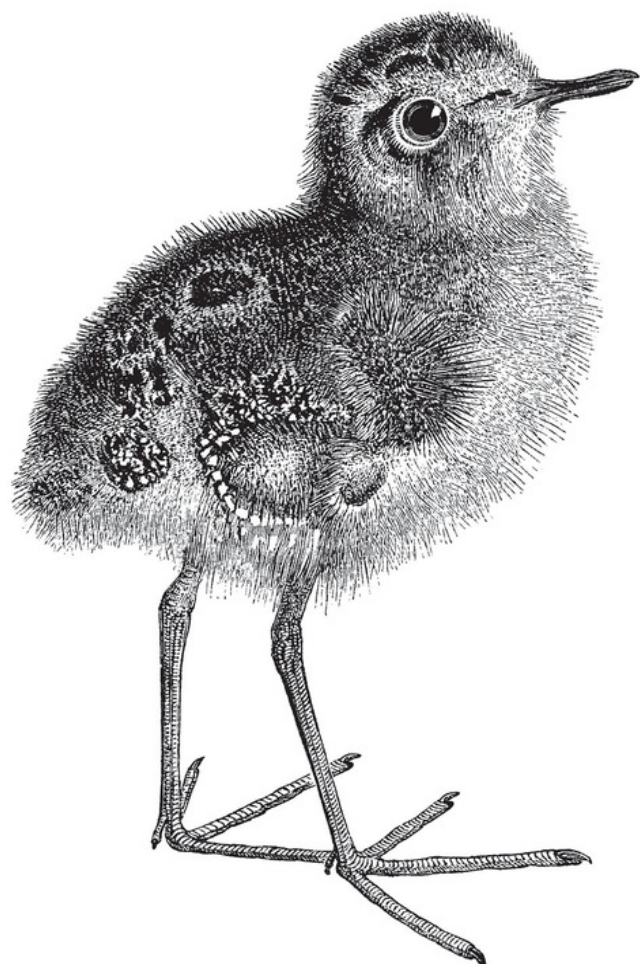


O'REILLY®

SQL Server Advanced Troubleshooting and Performance Tuning

Best Practices and Techniques



Early
Release

RAW &
UNEDITED

Dmitri Korotkevitch

www.datasense.ir

SQL Server Advanced Troubleshooting and Performance Tuning

Best Practices and Techniques

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

Dmitri Korotkevitch

SQL Server Advanced Troubleshooting and Performance Tuning

by Dmitri Korotkevitch

Copyright © 2021 Dmitri Korotkevitch. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://oreilly.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or corporate@oreilly.com.

Editors: Sarah Grey and Andy Kwan

Production Editor: Beth Kelly

Interior Designer: David Futato

Cover Designer: Karen Montgomery

Illustrator: Kate Dullea

May 2022: First Edition

Revision History for the Early Release

- 2020-12-18: First Release
- 2021-02-26: Second Release
- 2021-03-31: Third Release
- 2021-04-09: Fourth Release
- 2021-06-07: Fifth Release

- 2021-09-23: Sixth Release
- 2021-12-21: Seventh Release

See <http://oreilly.com/catalog/errata.csp?isbn=9781098101923> for release details.

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *SQL Server Advanced Troubleshooting and Performance Tuning*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

The views expressed in this work are those of the author, and do not represent the publisher's views. While the publisher and the author have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the author disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

978-1-098-10186-2

[LSI]

Chapter 1. SQL Server Setup and Configuration

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 1 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

Database servers never live in a vacuum. They belong to an ecosystem of one or more applications used by customers. Application databases are hosted on one or more instances of SQL Server, which, in turn, run on physical or virtual hardware. The data is stored on disks that are usually shared with other customers and database systems. Finally, all components use a network for communication and storage.

The complexity and internal dependencies of database ecosystems make troubleshooting a very challenging task. From the customers’ standpoint, most problems present themselves as general performance issues: applications might feel slow and unresponsive, database queries might time out, or applications might not connect to the database. The root cause of the issues could be anywhere. Hardware could be malfunctioning or incorrectly configured; the database might have inefficient schema, indexing, or code; SQL Server could be overloaded; client applications could have bugs or

design issues. This means you'll need to take a holistic view of your entire system in order to identify and fix problems.

This book is about troubleshooting SQL Server issues. However, we will always start this by analyzing your application ecosystem and SQL Server environment. This chapter will give you a set of guidelines on how to perform that validation and detect most common inefficiencies in SQL Server configuration. First, I'll discuss the hardware and operating system setup. Next, I'll talk about SQL Server and database configuration. I'll also touch on the topic of SQL Server consolidation and the overhead that monitoring can introduce into the system.

Hardware and Operating System Considerations

In most cases, troubleshooting and performance-tuning processes happen in production systems that host a lot of data and work under heavy loads. You have to deal with the issues and tune the live systems. Nevertheless, it is impossible to completely avoid discussion about hardware provisioning, especially because you may find that your servers cannot keep up with the load and need to be upgraded after troubleshooting.

I am not going to recommend particular vendors, parts or model numbers; computer hardware improves quickly and any such specific advice would be obsolete by the time the book is published. Instead, I'll focus on common-sense considerations with long-term relevance.

CPU

The license cost of a commercial database engine is, by far, the most expensive part in the system, and SQL Server is not an exception: you could build a decent server for less than the retail price of four cores in Enterprise Edition. You should buy the most powerful CPU your budget allows, especially if you are using non-Enterprise Editions, which limit the number of cores you can utilize.

Pay attention to CPU model. Each generation of CPUs will introduce performance improvements over the previous ones. You may get 10% to 15% performance improvements just by choosing newer CPUs, even when both generation of CPUs have the same clock speed.

In some cases, when licensing cost is not an issue, you may need to choose between slower CPUs with more cores and faster CPUs with fewer cores. In that case, the choice greatly depends on the system workload. In general, Online Transactional Processing (OLTP) systems, and especially In-Memory OLTP, would benefit from the higher single-core performance. A data warehouse and analytical workload, on the other hand, may run better with higher degree of parallelism and more cores.

Memory

There is a joke in the SQL Server community that goes like this:

Q. How much memory does SQL Server usually need?

A. More.

This joke has merits. SQL Server benefits from a large amount of memory, which allows it to cache more data. This, in turn, will reduce amount of disk input/output (I/O) activity and improve SQL Server's performance. Therefore, adding more memory to the server may be the cheapest and fastest way to address some performance issues.

For example, suppose the system suffers from non-optimized queries. You could reduce their impact by adding memory and eliminating the physical reads they introduce. This, obviously, does not solve the root cause of the problem. It is also dangerous, because as the data grows it eventually may not fit into the cache. However, in some cases it may be acceptable as a temporary Band-Aid solution.

The Enterprise Edition of SQL Server does not limit the amount of memory it can utilize. Non-Enterprise editions have limitations. In terms of memory utilization, Standard Edition of SQL Server 2016 and above can use up to 128GB of RAM for the buffer pool, 32GB of RAM per database for In-

Memory OLTP data, and 32GB of RAM for storing columnstore index segments. Web Edition memory usage is limited to half of what the Standard Edition provides. Factor those limits into your analysis when you are provisioning or upgrading non-Enterprise Edition instances of SQL Server. Don't forget to allocate some additional memory to other SQL Server components, for example plan cache and lock manager.

In the end, add as much memory as you can afford. It is cheap nowadays. There is no need to overallocate memory if your databases are small, but think about future data growth.

Disk Subsystem

A healthy, fast disk subsystem is essential for good SQL Server performance. SQL Server is very I/O intensive application - it is constantly reading from and writing data to disk.

There are many options for architecting the disk subsystem for SQL Server installations. The key is to build it in a way that provides low latency for I/O requests. For critical tier-1 systems, I recommend not exceeding 3 to 5 milliseconds (ms) of latency for data files reads and writes, and 1ms to 2ms of latency for transaction log writes. Fortunately, those numbers are now easily achieved with flash-based storage.

There's a catch, though: When you troubleshoot I/O performance in SQL Server, you need to analyze the latency metrics *within* SQL Server rather than on the storage level. It is common to see significantly higher numbers in SQL Server rather than in storage key performance indicators (KPIs), due to the queueing that may occur with I/O intensive workload. (Chapter 3 will discuss how to capture and analyze I/O performance data.)

If your storage subsystem provides multiple performance tiers, I recommend putting tempdb database on the fastest drive, followed by transaction log and data files. The tempdb is the shared resource on the server, and it is essential for this database to have good I/O throughput.

The writes to transaction log files are synchronous. It is critical to have low write latency for those files. The writes to transaction log are also sequential; however, remember that placing multiple log and/or data files to the same drive will lead to random I/O across multiple databases.

As a best practice, I'd put data and log files to the different physical drives for maintainability and recoverability reasons. You need to look at the underlying storage configuration though. In some cases, when disk arrays do not have enough spindles, splitting them across multiple LUNs may degrade disk array performance.

In my systems, I do not split clustered and nonclustered indexes across multiple filegroups by placing them on different drives. It rarely improves I/O performance unless you can completely separate storage path across the filegroups. On the other hand, this configuration can significantly complicate disaster recovery.

Finally, remember that some SQL Server technologies benefit from good sequential I/O performance. For example, In-Memory OLTP does not use random I/O at all, and performance of sequential reads usually becomes the limiting factor for database startup and recovery. Data warehouse scans would also benefit from sequential I/O performance when B-Tree and columnstore indexes are not heavily fragmented. The difference between sequential and random I/O performance is not very significant with flash-based storage; however, it may be a big factor with magnetic drives.

Network

SQL Server communicates with clients and other servers via the network. Obviously, it needs to provide enough bandwidth to support that communication. There are a couple items I want to mention, though.

First, you need to analyze entire network topology when you troubleshoot network-related performance. Remember that a network's throughput will be limited to the speed of its slowest component. You may have a 10 Gbps uplink from the server; however, if you have 1Gbps switch in network path, that would become the limiting factor. This is especially critical for

network-based storage: make sure that I/O path to disks is as efficient as possible.

Second, there is the common practice to build separate network for cluster heartbeat in AlwaysOn Failover Cluster and AlwaysOn Availability Groups. In some cases, people may also consider building separate network for all Availability Group traffic. This is the good approach that improves cluster reliability in simple configurations when all cluster nodes belong to the same subnet and may utilize Layer-2 routing. However, in complex, multi-subnet setup, multiple networks may lead to the routing issues. Be careful with that setup and make sure that networks are properly utilized in cross-node communication.

Virtualization adds another layer of complexity here. C

Consider a situation where you have a virtualized SQL Server cluster with nodes running on different hosts. You would need to check that the hosts can separate and route the traffic in the cluster network separately from the client traffic. Serving all vLan traffic through the single physical network card would defeat the purpose of a heartbeat network. (I will talk more about troubleshooting network-related issues in Chapter 13.)

Operating Systems and Applications

As a general rule, I suggest using the most recent version of your operating system that supports your version of SQL Server. Make sure that both the OS and SQL Server are patched, and implement a process to do patching regularly.

If you are using old version of SQL Server prior 2016, use 64-bit variant. In the most cases, the 64-bit version outperforms 32-bit version and scales better with the hardware.

Since SQL Server 2017, it's been possible to use Linux to host the database server. From a performance standpoint, Windows and Linux versions of SQL Server are very similar. The choice of operating system depends on enterprise ecosystem and on what your team is more comfortable to

support. Keep in mind, that Linux-based deployments may require a slightly different High Availability (HA) strategy compared to a Windows setup. For example, you may have to rely on Pacemaker instead of Windows Server Failover Cluster (WSFC) for automatic failovers.

Use a dedicated SQL Server host whenever possible. Remember that it's easier and cheaper to scale application servers—don't waste valuable resources on the database host!

On the same note, do not run nonessential processes on the server. I see database engineers running SQL Server Management Studio (SSMS) in remote desktop sessions all the time. It is always better to work remotely and not consume server resources.

Finally, if you are required to run antivirus software on the server, exclude any database folders from the scan.

Virtualization and Clouds

Modern IT infrastructure depends heavily on virtualization, which provides additional flexibility, simplifies management, and reduces hardware costs. As a result, more often than not, you'll have to work with virtualized SQL Server infrastructure.

There is nothing wrong with that. Properly implemented virtualization gives you many benefits, with negligible performance overhead. It adds another layer of High Availability with VMware vSphere vMotion or Hyper-V Live Migration. It allows you to seamlessly upgrade the hardware and simplifies database management. Unless you have the edge case when you need to squeeze the most from the hardware, I suggest virtualizing your SQL Server ecosystem.

NOTE

The overhead from virtualization increases on the large servers with many CPUs. However, it still may be acceptable in many cases.

Virtualization, however, adds another layer of complexity during troubleshooting. You need to pay attention to the host's health and load in addition to guest virtual machine (VM) metrics. To make matters worse, the performance impact of an overloaded host might not be clearly visible in standard performance metrics in guest OS.

I will discuss several approaches to troubleshooting the virtualization layer in Chapter 15 however, you can start by working with infrastructure engineers to confirm that the host is not overprovisioned. Pay attention to the number of physical CPUs and allocated vCPUs on the host along with physical and allocated memory. Mission-critical SQL Server VMs should have resources reserved for them to avoid performance impact.

Asides from the virtualization layer, troubleshooting virtualized SQL Server instances is the same as troubleshooting physical ones. The same applies to cloud installations when SQL Server is running within virtual machines. After all, the cloud is just a different datacenter managed by an external provider.

Configuring Your SQL Server

The SQL Server setup process's default configuration is relatively decent and may be suited to light and even moderate workloads. There are several things you need to validate and tune, however.

SQL Server Version and Patching Level

`SELECT @@VERSION` is the first statement I run during SQL Server system health checks. There are two reasons for that. First, it gives me a glimpse of the system's patching strategy, so I can potentially suggest some improvements. Second, it helps me to identify possible known issues that may exist in the system.

The latter one is very important. Many times, customers have asked me to troubleshoot problems that had already been resolved by service packs and

cumulative updates. Always look at the release notes to see if any of the issues mentioned look familiar; your problem may have already been fixed.

You might consider upgrading to the newest version of SQL Server when possible. Each version introduces performance, functional and scalability enhancements. This is especially true if you move to SQL Server 2016 or above from older versions. SQL Server 2016 was a milestone release that included many performance enhancements. In my personal experience, upgrading from SQL Server 2012 to 2016 and above can improve performance by 20 to 40% without any additional steps.

It is also worth noting that starting with SQL Server 2016 SP1, many former Enterprise Edition-only features became available in the lower-end editions of the product. Some of them, like data compression, allow SQL Server to cache more data in the buffer pool and improve performance of the system.

Obviously, you need to test the system prior to upgrading – there is always the chance of regressions. The risk is usually small with minor patching; however, it increases with the major upgrades. You can mitigate some risks with several database options, as you will see later in this chapter.

Instant File Initialization

Every time SQL Server grows data and transaction log files—either automatically or as part of ALTER DATABASE command—it fills the newly allocated part of the file with zeros. This process blocks all sessions that are trying to write to the corresponding file and, in case of transaction log, stops generating any log records. It may also generate the spike in I/O write workload.

That behavior cannot be changed for transaction log files – SQL Server always zeros them out. However, you can disable it for the data files by enabling *instant file initialization (IFI)*. This speeds up data file growth and reduces the time required to create or restore databases.

You can enable instant file initialization by giving an SA_MANAGE_VOLUME_NAME permission, also known as *Perform Volume Maintenance Task*, to the SQL Server startup account. This can be done in the *Local Security Policy* management application (secpol.msc). You will need to restart SQL Server for the change to take effect.

In SQL Server 2016 and above, you can also grant this permission as part of the SQL Server setup process (shown in **Figure 1-1**).

Server Configuration

Specify the service accounts and collation configuration.

Global Rules

Microsoft Update

Product Updates

Install Setup Files

Install Rules

Installation Type

Product Key

License Terms

Feature Selection

Feature Rules

Instance Configuration

Server Configuration

Database Engine Configuration

Feature Configuration Rules

Ready to Install

Installation Progress

Complete

Service Accounts Collation

Microsoft recommends that you use a separate account for each SQL Server service.

Service	Account Name	Password	Startup Type
SQL Server Agent	NT Service\SQLAgentSS...		Manual ▾
SQL Server Database Engine	NT Service\MSSQLSSQL...		Automatic ▾
SQL Server Browser	NT AUTHORITY\LOCAL...		Automatic ▾

Grant Perform Volume Maintenance Task privilege to SQL Server Database Engine Service

This privilege enables instant file initialization by avoiding zeroing of data pages. This may lead to information disclosure by allowing deleted content to be accessed.

[Click here for details](#)

< Back

Next >

Cancel

Figure 1-1. Enabling Instant File Initialization during SQL Server setup.

You can check if IFI is enabled by examining the `instant_file_initialization_enabled` column in the `sys.dm_server_services` data management view. This column is available in SQL Server 2012 SP4, SQL Server 2016 SP1, and above. In older versions, you can run the code shown in Listing 1-1.

Example 1-1. Checking if instant file initialization is enabled in old SQL Server versions

```
DBCC TRACEON(3004,3605,-1);
go
CREATE DATABASE Dummy;
go
EXEC sp_readerrorlog 0,1,N'Dummy';
go
DROP DATABASE Dummy;
go
DBCC TRACEOFF(3004,3605,-1);
go
```

If IFI is not enabled, the SQL Server log will indicate that SQL Server is zeroing out the mdf data file in addition to zeroing out the log ldf file, as shown in **Figure 1-2**. When IFI is enabled, it will only show zeroing out of the log ldf file.

	LogDate	ProcessInfo	Text
104	2020-10-26 15:57:45.370	spid32s	A connection timeout has occurred while attempting to establish a connection to availa..
105	2020-10-26 15:58:35.510	spid51	DBCC TRACEON 3004, server process ID (SPID) 51. This is an informational message only;..
106	2020-10-26 15:58:35.510	spid51	DBCC TRACEON 3605, server process ID (SPID) 51. This is an informational message only;..
107	2020-10-26 15:58:35.520	spid51	Zeroing C:\DB\Dummy.mdf from page 0 to 1024 (0x0 to 0x8000000)
108	2020-10-26 15:58:35.530	spid51	Zeroing completed on C:\DB\Dummy.mdf (elapsed = 2 ms)
109	2020-10-26 15:58:35.530	spid51	Zeroing C:\DB\Dummy_log.ldf from page 0 to 1024 (0x0 to 0x8000000)
110	2020-10-26 15:58:35.540	spid51	Zeroing completed on C:\DB\Dummy_log.ldf (elapsed = 2 ms)
111	2020-10-26 15:58:35.550	spid51	Starting up database 'Dummy'.

Figure 1-2. Checking if instant file initialization is enabled.

There is a small security risk associated with this setting. When IFI is enabled, the database administrators may see some data from previously deleted files in OS by looking at newly allocated data pages in the database. This is acceptable in most systems; if so, enable it.

Tempdb Configuration

Tempdb is the system database used to store temporary objects created by users and by SQL Server internally. This is a very active database and it often becomes a source of contention in the system. I will discuss how to troubleshoot tempdb-related issues in Chapter 9; in this chapter, I'll focus on configuration.

As already mentioned, you need to place tempdb on the fastest drive in the system. Generally speaking, this drive does not need to be redundant nor

persistent – the database is recreated at SQL Server startup, and local SSD disk or cloud ephemeral storage would work fine. Remember, however, that SQL Server will go down if tempdb is unavailable, so factor that into your design.

If you are using non-Enterprise Editions of SQL Server and the server has more memory than SQL Server can consume, you can put tempdb on the RAM drive. Don't do that with Enterprise Edition, though – you'll usually achieve better performance by using that memory for the buffer pool.

NOTE

Pre-allocate tempdb files to the maximum size of RAM drive and create additional small data and log files on disk to avoid running out of space. SQL Server will not use small on-disk files until RAM drive files are full.

The tempdb database should always have multiple data files. Unfortunately, default configuration created by SQL Server setup is not optimal, especially in the old versions of the product. We will discuss how to fine-tune the number of data files in tempdb later in the book, but you can use the following as the rule of thumb in initial configuration:

- If the server has 8 or fewer CPU cores, create the same number of data files.
- If the server has more than 8 CPU cores, use either 8 data files or 1/4 of the number of cores, whichever is greater, rounding up in batches of 4 files. For example, use 8 data files in the 24-core server and 12 data files in the 40-core server.

Finally, make sure that all tempdb data files have the same initial size and auto-growth parameters specified in megabytes (MB) rather than in percentages. This will allow SQL Server to better balance the usage of the data files, reducing possible contention in the system.

Trace Flags

SQL Server uses trace flags to enable or change the behavior of some product features. Although Microsoft has introduced more and more database and server configuration options in new versions of SQL Server, trace flags are still widely used. You will need to check any trace flags that are present in the system; you may also need to enable some of them.

You can get the list of enabled trace flags by running the DBCC TRACESTATUS command. You can enable them in SQL Server Configuration Manager and/or by using -T SQL Server startup option.

Let's look at some common trace flags.

T1118

This trace flag prevents usage of **mixed extents** in SQL Server. This will help to improve tempdb throughput in SQL Server 2014 and below by reducing the amount of changes and, therefore, contention in tempdb system catalogs. This trace flag is not required in SQL Server 2016 and above, where tempdb does not use mixed extents by default.

T1117

With this trace flag, SQL Server auto-grows all data files in the filegroup when one of the files is out of space. It provides more balanced I/O distribution across data files. You should enable it to improve tempdb throughput in old versions of SQL Server; however, check if any users' databases have filegroups with multiple unevenly sized data files. As with T1118, this trace flag is not required in SQL Server 2016 and above, where tempdb auto-grows all data files by default.

T2371

By default, SQL Server automatically updates statistics only after 20% of the data in the index has been changed. This means that with large tables, statistics are rarely updated automatically. The trace flag T2371 changes this behavior, making the statistics update threshold dynamic – the larger the table is, the lower the percentage of changes required to trigger the update. Starting with SQL Server 2016, you can also control this behavior via database compatibility level. I recommend enabling this trace flag unless all databases on the server have a compatibility level of 130 or above.

T3226

With this trace flag, SQL Server does not write information about successful database backups to the error log. This may help to reduce the size of the logs, making them more manageable.

T1222

This trace flag writes deadlock graphs to the SQL Server error log. This flag is benign; however, it makes SQL Server log harder to read and parse. It is also redundant – you can get deadlock graph from System_Health Extended Event session when needed. I usually remove this trace flag when I see it.

T4199

This trace flag and the QUERY_OPTIMIZER_HOTFIXES database option (in SQL Server 2016 and above) control the behavior of Query Optimizer hotfixes. When enabled, the hotfixes introduced in service packs and cumulative updates will be used. This may help to address

some of Query Optimizer bugs and improve query performance; however, it also increases the chance of plan regressions after patching. I usually do not enable it in production systems unless it is possible to perform thorough regression testing of the system before patching.

T7412

This trace flag enables lightweight execution profiling infrastructure in SQL Server 2016 and 2017. This allows you to collect execution plans and many execution metrics for the queries in the system with little CPU overhead. I am going to discuss it in more details in Chapter 5.

To summarize – in SQL Server 2014 and below, enable T1118, T2371 and, potentially, T1117. In SQL Server 2016 and above, enable T2371 unless all databases have compatibility level of 130 or above. After that – look at all other trace flags in the system and understand what they are doing. Some trace flags may be inadvertently installed by third-party tools and can negatively affect server performance.

Server Options

SQL Server provides many configuration settings. I'll cover many of them in depth later in the book; however, there are a few settings worth mentioning here.

Optimize for Ad-hoc Workloads

The first one is *Optimize for Ad-hoc Workloads*. This configuration option controls how SQL Server caches execution plans of ad-hoc (non-parameterized) queries. When it is disabled (by default), SQL Server caches full execution plans of those statements, which may significantly increase plan cache memory usage. As the opposite, when this setting is enabled, SQL Server starts by caching the small structure (just a few hundred bytes)

called *plan stub*, replacing it with the full execution plan if an ad-hoc query is executed the second time.

In majority of the cases, ad-hoc statements are not executed repeatedly, and it is beneficial to enable *Optimize for Ad-hoc Workloads* setting in every system. It could significantly reduce plan cache memory usage at cost of infrequent additional recompilations of ad-hoc queries. Obviously, this setting would not affect caching behavior of parameterized queries and T-SQL database code.

NOTE

Starting with SQL Server 2019 and in Azure SQL Database, you can control *Optimize for Ad-hoc Workload* behavior on the database level with the `OPTIMIZE_FOR_AD_HOC_WORKLOADS` database scoped configuration.

Max Server Memory

The second important setting is *Max Server Memory*, which controls how much memory SQL Server can consume. Database engineers love to debate how to properly configure it, and there are different approaches to calculating the proper value for the setting. Many engineers even suggest leaving the default value in place and allowing SQL Server to manage it automatically. In my opinion, it is best to fine-tune that setting, but it's important to do so correctly (Chapter 7 will discuss the details). An incorrect setting will impact SQL Server performance more than if you leave the default value in place.

One particular issue I often encounter during system health checks is severe underprovisioning of this setting. Sometimes people forget to change it after hardware or VM upgrades; other times, it's incorrectly calculated in nondedicated environments, where SQL Server is sharing the server with other applications. In both cases, you can get immediate improvements by increasing *Max Server Memory* or even setting it to the default value until you perform full analysis later.

Affinity Mask

You need to check SQL Server affinity and, potentially, set *affinity mask* if SQL Server is running on hardware with multiple non-uniform memory access (NUMA) nodes. In modern hardware, each physical CPU usually becomes a separate NUMA node. If you restrict SQL Server from using some of the physical cores, you need to balance SQL Server CPUs (or schedulers – see chapter 2) evenly across NUMAs. For example, if you are running SQL Server on a server with two 18-core Xeon processors and limiting SQL Server to 24 cores, you need to set affinity mask to utilize 12 cores from each physical CPU. This will give you better performance than having SQL Server use 18 cores from the first processor and 6 cores from the second.

Listing 1-2 shows how to analyze SQL Server schedulers (CPUs) distribution between NUMA nodes. Look at the count of schedulers for each `parent_node_id` column in the output.

Example 1-2. Checking the distribution of NUMA node schedulers (CPUs)

```
SELECT
    parent_node_id
    ,COUNT(*) as [Schedulers]
    ,SUM(current_tasks_count) as [Current]
    ,SUM(runnable_tasks_count) as [Runnable]
FROM sys.dm_os_schedulers
WHERE status = 'VISIBLE ONLINE'
GROUP BY parent_node_id;
```

Parallelism

It is important to check parallelism settings in the system. Default settings, like `MAXDOP = 0` and `Cost Threshold for Parallelism = 5`, do not work well in modern systems. As with Max Server Memory, it is better to fine-tune the settings based on the system workload (Chapter 6 will discuss this in detail). However, my rule of thumb for generic settings is:

- Set `MAXDOP` to 1/4 of the number of available CPUs in OLTP and half those in Data Warehouse systems. Do not exceed the number of CPUs in the NUMA node.

- Set Cost Threshold for Parallelism to 50.

Starting with SQL Server 2016 and in Azure SQL Server Database, you can set MAXDOP on the database level using the command ALTER DATABASE SCOPED CONFIGURATION SET MAXDOP. This is useful when the instance hosts databases that handle different workloads.

Configuration Settings

As with trace flags, analyze other changes in configuration settings that have been applied on the server. You can examine current configuration options using the sys.configurations view. Unfortunately, SQL Server does not provide a list of default configuration values to compare, so you need to hardcode it, as shown in Listing 1-3. I am including just a few configuration settings to save space, but you can download the full version of the script from this book's companion website.

Example 1-3. Detecting changes in server configuration settings

```
DECLARE
    @defaults TABLE
    (
        name SYSNAME NOT NULL PRIMARY KEY,
        def_value SQL_VARIANT NOT NULL
    )
INSERT INTO @defaults(name,def_value) VALUES('backup compression
default',0);
INSERT INTO @defaults(name,def_value) VALUES('cost threshold for
parallelism',5);
INSERT INTO @defaults(name,def_value) VALUES('max degree of
parallelism',0);
INSERT INTO @defaults(name,def_value) VALUES('max server memory
(MB)',2147483647);
INSERT INTO @defaults(name,def_value) VALUES('optimize for ad hoc
workloads',0);
/* Other settings are omitted in the book */
SELECT
    c.name, c.description, c.value_in_use, c.value
    ,d.def_value, c.is_dynamic, c.is_advanced
FROM
    sys.configurations c JOIN @defaults d ON
        c.name = d.name
WHERE
    c.value_in_use <> d.def_value OR
```

```

    c.value <> d.def_value
ORDER BY
    c.name;

```

Figure 1-3 shows the sample output of the code. The discrepancy between value and value_in_use columns indicates pending configuration changes that require restart to take an effect. The is_dynamic column shows if configuration option can be modified without restart.

	name	description	value_in_use	value	def_value	is_dynamic	is_advanced
1	max degree of parallelism	maximum degree of paralle...	1	1	0	1	1
2	optimize for ad hoc workloads	When this option is set, ...	1	1	0	1	1

Figure 1-3. Non-default server configuration options.

Configuring Your Databases

As the next step in analyzing your configuration, you'll need to validate several database settings and configuration options. Let's look at them.

Database Settings

SQL Server allows you to change multiple database settings, tuning behavior to workload and other requirements. I'll cover many of them later in the book; however, there are a few settings I would like to discuss here.

The first one is *Auto Shrink*. When this option is enabled, SQL Server periodically shrinks the database and releases unused free space from the files to the OS. While this looks appealing and promises to reduce disk space utilization, it may introduce issues.

Implementing this database shrink process, automatically or through the command DBCC SHRINKFILE, works on the physical level. It locates empty space in the beginning of the file and moves allocated extents from the end of the file there, without taking extent ownership into consideration.

This introduces noticeable load and lead to the serious index fragmentation. What's more, in many cases it's useless: the database files simply expand again as the data grows. It's always better to manage file space manually and disable Auto Shrink.

Another database option, *Auto Close*, controls how SQL Server caches data from the database. When it's enabled, SQL Server removes data pages from the buffer pool and execution plans from the plan cache when the database does not have any active connections. This will lead to performance impact with the new workload when data needs to be cached and queries need to be compiled again.

With very few exceptions, you should disable Auto Close. One such exception may be an instance that hosts a large number of rarely accessed databases. Even then, I would consider keeping this option disabled and allowing SQL Server to retire cached data in the normal way.

Make sure that *Page Verify* option is set to CHECKSUM. This will detect consistency errors more efficiently and helps to resolve database corruption cases.

Pay attention to the *database recovery model*. If the databases are in SIMPLE mode, in case of disaster or human error it would be impossible to recover the data beyond the last FULL database backup. If you find the database in this mode, immediately discuss it with the stakeholders, making sure that they understand the risk of data loss.

Database Compatibility Level controls SQL Server's compatibility and behavior on the database level. For example, if you are running SQL Server 2019 and have a database with a compatibility level of 130 (SQL Server 2016), SQL Server will behave as if the database is running on SQL Server 2016. Keeping the databases on the lower compatibility levels simplifies SQL Server upgrades by reducing possible regressions; however, it also blocks you from getting some new features and enhancements.

As a general rule, run databases on the latest compatibility level that matches the SQL Server version. Be careful when you change it: as with any version change, this may lead to regressions. Test the system before the

change and make sure you can roll back the change if needed, especially if the database has a compatibility level of 110 (SQL Server 2012) or below. Increasing compatibility level to 120 (SQL Server 2014) or above will enable a new *cardinality estimation model* and may significantly change execution plans for the queries. Test the system thoroughly to understand the impact of the change.

You can force SQL Server to use legacy cardinality estimation models with the new database compatibility levels by setting LEGACY_CARDINALITY_ESTIMATION database option to ON in SQL Server 2016 and above, or by enabling server-level trace flag T9481 in SQL Server 2014. This approach will allow you to perform upgrade or compatibility level changes in phases, reducing impact to the system. (Chapter 5 will cover cardinality estimation in more detail.)

Transaction Log Settings

SQL Server uses *write-ahead logging*, persisting information about all database changes in a transaction log. SQL Server works with transaction logs sequentially, in merry-go-round fashion. In most cases, you won't need multiple log files in the system – they make database administration more complicated and do not improve performance.

Internally, SQL Server splits transaction logs into chunks called *Virtual Log Files (VLF)* and manages them as single units. For example, SQL Server cannot truncate and reuse a VLF if it contains just a single active log record. Pay attention to the number of VLFs in the database. Too few of them will lead to very large VLFs, which make log management and truncation suboptimal. Too many small VLFs will degrade the performance of transaction log operations. Try not to exceed several hundred VLFs in production systems.

The number of VLFs SQL Server adds when it grows a log depends on SQL Server version and the size of the grows. In most cases, it creates 8 VLFs when the growth size is between 64MB and 1GB or 16 VLFs with above 1GB growth. Do not use percent-based auto-growth configuration

because it generates lots of unevenly sized VLFs. Change the log auto-growth setting to grow the file in chunks – I usually use chunks of 1,024 MB, which generates 128MB VLFs unless I need very large transaction log.

You can count the VLFs in the database with `sys.dm_db_log_info` data management **view** in SQL Server 2016 and above. In older versions of SQL Server, you can obtain the information by running `DBCC LOGINFO`. If the transaction log isn't configured well, consider rebuilding it. You can do this by shrinking the log to the minimal size and growing it in chunks of 1,024MB to 4,096 MB.

Do not auto-shrink transaction log files. They will grow again and affect performance when SQL Server zeroes out the file. It is better to pre-allocate the space and manage log file size manually. Do not restrict the maximum size and auto-growth, though – you want logs to grow automatically in case of emergencies. (Chapter 11 will provide more details on how to troubleshoot transaction-log issues.)

Data Files and Filegroups

By default, SQL Server creates new databases using the single-file PRIMARY filegroup and one transaction log file. Unfortunately, this configuration is suboptimal from performance, database management and High Availability standpoints.

SQL Server tracks space usage in the data files through system pages called *allocation maps*. In systems with highly volatile data, allocation maps can be a source of contention: SQL Server serializes access to them during their modifications (more about this in Chapter 10). Each data file has its own set of allocation map pages and you can reduce contention by creating multiple files in the filegroup with the active modifiable data.

Ensure that data is evenly distributed across multiple data files in the same filegroup. SQL Server uses an algorithm called *Proportional Fill*, which writes most data to the file that has the most free space. Evenly sized data files will help to balance those writes, reducing allocation maps contention.

Make sure that all data files in the filegroup have the same size and auto-growth parameters, specified in MBs.

You may also want to enable the `AUTOGROW_ALL_FILES` filegroup option (available in SQL Server 2016 and above), which triggers auto-growth for all files in the filegroup simultaneously. You can use trace flag T1117 for this in prior versions of SQL Server, but remember that this flag is set on the server level and will affect all databases and filegroups in the system.

It is often impractical or impossible to change the layout of existing databases. However, you may need to create new filegroups and move data around during performance tuning. Here are a few suggestions for doing so efficiently:

- Create multiple data files in filegroups with volatile data. I usually start with four files and increase the number if I see latching issues (see Chapter 10). Make sure that all data files have the same size and auto-growth parameters specified in MB; enable the `AUTOGROW_ALL_FILES` option. For filegroups with read-only data, one data file is usually enough.
- Do not spread clustered indexes, and nonclustered indexes, or large object (LOB) data across multiple filegroups. This rarely helps with performance and may introduce issues in cases of database corruption.
- Place related entities (for example, Orders and OrderLineItems) in the same filegroup. This will simplify database management and disaster recovery.
- Keep the PRIMARY filegroup empty if possible.

Figure 1-4 shows an example of a database layout for a hypothetical shopping-cart system. The data is partitioned and spread across multiple filegroups with the goal of minimizing downtime and utilizing *partial database availability* in case of disaster.¹ It will also improve your backup

strategy by implementing partial database backups and excluding read-only data from FULL backups.

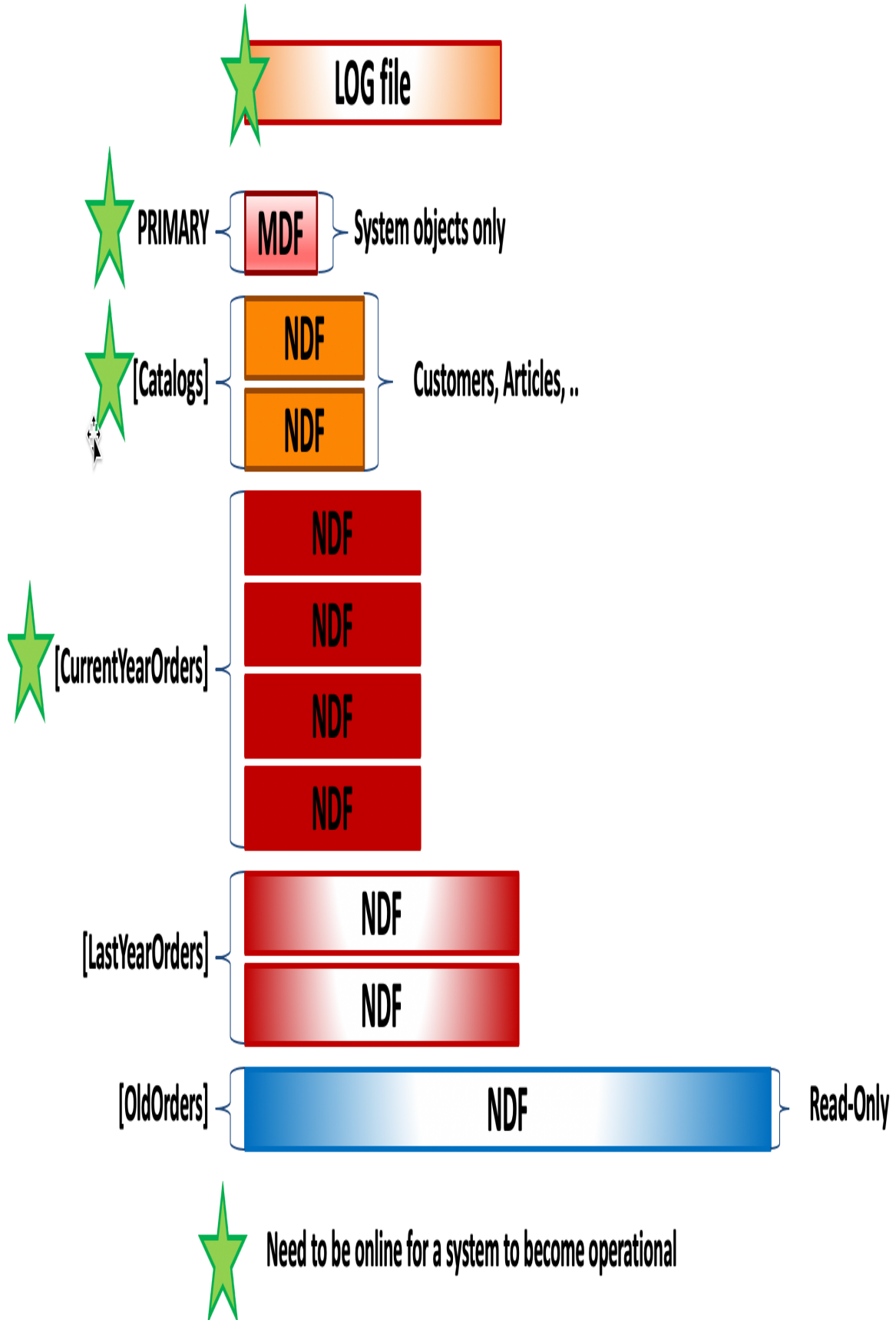


Figure 1-4. Database layout for a shopping cart system.

Analyzing SQL Server Error Log

SQL Server Error Log is another place I usually check at the beginning of troubleshooting. I like to see any errors it has, which can point to some areas to follow up. For example, errors 823 and 824 can indicate issues with disk subsystem and/or database corruption.

You can read the content of the log in SSMS. You can also get it programmatically using the `xp_readerrorlog` system stored procedure. The challenge here is the amount of data in the log: the noise from the information messages may hide useful data.

The code in Listing 1-4 helps you to address that problem. It allows you to filter out unnecessary noise and focus on the error messages. You can control the behavior of the code with the following variables:

`@StartDate` and `@EndDate`

- Define the time for analysis: `@NumErrorLogs`
- Specifies the number of log files to read if SQL Server rolls them over: `@ExcludeLogonErrors`
- Omits logon auditing messages: `@ShowSurroundingEvents` and `@ExcludeLogonSurroundingEvents`

These allow you to retrieve the information messages around the error entries from the log. The time window for those messages is controlled by the `@SurroundingEventsBeforeSeconds` and `@SurroundingEventsAfterSeconds` variables.

The script produces two outputs. The first one shows the entries from the error log that include word *error*. When `@ShowSurroundingEvents` parameter is enabled, it would also provide log entries around those error lines. You can also exclude some of log entries that contain the word *error* from the output by inserting them to `@ErrorsToIgnore` table.

Example 1-4. Analyzing SQL Server Error Log

```
IF OBJECT_ID('tempdb..#Logs',N'U') IS NOT NULL DROP TABLE #Logs;
IF OBJECT_ID('tempdb..#Errors',N'U') IS NOT NULL DROP TABLE
#Errors;
go
CREATE TABLE #Errors
(
    LogNum INT NULL,
    LogDate DATETIME NULL,
    ID INT NOT NULL identity(1,1),
    ProcessInfo VARCHAR(50) NULL,
    [Text] VARCHAR(MAX) NULL,
    PRIMARY KEY(ID)
);
CREATE TABLE #Logs
(
    [LogDate] DATETIME NULL,
    ProcessInfo VARCHAR(50) NULL,
    [Text] VARCHAR(max) NULL
);
DECLARE
    @StartDate DATETIME = DATEADD(DAY,-7,GETDATE())
    ,@EndDate DATETIME = GETDATE()
    ,@NumErrorLogs INT = 1
    ,@ExcludeLogonErrors BIT = 1
    ,@ShowSurroundingEvents BIT = 1
    ,@ExcludeLogonSurroundingEvents BIT = 1
    ,@SurroundingEventsBeforeSecond INT = 5
    ,@SurroundingEventsAfterSecond INT = 5    ,@LogNum INT = 0;

DECLARE
    @ErrorsToIgnore TABLE
    (
        ErrorText NVARCHAR(1024) NOT NULL
    );

INSERT INTO @ErrorsToIgnore(ErrorText)
VALUES
    (N'Registry startup parameters:%'),
    (N'Logging SQL Server messages in file%'),
    (N'CHECKDB for database%finished without errors%');

WHILE (@LogNum <= @NumErrorLogs)
BEGIN
    INSERT INTO #Errors(LogDate,ProcessInfo,Text)
    EXEC [master].[dbo].[xp_readerrorlog]
        @LogNum, 1, N'error', NULL, @StartDate, @EndDate, N'desc';
```

```

    IF @@ROWCOUNT > 0
        UPDATE #Errors SET LogNum = @LogNum WHERE LogNum IS NULL;
    SET @LogNum += 1;
END;

IF @ExcludeLogonErrors = 1
    DELETE FROM #Errors WHERE ProcessInfo = 'Logon';

DELETE FROM e
FROM #Errors e
WHERE EXISTS
(
    SELECT *
    FROM @ErrorsToIgnore i
    WHERE e.Text LIKE i.ErrorText
);

-- Errors only
SELECT * FROM #Errors ORDER BY LogDate DESC;

IF @@ROWCOUNT > 0 AND @ShowSurroundingEvents = 1
BEGIN
    DECLARE
        @LogDate DATETIME
        ,@ID INT = 0

    WHILE 1 = 1
    BEGIN
        SELECT TOP 1 @LogNum = LogNum, @LogDate = LogDate, @ID = ID
        FROM #Errors
        WHERE ID > @ID
        ORDER BY ID;

        IF @@ROWCOUNT = 0
            BREAK;

        SELECT
            @StartDate = DATEADD(SECOND, -@SurroundingEventsBeforeSecond,
@LogDate)
            ,@EndDate = DATEADD(SECOND, @SurroundingEventsAfterSecond,
@LogDate);

        INSERT INTO #Logs (LogDate, ProcessInfo, Text)
        EXEC [master].[dbo].[xp_readerrorlog]
            @LogNum, 1, NULL, NULL, @StartDate, @EndDate;
    END;

    IF @ExcludeLogonSurroundingEvents = 1

```

```

DELETE FROM #Logs WHERE ProcessInfo = 'Logon';

DELETE FROM e
FROM #Logs e
WHERE EXISTS
(
    SELECT *
    FROM @ErrorsToIgnore i
    WHERE e.Text LIKE i.ErrorText
);

SELECT * FROM #Logs ORDER BY LogDate DESC;
END

```

I am not going to put the full list of possible errors here – it may be excessive and, in many cases, is system specific. But you need to analyze any suspicious data from the output and understand its possible impact on the system.

Finally, I suggest setting up alerts for high-severity errors in SQL Server Agent, if this has not already been done. You can read [Microsoft documentation](#) on how to do that.

Consolidating Instances and Databases

You can't talk about SQL Server troubleshooting without discussing database and SQL Server instances consolidation. While consolidating often reduces hardware and licensing costs, it doesn't come for free; you need to analyze its possible negative impact on the current or future system performance.

There is no universal consolidation strategy that can be used with every project. You should analyze the amount of data, load, hardware configuration, and your business and security requirements when making this decision. However, as a general rule, avoid consolidating OLTP and Data Warehouse/Reporting databases on the same server when they are working under a heavy load (or, if they are consolidated, consider splitting them). Data warehouse queries usually process large amounts of data,

which leads to heavy I/O activity and flushes the content of the buffer pool. Taken together, this negatively affects the performance of other systems.

In addition, analyze your security requirements when consolidating databases. Some security features, such as Audit, affect the entire server and add performance overhead for all databases on the server. *Transparent Data Encryption (TDE)* is another example: even though it is a database-level feature, SQL Server encrypts tempdb when either of the databases on the server has TDE enabled. This leads to performance overhead for all other systems. As a general rule, do not keep databases with different security requirements on the same instance of SQL Server. Look at the trends and spikes in metrics and separate databases from each other when needed. (I will provide code to help you analyze CPU, I/O and Memory usage on a per-database basis later in the book.)

I suggest utilizing virtualization and consolidating multiple VMs on one or a few hosts, instead of putting multiple independent and active databases on a single SQL Server instance. This will give you much better flexibility, manageability, and isolation between the systems, especially if multiple SQL Server instances are running on the same server. It is much easier to manage their resource consumption when you virtualize them.

Observer Effect

The production deployment of every serious SQL Server system requires implementing a monitoring strategy. This may include third-party monitoring tools, code built based on standard SQL Server technologies, or both.

A good monitoring strategy is essential for SQL Server production support. It helps you to be more proactive and reduces incident detection and recovery times. Unfortunately, it does not come for free—every type of monitoring adds the overhead to the system. In some cases, this overhead may be negligible and acceptable; in others it may significantly affect server performance.

During my career as an SQL Server consultant, I've seen many cases of inefficient monitoring. For example, one client was using a tool that provided information about index fragmentation by calling the `sys.dm_db_index_physical_stats` function, in DETAILED mode, every four hours for every index in the database. This introduced huge spikes in I/O and cleared the buffer pool, leading to a noticeable performance hit. Another client used a tool that constantly polled various DMVs, adding significant CPU load to the server.

Fortunately, in many cases, you will be able to see those queries and evaluate their impact during system troubleshooting. This is not always the case, however, with other technologies, for example with monitoring based on Extended Events. (I will talk about methods for detecting inefficient queries in Chapter 4). Extended Events is a great technology that allows you to troubleshoot complex problems in SQL Server. It is not, however, the best choice as a profiling tool. Some events are heavy and may introduce large overhead in busy environments.

Let's look at the example and create an xEvents session that captures queries running in the system, as shown in Listing 1-5.

Example 1-5. Creating an xEvents session to capture queries in the system

```
CREATE EVENT SESSION CaptureQueries ON SERVER
ADD EVENT sqlserver.rpc_completed
(
    SET collect_statement=(1)
    ACTION
    (
        sqlserver.task_time,sqlserver.client_app_name
        ,sqlserver.client_hostname
        ,sqlserver.database_name
        ,sqlserver.nt_username
        ,sqlserver.sql_text
    )
),
ADD EVENT sqlserver.sql_batch_completed
(
    ACTION
    (
        sqlserver.task_time
        ,sqlserver.client_app_name
```

```

        ,sqlserver.client_hostname
        ,sqlserver.database_name
        ,sqlserver.nt_username
        ,sqlserver.sql_text
    )
),
ADD EVENT sqlserver.sql_statement_completed
ADD TARGET package0.event_file
(SET FILENAME=N'C:\PerfLogs\LongSql.xel',MAX_FILE_SIZE=(200))
WITH
(
    MAX_MEMORY =4096 KB
    ,EVENT_RETENTION_MODE=ALLOW_SINGLE_EVENT_LOSS
    ,MAX_DISPATCH_LATENCY=5 SECONDS
);

```

Next, deploy it to the server that operates under a heavy load with a large number of concurrent requests. Measure the throughput in the system, with and without xEvents session running. Obviously, be careful—and don't run it on the production server!

Figure 1-5 illustrates CPU load and number of batch requests per second in both scenarios on one of my servers. As you can see, enabling xEvents session decreased throughput by more than 20%. To make matters worse, it would be very hard to detect the existence of that session on the server.



Figure 1-5. Server throughput with and without an active xEvents session.

Obviously, the degree of impact would depend on the system's workload. In either case, check for any unnecessary monitoring or data-collection tools when you do the troubleshooting.

The bottom line: Evaluate the monitoring strategy and estimate its overhead as part of your analysis, especially when the server hosts multiple databases. For example, Extended Events work at the server level. While you can filter the events based on `database_id` field, the filtering occurs after an event has been fired. This can affect all databases on the server.

Summary

System troubleshooting is a holistic process that requires you to analyze your entire application ecosystem. You need to look at hardware, OS and

virtualization layers, and at SQL Server and database configuration and adjust them as needed.

SQL Server provides many settings that you can use to fine-tune the installation to the system workload. There are also best practices that apply to most systems, including enabling IFI and *Optimize for Ad-Hoc Workloads* settings, increasing the number of files in tempdb, turning on some trace flags, disabling Auto Shrink, and setting up correct auto-growth parameters for database files.

In the next chapter, we'll talk about one of the most important components in SQL Server, SQLOS, and a troubleshooting technique called Wait Statistics.

Troubleshooting Checklist

Troubleshoot for the following items:

- Perform a high-level analysis of hardware, network and disk subsystem
- Discuss host configuration and load in virtualized environments with infrastructure engineers
- Check OS and SQL Server versions, editions and patching level
- Check if *instant file initialization* is enabled
- Analyze trace flags
- Enable *Optimize for Ad-Hoc Workloads*
- Check memory and parallelism settings on the server
- Look at tempdb settings (including number of files); check for trace flag T1118 and potentially T1117, in SQL Server versions prior to 2016
- Disable *Auto Shrink* for databases

- Validate data and t-log file settings
- Check number of VLFs in transaction log files
- Check errors in SQL Server Log
- Check for unnecessary monitoring in the system

¹ For a deep dive into data partitioning and disaster recovery strategies, please see my book *Pro SQL Server Internals (2nd ed., Apress, 2016)*.

Chapter 2. SQL Server Execution Model and Wait Statistics

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 2 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

It is impossible to troubleshoot SQL Server instances without understanding its execution model. You need to know how SQL Server runs tasks and manages resources if you want to detect bottlenecks in the system. We will cover those questions in this chapter.

First, the chapter will describe SQL Server’s architecture and major components. Next, it will discuss SQL Server’s execution model and introduce you to the popular troubleshooting technique called Wait Statistics. It will also cover several data management views commonly used during troubleshooting. Finally, it will provide you an overview of Resource Governor, which you can configure to segregate different workloads in the system.

SQL Server: High-Level Architecture

As you know, SQL Server is a very complex product that consists of dozens of components and subsystems. It is impossible to cover all of them here, but in this section, you'll get a high-level overview. For the sake of understanding, we'll divide these components and subsystems into seven categories, as shown in Figure 2-1. Let's talk about them.

Protocol Layer (Client Communication)		Utilities (DBCC, Backup, Restore, BCP, etc)
Query Processor		
Query Optimization (Plan Generation, Costing, Statistics, etc)	Query Execution (Parallelism, Memory Grants, etc)	
Storage Engine (Data Access, Locking Manager, Tran Log Management, etc)	In-Memory OLTP Engine	
SQLOS / PAL (Scheduling, Resource Management, Deadlock Detection, etc)		

Figure 2-1. Major SQL Server Components

The *Protocol Layer* handles communication between SQL Server and client applications. It uses an internal format called *Tabular Data Stream (TDS)* to transmit data using network protocols such as TCP/IP or Name Pipes. If a client application and SQL Server are running on the same machine, you can use another protocol called *Shared Memory*.

NOTE

It is worth checking what protocols are enabled when you troubleshoot client connectivity issues. Some SQL Server editions, for example Express and Developer, disable TCP/IP and Name Pipes by default. They do not accept remote client connections until you enable network protocols in the *SQL Server Configuration Manager* utility.

The *Query Processor* layer is responsible for query optimization and execution. It parses, optimizes and manages compiled plans for the queries, and orchestrates all aspects of query execution.

The *Storage Engine* is responsible for data access and management in SQL Server. It works with the data on disk, manages transaction logs, and handles transactions, locking and concurrency along with a few other functions.

The *In-Memory OLTP Engine* supports In-Memory OLTP in SQL Server. It works with memory-optimized tables and is responsible for data management and access to those tables, native compilation, data persistence, and all other aspects of the technology.

There are layers of abstraction between the components. For example, *Query Interop* (not shown in Figure 2-1) allows the Query Processor to work with row-based and memory-optimized tables, transparently routing requests either to Storage or to In-Memory OLTP engines.

The most critical abstraction layer is *SQL Server Operating System (SQLOS)*, which isolates other SQL Server components from the operating systems and deals with scheduling, resource management and monitoring, exception handling, and many other aspects of SQL Server behavior. For example, when any SQL Server component needs to allocate memory, it does not call OS API functions: it requests memory from SQLOS. This allows SQL Server granular control over execution and internal resource usage without relying on the OS.

Finally, since the introduction of Linux support in SQL Server 2017, there is another component called *Platform Abstraction Layer (PAL)*, which serves as a proxy between SQLOS and operating systems. Except for few performance-critical use cases, SQLOS does not call OS API directly, relying on PAL instead. This architecture allows SQL Server's code to remain almost identical in Windows and Linux, which significantly speeds up development and product improvements.

From a troubleshooting standpoint, you'll see very little difference between SQL Server on Windows and on Linux. Obviously, you'll use different

techniques when analyzing the SQL Server ecosystem and OS configuration. However, both platforms behave the same when you start to analyze issues *inside* SQL Server, so I am not going to differentiate between them in this book.

Let's look at the layers in more detail, beginning with SQLOS.

SQLLOS and the Execution Model

Database servers are expected to handle a large number of user requests, and SQL Server is no exception. On a very high level, SQL Server assigns those requests to separate threads, executing the requests simultaneously. Except the cases when the server is idle, the number of active threads exceeds the number of CPUs in the system, and efficient scheduling is the key to good server performance.

Early versions of SQL Server relied on Windows scheduling. Unfortunately, Windows (and Linux) are general purpose OSs, which means they use *preemptive scheduling*. They allocate a time interval, or *time quantum*, to a thread to run, then switch to other threads when it expires. This is an expensive operation that requires switching between user and kernel modes, negatively affecting system performance.

In SQL Server 7.0, Microsoft introduced the first version of *User Mode Scheduler (UMS)*-a thin layer between Windows and SQL Server that was primarily responsible for scheduling. It used *cooperative scheduling*, with SQL Server threads coded to voluntarily yield every 4ms, allowing other threads to execute. This approach significantly reduced expensive context switching in the system.

NOTE

Some SQL Server processes, like extended stored procedures, CLR routines, external languages and a few others, may still run in preemptive scheduling mode.

Microsoft continued to make improvements in UMS in SQL Server 2000 and, finally, in SQL Server 2005 redesigned it to the much more robust SQLOS. In later versions of SQL Server, SQLOS is responsible for scheduling, memory and I/O management, exception handling, CLR and external languages hosting, and quite a few other functions.

When you start an SQL Server process, SQLOS creates a set of *schedulers* that manage workload across CPUs. The number of schedulers matches the number of logical CPUs in the system, with additional scheduler created for a *Dedicated Admin Connection (DAC)*. For example, if you have two quad-core physical CPUs with hyper-threading enabled, SQL Server will create 17 schedulers in the system. For all practical purposes, you can think of schedulers as the CPUs; I will use those terms interchangeably throughout the book.

NOTE

The Dedicated Admin Connection is your *last resort* troubleshooting connection. It allows you to access SQL Server if it becomes unresponsive and does not accept normal connections. I will talk about it in Chapter 13.

Each scheduler will be in an ONLINE or OFFLINE state, depending on its affinity mask setting and core-based licensing model. The schedulers usually do not migrate between CPUs; however, it is possible, especially under heavy load. Nevertheless, in most cases this behavior does not affect the troubleshooting process.

The schedulers are responsible to manage the set of *worker threads*, sometimes called *workers*. The maximum number of workers in a system is specified by the *Max Worker Thread* configuration option. The default value of zero indicates that SQL Server calculates the maximum number of worker threads based on number of schedulers in the system. In most cases, you do not need to change this default value—in fact, don't change it unless you know *exactly* what you are doing.

Each time there is a task to execute, it is assigned to an idle worker. When there are no idle workers, the scheduler creates a new one. It also destroys idle workers after 15 minutes of inactivity or in case of memory pressure. Each worker uses 512KB of RAM in 32-bit and 2MB of RAM in 64-bit SQL Server for the thread stack.

Workers do not move between schedulers; tasks do not move between workers. SQLOS, however, can create child tasks and assign them to different workers, for example in the case of parallel execution plans. This may explain situations when some schedulers are running under heavier loads than others – some workers could end up with more expensive tasks from time to time.

You can think about *workers* as the logical representation of OS threads, and *tasks* as the unit of works those threads handle.

In most cases, we focus on tasks during troubleshooting. There is an exception, however: when a task is in the PENDING state, which means that it is waiting for available worker after the task had been created. This is completely normal, and workers are usually assigned to tasks very quickly. However, it can also indicate a very dangerous condition when the system does not have enough workers to handle the requests. I will discuss how to detect and address that issue in Chapter 13.

Besides PENDING, a task may be in five other possible states:

RUNNING

The task is currently executing on the scheduler.

RUNNABLE

The task is waiting for the scheduler to be executed.

SUSPENDED

The task is waiting for an external event or resource.

SPINLOOP

The task is processing a *spinlock*. *Spinlocks* are synchronization objects that protect some internal objects. SQL Server may use them when it expects that access to the object will be granted very quickly, avoiding context switching for the workers. I will talk about troubleshooting spinlock issues in Chapter 13.

DONE

The task is complete.

The first three states are the most important and common. Each scheduler has at most one task in the RUNNING state. In addition, it has two different queues—one for RUNNABLE and one for SUSPENDED tasks. When the RUNNING task needs some resources—a data page from a disk, for example—it submits an I/O request and changes the state to SUSPENDED. It stays in the SUSPENDED queue until the request is fulfilled and the page has been read. After that, when it is ready to resume execution, the task is moved to the RUNNABLE queue.

Perhaps the closest real-life analogy to this process is a grocery-store checkout line. Think of cashiers as schedulers and customers as tasks in the RUNNABLE queue. A customer who is currently checking out is similar to a task in the RUNNING state.

If item is missing a UPC code, a cashier sends a store worker to do a price check. The cashier suspends the checkout process for the current customer, asking her or him to step aside (to the SUSPENDED queue). When the worker comes back with the price information, the customer moves to the end of the checkout line (the end of the RUNNABLE queue).

Of course, SQL Server's execution is much more efficient than a real-life store, where customers must wait patiently in line for the price check to complete. (A customer in the end of the RUNNABLE queue would probably wish for such efficiency!)

Wait Statistics

With exception of initialization and clean-up, a task spends its time switching between RUNNING, SUSPENDED and RUNNABLE states, as shown in Figure 2-2. The total execution time will include time in RUNNING state, when task actually executed; time in RUNNABLE state, when the task is waiting for scheduler (CPU) to execute; and time in SUSPENDED state, when task is waiting for resources.

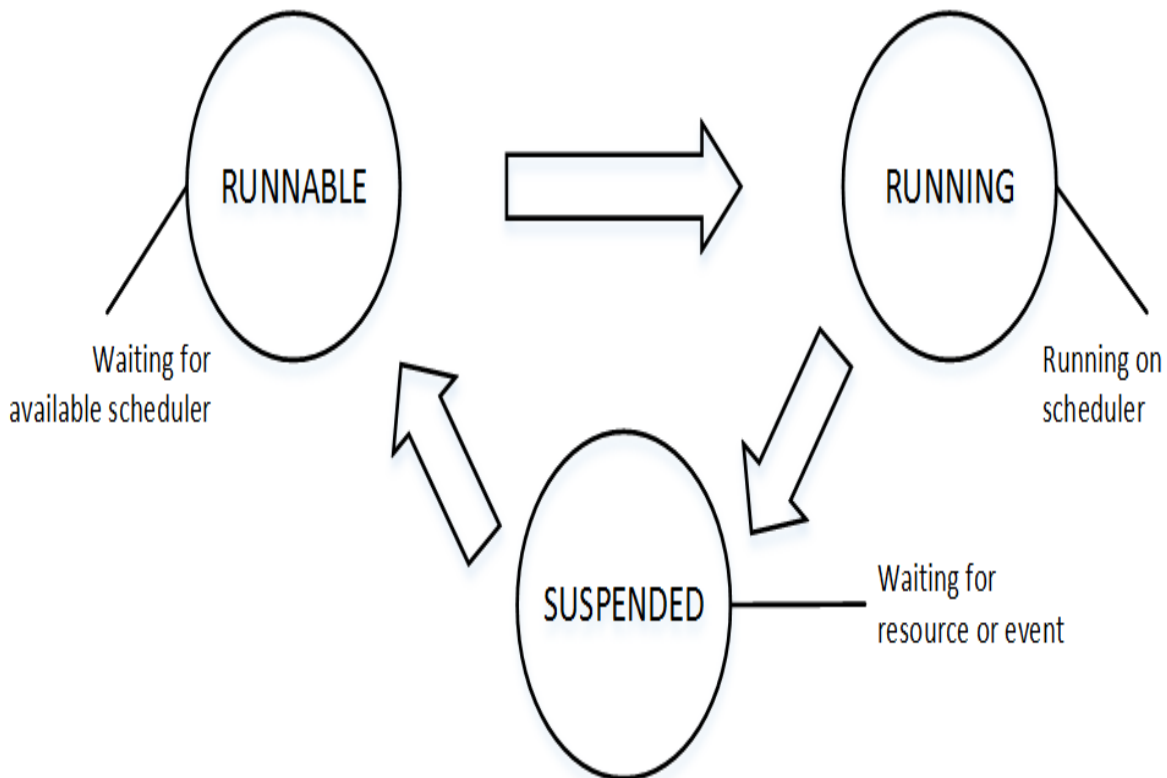


Figure 2-2. Task life cycle

In a nutshell, the goal of any performance-tuning process is improving system throughput by reducing query execution times. You can achieve this by reducing the time that query tasks spend in any of those states.

You can decrease query RUNNING time by upgrading hardware and moving to faster CPUs or by reducing amount of work tasks perform with query optimization.

You can shrink RUNNABLE time by adding more CPU resources or reducing the load on the system.

However, in most cases, you will get the most benefit by focusing on the time that tasks spend in SUSPENDED state while waiting for resources.

SQL Server tracks the cumulative time tasks spend in SUSPENDED state for different types of waits. You can view this data through the sys.dm_os_wait_stats view to get a quick sense of the main bottlenecks in your system and further fine-tune your troubleshooting strategy.

The code in Listing 2-1 shows you the wait types that take the most time in your system (filtering out some benign wait types, mainly related to internal SQL Server processes that spend most time waiting). The data is collected from the time of the last SQL Server restart, or since you last cleared it with the DBCC SQLPERF('sys.dm_os_wait_stats', CLEAR) command. Each new SQL Server version introduces new wait types. Some are useful for troubleshooting; others are benign and will need to be filtered out.¹

Example 2-1. Getting top wait types in the system (SQL Server 2012 and above)

```
;WITH Waits
AS
(
    SELECT
        wait_type, wait_time_ms, waiting_tasks_count, signal_wait_time_ms
        , wait_time_ms - signal_wait_time_ms AS resource_wait_time_ms
        , 100. * wait_time_ms / SUM(wait_time_ms) OVER() AS Pct
        , 100. * SUM(wait_time_ms) OVER(ORDER BY wait_time_ms DESC) /
            NULLIF(SUM(wait_time_ms) OVER(), 0) AS RunningPct
        , ROW_NUMBER() OVER(ORDER BY wait_time_ms DESC) AS RowNum
    FROM sys.dm_os_wait_stats WITH (NOLOCK)
    WHERE
        wait_type NOT IN /* Filtering out non-essential system waits */
        (N'BROKER_EVENTHANDLER',N'BROKER_RECEIVE_WAITFOR',N'BROKER_TASK_STOP
        ,
        ,N'BROKER_TO_FLUSH',N'BROKER_TRANSMITTER',N'CHECKPOINT_QUEUE',N'CHKP
        T'
        ,N'CLR_SEMAPHORE',N'CLR_AUTO_EVENT',N'CLR_MANUAL_EVENT'
        ,N'DBMIRROR_DBM_EVENT',N'DBMIRROR_EVENTS_QUEUE',N'DBMIRROR_WORKER_QU
        EUE'
        ,N'DBMIRRORING_CMD',N'DIRTY_PAGE_POLL',N'DISPATCHER_QUEUE_SEMAPHORE'
```

```

,N'EXECSYNC',N'FSAGENT',N'FT_IPTS_SCHEDULER_IDLE_WAIT',N'FT_IFTSHC_M
UTEX'
    ,N'HADR_CLUSAPI_CALL',N'HADR_FILESTREAM_IOMGR_IOCOMPLETION'
    ,N'HADR_LOGCAPTURE_WAIT',N'HADR_NOTIFICATION_DEQUEUE'

,N'HADR_TIMER_TASK',N'HADR_WORK_QUEUE',N'KSOURCE_WAKEUP',N'LAZYWRITE
R_SLEEP'
    ,N'LOGMGR_QUEUE',N'ONDEMAND_TASK_QUEUE'
    ,N'PARALLEL_REDO_WORKER_WAIT_WORK',N'PARALLEL_REDO_DRAIN_WORKER'
    ,N'PARALLEL_REDO_LOG_CACHE',N'PARALLEL_REDO_TRAN_LIST'
    ,N'PARALLEL_REDO_WORKER_SYNC'
,N'PREEMPTIVE_SP_SERVER_DIAGNOSTICS'
    ,N'PREEMPTIVE_OS_LIBRARYOPS'
,N'PREEMPTIVE_OS_COMOPS', N'PREEMPTIVE_OS_PIPEOPS'
    ,N'PREEMPTIVE_OS_GENERICOPS'
,N'PREEMPTIVE_OS_VERIFYTRUST'
    ,N'PREEMPTIVE_OS_FILEOPS'
,N'PREEMPTIVE_OS_DEVICEOPS'
    ,N'PREEMPTIVE_OS_QUERYREGISTRY'
,N'PREEMPTIVE_XE_CALLBACKEXECUTE'
    ,N'PREEMPTIVE_XE_DISPATCHER',N'PREEMPTIVE_XE_GETTARGETSTATE'
    ,N'PREEMPTIVE_XE_SESSIONCOMMIT',N'PREEMPTIVE_XE_TARGETINIT'

,N'PREEMPTIVE_XE_TARGETFINALIZE',N'PWAIT_ALL_COMPONENTS_INITIALIZED'

,N'PWAIT_DIRECTLOGCONSUMER_GETNEXT',N'PWAIT_EXTENSIBILITY_CLEANUP_TA
SK'
    ,N'QDS_PERSIST_TASK_MAIN_LOOP_SLEEP',N'QDS_ASYNC_QUEUE'
    ,N'QDS_CLEANUP_STALE_QUERIES_TASK_MAIN_LOOP_SLEEP'

,N'REQUEST_FOR_DEADLOCK_SEARCH',N'RESOURCE_QUEUE',N'SERVER_IDLE_CHEC
K'
    ,N'SLEEP_BPOOL_FLUSH',N'SLEEP_DBSTARTUP',N'SLEEP_DCOMSTARTUP'

,N'SLEEP_MASTERDBREADY',N'SLEEP_MASTERMDREADY',N'SLEEP_MASTERUPGRADE
D'
    ,N'SLEEP_MSDBSTARTUP',N'SLEEP_SYSTEMTASK',N'SLEEP_TASK'

,N'SLEEP_TEMPDBSTARTUP',N'SNI_HTTP_ACCEPT',N'SOS_WORK_DISPATCHER'
    ,N'SP_SERVER_DIAGNOSTICS_SLEEP',N'SQLTRACE_BUFFER_FLUSH'
    ,N'SQLTRACE_INCREMENTAL_FLUSH_SLEEP',N'SQLTRACE_WAIT_ENTRIES'
    ,N'STARTUP_DEPENDENCY_MANAGER',N'WAIT_FOR_RESULTS'
    ,N'WAITFOR',N'WAITFOR_TASKSHUTDOWN',N'WAIT_XTP_HOST_WAIT'

,N'WAIT_XTP_OFFLINE_CKPT_NEW_LOG',N'WAIT_XTP_CKPT_CLOSE',N'WAIT_XTP_
RECOVERY'

```

```

,N'XE_BUFFERMGR_ALLPROCESSED_EVENT',N'XE_DISPATCHER_JOIN',N'XE_DISPATCHER_WAIT'
,N'XE_LIVE_TARGET_TVF',N'XE_TIMER_EVENT')
)
SELECT
    w1.wait_type AS [Wait Type]
    ,w1.waiting_tasks_count AS [Wait Count]
    ,CONVERT(DECIMAL(12,3), w1.wait_time_ms / 1000.0) AS [Wait Time]
    ,CONVERT(DECIMAL(12,1), w1.wait_time_ms / w1.waiting_tasks_count)
      AS [Avg Wait Time]
    ,CONVERT(DECIMAL(12,3), w1.signal_wait_time_ms / 1000.0)
      AS [Signal Wait Time]
    ,CONVERT(DECIMAL(12,1), w1.signal_wait_time_ms /
w1.waiting_tasks_count)
      AS [Avg Signal Wait Time]
    ,CONVERT(DECIMAL(12,3), w1.resource_wait_time_ms / 1000.0)
      AS [Resource Wait Time]
    ,CONVERT(DECIMAL(12,1), w1.resource_wait_time_ms /
w1.waiting_tasks_count)
      AS [Avg Resource Wait Time]
    ,CONVERT(DECIMAL(6,3), w1.Pct)
      AS [Percent]
    ,CONVERT(DECIMAL(6,3), w1.RunningPct)
      AS [Running Percent]
FROM
    Waits w1

WHERE
    w1.RunningPct <= 99 OR w1.RowNum = 1
ORDER BY
    w1. RunningPct
OPTION (RECOMPILE, MAXDOP 1);

```

Figure 2-2 shows the output of this code from one of the production servers, early in the troubleshooting process. I can immediately see that majority of the waits in the system relate to blocking (LCK*) and I/O (PAGEIOLATCH*). This makes it much easier to decide where to focus my troubleshooting efforts.

	Wait Type	Wait Count	Wait Time	Avg Wait Time	Signal Wait Time	Avg Signal Wait Time
1	LCK_M_U	538312358	2952904.553	5.0	278904.170	1.0
2	PAGEIOLATCH_SH	132056495	730022.059	5.0	17938.737	0.0
3	LCK_M_S	196405075	379378.938	1.0	24706.314	0.0
4	ASYNC_NETWORK_IO	36665258	254793.758	6.0	100063.339	2.0
5	LOGBUFFER	11718571	165042.270	14.0	18931.562	1.0
6	PAGEIOLATCH_EX	153474407	133057.566	0.0	3225.941	0.0
7	LCK_M_IX	496185	98525.504	198.0	139.082	0.0
8	IO_COMPLETION	93217317	81833.420	0.0	3505.294	0.0
9	LATCH_EX	49863173	65876.146	1.0	10921.396	0.0
10	ASYNC_IO_COMPLETION	57845	56036.933	968.0	22.078	0.0
11	LCK_M_IS	57448	31694.644	551.0	9.403	0.0
12	LCK_M_SCH_M	2228	31016.126	13921.0	0.918	0.0
13	WRITELOG	1969821	26014.687	13.0	715.277	0.0
14	OLEDB	3041936	14911.992	4.0	6058.799	1.0
			Resource Wait Time	Avg Resource Wait Time	Percent	Running Percent
			2674000.376	4.0	58.316	58.316
			712083.322	5.0	14.417	72.733
			354672.624	1.0	7.492	80.226
			154730.419	4.0	5.032	85.258
			146110.708	12.0	3.259	88.517
			129831.625	0.0	2.628	91.145
			98386.422	198.0	1.946	93.091
			78328.126	0.0	1.616	94.707
			54954.750	1.0	1.301	96.008
			56014.855	968.0	1.107	97.114
			31685.241	551.0	0.626	97.740
			31015.208	13920.0	0.613	98.353
			25299.410	12.0	0.514	98.867
			8853.193	2.0	0.294	99.161

Figure 2-3. Example of sys.dm_os_wait_stats output

This troubleshooting approach is called *Wait Statistics Analysis*. It's the one of the most frequently used troubleshooting and performance-tuning

techniques in SQL Server. Figure 2-4 illustrates a typical troubleshooting cycle using Wait Statistics Analysis.

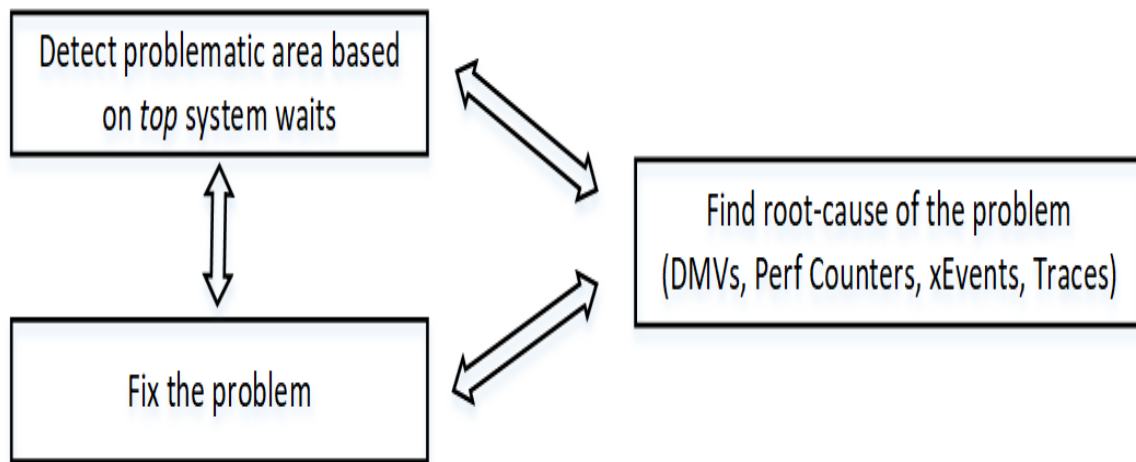


Figure 2-4. Typical Wait Statistics Analysis Troubleshooting Cycle

First, you identify the main bottleneck in the system by analyzing the top waits. Next, you confirm it with other tools and techniques and pinpoint the root cause of the problem. Finally, you fix it and repeat the cycle.

WARNING

A word of caution: This process may never end. While there are always opportunities to make things better, at some point further improvements become impractical. Remember the Pareto Principle – you will get 80% of improvements by spending 20% of your time – and don't waste time on nonessential tuning.

This looks very easy in theory; unfortunately, it is more complicated in real life. Many issues are related to each other, which can hide the real causes of bottlenecks. To choose a very common example: excessive disk waits are often triggered not by bad I/O performance, but by poorly optimized queries that constantly flush the buffer pool and overload the disk subsystem.

Figure 2-5 shows some of the high-level dependencies you might run into. This diagram is by no means exhaustive, but it illustrates the danger of tunnel vision during troubleshooting.

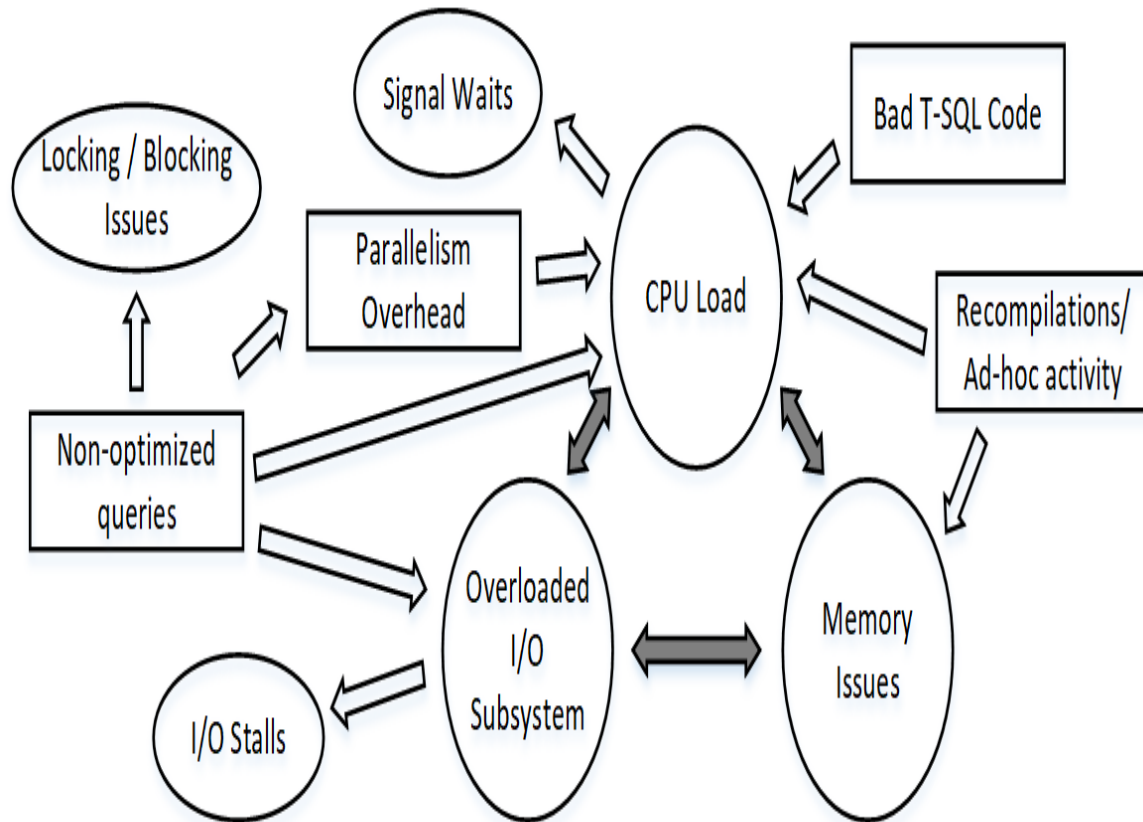


Figure 2-5. Dependencies and Issues

I considered listing the most common waits and possible root causes here, but I don't want you to start chasing symptoms rather than causes. Rather than jumping right to a list, you read the book first so that you can understand the possible dependencies involved.

I'll start going through specific issues and troubleshooting techniques in the upcoming chapters, but for now, let's cover important data management views in SQL Server related to SQLOS and the SQL Server execution model.

Execution Model–Related Data Management Views

SQL Server comes with a very large number of data management views (DMVs). For details on all of them, you can consult the [Microsoft](#)

Documentation. Here, I will focus on just a small subset that I regularly use during troubleshooting. We will look at many others later in the book.

sys.dm_os_wait_stats

As you saw earlier, the `sys.dm_os_wait_stats` **view** provides information about waits in the system. It will tell you how many times the wait occurs (`waiting_task_count`) along with cumulative times for resource (`resource_wait_time_ms`) and signal (`signal_wait_time_ms`) waits. The resource wait time indicates how long a task waited for the resource staying in `SUSPENDED` queue. The signal wait indicates the wait for the CPU in `RUNNABLE` queue after the resource wait was over.

For example, let's say a task is requested to read a data page from disk. The I/O request might take 6ms; then, the task might wait for another millisecond to resume the execution. If you view the wait data for this, you'll see 6ms of resource waits, 1ms of signal waits, and 7ms of total wait time.

Listing 2-2 shows you how to compare cumulative signal and resource waits in the system.

Example 2-2. Signal versus resource waits

```
SELECT
    SUM(signal_wait_time_ms) AS [Signal Wait Time (ms)]
    , CONVERT(DECIMAL(7,4), 100.0 * SUM (signal_wait_time_ms) /
        SUM(wait_time_ms)) AS [% Signal waits]
    , SUM(wait_time_ms - signal_wait_time_ms) AS [Resource Wait Time
(ms)]
    , CONVERT (DECIMAL(7,4), 100.0 * sum(wait_time_ms -
signal_wait_time_ms) /
        SUM(wait_time_ms)) AS [% Resource waits]
FROM
    sys.dm_os_wait_stats WITH (NOLOCK);
```

In most cases, signal waits should not exceed 10 to 15% of total wait time. A higher number may indicate a CPU bottleneck, with tasks spending a lot of time in the `RUNNABLE` queue. Do not jump to the conclusion that you need to add more CPUs, though—it may be entirely possible to address the problem with performance tuning.

Pay attention to how often waits occur. Sometimes, you'll see waits with a low `waiting_task_count` and high total wait time. Depending on the situation, you may or may not want to analyze them, especially during the initial phase of troubleshooting. Such waits are often triggered by production incidents or other atypical conditions.

Finally, make sure that you are working with representative data. As I mentioned, statistics are collected from the time of the last SQL Server restart, and workload on the server may change over time.

I usually ask customers to clear the waits a few days before starting the troubleshooting. It is safe to use the `DBCC SQLPERF('sys.dm_os_wait_stats', CLEAR)` command in production, although it may affect data collection in some third-party monitoring tools. As another option, you can collect two separate snapshots of wait statistics and calculate the delta between them. I am including the script to do that to the companion materials of the book.

sys.dm_exec_session_wait_stats

Starting with SQL Server 2016, you can look at waits on the session level, using the `sys.dm_exec_session_wait_stats` **view**. This is extremely useful when you need to troubleshoot performance of long-running queries or slow stored procedures in the system. The view will show you the waits that occurred during execution and help you pinpoint bottlenecks and areas to research.

The columns and data in this view are similar to those in `sys.dm_os_wait_stats`; you can easily adjust scripts to work in both scenarios. Remember that data in `sys.dm_exec_session_wait_stats` clears when a session opens and when the pooled connection resets.

You may notice that the data is not always updated for currently running statements. You need to wait until a query completes for the data to become available.

sys.dm_os_waiting_tasks

The `sys.dm_os_waiting_tasks` **view** shows you a list of tasks that are *currently* waiting in the SUSPENDED queue. This view is handy when the server is overloaded or unresponsive and you want to understand why sessions were suspended.

It is also very helpful when you troubleshoot concurrency issues and active blocking in the system, because it shows you the session ID of the blocker for the task (more in Chapter 8).

The most useful columns in this view are:

session_id

ID of the waiting session.

wait_type

Type of wait the session is waiting for.

wait_duration_ms

The duration of the wait.

blocking_session_id

The session blocking the current task. As I mention, this column is extremely useful when you troubleshoot active blocking in the system.

resource_address

Information on the resource the task is waiting for.

You may have more than one row per session in the output when you deal with parallel execution plans.

sys.dm_exec_requests

The `sys.dm_exec_requests` **view** provides detailed information on each request that is executing on the server. This gives you a great *at-a-glance*

snapshot of what is happening now and allows you to pinpoint most CPU or I/O intensive queries *currently* running in the system.

This view will return information for both user and system sessions. You can filter out most system sessions by using WHERE `session_id > 50` predicate, although you may have some system sessions with id greater than 50 nowadays.

The most useful columns in this view are:

session_id

The ID for the session. Unlike with `sys.dm_os_waiting_tasks`, you get a single row in the output per session unless you are using *Multiple Active Result Sets* (MARS) in your system.

start_time

The time when the request started.

total_elapsed_time

The request's duration.

status

The current request status (RUNNING, RUNNABLE, SUSPENDED, SLEEPING). SLEEPING status indicates an idle connection.

wait_type, wait_time, wait_resource, blocking_session_id

These appear if the request is currently suspended. Like `sys.dm_os_waiting_tasks`, the `blocking_session_id` column is very useful when you are troubleshooting active blocking in the system.

cpu_time, logical_reads, reads, writes, granted_query_memory, dop

These provide you with execution metrics.

sql_handle, plan_handle

These allow you to obtain the statement and its execution plan.

Listing 2-3 shows you the code that returns information about currently running CPU-intensive requests, along with connection information. This code requires SQL Server 2012 to run – you can remove TRY_CONVERT function if you have older version of SQL Server.

Example 2-3. Using sys.dm_exec_requests view

```
SELECT
    er.session_id
    ,er.request_id
    ,DB_NAME(er.database_id) as [database]
    ,er.start_time
    ,CONVERT(DECIMAL(21,3),er.total_elapsed_time / 1000.) AS
[duration]
    ,er.cpu_time
    ,SUBSTRING(
        qt.text,
        (er.statement_start_offset / 2) + 1,
        ((CASE er.statement_end_offset
            WHEN -1 THEN DATALENGTH(qt.text)
            ELSE er.statement_end_offset
            END - er.statement_start_offset) / 2) + 1
    ) AS [statement]
    ,er.status
    ,er.wait_type
    ,er.wait_time
    ,er.wait_resource
    ,er.blocking_session_id
    ,er.last_wait_type
    ,er.reads
    ,er.logical_reads
    ,er.writes
    ,er.granted_query_memory
    ,er.dop
    ,er.row_count
    ,er.percent_complete
    ,es.login_time
    ,es.original_login_name
    ,es.host_name
    ,es.program_name
    ,c.client_net_address
```



```

,ib.event_info AS [buffer]
,qt.text AS [sql]
,TRY_CONVERT(XML,p.query_plan) as [query_plan]
FROM
sys.dm_exec_requests er WITH (NOLOCK)
  OUTER APPLY sys.dm_exec_input_buffer
    (er.session_id, er.request_id) ib
  OUTER APPLY sys.dm_exec_sql_text(er.sql_handle) qt
  OUTER APPLY
    sys.dm_exec_text_query_plan
    (
      er.plan_handle
      ,er.statement_start_offset
      ,er.statement_end_offset
    ) p
LEFT JOIN sys.dm_exec_connections c WITH (NOLOCK) ON
  er.session_id = c.session_id
LEFT JOIN sys.dm_exec_sessions es WITH (NOLOCK) ON
  er.session_id = es.session_id
WHERE
  er.status <> 'background'
AND er.session_id > 50
ORDER BY
  er.cpu_time desc
OPTION (RECOMPILE, MAXDOP 1);

```

As a word of caution: getting a query execution plan with `sys.dm_exec_text_query_plan` function is expensive. Comment it out if your server is running under heavy CPU load.

sys.dm_os_schedulers

I do not use the `sys.dm_os_schedulers` **view** very often, only from time to time. As you can guess by the name, this view provides information about schedulers in the system. You can use it to get information about schedulers' distribution across NUMA nodes and to analyze metrics from individual schedulers.

I've already shown you the code for the first use case in Chapter 1, but it is worth repeating. Check the count of schedulers in each NUMA node to see if the CPU affinity has been set correctly.

Example 2-4. NUMA nodes schedulers statistics

```

SELECT
    parent_node_id
    ,COUNT(*) as [Schedulers]
    ,SUM(current_tasks_count) as [Current]
    ,SUM(runnable_tasks_count) as [Runnable]
FROM sys.dm_os_schedulers
WHERE status = 'VISIBLE ONLINE'
GROUP BY parent_node_id;

```

The `current_tasks_count` and `runnable_tasks_count` columns provide the number of tasks in the `RUNNING` and `RUNNABLE` queues in each node. A large `runnable_tasks_count` number may indicate a CPU bottleneck. Remember, however, that the numbers show what is happening in the system *now* and may not be representative over time. It is better to see cumulative information, for example the percentage of signal waits (see Listing 2-2) or CPU load overtime (see Chapter 6).

There are many other columns in the view that provide scheduler-specific statistics, such as status, number of workers and tasks in various states, number of context switches, CPU consumption and a few others. Check the documentation for more details.

Resource Governor Overview

Resource Governor is an Enterprise Edition feature that allows you to segregate and throttle different workloads on the server. Although it's been available for quite a long time, I consider Resource Governor a niche feature – I rarely see it in the field. (You may even consider skipping this section, and coming back if and when you have to deal with it.) Nevertheless, remember to check if Resource Governor is configured in the system you are troubleshooting – incorrect configuration can seriously impact server throughput.

When enabled, Resource Governor separates the sessions between different *workload groups* by calling *classifier function* at the time of the session's login. The classifier function is a simple user-defined function where you can use various connection properties (login name, application name, client IP address, etc.) to choose between workload groups.

Each workload group has several parameters, such as MAXDOP, maximum allowed CPU time for the request, and the maximum number of simultaneous requests allowed in the group. The workgroups are also associated with a *resource pool*, where you can customize resource usage for associated workload groups.

The SQL Server documentation refers to *resource pools* as “the virtual SQL Server instances inside of a SQL Server instance.” I do not think this is an accurate definition, though, because resource pools do not provide enough isolation from each other. However, you can control and limit CPU bandwidth and affinity, along with query memory grants (see Chapter 7).

Starting with SQL Server 2014, you can also control disk throughput by limiting resource pool IOPS. *You cannot, however, control buffer pool usage—it is shared across all pools.*

There are two system workload groups and resource pools: *internal* and *default*. As you can guess by the names, the first handles internal workload. The second is responsible for all non-classified workload. You can change the parameters of the *default* workload group without creating other user-defined workload groups and pools.

Figure 2-6 shows a Resource Governor configuration for an example scenario in which you want to separate OLTP and reporting workloads. This will reduce the impact of reporting queries on critical OLTP transactions, preventing them from saturating CPU and I/O.

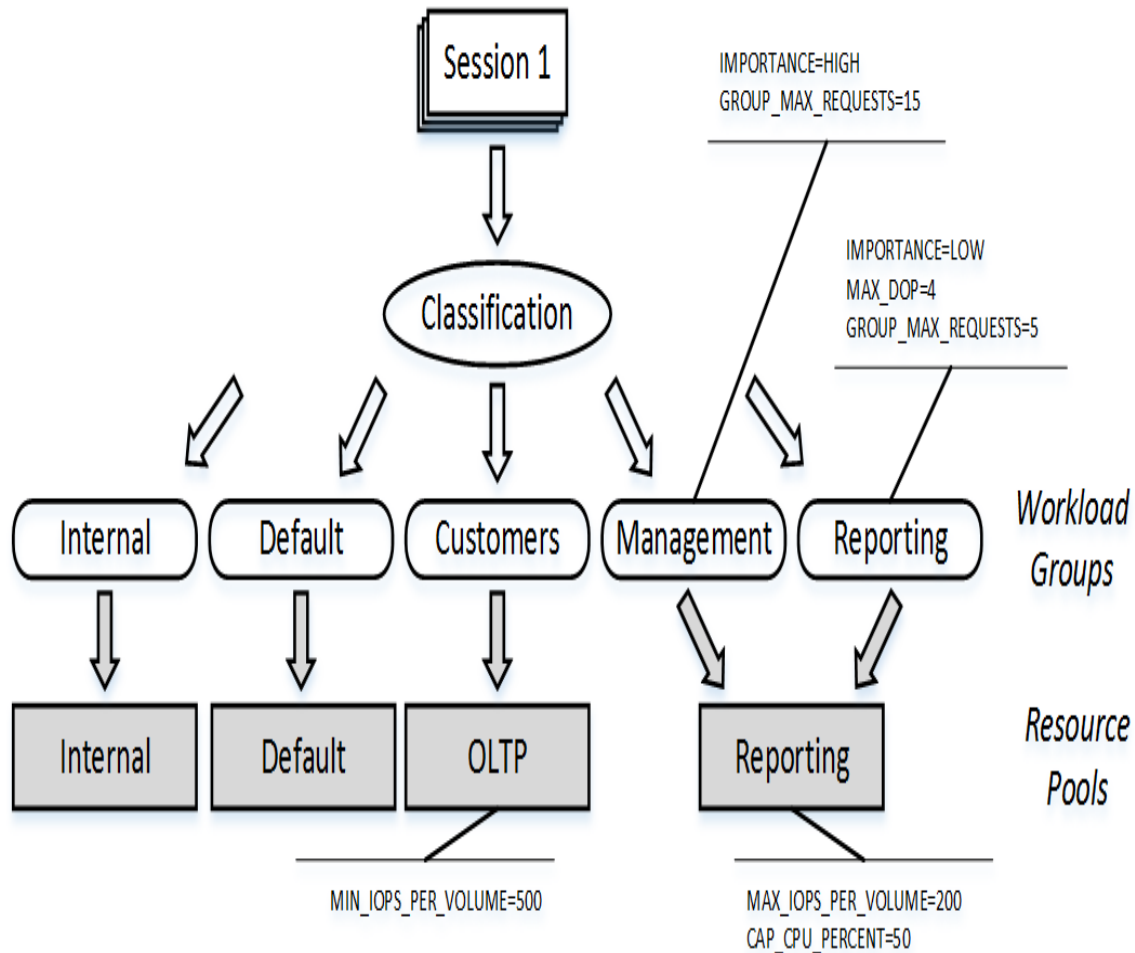


Figure 2-6. Example of Resource Governor Configuration

Resource Governor is useful, but it is not the easiest feature to configure and maintain. You need to do some planning and math when you want to configure resource throttling across multiple busy resource pools.

You also need to reevaluate the settings overtime, because hardware and workload requirements may change. I recently had to troubleshoot a case where a major disk subsystem upgrade did not improve system performance. We found that I/O in the system had been throttled by a MAX_IOPS_PER_VOLUME setting in the resource pool.

In conclusion, Resource Governor is good in use cases where you need to segregate different workloads in a single database on a standalone server or an instance that uses Failover Clustering. It is also useful for reducing the impact of database maintenance. For example, you can limit CPU utilized

by backup compression or I/O load from index maintenance by running them in a separate resource pool.

I recommend looking at different technologies when you need to segregate a different workload in the Always On Availability Groups setup. The readable secondaries may provide better scalability in the long term.

In addition, when you need to segregate workloads from multiple databases running on a single SQL Server instance, it's usually better to split the databases across multiple instances, and potentially virtualize them.

Summary

SQLOS is the vital subsystem responsible for scheduling and resource management in SQL Server. At startup, it creates schedulers—one per logical CPU—allocating the pool of worker threads to each scheduler to manage. User and system tasks are assigned to the worker threads, which perform the actual work.

SQL Server uses cooperative scheduling, with workers voluntarily yielding every 4ms. The tasks constantly migrate through the RUNNING, SUSPENDED, and RUNNABLE states while they are running on CPU or waiting for CPU and resources. SQL Server tracks the different type of waits and provides that information in `sys.dm_os_wait_tasks` view. You can analyze the most common waits and identify bottlenecks in the system with the troubleshooting process called Wait Statistics.

Be careful when analyzing waits; don't jump to immediate conclusions. Many performance issues may be related and can mask each other. You'll need to identify and confirm the root cause of the problem as part of your analysis.

In the next chapter, we will dive deeper into troubleshooting particular issues, starting with the disk, and learn how to diagnose and address them.

Troubleshooting Checklist

Troubleshoot as follows:

- Look at the waits in the system. Make sure that wait statistics are representative.
- Analyze percentages of signal and resource waits.
- Validate Resource Governor configuration when present.
- Triage the waits, looking for bottlenecks.

¹ The code in Listing 2-1 is good for versions up to SQL Server 2019. To exclude other wait types in future versions, see [Microsoft's documentation](#).

Chapter 3. Troubleshooting Disk Subsystem Issues

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 3 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

SQL Server is a very I/O intensive application: it is constantly reading data from and writing data to disk. Good I/O throughput is essential for SQL Server performance and health. Unfortunately, many SQL Server installations are I/O bound, even with modern flash-based storage.

In this chapter, I will show you how to analyze and troubleshoot disk subsystem performance issues. You will learn how SQL Server processes I/O requests internally and how to identify and detect possible bottlenecks through the entire I/O stack, on the SQL Server, OS, Virtualization and Storage levels.

Next, I will talk about checkpoint process tuning, a common source of I/O bottlenecks in busy OLTP systems.

Finally, I will cover the most common I/O-related waits you may encounter in your system.

Anatomy of the SQL Server I/O Subsystem

SQL Server never works with data pages directly in database files. Every time a data page needs to be read or modified, SQL Server reads that page to memory and caches it in the buffer pool. Each page in a buffer pool is referenced by a buffer structure, sometimes simply called buffer. It includes the page's address in the data file, a pointer to the data page in the memory, status information, and the page latching queue.

SQL Server uses latches to protect internal objects in memory preventing their corruptions when multiple threads modifying them simultaneously. The two most common types of latch are exclusive, which blocks any access to the object, and shared, which allows simultaneous reads but prevents modifications of the objects.

Conceptually, latches are similar to critical sections or mutexes in application development languages. We will talk about latches in detail in Chapter 10.

The location of data pages in a buffer pool does not represent the order in which they are stored in the database files. SQL Server, however, can efficiently locate the page in the buffer pool when needed. Every time SQL Server accesses the page there, it performs a logical read. When the page is not present in memory and needs to be read from disk, the physical read also occurred.

When data needs to be modified, SQL Server changes the pages in the buffer pool, marking them as dirty, then writes log records to the transaction log file. It saves dirty pages to the data files asynchronously in the Checkpoint and, sometimes, the Lazy Writer processes. We'll discuss both of those processes later in this chapter and transaction logs in Chapter 11. For now, remember that data modifications require SQL Server to read data pages from disk if they have not been already cached.

Now let's look at how SQL Server works with I/O in more detail.

Scheduling and I/O

As you remember from chapter 2, SQL Server uses cooperative scheduling, with multiple workers running on CPUs in a rotating fashion. The workers voluntarily yield when the short quantum expires, allowing other workers to proceed. This model requires SQL Server to use asynchronous I/O as much as possible – it is impossible for workers to wait until I/O request is completed, preventing other workers from executing.

By default, all SQL Server schedulers handle I/O requests. You can override this behavior and bind I/O to specific CPUs by setting the affinity I/O mask. In theory, this may help improve I/O throughput in very busy OLTP systems; however, I rarely find it necessary. In most cases, you'll achieve better results by performing optimizations and reducing CPU and I/O load.

You can read about affinity I/O masking in the [Microsoft documentation](#).

Every scheduler has a dedicated I/O queue. When a worker needs to perform an I/O operation, it creates an I/O request structure, puts it to the scheduler's queue and finally issues an asynchronous OS API I/O call. It does not wait until the request is completed; it either continues to run, doing other things, or suspends itself, moving to the SUSPENDED queue.

When a new worker starts to run on the scheduler (switching to RUNNING state), it goes through the scheduler's I/O queue. The I/O request structures contain enough information to check if the asynchronous OS API call has been completed, along with a pointer to callback function that the worker calls to complete the I/O request in SQL Server.

I know that this sounds complicated – please bear with me and we'll look at the details in the next section. The key things I'd like you to remember are:

All active schedulers are handling I/O requests by default.

Most I/O requests in SQL Server use asynchronous OS API calls. This is true even for write-ahead logging – the worker that issues the COMMIT statement may be suspended until the log record is written to disk; however, the OS API write command will be executed asynchronously.

The I/O request may be completed by a different worker than the one that issues it.

You can see a list of pending I/O requests in the `sys.dm_io_pending_io_requests` [view](#). The `io_pending_ms_ticks` column provides the duration of that request. The `io_pending` column indicates if the OS API call has been completed and if the request is waiting for a worker to finish it. This may help you to determine if request latency is being affected by CPU load in the system.

Now, as promised, let's look at that process again, with more concrete examples of reading data pages from disk.

Data Reads

When SQL Server needs to access a data page, it checks if the page already exists in buffer pool. If it does not, the worker allocates the buffer for the page, protecting it with an exclusive latch. This prevents workers from accessing the page until it is read – they will be blocked, waiting for the latch to clear.

Next, the worker creates the I/O request structure, puts it in the scheduler I/O queue, and initiates an OS API read request. Then it tries to acquire another shared latch on the buffer, which is blocked by the incompatible, exclusive latch held there. The worker then suspends itself with PAGEIOLATCH wait (Figure 3-1 illustrates that state).

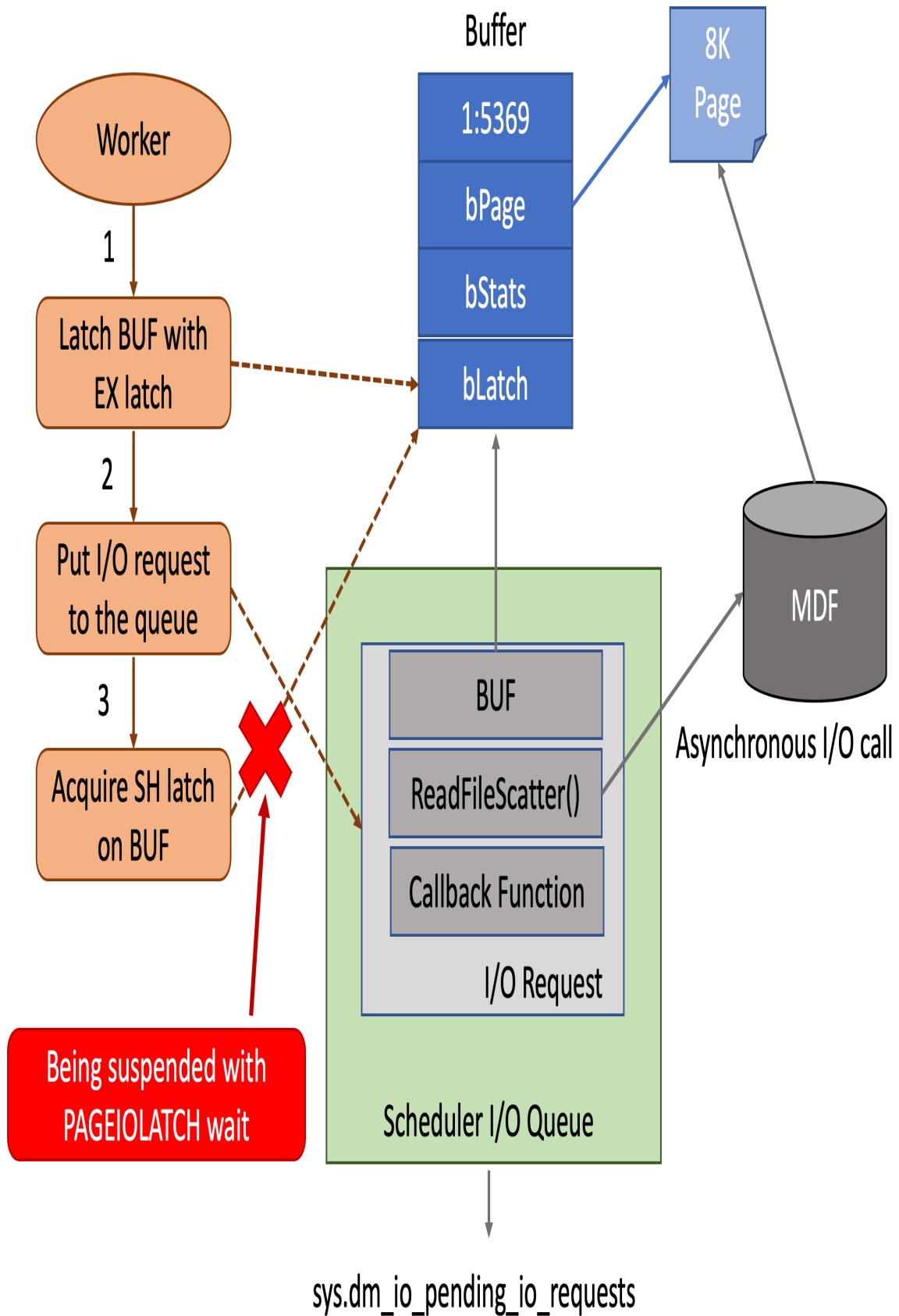


Figure 3-1. Reading data page from disk – initiating read

When another worker switches to a RUNNING state, it checks to see if any I/O requests in the scheduler's queue have been completed. If so, the worker calls the callback function to finalize the operation: this validates that page is not corrupted and removes the exclusive latch from the buffer. The worker that submitted the I/O request can then resume and access the data page (Figure 3-2).

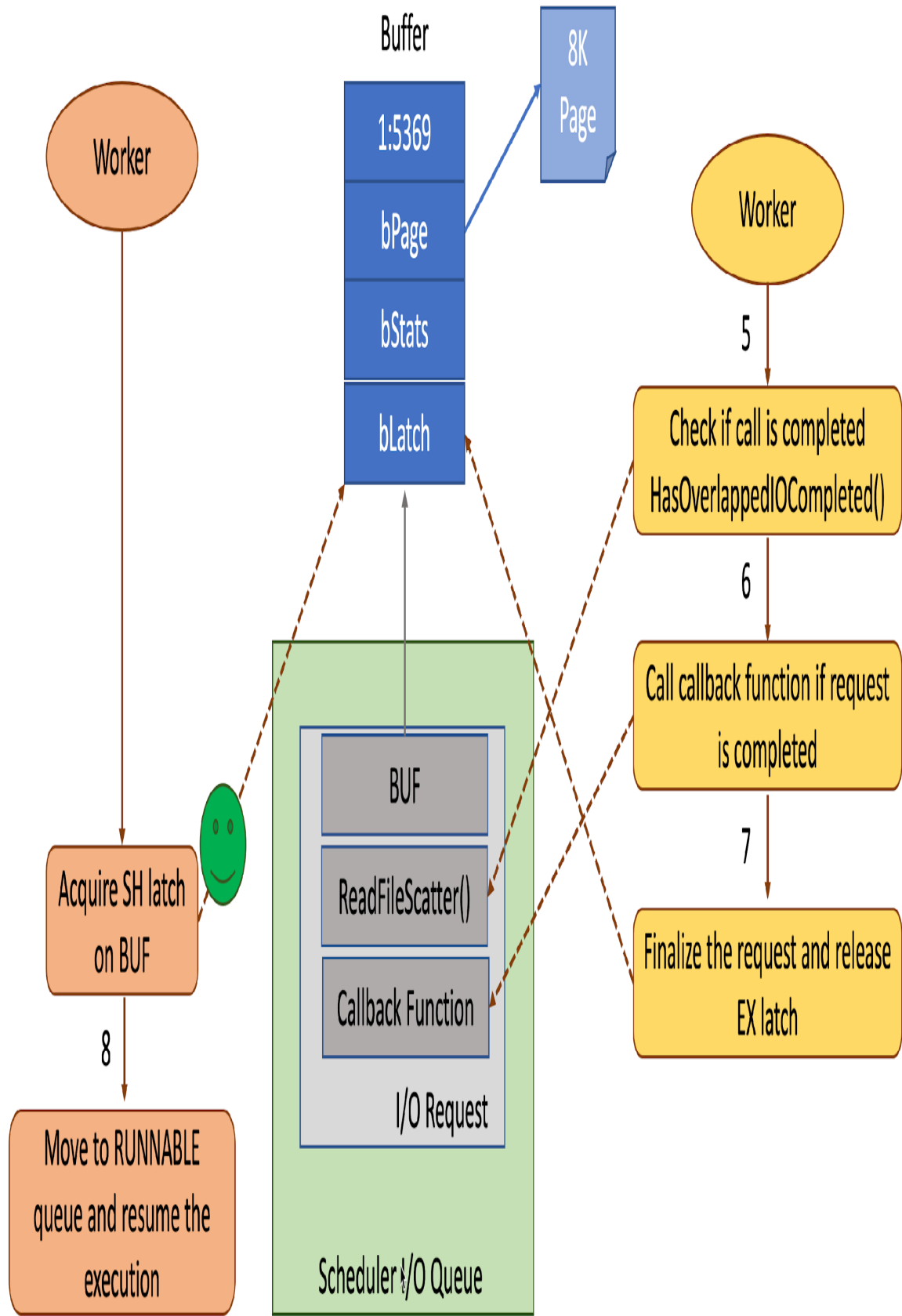


Figure 3-2. Reading data page from disk – completing the read

There are several errors that may occur during I/O requests. All of them are severe, and you need to set up alerts in the system for them.

- Error 823 indicates that the OS I/O API call was not successful. This is often a sign of hardware issues.
- Errors 605 and 824 indicate logical consistency issues with the data pages. When you encountered either of these errors, immediately check whether the database is corrupted, using the DBCC CHECKDB command. You may also encounter those errors in case of faulty I/O drivers, which can corrupt data pages during transfer.
- Error 833 tells you that an I/O request (OS API call) took longer than 15 seconds to return. This is abnormal; check the health of the disk subsystem when you see this error.
- Error 825 indicates that an I/O request failed and had to be retried in order to succeed. As with Error 833, check the health of the disk subsystem.

When troubleshooting those errors, you can look for the details in your SQL Server error log (use the code from Listing 1-4) and system event log.

It is very common for SQL Server to read multiple data pages in a single I/O request. For example, it uses read-ahead logic, reading multiple data pages during scans. As result, the query may perform thousands of logical reads with just a handful of physical reads. Another example is ramp-up reads, which is when SQL Server reads a large number of pages on each I/O request, trying to fill the buffer pool quickly on startup.

Data Writes

SQL Server handles data writes very similarly to data reads. In most cases, those writes are done asynchronously using a scheduler's I/O queues, as

you just saw in the previous examples. Obviously, the callback function will be implemented differently in different I/O operations.

When you change some data in the database, SQL Server modifies data pages in the buffer pool, reading pages from disk if needed. It generates log records for the modifications and saves them to the transaction log. The transaction is not considered to have been committed until the log records are hardened on disk. While, technically, you can treat write-ahead logging as synchronous writes, SQL Server uses an asynchronous I/O pattern for log writes.

SQL Server writes modified data pages in user databases asynchronously during checkpoint. This process finds dirty data pages in the buffer pool and saves them to disk. It tries to minimize the number of disk requests by combining and writing adjacent modified pages together in a single I/O operation when possible.

Another SQL Server process, called lazy writer, periodically sweeps the buffer pool to remove data pages that have not been recently accessed, freeing up the memory. In normal circumstances, lazy writer skips dirty data pages; however, it may also write them to disk if there is memory pressure in the system.

There are, as always, some exceptions. For example, during a bulk import operation, SQL Server allocates a set of buffers in the buffer pool and reuses them, writing data to the database outside of checkpoint. This preserves the content of the buffer pool, so it isn't flushed by massive data imports.

Checkpoint I/O may introduce issues on busy systems. I will talk about checkpoint tuning later in this chapter. But first, let's take a holistic look at the entire storage subsystem.

The Storage Subsystem: A Holistic View

Troubleshooting slow I/O performance in SQL Server is not an easy task. I've seen many heated discussions between database and infrastructure

teams. Database engineers generally complain about slow disk performance, while the storage engineers analyze the metrics from SAN devices with sub-millisecond latency and insist that all issues are on the SQL Server side. Neither team is right. They usually make the same mistake: oversimplifying the storage subsystem to just a couple of components. However, the storage subsystem isn't that simple.

Figure 3-3 shows a very high-level diagram of the network-based storage subsystem, with many details missing. (It also references some troubleshooting tools. We'll get to those, but don't focus on them yet.) The point here is that bad I/O performance can be caused by any component, so you need to analyze all layers in the stack.

Checkpoint Pages/sec, Background Writer Pages/sec, Lazy writer/sec,
Pages reads/sec, Pages writes/sec, Log Bytes Flushed/sec

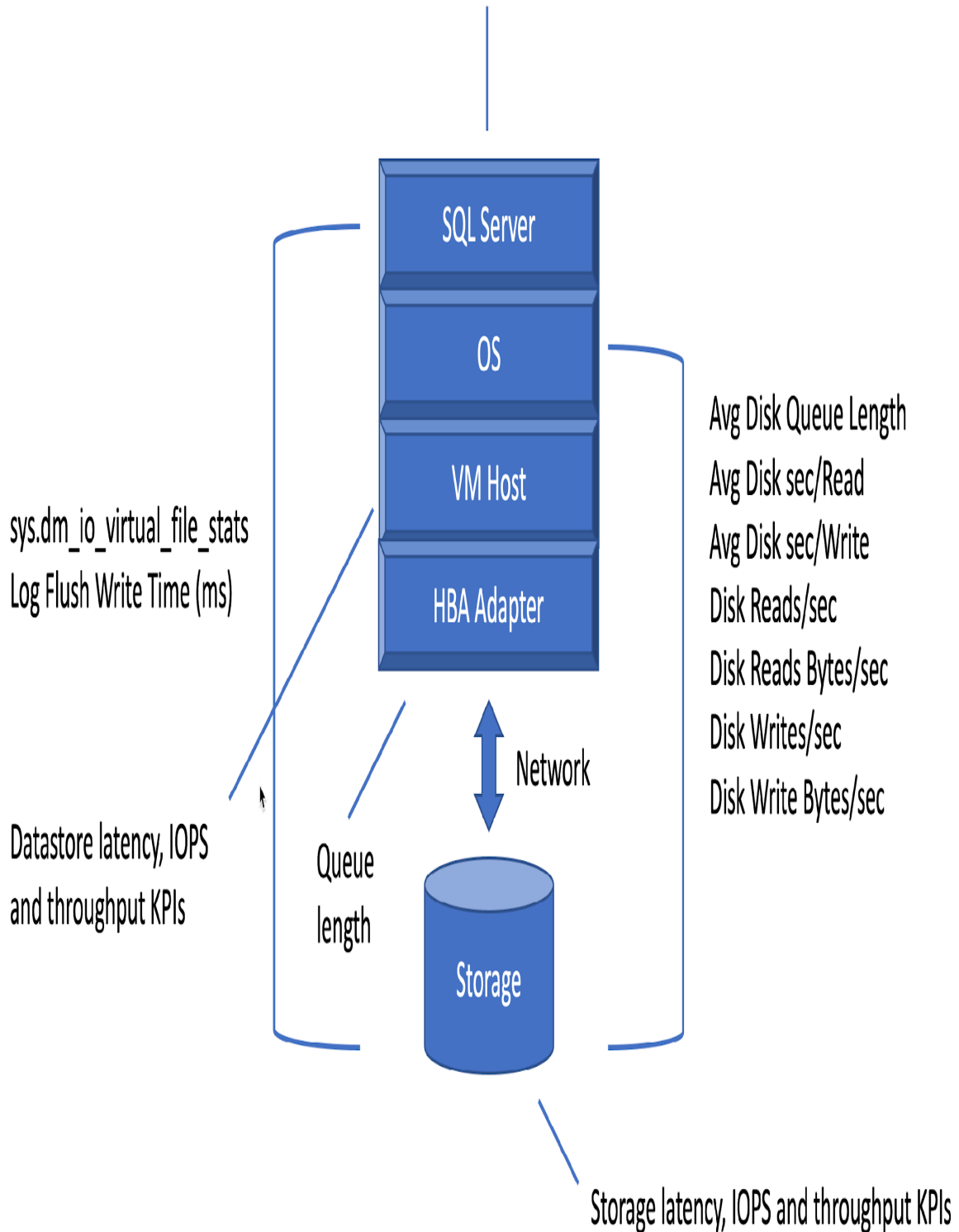


Figure 3-3. Storage Subsystem (Network-Based)

There is also an option of using direct-attached storage (DAS). In this configuration, the storage is either installed locally on the server (think about NVMe drives) or directly connected to it. This setup eliminates network from the storage path and may provide you better I/O performance in the system. As the downside, you'd lose the flexibility of external storage, where you can add additional space and perform maintenance on the fly, transparently to the server.

Every storage subsystem has a “tipping point” after which the latency of I/O requests will start to grow exponentially with increase in throughput and IOPS (I/O operations per second). For example, you may get a 1-millisecond response with an IOPS workload of 1,000 and a 3-millisecond response with an IOPS workload of 50,000. However, you might cross the tipping point at 100,000 IOPS and start to get double-digit or even triple-digit latency.

Every component in the stack will have its own tipping point. For example, low queue depth in the HBA adapter may lead to queuing on the controller level as the number of I/O requests increases. In this case SQL Server will suffer from high latency and poor I/O performance; however, all SAN metrics will be perfectly healthy, with no latency at all.

You can use the DiskSpd utility to test storage subsystem performance. That utility emulates SQL Server's workload in the system. You can download it from [GitHub](#).

As I've noted, you'll need to look at all storage subsystem components when you troubleshoot bad I/O performance. Nevertheless, the place to start is analyzing overall storage latency and the number of data SQL Server reads and writes. You can do this by looking at `sys.dm_io_virtual_file_stats` view.

sys.dm_io_virtual_file_stats view

The `sys.dm_io_virtual_file_stats` view is the most important tool in SQL Server I/O performance troubleshooting. This view provides I/O statistics by database file, including number of I/O operations, amount of data read and written, and information about stalls, or time for I/O requests to complete. (I use the terms latency and stalls interchangeably throughout this book.)

The data in this view is cumulative and is calculated from the time of the SQL Server restart. Take two snapshots of the data and calculate the delta between them (Listing 3-1 shows the code to do that). This code filters out database files with low I/O activity, since their metrics are usually skewed and not very useful.

Example 3-1. Using the `sys.dm_io_virtual_file_stats` view

```
CREATE TABLE #Snapshot
(
    database_id SMALLINT NOT NULL,
    file_id SMALLINT NOT NULL,
    num_of_reads BIGINT NOT NULL,
    num_of_bytes_read BIGINT NOT NULL,
    io_stall_read_ms BIGINT NOT NULL,
    num_of_writes BIGINT NOT NULL,
    num_of_bytes_written BIGINT NOT NULL,
    io_stall_write_ms BIGINT NOT NULL
);

INSERT INTO
#Snapshot(database_id,file_id,num_of_reads,num_of_bytes_read
,io_stall_read_ms,num_of_writes,num_of_bytes_written,io_stall_write
_ms)
    SELECT database_id,file_id,num_of_reads,num_of_bytes_read
,io_stall_read_ms,num_of_writes,num_of_bytes_written,io_stall_write
_ms
    FROM sys.dm_io_virtual_file_stats(NULL,NULL)
OPTION (RECOMPILE);

-- Set test interval (1 minute). Use larger intervals in production
WAITFOR DELAY '00:01:00.000';

;WITH Stats(db_id, file_id, Reads, ReadBytes, Writes
,WrittenBytes, ReadStall, WriteStall)
```

```

as
(
    SELECT
        s.database_id, s.file_id
        ,fs.num_of_reads - s.num_of_reads
        ,fs.num_of_bytes_read - s.num_of_bytes_read
        ,fs.num_of_writes - s.num_of_writes
        ,fs.num_of_bytes_written - s.num_of_bytes_written
        ,fs.io_stall_read_ms - s.io_stall_read_ms
        ,fs.io_stall_write_ms - s.io_stall_write_ms
    FROM
        #Snapshot s JOIN
sys.dm_io_virtual_file_stats(NULL, NULL) fs ON
        s.database_id = fs.database_id and
s.file_id = fs.file_id
)
SELECT
    s.db_id AS [DB ID], d.name AS [Database]
    ,mf.name AS [File Name], mf.physical_name AS [File Path]
    ,mf.type_desc AS [Type], s.Reads
    ,CONVERT(DECIMAL(12,3), s.ReadBytes / 1048576.) AS [Read
MB]
    ,CONVERT(DECIMAL(12,3), s.WrittenBytes / 1048576.) AS
[Written MB]
    ,s.Writes, s.Reads + s.Writes AS [IO Count]
    ,CONVERT(DECIMAL(5,2),100.0 * s.ReadBytes /
        (s.ReadBytes + s.WrittenBytes)) AS [Read %]
    ,CONVERT(DECIMAL(5,2),100.0 * s.WrittenBytes /
        (s.ReadBytes + s.WrittenBytes)) AS [Write
%]
    ,s.ReadStall AS [Read Stall]
    ,s.WriteStall AS [Write Stall]
    ,CASE WHEN s.Reads = 0
        THEN 0.000
        ELSE CONVERT(DECIMAL(12,3),1.0 * s.ReadStall /
s.Reads)
    END AS [Avg Read Stall]
    ,CASE WHEN s.Writes = 0
        THEN 0.000
        ELSE CONVERT(DECIMAL(12,3),1.0 * s.WriteStall /
s.Writes)
    END AS [Avg Write Stall]
FROM
    Stats s JOIN sys.master_files mf WITH (NOLOCK) ON
        s.db_id = mf.database_id and
        s.file_id = mf.file_id
    JOIN sys.databases d WITH (NOLOCK) ON
        s.db_id = d.database_id

```

```
WHERE -- Only display files with more than 20MB throughput
      (s.ReadBytes + s.WrittenBytes) > 20 * 1048576
ORDER BY
      s.db_id, s.file_id
OPTION (RECOMPILE);
```

Figure 3-4 shows the output from the view.

	DB ID	Database	File Name	File Path	Type	Reads	Read MB	Written MB	Writes
1	2	tempdb	tempdev	T:\SQLDa...	ROWS	279	17.375	17.563	281
2	2	tempdb	templog	T:\SQLDa...	LOG	0	0.000	76.727	1311
3	2	tempdb	temp2	T:\SQLDa...	ROWS	278	17.086	17.211	276
4	2	tempdb	temp3	T:\SQLDa...	ROWS	292	17.836	18.000	290
5	2	tempdb	temp4	T:\SQLDa...	ROWS	295	18.031	18.273	294
6	2	tempdb	temp5	T:\SQLDa...	ROWS	307	18.617	18.742	301
7	2	tempdb	temp6	T:\SQLDa...	ROWS	287	17.664	17.914	287
8	2	tempdb	temp7	T:\SQLDa...	ROWS	284	17.625	17.750	284
9	11				ROWS	76	1.578	88.719	7551
10	11				LOG	0	0.000	142.578	31166
11	11				ROWS	1111	9.453	38.422	4420
12	11				ROWS	2505	44.227	345.977	37251

	ID Count	Read %	Write %	Read Stall	Write Stall	Avg Read Stall	Avg Write Stall
1	560	49.73	50.27	205	322	0.735	1.146
2	1311	0.00	100.00	0	1191	0.000	0.908
3	554	49.82	50.18	195	308	0.701	1.116
4	582	49.77	50.23	206	322	0.705	1.110
5	589	49.67	50.33	206	323	0.698	1.099
6	608	49.83	50.17	225	343	0.733	1.140
7	574	49.65	50.35	214	308	0.746	1.073
8	568	49.82	50.18	208	299	0.732	1.053
9	7627	1.75	98.25	60	6847	0.789	0.907
10	31166	0.00	100.00	0	16944	0.000	0.544
11	5531	19.75	80.25	1994	3806	1.795	0.861
12	39756	11.33	88.67	2252	29908	0.899	0.803

Figure 3-4. Sample output from sys.dm_io_virtual_file_stats

The goal is to keep stalls/latency metrics as low as possible. It is impossible to define thresholds that can be applied to all systems, but my rule of thumb is not to exceed 1 to 2-millisecond write stalls for transaction logs and 5 to 7- millisecond read and write stalls for data files when network storage is used. The latency should be even lower, in the sub-millisecond range, when you are using modern direct-attached drives.

Next, analyze throughput in the system. High stalls with low throughput usually indicate performance issues outside of SQL Server. Don't forget to look at throughput across all files that share the same drive or controller. High throughput in some files may impact the metrics in others that share the same resource.

There is usually a correlation between throughput and stalls – the more data you are reading and writing, the higher latency you'll have. This correlation is usually linear until you reach the tipping point, after which latency increases very quickly.

A large amount of reads and read stalls in the data files is often accompanied by a significant percent of PAGEIOLATCH waits and a low Page Life Expectancy performance counter value. This indicates that a large amount of data is constantly being read from disk. You need to understand why that is happening. In most cases, it's due to nonoptimized queries that perform large scans reading data from disk. We will talk about how to detect those queries in the next chapter.

Don't discount the possibility, though, that the server is underprovisioned and doesn't have enough memory to accommodate an active dataset. That is also entirely possible. In either case, adding extra memory may be a completely acceptable solution that will reduce I/O load and improve performance of the system. It is, obviously, not the best solution, but in many cases it's easier and cheaper to use hardware to solve the problem.

In users' databases, large amount of writes and write stalls in data files often indicate inefficient checkpoint configuration. You may get some

improvements by tuning the checkpoint configuration, as I will show later in the chapter. In the longer term, you may need to analyze if it is possible to reduce the number of data pages SQL Server writes to disk. Some ways to do this include removing unnecessary indexes; reducing page splits by changing FILLFACTOR and tuning the index maintenance strategy; decreasing the number of data pages by implementing data compression; and, potentially, refactoring database schema and applications.

When you see large throughput and stalls in tempdb, identify what causes them. The three most common causes are version store activity, massive tempdb spills, and excessive usage of temporary objects. We will talk about these in Chapter 9.

Finally, you can also get an idea of I/O latency by analyzing resource wait time in PAGEIOLATCH and other I/O-related waits. This won't give you detailed information on a per-file basis, but it may be a good metric when you look at systemwide I/O performance.

Performance Counters and OS Metrics

The `sys.dm_io_virtual_file_stats` view provides useful and detailed information and points you in the right direction for further I/O troubleshooting, but it has one limitation: it averages data over the sampling interval.

This is completely acceptable when I/O latency is low. However, if latency numbers are high, you'll want to determine if performance is generally slow or if the numbers have been skewed by some bursts in activity. You can do this by looking at the performance counters correlating SQL Server and disk metrics.

The troubleshooting process will vary slightly between Windows and Linux. In Windows, the simplest way to analyze the metrics is using the well-known PerfMon (Performance Monitor) utility. You can look at the SQL Server and I/O performance counters together and correlate data from them.

*T
a
b
l
e*

*3
-*

*1
·
I
/*

*O
-*

*R
e
l
a
t
e
d*

*P
e
r
f
o
r
m
a
n
c
e*

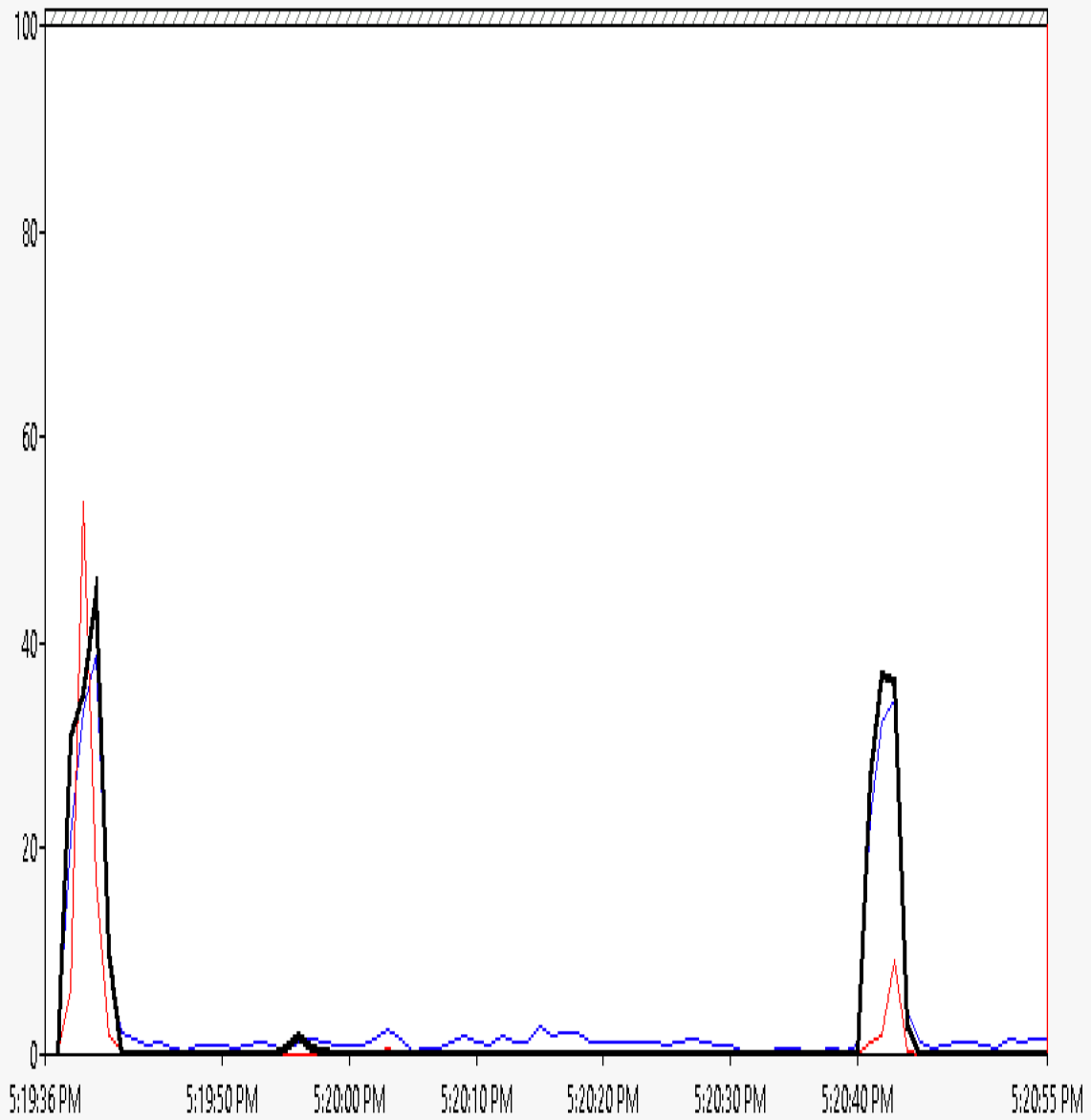
C
o
u
n
t
e
r
s

Performance Object	Performance Counters	Description
Physical Disk	Avg Disk Queue Length	Provides the average number of I/O requests (total, read, and write, respectively) queued during the sampling interval. Those numbers should be as low as possible. Spikes indicate that I/O requests are being queued at the OS level.
	Avg Disk Read Queue Length	
	Avg Disk Write Queue Length	
	Current Disk Queue Length	Gives you the size of the I/O request queue when the metric was collected.
	Avg Disk sec/Transfer	Indicates average latency for disk operations during the sampling interval. These numbers are usually similar to latency/stall metrics from the sys.dm_io_virtual_file_stats view when sampled over the same time period. However, because you typically measure sys.dm_io_virtual_file_stats over larger intervals, these counters will show you if I/O stalls were always high or if data has been affected by latency spikes.
Avg Disk sec/Read		
Avg Disk sec/Write		
Disk Transfers/sec	Displays the number of I/O operations and throughput at the time of the reading. Similar to latency counters, you can use them to analyze the uniformity of the disk workload.	
Disk Reads/sec		
Disk Writes/sec		

	Bytes/sec	Disk	
	Read Bytes/sec	Disk	
	Write Bytes/sec	Disk	
	Avg Disk Bytes/Transfer		Shows the average size of I/O requests, which can help you understand I/O patterns in the system.
	Disk Bytes/Read	Avg	
	Disk Bytes/Write	Avg	
SQL Server: Buffer Manager	Checkpoint pages/sec		Shows the number of dirty pages written by the checkpoint process.
	Background writer pages/sec		
	Lazy writer/sec		Provides number of pages written by lazy writer process
	Page reads/sec		Display the number of physical reads and writes
	Page writes/sec	Page	
	Readahead pages/sec		Shows the number of pages read by read-ahead process.
SQL Server: Databases	Log Bytes Flushed/sec		Provide you the data about throughput, latency and number of write requests for transaction log writes. Use those counters to understand uniformity of log generation when you troubleshoot high log write latency
	Log Flush Write Time (ms)	Log	
	Log Flushes/sec	Log	
SQL Server: SQL Statistics	Batch Requests/sec		While these two counters are not I/O-related, they can be used to analyze spikes in system workload that may lead to bursts in I/O activity.
SQL Server: Databases	Transactions/sec		

Usually, I start by looking at Avg Disk sec/Read and Avg Disk sec/Write latency counters, along with Avg Disk Queue Length. If I see any spikes in their values, I add SQL Server–specific counters to identify what processes may be leading to the bursts in activity.

Figure 3-5 illustrates one such example. You can see the correlation between Checkpoint pages/sec and high Avg Disk sec/Write and Avg Disk Queue Length values. This leads to the simple conclusion that the I/O subsystem cannot keep up with bursts of writes from the checkpoint process.



Last 0.000 Average 2,910.771 Minimum 0.000 Maximum 44,773.369 Duration 1:20

Show	Color	Scale	Counter	Instance	Parent	Object	Computer
<input checked="" type="checkbox"/>	Red	0.1	Avg. Disk Queue Length	_Total	---	PhysicalDisk	
<input checked="" type="checkbox"/>	Green	0.001	Checkpoint pages/sec	---	---	SQLServer:Buffer Manager	
<input checked="" type="checkbox"/>	Blue	0.001	Disk Writes/sec	_Total	---	PhysicalDisk	

Figure 3-5. Checkpoint and disk queueing

Pay attention to other applications installed on the server – it is possible that they are responsible for I/O activity bursts or other issues.

Linux doesn't offer the standard PerfMon utility; however, there are plenty of free and commercial monitoring tools available. You can also use tools like iostat, dstat, and iotop, which are included in major Linux distributions. They provide general disk performance metrics on a per-process or system level.

On the SQL Server side, you can access performance counters through `sys.dm_os_performance_counters` view. Listing 3-2 shows you how to do that.

Example 3-2. Using sys.dm_os_performance_counters view

```
CREATE TABLE #PerfCnters
(
    collected_time DATETIME2(7) NOT NULL DEFAULT SYSDATETIME(),
    object_name SYSNAME NOT NULL,
    counter_name SYSNAME NOT NULL,
    instance_name SYSNAME NOT NULL,
    cntr_value BIGINT NOT NULL,
    PRIMARY KEY (object_name, counter_name, instance_name)
);

;WITH Counters(obj_name, ctr_name)
AS
(
    SELECT C.obj_name, C.ctr_name
    FROM
    (
        VALUES
        ('SQLServer:Buffer Manager','Checkpoint
pages/sec')
        , ('SQLServer:Buffer Manager','Background
writer pages/sec')
        , ('SQLServer:Buffer Manager','Lazy
writes/sec')
        , ('SQLServer:Buffer Manager','Page
reads/sec')
        , ('SQLServer:Buffer Manager','Page
writes/sec')
```

```

pages/sec')
, ('SQLServer:Buffer Manager', 'Readahead
-- For all DBs
, ('SQLServer:Databases', 'Log Flushes/sec')
Flushed/sec') -- For all DBs
, ('SQLServer:Databases', 'Log Bytes
Time (ms)') -- For all DBs
, ('SQLServer:Databases', 'Log Flush Write
-- For all DBs
, ('SQLServer:Databases', 'Transactions/sec')
Requests/sec')
, ('SQLServer:SQL Statistics', 'Batch
) C(obj_name, ctr_name)
)
INSERT INTO
#PerfCntrs(object_name,counter_name,instance_name,cntr_value)
SELECT
pc.object_name, pc.counter_name, pc.instance_name,
pc.cntr_value
FROM
sys.dm_os_performance_counters pc WITH (NOLOCK)
JOIN Counters c ON
pc.counter_name = c.ctr_name AND
pc.object_name = c.obj_name;

WAITFOR DELAY '00:00:01.000';

;WITH Counters(obj_name, ctr_name)
AS
(
SELECT C.obj_name, C.ctr_name
FROM
(
VALUES
('SQLServer:Buffer Manager', 'Checkpoint
pages/sec')
, ('SQLServer:Buffer Manager', 'Background
writer pages/sec')
, ('SQLServer:Buffer Manager', 'Lazy
writes/sec')
, ('SQLServer:Buffer Manager', 'Page
reads/sec')
, ('SQLServer:Buffer Manager', 'Page
writes/sec')
, ('SQLServer:Buffer Manager', 'Readahead
pages/sec')
, ('SQLServer:Databases', 'Log Flushes/sec')
-- For all DBs

```

```

        , ('SQLServer:Databases','Log Bytes
Flushed/sec') -- For all DBs
        , ('SQLServer:Databases','Log Flush Write
Time (ms)') -- For all DBs
        , ('SQLServer:Databases','Transactions/sec')
-- For all DBs
        , ('SQLServer:SQL Statistics','Batch
Requests/sec')
    ) C(obj_name, ctr_name)
)
SELECT
    pc.object_name, pc.counter_name, pc.instance_name
    ,CASE pc.cntr_type
        WHEN 272696576 THEN
            (pc.cntr_value - h.cntr_value) * 1000 /
DATEDIFF(MILLISECOND,h.collected_time,SYSDATETIME())
        WHEN 65792 THEN
            pc.cntr_value
        ELSE NULL
    END as cntr_value
FROM
    sys.dm_os_performance_counters pc WITH (NOLOCK) JOIN
Counters c ON
        pc.counter_name = c.ctr_name AND pc.object_name =
c.obj_name
    JOIN #PerfCnters h ON
        pc.object_name = h.object_name AND
        pc.counter_name = h.counter_name AND
        pc.instance_name = h.instance_name
ORDER BY
    pc.object_name, pc.counter_name, pc.instance_name
OPTION (RECOMPILE);

```

You can also bring `sys.dm_io_virtual_file_stats` view to the analysis, sampling its data and performance counters together every second. The approach is the same one we just discussed – you’ll look at the correlation between disk latency and activity and evaluate the general performance of the I/O subsystem, identifying tipping points in the load.

Virtualization, HBA, and Storage Layers

There are several layers in the storage stack you may need to analyze in addition to OS. They include virtualization, HBA/SCSI controller

configuration, and the storage array itself.

I recommend working together with infrastructure and storage engineers during troubleshooting.

SQL Server mostly operates in shared environments. It shares storage and network infrastructure with other clients, and when virtualized, it runs on the same physical host with other VMs. As I said earlier in this book, when virtualization is being used, be sure to validate that the host is not overcommitted, which could lead to all sorts of performance issues.

Unless you have a very simple SQL Server setup that uses local storage, I/O requests will be serialized and sent over network. There are two typical problems here: insufficient queue depth and noisy neighbors.

Insufficient queue depth

The first is insufficient queue depth somewhere in the I/O path.

Unfortunately, the default query depth may not be enough for a highly demanding I/O workload. You'll need to check and potentially increase it in the datastore, vSCSI controller, and HBA adapter settings. The typical sign of insufficient queue depth is low latency on the storage combined with much higher latency in VM and/or OS, with disk queueing present.

Noisy neighbors

The second problem is noisy neighbors. Multiple I/O intensive VMs running on the same host may affect each other. Similarly, multiple high-throughput servers sharing the same network and storage may overload them. Unfortunately, troubleshooting noisy neighbor problem is never easy and you need to analyze multiple components in the infrastructure to detect it.

A word of caution – storage arrays can handle a limited number of outstanding requests. Increasing queue depth on a busy server could increase the number of outstanding requests on the storage. You might shift the bottleneck from the server to the storage layer, especially if the storage serves requests from many busy systems.

The virtualization host and storage both expose throughput, IOPS, and latency metrics for analysis. On virtualization layers, the metrics may vary based on technology. For example, in Hyper-V you can use regular disk performance counters on the host. In VMWare, you can get the data from ESXTOP utility. In either case, the troubleshooting approach is very similar to what we have already discussed. Look at the available metrics, correlate data from them, and detect the bottlenecks in the I/O path.

Finally, check the storage configuration. Storage vendors usually publish best practices for SQL Server workloads: they are a good starting point. Pay attention to the allocation unit size's alignment with the raid stripe size and partition offset, though.

For example, a 1024 MB partition offset, 4 KB disk block, 64 KB allocation unit, and 128 KB raid stripes are perfectly aligned, with each I/O request served by a single disk. On the other hand, 96 KB raid stripes will spread 64 KB allocation units across two disks, which leads to extra I/O requests and can seriously impact performance.

Again, it is always beneficial to work together with infrastructure and storage engineers. They are the subject matter experts and may help you to find the root cause of the problem faster than when you are working alone.

Finally, the best approach to get predictable performance in critical systems is to use a dedicated environment. Run SQL Server on dedicated hardware with direct-attached storage (DAS) to get the best performance possible.

Checkpoint Tuning

As we all know, SQL Server uses write-ahead logging. Transactions are considered to be committed only after the log records are hardened in the

transaction logs. SQL Server does not need to save dirty data pages to disk at the same time – it can reapply the changes by replaying log records if needed.

The checkpoint process saves data pages into the data files. The main goal of checkpoint is reducing recovery time in event of an SQL crash or failover: the fewer changes need to be replayed, the faster recovery will be. The maximum desired recovery time is controlled at either the server level or the database level. By default, both of them are 60 seconds.

NOTE

You should not consider the recovery target to be a hard value. In many cases, the database will recover much faster than that. It is also possible for bursts of activity and long running transactions to prolong recovery beyond the target time.

There are four different types of checkpoints:

Internal checkpoints

Internal checkpoints occur during some SQL Server operations, such as starting database backup or creating a database snapshot.

Manual checkpoints

Manual checkpoint occur manually, as the name indicates, when users trigger them with the CHECKPOINT command.

Automatic checkpoint

Historically, SQL Server used automatic checkpoints, with the recovery interval controlled at the server level. The checkpoint process wakes up once or few times each recovery interval and flushes dirty data pages to disk. Unfortunately, this approach can lead to bursts of data writes, which can be problematic in busy systems.

Indirect checkpoint

Starting with SQL Server 2012, you have another option: indirect checkpoint. With this method, SQL Server tries to balance I/O load by executing checkpoints much more frequently – in some cases, even continuously. This helps to mitigate bursts of data writes, making the I/O load much more balanced. Use it instead of automatic checkpoint whenever possible. Indirect checkpoint is controlled on a per-database basis and enabled by default in databases created in SQL Server 2016 and above. However, SQL Server does not enable indirect checkpoint automatically when you upgrade an SQL Server instance, or in SQL Server 2012 and 2014. You can do it manually by setting up a recovery target at the database level with the `ALTER DATABASE SET TARGET_RECOVERY_TIME` command.

Let me show you an example from one system I worked with. The sample of data from `sys.dm_io_virtual_file_stats` view over 1 minute had very high write latency for the data files. However, the smaller samples (1 to 3 seconds) rarely showed any activity at all.

Figure 3-6 shows the data, with the 1-minute sample at the top and the 1-second sample at the bottom.

	File Path	Reads	Read MB	Written MB	Writes
1	_BAK.NDF	3098	91.867	1404.727	138564
2	_Indexes.NDF	47398	2104.140	9648.797	1113960
	IO Count	Read Stall	Write Stall	Avg Read Stall	Avg Write Stall
	141662	11708	8287158	3.779	59.807
	1161098	466594	74303099	9.773	66.750

	File Path	Reads	Read MB	Written MB	Writes
1	_BAK.NDF	14	0.641	0.000	0
2	_Indexes.NDF	552	27.399	0.000	0
	IO Count	Read Stall	Write Stall	Avg Read Stall	Avg Write Stall
	14	14	1.000	0.000	0.000
	552	471	0.000	0.853	0.000

Figure 3-6. Sample sys.dm_io_virtual_file_stats with automatic checkpoint

This behavior led me to believe that the issue was related to checkpoint. I confirmed this hypothesis by looking at the Checkpoint pages/sec, Disk Writes/sec, and Avg Disk Queue Length performance counters. You can clearly see that burst of disk writes from the checkpoint process in Figure 3-5 earlier in the chapter, which shows the screenshot from PerfMon.

Although this instance ran SQL Server 2016, it used automatic checkpoint, because all databases had been upgraded from the earlier version of SQL Server. Enabling indirect checkpoint in the system immediately changed the I/O pattern, making it much more balanced.

You can see the performance counters in Figure 3-7. Notice that with indirect checkpoint, you should use Background writer pages/sec instead of the Checkpoint pages/sec counter.

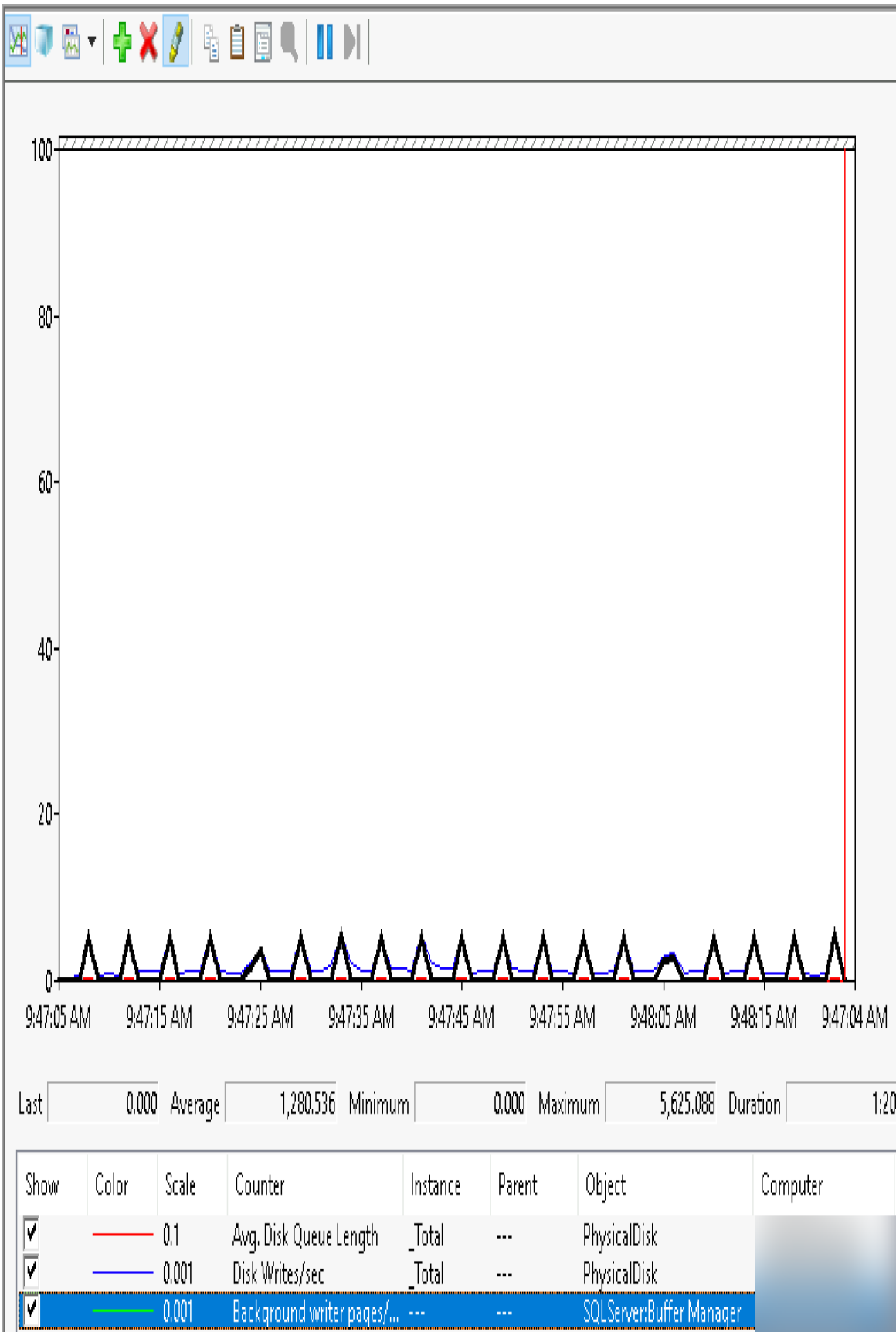


Figure 3-7. Indirect checkpoint performance counters

Figure 3-8 shows the output from a 1-minute sample in the `sys.dm_io_virtual_file_stats` view. As you can see, latency went back to normal.

	File Path	Reads	Read MB	Written MB	Writes
1	N:\Data3\AAD_BAK.NDF	2987	89.121	1367.542	135039
2	N:\Data3\AAD_Indexes.NDF	48212	2140.134	10944.122	1182783

IO Count	Read Stall	Write Stall	Avg Read Stall	Avg Write Stall
138026	3065	264675	1.026	1.960
1230995	66195	2280405	1.373	1.928

Figure 3-8. Sample `sys.dm_io_virtual_file_stats` with indirect checkpoint

Indirect checkpoints do not completely eliminate I/O bursts. You can still have them, especially if the system has some spikes in data modifications. However, they are less frequent than with automatic checkpoints.

You may also need to tune the recovery target to get the most balanced I/O load. In the case above, I got the best results with a 90-second target. Of course, high values may increase recovery time in the system.

I/O Waits

SQL Server uses several different wait types related to I/O operations. It is very common to see all of them present when the disk subsystem is not fast enough. Let's look at five of the most common:

`ASYNC_IO_COMPLETION`, `IO_COMPLETION`, `WRITELOG`, `WRITE_COMPLETION`, and `PAGEIOLATCH`.

ASYNC_IO_COMPLETION waits

This wait type occurs when SQL Server waits for asynchronous I/O operations (read or write) for non-buffer pool pages to complete. Examples include:

- Internal checkpoint when you start database backup or run DBCC CHECKDB
- Reading GAM pages from data files
- Reading data pages from database during database backup. (Unfortunately, this tends to skew the average wait time, making it harder to analyze.)

When I see significant presence of both ASYNC_IO_COMPLETION and PAGEIOLATCH waits in the system, I perform general I/O troubleshooting. If PAGEIOLATCH waits are not present, I look at how often ASYNC_IO_COMPLETION occurs. I may ignore that wait if its percentage is not very significant and disk latency is low.

IO_COMPLETION waits

The IO_COMPLETION wait type occurs during synchronous reads and writes in data files and during some read operations in transaction log. A few examples:

- Reading allocation map pages from the database
- Reading the transaction log during database recovery
- Writing data to tempdb during sort spills

When you see significant percentages of this wait in the system, perform general disk-performance troubleshooting. Pay specific attention to tempdb latency and throughput; in my experience, bad tempdb performance is the most common reason for this wait. We will talk more about tempdb troubleshooting in Chapter 9.

WRITELOG waits

As you can guess by the name, this wait occurs when SQL Server writes log records to the transaction log. It is normal to see this wait in any system; however, a large percentage may indicate a transaction-log bottleneck.

Look at average wait time and transaction log write latency in the `sys.dm_io_virtual_file_stats` view during troubleshooting. High numbers are impactful and may affect throughput in the system.

In addition to optimizing disk subsystem throughput, there are several other things you can do to reduce that wait. We will discuss them in Chapter 11.

WRITE_COMPLETION waits

This wait occurs during synchronous write operations in database and log files. In my experience, it is most common with database snapshots.

SQL Server maintains snapshot databases by persisting versions of data pages that existed at time the snapshot was created. At checkpoints after the snapshot was created, SQL Server reads old copies of data pages from data files and saves them into the snapshot before saving dirty pages to disk.

This can significantly increase the amount of I/O in the system.

When you see this wait in the system, check if there are database snapshots present. Remember that some internal processes, like DBCC CHECKDB, also create internal database snapshots.

When snapshots are present and their usage is legitimate, you may need to analyze how to improve disk performance to support them. In other cases, you may need to remove them from the system if storage cannot keep up.

PAGEIOLATCH waits

As you already know, PAGEIOLATCH waits occur when SQL Server reads data pages from disk. Those waits are very common and are present in any system. Technically, there are six such waits, but only three are typically present in the system:

PAGEIOLATCH_EX

Occurs when the worker wants to update the data page and is waiting for the page to be read from disk to the buffer pool.

PAGEIOLATCH_SH

Occurs when the worker wants to read the data page and is waiting for the page to be read from disk to the buffer pool.

PAGEIOLATCH_UP

Occurs when the worker wants to update a system page (for example, the allocation map) and is waiting for the page to be read from disk to the buffer pool.

Excessive amounts of PAGEIOLATCH waits show that SQL Server is constantly reading data from disk. This usually occurs under two conditions. The first is an underprovisioned SQL Server: when the active data does not fit into the memory. Second, and more often, it indicates the presence of nonoptimized queries that scan unnecessary data, flushing the contents of the buffer pool.

You can cross-check the data by looking at the Page Life Expectancy performance counter, which shows how long data pages stay in the buffer pool. As a baseline, you can generally use the value of 300 seconds per 4 GB of buffer pool memory: for example, 7,500 seconds on the server with 100 GB buffer pool.

You can see the value of Page Life Expectancy in the PerfMon utility or with the `sys.dm_os_performance_counters` view, as shown in Listing 3-3. It also returns values for individual NUMA nodes in the system.

Example 3-3. Getting Page Life Expectancy in the system

```
SELECT object_name, counter_name, instance_name, cntr_value as
```

```
[PLE(sec)]
FROM sys.dm_os_performance_counters WITH (NOLOCK)
WHERE counter_name = 'Page life expectancy'
OPTION (RECOMPILE);
```

A large percentage of PAGEIOLATCH waits always requires troubleshooting. While it does not always introduce customer-facing problems, especially with low-latency flash-based disk arrays, the data growth may push the disk subsystem over the limit, which can become a problem that quickly affects the entire system.

You can reduce the impact of PAGEIOLATCH waits by upgrading the disk subsystem or adding more memory to the server. However, the best approach is reducing the amount of data to read from disk by detecting and optimizing inefficient queries. We'll look at how to detect those queries in the next chapter.

Summary

SQL Server uses cooperative scheduling and, in the majority of cases, asynchronous I/O when it reads and writes data. By default, each scheduler has its own I/O queue and handles I/O in the system.

The `sys.dm_io_virtual_file_stats` view provides I/O throughput and latency metrics per database file. In a properly tuned system, the latency of transaction log writes should not exceed 1 to 2 milliseconds, and the latency of reads and writes to data files should not exceed 5 to 7 milliseconds with network-based storage and should be even lower with DAS.

Look at the entire I/O stack when troubleshooting bad I/O performance. The problem may be anywhere – in the OS, virtualization, network path, or storage layers.

In many cases, high I/O latency is introduced by bursts in I/O activity. Analyze and tune the checkpoint process – it is one of the most common offenders in busy systems.

In many cases, reducing disk activity will help you improve disk latency and system performance. Query optimization is one of the best ways to achieve that. We will look at how to detect non-optimized queries in the system in the next chapter.

Troubleshooting Checklist

Troubleshoot the following:

- Analyze disk subsystem latency with the `sys.dm_io_virtual_file_stats` view
- Check if high latency is caused by bursts in I/O activity by analyzing SQL Server and OS performance counters.
- Review I/O metrics at the VM and storage levels, paying attention to noisy neighbors in your setup.
- Check disk queue depth settings in the I/O stack.
- Troubleshoot SQL Server checkpoint performance and switch to indirect checkpoints.
- Troubleshoot log performance if you see significant `WRITELOG` waits (see Chapter 11).
- Troubleshoot tempdb performance if you see significant `IO_COMPLETION` waits and high tempdb usage and latency (see Chapter 9).
- Detect and optimize inefficient queries if you see high `PAGEIOLATCH` waits in the system.

Chapter 4. Detecting Inefficient Queries

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 4 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

Inefficient queries exist in every system. They impact performance in many ways, most notably by increasing I/O load, CPU usage and blocking in the system. It is essential to detect and optimize them. We’ll start with detection in this chapter, then move on to optimization strategies in subsequent chapters.

This chapter discusses inefficient queries and their potential impact on your system and provides guidelines for detecting them, starting with an approach that uses plan cache-based execution statistics. Next, I will talk about Extended Events and SQL Traces, and then cover Query Store. I’ll wrap up the chapter by sharing a few thoughts on third-party monitoring tools.

The Impact of Inefficient Queries

During my career as a database engineer, I have yet to see a system that wouldn't benefit from query optimization. I'm sure they exist – after all, no one calls me in to look at perfectly healthy systems. Nevertheless, those are few and far between, and there are always opportunities to improve and optimize.

Not every company prioritizes query optimization, though. It's a time-consuming and tedious process, and in many cases it's cheaper, given the benefits of speeding up development and time-to-market, to throw hardware at the problem than to invest hours in performance tuning.

At some point, however, that approach leads to scalability issues. Poorly optimized queries impact systems from many angles, but perhaps the most obvious is disk performance. If the I/O subsystem cannot keep up with the load of large scans, the performance of your entire system will suffer.

You can mask this problem, to a degree, by adding more memory to the server. This increases the size of the buffer pool and allows SQL Server to cache more data, reducing physical I/O. As amount of data in the system grows over time, however, this approach may become impractical or even impossible—especially in non-Enterprise editions of SQL Server that restrict the maximum buffer pool size.

Another effect to watch for is that nonoptimized queries burn CPU on the servers. The more data you process, the more CPU resources you consume. A server might spend just a few microseconds per logical read and in-memory data-page scan, but that quickly adds up as the number of reads increases.

Again, you can mask this by adding more CPUs to the server. (Note, however, that you will need to pay for additional licenses. In non-Enterprise editions, expect a cap on the number of CPUs.) Moreover, adding CPUs may not always solve the problem – nonoptimized queries will still contribute to blocking in the system. While there are ways to reduce blocking without performing query tuning, this can change system behavior and has performance implications.

The bottom line is: When you troubleshoot a system, always analyze whether queries in the system are poorly optimized. Once you've done that, estimate the impact of those inefficient queries.

While query optimization always benefits a system, it is not always simple, nor does it always provide the best ROI for your efforts. More often than not, you will at least need to tune some queries.

To put things in perspective: I perform query tuning when I see high disk throughput, blocking, or high CPU load in the system. However, I may initially focus my efforts elsewhere if data is cached in the buffer pool and the CPU load is acceptable. I have to be careful and think about data growth, though – it is possible that active data will one day outgrow the buffer pool, which could lead to sudden and serious performance issues.

Fortunately, query optimization does not require an all-or-nothing approach! You can achieve dramatic performance improvements by optimizing just a handful of frequently executed queries. Let's look at a few methods of how we can detect them.

Plan-Cache-Based Execution Statistics

In most cases, SQL Server caches and reuses execution plans for queries. For each plan in the cache, it also maintains execution statistics, including the number of times the query ran, cumulative CPU time, and I/O load. You can use this information to quickly pinpoint the most resource-intensive queries for optimization. (We will discuss plan caching in more details in Chapter 6.)

Analyzing plan cache-based execution statistics is not the most comprehensive detection technique; it has quite a few limitations. Nevertheless, it is very easy to use and, in many cases, *good enough*. It works in all versions of SQL Server and it is always present in the system. You don't need to set up any additional monitoring to collect the data.

NOTE

The code calls the `sys.dm_exec_query_plan` function for each cached plan in the system. This is a CPU-intensive operation, so remove it if your server is CPU-bound. You may also need to comment some of the columns in the statement, depending on the version and patching level of your SQL Server instance.

You can get execution statistics using the `sys.dm_exec_query_stats` view (as shown in Listing 4-1). The query there is a bit simplistic, but it demonstrates the view in action and shows you the list of metrics exposed in the view. We will use it to build a more sophisticated version of the code later in the chapter.

Example 4-1. Using the `sys.dm_exec_query_stats` view

```
SELECT TOP 50
    qs.creation_time AS [Cached Time]
    ,qs.last_execution_time AS [Last Exec Time]
    ,SUBSTRING(qt.text, (qs.statement_start_offset/2)+1,
    ((
        CASE qs.statement_end_offset
            WHEN -1 THEN DATALENGTH(qt.text)
            ELSE qs.statement_end_offset
        END - qs.statement_start_offset)/2)+1) AS SQL
    ,qp.query_plan AS [Query Plan]
    ,qs.execution_count AS [Exec Cnt]
    ,CONVERT(DECIMAL(10,5),
        IIF(datediff(second,qs.creation_time,
qs.last_execution_time) = 0,
        NULL,
        1.0 * qs.execution_count /
            datediff(second,qs.creation_time,
qs.last_execution_time)
        )
    ) AS [Exec Per Second]
    ,(qs.total_logical_reads + qs.total_logical_writes) /
        qs.execution_count AS [Avg IO]
    ,(qs.total_worker_time / qs.execution_count / 1000)
        AS [Avg CPU(ms)]
    ,qs.total_logical_reads AS [Total Reads]
    ,qs.last_logical_reads AS [Last Reads]
    ,qs.total_logical_writes AS [Total Writes]
    ,qs.last_logical_writes AS [Last Writes]
    ,qs.total_worker_time / 1000 AS [Total Worker Time]
```

```

,qs.last_worker_time / 1000 AS [Last Worker Time]
,qs.total_elapsed_time / 1000 AS [Total Elapsed Time]
,qs.last_elapsed_time / 1000 AS [Last Elapsed Time]
,qs.total_rows AS [Total Rows]
,qs.last_rows AS [Last Rows]
,qs.total_rows / qs.execution_count AS [Avg Rows]
,qs.total_physical_reads AS [Total Physical Reads]
,qs.last_physical_reads AS [Last Physical Reads]
,qs.total_physical_reads / qs.execution_count
    AS [Avg Physical Reads]
,qs.total_grant_kb AS [Total Grant KB]
,qs.last_grant_kb AS [Last Grant KB]
,(qs.total_grant_kb / qs.execution_count)
    AS [Avg Grant KB]
,qs.total_used_grant_kb AS [Total Used Grant KB]
,qs.last_used_grant_kb AS [Last Used Grant KB]
,(qs.total_used_grant_kb / qs.execution_count)
    AS [Avg Used Grant KB]
,qs.total_ideal_grant_kb AS [Total Ideal Grant KB]
,qs.last_ideal_grant_kb AS [Last Ideal Grant KB]
,(qs.total_ideal_grant_kb / qs.execution_count)
    AS [Avg Ideal Grant KB]
,qs.total_columnstore_segment_reads
    AS [Total CSI Segments Read]
,qs.last_columnstore_segment_reads
    AS [Last CSI Segments Read]
,(qs.total_columnstore_segment_reads / qs.execution_count)
    AS [AVG CSI Segments Read]
,qs.max_dop AS [Max DOP]
,qs.total_spills AS [Total Spills]
,qs.last_spills AS [Last Spills]
,(qs.total_spills / qs.execution_count) AS [Avg Spills]
FROM
    sys.dm_exec_query_stats qs WITH (NOLOCK)
        CROSS APPLY sys.dm_exec_sql_text(qs.sql_handle) qt
        CROSS APPLY sys.dm_exec_query_plan(qs.plan_handle) qp
ORDER BY
    [Avg IO] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

You will likely sort data differently based on your tuning goals: by I/O when you need to reduce disk load; by CPU on CPU-bound systems, and so on.

Figure 4-1 shows a partial output of the query from one of the servers. As you can see, it is very easy to choose queries to optimize based on frequency of query executions and resource consumption data in the output.

	Cached Time	Last Exec Time	SQL	Query Plan	Exec Cnt
1	2021-01-12 15:30:57.623	2021-01-12 18:11:05.270	UPDATE ls...	<ShowPlanXML...	18
2	2021-01-12 16:19:38.383	2021-01-12 18:14:00.007	merge bi...	<ShowPlanXML...	13
3	2021-01-12 17:49:29.357	2021-01-12 17:49:41.067	SELECT th...	<ShowPlanXML...	1
4	2021-01-11 12:02:01.930	2021-01-12 18:11:04.333	UPDATE bi...	<ShowPlanXML...	6
5	2021-01-12 15:25:10.940	2021-01-12 18:15:42.060	select * ...	<ShowPlanXML...	19
6	2021-01-03 09:36:06.373	2021-01-03 09:36:22.527	SELECT th...	<ShowPlanXML...	1
7	2021-01-12 18:12:11.743	2021-01-12 18:12:11.770	select * ...	<ShowPlanXML...	1

	Exec Per Sec..	Avg IO	Avg CPU(ms)	Total Reads	Last Reads	Total Writes
	0.00187	38737113	51757	696209779	38690449	1058256
	0.00189	23303547	36656	302680027	23289669	266093
	0.08333	19881441	183527	19836720	19836720	44721
	0.00006	13916838	44736	83500847	13917209	182
	0.00186	10609666	36899	201079410	10584684	504246
	0.06250	8505727	157595	8491099	8491099	14628
	NULL	7271806	8437	7271796	7271796	10

Figure 4-1. Fig 4-1 Partial output from sys.dm_exec_query_stats view

The execution plans you get in the output do not include actual execution metrics. In this respect, they are similar to estimated execution plans. You'll need to take this into consideration during optimization (we'll talk more about that in chapter 5).

There are several other important limitations to remember.

First and foremost, you won't see any data for the queries that do not have execution plans cached. You may miss some infrequently executed queries with plans evicted from the cache. Usually, this is not a problem – infrequently executed queries rarely need to be optimized at the beginning of tuning.

There is another possibility, however. SQL Server won't cache execution plans if you are using statement-level recompile or executing stored procedures with a RECOMPILE clause. You need to capture those queries using Query Store or Extended Events, which we will discuss later in the chapter.

The second problem is related to how long plans stay cached. This varies by plan, which may skew the results when you sort data by *total* metrics. For example, a query with lower *average* CPU time may show a higher *total* number of executions and CPU time than a query with higher *average* CPU time, depending on the time when both plans were cached.

You can look at the `creation_time` and `last_execution_time` columns, which show the last time when plans were cached and executed, respectively. I usually look at the data sorted based on both *total* and *average* metrics, taking the frequency of executions into consideration.

The final problem is more complicated: it is possible to get multiple results for the same or similar queries. This can happen with ad-hoc workloads, with clients that have different SET settings in their sessions, when users run the same queries with slightly different formatting, or in many other cases.

Fortunately, you can address that problem by using two columns, `query_hash` and `query_plan_hash`, both exposed in the `sys.dm_exec_query_stats` view. The same values in those columns would

indicate similar queries and execution plans. You can use those columns to aggregate data.

WARNING

The DBCC FREEPROCCACHE statement clears the plan cache to reduce the size of the output in the demo. Do not run it on production servers!

Let me demonstrate with a simple example. Listing 4-2 runs three queries and then examines the content of the plan cache. The first two are the same—they just have different formatting. The third one is different.

Example 4-2. Query_hash and query_plan_hash in action

```
DBCC FREEPROCCACHE -- Do not run in production!
GO
SELECT /*V1*/ TOP 1 object_id FROM sys.objects WHERE object_id = 1;
GO
SELECT /*V2*/ TOP 1 object_id
FROM sys.objects
WHERE object_id = 1;
GO
SELECT COUNT(*) FROM sys.objects
GO
SELECT
    qs.query_hash, qs.query_plan_hash, qs.sql_handle,
    qs.plan_handle,
    SUBSTRING(qt.text, (qs.statement_start_offset/2)+1,
        ((
            CASE qs.statement_end_offset
                WHEN -1 THEN DATALENGTH(qt.text)
                ELSE qs.statement_end_offset
            END - qs.statement_start_offset)/2)+1
        ) as SQL
FROM
    sys.dm_exec_query_stats qs
    CROSS APPLY sys.dm_exec_sql_text(qs.sql_handle) qt
ORDER BY query_hash
OPTION (MAXDOP 1, RECOMPILE);
```

You can see the results in Figure 4-2. There are three execution plans in the output. The last two rows have the same query_hash and query_plan_hash and different sql_handle and plan_handle values.

	query_hash	query_plan_hash	sql_handle
1	0x927A56E8332165F5	0xF5E14111A9DA50EE	0x02000000983965123CEF4131822A91507E192F...
2	0xB0B66B1CD33AF4F7	0x3EE0B044F8702057	0x020000008B28A425CDCE657B1DA2E89083EB2F...
3	0xB0B66B1CD33AF4F7	0x3EE0B044F8702057	0x02000000E145B236E584443C694C2299F79EE5...

plan_handle	SQL
0x0600010098396512003120340200000001000000000000...	SELECT COUNT(*) FROM sys.objects
0x060001008B28A425C02920340200000001000000000000...	SELECT /*V2*/ TOP 1 object_id FROM ...
0x06000100E145B236802220340200000001000000000000...	SELECT /*V1*/ TOP 1 object_id FROM s...

Figure 4-2. Fig. 4-2. Multiple plans with the same query_hash and query_plan_hash

Listing 4-3 provides a more sophisticated version of the script from Listing 4-1 by aggregating statistics from similar queries. The statement and execution plans are picked up randomly from the first query in each group, so factor that into your analysis.

Example 4-3. Using the sys.dm_exec_query_stats view with query_hash aggregation

```

;WITH Data
AS
(
    SELECT TOP 50
        qs.query_hash
        ,COUNT(*) as [Plan Count]
        ,MIN(qs.creation_time) AS [Cached Time]
        ,MAX(qs.last_execution_time) AS [Last Exec Time]
        ,SUM(qs.execution_count) AS [Exec Cnt]
        ,SUM(qs.total_logical_reads) AS [Total Reads]
        ,SUM(qs.total_logical_writes) AS [Total Writes]
        ,SUM(qs.total_worker_time / 1000) AS [Total Worker Time]
        ,SUM(qs.total_elapsed_time / 1000) AS [Total Elapsed Time]
        ,SUM(qs.total_rows) AS [Total Rows]

```

```

    ,SUM(qs.total_physical_reads) AS [Total Physical Reads]
    ,SUM(qs.total_grant_kb) AS [Total Grant KB]
    ,SUM(qs.total_used_grant_kb) AS [Total Used Grant KB]
    ,SUM(qs.total_ideal_grant_kb) AS [Total Ideal Grant KB]
    ,SUM(qs.total_columnstore_segment_reads)
        AS [Total CSI Segments Read]
    ,MAX(qs.max_dop) AS [Max DOP]
    ,SUM(qs.total_spills) AS [Total Spills]
FROM
    sys.dm_exec_query_stats qs WITH (NOLOCK)
GROUP BY
    qs.query_hash
ORDER BY
    SUM((qs.total_logical_reads + qs.total_logical_writes) /
        qs.execution_count) DESC
)
SELECT
    d.[Cached Time]
    ,d.[Last Exec Time]
    ,d.[Plan Count]
    ,sql_plan.SQL
    ,sql_plan.[Query Plan]
    ,d.[Exec Cnt]
    ,CONVERT(DECIMAL(10,5),
        IIF(datediff(second,d.[Cached Time], d.[Last Exec Time]) =
0,
            NULL,
            1.0 * d.[Exec Cnt] /
                datediff(second,d.[Cached Time], d.[Last Exec Time])
        )
    ) AS [Exec Per Second]
    ,(d.[Total Reads] + d.[Total Writes]) / d.[Exec Cnt] AS [Avg IO]
    ,(d.[Total Worker Time] / d.[Exec Cnt] / 1000) AS [Avg CPU(ms)]
    ,d.[Total Reads]
    ,d.[Total Writes]
    ,d.[Total Worker Time]
    ,d.[Total Elapsed Time]
    ,d.[Total Rows]
    ,d.[Total Rows] / d.[Exec Cnt] AS [Avg Rows]
    ,d.[Total Physical Reads]
    ,d.[Total Physical Reads] / d.[Exec Cnt] AS [Avg Physical Reads]
    ,d.[Total Grant KB]
    ,d.[Total Grant KB] / d.[Exec Cnt] AS [Avg Grant KB]
    ,d.[Total Used Grant KB]
    ,d.[Total Used Grant KB] / d.[Exec Cnt] AS [Avg Used Grant KB]
    ,d.[Total Ideal Grant KB]
    ,d.[Total Ideal Grant KB] / d.[Exec Cnt] AS [Avg Ideal Grant KB]
    ,d.[Total CSI Segments Read]

```

```

    ,d.[Total CSI Segments Read] / d.[Exec Cnt] AS [AVG CSI Segments
Read]
    ,d.[Max DOP]
    ,d.[Total Spills]
    ,d.[Total Spills] / d.[Exec Cnt] AS [Avg Spills]
FROM
    Data d
    CROSS APPLY
    (
        SELECT TOP 1
            SUBSTRING(qt.text, (qs.statement_start_offset/2)+1,
                ((
                    CASE qs.statement_end_offset
                        WHEN -1 THEN DATALENGTH(qt.text)
                        ELSE qs.statement_end_offset
                    END - qs.statement_start_offset)/2)+1
            ) AS SQL
        ,qp.query_plan AS [Query Plan]
        FROM
            sys.dm_exec_query_stats qs
            CROSS APPLY sys.dm_exec_sql_text(qs.sql_handle)
        qt
            CROSS APPLY
            sys.dm_exec_query_plan(qs.plan_handle) qp
        WHERE
            qs.query_hash = d.query_hash AND ISNULL(qt.text, '')
    ) sql_plan
ORDER BY
    [Avg IO] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

Starting with SQL Server 2008, you can get execution statistics for stored procedures through the `sys.dm_exec_procedure_stats` view. You can use the code from Listing 4-4 to do that. As with the `sys.dm_exec_query_stats` view, you can sort data by various execution metrics, depending on your optimization strategy.

Example 4-4. Using the `sys.dm_exec_procedure_stats` view

```

SELECT TOP 50
    DB_NAME(ps.database_id) AS [DB]
    ,OBJECT_NAME(ps.object_id, ps.database_id) AS [Proc Name]
    ,ps.type_desc AS [Type]
    ,ps.cached_time AS [Cached Time]
    ,ps.last_execution_time AS [Last Exec Time]
    ,qp.query_plan AS [Plan]

```



```

    ,ps.execution_count AS [Exec Count]
    ,CONVERT(DECIMAL(10,5),
        IIF(datediff(second,ps.cached_time, ps.last_execution_time)
= 0,
        NULL,
        1.0 * ps.execution_count /
            datediff(second,ps.cached_time,
ps.last_execution_time)
        )
    ) AS [Exec Per Second]
    ,(ps.total_logical_reads + ps.total_logical_writes) /
        ps.execution_count AS [Avg IO]
    ,(ps.total_worker_time / ps.execution_count / 1000)
        AS [Avg CPU(ms)]
    ,ps.total_logical_reads AS [Total Reads]
    ,ps.last_logical_reads AS [Last Reads]
    ,ps.total_logical_writes AS [Total Writes]
    ,ps.last_logical_writes AS [Last Writes]
    ,ps.total_worker_time / 1000 AS [Total Worker Time]
    ,ps.last_worker_time / 1000 AS [Last Worker Time]
    ,ps.total_elapsed_time / 1000 AS [Total Elapsed Time]
    ,ps.last_elapsed_time / 1000 AS [Last Elapsed Time]
    ,ps.total_physical_reads AS [Total Physical Reads]
    ,ps.last_physical_reads AS [Last Physical Reads]
    ,ps.total_physical_reads / ps.execution_count AS [Avg Physical
Reads]
    ,ps.total_spills AS [Total Spills]
    ,ps.last_spills AS [Last Spills]
    ,(ps.total_spills / ps.execution_count) AS [Avg Spills]
FROM
    sys.dm_exec_procedure_stats ps WITH (NOLOCK)
        CROSS APPLY sys.dm_exec_query_plan(ps.plan_handle) qp
ORDER BY
    [Avg IO] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

Figure 4-3 shows partial output of the code.

DB	Proc Name	Type	Cached Time	Last Exec Time
1	IndexOptimize	SQL_STORED_PROCEDURE	2021-01-03 01:05:01.067	2021-01-10 01:05:00.640
2	archive_data	SQL_STORED_PROCEDURE	2021-01-10 05:45:00.503	2021-01-12 05:45:00.367
3	archive_misc	SQL_STORED_PROCEDURE	2021-01-03 05:47:42.547	2021-01-12 05:49:36.920
4	AGGREGATE_LI...	SQL_STORED_PROCEDURE	2021-01-10 03:01:48.407	2021-01-12 18:19:46.477
5	archive_orde...	SQL_STORED_PROCEDURE	2021-01-10 05:45:03.197	2021-01-12 05:45:06.397
6	archive_inte...	SQL_STORED_PROCEDURE	2021-01-11 05:47:20.680	2021-01-12 05:48:32.947
7	GENERATE_DB...	SQL_STORED_PROCEDURE	2021-01-10 03:01:48.407	2021-01-12 18:19:46.477

Plan	Exec Count	Exec Per Second	Avg IO	Avg CPU(ms)
NULL	2	0.00000	997605286	3231184
<ShowPlanXML xmlns="http://...	3	0.00002	177832619	940548
<ShowPlanXML xmlns="http://...	10	0.00001	80059850	677919
<ShowPlanXML xmlns="http://...	404	0.00177	44012127	117734
NULL	3	0.00002	43522659	73929
<ShowPlanXML xmlns="http://...	2	0.00002	23452696	40479

Figure 4-3. Fig. 4-3 Partial output of sys.dm_exec_procedure_stats view

As you can see in the output, you can get execution plans for the stored procedures. Internally, the execution plans of stored procedures and other T-SQL modules are just collections of each statement's individual plan. In

some cases—for example, when a stored procedure involves dynamic SQL—the script will not return a plan in the output.

Listing 4-5 helps to address this. You can use it to get cached execution plans and their statistics (for stored procedure statements that have plans cached).

Example 4-5. Getting execution plan and statistics for stored procedure statements

```
SELECT
    qs.creation_time AS [Cached Time]
    ,qs.last_execution_time AS [Last Exec Time]
    ,SUBSTRING(qt.text, (qs.statement_start_offset/2)+1,
        ((
            CASE qs.statement_end_offset
                WHEN -1 THEN DATALENGTH(qt.text)
                ELSE qs.statement_end_offset
            END - qs.statement_start_offset)/2)+1) AS SQL
    ,qp.query_plan AS [Query Plan]
    ,CONVERT(DECIMAL(10,5),
        IIF(datediff(second,qs.creation_time,
qs.last_execution_time) = 0,
            NULL,
            1.0 * qs.execution_count /
                datediff(second,qs.creation_time,
qs.last_execution_time)
        )
    ) AS [Exec Per Second]
    ,(qs.total_logical_reads + qs.total_logical_writes) /
        qs.execution_count AS [Avg IO]
    ,(qs.total_worker_time / qs.execution_count / 1000)
        AS [Avg CPU(ms)]
    ,qs.total_logical_reads AS [Total Reads]
    ,qs.last_logical_reads AS [Last Reads]
    ,qs.total_logical_writes AS [Total Writes]
    ,qs.last_logical_writes AS [Last Writes]
    ,qs.total_worker_time / 1000 AS [Total Worker Time]
    ,qs.last_worker_time / 1000 AS [Last Worker Time]
    ,qs.total_elapsed_time / 1000 AS [Total Elapsed Time]
    ,qs.last_elapsed_time / 1000 AS [Last Elapsed Time]
    ,qs.total_rows AS [Total Rows]
    ,qs.last_rows AS [Last Rows]
    ,qs.total_rows / qs.execution_count AS [Avg Rows]
    ,qs.total_physical_reads AS [Total Physical Reads]
    ,qs.last_physical_reads AS [Last Physical Reads]
    ,qs.total_physical_reads / qs.execution_count
```

```

        AS [Avg Physical Reads]
,qs.total_grant_kb AS [Total Grant KB]
,qs.last_grant_kb AS [Last Grant KB]
,(qs.total_grant_kb / qs.execution_count)
    AS [Avg Grant KB]
,qs.total_used_grant_kb AS [Total Used Grant KB]
,qs.last_used_grant_kb AS [Last Used Grant KB]
,(qs.total_used_grant_kb / qs.execution_count)
    AS [Avg Used Grant KB]
,qs.total_ideal_grant_kb AS [Total Ideal Grant KB]
,qs.last_ideal_grant_kb AS [Last Ideal Grant KB]
,(qs.total_ideal_grant_kb / qs.execution_count)
    AS [Avg Ideal Grant KB]
,qs.total_columnstore_segment_reads
    AS [Total CSI Segments Read]
,qs.last_columnstore_segment_reads
    AS [Last CSI Segments Read]
,(qs.total_columnstore_segment_reads / qs.execution_count)
    AS [AVG CSI Segments Read]
,qs.max_dop AS [Max DOP]
,qs.total_spills AS [Total Spills]
,qs.last_spills AS [Last Spills]
,(qs.total_spills / qs.execution_count) AS [Avg Spills]
FROM
    sys.dm_exec_query_stats qs WITH (NOLOCK)
        CROSS APPLY sys.dm_exec_sql_text(qs.sql_handle) qt
        CROSS APPLY sys.dm_exec_text_query_plan

(qs.plan_handle,qs.statement_start_offset,qs.statement_end_offset)
qp
WHERE
    OBJECT_NAME(qt.objectid, qt.dbid) = <SP Name>
ORDER BY
    qs.statement_start_offset, qs.statement_end_offset
OPTION (RECOMPILE, MAXDOP 1);

```

Finally, starting with SQL Server 2016, you can get execution statistics for triggers and scalar user-defined functions, using `sys.dm_exec_trigger_stats` and `sys.dm_exec_function_stats`, respectively. Listing 4-6 shows the code to do that.

Example 4-6. Getting execution statistics for user-defined functions and triggers

```

SELECT TOP 50
    DB_NAME(fs.database_id) AS [DB]
    ,OBJECT_NAME(fs.object_id, fs.database_id) AS [Function]

```

```

,fs.type_desc AS [Type]
,fs.cached_time AS [Cached Time]
,fs.last_execution_time AS [Last Exec Time]
,qp.query_plan AS [Plan]
,fs.execution_count AS [Exec Count]
,CONVERT(DECIMAL(10,5),
        IIF(datediff(second,fs.cached_time, fs.last_execution_time)
= 0,
        NULL,
        1.0 * fs.execution_count /
        datediff(second,fs.cached_time,
fs.last_execution_time)
        )
) AS [Exec Per Second]
,(fs.total_logical_reads + fs.total_logical_writes) /
    fs.execution_count AS [Avg IO]
,(fs.total_worker_time / fs.execution_count / 1000) AS [Avg
CPU(ms)]
,fs.total_logical_reads AS [Total Reads]
,fs.last_logical_reads AS [Last Reads]
,fs.total_logical_writes AS [Total Writes]
,fs.last_logical_writes AS [Last Writes]
,fs.total_worker_time / 1000 AS [Total Worker Time]
,fs.last_worker_time / 1000 AS [Last Worker Time]
,fs.total_elapsed_time / 1000 AS [Total Elapsed Time]
,fs.last_elapsed_time / 1000 AS [Last Elapsed Time]
,fs.total_physical_reads AS [Total Physical Reads]
,fs.last_physical_reads AS [Last Physical Reads]
,fs.total_physical_reads / fs.execution_count AS [Avg Physical
Reads]
FROM
    sys.dm_exec_function_stats fs WITH (NOLOCK)
    CROSS APPLY sys.dm_exec_query_plan(fs.plan_handle) qp
ORDER BY
    [Avg IO] DESC
OPTION (RECOMPILE, MAXDOP 1);
SELECT TOP 50
    DB_NAME(ts.database_id) AS [DB]
    ,OBJECT_NAME(ts.object_id, ts.database_id) AS [Function]
    ,ts.type_desc AS [Type]
    ,ts.cached_time AS [Cached Time]
    ,ts.last_execution_time AS [Last Exec Time]
    ,qp.query_plan AS [Plan]
    ,ts.execution_count AS [Exec Count]
    ,CONVERT(DECIMAL(10,5),
        IIF(datediff(second,ts.cached_time, ts.last_execution_time)
= 0,
        NULL,

```

```

        1.0 * ts.execution_count /
            datediff(second,ts.cached_time,
ts.last_execution_time)
    )
    ) AS [Exec Per Second]
    , (ts.total_logical_reads + ts.total_logical_writes) /
        ts.execution_count AS [Avg IO]
    , (ts.total_worker_time / ts.execution_count / 1000) AS [Avg
CPU (ms)]
    , ts.total_logical_reads AS [Total Reads]
    , ts.last_logical_reads AS [Last Reads]
    , ts.total_logical_writes AS [Total Writes]
    , ts.last_logical_writes AS [Last Writes]
    , ts.total_worker_time / 1000 AS [Total Worker Time]
    , ts.last_worker_time / 1000 AS [Last Worker Time]
    , ts.total_elapsed_time / 1000 AS [Total Elapsed Time]
    , ts.last_elapsed_time / 1000 AS [Last Elapsed Time]
    , ts.total_physical_reads AS [Total Physical Reads]
    , ts.last_physical_reads AS [Last Physical Reads]
    , ts.total_physical_reads / ts.execution_count AS [Avg Physical
Reads]
FROM
    sys.dm_exec_trigger_stats ts WITH (NOLOCK)
        CROSS APPLY sys.dm_exec_query_plan(ts.plan_handle) qp
ORDER BY
    [Avg IO] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

Troubleshooting based on plan cache-based execution statistics has several limitations, and you may miss some queries. Nevertheless, it is a great starting point. Most importantly, the data is collected automatically and you can access it immediately, without setting up additional monitoring tools.

SQL Traces and Extended Events

I am sure that every SQL Server engineer is aware of SQL Traces and Extended Events. They allow you to capture various events in a system for analysis and troubleshooting in real time. You can also use them to capture long-running and expensive queries, including those that don't cache execution plans and are therefore missed by the sys.dm_exec_query_stats view.

I'd like to start this section with a warning, though: Do *not* use SQL Traces and xEvents for this purpose unless it is *absolutely necessary*. Capturing executed statements is an expensive operation that may introduce significant performance overhead in busy systems. (You saw one such example in Chapter 1.)

It does not matter how much data you collect. You can exclude most statements from the output by filtering out queries with low resource consumption. But SQL Server will still have to capture *all* statements to evaluate, filter, and discard unnecessary events.

Don't collect unnecessary information in events you are collecting or in xEvent actions you are capturing. Some actions—for example, callstack—are expensive and lead to a serious performance hit when enabled.

I do not want to beat a dead horse, but I have no choice: use Extended Events instead of SQL Traces. They are lighter and introduce less overhead in the system. Choose an in-memory ring_buffer target and allow event loss in configuration when possible.

Table 4-1 shows several Extended and SQL Trace Events that can be used to detect inefficient queries.

*T
a
b
l
e
4
-
l
.
E
x
t
e
n
d
e
d
a
n
d
S
Q
L*

*T
r
a
c
e
E
v
e
n*

*t
s
t
o
d
e
t
e
c
t
i
n
e
f
f
i
c
i
e
n
t
q
u
e
r
i
e
s*

SQL Trace Event

xEvent

Comments

Fired when statement starts the

SQL:StmtStarting	sqlserver.sql_statement_starting	execution.
------------------	----------------------------------	------------

SQL:StmtCompleted	sqlserver. sql_statement_completed	Fired when statement finishes the execution.
-------------------	---------------------------------------	--

SP:StmtStarting	sqlserver.sp_statement_starting	Fired when SQL statement within T-SQL module (stored procedure, user-defined function, etc.) starts the execution.
-----------------	---------------------------------	--

SP:StmtCompleted	sqlserver.sp_statement_completed	Fired when SQL statement within T-SQL module completes the execution.
------------------	----------------------------------	---

RPC:Starting	sqlserver.rpc_starting	Fired when remote procedure call (RPC) is starting. Those calls are parameterized SQL requests, such as calls of stored procedures or parameterized batches, sent from applications. Many client libraries will run queries via sp_executesql calls, which can be captured by that event.
--------------	------------------------	---

RPC:Completed	sqlserver.rpc_completed	Fired when RPC completes.
---------------	-------------------------	---------------------------

SP:Starting	sqlserver.module_start	Fired when T-SQL module starts execution.
-------------	------------------------	---

SP:Completed	sqlserver.module_completed	Fired when T-SQL module completes the execution.
--------------	----------------------------	--

Error:Attention	sqlserver.attention	Occurs when client terminates query execution due to timeout or connection loss.
-----------------	---------------------	--

Usually, when I need to capture inefficient queries, I set up an xEvents session capturing `sqlserver.rpc_completed`, `sqlserver.sql_completed`, and `sqlserver.attention` events and filtering data by execution metrics, such as `cpu_time` or `logical_reads`. As part of the event, I capture several actions, most notably `sqlserver.sql_text` and client information.

Listing 4-7 shows code to capture queries that consume more than 3,000ms of CPU time or produce more than 10,000 logical reads or writes. This code will work in SQL Server 2012 and above; it may require small modifications in SQL Server 2008 due to the different way it works with the file target.

Example 4-7. Capturing CPU- and I/O intensive queries

```
CREATE EVENT SESSION [Expensive Queries]
ON SERVER
ADD EVENT
    sqlserver.sql_statement_completed
    (
        ACTION
        (
            sqlserver.client_app_name
            ,sqlserver.client_hostname
            ,sqlserver.database_id
            ,sqlserver.plan_handle
            ,sqlserver.sql_text
```

```

        ,sqlserver.username
    )
WHERE
    (
        (
            cpu_time >= 3000000 or -- Time in microseconds
            logical_reads >= 10000 or
            writes >= 10000
        ) AND
        sqlserver.is_system = 0
    )
),
ADD EVENT
    sqlserver.rpc_completed
    (
        ACTION
            (
                sqlserver.client_app_name
                ,sqlserver.client_hostname
                ,sqlserver.database_id
                ,sqlserver.plan_handle
                ,sqlserver.sql_text
                ,sqlserver.username
            )
        WHERE
            (
                (
                    cpu_time >= 3000000 or
                    logical_reads >= 10000 or
                    writes >= 10000
                ) AND
                sqlserver.is_system = 0
            )
    )
)
ADD TARGET
    package0.event_file
    (
        SET FILENAME = 'c:\ExtEvents\Expensive Queries.xel'
    )
WITH
    (
        event_retention_mode=allow_single_event_loss
        ,max_dispatch_latency=30 seconds
    );

```

You can parse the captured results with the code from Listing 4-8.

Example 4-8. Parsing collected xEvent data

```

;WITH TargetData(Data, File_Name, File_Offset)
AS
(
    SELECT CONVERT(xml,event_data) AS Data, file_name, file_offset
    FROM
        sys.fn_xe_file_target_read_file
            ('c:\extevents\Expensive Queries*.xel',NULL,NULL,NULL)
)
,EventInfo([Event],[Event Time],[DB],[Statement],[SQL],[User Name]
,[Client],[App],[CPU Time],[Duration],[Logical Reads]
,[Physical Reads],[Writes],[Rows],[PlanHandle]
,[File_Name],[File_Offset])
as
(
    SELECT
        Data.value('/event[1]/@name','sysname') AS [Event]
        ,Data.value('/event[1]/@timestamp','datetime') AS [Event Time]
        ,Data.value('((/event[1]/data[@name="database_id"]/value/text())
[1]))','INT')
        AS [DB]
        ,Data.value('((/event[1]/data[@name="statement"]/value/text())
[1]))'
        , 'nvarchar(max)') AS [Statement]
        ,Data.value('((/event[1]/data[@name="sql_text"]/value/text())
[1]))'
        , 'nvarchar(max)') AS [SQL]
        ,Data.value('((/event[1]/data[@name="username"]/value/text())
[1]))'
        , 'nvarchar(255)') AS [User Name]

        ,Data.value('((/event[1]/data[@name="client_hostname"]/value/text())
[1]))'
        , 'nvarchar(255)') AS [Client]

        ,Data.value('((/event[1]/data[@name="client_app_name"]/value/text())
[1]))'
        , 'nvarchar(255)') AS [App]
        ,Data.value('((/event[1]/data[@name="cpu_time"]/value/text())
[1]))'
        , 'bigint') AS [CPU Time]
        ,Data.value('((/event[1]/data[@name="duration"]/value/text())
[1]))'
        , 'bigint') AS [Duration]

        ,Data.value('((/event[1]/data[@name="logical_reads"]/value/text())
[1]))'
        , 'int') AS [Logical Reads]

```

```

,Data.value('(/event[1]/data[@name="physical_reads"]/value/text())[1])'
    , 'int') AS [Physical Reads]
    ,Data.value('(/event[1]/data[@name="writes"]/value/text())[1])'
    , 'int') AS [Writes]
    ,Data.value('(/event[1]/data[@name="row_count"]/value/text())[1])'
    , 'int') AS [Rows]
    ,Data.value(
'xs:hexBinary((/event[1]/action[@name="plan_handle"]/value/text())[1]))'
    , 'varbinary(64)') AS [PlanHandle]
    ,File_Name
    ,File_Offset
FROM
    TargetData
)
SELECT
    ei.*, qp.Query_Plan
FROM
    EventInfo ei
    OUTER APPLY sys.dm_exec_query_plan(ei.PlanHandle) qp
OPTION (MAXDOP 1, RECOMPILE);

```

When you work with SQL Traces and xEvents, you have to deal with raw data. You'll need to aggregate it to determine which queries introduce the most cumulative impact.

Again: beware of the overhead that xEvents and SQL Traces introduce in systems. Do not create and run those sessions permanently. In many cases you can get enough troubleshooting data by enabling the session or trace for just a few minutes.

For more extensive examples on how to work with different xEvents targets, see the code repository that accompanies this book. You can also read more about Extended Events in my book *Pro SQL Server Internals*.

Query Store

So far in this chapter, we have discussed two approaches to detecting inefficient queries. Both have limitations. Plan-cache-based data may miss

some queries; SQL Traces and xEvents require you to perform complex analysis of the output and may have significant performance overhead in busy systems.

The Query Store, introduced in SQL Server 2016, helps to address those limitations. You can think of it as something like the flight data recorders (or “black boxes”) in airplane cockpits, but for SQL Server. When the Query Store is enabled, SQL Server captures and persists runtime statistics and execution plans of the queries in the database. It shows how the execution plans perform and how they evolve over time. Finally, it allows you to force specific execution plans to queries addressing parameter-sniffing issues, which we will discuss in Chapter 6.

NOTE

The Query Store is disabled by default in the on-premises version of SQL Server. It is enabled by default in Azure SQL Databases and Azure SQL Managed Instances.

The Query Store is fully integrated into the query processing pipeline, as illustrated by the high-level diagram in Figure 4-4.

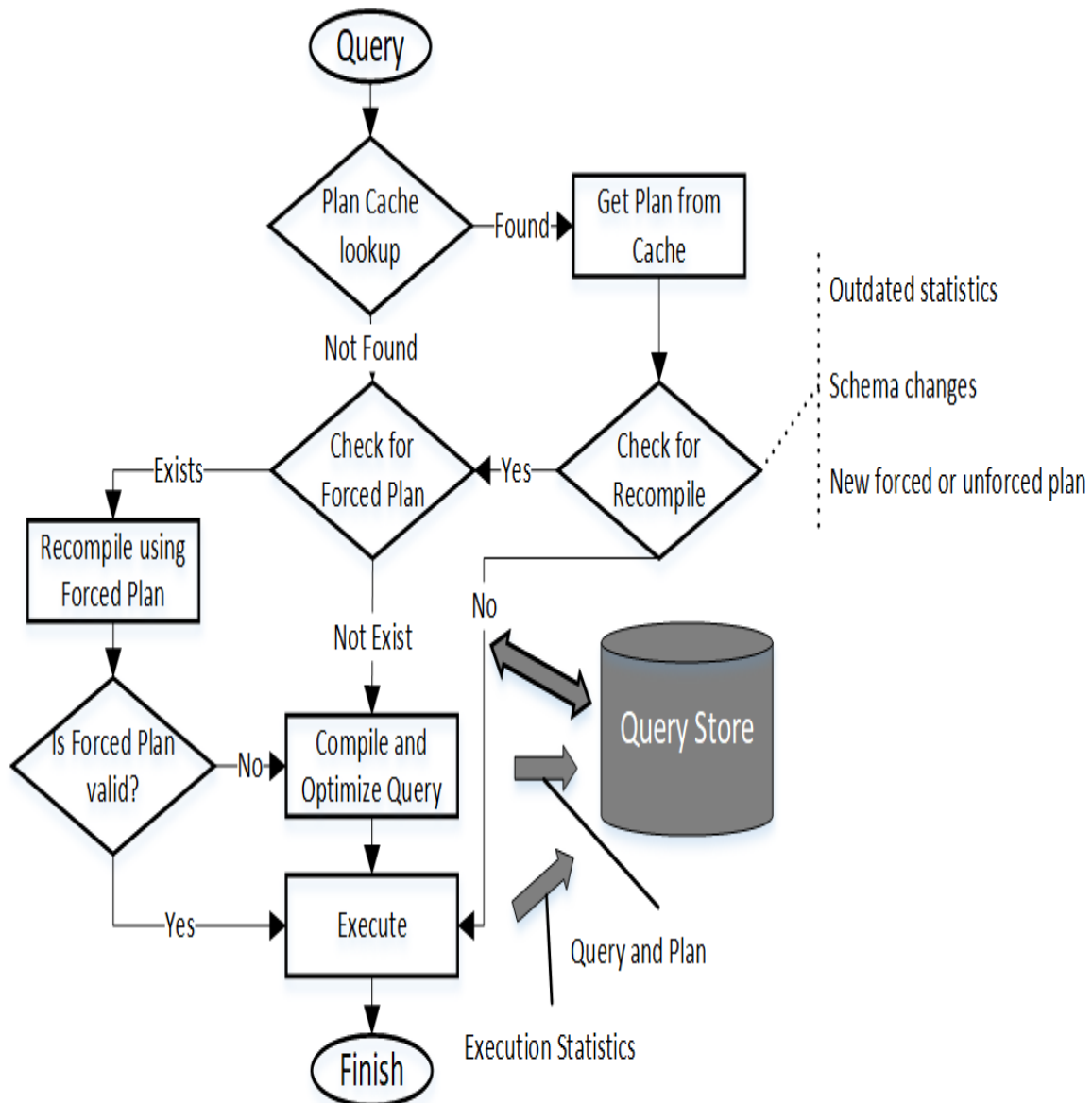


Figure 4-4. Fig. 4-4. Query processing pipeline

When a query needs to be executed, SQL Server looks up the execution plan from the plan cache. If it finds a plan, SQL Server checks if the query needs to be recompiled (due to statistics updates or other factors), if a new forced plan has been created, and if an old forced plan has been dropped from the Query Store.

During the compilation, SQL Server checks if the query has a forced plan available. When that happens, the query essentially gets compiled with the

forced plan, much like when the USE PLAN hint is used. If the resulting plan is valid, it is stored in the plan cache for reuse.

If the forced plan is no longer valid (for example, when a user drops an index referenced in the forced plan), SQL Server does not fail the query. Instead, it compiles the query again without the forced plan and without caching it afterwards. The Query Store, on the other hand, persists both plans, marking the forced plan as invalid. All of that happens without affecting the applications.

Despite its tight integration with the query processing pipeline and various internal optimizations, Query Store still adds overhead to the system. Just how much overhead depends on two main factors: the number of compilations and the data collection settings.

The more compilations SQL Server performs, the more load the Query Store must handle. In particular, the Query Store may not work very well in systems that have a very heavy, ad-hoc, non-parameterized workload.

Query Store's configurations allow you to specify if you want to capture all queries or just expensive ones, along with aggregation intervals and data retention settings. If you collect more data and/or use smaller aggregation intervals, you'll have more overhead.

The overhead introduced by the Query Store is usually relatively small. However, it may be significant in some cases. For example, I've been using the Query Store to troubleshoot the performance of one process that consists of a very large number of small ad-hoc queries. I captured all queries in the system using QUERY_CAPTURE_MODE=ALL mode, collecting almost 10GB of data in the Query Store. The process took 8 hours to complete with the Query Store enabled, comparing to 2.5 hours without it.

Nevertheless, I suggest enabling Query Store if your system can handle the overhead. Some SQL Server features, such as Intelligent Query Processing, rely on Query Store data and will benefit from it.

NOTE

Monitor QDS* waits when you enable Query Store. Excessive QDS* waits may be a sign of higher Query Store overhead in the system. Ignore QDS_PERSIST_TASK_MAIN_LOOP_SLEEP and QDS_ASYNC_QUEUE waits – they are benign.

You can work with the Query Store in two ways – through the graphics UI in SSMS or by querying data management views directly. Let's look at the UI first.

Query Store SSMS Reports

After you enable the Query Store in the database, you'll see a *Query Store* folder in the *Object Explorer* (Figure 4-5). The number of reports in the folder will depend on the versions of SQL Server and SSMS in your system. The rest of this section will walk you through the seven reports shown in Figure 4-5.

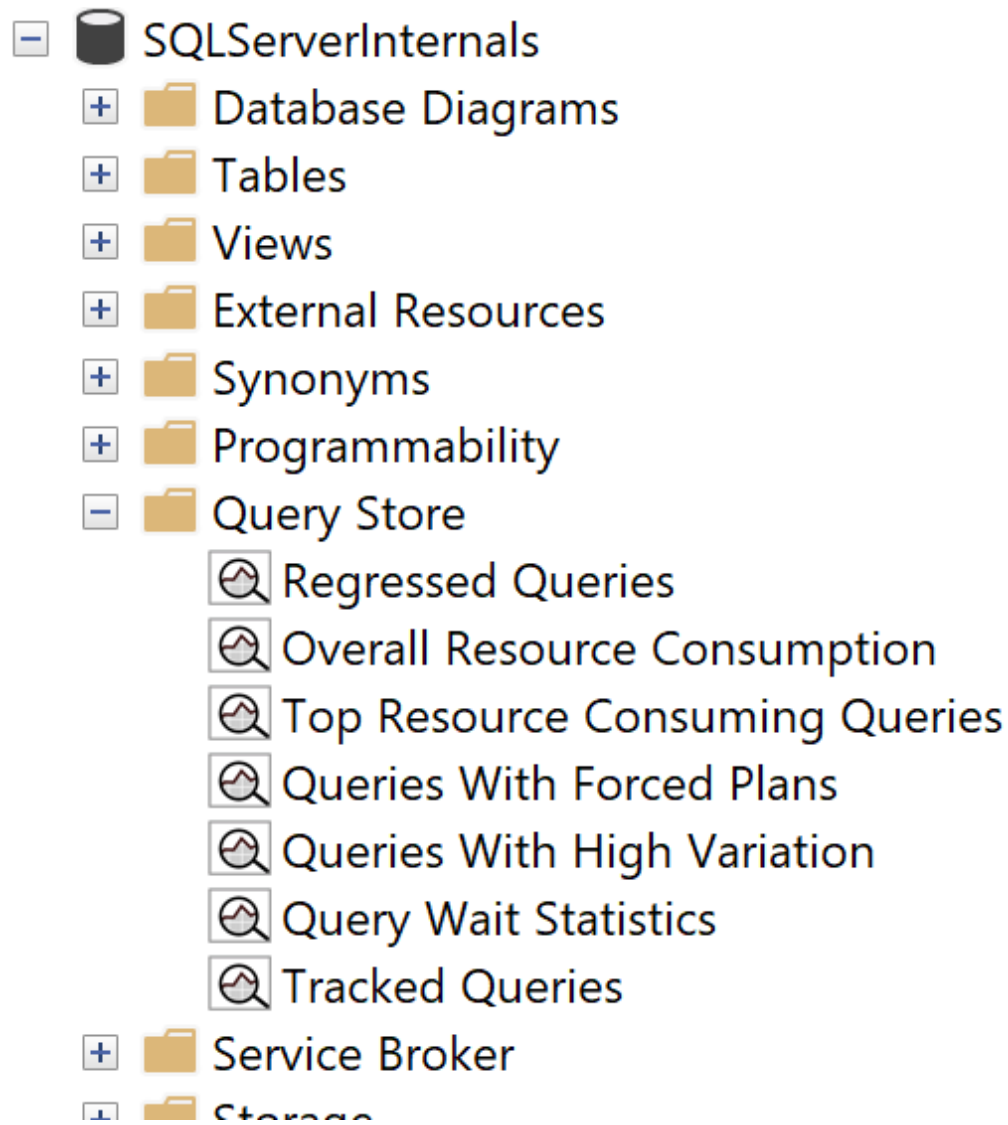


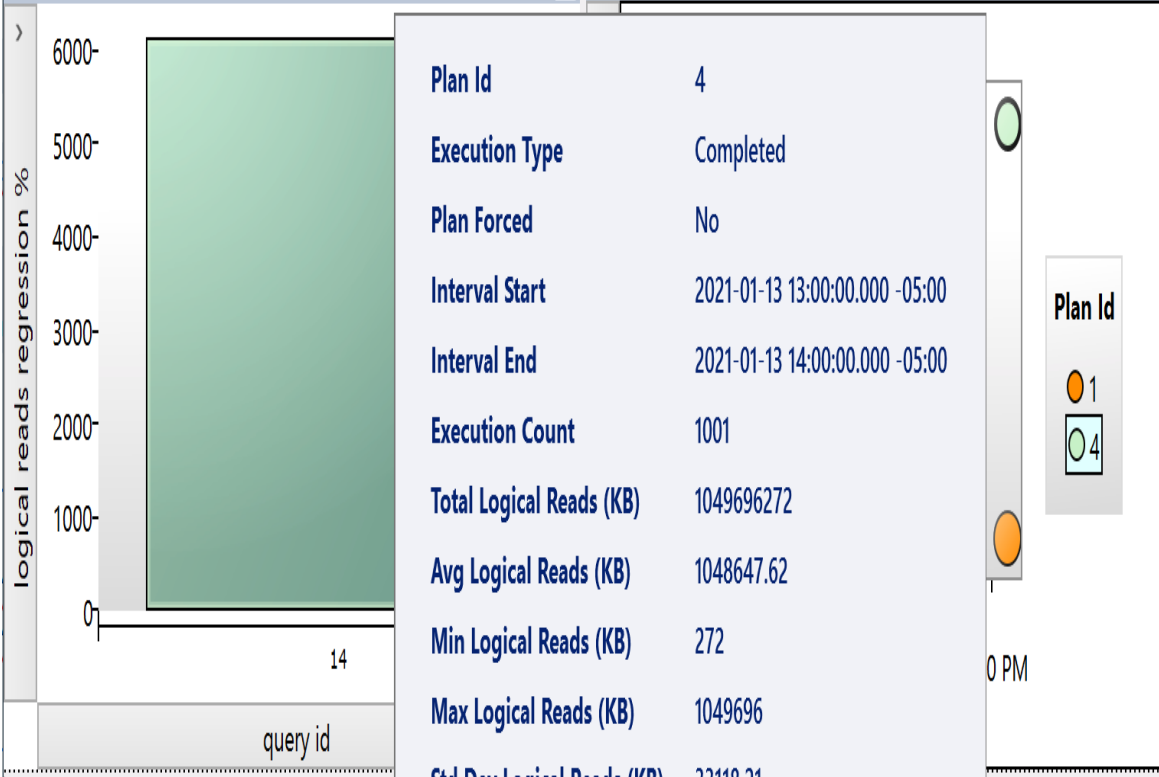
Figure 4-5. Fig. 4-5. Query Store reports in SSMS

Regressed Queries

This report, shown in Figure 4-6, shows queries whose performance has regressed overtime. You can configure the time frame and regression criteria (such as disk operations, CPU consumption, and number of executions) for analysis.

Top 25 regressed queries for database SQLServerInternals. Time period: 1/12/2021 1:40 PM - 1/13/2021 1:55 PM

Metric Logical Reads (KB) Statistic Avg Plan summary for query 14



Plan Id	4
Execution Type	Completed
Plan Forced	No
Interval Start	2021-01-13 13:00:00.000 -05:00
Interval End	2021-01-13 14:00:00.000 -05:00
Execution Count	1001
Total Logical Reads (KB)	1049696272
Avg Logical Reads (KB)	1048647.62
Min Logical Reads (KB)	272
Max Logical Reads (KB)	1049696
Std Dev Logical Reads (KB)	33118.21
Variation Logical Reads (KB)	0.03

Plan 4 [not forced]

Query 1: Query cost (relative to the query plan) 100%
 select @Salary = Avg(Salary) from dbo.Employees where Country = @Country

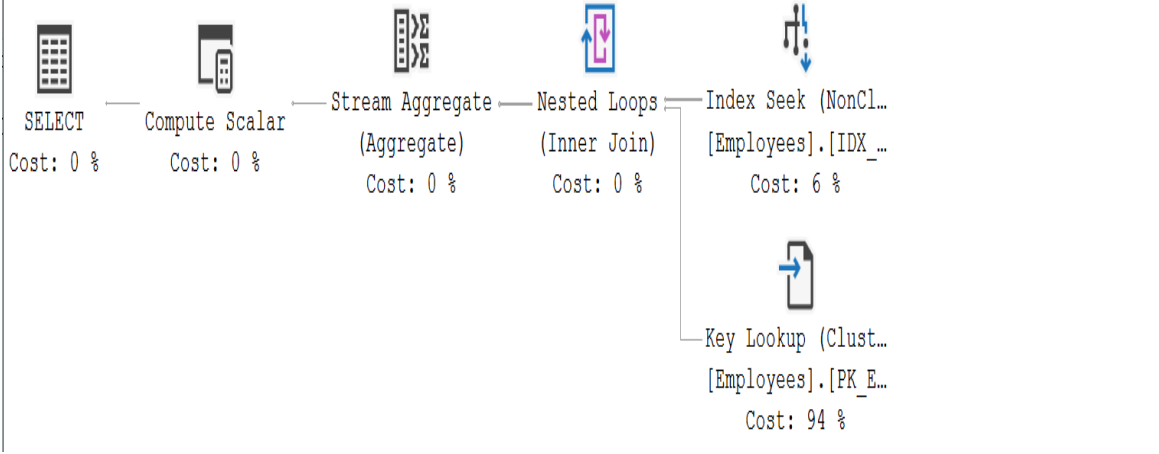


Figure 4-6. Fig 4-6. Regressed Queries report

Choose the query in the graph on the top left. The top right portion of the report illustrates collected execution plans for the selected query. You can click on the dots, which represent different execution plans, and see the plans at the bottom. You can also compare different execution plans.

The *Force Plan* button allows you to force a selected plan for the query. It calls the `sys.sp_query_store_force_plan` stored procedure internally. Similarly, the *Unforce Plan* button removes a forced plan by calling the `sys.sp_query_store_unforce_plan` stored procedure.

The Regressed Queries report is a great tool for troubleshooting issues related to parameter sniffing, which we will discuss in Chapter 6, and fixing them quickly by forcing specific execution plans.

Top Resource Consuming Queries

This report (Figure 4-7) allows you to detect the most resource-intensive queries in the system. While it works similarly to the data provided by `sys.dm_exec_query_stats` view, it does not depend on the plan cache. You can customize the metrics used for data sorting and the time interval.

Top Resource Consumers [SQLServerInternals]

Top 25 resource consumers for database SQLServerInternals. Time period: Last hour ending at 1/13/2021 1:50 PM

Metric: CPU Time (ms) Statist: Plan summary for query 14

plan id	plan forced	execution type	execution count	min cpu time	max cpu time	avg cpu time	std dev cpu time	va	
1	4	0	0	1001	0.1	101.31	61.47	3.45	0.
2	1	0	0	2000	13.47	21.05	15.82	0.54	0.

total cpu time

Query Id: 14
 Object Id: 1605580758
 Object Name: GetAverageSalary
 Total CPU Time (ms): 93168.22
 Execution Count: 3001
 Plan Count: 2

Plan 4 [not forced] (@Country varchar(64),@Salary money)select @Salary = Avg(Salary)
 from dbo.Employees
 where Country = @Country

Query 1: Quer
 select @Salari
 ry = @Country

Force Plan Unforce Plan

SELECT
 Cost: 0 %

Compute Scalar
 Cost: 0 %

Stream Aggregate (Aggregate)
 Cost: 0 %

Nested Loops (Inner Join)
 Cost: 0 %

Index Seek (NonCl...
 [Employees].[IDX_...]
 Cost: 6 %

Key Lookup (Clust...
 [Employees].[PK_E...]
 Cost: 94 %

Figure 4-7. Fig. 4-7. Top Resource Consuming Queries report

Overall Resource Consumption

This report shows you the workload's statistics and resource usage over the time intervals you specify. It will allow you to detect and analyze spikes in resource usage and drill down to the queries that introduce such spikes.

Figure 4-8 shows the output of the report.

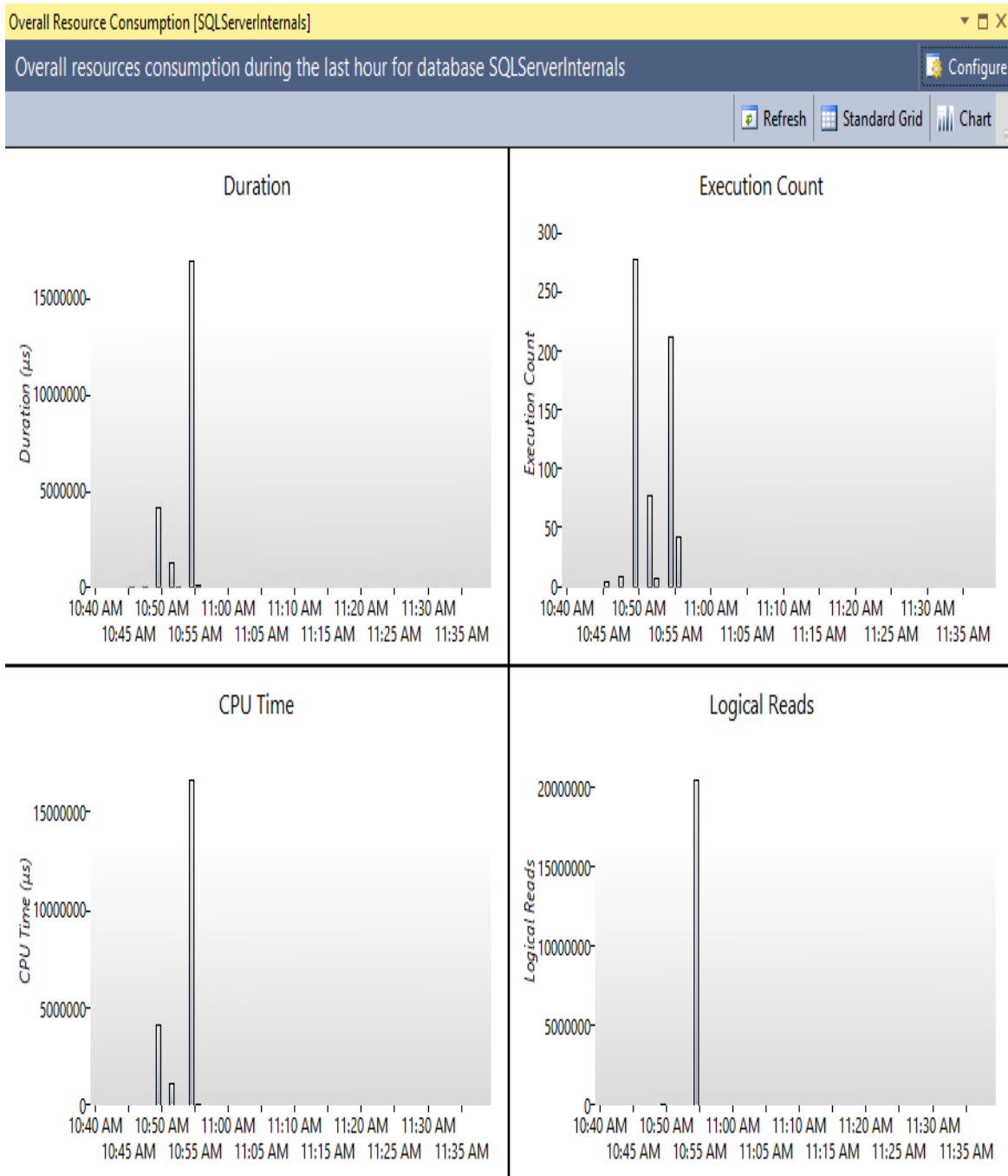


Figure 4-8. Fig. 4-8. Overall Resource Consumption report

Queries With High Variation

This report allows you pinpoint queries with high performance variation. You can use it to detect anomalies in the workload, along with possible performance regressions. (For the sake of space, I'll skip the screenshots here.)

Queries With Forced Plan

This report shows you the queries that have an execution plan forced in the system.

Query Wait Statistics

This report allows you to detect queries with high waits. The data is grouped by several categories (such as CPU, disk, and blocking), depending on wait type. You can see details on wait mapping in the [Microsoft Documentation](#).

Tracked Queries

Finally, the Tracked Queries report allows you to monitor execution plans and statistics for individual queries. It provides similar information to the *Regressed Queries* and *Top Resource Consuming Queries* reports, at the scope of individual queries.

These reports will give you a large amount of data for analysis. However, in some cases, you'll want to use T-SQL and work with the Query Store data directly. Let's look at how you can accomplish that.

Working with Query Store DMVs

The Query Store data management views (DMVs) are highly normalized, as shown in Figure 4-9. Execution statistics are tracked for each execution plan and grouped by collection intervals, which are defined by the `INTERVAL_LENGTH_MINUTES` setting.

As I've noted, the smaller the intervals you use, the more data will be collected and persisted in the Query Store. The same applies to the system workload: an excessive number of ad-hoc queries may balloon the Query Store's size. Keep this in mind when you configure the Query Store in your system.

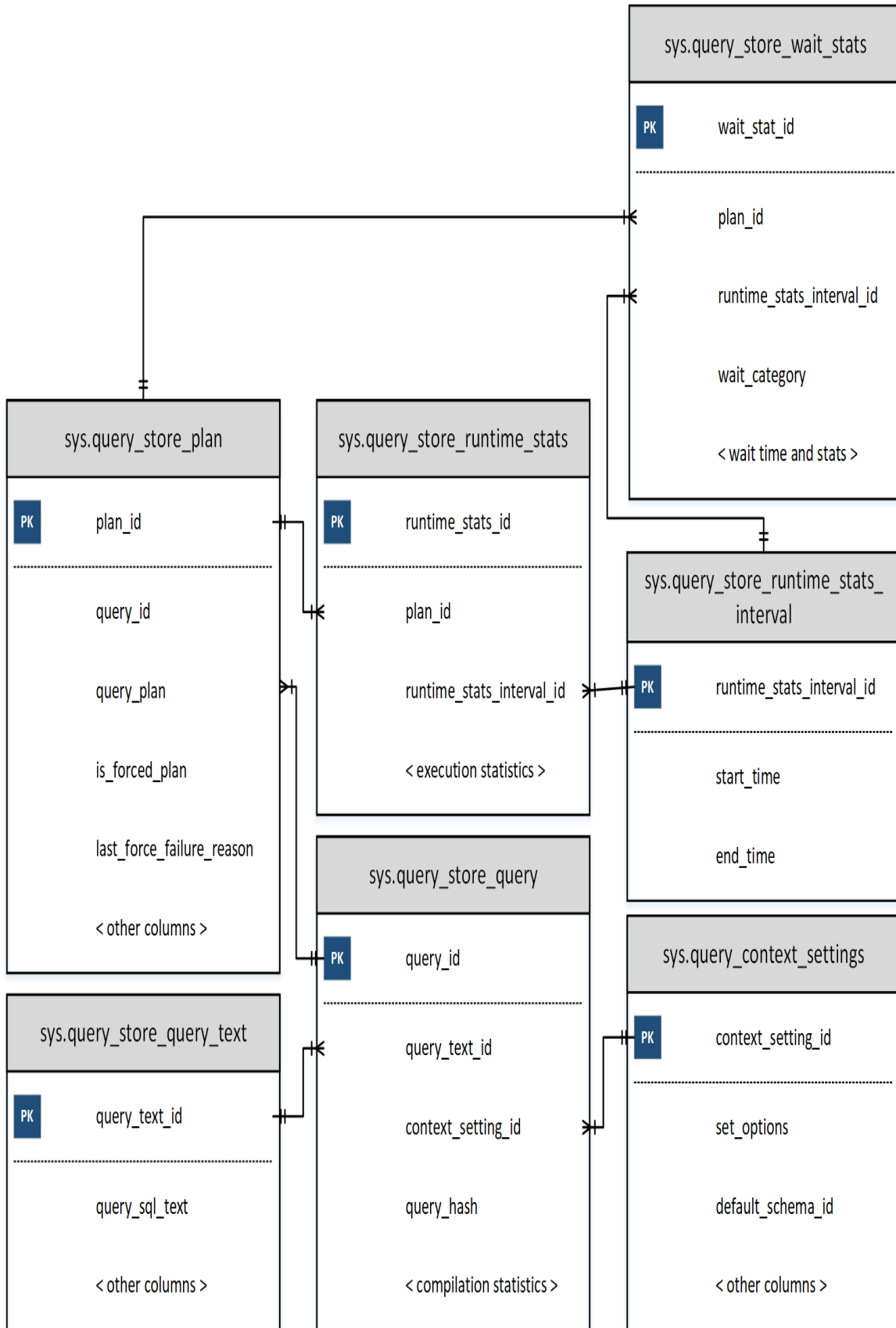


Figure 4-9. Fig. 4-9. Query Store DMVs

You can logically separate DMVs into two categories: plan store and runtime statistics. The former ones include the following views:

sys.query_store_query

The `sys.query_store_query` **view** provides information about queries and their compilation statistics, and last execution time.

sys.query_store_query_text

The `sys.query_store_query_text` **view** shows information about query text.

sys.query_context_setting

The `sys.query_context_setting` **view** contains information about context settings associated with the query. It includes SET options, default schema for the session, language, and other attributes. SQL Server may generate and cache separate execution plans for the same query when those settings are different.

sys.query_store_plan

The `sys.query_store_plan` **view** provides information about query execution plans. The `is_forced_plan` column indicates whether the plan is forced. The `last_force_failure_reason` tells you why a forced plan was not applied to the query.

As you can see, each query can have multiple entries in the `sys.query_store_query` and `sys.query_store_plan` views. This will vary based on your session context options, recompilations, and other factors.

Three other views represent runtime statistics:

sys.query_store_runtime_stats_interval

The `sys.query_store_runtime_stats_interval` **view** contains information about statistics collection intervals.

sys.query_store_runtime_stats

The `sys.query_store_runtime_stats` **view** references the `sys.query_store_plan` view and contains information about runtime statistics for a specific plan during a particular `sys.query_store_runtime_stats_interval` interval. It provides information about execution count, CPU time and call durations, logical and physical I/O statistics, transaction log usage, degree of parallelism, memory grant size, and a few other useful metrics.

sys.query_store_wait_stats

Starting with SQL Server 2017, you can get information about query waits with the `sys.query_store_wait_stats` **view**. The data is collected for each plan and time interval and grouped by several wait categories, including CPU, memory, and blocking.

Let's look at a few scenarios for working with Query Store data.

Listing 4-9 provides code that returns information about the system's 50 most I/O-intensive queries. Because the Query Store persists execution statistics over time intervals, you'll need to aggregate data from multiple `sys.query_store_runtime_stats` rows. The output will include data for all intervals that ended within the last 24 hours, grouped by queries and their execution plans.

Example 4-9. Getting information about expensive queries from Query Store

```
SELECT TOP 50
    q.query_id, qt.query_sql_text, qp.plan_id, qp.query_plan
    ,SUM(rs.count_executions) AS [Execution Cnt]
    ,CONVERT(INT,SUM(rs.count_executions *
        (rs.avg_logical_io_reads + avg_logical_io_writes)) /
        SUM(rs.count_executions)) AS [Avg IO]
    ,CONVERT(INT,SUM(rs.count_executions *
        (rs.avg_logical_io_reads + avg_logical_io_writes))) AS [Total
IO]
    ,CONVERT(INT,SUM(rs.count_executions * rs.avg_cpu_time) /
        SUM(rs.count_executions)) AS [Avg CPU]
    ,CONVERT(INT,SUM(rs.count_executions * rs.avg_cpu_time)) AS [Total
CPU]
```

```

, CONVERT(INT, SUM(rs.count_executions * rs.avg_duration) /
SUM(rs.count_executions)) AS [Avg Duration]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_duration))
AS [Total Duration]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_physical_io_reads) /
SUM(rs.count_executions)) AS [Avg Physical Reads]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_physical_io_reads))
AS [Total Physical Reads]
, CONVERT(INT, SUM(rs.count_executions *
rs.avg_query_max_used_memory) /
SUM(rs.count_executions)) AS [Avg Memory Grant Pages]
, CONVERT(INT, SUM(rs.count_executions *
rs.avg_query_max_used_memory))
AS [Total Memory Grant Pages]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_rowcount) /
SUM(rs.count_executions)) AS [Avg Rows]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_rowcount)) AS [Total
Rows]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_dop) /
SUM(rs.count_executions)) AS [Avg DOP]
, CONVERT(INT, SUM(rs.count_executions * rs.avg_dop)) AS [Total DOP]
FROM
sys.query_store_query q WITH (NOLOCK)
JOIN sys.query_store_plan qp WITH (NOLOCK) ON
q.query_id = qp.query_id
JOIN sys.query_store_query_text qt WITH (NOLOCK) ON
q.query_text_id = qt.query_text_id
JOIN sys.query_store_runtime_stats rs WITH (NOLOCK) ON
qp.plan_id = rs.plan_id
JOIN sys.query_store_runtime_stats_interval rsi WITH (NOLOCK) ON
rs.runtime_stats_interval_id = rsi.runtime_stats_interval_id
WHERE
rsi.end_time >= DATEADD(DAY, -1, GETDATE())
GROUP BY
q.query_id, qt.query_sql_text, qp.plan_id, qp.query_plan
ORDER BY
[Avg IO] DESC
OPTION (MAXDOP 1, RECOMPILE);

```

Obviously, you can sort data by different criteria than average I/O. You can also add predicates to the WHERE and/or HAVING clauses of the query to narrow down the results. For example, you can filter by *DOP* columns if you want to detect queries that use parallelism in an OLTP environment and fine-tune the *Cost Threshold for Parallelism* setting.

Another example is for detecting queries that balloon the plan cache. The code in Listing 4-10 provides information about queries that generate multiple execution plans due to different context settings. The two most common reasons for this are sessions that use different SET options and queries that reference objects without schema names.

Example 4-10. Queries with different context settings

```
SELECT
    q.query_id, qt.query_sql_text
    ,COUNT(DISTINCT q.context_settings_id) AS [Context Setting Cnt]
    ,COUNT(DISTINCT qp.plan_id) AS [Plan Count]
FROM
    sys.query_store_query q WITH (NOLOCK)
    JOIN sys.query_store_query_text qt WITH (NOLOCK) ON
        q.query_text_id = qt.query_text_id
    JOIN sys.query_store_plan qp WITH (NOLOCK) ON
        q.query_id = qp.query_id
GROUP BY
    q.query_id, qt.query_sql_text
HAVING
    COUNT(DISTINCT q.context_settings_id) > 1
ORDER BY
    COUNT(DISTINCT q.context_settings_id)
OPTION (MAXDOP 1, RECOMPILE);
```

Listing 4-11 shows you how to find similar queries based on query_hash value. Usually, those queries belong to a non-parameterized ad-hoc workload in the system. You can parameterize those queries in the code. If that's not possible, consider using forced parameterization, which we will discuss in Chapter 6.

Example 4-11. Detecting queries with duplicated query_hash value

```
SELECT TOP 100
    q.query_hash
    ,COUNT(*) AS [Query Count]
    ,AVG(rs.count_executions) AS [Avg Exec Count]
FROM
    sys.query_store_query q WITH (NOLOCK)
    JOIN sys.query_store_plan qp WITH (NOLOCK) ON
        q.query_id = qp.query_id
    JOIN sys.query_store_runtime_stats rs WITH (NOLOCK) ON
        qp.plan_id = rs.plan_id
GROUP BY
    q.query_hash
```

```
HAVING
    COUNT(*) > 1
ORDER BY
    [Avg Exec Count] ASC, [Query Count] DESC
OPTION (MAXDOP 1, RECOMPILE);
```

You can view additional examples in the book's code repository.

As you can see, the possibilities are endless. Use the Query Store if you can afford its overhead in your system.

Third-Party Tools

As you've now seen, SQL Server provides a very rich and extensive set of tools to locate inefficient queries. Nevertheless, you may also benefit from monitoring tools developed by other vendors. Most will provide you with a list of most resource-intensive queries for analysis and optimization. Many will also give you the baseline, which you can use to analyze trends and detect regressed queries.

I am not going to discuss specific tools; instead, I want to offer you a few tips for choosing and using these tools.

The key to using any tool is understanding it. Research how it works and analyze its limitations and what data it may miss. For example, if a tool gets data by polling the `sys.dm_exec_requests` view on schedule, it may miss a big portion of small but frequently executed queries that run in between polls. Alternatively, if a tool determines inefficient queries by session waits, the results will greatly depend on your system's workload, the amount of data cached in the buffer pool, and many other factors.

Depending on your specific needs, these limitations might be acceptable. Remember the **Pareto principle** (also known as the "80/20 rule"): you don't need to optimize all inefficient queries in the system. Nevertheless, you may benefit from a holistic view and from multiple perspectives. For example, it is very easy to cross-check a tool's list of inefficient queries against the plan-cache-based execution statistics for a more complete list.

There is another important reason to understand your tool, though: estimating the amount of overhead it could introduce. Some DMVs are very expensive to run. For example, if a tool calls the `sys.dm_exec_query_plan` function during each `sys.dm_exec_requests` poll, it may lead to a measurable increase in overhead in busy systems. It is also not uncommon for tools to create traces and xEvent sessions without your knowledge.

In the end, choose the approach that best allows you to pinpoint inefficient queries and that works best with your system. Remember that query optimization will help in any system.

Summary

Inefficient queries impact SQL Server's performance and can overload the disk subsystem. Even in systems that have enough memory to cache data in the buffer pool, those queries burn CPU, increase blocking, and affect the customer experience.

SQL Server keeps track of execution metrics for each cached plan and exposes them through the `sys.dm_exec_query_stats` view. You can also get execution statistics for stored procedures, triggers, and scalar user-defined functions with `sys.dm_exec_procedure_stats`, `sys.dm_exec_trigger_stats`, and `sys.dm_exec_function_stats` views, respectively.

Your plan-cache-based execution statistics will not track runtime execution metrics in execution plans, nor will it include queries that do not have plans cached. Make sure to factor this to your analysis and query-tuning process.

You can capture inefficient queries in real time with Extended Events and SQL Traces. Both approaches introduce overhead, especially in busy systems. They also provide raw data, which you'll need to process and aggregate for further analysis.

In SQL Server 2016 and above, you can utilize the Query Store. This is a great tool that does not depend on the plan cache and allows you to quickly pinpoint plan regressions. The Query Store adds some overhead; this may be acceptable in many cases, but monitor it when you enable the feature.

Finally, I discussed how you can use third-party monitoring tools to find inefficient queries. Remember to research how a tool works and understand its limitation and overhead.

In the next chapter, we will discuss a few common techniques that you can use to optimize inefficient queries.

Troubleshooting Checklist

- Get the list of inefficient queries from the `sys.dm_exec_query_stats` view. Sort the data according to your troubleshooting strategy (CPU, I/O, and so forth).
- Detect the most expensive stored procedures with the `sys.dm_exec_procedure_stats` view.
- Consider enabling the Query Store in your system and analyzing the data you collect. (This may or may not be feasible if you already use external monitoring tools.)
- Analyze data from third-party monitoring tools and cross-check it with SQL Server data.
- Analyze the overhead that inefficient queries introduce in the system. Correlate queries' resource consumption with wait statistics and server load.
- Optimize queries if you determine this is needed.

Chapter 5. Intro to Query Tuning

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 5 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

The topic of query optimization and tuning is easily worth another book. Indeed, there are many books available already and I encourage you to read them and master your skills. I will not try to duplicate them here; instead, this chapter will cover some of the most important concepts you need to understand to tune the queries.

You cannot master the process of query optimization without understanding the internal index structure and patterns that SQL Server uses to access data. This chapter thus begins with a high-level overview of B-Tree indexes and seek-and-scan operations.

Next, I discuss statistics and cardinality estimations, along with ways to read and analyze execution plans.

Finally, I cover several common issues you might encounter during the query tuning process, offering advice on how to address them and index the data.

Data Storage and Access Patterns

Modern SQL Server versions support three data storage and processing technologies. The oldest and most commonly used one is *row-based storage*. With row-based storage, all table columns are combined together into the data rows that reside on 8KB data pages. Logically, those data rows belong to *B-Tree indexes* or *heaps* (which we'll discuss in a moment).

Starting with SQL Server 2012, you can store some indexes or entire tables in columnar format using *column-based storage*. The data in such indexes is heavily compressed and stored on a per-column basis. This technology is optimized and provides great performance for read-only analytical queries that scan large amounts of data. Unfortunately, it does not scale well in an OLTP workload.

Finally, starting with SQL Server 2014, you can use *In-Memory OLTP* and store data in *memory-optimized tables*. The data in such tables resides completely in memory and is great for heavy OLTP workloads.

NOTE

You can use all three technologies—row-based, column-based, and memory-optimized tables—together, partitioning data between them. This approach is extremely useful when you need to support heavy OLTP and analytical workloads in the same system. I cover that architecture pattern in detail in my book *Pro SQL Server Internals*.

Row-based storage is the default and by far most common storage technology in SQL Server. The `CREATE TABLE` and `CREATE INDEX` statements will store data in a row-based format unless you specify otherwise. It can handle moderate OLTP and analytical workloads and introduces less database administration overhead than columnstore indexes and In-Memory OLTP.

In this chapter, I will focus on row-based storage and queries that work with B-Tree indexes. I will discuss troubleshooting aspects of columnstore indexes and In-Memory OLTP in Chapters 8 and 14.

Let's look at how SQL Server stores data in row-based storage.

Row-Based Storage Tables

Internally, the structure of a row-based table consists of multiple elements and internal objects, as shown in Figure 5-1.

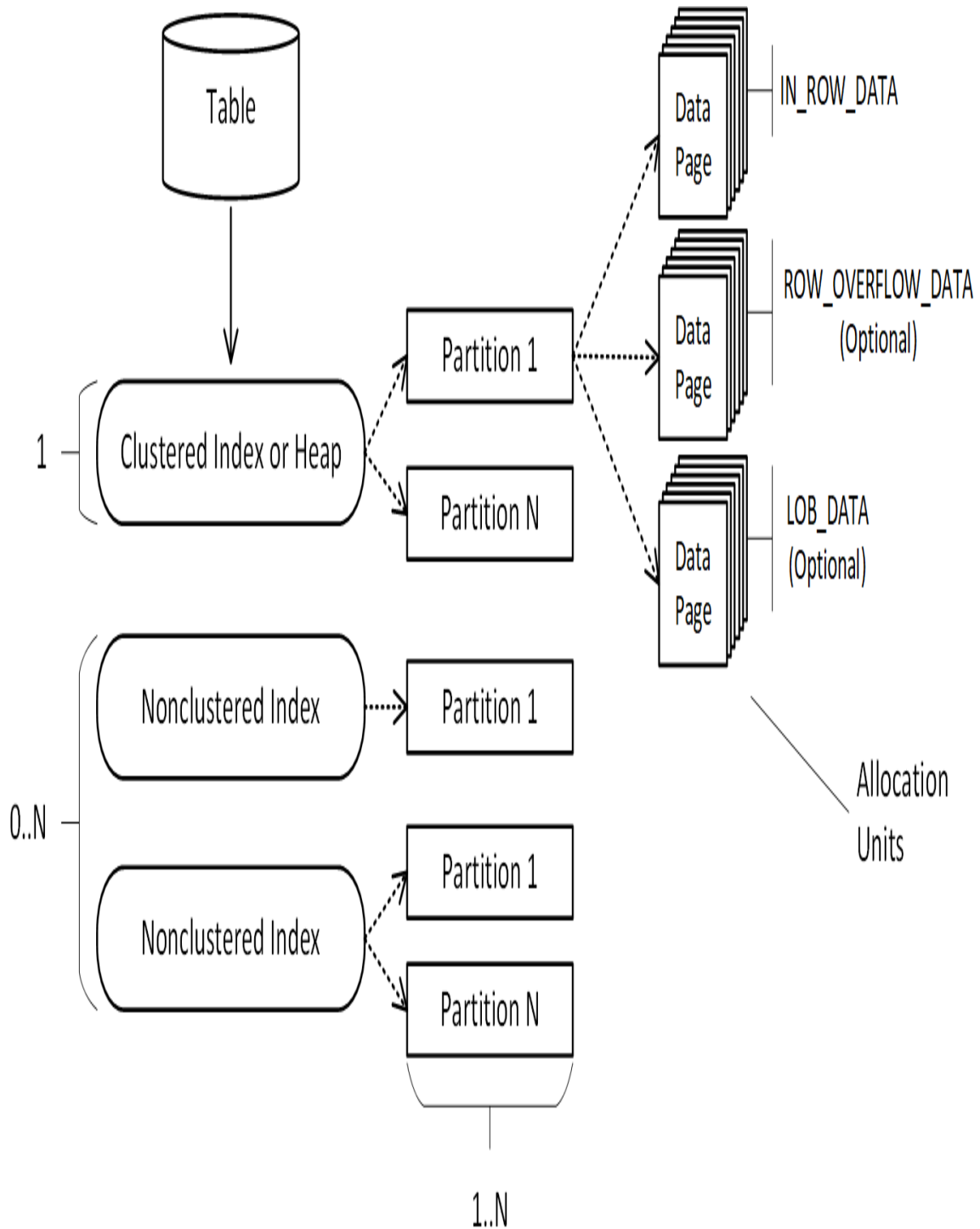


Figure 5-1. Internal table structure

The data in the tables is stored either completely unsorted (those tables are called *heap tables* or *heaps*) or sorted based on the value of a clustered index key, when such an index is defined.

I am not going to dive deep into detail, but as a general rule, it is better to avoid heaps and define clustered indexes on your tables. There are some edge cases when heap tables may outperform tables with clustered indexes; nevertheless, heaps have several shortcomings. In most cases, you'll get better performance when tables have clustered indexes.

In addition to a single clustered index, every table may have a set of *nonclustered indexes*: separate data structures that store copies of the data from a table sorted according to index key columns. For example, if a column is included in two nonclustered indexes, SQL Server would store that data three times - once in a clustered index or heap, and once in each nonclustered index.

While SQL Server allows you to create large numbers of nonclustered indexes, doing so is not a good idea. In addition to storage overhead, SQL Server needs to insert, update, or delete data in each nonclustered index during data modifications maintaining multiple copies of the data.

Internally, each index (and heap) consists of one or more partitions. You can think of each partition as an internal data structure (index or heap) that is independent from other partitions in the table. You can use a different partition strategy for every index in the table; however, it is usually beneficial to partition all indexes in the same way, aligning them with each other.

As I mentioned above, the actual data is stored in data rows on 8KB data pages with 8,060 bytes available to users. The data from all columns is stored together with exception when column data does not fit on the data page.

The data pages combine into three different categories called *allocation units*.

IN_ROW_DATA allocation unit pages store the main data row objects, which consist of internal attributes and the data from fixed-length columns (such as int, datetime, float, etc.). The in-row part of a data row must fit on a single data page, so it cannot exceed 8,060 bytes. The data from variable-length columns, such as (n)varchar(max), (n)varbinary(max), xml and

others, may also be stored in-row in the main row object when it fits into this limit.

When variable-length data does not fit in-row, SQL Server stores it off-row on different data pages, referencing them through in-row pointers. Variable-length data that exceeds 8,000 bytes is stored on LOB_DATA allocation unit data pages (LOB stands for large objects). Otherwise, the data is stored in ROW_OVERFLOW_DATA allocation unit pages.

I'd like to repeat a well-known piece of advice here: Do *not* use retrieve unnecessary columns in SELECT statements, especially with the SELECT * pattern. This may lead to additional I/O operations to get data from off-row pages, and may also defer usage of covered indexes, as you'll see later in the chapter.

Finally, SQL Server logically groups sets of eight pages into 64KB units called *extents*. There are two types of extents available. *Mixed extents* store data that belongs to different objects. *Uniform extents* store the data for the same object. By default, when a new object is created, SQL Server stores the first eight object pages in mixed extents. After that, all subsequent space allocation for that object is done with uniform extents.

You can disable mixed extents allocation with server-level trace flag T1118. In SQL Server 2016 and above, you can control it on the database level with MIXED_PAGE_ALLOCATION database option. Turning mixed extents off will reduce the number of modifications in the system tables when a new table is created. In users' databases, doing so rarely gives you noticeable benefits; however, it may significantly improve tempdb throughput in busy OLTP systems. You can disable mixed extents allocation with trace flag T1118 in old versions of SQL Server (prior to 2016). From SQL Server 2016 on, tempdb stopped using mixed extents, so you don't need to enable that trace flag in the system.

Next, let's look at the structure of the indexes.

B-Tree Indexes

Clustered and nonclustered indexes have a very similar internal format called *B-Tree*. Let's create an example table called Customers, defined in Listing 5-1. The table has the clustered index defined on CustomerId and nonclustered index on Name columns.

Example 5-1. Customers table

```
CREATE TABLE dbo.Customers
(
    CustomerId INT NOT NULL,
    Name NVARCHAR(64) NOT NULL,
    Phone VARCHAR(32) NULL,
    /* Other Columns */
);
CREATE UNIQUE CLUSTERED INDEX IDX_Customers_CustomerId
ON dbo.Customers(CustomerId);
CREATE NONCLUSTERED INDEX IDX_Customers_Name
ON dbo.Customers(Name);
```


CONSTRAINTS VERSUS INDEXES

As you may have noticed, I defined the clustered index on the table instead of creating the primary key constraint. I did this on purpose. I always consider constraints to be the part of a logical database design, which defines entities and their key attributes. Indexes, on the other hand, belong to the physical database design and physical data structures of the database.

By default, SQL Server creates unique clustered indexes for primary key constraints. However, you can—and in many cases should—mark primary keys as nonclustered, which will make them unique nonclustered indexes.

With the exception of a few SQL Server features that require you to define primary keys, the choice between constraints and indexes is a matter of personal preference. Primary and unique constraints are implemented as indexes internally and behave the same. During performance tuning, you'll work with indexes, so I am not going to reference constraints in this book. Their absence from the discussion does not mean that you should not use constraints in your databases.

The logical structure of the clustered index is shown in Figure 5-2.

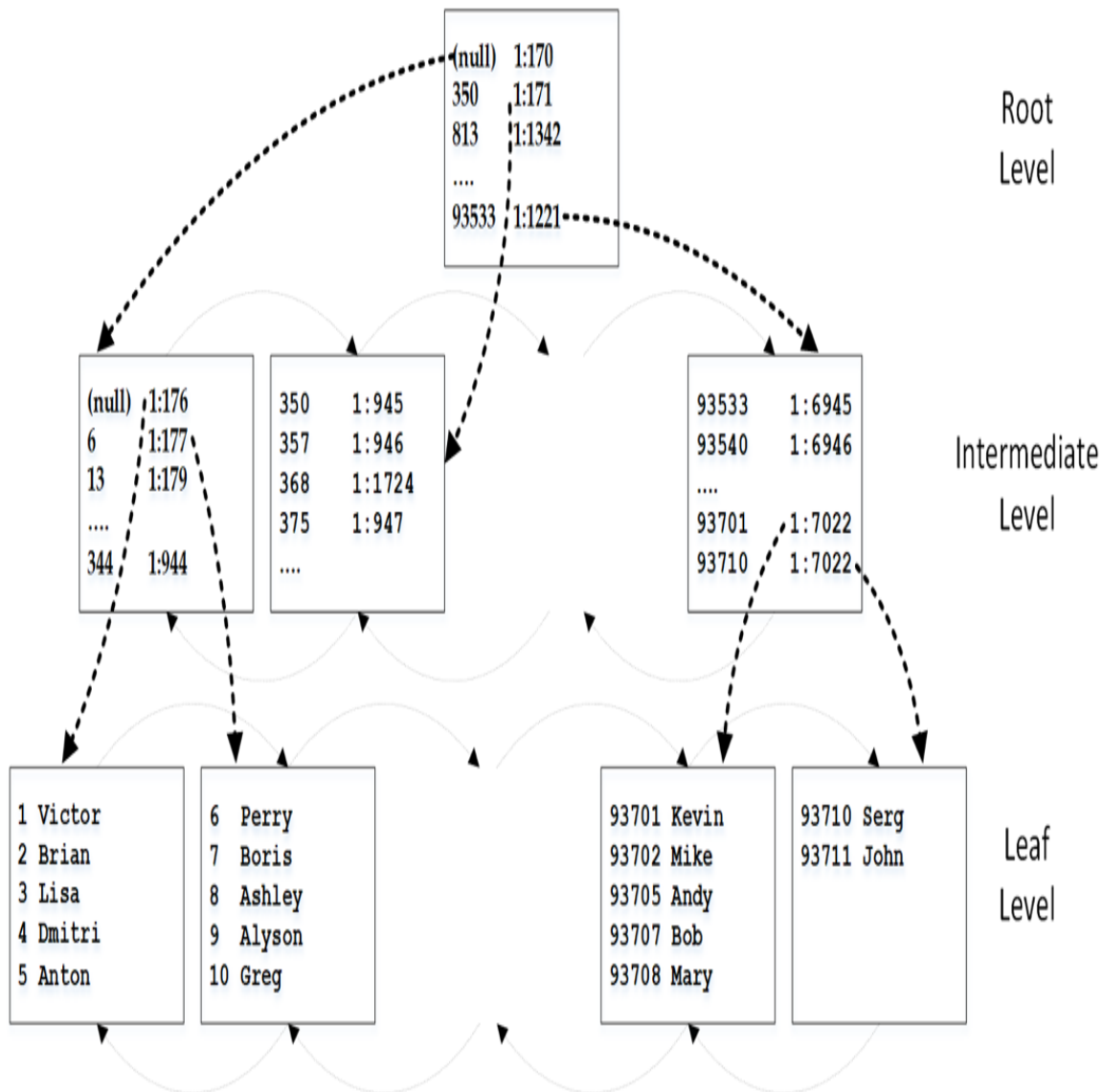


Figure 5-2. B-Tree index

The bottom level of the index is called the *leaf level*. It stores the data sorted according to the index key value. If it is a clustered index, the leaf level stores all data from the table sorted based on the clustered key. To be exact, the leaf level includes IN_ROW data only, which may reference off-row column data on the other pages.

If all data in the index fits into a single data page, the index will consist of that single leaf page. Otherwise, SQL Server will start to build *intermediate levels* of the index. Each row on an intermediate level page references the page from the level below and contains the minimal value of the key in the

referenced page, as well as its physical address (FileId:PageId) in the database. The only exception is the very first row, which stores NULL instead of the minimal value of the key.

SQL Server continues to build intermediate levels until it ends with a level with a single page. This level is called the *root level*; it is the entry point to the index.

The pages on each level of the index are linked into the double-linked list. Each page knows the previous and the next page in the index. This allows SQL Server to scan the indexes forward and backward. (Keep in mind, however, that the backward scan may be less efficient, since SQL Server does not use parallelism during that operation.)

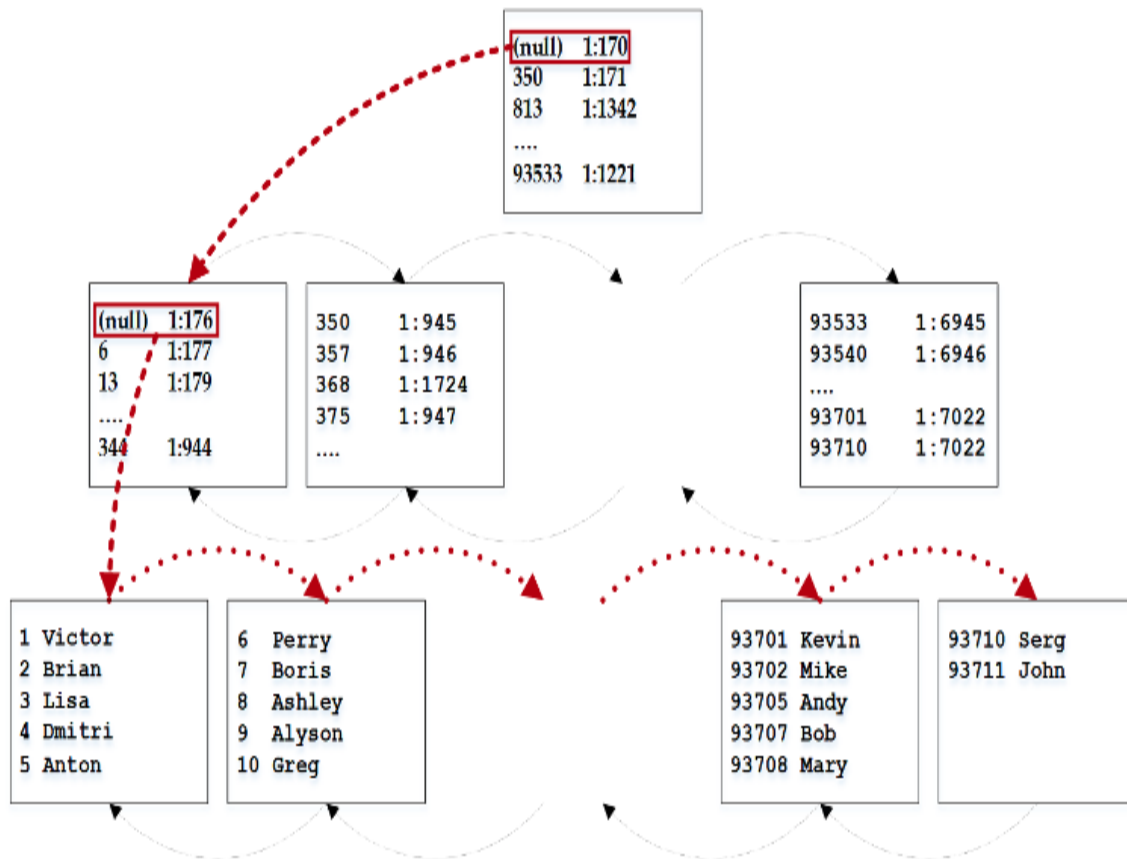
SQL Server can access data in the index either through *index scan* or *index seek*. With scans, there are two ways SQL Server can do that.

The first is an *allocation order scan*. SQL Server tracks the extents that belong to each index in the database through system pages called *Index Allocation Maps* (IAM). It reads the data pages from the index in random order according to index allocation data. This method, however, could introduce data consistency problems and is rarely used.

The second, more common method is called an *ordered scan*. Let's assume that you want to run the SELECT Name FROM dbo.Customers query. All data rows reside on the leaf level of the index and SQL Server can scan it and return the rows to the client.

SQL Server starts with the root page of the index and reads the first row from there. That row references the intermediate page with the minimum key value from the table. SQL Server reads that page and repeats the process until it finds the first page on the leaf level. Then SQL Server starts to read rows one by one, moving through the linked list of the pages until all rows have been read (Figure 5-3).

Figure 5-3. Index scan



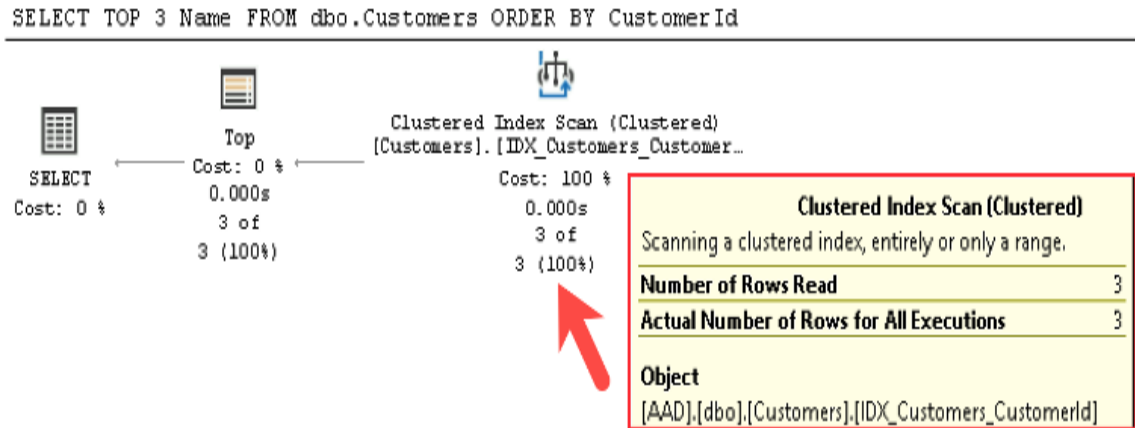
Obviously, in real life, it may become more complicated. In some cases, a query may simultaneously scan multiple parts of the index with parallel execution plans. In others, SQL Server may combine multiple index scans of simultaneously running queries together into the single physical index scan. Nevertheless, when you see the *Index Scan* operator in the execution plan, you can assume that this operator will access all data from the index.

There is one exception, however: when the plan has an index scan immediately following the *Top* operator. In that case, the scan operator will stop after it returns the number of rows requested by TOP and will not access the entire table. Usually, this happens if your query does not have an ORDER BY clause, or if the ORDER BY clause matches the index key.

Figure 5-4 shows part of the execution plan of SELECT TOP 3 Name FROM dbo.Customers ORDER BY CustomerId query. The *Number of*

Rows Read and *Actual Number of Rows* properties in the *Index Scan* operator indicate that the scan stopped after it read three rows.

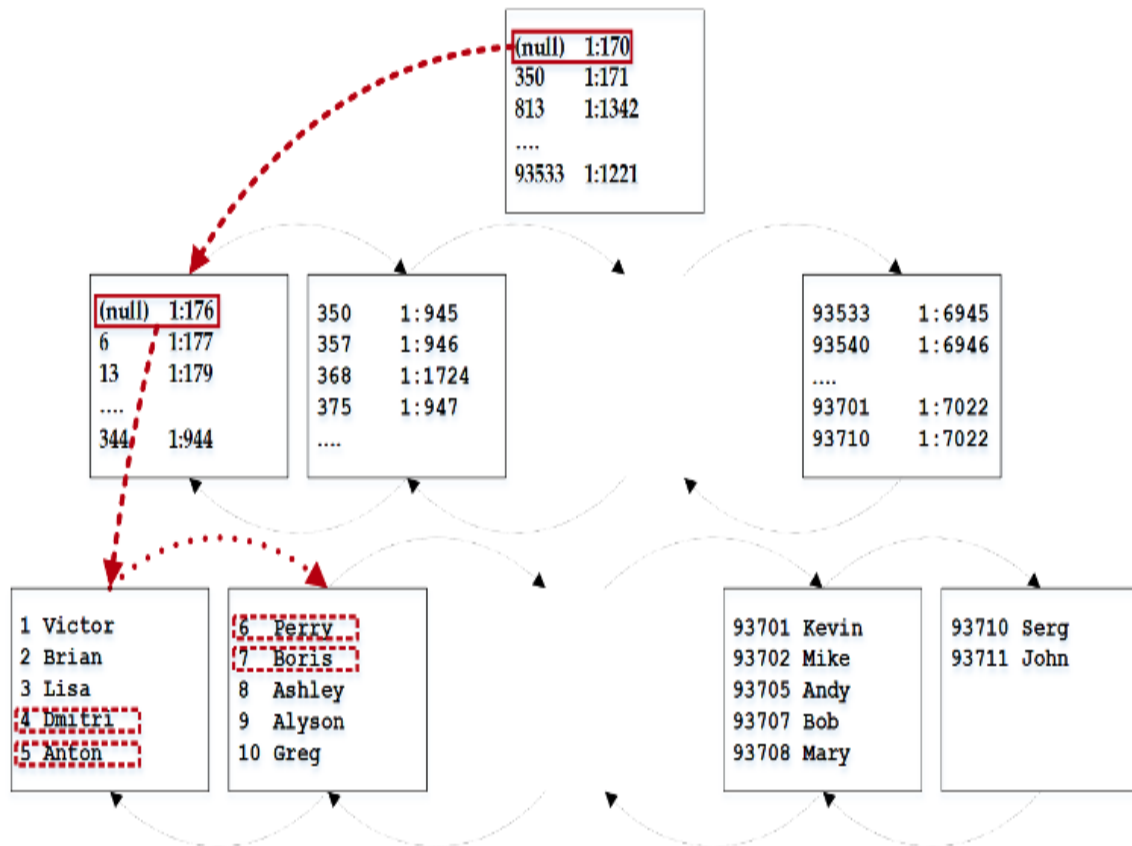
Figure 5-4. Top and Index Scan operators



Top and Index Scan operators

As you can imagine, reading all data from the large index is an expensive operation. Fortunately, SQL Server can access subset of the data by using the *index seek* operation. Say you want to run the following query: `SELECT Name FROM dbo.Customers WHERE CustomerId BETWEEN 4 AND 7`. Figure 5-5 illustrates how SQL Server might process it.

Figure 5-5. Index seek



In order to read the range of rows from the table, SQL Server needs to find the row with the minimum value of the key from the range, which is 4. SQL Server starts with the root page, where the second row references a page with a minimum key value of 350. That is greater than the key value you're looking for, so SQL Server reads the intermediate-level data page (1:170) referenced by the first row on the root page.

Similarly, the intermediate page leads SQL Server to the first leaf-level page (1:176). SQL Server reads that page, then it reads the rows with CustomerId equal 4 and 5, and finally, it reads the two remaining rows from the second page.

Technically speaking, there are two kinds of index seek operations:

Point-lookup

The first is called a *point-lookup* (or, sometimes, *singleton lookup*) where SQL Server seeks and returns a single row: for example, the WHERE CustomerId = 2 predicate is point-lookup operation.

Range scan

The other type is called a *range scan*. It requires SQL Server to find the lowest or highest value of the key and scan the set of rows (either forward or backward) until it reaches the end of the scan range. The predicate WHERE CustomerId BETWEEN 4 AND 7 leads to the range scan. Both cases are shown as *Index Seek* operators in the execution plans.

As you can guess, index seek is more efficient than index scan because SQL Server usually processes just a subset of rows and data pages, rather than scanning the entire index. However, the *Index Seek* operator in the execution plan may be misleading and represent an inefficient range scan that reads a large number of rows, or even the entire index. I will talk about this condition later in the chapter.

There is a concept in relational databases called *SARGable predicates*, which stands for *Search Argument-able*. SARGable predicates allow SQL Server to isolate a subset of the index key to process. In a nutshell, with SARGable predicate, SQL Server can determine a single key value or a range of index key values to read during a predicate evaluation and utilize Index Seek operation when the index exists.

Obviously, it is beneficial to write queries using SARGable predicates and utilize index seek whenever possible. This is done using operators, which include =, >, >=, <, <=, IN, BETWEEN, and LIKE (for prefix matching). Non-SARGable operators include NOT, <>, LIKE (when not prefix matching), and NOT IN.

Predicates are also non-SARGable when using functions (system or non-inlined user-defined) against the table columns. SQL Server must call the function for every row it processes to evaluate the predicate. This prevents you from using an index seek.

The same applies to data type conversions where SQL Server uses the `CONVERT_IMPLICIT` internal function. One common example is using the Unicode `nvarchar` parameter in the predicate with a `varchar` column. Another case is when you have different data types in the columns that participate in the join predicate. Both cases could lead to an index scan, even when the predicate operator appears to be SARGable.

Composite Indexes

Indexes with multiple key columns are called *composite indexes*. The data in the composite indexes is sorted per column, from left to right. Figure 5-6 shows the structure of a composite index defined on `LastName` and `FirstName` columns in the table. The data is first sorted based on `LastName` (the leftmost column), then on `FirstName` within each `LastName` value.

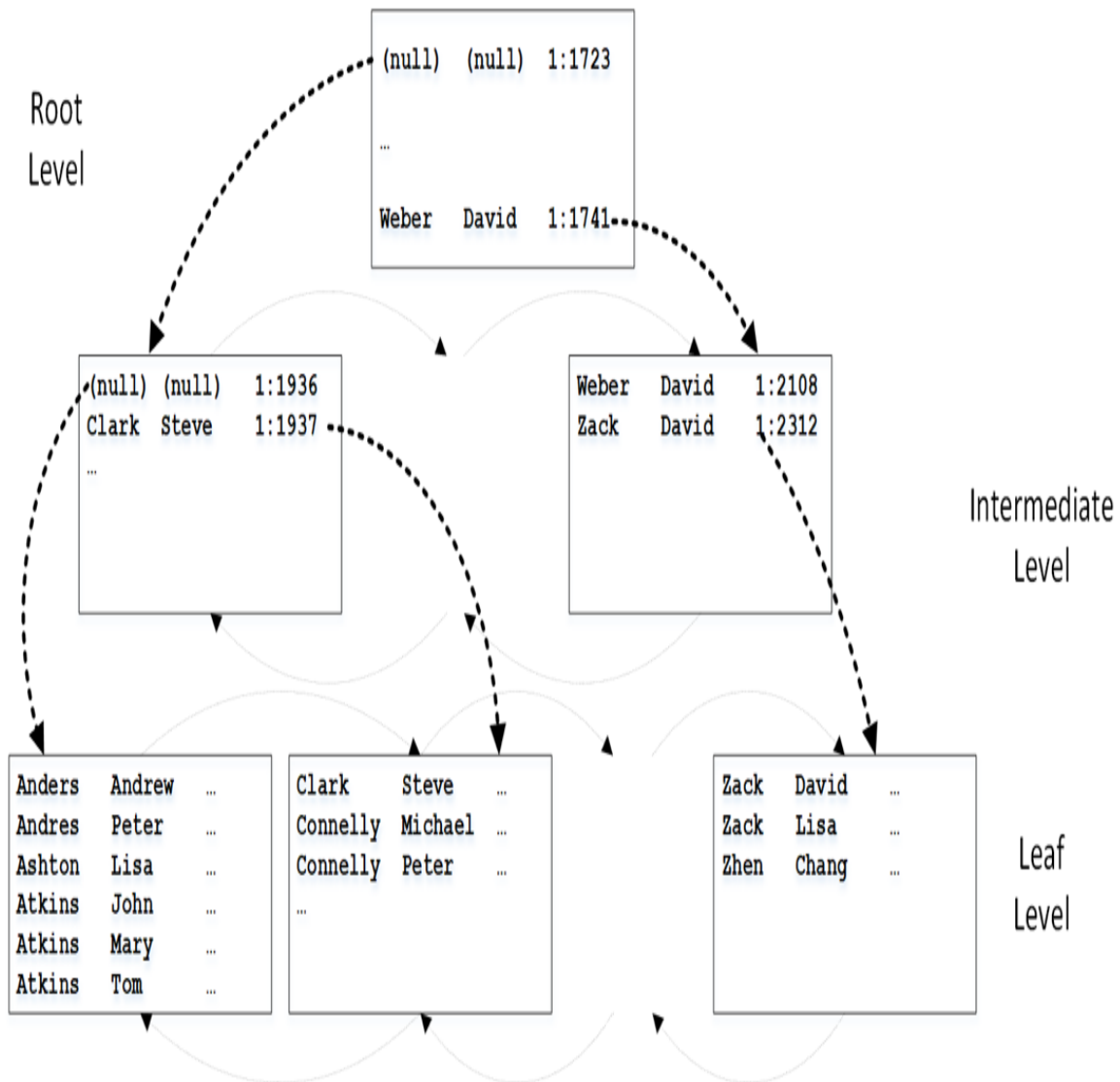


Figure 5-6. Composite indexes

The SARGability of a composite index depends on the SARGability of the predicates on the leftmost index columns, which allow SQL Server to determine and isolate the range of the index keys to process.

Table 5-1 shows examples of SARGable and non-SARGable predicates, using the index from Figure 5-6.

*T
a
b
l
e
5
-
1
.
T
a
b
l
e
5
-
1
.
C
o
m
p
o
s
i
t
e
i
n
d
e
x
a*

n
d
S
A
R
G
a
b
i
l
i
t
y

SARGable predicates Non-SARGable predicates

```
LastName  
= 'Clark'  
AND  
FirstName = 'Steve'
```

```
LastName  
<> 'Clark'  
  
AND  
  
FirstName = 'Steve'
```

```
LastName  
= 'Clark'  
AND  
FirstName <> 'Steve'
```

```
LastName  
LIKE '%  
ar  
%'  
AND  
FirstName = 'Steve'
```

```
LastName      FirstName = 'Steve'
= 'Clark'
```

```
LastName
LIKE 'C1%'
```

Nonclustered Indexes

While a clustered index specifies how data rows are sorted in a table, nonclustered indexes define a separate sorting order for a column or set of columns, storing them as separate data structures.

Think about a book, for example. Its page numbers represent the book's *clustered* index. The term index—that's the one labeled "Index" at the end of the book—lists terms from the book in alphabetical order. Each term references the numbers of each page where the term is mentioned. It is thus a *nonclustered* index of the terms.

When you need to find a term in the book, you can look it up in the term index. It is a fast and efficient operation, because terms are sorted in alphabetical order. Next, you can quickly find the pages on which the terms are mentioned using the page numbers specified there. Without the term index, your only choice would be to read the entire book, page by page, until you find all references to the term.

As I have noted, clustered and nonclustered indexes use a similar B-Tree structure. Figure 5-7 shows the structure of a nonclustered index (right) on the Name column we created in Listing 5-1. It also shows the clustered index (left), for reference.

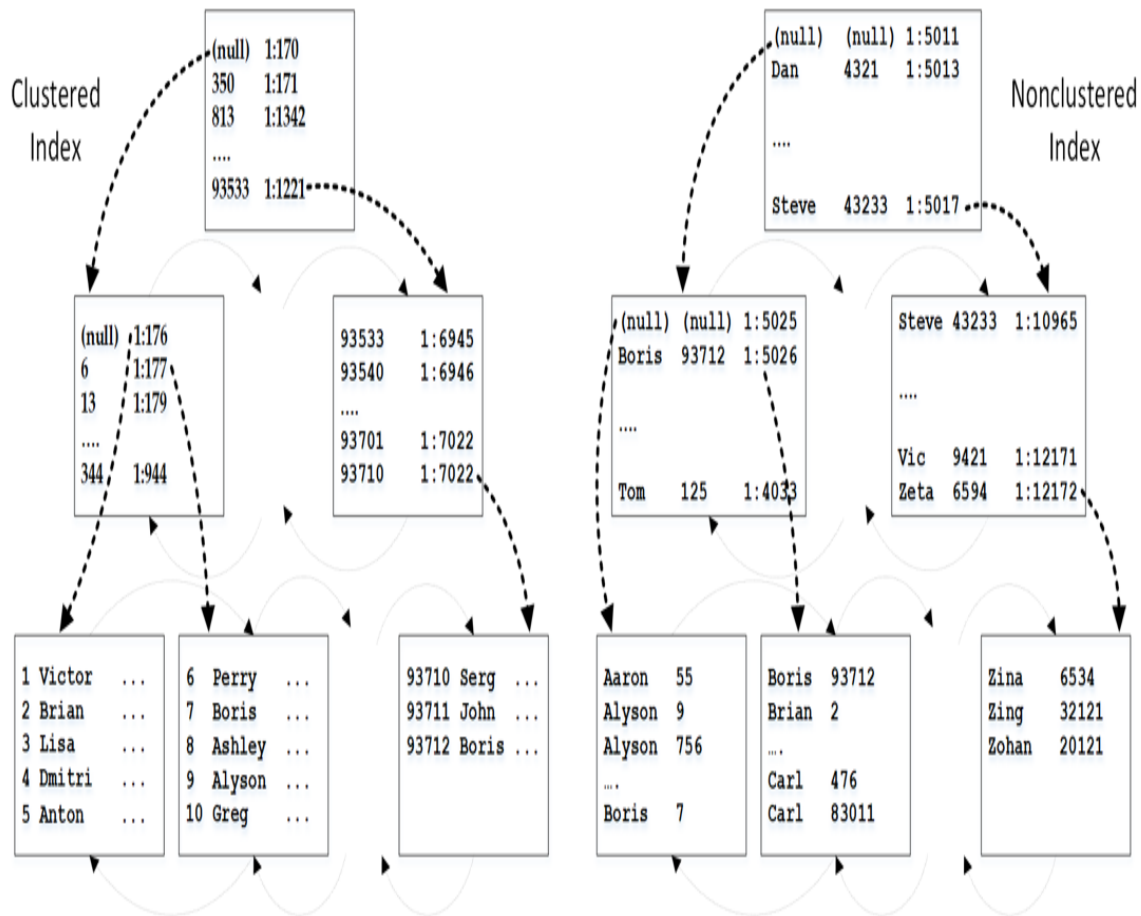


Figure 5-7. Clustered and nonclustered indexes in Customers table

The leaf level of the nonclustered index is sorted based on the value of the index key (Name). Every row on the leaf level includes the key value and *row-id* value. For tables with a clustered index, *row-id* represents the value of the clustered index key of the row.

This is a very important thing to remember: nonclustered indexes do *not* store information about physical row location when a table has a clustered index defined. They store the *value of the clustered index key* instead. This also means that nonclustered indexes include the data from clustered index key columns *even if you don't explicitly add those columns to the index definition*.

Like clustered indexes, the intermediate and root levels of nonclustered indexes store one row per page from the level they reference. That row

consists of the physical address and the minimum value of the key from the page. In non-unique indexes, it also stores the row-id of such a row.

Let's look at how SQL Server uses nonclustered indexes. I'll run the following query: `SELECT Name, Phone FROM dbo.Customers WHERE Name = 'Boris'`. Figure 5-8 shows that process.

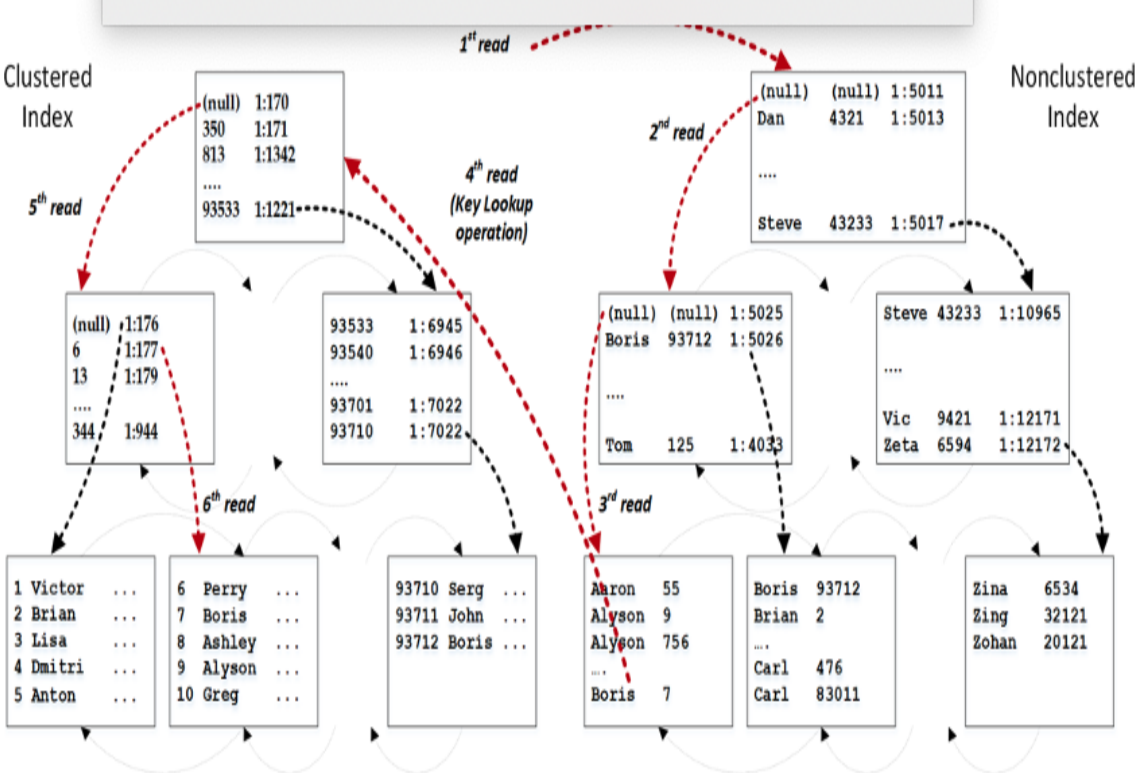


Figure 5-8. Nonclustered Index Usage: Part 1

Similar to the clustered index, SQL Server starts with the root page of the nonclustered index. The key value *Boris* is less than *Dan*, so SQL Server goes to the intermediate page referenced from the first row in the root-level page.

The second row of the intermediate page indicates that the minimum key value on the page is *Boris*, although the index had not been defined as unique and SQL Server does not know if there are other *Boris* rows stored on the first page. As a result, it goes to the first leaf page of the index and finds the row with the key value *Boris* and row-id equaling 7 there.

In our case, the nonclustered index does not store any data besides CustomerId and Name, and SQL Server needs to traverse the clustered index tree and obtain the data for Phone column from there. This operation is called *key lookup* (*RID lookup* in heap tables).

In the next step shown in Figure 5-9, SQL Server comes back to the nonclustered index and reads the second page from the leaf level. It finds another row with the key value *Boris* and row-id of 93712, and it performs key lookup again.

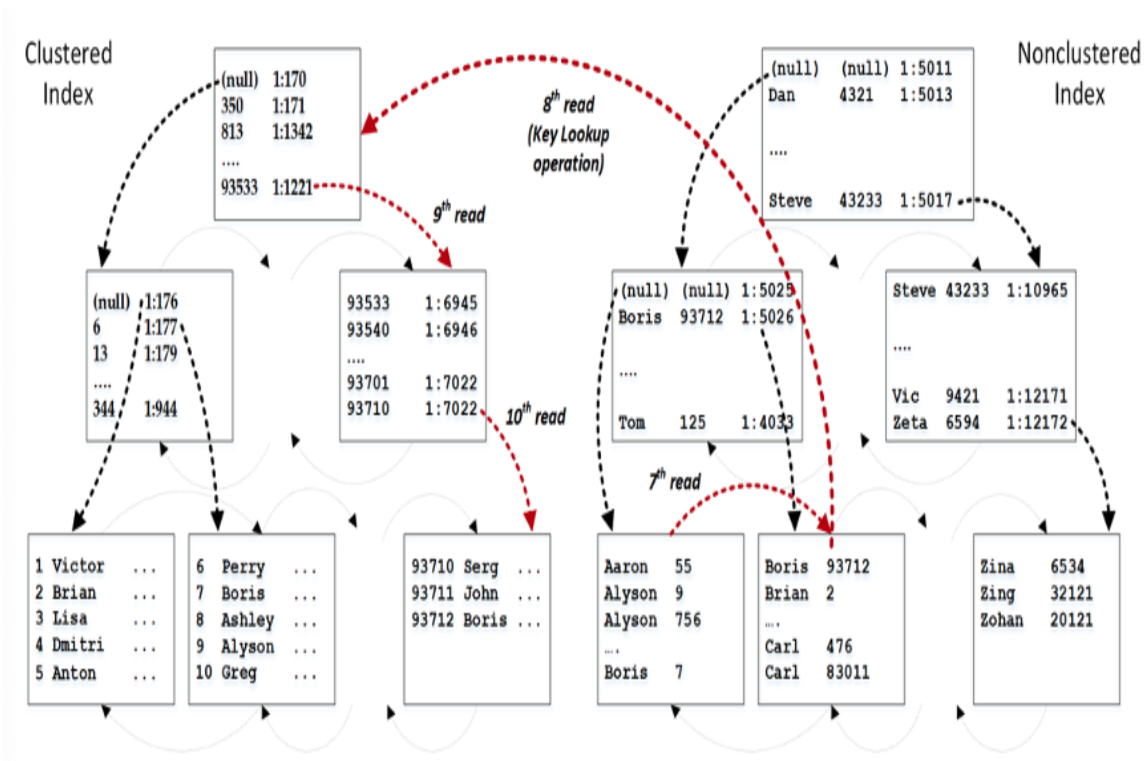


Figure 5-9. Nonclustered Index Usage: Part 2

As you can see from Figure 5-9, SQL Server had to perform 10 reads even though query returned just two rows. The number of I/O operations can be calculated based on the following formula:

(# of levels in nonclustered index) + (number of pages read from the leaf level of nonclustered index) + (number of rows found) * (# of levels in clustered index)

A large number of rows found (key lookup operations) leads to a large number of I/O operations, which makes using a nonclustered index inefficient.

As a result, SQL Server is very conservative in choosing nonclustered indexes when it expects that a large number of key lookup operations will be required. It may choose to scan a clustered or another nonclustered index instead. The threshold when SQL Server decides not to use a nonclustered index with key lookup varies, but it is very low – often a fraction of a percent of the total number of rows in the table.

The same applies to RID lookup operations. Nonclustered indexes in heap tables store physical address of the row in row-id. Technically, SQL Server can access the row in a heap through the single read operation; however, it is still expensive. Moreover, if the new version of the row does not fit into the old data page during an update, SQL Server will move it to another place, referencing it through another structure called *forwarding pointer*, which contains the address of the new (updated) version of the row. Nonclustered indexes will continue to reference forwarding pointers in row-id, and the RID lookup may lead to multiple read operations to access the row.

Index Fragmentation

SQL Server always maintains the order of the data in the index, inserting new rows on the data pages to which they belong. If the data page does not have enough free space, SQL Server allocates a new page and places the row there, adjusting the pointers in the double-linked page list to maintain logical sorting order in the index. This operation is called *page split*, and it leads to index fragmentation, as you'll see in this section.

Figure 5-10 illustrates this condition. When the original page does not have enough space to accommodate the new row, SQL Server performs a page split, moving about half of the data from original page to the new page and adjusting page pointers afterward.

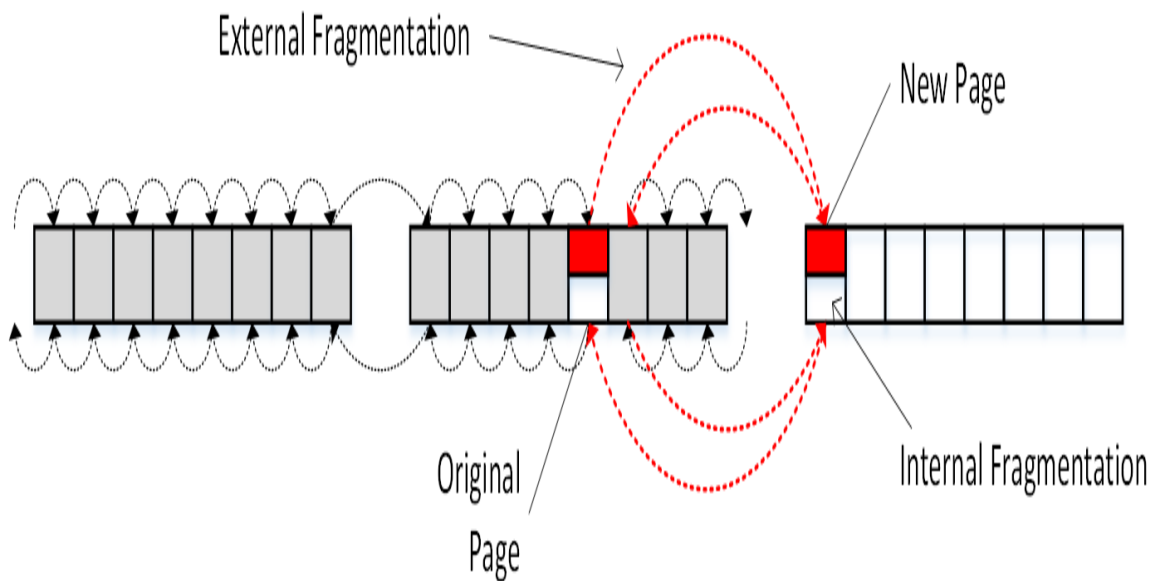


Figure 5-10. Page split

Page splits can also occur during data modifications. When an update cannot be done in place—for example, during a data row size increase—SQL Server performs a page split and moves updated and subsequent rows from that page to another page. It maintains the index sorting order through the page pointers.

There are two kinds of index fragmentation – internal and external.

External fragmentation

External fragmentation means that the logical order of the pages does not match their physical order in the data files, and/or that logically subsequent pages are not located in the same or adjacent extents.

External fragmentation forces SQL Server to jump around while reading the data from the disk, which makes read-ahead less efficient and increases the number of physical reads required. The impact is higher with magnetic drives where random I/O is less efficient than sequential I/O.

Internal fragmentation

Internal fragmentation, on the other hand, means that data pages in the index have free space. As a result, the index uses more data pages to store data, which in turn increases the number of logical reads during query execution. In addition, SQL Server uses more memory in the buffer pool to cache index pages.

A small degree of internal fragmentation is not always bad. It reduces page splits during insert and update operations, when data is inserted into or updated in different pages in the index. A large degree of internal fragmentation, however, wastes index space and reduces the performance of the system.

You can analyze index fragmentation in the system with `sys.dm_db_index_physical_stats` **view**. The three most important columns from the result set are:

avg_page_space_used_in_percent

`avg_page_space_used_in_percent` shows the average percentage of the data storage space used on the page. This value shows you the internal index fragmentation.

avg_fragmentation_in_percent

`avg_fragmentation_in_percent` provides you with information about external index fragmentation. For tables with clustered indexes, it indicates the percent of out-of-order pages when the next physical page allocated in the index is different from the page referenced by the next-page pointer of the current page. For heap tables, it indicates the percent

of out-of-order extents, when extents are not residing continuously in data files.

fragment_count

`fragment_count` indicates how many continuous data fragments the index has. Every fragment constitutes the group of extents adjacent to each other. Adjacent data increases the chances that SQL Server will use sequential I/O and Read-Ahead while accessing the data.

The impact of index fragmentation can be offset by modern hardware, when servers have enough memory to cache the data in the buffer pool and fast flash-based I/O subsystems to read the data. While it is always beneficial to reduce fragmentation in the system, you need to analyze its impact when designing your index maintenance strategy.

To put things in perspective: If your system has low-activity hours during nights or weekends, use them for index maintenance. However, if your system handles thousands of transactions per second around the clock, do the analysis and estimate the benefits and downsides of different index maintenance strategies. Remember that index maintenance is an expensive operation and will add overhead to the system while it's running.

There are two index maintenance methods that reduce fragmentation: index reorganize and index rebuild. Let's look at each in turn.

Index reorganize

Index reorganize, often called *index defragmentation*, reorders leaf-level data pages into their logical order. It also tries to compress pages, reducing their internal fragmentation. This is an online operation that can be interrupted at any time without losing the operation's progress up

to the point of interruption. You can also reorganize indexes with the ALTER INDEX REORGANIZE command.

Index rebuild

Index rebuild (ALTER INDEX REBUILD), on the other hand, creates another copy of the index in the table. It is an offline operation, which will lock the table in non-Enterprise editions of SQL Server. In the Enterprise edition it can be done online, though it will still require a short table-level lock at the beginning and end of execution.

Microsoft [documentation](#) recommends rebuilding indexes if their external fragmentation (`avg_fragmentation_in_percent`) exceeds 30% and reorganize indexes for fragmentation between 5% and 30%. You can use those values as a rule of thumb; however, as I mentioned, it may be better to analyze and tune for your own use-cases.

Pay attention to the FILLFACTOR index property, which allows you to reserve some free space during index creation or rebuild, reducing page splits afterwards. Unless you have an ever-increasing append-only index, you should set FILLFACTOR below 100%. I usually start with 85 or 90% and fine-tune the values to get the least internal and external fragmentation in the index.

Finally, in heap tables, `sys.dm_db_index_physical_stats` view provides the information about forwarding pointers with the `forwarded_record_count` column. Tables with a large number of forwarding pointers are inefficient and need to be rebuilt with the ALTER TABLE REBUILD operation. However, the better option in most cases is converting them into clustered index tables.

NOTE

Ola Hallengren has provided a set of *scripts* that have become the de facto standard for database maintenance tasks. Consider using them in your systems.

Statistics and Cardinality Estimation

SQL Server stores information about data distribution in the index in internal objects called *statistics*. By default, SQL Server creates statistics for each index in the database and uses it during query optimization. Let's look what information is stored in the statistics first.

Listing 5-2 creates a table with clustered and nonclustered indexes and populates it with some data. Finally, it provides information about the statistics using the DBCC SHOW_STATISTICS command.

Example 5-2. Examining statistics

```
CREATE TABLE dbo.DBObjects
(
    ID INT NOT NULL IDENTITY(1,1),
    Name SYSNAME NOT NULL,
    CreateDate DATETIME NOT NULL
);
CREATE UNIQUE CLUSTERED INDEX IDX_DBObjects_ID
ON dbo.DBObjects(ID);

INSERT INTO dbo.DBObjects(Name,CreateDate)
    SELECT name, create_date FROM sys.objects ORDER BY name;

-- Creating some duplicate values
INSERT INTO dbo.DBObjects(Name, CreateDate)
    SELECT t1.Name, t1.CreateDate
    FROM dbo.DBObjects t1 CROSS JOIN dbo.DBObjects t2
    WHERE t1.ID = 5 AND t2.ID between 1 AND 20;

CREATE NONCLUSTERED INDEX IDX_DBObjects_Name_CreateDate
ON dbo.DBObjects(Name, CreateDate);

DBCC
SHOW_STATISTICS('dbo.DBObjects','IDX_DBObjects_Name_CreateDate');
```

Figure 5-11 shows the output of the listing (you may get different results in your system).

Figure 5-11. Statistics

	Name	Updated	Rows	Rows Sampled	Steps	Density
1	IDX_DBOObjects_Name_CreateDate	Mar 6 2021 9:40AM	137	137	102	1
	Average key length	String Index	Filter Expression	Unfiltered Rows	Persisted Sample Percent	
	42.49635	YES	NULL	137	0	

	All density	Average Length	Columns
1	0.008547009	30.49635	Name
2	0.008547009	38.49635	Name, CreateDate
3	0.00729927	42.49635	Name, CreateDate, ID

	RANGE_HI_KEY	RANGE_ROWS	EQ_ROWS	DISTINCT_RANGE_ROWS	AVG_RANGE_ROWS
1	CommandExecute	0	1	0	1
2	CommandLog	0	1	0	1
3	DBObjects	0	1	0	1
4	DatabaseBackup	0	1	0	1
5	DatabaseIntegrityCheck	0	21	0	1
6	EventNotificationError...	0	1	0	1
7	MSreplication_options	1	1	1	1
8	PK_C.....	0	1	0	1

As you can see, the DBCC SHOW_STATISTICS command returns three result sets. The first contains general metadata information about statistics,

such as name, update date, number of rows in the index at the time when the statistics were updated, and so on.

The second result set, called *density vector*, contains information about density for the combination of key values from the statistics (index). It is calculated based on the formula $(1 / \text{number of distinct values})$, and it indicates how many rows on average every combination of key values has. It is worth noting that although `IDX_DBOObjects_Name_CreateDate` index has two index keys, row-id (clustered index column)- ID- also presents in the index and is returned in the density vector.

The last and most important result set is called the *histogram*. It provides information about data distribution in the index. Each record in the histogram, called a *histogram step*, includes the sample key value from the left-most column from the statistics (index) and information about data distribution in the interval of values from the preceding to the current `RANGE_HI_KEY` value. It also includes the estimated number of rows in the interval (`RANGE_ROWS`), number of rows with key value equal to `RANGE_HI_KEY` (`EQ_ROWS`), number of distinct key values in the interval (`DISTINCT_RANGE_ROWS`), and average of rows per distinct key values (`AVG_RANGE_ROWS`).

SQL Server uses statistics information during query optimization estimating the number of rows that each operator in the execution plan would process and return to the next operator there. That process is called *cardinality estimation*.

Cardinality estimation greatly affects the execution plan. SQL Server uses it to choose the sequence of operators in the plan, indexes to access the data, type of join operators, and many other things. The efficiency of its execution plans greatly depends on the correctness of its cardinality estimation, and therefore on accurate statistics in the system.

There are three things you need to remember about statistics. First, and most important, SQL Server maintains the histogram and has information about data distribution only for left-most column of the index. There is no

information about data distribution for other index columns or for combinations of index column values.

The common advice you'll hear suggests using the most selective column as the left-most column in the composite indexes. While following this advice may sometimes improve the quality of cardinality estimations, don't follow it blindly. You need to analyze the queries, making sure that the predicates in left-most columns are SARGable and support efficient index seek operations.

The second important thing to remember about statistics is that the histogram stores, at most, 200 steps, regardless of the table size and whether the table is partitioned or not. This can affect cardinality estimations in large tables with uneven data distribution, since each step stores information about larger key intervals.

Finally, you need to know how SQL Server updates statistics. In databases with a compatibility level below 130 (as of SQL Server 2016), statistics are *only* updated automatically after 20 percent of the data in the index has changed. For example, in a table with 100 million rows, you would need to insert, delete or update index key columns in 20 million rows before an automatic update is triggered. This means that in large tables, statistics are rarely updated automatically and tend to become inaccurate over time.

Starting with a database compatibility level of 130, the statistics update threshold becomes dynamic. The percentage of changes that triggers the statistics update becomes smaller as the amount of the data in the table grows. You can force the same behavior for databases with older compatibility levels and in old versions of SQL Server with trace flag T2371. This is one of the trace flags I enable in every system.

Statistics Maintenance

Accurate, up-to-date statistics improve system performance. Analyze your statistics maintenance strategy when you perform system troubleshooting, and validate whether it provides you with accurate information.

You can rely on automatic statistics updates, maintain statistics manually, or combine both approaches. Index maintenance also affects statistics maintenance strategy, since index rebuild automatically updates statistics in the index. Index reorg, on the other hand, does not update it.

You can control whether SQL Server creates and updates statistics automatically at the database level with the *Auto Create Statistics* and *Auto Update Statistics* database options. When these are enabled, SQL Server automatically maintains statistics on all indexes except those that have `STATISTICS_NORECOMPUTE` enabled (it is disabled by default).

SQL Server may use different methods to update statistics. By default, it just samples the data from the index. This approach is lightweight, but it does not always provide accurate results. Alternatively, you can update statistics using the `UPDATE STATISTICS WITH FULLSCAN` statement, which will read the entire index.

You can also update the statistics specifying percent or number of rows to sample with `UPDATE STATISTICS WITH SAMPLE` statement. Obviously, the more data you read, the more I/O overhead you'll have on large indexes.

During query compilation, SQL Server detects whether statistics are outdated and may update them synchronously or asynchronously, based on the selected *Auto Update Statistics Asynchronously* database option. With synchronous updates, Query Optimizer defers query compilation until the update is done. With asynchronous updates, the query is optimized using old statistics while statistics are updated in background. You can keep your default synchronous statistics update unless your system requires extremely low response time from the queries.

The default automatic statistics maintenance is acceptable in many cases, as long as the database has a compatibility level of 130 or above, or T2371 is set. However, in some cases, you can also update statistics of the key indexes manually and/or run statistics update with `FULLSCAN` after hours.

It is usually beneficial to update statistics on the filtered indexes manually. Modifications of filtered columns do not count towards the statistics update

threshold, which may make automatic statistics maintenance inefficient.

NOTE

Filtered indexes allow you to filter subsets of data in the table, reducing index size and index maintenance cost. Read about them in the Microsoft [documentation](#).

Listing 5-3 shows you how to view statistics properties, such as when statistics were last updated and how many changes in the data have occurred since the last update. You can use it as part of your custom statistics maintenance in the system, if needed.

Example 5-3. Analyzing statistics properties

```
SELECT
    s.stats_id AS [Stat ID]
    ,sc.name + '.' + t.name AS [Table]
    ,s.name AS [Statistics]
    ,p.last_updated
    ,p.rows
    ,p.rows_sampled
    ,p.modification_counter AS [Mod Count]
FROM
    sys.stats s JOIN sys.tables t ON
        s.object_id = t.object_id
    JOIN sys.schemas sc ON
        t.schema_id = sc.schema_id
    OUTER APPLY
        sys.dm_db_stats_properties(t.object_id,s.stats_id) p
ORDER BY
    p.last_updated
```

NOTE

This section barely scratches the surface of statistics and their maintenance—I strongly recommend reading the Microsoft [documentation](#) to learn more about it.

Cardinality Estimation Models

As you already know, the quality of query optimization depends on accurate cardinality estimations. SQL Server must correctly estimate the number of rows in each step of query execution to generate an efficient execution plan. Accurate statistics go a long way in improving the estimations; however, they are just part of the picture.

During the cardinality estimation process, Query Optimizer relies on a set of assumptions that cover, among other things:

- data distribution in the tables
- the impact of different operators and predicates on the size of the output
- relations between multiple predicates in a single table
- correlation of the data in multiple tables during joins

These assumptions, along with cardinality estimation algorithms, define the cardinality estimation model used during optimization.

The original (legacy) cardinality estimation model was initially developed for SQL Server 7.0 and used exclusively until the release of SQL Server 2014. Aside from some minor improvements across versions, the model remained conceptually the same.

In SQL Server 2014, Microsoft released a new cardinality estimation model enabled in the databases with compatibility level of 120. That model uses different assumptions, which lead to different cardinality estimations and execution plans.

It is impossible to tell which model is better. Some queries behave better with the new model; others may regress when you upgrade. You can continue to use the legacy cardinality estimation model with new versions of SQL Server; however, it may be beneficial to upgrade at some point. Microsoft says it is not going to remove legacy model from SQL Server in the future, but it won't be enhanced, either.

Unfortunately, that's easier said than done, especially in large and complex systems. Changing the model may lead to massive changes in the execution plans, so you need to be prepared to detect and address regressions quickly. Fortunately, Query Store can simplify the process. You can collect the data before the change and force SQL Server to use old execution plans for those queries that regressed under the new model. Obviously, you'll still need to analyze and optimize them later.

You can control the cardinality estimation model with the database compatibility level. Keep in mind that new model may behave slightly differently in each compatibility level, starting with 120 (SQL Server 2014). Legacy models, on the other hand, will behave the same in each SQL Server version. Enabling the `QUERY_OPTIMIZER_HOTFIXES` database setting or setting `T4199` may also affect the estimations.

In SQL Server 2014, you can control the model with database compatibility levels or with trace flags. `T2312` and `T9481` force SQL Server to use new and legacy models respectively ignoring database compatibility level. In SQL Server 2016 and above, you can choose the model by setting the `LEGACY_CARDINALITY_ESTIMATION` database option.

When you perform the SQL Server version upgrade, I recommend doing it in phases to reduce the risk of regression. First, upgrade the server version, keeping the old cardinality estimation model in place. Validate that everything works as expected after the upgrade. Then you can consider changing the model. As mentioned, use Query Store as part of that process.

As a general rule, I do not recommend switching to the new cardinality estimation model in SQL Server 2014. I encountered several bugs in early builds of this version, which led to more regressed queries. If you do switch, install the latest service pack, enable `T4199` and carefully test the system.

Analyzing Your Execution Plan

The query optimization process in SQL Server, done by the Query Optimizer, generates a query execution plan. This plan consists of multiple operators that access and manipulate the data, achieving results for the query. The query tuning process, in a nutshell, requires us to analyze and improve execution plans for the queries.

Even though every database engineer is familiar with execution plans, I'd like to discuss several things related to query tuning. First, we need to look how SQL Server executes operators in the plan.

Row Mode and Batch Mode Execution

SQL Server has two processing methods for queries. The default, *row mode*, is traditionally used with row-based storage and B-Tree indexes. In this mode, each operator in the execution plan processes data rows one at a time, requesting them from child operators when needed.

Let's look at the simple example query shown in Listing 5-4.

Example 5-4. Row mode execution: Sample query

```
SELECT TOP 10 c.CustomerId, c.Name, a.Street, a.City, a.State,
a.ZipCode
FROM
    dbo.Customers c JOIN dbo.Addresses a ON
        c.PrimaryAddressId = a.AddressId
ORDER BY
    c.Name
```

This query would produce the execution plan shown in Figure 5-12. SQL Server selects all of the data from the Customers table, sorts it based on the Name column, gets the first 10 rows, joins it with the Addresses data, and returns it to the client.

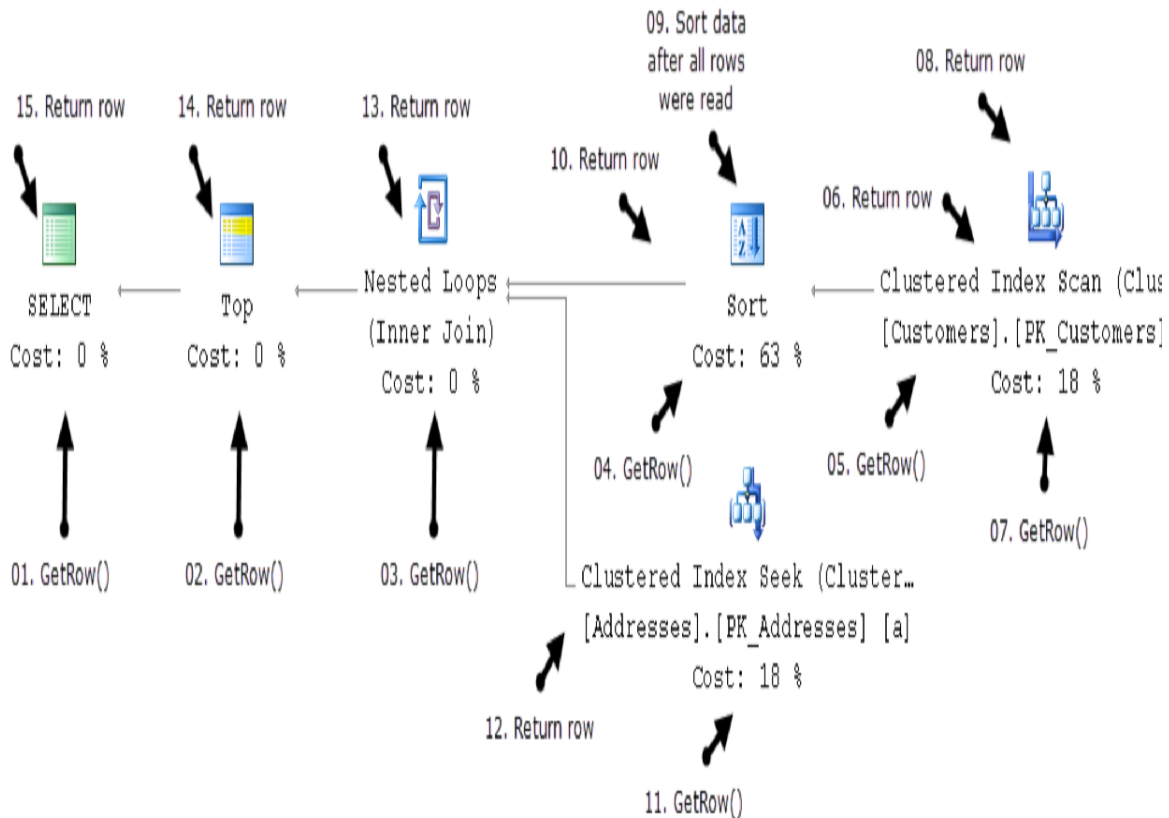


Figure 5-12. Row mode execution: Getting the first row

Let's analyze how SQL Server executes a query. The *Select* operator, which is the parent operator in the execution plan, calls the `GetRow()` method of the *Top* operator. The *Top* operator, in turn, calls the `GetRow()` method of the *Nested Loop Join*.

A *Join* operator gets the data from two different inputs. First, it calls the `GetRow()` method of the *Sort* operator. In order to sort, SQL Server needs to read all of the rows first. So the *Sort* operation calls the `GetRow()` method of the *Clustered Index Scan* operator multiple times, accumulating the results. The *Scan* operator, which is the lowest operator in the execution plan tree, returns one row from the *Customers* table per call. Figure 5-12 shows just two `GetRow()` calls, for simplicity's sake.

When all of the data from the *Customers* table has been read, the *Sort* operator performs sorting and returns the first row back to the *Join* operator, which calls the `GetRow()` method of the *Clustered Index Seek* operator on

the *Addresses* table after that. If there is a match, the *Join* operator concatenates data from both inputs and passes the resulting row back to the *Top* operator, which, in turn, passes it to *Select*.

The *Select* operator returns a row to the client and requests the next row by calling the *GetRow()* method of the *Top* operator again. The process repeats until the first 10 rows are selected. All operators keep their state and the *Sort* operator preserves the sorted data. It does not need to access the *Clustered Index Scan* operator again, as shown in Figure 5-13.

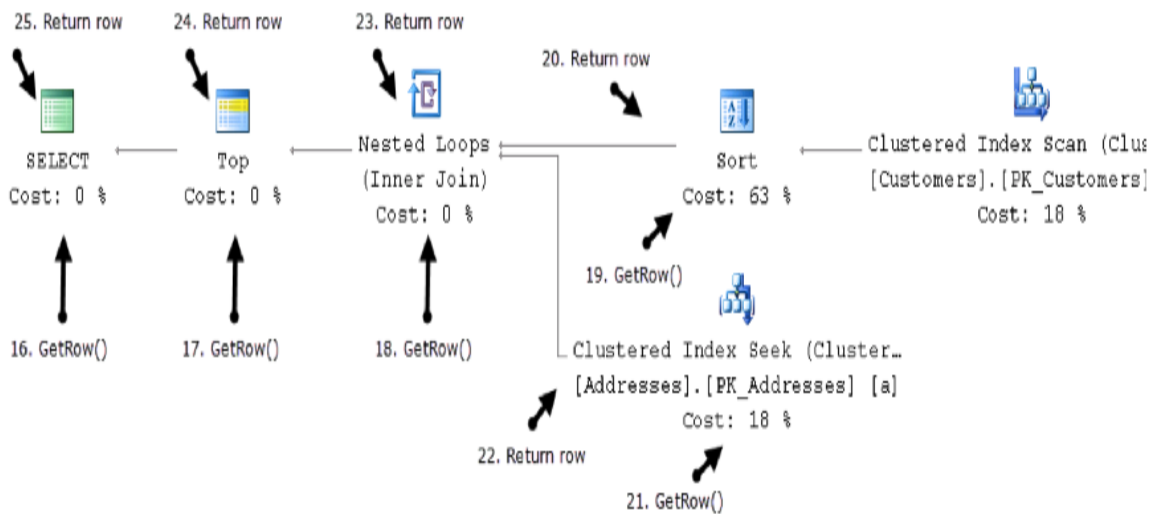


Figure 5-13. Row mode execution: Getting the next row

Each operator in the execution plan has multiple properties, the names of which may vary slightly in different versions of SSMS and in other applications. Let's look at the most important ones.

Actual Number of Rows and Number of Rows Read

Those two properties illustrate how many rows were returned by the operator and how many rows were processed during execution. For example, *Index Scan* operator with a predicate may process 1,000 rows, filtering out 950 of them. In that case, the property would show 50 and 1,000 rows, respectively.

Estimated Number of Rows and Estimated Number of Rows Read

Those two properties provide cardinality estimation data and indicate how many rows Query Optimizer expected the operator to return and process. The large discrepancy between *estimated* and *actual* metrics indicates a cardinality estimation error, which could lead to a suboptimal execution plan.

Number of Executions and Estimated Number of Executions

The *Number of Executions* metric indicates how many times the operator was executed. It does not correspond to the number of GetRow() calls but rather indicates how many times this part of the execution plan was processed. For example, in the plan shown in Figures 5-12 and 5-13, *Clustered Index Scan* in the Customers table would be executed once, while *Clustered Index Seek* in the Addresses table would be executed 10 times.

The *Estimated Number of Executions* metric shows the estimate used by Query Optimizer.

Startup Predicate

In some cases, operators may have a *Startup Predicate* property, which indicates the condition that needs to be met for the operator to execute. For example, a WHERE @ProvideDetails = 1 clause may generate *Filter* operator with *Startup Predicate* @ProvideDetails = 1. The execution plan subtree after the *Filter* operator may or may not be executed, depending on the @ProvideDetails parameter in runtime.

Unfortunately, row mode execution and per-row processing do not scale well with large analytical queries that process millions or even billions of rows. To address this, SQL Server 2012 introduced another execution model, called *batch mode execution*. This allows operators in the execution plans to process rows in batches. The process is optimized for large amounts of data and parallel execution plans.

Until SQL Server 2019, Query Optimizer did not consider batch mode execution unless at least one of the tables in the query had a columnstore index. This restriction was removed in SQL Server 2019, where batch mode can be used with row-based B-Tree indexes in databases with a compatibility level of 150. This does not mean that all execution plans will use batch mode; however, Query Optimizer will consider batch mode during optimization.

As with any feature that affects execution plans, batch mode can introduce regressions in some cases. You can enable and disable it on the database level with the `BATCH_MODE_ON_ROWSTORE` database option or on the query level with the `ALLOW_BATCH_MODE` and `DISALLOW_BATCH_MODE` query hints.

Finally, there is a trick that may enable batch mode execution on B-Tree tables in SQL Server 2016 and 2017: You can create empty and filtered nonclustered columnstore indexes on B-Tree tables that run large analytical queries. For example, if the table has ID column that stores only positive values, the following index will allow Query Optimizer to consider batch mode during optimization: `CREATE NONCLUSTERED COLUMNSTORE INDEX NCCI ON T WHERE ID < 0`.

You can see the operator's execution mode with *Actual Execution Mode* property in the execution plan. *Actual Number of Batches* will tell you how many batches were processed. However, the query tuning strategy would be the same regardless of the execution mode.

Live Query Statistics and Execution Statistics Profiling

There are several tools that allow you to analyze execution plans. In addition to the well-known SSMS, you can use another freeware tool from **Microsoft**—Azure Data Studio. Despite the name, it works perfectly well with on-prem instances of SQL Server and can be installed on other operating systems besides Windows.

I consider Azure Data Studio to be targeted to developers rather than database administrators. Nevertheless, it provides basic database administration and tuning features and can be expanded with multiple third-party extensions. Some extensions will even bring support of other database platforms in addition to SQL Server.

I consider **SentryOne Plan Explorer** a must-have freeware tool for query tuning. This tool focuses on execution plan analysis. I find it more advanced and easier to use than SSMS. I suggest you download and test it if you have not done so already.

SSMS has another very useful feature called *Live Query Statistics*. This feature allows you to monitor query execution in runtime, detecting possible inefficiencies in the execution plan.

Figure 5-14 shows an example of the Live Query Statistics window in SSMS (screenshot is copied from Microsoft **documentation**). The operators with solid lines have been completed. The dotted lines represent the tree of the operators that are currently executing. You can also see the estimated progress of each active operator, along with the actual and estimated numbers of rows. All metrics are updating during query execution.

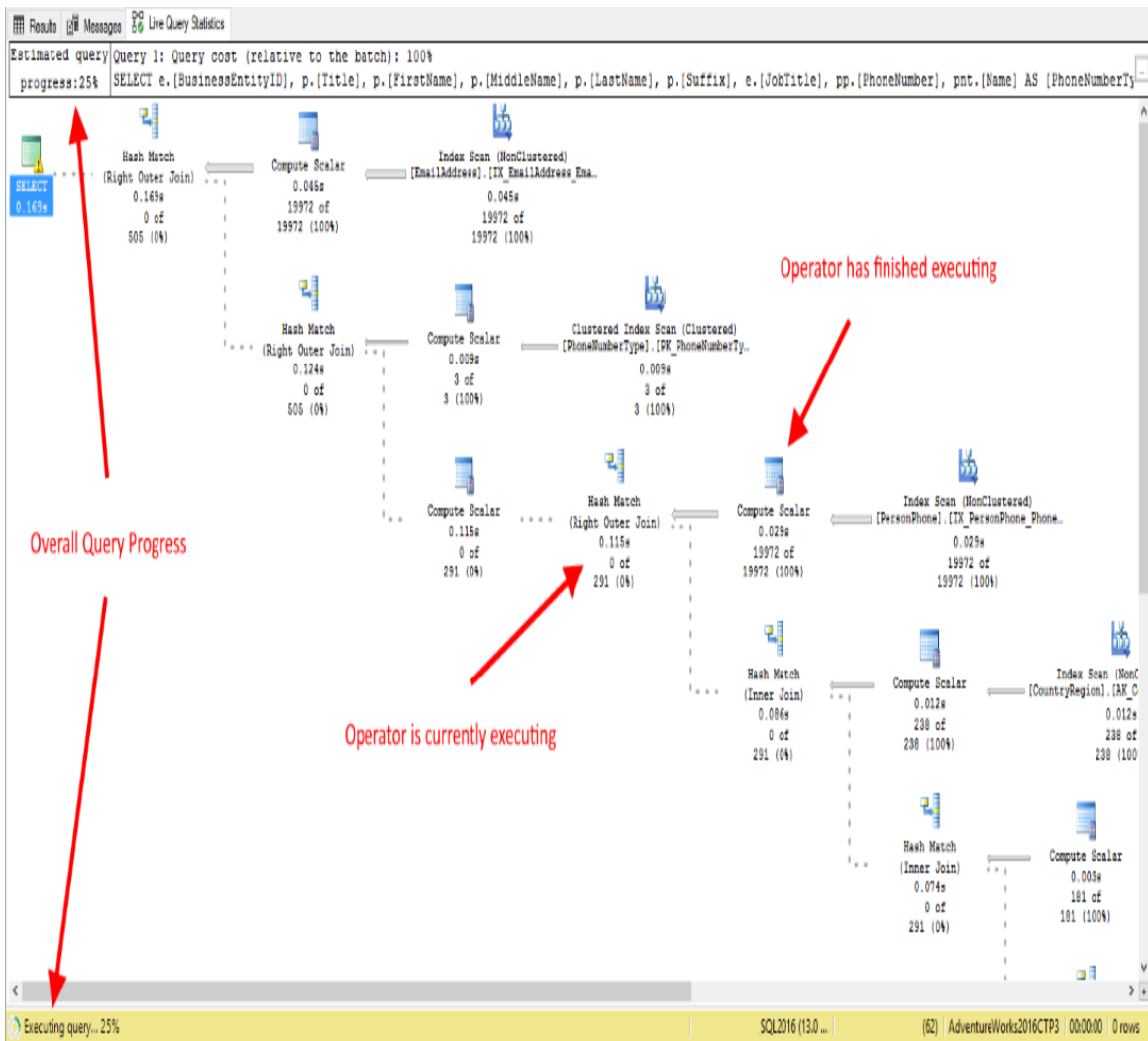


Figure 5-14. Live Query Statistics (Source: *Microsoft documentation*)

Live Query Statistics is very useful when you need to debug long-running queries. It allows you to pinpoint inefficiencies in the execution plans and speed up further query tuning. You can enable Live Query Statistics for queries you run in SSMS. You can also access it from the Active Expensive Query section of the Activity Monitor window.

Live Query Statistics collects data based on *query execution statistics profiling*. It uses two different methods. *Standard profiling* exists in all versions of SQL Server and has historically been used to obtain the actual execution plan for the queries. Unfortunately, this method introduces significant overhead.

Starting with SQL Server 2014 SP2, there is another option called *lightweight profiling*. With this method, the overhead is significantly smaller; however, it does not collect runtime CPU information.

Table 5-2 illustrates how you can enable profiling in different versions of SQL Server. It also shows xEvents that enable profiling globally in the system. Live Query Statistics integrates with the latest version of profiling supported by the SQL Server instance where it runs.

*T
a
b
l
e*

*5
-
2*

*·
T
a
b
l
e*

*5
-
2*

*·
Q
u
e
r
y*

*e
x
e
c
u
t
i
o*

*n
s
t
a
t
i
s
t
i
c
s

p
r
o
f
i
l
i
n
g*

	Type	How to enable	xEvent
Prior SQL Server 2014 SP2	Standard	SET STATISTICS XML SET STATISTICS PROFILE	post_query_execution_showplan
SQL Server 2014 SP2 – SQL Server 2016 RTM		Lightweight v1	Live Query Statistics
			post_query_execution_showplan (less

overhead)
query_thread_profile

SQL Server 2016 SP1 – SQL Server 2017	Lightweight v2	T7412 QUERY_PLAN_PROFILE query hint	query_plan_profile
---------------------------------------	----------------	--	--------------------

SQL Server 2019	Lightweight v3	Enabled by default LIGHTWEIGHT_QUERY_PROFILING database option	query_post_execution_plan_profile
-----------------	----------------	---	-----------------------------------

The overhead of standard profiling is significant. With lightweight profiling, on the other hand, it is very low. According to Microsoft, starting with SQL Server 2016 SP1, the overhead of continuously running lightweight profiling is about 2 to 4%. Technically speaking, you can run an xEvents session and collect profiling information for all queries in the system if the server is not CPU bound. Nevertheless, be careful, and measure the impact of this monitoring in your system.

There is another useful new function, `sys.dm_exec_query_statistics_xml`, that utilizes lightweight profiling. It provides an in-flight execution plan for the currently running request. The result looks like the snapshot of Live Query Statistics. You can use this function together with the `sys.dm_exec_requests` view, as shown in Listing 5-5.

Example 5-5. Using `sys.dm_exec_query_statistics_xml`

```
SELECT
    er.session_id
    ,er.request_id
    ,DB_NAME(er.database_id) as [database]
    ,er.start_time
    ,CONVERT(DECIMAL(21,3),er.total_elapsed_time / 1000.) AS
```

```

[duration]
    ,er.cpu_time
    ,SUBSTRING(
        qt.text,
        (er.statement_start_offset / 2) + 1,
        ((CASE er.statement_end_offset
            WHEN -1 THEN DATALENGTH(qt.text)
            ELSE er.statement_end_offset
            END - er.statement_start_offset) / 2) + 1
    ) AS [statement]
    ,er.status
    ,er.wait_type
    ,er.wait_time
    ,er.wait_resource
    ,er.blocking_session_id
    ,er.last_wait_type
    ,er.reads
    ,er.logical_reads
    ,er.writes
    ,er.granted_query_memory
    ,er.dop
    ,er.row_count
    ,er.percent_complete
    ,es.login_time
    ,es.original_login_name
    ,es.host_name
    ,es.program_name
    ,c.client_net_address
    ,ib.event_info AS [buffer]
    ,qt.text AS [sql]
    ,p.query_plan
FROM
    sys.dm_exec_requests er WITH (NOLOCK)
        OUTER APPLY sys.dm_exec_input_buffer(er.session_id,
er.request_id) ib
        OUTER APPLY sys.dm_exec_sql_text(er.sql_handle) qt
        OUTER APPLY
sys.dm_exec_query_statistics_xml(er.session_id) p
        LEFT JOIN sys.dm_exec_connections c WITH (NOLOCK)
ON
        er.session_id = c.session_id
        LEFT JOIN sys.dm_exec_sessions es WITH (NOLOCK) ON
        er.session_id = es.session_id
WHERE
    er.status <> 'background'
    AND er.session_id > 50
ORDER BY

```



```
er.cpu_time desc  
OPTION (RECOMPILE, MAXDOP 1);
```

Common Issues and Inefficiencies

The query optimization and tuning process may touch multiple layers in the system. For example, with third-party applications, you might not have access to the source code and must deal with predefined set of queries. This, perhaps, is the most challenging case—where your optimization is limited to creating and modifying indexes.

The situation is much better when you are able to change queries and database code. Those changes can be time-consuming and require testing, but they yield better results and significant performance improvements.

In some cases, you need to go beyond the database code. You might need to change the database schema, the application architecture, and sometimes even the technology to scale the system. While this is an extremely complex process, it may provide you the best results long-term.

I will not go that deep in this section, however. Instead, I will cover several common inefficiencies that can be addressed with indexing and code changes. Just remember that your options are not as limited in real life.

Inefficient Code

With exception of *black-box* optimizations, where you don't have access to the code, you should start the tuning by reviewing and, potentially, refactoring the queries. There are several anti-patterns and issues to detect.

Non-SARGable predicates

SARGable predicates allow SQL Server to utilize the *Index Seek* operator by limiting the range of the key values to process. You need to analyze the query detecting and removing non-SARGable predicates when possible.

One very common case that leads to non-SARGable predicates is functions. SQL Server calls the function for every row it is processing in order to

evaluate the predicate. This applies to both system and scalar user-defined functions (UDFs).

NOTE

SQL Server 2019 may inline some scalar UDFs into the query statement. However, there are many limitations that prevent inlining, so it is better to avoid them when possible.

Table 5-3 shows several examples of how you can refactor some predicates to make them SARGable.

*T
a
b
l
e*

*5
-
3*

*.
E
x
a
m
p
l
e
s*

*o
f
r
e
f
a
c
t
o
r
i
n
g*

*N
o
n
-
S
A
R
G
a
b
l
e*

*p
r
e
d
i
c
a
t
e
s*

*i
n
t
o*

*S
A
R
G
a
b
l*

e
o
n
e
s

Operation	Non-SARGable implementation	SARGable implementation
-----------	-----------------------------	-------------------------

Mathematical calculations		
---------------------------	--	--

	Column - 1 = @Value	
--	---------------------	--

		Column = @Value + 1
--	--	---------------------

	ABS(Column) = 1	
--	-----------------	--

		Column IN (-1, 1)
--	--	-------------------

Date manipulation		
-------------------	--	--

	CONVERT(DATETIME, CONVERT(VARCHAR(10),Column,121)) = @Date	
--	---	--

		Column >= @Date AND Column < DATEADD(DAY,1,@ Date)
--	--	--

	DATEPART(YEAR,Column) = @Year	
--	----------------------------------	--

		Column >= @Year AND Column < DATEADD(YEAR,1,@Year)
--	--	--

DATEADD(DAY,7,Column) >
GETDATE()

Column >
DATEADD(DAY,-7,GETDATE())

Prefix search

LEFT(Column,3) = 'ABC'

Column LIKE 'ABC%'

Substring search

Column LIKE
'%ABC%'

Use Full-Text Search or other
technologies

Pay attention to the data types used in the predicates. Implicit conversion operation is, in a nutshell, a call of the system `CONVERT_IMPLICIT` function, which in many cases would prevent index seek. Remember to analyze JOIN predicates in addition to WHERE clauses. Data type mismatches in join columns are a very common source of problems.

User-Defined Functions

I noted just now that system and non-inlined scalar UDFs may prevent SQL Server from using the index seek operation. Moreover, they introduce

significant performance overhead, especially in case of user-defined functions. SQL Server calls them for every row it processes, and those calls are similar to stored procedure calls (you can confirm this by capturing an `rpc_starting` xEvent or `SP:Starting` trace event).

SQL Server optimizes the code in multi-statement UDFs (scalar and table-valued) separately from the caller query. This usually leads to less efficient execution plans. More importantly, depending on the version of SQL Server, database compatibility level and configuration settings, SQL Server estimates that multi-statement table-valued functions return either 1 or 100 rows. This could completely invalidate cardinality estimations and lead to highly inefficient plans.

The latter situation is improved in SQL Server 2017 with *Adaptive Query Processing*. One of its features, Interleaved Execution, defers final compilation of the query until run-time, when SQL Server can measure the actual number of rows returned by the function and finish optimization using that data.

Nevertheless, it is better to avoid multi-statement functions and use inline table-valued functions when possible. SQL Server embeds and optimizes them together with the caller queries. Fortunately, in many cases, scalar and multi-statement table-valued functions can be converted to inline table-valued functions with very little effort.

Temporary Tables and Table Variables

Temporary tables and table variables are very valuable during query optimization. You can use them to persist the intermediate results of queries. This allows you to simplify queries, which may improve cardinality estimations and generate more efficient execution plans.

There are two common mistakes associated with temporary tables and table variables.

First, some people choose table variables over temporary tables because of the common misconception that table variables are in-memory objects that don't use tempdb and are therefore more efficient than temporary tables.

This is not the case. Both objects rely on tempdb. Even though table variables are slightly more efficient than temporary tables, this efficiency comes from an important limitation: they do not maintain statistics on primary keys and unique constraints.

Similar to multi-statement table-valued functions, Query Optimizer estimates that a table variable stores either 1 or 100 rows. In SQL Server 2019, you can use Adaptive Query Processing to defer compilation of the query and allow Query Optimizer to use the actual number of rows to generate the execution plan. You can achieve the same results in older versions of SQL Server with a statement-level recompile and `OPTION (RECOMPILE)` clause. However, SQL Server does not maintain statistics histograms and/or information about actual data distribution in table variables.

Temporary tables, on the other hand, behave like regular tables. They maintain index statistics and allow SQL Server to use them during optimization. While there are some edge cases when table variables could be the better choice, in most cases it is safer to use temporary tables. In many cases throughout my career, I've achieved great results by replacing table variables with temporary tables, without any further code or indexing changes.

The second common mistake is not indexing temporary tables, which negatively affects cardinality estimations and can lead to inefficient table scans. Treat temporary tables as regular tables and index them to support efficient querying, especially when they store significant amounts of data.

You can use a properly indexed temporary table to persist results from multi-statement table-valued functions. This will improve cardinality estimations, especially in the old versions of SQL Server without Adaptive Query Processing.

Obviously, temporary tables and table variables come with a price. There is overhead involved in creating and populating them. When they are used wisely, the benefits may outweigh the downsides, but they are generally not a good choice to store millions of rows.

Stored Procedures and ORM Frameworks

While this topic is not directly related to query tuning, it is impossible to avoid mentioning Object Relational Mapping (ORM) frameworks. They are extremely common nowadays and saying that all database engineers hate them would not be exaggerating. The queries generated by ORM frameworks are extremely complex and hard to optimize.

Unfortunately, we need to accept that those frameworks simplify development and reduce its time and cost. In most cases, it is unrealistic and unreasonable to insist that application developers not use them. More importantly, in many cases, the less efficient queries generated by frameworks are totally acceptable.

Performance-critical queries are different, though. You may not have many of them, but there are always some that will require extensive tuning and optimization. In those cases, autogenerated and/or ad-hoc queries are not the best choice. It is preferable to switch to stored procedures, which provide full flexibility and a larger set of techniques for optimization.

While this switch may require changes in the application code, in many cases it will reduce tuning time and cost. Remember that when you are choosing your tuning approach.

Inefficient Index Seek

As you already know, an index seek operation is usually more efficient than an index scan. This does not mean, however, that every index seek is efficient. SQL Server uses an index seek when query predicates allow it to isolate the range of data rows from the index during query execution. If this range is very large, this can reduce the efficiency of the operation.

Let's look at a simple example: we'll create a table and populate it with some data. Then we'll run two SELECT statements – with and without a WHERE clause — as shown in Listing 5-6.

Example 5-6. Index seek inefficiency

```

CREATE TABLE dbo.T1
(
    IndexedCol INT NOT NULL,
    NonIndexedCol INT NOT NULL
);
CREATE UNIQUE CLUSTERED INDEX IDX_T1
ON dbo.T1(IndexedCol);

;WITH N1(C) AS (SELECT 0 UNION ALL SELECT 0) -- 2 ROWS
,N2(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N1 AS T2) -- 4 ROWS
,N3(C) AS (SELECT 0 FROM N2 AS T1 CROSS JOIN N2 AS T2) -- 16 ROWS
,N4(C) AS (SELECT 0 FROM N3 AS T1 CROSS JOIN N3 AS T2) -- 256 ROWS
,N5(C) AS (SELECT 0 FROM N4 AS T1 CROSS JOIN N4 AS T2) -- 65,536
ROWS
,N6(C) AS (SELECT 0 FROM N3 AS T1 CROSS JOIN N5 AS T2) -- 1,048,576
ROWS
,IDs(ID) AS (SELECT ROW_NUMBER() OVER (ORDER BY (SELECT NULL)) FROM
N6)
INSERT INTO dbo.T1(IndexedCol, NonIndexedCol)
    SELECT ID, ID FROM IDs;

SET STATISTICS IO ON
SELECT COUNT(*) FROM dbo.T1;
SELECT COUNT(*) FROM dbo.T1 WHERE IndexedCol > 0;

```

Figure 5-15 shows the execution plans of both queries along with their I/O statistics. All rows in the table have positive IndexedCol values, so both queries must scan an entire index. In short, an index seek operation is identical to an index scan.

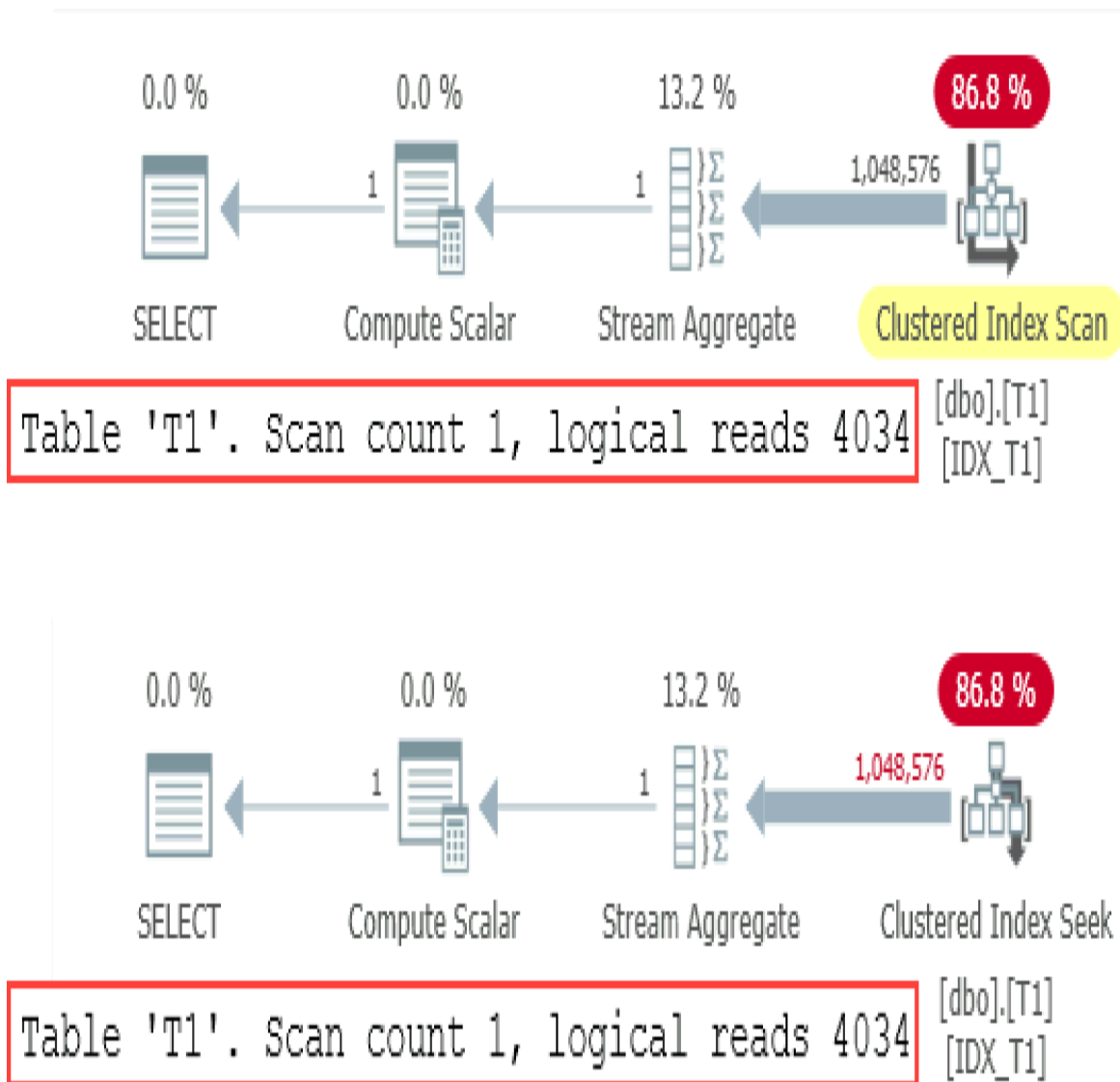


Figure 5-15. Inefficient index seek

It is very common to see inefficient index seeks in multi-tenant systems. Take, for example, order fulfillment software, where data is generally spread across a relatively small number of warehouses. It is common to see the tenant-id (or warehouse_id, in this example) as the left-most column in the index keys.

Queries in those systems usually process data from a single tenant, using it as the predicate in the WHERE clause, which will rightfully lead to index seek operations in the execution plan. But if each tenant (or warehouse) stores very large amounts of data, you may get perfectly looking and

inefficient execution plan even without any scans present. You'll need to use other predicates to make index seeks selective and improve performance.

You can analyze the efficiency of index seeks in the execution plan by looking at the operator's properties. Let's run the query shown in Listing 5-7 against the table we created in Listing 5-6.

Example 5-7. Sample query

```
SELECT IndexedCol, NonIndexedCol
FROM dbo.T1
WHERE
    IndexedCol BETWEEN 100 AND 150 AND NonIndexedCol % 2 = 0;
```

Figure 5-16 illustrates several key properties for the analysis. The screenshot here was captured in SentryOne Plan Explorer; however, you'd see the same data in SSMS.

Clustered Index Seek
Scans a particular range of rows from a clustered index.
Actual Rows: 26
Actual Rows Read: 51
Estimated Rows: 2,949
Estimated Rows To Be Read: 94,372
Table: [dbo].[T1]
Clustered Index: [IDX_T1]
Seek Predicates:
[AAD].[dbo].[T1].[IndexedCol] >= [@0]
[AAD].[dbo].[T1].[IndexedCol] <= [@1]
Predicate:
[AAD].[dbo].[T1].[NonIndexedCol]%(2)=[@3]
Output List:
IndexedCol
NonIndexedCol

Figure 5-16. Index Seek operator properties

Let's look at the most important properties of the operator.

Seek Predicate

This property shows the predicate(s) that SQL Server uses to limit the range of rows during the index seek. The more selective this predicate is, the more efficient index seek will be.

Predicate

The predicate illustrates additional filter criteria that SQL Server applies to every row read by an *Index Seek* operator. It does not reduce the size of the data the operator must process; however, it may reduce the number of rows the operator returns in the execution plan. In our case, the operator read 51 rows from the index and returned 26 rows to the next operator in the execution plan.

It is always more efficient to reduce the size of the data with an efficient seek predicate. When seek predicates are not selective enough, consider restructuring your index in a way that allows SQL Server to use regular predicates as seek predicates.

Actual Rows and Actual Rows Read

These two properties are called *Actual Number of Rows* and *Number of Rows Read* in SSMS. They illustrate how many rows were returned by the operator and how many rows were processed during execution. The large value in *Number of Rows Read* indicates that index seek processed a large amount of data and may require further investigation. The large discrepancy between those two values shows potential index inefficiency, with a significant portion of the data being filtered out by the *Predicate* rather than the *Seek Predicate* of the operator.

Estimated Rows and Estimated Rows to be Read

Those two properties are called *Estimated Number of Rows* and *Estimated Number of Rows to be Read* in SSMS. As noted, you can compare *estimated* and *actual* metrics in the execution plan to estimate

the quality of the cardinality estimation data. A noticeable cardinality estimation error may indicate a wrong choice of index and/or join type (more on that later), which you will likely want to address.

Obviously, it is much easier to analyze an execution plan when you have the actual execution metrics available. While estimated metrics can be useful for initial analysis, cardinality estimation errors may provide a very wrong or incomplete picture. Do additional analysis and look at the data distribution in the database when dealing with estimated execution plans.

Incorrect Join Type

SQL Server uses many physical join operators during query execution. These belong to one of the three logical join types: *loop*, *hash* and *merge*. Each is optimized for specific conditions, and the incorrect choice may have a serious negative impact on the query performance. Unfortunately, people often don't pay attention to the join type chosen by SQL Server, overlooking opportunities for optimization.

Let's look at all three types in more detail.

Loop Join

A loop join (or nested loop join) is the simplest join algorithm. As with any join type, it accepts two inputs, which are called *outer* and *inner* tables. The algorithm for the join is very simple (Listing 5-8). Briefly, SQL Server goes through the outer table looking up rows to join in the inner table for each outer row.

Example 5-8. Loop join algorithm (pseudo code)

```
/* Inner join */
for each row R1 in outer table
    find row(s) R2 in inner table
        if R1 joins with R2
            return join (R1, R2)
/* Outer join */
for each row R1 in outer table
```

```
find row(s) R2 in inner table
  if R1 joins with R2
    return join (R1, R2)
  else
    return join (R1, NULL)
```

The cost of the join depends on two factors. The first is the size of the outer table. SQL Server goes through each row there, locating corresponding rows in the inner table to join. The more data it needs to process, the less efficient it will be.

The second factor is the efficiency of the inner table search. When the join column(s) in the inner table are properly indexed, SQL Server can utilize the efficient index seek operation. In that case, the cost of the inner table search on each iteration will be relatively low. Without the index, SQL Server might have to scan the inner table multiple times – once for each row from the outer table. As you can guess, that is extremely inefficient.

The loop join is optimized for conditions in which one of the tables is small and the other has an index to support an index seek operation for the join. It is impossible to define the hard threshold after which the join becomes inefficient. It may perform well with thousands and sometimes tens of thousands of rows in the outer input; however, it would not scale well with millions of rows. Nevertheless, in proper conditions, this type of join is extremely efficient. It has very little startup cost, does not use tempdb, and does not consume large amount of memory.

Finally, loop join is the only join type that does not require an equality predicate. SQL Server may evaluate a join predicate between every row from both inputs. It does not require a join predicate at all. For example, the CROSS JOIN operator would lead to a nested loop physical join when every row from both inputs has been joined together. Obviously, SQL Server cannot use index seek if the join predicate is not SARGable, which would lead to extremely inefficient operation with large inputs.

Merge Join

The merge join works with two sorted inputs. It compares two rows, one at time, and returns their join to the client if they are equal. If they are not, it

discards the lesser value and moves on to the next row in the input. The algorithm for the join is shown in Listing 5-9.

Example 5-9. Inner merge join algorithm (pseudocode)

```
/* Prerequisites: Inputs I1 and I2 are sorted */
get first row R1 from input I1
get first row R2 from input I2
while not end of either input
begin
    if R1 joins with R2
    begin
        return join (R1, R2)
        get next row R2 from I2
    end
    else if R1 < R2
        get next row R1 from I1
    else /* R1 > R2 */
        get next row R2 from I2
end
```

The merge join is optimized for medium and large inputs, when both of those inputs are sorted. This means that inputs need to be indexed on the join predicate columns. However, in practice, SQL Server may decide to sort inputs during query execution; the cost of the sort may thus far exceed the cost of the merge join itself. Check if that is the case and factor the cost of the Sort operator into your analysis.

Hash Join

A *hash join* is designed to handle large unsorted inputs. Its algorithm consists of two different phases.

During the first, or *build*, phase, a hash join scans one of the inputs (usually the smaller one), calculates the hash values of the join key, and places them into the hash table. In the second, or *probe*, phase, it scans the second input, and checks (probes) to see if the hash value of the join key from the second input exists in the hash table. If so, SQL Server evaluates the join predicate for the row from the second input and all rows from the first input that belong to the same hash bucket. The algorithm is shown in Listing 5-10.

Example 5-10. Inner hash join algorithm (pseudocode)

```

/* Build Phase */
for each row R1 in input I1
begin
    calculate hash value on R1 join key
    insert hash value to appropriate bucket in hash table
end
/* Probe Phase */
for each row R2 in input I2
begin
    calculate hash value on R2 join key
    for each row R1 in hash table bucket
    if R1 joins with R2
        return join (R1, R2)
end

```

A hash join requires memory to store the hash table. When there is not enough memory, the join stores some hash table buckets in tempdb. This condition is called a *spill* and can greatly affect join performance, since tempdb access is significantly slower.

Spills often happen due to incorrect memory grant estimation, which in turn may be triggered by incorrect cardinality estimation. When this is the case, make sure that the statistics are up to date; consider simplifying or refactoring the query if this does not help.

Adaptive Query Processing in SQL Server 2017 introduced a new feature called *memory grant feedback*, which increases or decreases memory grants for a query based on memory usage in previous executions. In SQL Server 2017, that feature is limited to batch mode execution. Starting with SQL Server 2019, it is also enabled in row mode execution. Read the [Microsoft documentation](#) for more information, and consider enabling it. This may reduce tempdb spills in the system.

Comparing Join Types

Table 5-4 summarizes the behavior of different join types and the use cases for which they are optimized.

*T
a
b
l
e

5
-
4
.
J
o
i
n

c
o
m
p
a
r
i
s
o
n*

Nested Loop Join Merge Join Hash Join

Best use case	At least input is small; index on the join column(s) in another input	Medium to large inputs, sorted on index key	Medium to large inputs
---------------	---	---	------------------------

Requires sorted input	No	Yes	No
-----------------------	----	-----	----

Requires equality predicate	No	Yes	Yes
-----------------------------	----	-----	-----

Blocking operator	No	No	Yes (Build phase only)
-------------------	----	----	------------------------

Uses memory	No	No	Yes
-------------	----	----	-----

Uses tempdb	No	No (Sort may spill to tempdb)	Yes, in case of spills
-------------	----	-------------------------------	------------------------

Preserves order	Yes (outer input)	Yes	No
-----------------	-------------------	-----	----

Adaptive Query Processing in SQL Server 2017 also introduced the concept of the *adaptive join*. With this join, SQL Server chooses to use either a loop or a hash join based on the size of the inputs in runtime. Unfortunately, in SQL Server 2017 and 2019, this works only in batch mode execution, which, in most cases, is triggered by columnstore indexes. You need to enable Live Query Statistics in SSMS to see adaptive join in the execution plan.

I mentioned just now that each type of join is optimized for specific use cases and may not perform well in other cases. Let's look at a simple example and compare the performance of different join types. Listing 5-11 creates another table (similar to the one in Listing 5-6) and populates it with the same data. Both tables have two columns each and a clustered index defined on one of the columns.

Example 5-11. Join performance: Table creation

```
CREATE TABLE dbo.T2
(
    IndexedCol INT NOT NULL,
    NonIndexedCol INT NOT NULL
);
CREATE UNIQUE CLUSTERED INDEX IDX_T2
ON dbo.T2(IndexedCol);
INSERT INTO dbo.T2(IndexedCol, NonIndexedCol)
    SELECT IndexedCol, NonIndexedCol FROM dbo.T1;
```

Next, let's compare the performance of different join types using the code in Listing 5-12. Here, I am forcing different join types with join hints (more on those later). I put the execution time of the statements in my test environment into code comments.

Example 5-12. Join performance: Test cases

```
-- Loop join with index seek in inner table
-- Elapsed time: 137ms.
SELECT COUNT(*)
FROM dbo.T1 INNER LOOP JOIN dbo.T2 ON    T1.IndexedCol =
T2.IndexedCol
WHERE
    T1.NonIndexedCol <= 100;

-- Loop join with inefficient index scan in inner table
-- Elapsed time: 16,732ms
SELECT COUNT(*)
FROM dbo.T1 INNER LOOP JOIN dbo.T2 ON
    T1.IndexedCol = T2.NonIndexedCol
WHERE
    T1.NonIndexedCol <= 100;

-- Hash join. Slower than loop join on small inputs
-- Elapsed time: 411ms.
SELECT COUNT(*)
FROM dbo.T1 INNER HASH JOIN dbo.T2 ON
    T1.IndexedCol = T2.IndexedCol
WHERE
    T1.NonIndexedCol <= 100;

-- Loop join with index seek in inner table with large input
-- Elapsed time: 1,514ms
SELECT COUNT(*)
FROM dbo.T1 INNER LOOP JOIN dbo.T2 ON
```

```

T1.IndexedCol = T2.IndexedCol;

-- Hash join using indexed columns to join
-- Faster than loop join on large input
-- Elapsed time: 1,215ms
SELECT COUNT(*)
FROM dbo.T1 INNER HASH JOIN dbo.T2 ON
    T1.IndexedCol = T2.IndexedCol;

-- Hash join using non-indexed columns to join
-- Performance does not depend on if join columns were indexed
-- Elapsed time: 1,235ms
SELECT COUNT(*)
FROM dbo.T1 INNER HASH JOIN dbo.T2 ON
    T1.IndexedCol = T2.NonIndexedCol;

-- Merge join with pre-sorted inputs
-- Elapsed time: 440ms
SELECT COUNT(*)
FROM dbo.T1 INNER MERGE JOIN dbo.T2 ON
    T1.IndexedCol = T2.IndexedCol;

-- Merge join without pre-sorted inputs
-- Elapsed time: 774ms
SELECT COUNT(*)
FROM dbo.T1 INNER MERGE JOIN dbo.T2 ON
    T1.IndexedCol = T2.NonIndexedCol;

```

Loop join is faster than hash join on the small inputs; however, hash join becomes more efficient as size of the input grows. Merge join, on the other hand, is great when inputs are sorted. Otherwise, it adds a *Sort* operator to the execution plan. While this might work fine with small inputs, sorting very large inputs would not work well.

As these examples illustrate, an incorrect choice of join type can reduce query performance dramatically. In most cases, this happens due to incorrect cardinality estimations, especially when SQL Server seriously underestimates the size of join inputs.

With hash and merge joins, this may lead to tempdb spills and slower join performance. However, that situation is most dangerous with the loop join, especially when the cost of inner input processing is high. (Recall that SQL

Server processes the inner input for each row from the outer input, so costs add up quickly with each iteration.

You can detect that condition by comparing actual to estimated rows in the outer input or by looking at actual versus estimated number of executions in the first operator from inner input (see Figure 5-17). A large discrepancy would indicate a cardinality estimation error. When a cardinality estimation error leads to a high number of executions in the inner input, a loop join may be the wrong choice.

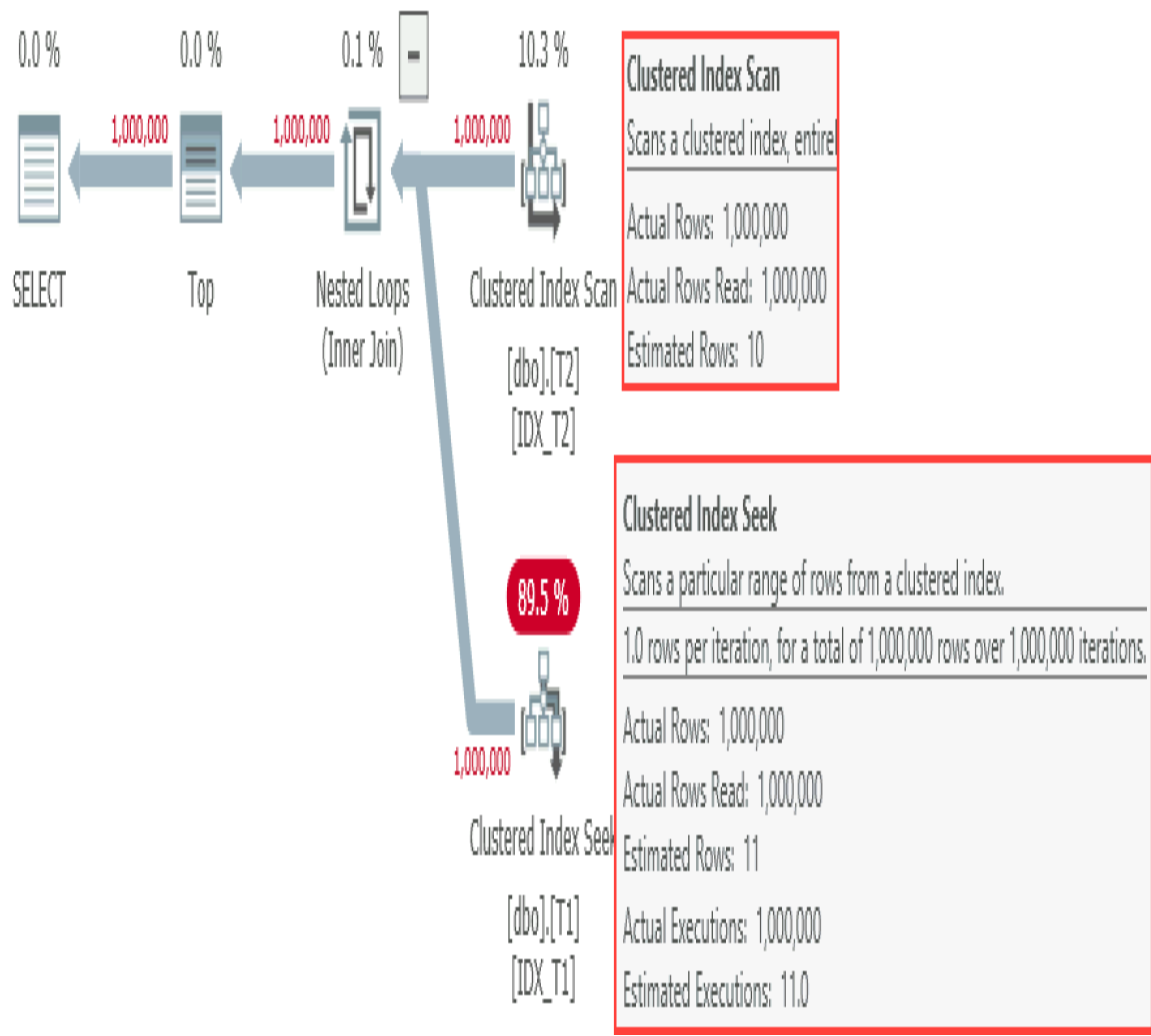


Figure 5-17. Cardinality estimation error in loop join

You have several options for addressing this problem. As a first step, review the query, looking for opportunities to refactor. Remove code patterns that could affect cardinality estimations, like multi-statement functions and table

variables. See if there is a possibility of better indexing, which could improve the execution plan.

It's worth checking if you are being affected by parameter sniffing in parameter-sensitive plans. SQL Server sometimes caches and reuses execution plans compiled for atypical parameter values, especially when data is distributed very unevenly in the table. Think about multi-tenant systems where some tenants may have very little data and others a great deal. Execution plans generated for the former group of tenants would be inefficient for the latter.

When this is the case, consider using statement-level recompilation with `OPTION (RECOMPILE)` or disabling parameter sniffing in the database. I will discuss those and other options in more detail in the next chapter.

You can also try updating the statistics in the tables from the query. Unfortunately, this may not always help. Query Optimizer in SQL Server does not always make the right assumptions when estimating cardinality in complex queries with multiple joins.

I have mentioned query simplification, which may help to address that problem. Consider splitting the query and persist intermediate results in the properly indexed temporary tables. SQL Server will be able to see the data distribution there and accurately estimate cardinality in the queries. Obviously, remember the overhead that temporary tables may introduce and do not use them to cache very large datasets.

As a last resort, you can force join types with query hints. This is a dangerous method, and you need to be very careful. The hints will force Query Optimizer to perform optimization in a specific way, which may be inefficient in the future if the data distribution changes. If you do this, remember it and re-evaluate periodically if hints are still required.

There are two ways to use join hints. The first is specifying a list of allowed join types in query options. The first statement in Listing 5-13 shows how to use this method to prevent Query Optimizer from using loop joins anywhere in the query. The second option is specifying the type of specific

join between the tables. The second statement in the listing forces SQL Server to use a hash join between tables A and B.

Example 5-13. Forcing join types

```
SELECT A.Col1, B.Col2
FROM
    A JOIN B ON A.ID = B.ID
    JOIN C ON B.CID = C.ID
OPTION (MERGE JOIN, HASH JOIN);
SELECT A.Col1, B.Col2
FROM
    A INNER HASH JOIN B ON A.ID = B.ID
    JOIN C ON B.CID = C.ID;
```

Unfortunately, the second approach also forces join orders for *all* joins in the query. SQL Server will always join the tables in the order they were specified in the query without trying to reorder the joins. In Listing 5-13, for example, table A will be always be joined with table B first using hash join, and the result of their join will be joined with table C. This makes it dangerous for queries with multiple joins by preventing possible join-order optimizations. Be careful when you use them!

Excessive Key Lookups

As you already know, key and RID lookups¹ become extremely inefficient on a large scale. SQL Server does not use indexes for seek operations when it estimates that large number of key or RID lookups will be required. However, you may still encounter excessive lookups in the execution plans.

This situation often occurs due to incorrect cardinality estimations. If SQL Server estimates that just a handful of key lookups will be required, it might decide to use nonclustered index seek, and the error could lead to a large number of key lookups. That error usually happens due to cardinality estimation model limitations or parameter sniffing (in parameter-sensitive plans, about which you'll learn more in the next chapter).

You can detect this issue by analyzing estimated and actual numbers of rows in the execution plan, as shown in Figure 5-18.

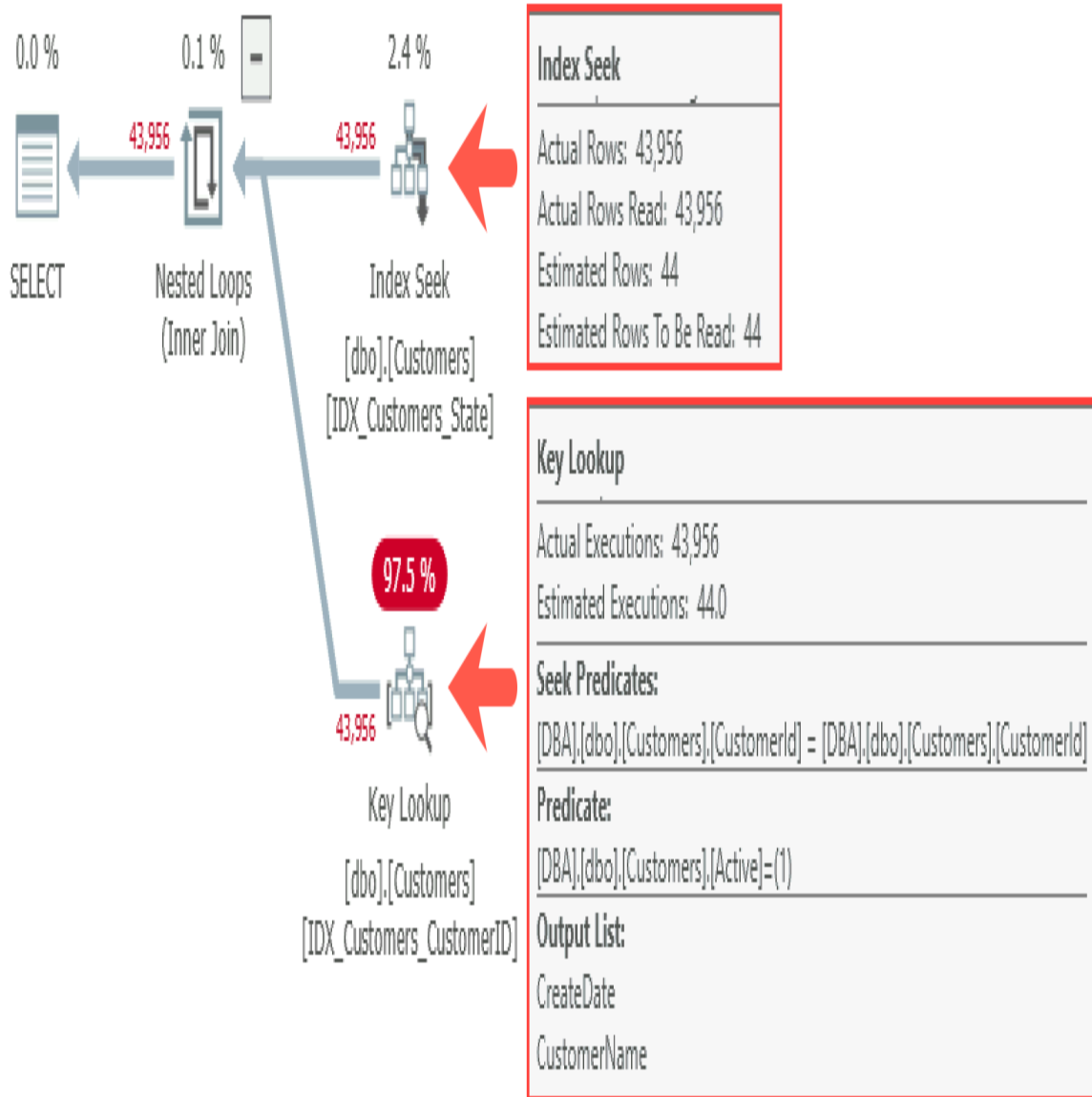


Figure 5-18. Incorrect cardinality estimation and key lookups

SQL Server may also decide to use key lookups, choosing between bad data access strategies and worse ones. Running millions of key lookups is extremely inefficient; however, it could be better than scanning a table with billions of rows.

In many cases, you can remove key lookups with covering indexes. If all columns required for the query are present in a nonclustered index, SQL Server doesn't need to access the main data row in a clustered index or heap. By definition, it includes all columns in the nonclustered index key along with the clustered index columns in row-id.

SQL Server allows you to include other columns in the index using the INCLUDE index clause. Data from these columns is stored on the leaf level only. It is not added to the index key and does not affect the sorting order of the index rows. Figure 5-19 illustrates the structure of an index with included columns, defined as `CREATE INDEX IDX_Customers_Name ON dbo.Customers(Name) INCLUDE(DateOfBirth)` on the table, which has `CustomerId` as the clustered index column.

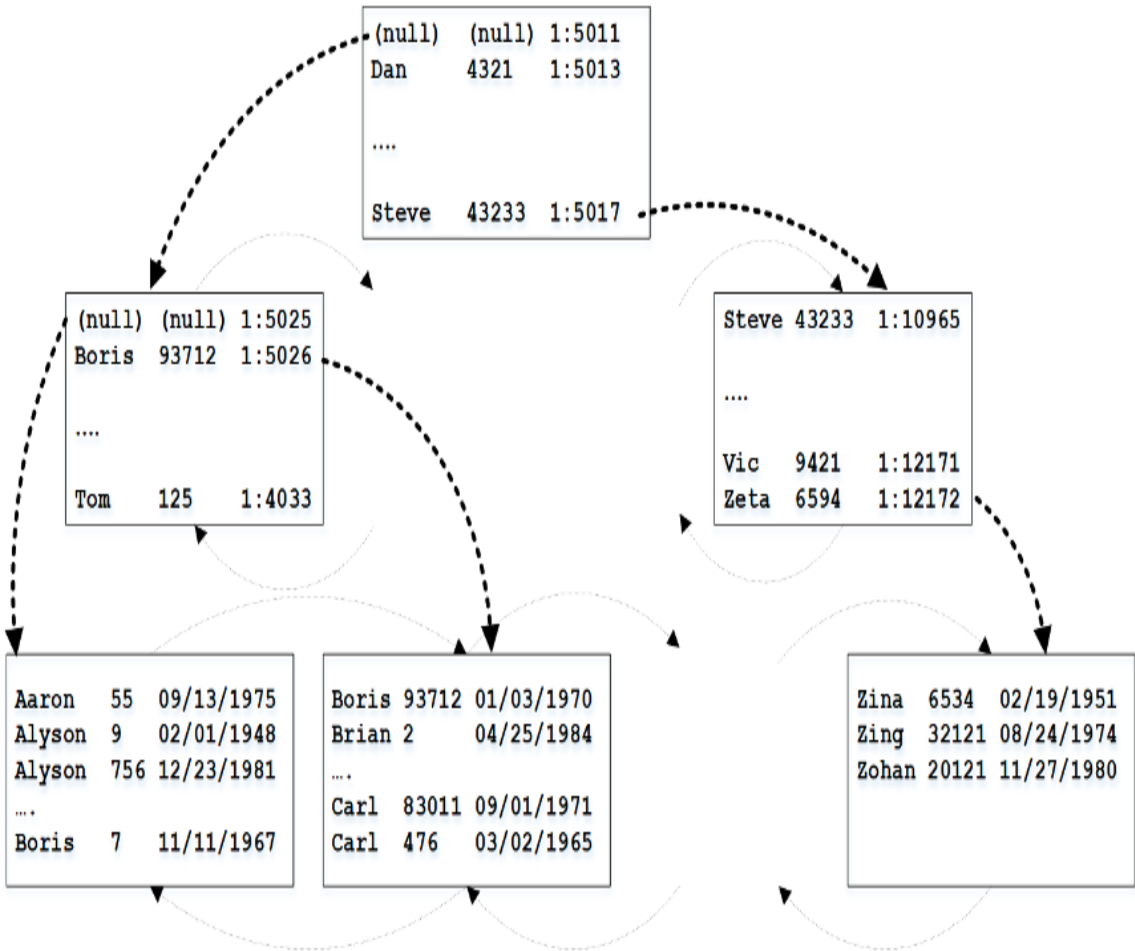


Figure 5-19. Index with included columns

Now, if the only columns the query references are present in the index, SQL Server can obtain all data from the leaf level of the nonclustered index B-Tree without performing key lookups. It can use the index regardless of how many rows would be selected from there.

Making nonclustered indexes covering is one of the most commonly used query optimization techniques. You can open the properties of Key Lookup operator in the execution plan (see Figure 5-18 above) and get a list of columns from the operator's output and filter predicate. Including those columns in the nonclustered index would eliminate the need to do lookups.

Although covering indexes are a great tool for optimizing queries, they come at a cost. Every column in the index increases its row size, as well as the number of data pages it uses on disk and in memory. That introduces additional overhead during index maintenance and increases the database size. Moreover, queries need to read more pages when scanning all or part of the index. This doesn't necessarily introduce a noticeable performance impact during small range scans, when reading a few extra pages is far more efficient than key lookups, but it can degrade the performance of queries that scan a large number of data pages or the entire index.

Covering indexes also add update overhead. By adding a column to nonclustered indexes, you store the data in multiple places. This improves the performance of queries that select the data. However, during updates, SQL Server needs to change the rows in every index where updated columns are present.

It is not a good idea to create very wide nonclustered indexes that include majority of the columns in the table. Nor do you want to have very large number of indexes. Both conditions will increase the size of the database and lead to excessive update overhead, especially in OLTP environments. (I will talk more about index analysis and consolidation in chapter 14.)

Finally, I would like to state a very important and obvious thing. While excessive key lookups are bad for performance, having key lookups is *completely normal*. Doing lookups of hundreds or even thousands of rows may be a better option than creating large covering indexes.

Think about it from a different angle: the key lookup is basically a loop join with efficient index seek in the inner table (clustered index). It is very fast and efficient on a small scale.

Remember, you don't have eradicate all key lookups from the execution plan. Just analyze their efficiency and take care of the inefficient ones.

Indexing the Data

Designing proper indexes is both an art and a science. You will master this skill over time. Let me give you a few tips on how to start.

Most queries in the system have some parameters. The most selective SARGable predicates filter out most of the data; the columns in those predicates are the best candidates for the index.

Let's look at a hypothetical query that returns a list of one customer's orders (Listing 5-14).

Example 5-14. Selecting a list of a customer's orders

```
SELECT c.CustomerName, c.CustomerNumber, o.OrderId, o.OrderDate,
o.Amount
FROM
    dbo.Customers c JOIN dbo.Orders o ON
        c.CustomerId = o.CustomerId
WHERE
    c.CustomerNumber = @CustNum AND
    c.Active = 1 AND
    o.OrderDate between @StartDate AND @EndDate AND
    o.Fulfilled = 1;
```

Let's assume that you run this query against tables that have no indexes except clustered indexes on CustomerId and OrderId columns. The execution plan for the query consists of two clustered index scans and a hash join between the tables (Figure 5-20). As you can see, it is inefficient.

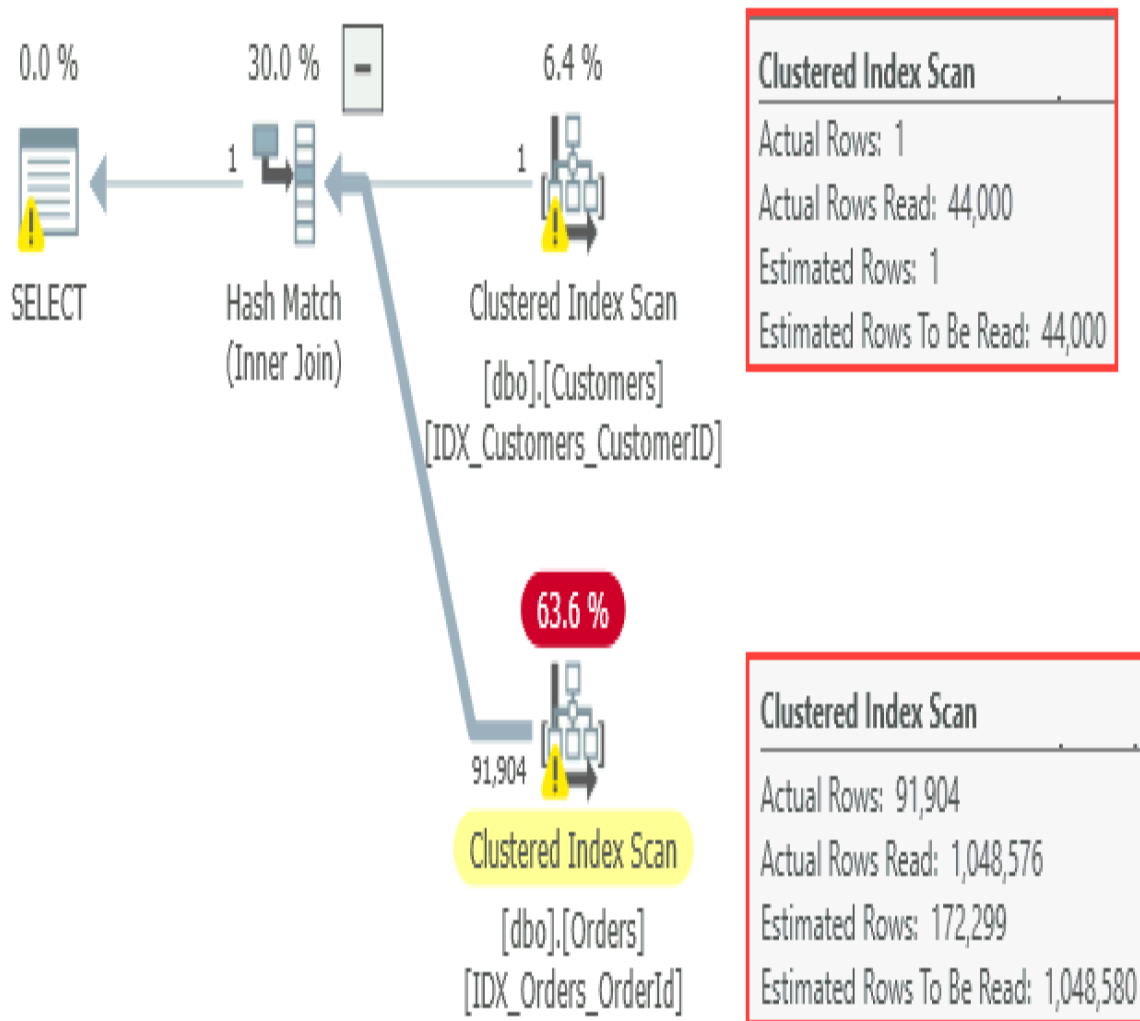


Figure 5-20. Indexing example: Initial execution plan

The natural place to start is the CustomerNumber column. The data here is likely to be unique, and the index on this column will be highly selective. The predicate `Active=1` checks the status of the customer. You can expect it to be commonly used in the queries. It is a good idea to add the Active column to the index as an included column.

I'd also expect CustomerName to often be selected alongside CustomerNumber data and added to the index as another included column, eliminating the need for a key lookup operation. Still, key lookups may be completely acceptable on a small scale with highly selective indexes.

Figure 5-21 shows the execution plan after you create the following index:
`CREATE UNIQUE INDEX IDX_Customers_CustomerNumber ON`

dbo.Customers(CustomerNumber) INCLUDE (Active, CustomerName).

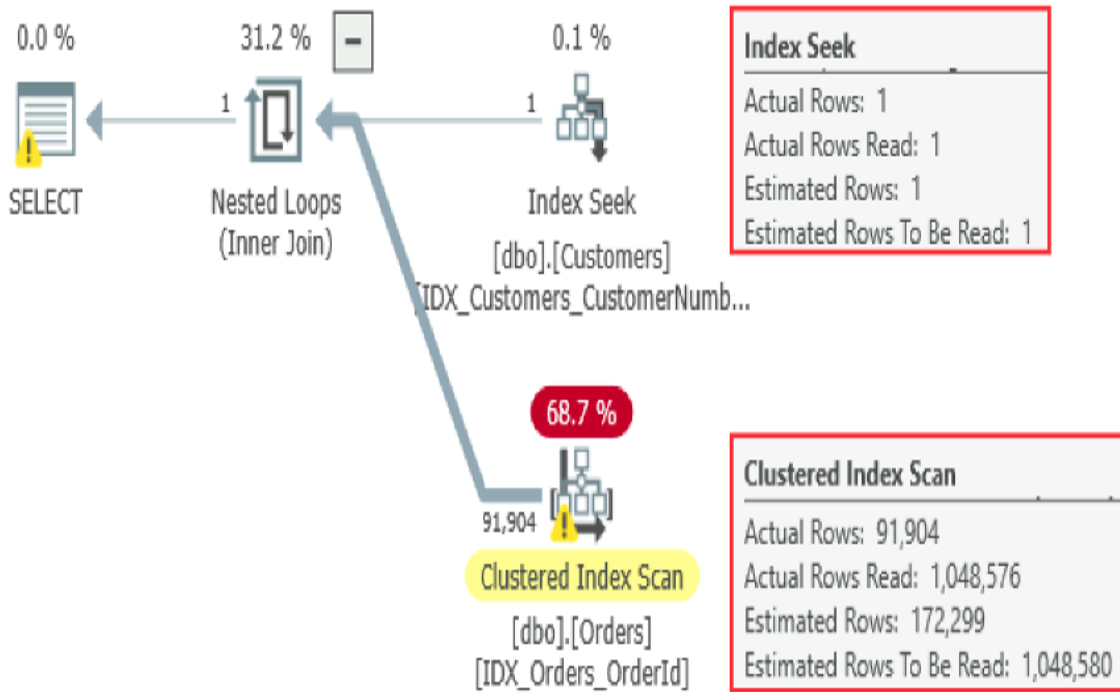


Figure 5-21. Indexing example: Plan with the index on Customers table

There are three predicates with Orders table columns. Two of them—on the CustomerId and OrderDate columns—are selective, and thus good candidates for the index. You can define the index with either (CustomerId, OrderDate) or (OrderDate, CustomerId) column order.

To choose, consider how data will be sorted in the index. With the first option, (CustomerId, OrderDate), SQL Server sorts the data by CustomerId first. Then the orders for each customer are sorted by OrderDate. With the second option, the data will be sorted by OrderDate across all customers.

Both indexes will allow index seek in the Orders table. However, the first index is more efficient for our query. SQL Server will be able to do a range scan for orders that belong to this single customer for the time interval defined by @StartDate and @EndDate. With the second index, SQL Server would have to read all orders in that time interval for all customers, which would force it to scan more data.

It's good to add the Fulfilled column to the index as an included column, to evaluate the predicate as part of the seek operation. In this case, I'd also include the Amount column – it is small enough and would not increase the size of the index much.

Figure 5-22 shows the final execution plan after the following index has been created: `CREATE INDEX IDX_Orders_CustomerId_OrderDate ON dbo.Orders(CustomerId, OrderDate) INCLUDE (Fulfilled, Amount)`.

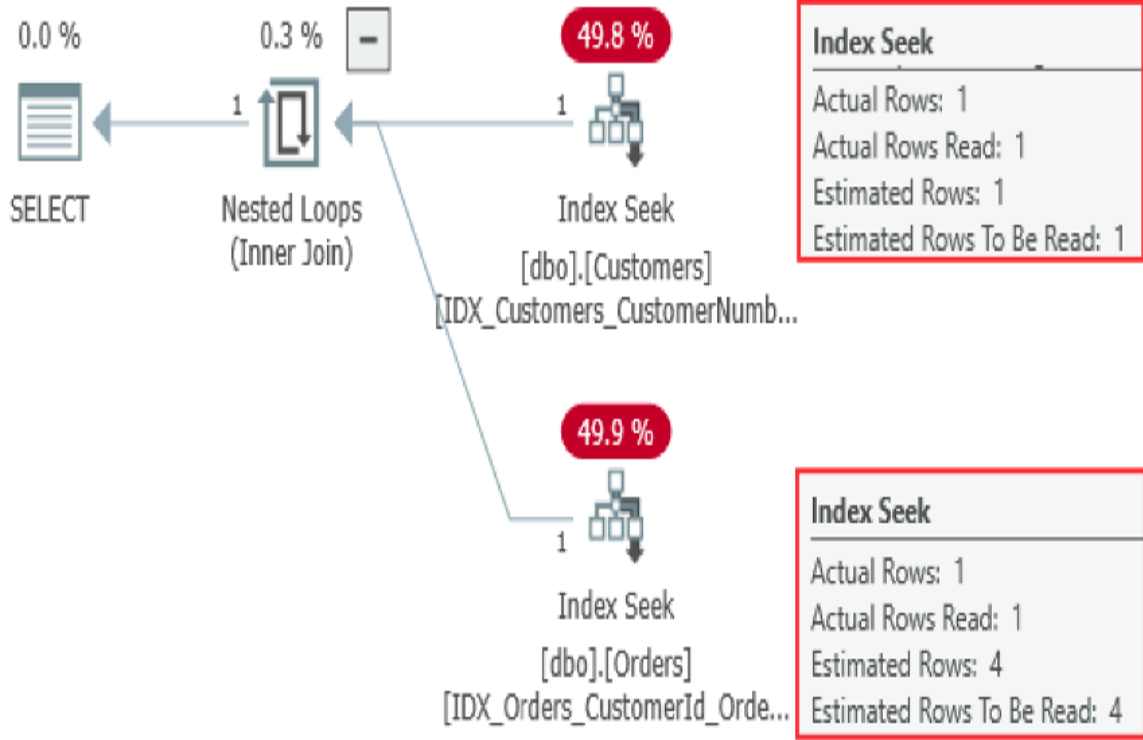


Figure 5-22. Indexing example: Final execution plan

Query optimization is never boring. It constantly challenges you and helps you to learn. I hope, this chapter gives you some tips on where to start and encourage you to practice and learn more. After all, query tuning is the most efficient way to improve performance of the system.

A few more words of advice: Don't create separate indexes for each query. Instead, analyze the workload in the system and create indexes that can be useful in multiple queries. When you start optimization, review the least efficient queries and identify common access patterns.

Be careful with covering indexes. Making them very wide is not a good idea. Again, there is nothing wrong with key lookups if they don't introduce a significant impact on the system.

In the next chapter, we will discuss high CPU load and options how to reduce it.

Summary

SQL Server supports three data storage and processing technologies. Row-based storage, which is the most common, stores the data from all columns together in unsorted heaps or sorted B-Tree indexes.

B-Tree indexes sort data based on index keys. Clustered indexes store the data from all table columns. Nonclustered indexes store another copy of the data in separate physical indexes. They reference clustered indexes through row-id, which is the clustered index key values. When data is not present in the nonclustered index, SQL Server goes through the clustered index B-Tree using key lookup operation. This operation is expensive at scale.

SQL Server accesses data in two ways. An *index scan* usually reads all rows from the index. *Index seek* isolates and processes a subset of the index rows, which is usually more efficient than an index scan. Write your queries to allow SQL Server to utilize index seeks. You should also analyze the efficiency of your index seeks, making sure that they don't process large amounts of data.

SQL Server stores information about indexes and data distribution in statistics, which it uses to estimate how many rows each operator will need to process in the execution plan. Accurate cardinality estimation helps SQL Server generate efficient execution plans. Up-to-date statistics are a key element in correct cardinality estimation.

When you analyze execution plans, pay attention to the efficiency of the Index Seek operators, the choices of join type, and the number of key and RID lookups. Check cardinality estimations, too: improper cardinality estimation is one of the most common problems that lead to poor execution

plans. Make sure that statistics are up to date, add the required indexes, and simplify and refactor queries to address issues.

Troubleshooting Checklist

Troubleshoot the following:

- Analyze statistics maintenance strategy in the system. Add a T2371 trace flag if the system has databases with compatibility level below 130 (SQL Server 2016).
- Analyze and improve index maintenance strategies.
- Make sure that statistics on filtered indexes are frequently updated. Optionally, consider rebuilding them frequently.
- Identify and optimize inefficient queries.

¹ For the sake of brevity, from this point on, I will stop referring to “key and RID lookups”; everything I say about “key lookups” applies to RID lookups as well.

Chapter 6. High CPU Load

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 6 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

My first SQL Server tuning project happened more than 20 years ago, and I’ve been dealing with many systems ever since. Over the years, I’ve noticed an interesting trend. Most of the systems I optimized in the past were I/O bound. Of course there were other problems, but reducing I/O load through query tuning and code refactoring was usually enough to get the job done.

This started to change several years ago. While I still see non-optimized and I/O intensive queries, their impact is masked by high-performance, low-latency flash-based drives. Moreover, the availability of cheap hardware allows for bigger servers that can handle the load from more users. The need to reduce high CPU load is quite common nowadays.

In this chapter, I will talk about several common patterns that increase CPU load and options to address it. I will start with non-optimized queries and inefficient database code. Next, I will cover query compilation overhead, along with plan caching, and the issues they can introduce. Finally, I will

discuss the benefits and downsides of parallelism in systems and ways to tune your parallelism settings.

Non-Optimized Queries and T-SQL Code

Why does your server have a high CPU load? There are several possibilities, but I'll start with the most obvious and common one: non-optimized queries. It does not matter how fast the disk subsystem is. Nor does it matter whether the servers have enough memory to cache all data in the buffer pool and eliminate all physical I/O. Non-optimized queries *will* increase CPU load.

To put things in perspective: a query that scans 10 million data pages uses a million times more CPU resources than a query that scans just 10 pages. It does not matter that each logical read takes just a few microseconds of CPU time: that adds up quickly when multiple users are running those queries in parallel.

You can detect CPU-intensive queries using the techniques I discussed in Chapter 4, such as sorting data by CPU (worker) time while choosing targets for optimization. Optimizing those queries will decrease CPU load.

Don't confuse *CPU time* with *duration*, though. While queries with higher CPU time usually take longer to complete, the opposite is not true. A query may be blocked and use little CPU but still take a long time.

NOTE

Reducing query duration will improve users' experience, but I rarely choose optimization targets based on this factor. Optimizing queries with high resource usage usually reduces duration as well.

Inefficient T-SQL Code

Inefficient T-SQL code also contributes to the problem. Except for natively compiled In-Memory OLTP modules, SQL Server interprets T-SQL code. This leads to additional CPU overhead.

Don't get me wrong: I don't want to discourage you from using stored procedures and T-SQL code. The benefits of properly designed and implemented T-SQL modules greatly outweigh CPU overhead. But there's one case I need to mention specifically – row-by-row processing.

Regardless of how you implement row-by-row processing – with cursors or with loops – it is inefficient. Imperative row-by-row execution will be slower and more CPU-intensive than declarative set-based logic. There are some rare cases when you absolutely have to implement row-by-row processing; however, avoid it when possible.

Statements that perform row-by-row processing may not always appear to be the most resource-intensive statements. You can look at plan cache-based execution statistics for T-SQL modules with `sys.dm_exec_procedure_stats`, `sys.dm_exec_function_stats`, and `sys.dm_exec_trigger_stats` views (discussed in Chapter 4) to detect the modules with the most cumulative resource usage. Analyze what they are doing, keeping an eye on row-by-row logic.

Other T-SQL constructs contribute to CPU load, too: for example, JSON and (especially) XML support are CPU-intensive. It is better to parse semi-structured data on the client side, rather than in SQL Server. It's also easier and cheaper to scale application servers since you don't need to pay the SQL Server licensing cost.

Be aware of CLR, **external languages** code, and extended stored procedures with complex logic. Avoid extensive function calls, especially with user-defined functions. They add overhead and may lead to less efficient execution plans when they are not inlined.

Pay attention to views usage. Depending on the database schema and definition, views may introduce unnecessary joins, accessing tables the queries do not need to access. This is especially common if the tables do not have proper foreign keys defined.

Scripts for Troubleshooting High CPU Load

I'd like to provide you with a couple of scripts that are helpful when troubleshooting high CPU load. The first, in Listing 6-1, shows you CPU load on the server during the last 256 minutes. The data is measured once per minute, so it may miss short CPU load bursts that occur in between measurements.

Example 6-1. Getting CPU Load History

```

DECLARE
    @now BIGINT;

SELECT @now = cpu_ticks / (cpu_ticks / ms_ticks)
FROM sys.dm_os_sys_info WITH (NOLOCK);
;WITH RingBufferData([timestamp], rec)
AS
(
    SELECT [timestamp], CONVERT(XML, record) AS rec
    FROM sys.dm_os_ring_buffers WITH (NOLOCK)
    WHERE
        ring_buffer_type = N'RING_BUFFER_SCHEDULER_MONITOR' AND
        record LIKE N'%<SystemHealth>%'
)
,Data(id, SystemIdle, SQLCPU, [timestamp])
AS
(
    SELECT
        rec.value('(/Record/@id)[1]', 'int')
        ,rec.value
            ('(/Record/SchedulerMonitorEvent/SystemHealth/SystemIdle)
[1]', 'int')
        ,rec.value
            ('(/Record/SchedulerMonitorEvent/SystemHealth/ProcessUtilization)
[1]', 'int')
        , [timestamp]
    FROM RingBufferData
)
SELECT TOP 256
    dateadd(MS, -1 * (@now - [timestamp]), getdate()) AS [Event Time]
    ,SQLCPU AS [SQL Server CPU Utilization]
    ,SystemIdle AS [System Idle]
    ,100 - SystemIdle - SQLCPU AS [Other Processes CPU Utilization]
FROM Data
ORDER BY id desc
OPTION (RECOMPILE, MAXDOP 1);

```

Figure 6-1 illustrates the output of the code. Data in the [Other Processes CPU Utilization] column shows CPU load in the system outside of SQL Server. If that load is significant, analyze what processes are running on the server and generating it.

	Event Time	SQL Server CPU Utilization	System Idle	Other Processes CPU Utilization
1	2021-03-27 09:44:52.903	61	35	4
2	2021-03-27 09:43:52.877	62	36	2
3	2021-03-27 09:42:52.863	56	42	2
4	2021-03-27 09:41:52.847	68	30	2
5	2021-03-27 09:40:52.833	77	20	3
6	2021-03-27 09:39:52.813	63	34	3
7	2021-03-27 09:38:52.800	64	33	3
8	2021-03-27 09:37:52.787	62	35	3

Figure 6-1. Script output showing CPU load history

Listing 6-2 helps you analyze CPU load per database. This may be beneficial when your server hosts multiple databases, and you are considering splitting the busy ones between different servers. (Please note: This script uses plan cache data, so the output is imprecise.)

Example 6-2. Per-database CPU Load

```

;WITH DBCPU
AS
(
    SELECT
        pa.DBID, DB_NAME(pa.DBID) AS [DB]
        ,SUM(qs.total_worker_time/1000) AS [CPUTime]

```

```

FROM
    sys.dm_exec_query_stats qs WITH (NOLOCK)
CROSS APPLY
    (
        SELECT CONVERT(INT, value) AS [DBID]
        FROM sys.dm_exec_plan_attributes(qs.plan_handle)
        WHERE attribute = N'dbid'
    ) AS pa
GROUP BY pa.DBID
)
SELECT
    [DB]
    ,[CPUTime] AS [CPU Time (ms)]
    ,CONVERT(decimal(5,2), 1. * [CPUTime] /
        SUM([CPUTime]) OVER() * 100.0) AS [CPU Percent]
FROM DBCPU
WHERE DBID <> 32767 -- ResourceDB
ORDER BY [CPUTime] DESC;

```

Figure 6-2 shows the output from a production server.

	DB	CPU Time (ms)	CPU Percent
1	Appdb	55332331	73.69
2	AppDW	16190787	21.56
3	master	2264420	3.02
4	DBA	918845	1.22

Figure 6-2. Script output showing CPU load per database

Non-Optimized Query Patterns to Watch For

In case of non-optimized queries, there are two distinct patterns that can trigger high CPU load. I call them “the worst offenders” and “death by a thousand cuts.” (This terminology is by no means standard – it’s just how I like to differentiate between them.)

The worst offenders

The “worst offenders” occur when you have one or more expensive, long-running queries that generate heavy CPU load: think of non-optimized queries with parallel execution plans that scan millions of rows and perform sorting and aggregation. They can bring a server to its knees, especially if you have several running simultaneously.

Fortunately, it is easy to detect the worst offenders in real time by querying the `sys.dm_exec_requests` view and analyzing the `cpu_time` column. (You can use the code from Listing 2-3 in Chapter 2 to do that.)

WARNING BOX

A word of caution: The code in Listing 2-3 filters out system processes with a `session_id` below 50. In some cases, you might want to remove this filter and analyze all sessions running on the server. Keep in mind that some sessions may have been running since SQL Server startup and have high cumulative `cpu_time`: pay attention to the request start time.

Death by a thousand cuts

With the second pattern, “death by a thousand cuts,” the load on the server is generated by a large number of simultaneously running requests. Each request may be relatively small and even optimized; however, the sheer number of requests drives CPU usage and server load up.

This case is more challenging to handle. While query optimization (covered in Chapter 5) may help, you’ll likely have to optimize a large number of queries consuming significant time and effort. It often requires refactoring database schemas, code, and applications on a massive scale to achieve results.

In the end, you have to reduce the load on the server to address the problem. Let’s talk about several other factors that can increase that load, starting with the query compilation and plan caching processes.

Query Compilation and Plan Caching

Every time you submit a query to the system, SQL Server needs to compile and optimize it. This process is resource intensive, so SQL Server tries to minimize the number of compilations by caching execution plans for later reuse. In addition to regular client queries and batches, it caches plans of various objects, such as stored procedures, triggers, and user-defined functions. The memory area where these are stored is called the *plan cache*.

SQL Server uses different algorithms to determine which plans to remove from the cache in case of memory pressure. For ad-hoc queries, this selection is based on how often the plan is reused. For other types of plans, the cost of plan generation is also factored into the decision.

SQL Server recompiles queries when it suspects that currently cached plans are no longer valid. This may happen if the plan references objects whose schemas have changed, or because of stale statistics. SQL Server checks to see if the statistics are outdated when it looks up a plan from the cache, and it recompiles the query if they are. That recompilation, in turn, triggers a statistics update.

Plan caching and reuse can significantly reduce the number of compilations and the CPU load, as I will demonstrate later in the chapter. However, it can also introduce problems. Let's look at some of the most common issues that arise, starting with *parameter sensitivity* in *parameter-sensitive plans*. (This is sometimes called *parameter sniffing*, which just describes the SQL Server behavior that leads to that issue).

Parameter-Sensitive Plans

Except for some trivial queries, SQL Server always offers multiple options for generating the execution plan for the query. It can use different indexes to access data, select different join types, and choose among operators and execution strategies.

By default, SQL Server analyzes (*sniffs*) parameter values at the time of optimization and generates and caches an optimal plan for those values. Nothing is wrong with this behavior—though it can, if your data is unevenly distributed, lead to a cached plan that is optimal for atypical, rarely used

parameter values but highly inefficient for queries with more common parameters.

I'm sure we're all experienced a situation where some queries or stored procedures suddenly started taking much longer to complete, even though there were no recent changes in the system. In most cases, these situations happen when queries are recompiled after a statistics update, due to parameter sniffing.

Let me show you an example. The script in Listing 6-3 creates a table and populates it with 1 million rows, evenly distributed across 10 StoreId values (a little more than 100,000 rows per StoreId), along with 10 rows with a StoreId of 99.

Example 6-3. Parameter-sensitive plans: Table creation

```
CREATE TABLE dbo.Orders
(
    OrderId INT NOT NULL IDENTITY(1,1),
    OrderNum VARCHAR(32) NOT NULL,
    CustomerId UNIQUEIDENTIFIER NOT NULL,
    Amount MONEY NOT NULL,
    StoreId INT NOT NULL,
    Fulfilled BIT NOT NULL
);
;WITH N1(C) AS (SELECT 0 UNION ALL SELECT 0) -- 2 rows
,N2(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N1 AS T2) -- 4 rows
,N3(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N2 AS T2) -- 16 rows
,N4(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN T2 CROSS JOIN N2 AS T3)
-- 1024 rows
,N5(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N4 AS T2 ) -- 1,048,576
rows
,IDs(ID) AS (SELECT ROW_NUMBER() OVER (ORDER BY (SELECT NULL)) FROM
N5)
INSERT INTO dbo.Orders(OrderNum, CustomerId, Amount, StoreId,
Fulfilled)
select
    'Order: ' + convert(varchar(32),ID)
    ,newid()
    ,ID % 100
    ,ID % 10
    ,1
from IDs;
INSERT INTO dbo.Orders(OrderNum, CustomerId, Amount, StoreId,
Fulfilled)
```

```

        select top 10 OrderNum, CustomerId, Amount, 99, 0
        from dbo.Orders
        order by OrderId;
CREATE UNIQUE CLUSTERED INDEX IDX_Orders_OrderId
ON dbo.Orders(OrderId);
CREATE NONCLUSTERED INDEX IDX_Orders_CustomerId
ON dbo.Orders(CustomerId);
CREATE NONCLUSTERED INDEX IDX_Orders_StoreId
ON dbo.Orders(StoreId);

```

Next, let's create a stored procedure that calculates the total sales amount for a specific store (Listing 6-4). I'm using a stored procedure in this example; however, parameterized queries called from client applications would behave the same way.

Listing 6-4. Parameter-sensitive plans: Stored procedure

```

CREATE PROC dbo.GetTotalPerStore(@StoreId int)
AS
    SELECT SUM(Amount) as [Total Amount]
    FROM dbo.Orders
    WHERE StoreId = @StoreId;

```

With the current data distribution, when the stored procedure is called with any `@StoreId` other than 99, the optimal execution plan involves scanning the clustered index in the table. However, if `@StoreId=99`, a better execution plan would be to use an index seek on `IDX_Orders_StoreId` index, with the key lookup afterwards.

Let's call the stored procedure twice: the first time with `@StoreId=5` and the second time with `@StoreId=99`, as shown in Listing 6-5.

Example 6-5. Parameter-sensitive plans: Calling the procedure (Test 1)

```

EXEC dbo.GetTotalPerStore @StoreId = 5;
EXEC dbo.GetTotalPerStore @StoreId = 99;

```

As you can see from the execution plan in Figure 6-3, SQL Server compiles the stored procedure, caches the plan with the first call, and reuses the plan later. Even though this plan is less efficient for the second call with `@StoreId=99`, it may be acceptable when those calls are rare, which is expected with such a data distribution.

Query 1: Query cost (relative to the batch): 50%

```
SELECT SUM(Amount) as [Total Amount] FROM dbo.Orders WHERE StoreId = @StoreId
```

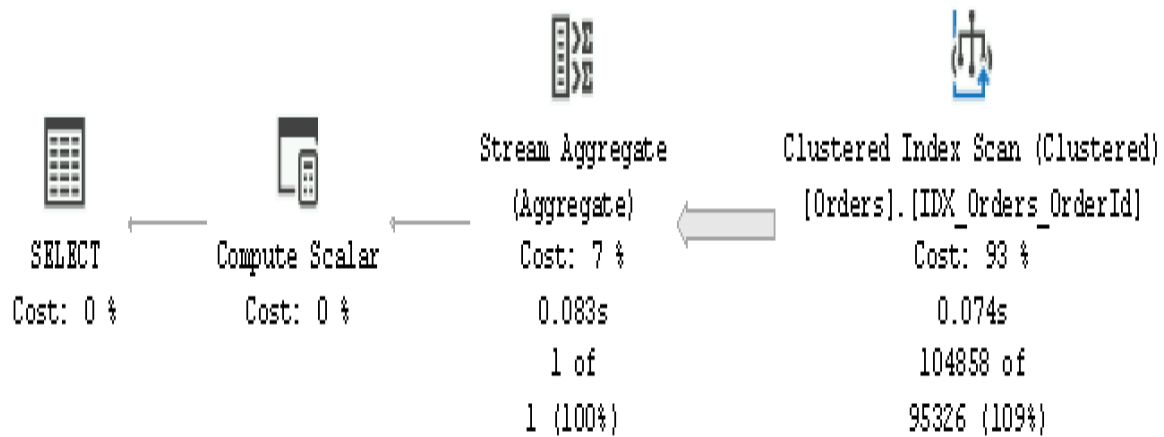


Table 'Orders'. Scan count 1, logical reads 8468

Query 2: Query cost (relative to the batch): 50%

```
SELECT SUM(Amount) as [Total Amount] FROM dbo.Orders WHERE StoreId = @StoreId
```

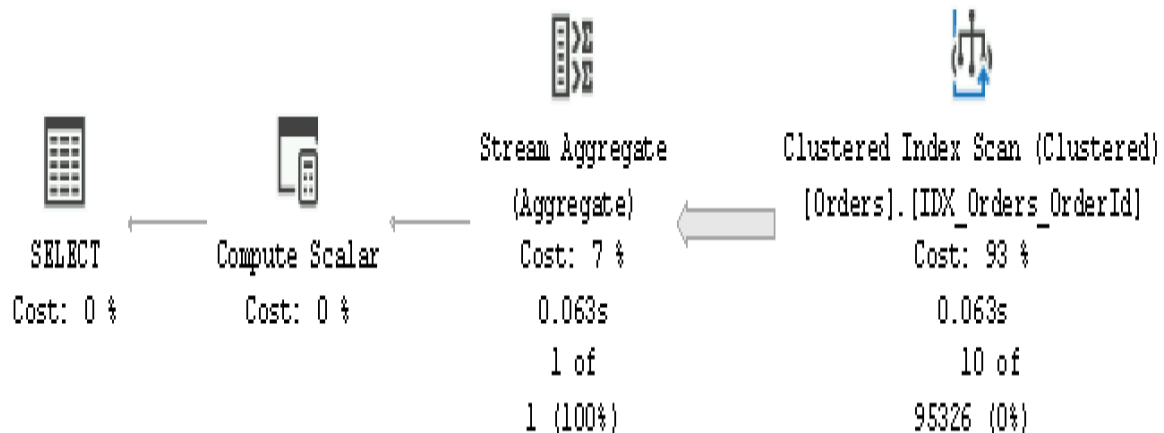


Table 'Orders'. Scan count 1, logical reads 8468

Figure 6-3. Execution plans of the queries (Test 1).

Now let's take a look at what happens if we swap those calls when the plan is *not* cached (Listing 6-6). I am clearing the plan cache with the DBCC FREEPROCCACHE command – do not run this demo on a production

server! Note, however, that the same thing can happen when a statistics update triggers the query to recompile.

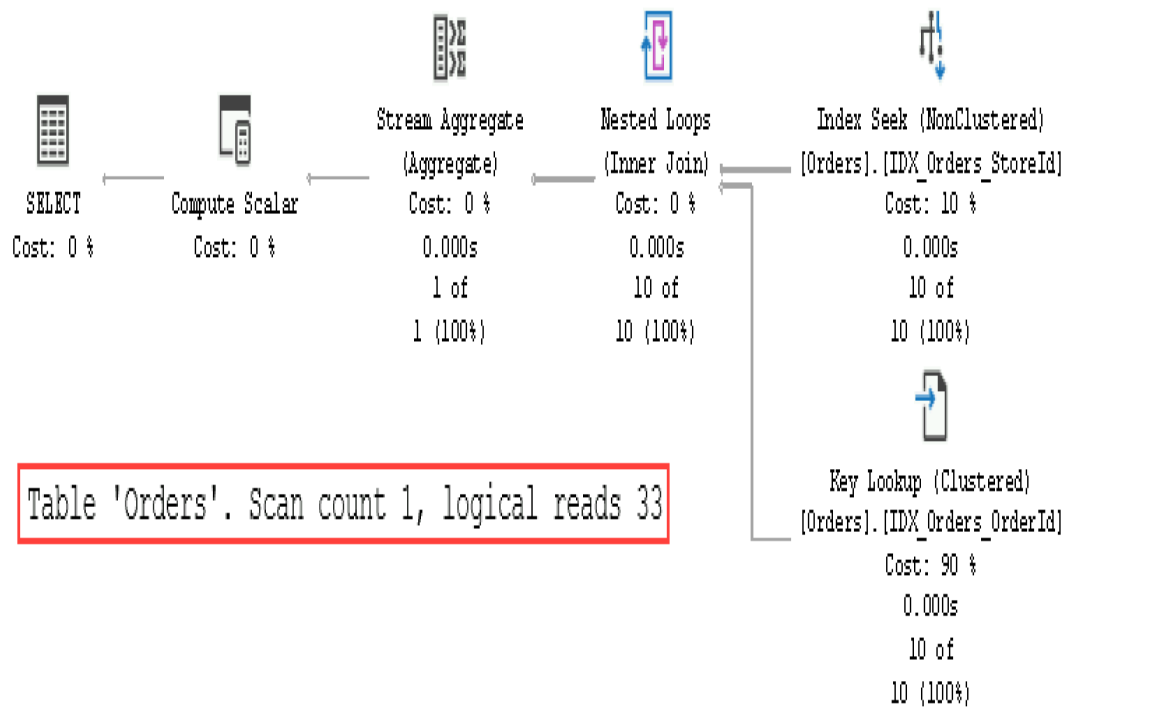
Example 6-6. Parameter-sensitive plans: Calling the procedure (Test 2)

```
DBCC FREEPROCCACHE;  
EXEC dbo.GetTotalPerStore @StoreId = 99;  
EXEC dbo.GetTotalPerStore @StoreId = 5;
```

As you can see in Figure 6-4, SQL Server now caches the plan compiled for the `@StoreId=99` parameter value. Even though this plan is more efficient when the stored procedure is called with this parameter, it is highly inefficient for other `@StoreId` values.

Query 1: Query cost (relative to the batch): 50%

SELECT SUM(Amount) as [Total Amount] FROM dbo.Orders WHERE StoreId = @StoreId



Query 2: Query cost (relative to the batch): 50%

SELECT SUM(Amount) as [Total Amount] FROM dbo.Orders WHERE StoreId = @StoreId

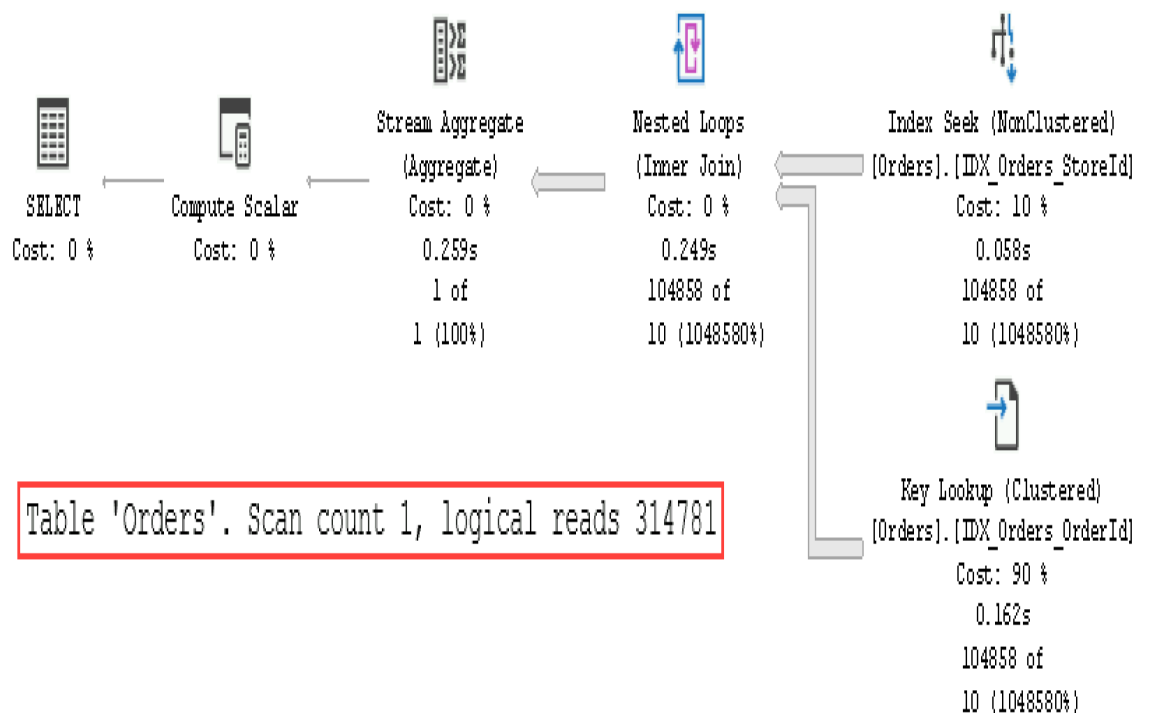


Figure 6-4. Execution plans of the queries (Test 2).

Inefficient parameter-sensitive plans often become the “worst offender” queries that drive CPU load up. As I mentioned, you can detect those queries with the `sys.dm_exec_requests` view (Listing 2-3) and recompile them to remediate the issue.

You can force stored procedures and other T-SQL modules to recompile with the `sp_recompile` stored procedure. For ad-hoc queries, you can call `DBCC FREEPROCCACHE`, providing `plan_handle` or `sql_handle` as the parameter. Finally, if you have Query Store enabled, you can force a more efficient query execution plan there.

Obviously, it is better to address the root cause of the issue. First, see if there are any opportunities for query tuning, which would eliminate the plans’ parameter sensitivity. We usually end up with parameter-sensitive plans because there are no efficient plans that do not depend on parameter values. For example, if you have the Amount column included to the `IDX_Orders_StoreId` index, that index would become covering. SQL Server can use it for all parameter values regardless of how many rows will be read, because the *Key Lookup* operation will no longer be required.

If you are using SQL Server 2017 or above, you can benefit from *automatic plan correction*, which is part of the *automatic tuning* technology. When this feature is enabled, SQL Server can detect parameter sniffing issues and automatically force the last known good plan that was used before regression occurred.

Automatic plan correction relies on the *Force Plan* feature of Query Store and, as you can guess, requires Query Store to be enabled in the database. Moreover, you need to enable it in the database with `ALTER DATABASE SET AUTOMATIC_TUNING (FORCE_LAST_GOOD_PLAN = ON)` statement. You can read more about it in the Microsoft *documentation*.

If neither of those options works, you can force the recompilation of either stored procedure using `EXECUTE WITH RECOMPILE` or a statement-level recompile with `OPTION (RECOMPILE)` clauses. Listing 6-7 shows the latter approach.

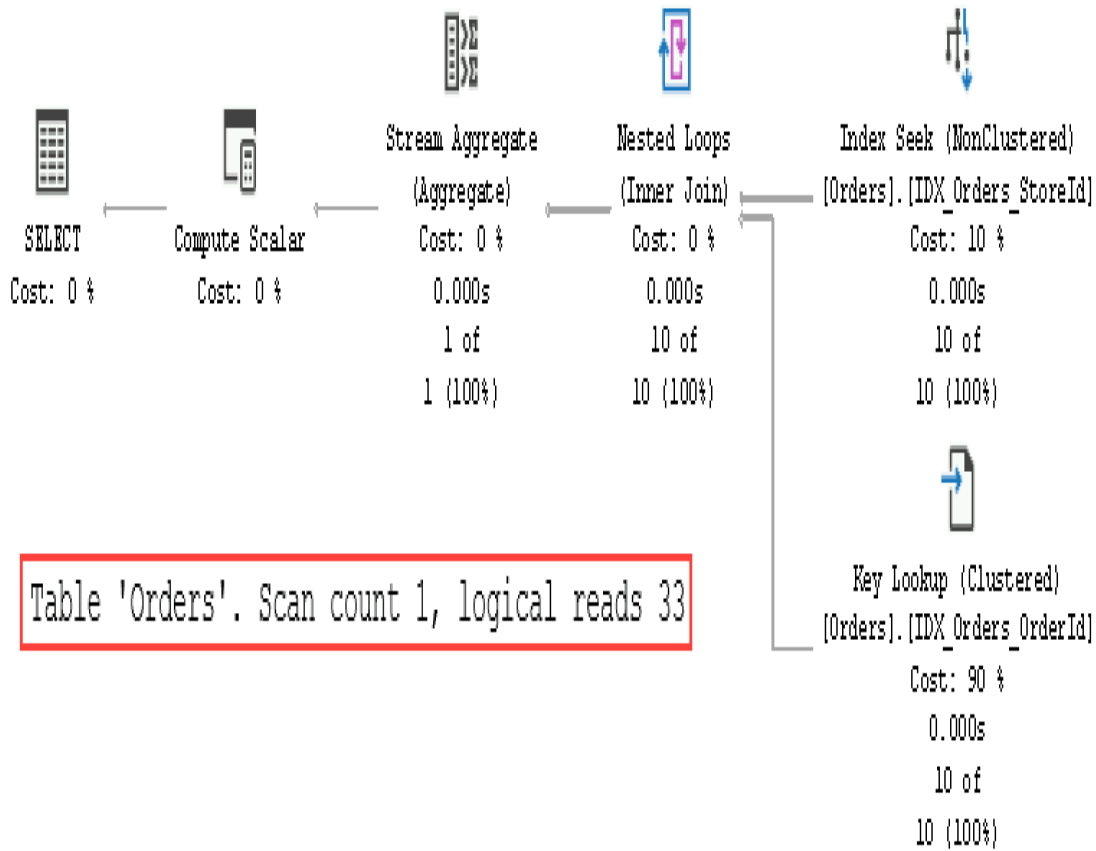
Example 6-7. Parameter-sensitive plans: Statement-level recompile

```
ALTER PROC dbo.GetTotalPerStore(@StoreId int)
AS
    SELECT SUM(Amount) as [Total Amount]
    FROM dbo.Orders
    WHERE StoreId = @StoreId
        OPTION (RECOMPILE);
GO
EXEC dbo.GetTotalPerStore @StoreId = 99;
EXEC dbo.GetTotalPerStore @StoreId = 5;
```

Figure 6-5 shows that SQL Server recompiles the query sniffing parameters during each call.

Query 1: Query cost (relative to the batch): 0%

SELECT SUM(Amount) as [Total Amount] FROM dbo.Orders WHERE StoreId = @StoreId OPTION (RECOMPILE)



Query 2: Query cost (relative to the batch): 100%

SELECT SUM(Amount) as [Total Amount] FROM dbo.Orders WHERE StoreId = @StoreId OPTION (RECOMPILE)

Missing Index (Impact 96.7039): CREATE NONCLUSTERED INDEX [<Name of Missing Index, sysname,>] ON |

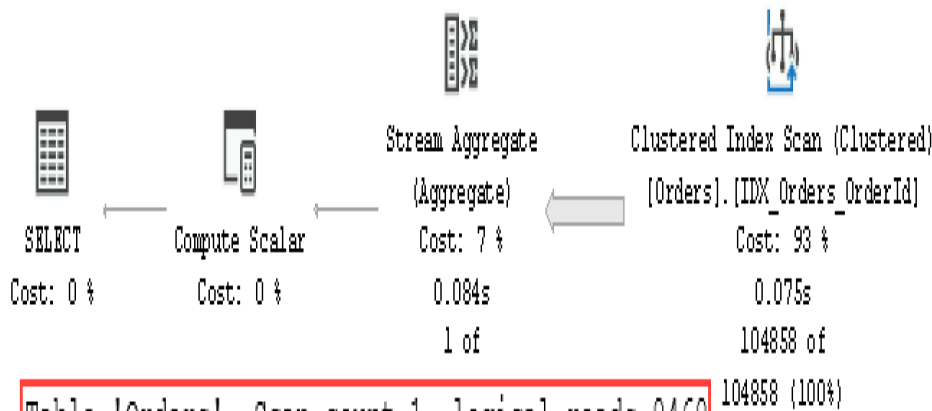


Figure 6-5. Execution plans of the queries (statement-level recompile).

Forcing the recompile will allow you to get the most efficient execution plans on each call—at the cost of constant recompilation overhead. This approach may be completely acceptable with infrequently executed queries; however, with frequently executed ones it may lead to noticeable CPU increase, as I’ll show later in the chapter.

You can address this by utilizing another hint – OPTIMIZE FOR. This hint allows you to specify parameter values for Query Optimizer to use during optimization. For example, with OPTIMIZE FOR (@StoreId=5) hint, Query Optimizer will not sniff @StoreId, instead optimizing it for the value of 5 all the time.

As you can guess, the danger of using the OPTIMIZE FOR hint is that data distribution changes. For example, if the store with @StoreId=5 went out of business, you’d end up with highly inefficient execution plans.

Fortunately, there is another form of this hint: OPTIMIZE FOR UNKNOWN. With this hint, SQL Server performs an optimization based on the most statistically common value in the table. In our case, this hint would lead to the plan with clustered index scan, which is expected with data distribution in the table.

You can use the hints OPTIMIZE FOR UNKNOWN (in all SQL Server versions after 2008) or DISABLE_PARAMETER_SNIFFING (in SQL Server 2016 and up) on the query level. Both hints are essentially the same. In SQL Server 2016, you can also control it at the database level with the PARAMETER_SNIFFING database option. Finally, you can disable parameter sniffing on the server-level with trace flag T4136. This trace flag also works in SQL Server versions prior to 2016.

NOTE

In my experience, disabling parameter sniffing leads to better and more stable execution plans in multi-tenant systems and systems with very uneven data distribution. Your mileage may vary, but it’s worth trying if you see a large number of parameter-sensitive plans in your system.

Caching inefficient parameter-sensitive plans increases CPU load. Unfortunately, that's not the only issue you can encounter with plan caching.

Parameter-Value Independence

Cached execution plans need to be valid for all possible combinations of parameters in future calls. As a result, even with parameter sniffing, SQL Server will not generate an execution plan that cannot be used with some parameters in the future when SQL Server expects to cache it.

This sounds a bit confusing, so let me demonstrate it with a simple example. Listing 6-8 shows a very common (and very bad) pattern: the stored procedure accepts optional parameters, using a single query to cover them all.

Example 6-8. Parameter-Value Independence

```
CREATE PROC dbo.SearchOrders
(
    @StoreId INT
    ,@CustomerId UNIQUEIDENTIFIER
)
AS
    SELECT OrderId, CustomerId, Amount, Fulfilled
    FROM dbo.Orders
    WHERE
        ((@StoreId IS NULL) OR (StoreId = @StoreId)) AND
        ((@CustomerId IS NULL) OR (CustomerId = @CustomerId));
GO
EXEC dbo.SearchOrders
    @StoreId = 99
    ,@CustomerId = 'A65C047D-5B08-4041-B2FE-8E3DD6570B8A';
```

Regardless of what parameters you are using at the time of compilation, you will get a plan similar to the one shown in Figure 6-6. Even though there are indexes and SARGable predicates on both the CustomerId and StoreId columns, SQL Server uses the *Index Scan* operation instead of *Index Seek*. Unfortunately, SQL Server cannot use index seek because the plan needs to be cached and reused in the future; this plan would not be valid if the *seek predicate* (@StoreId parameter in the plan below) was not provided.

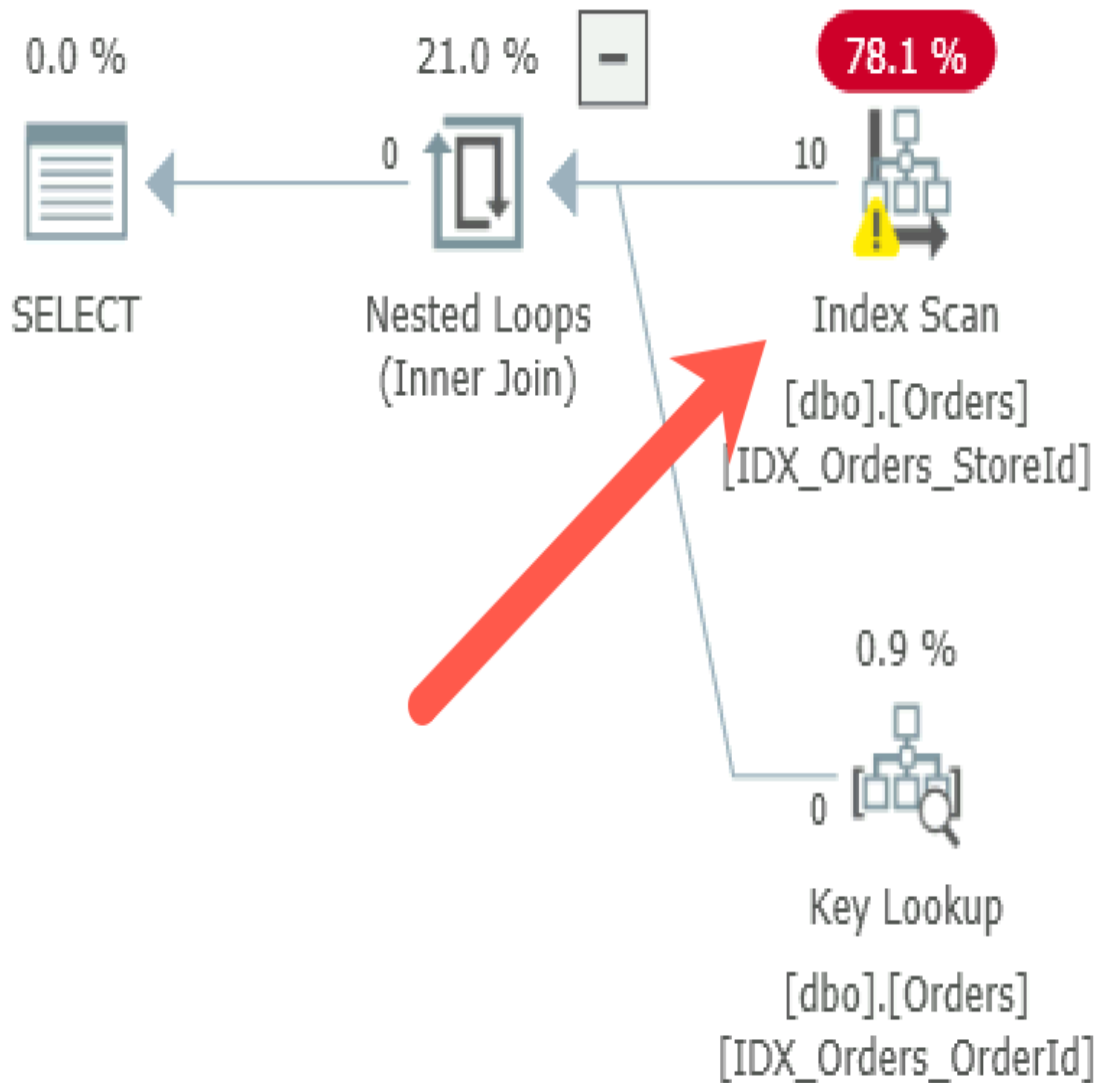


Figure 6-6. Parameter-value independence: Execution plan of the stored procedure

The statement-level recompile addresses the problem, again at the cost of additional compilation overhead. As noted, this overhead may be acceptable for infrequently executed queries.

As another option, you could rewrite the code using IF statements that cover all possible combinations of parameters. If you did, SQL Server would cache the plan for each statement. It would work in simple cases, but would quickly become unmanageable as the number of parameters grows.

Finally, writing the code using dynamic SQL is a completely valid option in many cases. (Listing 6-9 shows how to do that.) Be careful, of course, and

utilize parameters to prevent SQL injection.

Example 6-9. Dynamic SQL implementation

```
ALTER PROC dbo.SearchOrders
(
    @StoreId INT
    ,@CustomerId UNIQUEIDENTIFIER
)
AS
BEGIN
    DECLARE
        @SQL nvarchar(max) =
        N'SELECT OrderId, CustomerId, Amount, Fulfilled
        FROM dbo.Orders
        WHERE
            (1=1)' +
            IIF(@StoreId IS NOT NULL, N'AND (StoreId = @StoreId)', '') +
            IIF(@CustomerId IS NOT NULL, N'AND (CustomerId =
@CustomerId)', '');
        EXEC sp_executesql
            @SQL = @SQL
            ,@Params = N'@StoreId INT, @CustomerId UNIQUEIDENTIFIER'
            ,@StoreId = @StoreId, @CustomerId = @CustomerId;
END
```

There are other times when caching and reusing plans may lead to inefficient plans. One case, which is often overlooked, involves filtered indexes.

The query in Listing 6-10 will not use a filtered index even if you call it with `@Fulfilled = 0` value. This happens because the cached execution plan that uses the filtered index will not be valid for `@Fulfilled = 1` calls.

Example 6-10. Query that would not use filtered index

```
CREATE NONCLUSTERED INDEX IDX_Orders_ActiveOrders_Filtered
ON dbo.Orders(OrderId)
INCLUDE(Fulfilled)
WHERE Fulfilled = 0;
GO
DECLARE
    @Fulfilled bit = 0;
SELECT COUNT(*) AS [Active Order Count]
FROM dbo.Orders
WHERE Fulfilled = @Fulfilled;
```

NOTE

Always add *all* columns from the filter to either key or included columns in filtered indexes. This leads to more efficient execution plans.

Unfortunately, this problem can also occur due to auto-parameterization, which I will discuss later in the chapter. But first, let's look at compilations and the overhead they introduce.

Compilation and Parameterization

As you know, SQL Server caches and reuses execution plans for T-SQL modules and ad-hoc client queries and batches. For ad-hoc queries, however, the plans are reused only for *identical* queries. There are a few factors that dictate that.

First, identical queries need to be *exactly* the same: a complete character-for-character match. Look at the queries in Listing 6-11. Only two are identical (the first and second), even though all of these queries are logically the same.

Example 6-11. Identical queries

```
SELECT COUNT(*) FROM dbo.Orders WHERE StoreId = 99;  
SELECT COUNT(*) FROM dbo.Orders WHERE StoreId = 99;  
SELECT COUNT(*) FROM dbo.Orders WHERE StoreId=99;  
select count(*) from dbo.Orders where StoreId = 99;
```

In addition, some of the SET options affect plan reuse, including ANSI_NULL_DLFT_OFF, ANSI_NULL_DLFT_ON, ANSI_NULL, ANSI_PADDING, ANSI_WARNING, ARITHABORT, CONCAT_NULL_YELDS_NULL, DATEFIRST, DATEFORMAT, FORCEPLAN, DATEFORMAT, LANGUAGE, NO_BROWSETABLE, NUMERIC_ROUNDABORT, and QUOTED_IDENTIFIER. Plans generated with one set of SET options cannot be reused by sessions that use a different set of SET options.

You've probably noticed that I keep emphasizing the point that query compilation and optimization processes are resource intensive and may introduce significant CPU load with a heavy ad-hoc workload. To demonstrate this, I have created a small application that runs simple queries from Listing 6-12 in a loop in multiple threads. You can download it from the companion materials of the book.

In the first test case, the application runs ad-hoc queries using non-parameterized CustomerId values (the queries are constructed in the application). Each query in the call is unique and needs to be compiled. The second test, on the other hand, uses parameterized query. The plan for this query can be reused across calls.

Example 6-12. Ad hoc versus parameterized workload

```
-- Test Case 1
SELECT TOP 1 OrderId
FROM dbo.Orders
WHERE CustomerId = '<ID Generated in the app>';
-- Test Case 2
SELECT TOP 1 OrderId
FROM dbo.Orders
WHERE CustomerId = @CustomerId;
```

Both of the queries are extremely light. Moreover, I ran them in the test environment with enough memory to cache the entire table and eliminate physical I/O.

Figure 6-7 illustrates the performance metrics collected during the tests. As you can see, during the second test (on the right), the system was able to handle almost 6 times more requests per second than during the first test.

\\SQL2019-1		\\SQL2019-1	
Processor Information		Processor Information	
_Total		_Total	
% Processor Time	99.206	% Processor Time	98.413
SQLServer:Plan Cache		SQLServer:Plan Cache	
Total		Total	
Cache Object Counts	41,892.000	Cache Object Counts	3.000
SQLServer:SQL Statistics		SQLServer:SQL Statistics	
Batch Requests/sec	3,898.249	Batch Requests/sec	22,584.970
SQL Compilations/sec	3,897.233	SQL Compilations/sec	0.000

Figure 6-7. Ad hoc versus parameterized workload throughput

Obviously, this scenario is completely synthetic; in real life, you are unlikely to see a situation where SQL Server has to spend majority of its time compiling queries. Nevertheless, in systems with heavy ad-hoc workloads, the impact of compilations can be very significant. In addition to CPU load, there is also an impact on memory, which I will discuss in the next chapter.

There are three SQL Server: SQL Statistics performance counters that can help you to see system throughput and number of compilations.

Batch Requests/sec

Batch Requests/sec counter shows the number of batches SQL Server receives per second. Higher values indicate higher system load and throughput.

SQL Compilations/sec

SQL Compilations/sec counter shows how many compilations SQL Server performs every second. The higher this number is, the more compilations and, therefore, the more overhead you have.

SQL Re-Compilations/sec

SQL Re-Compilations/sec counter gives you the number of recompilations for already cached execution plans. This may happen due to frequent changes in underlying data in both users and temporary tables.

In a properly tuned OLTP system, the number of compilations and recompilations should be just a fraction of the total number of batch requests. If that is not the case, analyze and reduce the compilations. (We will talk about how to analyze plan cache data in the next chapter.)

Non-parameterized, ad-hoc client queries are the most common cause of compilations. As you can guess, the best approach is changing the queries and parameterizing them. Unfortunately, this usually requires you to change the client code, which is not always possible.

Fortunately, there is another option: *auto-parameterization*.

Auto-Parameterization

SQL Server tries to reduce compilation overhead by replacing constants in ad-hoc queries with parameters and cache compiled plans as if the queries were parameterized. When this happens, similar ad-hoc queries that use different constants can reuse cached plans.

Let's look at the example and run the queries in Listing 6-13. As before, I am clearing the plan cache with the DBCC FREEPROCCACHE command to reduce the size of the output.

Example 6-13. Auto-parameterization

```
DBCC FREEPROCCACHE
GO
SELECT * FROM dbo.Orders WHERE OrderId = 1;
GO
SELECT * FROM dbo.Orders WHERE OrderId = 2;
GO
SELECT
    p.usecounts, p.cacheobjtype, p.objtype, p.size_in_bytes, t.
```

```

[text]
FROM
    sys.dm_exec_cached_plans p CROSS APPLY
        sys.dm_exec_sql_text(p.plan_handle) t
WHERE
    p.cacheobjtype LIKE 'Compiled Plan%' AND
    t.[text] LIKE '%Orders%'
ORDER BY
    p.objtype DESC
OPTION (RECOMPILE);

```

Figure 6-8 shows the output of the last statement from the code. As you can see, there are three entries in the plan cache: a compiled plan used for both auto-parameterized ad-hoc queries, and two other objects called *shell queries*. Each shell query uses about 16KB of memory and stores information about the original ad-hoc query and links it to the compiled plan.

	usecounts	cacheobjtype	objtype	size_in_bytes	text
1	2	Compiled Plan	Prepared	40960	(@1 tinyint)SELECT * FROM [dbo].[Orders] WHERE [OrderId]=@1
2	1	Compiled Plan	Adhoc	16384	SELECT * FROM dbo.Orders WHERE OrderId = 2;
3	1	Compiled Plan	Adhoc	16384	SELECT * FROM dbo.Orders WHERE OrderId = 1;

Figure 6-8. Plan cache after auto-parameterization

Simple parameterization

By default, SQL Server uses SIMPLE parameterization, and it is very conservative in parameterizing queries. Simple parameterization only happens when a cached plan is considered *safe to parameterize*. This means that the plan would have the same shape and cardinality estimations, even when constant or parameter values change.

For example, a plan with a nonclustered index seek and key lookup on a unique index is safe because nonclustered index seek would never return more than one row, regardless of parameter value. However, the same

operation on a non-unique index is not safe. Different parameter values would lead to different cardinality estimations; this could make a clustered index scan the better option for some parameter values.

Moreover, there are many language constructs that prevent simple parameterization, including IN, TOP, DISTINCT, JOIN, UNION, and subqueries. In practice, this means the majority of queries will not be auto-parameterized.

Forced Parameterization

Alternatively, SQL Server can use FORCED parameterization. This can be enabled at the database level with the ALTER DATABASE SET PARAMETRIZATION FORCED command, or at the query level with a PARAMETRIZATION FORCED hint. In this mode, SQL Server auto-parameterizes most ad-hoc queries (with very few exceptions).

Figure 6-9 shows results of the first test case (non-parameterized ad-hoc queries) from Listing 6-12 after I enable forced parameterization in the database. While the system throughput is still significantly lower than with a properly parameterized workload, it is much better than with simple parameterization (Figure 6-7) that did not auto-parameterize the query. SQL Server still needs to spend CPU time to auto-parameterize queries; however, it can reuse the cached execution plan and does not need to optimize all queries.

\\SQL2019-1

Processor Information	_Total
% Processor Time	100.000
SQLServer:Plan Cache	_Total
Cache Object Counts	160,039.000
SQLServer:SQL Statistics	
Batch Requests/sec	12,887.195
Forced Parameterizations/sec	12,887.195
SQL Compilations/sec	12,887.195

Figure 6-9. Throughput with forced parameterization

Enabling forced parameterization may significantly reduce compilation overhead and CPU load in systems with heavy ad-hoc workload. You will of course get different results in different systems, but I've had a few cases where enabling forced parameterization reduced CPU load by as much as 25 to 30%.

Forced parameterization has its downsides, though. When it is enabled, SQL Server starts to auto-parameterize the majority of ad-hoc queries, which will open the door to parameter-sensitive plans and parameter-sniffing-related issues. You can expect some ad-hoc queries to regress because of that.

NOTE

I recommend that you consider disabling parameter sniffing after you enable forced parameterization. While this may not be the best option for *every* system, I have found it helpful in most cases.

Fortunately, you are not always forced to take an all-or-nothing approach. As I mentioned, you can enable forced parameterization on the query level with the `PARAMETERIZATION FORCED` query hint. This is useful when you have just a handful of non-parameterized ad-hoc queries and do not want to enable forced parameterization globally.

If you don't have access to the source code, you can force the hint through plan guides. Listing 6-14 shows how to do that. It uses two stored procedures. The first, `sp_get_query_template`, creates the query template based on the sample query provided as parameter. You can use any constant values in an ad-hoc query for template creation. The second procedure, `sp_create_plan_guide`, creates the plan guide.

Example 6-14. Applying forced parameterization through a plan guide

```
DECLARE
    @stmt nvarchar(max)
    ,@params nvarchar(max)
    ,@query nvarchar(max) =
N'SELECT TOP 1 OrderId FROM dbo.Orders WHERE CustomerId =
''B970D68B-F88E-438B-9B04-6EDE47CC1D9A''';
EXEC sp_get_query_template
    @querytext = @query
    ,@temptatetext = @stmt output
    ,@params = @params output;

EXEC sp_create_plan_guide
    @type = N'TEMPLATE'
    ,@name = N'forced_parameterization_plan_guide'
    ,@stmt = @stmt
    ,@module_or_batch = null
    ,@params = @params
    ,@hints = N'OPTION (PARAMETERIZATION FORCED)';
```

You can download the test application from this book's companion materials and repeat the load tests to validate that the plan guide is working in a database that uses `SIMPLE` parameterization.

In some cases, you need to do the opposite and force simple parameterization for specific queries in the database that use forced parameterization. This can happen when some ad-hoc queries have parameter-sensitive plans. Listing 6-15 shows how you can force simple

parameterization through the plan guide, allowing the ad-hoc query to utilize filtered index. You need to provide the statement to `sp_create_plan_guide` stored procedure as if it already had been auto-parameterized. You can obtain it from the plan cache, as shown in the listing, along with the parameters of the statement.

The first query in Listing 6-15 is the one to which I am applying the plan guide.

Example 6-15. Applying simple parameterization through a plan guide

```
SELECT OrderId
FROM dbo.Orders
WHERE Fulfilled = 0;
GO
SELECT
    SUBSTRING(qt.text, (qs.statement_start_offset/2)+1,
    ((
        CASE qs.statement_end_offset
            WHEN -1 THEN DATALENGTH(qt.text)
            ELSE qs.statement_end_offset
        END - qs.statement_start_offset)/2)+1) AS SQL
    ,qt.text AS [Full SQL]
FROM
    sys.dm_exec_query_stats qs with (nolock)
    CROSS APPLY sys.dm_exec_sql_text(qs.sql_handle) qt
WHERE
    qt.text like '%Fulfilled%'
OPTION(RECOMPILE, MAXDOP 1);
DECLARE
    @stmt nvarchar(max) =
        N'select OrderId from dbo . Orders where Fulfilled = @0'
    ,@params nvarchar(max) = N'@0 int'

-- Creating plan guide
EXEC sp_create_plan_guide
    @type = N'TEMPLATE'
    ,@name = N'simple_parameterization_plan_guide'
    ,@stmt = @stmt
    ,@module_or_batch = null
    ,@params = @params
    ,@hints = N'OPTION (PARAMETERIZATION SIMPLE)';
```

SQL Server does not auto-parameterize queries in stored procedures and other T-SQL modules. You can move some ad-hoc queries to stored

procedures, avoiding parameter-sensitivity issues when forced parameterization is enabled.

Finally, I'd like to reiterate: *recompilations may lead to significant CPU overhead in systems with heavy ad-hoc workloads*. Pay attention to it!

Parallelism

SQL Server uses parallelism to speed up execution of complex queries by splitting the queries across multiple CPUs (workers). It improves user experience by completing queries faster. However, there ain't such a thing as a free lunch: parallelism always comes with overhead. With parallel execution plans, SQL Server needs to do additional work, splitting and merging the data across multiple workers and managing their execution.

Assume that a query finishes in 1 second, with the serial execution plan using the same 1,000ms of worker time. The same query may complete in 300ms with a parallel four-CPU plan, consuming 1,050ms of worker time in total. Managing parallelism requires SQL Server to perform extra work, and cumulative CPU time will always be higher than in the serial plan.

That overhead may impact throughput in busy OLTP systems. Faster execution time for a single query does not matter much, since there are many other queries waiting for CPUs in the queue. The overhead of parallelism forces them to wait longer for a CPU to become available. Although parallelism is good in complex reporting and analytical workloads, it can become a problem in OLTP systems, especially when the server operates under high CPU loads.

Unfortunately, it is extremely hard to find a system that does not mix both workloads. Even when you implement a dedicated data warehouse and operational data store (ODS), there will still be some reports and complex queries running in the source OLTP systems. Ideally, you want to separate those workloads, running them with different parallelism settings.

To make matters worse, SQL Server's default parallelism configuration is far from optimal. It allows SQL Server to utilize all CPUs in parallel

execution plans (MAXDOP=0 setting) and generate parallel execution plans when the cost of queries is equal to or greater than 5 (in technical terms, when the *cost threshold for parallelism*, or CTFP, is 5 or more). The meaning of *cost* is synthetic: it does not represent anything meaningful and is used as a baseline metric during query optimization. Nevertheless, the value of 5 is extremely low nowadays as amounts of data grow; this value allows parallel execution plans for many queries.

Parallelism presents itself in the system with CXPACKET, CXCONSUMER and EXCHANGE waits. It is very important to remember, however, that parallelism is not the root cause but a *symptom* of the issue. A high percentage of parallelism waits merely indicates a large number of expensive queries, which could be completely normal for a given reporting workload. In OLTP systems, on the other hand, such a figure usually means that queries are not properly optimized (optimized queries would have a lower cost).

NOTE

You can see the cost of an individual statement by examining the property of the root operator in the execution plan.

When I see substantial parallelism waits in OLTP systems, I adjust the parallelism settings and continue troubleshooting and query tuning. Optimized queries have a lower cost and therefore reduce parallelism. In some cases, I even filter out parallelism waits from the wait statistics output, to get more a detailed picture of other waits.

There are several approaches to tuning parallelism settings. In OLTP systems I start by setting MAXDOP to one-fourth of the number of available CPUs. If the server has a large number of CPUs or handles lots of OLTP requests, I may decrease the number to one-eighth or even lower. In data warehouse systems I might use half of the available CPUs instead.

More importantly, I increase the CTFP. I often start with a CTFP of 50, but you can examine the cost of the queries to analyze if other thresholds would

work better. You can run the code from Listing 4-X, uncommenting the [Query Cost] column to see the cost of cached execution plans.

After the change is done, I monitor CPU load, percent of signal waits, and parallelism waits and adjust the settings. One goal for these adjustments is finding the right CFTP value, which will allow SQL Server to separate different workloads and reduce or even prevent parallelism in OLTP queries.

There are other, more granular options to control parallelism. For example, you can separate OLTP and reporting workloads with Resource Governor and set different MAXDOP options for different workload groups. You could also consider setting MAXDOP to 1 in OLTP systems, enabling parallelism for reporting queries with a MAXDOP query hint. Either of those options would require you to monitor the system constantly and work closely with development teams.

Whatever you do, do *not* set MAXDOP to 1 for all system workloads. This just hides the problem. Remember that parallelism is normal – you just need to make sure it is used legitimately.

Summary

Issues with high CPU load are common nowadays, as fast disk subsystems and large amount of memory hide the impact of non-optimized queries. Reducing CPU load often becomes the goal of the performance tuning process.

Nevertheless, non-optimized queries are still a major factor in increasing CPU load on the server. The more data SQL Server needs to scan, the more CPU resources it uses. General query optimization helps to reduce that.

You also learned in this chapter about the overhead of query compilation. In systems with heavy ad-hoc workloads, query compilation may lead to very significant CPU usage. Query parameterization in the code is the best option to address the issue. Alternately, consider enabling forced parameterization for some queries or at the database level.

Unfortunately, parameterization may lead to issues with parameter-sensitive plans, where SQL Server compiles and caches plans for atypical parameter values. Those plans may be highly inefficient for other combinations of parameters. In many cases, disabling parameter sniffing improves the situation.

Pay attention to the amount of parallelism in your system. Parallelism is completely normal for reporting and analytical workloads; however, it is not desirable in OLTP systems, because parallelism management always adds overhead.

Remember that parallelism indicates the existence of expensive queries. You need to optimize them, instead of disabling parallelism and hiding the problem. Nevertheless, SQL Server's default parallelism settings are suboptimal and need to be tuned.

In the next chapter, you'll learn how to troubleshoot memory-related issues in SQL Server.

Troubleshooting Checklist

- Analyze and reduce CPU load from the processes outside of SQL Server.
- Detect and optimize the *worst offenders* – the queries that use the most worker time.
- Detect and optimize the most resource-intensive stored procedures and T-SQL modules.
- Review the impact of compilations. Plan cache metrics (which we'll discuss in the next chapter) may be useful in cross-checking the data.
- Parameterize critical queries. In the most severe cases of heavy, non-parameterized ad-hoc workloads, consider enabling forced parameterization and, potentially, disabling parameter sniffing.
- Tune your parallelism settings.

Chapter 7. Troubleshooting Memory Issues

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 7th chapter of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at dmitri@aboutsqlserver.com.

SQL Server is a memory-intensive application: it can consume hundreds of gigabytes or even terabytes of memory. This is completely normal and often a good thing – using this much memory reduces the need for physical I/O and recompilations, improving server performance.

In this chapter, I will discuss how SQL Server works with memory. I will start with an overview of how SQL Server uses memory and give you a few tips on memory configuration. Next, I’ll discuss the memory allocation process and show you how to analyze the memory usage of internal SQL Server components. Then I’ll talk about query memory grants and the ways to troubleshoot extensive query memory usage. Finally, I’ll briefly discuss memory management and potential issues in In-Memory OLTP implementation.

SQL Server Memory Usage and Configuration

By default, SQL Server tries to allocate as much memory as possible if memory is required for the operation. It does not allocate all memory at start time; the allocation occurs as needed, for example when SQL Server reads data pages to the buffer pool or stores compiled plans in the cache.

You can often see SQL Server consuming most of the OS memory. This is completely normal – when properly configured, SQL Server responds to OS requests and deallocates some process memory when needed. This condition is called *external memory pressure*. It usually occurs when the OS does not have enough memory for other applications. In small amounts, external memory pressure is not necessarily dangerous; however, deallocating a large amount of memory can significantly impact SQL Server's performance.

You can detect those events by setting up an alert on error 17890. The error generates the following message in the error log providing the information on how much memory had been trimmed - *A significant part of SQL Server process memory has been paged out. This may result in a performance degradation*. I will discuss how to avoid that later in that section.

Another condition, *internal memory pressure*, may occur when some SQL Server components consume large amounts of memory, impacting other components on the server. In most cases, SQL Server handles those cases gracefully, dynamically adjusting internal memory usage; however, it may lead to problems. I'll show you how to troubleshoot those later in the chapter.

There are several ways to monitor memory usage on the server. You can look at the following performance counters in the *Memory Manager* object. As a reminder, you can see them in the *Performance Monitor* utility or query the `sys.dm_os_performance_counters` view.

Target Server Memory (KB)

This performance counter indicates the ideal amount of memory SQL Server should consume. It depends on configuration settings, the total amount of memory available to the OS, and a few other factors.

Total Server Memory (KB)

This performance counter shows the amount of memory SQL Server currently uses.

In normal circumstances, the values of the two counters should stay very close. There are three cases, however, when Total Server Memory (KB) can become significantly lower than Target Server Memory (KB):

- When the server is ramping up shortly after startup; this is completely normal behavior
- When hardware is overallocated and SQL Server does not need all the available memory; this is also normal, but may indicate inefficient capacity planning
- During a memory pressure event, when SQL Server responds by trimming the memory; this condition requires further troubleshooting

You can also get memory metrics from the `sys.dm_os_sys_memory` and `sys.dm_os_process_memory` views, which provide information about OS and SQL Server memory usage, respectively. Listing 7-1 shows the code that uses them.

Listing 7-1. Analyzing OS and SQL Server memory usage

```
SELECT total_physical_memory_kb / 1024 AS [Physical Memory (MB)]
, available_physical_memory_kb / 1024 AS [Available Memory (MB)]
, total_page_file_kb / 1024 AS [Page File Commit Limit (MB)]
, available_page_file_kb / 1024 AS [Available Page File (MB)]
, (total_page_file_kb - total_physical_memory_kb) / 1024
AS [Physical Page File Size (MB)]
, system_cache_kb / 1024 AS [System Cache (MB)]
```

```

/* Values: LOW/HIGH/STEADY */
,system_memory_state_desc AS [System Memory State]
FROM sys.dm_os_sys_memory WITH (NOLOCK);

SELECT
    physical_memory_in_use_kb / 1024
    AS [SQL Server Memory Usage (MB)]
,locked_page_allocations_kb / 1024
    AS [SQL Server Locked Pages Allocation (MB)]
,large_page_allocations_kb / 1024
    AS [SQL Server Large Pages Allocation (MB)]
,memory_utilization_percentage
,available_commit_limit_kb
,process_physical_memory_low /* May indicate memory pressure */
,process_virtual_memory_low
FROM sys.dm_os_process_memory WITH (NOLOCK);

```

You can get historical information about Target and Total Server Memory from a `system_health` xEvent session. You can also see it in the `sp_server_diagnostics_component_result` event on the target (the partial output is shown in Listing 7-2). This information may be useful when you are troubleshooting unexplained performance issues and need to check whether the server experienced memory pressure during the time the problem occurred.

This data is also captured by a hidden xEvent session and stored in XEL files in SQL Server Log folder. The names of those files consist of the server and instance name followed by an SQLDIAG string.

Listing 7-2. `sp_server_diagnostics_component_result` event in `system_health` session (partial)

```

<resource lastNotification="RESOURCE_MEM_STEADY"
outOfMemoryExceptions="0" isAnyPoolOutOfMemory="0"
processOutOfMemoryPeriod="0">
  <memoryReport name="Process/System Counts" unit="Value">
    <entry description="Available Physical Memory"
value="65669554176" />
    <entry description="Available Virtual Memory"
value="138792447782912" />
    <entry description="Available Paging File" value="67695706112"
/>
  <..>

```

```
</memoryReport>
<memoryReport name="Memory Manager" unit="KB">
<entry description="Locked Pages Allocated" value="641593188" />
<entry description="Large Pages Allocated" value="3248128" />
<entry description="Target Committed" value="653261832" />
<entry description="Current Committed" value="653263320" />
    <..>
</memoryReport>
</resource>
```

Let's look at several options you can use to configure SQL Server memory.

Configuring SQL Server Memory

There are two well-known configuration settings that control SQL Server memory usage: Maximum and Minimum Server Memory.

Maximum Server Memory

The maximum amount of memory SQL Server can allocate. There are some cases when it can allocate memory beyond this amount; if so, it will detect that condition and deallocate excess memory.

Minimum Server Memory

The minimum amount of memory reserved for an SQL Server instance. SQL Server does not pre-allocate memory to match Minimum Server Memory value on startup. However, SQL Server would not deallocate the memory below it once the threshold is reached.

The default settings allow SQL Server to allocate all available memory without reserving any memory for the instance. This behavior may be sufficient in many systems with low or even mid-size loads. You may benefit, however, from tuning them (especially Maximum Server Memory) in your environment. Keep in mind that incorrect memory configuration may harm your system more than if you just keep default values.

Setting Maximum Server Memory usually requires some tuning. You can start with the base value and then adjust it by monitoring available physical memory on the server, using the `available_physical_memory_kb` column in the `sys.dm_os_sys_memory` view or the `Memory\Available MBytes` performance counter. Keep at least 512MB of memory reserved on the small servers, and 1GB or more reserved on servers with 128GB of RAM or more.

You can calculate the base value to start with the following formula:

$$\text{Total_Physical_Memory} - (4\text{GB} + 1\text{GB} * (\text{Total_Physical_Memory} - 16\text{GB}) / 8) - \text{Memory_For_Other_Apps}.$$

In case you are using an older version of SQL Server than 2012, you need to reserve additional memory, since the Maximum Server Memory setting in those versions controls memory usage of the buffer pool only.

NOTE

I am not covering memory configuration in the 32-bit version of SQL Server. If you are still using that, it's time to upgrade!

You need to properly estimate memory usage for other applications and reserve some memory for them. Those applications make memory management more complicated and can also impact system performance. I strongly recommend using dedicated SQL Servers in mission critical systems and not running any applications (including SSRS and SSAS) there.

Setting Maximum Server Memory does not prevent SQL Server from responding to memory pressure. In some extreme cases, Windows can even page some of SQL Server's physical memory to a page file. You can prevent this by granting SQL Server *Lock Pages In Memory (LPIM)* permission in Group Policy.

Using LPIM may help to improve SQL Server's responsiveness in the event of extreme external memory pressure. Be very careful with this setting,

though, especially in non-dedicated environments. It requires you to configure Maximum Server Memory properly and may lead to OS stability issues and even crashes if you over-allocate it.

The overhead of memory management on large servers can also increase shutdown time, which, in turn, can impact failover duration in SQL Server Failover Cluster. You can improve it by enabling *Large Page Allocations* with trace flag T834. In this mode, SQL Server allocates memory in the larger chunks, which speeds up the process and reduces memory management overhead. This setting requires you to enable LPIM and forces SQL Server to pre-allocate all memory up to Maximum Server Memory value on startup. This may increase SQL Server startup time, especially on servers with a large amount of memory.

Test Large Page Allocations carefully before you enable it. You are unlikely to benefit from it unless the server has at least 384GB of RAM. Do not enable it in non-dedicated environments, nor in environments that use columnstore indexes. Unfortunately, those two technologies do not work well with each other. That problem has been partially addressed in SQL Server 2019, where you can utilize some of the Large Page Allocations features in environments with columnstore indexes. You need to use another trace flag T876 instead of T834 to enable this. Nevertheless, carefully test the system before switching it on.

How Much Memory Is Enough?

I discussed hardware in Chapter 1, but I'd like to repeat a few things here. Memory is the key resource in SQL Server. *Adding more memory to the servers is often the fastest and cheapest way to improve system performance.*

There are no limitations on how much memory SQL Server can utilize with Enterprise Edition. Add as much memory as your server can support and use the fastest memory possible. Pay attention to the size of your active data – there's no need to build a server with terabytes of RAM for a 100GB

database, but keep future growth in mind and add some memory to support it.

In Standard Edition, the buffer pool size is limited to 128GB, but you'll need additional memory for other SQL Server components and the OS. I recommend provisioning the servers with 192GB of RAM to be on the safe side.

You will need even more memory if you are using columnstore indexes or In-Memory OLTP. The former uses an additional 32GB per Standard Edition instance to store segment data. The latter uses up to 32GB of RAM *per database*.

You can improve memory utilization by decreasing internal index fragmentation (the `avg_page_space_used_in_percent` column in the `sys.dm_db_index_physical_stats` view) and/or applying data compression. That will reduce the number of data pages and allow SQL Server to cache more data in the buffer pool.

Fortunately, memory is cheap nowadays – benefit from it!

Memory Allocations

As I mentioned in Chapter 2, almost all memory allocations in SQL Server are done through SQLOS. Extended stored procedures and linked server providers can perform memory allocation outside of SQLOS and, therefore, would not be controlled by the Maximum Server Memory setting.

Internally, SQLOS partitions the memory into *memory nodes* based on the server's NUMA configuration – one memory node per NUMA node. Each memory node has a *memory allocator*, which uses various Windows and Linux API methods to allocate and deallocate the memory.

In earlier versions of SQL Server, memory allocation was carried out with single-page allocators for memory allocations less than 8 KB and multi-page allocators for any allocations greater than 8KB. Starting with SQL Server 2012, the memory allocators were consolidated into one allocator for

any size page. You can track memory usage and allocations per memory node with the `sys.dm_os_memory_nodes` [view](#).

Internally, allocations become *memory objects*. Each memory object stores an allocated memory along with its metadata (size, owner, etc.). You can analyze memory objects with the `sys.dm_os_memory_objects` [view](#); although, I rarely use it during troubleshooting.

Another key element of SQL Server memory architecture is called *memory clerks*. Each major component of SQL Server has its own memory clerk, which works as the proxy between component and memory allocator. When a component needs memory, it sends a request to the corresponding memory clerk, which in turn, gets the memory object from the memory allocator.

The last major component in SQL Server dynamic memory management is called the *memory broker*. Memory brokers supervise memory clerks by adjusting their memory usage based on available process memory, memory pressure, and other conditions. Memory brokers do not allocate or deallocate memory by themselves; however, they can send the signal to memory clerks to shrink or grow. You can look at the `sys.dm_os_memory_brokers` [view](#) to see memory broker state and amount of memory they allocate.

Figure 7-1 shows dependencies between memory management components in SQL Server.

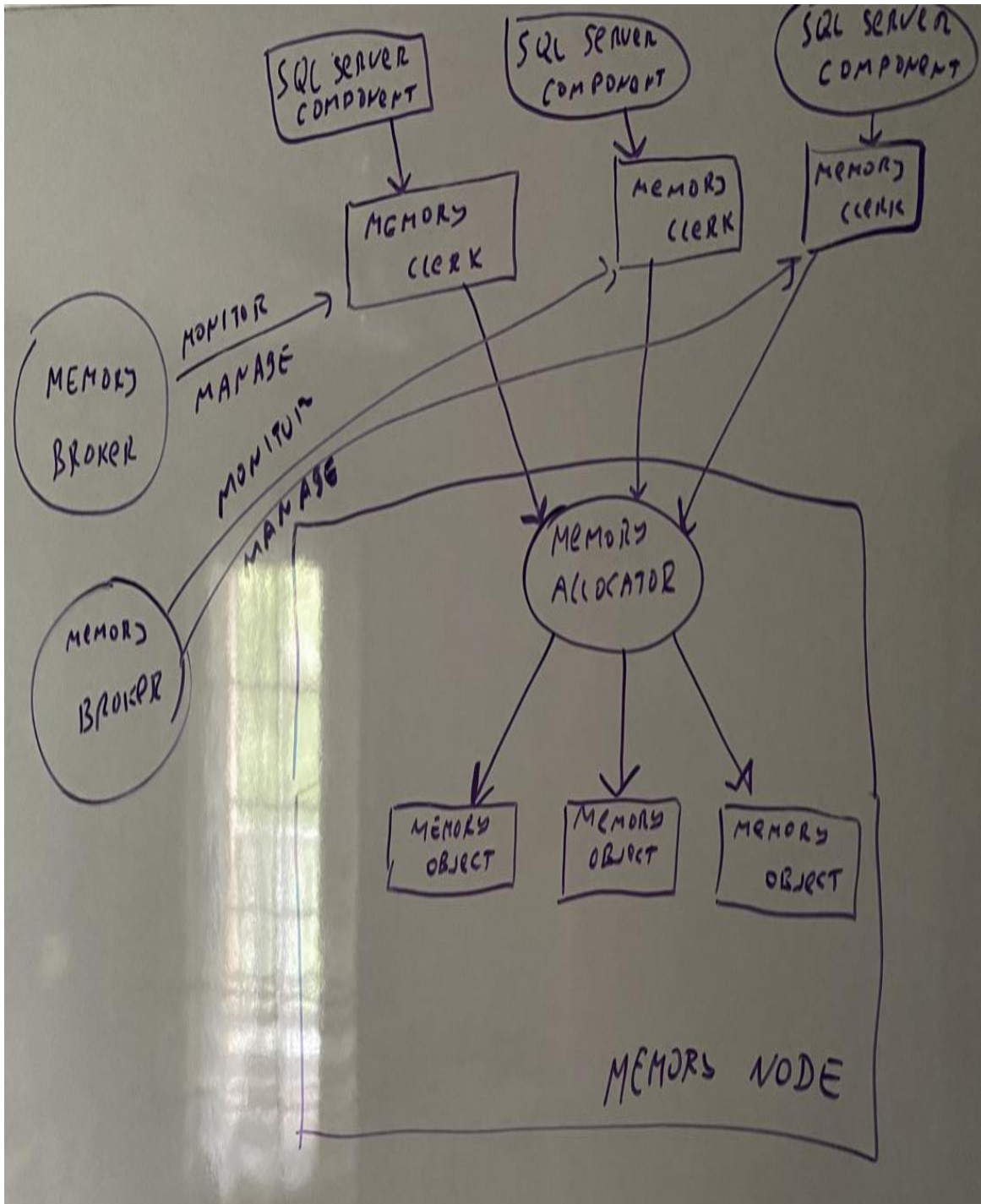


Figure 7-1. Memory management components

I usually start analyzing memory usage by looking at memory clerks using the `sys.dm_os_memory_clerks` [view](#). Listing 7-3 shows how to do that. The code will work in SQL Server 2012 and above. If you are using an older

version of SQL server, you should replace the single_page_kb column with the sum of pages_kb and multi_pages_kb columns.

Listing 7-3. Analyzing memory usage

```
SELECT TOP 15
  [type] AS [Memory Clerk]
, CONVERT(DECIMAL(16,3),SUM(pages_kb) / 1024.0) AS [Memory
Usage(MB)]
FROM sys.dm_os_memory_clerks WITH (NOLOCK)
GROUP BY [type]
ORDER BY sum(pages_kb) DESC
```

Figure 7-2 shows the output of the script from one of the production SQL Servers that use 640GB of RAM. The sys.dm_os_memory_clerks view provides you enough information to quickly evaluate memory usage and detect potential abnormalities.

	Memory Clerk	Memory Usage(MB)
1	MEMORYCLERK_SQLBUFFERPOOL	583425.398
2	CACHESTORE_OBJCP	10116.328
3	OBJECTSTORE_LOCK_MANAGER	7620.531
4	CACHESTORE_SQLCP	4388.898
5	USERSTORE_SCHEMAMGR	1542.930
6	OBJECTSTORE_XACT_CACHE	1083.602
7	MEMORYCLERK_SOSNODE	534.867
8	OBJECTSTORE_SERVICE_BROKER	423.273
9	MEMORYCLERK_SQLSTORENG	315.492
10	USERSTORE_TOKENPERM	275.320
11	MEMORYCLERK_SQLQERESERVATI...	218.430
12	CACHESTORE_PHDR	193.898
13	OBJECTSTORE_SNI_PACKET	179.516
14	MEMORYCLERK_XE	146.688

Figure 7-2. Memory usage on one of production servers

You can see a complete list of memory clerks in the [Microsoft Documentation](#).

Memory Clerks

Let's look at the most common memory clerks you may encounter.

MEMORYCLERK_SQLBUFFERPOOL

As you can guess by its name, MEMORYCLERK_SQLBUFFERPOOL controls memory allocation of the buffer pool. Usually, it is one of the largest memory consumers, especially when the server works with large databases.

There are no specific thresholds for *ideal* buffer pool size - it depends on your system. Large buffer pools are completely normal and mean that SQL Server just caches more data. There is nothing to worry about - when other SQL Server components need more memory, SQL Server trims the size of the buffer pool if there is no OS memory left to allocate.

Small buffer pools, on the other hand, may require some investigation. It does not necessarily present a problem – remember that SQL Server caches the *active* data the system is using. For example, even if you have a large multi-terabyte database, your buffer pool will be small if you work with a small subset of data.

You can check the *Page Life Expectancy* performance counter (Listing 3-3), which shows the number of reads in data files (Listing 3-1) and percentage of PAGEIOWAITS, to determine if the buffer pool is constantly flushing. If this is the case, especially in OLTP systems, you should analyze the memory usage of other clerks and potential memory pressure conditions. Obviously, you should detect and optimize inefficient queries – the less data you need to scan, the less data you'd bring to the buffer pool.

Remember that non-Enterprise editions of SQL Server have limitations on the maximum buffer pool size. For example, the Standard Edition is limited to 128GB of RAM in SQL Server 2016 and above, and to 64GB in older versions.

Listing 7-4 shows a query that displays your buffer pool memory usage on a per-database basis. This may help you analyze when the server uses

multiple databases. As a bonus, the code also returns the average time of the physical reads for data pages from disk.

Listing 7-4. Buffer pool usage on per-database basis

```
;WITH BufPoolStats
AS
(
  SELECT
    database_id
  ,COUNT_BIG(*) AS page_count
  ,CONVERT(DECIMAL(16,3),COUNT_BIG(*) * 8 / 1024.) AS size_mb
  ,AVG(read_microsec) AS avg_read_microsec
  FROM
    sys.dm_os_buffer_descriptors WITH (NOLOCK)
  GROUP BY
    database_id
)
SELECT
  DB_NAME(database_id) AS [DB]
  ,size_mb
  ,page_count
  ,avg_read_microsec
  ,CONVERT(DECIMAL(5,2), 100. * (size_mb / SUM(size_mb) OVER()))
  AS [Percent]
FROM
  BufPoolStats
ORDER BY
  size_mb DESC
OPTION (MAXDOP 1, RECOMPILE);
```

CACHESTORE_OBJCP, CACHESTORE_SQLCP, and CACHESTORE_PHDR Memory Clerks

CACHESTORE and USERSTORE memory clerks are, in a nutshell, the caches for different types of data. The CACHESTORE_OBJCP, CACHESTORE_SQLCP, and CACHESTORE_PHDR memory clerks store plan cache–related objects.

CACHESTORE_PHDR

This memory clerk caches internal objects used during query compilations.

CACHESTORE_OBJCP

This memory clerk stores compiled execution plans for stored procedures, functions, triggers, and other SQL objects.

CACHESTORE_SQLCP

This memory clerk stores compiled execution plans for ad-hoc queries, prepared statements, and server-side cursors.

CACHESTORE_PHDR clerks are most commonly used for compiling queries that reference complex views, large batches, and statements with a large number of constants in the IN clause. The objects cached by them are short-lived and are cached only during query compilation.

The *CACHESTORE_PHDR* clerk rarely consumes large amounts of memory. When you see high memory usage from this clerk, analyze the code and database schema to see if there are opportunities for refactoring. Don't do unnecessary refactoring though. For example, an application that passes a large amount of data to a stored procedure through table-valued parameters (TVP) may generate a large batch of individual insert statements to populate them. Although refactoring and switching to per-row processing may reduce memory usage during compilation, it could seriously impact system throughput as per-row processing will be slower than set-based operations with TVPs.

In most systems, the majority of plan cache memory will be allocated by *CACHESTORE_OBJCP* and *CACHESTORE_SQLCP* clerks. Memory consumption of *CACHESTORE_OBJCP* depends on the database schema and data tier architecture. Systems with a large number of actively used stored procedures and other T-SQL modules use more memory to store their execution plans. A *reasonably* high memory usage of *CACHESTORE_OBJCP* would not introduce any issues as long as it does not impact other components.

The situation is different with `CACHESTORE_SQLCP` clerks. Large memory consumption there usually indicates an excessive ad-hoc workload, which is CPU intensive and consumes plan cache memory.

You have already seen CPU overhead from ad-hoc queries in the previous chapter. Now, let me show you an example of memory overhead introduced by it. Listing 7-5 runs 1,000 simple ad-hoc queries and checks the plan cache state afterward using `sys.dm_exec_cached_plans` [view](#). This script clears the content of the cache with the `DBCC FREEPROCCACHE` command – do not run it on the production server!

Listing 7-5. Running 1,000 ad-hoc queries and examining plan cache content

```
DBCC FREEPROCCACHE
GO

DECLARE
    @SQL NVARCHAR(MAX)
    ,@I INT = 0

WHILE @I < 1000
BEGIN
    SELECT @SQL =
N'DECLARE @C INT;SELECT @C=object_id FROM sys.objects WHERE
object_id='
+ CONVERT(NVARCHAR(10),@I);
    EXEC(@SQL);
    SELECT @I += 1;
END;

SELECT
    p.usecounts, p.cacheobjtype, p.objtype, p.size_in_bytes, t.
[text]
FROM
    sys.dm_exec_cached_plans p WITH (NOLOCK)
    CROSS APPLY sys.dm_exec_sql_text(p.plan_handle) t
WHERE
    p.objtype = 'Adhoc'
ORDER BY
    p.objtype DESC
OPTION (RECOMPILE);

SELECT
```

```

CONVERT(DECIMAL(12,3),SUM(1. * p.size_in_bytes)/1024.) AS [Size
(KB) ]
FROM
  sys.dm_exec_cached_plans p WITH (NOLOCK)
WHERE
  p.objtype = 'Adhoc'
OPTION (RECOMPILE);

```

Figure 7-3 shows the content of the plan cache and its memory usage when I run a query with the *Optimize for ad-hoc workloads* configuration setting disabled. There are 1,000 plans cached, each of which uses 48KB of memory, or 48MB total. Keep in mind, I ran a very simple query. Large and complex queries will use significantly more memory for the execution plans.

	usecounts	cacheobjtype	objtype	size_in_bytes	text
1	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=999
2	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=998
3	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=997
4	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=996
5	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=995
6	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=994
7	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=993
8	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=992
9	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=991
10	1	Compiled Plan	Adhoc	49152	DECLARE @C INT;SELECT @C= t_id=990

	Size (KB)
1	48000.000

Figure 7-3. Plan cache content when *Optimize for ad-hoc workloads* is disabled

Figure 7-4 shows plan cache statistics after I repeat the test with the *Optimize for ad-hoc workloads* configuration setting enabled. I still have

1,000 plans in the cache; however, memory usage is significantly lower. SQL Server caches small 400-byte structures, called *compiled plan stubs*, instead of actual compiled plans. These structures are the placeholders that are used to keep track of which ad-hoc queries were executed. When the same query runs a second time, SQL Server replaces compiled plan stub with the actual compiled plan and reuses it going forward.

	usecounts	cacheobjtype	objtype	size_in_bytes	text
1	1	Compiled Plan Stub	Adhoc	408	DECLARE @C INT;SELECT @C=0 t_id=999
2	1	Compiled Plan Stub	Adhoc	408	DECLARE @C INT;SELECT @C=0 t_id=998
3	1	Compiled Plan Stub	Adhoc	408	DECLARE @C INT;SELECT @C=0 t_id=997
4	1	Compiled Plan Stub	Adhoc	408	DECLARE @C INT;SELECT @C=0 t_id=996
5	1	Compiled Plan Stub	Adhoc	408	DECLARE @C INT;SELECT @C=0 t_id=995
...					
Size (KB)					
1	398.438				

Figure 7-4. Plan cache content when *Optimize for ad-hoc workload* is enabled

As I already mentioned multiple times, the *Optimize for ad-hoc workloads* setting should be enabled in most systems. It will reduce plan cache memory consumption and improve system performance.

Finally, let's repeat the same test using parameterized query as shown in Listing 7-6.

Listing 7-6. Running parameterized query and examine plan cache content

```

DBCC FREEPROCCACHE
GO

DECLARE
    @SQL NVARCHAR(MAX),@I INT = 0

WHILE @I < 1000
BEGIN
    SELECT @SQL =
N'DECLARE @C INT;SELECT @C=object_id FROM sys.objects WHERE
object_id=@P';
    EXEC sp_executesql @SQL=@SQL,@Params=N'@P INT',@P = @I;
    SELECT @I += 1;
END;

SELECT
    p.usecounts, p.cacheobjtype, p.objtype, p.size_in_bytes, t.
[text]
FROM
    sys.dm_exec_cached_plans p WITH (NOLOCK)
    CROSS APPLY sys.dm_exec_sql_text(p.plan_handle) t
WHERE
    p.objtype = 'Adhoc'
ORDER BY
    p.objtype DESC
OPTION (RECOMPILE);

SELECT
    CONVERT(DECIMAL(12,3),SUM(1. * p.size_in_bytes)/1024.) AS [Size
(KB)]
FROM
    sys.dm_exec_cached_plans p WITH (NOLOCK)
WHERE
    p.objtype = 'Adhoc'
OPTION (RECOMPILE);

```

As you can see in Figure 7-5, there is a single execution plan cached. This reduces memory consumption and CPU load on the server.

	usecounts	cacheobjtype	objtype	size_in_bytes	text
1	1000	Compiled Plan	Prepared	49152	(@P INT)DECLARE @C INT;SELECT @C=r _id=@P

	Size (KB)
1	48.000

Figure 7-5. Plan cache content with parameterized query

The size and amount of memory consumed by plan cache depends on the system workload, complexity of the execution plans, and data-tier design. It is not uncommon to see memory clerks using gigabytes or even tens of gigabytes of memory. Nevertheless, pay attention to it and analyze the content of the plan cache when you see the large numbers.

This is especially important if you see high memory usage by the `CACHESTORE_SQLCP` clerk. More often than not, it is the sign of a heavy ad-hoc workload and large number of single-use plans. As you saw earlier in the chapter, this impacts system performance.

There is another caveat regarding the ad-hoc workload. By default, SQL Server can cache about 160,000 objects in the plan cache. Reaching this limit may lead to additional CPU contention.

You may increase this number to about 640,000 objects by enabling trace flag T174. This may reduce CPU load in systems with very heavy ad-hoc workloads; however, it may also negatively impact performance as SQL Server would use more memory for plan cache. Test it before deploying to production.

Listing 7-7 shows three queries you can use for plan cache analysis and troubleshooting. The first gives you the information about cached plan-cache objects grouped by their types. The data is collected from all plan-cache memory clerks.

The second query gives you the number of single-used plans along with the memory they consume in plan cache. Large numbers, and especially large memory consumption, may warrant further investigation. The third query helps you detect the most memory-intensive single-used plans.

Listing 7-7. Analyzing plan cache

```
-- Number of cached object and their memory usage grouped by type
SELECT
    cacheobjtype
    ,objtype
    ,COUNT(*) AS [Count]
    ,CONVERT(DECIMAL(12,3),SUM(1.*size_in_bytes)/1024./1024.)
    AS [Size (MB)]
FROM
    sys.dm_exec_cached_plans WITH (NOLOCK)
GROUP BY
    cacheobjtype, objtype
ORDER BY
    [Size (MB)] DESC
OPTION (RECOMPILE);
-- Statistics on single-used execution plans
SELECT
    COUNT(*) AS [Single-used plan count]
    ,CONVERT(DECIMAL(10,3),SUM(cp.size_in_bytes)/1024./1024.)
    AS [Size (MB)]
FROM
    sys.dm_exec_cached_plans cp WITH (NOLOCK)
WHERE
    cp.objtype in (N'Adhoc', N'Prepared') AND
    cp.usecounts = 1
OPTION (RECOMPILE);
-- 25 most memory-intensive single-used plans
SELECT TOP 25
    DB_NAME(t.dbid) as [DB]
    ,cp.usecounts
    ,cp.plan_handle
    ,t.[text]
    ,cp.objtype
    ,cp.size_in_bytes
    ,CONVERT(DECIMAL(12,3),cp.size_in_bytes/1024.) as [Size (KB)]
FROM
    sys.dm_exec_cached_plans cp WITH (NOLOCK)
    CROSS APPLY sys.dm_exec_sql_text(cp.plan_handle) t
WHERE
    cp.cacheobjtype = N'Compiled Plan'
```

```
AND cp.objtype in (N'Adhoc', N'Prepared')
AND cp.usecounts = 1
ORDER BY
  cp.size_in_bytes DESC
OPTION (RECOMPILE);
```

You can remove individual execution plans from the cache by calling the `DBCC FREEPROCCACHE` statement and providing the plan handle as the parameter. This can be useful if you want to remove specific large single-use plans without affecting others. You can also do that to remove regressed parameter-sensitive plans affected by parameter sniffing.

You can also clear all plans stored by `CACHESTORE_SQLCP` memory clerk with the `DBCC FREESYSTEMCACHE('SQL Plans') WITH MARK_IN_USE_FOR_REMOVAL` command. This removes all ad-hoc and prepared plans from the cache regardless of how often they were re-used keeping plans from stored procedures and other T-SQL modules intact.

Obviously, address the root-cause of the issue when it is possible. Enable *Optimize for ad-hoc workloads* and parameterize the queries as I discussed in the previous chapter.

Finally, pay attention when the plan cache is small and does not use much memory. It could be completely normal when the system uses parameterized queries and stored procedures and has little ad-hoc activity. On the other hand, it may be a sign of memory pressure when the plan cache is constantly shrinking. Analyze compilation and recompilation performance counters when this is the case.

OBJECTSTORE_LOCK_MANAGER Memory Clerk

The `OBJECTSTORE_LOCK_MANAGER` clerk stores lock structures that SQL Server uses to support concurrency. High memory usage of this memory clerk indicates a large number of active locks being currently held. Usually, you'd try to keep the number of active locks as small as possible; however, there are the cases when you may need to acquire many row-level locks and prevent their escalation to object level.

I will talk about troubleshooting the SQL Server concurrency issues in the next chapter.

MEMORYCLERK_SQLQERESERVATIONS Memory Clerk

The MEMORYCLERK_SQLQERESERVATIONS clerk manages *memory grants* - memory allocated to queries during their execution. It is common to see this clerk in the list of top memory consumers, especially in systems with data warehouse and reporting workloads.

I will talk about memory management during query execution and troubleshooting of extensive memory grants later in this chapter.

USERSTORE_TOKENPERM Memory Clerk

The USERSTORE_TOKENPERM clerk provides memory for a security token store that is used to track user permissions and various other security objects. A large token store may introduce performance issues by increasing CPU load and triggering internal memory pressure by stealing the memory from other SQL Server components.

Unfortunately, this is a hard problem to address. To make it even worse, there are many known issues related to the token store in SQL Server. Some of them may be triggered by extensive usage of application roles and heavy ad-hoc workload.

Keep an eye on the USERSTORE_TOKENPERM memory clerk. If the size is more than a few gigabytes, it may indicate that you have a problem, especially if it continues to grow. In this case, apply the latest service pack and/or cumulative update. If that does not help, consider opening a support case with Microsoft CSS.

As a temporary solution, you can clear the token store by using the DBCC FREESYSTEMCACHE ('TokenAndPermUserStore') command. In some cases, you may need to clear the token store on a regular basis by using the SQL Server Agent job, for example, until you have a permanent solution.

MEMORYCLERK_SQLCONNECTIONPOOL Memory Clerk

The MEMORYCLERK_SQLCONNECTIONPOOL clerk provides memory for connection-specific objects that the client needs the server to maintain. The most common case is to use prepared statement handles that are generated by SP_PREPEXECRPC stored procedure calls during some RPC calls.

Memory consumption by this clerk is rarely a problem. It may grow, however, when client applications do not properly discard prepared handles while maintaining open connections to the database. You may have to restart the application and, in some cases, restart SQL Server to clear up the memory.

Obviously, it is better to address the root-cause of the issue in the application code and close the handles after execution.

MEMORYCLERK_SQLCLR, MEMORYCLERK_SQLCLRASSEMBLY, and MEMORYCLERK_SQLEXTENSIBILITY Memory Clerks

MEMORYCLERK_SQLCLR, MEMORYCLERK_SQLCLRASSEMBLY and MEMORYCLERK_SQLEXTENSIBILITY clerks are used for memory allocations in CLR and other supported language extensions (Java, R, and Python). All external languages will manage the memory and perform garbage collection automatically; however, it is possible to write code that consumes a large amount of memory during the execution. Think about processing large files or working with large documents, for example.

When you see high memory usage in those memory clerks, analyze the usage of CLR and/or external languages. You may need to work with developers and refactor or even migrate some code to the application servers.

MEMORYCLERK_XTP memory clerk

The MEMORYCLERK_XTP clerk controls memory allocations for In-Memory OLTP technology. When you see high memory usage in this clerk, you need to look at memory-optimized tables memory consumption along with a few other things I will discuss later in that chapter.

The high memory usage of In-Memory OLTP may be completely legitimate when memory-optimized tables store large amounts of data. More importantly, that memory is not going to be released in case of memory pressure when In-memory OLTP is in use. You should take this into consideration when planning hardware capacity for servers (secondary AG replicas, DR servers, lower environments, etc.). The database will not start up and/or data in memory-optimized tables will become read-only if the server does not have enough memory to accommodate In-Memory OLTP data.

Wrapping Up

It is impossible and unnecessary to cover all memory clerks in this section. Fortunately, Microsoft provides a comprehensive list of all memory clerks in their [documentation](#). Even though it does not include troubleshooting guidelines, it will help you understand what each memory clerk is responsible for and will point you in the right direction for troubleshooting.

As a word of caution - don't have tunnel vision and jump to immediate conclusions strictly based on clerks' memory usage. With very few exceptions, there are no specific guidelines on how much memory each memory clerk should use. You need to look at memory usage holistically and understand how different components impact each other.

For example, having a lot of memory reserved for query execution (MEMORYCLERK_SQLQRESERVATIONS) may be completely normal if it does not impact buffer pool, plan cache and other SQL Server components. On the other hand, it may be dangerous when it introduces internal memory pressure. Remember to use memory clerk information together with other metrics in your analysis.

DBCC MEMORYSTATUS command

If you have worked with SQL Server long enough, you are probably familiar with the DBCC MEMORYSTATUS command, which provides a snapshot of the current memory usage in SQL Server. Personally, I consider

this command to be a mixed bag. Although it consolidates all memory usage information into one multi-result-set output, it has a limited ability to filter and aggregate data. More often than not, I collect the metrics using data management views instead.

I am not going to cover the DBCC MEMORYSTATUS command in this book, but I'd encourage you to run it and see if you like how it presents the information and if you find it easy to interpret. It is just a different and consolidated projection of the data I've already discussed in this chapter.

Query Execution and Memory Grants

Every query in SQL Server needs to use some memory in order to run. This memory is called *memory grant* and it is assigned to the query before it starts executing. The query will not start until the memory is available and may eventually time out if it cannot start.

The size of the memory grant depends on the operators in the execution plan, cardinality estimations, degree of parallelism, execution mode (row- or batch-mode), and a few other factors. For example, *Sort* or *Hash* operators need additional memory to support internal structures. They also benefit from extra memory to store all or a subset of the data in memory rather than spilling it to tempdb.

You can see memory grant information in the query execution plan. Figure 7-6 shows the information available in the query plan window in the SELECT (top plan operator) pop up and SSMS properties window. You can also get the same data from the XML representation of the execution plan. The numbers are in KB.

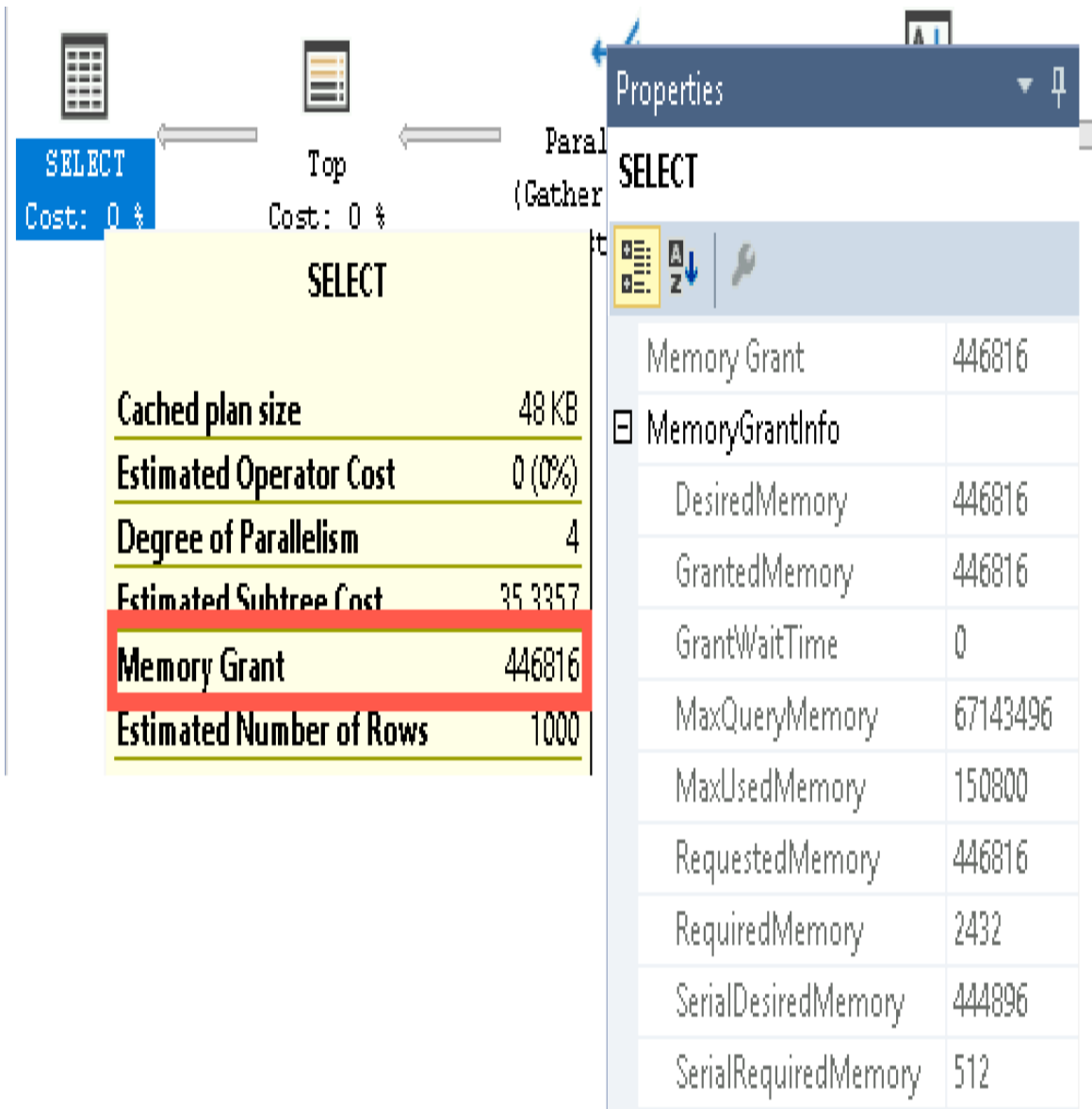


Figure 7-6. Memory grant information in SSMS

Let's look at memory grant properties.

Required Memory

Absolute minimum amount of memory required for a query to execute.

A query will not start until this memory is available.

Serial Required Memory

Absolute minimum amount of memory required for query to execute in case of serial execution plan. This will match Required Memory if a query runs serially.

Desired Memory

Amount of memory the query wants in the perfect scenario. For example, if a query plan contains a Sort operator, Desired Memory may include enough memory to sort all data in memory based on cardinality estimations.

Serial Desired Memory

Desired amount of memory if queries are executed serially.

Requested Memory

Amount of memory query requested from SQL Server when asking for a memory grant.

Granted Memory

Memory granted to a query.

Max Used Memory

Amount of memory used by a query during the execution.

Max Query Memory

Maximum possible size of a memory grant for queries.

Grant Wait Time

Amount of time that a query waited for a memory grant.

The size of a memory grant is calculated at query optimization time and cached with the execution plan. Further executions of the queries will reuse the same grant sizes; however, SQL Server 2017 and above can recalculate the size of the memory grant based on actual memory usage from previous execution. I will cover this behavior later in this chapter.

As I mentioned earlier, memory grants are managed by a `MEMORYCLERK_SQLQERESERVATIONS` clerk, which uses a thread synchronization object called *resource semaphore* to allocate the memory.

When memory cannot be allocated, the resource semaphore puts queries into the wait queue, generating `RESOURCE_SEMAPHORE` waits. Internally, the resource semaphore uses two wait queues and ranks queries based on their memory grant size and query cost. The first queue, called the *small-query resource semaphore*, stores queries that require less than 5MB and costs less than 3 cost units. The second queue stores all other queries.

The resource semaphore processes requests on a first-come-first-serve basis. It favors the small-query resource queue over the other one, and reduces the waiting time for the small queries that do not require a large amount of memory.

You need to monitor the situation when queries are waiting for memory grants to execute. This is unhealthy and requires investigation. The presence of `RESOURCE_SEMAPHORE` in noticeable amounts indicates a problem.

There are several performance counters in the Memory Management object that you can use for the monitoring and troubleshooting.

Memory Grants Pending

Shows the number of memory grant requests that are currently pending.

Ideally, this counter should show 0 all the time, which indicates that there are no queries waiting for memory grants.

Memory Grants Outstanding

Shows the number of currently running queries with fulfilled memory grant requests. Large values in this counter indicates a memory-intensive workload (usually queries with Sort and Hash operators). Although this can be normal in data warehouses, you need to investigate this condition in OLTP systems.

Maximum Workspace Memory

Provides total amount of workspace memory in KB.

Granted Workspace Memory

Indicates how much workspace memory in KB is currently in use.

You can get more detailed information about total and granted sizes of workspace memory with the `sys.dm_exec_query_resource_semaphores` view. It provides you with statistics for both resource semaphore queues including memory information, number of queries in the waiting queue, and a few other metrics.

You can obtain information about pending and outstanding memory grants from `sys.exec_query_memory_grants` **view** as shown in Listing 7-8. The `grant_time` column shows the time when the grant was fulfilled. A NULL value in this column indicates that grant is pending.

Listing 7-8. Memory grant information

```
SELECT mg.session_id
, t.text AS [sql]
, qp.query_plan AS [plan]
, mg.is_small /* Resource Semaphore Queue information */
, mg.dop
, mg.query_cost
, mg.request_time
, mg.grant_time
, mg.wait_time_ms
, mg.required_memory_kb
, mg.requested_memory_kb
```



```

,mg.granted_memory_kb
,mg.used_memory_kb
,mg.max_used_memory_kb
,mg.ideal_memory_kb
FROM
  sys.dm_exec_query_memory_grants mg WITH (NOLOCK)
  CROSS APPLY sys.dm_exec_sql_text(mg.sql_handle) t
  CROSS APPLY sys.dm_exec_query_plan(mg.plan_handle) qp
--WHERE -- Uncomment to see only pending memory grants
-- mg.grant_time IS NULL
ORDER BY
  mg.requested_memory_kb DESC
OPTION (RECOMPILE, MAXDOP 1);

```

The `sys.dm_exec_query_memory_grants` view shows you the current status of memory grants. However, you may need to look at historical memory usage and detect the most memory-intensive queries. You can use the methods discussed in Chapter 4 to analyze memory grants data during troubleshooting.

The simplest approach, perhaps, is using query execution statistics and `sys.dm_exec_query_stats` view. It includes several columns to track query memory grants and memory usage. You can use the code from Listing 4-1 sorting data by `total_grant_kb` and/or `[avg grant kb]` columns and detect most memory intensive queries. Alternatively, you can get the data from the Query Store when it is enabled.

Finally, you can track memory grants requests in run-time with Extended Events. In SQL Server 2014 SP2 and above, you can use lightweight query profiling, which I discussed in Chapter 5. In older versions, you can use `query_memory_grant_usage` event instead.

Optimizing Memory Intensive Queries

When you detect a lot of memory-intensive queries, analyze their execution plans. In most cases, memory grants are driven by the memory usage of *Hash and Sort* operators. If possible, capture actual execution plans, as they will provide you with the actual number of rows processed by operators along with their memory usage.

Unfortunately, there is no silver bullet that can magically optimize queries and reduce their memory consumption. However, there are a few things you can do instead. First, analyze if there is an opportunity for better indexing. This may eliminate unnecessary sorts and, in some cases, can change hash join to loop join, which does not use much memory.

Let's look at an example and create a table as shown in Listing 7-9.

Listing 7-9. Optimizing memory intensive queries: Table creation

```
CREATE TABLE dbo.Orders
(
  OrderID INT NOT NULL,
  OrderDate DATETIME2(0) NOT NULL,
  Placeholder CHAR(8000) NULL,
  CONSTRAINT PK_Orders PRIMARY KEY CLUSTERED(OrderID)
);

;WITH N1(C) AS (SELECT 0 UNION ALL SELECT 0) -- 2 ROWS
,N2(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N1 AS T2) -- 4 ROWS
,N3(C) AS (SELECT 0 FROM N2 AS T1 CROSS JOIN N2 AS T2) -- 16 ROWS
,N4(C) AS (SELECT 0 FROM N3 AS T1 CROSS JOIN N3 AS T2) -- 256
ROWS
,N5(C) AS (SELECT 0 FROM N4 AS T1 CROSS JOIN N4 AS T2) -- 65,536
ROWS
,IDs(ID) AS (SELECT ROW_NUMBER() OVER (ORDER BY (SELECT NULL))
FROM N5)
INSERT INTO dbo.Orders(OrderID, OrderDate)
  SELECT ID, DATEADD(day, ID % 365, '2021-01-01')
FROM IDs;
```

Next, let's run a query that returns the 200 most recent orders, as shown in Listing 7-10.

Listing 7-10. Optimizing memory intensive queries: Test query 1

```
SELECT TOP 200 OrderID, OrderDate, Placeholder
FROM dbo.Orders
ORDER BY OrderDate DESC
```

The execution plan and memory grant metrics for the query are shown in Figure 7-7. The query uses a 630MB memory grant, which is driven by a

Sort TOP N operator. This operator caches rows in memory before performing sorting.

It is worth noting that the Sort TOP N operator has internal and undocumented optimizations for the cases when the TOP condition does not exceed 100 rows. In that mode, the operator uses very little memory during execution.

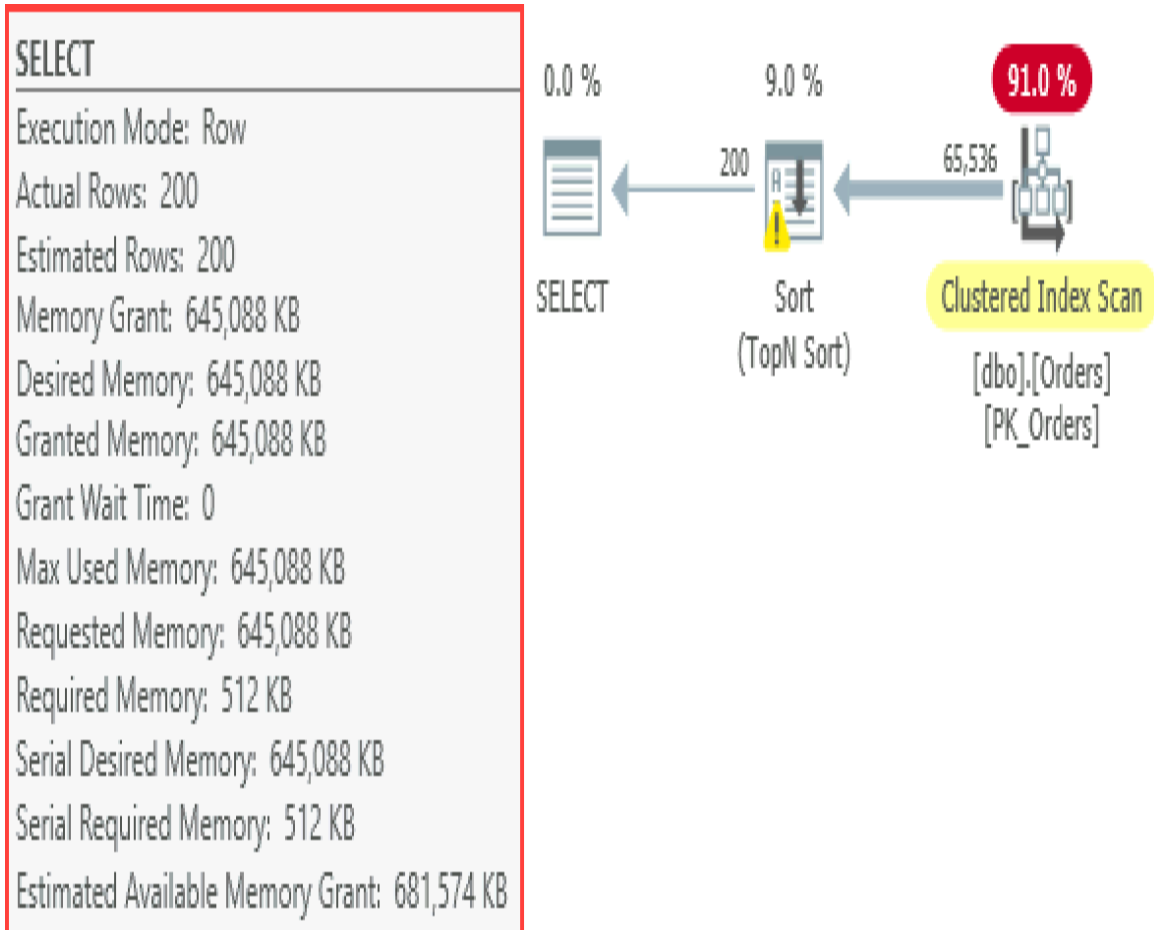


Figure 7-7. Optimizing memory intensive queries: Execution plan of the query

As I mentioned, you can often eliminate sorting with the proper indexing. For example, let's create the index with `CREATE INDEX IDX_Orders_OrderDate ON dbo.Orders(OrderDate)` command and run the query from Listing 7-10 again.

Figure 7-8 shows the new execution plan. The *Sort* operator is no longer required, and query does not require a memory grant to support it.

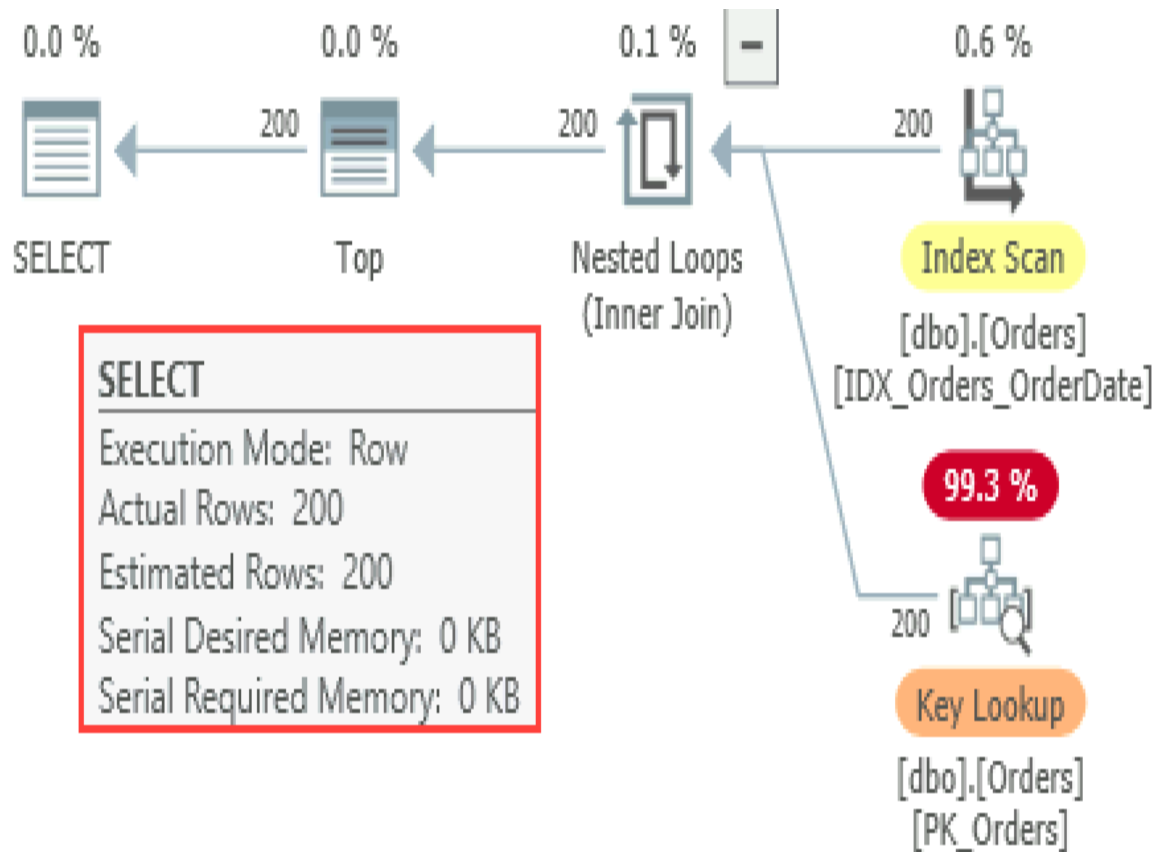


Figure 7-8. Optimizing memory intensive queries: Execution plan of the query after index has been created

The size of the requested memory grant depends on the estimated number of rows and the size of the rows the operator needs to process. For example, if Query Optimizer expects to sort 10,000 rows of 100 bytes each, it would need about 10MB to accommodate the data in memory. Both cardinality estimation errors and row size estimation errors may lead to incorrect memory grants.

Let's look at the impact of cardinality estimation errors first and run another query in Listing 7-11.

Listing 7-11. Optimizing memory intensive queries: Test query 2

```
SELECT TOP 200 OrderID, OrderDate, Placeholder
FROM dbo.Orders
WHERE OrderDate BETWEEN '2021-07-01' AND '2021-08-01'
ORDER BY Placeholder;
```

Figure 7-9 shows the partial execution plan for a query along with the memory grant statistics. In this example, I just created the index and, therefore, the statistics are up to date.

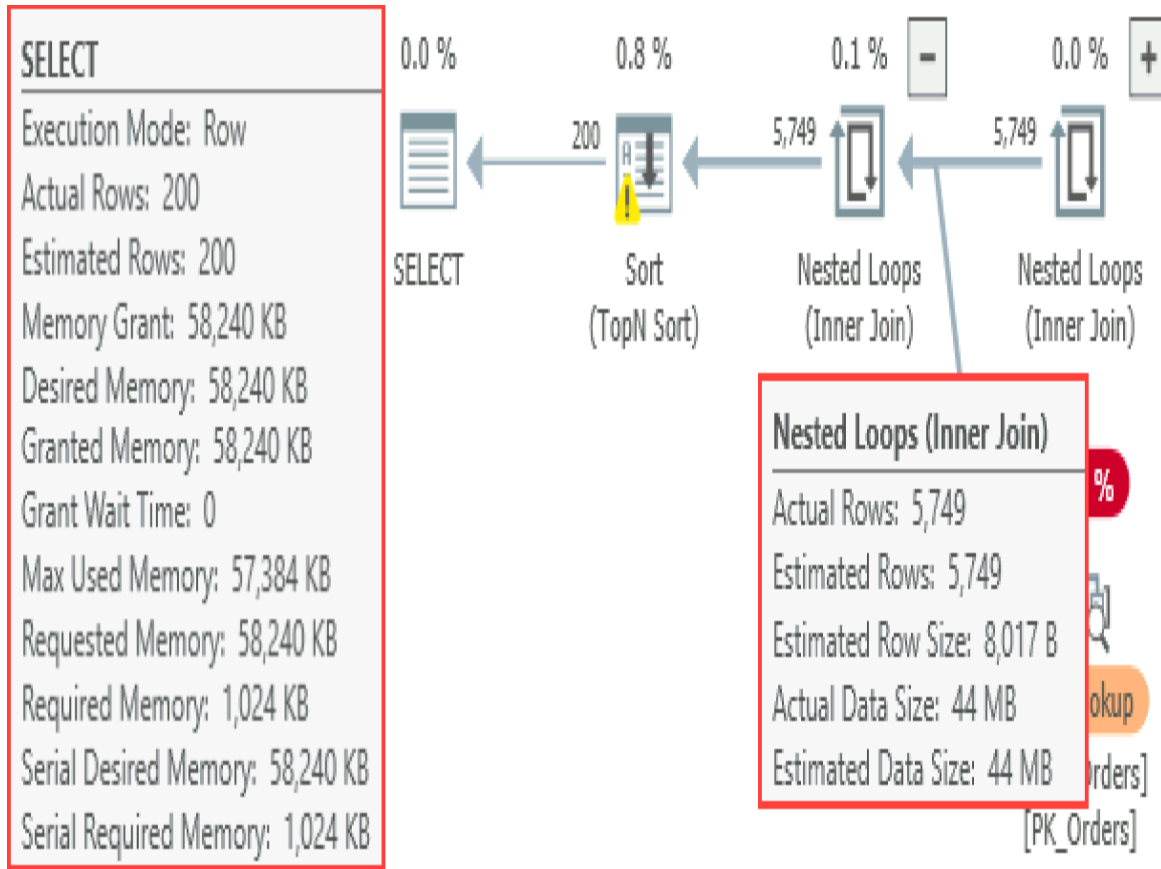


Figure 7-9. Optimizing memory intensive queries: Execution plan with up-to-date statistics

Now, let's run the code from Listing 7-12. It disables automatic statistics updates in the index and then deletes a large number of rows and clears the plan cache.

Listing 7-12. Optimizing memory intensive queries: Outdating statistics

```
ALTER INDEX IDX_Orders_OrderDate ON dbo.Orders
SET (STATISTICS_NORECOMPUTE = ON);
DELETE FROM dbo.Orders
WHERE OrderDate BETWEEN '2021-07-02' AND '2021-09-01';
DBCC FREEPROCCACHE;
```

Now let's repeat the test and run the query from Listing 7-11 again. As you can see in Figure 7-10, cardinality estimation error led to an incorrect and excessive memory grant for the query.

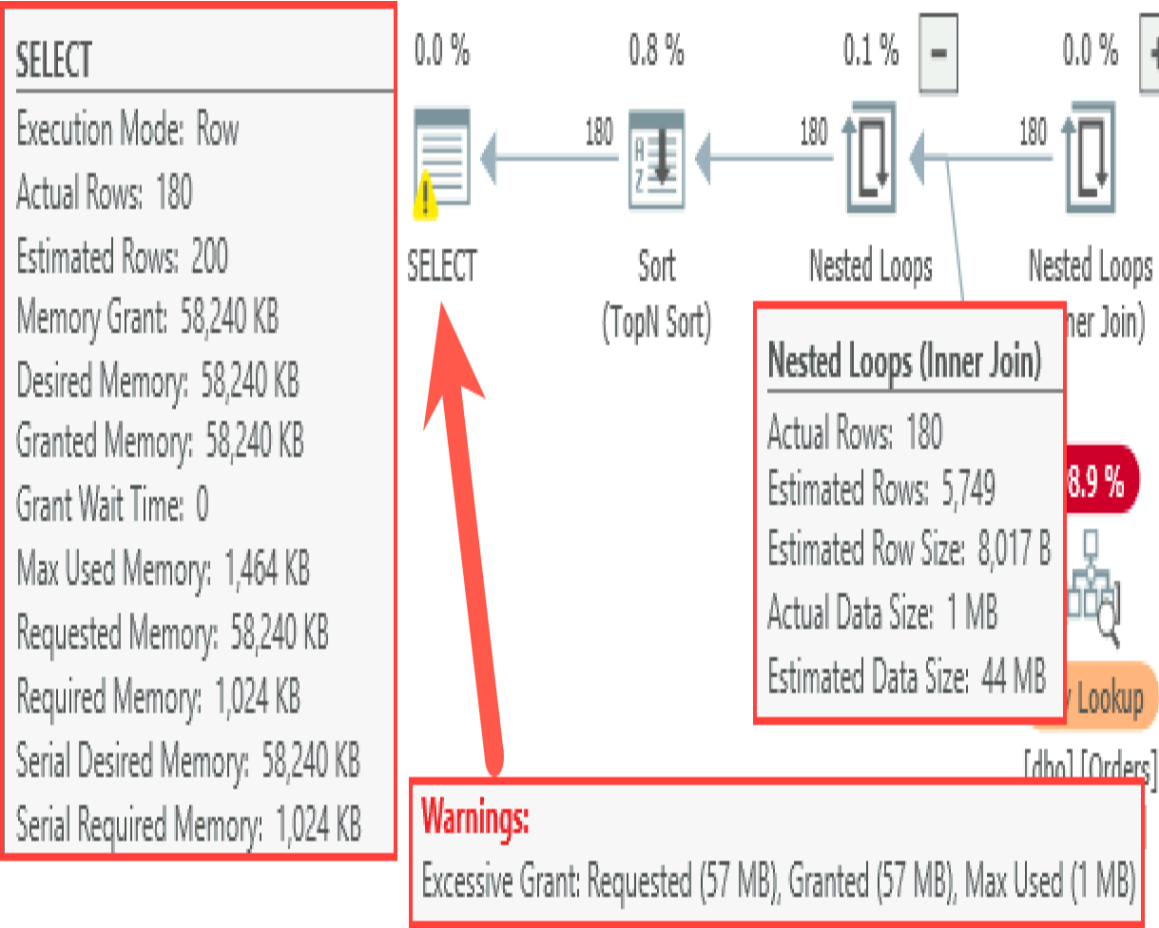


Figure 7-10. Optimizing memory intensive queries: Cardinality estimation error and memory grant

Let's update the statistics with the UPDATE STATISTICS dbo.Orders IDX_Orders_OrderDate WITH FULLSCAN command and run the query again. As you can see in Figure 7-11, correct cardinality estimation led to a significantly smaller memory grant for the query.

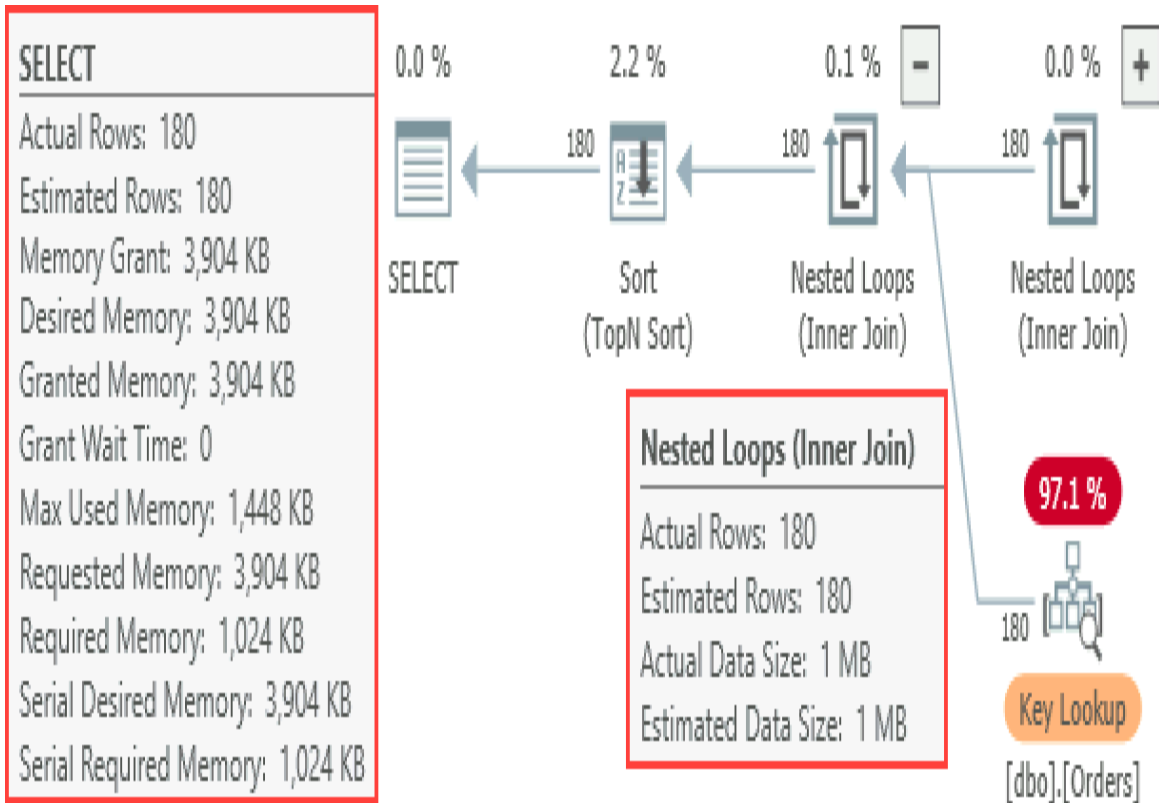


Figure 7-11. Optimizing memory intensive queries: Memory grant after UPDATE STATISTICS statement

Cardinality estimation errors are hard to deal with. Make sure that the statistics are up to date and avoid any constructs that may affect cardinality estimations (table variables, multi-statement table-value functions, etc.). In some cases, you can split complex queries into the smaller ones and persist intermediate results in temporary tables. This will lead to the overhead, which I will discuss in Chapter 9; however, it may be the small price to pay for the better execution plans in some cases.

Finally, let's look at another factor that contributes to the size of memory grant, which is the data row size. SQL Server calculates it based on data types of the columns processed by operators. For fixed-length columns the size is predefined. For example, tinyint will use one byte, int – 4 bytes, and so on.

The estimations for variable-length columns, on the other hand, depend on their length in table definition. SQL Server estimates them to be populated by 50%. For example, the column defined as varchar(100) will have 50-

byte estimation and nvarchar(200) will have 200 bytes as unicode characters are using 2 bytes to store. Finally, columns defined as (MAX), will have 4,000 bytes.

Do not select unnecessary columns and avoid using large fixed-length (n)char(n) and binary(n) data types, as they increase row size estimations and size of memory grants. You can see the impact of having a large CHAR(8000) column in the previous examples. SQL Server estimated data rows being 8,017 bytes each despite that Placeholder column stored NULL value in all rows.

Let's change the Placeholder column data type with the ALTER TABLE dbo.Orders ALTER COLUMN Placeholder VARCHAR(32) command and run the query from Listing 7-11 again. As you can see in Figure 7-12, it changed the row size estimation from 8,017 to just 37 bytes, making the memory grant significantly smaller.

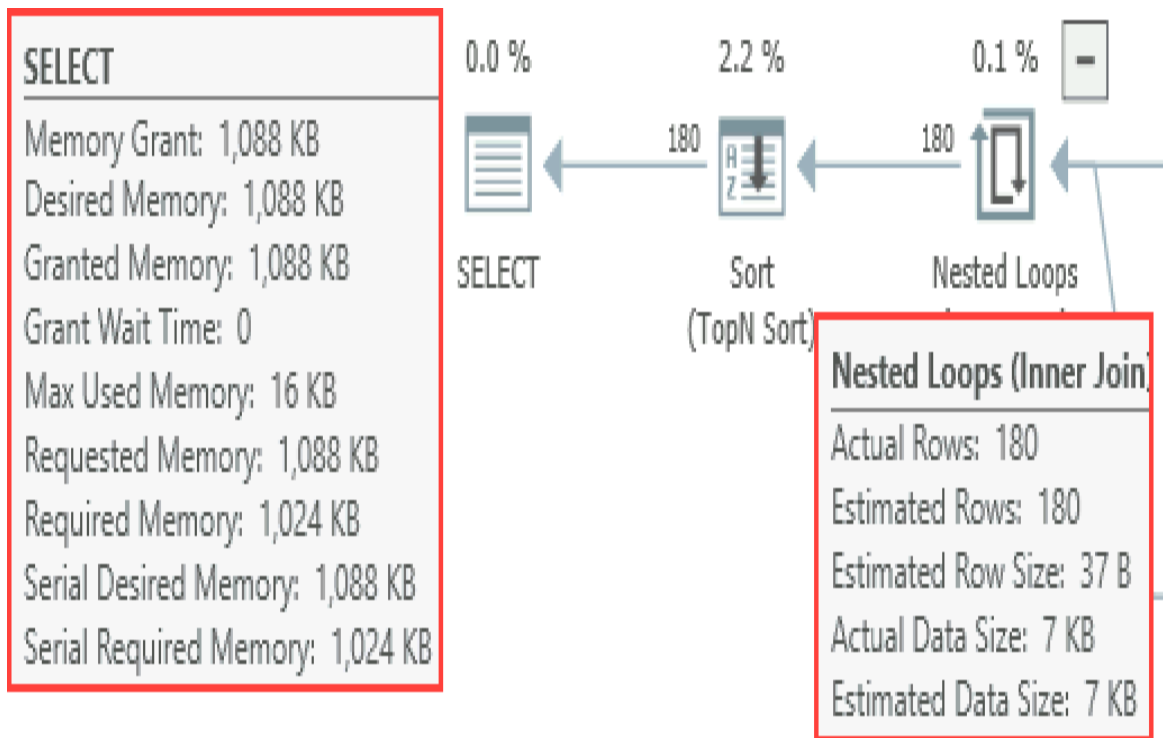


Figure 7-12. Optimizing memory intensive queries: Memory grant with VARCHAR data type

Memory Grant Feedback

The *memory grant feedback* feature, introduced in SQL Server 2017, allows SQL Server to dynamically adjust memory grants for cached execution plans based on memory usage from previous query executions. In SQL Server 2017, memory grant feedback works only with the queries that utilize batch-mode execution, which, in most cases, limits that to the queries that work with columnstore indexes. In SQL Server 2019, memory grant feedback was also extended to row-mode execution queries.

The memory grant feedback corrects both excessive and insufficient memory grants. For excessive memory grants, recalculation is triggered if the query uses less than 50% of granted memory. For insufficient grants, recalculation is triggered in the event of tempdb spill (more on this in Chapter 9). After recalculation, SQL Server updates the memory grant parameters in the cached execution plan and uses the new values going forward.

The batch-mode memory grant feedback is enabled in databases with compatibility level 140 (SQL Server 2017) and above. In the row-mode, it requires compatibility level 150 (SQL Server 2019) to be set. It also works in Azure SQL Databases.

Both, row- and batch-mode feedbacks may be disabled through database scope configuration with ALTER DATABASE SCOPED CONFIGURATION command. Row mode feedback is controlled by the ROW_MODE_MEMORY_GRANT_FEEDBACK setting. For batch-mode, however, the settings are different between SQL Server 2017 and 2019.

In SQL Server 2017, you can disable batch-mode feedback by setting the DISABLE_BATCH_MODE_MEMORY_GRANT_FEEDBACK to ON. In SQL Server 2019, you should set BATCH_MODE_MEMORY_GRANT_FEEDBACK to OFF.

You can also disable memory grant feedback on the query level with DISABLE_BATCH_MODE_MEMORY_GRANT_FEEDBACK and DISABLE_ROW_MODE_MEMORY_GRANT_FEEDBACK query hints. It may be useful with parameter-sensitive plans that suffer from parameter sniffing issues. Memory grant feedback should detect those conditions and

stop adjusting memory grants; however, in some cases you may decide to use hints to prevent any adjustments from happening.

Finally, adjusted memory grants are not persisted and are applied to cached plans only. The changes will be lost if a plan is evicted from the cache or in the event of a SQL Server restart or failover.

Controlling Memory Grant Size

SQL Server provides you two query hints- `MIN_GRANT_PERCENT` and `MAX_GRANT_PERCENT`- that allow you to specify the minimum and maximum *percentage* of workspace memory that can be granted to a query. Unfortunately, those hints are tricky to deal with. You cannot specify memory grant size in absolute units, such as KB or MB, and you must deal with a percentage of available workspace memory, which is based on hardware and SQL Server configuration.

Workspace memory can use up to 75% of target server memory. You can track both of those metrics through the Maximum Workspace Memory (KB) and Target Server Memory (KB) performance counters. By default, the maximum size of the memory grant for individual queries is 25% of workspace memory.

Let's assume that you have the query shown in Listing 7-13.

Listing 7-13. Hypothetical query

```
SELECT Col1, Col2
FROM T1
ORDER BY Col3
OPTION (MIN_GRANT_PERCENT=0.5, MAX_GRANT_PERCENT=3) ;
```

In a SQL Server instance with a target server memory of 100GB, the maximum workspace memory will be set to 75GB by default. In that configuration, the query would get at minimum 0.5% of 75GB, which is 0.375GB (384MB), and at maximum 3% of 75GB, which is 2.25GB, as the memory grant. It is worth noting that the size may be adjusted based on *required memory* that is needed to start the query.

The situation becomes even more complicated when you are using Resource Governor. It allows you to separate workspace memory across multiple resource pools using `MIN_MEMORY_PERCENT` and `MAX_MEMORY_PERCENT` resource pool settings. Moreover, you can further limit memory grants of individual requests in a resource pool's workload group by setting `REQUEST_MAX_MEMORY_GRANT_PERCENT` property.

As an example, assume that you have a resource pool with `MAX_MEMORY_PERCENT` set to 60% and a workload group with `REQUEST_MAX_MEMORY_GRANT_PERCENT` set to 10%. In that configuration, the memory grant sizes for the query from Listing 7-13 will be:

- Minimum size: $75\text{GB workspace memory} * 60\% \text{ resource pool limit} * 10\% \text{ workload group limit} * 0.5\% \text{ query limit} = 0.0225\text{GB}$ (23.04MB)
- Maximum size: $75\text{GB workspace memory} * 60\% \text{ resource pool limit} * 10\% \text{ workload group limit} * 3\% \text{ query limit} = 0.135\text{GB}$ (138.24MB)

While the Resource Governor and query hints allow you some control over memory grant size, this solution is extremely fragile. Any changes that impact memory configuration on SQL Server would lead to the different, and often unexpected, memory grant calculations. Use it with extreme care and only as the last resort when query optimization and memory grant feedback did not help. (You can read more about Resource Governor in the Microsoft [documentation](#).)

In-Memory OLTP Memory Usage and Troubleshooting

Our discussion about SQL Server memory management and troubleshooting would not be complete without covering In-Memory OLTP.

This technology relies on memory-optimized tables. The data in those tables may be persisted on-disk for durability purposes; however, SQL Server loads entire tables into memory on database startup.

This behavior is very different when compared to regular disk-based tables. With them, SQL Server always loads the data to the buffer pool; however, it does not need to load and cache the entire table. Only the active portion of the data from the table needs to be in memory.

More importantly, in case of memory pressure, SQL Server can shrink the size of the buffer pool and cache less data. It may impact performance of the system by increasing the amount of physical I/O; nevertheless, the system will be operational if it happens.

With In-Memory OLTP, on the other hand, SQL Server loads all memory-optimized data into memory on database startup. The database will not come online if the server does not have enough memory to store the data. Moreover, if at any point of time SQL Server does not have enough memory to support data growth, the memory-optimized tables become read-only.

As the amount of memory consumed by In-Memory OLTP grows, it may start to impact other SQL Server components, which would then have less memory to utilize. In the Standard Edition of SQL Server, In-Memory OLTP can utilize at most 32GB per database. In theory, the Enterprise Editions of SQL Server 2016 and above do not have any limits. However, in practice, In-Memory OLTP memory is limited to about 80% of the Resource Governor resource pool memory, the database is bound to (or DEFAULT resource pool if not bound).

You can use this behavior to limit the amount of memory available to In-Memory OLTP. Listing 7-14 shows how to do that. The `sys.sp_xtp_bind_db_resource_pool` and `sys.sp_xtp_unbind_db_resource_pool` stored procedures bind and unbind the database to and from the resource pool. You may need to restart the database for the change to take effect.

Obviously, be careful and remember that In-Memory OLTP data will become read-only if you reach the limits.

Listing 7-14. Limiting amount of memory for In-Memory OLTP

```
CREATE RESOURCE POOL InMemoryDataPool
WITH (MIN_MEMORY_PERCENT=40,MAX_MEMORY_PERCENT=40);
ALTER RESOURCE GOVERNOR RECONFIGURE;
EXEC sys.sp_xtp_bind_db_resource_pool
    @database_name = 'InMemoryOLTPDemo'
    ,@pool_name = 'InMemoryDataPool';

-- You need to take DB offline and bring it back online
-- for the changes to take effect
ALTER DATABASE MyDB SET OFFLINE;
ALTER DATABASE MyDB SET ONLINE;
```

You can monitor how much memory is consumed by In-Memory OLTP through MEMORYCLERK_XTP clerk memory consumption. In case of high memory usage, you can analyze the per-object memory consumption with the `sys.dm_db_xtp_table_memory_stats` [view](#). Listing 7-15 shows you the code to do that. You can also use the *Memory Usage by Memory Optimized Objects* report in SSMS that provides similar output.

Listing 7-13. Analyzing memory consumption of memory-optimized tables

```
SELECT
    ms.object_id
    ,s.name + '.' + t.name AS [table]
    ,ms.memory_allocated_for_table_kb
    ,ms.memory_used_by_table_kb
    ,ms.memory_allocated_for_indexes_kb
    ,ms.memory_used_by_indexes_kb
FROM sys.dm_db_xtp_table_memory_stats ms WITH (NOLOCK)
LEFT OUTER JOIN sys.tables t WITH (NOLOCK) ON
    ms.object_id = t.object_id
LEFT OUTER JOIN sys.schemas s WITH (NOLOCK) ON
    t.schema_id = s.schema_id
ORDER BY
    ms.memory_allocated_for_table_kb DESC;
```

Analyze the amount of data stored in large memory-optimized tables. In most cases, you should not retain a lot of historical data in memory – it is better to partition it between memory-optimized and disk-based tables.

Pay attention to the schema of the tables and (n)varchar(max) and varbinary(max) columns. In-Memory OLTP works very differently than disk-based tables. LOB columns introduce significant storage and performance overhead even when they are empty.

Most importantly, you need to make sure that the system does not have long-running or run-away transactions. In-Memory OLTP uses row-versioning. Data modifications generate the new versions of data rows that consume memory. Old versions and deleted data rows are eventually deallocated by garbage collection process; however, it will not process the data generated after start time of the oldest active transaction. The memory usage will continue to grow, and it eventually may put the system down.

Listing 7-16 shows the code that detects the ten oldest In-Memory OLTP transactions. You can use it for troubleshooting and you can build monitoring and alerting around it.

Listing 7-16. Detecting 10 oldest In-Memory OLTP transactions

```
SELECT TOP 10
    t.session_id
  ,t.transaction_id
  ,t.begin_tsn
  ,t.end_tsn
  ,t.state_desc
  ,t.result_desc
  ,SUBSTRING(
    qt.text
  ,er.statement_start_offset / 2 + 1
  ,(CASE er.statement_end_offset
    WHEN -1 THEN datalength(qt.text)
    ELSE er.statement_end_offset
    END - er.statement_start_offset
  ) / 2 + 1
  ) AS SQL
FROM
    sys.dm_db_xtp_transactions t WITH (NOLOCK)
  LEFT OUTER JOIN sys.dm_exec_requests er ON
```

```
t.session_id = er.session_id
CROSS APPLY sys.dm_exec_sql_text(er.sql_handle) qt
WHERE
t.state IN (0,3) /* ACTIVE/VALIDATING */
ORDER BY
t.begin_tsn
OPTION (RECOMPILE, MAXDOP 1);
```

In-Memory OLTP is a great technology that can significantly improve the throughput of OLTP systems. It is not, however, a set-it-and-forget-it type of technology. It requires proper system and database design and adequate monitoring in production. Consider reading my book *Expert SQL Server In-Memory OLTP* (2nd edition, Apress, 2017) if you want to learn more.

Summary

SQL Server is a memory-intensive application that may consume hundreds of gigabytes or terabytes of memory. This is completely normal and it can help improve SQL Server performance. Nevertheless, it is important to properly configure server memory, especially in non-dedicated environments.

Set and tune the *Maximum Server Memory* setting, leaving enough memory for the OS and other applications. Consider granting the *Lock Pages In Memory* privilege to SQL Server accounts; however, remember that it may lead to system stability issues if the system is not properly configured.

You may also consider enabling *Large Page Allocations* on servers with a large amount of memory. This feature does not work well with columnstore indexes in SQL Server versions prior to 2019.

You can analyze the current memory usage by looking at the memory consumption of various memory clerks. Detect anomalies and address root-causes of the issues.

Monitor the status of memory grants and presence of RESOURCE_SEMAPHORE waits. Optimize queries with high memory

grants when possible. Enable memory grant feedback feature if it is available in your version of SQL Server.

In the next chapter, I will talk about the QL Server concurrency model and explain how to troubleshoot blocking issues and deadlocks.

Troubleshooting Checklist

Check and adjust memory configuration.

Analyze memory usage with `sys.dm_os_memory_clerks` view. Address possible issues.

Analyze plan cache memory usage.

Analyze memory usage from single-use ad-hoc execution plans.

Check for `RESOURCE_SEMAPHORE` waits. Detect and optimize for the most memory-intensive queries.

Chapter 8. Troubleshooting TempDB Usage and Performance

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 9th chapter of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at dmitri@aboutsqlserver.com.

The tempdb is the system database, which is shared across all user and system sessions. It stores user-created and internal temporary objects and data, and is used by many processes. High tempdb performance and throughput are essential for good server performance.

I will start this chapter with an overview of tempdb consumers and usage patterns and share several best practices related to the usages of temporary objects. Next, I will show how to diagnose and address common tempdb issues. Finally, I’ll provide you with several tempdb configuration tips.

Temporary Objects: Usage and Best Practices

Tempdb performance tuning is a complex topic. I usually like to start with an overview of usage patterns. After all, tempdb is just another database, and reducing its load usually improves its throughput. There are some internal optimizations in tempdb behavior, which I will discuss later in that chapter. But for all practical purposes, you can consider tempdb similar to other user databases.

The tempdb database stores temporary objects created by users, internal record sets generated during query executions, the version store, and a few other objects. You can look at how much space different object types are using by running the code from Listing 9-1.

Listing 9-1. Tempdb space usage

```
SELECT
    CONVERT(DECIMAL(12,3),
    SUM(user_object_reserved_page_count) / 128.)
        AS [User Objects (MB)]
    ,CONVERT(DECIMAL(12,3),
    SUM(internal_object_reserved_page_count) / 128.
    ) AS [Internal Objects (MB)]
    ,CONVERT(DECIMAL(12,3),
    SUM(version_store_reserved_page_count) / 128.
    ) AS [Version Store (MB)]
    ,CONVERT(DECIMAL(12,3),
    SUM(unallocated_extent_page_count) / 128.
    ) AS [Free Space (MB)]
FROM
    tempdb.sys.dm_db_file_space_usage WITH (NOLOCK);
```

Let's start our discussion with the first group in that list – objects created by users.

Temporary Tables and Table Variables

Temporary tables and table variables store short-lived information and intermediate results during data processing. For most part, temporary tables behave similarly to regular user tables. They don't support triggers and cannot be included in views; however, they support indexes and constraints

and can be altered like regular tables. Altering them is not a good idea, though – I'll explain why later in the chapter.

There are two kind of temporary tables – global and local. They differ in lifespan and visibility. *Global temporary tables* are created with names that start with two hash symbols (##) and are visible to all sessions. They are dropped when the session in which they were created disconnects and other sessions stop referencing them.

You can use global temporary tables to store and share temporary data between sessions. This approach, however, is fragile and prone to errors. For example, if the session that created the global temporary table loses connection to the database, the table may be dropped at an unpredictable time. You may get better results by creating regular tables in tempdb instead.

Local temporary tables are named starting with one hash symbol (#) and are visible only in the session in which they were created. When multiple sessions simultaneously create local temporary tables with the same name, every session will have its own instance of the table.

Local temporary tables are visible in the module in which they were created and in all other modules called from that module. For example, if you open a connection and create a temporary table in that session, the table will be visible everywhere in that session and live while the session is open.

Alternatively, if you create a temporary table in the stored procedure, it will be visible in that stored procedure and all other T-SQL modules or dynamic SQL called from there. It will be dropped automatically when that stored procedure completes.

You can use this behavior to pass data between T-SQL modules.

Nevertheless, it has a couple of downsides: First, it increases the number of compilations and CPU load. SQL Server will need to recompile the inner (called) module as it does not know anything about external table until the module was called.

Second, this approach is also extremely fragile. Any changes in temporary table schema in outer (calling) modules can break the inner ones. The

situation will become even worse if the inner modules are executed by multiple callers – it quickly gets very hard to support. Use this approach with extreme care and only when absolutely necessary.

In contrast – table variables are visible only in the module where they were defined. You can pass them as parameters to other modules (more about this later in this chapter).

Despite the old myth, table variables are not in-memory objects. They use tempdb similarly to temporary tables. They introduce less overhead than temporary tables; however, that benefit comes with a major limitation: They don't support indexes, except for primary key and unique constraints. More importantly, table variables do not maintain statistics on those constraints. This can lead to significant cardinality estimation errors and highly inefficient execution plans.

Let's look at the example in Listing 9-2. Here you'll create a temporary table and populate it with some data.

Listing 9-2. Cardinality estimations: Temporary table creation

```
CREATE TABLE #TT(ID INT NOT NULL PRIMARY KEY);
;WITH N1(C) AS (SELECT 0 UNION ALL SELECT 0) -- 2 rows
,N2(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N1 AS T2) -- 4 rows
,N3(C) AS (SELECT 0 FROM N2 AS T1 CROSS JOIN N2 AS T2) -- 16
rows
,N4(C) AS (SELECT 0 FROM N3 AS T1 CROSS JOIN N3 AS T2) -- 256
rows
,IDs(ID) AS (SELECT ROW_NUMBER() OVER (ORDER BY (SELECT NULL))
FROM N4)
INSERT INTO #TT(ID)
SELECT ID FROM IDs;
```

As the next step, run the code from Listing 9-3. This will select the data from the temporary table and table variable and compare cardinality estimations in the queries. (Note that I am running that demo in the database with a compatibility level of 140 in SQL Server 2017. It will behave slightly differently at compatibility level 150 or later, as in SQL Server 2019 – I will cover this shortly.)





Listing 9-3. Cardinality estimations: Selecting data from temporary objects

```
DECLARE
    @TTV TABLE (ID INT NOT NULL PRIMARY KEY);
INSERT INTO @TTV (ID)
    SELECT ID FROM #TT;
SELECT COUNT(*) FROM #TT;
SELECT COUNT(*) FROM @TTV;
SELECT COUNT(*) FROM @TTV OPTION (RECOMPILE);
```

Figure 9-1 shows cardinality estimations for SELECT queries. As you can see, the estimation is correct for the temporary table. However, unless you are using a statement-level recompile, SQL Server estimates that a table variable has only one row. Cardinality estimation errors can progress quickly through the execution plan, which means that using table variables can lead to highly inefficient plans.

Query 2: Query cost (relative to the batch): 15%





SELECT COUNT(*) FROM #TT

		
	Stream Aggregate (Aggregate)	Clustered Index S...
	[#TT].[PK_#TT_...	
 SELECT	 Compute Scalar	
Cost: 0 %	Cost: 0 %	Cost: 4 %
		0.000s
		1 of
		1 (100%)
		256 of
		256 (100%)

Clustered Index Scan	
Number of Rows Read	256
Estimated Number of Rows to be Read	256

Query 3: Query cost (relative to the batch): 13%





SELECT COUNT(*) FROM @TTV

		
	Stream Aggregate (Aggregate)	Clustered Index S...
	[@TTV].[PK_#B37A...	
 SELECT	 Compute Scalar	
Cost: 0 %	Cost: 0 %	Cost: 0 %
		0.000s
		1 of
		1 (100%)
		256 of
		1 (25600%)

Clustered Index Scan	
Number of Rows Read	256
Estimated Number of Rows to be Read	1

Query 4: Query cost (relative to the batch): 15%

SELECT COUNT(*) FROM @TTV OPTION (RECOMPILE)

		
	Stream Aggregate (Aggregate)	Clustered Index S...
	[@TTV].[PK_#B562...	
 SELECT	 Compute Scalar	
Cost: 0 %	Cost: 0 %	Cost: 4 %
		0.000s
		1 of
		1 (100%)
		256 of
		256 (100%)

Clustered Index Scan	
Number of Rows Read	256
Estimated Number of Rows to be Read	256

Figure 8-1. Cardinality Estimations – Temporary Tables vs. Table Variables

A statement-level recompile provides Query Optimizer information about the total number of rows (note, however, that table variables do not keep statistics or information about data distribution).

Let's repeat the test by adding WHERE clause to the queries, as shown in Listing 9-4. All rows in the tables have positive ID values.

Listing 9-4. Cardinality estimations: Selecting data with WHERE clause





```
DECLARE
    @TTV TABLE(ID INT NOT NULL PRIMARY KEY);
INSERT INTO @TTV(ID)
    SELECT ID FROM #TT;
SELECT COUNT(*) FROM #TT WHERE ID > 0;
SELECT COUNT(*) FROM @TTV WHERE ID > 0;
SELECT COUNT(*) FROM @TTV WHERE ID > 0 OPTION (RECOMPILE);
```

Figure 9-2 shows the cardinality estimations for the scenario in Listing 9-4. Temporary tables maintain statistics on indexes, so SQL Server was able to estimate the number of rows in the first SELECT correctly.

As before, without a statement-level recompile, SQL Server always assumes that the table variable has only a single row. But even with a statement-level recompile, the estimations are way off. There are no statistics, and SQL Server assumes that the greater-than (>) operator will filter out about 70% of the rows from the table, which is completely incorrect.

Query 2: Query cost (relative to the batch): 15%





SELECT COUNT(*) FROM #TT WHERE ID > 0

			
		Stream Aggregate (Aggregate)	Clustered Index S...
SELECT	Compute Scalar	Cost: 4 %	Cost: 96 %
Cost: 0 %	Cost: 0 %	0.000s	0.000s
		1 of	256 of
		1 (100%)	256 (100%)

Clustered Index Seek	
Number of Rows Read	256
Estimated Number of Rows to be Read	256

Query 3: Query cost (relative to the batch): 14%





SELECT COUNT(*) FROM @TTV WHERE ID > 0

			
		Stream Aggregate (Aggregate)	Clustered Index S...
SELECT	Compute Scalar	Cost: 0 %	Cost: 100 %
Cost: 0 %	Cost: 0 %	0.000s	0.000s
		1 of	256 of
		1 (100%)	1 (25600%)

Clustered Index Seek	
Number of Rows Read	256
Estimated Number of Rows to be Read	1

Query 4: Query cost (relative to the batch): 14%

SELECT COUNT(*) FROM @TTV WHERE ID > 0 OPTION (RECOMPILE)

			
		Stream Aggregate (Aggregate)	Clustered Index S...
SELECT	Compute Scalar	Cost: 1 %	Cost: 99 %
Cost: 0 %	Cost: 0 %	0.000s	0.000s
		1 of	256 of
		1 (100%)	77 (332%)

Clustered Index Seek	
Number of Rows Read	256
Estimated Number of Rows to be Read	76.8

Figure 8-2. Cardinality Estimations – Table Variables and WHERE clause

If you run the scripts in a database with compatibility level 150 or above (SQL Server 2019), you'll get slightly different estimations in the second query. SQL Server defers compiling statements with table variables and uses the number of rows there at the moment of compilation, caching the generated plan afterwards. Nevertheless, there are still no statistics, and the estimation will not be correct when a WHERE clause is used. Remember this behavior.

WARNING BOX

WARNING

Table variables are *not* a good choice when cardinality estimation errors could impact execution plans—for example, if they store large amounts of data and participate in joins with other tables. In most cases, the Query Optimizer will choose a loop join, which is inefficient with large inputs. Use temporary tables instead – they are a *much* safer choice than table variables. In many cases, simply switching table variables to temporary tables has improved my query performance significantly.

Do not forget to index temporary tables properly when you work with them. All of the optimization rules you learned in Chapter 5 apply, whether you're dealing with regular or temporary tables.

Temporary tables are a great tool to improve cardinality estimations and optimize complex queries. You can split a complex query into a few simpler ones and store the intermediate results in a temporary table. Smaller queries are easier to optimize, especially with up-to-date statistics generated in temporary tables.

As usual, this approach comes with a downside: the overhead of creating and populating temporary tables. Figure 9-3 gives a high-level overview of that process. With one exception, which I'll mention, the same process also applies to table variables.

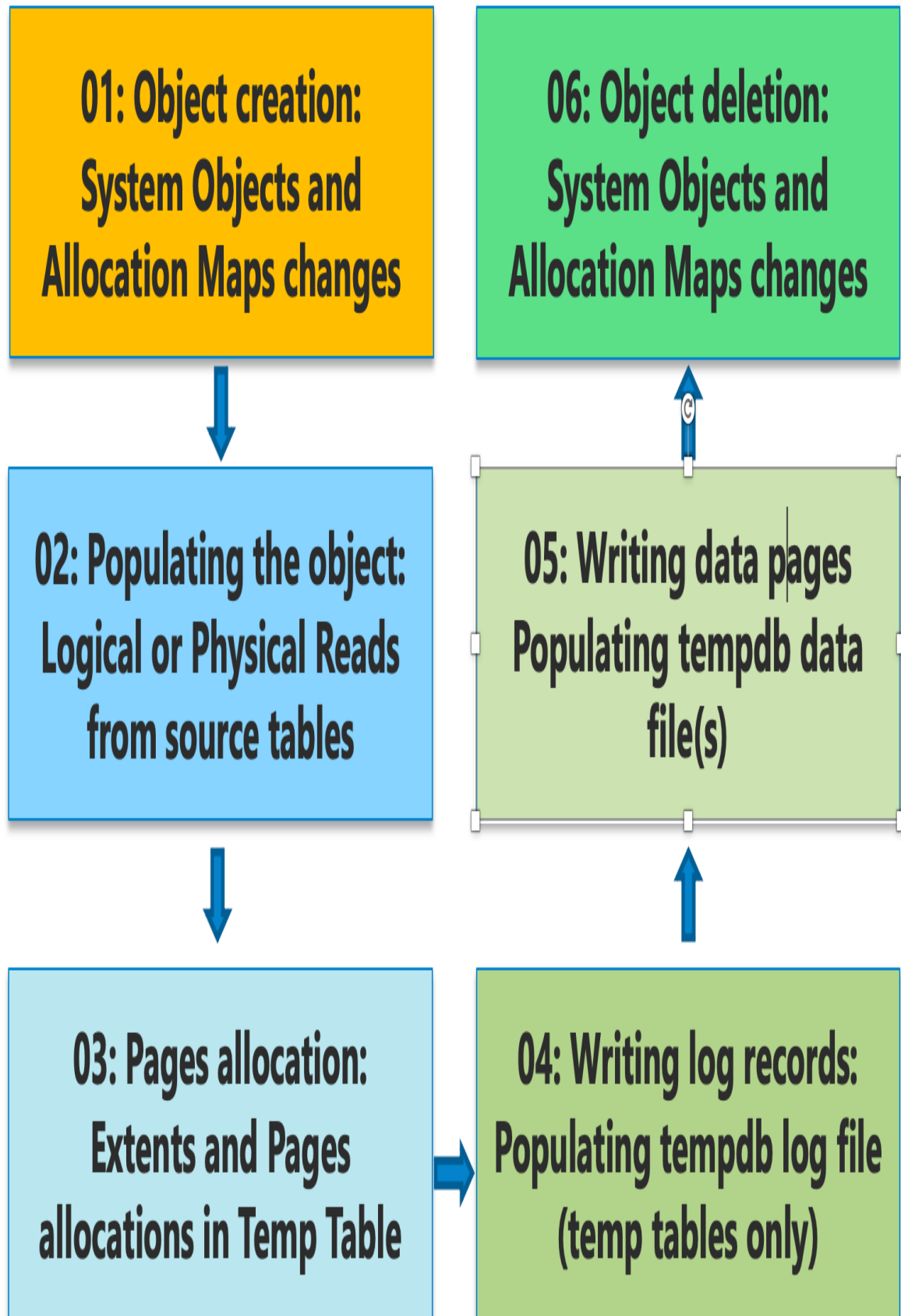


Figure 8-3. Using temporary table to store intermediate data

Let's discuss what happens during the process pictured in Figure 9-3.

Step 1: Object Creation

When SQL Server creates a temporary table or table variable, it makes several modifications in the system catalogs and allocation map pages in tempdb. This is a very fast process; however, it may become a source of contention when multiple sessions are updating system pages simultaneously.

This condition can be reduced by proper tempdb configuration, which I will cover later in this chapter. I'll also explain how to properly diagnose it in this and the next chapter of the book.

Step 2: Populating the temporary object

Next, SQL Server populates the temporary object by performing logical or physical reads from the source tables.

Step 3: Pages allocation

SQL Server allocates data pages for the new object in the buffer pool and modifies them, marking as dirty.

Step 4: Writing log records

For temporary tables, SQL Server logs the above actions in the tempdb transaction log.

Changes in table variables are not logged: they are not transaction-aware and a ROLLBACK operation would not undo any actions against them. (This is another performance advantage they offer over temporary tables.)

Step 5: Writing data pages

In a separate, asynchronous process, dirty data pages will eventually be written to the data files. This might even be done after the object is

dropped from the database.

Data pages that belong to temporary objects stay in the buffer pool and behave similarly to data pages from the regular tables. They consume buffer pool space and reduce the amount of memory needed to cache the data from the regular tables. You can use Listing 7-4 to analyze tempdb buffer pool usage.

Step 6: Object deletion

Eventually, the temporary objects need to be deallocated. While this operation is relatively light, it also modifies the system catalogs and may lead to contention in busy systems.

As you can see, creating and populating temporary objects can be expensive, especially when you're dealing with large amounts of data. While temporary tables are a great query optimization tool for *reasonably* small intermediate results, it is usually not a good idea to store millions of rows in them.

Moreover, reading the data is significantly less expensive than writing it. It may be cheaper and faster to read large amounts of data from the regular table a few times than to write significant portions of it into the temporary table.

There are no hard thresholds for when you should or should not use temporary tables. It depends on your workload, the amount of data, your hardware configuration, and what problems you are trying to solve. Just remember the overhead they introduce and make sure it doesn't outweigh the benefits of using temporary tables.

Temporary Object Caching

I mentioned just now that creating and deallocating temporary objects requires SQL Server to modify system catalogs and allocation map pages in tempdb. One of SQL Server's optimizations, *temporary object caching*, helps to reduce that overhead.

The name of this feature is a bit confusing. It relates to caching temporary objects' *allocation* pages, not data pages. When you use temporary object caching, instead of dropping the table, SQL Server truncates it. It keeps two pages per index pre-allocated: one IAM and one data page. The next time the table is created, SQL Server reuses these pages, which helps reduce the number of modifications required in the allocation map pages.

Let's look at Listing 9-5 and define the stored procedure that creates and drops the temporary table.

Listing 9-5. Temporary object caching: Stored procedure

```
CREATE PROC dbo.TempTableCaching
AS
    CREATE TABLE #T(C INT NOT NULL PRIMARY KEY);
    DROP TABLE #T;
```

Next, run the stored procedure and examine the transaction log records it generates. You can see the code in Listing 9-6.

Listing 9-6. Temporary object caching: Running stored procedure

```
CHECKPOINT;
GO

EXEC dbo.TempTableCaching;

SELECT
    Operation, Context, AllocUnitName
    , [Transaction Name], [Description]
FROM
    tempdb.sys.fn_dblog(null, null);
```

You can see the output of the code in Figure 9-4. Here, the first stored procedure call produced 45 log records, most related to updating the allocation map pages and system tables while creating the temporary table.

	Operation	Context	AllocUnitName	Transaction Name	Description
8	LOP_INSERT_ROWS	LCX_INDEX_LEAF	sys.syssschobjs.nc2	NULL	
9	LOP_INSERT_ROWS	LCX_INDEX_LEAF	sys.syssschobjs.nc3	NULL	
10	LOP_MODIFY_ROW	LCX_CLUSTERED	sys.syssschobjs.clst	NULL	
11	LOP_INSERT_ROWS	LCX_CLUSTERED	sys.syscolpars.clst	NULL	
12	LOP_INSERT_ROWS	LCX_INDEX_LEAF	sys.syscolpars.nc	NULL	
13	LOP_MODIFY_ROW	LCX_CLUSTERED	sys.syssschobjs.clst	NULL	
14	LOP_MODIFY_ROW	LCX_CLUSTERED	sys.syssschobjs.clst	NULL	
15	LOP_INSERT_ROWS	LCX_CLUSTERED	sys.sysrowsets.clu...	NULL	
16	LOP_INSERT_ROWS	LCX_CLUSTERED	sys.sysallocunits...	NULL	
17	LOP_INSERT_ROWS	LCX_INDEX_LEAF	sys.sysallocunits...	NULL	
18	LOP_INSERT_ROWS	LCX_CLUSTERED	sys.sysrscols.clst	NULL	
19	LOP_HOBT_DDL	LCX_NULL	NULL	NULL	Action 1 (CREATE_HOBT) on...
20	LOP_INSERT_ROWS	LCX_CLUSTERED	sys.sysidxstats.cl...	NULL	

↙

Ln 17, Col 1 Spaces: 4 UTF-8 LF SQL SentryOne Plan Explorer : On MSSQL 45 rows 00:00:00 19

Figure 8-4. Temporary object caching: Log records after the first call

The situation changes when you run the code from Listing 9-6 a second time. Now, when the temporary table is cached, table creation introduces just a few log records, all of which are against the system table, with no allocation map pages involved. You can see this in Figure 9-5.

	Operation	Context	AllocUnitName	Transaction Name	Description
2	LOP_XACT_CKPT	LCX_BOOT_PAGE_CKPT	NULL	NULL	
3	LOP_END_CKPT	LCX_NULL	NULL	NULL	2021/06/29 10:10:53:063;L...
4	LOP_BEGIN_XACT	LCX_NULL	NULL	CREATE TABLE	2021/06/29 10:10:55:813;C...
5	LOP_SHRINK_NOOP	LCX_NULL	NULL	NULL	
6	LOP_MODIFY_ROW	LCX_CLUSTERED	sys.sysschobjs.clst	NULL	
7	LOP_DELETE_ROWS	LCX_MARK_AS_GHOST	sys.sysschobjs.nc1	NULL	
8	LOP_INSERT_ROWS	LCX_INDEX_LEAF	sys.sysschobjs.nc1	NULL	
9	LOP_DELETE_ROWS	LCX_MARK_AS_GHOST	sys.sysschobjs.nc2	NULL	
10	LOP_INSERT_ROWS	LCX_INDEX_LEAF	sys.sysschobjs.nc2	NULL	
11	LOP_MODIFY_ROW	LCX_CLUSTERED	sys.sysschobjs.clst	NULL	
12	LOP_COMMIT_XACT	LCX_NULL	NULL	NULL	2021/06/29 10:10:55:813
13	LOP_BEGIN_XACT	LCX_NULL	NULL	DROPOBJ	2021/06/29 10:10:55:813;D...
14	LOP_SHRINK_NOOP	LCX_NULL	NULL	NULL	

Ln 10, Col 1 (156 selected) Spaces: 4 UTF-8 LF SQL SentryOne Plan Explorer : On MSSQL 24 rows 00:00:00 192.168.250.100 : t

Figure 8-5. Temporary object caching: Log records after the second call

Temporary object caching is enabled by default for all temporary objects created in stored procedures and triggers (session-level objects are not cached). However, there are a few conditions:

- The table needs to be smaller than 8MB. Large tables are not cached.
- There can be no DDL statements that change the table structure. Any schema modification, with exception of DROP TABLE, will prevent temporary object caching. However, you can create indexes on the table and SQL Server will cache them.
- No named constraints can be defined in the table. Unnamed constraints will not prevent caching.

Consider these guidelines as you write your code. Temporary object caching is a useful performance feature, and you will benefit from it.

Table-Valued Parameters

SQL Server allows you to define table types in the database. When you declare the variable of the table type in the code, it works the same way as table variables. You can also pass the variables of the table types as parameters to T-SQL modules. Those parameters are called *table-valued parameters*.

Table-valued parameters are implemented as table variables under the hood and inherit all their benefits and limitations, most notable around missing statistics and cardinality estimation errors they may introduce. In addition, those parameters are read-only. You cannot insert, update, or delete the data from table-valued parameters in the modules they passed to.

Listing 9-7 shows how you can use table-valued parameters. (This is just an example of a possible usage scenario, not a reference implementation for any use cases!) It also demonstrates that table variables are not transaction-aware. You can use them to pass the information outside of the transaction you are rolling back.

Listing 9-7. Using table-valued parameters

```
CREATE TYPE dbo.tvpTransfers AS TABLE
(
    FromAccount BIGINT NOT NULL,
    ToAccount BIGINT NOT NULL,
    ADate DATETIME2(0) NOT NULL,
    Amount MONEY NOT NULL,
    PRIMARY KEY(FromAccount, ToAccount)
);
GO

CREATE PROC dbo.ProcessRejectedTransfers
(
    @RejectedTransfers dbo.tvpTransfers READONLY
)
AS
    SELECT FromAccount, ToAccount, ADate, Amount
    FROM @RejectedTransfers;
GO

CREATE PROC dbo.DoTransfers
```



```

(
    @Transfers dbo.tvpTransfers READONLY
)
AS
    DECLARE
        @RejectedTransfers dbo.tvpTransfers

    BEGIN TRAN
        INSERT INTO @RejectedTransfers
            (FromAccount, ToAccount, ADate, Amount)
            SELECT FromAccount, ToAccount, ADate, Amount
            FROM @Transfers
            WHERE Amount > 10000;
        /* ... */
    ROLLBACK -- Table variables are not transaction-aware.
    EXEC sp_executesql
        N'EXEC dbo.ProcessRejectedTransfers @Transfers;'
        ,N'@Transfers dbo.tvpTransfers READONLY'
        ,@Transfers = @RejectedTransfers;
GO

DECLARE
    @Transfers dbo.tvpTransfers

INSERT INTO @Transfers
    (FromAccount, ToAccount, ADate, Amount)
VALUES
    (1,2,'2021-08-01',100)
    , (3,4,'2021-08-02',15000)
    , (5,6,'2021-08-03',20000);

EXEC dbo.DoTransfers @Transfers;

```

Table-valued parameters are one of the fastest ways to pass a batch of rows from a client application to a T-SQL routine. They are an order of magnitude faster than separate DML statements and, in some cases, can even outperform bulk operations. Look for opportunities to use them – they can provide significant performance improvements.

NOTE

I am including the demo application that compares performance of different techniques to pass a batch of rows from client application to the database in companion materials of the book.

Regular Tables in TempDb and Transaction Logging

You can create and use regular tables in tempdb. Those tables are visible to all sessions and behave the same way as in the other databases. Obviously, the system needs to be ready in case those tables disappear during SQL Server restarts, when tempdb is recreated. You have two options: you can recreate them by defining them in a model database or you can use startup stored procedures.

Tempdb can be a good choice as the staging area for ETL processes in cases where you need to load and process a large amount of data *and* the process does not have High Availability requirements. Tempdb uses the SIMPLE recovery model and more efficient transaction logging. It does not need to support crash recovery and, therefore, does not need to keep the REDO portion of transaction log record, which is what allows SQL Server to reapply changes from committed transactions if it crashes before the checkpoint.

Transaction logging of local temporary tables is even more efficient. Those tables are visible in the scope of a single session, which allows SQL Server to use minimally logged operations in more cases (I'll cover this in more detail in Chapter 11). You can use use minimally logged operations to improve ETL performance even further. Remember, however, that temporary tables will disappear if the client loses connection to the database.

Table 9-1 shows examples of when operations are minimally logged for regular and temporary tables in tempdb. You can use it as a reference to speed up ETL processes and initial data load to the tables. As I stated

earlier, table variables do not use logging at all; however, they open the door to cardinality estimation errors during query optimization.

*T
a
b
l
e*

*8
-*

*l
·
M
i
n
i
m
a
l
l
y*

*l
o
g
g
e
d*

*o
p
e
r
a
t*

i
o
n
s

i
n

t
e

m

p

d

b

Operation	Minimally logged?
-----------	-------------------

SELECT INTO dbo.RegularTable SELECT INTO #tempTable	Yes—although I don't recommend this pattern, because it does not create clustered indexes in the table
--	--

Yes

INSERT INTO dbo.RegularTable WITH (TABLOCK) SELECT

No

INSERT INTO dbo.RegularTable SELECT

Yes

INSERT INTO #tempTable SELECT

Using tempdb for ETL processes may improve their performance and reduce the load on the server. In many cases, however, you can achieve even better results by using durable or non-durable memory-optimized tables and In-Memory OLTP. Consider them as options, but test implementation carefully: In-Memory OLTP behaves differently than the classic Storage Engine.

Internal TempDB Consumers

In addition to objects created by users, SQL Server uses tempdb to store internal objects. The two most common space consumers in that group are version store and internal row sets, which are generated by Sort, Hash and Exchange operations that spill over to tempdb.

Let's look at both in detail.

Version Store

Several SQL Server features rely on row versioning. In addition to optimistic isolation levels, such as RCSI and SNAPSHOT, row versioning is used by online index rebuild, triggers and Multiple Active Result Sets (MARS). As you learned in chapter 8, the old versions of the rows are stored in tempdb version store.

Figure 9-6 shows version store behavior with a large row-versioning transaction running in the system. You can see the growth of version store size (Version Store Size (KB) performance counter), which triggers tempdb auto-growth events.

In this case, SQL Server cleared up the version store after the transaction completed. It took some time, however, for clean-up to occur. The clean-up process is asynchronous and runs on a schedule.

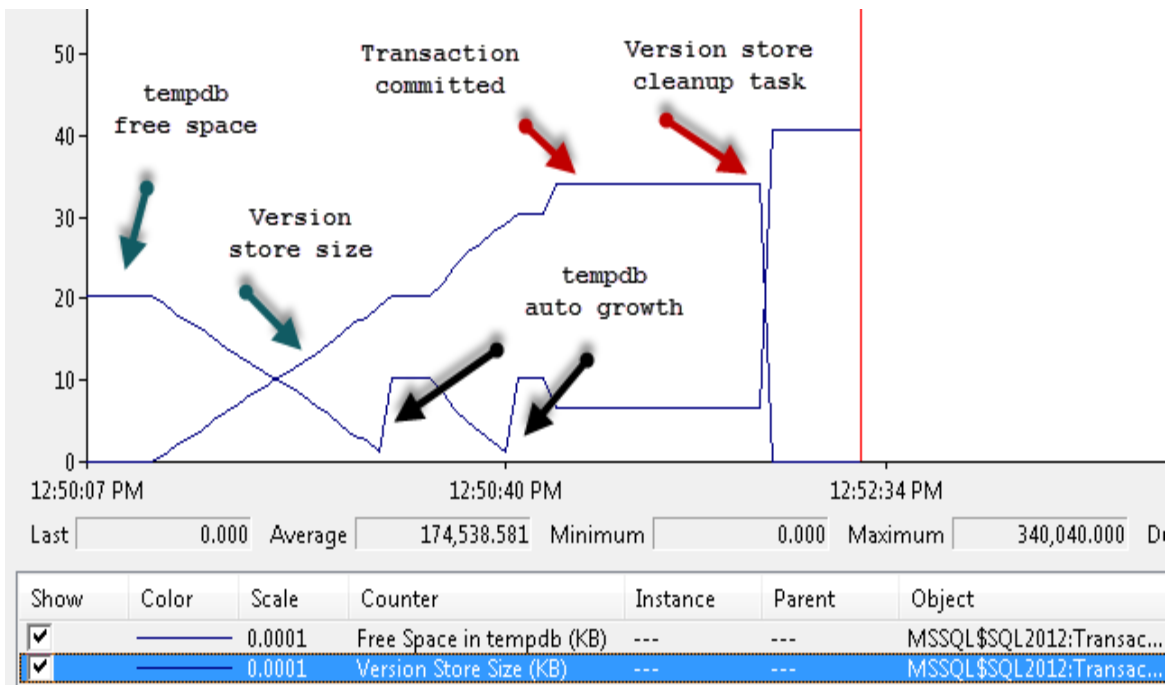


Figure 8-6. Version store growth and clean-up

Version store growth is one of the common reasons of excessive growth of the tempdb. SQL Server does not clean up the version store beyond the starting point of the oldest active row-versioning transaction. Uncommitted run-away transactions can lead to extensive version store and tempdb growth, even if tempdb does not generate row versions by itself.

Listing 9-8 shows you how to detect the five oldest row-versioning transactions using `sys.dm_tran_active_snapshot_database_transactions` view. You can kill the session that prevents clean-up of the version store in case of emergency.

Listing 9-8. Detecting the five oldest row-versioning transactions

```

SELECT TOP 5
    at.transaction_id
    ,at.elapsed_time_seconds
    ,at.session_id
    ,s.login_time
    ,s.login_name
    ,s.host_name
    ,s.program_name
    ,s.last_request_start_time
    ,s.last_request_end_time

```

```

,er.status
,er.wait_type
,er.wait_time
,er.blocking_session_id
,er.last_wait_type
,st.text AS [SQL]
FROM
    sys.dm_tran_active_snapshot_database_transactions at WITH
(NOLOCK)
    JOIN sys.dm_exec_sessions s WITH (NOLOCK) on
        at.session_id = s.session_id
    LEFT JOIN sys.dm_exec_requests er WITH (NOLOCK) on
        at.session_id = er.session_id
    CROSS APPLY
        sys.dm_exec_sql_text(er.sql_handle) st
ORDER BY
    at.elapsed_time_seconds DESC;

```

NOTE

Long-running transactions on readable secondaries in Availability Groups can defer version store clean-up. I'll discuss this in Chapter 12.

There are several performance counters in the Transactions performance object that you can use to monitor version store behavior:

Version Store Size (KB)

Current size of the version store.

Version Generation rate (KB/s)

Growth rate of the version store.

Version Cleanup rate (KB/s)

Cleanup rate of the version store.

Longest Transaction Running Time

Duration in seconds of the oldest active transaction that is using row versioning.

Free Space in tempdb (KB)

Amount of available space in tempdb. While this counter is not related to the version store, you can use it for general tempdb monitoring.

Listing 9-9 shows the queries that you can use to analyze version store usage per database. This may be useful when you troubleshoot excessive version-store growth on servers that host multiple databases.

The first query would work on SQL Server 2016 and above. The second can be used on the older versions of SQL Server. It is more resource intensive, however, and provides slightly less accurate results.

Listing 9-9. Version store usage per database

```
-- SQL Server 2016 SP2 and above
SELECT
    DB_NAME(database_id) AS [DB]
    ,database_id
    ,reserved_page_count
    ,CONVERT(DECIMAL(12,3),reserved_space_kb / 1024.)
        AS [Reserved Space (MB)]
FROM
    sys.dm_tran_version_store_space_usage WITH (NOLOCK)
OPTION (RECOMPILE);

-- SQL Server 2014 and below. Less accurate.
SELECT
    DB_NAME(database_id) AS [DB]
    ,database_id
    ,CONVERT(DECIMAL(12,3),
        SUM(record_length_first_part_in_bytes +
            record_length_second_part_in_bytes) / 1024. / 1024.
    ) AS [Version Store (MB)]
FROM
    sys.dm_tran_version_store WITH (NOLOCK)
GROUP BY
    database_id
OPTION (RECOMPILE, MAXDOP 1);
```

Version store is a key component for multiple SQL Server features. Monitor its size and behavior, especially looking for long-running and runaway row-

versioning transactions.

Spills

As you remember from Chapter 7, Sort, Hash and Exchange operators require memory to store the data. When query memory grant is insufficient, they *spill* data to tempdb and perform the operation there. The spills impact query performance, since in-database data access is significantly slower than doing the operation in memory. It also adds to the tempdb load and may increase its size.

In recent versions of SQL Server, memory grant feedback can reduce the number of spills. However, it does not solve the problem completely. You may need to detect and optimize queries that have incorrect or excessive memory grants, using the techniques discussed in Chapter 7.

You can detect queries that spill to tempdb by using `sort_warning`, `hash_warning` and `exchange_spill` xEvents. Listing 9-10 shows the code that creates an xEvent session and the query you can use to parse the results.

Listing 9-10. Using xEvents to detect queries that spill to tempdb

```
CREATE EVENT SESSION [Spills]
ON SERVER
ADD EVENT sqlserver.hash_warning
(
    ACTION
    (
        sqlserver.database_id
        ,sqlserver.plan_handle
        ,sqlserver.session_id
        ,sqlserver.sql_text
        ,sqlserver.query_hash
        ,sqlserver.query_plan_hash
    )
    WHERE ([sqlserver].[is_system]=0)
),
ADD EVENT sqlserver.sort_warning
(
    ACTION
    (
```

```

        sqlserver.database_id
        ,sqlserver.plan_handle
        ,sqlserver.session_id
        ,sqlserver.sql_text
        ,sqlserver.query_hash
        ,sqlserver.query_plan_hash
    )
    WHERE ([sqlserver].[is_system]=0)
),
ADD EVENT sqlserver.exchange_spill
(
    ACTION
    (
        sqlserver.database_id
        ,sqlserver.plan_handle
        ,sqlserver.session_id
        ,sqlserver.sql_text
        ,sqlserver.query_hash
        ,sqlserver.query_plan_hash
    )
    WHERE ([sqlserver].[is_system]=0)
)
ADD TARGET package0.ring_buffer;
GO

-- Analyze the results
DROP TABLE IF EXISTS #tmpXML;
CREATE TABLE #tmpXML
(
    EventTime DATETIME2(7) NOT NULL,
    [Event] XML
);
DECLARE
    @TargetData XML;
SELECT
    @TargetData = CONVERT(XML,st.target_data)
FROM
    sys.dm_xe_sessions s WITH (NOLOCK)
        JOIN sys.dm_xe_session_targets st WITH(NOLOCK) ON
            s.address = st.event_session_address
WHERE
    s.name = 'Spills' and st.target_name = 'ring_buffer';
INSERT INTO #tmpXML(EventTime, [Event])
SELECT
    t.e.value('@timestamp','datetime'), t.e.query('.')
FROM
    @TargetData.nodes('/RingBufferTarget/event') AS t(e);
;WITH EventInfo

```

```

AS
(
    SELECT
        t.EventTime
        ,t.[Event].value('/event[1]/@name','sysname') AS [Event]
        ,t.
[Event].value('/event[1]/action[@name="session_id"]/value/text()
)[1]'
        , 'smallint') AS [Session]
        ,t.
[Event].value('/event[1]/action[@name="database_id"]/value/text(
))[1]'
        , 'smallint') AS [DB]
        ,t.
[Event].value('/event[1]/action[@name="sql_text"]/value/text()
)[1]'
        , 'nvarchar(max)') AS [SQL]
        ,t.[Event]

.value('/event[1]/data[@name="granted_memory_kb"]/value/text()
)[1]'
        , 'bigint') AS [Granted Memory (KB)]
        ,t.[Event]

.value('/event[1]/data[@name="used_memory_kb"]/value/text())[1]'
        , 'bigint') AS [Used Memory (KB)]
        ,t.[Event]

.value('xs:hexBinary(/event[1]/action[@name="plan_handle"]/value
/text())[1])'
        , 'varbinary(64)') AS [PlanHandle]
        ,t.
[Event].value('/event[1]/action[@name="query_hash"]/value/text()
)[1]'
        , 'nvarchar(64)') AS [QueryHash]
        ,t.[Event]

.value('/event[1]/action[@name="query_plan_hash"]/value/text()
)[1]'
        , 'nvarchar(64)') AS [QueryPlanHash]
FROM
    #tmpXML t
)
SELECT
    ei.*, qp.query_plan
FROM
    EventInfo ei

```

```

OUTER APPLY sys.dm_exec_query_plan(ei.PlanHandle) qp
OPTION (RECOMPILE, MAXDOP 1);

```

You can group the results by query hash or plan hash columns to find the queries that spill most often.

SSMS and other tools display warnings about spills in the execution plans. Figure 9-7 shows an example in SSMS. Do not ignore these warnings when you tune your queries.

The image shows a SQL Server execution plan with a Sort operator. A yellow warning icon is present above the operator, and a red arrow points to it. The operator's cost is 86%. The Properties window is open, showing the following details:

Sort	
Sort the input.	
Warnings	
Operator used tempdb to spill data during execution with spill level 1 and 1 spilled thread(s), Sort wrote 585 pages to and read 585 pages from tempdb with granted memory 8472KB and used memory 8472KB	
Warnings	Operator used
Operator used tempdb to spill	
SortSpillDetails	
GrantedMemoryKb	8472
ReadsFromTempDb	585
UsedMemoryKb	8472
WritesToTempDb	585

Figure 8-7. Sort warning in SSMS

Unfortunately, it is impossible to eliminate all spills. Make sure, however, that critical queries do not spill, and evaluate the impact of any spills on tempdb performance.

Common TempDB Issues

Issues with tempdb performance and throughput will affect the entire system. They degrade query performance, slow down T-SQL code, and introduce other issues and side effects. Unfortunately, their impact is not easy to estimate.

I usually start my analysis by looking at the general throughput and latency/stall metrics of tempdb database files using the `sys.dm_io_virtual_file_stats` view and the code from Listing 3-1. High read and write throughput indicates heavy tempdb usage. High latency may indicate issues with tempdb configuration (more on that later) or that the disk subsystem is overloaded, potentially due to high throughput.

When I see high tempdb throughput, I try to pinpoint its root cause. I look at what consumes space in tempdb by running the query from Listing 9-1. The space usage of each type of object does not always correlate with the throughput it generates, but I can often get useful information by analyzing trends in how space is used over time.

I can estimate version store throughput by looking at the Version Generation rate and Version Cleanup rate performance counters. Usually, there is little you can do about load it generates. The “solution” of switching off optimistic isolation levels is rarely appropriate. Although, remember that other SQL Server features—such as triggers, online index rebuild, and MARS—use row versioning; see if you can find any opportunities to reduce the load they generate.

There are two counters in General Statistics performance objects that help me analyze the usage of user-created temporary objects.

Temp Tables Creation Rate

Shows the count of temporary tables and table variables created per second

Active Temp Tables

Provides the number of temporary tables and table variables currently in use

High numbers in those counters indicate excessive usage of temporary objects, especially when you also see high space usage of user objects in the output from Listing 9-1. You might look for opportunities to reduce usage of temporary objects: for example, by optimizing a few frequently executed stored procedures.

You can look at the following performance counters to track how often some internal objects are created in tempdb:

SQLServer:Access Methods\ Worktables Created/sec

Provides the count of internal worktables SQL Server created to support spools, cursors, and LOB and XML variables

SQLServer:Cursor Manager By Type\Cursor worktable usage

Number of worktables used by cursors

Unfortunately, SQL Server does not provide performance counters to track spills. However, I can use an xEvent session to capture `sort_warning`, `hash_warning` and `exchange_spill` events using the `event_counter` target as shown in Listing 9-11.

Listing 9-11. Counting the number of spills

```
CREATE EVENT SESSION [Spill_Count]
ON SERVER
ADD EVENT sqlserver.exchange_spill,
ADD EVENT sqlserver.hash_warning,
ADD EVENT sqlserver.sort_warning
ADD TARGET package0.event_counter;

DECLARE
    @TargetData XML
SELECT
    @TargetData = CONVERT(XML,st.target_data)
FROM
    sys.dm_xe_sessions s WITH (NOLOCK)
    JOIN sys.dm_xe_session_targets st WITH(NOLOCK) ON
        s.address = st.event_session_address
WHERE
    s.name = 'Spill Count' and st.target_name = 'event_counter';
```

```

;WITH EventInfo
AS
(
    SELECT
        t.e.value('@name','sysname') AS [Event]
        ,t.e.value('@count','bigint') AS [Count]
    FROM
        @TargetData.nodes

    ('/CounterTarget/Packages/Package[@name="sqlserver"]/Event')
        AS t(e)
)
SELECT [Event], [Count]
FROM EventInfo
OPTION (RECOMPILE, MAXDOP 1);

```

You can achieve some benefits by optimizing queries that spill to tempdb and reducing cursor usage. Focus on the most critical and frequently executed queries and stored procedures.

That said, you can often get better ROI by scaling up your hardware. It is cheap nowadays. However, there are a couple of other common problems with tempdb that you need to be aware of.

System Pages Contention

Tempdb is a busy database – multiple sessions create and drop objects there all the time. During those processes, sessions make several changes in the system pages that track object allocations and metadata.

To protect the integrity of those pages, SQL Server serializes access to them – only one session may make modifications at any given time. This can lead to contention and reduce tempdb throughput in busy systems.

SQL Server enforces that serialization with internal objects called *latches*. They maintain the consistency of various internal data structures in SQL Server’s memory, allowing only one thread to change the object at a time. (I will discuss latches in depth in the next chapter.)

This contention presents itself with PAGELATCH waits, which are very short-lived: their duration is usually measured in microseconds.

Nevertheless, they may become very noticeable in systems with large numbers of concurrent users.

Contention during system pages modifications in tempdb is a common source of PAGELATCH waits. There are other cases, however, when you can see page latches and need to analyze wait resources to understand if those waits were tempdb-related. You can do this with the code from Listing 9-12, which captures current PAGELATCH waits using the sys.dm_os_waiting_tasks view.

Listing 9-12. Capturing currently waiting sessions

```
-- SQL Server 2005-2017
SELECT
    wt.session_id
    ,wt.wait_type
    ,er.wait_resource
    ,er.wait_time
FROM
    sys.dm_os_waiting_tasks wt WITH (NOLOCK)
    JOIN sys.dm_exec_requests er WITH (NOLOCK) ON
        wt.session_id = er.session_id
WHERE
    wt.wait_type LIKE 'PAGELATCH%'
OPTION (MAXDOP 1, RECOMPILE);
-- SQL Server 2019
SELECT
    wt.session_id
    ,wt.wait_type
    ,er.wait_resource
    ,er.wait_time
    ,pi.database_id
    ,pi.file_id
    ,pi.page_id
    ,pi.object_id
    ,OBJECT_NAME(pi.object_id,pi.database_id) as [object]
    ,pi.index_id
    ,pi.page_type_desc
FROM
    sys.dm_os_waiting_tasks wt WITH (NOLOCK)
    JOIN sys.dm_exec_requests er WITH (NOLOCK) ON
        wt.session_id = er.session_id
    CROSS APPLY
        sys.fn_PageResCracker(er.page_resource) pc
    CROSS APPLY
```

```

sys.dm_db_page_info(pc.db_id,pc.file_id
,pc.page_id,'DETAILED') pi
WHERE
    wt.wait_type LIKE 'PAGELATCH%'
OPTION (MAXDOP 1, RECOMPILE);

```

Figure 9-8 shows an example of the output. SQL Server 2019 gives you a couple of additional functions to get more detailed information. In older versions of SQL Server, you can look at wait_resource column, which shows database_id (2) as the first part of the value.

	session_id	wait_type	wait_resource	wait_time	database_id	file_id	page_id	object_id	object	index_id	page_type_desc
1	60	PAGELATCH_EX	2:1:118	3	2	1	118	34	syssschobjs	2	INDEX_PAGE
2	65	PAGELATCH_EX	2:1:307	0	2	1	307	34	syssschobjs	1	DATA_PAGE
3	69	PAGELATCH_EX	2:1:307	2	2	1	307	34	syssschobjs	1	DATA_PAGE
4	70	PAGELATCH_EX	2:1:118	4	2	1	118	34	syssschobjs	2	INDEX_PAGE
5	75	PAGELATCH_EX	2:1:118	2	2	1	118	34	syssschobjs	2	INDEX_PAGE
6	76	PAGELATCH_SH	2:1:118	3	2	1	118	34	syssschobjs	2	INDEX_PAGE
7	77	PAGELATCH_EX	2:3:830	0	2	3	830	34	syssschobjs	1	DATA_PAGE
8	78	PAGELATCH_EX	2:1:118	5	2	1	118	34	syssschobjs	2	INDEX_PAGE
9	80	PAGELATCH_EX	2:3:2052	3	2	3	2052	34	syssschobjs	1	DATA_PAGE
10	81	PAGELATCH_EX	2:3:830	6	2	3	830	34	syssschobjs	1	DATA_PAGE
11	82	PAGELATCH_EX	2:3:830	3	2	3	830	34	syssschobjs	1	DATA_PAGE
12	83	PAGELATCH_EX	2:3:2052	4	2	3	2052	34	syssschobjs	1	DATA_PAGE
13	84	PAGELATCH_EX	2:3:830	5	2	3	830	34	syssschobjs	1	DATA_PAGE

Figure 8-8. PAGELATCH waits

Because individual PAGELATCH waits are usually so short, you may miss them when you query the sys.dm_os_waiting_tasks view. You can run the query multiple times or use xEvents to track them.

Listing 9-13 shows the xEvents session to capture latch waits per database. This session will introduce overhead, so don't keep it running outside of troubleshooting. Here, to reduce overhead, I am limiting the number of events to collect.

Listing 9-13. Capturing latch waits

```
CREATE EVENT SESSION [Latch Waits] ON SERVER
ADD EVENT sqlserver.latch_suspend_end
ADD TARGET package0.ring_buffer
  (SET max_events_limit=2000);
GO
-- The code below parses collected results
DROP TABLE IF EXISTS #tmpXML;
CREATE TABLE #tmpXML
(
    EventTime DATETIME2(7) NOT NULL,
    [Event] XML
);
DECLARE
    @TargetData XML;
SELECT
    @TargetData = CONVERT(XML,st.target_data)
FROM
    sys.dm_xe_sessions s WITH (NOLOCK)
    JOIN sys.dm_xe_session_targets st WITH(NOLOCK) ON
        s.address = st.event_session_address
WHERE
    s.name = 'Latch Waits' and st.target_name = 'ring_buffer';
INSERT INTO #tmpXML(EventTime, [Event])
    SELECT t.e.value('@timestamp','datetime'), t.e.query('.')
    FROM @TargetData.nodes('/RingBufferTarget/event') AS t(e);
;WITH EventInfo
AS
(
    SELECT
        t.[EventTime] as [Time]
        ,t.
[Event].value('(//event[1]/data[@name="database_id"]/value/text())
[1]'
        , 'smallint') AS [DB]
        ,t.
[Event].value('(//event[1]/data[@name="duration"]/value/text())
[1]'
        , 'bigint') AS [Duration]
    FROM
        #tmpXML t
)
SELECT
    MONTH([Time]) as [Month]
    ,DAY([Time]) as [Day]
    ,DATEPART(hour,[Time]) as [Hour]
```

```

,DATEPART(minute,[Time]) as [Minute]
,[DB]
,COUNT(*) as [Latch Count]
,CONVERT(DECIMAL(15,3),SUM(Duration / 1000.)) as [Duration
(MS)]
FROM
    EventInfo ei
GROUP BY
    MONTH([Time]),DAY([Time]),DATEPART(hour,
[Time]),DATEPART(minute,[Time]),[DB]
ORDER BY
    [Month],[Day],[Hour],[Minute],[DB]
OPTION (RECOMPILE, MAXDOP 1);

```

There are a few things you can do in cases of tempdb system pages contention. In SQL Server 2019, you can benefit from a new feature called *memory-optimized tempdb metadata*, which converts system tables in tempdb to latch-free, non-durable memory-optimized tables.

You can enable this feature by running the ALTER SERVER CONFIGURATION SET MEMORY_OPTIMIZED TEMPDB_METADATA = ON command. You'll need to restart the server for the change to take effect. You can check if this feature is enabled by running the SELECT SERVERPROPERTY('IsTempdbMetadataMemoryOptimized') command.

There are a few limitations when this feature is enabled: most notably, you cannot create columnstore indexes on tempdb tables. Nevertheless, it is beneficial – use it unless your server has very little memory provisioned. (You can read more about the limitations of memory-optimized tempdb metadata in the Microsoft [documentation](#).)

Unfortunately, this contention may be harder to address in previous versions of SQL Server. There are a few things you can do, however. First, check tempdb configuration. Since allocations are done per file, you can reduce contention by creating multiple data files. The benefit, however, will quickly diminish as the number of files grows.

In old versions of SQL Server (prior 2016), enable the T1118 trace flag, which disables mixed-extent allocations. This will reduce the number of

changes to the system pages during allocation and deallocation operations. In the end, you may have to look at reducing tempdb usage. The less work the database performs, the less chance of contention.

Running out of Space

Running out of space in tempdb is never a good thing. It usually leads to production incidents, when a query that needs to write to the database fails. This can even impact sessions that don't explicitly use tempdb – for example, if you have optimistic isolation levels enabled, the UPDATE and DELETE statements will fail when they try to write to the version store.

You may have noticed that I keep repeating this, but for good reason: it is essential to implement free space monitoring. Depending on your configuration and maintenance practices, you might monitor free space on disk, available space in the database, or both. Set alerts to tell you when tempdb space utilization becomes critical.

If you place the tempdb log file on the same drive with data files, you might also need to look at its size and the available space there. You can do this by running the DBCC SQLPERF(LOGSPACE) command or using sys.database_files view. I'll provide additional scripts and discuss how to troubleshoot log growth issues in Chapter 11.

You have three data management views to use. The first, sys.dm_db_file_space_usage, provides information about space usage and available space in tempdb database files. Run the code from Listing 9-1 as your first step in troubleshooting.

You can follow the techniques you learned earlier in the chapter to analyze version store growth problems. Alternatively, if you see that space is consumed by users' or internal temporary objects, you may need to find the sessions that consume the most space by using two other views:

sys.dm_db_session_space_usage

The sys.dm_db_session_space_usage **view** returns the number of pages allocated and deallocated by each session for the database.

sys.dm_db_task_space_usage

The `sys.dm_db_task_space_usage` **view** provides you the allocation information for currently running tasks. That data is not reflected in `sys.dm_db_session_space_usage` view until the task is completed.

Listing 9-14 shows you code that combines the data from both views, allowing you to detect the sessions that consume the most space in tempdb.

Listing 9-14. Detecting the sessions that consume the most tempdb space

```
;WITH SpaceUsagePages
AS
(
    SELECT
        ss.session_id
        ,ss.user_objects_alloc_page_count +
          ISNULL(SUM(ts.user_objects_alloc_page_count),0)
          AS [user_alloc_page_count]
        ,ss.user_objects_dealloc_page_count +
          ISNULL(SUM(ts.user_objects_dealloc_page_count),0)
          AS [user_dealloc_page_count]
        ,ss.user_objects_deferred_dealloc_page_count
          AS [user_deferred_page_count]
        ,ss.internal_objects_alloc_page_count +
          ISNULL(SUM(ts.internal_objects_alloc_page_count),0)
          AS [internal_alloc_page_count]
        ,ss.internal_objects_dealloc_page_count +
          ISNULL(SUM(ts.internal_objects_dealloc_page_count),0)
          AS [internal_dealloc_page_count]
    FROM
        sys.dm_db_session_space_usage ss WITH (NOLOCK) LEFT JOIN
        sys.dm_db_task_space_usage ts WITH (NOLOCK) ON
        ss.session_id = ts.session_id
    GROUP BY
        ss.session_id
        ,ss.user_objects_alloc_page_count
        ,ss.user_objects_dealloc_page_count
        ,ss.internal_objects_alloc_page_count
        ,ss.internal_objects_dealloc_page_count
        ,ss.user_objects_deferred_dealloc_page_count
)
,SpaceUsage
AS
```

```

(
    SELECT
        session_id
        ,CONVERT(DECIMAL(12,3),
            ([user_alloc_page_count] - [user_dealloc_page_count])
/ 128.
        ) AS [user_used_mb]
        ,CONVERT(DECIMAL(12,3),
            ([internal_alloc_page_count] -
[internal_dealloc_page_count]) / 128.
        ) AS [internal_used_mb]
        ,CONVERT(DECIMAL(12,3),user_deferred_page_count / 128.)
            AS [user_deferred_used_mb]
    FROM
        SpaceUsagePages
)
SELECT
    su.session_id
    ,su.user_used_mb
    ,su.internal_used_mb
    ,su.user_deferred_used_mb
    ,su.user_used_mb + su.internal_used_mb AS [space_used_mb]
    ,es.open_transaction_count
    ,es.login_time
    ,es.original_login_name
    ,es.host_name
    ,es.program_name
    ,er.status as [request_status]
    ,er.start_time
    ,CONVERT(DECIMAL(21,3),er.total_elapsed_time / 1000.) AS
[duration]
    ,er.cpu_time
    ,ib.event_info as [buffer]
    ,er.wait_type
    ,er.wait_time
    ,er.wait_resource
    ,er.blocking_session_id
FROM
    SpaceUsage su
    LEFT JOIN sys.dm_exec_requests er WITH (NOLOCK) ON
        su.session_id = er.session_id
    LEFT JOIN sys.dm_exec_sessions es WITH (NOLOCK) ON
        su.session_id = es.session_id
    OUTER APPLY
        sys.dm_exec_input_buffer(es.session_id,
er.request_id) ib
WHERE
    su.user_used_mb + su.internal_used_mb >= 50

```

```
ORDER BY
    [space_used_mb] DESC
OPTION (RECOMPILE)
```

You can kill the sessions that consume the most space in tempdb to alleviate the problem. Obviously, it is a good idea to troubleshoot what caused the high space usage and address the root cause of the issue.

TempDB Configuration

Let's finish the chapter with a few tips on tempdb configuration. I've covered many aspects of it already; nevertheless, I'd like to reiterate them here.

First of all, place tempdb on the fastest drive possible. Use fast local storage (DAS) if it's an option, to give you lower latency and higher throughput compared to network-based storage (NAS). (See chapter 3.)

In non-Enterprise editions, you can even consider putting tempdb on the RAM drive if the server has enough memory available. In Enterprise edition, it is usually better to leave the memory to SQL Server.

Obviously, perform capacity planning and allocate enough storage space for tempdb workload. Set monitoring on available space and set alerts for when you are low on space. Running out of space in tempdb will lead to production incidents and potential downtime.

There is a trick you can use when you provision tempdb using fast storage with limited capacity: you can split tempdb files between fast and slow drives. Pre-allocate the files on the fast drive to the maximum possible size and disable auto-growth for them. At the same time, leave files on the slow drives very small, but keep auto-growth enabled.

In that configuration, SQL Server will favor the files on the fast drives during normal operations. However, in extreme condition, the files on the slow drive will start to grow, which helps avoid out-of-space incidents. I do not use this configuration often; however, it may be beneficial in rare

circumstances, when fast storage is big enough to handle regular tempdb workload but doesn't not have enough capacity for the spikes.

Create multiple data files. The old well-known advice that you should match the number of files with the number of CPUs may not be relevant anymore. My rule of thumb is:

- Match the number of files with number of CPUs if the server has eight or fewer CPUs.
- If server has more than eight CPU cores, use either eight data files or 1/4 of the number of cores, whichever is greater, rounding up in batches of four files. For example, use eight data files in the 24-core server and twelve data files in the 40-core server.
- Add files in the batches of four if you see any allocation contention or bottlenecks.

Use the same initial size and auto-growth parameters specified in MB for all files. If you are using SQL Server prior 2016, set the trace flag T1118 (disable mixed extents allocation) and, potentially, T1117 (grow all files in the filegroup simultaneously). Both of those flags are server-wide and impact user databases. This is not a problem with the T1118 flag; however, do not enable T1117 if users' databases have unevenly sized data files.

Finally, if you are using SQL Server 2019 or above, consider enabling memory-optimized catalogs in tempdb, especially if you see PAGELATCH waits over tempdb allocation pages.

Summary

Tempdb is a busy database that is shared among all user and system sessions. High tempdb performance and throughput are essential for good system performance.

Place tempdb on the fastest drive you have. Make sure that the database configuration is correct. Create multiple data files using the same auto-

growth parameters specified in MB. Enable the T1118 and, potentially, T1117 flags in SQL Server versions prior to 2016.

Analyze PAGELATCH waits if you see them and detect if they are tempdb-related. In SQL Server 2019, you can mitigate them by enabling memory-optimized tempdb catalogs. In older versions of SQL Server, analyze the number of tempdb files, use temporary object caching, and reduce tempdb usage.

Troubleshoot and address excessive tempdb load. Refactor code that does not use temporary objects legitimately, and optimize queries with excessive spills.

Be aware that table variables do not store statistics; doing so leads to cardinality estimation errors and inefficient query plans. Consider using temporary tables as a safer choice.

In the next chapter, I'll talk about latches – synchronization objects used by SQL Server to protect internal in-memory structures.

Troubleshooting Checklist

Check tempdb configuration (number of files, auto-growth settings, T1118 in old versions of SQL Server, etc.).

Analyze performance of tempdb storage.

Review tempdb space consumption and usage. Address issues if feasible.

Troubleshoot if PAGELATCH waits are tempdb-related.

Consider enabling memory-optimized tempdb metadata in SQL Server 2019, especially if you notice tempdb-related PAGELATCH waits.

Setup monitoring on tempdb drive space utilization.

Chapter 9. Latches

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 10th chapter of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at dmitri@aboutsqlserver.com.

Latches are lightweight synchronization objects that protect the consistency of SQL Server internal data structures. As the opposite of locks, which protect transactional data consistency, latches prevent corruption of the data structures in memory.

In most cases, latches are short-lived and may be unnoticeable in systems with light loads. However, as loads grow, latch contention may become an issue and can limit system scalability and throughput. In this chapter, I will discuss how to detect and address those situations.

I will start this chapter with an overview of latches and their categories and types. Next, I will discuss page latches and several mitigation techniques to address their contention. Finally, I will cover other ways to deal with common latch types.

Introduction to Latches

There is a concept in computer science called *mutual exclusion*, which means that multiple threads or processes cannot execute critical code simultaneously. Think about the multi-threaded application in which threads use the shared objects. In those systems, you often need to serialize the code that accesses those objects to avoid creating race conditions when multiple threads read and update them simultaneously.¹

Internally, SQL Server enforces mutual exclusion and protects in-memory data structures by using latches. People are often confused by latches' similarities to locks: both types of objects affect concurrency and may prevent simultaneous access to the same data. There is a subtle difference, however.

SQL Server uses locks to enforce *logical* data consistency, preventing sessions from working with transactionally inconsistent data. Latches, on the other hand, enforce the *physical* consistency of in-memory data structures, preventing corruption by multiple workers that access them simultaneously. Latches do not work in transaction boundaries; they are acquired when the worker needs to access an in-memory object and released after the operation is done.

Consider a situation where multiple sessions need to update different rows on the same data page. Those sessions would not block each other with locks, as long as they don't acquire incompatible locks on the same rows. However, they might block each other with latches to prevent simultaneous updates of the data page object in-memory, which can corrupt it.

There are five different latch types in SQL Server. In terms of compatibility, they are similar to locks, but don't confuse the two – they are different beasts!

Keep latch (KP)

Keep latch (KP) ensures that the referenced structure cannot be destroyed. It is compatible with all other latch types except the destroy (DT) latch. This latch type has similarities to schema stability (Sch-S) locks.

Shared latch (SH)

A shared latch (SH) is required when a thread needs to read the data structure. Shared (SH) latches are compatible with each other and with the keep (KP) and update (UP) latches.

Update latch (UP)

Update latches (UP) allow other threads to read the structure but prevents them from updating it. SQL Server uses them in some scenarios to improve concurrency, similar to update (U) locks. Update (UP) latches are compatible with keep (KP) and shared (SH) latches and incompatible with all other types.

Exclusive latch (EX)

An exclusive latch (EX) is required when a thread modifies the data structure. Conceptually, exclusive (EX) latches are similar to exclusive (X) locks; they are incompatible with all other latch types except keep (KP) latches.

Destroy latch (DT)

A destroy latch (DT) is required to destroy a data structure. For example, a destroy (DT) latch would be acquired when the lazy writer process removes a data page from the buffer pool. These latches are incompatible with other latch types.

When the worker cannot obtain a latch on the data structure, it becomes suspended and generates a latch-related wait type. Those wait types can belong to one of three categories:

PAGEIOLATCH

The PAGEIOLATCH waits indicate I/O related latches. SQL Server uses these latches and wait types while waiting for data pages to be read from disk to the buffer pool. A large percentage of such waits could indicate a large number of non-optimized queries and/or suboptimal

disk subsystem performance. Chapter 3 covers how to troubleshoot those conditions.

PAGELATCH

PAGELATCH waits indicate buffer-pool-related latches, which occur when threads need to access or modify data and allocation map pages in the buffer pool. As you know from chapter 9, those waits can be triggered by contention in tempdb system catalogs. I'll talk about other conditions that can trigger them later in this chapter.

LATCH

All other non-buffer-pool-related latches. I'll talk about them later in the chapter.

SQL Server uses different wait types for each latch type. For example, PAGELATCH_EX indicates the wait for an exclusive (EX) latch on the buffer pool page. The LATCH_SH wait type shows a shared (SH) non-buffer pool latch wait.

Let's look at the latch categories in depth.

Page Latches

SQL Server uses page latches to protect the consistency of the buffer pool pages in memory. When a worker needs to change anything on the page, it obtains an exclusive (EX) page latch. Similarly, when a worker needs to read something the page, it acquires a shared (SH) page latch. The workers can read the data simultaneously; however, only one worker can modify the page at any given time, and only when no other workers are accessing it.

The impact of page latches greatly depends on the system load and workload patterns. SQL Server reads data from and writes data to in-memory pages very quickly; you might not notice it in a system with a light load and a small number of concurrent users. However, page latch

contention increases with the number of active sessions simultaneously accessing the same data pages.

There are two common reasons for page latches. They are not mutually exclusive and both can contribute to a high percentage of PAGELATCH waits on the server. You need to address both when troubleshooting.

The first is related to tempdb. A heavy concurrent tempdb workload may introduce latch contention on system pages (chapter 9 discusses how to detect and mitigate that condition).

The second most common case for page latching usually occurs in users' tables. It is called the *hotspot*. Hotspots usually occur in indexes with ever-increasing (or ever-decreasing) values, such as identities, sequences or datetime columns that populate with the time when rows are inserted. When multiple sessions simultaneously insert data into those indexes, all rows are placed on the same page (the last one), which leads to page latch contention as workers start blocking each other.

Let's look at a hypothetical scenario in which you want to log application-request information in the database. Listing 10-1 creates a few tables to store the data.

NOTE

Do not consider this implementation a real-life example. In fact, relational databases are usually the worst place to store logs. I am using it only for the sake of demonstrating page latch contention.

Listing 10-1. Creating tables to store the logs

```
CREATE TABLE dbo.WebRequests_Disk
(
    RequestId INT NOT NULL identity(1,1),
    RequestTime DATETIME2(4) NOT NULL
        CONSTRAINT DEF_WebRequests_Disk_RequestTime
            DEFAULT SYSUTCDATETIME(),
    URL VARCHAR(255) NOT NULL,
    RequestType TINYINT NOT NULL,
```

```

ClientIP VARCHAR(15) NOT NULL,
BytesReceived INT NOT NULL,

CONSTRAINT PK_WebRequests_Disk
PRIMARY KEY NONCLUSTERED(RequestID)
ON [LOGDATA]
);
CREATE UNIQUE CLUSTERED INDEX
IDX_WebRequests_Disk_RequestTime_RequestId
ON dbo.WebRequests_Disk(RequestTime,RequestId)
ON [LOGDATA];
CREATE TABLE dbo.WebRequestHeaders_Disk
(
    RequestId INT NOT NULL,
    HeaderName VARCHAR(64) NOT NULL,
    HeaderValue VARCHAR(256) NOT NULL,

    CONSTRAINT PK_WebRequestHeaders_Disk
    PRIMARY KEY CLUSTERED(RequestID,HeaderName)
    ON [LOGDATA]
);
CREATE TABLE dbo.WebRequestParams_Disk
(
    RequestId INT NOT NULL,
    ParamName VARCHAR(64) NOT NULL,
    ParamValue nVARCHAR(256) NOT NULL,

    CONSTRAINT PK_WebRequestParams_Disk
    PRIMARY KEY CLUSTERED(RequestID,ParamName)
    ON [LOGDATA]
);

```

Next, I will run the application, which will insert data into those tables from multiple threads simultaneously. You can get the application code from the book's companion material.

Figure 10-1 shows the metrics from my test server when the application was running. Although the duration of individual latch waits was very short—just a fraction of a millisecond—cumulatively they contributed to very large numbers and limited application throughput. CPU load had *not* been maxed out when that happened.

\\SQL2019-1

Processor Information **_Total**
% Processor Time 66.646

SQLServer:Latches
Average Latch Wait Time (ms) 0.294
Latch Waits/sec 41,990.208

SQLServer:SQL Statistics
Batch Requests/sec 5,049.916

Figure 9-1. Fig. 10-1. Performance metrics during page latch contention

Figure 10-2 shows wait statistics, obtained with the code in Listing 2-1. As you can see, PAGELATCH_EX and PAGELATCH_SH waits contribute to more than 60% of total waits.

	Wait Type	Wait Count	Wait Time	Avg Wait Time	Avg Signal Wait Time	Avg Resource Wait Time	Percent	Running Percent
1	PAGELATCH_EX	1951297	749.168	0.0	0.0	0.0	45.332	45.332
2	WRITELOG	389893	488.770	1.0	0.0	1.0	29.575	74.908
3	PAGELATCH_SH	672816	303.111	0.0	0.0	0.0	18.341	93.249

Figure 9-2. Fig. 10-2. Wait statistics with page latch contention

You can detect the indexes that introduce the most page latch waits using the `sys.dm_db_index_operational_stats` **function**. As you can guess by the name, this function tracks indexes' operational metrics, including the number of latches and their wait times.

The code in Listing 10-2 detects indexes that contribute to hotspots. This is a very simplified implementation, however; I'll provide a more sophisticated query to analyze indexes' health and usage in Chapter 14.

Listing 10-2. Analyzing page latch index statistics

```
SELECT
    s.name + '.' + t.name as [table]
    ,i.index_id
    ,i.name as [index]
    ,SUM(os.page_latch_wait_count) AS [latch count]
    ,SUM(os.page_latch_wait_in_ms) as [latch wait (ms)]
FROM
    sys.indexes i WITH (NOLOCK) JOIN sys.tables t WITH (NOLOCK)
on
    i.object_id = t.object_id
JOIN sys.schemas s WITH (NOLOCK) ON
    t.schema_id = s.schema_id
CROSS APPLY
    sys.dm_db_index_operational_stats
    (
        DB_ID()
        ,t.object_id
        ,i.index_id
        ,0
    ) os
GROUP BY
    s.name, t.name, i.name, i.index_id
ORDER BY
    SUM(os.page_latch_wait_in_ms) DESC;
```

Figure 10-3 shows the output from the code. As you can see, the data allows you to detect problematic indexes in the database quickly.

	table	index_id	index	latch count	latch wait (ms)
1	dbo.WebRequestHeaders_Disk	1	PK_WebRequestHeaders_Disk	2230279	933152
2	dbo.WebRequestParams_Disk	1	PK_WebRequestParams_Disk	356426	112319
3	dbo.WebRequests_Disk	1	IDX_WebRequests_Disk_Requ...	36621	7568
4	dbo.WebRequests_Disk	2	PK_WebRequests_Disk	984	722

Figure 9-3. Fig. 10-3. Output from sys.dm_db_index_operational_stats function

Addressing Hotspots: OPTIMIZE_FOR_SEQUENTIAL_KEY Index Option

Unfortunately, addressing page latching due to hotspots is not an easy task, especially in old versions of SQL Server. In SQL Server 2019 and above, you can enable the `OPTIMIZE_FOR_SEQUENTIAL_KEY` index property. This setting will improve throughput in hotspot scenarios, but will not solve the problem completely.

Let's enable that setting, as shown in Listing 10-3, and repeat the test. I also recommend clearing wait statistics before restarting the application.

Listing 10-3. Enable `OPTIMIZE_FOR_SEQUENTIAL_KEY` option (SQL Server 2019 and above)

```
ALTER INDEX PK_WebRequestHeaders_Disk
ON dbo.WebRequestHeaders_Disk
SET (OPTIMIZE_FOR_SEQUENTIAL_KEY = ON);
ALTER INDEX PK_WebRequestParams_Disk
ON dbo.WebRequestParams_Disk
SET (OPTIMIZE_FOR_SEQUENTIAL_KEY = ON);
```

Figure 10-4 shows performance metrics when the `OPTIMIZE_FOR_SEQUENTIAL_KEY` is enabled. As you can see, that option improves throughput; however, the latch waits are still significant.

\\SQL2019-1

Processor Information	_Total
% Processor Time	83.373
SQLServer:Latches	
Average Latch Wait Time (ms)	0.077
Latch Waits/sec	15,173.777
SQLServer:SQL Statistics	
Batch Requests/sec	8,326.207

Figure 9-4. Fig.10-4. Performance metrics with OPTIMIZE_FOR_SEQUENTIAL_KEY enabled

Figure 10-5 shows the wait statistics. You can see the new wait type, BTREE_INSERT_FLOW_CONTROL, in the output. This wait type is specific to OPTIMIZE_FOR_SEQUENTIAL_KEY implementation and, for all practical purposes, masks PAGELATCH waits and latch contention in wait statistics.

	Wait Type	Wait Count	Wait Time	Avg Wait Time	Avg Signal Wait Time	Avg Resource Wait Time	Percent	Running Percent
1	WRITELOG	609044	1020.615	1.0	0.0	1.0	60.430	60.430
2	BTREE_INSERT_FLOW_CONTROL	1477495	554.722	0.0	0.0	0.0	32.845	93.275
3	PAGELATCH_EX	1128697	99.052	0.0	0.0	0.0	5.865	99.140

Figure 9-5. Fig. 10-5. Wait statistics with OPTIMIZE_FOR_SEQUENTIAL_KEY enabled

While enabling OPTIMIZE_FOR_SEQUENTIAL_KEY may help improve throughput in hotspot situations, it does not completely eliminate latching. Moreover, this option is not available in older versions of SQL Server (prior

to 2019). In those cases, you have very few options besides refactoring database schema and applications.

You need to detect and analyze indexes that contribute the most to latching. Review index usage (I'll share a few helpful techniques in Chapter 14) and determine if problematic indexes can be dropped or altered in ways that prevent them from constantly increasing and contributing to hotspots.

NOTE

Clustered indexes defined on identity columns in active tables are one of the most common places for hotspots. To address that, use clustered indexes that spread insert activity across the table.

Addressing Hotspots: Hash Partitioning

One technique that allow you to spread inserts across a table is called *hash partitioning*. You can partition the table and use random hash values to spread rows across multiple partitions.

Listing 10-4 shows such an example. It redefines two of the tables from Listing 10-1, partitioning them with a new HashVal column calculated as CHECKSUM(RequestId). This randomly distributes the data across multiple partitions, reducing insertion rates and latch contention on each individual partition.

Note that the HashVal column is defined as the rightmost column in the indexes, to preserve the sorting order on each individual partition.

Listing 10-4. Implementing hash partitioning

```
-- For demo purposes
TRUNCATE TABLE dbo.WebRequests_Disk;
DROP TABLE dbo.WebRequestHeaders_Disk;
DROP TABLE dbo.WebRequestParams_Disk;
GO
CREATE PARTITION FUNCTION pfHash(int)
AS RANGE LEFT FOR VALUES
(-1847483647, -1547483647, -1247483647, -947483647, -647483647, -34748
```

```

3647
,-47483647,252516353,552516353,852516353,1152516353,1452516353,17
52516353);
CREATE PARTITION SCHEME psHash
AS PARTITION pfHash
ALL TO ([LOGDATA]);
CREATE TABLE dbo.WebRequestHeaders_Disk
(
    RequestId INT NOT NULL,
    HeaderName VARCHAR(64) NOT NULL,
    HeaderValue VARCHAR(256) NOT NULL,
    HashVal AS CHECKSUM(RequestId) PERSISTED,
    CONSTRAINT PK_WebRequestHeaders_Disk
    PRIMARY KEY CLUSTERED(RequestID,HeaderName,HashVal)
    ON psHash(HashVal)
);
CREATE TABLE dbo.WebRequestParams_Disk
(
    RequestId INT NOT NULL,
    ParamName VARCHAR(64) NOT NULL,
    ParamValue nVARCHAR(256) NOT NULL,
    HashVal AS CHECKSUM(RequestId) PERSISTED,

    CONSTRAINT PK_WebRequestParams_Disk
    PRIMARY KEY CLUSTERED(RequestID,ParamName,HashVal)
    ON psHash(HashVal)
);

```

Figure 10-6 shows the performance metrics on my test server with hash partitioning implemented. The insert throughput is better than without the partitioning; however, using the `OPTIMIZE_FOR_SEQUENTIAL_KEY` approach would outperform it. Your mileage may vary, of course, and results may differ in other use cases and workloads.

\\SQL2019-1

Processor Information	_Total
% Processor Time	70.444
SQLServer:Latches	
Average Latch Wait Time (ms)	0.219
Latch Waits/sec	32,993.055
SQLServer:SQL Statistics	
Batch Requests/sec	7,379.357

Figure 9-6. Fig.10-6. Performance metrics with hash partitioning

Figure 10-7 shows the wait statistics. As you see, the page latch wait percentage is lower comparing than without the partitioning (shown in Figure 10-2).

	Wait Type	Wait Count	Wait Time	Avg Wait Time	Avg Signal Wait Time	Avg Resource Wait Time	Percent	Running Percent
1	WRITELOG	331696	722.286	2.0	0.0	1.0	66.211	66.211
2	PAGELATCH_EX	1248937	291.905	0.0	0.0	0.0	26.759	92.970
3	PAGELATCH_SH	246413	68.703	0.0	0.0	0.0	6.298	99.268

Figure 9-7. Fig. 10-7. Wait statistics with hash partitioning

While hash partitioning may help reduce latch contention, it is dangerous. As with any partitioning, the data will be spread across multiple internal physical tables and partitions. This will change execution plans and may lead to regressions in query performance. Test the system carefully when you implement this.

In some cases, you can reduce the side effects of hash partitioning by using staging tables. In that scenario, the application may insert data into hash-partitioned staging tables and have another process running on schedule, copying data from staging to main tables. While this approach adds some overhead from staging tables management, it protects against the unanticipated performance regressions that hash partitioning can introduce.

Addressing Hotspots: In-Memory OLTP

It is impossible to discuss performance issues introduced by latching without mentioning In-Memory OLTP. After all, one of the main goals of that technology is addressing latching and locking challenges with disk-based tables under heavy loads. Memory-optimized tables are latch- and lock-free and would scale perfectly under a concurrent OLTP workload.

Those benefits don't come for free, though. In-Memory OLTP and memory-optimized tables behave differently than classic Storage Engine and disk-based tables. You need to design the system properly and utilize technology to avoid possible side effects and issues. As with hash partitioning, you may get performance regressions if you simply switch disk-based tables to become memory-optimized without taking different technology behavior into consideration.

Fortunately, memory-optimized tables are perfect candidates for staging tables. You can completely eliminate hotspots and latch contention and reduce any possible side effects In-Memory OLTP may introduce. Consider it a valid implementation option to solve hotspot problems.

I will repeat my word of caution, though: you need to know how to deploy and maintain In-Memory OLTP properly before using it. You can read my book *Expert SQL Server In-Memory OLTP* to learn more.

Other Latch Types

The third latch category, non-buffer-pool latches, presents itself with generic LATCH wait types. Similar to wait statistics, you can get

information about individual latch types using the `sys.dm_os_latch_stats` view shown in Listing 10-5. You can also clear latch statistics with the command `DBCC SQLPERF('sys.dm_os_latch_stats', CLEAR)`.

Listing 10-5. Analyzing latch statistics

```
;WITH Latches
AS
(
    SELECT
        latch_class, wait_time_ms, waiting_requests_count
        ,100. * wait_time_ms / SUM(wait_time_ms) OVER() AS Pct
        ,100. * SUM(wait_time_ms) OVER(ORDER BY wait_time_ms
DESC) /
            NULLIF(SUM(wait_time_ms) OVER(), 0) AS RunningPct
        ,ROW_NUMBER() OVER(ORDER BY wait_time_ms DESC) AS RowNum
    FROM
        sys.dm_os_latch_stats WITH (NOLOCK)
    WHERE
        wait_time_ms > 0 AND
        latch_class NOT IN (N'BUFFER',N'SLEEP_TASK')
)
SELECT
    l1.latch_class AS [Latch Type]
    ,l1.waiting_requests_count AS [Latch Count]
    ,CONVERT(DECIMAL(12,3), l1.wait_time_ms / 1000.0)
        AS [Wait Time]
    ,CONVERT(DECIMAL(12,1), l1.wait_time_ms /
l1.waiting_requests_count)
        AS [Avg Wait Time]
    ,CONVERT(DECIMAL(6,3), l1.Pct)
        AS [Percent]
    ,CONVERT(DECIMAL(6,3), l1.RunningPct)
        AS [Running Percent]
FROM
    Latches l1
WHERE
    l1.RunningPct <= 99 OR l1.RowNum = 1
ORDER BY
    l1.RunningPct
OPTION (RECOMPILE, MAXDOP 1);
```

Figure 10-8 shows the output from one of the servers.

	Latch Type	Latch Count	Wait Time	Avg Wait Time	Percent	Running Percent
1	ACCESS_METHODS_DATASET_PARENT	3858110729	1249807.056	0.0	47.031	47.031
2	NESTING_TRANSACTION_FULL	766539824	972751.713	1.0	36.605	83.635
3	TRACE_CONTROLLER	38694619	270191.197	6.0	10.167	93.803

Figure 9-8. Fig.10-8. Latch statistics

Unfortunately, latch types are poorly documented. You often need to search multiple sources to understand the meaning of each latch. Nevertheless, let's look at a few common types you will encounter.

Parallelism-related latches

There are multiple parallelism-related latch types. The most common are ACCESS_METHOD_DATASET_PARENT, ACCESS_METHODS_SCAN_RANGE_GENERATOR, ACCESS_METHODS_SCAN_KEY_GENERATOR, and NESTING_TRANSACTION_FULL. In my experience, those latches usually appear at the top of the list in the sys.dm_os_latch_stats view output.

Treat these latches as you would parallelism waits (CXPACKET, CXCONSUMER and EXCHANGE), which often appear alongside parallelism-related latches. Analyze and tune parallelism usage to address this, as discussed in Chapter 6.

LOG_MANAGER

The LOG_MANAGER latch type indicates growth in the transaction log. You may see that latch in situations where *something* is regularly preventing the transaction log from being truncated. Monitor the log_reuse_wait_desc column in the sys.databases view to troubleshoot it (more in the next chapter).

I have also encountered this latch type in systems that regularly shrink the transaction log after each log backup. This is a bad practice, since the log file is zero-initializing at the time of growth.

ACCESS_METHODS_HOBT_VIRTUAL_ROOT

The ACCESS_METHODS_HOBT_VIRTUAL_ROOT latch is used during access to the index metadata. Significant presence of this latch indicates a large number of page splits of the root pages in B-Tree indexes. Usually this happens in small indexes with very volatile data, for example tables that operate as in-database queues.

You can detect potentially problematic indexes with tree_page_latch_wait_count and tree_page_latch_wait_time_ms columns in the sys.db_db_index_operational_stats view (more in Chapter 14).

ACCESS_METHODS_HOBT_COUNT

The ACCESS_METHODS_HOBT_COUNT latch is used to update page and row count information in table metadata. Contention with this latch indicates lots of small concurrent data modifications in some tables.

Chapter 14 will discuss an approach you can use to detect such tables; however, it may be easier to work with the development team instead of using SQL Server views, since addressing this latch contention usually requires application changes.

FGCB_ADD_REMOVE

The FGCB_ADD_REMOVE latch occurs during adding, removing, growing and shrinking files in the filegroup. Check if *Instant File Initialization* is enabled and *Auto Shrink* database option is disabled. Also check that auto-growth parameters are not growing files in very small chunks.

You may also see that latch is the databases with very large number of data files, especially in tempdb. Re-evaluate database configuration to address it.

TRACE_CONTROLLER

As you can guess by the name, the TRACE_CONTROLLER latch indicates large number of traces . Evaluate monitoring strategy and remove unnecessary monitoring when you see that.

Similar to waits, latches are unavoidable. They are completely acceptable in small numbers; however, excessive latching is usually a sign of the issues in the system. Don't jump to immediate action, though – identify and address the root causes!

Summary

Latches are lightweight synchronization objects that protect the consistency of SQL Server internal data structures. They are usually short-lived and may be unnoticeable in systems with light loads. However, as loads grow, latch contention may become an issue, limiting system scalability and throughput.

There are three categories of latches and wait types associated with them; each latch type (shared (SH), exclusive (EX), etc.) has corresponding wait type in wait statistics.

PAGEIOLATCH waits occur when SQL Server waits for the data page to be read to the buffer pool. A significant percentage of those waits require you to troubleshoot disk subsystem load (see Chapter 3).

Page latches and corresponding PAGELATCH waits occur when multiple workers are simultaneously accessing and updating data pages in memory. They are usually triggered by tempdb system object contention or by hotspots in ever-increasing indexes in users' databases.

In SQL Server 2019 and above, you can reduce the impact of hotspots by enabling the `OPTIMIZE_FOR_SEQUENTIAL_KEY` index option. However, in many cases, you need to drop or change the index, or implement workarounds with hash partitioning and/or staging tables.

Non-buffer-pool-related latches present themselves with LATCH wait types. You can obtain latch statistics with the `sys.dm_os_latch_stats` view. Identify and address the root cause of these issues when you troubleshoot them.

Now it is time to look at transaction log problems you may encounter.

Troubleshooting Checklist

Analyze the impact of page latches with PAGELATCH waits.

Detect if PAGELATCH waits are coming from tempdb or hotspots in users' databases.

Address tempdb system object contention (chapter 9).

Identify indexes that contribute to hotspots with `sys.dm_db_index_operational_stats` view. Analyze if indexes can be refactored or dropped (more in Chapter 14).

Enable `OPTIMIZE_FOR_SEQUENTIAL_KEY` option in SQL Server 2019 and above

Consider refactoring the application using In-Memory OLTP, staging tables, and/or hash partitioning if neither of the other options helps with hotspots.

Review latch statistics with `sys.dm_os_latch_stats` view; troubleshoot and address issues if needed.

¹ Every development language has a set of synchronization primitives: for example, mutexes and critical sections. In T-SQL, you can use application locks to serialize access to some code. I am not covering them in this book; however, you can read about them in my Expert SQL Server Transactions and Locking book or in the Microsoft [Documentation](#).

Chapter 10. Transaction Log

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be Chapter 11 of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the editor at sgrey@oreilly.com.

Every database in SQL Server has one or more transaction log files in addition to data files. The transaction log stores information about the changes made in the database and allows SQL Server to recover databases to transactionally consistent states in case of an unexpected shutdown or crash. Every data modification in the database is stored there, and high log file throughput is essential for good system performance.

In this chapter, I’ll explain how SQL Server logs transactions and how the transaction log works internally. Next, I’ll cover several best practices for transaction log configuration and talk about how to address “Transaction log full” situation. Finally, I’ll discuss how to troubleshoot insufficient transaction log throughput.

Transaction Log Internals

SQL Server uses a transaction log to keep each database in a *transactionally consistent* state, meaning that data modifications done from

within transactions must either be committed or rolled back in full. SQL Server never allows data to be transactionally inconsistent by applying just a subset of the changes from uncommitted transactions.

The transaction log guarantees consistency. It stores the stream of *log records* generated by data modifications and some internal operations. Every log record has a unique, auto-incrementing *Log Sequence Number* (LSN) and describes the data change. It includes information about the affected row, the old and new versions of the data, the transaction that performed the modification, and so forth.

Every data page keeps the LSN of the last log record that modified it. During recovery, SQL Server can compare the LSNs of the log records and data pages and find out if the most recent changes were saved to the data files. There is enough information stored in a log record to undo or redo the operation if needed.

SQL Server uses *Write-Ahead Logging (WAL)*, which guarantees that log records are always written to the log file *before* dirty data pages are saved to the database. Sharp-eyed readers may notice that in Chapter 3, I mentioned that log records are saved synchronously with data modifications, while data pages are saved asynchronously during the checkpoint process. While that is conceptually correct, I will be more precise here: SQL Server caches the log records in small memory caches called *log buffers*, writing them in batches to reduce the number of log-write I/O operations.

Each database has its own log buffer, consisting of 60KB structures called *log blocks*. Each log buffer (and database) can have up to 128 log blocks. SQL Server writes log blocks to the log file in a single I/O operation. However, it does not always wait until the log block is full. The typical size of a log-writing I/O operation varies, from 512 bytes to 60KB.

Unfortunately, the SQL Server documentation is inconsistent in its terminology, and often references *log blocks* as *log buffers*. Just remember that SQL Server caches the log records in memory before writing them on disk.

Data Modifications and Transaction Logging

Let's look at how SQL Server modifies data in more detail. Figure 11-1 shows a database with an empty log buffer and transaction log. The last transaction in the log has an LSN of 7314.

Let's assume that there are two active transactions: T1 and T2. The BEGIN TRAN log records for both of those transactions have already been saved in the log and are not shown in the diagram.

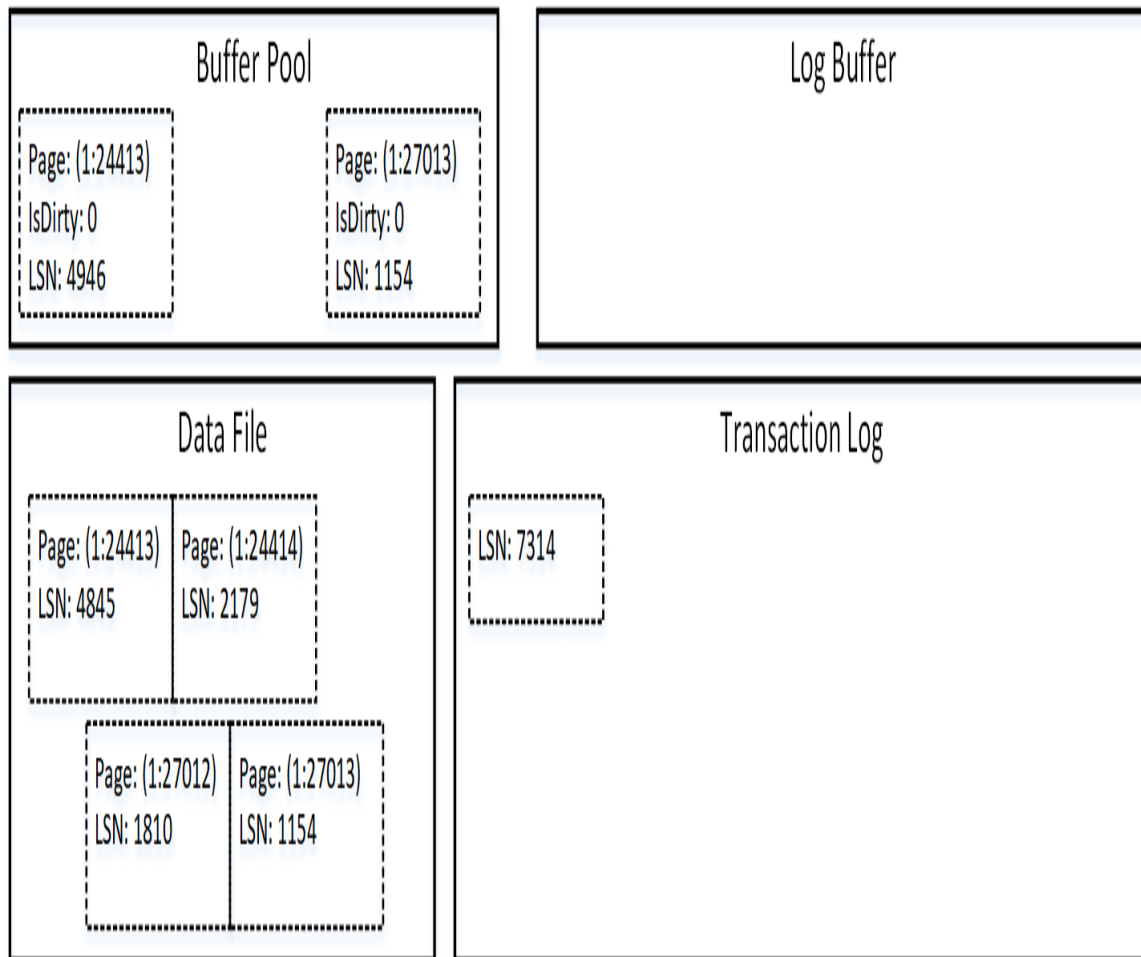


Figure 10-1. Data modifications and transaction logging: Initial state

Let's assume that transaction T1 updates one of the rows on page (1:24413). This operation generates a new log record, which will be placed into the log buffer. It will also update the data page, marking it as dirty and changing the LSN in the page header. Figure 11-2 illustrates that.

At this point, the log record has not been saved to the log file (this is often called *hardened*). It does not introduce any issues as long as the data page has not been saved in the data file. In event of an SQL Server crash, both the log record and the modifications on the data page will be gone—which is fine, because the transaction has not been committed.

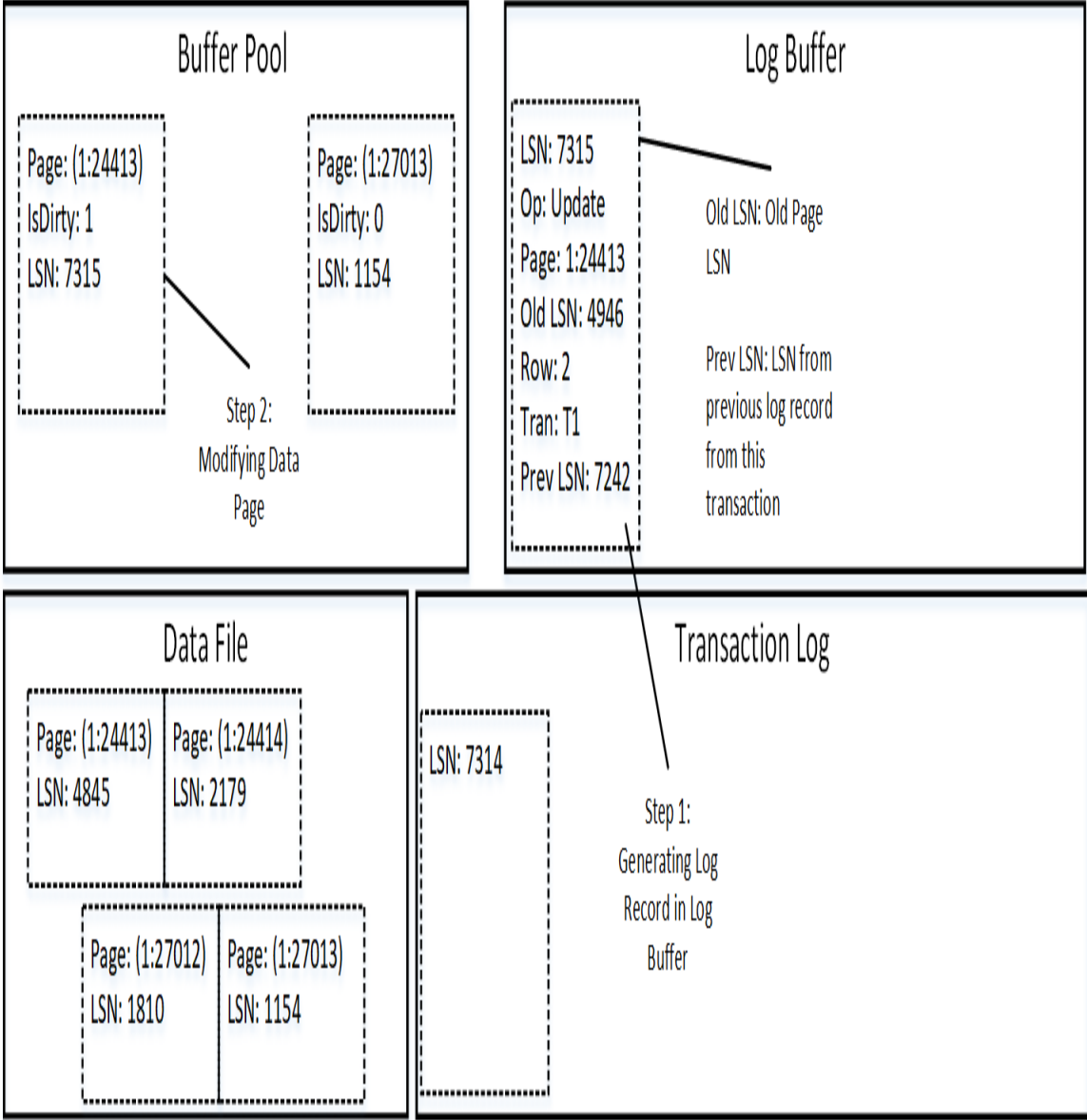


Figure 10-2. Data modifications and transaction logging: State after the first update

Next, let's assume that transaction T2 inserts a new row into page (1:27013), while transaction T1 deletes another row on the same page. Those operations generate two log records. These are placed into the log

buffer, as shown in Figure 11-3. Right now, all log records are still in the log buffer.

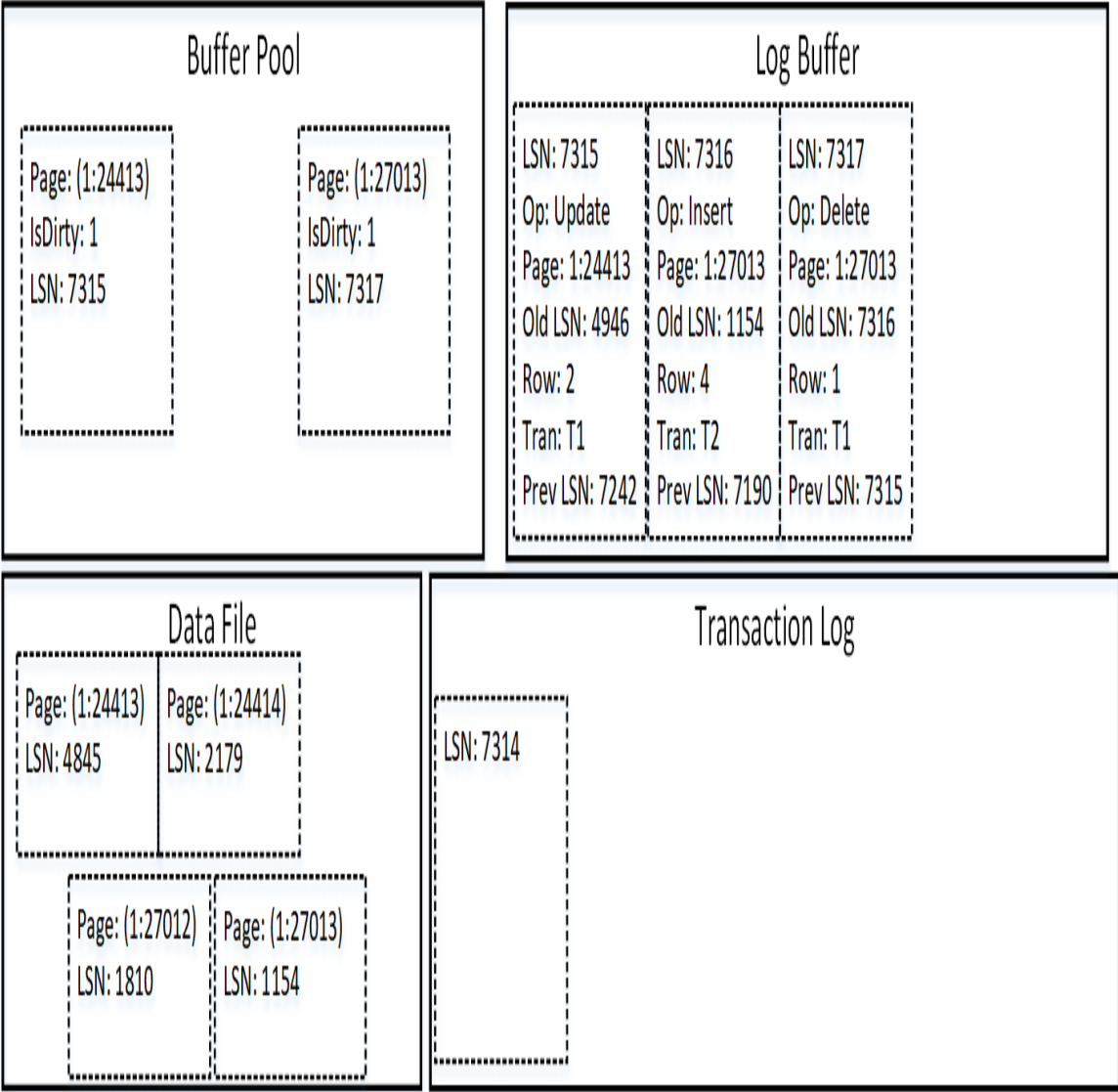


Figure 10-3. Data modifications and transaction logging: State after two data modifications

Now let's assume that an application commits transaction T2. This action generates a COMMIT log record and forces SQL Server to write (harden) the content of the log block to disk. It writes *all* log records from the buffer to disk, regardless of what transaction generated them (Figure 11-4).

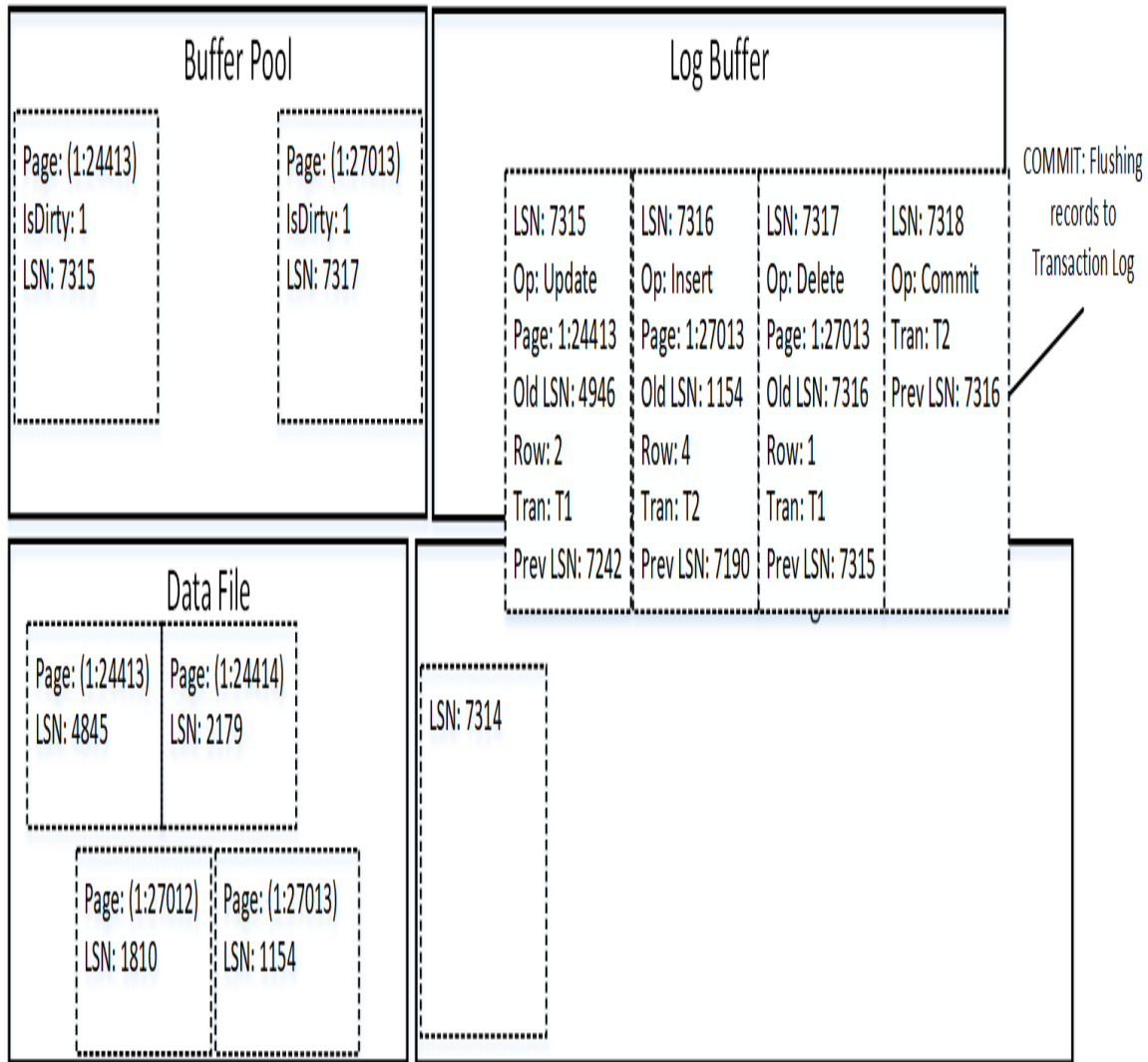


Figure 10-4. Data modifications and transaction logging: Commit operation

The applications receive confirmation that the transaction has been committed *only* after all log records are hardened. Even though the data page (1:27013) is still dirty and has not been saved into the data file, the hardened log records on the disk have enough information to re-apply the changes made by the committed T2 transaction if needed.

The dirty pages from the buffer pool will be saved to data files on the checkpoint. This operation also generates a CHECKPOINT log record and immediately hardens it into the log. Figure 11-5 shows that state.

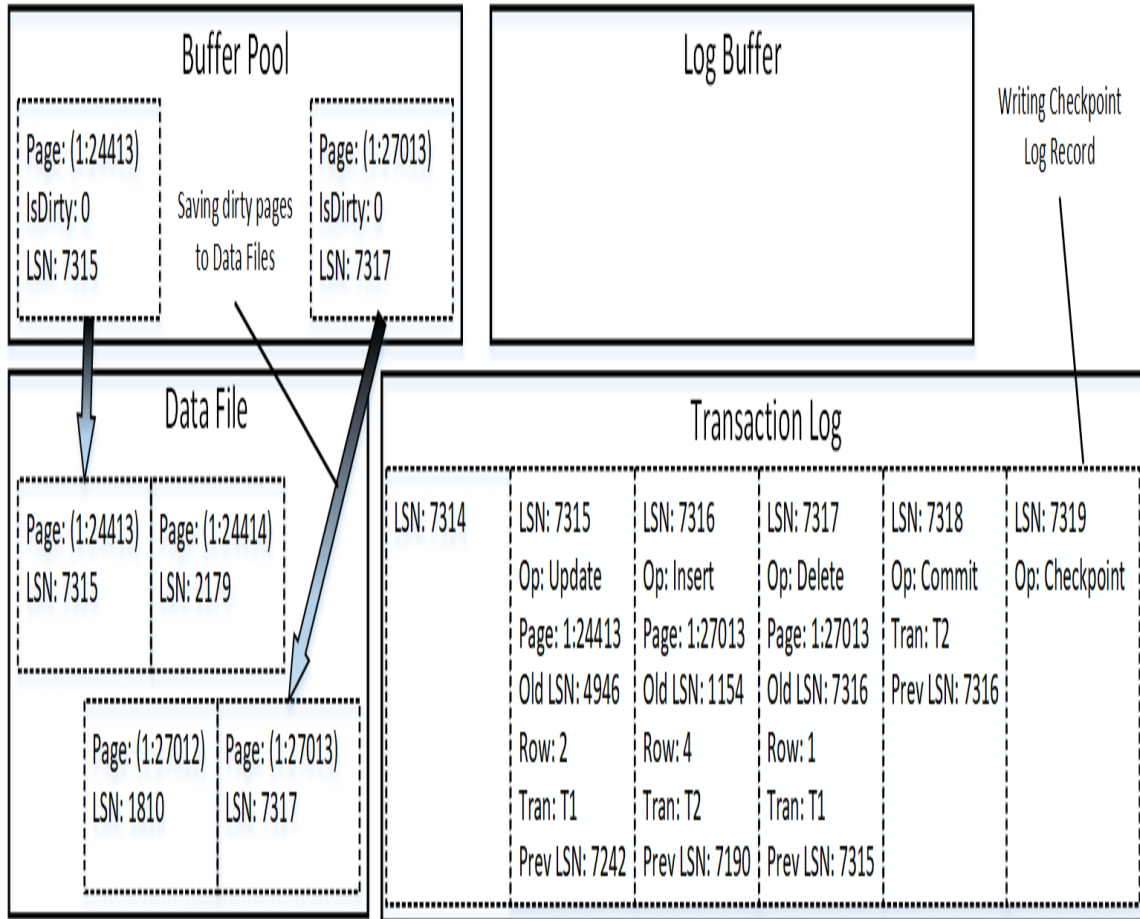


Figure 10-5. Data modifications and transaction logging: Checkpoint

After the checkpoint occurs, the pages in the data file may store data from uncommitted transactions (T1, in our example). However, log records in the transaction log include enough information to undo the changes if needed. When this is the case, SQL Server performs *compensation operations*, executing actions opposite to those that made the original data modifications and generating compensation log records.

Figure 11-6 shows such an example, rolling transaction T1 back. Here, SQL Server has performed a compensation update, generating a compensation log record (LSN: 7320) to reverse the changes of the original update operation (LSN: 7315). It has also generated a compensation insert (LSN: 7321) to compensate for the delete operation (LSN: 7317).

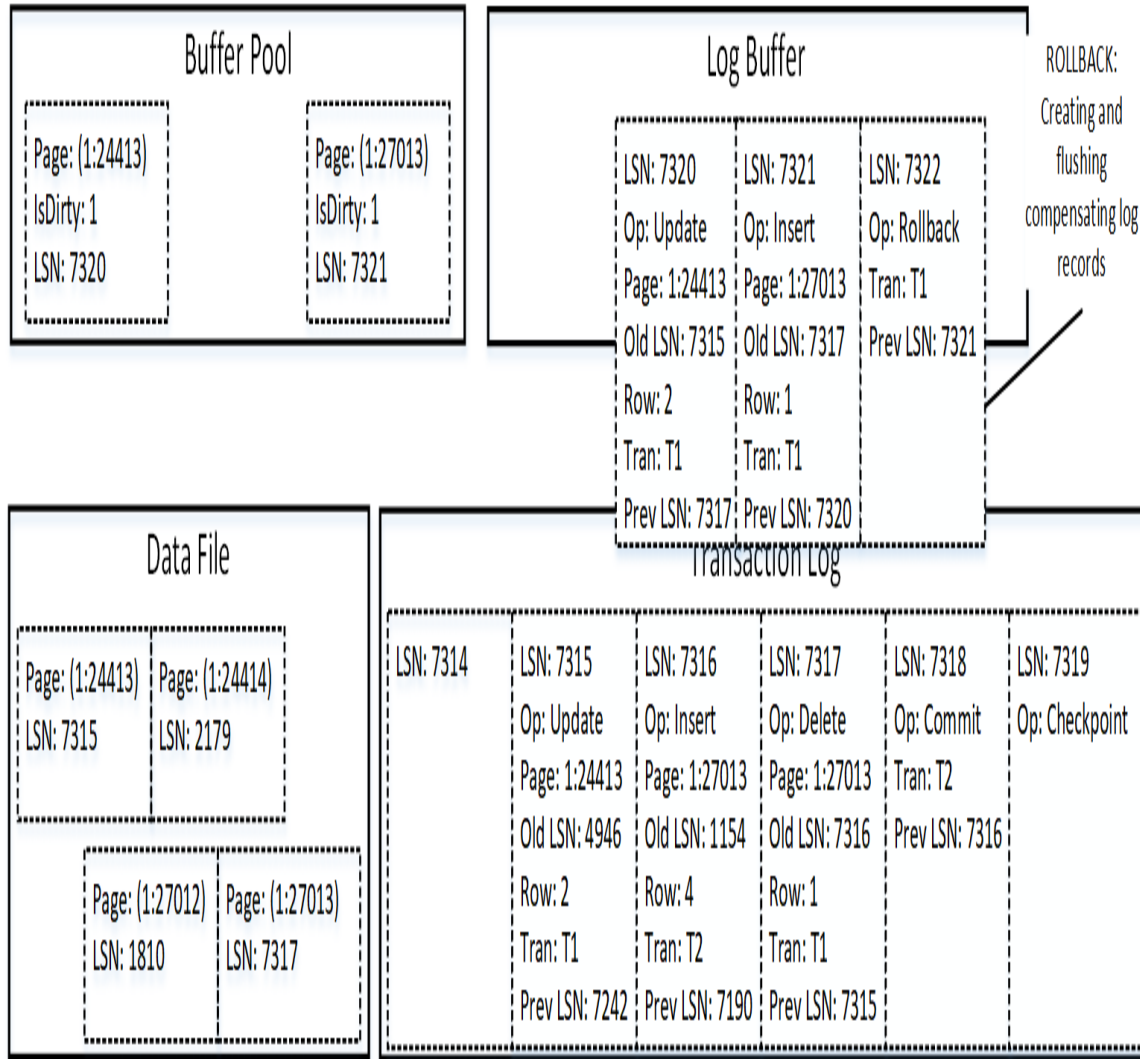


Figure 10-6. Data modifications and transaction logging: Rollback

There are two transaction logging-related waits to monitor.

WRITELOG

WRITELOG waits occur when SQL Server is waiting for the completion of an I/O operation that writes a log block to disk. With the exception of delayed durability (covered later in this chapter), this type of wait is synchronous, since it prevents transactions from committing while the write I/O is in progress. Your goal should be to minimize that wait and improve transaction-log throughput as much as possible.

LOGBUFFER

LOGBUFFER waits occur when SQL Server is waiting for an available log block to save the log records. In most cases, this happens due to insufficient I/O throughput, when SQL Server cannot write log blocks to disk fast enough. Usually, when LOGBUFFER waits are present, you'll also see WRITELOG waits. Improving transaction log throughput would help to address that.

I will talk how to troubleshoot and improve log-file throughput later in this chapter. However, you can also improve transaction-log performance by reducing the amount of logging. You can do this by removing unnecessary and unused indexes (more on those in Chapter 14), tuning your index-maintenance strategy to reduce page splits, and reducing the row size in frequently modified indexes.

You can also improve your transaction-management strategy by avoiding autocommitted transactions. This greatly reduces the amount of logging and the number of write log I/O requests in the system. Let's look at that in more detail.

Explicit and Autocommitted Transactions and Log Overhead

As you learned in Chapter 8, SQL Server always executes statements in the context of a transaction. If you don't have any explicit or implicit transactions started, SQL Server runs the statement in an autocommitted transaction, as if that statement was wrapped into a BEGIN TRAN ... COMMIT block.

Logging autocommitted transactions means including BEGIN XACT and COMMIT XACT transaction log records, which can significantly increase the amount of logging in the system. More importantly, it also decreases log performance, since SQL Server has to flush the log blocks after each statement on every COMMIT operation.

Figure 11-7 illustrates this. INSERT_1, UPDATE_1 and DELETE_1 operations run in autocommitted transactions, generating additional log

records and forcing the log buffer to flush on each COMMIT. Alternatively, running those operations in explicit transactions leads to more efficient logging.

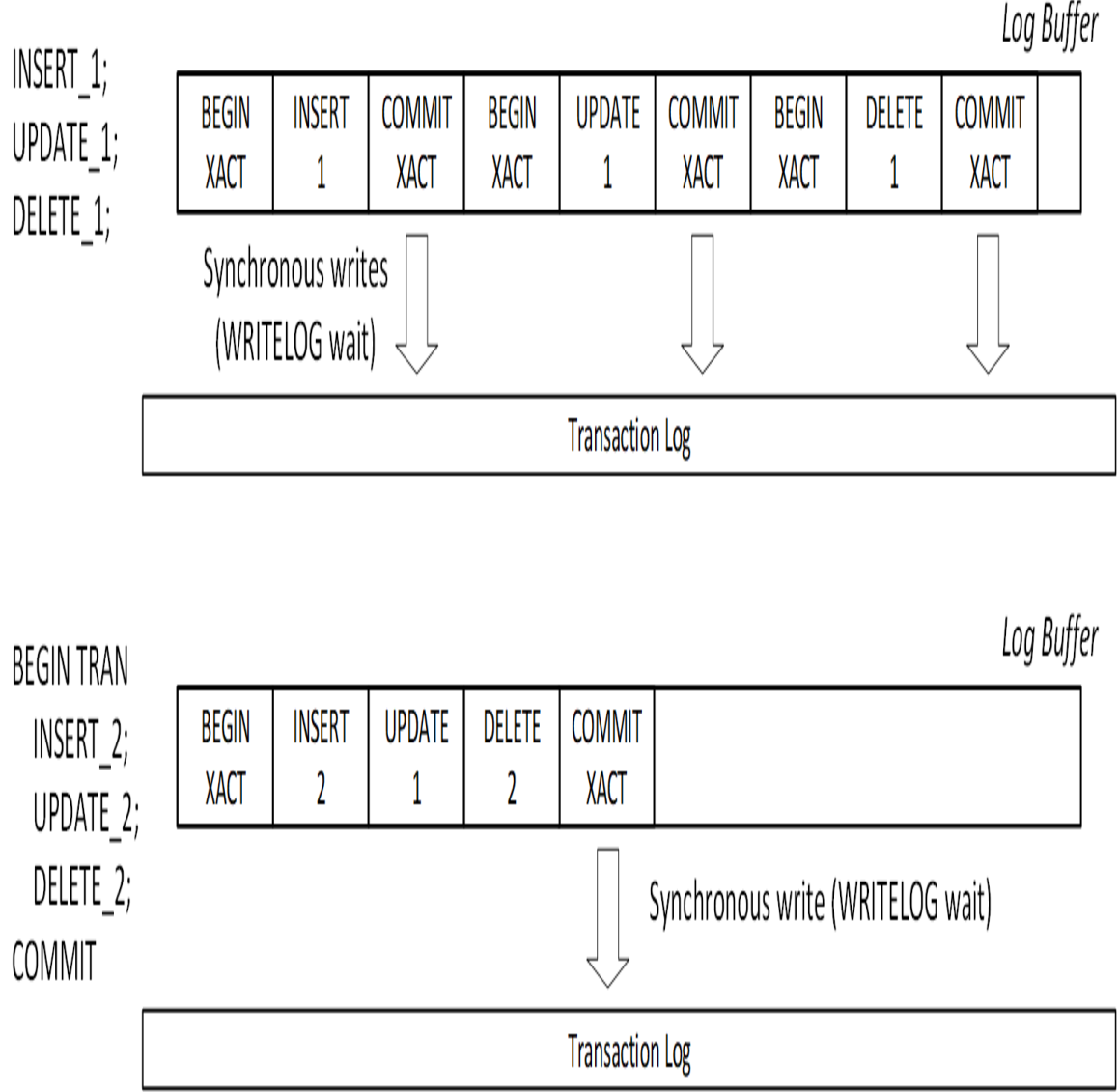


Figure 10-7. Explicit and autocommitted transactions

The code in Listing 11-1 shows the overhead involved in autocommitted transactions as compared to explicit transactions. It performs an INSERT/UPDATE/DELETE sequence 10,000 times in the loop, in autocommitted and explicit transactions, respectively. It then compares their

execution time and transaction log throughput using the sys.dm_io_virtual_file_stats view.

Example 10-1. Explicit and autocommitted transactions

```
CREATE TABLE dbo.TranOverhead
(
    Id INT NOT NULL,
    Col CHAR(50) NULL,
    CONSTRAINT PK_TrانOverhead
    PRIMARY KEY CLUSTERED(Id)
);

-- Autocommitted transactions
DECLARE
    @Id INT = 1
    ,@StartTime DATETIME = GETDATE()
    ,@num_of_writes BIGINT
    ,@num_of_bytes_written BIGINT

SELECT @num_of_writes = num_of_writes, @num_of_bytes_written =
num_of_bytes_written
FROM sys.dm_io_virtual_file_stats(db_id(),2);

WHILE @Id <= 10000
BEGIN
    INSERT INTO dbo.TrانOverhead(Id, Col) VALUES(@Id, 'A');
    UPDATE dbo.TrانOverhead SET Col = 'B' WHERE Id = @Id;
    DELETE FROM dbo.TrانOverhead WHERE Id = @Id;

    SET @Id += 1;
END;

SELECT
    DATEDIFF(MILLISECOND,@StartTime,GETDATE())
        AS [Time(ms): Autocommitted Tran]
    ,s.num_of_writes - @num_of_writes
        AS [Number of writes]
    ,(s.num_of_bytes_written - @num_of_bytes_written) / 1024
        AS [Bytes written (KB)]
FROM
    sys.dm_io_virtual_file_stats(db_id(),2) s;
GO

-- Explicit Tran
DECLARE
    @Id INT = 1
    ,@StartTime DATETIME = GETDATE()
```



```

        ,@num_of_writes BIGINT
        ,@num_of_bytes_written BIGINT

SELECT @num_of_writes = num_of_writes, @num_of_bytes_written =
num_of_bytes_written
FROM sys.dm_io_virtual_file_stats(db_id(),2);

WHILE @Id <= 10000
BEGIN
    BEGIN TRAN
        INSERT INTO dbo.TranOverhead(Id, Col) VALUES(@Id, 'A');
        UPDATE dbo.TranOverhead SET Col = 'B' WHERE Id = @Id;
        DELETE FROM dbo.TranOverhead WHERE Id = @Id;
    COMMIT
    SET @Id += 1;
END;

SELECT
    DATEDIFF(MILLISECOND,@StartTime,GETDATE())
        AS [Time(ms): Explicit Tran]
    ,s.num_of_writes - @num_of_writes
        AS [Number of writes]
    ,(s.num_of_bytes_written - @num_of_bytes_written) / 1024
        AS [Bytes written (KB)]
FROM
    sys.dm_io_virtual_file_stats(db_id(),2) s;

```

You can see the output from the code in my environment in Figure 11-8. Explicit transactions were about three times faster and generated three times less log activity than autocommitted ones.

	Time(ms): Autocommitted Tran	Number of writes	Bytes written (KB)
1	13117	30000	15670
	Time(ms): Explicit Tran	Number of writes	Bytes written (KB)
1	4500	10000	10658

Figure 10-8. EPerformance of explicit and autocommitted transactions

As I have stated, proper transaction management with explicit transactions can significantly improve your transaction-log throughput. Remember, however, the impact of long-running transactions on blocking. Exclusive (X) locks are held until the end of the transaction. Keep this locking behavior in mind as you design your transaction strategy and write your code, as well as considering the factors discussed in Chapter 8.

Unfortunately, changing transaction strategy in existing systems is not always possible. If your system suffers from a large number of autocommitted transactions and can tolerate a small amount of data loss, consider using another feature: *Delayed Durability*. This feature is available in SQL Server 2014 and above.

Delayed Durability

As you already know, SQL Server flushes the contents of the log block into a log file at the time of commit. It sends a confirmation to the client only after a commit record has been hardened to disk. This may lead to a large number of small log-write I/O requests with autocommitted transactions.

Delayed durability changes this behavior, making commit operations asynchronous. The client receives confirmation that the transaction has been committed immediately, without having to wait for the commit record to be hardened to disk. The commit record stays in the log buffer until one or more of the following conditions occur:

- The log block is full
- A fully durable transaction in the same database is committed, and its commit record flushes the contents of the log buffer to disk
- A CHECKPOINT operation occurs
- A `sp_flush_log` stored procedure is called
- A log buffer flush operation is triggered based on the log generation rate and/or timeout thresholds

There is obviously some risk here. If SQL Server crashes before the commit record is hardened, the data modifications from that transaction will be rolled back at recovery, as if the transaction had never been committed at all. However, other transactions would be able to see the data modifications made by the delayed durability transaction between the commit and the crash.

You can control delayed durability on both the database and transaction levels. The database option `DELAYED_DURABILITY` supports three different values:

DISABLED

This is the default option. It disables delayed durability in the database regardless of the transaction durability mode. All transactions in the database are always fully durable.

FORCED

This option forces delayed durability for all database transactions regardless of the transaction durability mode.

ALLOWED

With this option, delayed durability is controlled at the transaction level. Transactions are fully durable unless delayed durability is specified. Listing 11-2 shows how to specify it at the transaction level.

Example 10-2. Controlling delayed durability on transaction level

```
BEGIN TRAN
/* Do the work */
COMMIT WITH (DELAYED_DURABILITY=ON);
```

Delayed durability may be used in chatty systems with large numbers of autocommitted transactions and insufficient log throughput. In most cases, however, I prefer to avoid it. I use it as a last resort, only when all other log throughput improvement techniques have been unsuccessful and only when data loss is acceptable. Use with care!

In-Memory OLTP Transaction Logging

Although covering In-Memory OLTP in detail would be outside the scope of this book, I need to mention In-Memory OLTP transaction logging. As opposed to row-based and column-based technologies, In-Memory OLTP generates transaction log records at the time of the COMMIT operation and only when the transaction is successfully committed. Logging is also optimized. Transactions usually generate just one transaction log record, or a few large ones, even when they modify large amounts of data. Those records are stored in the regular log file and backed up with all the other log records.

This behavior may change I/O patterns for log operations. In-Memory OLTP log writes may lead to larger write requests, especially with large In-Memory OLTP transactions. Moreover, the log files are continuously read by In-Memory OLTP's continuous checkpoint process, which parses log records and updates In-Memory OLTP data persisted on disk.

You don't need to worry about those details in most cases; however, remember about I/O patterns when you design an I/O subsystem for databases that use In-Memory OLTP.

VIRTUAL LOG FILES

Internally, SQL Server divides physical log files into smaller parts called *Virtual Log Files (VLF)*. SQL Server uses them as a unit of management, and they can be active or inactive.

Active VLFs store the active portion of the transaction log, which contains log records required to keep the database transactionally consistent, provide point-in-time recovery, and support active SQL Server processes such as transactional replication and AlwaysOn Availability Groups. An inactive VLF contains the truncated (inactive) and unused parts of the transaction log.

Figure 11-9 shows an example transaction log file and VLFs. The active portion of the log starts with VLF3, the oldest active transaction in the

system. In case of a rollback, SQL Server would need to access log records generated by that transaction.

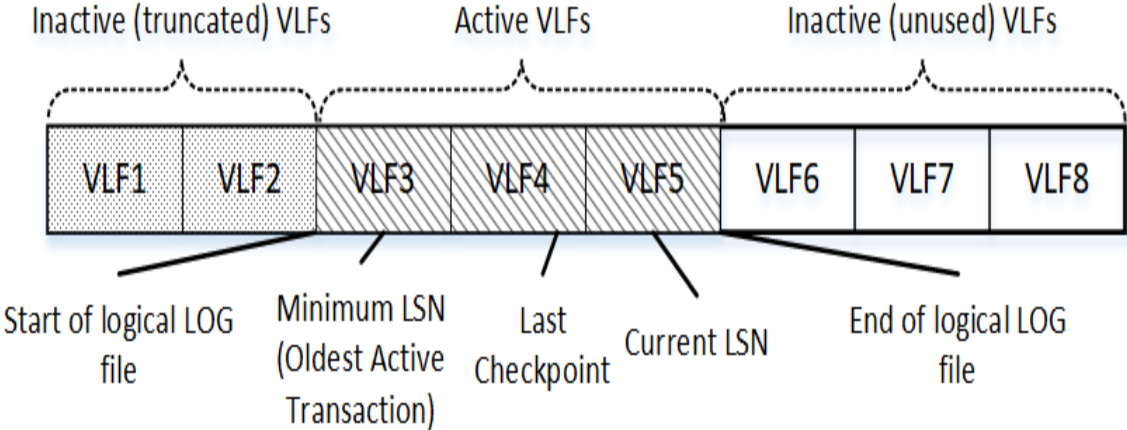


Figure 10-9. Transaction log and VLFs

In Figure 11-9, the only process that keeps VLF3 active is the active transaction. When this transaction commits, SQL Server truncates the log, marking VLF3 as inactive (Figure 11-10). Truncating the transaction log does not reduce the size of the log file on disk; it just means that parts of the transaction log (one or more VLFs) are marked as inactive and ready for reuse.

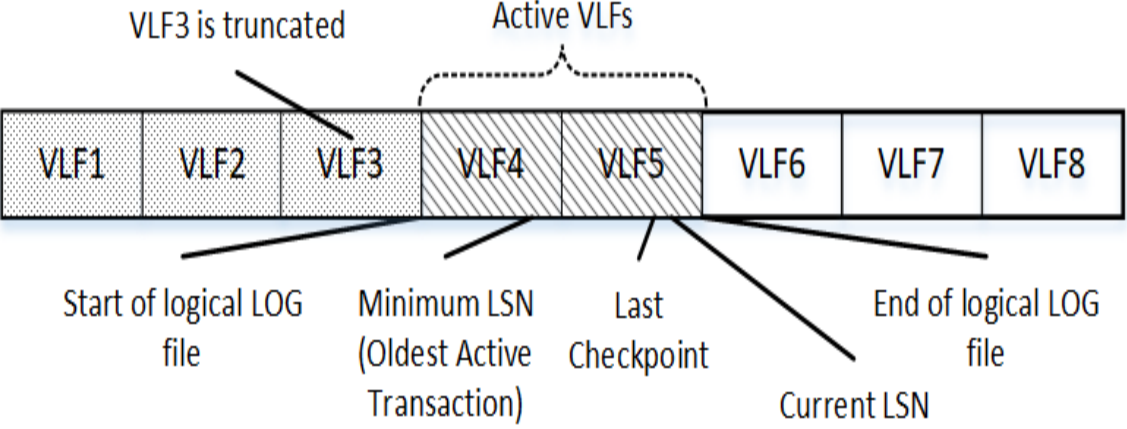


Figure 10-10. Transaction log and VLFs – After commit

SQL Server uses VLFs as the unit of truncation. A VLF cannot be marked as inactive if it contains the single log record from the active portion of the

log. This is one reason why having very large VLFs is not a good idea.

Transaction-log truncation behavior is controlled by the *Database Recovery Model* setting. There are three recovery models. Each guarantees that the active portion of the log has enough data to keep the database consistent; however, the models may provide different levels of recoverability in case of disaster, and SIMPLE and BULK-LOGGED models may prevent you from using some SQL Server technologies.

SIMPLE

In the SIMPLE recovery model, the log is truncated at checkpoint. All data pages with LSNs prior to checkpoint LSN are saved on-disk. SQL Server does not need to access log records prior to the checkpoint to re-apply them to data pages in the event of an unexpected shutdown or crash. Old active transactions and transaction replication may defer truncation, keeping VLFs active.

In this recovery model, SQL Server does not use transaction log backups. It prevents you from performing a point-in-time recovery and may lead to data loss if either of the database files (data or log) becomes corrupted. The *recovery point (RPO)* for the database, in this model, becomes the time of the last full backup.

In some cases, such as when data is static or can be recreated from other sources, the SIMPLE recovery model may be completely acceptable. However, when you encounter this during a system health check, you need to confirm this with the system stakeholders and discuss the possibility of data loss.

FULL

In the FULL recovery model, SQL Server fully logs all operations in the database and requires you to perform transaction log backups to truncate the transaction log. Because the transaction log backups store all log records in the database, this mode supports point-in-time recovery, as long as the sequence of backup files (*backup chain*) is

available. With a few exceptions, you will probably want to use the FULL recovery model in production databases.

To support various SQL Server features and technologies that rely on transaction log records, the FULL recovery model is required. Those technologies may also defer truncation of the log, even when log backups are taken. For example, if an AlwaysOn Availability Group node goes offline, SQL Server will be unable to truncate the log until the node is back and catches up with the replication.

BULK-LOGGED

The BULK-LOGGED recovery model works similarly to the FULL model, except that some operations are minimally logged: for example, index creation or BULK INSERT statements. With minimally logged operations, SQL Server logs only the page allocation information in the log file. While this reduces the log file usage, you cannot perform a point-in-time recovery when minimally logged operations are present. As with the SIMPLE recovery model, analyze the frequency of bulk-logged operations and the risks of potential data loss when you see databases with this recovery model in production.

You can analyze VLFs by using the `sys.dm_db_log_info` data management view in SQL Server 2016 and above or with the `DBCC LOGINFO` command in older versions of SQL Server. Listing 11-3 shows the code that uses this view against one of the databases.

Example 10-3. Analyzing VLFs in the database

```
SELECT *
FROM sys.dm_db_log_info(DB_ID());

SELECT
    COUNT(*) as [VLF Count]
    ,MIN(vlf_size_mb) as [Min VLF Size (MB)]
    ,MAX(vlf_size_mb) as [Max VLF Size (MB)]
    ,AVG(vlf_size_mb) as [Avg VLF Size (MB)]
FROM sys.dm_db_log_info(DB_ID());
```

Figure 11-11 shows the output from the view for a database with an incorrect log file configuration that uses a 10% auto-growth setting. As you can see, the database has a large number of unevenly sized VLFs.

	database_id	file_id	vlf_begin_offset	vlf_size_mb	vlf_sequence_number	vlf_active	vlf_st
35	16	2	27459584	2.62	281	1	2
36	16	2	30212096	2.87	282	1	2
37	16	2	33226752	3.18	283	1	2
38	16	2	36569088	3.5	284	1	2
39	16	2	40239104	3.87	285	1	2
40	16	2	44302336	4.25	286	1	2
41	16	2	48758784	4.68	287	1	2
42	16	2	53673984	5.12	288	1	2
43	16	2	59047936	5.62	289	1	2
	VLF Count	Min VLF Size (MB)	Max VLF Size (MB)	Avg VLF Size (MB)			
1	105	0.24	1889.87	197.98695238095235			

Figure 10-11. Inefficient VLF configuration

Transaction Log Configuration

SQL Server works with transaction logs sequentially while writing and reading a stream of log records. Even though the log may have multiple

physical files, SQL Server does not usually benefit from them; in most cases, a single log file is easier to maintain and manage.

There are a couple of edge cases when multiple log files can be beneficial:

SQL Server may zero-initialize log file in parallel. This could speed up database creation or restoration if it uses large (multi-terabyte) log files.

If you want to place the transaction log on a fast but small drive, you can create a file pre-allocating the size to fill the fast drive and add another small log file on a larger and slower drive. SQL Server will use the file on the fast drive most of the time; however, the small file will protect you if the transaction log is full and not truncating.

It is better to manage transaction log size manually, avoiding the overhead of zero-initializing at the time of auto-growth. You can analyze a transaction log's size and recreate the log, pre-allocating the size as needed. Be sure to take log-intensive operations such as index maintenance into account when you do the analysis.

As you learned stated in Chapter 1, you can rebuild the log by shrinking it to the minimal size and then pre-allocate the space using chunks of 1,024 to 4,096 MB. I usually use 1,024MB chunks, which will create 128MB VLFs. If I need *very* large log files – hundreds of gigabytes or even terabytes—I might use larger chunks.

Do not restrict the log's maximum size and auto-growth: you need to be able to grow the log in case of emergencies.

Log Truncation Issues

Excessive transaction log growth is a common problem that junior or accidental DBAs should handle. It happens when SQL can't truncate the transaction log and reuse the space in the log file. In such cases, the log file continues to grow until it fills the entire disk, switching the database to read-only mode and raising a "*Transaction log full*" error (error code 9002).

The best way to handle this condition is to avoid it in the first place. As I discussed in the *Monitoring Disk Space Usage* section in Chapter 9, it is essential to monitor for low disk space condition and set up alerts. If you pre-allocate the log file to fit the entire drive, monitor the amount of free space in the log file and have it send an alert when it is low.

If you end up in a *Transaction Log Full* situation, my first and most important advice is: Don't panic. First, you need to analyze the root cause of the issue and see if you can mitigate it quickly. You can do this by looking at the `log_reuse_wait_desc` column in the `sys.databases` view, either querying it directly or using the more sophisticated version shown in Listing 11-4. This column shows you why the log is not truncated.

Example 10-4. Analyzing the log_reuse_wait_desc column in the sys.databases view

```
CREATE TABLE #SpaceUsed
(
    database_id SMALLINT NOT NULL,
    file_id SMALLINT NOT NULL,
    space_used DECIMAL(15,3) NOT NULL,
    PRIMARY KEY(database_id, file_id)
);

EXEC master..sp_MSforeachdb
N'USE[?];
INSERT INTO #SpaceUsed(database_id, file_id, space_used)
    SELECT DB_ID('?'), file_id,
        (size - CONVERT(INT,FILEPROPERTY(name, 'SpaceUsed')) /
128.
FROM sys.database_files
WHERE type = 1;';

SELECT
    d.database_id, d.name, d.recovery_model_desc
    ,d.state_desc, d.log_reuse_wait_desc, m.physical_name
    ,m.is_percent_growth
    ,IIF(m.is_percent_growth = 1
        ,m.growth
        ,CONVERT(DECIMAL(15,3),m.growth / 128.0)
    ) AS [Growth (MB or %)]
    ,CONVERT(DECIMAL(15,3),m.size / 128.0) AS [Size (MB)]
    ,IIF(m.max_size = -1
        ,-1
```

```

        ,CONVERT(DECIMAL(15,3),m.max_size / 128.0)
    ) AS [Max Size(MB)]
    ,s.space_used as [Space Used(MB)]
FROM
    sys.databases d WITH (NOLOCK)
        JOIN sys.master_files m WITH (NOLOCK) ON
            d.database_id = m.database_id
    LEFT OUTER JOIN #SpaceUsed s ON
        s.database_id = m.database_id AND
        s.file_id = m.file_id
ORDER BY
    d.database_id;

```

Figure 11-12 shows example output from Listing 11-4.

database_	name	recovery_	state_desc	log_reuse_wait_desc	physical_	is_percent_growth	Growth (MB or %)	Size (MB)	Max Size(MB)	Space Used(MB)
1	master	SIMPLE	ONLINE	NOTHING	M:\SQLDa...	1	10.000	2.250	-1.000	1.328
2	tempdb	SIMPLE	ONLINE	NOTHING	T:\SQLDa...	0	1024.000	1024.000	-1.000	966.094
2	tempdb	SIMPLE	ONLINE	NOTHING	L:\SQLDa...	0	1024.000	1024.000	2097152.000	966.094
3	model	FULL	ONLINE	LOG_BACKUP	M:\SQLDa...	0	1024.000	1024.000	-1.000	1015.336
4	msdb	SIMPLE	ONLINE	NOTHING	M:\SQLDa...	1	10.000	19.625	2097152.000	18.133
5	DBA	SIMPLE	ONLINE	NOTHING	L:\SQLDa...	0	1024.000	1024.000	2097152.000	937.281
16	LogDemo2	FULL	ONLINE	LOG_BACKUP	L:\SQLDa...	0	0.000	20789.0...	21000.000	-16.281

Figure 10-12. Analysing log_reuse_wait_desc data

Let's look at the most common reasons for deferred log truncation and log_reuse_wait_desc values.

LOG_BACKUP Log Reuse Wait

The LOG_BACKUP log reuse wait is one of the most common waits for databases in the FULL and BULK-LOGGED recovery models. It indicates that the log cannot be truncated due to a lack of recent transaction log backup.

When you see this log reuse wait, check the status of the transaction log backup job. Make sure it is not failing due to a lack of space in the backup destination, or for any other reasons. It is also possible that log backup frequency is not fast enough during log-intensive operations. For example, index maintenance can generate enormous amounts of transaction log records in a very short time.

You can mitigate the issue by performing a transaction log backup. Remember to keep the backup file if you run this operation manually using a non-standard backup destination. The file becomes part of your backup chain and will be required for database recovery.

If you don't use any technologies that rely on the FULL recovery model, you can temporarily switch the database to SIMPLE mode, which will truncate the transaction log. Remember that this will leave you exposed to data loss. Switch back to the FULL recovery model and reinitialize the backup chain by performing FULL and LOG backups as quickly as possible.

Finally, in many cases, this situation may be avoided by properly monitoring the health of the backup jobs. Set up alerts for continuous log backup failures.

ACTIVE_TRANSACTION Log Reuse Wait

The ACTIVE_TRANSACTION log reuse wait indicates that the log cannot be truncated due to the presence of the old active transaction. The most common case for that is incorrect transaction management in the application, which leads to runaway uncommitted transactions. For example, the application may issue multiple BEGIN TRAN statements without corresponding COMMIT for each of them.

You can see the list of active transactions using the code from Listing 11-5. The code may provide you multiple rows for each transaction, because it gets log usage information on a per-database basis.

Example 10-5. Getting active transactions

```
SELECT
    dt.database_id
    ,DB_NAME(dt.database_id) as [DB]
    ,st.session_id
    ,CASE at.transaction_state
        WHEN 0 THEN 'Not Initialized'
        WHEN 1 THEN 'Not Started'
        WHEN 2 THEN 'Active'
        WHEN 3 THEN 'Ended (R/O)'
        WHEN 4 THEN 'Commit Initialize'
        WHEN 5 THEN 'Prepared'
        WHEN 6 THEN 'Committed'
        WHEN 7 THEN 'Rolling Back'
        WHEN 8 THEN 'Rolled Back'
    END AS [State]
    ,at.transaction_begin_time
    ,es.login_name
    ,ec.client_net_address
    ,ec.connect_time
    ,dt.database_transaction_log_bytes_used
    ,dt.database_transaction_log_bytes_reserved
    ,er.status
    ,er.wait_type
    ,er.last_wait_type
    ,sql.text AS [SQL]
FROM
    sys.dm_tran_database_transactions dt WITH (NOLOCK)
        JOIN sys.dm_tran_session_transactions st WITH (NOLOCK) ON
            dt.transaction_id = st.transaction_id
        JOIN sys.dm_tran_active_transactions at WITH (NOLOCK) ON
            dt.transaction_id = at.transaction_id
        JOIN sys.dm_exec_sessions es WITH (NOLOCK) ON
            st.session_id = es.session_id
        JOIN sys.dm_exec_connections ec WITH (NOLOCK) ON
            st.session_id = ec.session_id
        LEFT OUTER JOIN sys.dm_exec_requests er WITH (NOLOCK) ON
            st.session_id = er.session_id
        CROSS APPLY
            sys.dm_exec_sql_text(ec.most_recent_sql_handle) sql
ORDER BY
    dt.database_transaction_begin_time
```

You can kill the session that holds active transaction using the KILL command. Later, you can analyze why the transaction was not properly managed.

AVAILABILITY_REPLICA Log Reuse Wait

As you can guess by the name, the AVAILABILITY_REPLICA log reuse wait may occur in systems that use AlwaysOn Availability Groups. In that technology, the primary node communicates with the secondary nodes by sending them the stream of transaction log records. The log cannot be truncated until those records have been sent and replayed on the secondaries.

The AVAILABILITY_REPLICA log reuse waits usually occur during certain issues in Availability Groups: most commonly the secondary node being unavailable, replication between the nodes falling behind, or secondaries being unable to replay changes fast enough to keep up with the load.

When you see this wait, check the health of the Availability Group. In most cases, the only quick option to address the issue, besides adding more space to the log, is removing a problematic secondary node from the Availability Group. I'll discuss Availability Groups in more detail in the next chapter.

DATABASE_MIRRORING Log Reuse Wait

The DATABASE_MIRRORING log reuse wait occurs in systems that use database mirroring technology. This technology was a predecessor of AlwaysOn Availability Groups and behaves similarly to it, communicating through the stream of log records.

I'm not going to discuss how to troubleshoot database mirroring in this book, since that technology has long been obsolete. Conceptually, it is similar to troubleshooting Availability Groups issues.

As with AVAILABILITY_REPLICA reuse waits, analyze the health of database mirroring when you see a DATABASE_MIRRORING reuse wait.

REPLICATION Log Reuse Wait

The REPLICATION log reuse wait occurs when the Log Reader agent falls behind while harvesting log records for transactional replication or Change Data Capture (CDC) processes. When you see this log reuse wait, check the status of the Log Reader agent and address the issues you discover.

One particular issue you may experience is the Log Agent query timing out. By default, it is 30 minutes, which is sufficient in most cases. However, it may not be enough when system processes very large data modifications (millions of rows) in the replicated tables. You can **increase the QueryTimeout in the Log Agent profile** if this is the case.

There are two “nuclear” options: you can remove the replication or mark log records as harvested by using the `sp_repldone` command. Both approaches may require you to reinitialize the replication later.

ACTIVE_BACKUP_OR_RESTORE Log Reuse Wait

The ACTIVE_BACKUP_OR_RESTORE reuse wait indicates that the log cannot be truncated due to active database backup or restore processes, regardless of what type of backup or restore is running.

One of the common cases for this wait is degraded performance of network or backup storage during the large FULL backup. Check the status of the active backup and restore jobs when you see this reuse wait.

Other Mitigation Strategies

It isn't always possible to mitigate the root cause of an issue quickly. Sometimes it isn't even feasible. For example, removing unavailable replica from Availability Group or disabling replication may lead to significant work rebuilding them later.

You can add another log file or expand the size of the log drive as a temporary solution. This allows the database to operate and gives you some time to mitigate the root cause of the issue.

Remember that you could have multiple issues preventing transaction log truncation simultaneously. For example, a network outage could prevent the server from communicating with Availability Group replicas and from accessing the backup destination. Check the `log_reuse_wait_desc` value and the amount of free space in the log file after you address each issue.

Finally, learn from the experience. *Transaction Log Full* issues are serious, and you should avoid it at all costs. Do the root-cause analysis, evaluate your monitoring strategy, and perform capacity planning to reduce the possibility of this happening again.

Transaction Log Throughout

The impact of bad transaction log throughput is not always visible. While any operation that changes something in the database writes to transaction log, those writes are considered part of the operation. Engineers tend to look at the “big picture” and performance of an entire operation as a whole, overlooking the impact of individual components.

Think about index maintenance, for example. This is an extremely log-intensive operation and bad log throughput greatly affects its performance. Nevertheless, database administrators usually try to reduce index maintenance impact and duration by adjusting its schedule or excluding indexes from the maintenance, overlooking slow log writes. (There are some exceptions; however, they are few and far between.)

It is also easy to overlook the impact of bad transaction log throughput on the regular workload in OLTP systems. High log write latency increases queries execution time, though people rarely look at it during query tuning. In either case, improving transaction log performance always improves system performance.

A word of caution, though: While improving transaction log performance is always beneficial, it is not a magic solution to every problem. Nor will impact always be visible to users. You may get better ROI from addressing other bottlenecks first.

From a wait statistics standpoint, you can look at WRITELOG and LOGBUFFER waits: as I've mentioned, WRITELOG occurs when SQL Server waits for the log write to complete. LOGBUFFER happens when SQL Server does not have the available log buffer to cache the log records or can't flush log blocks to disk fast enough.

I wish I could tell you an exact percentage threshold when those waits start to represent a problem, but that would be impossible. You'll always see them in the system, but you need to look at general I/O health and throughput to estimate their impact. You'll also often see WRITELOG wait together with other I/O waits (PAGEIOLATCH, etc), especially when log and data files are sharing the same *physical* storage and network I/O path. In most cases, this is the sign that the I/O subsystem is either overloaded or configured incorrectly.

As Chapter 3 discussed, the sys.dm_io_virtual_file_stats view provides information about database files latency and throughput. I generally like to see an average write latency in the log files of within 1 to 3 milliseconds when network-based storage is used.

In high-end OLTP systems, you can place log files to DAS NVMe drives, which should bring the latency into the sub-millisecond range. In some edge cases, you can also utilize persistent memory technologies to reduce latency even further.

Pay attention to the average size of log writes. Small writes are less efficient and may impact log throughput. In most cases, those writes are a sign of autocommitted transactions. You can reduce them with proper transaction management or by enabling delayed durability (in cases when you can tolerate a small amount of data loss).

There are several performance counters in the Databases object that you can use to monitor transaction log activity in real time. These include:

Log Bytes Flushed/sec

The Log Bytes Flushed/sec counter shows how much data has been written to the log file.

Log Flushes/sec

The Log Flushes/sec counter indicates how many log write operations have been performed every second. You can use it with the Log Bytes Flushed/sec counter to estimate the average log write size.

Log Flush Write Time(ms)

The Log Flush Write Time(ms) counter shows the average time of a log write operation. You can use it to see log write I/O latency in real time.

Log Flush Waits/sec

The Log Flush Waits/sec counter provides the number of commit operations per second while waiting for the log records to be flushed. Ideally, this number should be very low. High values indicate a log-throughput bottleneck.

Keep in mind that log-intensive operations may dramatically change the numbers. For example, an unthrottled index rebuild can generate an enormous number of log records very quickly, saturating your I/O subsystem and log throughput.

I like to run a SQL Agent job and collect data from the `sys.dm_io_virtual_file_stats` view every 5 to 15 minutes, then persist it in the DBA utility database. This provides detailed information about I/O workload overtime. This information is extremely useful when you're analyzing transaction-log throughput.

In either case, you may need to look at a few different areas to improve transaction log performance and throughput. First, analyze the hardware – it is essential to use a fast disk subsystem with the log files. Next, look at the log configuration, and try to reduce the overhead of large number of VLFs and growth management. Finally, look for opportunities to reduce log generation through proper transaction management and efficient database maintenance jobs.

Summary

SQL Server uses Write-Ahead Logging to support database consistency requirements. Log records are hardened to the log file before transactions are committed. Insufficient transaction log throughput impacts system performance. Put transaction logs to the fast low-latency storage when possible.

Transaction log throughput issues present themselves with WRITELOG and LOGBUFFER waits. Troubleshoot log performance when these waits become noticeable. You can use `sys.dm_io_virtual_file_stats` view and performance counters for troubleshooting.

You can improve transaction log performance by reducing the databases' log generation rate. It can be done by tuning index maintenance strategy, removing unnecessary indexes, refactoring database schema, and improving transaction management. You can enable delayed durability for databases that handle large numbers of autocommitted transactions and can sustain small data losses.

Make sure that the log file is properly configured and the number of VLFs in the file is manageable. Consider rebuilding the log if you detect suboptimal configuration.

In the next chapter, we'll talk about AlwaysOn Availability Groups and the issues you may encounter when using them.

Troubleshooting Checklist

- Review and adjust transaction log configuration for the databases
- Analyze number of VLFs and rebuild logs if needed
- Check databases' recovery model and discuss disaster recovery strategy with stakeholders
- Reduce transaction log generation rate when possible
- Analyze and improve transaction log throughput

Chapter 11. AlwaysOn Availability Groups

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 12th chapter of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at dmitri@aboutsqlserver.com.

AlwaysOn Availability Groups is the most common High Availability technology used in SQL Server. It persists multiple copies of databases, eliminating storage from being the single point of failure. It also allows you to scale read workload through multiple readable secondary nodes.

In this chapter, I’ll provide an overview of how Availability Groups work internally and explain how to troubleshoot common issues. You will learn about the overhead introduced by synchronous replicas and readable secondaries. Finally, I will discuss Availability Groups monitoring and cover how to troubleshoot failover events.

AlwaysOn Availability Groups Overview

Perhaps the easiest way to explain how AlwaysOn Availability Groups work is to look at the history of this technology. It was introduced in SQL Server

2012 as the replacement and successor to Database Mirroring, which used to have the internal Microsoft code name *Real Time Log Shipping*. That name is a good description: both Database Mirroring and Availability Groups rely on the stream of transaction log records to communicate. The primary node sends log records to secondaries, which harden them in their transaction logs and replay (redo) the changes in the databases.

Both technologies support the single primary node that handles write workload. The Availability Group in Enterprise edition could handle read-only workload, which you can scale across multiple secondary replicas. In Standard Edition, the technology is limited to Basic Availability Groups, which support just a single secondary replica. This replica is used strictly for High Availability (HA) and/or Disaster Recovery (DR) purposes. Clients cannot connect and read data from there.

Availability Groups support automatic failover for HA. They rely on Windows Server Failover Cluster (WSFC) in Windows and on Pacemaker in Linux. Prior to SQL Server 2017, you had to have WSFC to set up Availability Groups. Starting with SQL Server 2017, you can use Availability Groups without WSFC. Automatic failover is not supported in that configuration.

Figure 12-1 shows an example of a three-node Availability Group setup in a Windows environment. Availability Group Listener is the name of the virtual network (similar to WSFC Cluster Endpoint) that provides a level of abstraction for clients' connectivity. In this configuration, clients can connect to Availability Group without knowing which node is working as the primary at that moment.

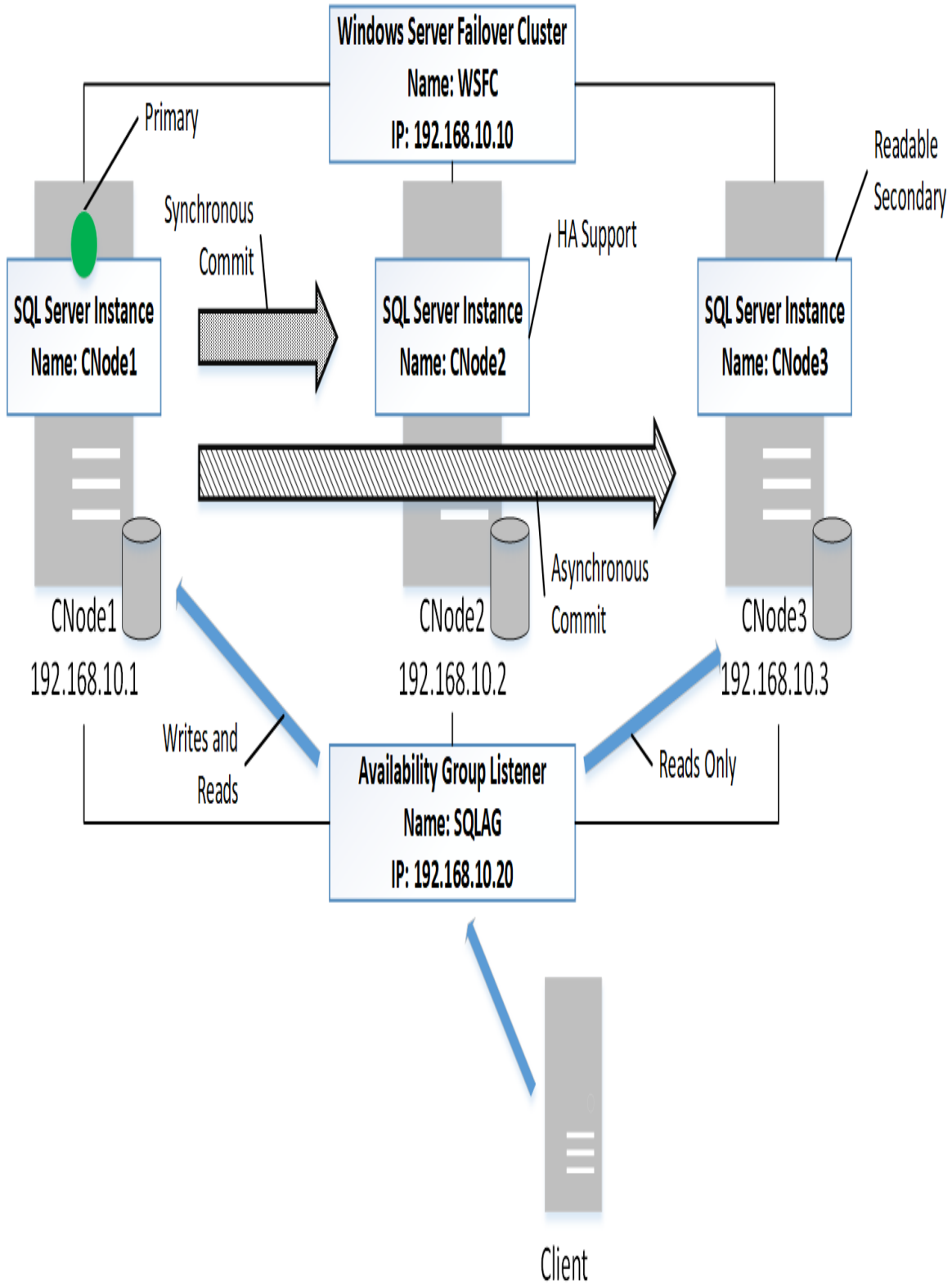


Figure 11-1. Example of Availability Group setup

Availability Groups work on the database-group level. They replicate the group of databases, failing them over together to another node in the event of failover. Each node in the topology stores its own copies of the data. This avoids a situation where storage can be the single point of failure.

Unfortunately, Availability Groups do not replicate instance-level objects. You need matching sets of logins, jobs, and all other instance objects on all nodes to support HA properly and ensure that it can operate after failover. Make sure you validate the configuration as part of any system health check, and regularly test HA implementation in production systems.

NOTE

The dbatools [open source library](#) provides many PowerShell cmdlets to replicate instance-level objects and validate the setup.

Now let's look at the key components of Availability Group technology – send and redo queues.

Availability Group Queues

Availability Group nodes communicate through the stream of transaction log records. Secondary nodes write those records to transaction logs and asynchronously apply the changes into the databases they host.

Figure 12-2 shows a high-level view of this process. The key components here are *send queues and redo queues*.

Send Queue

Send queues exist on the primary node. They store log records that need to be sent to the secondary nodes. There are separate queues per database for each secondary node in the Availability Group.

Redo Queue

Redo queues exist on secondary nodes. They store the log records for changes that need to be applied to the databases by an asynchronous *redo process*. Each database on each secondary node has its own redo queue.

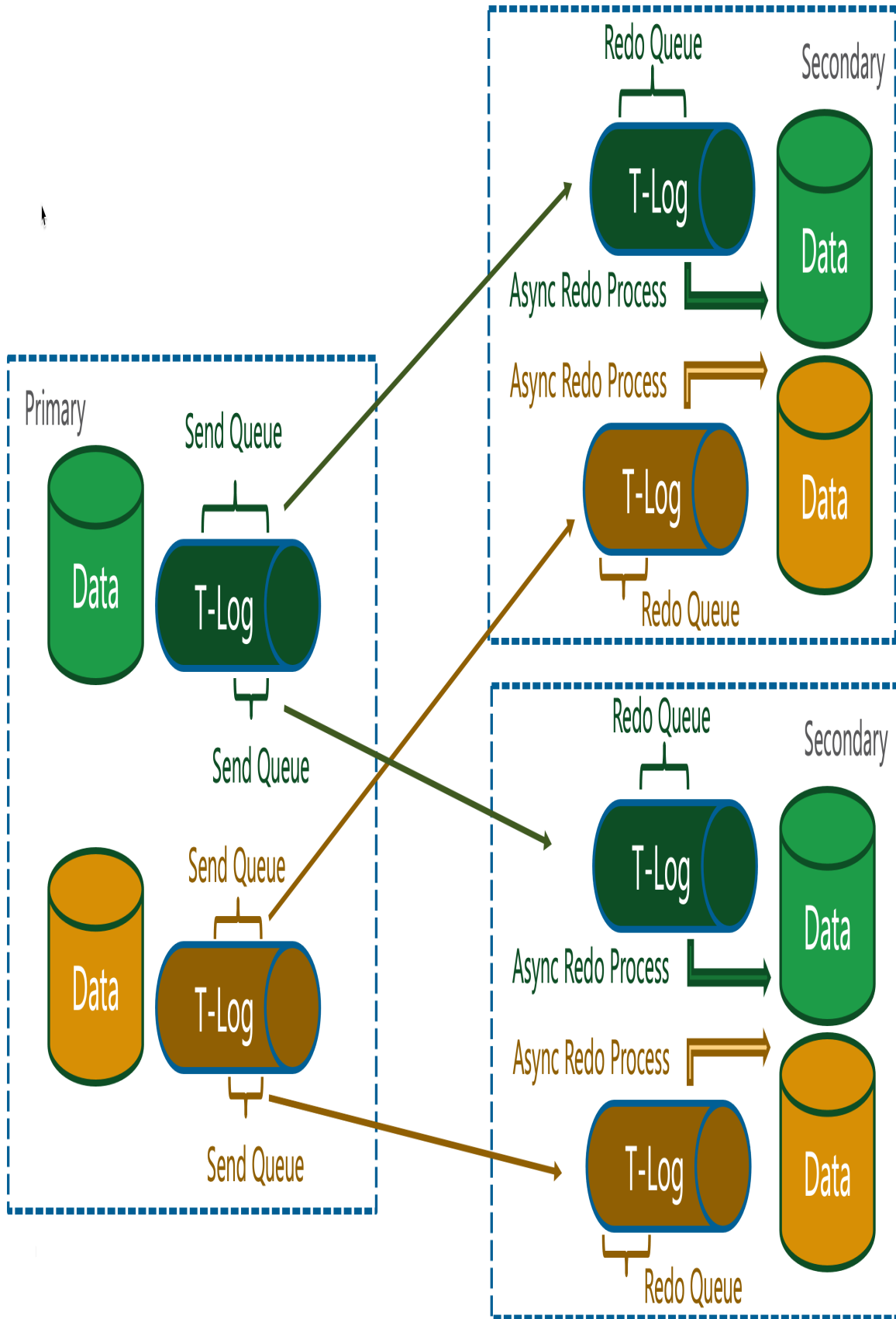


Figure 11-2. Availability Group queues

In a healthy Availability Group, send and redo queues stay as small as possible. A large send queue increases the amount of possible data loss in asynchronous replicas and may significantly increase blocking with synchronous ones (more on that later in this chapter). A large redo queue will impact database recovery time and increase replication latency, which in turn can affect queries on readable secondary nodes.

Finally, large send and redo queues also impact transaction-log truncation. SQL Server will not truncate a log beyond the largest send-queue starting point. Nor (though this is not documented) will it truncate the log beyond the oldest redo starting point across all secondary nodes.

Listing 12-1 shows you the code you can use to monitor the health of Availability Groups. You need to run it on the primary node to get the right results.

Listing 12-1. Availability Group monitoring code

```
SELECT
    ar.replica_server_name as [Replica]
    ,DB_NAME(drs.database_id) AS DB
    ,drs.synchronization_state_desc as [Sync State]
    ,ars.synchronization_health_desc as [Health]
    ,ar.availability_mode as [Synchronous]
    ,drs.log_send_queue_size
    ,drs.redo_queue_size
    ,ISNULL(
        GhostReplicaState.max_low_water_mark_for_ghosts -
        drs.low_water_mark_for_ghosts,0
    ) AS [water_mark_diff]
    ,drs.log_send_rate
    ,drs.redo_rate
    ,pri.last_commit_time AS primary_last_commit_time
    ,IIF(drs.is_primary_replica = 1
        ,pri.last_commit_time
        ,drs.last_commit_time
    ) AS node_last_commit_time
    ,IIF(drs.is_primary_replica = 1
        ,0
        ,DATEDIFF(ms,drs.last_commit_time,pri.last_commit_time)
    ) AS commit_latency
```

```

FROM
    sys.availability_groups ag WITH (NOLOCK)
    JOIN sys.availability_replicas ar WITH (NOLOCK) ON
        ag.group_id = ar.group_id
    JOIN sys.dm_hadr_availability_replica_states ars WITH
(NOLOCK) ON
        ar.replica_id = ars.replica_id
    JOIN sys.dm_hadr_database_replica_states drs WITH (NOLOCK)
ON
        ag.group_id = drs.group_id AND
        drs.replica_id = ars.replica_id
    LEFT JOIN sys.dm_hadr_database_replica_states pri WITH
(NOLOCK) ON
        pri.is_primary_replica = 1 AND
        drs.database_id = pri.database_id
    OUTER APPLY
        (
            SELECT MAX(drs2.low_water_mark_for_ghosts) AS
                max_low_water_mark_for_ghosts
            FROM sys.dm_hadr_database_replica_states drs2 WITH
(NOLOCK)
                WHERE drs.database_id = drs2.database_id
        ) GhostReplicaState
WHERE
    ars.is_local = 0
ORDER BY
    replica_server_name, DB;

```

Figure 12-3 shows the output from Listing 12-1 on one of the production servers.

	Replica	DB	Sync State	Health	Synchronous	log_send_queue_size	redo_queue_size
1			SYNCHRONIZED	HEALTHY	1	0	3605
2			SYNCHRONIZED	HEALTHY	1	0	81
3			SYNCHRONIZED	HEALTHY	1	0	0
4			SYNCHRONIZED	HEALTHY	1	0	0
5			SYNCHRONIZED	HEALTHY	1	0	0
6			SYNCHRONIZED	HEALTHY	1	0	0
7			SYNCHRONIZED	HEALTHY	1	0	0
8			SYNCHRONIZING	HEALTHY	0	60	2898
9			SYNCHRONIZING	HEALTHY	0	0	128
	water_mark_diff	log_send_rate	redo_rate	primary_last_commit_time	node_last_commit_time	commit_latency	
	985	236939	34739	2021-05-13 07:15:08.397	2021-05-13 07:15:05.960	2436	
	53	0	85281	2021-05-13 07:15:08.340	2021-05-13 07:15:07.680	660	
	1	0	22799	2021-05-13 07:00:33.760	2021-05-13 07:00:33.760	0	
	1	0	0	2021-05-13 07:00:34.220	2021-05-13 07:00:34.220	0	
	1	0	0	2021-05-13 07:00:34.493	2021-05-13 07:00:34.493	0	
	1	0	0	2021-05-13 07:00:34.777	2021-05-13 07:00:34.777	0	
	1	0	1	2021-05-13 07:00:35.050	2021-05-13 07:00:35.050	0	
	38251	0	10584	2021-05-13 07:15:08.397	2021-05-13 07:15:06.483	1913	
	117	0	50990	2021-05-13 07:15:08.340	2021-05-13 07:15:07.510	830	

Figure 11-3. Output from Availability Group monitoring script

You need to monitor several things:

Synchronization Health and State

Synchronization health and state, provided by the `synchronization_health_desc` and `synchronization_state_desc` columns, indicate if the Availability Group is healthy and if the data is synchronized. Those are the key metrics to monitor.

Send and Redo Queue Sizes

You can see send and redo queue sizes in the `log_send_queue_size` and `redo_queue_size` columns output. Both queues should be as small as possible.

Replication Lag

You can monitor replication lag by comparing the last commit times on the primary and secondary nodes. The data is available in the `last_commit_time` columns (these appear as `primary_last_commit_time` and `node_last_commit_time` columns in the script output). Both the send and redo queues affect the lag: the more data you have in the queues, the higher the lag will be. Obviously, the lag should be as small as possible, especially if you are using readable secondaries.

The `sys.dm_hadr_database_replica_states` view has a `secondary_lag_seconds` column; however, I have found it less accurate than calculating lag based on `last_commit_time` data.

Ghost Cleanup Lag

Large send and redo queues and long active transactions on the readable secondaries will defer the ghost cleanup and version store cleanup processes on the primary node, impacting system performance. From a monitoring standpoint, you can detect that condition by analyzing the difference in `low_water_mark_for_ghosts` values between the primary

and secondary nodes. That data is exposed by the `water_mark_diff` column in the script.

The impact of readable secondaries on system performance is an important topic which I'll cover later in this chapter.

I'd like to repeat: it is *extremely important* to monitor the health and performance of the Availability Groups in your setup. There are many things that can go wrong. I'll talk about a few of them in this chapter.

You can build the monitoring code using the script from Listing 12-1. You can compare the metrics returned by the script against pre-defined thresholds, triggering alerts as needed. Tune the thresholds for your specific workload and infrastructure. For example, asynchronous off-site replicas may need higher send queue alert thresholds than synchronous on-site HA replicas.

You can also obtain queue sizes and several other Availability Group performance metrics through performance counters in the *Database Replica performance object*. You can build alerting around them; however, this method is more susceptible to load spikes. For example, large batch operations may trigger short spikes in log generation, leading to unnecessary alerts.

Let's look at a few common issues you may encounter with Availability Groups.

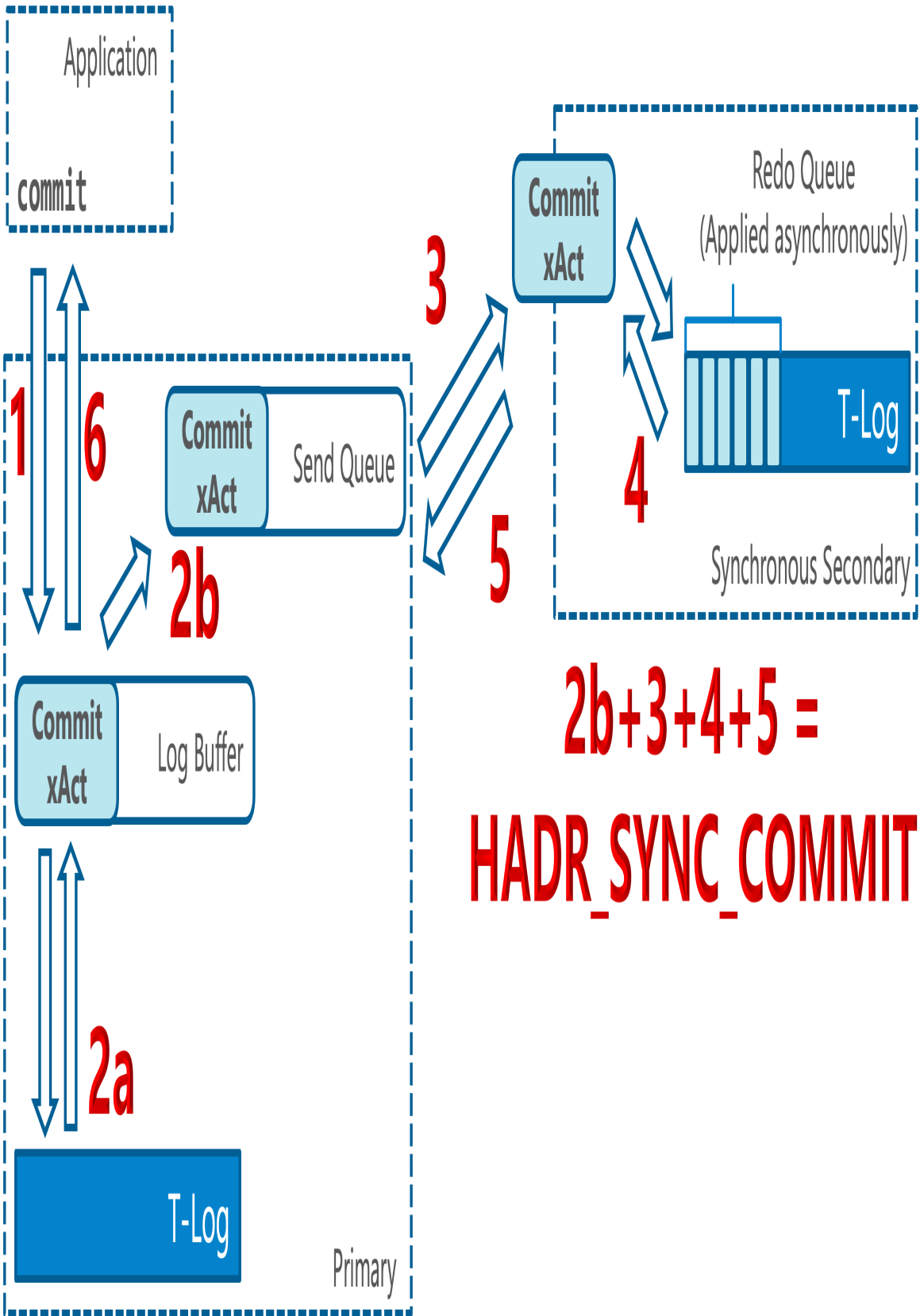
Synchronous Replication and the Danger of the `HADR_SYNC_COMMIT` Wait

Availability Groups allow you to configure replication using either synchronous or asynchronous commit modes. The synchronous mode allows you to avoid data loss, but at the cost of additional commit latency.

There is a common misconception that, in synchronous mode, data on the secondaries is updated synchronously with the primary node. This is not the

case. Only the log records are hardened synchronously. The redo process will still be asynchronous and can fall behind.

Figure 12-4 shows the replication data flow in synchronous commit mode. As you can see, the client does not receive confirmation that the transaction is committed until the primary node gets acknowledgement that commit log records have been hardened on the secondaries. The primary node waits for the confirmation, generating an HADR_SYNC_COMMIT wait.



$$2b+3+4+5 = \text{HADR_SYNC_COMMIT}$$

Figure 11-4. Synchronous commit data flow

This behavior introduces subtle, hard-to-understand concurrency issues. SQL Server keeps transactions active and does not release locks until it receives commit acknowledgements. This increases the chances of competing lock requests and blocking.

That increase in transaction duration may not always be visible in normal circumstances. However, it may become an issue during log-intensive operations, when the primary node does not have enough throughput to send log records to secondaries and send queues start to grow.

You can duplicate this problem by running an unthrottled clustered index rebuild on a large table with LOB columns (tables with LOB columns will produce more log records). If your server is powerful enough, you'll notice that the send queues for synchronous replicas start growing. This will lead to extensive blocking in busy OLTP systems, consume available workers and likely bring the system down in short amount of time. To make matters worse, cancelling the index rebuild will not solve the issue immediately, because the primary will still need to transmit all log records from the send queue.

NOTE

You can see the short YouTube [video](#) I published to demonstrate the problem

It is common to see HADR_SYNC_COMMIT become one of the top waits in busy OLTP systems. This situation may be legitimate and does not always represent the problem. Nevertheless, you need to estimate the impact on your system of this wait and the overhead of a synchronous commit.

Look at the average resource wait time of an HADR_SYNC_COMMIT wait in the output from sys.dm_os_wait_stats view (Listing 2-1). That value represents the average time SQL Server waits for acknowledgement that log records have been hardened on the secondary nodes. It should be as low as

possible, ideally no more than a couple milliseconds. Its duration depends on three key factors:

Network performance

Both log records and acknowledgement messages go through the network, so good network throughput is *essential* to support Availability Group replication. You can look at network performance counters, such as Bytes Received/sec, Bytes Sent/sec and Current Bandwidth, to analyze and troubleshoot possible network issues.

In some cases, you might build a separate network for Availability Groups, segregating replication from client traffic. If so, make sure that these networks are physically separated from each other; if they share the same physical LAN adapters, this topology would not provide you many benefits as the traffic from both networks will go to the same physical infrastructure.

I/O performance on secondary nodes

Synchronous replicas harden the log records in the log files before sending an acknowledgement back to primary. Insufficient I/O performance there will increase commit latency. Look at the log-write stalls in the `sys.dm_io_virtual_file_stats` view. You can troubleshoot I/O performance using the methods discussed in Chapter 3.

CPU bandwidth

Both the primary and secondary nodes need enough CPU bandwidth to handle replication. Make sure that servers are not overloaded and that schedulers are evenly balanced across NUMA nodes (Listing 2-4).

As a general rule, I do not recommend using readable synchronous replicas. Client queries add additional load and can impact replication throughput. Think about non-optimized queries that overload the I/O subsystem and thus increase log-write latency. It is usually better to build separate asynchronous replicas to scale the read workload instead.

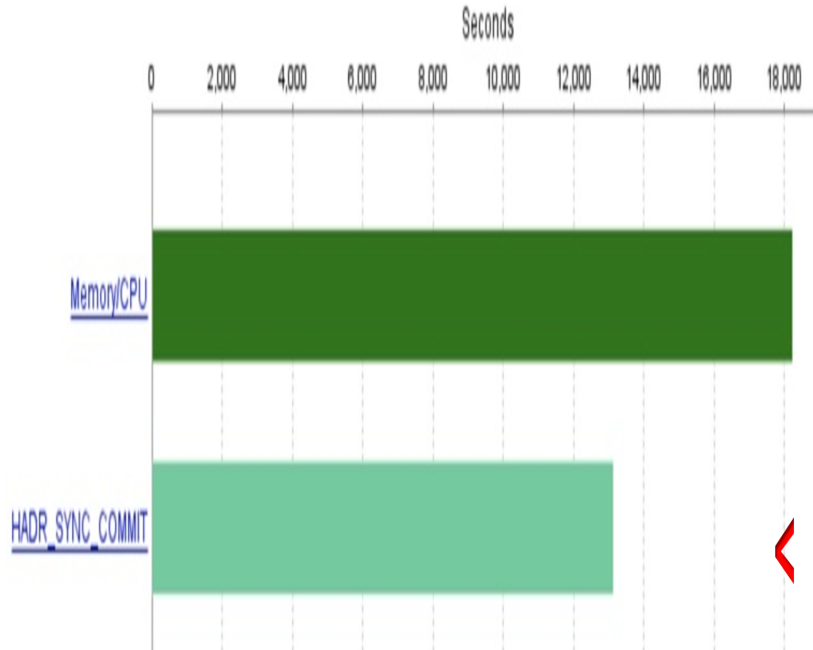
You can improve Availability Group performance by reducing the number of log records to process. I mentioned a few ways to do that in the previous chapter – all of them are applicable here.

Finally, if you are still using SQL Server 2012 or 2014, consider upgrading to the newer versions. SQL Server 2016 introduced performance enhancements in many areas, including Availability Groups. Newer versions of SQL Server will provide even better results.

Figure 12-5 shows the waits on one production server before and after SQL Server 2016 upgrade. Both snapshots were taken in the SolarWinds DPA application, with the server handling the same workload. As you can see, SQL Server 2016 reduced HADR_SYNC_COMMIT waits to less than a third of their previous levels. The upgrade also reduced CPU load by 35 percent without any further changes to the applications.

Top Waits |

| January 14, 2019 - 8:00AM to 9:00AM



< SQL Server 2012

Top Waits |

| January 21, 2019 - 8:00AM to 9:00AM

SQL Server 2016 >

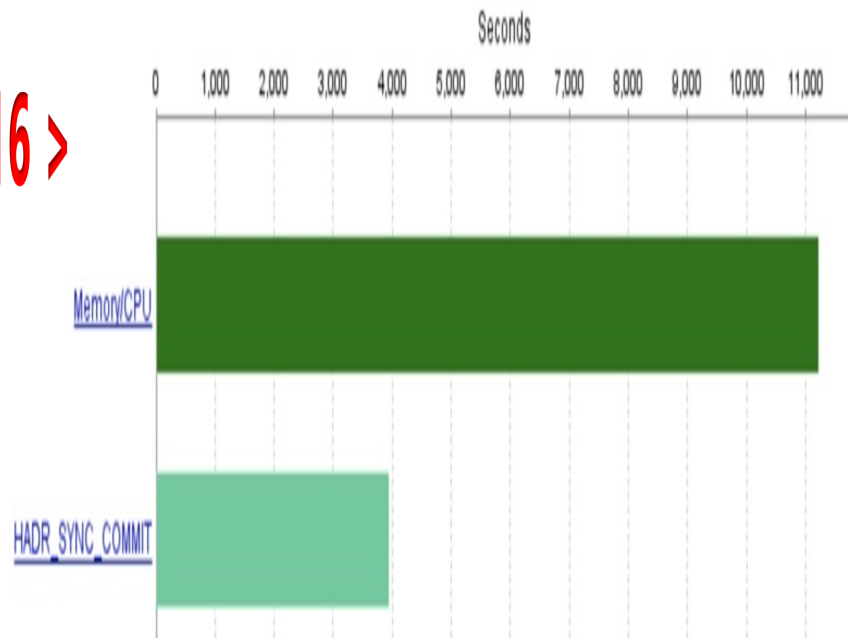


Figure 11-5. Waits before and after SQL Server 2016 upgrade

Usually, analyzing your network and I/O performance and your CPU load provides you enough information to troubleshoot bad Availability Group throughput. In some cases, however, you need to go further and look at the performance of individual operations. You can do this by analyzing Availability Group xEvents.

Availability Group Extended Events

SQL Server exposes large number of Extended Events (xEvents) you can use while troubleshooting Availability Group performance. Table 12-1 shows the most important ones.

*T
a
b
l
e*

*l
l
-
l
.
E
x
t
e
n
d
e
d*

*E
v
e
n
t
s*

*f
o
r*

*A
v*

*a
i
l
a
b
i
l
i
t
y*

*G
r
o
u
p*

*p
e
r
f
o
r
m
a
n
c
e*

*t
r
o
u
b
l
e*

*s
h
o
o
t
i
n
g*

xEvent	Location	Description
log_flush_start	Primary Secondary	Node starts hardening log records in transaction log files
log_flush_complete	Primary Secondary	Node finishes hardening log records in transaction log files
hadr_log_block_compression	Primary	Primary node compresses a log block
hadr_log_block_decompression	Secondary	Secondary node decompresses a log block
	Primary	Primary node captures log block for the replication. Mode column indicates the action:

hadr_capture_log_block		1: Log block is captured 2: Log block is enqueued into a send queue 3: Log block is dequeued and ready to be send 4: Log block is routed to proper replica
------------------------	--	---

ucs_connection_send_msg	Primary Secondary	A log block is sent to transport
-------------------------	----------------------	----------------------------------

hadr_transport_receive_log_block_message	Secondary	Secondary receives the log block. Mode column indicates the action: 1: Log block is received 2: Log block is enqueued in a working queue
--	-----------	--

hadr_send_harden_lsn_message	Secondary	Secondary is sending acknowledgement that a log block is hardened. Mode column indicates: 1: Message is created 2: Message is ready to be sent 3: Message is routed to primary
------------------------------	-----------	---

hadr_lsn_send_complete	Secondary	Acknowledgement has been sent
------------------------	-----------	-------------------------------

hadr_receive_harden_lsn_message	Primary	Primary receives an acknowledge message. Mode field indicates: 1: Message is received 2: Message is dequeued for processing
---------------------------------	---------	---

hadr_db_commit_msg_harden	Primary	Acknowledgement has been processed
---------------------------	---------	------------------------------------

hadr_db_commit_ mgr_harden_still_ waiting	Primary	Indicates that a commit acknowledgement has not been received after 2 seconds. This should not happen in normal circumstances and usually indicates an issue with Availability Groups.
---	---------	--

Figure 12-6 shows the sequence of events that occur during synchronous Availability Groups communication. On the primary node, SQL Server starts the log flush and simultaneously sends the log block to the secondary node. The secondary node decodes the block, hardens it in the log file, and then constructs and sends an acknowledgement message back to the primary.

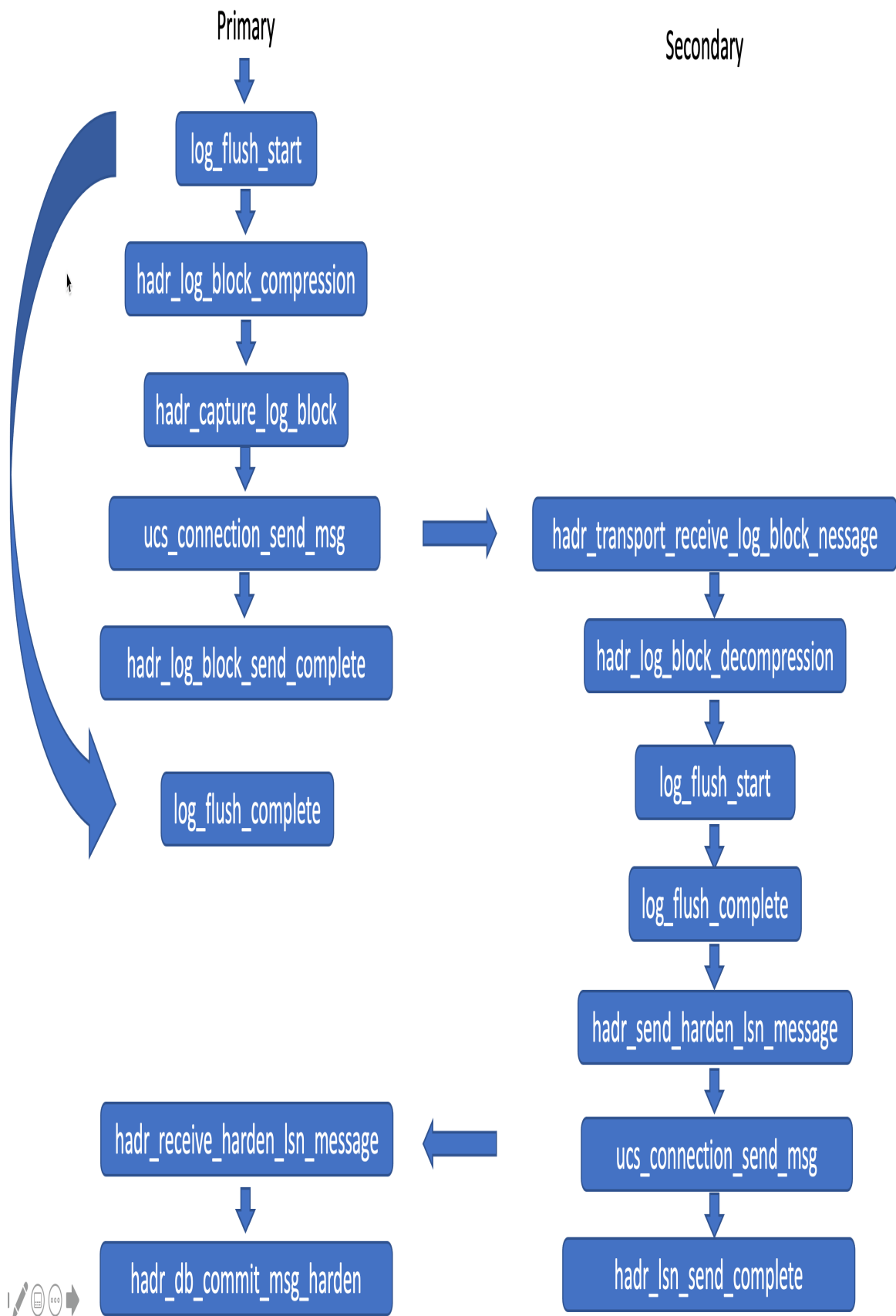


Figure 11-6. Availability Groups communication flow in synchronous mode

Listing 12-2 shows two xEvent sessions you can run to capture xEvents from Table 12-1. A word of caution: those sessions can collect a large number of events. Do not keep them running outside of the troubleshooting.

I am not including an `ucs_connection_send_msg` event here, because it introduces a lot of noise. Consider including it if you troubleshoot possible network performance issues.

Listing 12-2. Creating xEvent sessions for Availability Group performance troubleshooting.

```
-- Create on primary node
CREATE EVENT SESSION [AlwaysOn_Tracing_Primary] ON SERVER
ADD EVENT sqlserver.hadr_capture_log_block,
ADD EVENT sqlserver.hadr_db_commit_mgr_harden,
ADD EVENT sqlserver.hadr_db_commit_mgr_harden_still_waiting,
ADD EVENT sqlserver.hadr_log_block_compression,
ADD EVENT sqlserver.hadr_log_block_send_complete,
ADD EVENT sqlserver.hadr_receive_harden_lsn_message,
ADD EVENT sqlserver.log_flush_complete,
ADD EVENT sqlserver.log_flush_start
ADD TARGET package0.ring_buffer (SET max_events_limit=
(0),max_memory=(16384));
GO

-- Create on secondary node
CREATE EVENT SESSION [AlwaysOn_Tracing_Secondary] ON SERVER
ADD EVENT sqlserver.hadr_apply_log_block,
ADD EVENT sqlserver.hadr_log_block_decompression,
ADD EVENT sqlserver.hadr_lsn_send_complete,
ADD EVENT sqlserver.hadr_send_harden_lsn_message,
ADD EVENT sqlserver.hadr_transport_receive_log_block_message,
ADD EVENT sqlserver.log_flush_complete,
ADD EVENT sqlserver.log_flush_start
ADD TARGET package0.ring_buffer (SET max_events_limit=
(0),max_memory=(16384));
```

You can use the `log_block_id` and `database_id` fields in both sessions to correlate session data. There is a catch, though: in `hadr_send_harden_lsn_message`, `hadr_receive_harden_lsn_message`, and `hadr_lsn_send_complete` events, the `log_block_id` will be higher than in

previous events. This is due to how xEvents collects the data. The difference in values depends on the load in the database; however, it won't exceed 120.

Figure 12-7 shows the events collected in my test environment on the primary node. You may notice that, first, a few events were collected out of order on the target. They were all fired very rapidly and have the same timestamp.

name	timestamp	mode	log_block_id
hadr_capture_log_block	2021-05-10 19:09:31.4952867	1	158913901495
hadr_capture_log_block	2021-05-10 19:09:31.4952867	2	158913901495
log_flush_start	2021-05-10 19:09:31.4952867	NULL	158913901495
hadr_log_block_compression	2021-05-10 19:09:31.4952867	NULL	158913901495
hadr_log_block_send_complete	2021-05-10 19:09:31.4952867	NULL	158913901495
log_flush_complete	2021-05-10 19:09:31.4952867	NULL	158913901495
hadr_receive_harden_lsn_message	2021-05-10 19:09:31.4962860	1	158913901496
hadr_receive_harden_lsn_message	2021-05-10 19:09:31.4962860	2	158913901496
hadr_db_commit_mgr_harden	2021-05-10 19:09:31.4962860	NULL	NULL

Figure 11-7. xEvents from primary node

Figure 12-8 shows the events collected on the secondary node. There is a small gap in timestamps with the primary node, because I was unable to synchronize time perfectly between the servers.

name	timestamp	mode	log_block_id
hadr_transport_receive_log_block_messa...	2021-05-10 19:09:31.5116162		1 158913901495
hadr_transport_receive_log_block_messa...	2021-05-10 19:09:31.5116664		2 158913901495
hadr_log_block_decompression	2021-05-10 19:09:31.5116888	NULL	158913901495
hadr_log_block_decompression	2021-05-10 19:09:31.5116967	NULL	158913901495
log_flush_start	2021-05-10 19:09:31.5117582	NULL	158913901495
log_flush_complete	2021-05-10 19:09:31.5123777	NULL	158913901495
hadr_send_harden_lsn_message	2021-05-10 19:09:31.5124051		1 158913901496
hadr_send_harden_lsn_message	2021-05-10 19:09:31.5124295		2 158913901496
hadr_send_harden_lsn_message	2021-05-10 19:09:31.5124316		3 158913901496
hadr_lsn_send_complete	2021-05-10 19:09:31.5125995	NULL	158913901496

Figure 11-8. xEvents from secondary node

The timing of individual operations can help you identify possible bottlenecks. For example, a slow disk subsystem on the secondary would introduce a delay between the `log_flush_start` and `log_flush_complete` events. Insufficient CPU throughput may prolong `hadr_log_block_compression` and `hadr_log_block_decompression` events when compression is used. Analyze the data and cross-check it with other metrics.

Compression behavior varies in different versions of SQL Server. SQL Server 2012 and 2014 compressed all Availability Group traffic by default. In SQL Server 2016 and above, however, compression is used only in asynchronous communication. Synchronous commit, on the other hand, does

not compress the traffic. Nevertheless, you can still see compression and decompression xEvents generated even when the action is not performed.

There are three trace flags that you can use to change compression behavior in SQL Server 2016 and above. They may shift the workload and bottlenecks between CPU and network, so use them with care!

T9592

This trace flag enables compression of the traffic between synchronous replicas. Consider enabling it if you have to support synchronous replicas over a slow network. Keep in mind that compression adds CPU overhead on both nodes and may increase HADR_SYNC_COMMIT latency in fast networks.

T1462

This trace flag disables compression of the traffic between asynchronous replicas. This may reduce CPU load on the nodes in very busy OLTP systems, at the cost of additional network traffic.

T9567

SQL Server does not use compression when it performs automatic seeding of the new nodes in Availability Groups. The T9567 flag allows you to enable it. This can speed the up automatic seeding process, but at the cost of additional CPU load on the primary node.

Finally, modern versions of SSMS provide tools to get similar latency data. You can do this by clicking *Collect Latency Data* in the Availability Group dashboard. This action creates and runs xEvent sessions on the nodes that collect Availability Group events.

The sessions will run for two minutes, using a file target to store the data. After that, SQL Server processes the data and allows you to access it through the standard reports in the Availability Group popup menu in SSMS Object Explorer.

There are some limitations to SSMS data collection: First, it works only when you connect to the servers using Windows Authentication and have sysadmin permissions on all replicas. It also relies on SQL Server Agent running on all servers. Finally, it uses file targets, which may add an overhead on busy servers with heavy Availability Groups traffic. Nevertheless, SSMS data collection is often easier than manual xEvent analysis.

Asynchronous Replication and Readable Secondaries

In contrast to synchronous commit mode, in asynchronous mode, primary does not wait for confirmation that the log records have been hardened on the secondaries. A transaction becomes committed when the commit log record is saved in the primary's transaction log

Figure 12-9 shows replication data flow in asynchronous mode. The size of the send queue on primary dictates the possible data loss and recovery point objective (RPO) in the event of an SQL Server crash.

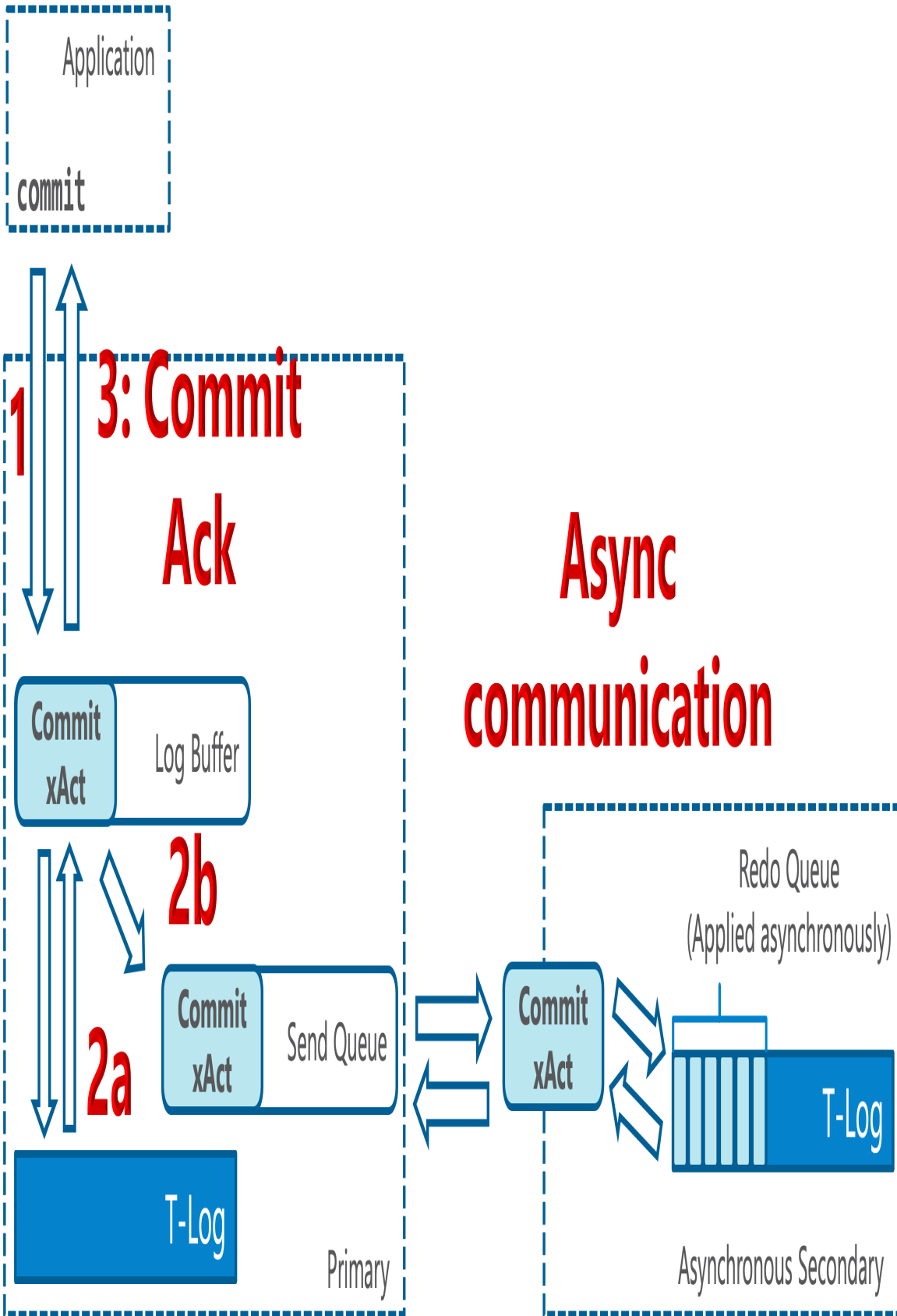


Figure 11-9. Asynchronous commit data flow

The Enterprise Edition of SQL Server allows you to scale the read workload by running queries on the secondary nodes. This is a great feature that helps to improve system throughput and reduces the load on the primary node. There are a few factors you need to consider, however.

As I have mentioned, be extremely careful when querying nodes that use synchronous replication. The overhead of the queries may increase replication latency and HADR_SYNC_COMMIT waits, impacting the primary node. It is always better to leave synchronous nodes alone and build asynchronous nodes to scale the read workload.

Regardless of the replication mode, the data on secondaries will always be behind primary, because the redo process that applies changes to the database is asynchronous. In normal circumstances the lag may be very small—measured in milliseconds or even microseconds; however, it may grow during log-intensive operations such as index maintenance or large batch processing.

Do not use secondary nodes for critical queries that need up-to-date data. Remember that you cannot guarantee that the lag is always going to be small. Neither should you split the read and write queries across the nodes within a single business transaction. This will lead to inconsistent results.

Next, let's look at another, less well known issue related to readable secondaries that can affect the primary node in quite unexpected ways.

The Impact of the Readable Secondaries

SQL Server always uses the SNAPSHOT isolation level for queries on the secondary nodes, ignoring SET TRANSACTION ISOLATION LEVEL statements and isolation-level table hints. It allows you to eliminate the possibility of readers being blocked by writers. This happens even if you do not enable the ALLOW_SNAPSHOT_ISOLATION database option.

Using the SNAPSHOT isolation level also means that SQL Server will use row versioning on the secondary nodes. As you remember from Chapter 8,

SQL Server will start using version store in tempdb and also maintain 14-byte version store pointers in modified data rows.

In the Availability Group infrastructure, the databases on the primary and secondary nodes are exactly the same, so it is impossible to maintain row versioning on secondary nodes only. SQL Server needs to allocate space for 14-byte version store pointers on the primary node for the databases to match, even if optimistic isolation levels are not enabled in the database.

SQL Server does not use tempdb version store on primary nodes if you don't enable the `READ_COMMITTED_SNAPSHOT` or `ALLOW_SNAPSHOT_ISOLATION` options. Nevertheless, it adds extra 14 bytes to the data rows during data modifications, which may lead to additional page splits and index fragmentation.

Unfortunately, readable secondaries also introduce another phenomenon that is less well known: long-running `SNAPSHOT` transactions on secondary nodes may defer the ghost and version store cleanup tasks on the primary. Such transactions work with a snapshot of the data at the time the transaction started. SQL Server cannot remove deleted rows and reuse the space because of the possibility that `SNAPSHOT` transaction will need to access the old versions of the rows.

The same applies to large send and redo queues. SQL Server cannot clean up deleted rows, because the secondaries could start a `SNAPSHOT` transaction before replaying the ghost cleanup log records. This can become an issue if the replica goes offline or constantly falls behind in applying the changes.

Take a look at the example in Listing 12-3. This will create two tables in the database. Table T1 will have 65,536 rows and use 65,536 pages: one row per data page.

Listing 12-3. Readable secondaries: Table creation

```
CREATE TABLE dbo.T1
(
    ID INT NOT NULL,
    Placeholder CHAR(8000) NULL,
    CONSTRAINT PK_T1 PRIMARY KEY CLUSTERED(ID)
);
```

```

CREATE TABLE dbo.T2
(
    Col INT
);

;WITH N1(C) AS (SELECT 0 UNION ALL SELECT 0) -- 2 rows
,N2(C) AS (SELECT 0 FROM N1 AS T1 CROSS JOIN N1 AS T2) -- 4 rows
,N3(C) AS (SELECT 0 FROM N2 AS T1 CROSS JOIN N2 AS T2) -- 16 rows
,N4(C) AS (SELECT 0 FROM N3 AS T1 CROSS JOIN N3 AS T2) -- 256 rows
,N5(C) AS (SELECT 0 FROM N4 AS T1 CROSS JOIN N4 AS T2 ) -- 65,536
rows
,IDS(ID) AS (SELECT ROW_NUMBER() OVER (ORDER BY (SELECT NULL))
FROM N5)
INSERT INTO dbo.T1(ID)
    SELECT ID FROM IDS;

```

Next, let's start a transaction on the secondary node and run the query against table T2, as shown in Listing 12-4. Even though I am using explicit transactions, the same behavior will occur if there is a long-running statement in an autocommitted transaction.

Listing 12-4. Readable Secondaries: Starting transaction on secondary node

```

BEGIN TRAN
    SELECT * FROM dbo.T2;

```

Next, using the code in Listing 12-5, delete all data from T1 table and then run the query that will do the clustered index scan on the primary node.

Listing 12-5. Readable Secondaries: Deleting data and performing clustered index scan

```

DELETE FROM dbo.T1;
-- Wait 1 minute
WAITFOR DELAY '00:01:00.000';

SET STATISTICS IO ON
SELECT COUNT(*) FROM dbo.T1;
SET STATISTICS IO OFF
--Output: Table 'T1'. Scan count 1, logical reads 65781

```

Even though the table is empty, the data pages have not been deallocated. This leads to significant I/O overhead on the primary node.

Finally, let's look at the index statistics, using the code from Listing 12-6.

Listing 12-6. Readable Secondaries: Checking index stats

```
SELECT
    index_id, index_level, page_count
    , record_count, version_ghost_record_count
FROM
    sys.dm_db_index_physical_stats
    (
        DB_ID()
        , OBJECT_ID(N'dbo.T1')
        , 1
        , NULL
        , 'DETAILED'
    );
```

Figure 12-10 shows the output of the query. The leaf index level shows 65,536 rows in the `version_ghost_record_count` column. This column contains the number of ghosted rows that cannot be removed due to active transactions that rely on row versioning. In our case, this transaction runs on a different (secondary) node.

	index_id	index_level	page_count	record_count	version_ghost_record_count
1	1	0	65536	0	65536
2	1	1	243	65536	0
3	1	2	1	243	0

Figure 11-10. Index Statistics

This behavior may increase I/O and CPU load, because SQL Server needs to scan ghosted data rows during query execution. It also increases the size of tempdb since the version store there is not being cleaned up.

The impact on I/O and CPU load may be especially high if the system is implementing message processing based on in-database queuing. The tables that store the messages are usually small, but the data there is extremely volatile. Ghost records accumulate quickly, driving up CPU and I/O load.

I encountered this issue for the first time in such a scenario. One of the tables was used for message processing, handling about 100 inserts and

deletes every second. Runaway reporting transaction on the secondary node increased CPU load on primary by 30 percent overnight without any changes in the system workload. By the time I detected the issue, the table had more than a million data pages storing just a handful of rows. The queries that read the messages scanned all those data pages, driving the CPU load up.

Unfortunately, it is very common for people to offload non-optimized reporting queries to secondary nodes without understanding the potential consequences. Remember this and monitor ghost cleanup lag with the `water_mark_diff` column from Listing 12-1. The alert threshold will depend on your system workload. Analyze the overhead of deferred ghost cleanup and set it accordingly.

Last but not least, do not enable readable secondaries unless you are querying them. They introduce performance impact and increase licensing cost. There is no need to incur that unless you are using the feature.

Parallel Redo

In SQL Server 2012 and 2014, the redo process in Availability Groups used one thread per database. In high-end OLTP systems with extremely high rate of modifications, that led to insufficient redo throughput, making it hard to use Availability Groups at all.

Starting with SQL Server 2016, the redo process may become parallel with SQL Server using multiple threads to replay the log records. The number of threads that perform redo process depends on the number of CPUs on the server and the number of databases in Availability Group.

In Availability Group with multiple databases, only the first six use parallel redo. Their order usually depends on when the databases joined the Availability Group. Unfortunately, this feature is poorly documented and may change in future versions of SQL Server.

NOTE

SQL Server 2016 and 2017 does not use parallel redo with the databases that use In-Memory OLTP. This restriction has been removed in SQL Server 2019.

Parallel redo is a great feature—*when it works*. Unfortunately, it doesn't always work. I've experienced quite a few cases it suddenly *stopped* working, leading to significant performance degradation of redo process. Usually it occurs in busy OLTP systems that use readable secondaries.

The issue usually presents as elevated CPU usage, with constantly growing redo queues. You are also likely to see multiple waiting tasks, usually including one or more of the following wait types in the `sys.dm_os_waiting_tasks` and the `sys.dm_exec_requests` views output:

- DIRTY_PAGE_TABLE_LOCK
- DPT_ENTRY_LOCK
- PARALLEL_REDO_TRAN_TURN
- PARALLEL_REDO_FLOW_CONTROL

You can disable parallel redo with server-level trace flag T3459, using the command `DBCC TRACEON(3459,-1)`. You can make this change online without restarting the server. However, disabling the trace flag and switching back to parallel redo does requires restarting in the builds prior to SQL Server 2017 CU9, SQL Server 2016 SP2 CU2 and SQL Server 2016 SP1 CU10.

As usual, apply the latest cumulative update to SQL Server. There are many fixes related to parallel redo, especially in SQL Server 2016 and 2017.

Troubleshooting Failover Events

Automatic failover is a great feature that improves High Availability. It also introduces a few tasks. When an unintended failover occurs, you'll need to

find the root cause. In other cases, you may need to troubleshoot why failover did *not* occur as expected.

The troubleshooting strategy is the same in both cases: collect information, then analyze it afterwards. I will discuss the sources of that information later in this section. First, though, let's get a high-level overview of how SQL Server interacts with Windows Failover Cluster.

Availability Groups and Windows Server Failover Cluster

Availability Groups rely on WSFC services to support automatic failovers. They become the cluster resource in the cluster, which manages them similarly to other services. This means that if WSFC has a problem—for example, if it loses the quorum—underlying Availability Groups will also be affected.

There are two key checks in cluster resource management: `IsAlive` and `IsHealthy`. The cluster executes `IsAlive` and `IsHealthy` checks frequently, validating that the resource is online and it is healthy, operating as expected.

To put things in perspective, if an `IsAlive` check fails, the cluster may initiate the failover or shut down the SQL Server instance. It could also trigger those actions if it decides that SQL Server is unhealthy, based on the results of the `IsHealthy` check; if it finds a large number of access violation errors; or if the server shows extremely high and prolonged resource consumption.

Both checks are done by SQL Server resource DLL, which constantly communicates with the SQL Server instance. The `IsAlive` check is done through the Shared Memory protocol, which allows two processes to share memory for communication. The frequency of communication is controlled by the `LeaseTimeout` property in the Availability Group cluster resource. By default, `LeaseTimeout` is set to 20,000 milliseconds; the `IsAlive` check runs every 5 seconds, just 25% of that value.

The lease mechanism exists only on the primary node, and you can think of it as the Availability Group's heartbeat. When the cluster does not receive confirmation that the lease is active, it considers the lease to be expired and considers Availability Group unhealthy. It stops accepting write requests in

order to avoid *split brain*, a condition when multiple nodes behave as the primary replica, accepting write requests. Next, the Availability Group tries to failover, assuming WSFC itself is healthy.

The IsHealthy check, on the other hand, relies on a stored procedure in the sp_server_diagnostics system that provides health information about the general system and Availability Group, as well as SQL Server component status and resource utilization. The frequency of the check is controlled by the HealthCheckTimeout property in the Availability Group cluster resource. By default, it is 30,000 milliseconds; the frequency of the IsHealthy check is one-third of that value, or 10 seconds.

Another property, FailureConditionLevel, specifies the failure condition for the health check. It usually has the following values:

1. OnServerDown: The health check validates that Availability Group is online. The validation succeeds if the IsAlive check passes.
2. OnServerUnresponsive: The health check fails if no data is received from sp_server_diagnostics procedure within the time specified in HealthCheckTimeout.
3. OnCriticalServerError: The health check fails if any of the major SQL Server components returns an error. This is the default setting.
4. OnModerateServerError: The health check fails if the resource component returns an error.
5. OnAnyQualifiedFailureConditions: The health check fails if the query processing component returns an error.

As you can guess, the higher the value, the more often failover will be triggered. In most cases, you don't need to change the FailureConditionLevel parameters; however, you may decide to temporarily decrease them when you are troubleshooting failovers due to extremely high resource utilization. I do not recommend increasing the parameters unless you want to failover at the first sign of a possible issue, even when the issue could resolve itself.

You can increase LeaseTimeout if you are experiencing connectivity issues between cluster nodes. Use the following formula to determine the maximum possible value for the property:

$$\text{LeaseTimeout} \leq (\text{SameSubnetDelay} * \text{SameSubnetThreshold}) * 2$$

SameSubnetDelay and SameSubnetThreshold are cluster-level properties that indicate the cluster's heartbeat frequency and the number of times it can be missed before the cluster considers it unhealthy. You can tune those parameters, as well as LeaseTimeout.

NOTE

The CrossSubnetDelay and CrossSubnetThreshold parameters control the heartbeat parameters in the cross-subnet cluster. Make sure that the SameSubnetDelay and SameSubnetThreshold values do not exceed them.

Increase HealthTimeout if you do not want to failover when short spikes in resource utilization cause SQL Server to become unresponsive. Using virtualization and VM backup software to pause virtual machines while the backup is in progress can also help. Be careful: this may delay the failover in the cluster.

As I have noted, SQL Server on Linux relies on an external service, Pacemaker, for failover support. Conceptually, Pacemaker behaves similarly to WSFC; however, SQL Server cannot communicate with Pacemaker directly. The implementation relies on a polling mechanism, with Pacemaker regularly querying SQL Server and database states.

Troubleshooting Failovers

So why do failovers occur? You may have already guessed some of the reasons. The most common are:

Hardware issues

WSFC issues, such as cluster quorum loss due to network or connectivity problems

Lease expiration, which can be triggered by WSFC issues like extremely high load on the server or an unresponsive OS

Health check timeouts, often triggered by high SQL Server load (CPU and/or disk), high numbers of access violation errors and thread dumps, an unresponsive OS, a frozen VM, or other factors

Fortunately, SQL Server and your OS provide a lot of information to help you research the cause. Let's look at several of them now.

AlwaysOn_health xEvent Session

The AlwaysOn_health xEvent session is created when you provision Availability Groups. This session contains the set of events that track the health of Availability Groups, the state and role changes of the nodes, and Availability Group-related DDL operations. It also includes several events that track high-severity errors: for example, the `availability_group_lease_expired` event is generated when the lease expires.

The AlwaysOn_health session is the great place to start troubleshooting. It may not show you the root cause of the incident; however, it gives you information about what happened with Availability Group and when, as well as what actions were taken. Finally, it allows you to pinpoint conditions when manual failover was triggered either advertently or by human error.

SQLDIAG files

The WSFC, as noted, uses a `sp_server_diagnostics` stored procedure as part of IsHealthy cluster checks. This stored procedure provides information about SQL Server component health and resource utilization, such as memory and CPU usage, the state of the deadlock monitor, access violations and thread dumps, and many others. It also provides general health information on all Availability Groups on the server.

The data from `sp_server_diagnostics` is captured by a hidden xEvent session and stored in XEL files in the SQL Server Log folder. These files

are extremely useful during failover troubleshooting, because they can tell you if some of the server components were overloaded or not healthy. The naming convention for the files consists of the server and instance name followed by an SQLDIAG string.

system_health xEvent Session

The system_health xEvent session is another session that is created and runs by default in SQL Server. It provides information about general SQL Server health and high-severity errors along with component health and diagnostics data. The latter is similar to the information in the SQLDIAG files; however, the metrics are aggregated in larger intervals of 5 minutes.

Cluster Log

The cluster log is a key source of data for the troubleshooting. It contains detailed information about the errors that led to failover conditions and can help you pinpoint quorum and other cluster-related issues.

SQL Server Log

The SQL Server log provides information about critical errors, the status of failovers, and the state of Availability Group replicas. Looking at the events around the time of incident can be useful when troubleshooting.

System Logs

System logs and the Windows Event Log can provide details about critical system conditions and errors, such as hardware failures.

Keep in mind that SQL Server and Windows Failover Cluster roll the logs and xEvent files over. You need to collect the data shortly after failover occurs, or it may disappear. When you have collected the data, analyze it holistically to identify the problem.

When Failover Does Not Occur

In some cases, you need to troubleshoot the opposite problem: why Availability Group did not failover automatically when it should have. Here are a few common reasons this happens.

Incorrect Configuration

The Availability Group has not been configured correctly: for example, it does not have Automatic Failover set or it has databases that joined Availability Group on a single node only.

Communication Issues

Failover may not occur if primary node loses connectivity to the secondary node to which it is supposed to failover. Failover also requires the databases to be fully synchronized, so it might not happen if the secondary has not caught up with replication after being offline.

Cluster Issues

Some WSFC issues may prevent automatic failover, such as the SQL Server resource DLL losing its connection to SQL Server.

Availability Group Exceeded Maximum Failover Threshold

Each cluster resource has a set of properties that specify the maximum number of failovers that can occur within a specific time period. In Availability Groups, the default maximum number of failovers is $N - 1$, where N is the number of nodes in the cluster. The default time period is 6 hours.

You can change these values in the properties of Availability Group cluster resource.

The troubleshooting approach is very similar to what you did to identify an unwanted failover: look at the data from the Cluster and SQL Server logs, SQLDIAG files, AlwaysOn_health sessions, and system_health sessions and reconstruct what happened during the incident. This usually provides enough data to help you understand the cause of the problem.

Summary

AlwaysOn Availability Groups is the most popular High Availability technology used in SQL Server. It avoids a situation where storage can be the single point of failure and, in Enterprise Edition, allows you to scale read workload across multiple replicas.

Availability Groups do not replicate instance-level objects, such as logins and jobs. Synchronize them across the nodes to allow systems to operate after failovers.

Replication in Availability Groups is based on the stream of transaction log records. Monitor send and redo queues and set alerts to notify you when they grow. Large queues may lead to data loss and prolong recovery time. They can also prevent log truncation and introduce other performance issues.

Synchronous commit mode helps avoid data loss, but the tradeoff is additional commit latency during replication. Analyze the HADR_SYNC_COMMIT wait time and reduce it as much as possible.

Readable secondaries allow you to scale read workload. However, they may defer ghost and version store cleanup tasks on the primary node, increasing CPU and I/O load there. Avoid long-running transactions on secondaries and monitor workload there.

When you troubleshoot failover events, look at the data from AlwaysOn_health and system_health xEvent sessions; SQLDIAG files; and OS, SQL Server, and Cluster logs. They usually contain enough information to diagnose the problem.

In the next chapter, I'll discuss several other common wait types.

Troubleshooting Checklist

Make sure that instance-level objects are synchronized across Availability Group nodes

Perform failover/HA testing if possible

Set up Availability Group monitoring and alerting to cover queue sizes, replication latency, ghost cleanup lag, and failover events

Analyze the impact of commit latency if synchronous commit is used; review the HADR_SYNC_COMMIT resource wait time and troubleshoot replication performance if needed

Check if readable secondaries are enabled. Evaluate the impact of read-only queries on synchronous readable replicas.

Disable readable secondaries if they are not being used.

Chapter 12. Other Notable Wait Types

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 13th chapter of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at dmitri@aboutsqlserver.com.

This rather small chapter covers several wait types I have not discussed yet. I will start with the `ASYNC_NETWORK_IO` wait, which occurs when the client does not consume data from SQL Server fast enough. Next, I will talk about the `THREADPOOL` wait and the dangerous condition of worker thread starvation. After that, I will address backup-related wait types and ways to improve backup performance.

I will conclude the chapter with an overview of `OLEDB` and a few other preemptive wait types that occur when SQL Server calls OS API switching to preemptive execution mode.

ASYNC_NETWORK_IO Waits

The `ASYNC_NETWORK_IO` is a common wait I see in nearly every system. Inexperienced engineers usually guess by the wait type name,

associating `ASYNC_NETWORK_IO` with bad network performance. This wait, however, indicates a much broader condition that occurs when SQL Server has to wait for the client application to consume data.

Slow networks can definitely trigger that condition, but more often than not, the cause is inefficient client application design. If the application reads and processes data row by row, this forces SQL Server to wait during processing.

Listing 13-1 shows the code pattern that will trigger the issue. The client application consumes and processes rows one by one, keeping `SqlDataReader` open. The SQL Server worker waits for the client to consume all rows and generates the `ASYNC_NETWORK_IO` wait in the meantime.

Example 12-1. Listing 13-1. Code pattern that triggers `ASYNC_NETWORK_IO` wait

```
using (SqlConnection connection = new
SqlConnection(connectionString))
{
    SqlCommand command = new SqlCommand(cmdText, connection);
    connection.Open();
    using (SqlDataReader reader = command.ExecuteReader())
    {
        while (reader.Read())
            ProcessRow((IDataRecord) reader);
    }
}
```

The right approach here is to read all rows first as quickly as possible, then process them afterward, as shown in Listing 13-2. The client may need to batch the operation, if the size of the data is very large and the clients do not have enough memory to cache it.

Example 12-2. Listing 13-2. Code pattern that removes `ASYNC_NETWORK_IO` waits

```
List<Orders> orderRows = new List<Orders>();
using (SqlConnection connection = new
SqlConnection(connectionString))
{
    SqlCommand command = new SqlCommand(cmdText, connection);
```

```

connection.Open();
using (SqlDataReader reader = command.ExecuteReader())
{
    while (reader.Read())
        orderRows.Add(ReadOrderRow((IDataRecord) reader));
}
}
ProcessAllOrderRows(orderRows);

```

My favorite way to prove that row-by-row processing triggers an `ASYNC_NETWORK_IO` wait is running a small demo in SSMS. Connect to the local SQL Server instance—SSMS will use the *shared memory* protocol, which does not involve any network traffic. Next, clear the wait statistics and run the `SELECT *` statement from the table, selecting several thousand rows. When you check the waits after the execution, you'll see `ASYNC_NETWORK_IO` at the top of the list despite the total absence of network traffic.

Next, enable the *Discard results after execution* setting in the Tools/Options/Query Results/SQL Server/Results to Grid options form and repeat the test. You'll see that the `ASYNC_NETWORK_IO` wait is no longer present. The reason you saw the wait in the first test is SSMS inefficiency, which populates the results grid on a row-by-row basis. This implementation is slow, and SQL Server waits for each row to be displayed in the grid, generating the wait.

When this wait is present in significant amounts, analyze network performance first. Review network topology (remember, network throughput depends on the slowest component) and check network-performance counters and metrics.

When you are confident that network performance is not an issue, analyze the situation with client applications. You might detect inefficient code with the `sys.dm_os_waiting_tasks`, `sys.dm_exec_requests`, `sys.dm_exec_sessions`, and `sys.dm_exec_connections` views. However, in some cases, it may not be feasible or even possible to change the client code.

In my list, the `ASYNC_NETWORK_IO` waits belong to the “it depends” category. I sometimes ignore them if they are not very significant and the system is not operating under an extremely heavy load.

Don’t take this the wrong way, though: this wait is not benign. It consumes workers on the server and introduces memory and CPU overhead. Nevertheless, it is not the biggest fish to fry. I usually get better ROI by focusing on other areas.

THREADPOOL Waits

In contrast to `ASYNC_NETWORK_IO` waits, `THREADPOOL` waits need to be investigated even when you see only a few. They indicate that SQL Server does not have enough worker threads to assign to tasks. When this happens, clients are unable to connect to SQL Server and will get *connection timeout* errors, as if SQL Server was down.

One situation where this may happen is when you have long blocking chains, usually with schema modification (Sch-M) locks involved. Follow along with me as we go through an example that emulates this condition.

First, let’s reduce the number of workers on the server by running the code from Listing 13-3. This code also enables a remote *dedicated admin connection* (DAC) – more about that later.

WARNING

Do *not* execute this scenario on production servers. It will bring the server down! Also, reset the *max worker threads* value back to 0 after the test is complete.

Example 12-3. Listing 13-3. Enabling a remote dedicated admin connection and reducing the number of workers

```
EXEC sys.sp_configure N'max worker threads', N'128';  
GO  
EXEC sp_configure 'remote admin connections', 1;  
GO
```

```
RECONFIGURE
GO
```

Next, let's create a small table and insert a row into it acquiring an intent-exclusive (IX) lock there. Run the code in Listing 13-4; do not commit the transaction and do not close the session afterwards.

Example 12-4. Listing 13-4. THREADPOOL wait: Creating a test table and acquiring an intent-exclusive lock

```
CREATE TABLE dbo.ThreadPoolDemo
(
    Col1 INT
);
GO
BEGIN TRAN
    INSERT INTO dbo.ThreadPoolDemo VALUES(0);
```

Next, let's open another session and run the ALTER TABLE statement shown in Listing 13-5. This statement will be blocked because Sch-M and IX locks are incompatible.

Example 12-5. Listing 13-5. THREADPOOL wait: Altering the table

```
-- Run in the different session from Listing 13-4
ALTER TABLE dbo.ThreadPoolDemo ADD Col2 INT;
```

Now, any other sessions trying to access the table will be blocked due to the presence of the Sch-M lock request in the queue. They will be suspended and wait, consuming workers on the server.

Let's emulate this condition by running the Windows Batch script shown in Listing 13-6. This script will open multiple sessions that try to select data from the table. You may need to change the server and database names in your environment.

Example 12-6. Listing 13-6. THREADPOOL wait: Generating multiple sessions and blocking

```
@ECHO OFF
SET query="SELECT * FROM SQLServerInternals.dbo.T"
SET p1=-S
SET server=.
SET p3=-E
SET p4=-Q
```

```
FOR /L %%I IN (1,1,150) DO (  
    START "" "sqlcmd.exe" %p1% %server% %p3% %p4% %query%  
)
```

This will consume all available workers. Now, if you try to connect to the server from SSMS, you'll get the connection error shown in Figure 13-1 because there are no available workers in the system to pick up the connection.

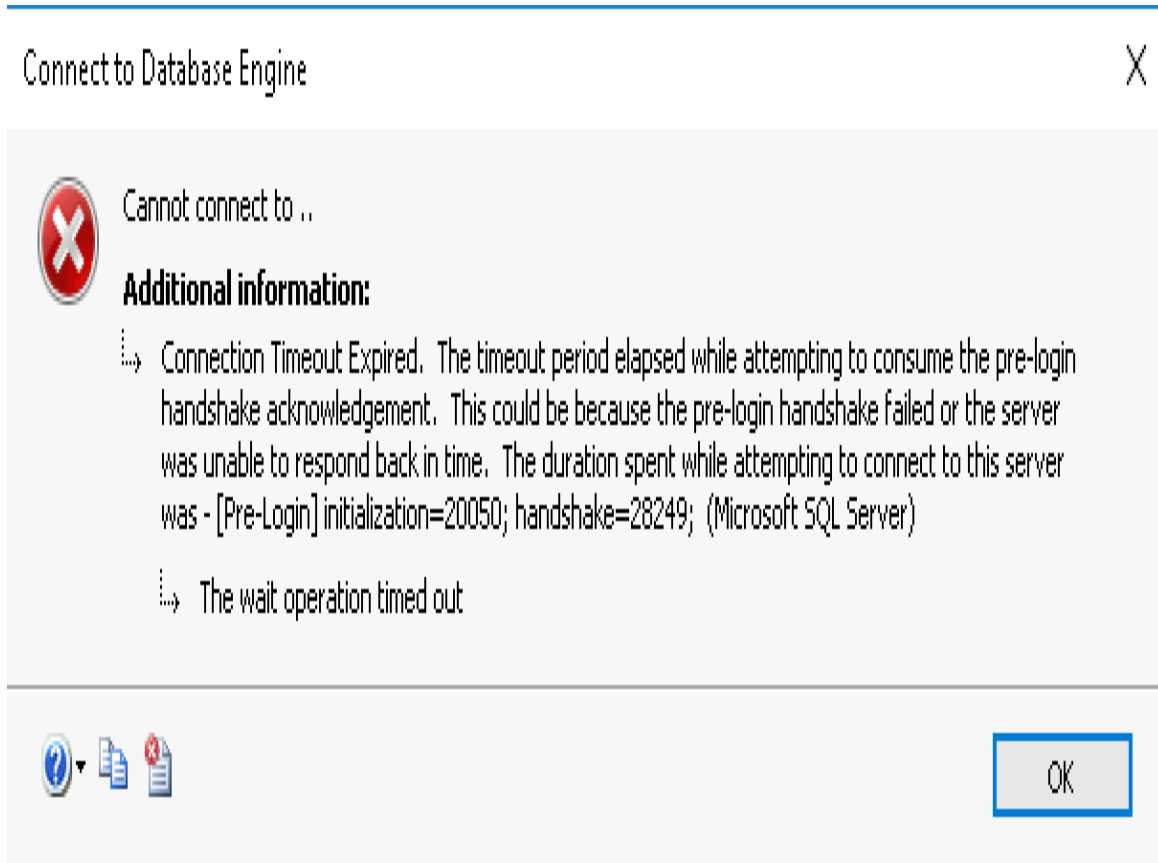


Figure 12-1. Connection error when worker pool is exhausted

When this condition occurs (and in other cases when the server stops responding), you can connect to it through a dedicated admin connection. SQL Server reserves a private scheduler and a small amount of memory for DAC to support those troubleshooting scenarios.

You can connect to SQL Server with DAC by using the ADMIN: server name prefix in the SSMS connection box or with the -A option in sqlcmd utility. Only members of the sysadmin server role are allowed to connect,

and only one session at a time can use a DAC connection. Do not connect to DAC in SSMS Object Explorer, which requires its own connection. Similarly, disable the Intellisense feature in the Query Window before connecting.

By default, DAC is available only locally. In some cases, when a server is completely overloaded, OS may not respond, preventing you from utilizing it. You always need to enable remote access to DAC during the initial server setup, as I did in Listing 13-3.

Now, connect to the server through DAC and query the `sys.dm_os_waiting_tasks` view with the code shown in Listing 13-7. Note that I filtered out system sessions from the output to focus on the blocking chain and THREADPOOL waits. I would not do that if I were troubleshooting an ongoing issue.

Example 12-7. Listing 13-6. THREADPOOL wait: Querying sys.dm_os_waiting_tasks view

```
SELECT
    session_id
    ,wait_type
    ,wait_duration_ms
    ,blocking_session_id
    ,resource_description
FROM
    sys.dm_os_waiting_tasks WITH (NOLOCK)
WHERE
    (session_id > 50 OR session_id IS NULL) AND
    wait_type NOT IN (N'CLR_AUTO_EVENT',N'QDS_ASYNC_QUEUE'
        ,N'XTP_PREEMPTIVE_TASK',N'BROKER_RECEIVE_WAITFOR'
        ,N'QDS_PERSIST_TASK_MAIN_LOOP_SLEEP')
ORDER BY
    session_id,
    wait_duration_ms DESC;
```

Figure 13-2 shows the output of the code. You can see the large number of sessions waiting for a schema stability lock on the table, generating a LCK_M_SCH_S wait type. They are blocked by the session with an ALTER TABLE statement, due to incompatibility between the Sch-S and Sch-M locks. The ALTER TABLE, in turn, is blocked by the session with SPID=55 that executed an INSERT statement in the uncommitted

transaction. You can terminate either the root blocker or the ALTER TABLE session with the KILL command to solve the problem.

	session_id	wait_type	wait_duration_ms	blocking_session_id	resource_description
9	NULL	THREADPOOL	95069	NULL	threadpool id=scheduler1d7..
10	NULL	THREADPOOL	92985	NULL	threadpool id=scheduler1d7..
11	NULL	THREADPOOL	91421	NULL	threadpool id=scheduler1d7..
12	NULL	THREADPOOL	82514	NULL	threadpool id=scheduler1d7..
13	52	LCK_M_SCH_S	106991	64	objectlock lockPartition=0..
14	60	LCK_M_SCH_S	106991	64	objectlock lockPartition=0..
15	61	LCK_M_SCH_S	106928	64	objectlock lockPartition=0..
16	63	LCK_M_SCH_S	106843	64	objectlock lockPartition=0..
17	64	LCK_M_SCH_M	106772	55	objectlock lockPartition=0..
18	76	LCK_M_SCH_S	106586	64	objectlock lockPartition=0..
19	80	LCK_M_SCH_S	106410	64	objectlock lockPartition=0..
20	81	LCK_M_SCH_S	106388	64	objectlock lockPartition=0..

Figure 12-2. Output of sys.dm_os_waiting_tasks view

As you may have noticed, the root-blocker session with SPID=55 is not present in the output. That session has the runaway uncommitted transaction and had not been blocked in our example. You can see it in *sleeping* status with `open_transaction_count=1` in the `sys.dm_exec_sessions`

view output. You can use Listing 2-3 to obtain information about the session and client application for further troubleshooting.

You can also see the tasks with THREADPOOL waits in the `sys.dm_os_waiting_tasks` view output. Those tasks belong to connection requests from the clients. SQL Server does not have available workers to handle them, which will eventually trigger connection errors on the client side.

It is also worth noting that those tasks have empty `session_id` values, as the sessions have not been established yet. Thus, they would not appear in the `sys.dm_exec_requests` view output.

There are other cases that may lead to THREADPOOL waits. For example, they may occur in low memory conditions, due to inadequately provisioned servers or memory pressure. The number of workers in the system depends on the amount of memory and other factors; you may have insufficient workers to handle the load. Look at the error log for signs of memory pressure and SQL Server dumps.

It is also possible that server cannot handle the workload because there are too many active sessions or too many concurrent queries with the parallel plans. You need to troubleshoot those conditions and potentially tune the system to reduce the overall load.

Check the *max worker thread* setting. It is sometimes misconfigured and set to a low number. Reset it to the default value and see how that impacts the system.

You can consider increasing the *max worker thread* setting; however, this rarely addresses the problem. For example, massive blocking may quickly reappear, with the new workers also blocked. Don't panic! Just troubleshoot the root cause of the issue – that's always the better option.

Backup-Related Waits

As you can guess by their names, BACKUIO and BACKUPBUFFER waits occur during backup and restore. They indicate insufficient operations throughput when SQL Server cannot write to or read from backup files fast enough. I rarely focus on those waits in the beginning of the performance tuning process; however, there are a couple of things worth analyzing if you see those waits in noticeable amounts.

It is common to see backup, I/O, and network-related waits together when they share the same resources. For example, a slow or overloaded disk subsystem may lead to BACKUIO, PAGEIOLATCH, WRITELOG and other I/O-related waits. Backup-related waits, in such cases, may become another element confirming the problem.

However, when resources are shared, check if the backup operation impacts other components when it is running. For example, is there an increase in the number of occurrences and average wait times of PAGEIOLATCH, WRITELOG or HADR_SYNC_COMMIT waits during that time?

Many monitoring tools collect information about waits that occur in specific time intervals, and you can use this in your analysis. Alternatively, you can build a solution by persisting sys.dm_os_wait_stats data in utility database at regular intervals, perhaps using SQL Server Agent Job. The companion materials for this book include a script that provides a snapshot of the wait statistics for the time interval, which you can use to build your implementation.

In the end, redesign the backup process if you see an impact. Utilize differential backups, use different targets to store files, and adjust the schedule. Your strategy and options will depend on the infrastructure, your RTO and RPO requirements, version and edition of SQL Server, and other factors.

Improving Backup Performance

The native SQL Server backup does not provide a lot of configuration options. Nevertheless, there are some knobs you can turn to tune and improve backup performance. Let's look at them.

Backup Compression

SQL Server allows you to compress backup files. In most cases, it reduces the size of the backup files, at the cost of extra CPU load introduced by the compression.

Usually, this is a good trade-off. Smaller backup files will take less time to transmit over the network, which improves recovery time in the event of disaster. They also reduce the load and storage size on the disk subsystem.

As a general rule, unless the server is CPU-bound, enable backup compression. There are a few considerations related to *Transparent Data Encryption* (TDE), which I will cover shortly.

Striped Backups

You can split the database backup into multiple files and create a *striped backup*. This will parallelize backup and restore operations, allowing SQL Server to use multiple threads to perform them.

This feature is very beneficial for large databases and may significantly reduce backup and restore times. Keep in mind that it is resource intensive and will add load to the infrastructure.

You also need to analyze bottlenecks during the backup process. For example, if you are backing up the database to a network location and don't have sufficient network throughput, striping the backups and parallelizing the process won't help you much. In that case, you might consider striping backups to local DAS storage and copying the files to network location later, when the backup is complete.

BUFFERCOUNT and MAXTRANSFERSIZE Options

SQL Server allows you to control the number of I/O buffers (BUFFERCOUNT option) and the maximum size of transfer block (MAXTRANSFERSIZE option) for your backup operation. Usually, increasing both speeds up the process. You need to carefully tune them in your system, though. After a speed increase, the backup operation will

consume more memory, which may impact other SQL Server components. It may even lead to an *out of memory* condition if you the set parameters incorrectly.

If you want to compress a TDE-enabled database in SQL Server 2016 or 2017, set MAXTRANSFERSIZE larger than 65,536 (64KB). This will switch SQL Server to the new and improved compression algorithm, which works with encryption. Without that, compressing an encrypted database would not save much space. This is not required in SQL Server 2019, which will adjust MAXTRANSFERSIZE automatically while it compresses the TDE-enabled database.

Partial Database Backups

When you deal with large databases, it is very common for a large portion of the data to become static over time: think of scenarios with append-only tables or when the data becomes read-only after some time.

When this is the case, you can partition the data, utilizing partitioned tables and/or views. You can place the static portion of the data into separate filegroups, marking them as read-only. Those filegroups can be backed up just once and then excluded from the regular FULL backups, saving time and significantly reducing the backup's size on disk. (You can read about this implementation in more detail on my [blog](#) or in my book *Pro SQL Server Internals*.)

In the end, the way you tune backup and restore processes should be aligned with your organization's disaster-recovery strategy. Review the requirements and the RTO and RPO metrics and design an implementation that can support them.

HTBUILD and Other HT* Waits

The HTBUILD, HTDELETE, HTMEMO, HTREINIT, and HTREPARTITION waits occur during management of internal hash tables in batch-mode execution. Prior to SQL Server 2019, batch-mode execution

was almost exclusively used with columnstore indexes. In SQL Server 2019 and above, it can also be used with row-based tables.

Having those waits present in small amounts does not necessarily indicate a problem. They may, however, be a sign of poorly maintained columnstore indexes. In particular, they sometimes indicate the presence of large uncompressed delta stores or of many small and unevenly-sized rowgroups. Review the columnstore indexes and rebuild any partitions that have inefficiencies. You can use the `sys.column_store_row_group` view for the analysis.

I have yet to see those waits become an issue with batch-mode execution on row-based tables in SQL Server 2019. I'd guess, however, that inaccurate statistics may introduce them.

Finally, I'd like to note that Microsoft Documentation suggests reducing `MAXDOP` or increasing the cost threshold for parallelism as ways to mitigate these waits. This is not the right approach, in my opinion: it will mask the problem and disable batch-mode execution for some queries, degrading performance.

Preemptive Waits

As you remember from Chapter 2, SQL Server OS uses a cooperative execution model. The workers voluntarily yield when their CPU time quantum expires, allowing other workers to execute.

There are some exceptions, however, when SQL Server needs to call external functions it does not control. Think about OS API calls to authenticate the user against a domain controller or extended stored procedures calls. When it happens, SQL Server switches the worker to preemptive execution mode. SQL Server does not control its scheduling anymore. The worker continues to show the `RUNNING` state; however, it also generates a wait of one of the preemptive types, of which there are more than 200 in SQL Server 2019. Most of them do not represent any

contention and can be ignored. There are a few, however, you need to be aware of. Let's look at them.

PREEMPTIVE_OS_WRITEFILEGATHER Wait Type

PREEMPTIVE_OS_WRITEFILEGATHER waits occur during the zero-initializing process. As you'll recall, SQL Server always zero-initializes log files and may also zero-initialize data files if *instant file initialization* is not enabled.

When you see that wait in the system, check and enable instant file initialization by granting *Perform volume management tasks* (SE_MANAGE_VOLUME_NAME) permission to the SQL Server startup account. Remember that there is a small security risk (discussed in Chapter 1), although it is usually not an issue in most systems.

Also, review the transaction log's auto-growth parameters. Growing the log files in large chunks can prolong zero-initializing time and lead to inefficient VLF structure. As you learned in Chapter 11, it is better to manage transaction-log size manually.

PREEMPTIVE_OS_WRITEFILE Wait Type

PREEMPTIVE_OS_WRITEFILE waits may indicate a bottleneck during synchronous writes to the files. This wait type rarely becomes the issue; however, it may be a sign of a slow or overloaded disk subsystem on the server.

When I see this wait type present in significant amounts, I check whether the server runs multiple SQL Traces using files as the targets to save data. Another occurrence of this wait may be related to database-snapshot writes – either snapshots created by users or to internal snapshots created by the DBCC CHECKDB operation.

Authentication-Related Wait Types

There are several wait types that occur during user authentication calls when SQL Server needs to wait for responses from the Domain Controller. Their names start with `PREEMPTIVE_OS_AUTH*`; and they also include `PREEMPTIVE_OS_LOOKUPACCOUNTSID` and `PREEMPTIVE_OS_ACCEPTSECURITYCONTEXT` waits.

These waits may be infrastructure-related. One typical example involves Cloud-based SQL Servers authenticating against remote on-premises Active Directory controllers. The latency of the calls can prolong the authentication process, leading to a high percentage of authentication-related waits.

Another possible reason is code that runs under a different execution context than the calling session. The `EXECUTE AS OWNER` or `EXECUTE AS USER` modules may require constant authentication calls, which can be expensive under a heavy load.

When you see authentication-related waits in the system, check the health of the AD infrastructure and the latency of authentication calls, then review how often those calls have been performed.

OLEDB Waits

OLEDB is another preemptive wait type that occurs when SQL Server is waiting for data from an OLE DB provider. Most often it happens in the following cases:

- Calls to linked servers

- Execution of some SSIS packages

- Operations during DBCC CHECKDB execution

- Queries against DMVs

Waits from the first two categories usually have relatively high wait times. Such cases indicate long-running remote queries and prolonged SSIS package execution. Consider troubleshooting the performance of those calls on remote servers and/or reviewing the SSIS package logic.

Waits from DBCC CHECKDB operation and from DMV access are usually short, no more than a few milliseconds on average. The cumulative numbers will give you an idea of the overhead of the operations. You can address the DBCC CHECKDB overhead by offloading it to the secondary Availability Group replicas or to a backup validation server.

DMV access waits can give you an idea of the overhead of the monitoring tools that are constantly querying those views. If they produce excessive waits, you might to decide to redesign your monitoring strategy and/or switch to different tools.

Finally, another known issue for OLEDB waits sometimes occurs when a remote non-SQL Server linked server does not terminate the connection properly. SQL Server will keep connection open, generating a never-ending OLEDB wait for the session. Unfortunately, there is no easy way to fix this besides restarting SQL Server.

Wait Types: Wrapping Up

There are hundreds of wait types in SQL Server and it is impossible to cover all of them. Nor is it practical – in all likelihood, you’ll never deal with most of them during system troubleshooting.

This book has so far covered the most common wait types, enough to troubleshoot most issues; however, you may encounter other wait types. Don’t be confused when it happens!

Start by researching the wait type. The Microsoft [documentation](#) is a good place to start. Another useful resource is the SQLSkills [Wait Types Library](#). Obviously, use Google or Bing as well – sometimes you’ll find useful information there.

The conditions when the wait type is generated will point you to the problematic areas. By now, you know enough about SQL Server’s execution model and components to define your troubleshooting strategy and address issues. Just remember to look at the problem holistically and

avoid tunnel vision. All of SQL Server's components work together, and issues that arise can impact many of them.

Summary

ASYNC_NETWORK_IO waits occur when SQL Server has to wait for a client application to consume the data. It may indicate inefficient network throughput or issues with client applications that read and process data row by row.

The presence of THREADPOOL waits indicates the dangerous condition of worker thread starvation. When this happens, SQL Server does not have enough workers to execute client requests and becomes unresponsive. Typical issues that trigger that condition include long blocking chains, insufficient memory (due to incorrect configuration or extreme memory pressure, and heavy concurrent (and often non-optimized) load.

You can use dedicated admin connection to troubleshoot ongoing issues when SQL Server is unresponsive and does not accept regular connections. Make sure to enable remote access to DAC during server provisioning.

BACKUPIO and BACKUPBUFFER waits occur if SQL Server does not have enough throughput to write to or read from backup files. Analyze the infrastructure and tune the backup process when you see them.

Preemptive waits occur when SQL Server needs to call external functions by switching to a preemptive execution model. You can ignore most of those waits; however, pay attention to authentication-related waits, OLEDB waits, and I/O-related preemptive wait types.

In the next chapter we'll switch gears, discussing how to detect inefficiencies in the database schema and indexing.

Troubleshooting Checklist

Review the network topology and client implementation when you see large amount of ASYNC_NETWORK_IO waits.

Investigate THREADPOOL waits if you see them.

Review Disaster Recovery and Backup Strategies and tune backup process if needed.

Analyze the size of delta stores and state of rowgroups in columnstore indexes if you see HTBUILD or other batch-mode execution-related waits.

Troubleshoot preemptive and OLEDB waits if you see a noticeable number of them.

Chapter 13. Database Schema and Index Analysis

A NOTE FOR EARLY RELEASE READERS

With Early Release ebooks, you get books in their earliest form—the author’s raw and unedited content as they write—so you can take advantage of these technologies long before the official release of these titles.

This will be the 13th chapter of the final book. Please note that the GitHub repo will be made active later on.

If you have comments about how we might improve the content and/or examples in this book, or if you notice missing material within this chapter, please reach out to the author at dmitri@aboutsqlserver.com.

Until this point, most of the troubleshooting efforts described in this book have treated users’ databases and applications as *black boxes*. I’ve focused on performance improvements that do not require any changes in the databases and applications beyond indexing and simple T-SQL code changes. This approach provides easier and faster ROI; however, it also limits the results you can accomplish.

Don’t take this the wrong way: in many cases, you can achieve *good enough* results without needing to make significant database and application changes. Nevertheless, it may be beneficial to perform a high-level review of your database schema and index usage and address some of the problems you find.

I’ll start this chapter with an overview of several SQL Server catalog views and show you how to detect a few database-design issues. Next, I’ll demonstrate how to identify inefficient indexing through index usage and operational statistics and provide a handful of scripts for this analysis,

including one that lets you view various index metrics consolidated together at a glance.

Database Schema Analysis

SQL Server provides quite a few **catalog views** that expose information about server- and database-level objects. They are extremely useful when you need to analyze and detect inefficiencies in the database schema.

Figure 14-1 shows several database-object-related catalog views with their dependencies and key attributes. This is just a small subset of the available views, to give you an idea of how much information you can obtain.

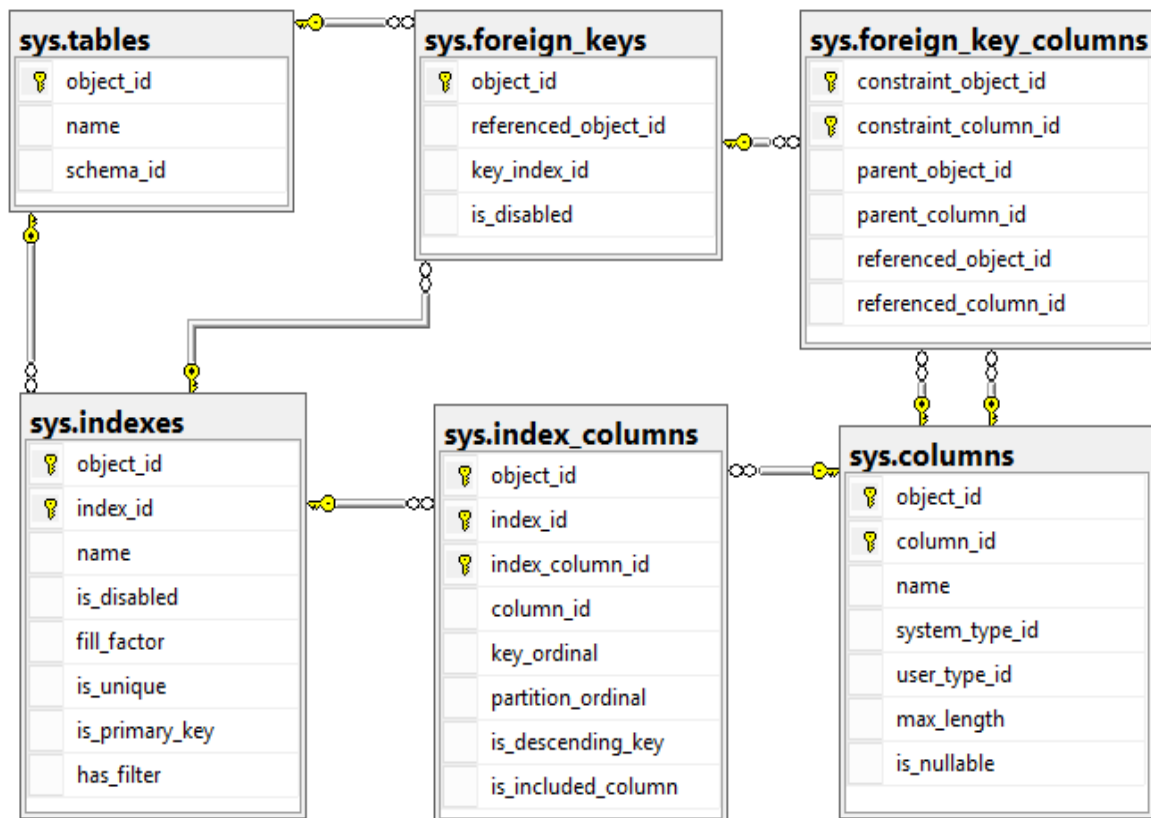


Figure 13-1. Catalog views

Let's look at seven common design issues you can identify with catalog views:

- Heap tables
- Indexes on uniqueidentifier columns

- Wide and nonunique clustered indexes
- Untrusted foreign keys
- Non-indexed foreign keys
- Redundant indexes
- High identity values

Heap Tables

Heap tables are a somewhat controversial subject. They certainly have their uses – for example, they may be a good choice in some ETL processes where fast insert throughput is critical and you cannot use memory-optimized tables. However, as I briefly mentioned in Chapter 5, it is better to avoid heap tables; tables with clustered indexes outperform them in most workloads.

When I see heap tables during my analysis, I usually review them and consider creating clustered indexes on them. If that’s impossible, I look for inefficient heap tables and rebuild them.

The code in Listing 14-1 provides a list of heap tables in the database. You may want to review index usage statistics in those tables to see if there is a natural candidate for a clustered index. (I’ll talk more about this later in the chapter.)

Example 13-1. Listing 14-1. Obtaining heap tables in the database

```

SELECT
    t.object_id
    ,s.name + '.' + t.name AS [table]
    ,p.rows
FROM
    sys.tables t WITH (NOLOCK)
    JOIN sys.schemas s WITH (NOLOCK) ON
        t.schema_id = s.schema_id
    CROSS APPLY
        (
            SELECT SUM(p.rows) AS [rows]
            FROM sys.partitions p WITH (NOLOCK)
            WHERE t.object_id = p.object_id AND p.index_id = 0
        ) p
WHERE
    t.is_memory_optimized = 0 AND

```

```

t.is_ms_shipped = 0 AND
EXISTS
(
    SELECT *
    FROM sys.indexes i WITH (NOLOCK)
    WHERE t.object_id = i.object_id AND i.index_id = 0
)
ORDER BY
    p.rows DESC
OPTION (RECOMPILE, MAXDOP 1);

```

There are two main metrics to monitor with heap tables. The first is the number of *forwarded records*. In contrast to B-Tree indexes, when a heap table page does not have enough space to accommodate the new version of a row, SQL Server does not perform a page split. Instead, it puts the new row on another page and replaces the original row with a small structure called a *forwarding pointer*. Large numbers of forwarding pointers and forwarded records will reduce the performance of I/O operations on the table.

The second metric is *internal fragmentation*, which is the amount of free space available on the data pages. A large degree of internal fragmentation will increase the size of the table on disk and in memory, which decreases performance.

Listing 14-2 shows how to detect inefficient heap tables. The code uses the expensive `sys.dm_db_index_physical_stats` function, which scans the entire table, so it's better to run it in a non-production environment—for example, against a recent database backup restored on a non-production server.

You can rebuild inefficient heap tables with the `ALTER TABLE REBUILD` statement after you detect them.

Example 13-2. Listing 14-2. Detecting inefficient heap tables

```

SELECT TOP 25
    t.object_id
    ,s.name + '.' + t.name AS [table]
    ,SUM(ips.record_count) AS [rows]
    ,SUM(ips.forwarded_record_count)
        AS [forwarding pointers]
    ,SUM(ips.avg_page_space_used_in_percent * ips.page_count) /
        NULLIF(SUM(ips.page_count),0)
        AS [internal fragmentation %]
FROM
    sys.tables t WITH (NOLOCK)

```

```

        JOIN sys.schemas s WITH (NOLOCK) ON
            t.schema_id = s.schema_id
        CROSS APPLY
            sys.dm_db_index_physical_stats
                (DB_ID(),t.object_id,0,NULL,'DETAILED') ips
WHERE
    t.is_memory_optimized = 0 AND
    t.is_ms_shipped = 0 AND
    EXISTS
    (
        SELECT *
        FROM sys.indexes i WITH (NOLOCK)
        WHERE t.object_id = i.object_id AND i.index_id = 0
    )
GROUP BY
    t.object_id, s.name, t.name
ORDER BY
    [forwarding pointers] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

Indexes with the Uniqueidentifier Data Type

B-Tree indexes that use randomly generated values for keys often cause performance problems. Random values introduce very significant index fragmentation, and they are slow during large batch operations.

One of the most common approaches to generating random values is using the uniqueidentifier data type. The code in Listing 14-3 detects indexes that have uniqueidentifier as the leftmost key column in the index. When you detect this, you can also look at index fragmentation with the sys.dm_db_index_physical_stats view.

Example 13-3. Listing 14-3. Getting indexes with uniqueidentifiers in the leftmost key column

```

SELECT
    t.object_id
    ,s.name + '.' + t.name AS [table]
    ,i.name AS [index]
    ,i.is_disabled
    ,p.rows
FROM
    sys.tables t WITH (NOLOCK)
        JOIN sys.schemas s WITH (NOLOCK) ON
            t.schema_id = s.schema_id
        JOIN sys.indexes i WITH (NOLOCK) ON

```

```

        t.object_id = i.object_id
CROSS APPLY
    (
        SELECT SUM(p.rows) AS [rows]
        FROM sys.partitions p WITH (NOLOCK)
        WHERE i.object_id = p.object_id AND i.index_id =
p.index_id
    ) p
WHERE
    t.is_memory_optimized = 0 AND
    i.type in (1,2) AND /* CI and NCI */
    i.is_hypothetical = 0 AND
    EXISTS
    (
        SELECT *
        FROM
            sys.index_columns ic WITH (NOLOCK)
            JOIN sys.columns c WITH (NOLOCK) ON
                ic.object_id = c.object_id AND
                ic.column_id = c.column_id
        WHERE
            ic.object_id = i.object_id AND
            ic.index_id = i.index_id AND
            ic.key_ordinal = 1 AND
            c.system_type_id = 36 /* uniqueidentifier */
    )
ORDER BY
    p.[rows] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

You'll need to analyze how new uniqueidentifier values are generated. Random values generated in the application or with the NEWID() function will lead to problems; however, the NEWSEQUENTIALID() function introduces pseudo-ever-increasing values that behave similarly to identity columns (except less efficient, due to their larger data type size). Keep in mind that the NEWSEQUENTIALID() function may generate values lower than it previously generated after OS restart or failover.

Unfortunately, there is no simple way to address the problems introduced by indexed random keys. I recommend that you analyze the amount of fragmentation, tune the FILLFACTOR index property to minimize page splits, and consider code refactoring when it is an option.

Wide and Non-Unique Clustered Indexes

For good system performance, you need efficient clustered indexes. In addition to supporting critical queries, the ideal clustered index has three characteristics: it should be static, narrow, and unique.

Static

First, an ideal clustered index should be static. The goal is generally to avoid updating clustered keys, a very expensive operation that requires moving data rows to another place in the clustered index B-Tree and then updating row-ids in all nonclustered index rows that reference the data row.

Narrow

Second, an ideal clustered index should be narrow. Because clustered index key columns present as row-id in all nonclustered indexes, wide clustered index keys lead to wide and therefore less efficient nonclustered indexes.

Unique

Finally, an ideal clustered index should be defined as unique. When you don't define this, SQL Server adds another 4-byte internal column called uniquefier (not to be confused with uniqueidentifier!) that enforces the uniqueness of the clustered index keys. This column increases the size of clustered and nonclustered index keys and should be avoided when possible.

Listing 14-4 shows two queries. The first provides the 25 tables with the widest clustered index keys, based on key column data types. The second produces a list of tables with non-unique clustered indexes.

Example 13-4. Listing 14-4. Detecting inefficient clustered indexes

```
SELECT TOP 25
    t.object_id
    ,s.name + '.' + t.name AS [table]
    ,p.rows
    ,ic.[max length]
FROM
    sys.tables t WITH (NOLOCK)
    JOIN sys.schemas s WITH (NOLOCK) ON
```



```

        t.schema_id = s.schema_id
    CROSS APPLY
    (
        SELECT SUM(p.rows) AS [rows]
        FROM sys.partitions p WITH (NOLOCK)
        WHERE i.object_id = p.object_id AND p.index_id = 1
    ) p
    CROSS APPLY
    (
        SELECT SUM(c.max_length) as [max length]
        FROM
            sys.indexes i
            JOIN sys.index_columns ic WITH (NOLOCK) ON
                i.object_id = ic.object_id AND
                i.index_id = ic.index_id AND
                ic.is_included_column = 0
            JOIN sys.columns c WITH (NOLOCK) ON
                ic.object_id = c.object_id AND
                ic.column_id = c.column_id
        WHERE
            i.object_id = t.object_id AND
            i.index_id = 1 AND
            i.type = 1
    ) ic
WHERE
    t.is_memory_optimized = 0
ORDER BY
    ic.[max length] DESC
OPTION (RECOMPILE, MAXDOP 1);
-- Non-unique CI
SELECT
    t.object_id
    ,s.name + '.' + t.name AS [table]
    ,p.rows
FROM
    sys.tables t WITH (NOLOCK)
    JOIN sys.schemas s WITH (NOLOCK) ON
        t.schema_id = s.schema_id
    CROSS APPLY
    (
        SELECT SUM(p.rows) AS [rows]
        FROM sys.partitions p WITH (NOLOCK)
        WHERE t.object_id = p.object_id AND p.index_id = 1
    ) p
WHERE
    t.is_memory_optimized = 0 AND
    EXISTS
    (
        SELECT *

```

```

FROM sys.indexes i WITH (NOLOCK)
WHERE
    t.object_id = i.object_id AND
    i.index_id = 1 AND
    i.is_unique = 0 AND
    i.type = 1 /* CI */
)
ORDER BY
    p.[rows] DESC
OPTION (RECOMPILE, MAXDOP 1);

```

Do not jump to conclusions based strictly on query outputs. The benefits of having the right clustered indexes to support critical queries may easily offset the overhead introduced by wide and non-unique indexes. Look at the queries and indexing holistically before making the decision to refactor database schema.

Nevertheless, it is always better to recreate clustered indexes as unique if data in the index key is unique. The overhead of a uniqueifier column is completely unnecessary.

Untrusted Foreign Keys

Aside from a very few edge cases, it's always good to define foreign key constraints in the database. They improve data quality and reduce errors and bugs. Moreover, they often improve performance—for example, Query Optimizer can remove unnecessary joins between tables with foreign keys present.

There are two ways to create foreign key constraints. By default, they are created as *trusted* with SQL Server validating that existing data in the tables do not violate the constraint. You can also create the constraint as *untrusted* using WITH NOCHECK clause in ALTER TABLE statement. In that case, SQL Server enforces constraint going forward; however, it does not validate existing data.

Unfortunately, untrusted foreign key constraints limit possible optimizations for Query Optimizer, and it is better to validate untrusted foreign keys. Listing 14-5 allows you to detect untrusted constraints in the database.

Example 13-5. Listing 14-5. Selecting untrusted foreign key constraints

```

SELECT
    fk.is_disabled
    ,fk.is_not_trusted
    ,fk.name AS [FK]
    ,ps.name + '.' + pt.name AS [Referencing Table / Detail]
    ,rs.name + '.' + rt.name AS [Referenced Table / Master]
    ,fk.update_referential_action_desc
    ,fk.delete_referential_action_desc
FROM
    sys.foreign_keys fk WITH (NOLOCK)
    JOIN sys.tables pt WITH (NOLOCK) ON
        fk.parent_object_id = pt.object_id
    JOIN sys.schemas ps WITH (NOLOCK) ON
        pt.schema_id = ps.schema_id
    JOIN sys.tables rt WITH (NOLOCK) ON
        fk.referenced_object_id = rt.object_id
    JOIN sys.schemas rs WITH (NOLOCK) ON
        rt.schema_id = rs.schema_id
WHERE
    fk.is_not_trusted = 1 OR fk.is_disabled = 1
OPTION (RECOMPILE, MAXDOP 1);

```

You can validate constraints with **ALTER TABLE CHECK CONSTRAINT** command. The word of caution though. This operation scans the table acquiring schema modification (Sch-M) lock for duration of the scan. Schedule it during the downtime, especially with the large tables.

Non-Indexed Foreign Keys

Foreign key constraints require you to have an index on constraint columns in a referenced (master) table. You are not, however, required to have an index in a referencing (detail) table. This may lead to very serious issues during referential integrity checks.

Consider the situation when you delete the row in referenced table. SQL Server needs to check if this operation violated the constraint and potentially perform cascading action in referencing table. Without the index, it would lead to the referencing table scan, which introduce performance and potential blocking issues.

Listing 14-6 returns the list of foreign key constraints that don't have corresponding indexes in referencing tables defined. With very few

exceptions, you want to create the indexes to support referential integrity operations.

Example 13-6. Listing 14-6. Getting non-indexed foreign key constraints

```
SELECT
    fk.is_disabled
    ,fk.is_not_trusted
    ,fk.name as [FK]
    ,ps.name + '.' + pt.name AS [Referencing Table / Detail]
    ,rs.name + '.' + rt.name AS [Referenced Table / Master]
    ,fk.update_referential_action_desc
    ,fk.delete_referential_action_desc
    ,fk_cols.cols as [fk columns]
FROM
    sys.foreign_keys fk WITH (NOLOCK)
    JOIN sys.tables pt WITH (NOLOCK) ON
        fk.parent_object_id = pt.object_id
    JOIN sys.schemas ps WITH (NOLOCK) ON
        pt.schema_id = ps.schema_id
    JOIN sys.tables rt WITH (NOLOCK) ON
        fk.referenced_object_id = rt.object_id
    JOIN sys.schemas rs WITH (NOLOCK) ON
        rt.schema_id = rs.schema_id
    CROSS APPLY
        (
            SELECT
                (
                    SELECT
                        UPPER(col.name) AS [text()]
                        ,',' AS [text()]
                    FROM
                        sys.foreign_key_columns fkc WITH (NOLOCK)
                        JOIN sys.columns col WITH (NOLOCK) ON
                            fkc.parent_object_id = col.object_id
                )
            AND
                fkc.parent_column_id = col.column_id
            WHERE
                fkc.constraint_object_id = fk.object_id
            ORDER BY
                fkc.constraint_column_id
            FOR XML PATH('')
        ) as cols
    ) fk_cols
WHERE
    NOT EXISTS
        (
            SELECT *
            FROM
```

```

sys.indexes i WITH (NOLOCK)
  CROSS APPLY
  (
    SELECT
      (
        SELECT
          UPPER(col.name) AS [text()]
          ,',' AS [text()]
        FROM
          sys.index_columns ic WITH (NOLOCK)
          JOIN sys.columns col WITH (NOLOCK)
            ON
              ic.object_id = col.object_id
            AND
              ic.column_id = col.column_id
        WHERE
          i.object_id = ic.object_id AND
          i.index_id = ic.index_id AND
          ic.is_included_column = 0
        ORDER BY
          ic.partition_ordinal
        FOR XML PATH('')
      ) AS cols
    ) idx_col
  WHERE
    i.object_id = fk.parent_object_id AND
    CHARINDEX(fk_cols.cols,idx_col.cols) = 1 AND
    i.is_disabled = 0 AND
    i.is_hypothetical = 0 AND
    i.has_filter = 0 AND
    i.type IN (1,2)
  )
OPTION (RECOMPILE, MAXDOP 1);

```

Redundant Indexes

As you remember from Chapter 5, the data in the composite B-Tree indexes is sorted starting from left-most to right-most key columns. SQL Server may use indexes for *Index Seek* operation as long as it has SARGable predicates on left-most columns of the index.

Think about two indexes: IDX1(LastName) and IDX2(LastName, FirstName) as the example. The data in both indexes is sorted based on LastName. Next, IDX2 has the data sorted by FirstName within each LastName, while IDX1 would not have FirstName data sorted. Nevertheless, both indexes would

support index seek on LastName data, which makes IDX1 redundant and unnecessary.

Listing 14-7 shows you the code that you can use to detect potentially redundant indexes with matching left-most column.

Example 13-7. Listing 14-7. Locating potentially redundant indexes

```
SELECT
    s.name + '.' + t.name AS [Table]
    ,i1.index_id AS [I1 ID]
    ,i1.name AS [I1 Name]
    ,dupIdx.index_id AS [I2 ID]
    ,dupIdx.name AS [I2 Name]
    ,LEFT(i1_col.key_col,LEN(i1_col.key_col) - 1) AS [I1 Keys]
    ,LEFT(i1_col.included_col,LEN(i1_col.included_col) - 1) AS [I1
Included Col]
    ,i1.filter_definition AS [I1 Filter]
    ,LEFT(i2_col.key_col,LEN(i2_col.key_col) -1) AS [I2 Keys]
    ,LEFT(i2_col.included_col,LEN(i2_col.included_col) - 1) AS [I2
Included Col]
    ,dupIdx.filter_definition AS [I2 Filter]
FROM
    sys.tables t WITH (NOLOCK)
    JOIN sys.indexes i1 WITH (NOLOCK) ON
        t.object_id = i1.object_id
    JOIN sys.index_columns ic1 WITH (NOLOCK) ON
        ic1.object_id = i1.object_id AND
        ic1.index_id = i1.index_id AND
        ic1.index_column_id = 1
    JOIN sys.columns c WITH (NOLOCK) ON
        c.object_id = ic1.object_id AND
        c.column_id = ic1.column_id
    JOIN sys.schemas s WITH (NOLOCK) ON
        t.schema_id = s.schema_id
CROSS APPLY
    (
        SELECT i2.index_id, i2.name, i2.filter_definition
        FROM
            sys.indexes i2 WITH (NOLOCK)
            JOIN sys.index_columns ic2 WITH (NOLOCK) ON
                ic2.object_id = i2.object_id AND
                ic2.index_id = i2.index_id AND
                ic2.index_column_id = 1
        WHERE
            i2.object_id = i1.object_id AND
            i2.index_id > i1.index_id AND
            ic2.column_id = ic1.column_id AND
            i2.type in (1,2) AND
```

```

i2.is_disabled = 0 AND
i2.is_hypothetical = 0 AND
(
    i1.has_filter = i2.has_filter AND
    ISNULL(i1.filter_definition, '') =
        ISNULL(i2.filter_definition, '')
)
) dupIdx
CROSS APPLY
(
    SELECT
        (
            SELECT
                col.name AS [text()]
                , IIF(icol_meta.is_descending_key = 1, '
DESC', '')
                AS [text()]
                , ',' AS [text()]
            FROM
                sys.index_columns icol_meta WITH (NOLOCK)
                JOIN sys.columns col WITH (NOLOCK) ON
                    icol_meta.object_id = col.object_id
            AND
                    icol_meta.column_id = col.column_id
            WHERE
                icol_meta.object_id = i1.object_id AND
                icol_meta.index_id = i1.index_id AND
                icol_meta.is_included_column = 0
            ORDER BY
                icol_meta.key_ordinal
            FOR XML PATH('')
        ) AS key_col
    , (
        SELECT
            col.name AS [text()]
            , ',' AS [text()]
        FROM
            sys.index_columns icol_meta WITH (NOLOCK)
            JOIN sys.columns col WITH (NOLOCK) ON
                icol_meta.object_id = col.object_id
            AND
                icol_meta.column_id = col.column_id
        WHERE
            icol_meta.object_id = i1.object_id AND
            icol_meta.index_id = i1.index_id AND
            icol_meta.is_included_column = 1
        ORDER BY
            icol_meta.key_ordinal
        FOR XML PATH('')
    )
)

```

```

        ) AS included_col
    ) i1_col
CROSS APPLY
(
    SELECT
        (
            SELECT
                col.name AS [text()]
                ,IIF(icol_meta.is_descending_key = 1, '
DESC', '')
                AS [text()]
                ,',' AS [text()]
            FROM
                sys.index_columns icol_meta WITH (NOLOCK)
                JOIN sys.columns col WITH (NOLOCK) ON
                    icol_meta.object_id = col.object_id
            AND
                icol_meta.column_id = col.column_id
            WHERE
                icol_meta.object_id = t.object_id AND
                icol_meta.index_id = dupIdx.index_id AND
                icol_meta.is_included_column = 0
            ORDER BY
                icol_meta.key_ordinal
            FOR XML PATH('')
        ) AS key_col
    , (
        SELECT
            col.name AS [text()]
            ,',' AS [text()]
        FROM
            sys.index_columns icol_meta WITH (NOLOCK)
            JOIN sys.columns col WITH (NOLOCK) ON
                icol_meta.object_id = col.object_id
        AND
            icol_meta.column_id = col.column_id
        WHERE
            icol_meta.object_id = t.object_id AND
            icol_meta.index_id = dupIdx.index_id AND
            icol_meta.is_included_column = 1
        ORDER BY
            icol_meta.key_ordinal
        FOR XML PATH('')
    ) AS included_col
    ) i2_col
WHERE
    i1.is_disabled = 0 AND
    i1.is_hypothetical = 0 AND
    i1.type in (1,2)

```



```
ORDER BY
    s.name, t.name, i1.index_id
OPTION (RECOMPILE, MAXDOP 1);
```

You can usually drop redundant indexes and perform additional analysis looking for the candidates for index consolidation. For example, you can consolidate IDX3 and IDX4 indexes defined in Listing 14-8 into IDX5. Pay attention to index usage statistics, as it may provide you the information what indexes are rarely used.

Example 13-8. Listing 14-8. Examples of index consolidation

```
CREATE INDEX IDX3 ON T(LastName, FirstName) INCLUDE(Phone);
CREATE INDEX IDX4 ON T(LastName) INCLUDE (SSN);
--IDX3 and IDX4 can be consolidated to IDX5
CREATE INDEX IDX3 ON T(LastName, FirstName) INCLUDE(Phone,SSN);
```

High Identity Values

One of silly, and the same time, very dangerous conditions you may encounter is running out-of-capacity for integer-based columns. Think about the table with INT IDENTITY PRIMARY KEY, which reaches the value of 2,147,483,647, which is the maximum INT data type allows you to store. All further INSERT operations to the table would fail, which may lead to production outage.

To make matter worse, this condition is hard to recover from. While you can change the data type from INT to BIGINT, the ALTER TABLE ALTER COLUMN statement may take hours or even days to complete on large tables. Moreover, it would lock the table with schema modification (Sch-M) lock during the execution. It is sad, but it is not uncommon to see people losing their jobs when it happened.

Listing 14-8 shows the code you can use to detect identity columns that are running out of capacity. You need to regularly run it on your systems and proactively address the risks when you detect the issues.

Example 13-9. Listing 14-8. Identifying identity columns with high utilization

```
DECLARE
    @Types table
    (
        type_id int not null primary key,
```

```

        name sysname not null,
        max_val bigint not null
    )
INSERT INTO @Types(type_id, name, max_val)
VALUES
    (48,'tinyint',255)
    ,(52,'smallint',32767)
    ,(56,'int',2147483647)
    ,(127,'bigint',9223372036854775807);
DECLARE
    @percentThreshold int = 50;
SELECT
    s.name + '.' + t.name AS [table]
    ,c.name AS [column]
    ,tp.name AS [type]
    ,IDENT_CURRENT(t.name) AS [identity]
    ,CONVERT(DECIMAL(6,3),
        100. * IDENT_CURRENT(t.name) / tp.max_val
    ) AS [percent full]
FROM
    sys.tables t WITH (NOLOCK)
    JOIN sys.schemas s WITH (NOLOCK) ON
        t.schema_id = s.schema_id
    JOIN sys.columns c WITH (NOLOCK) ON
        c.object_id = t.object_id
    JOIN @Types tp ON
        tp.type_id = c.system_type_id
WHERE
    c.is_identity = 1 AND
    100. * IDENT_CURRENT(t.name) / tp.max_val > @percentThreshold
ORDER BY
    [percent full] DESC;

```

Obviously, that script would not catch the situation when the values are generated in the code. You may also consider checking sys.sequences view and see current sequence values as another layer of protection.

There are just a few examples of what you can do with catalog views to detect inefficiencies and issues with the database schema. Explore other views and be creative!

Now, let's look at another goldmine of the information – index usage and operational statistics.

Index Analysis

As we all know, indexes help to improve performance of the queries. Unfortunately, they come at cost – they increase amount of data in the database, consume additional memory and add an overhead during data modifications. The large number of inefficient and/or unused indexes may significantly degrade system performance.

Query optimization process is essential part of performance tuning and more often than not leads to creation of new indexes in the database. However, unless you are operating in emergency situation, I'd always advocate to spend time analyzing and removing inefficient indexes *before* starting to create the new ones. It will remove the overhead they introduce and make the process of optimization simpler.

SQL Server comes with two data management views for index usage analysis. The first, `sys.dm_db_index_usage_stats` **view** provides you the information how many queries used the index for seek, scan, update and lookup operations. The second, `sys.dm_db_index_operational_stats` **view** gives you row-level access and operational metrics including information about I/O, locking, latching and a few others.

There is the key difference how the metrics are collected. The `sys.dm_db_index_usage_stats` view provides the information on query level, while `sys.dm_db_index_operational_stats` view works on the row level. For example, if you run the query that inserts 1,000 rows to the table, the first view will have `user_updates` value increased by one. In contrast, the `leaf_insert_count` column in the second view will be incremented by 1,000.

SQL Server does not persist index usage statistics. The metrics in the views will be cleared at time of SQL Server restart or when database becomes offline. Moreover, in some old SQL Server builds (SQL Server 2012 prior SP2-CU12 and SP3-CU3; SQL Server 2014 prior SP2), the metrics will be cleared at time of index rebuild operation. You need to factor it into your analysis making sure that statistics is representative and does not miss important but unfrequently executed queries. Think about critical process in accounting system that runs on schedule once per month as an example.

You also need to look at index usage on readable secondaries in Availability Group setup. It is common to have some indexes to support queries on

secondary nodes and those indexes may appear as unused on primary node.

Let's look at both views in more details and see how we can interpret information from them.

The sys.dm_db_index_usage_stats view

The sys.dm_db_index_usage_stats view is one of the main tools during index usage analysis. It provides you the data on how often index had been used or, to be precise, how many times queries utilized the index and index appeared in the execution plans.

The data is grouped by access methods, and you'd see the separate metrics for seek, scan and lookup operations. Finally, the view also shows how often the index had been updated, which helps you to estimate update overhead it introduces.

Listing 14-9 shows you the code that utilizes the view. I am renaming some of the view columns in the output to make them more compact and will use both, view column names and my aliases interchangeably in that chapter.

Example 13-10. Listing 14-9. Using sys.dm_db_index_usage_stats view

```
SELECT
    t.object_id
    ,i.index_id
    ,s.name + '.' + t.name AS [Table]
    ,i.name AS [Index]
    ,i.type_desc
    ,i.has_filter AS [Filtered]
    ,i.is_unique AS [Unique]
    ,p.rows AS [Rows]
    ,ius.user_seeks AS [Seeks]
    ,ius.user_scans AS [Scans]
    ,ius.user_lookups AS [Lookups]
    ,ius.user_seeks + ius.user_scans + ius.user_lookups AS [Reads]
    ,ius.user_updates AS [Updates]
    ,ius.last_user_seek AS [Last Seek]
    ,ius.last_user_scan AS [Last Scan]
    ,ius.last_user_lookup AS [Last Lookup]
    ,ius.last_user_update AS [Last Update]
FROM
    sys.tables t WITH (NOLOCK)
    JOIN sys.indexes i WITH (NOLOCK) ON
        t.object_id = i.object_id
```

```

JOIN sys.schemas s WITH (NOLOCK) ON
    t.schema_id = s.schema_id
CROSS APPLY
(
    SELECT SUM(p.rows) AS [rows]
    FROM sys.partitions p WITH (NOLOCK)
    WHERE
        i.object_id = p.object_id AND
        i.index_id = p.index_id
) p
LEFT OUTER JOIN sys.dm_db_index_usage_stats ius ON
    ius.database_id = DB_ID() AND
    ius.object_id = i.object_id AND
    ius.index_id = i.index_id
WHERE
    i.is_disabled = 0 AND
    i.is_hypothetical = 0 AND
    t.is_memory_optimized = 0 AND
    t.is_ms_shipped = 0
ORDER BY
    s.name, t.name, i.index_id
OPTION (RECOMPILE, MAXDOP 1);

```

Figure 14-2 shows the output of the code.

	object_id	index_id	Table	Index	type_desc	Filtered	Unique	Rows	Seeks	Scans	Lookups
1	5801971...	1	Table_5801	Index_1	CLUSTERED	0	1	920095	3003390844	442375	35632327
2	5801971...	3	Table_5801	Index_3	NONCLUSTERED	0	1	920095	420014	0	0
3	5801971...	7	Table_5801	Index_7	NONCLUSTERED	0	0	920095	2040257472	7	0
4	5801971...	20	Table_5801	Index_20	NONCLUSTERED	0	0	920095	1805673313	43	0
5	5801971...	22	Table_5801	Index_22	NONCLUSTERED	0	0	920095	1288379097	0	0
6	2107310...	1	Table_2107	Index_1	CLUSTERED	0	1	2701711053	136545557	44714994	104653148
7	2107310...	2	Table_2107	Index_2	NONCLUSTERED	0	0	2701711053	2090315	0	0
8	2107310...	3	Table_2107	Index_3	NONCLUSTERED	0	0	2701711053	10	0	0
	Reads	Updates	Last Seek	Last Scan	Last Lookup	Last Update					
	3039465546	1429545988	2021-08-22 09:55:24.943	2021-08-22 09:55:14.567	2021-08-22 09:55:23.983	2021-08-22 09:55:24.923					
	420014	673402076	2021-08-22 09:54:39.727	NULL	NULL	2021-08-22 09:55:24.923					
	2040257479	1344552091	2021-08-22 09:55:24.963	2021-02-10 17:43:00.593	NULL	2021-08-22 09:55:24.923					
	1805673356	1344552091	2021-08-22 09:55:24.990	2021-08-20 10:56:54.933	NULL	2021-08-22 09:55:24.923					
	1288379097	1344552091	2021-08-22 09:55:25.027	NULL	NULL	2021-08-22 09:55:24.923					
	285913699	44717469	2021-08-22 09:55:21.800	2021-08-22 09:55:20.300	2021-08-22 09:55:23.583	2021-08-22 09:55:20.297					
	2090315	44717461	2021-08-22 09:55:04.730	NULL	NULL	2021-08-22 09:55:20.297					
	10	44717461	2021-07-22 13:15:04.757	NULL	NULL	2021-08-22 09:55:20.297					

Figure 13-2. Sys.dm_db_index_usage_stats view output

Let's look at the columns in the output

database_id , object_id and index_id

The `database_id`, `object_id` and `index_id` columns reference database, table and index respectively. You can use them to filter the output and join them with `sys.databases`, `sys.tables` and `sys.indexes` catalog views.

user_seek , user_scan and user_lookup

The `user_seek`, `user_scan` and `user_lookup` columns provide you the information of how often queries used the index in Index Seek, Index Scan, Key Lookup and RID Lookup operators. An efficient index B-Tree index in OLTP system should primarily use an index seek.

Remember, however, as I discussed in Chapter 5, the seek operation is not always efficient and may scan very large range of rows.

user_update

The `user_update` data shows how many times the index had been modified by insert, update, delete or merge operations. In the nutshell, it allows you to estimate the overhead required for index maintenance during data modifications.

It is very important to remember that `user_update` represents how many times operation occurred rather than how many rows it changed. For example, the single DELETE call would always increment `user_update` value by 1 regardless of how many rows were deleted.

last_user_seek , last_user_scan , last_user_lookup and last_user_update

The `last_user_*` columns provide you the time of the last corresponding operation in the index. That data is useful when SQL Server had not been restarted and usage metrics have been collected for a long time.

system_seek , system_scan , system_lookup , system_update , last_system_seek , last_system_scan , last_system_lookup and last_system_update

The `system_*` columns (not shown in the script and Figure 14-2) provides the statistics on index usage by system processes. It includes statistics

update, index maintenance and a few others. In most cases, you don't need to worry about those metrics.

Let's look at a few common patterns in the view output and conclusions you can draw from them.

Unused indexes

I've performed many SQL Server health checks in the past and I rarely saw the systems that don't have unused and unnecessary indexes present. Those indexes are easy to detect as `sys.dm_db_index_usage_stats` view does not show any read (`user_seek`, `user_scan` and `user_lookup`) activity in them.

Figure 14-3 shows an example of such output (I am removing some columns from Listing 14-9 results to make output more concise). As you can see, the indexes with `index_id` of 43 and 47 do not perform any reads.

	index_id	Table	Index	Filtered	Unique	Rows	Seeks	Scans	Lookups	Reads	Updates
1	1	Table_1888	Index_1	0	1	38128199	6464388153	248501	1597518500	8062155154	2169976596
2	21	Table_1888	Index_21	0	0	38128205	4604192	0	0	4604192	153890418
3	40	Table_1888	Index_40	0	0	38128205	0	11	0	11	186522085
4	41	Table_1888	Index_41	0	0	38128205	2263209	40700	0	2303909	288356234
5	42	Table_1888	Index_42	0	0	38128205	561352904	3	0	561352907	186829746
6	43	Table_1888	Index_43	0	0	38128205	0	0	0	0	186526890
7	44	Table_1888	Index_44	0	0	38128205	1681656033	0	0	1681656033	389900596
8	47	Table_1888	Index_47	0	0	38128205	0	0	0	0	1248665902
9	53	Table_1888	Index_53	0	0	38128205	11859	3	0	11862	428364779
10	64	Table_1888	Index_64	1	1	298499	23565855	6193576	0	29759431	1922246283

Figure 13-3. Unused indexes

In most cases, those indexes are the easiest to deal with. You can disable and/or drop them with very little risk involved. Just remember and analyze a few things before you do that.

First, as I've mentioned, make sure that those indexes are not used across all Availability Group nodes. It is very common to see unused indexes on primary node being heavily utilized for reporting on readable secondaries.

Second, make sure that usage statistics are representative and that you don't overlook infrequent processes the index supports. That condition is tricky. You need to analyze the benefits and downsides of keeping the index. In some cases, you may need to consider removing the index and its overhead and

running infrequent processes inefficiently. You can also consider disabling the index, re-enabling it on schedule when needed.

Finally, pay extra attention to unique indexes. They may support unique constraints and removing them may lead to data quality issues.

Nevertheless, in many cases unused indexes are low-hanging fruit, and removing them can immediately benefit the system.

Indexes with High Maintenance Cost

I'd consider indexes that have significantly higher number of updates over the reads being the more complex case of unused indexes. You may benefit from removing them; however, you need to analyze their usage and estimate negative impact when infrequently executed queries run without them.

I don't have formal threshold after which indexes should be analyzed. Personally, I use several criteria's. First, I'd check indexes with low number of reads, especially those with low number of seeks (see index with index_id=40 on Figure 14-3). Next, I'd look at the indexes where number of updates is significantly, order or orders of magnitude, higher than the number of reads (index_id=53 on Figure 14-3).

You can get some of the queries that utilize the index using the code in Listing 14-10. This code analyzes plan cache data and may miss queries that don't have execution plan cached. You can adjust the code to run against sys.query_store_plan and other Query Store catalog **views** if you have Query Store enabled.

As the word of caution, the code is also slow – you may consider dumping content of plan cache into the tables in utility database and do the analysis in non-production environment.

Example 13-11. Listing 14-10. Detecting queries that use the index

```
DECLARE
    @IndexName SYSNAME = QUOTENAME('<INDEX NAME>');
;WITH XMLNAMESPACES
(DEFAULT 'http://schemas.microsoft.com/sqlserver/2004/07/showplan')
,CachedData
AS
(
    SELECT DISTINCT
```

```

        obj.value('@Database','SYSNAME') AS [Database]
        ,obj.value('@Schema','SYSNAME') + '.' +
obj.value('@Table','SYSNAME')
        AS [Table]
        ,obj.value('@Index','SYSNAME') AS [Index]
        ,obj.value('@IndexKind','VARCHAR(64)') AS [Type]
        ,stmt.value('@StatementText','NVARCHAR(MAX)') AS [Statement]
        ,CONVERT(NVARCHAR(MAX),qp.query_plan) AS query_plan
        ,cp.plan_handle
FROM
    sys.dm_exec_cached_plans cp WITH (NOLOCK)
        CROSS APPLY sys.dm_exec_query_plan(plan_handle) qp
        CROSS APPLY query_plan.nodes

('/ShowPlanXML/BatchSequence/Batch/Statements/StmtSimple') batch(stmt)
        CROSS APPLY stmt.nodes

('..//IndexScan/Object[@Index=sql:variable("@IndexName")]') idx(obj)
)
SELECT
    cd.[Database]
    ,cd.[Table]
    ,cd.[Index]
    ,cd.[Type]
    ,cd.[Statement]
    ,CONVERT(XML,cd.query_plan) AS query_plan
    ,qs.execution_count
    ,(qs.total_logical_reads + qs.total_logical_writes) /
qs.execution_count
    AS [Avg IO]
    ,qs.total_logical_reads
    ,qs.total_logical_writes
    ,qs.total_worker_time
    ,qs.total_worker_time / qs.execution_count / 1000 AS [Avg Worker
Time (ms)]
    ,qs.total_rows
    ,qs.creation_time
    ,qs.last_execution_time
FROM
    CachedData cd
        OUTER APPLY
            (
                SELECT
                    SUM(qs.execution_count) AS execution_count
                    ,SUM(qs.total_logical_reads) AS total_logical_reads
                    ,SUM(qs.total_logical_writes) AS total_logical_writes
                    ,SUM(qs.total_worker_time) AS total_worker_time
                    ,SUM(qs.total_rows) AS total_rows
                    ,MIN(qs.creation_time) AS creation_time

```

```

        ,MAX(qs.last_execution_time) AS last_execution_time
    FROM sys.dm_exec_query_stats qs WITH (NOLOCK)
    WHERE qs.plan_handle = cd.plan_handle
) qs
OPTION (RECOMPILE, MAXDOP 1);

```

The further actions would depend on number of queries that utilize the index. In some cases, you can refactor them switching to another index. In other cases, you may drop the index or leave it in the database.

Inefficient Reads

Another target for analysis in OLTP environments is the indexes with large number of scans, especially when that number is significantly higher than number of seeks. Figure 14-4 shows such an example.

	index_id	Table	Index	Filtered	Unique	Rows	Seeks	Scans	Lookups	Reads	Updates
1	1	Table_2887	Index_1	0	1	76707068	2616312067	157348	4961245333	7577714748	1398991251
2	3	Table_2887	Index_3	1	0	573442	231677400	0	0	231677400	1387934116
3	4	Table_2887	Index_4	0	0	76707068	21	59382702	0	59382723	854126441
4	5	Table_2887	Index_5	0	0	76707067	3393509	0	0	3393509	824102200
5	6	Table_2887	Index_6	0	0	76707068	1129131537	1	0	1129131538	868627246
6	8	Table_2887	Index_8	0	0	76707068	175291797	3	0	175291800	817555464
7	9	Table_2887	Index_9	0	0	76707068	0	3124770	0	3124770	857147269
8	35	Table_2887	Index_35	0	0	76707068	568834597	0	0	568834597	719062166

Figure 13-4. Indexes with inefficient reads

You can use similar approach as I discussed in previous section and find queries that use those indexes. I usually start with indexes that have very little

or no seeks at all. Queries that scan them usually need to be optimized, which may allow to remove inefficient usage of the indexes.

Pay attention if index is filtered. It is common to scan filtered indexes and it may be completely normal when index is small.

Inefficient Clustered Indexes and Heaps

In the tables with clustered indexes (index_id=1), large user_lookup value indicates excessive *Key Lookup* operations. This is usually the sign of either inefficient clustered indexes or presence of nonclustered indexes that do not cover frequently-executed queries.

Figure 14-5 shows one of such examples. The clustered index shows large number of lookups with no seeks at all. This pattern is very common when a table has synthetic CLUSTERED PRIMARY KEY defined on IDENTITY column and queries use different columns and indexes to access the data.

index_id	Table	Index	Filtered	Unique	Rows	Seeks	Scans	Lookups	Reads	Updates
1	Table_2067	Index_0	0	0	8362494	0	1	364167805	364167806	286971064
2	Table_2067	Index_2	0	1	8362494	364167805	478	0	364168283	195902738
3	Table_2067	Index_2	0	1	8362494	27121	315	0	27436	195902738
4	Table_2067	Index_2	0	1	8362494	165623	0	0	165623	195902738

Figure 13-5. Inefficient clustered index

Some of those cases are very easy to analyze and address. For example, in the case shown in Figure 14-5, the index with index_id=2 is clearly the one that does not cover frequently-executed queries and lead to *Key Lookup* operations. You can consider making it clustered index, especially if it is not very wide. Alternatively, you can include additional columns to the index to cover the queries.

Unfortunately, not all cases are easy and straightforward. For example, if you look at the statistics shown in Figure 14-4, you'd see that the clustered index is

used for both, *Index Seek* and *Key Lookup* operations. Neither does it clearly show nonclustered index that may be responsible for lookups. In that case, you'd likely want to keep existing clustered index, analyze queries against nonclustered indexes and potentially make those indexes covering.

Finally, the large number of lookups in heap tables (`index_id=0`) indicates excessive *RID Lookup* operations. The most common option to address it would be changing one of frequently used nonclustered indexes to become clustered.

Wrapping Up

The `sys.dm_db_index_usage_stats` view is the great tool to detect and remove inefficient indexes. However, and I am repeating it over and over again, make sure that you are working with the data collected from all servers in Availability Group and over representative period of time.

It is also very beneficial to include the data provided by `sys.dm_db_operational_index_stats` view into analysis. Let's look at that view in more details.

The `sys.dm_db_index_operational_stats` view

The `sys.dm_db_index_operational_stats` view provides low-level statistics on index access methods, locking, latching, I/O and a few other areas. This data is incredibly useful to troubleshoot index performance and identify locking and latching bottlenecks.

Listing 14-11 shows the simple code that utilizes that view. I show it for demo purposes here – you may benefit from combining the data together with `sys.dm_db_index_usage_stats` and other views, as I'll discuss later in the chapter.

Example 13-12. Listing 14-11. Using `sys.dm_db_index_operational_stats` view

```
SELECT
    t.object_id
    ,i.index_id
    ,s.name + '.' + t.name AS [Table]
    ,i.name AS [Index]
    ,i.type_desc
    ,i.has_filter AS [Filtered]
```

```

,i.is_unique AS [Unique]
,p.rows AS [Rows]
,ous.*
FROM
sys.tables t WITH (NOLOCK)
    JOIN sys.indexes i WITH (NOLOCK) ON
        t.object_id = i.object_id
    JOIN sys.schemas s WITH (NOLOCK) ON
        t.schema_id = s.schema_id
CROSS APPLY
(
    SELECT SUM(p.rows) AS [rows]
    FROM sys.partitions p WITH (NOLOCK)
    WHERE
        i.object_id = p.object_id AND
        i.index_id = p.index_id
) p
OUTER APPLY sys.dm_db_index_operational_stats
(DB_ID(),i.object_id,i.index_id,NULL) ous
WHERE
    i.is_disabled = 0 AND
    i.is_hypothetical = 0 AND
    t.is_memory_optimized = 0 AND
    t.is_ms_shipped = 0
ORDER BY
    s.name, t.name, i.index_id
OPTION (RECOMPILE, MAXDOP 1);

```

Figure 14-6 shows the output of the code. You can look at the name of columns in the output to get some sense on the information the view provides.

object_id	index_id	Table	Index	type_desc	Filtered	Unique	Rows	database_id	object_id	index_id	partition_number	
1	1946490013	1	Table_0013	Index_1	CLUSTERED	0	1	0	11	1946490013	1	1
2	274100017	1	Table_0017	Index_1	CLUSTERED	0	1	134	11	274100017	1	1
3	1397580017	1	Table_0017	Index_1	CLUSTERED	0	1	0	NULL	NULL	NULL	NULL
hobt_id		leaf_insert_count	leaf_delete_count	leaf_update_count	leaf_ghost_count	nonleaf_insert_count	nonleaf_delete_count					
72057597802774528		0	0	0	0	0	0					
72057597681336320		61	0	234	0	3	0					
NULL		NULL	NULL	NULL	NULL	NULL	NULL					
nonleaf_update_count		leaf_allocation_count	nonleaf_allocation_count	leaf_page_merge_count	nonleaf_page_merge_count	range_scan_count						
0		0	0	0	0	4214002						
0		3	1	0	0	20302445						
NULL		NULL	NULL	NULL	NULL	NULL	NULL					
singleton_lookup_count		forwarded_fetch_count	lob_fetch_in_pages	lob_fetch_in_bytes	lob_orphan_create_count	lob_orphan_insert_count						
0		0	0	0	0	0						
3531733818		0	0	0	0	0						
NULL		NULL	NULL	NULL	NULL	NULL	NULL					
row_overflow_fetch_in_pages		row_overflow_fetch_in_bytes	column_value_push_off_row_count	column_value_pull_in_row_count	row_lock_count							
0		0	0	0	0							
0		0	0	0	594							
NULL		NULL	NULL	NULL	NULL	NULL						
row_lock_wait_count		row_lock_wait_in_ms	page_lock_count	page_lock_wait_count	page_lock_wait_in_ms	index_lock_promotion_attempt_count						
0		0	0	0	0	0						
0		0	377	0	0	4045						
NULL		NULL	NULL	NULL	NULL	NULL						
index_lock_promotion_count		page_latch_wait_count	page_latch_wait_in_ms	page_io_latch_wait_count	page_io_latch_wait_in_ms							
0		0	0	1	1							
0		1	27	1	1							
NULL		NULL	NULL	NULL	NULL	NULL						
tree_page_latch_wait_count		tree_page_latch_wait_in_ms	tree_page_io_latch_wait_count	tree_page_io_latch_wait_in_ms								
0		0	0	0								
0		0	0	0								
NULL		NULL	NULL	NULL	NULL							
page_compression_attempt_count		page_compression_success_count										
0		0										
0		0										
NULL		NULL										

Figure 13-6. Sys.dm_db_index_operational_stats view output

I am not going to cover all columns in the output and encourage you to read the [documentation](#). Nevertheless, let's look at a few categories of the columns and discuss how to use them.

Data modification statistics

The `leaf_insert_count`, `leaf_update_count`, `leaf_delete_count`, `leaf_ghost_count`, and `leaf_allocation_count` columns provide you information about data modifications in the index. As the opposite to `sys.dm_db_index_usage_stats` view, which counts number of operations, `sys.dm_db_index_operational_stats` data gives you the number of affected rows. You can correlate the information from both views to get better understanding of index maintenance overhead.

In addition, `sys.dm_db_index_operational_stats` view provides the same metrics for intermittent and root levels of B-Tree indexes though the set of `nonleaf_*` columns. You can use them for troubleshooting of `ACCESS_METHODS_HOBT_VIRTUAL_ROOT` latch caused by root-page splits and contention.

Data access statistics

The `singleton_lookup` and `range_scan_count` columns provide you access method data. The first column counts seeks and lookup operations that return single rows. The second column counts *Index Seek* operations that perform the range scan of multiple rows along with index scans. You can use that data to estimate efficiency of index seeks based on `singleton_lookup` value. However, it is impossible to estimate how large are the range scans based on `range_scan_count` value alone.

The `forwarded_fetch_count` column gives you the number of forwarding-pointer reads in heap tables. The heap tables with high value in that column are inefficient and need to be rebuilt. You can use it together with `forwarded_record_count` data in `sys.dm_db_index_physical_stats` view as I discussed earlier in the chapter.

The `lob_fetch_in_pages`, `lob_fetch_in_bytes`, `row_overflow_fetch_in_pages` and `row_overflow_fetch_in_bytes` columns give you statistics on off-row

column access. The queries against those tables may select unnecessary columns, perhaps using `SELECT *` antipattern.

Locking information

The `row_lock_count`, `row_lock_wait_count`, `row_lock_wait_in_ms`, `page_lock_count`, `page_lock_wait_count` and `page_lock_wait_in_ms` columns give you row- and page-level locking statistics. You can use this information to detect the indexes and tables that suffers the most from the locking issues.

The `index_lock_promotion_attempt_count` and `index_lock_promotion_count` columns contain the number of lock escalation attempts and successful lock escalations on the index. The latter column is useful when you see high percent of intent lock waits (`LCK_M_I*`), which are often triggered by lock escalations.

Usually, I do not use `sys.dm_db_index_operational_stats` view as the main tool for locking and blocking troubleshooting. Nevertheless, it is useful tool to cross-check data collected from other venues I discussed in Chapter 8.

Latching information

The `page_latch_wait_count`, `page_latch_wait_in_ms`, `tree_page_latch_wait_count` and `tree_page_latch_wait_in_ms` columns give you page-latch statistics for the index. The former two columns show the data for leaf-level index pages; the latter two – for intermediate and root pages.

Those metrics are extremely useful when you see high percent of `PAGELATCH` waits generated in users' databases. Indexes with highest page latch waits are likely the ones that lead to hotspots and bottleneck.

I/O information

Similarly, `page_io_latch_wait_count`, `page_io_latch_wait_in_ms`, `tree_page_io_latch_wait_count` and `tree_page_io_latch_wait_in_ms` columns point you to the indexes that experienced the most `PAGEIOLATCH` waits.

As with locking troubleshooting, you should not use `sys.dm_db_index_operational_stats` data as the main venue for I/O performance troubleshooting. Nevertheless, it is useful to analyze indexes with

the highest PAGEIOLATCH waits. You can drop them if they are not in use or compress them and get immediate relief by decreasing I/O they introduce.

Wrapping Up

The `sys.dm_db_index_operational_stats` view provides you *a lot* of useful information for the troubleshooting. The ability to look at low-level statistics on per-index basis gives you another perspective to validate your assumptions and confirm potential root-causes of the issues. Do not use it as the main source to drive conclusions though.

Think about the situation when you troubleshoot disk subsystem performance. You may detect and drop nonclustered indexes with highest `page_io_latch_wait_in_ms` values; however, non-optimized queries will just start scanning other indexes. Moreover, it could lead to even higher I/O throughput if SQL Server start to scan larger clustered indexes.

The proper approach during the troubleshooting would be confirming non-optimized queries as the source of the I/O bottleneck and then detect and optimize most I/O intensive queries. You may not even need `sys.dm_db_index_operational_stats` data during query optimization; however, in some cases, it may help to pinpoint I/O intensive operators in the execution plan.

There is one exception though – page latching. That `sys.dm_db_index_operational_stats` view can be the key tool to identify hotspots in the indexes. Obviously, you will need to analyze the schema and workload, and confirm the assumption. Nevertheless, the view usually points you in the right direction and speed up the troubleshooting.

As I said multiple times, look at the system holistically and utilize all tools available for you. The index usage and operational statistics views are the great tools to have in the toolbox.

Holistic View : `sp_Index_Stats`

The catalog and data management views provide you a lot of information for analysis. Unfortunately, there are quite a few of them and you need to look at

many different places to get the full picture. It becomes inconvenient and slows down the process.

In my workflow, I addressed this by writing the stored procedure- `sp_Index_Stats`- which combines the information from various views and return it in the single output. I will share this code with you – you can download it as part of this book’s companion materials or from my [blog](#).

NOTE

The code will be available prior the book is published. It is not available at time of Early Release of the chapter.

The stored procedure provides large amount of information including:

- Index definition and metadata
- Size of the index on-disk and in buffer pool
- Index usage statistics
- Index operational statistics
- Statistics information

You can see the sample output in Figure 14-7.

Server	Database	OnTime	GuidThere	Max Compress...	ObjectId	Index...	Table	Index	Rows	Index U...	fill_factor	
1	PROD-S...	PROD-DB	2021-08-2...	No	PAGE	21073...	1	Table_0717	Index_1	2713565269	1	97
2	PROD-S...	PROD-DB	2021-08-2...	No	PAGE	21073...	2	Table_0717	Index_2	2713565269	0	90
3	PROD-S...	PROD-DB	2021-08-2...	No	PAGE	21073...	3	Table_0717	Index_3	2713565269	0	90
is_disabled	Filtered	Lock Escalation	TotalPages	UsedPages	DataPages	TotalSpaceMB	UsedSpaceMB	DataSpaceMB	BufferPoolSizeMb	Statistics Date		
0	0	DISABLE	41347098	41341938	41159291	323024	322983	321556	136093.492	2021-08-28 06:21		
0	0	DISABLE	13685014	13662719	13556170	106914	106739	105907	182.664	2021-08-28 08:4		
0	0	DISABLE	16467630	16457421	16304083	128653	128573	127375	24.555	2021-08-28 06:2		
Seeks	Scans	Lookups	Reads	Updates	Last Seek	Last Scan	Last Lookup	Last Update	range_scan_count	singleton_lookup_count		
198203	10836	260686	469725	0	2021-08-28	2021-08-28	2021-08-28	NULL	328907	108117208696		
3667	3	0	3670	0	2021-08-28	2021-08-28	NULL	NULL	29075939	0		
13413	0	0	13413	0	2021-08-19	NULL	NULL	NULL	633286	0		
row_lock_wait_count	row_lock_wait_in_ms	page_lock_wait_count	page_lock_wait_in_ms	page_latch_wait_count								
0	0	0	0	5687								
0	0	0	0	2								
0	0	0	0	2								
page_latch_wait_in_ms	page_io_latch_wait_count	page_io_latch_wait_in_ms										
8491	248821989	962889860										
1	322655	580721										
0	583154	1648198										

**Early Release - Will be
changed**

Figure 13-7. sp_Index_Stats output

The stored procedure collects the information from one or multiple databases and allows you to persist output in the table for further analysis. I usually consolidate collected data from all Availability Group nodes in the single table before I start researching it.

Obviously, you are not obligated to use that stored procedure. However, it would speed up the process and simplify your work.

Summary

SQL Server catalog and data management views are the goldmines in identifying inefficiencies in database design and indexing strategy. You need to perform database schema and index usage analysis as part of system health check.

You can use key system catalogs, such as `sys.tables`, `sys.indexes`, `sys.index_columns`, `sys.foreign_keys` and others to identify database design inefficiencies. The possible issues include inefficient heap tables and clustered indexes, non-indexed foreign key constraints, redundant indexes, and many others.

Always check if `IDENTITY` columns are reaching maximum data type capacity. Reaching data type capacity may lead to prolonged downtime.

You can use the `sys.dm_db_index_usage_stats` view to analyze and detect inefficient indexes in the database. It includes unused and suboptimal indexes, indexes with high maintenance cost, and inefficient clustered indexes, among other issues.

Another view, `sys.dm_db_index_operational_stats`, provides low-level statistics on access, locking, latching and I/O. It may help you to detect hotspots during page latching and find indexes that contribute to other performance issues.

Both views are cleared when SQL Server restarts and when the database goes offline. They do not represent index usage on secondary Availability Group replicas. Make sure you are dealing with representative information during your analysis.

Finally, the `sp_Index_Stats` stored procedure provides holistic information about indexes, including their metadata, size on disk and in buffer pool, and usage and operational metrics. You can download it from the companion materials.

Now, it's time for us to discuss virtualizing, as well as troubleshooting SQL Server in a virtualized environment.

Troubleshooting Checklist

- Detect potential database-schema inefficiencies with catalog views
- Review the current values and remaining capacity of IDENTITY columns
- Analyze index usage with the `sys.dm_db_index_usage_stats` and `sys.dm_db_index_operational_stats` views and address possible issues

About the Author

Dmitri Korotkevitch is a Microsoft Data Platform MVP and the Director of Database Services at Chewy.com. He specializes in the design, development, and performance tuning of complex OLTP systems that handle thousands of transactions per second, and has years of experience working with Microsoft SQL Server as an Application and Database Developer, Database Administrator, and Database Architect. Dmitri also provides SQL Server consulting services and training to clients around the world.