# Query Processing Bottlenecks

*by Sophia*

# 1. Introduction

A bottleneck is a degradation in performance that occurs when demand for services or resources exceeds what is available. There are a variety of reasons why bottlenecks occur. These include hardware limitations, inefficient query processing, inadequate indexing, locking and concurrency issues, disk I/O constraints, and network latency. A bottleneck can have a negative impact on the overall responsiveness, speed, and scalability of a database, resulting in slow queries, increased response times, and poor system performance. The ability of the database system to handle growing data volumes, user loads, and transactional demands without compromising its efficiency requires identifying and resolving database bottlenecks.

Addressing database bottlenecks can involve performance tuning, query optimization, hardware upgrades, and proper database resource configuration. Database administrators and developers must regularly monitor and profile their systems to identify potential bottlenecks and implement suitable solutions. There are several ways to alleviate bottlenecks and enhance database performance, including indexing, caching, partitioning, and load balancing. Management and mitigation of database bottlenecks can help organizations maintain an efficient database system that is reliable and responsive to the demands of applications and users.

# 2. Types of Bottlenecks

There are many types of bottlenecks that affect the performance of a database. These include memory, input/output, processing, and indexing bottlenecks. We'll first learn about memory bottlenecks.

## 2a. Memory Bottlenecks

**Memory bottlenecks** occur when the database system cannot access or manage memory resources effectively. It is common to encounter a memory bottleneck when insufficient memory is allocated to the database. This results in frequent disk writes and reads, which causes query processing to be significantly slowed down. Data paging in and out of disk storage is a means of simulating additional memory using hard disk space, at the expense of performance, when the current working set (data actively used by queries and operations) exceeds the available memory. The inefficient use of memory can also cause memory bottlenecks, causing certain database components to consume excessive amounts of memory and starving other components. A database's overall responsiveness may be negatively impacted by **memory leaks**, suboptimal caching strategies, or excessive temporary table usage.

Management and allocation of memory resources are essential to solving database memory bottlenecks. A server's memory settings need to be configured to optimize query performance, and the database must have enough memory to accommodate the working set and anticipated growth. Caching and buffering techniques can be implemented to minimize disk I/O and improve memory utilization by actively monitoring memory usage and identifying memory-intensive queries or processes. Database systems must be tuned to efficiently utilize memory and provide a high-performance and responsive environment for critical data operation by profiling and tuning memory.

It is usually a good indication of an issue when we see longer query execution times, excessive input/output on the system, or even out-of-memory messages in the logs. A **query optimizer** is like a smart assistant inside a DBMS that determines the most efficient way to execute a query to retrieve or modify data, reducing the amount of memory and processing required. Query optimizers are commonly used to improve query efficiency and avoid memory bottlenecks. It is also possible to resolve a bottleneck by adding more physical memory. If the system does not have enough memory (RAM), all processes running on it would compete for this resource. This can result in performance degradation and could cause the system to become unstable. Adding more physical memory can help to fix the bottleneck and improve system performance.

> 📄 **TERMS TO KNOW**
>
> **Memory Bottleneck**
> Occurs when the database system either does not have enough physical memory allocated for its use or is unable to access or manage memory resources effectively.
>
> **Memory Leak**
> A situation where a program allocates memory but then fails to release it when the memory is no longer needed.
>
> **Query Optimizer**
> A software that determines the most efficient way to execute a query.

## 2b. Input/Output Bottlenecks

**Input/output (I/O) bottlenecks** occur due to a bandwidth limitation in reading from or writing to the storage device. I/O bottlenecks are generally associated with physical hardware limitations or latency issues. **Latency** is a delay between when a request is issued and its fulfillment. For example, when an application requests a file from a hard disk, the time it must wait for that file to be read and delivered is latency.

> ✏️ KEY CONCEPT
>
> I/O bottlenecks can be caused by slow disk drives or outdated storage technology. You can significantly improve I/O performance by upgrading to faster and more modern storage solutions, such as solid-state drives (SSDs).

When multiple database operations vie for the same disk resources, the throughput rate available may be insufficient to handle all requests. One way to increase throughput is to implement a RAID. RAID stands for Redundant Array of Independent Disks. It involves combining multiple physical disks into a single logical unit. Some types of RAID are designed to increase I/O performance by spreading out the data across multiple disks. For example, if a file is written across four different disks, with each disk storing just a portion of the file, the amount of data written per disk is 1/4 the original amount and takes only 1/4 as long to write or read.

Fragmentation can slow down data retrieval on a traditional magnetic hard disk drive. **Fragmentation** refers to a condition in which a file is not stored contiguously on the disk, so the read/write head has to move around, picking up all the pieces of the file when it is requested. You can reduce fragmentation by running a disk defragmenter or optimization utility. Such a utility rearranges the content on the disk so that all the parts of each file are stored in adjacent sectors. This reduces the time it takes for the disk's read/write head to move to the right spots to read all the parts of the file. Defragmentation is not necessary on a solid-state drive (SSD).

Query pagination can also help. It retrieves and displays a large set of data in smaller, manageable chunks or pages instead of fetching and displaying it all at once. You can employ pagination by adding a page size parameter to the query.

Multiple queries accessing the same data simultaneously can strain disk I/O with a high volume of concurrent transactions. You can manage concurrency and I/O contention using transaction isolation levels and database locking. **Isolation levels** define the degree to which one transaction must be isolated from the effects of other concurrent transactions. Lower isolation levels offer better performance but may allow for anomalies between concurrent transactions.

**Database locking** is a mechanism for controlling access to a database in a multiuser environment. When a record is locked by one user, no other user may edit it.

Database **disk caching** is a way to improve performance by reducing the need to access physical storage for frequently accessed data. It involves storing a portion of the database in the computer's memory to speed up read operations. This works because accessing data from memory is significantly faster than reading it from the disk. Using **in-memory databases** (that is, databases that rely primarily on memory for data storage and retrieval) can increase performance, as can increasing the size of a disk's cache.

> 🚩 HINT
>
> Distributed databases are prone to I/O bottlenecks due to network latency between nodes. You can mitigate this issue by reducing network latency or using distributed caching.

Database administrators can monitor and profile the system, looking for ways to improve I/O speed. Some of the methods they can use include upgrading storage hardware, optimizing queries, creating indexes, and adjusting disk caching settings. Administrators should also regularly tune and optimize the database configuration, run maintenance utilities such as a disk defragmenter (if appropriate for the storage type), and monitor disk usage to ensure that I/O is not a bottleneck to optimal performance.

📄 TERMS TO KNOW

**Input/Output (IO) Bottleneck**
Occurs when there is a bandwidth limitation in reading from or writing to the storage device.

**Latency**
A delay between when a request is issued and its fulfillment.

**Fragmentation**
A condition in which a file is not stored contiguously on the disk, so the read/write head has to move around, picking up the pieces of the file when it is requested.

**Isolation Level**
The degree to which one transaction must be isolated from the effects of other concurrent transactions.

**Database Locking**
A mechanism for controlling access to a database in a multiuser environment.

**Disk Caching**
A way to improve performance by reducing the need to access physical storage for frequently accessed data.

**In-Memory Databases**
Databases that rely primarily on memory for data storage and retrieval.

## 2c. Processing Bottlenecks

Many factors can affect query execution and overall system performance, leading to processing bottlenecks—in other words, bottlenecks due to the CPU not being able to keep up with the demands on it. It is possible that some of these problems have more than one cause and solution.

🖊 KEY CONCEPT

One way to resolve a processing bottleneck is to make queries more efficient so they do not require as much processing time. For example, a query can be optimized and rewritten to use more efficient join strategies.

Indexing problems can also lead to database issues. A lack of indexes for frequently queried columns can lead to full table scans, which can cause processing bottlenecks. Creating indexes on key columns can significantly speed up data retrieval.

Suboptimal database design can also cause processing bottlenecks. Evaluating database design and considering denormalization, data partitioning, and proper table structures can improve data access and processing.

A high level of concurrency can lead to high levels of contention and performance bottlenecks as a result of locking. You can mitigate concurrency issues by implementing appropriate transaction isolation levels, optimizing locking strategies, and employing techniques such as row-level locking.

On the hardware side of things, insufficient CPU, memory, and storage resources can limit a database's processing capacity. You can improve processing capabilities by upgrading hardware components or sharding databases to distribute workloads across multiple servers. **Sharding** is a database architecture strategy that breaks up a large database into smaller, more manageable parts called shards. Each shard is an independent database that stores a subset of the overall data.

The absence of query caching can result in significant processing overhead for frequently executed queries. Implementing query caching mechanisms can improve response times for repetitive queries.

A database that supports multiple concurrent connections can experience a processing bottleneck when many users access it at once. You can optimize resource utilization via connection pooling and connection management techniques. **Connection pooling** is a technique for managing and reusing database connections. Establishing and closing database connections can be resource intensive; connection pooling aims to improve performance by reusing existing connections rather than creating new ones for each interaction with the database.

⑦ DID YOU KNOW

Process bottlenecks can also be caused by inadequate or irregular database maintenance. Database performance can be enhanced by regularly rebuilding indexes, updating statistics, and cleaning up data.

📄 TERMS TO KNOW

**Sharding**
A database architecture strategy that breaks up a large database into smaller, more manageable parts called shards.

**Connection Pooling**
A technique for managing and reusing database connections.

## 2d. Indexing Bottlenecks

As you learned earlier in the course, indexing involves making a list of all the values in a certain column to help speed up searches. If there is poor or no indexing for a frequently accessed column, performance can suffer, especially in a table with many records. Indexing columns that are frequently used to filter query records can help.

An **indexing bottleneck** occurs when the process of maintaining and updating indexes becomes a performance bottleneck. There are many reasons why indexing bottlenecks can occur, impacting the performance of database queries. Examples include:

| There are too many indexes on | This can cause data modification operations to run slowly. This is because the database maintains each index during INSERT, UPDATE, and DELETE operations. Identify any |

| a table. | redundant or unused indexes and remove them to improve the speed of data modification. |
|---|---|
| Query execution may be slowed. | This is caused by missing indexes on frequently queried columns since the database may need to perform full table scans. Create indexes on frequently used columns in WHERE, JOIN, and ORDER BY clauses to speed up data retrieval. |
| Query execution performance degrades over time. | As data changes over time, indexes can become fragmented, resulting in decreased performance. Maintain index efficiency by regularly defragmenting or rebuilding them. |
| A wrong index type is chosen. | This is based on query patterns and data characteristics (e.g., clustered, non-clustered, covering). Assess how queries are accessed and how data is distributed to determine the right index type for each query. |
| Composite indexes have an inefficient design. | Composite indexes can be useful, but if they are not designed properly, they may cause bottlenecks. Make sure the most selective columns are at the beginning of composite indexes to match query predicates. |
| There are issues with index cardinality. | The cardinality of an index is a measure of how many unique values there are in the indexed columns compared to the total number of rows in the table. It is often expressed as a ratio or percentage. When the cardinality of an index is high, the database may ignore it, and when it is low, the search results may need to be narrowed down more effectively. Data distribution and query patterns should be aligned with index cardinality. Out-of-date statistics can lead to suboptimal query execution plans caused by stale statistics. Update statistics regularly so that the optimizer can choose the most efficient query plan. |
| There are issues with index overhead. | Maintaining indexes can have a significant impact on database performance. They consume storage space. It is important to balance query optimization with storage overhead carefully. |

Index maintenance, query performance monitoring, and query analysis are key to solving indexing bottlenecks. You should periodically review index usage, identify redundant or missing indexes, and optimize the indexes you want to keep. In order to ensure optimal query performance, you must strike a balance between the number and types of indexes. A well-performing database system can be maintained by regularly evaluating and tuning indexes.

📄 **TERM TO KNOW**

**Indexing Bottleneck**

Occurs when the process of maintaining and updating indexes becomes a performance bottleneck.

📋 **SUMMARY**

In this lesson, you learned about different **types of bottlenecks**, which refers to a degradation in performance that occurs when demand for services or resources exceeds what is available. You learned that database bottlenecks affect databases' smooth and efficient operation. These include **memory bottlenecks, input/output bottlenecks, processing bottlenecks, and indexing bottlenecks**.

Various factors may contribute to these bottlenecks, such as hardware limitations, inefficient query processing, inadequate indexing, locking or concurrency issues, and network latency. Bottlenecks negatively impact the database's responsiveness, speed, and scalability, resulting in slow query execution, increased response times, and degraded system performance. A database system's ability to handle growing data volumes, user loads, and transaction demands depends on identifying and resolving database bottlenecks.

Tuning database performance, optimizing queries, upgrading hardware, and properly configuring database resources are key to addressing database bottlenecks. To identify potential bottlenecks and implement appropriate solutions, database administrators and developers must monitor and profile the database system regularly. By employing indexing, caching, partitioning, and load balancing, bottlenecks can be alleviated, and database performance can be improved. Proactively managing and mitigating database bottlenecks can ensure that organizations maintain a responsive and reliable database system that meets the demands of their applications and users.

Source: THIS TUTORIAL WAS AUTHORED BY DR. VINCENT TRAN, PHD (2020) AND Faithe Wempen (2024) FOR SOPHIA LEARNING. PLEASE SEE OUR **TERMS OF USE**.

---

📄 **TERMS TO KNOW**

**Connection Pooling**
A technique for managing and reusing database connections.

**Database Locking**
A mechanism for controlling access to a database in a multiuser environment.

**Disk Caching**
A way to improve performance by reducing the need to access physical storage for frequently accessed data.

**Fragmentation**
A condition in which a file is not stored contiguously on the disk, so the read/write head has to move around, picking up the pieces of the file when it is requested.

**In-Memory Databases**
Databases that rely primarily on memory for data storage and retrieval.

**Indexing Bottleneck**
Occurs when the process of maintaining and updating indexes becomes a performance bottleneck.

**Input/Output (I/O) Bottleneck**
Occurs when there is a bandwidth limitation in reading from or writing to the storage device.

**Isolation Level**

The degree to which one transaction must be isolated from the effects of other concurrent transactions.

**Latency**

A delay between when a request is issued and its fulfillment.

**Memory Bottleneck**

Occurs when the database system either does not have enough physical memory allocated for its use or is unable to access or manage memory resources effectively.

**Memory Leak**

A situation where a program allocates memory but then fails to release it when the memory is no longer needed.

**Query Optimizer**

A software that determines the most efficient way to execute a query.

**Sharding**

A database architecture strategy that breaks up a large database into smaller, more manageable parts called shards.