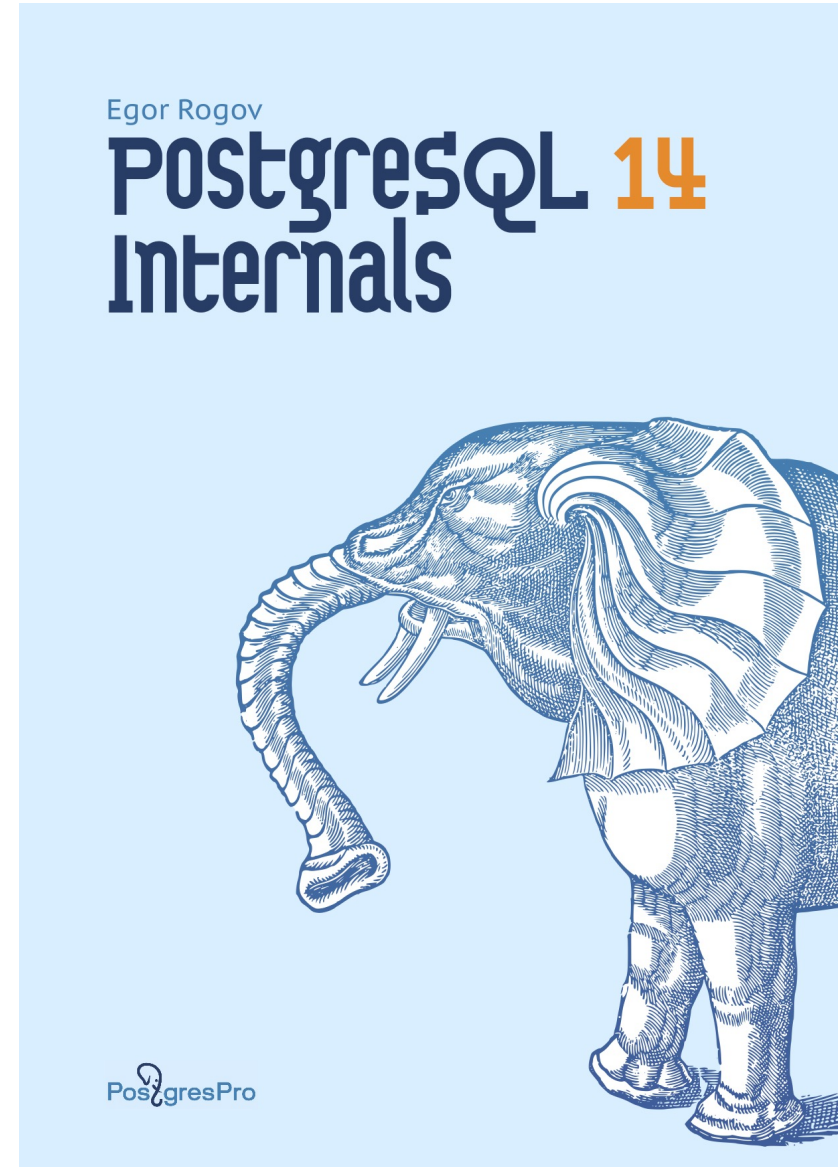
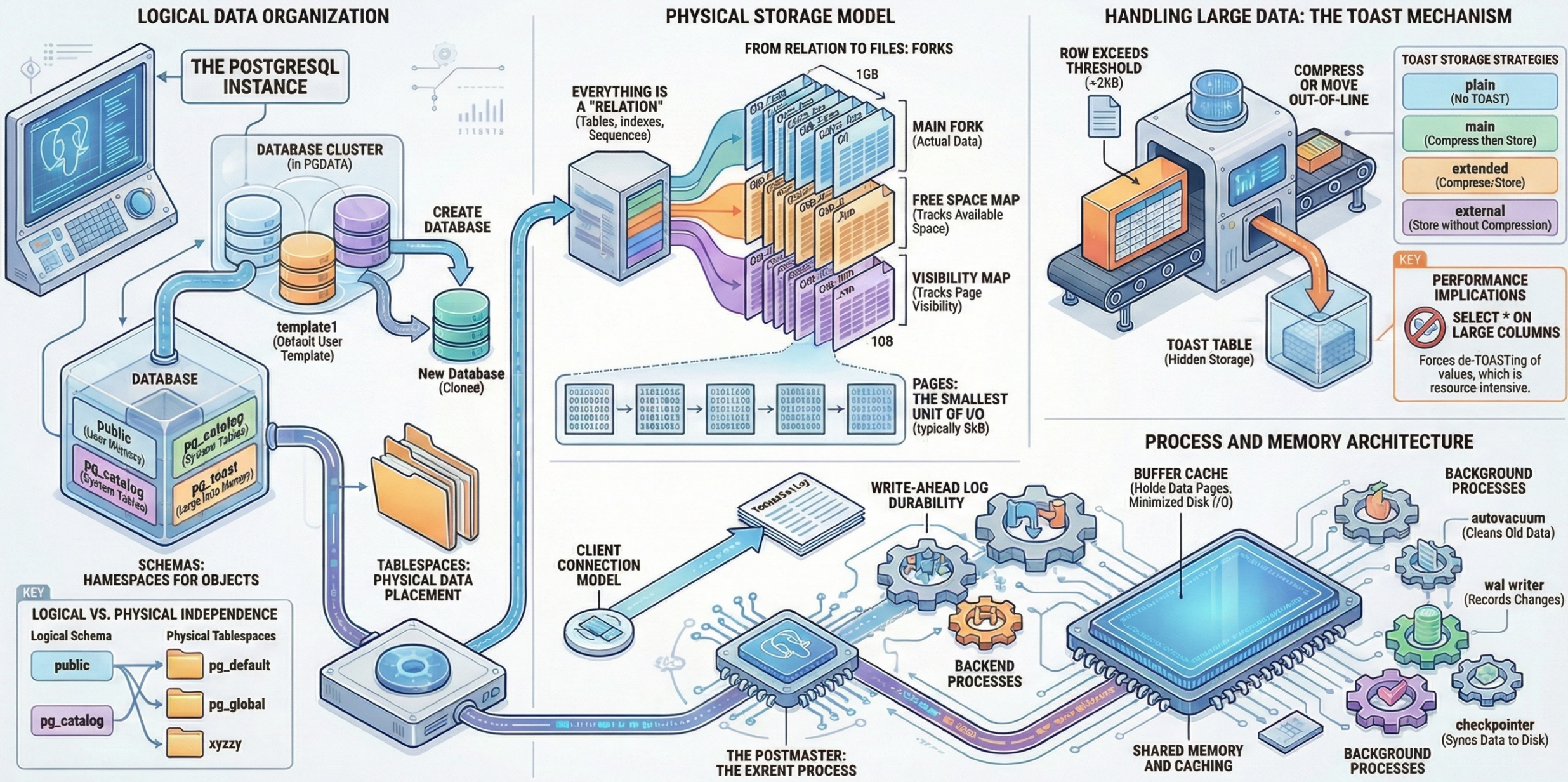


InfoGraph




Under the Hood: The Architecture of PostgreSQL




A Visual Guide to Database Transaction Isolation


1. The Core Challenge: Consistency vs. Concurrency




DATA CONSISTENCY
Assurance of data accuracy and correctness, correctness, defined by application rules (e.g., total balance remains constant during a transfer).



TRANSACTIONS
Bundles multiple operations into a single all-or-nothing unit.

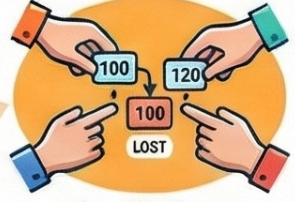


THE CONCURRENCY PROBLEM
Multiple transactions running simultaneously can interfere, violating consistency rules.

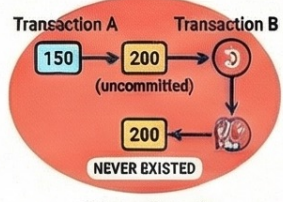


THE NEED FOR ISOLATION
To guarantee consistency, transactions must be isolated, ensuring concurrent execution equals sequential execution.

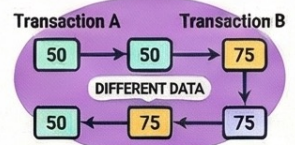
2. Concurrency Anomalies: What Can Go Wrong?



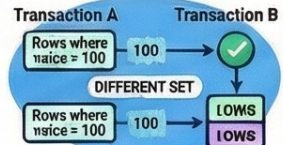
Lost Update
Two transactions read and update the same value, overwriting and "losing" the first's change.



Dirty Read
Transaction A reads uncommitted data from Transaction B. If B rolls back, A reads invalid data.

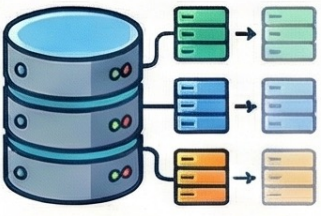


Non-Repeatable Read
A transaction reads a row twice but gets different data each time because another transaction modified and committed that row in between the two reads.




Phantom Read
A transaction runs the same query twice but gets a different set of rows because another transaction inserted or deleted rows that match the query's criteria.


4. The PostgreSQL Way: Snapshot Isolation




PostgreSQL Uses Snapshot Isolation
Employ Multi-Version Concurrency Control (MVCC). Each transaction operates on a consistent "snapshot" from its start, reducing locking.



Stricter Than The Standard
In PostgreSQL, Dirty Reads are never permitted. Read Uncommitted behaves exactly like Read Committed.



Beware of New Anomalies
Prevents standard anomalies but can suffer from Read Skew (in Read Committed) and Write Skew (in Repeatable Read).



Serialization Failures Instead of Anomalies
At higher levels (Repeatable Read, Serializable), PostgreSQL proactively aborts conflicting transactions.

5. Practical Guidance: Which Level Should You Use?

Read Committed (The Default)

PROS:
Good performance, no forced transaction retries.

CONS:
Developer is responsible for preventing anomalies like read skew and lost updates using explicit locks (SELECT FOR UPDATE) or single atomic SQL statements.

Repeatable Read

PROS:
Protects against more anomalies, great for complex, multi-statement read-only reports.

CONS:
Application must have logic to retry transactions that fail due to serialization errors. Still vulnerable to write skew.

Serializable (The Safest)

PROS:
Guarantees full data consistency, simplifying application logic.

CONS:
Highest performance overhead, requires robust transaction retry logic, and cannot be used with read replicas.

3. The SQL Standard: Four Levels of Isolation

Read Uncommitted
Most lenient. Choosing a level is a trade-off between consistency guarantees and performance.

Read Committed

Repeatable Read

Serializable
Stricter levels often require more locking, reducing concurrency.

Isolation Level	Lost Update	Dirty Read	Non-Repeatable Read	Phantom Read	Other Anomalies
Read Uncommitted	✓	✓	✓	✓	Yes
Read Committed	✓	✗	✓	✓	Yes
Repeatable Read	✓	✗	✗	✓	Yes
Serializable	✗	✗	✗	✗	—

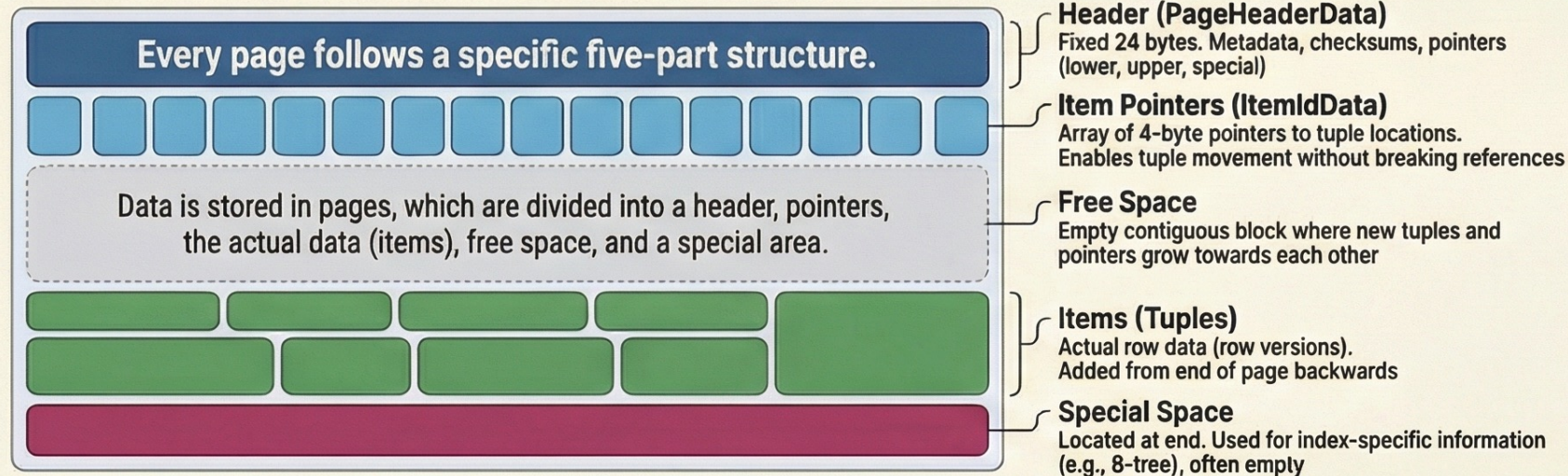
Note: "—" indicates the anomaly is prevented. "Yes" indicates it is allowed.

PostgreSQL Level	Lost Updates	Dirty Reads	Non-Repeatable Reads	Phantom Reads	Other Anomalies
Read Committed	✓	✗	✓	✓	Yes
Repeatable Read	✗	✗	✓	✗	Yes
Serializable	✗	✗	✗	✗	—

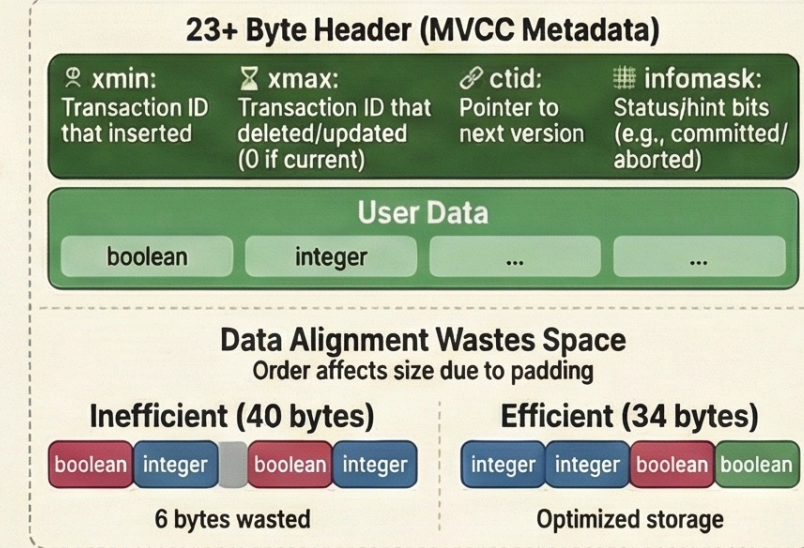
Note: "Lost Updates" refers to application logic reading, calculating, and writing back, not prevented by default.

Inside a PostgreSQL Page: A Visual Guide to Data Storage and Transactions

Anatomy of a PostgreSQL Page (8KB Block)



The Structure of a Row Version (Tuple)



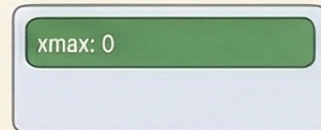
The Lifecycle of a Tuple: MVCC in Action

1. INSERT



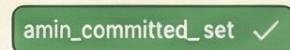
New tuple added. xmax set to 0, indicating current, valid version

2. COMMIT



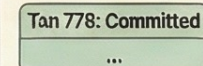
Status recorded in external Commit Log. Tuple on page is not immediately changed

Hint Bits Are Set Lazily



Next reader checks CLOG, then sets xmin_committed hint bit to avoid future lookups

Commit Log (CLOG)



3. UPDATE



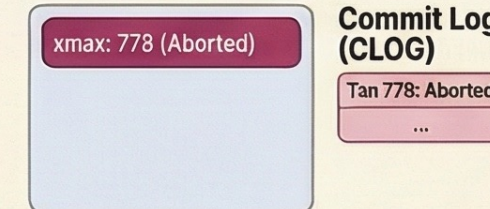
Two-step process: Old tuple xmax set to current transaction ID. New tuple created with xmin as current ID

4. DELETE



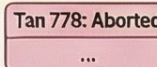
Tuple is not physically removed. xmax set to ID of deleting transaction. Becomes invisible after commit

5. ROLLBACK (Abort)

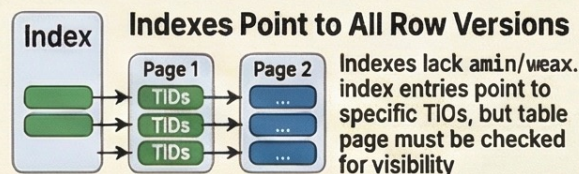


Transaction marked aborted in CLOG. Changes remain on page; xmax_aborted hint bit will be set by next observer, making change invisible

Commit Log (CLOG)



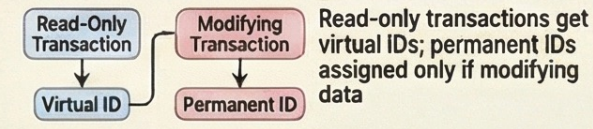
Supporting Structures & Concepts



TOAST for Large Data



Virtual Transactions Optimize Reads



Subtransactions for Savepoints



A Deep Dive into PostgreSQL Snapshots: Understanding Transaction Visibility

What is a Snapshot?



A Consistent View of the Database

A snapshot represents the state of committed data at a specific moment, providing a stable and consistent view for a transaction.



Ensure Transaction Isolation

Each transaction uses its own snapshot, meaning different transactions can see different states of the data simultaneously without interfering with each other.



Not a Physical Copy

Instead of copying data, a snapshot is a logical construct defined by a set of numbers that helps apply visibility rules to existing raw versions.

Snapshots and Isolation Levels

Read Committed: One Snapshot Per Statement

In this isolation level, a new snapshot is created at the beginning of each individual SQL statement within a transaction.



Snapshot



Snapshot

Statement 1

Statement 2

Visualize the Difference



Snapshot

Statement 1

Statement 2

Repeatable Read & Serializable: One Snapshot Per Transaction

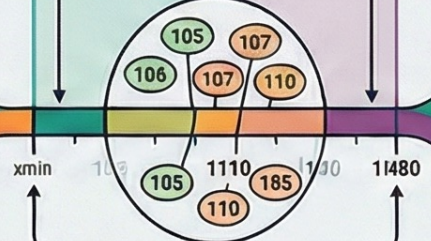
In these isolation levels, a single snapshot is created at the start of the first statement and is reused for the entire duration of the transaction.

The Anatomy of a Snapshot

xmin (Lower Bound)
The ID of the oldest active transaction. All transactions with a lower ID are guaranteed to be either committed or rolled back.

xmax (Upper Bound)
The first unused transaction ID. All transactions with an ID greater than or equal to xmax started after the snapshot was taken and are invisible.

Snapshot Range



Transaction ID (xid) Axis

xmin (Lower Bound)
The ID of the oldest active transaction. All transactions with a lower ID are guaranteed to be either committed or rolled back.

xmax (Upper Bound)
The first unused transaction ID. All transactions with an ID greater than or equal to xmax started after the snapshot was taken and are invisible.

xip_list (In-Progress List)

A list of active transaction IDs that were running between xmin and xmax when the snapshot was created. Their changes are not visible.

The Rules of Visibility

Visible: Committed Before Snapshot
(xid < xmin)



Visible Data

Changes from transactions that finished long before the snapshot was taken are always visible.

Conditionally Visible: Committed Within Snapshot Range
(xmin ≤ xid < xmax)



Visible Data

Changes are visible ONLY if the transaction ID is NOT in the xip_list (meaning it committed before the snapshot was taken).

Is NOT in xip_list?



Invisible: Started After Snapshot
(xid ≥ xmax)

Invisible Data

Changes from transactions that began after the snapshot was created are never visible.

Exception: Your Own Changes Are Always Visible

A transaction can always see its own uncommitted modifications.

The Database Horizon and VACUUM



Outdated Tuples (Dead Data)

The Database Horizon

The object xmin among all currently active snapshots in the database. It defines the object point to history any active transaction needs to see.

Horizon Enables Cleanup

Old raw versions (outdated tuples) no longer visible to any active transaction (i.e., those whose xmin is older than the horizon) can be safely removed by the VACUUM process.



Long-Running Transactions Cause Bloat

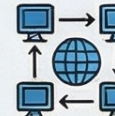
A long-running transaction with an old snapshot holds the database between back, preventing VACUUM from cleaning up dead tuples and causing tables and indexes to grow in size (bloat).

Special Cases



System Catalogs Use Up-to-Date Snapshots

To ensure correctness, queries on system tables (e.g., checking constraints) use the most recent data, ignoring the transaction's primary snapshot.

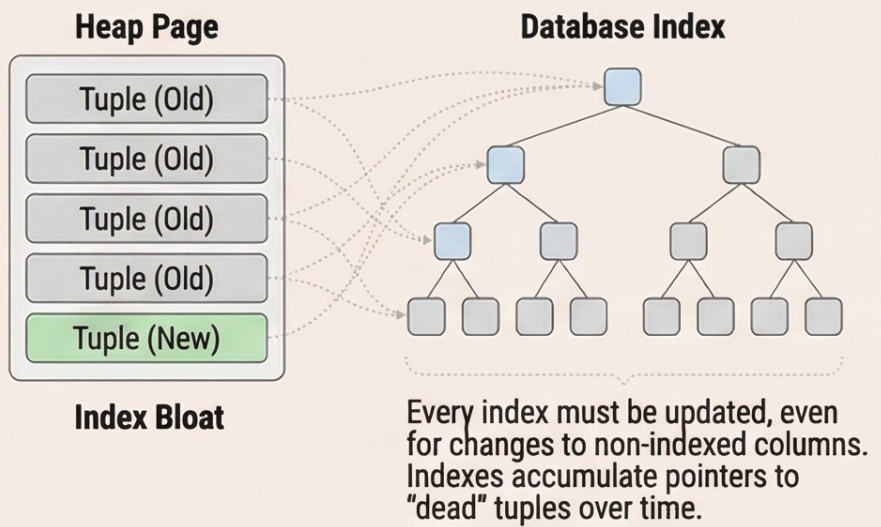


Exporting Snapshots for Consistency

The pg_export_snapshot function allows multiple concurrent transactions to export and share the exact same snapshot, ensuring they all see an identical state of the database (critical for tools like parallel pg_dump).

The Problem: Standard Updates Are Expensive

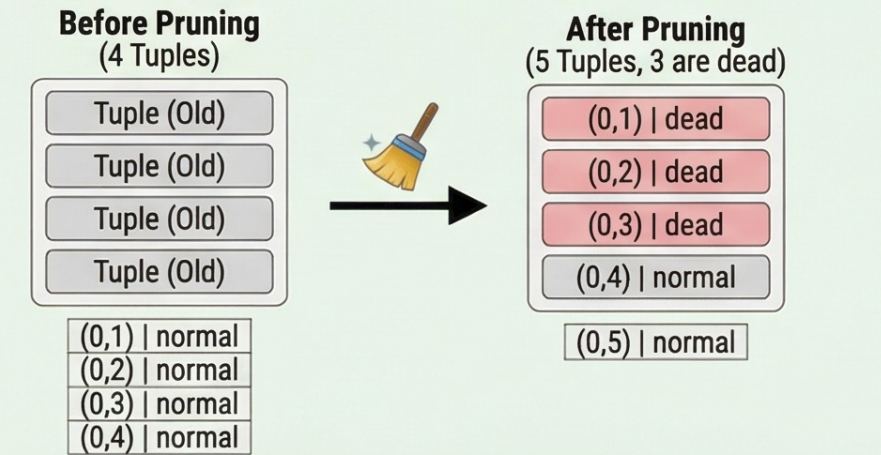
Every UPDATE creates a new tuple version, leaving the old one behind. Every index accumulates pointers to "dead" but tuples over time.



First Line of Defense: Page Pruning

What is Page Pruning?

A fast, automatic cleanup that removes dead tuples (no longer visible to any transaction) from a single heap page.

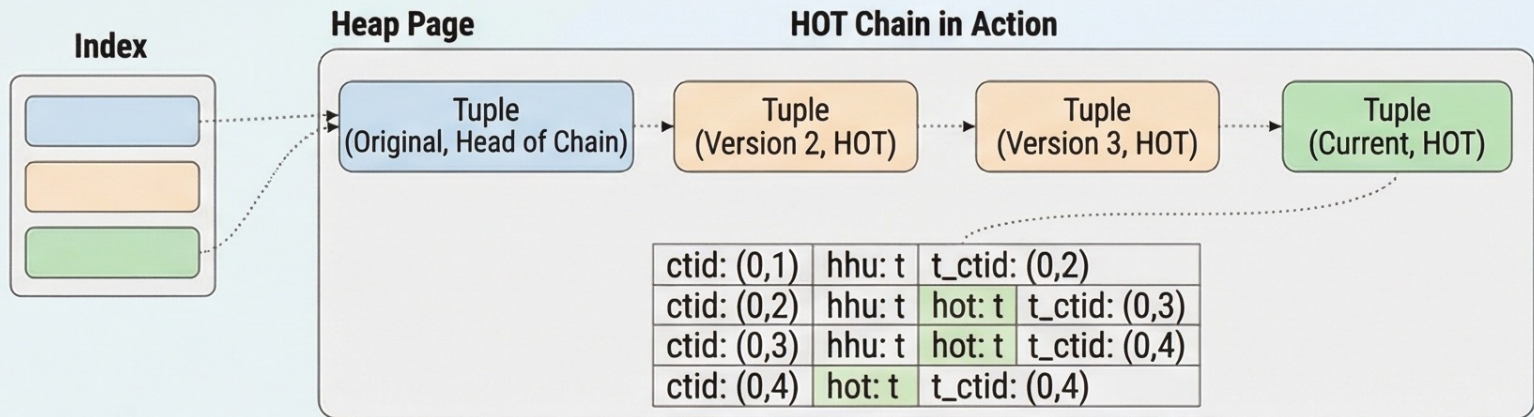


Anatomy of a PostgreSQL Update: How HOT & Page Pruning Optimize Performance

The Ultimate Optimization: Heap-Only Tuple (HOT) Updates

What are HOT Updates?

An optimization that avoids creating new index entries when an UPDATE does not modify any indexed columns.



How a HOT chain works: An index scan finds the first tuple. If it's marked "Heap Hot Updated," the scan follows the pointer chain on the heap page to find the current, visible version.

Managing HOT Chains: Pruning & Splitting

Pruning a HOT Chain

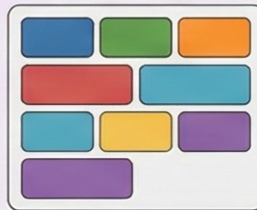
Page pruning is smarter with HOT chains. Since intermediate tuples are not referenced by indexes, they can be completely removed. The first tuple in the chain is kept but changed to a "redirect" state, pointing to the current head of the shortened chain.

fillfactor = 75%

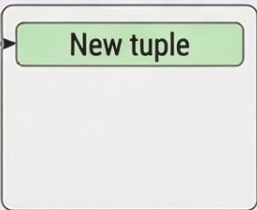
Splitting a HOT Chain

If a page becomes full and a new tuple version cannot be stored, the chain is "split." A new tuple is created on a different page, and a new, second entry must be added to the index. This breaks the HOT optimization for that chain.

Heap Page (Full)



Heap Page (New)

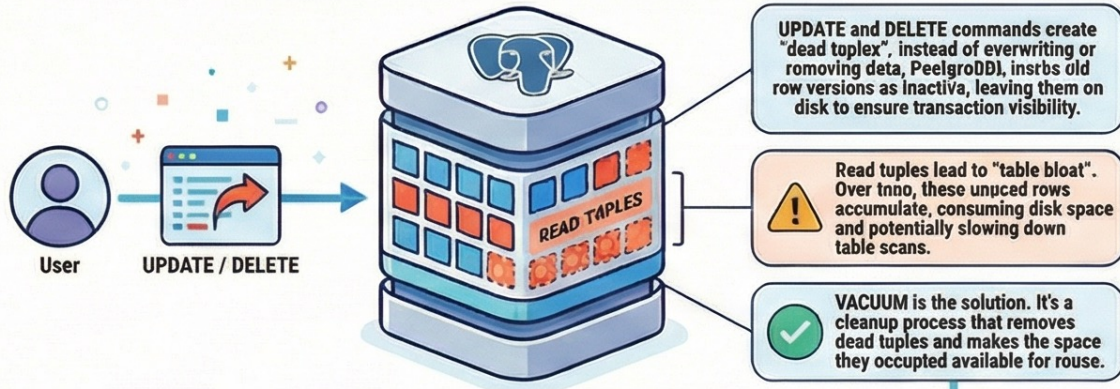


Use fillfactor to reserve space for HOT updates.

For tables with frequent updates to non-indexed columns, lowering fillfactor (e.g., to 75%) leaves empty space on each page, allowing HOT chains to grow without splitting. The trade-off is a larger overall table size.

PostgreSQL's Cleanup Crew: A Guide to VACUUM & AUTOVACUUM

The Problem: "Dead Tuples" in PostgreSQL



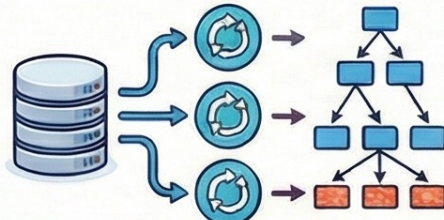
The 4 Stages of a Manual VACUUM

1. Heap (Table) Scan



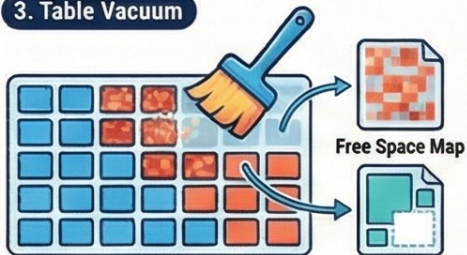
Scans the table to find dead tuples, using the Visibility Map to skip pages known to have no dead tuples. Their IDs are collected in memory (maintenance_work_mem).

2. Index Vacuum



Scans all of the table's indexes to find and remove entries that point to the collected dead tuples. This can run in parallel for multiple indexes.

3. Table Vacuum



Re-scans the table to remove the actual dead tuple versions, now that they are no longer referenced by any index. Updates the Free Space Map.

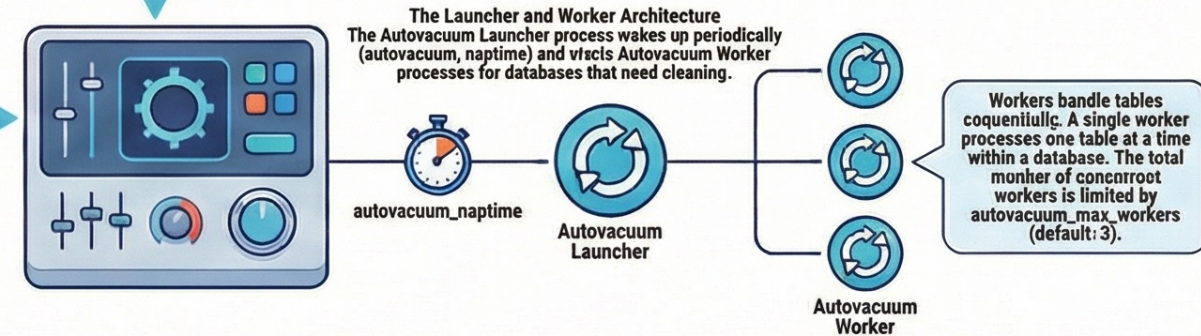
4. Table Truncation



If enough empty pages exist at the end of the table file, this final step may lock the table briefly to return the empty space to the operating system.

The Automation Engine: How Autovacuum Works

Autovacuum solves the scheduling problem. Running VACUUM manually is inefficient, too often wastes resources, too rarely causes bloat. Autovacuum automates this based on actual table activity.



Autovacuum & Autoanalyze Trigger Conditions

Autovacuum is triggered by thresholds. Autovacuum runs on a table when the number of modified or dead tuples exceeds a calculated threshold.

AUTOVACUUM Trigger

VACUUM runs if:

$$\begin{aligned} & \text{autovacuum_vacuum_threshold (50)} > \text{autovacuum_vacuum_scale_factor (50)} \\ & + \text{table_rows} \times (0.2) \end{aligned}$$

Insert-Only VACUUM Trigger

VACUUM runs to update the visibility map if:

$$\begin{aligned} & \text{autovacuum_vacuum_insert_threshold (1000)} > \text{autovacuum_vacuum_insert_scale_factor (1000)} \\ & + \text{table_rows} \times (0.2) \end{aligned}$$

AUTOANALYZE Trigger

ANALYZE (which gathers statistics for the query planner) runs if:

$$\begin{aligned} & \text{autovacuum_analyze_threshold (50)} > \text{autovacuum_analyze_scale_factor (50)} \\ & + \text{table_rows} \times (0.1) \end{aligned}$$

Default values are shown in parentheses. These server-level settings can be overridden for individual tables to fine tune behavior.

Tuning and Monitoring Vacuum Operations



Manage system load with cost-based throttling.
Autovacuum pauses periodically to avoid accumulating too many resources. It does a set amount of "work" (autovacuum_vacuum_scale_factor), then pauses for a duration (autovacuum_vacuum_scale_delay).



Default autovacuum_vacuum_scale_delay is 2ms. This is a significant change from older versions (which used 20ms) and is better suited for modern hardware.



Monitor active VACUUM jobs with pg_stat_progress_vacuum.
This system view shows the current phase, pages scanned, and number of index scans (index_vacuum_count).



Multiple Index scans indicate insufficient memory.
If index_vacuum_count is greater than 1, it means maintenance_work_mem was too small to hold all dead tuple files, forcing multiple, expensive passes over the indexes.



Log autovacuum activity for long-term analysis.
Set log_autovacuum_min_duration to 0 to log every autovacuum run, helping you identify tables that are frequently vacuumed or take a long time to process.

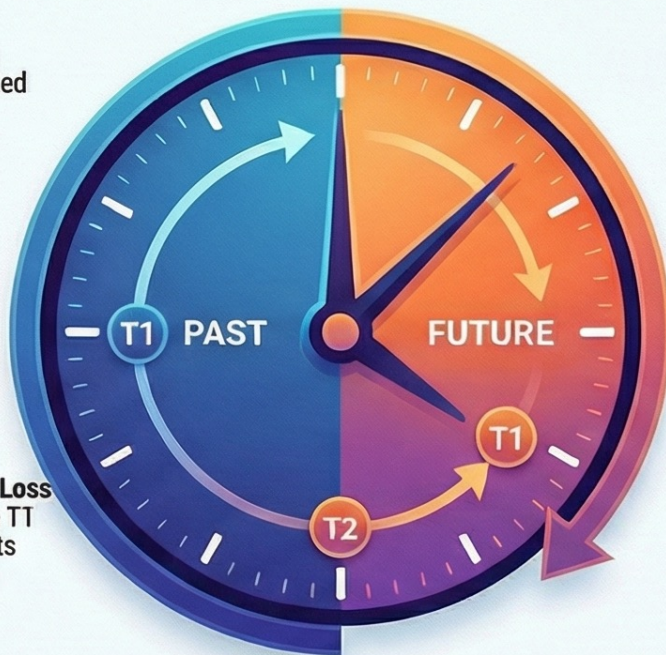
PostgreSQL Freezing: Preventing Transaction ID Wraparound

THE PROBLEM: TXID WRAPAROUND

32-Bit Limit Reached in Weeks:
4 billion Transaction IDs exhausted in about 6 weeks at 1,000 transactions/second.

Age, Not ID, Determines Order:
PostgreSQL compares "age" (transactions occurred since) rather than absolute ID number.

The Danger: "Time Travel" Data Loss
New transactions incorrectly see T1 as being in the "future", making its changes invisible and leading to catastrophic data inconsistency.



THE SOLUTION: TUPLE FREEZING

What is Freezing?
VACUUM identifies old tuples and marks them as "frozen", meaning they are universally visible to all transactions, past and future.

Modern Freezing Uses Hint Bits:
Previously changed tuple's `xltn` to special value (2). Now, sets two "hint bits" in the header, preserving original `xltn` for debugging.

Frozen Tuples are Infinitely Old:
A frozen tuple is treated as if its creation transaction is in the distant past for every other transaction, safely removing it from wraparound risk.



HOW FREEZING IS MANAGED AND AUTOMATED



Step 1: Routine Freezing
`'vacuum_freeze_min_age'`
(Default: 50 million transactions)
Standard VACUUM freezes tuples older than this age on processed pages.



Step 2: Aggressive Freezing
`'vacuum_freeze_table_age'`
(Default: 150 million transactions)
If table's oldest unfrozen TXID (`relfrozenxid`) exceeds this, VACUUM scans **all** pages.



Step 3: Forced Autovacuum
`'autovacuum_freeze_max_age'`
(Default: 200 million transactions)
A safety net. Autovacuum is forced to run on a table if it gets too old, even if disabled.



Step 4: Emergency Failsafe
`'vacuum_failsafe_age'`
(Default: 1.6 billion transactions)
The final defense. High-priority VACUUM runs, skipping non-essential work to freeze tuples as fast as possible to prevent shutdown.

TAKING MANUAL CONTROL OF FREEZING



"VACUUM FREEZE": This command performs an aggressive freeze on an entire table immediately, freezing all tuples regardless of their age.



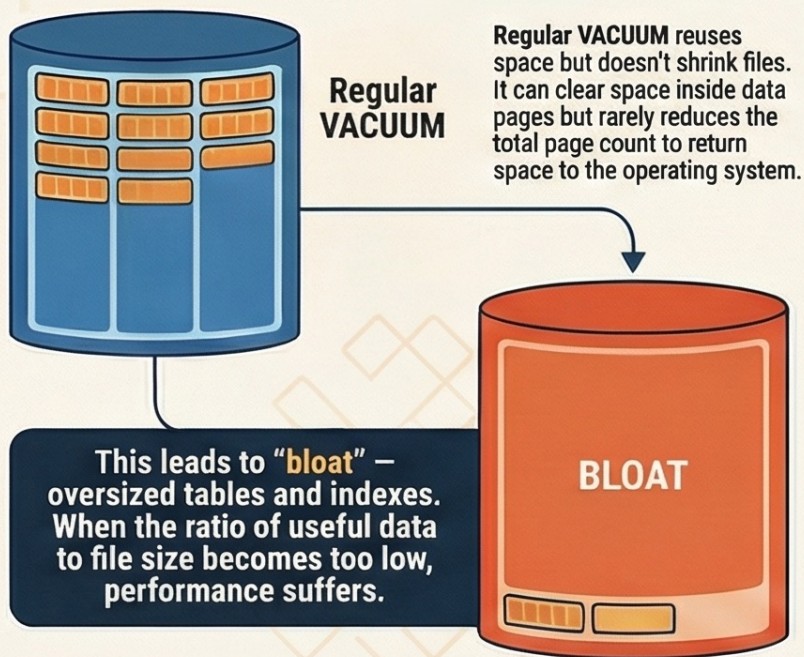
"COPY ... WITH FREEZE": For bulk-loading static data, this command freezes the rows as they are inserted, preventing future VACUUM on them (as long as they don't change).



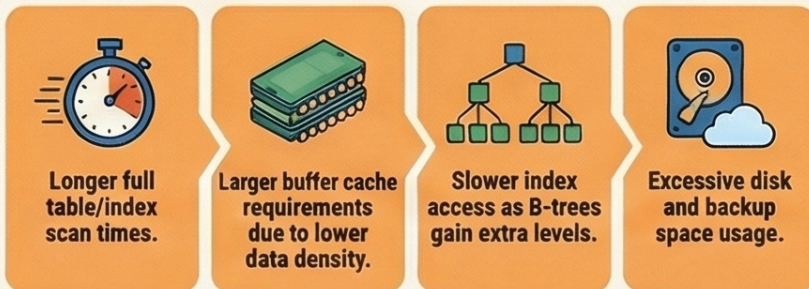
Ideal for Unchanging Data: Manually freezing is most useful for tables that are loaded once and rarely or never updated, avoiding unnecessary vacuum overhead.

Taming the Bloat: A Guide to PostgreSQL Table Reorganization

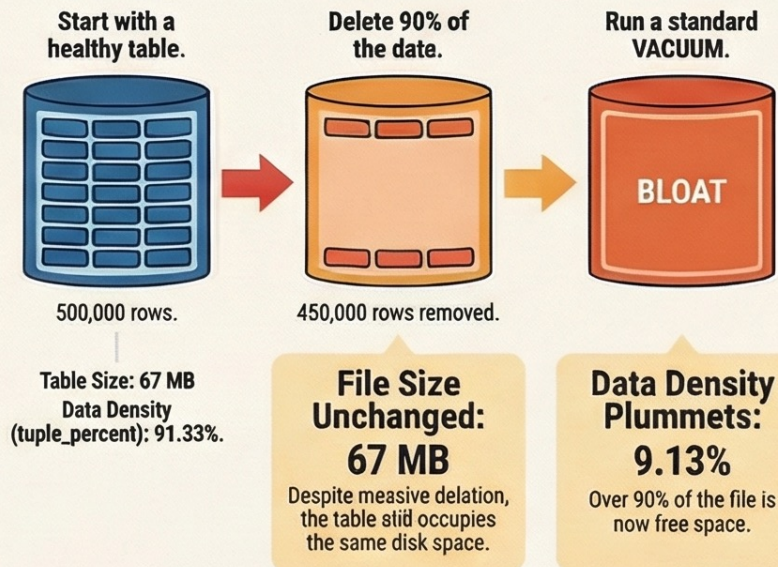
The Problem: Why Regular VACUUM Isn't Enough



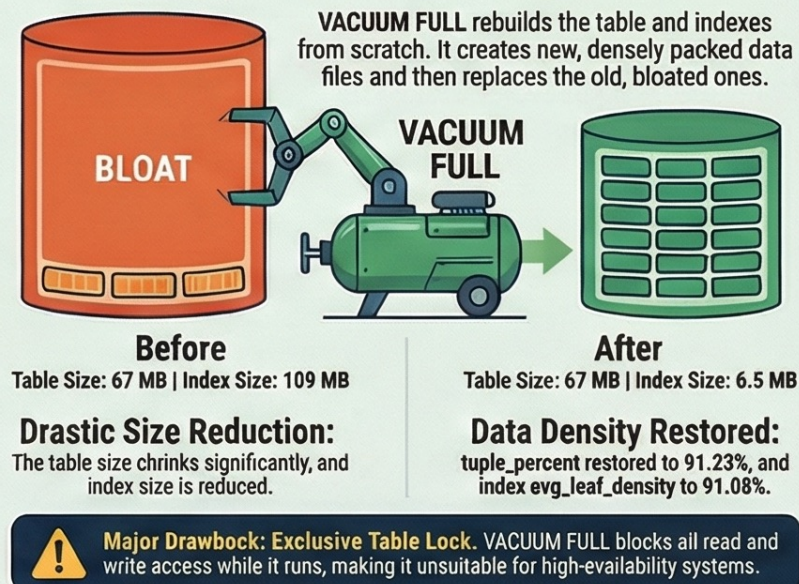
The Consequences of Bloat



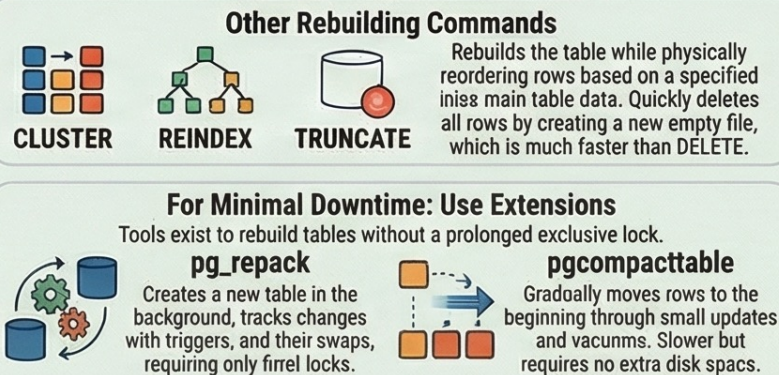
A Practical Demonstration of Bloat



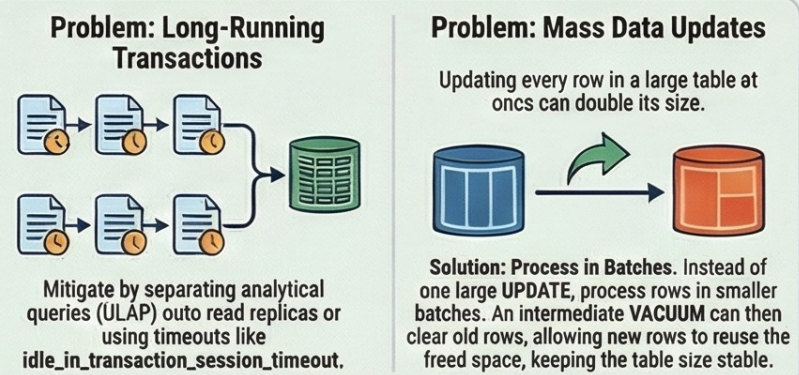
The Solution: Rebuilding with VACUUM FULL



Alternatives to VACUUM FULL



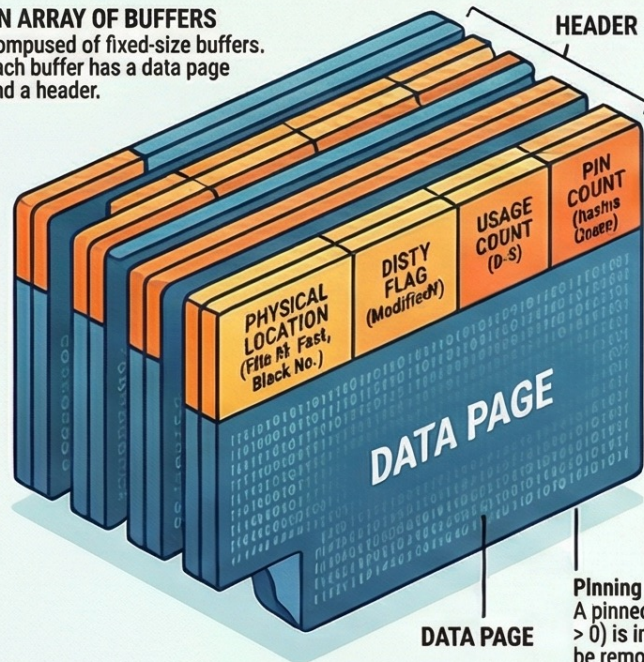
Preventing Bloat Before It Starts



PostgreSQL's Brain: A Deep Dive into the Buffer Cache

ANATOMY OF THE BUFFER CACHE

AN ARRAY OF BUFFERS
Composed of fixed-size buffers. Each buffer has a data page and a header.



THE HEADER STORES METADATA
Each header holds vital information about its data page.

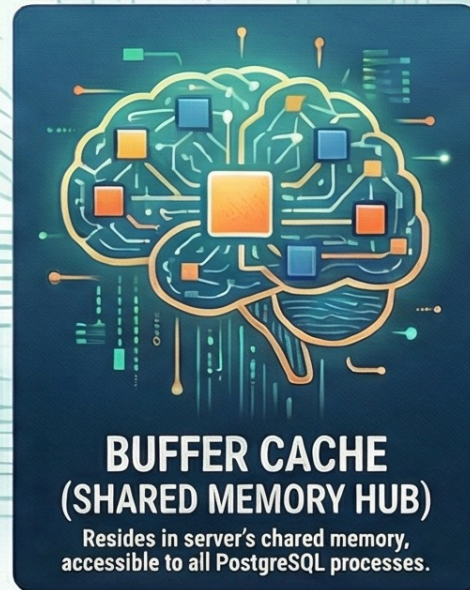
Pinning Locks the Buffer:
A pinned buffer (pin count > 0) is in-use and cannot be removed.

TUNING AND ADMINISTRATION

KEY FINDING
SIZING YOUR CACHE (shared_buffers)
Default size (125MB) is too low. Start with 25% of total RAM.

HOW TO MONITOR WITH pg_buffercache
Inspect cache contents to see cached relations and 'hot' pages (based on usage count).

HOW TO PRE-WARM THE CACHE WITH pg_preware
Load important tables into cache after restart or save/restore cache state.



HANDLING BULK OPERATIONS

THE THREAT OF CACHE POLLUTION
Large operations can flood the buffer cache, pushing out hot pages.

THE BUFFER RING SOLUTION
PostgreSQL uses small, dedicated 'buffer rings' for bulk operations to contain eviction.

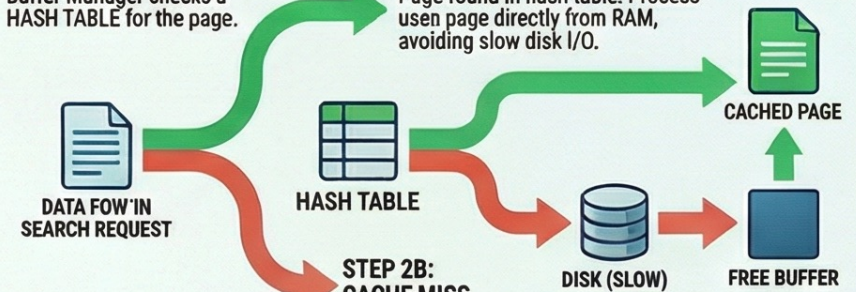


THREE EVICTION STRATEGIES:
Different strategies depending on the operation type.

THE PAGE REQUEST LIFECYCLE: HIT OR MISS?

STEP 1: THE SEARCH BEGINS
Buffer Manager checks a HASH TABLE for the page.

STEP 2A: CACHE HIT!
Page found in hash table. Process uses page directly from RAM, avoiding slow disk I/O.



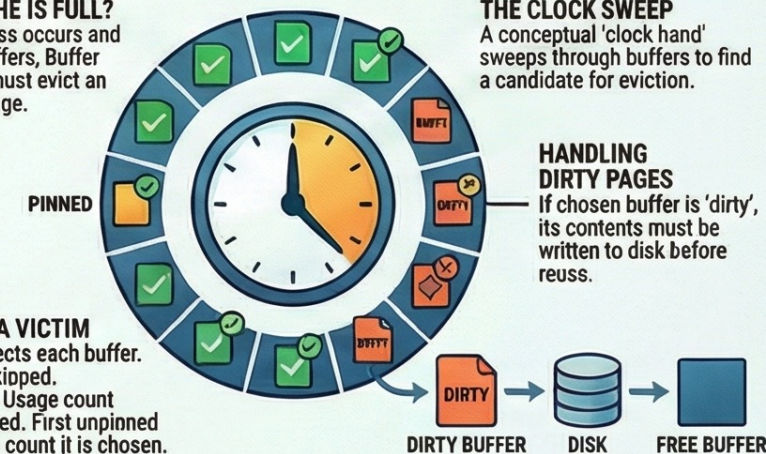
STEP 2B: CACHE MISS...
Page not in hash table. System must read from disk into a free buffer.

UPDATING THE CACHE
After a miss, page is loaded, metadata updated in header, reference added to hash table, usage count set to 1.

MAKING SPACE: THE EVICTION STRATEGY

WHAT HAPPENS WHEN THE CACHE IS FULL?
When a miss occurs and no free buffers, Buffer Manager must evict an existing page.

THE CLOCK SWEEP
A conceptual 'clock hand' sweeps through buffers to find a candidate for eviction.



FINDING A VICTIM
Hand inspects each buffer. Pinned? Skipped. Unpinned? Usage count decremented. First unpinned with usage count it is chosen.

HANDLING DIRTY PAGES
If chosen buffer is 'dirty', its contents must be written to disk before reuse.

THE EXCEPTION: LOCAL CACHE FOR TEMPORARY TABLES

TEMPORARY TABLES GET THEIR OWN CACHE
Temporary data is session-private, uses a simpler, local cache instead of the shared buffer cache.

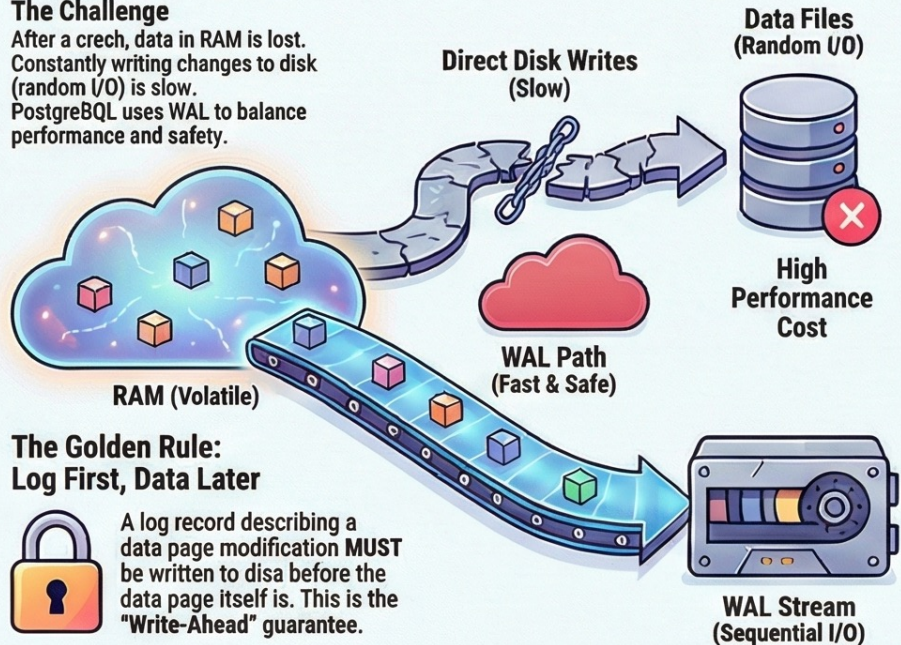
SIMPLER AND MORE EFFICIENT
Local cache doesn't need complex locking. Size controlled by temp_buffers parameter (default 8MB).

A Visual Guide to PostgreSQL's Write-Ahead Log (WAL)

1. The "Why": Durability vs. Performance

The Challenge

After a crash, data in RAM is lost. Constantly writing changes to disk (random I/O) is slow. PostgreSQL uses WAL to balance performance and safety.



The Golden Rule: Log First, Data Later



A log record describing a data page modification **MUST** be written to disk before the data page itself is. This is the "Write-Ahead" guarantee.

What Gets Logged?

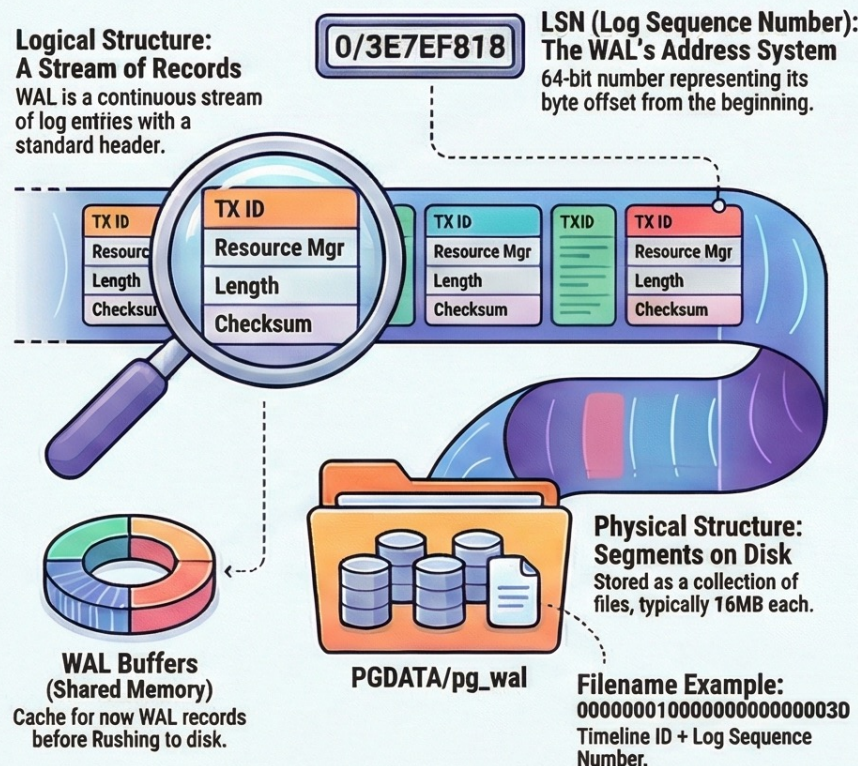
- Page modifications in buffer cache
- Transaction commits/rollbacks
- File operations

What is NOT Logged?

- Operations on UNLOGGED & temporary tables

2. The "What": Anatomy of the WAL

**Logical Structure:
A Stream of Records**
WAL is a continuous stream of log entries with a standard header.



3. The "How": The Checkpoint Process

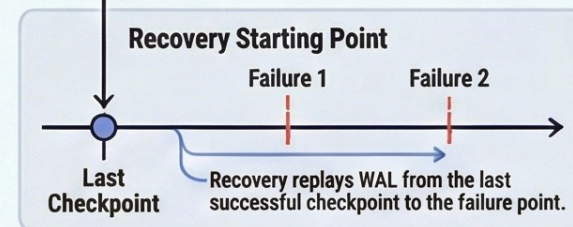


Checkpoint: Safe Starting Point

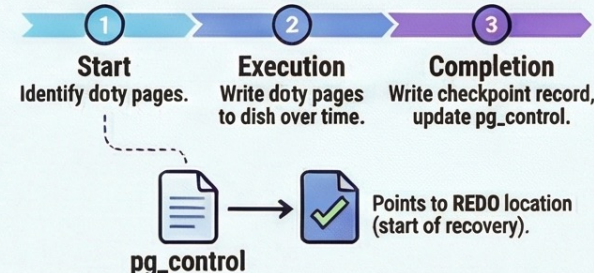
A point in the WAL sequence where all modified data pages are flushed to disk.

Why Necessary?

Create a known safe point for recovery, allowing old WAL files to be recycled.



The 3 Phases of a Checkpoint (checkpointer process)



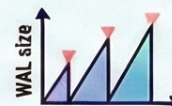
4. Putting It All Together: Crash Recovery



Crash Detected



5. Tuning and Monitoring WAL



Balancing Checkpoint Triggers

Triggers: checkpoint_timeout (time)
OR max_wal_size (size).
Goal: Most checkpoints time-based.

Warning Signs

If size-based checkpoints are much more frequent, max_wal_size may be too low.

Role of Background Writer (bgwriter)

Proactively writes pages to disk, reducing checkpointer work.

Key Metric: buffers_backend (in pg_stat_bgwriter) should be LOW.

Key Configuration Parameters

Parameter	Default	What it Does	Value
checkpoint_timeout	5 min	Max time between automatic checkpoints	
max_wal_size	1 GB	WAL size triggering a checkpoint	
checkpoint_completion_target	0.9	Spreads checkpoint I/O over time	
min_wal_size	80 MB	Min WAL size to keep for reuse	

The Performance vs. Durability Trade-off: Inside PostgreSQL's WAL Mode

Synchronous Commit: The Safety-First Approach

Transaction only 'committed' after WAL records physically written to disk.
(Default: `synchronous_commit = on`)



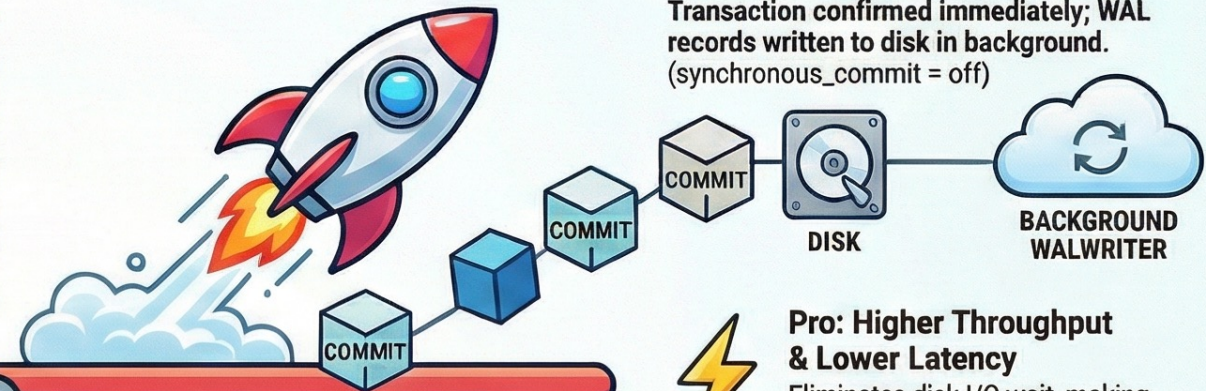
Pro: Maximum Reliability
Ensures ACID durability requirements are met. Once a commit is acknowledged, the data is safe from crashes.

Con: Slower Performance
Waiting for disk I/O increases latency and reduces throughput.

Benchmark Data
Synchronous Mode
Transactions (30s): 20,123
Avg Latency: 1.491 ms
TPS: 670.8

Asynchronous Commit: The Speed-Focused Alternative

Transaction confirmed immediately; WAL records written to disk in background.
(`synchronous_commit = off`)



Pro: Higher Throughput & Lower Latency
Eliminates disk I/O wait, making commit significantly faster.

Con: Risk of Data Loss
...recently committed transactions can be lost (a window of up to 0.6 seconds by default).

Benchmark Data
Asynchronous Mode
Transactions (30s): 61,809
Avg Latency: 0.485 ms
TPS: 2066.4

Ensuring Fault Tolerance

Challenge 1: Non-Atomic Writes
Database page (8KB) written in smaller blocks (4KB). Crash can leave corrupted, partial page.

Solution: Full Page Writes (FPI)
PostgreSQL writes a full copy (FPI) to WAL on first modification after checkpoint for recovery.

FPIs increase WAL size (e.g., 71.5% of data). Enable '`wal_compression = on`' to reduce size (e.g., 29 MB to 10 MB).

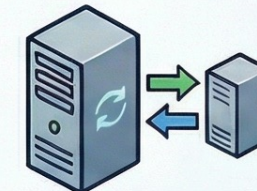
Challenge 2: Data Corruption
Hardware failures can silently corrupt data in memory, transfer, or on disk, spreading to backups.

Solution: Checksums
Enable '`data_checksums`' to verify page checksum on read. WAL records always protected by checksums.

Understanding WAL Levels



Level 1: Minimal
(`wal_level = minimal`)
Logs essential crash recovery info. Skips bulk operations to save space. No backups/replication.



Level 2: Replica
This is the default level.
Logs enough for Point-in-Time Recovery and Physical Streaming Replication. All data changes logged.



Level 3: Logical
(`wal_level = logical`)
Includes Replica info plus data for Logical Decoding. Required for Logical Replication to other systems.

Understanding Database Locking in PostgreSQL

Database locking is a fundamental mechanism for managing concurrency, preventing data corruption by controlling simultaneous access to shared resources. PostgreSQL employs a sophisticated system to balance performance and data integrity.

The Fundamentals of Locking



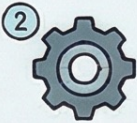
What is a Lock?

A lock is a mechanism that controls concurrent access to a shared resource, ensuring that multiple processes don't interfere with each other.

The Lock Lifecycle



1 Acquire a lock on a resource



2 Perform its operation



3 Release the lock so other processes can use the resource.



The Trade-off: Finer-grained locks increase concurrency but also increase the number of locks to manage. Coarse-grained locks are simpler but limit concurrency.

Common Heaverlyweight Lock Types

- relation:** Table-level locks
- tuple:** Row-level locks
- transactionid:** Locks on a transaction itself
- page:** Locks on data pages (used by some indexes)
- advisory:** User-managed locks
- object:** Locks on non table database objects



No Concurrency = No Locks

If a resource is not accessed simultaneously by processes it doesn't require a lock.

Key Characteristics of Locks

Granularity: The Scope of the Lock

Coarse-grained



A table-level lock is coarse, preventing all concurrent access to that table.

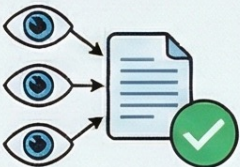
Fine-grained



A row-level lock is fine, allowing processes to work on different rows of the same table simultaneously.

Lock Modes: The Type of Access

Shared Mode



Allows multiple processes to read a resource simultaneously.

Exclusive Mode



Prevents any other process from accessing the resource; used for writing.

Basic Compatibility (Shared vs. Exclusive)

	Shared	Exclusive
Shared	✓	✗
Exclusive	✓	✗
Exclusive	✗	✗

Classifying Locks by Duration



Long-term Locks

Held for a potentially long time, often until a transaction ends. Protects resources like tables and rows. Includes advanced features like wait queues and deadlock detection.



Short-term Locks

Acquired for very brief periods (microseconds) to protect data structures in shared memory. Managed automatically with simple infrastructure.

A Closer Look at PostgreSQL's Heavyweight Locks



What are Heavyweight Locks?

These are the long-term, object-level locks in PostgreSQL, visible in 'pg_locks' view, managed in a shared memory pool.

Common Heavyweight Lock Types ("locktype")



relation: Table-level locks



tuple: Row-level locks



transactionid: Locks on a transaction itself



page: Locks on data pages



advisory: User-managed locks



object: Locks on non-table objects

	AS	RS	SUE	E	AE
AS	AS	RS	SUE	E	AE
RS	AS	RS	SUE	E	AE
SUE	AS	RS	SUE	E	AE
E	AS	RS	SUE	E	AE
AE	AS	RS	SUE	E	AE

Table-Level Lock Modes & Compatibility

8 Modes for Maximum Concurrency

PostgreSQL provides eight table level lock modes to allow the maximum number of operations to run in parallel without conflict.



Most Permissive: Access Share

Used by 'SELECT' queries. Compatible with all modes except the most restrictive, allowing reads alongside almost any other operation.



Most Restrictive: Access Exclusive

Used by commands like 'DROP TABLE' or 'VACUUM FULL'. Incompatible with all other lock modes, ensuring no other process can access the table.

The Wait Queue in Action



What is a Wait Queue?

When a process tries to acquire a lock that conflicts with an existing lock, it enters a "first-in, first-out" wait queue until the resource is freed.

Table-Level Lock Compatibility Matrix

	AS	RS	RE	SRE	S	X	AE
AS	AS	RS	RE	SRE	S	X	AE
RS	AS	RS	RE	SRE	S	X	AE
RE	AS	RS	RE	SRE	S	X	AE
SRE	AS	RS	RE	SRE	S	X	AE
S	AS	RS	RE	SRE	S	X	AE
X	AS	RS	RE	SRE	S	X	AE
AE	AS	RS	RE	SRE	S	X	AE



Deadlock Detection

If two or more transactions waiting for each other in a circular chain, PostgreSQL automatically detects this "deadlock" and terminates one of the transactions.

A Deep Dive into PostgreSQL Row-Level Locking

1. The PostgreSQL Approach: "Virtual" Locks in the Row Header

No Heavyweight Locks in Memory

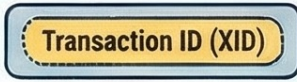


Heavyweight lock structures avoided



Heavyweight lock structures avoided

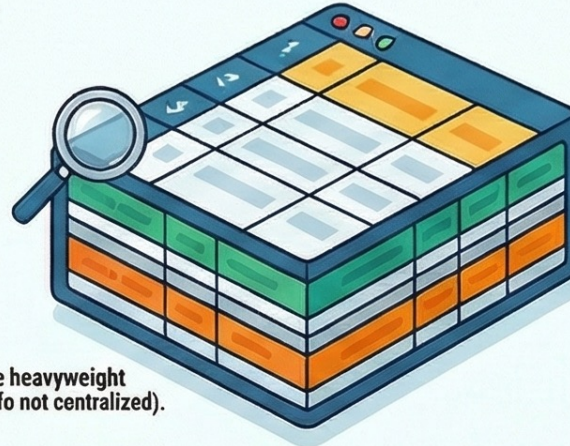
The 'xmax' Field is the Lock Holder



Transaction ID (XID) written here marks row as locked.

The Trade-Off: Efficiency vs. Complexity

Efficiency: Lock countless rows with no extra RAM cost. Complexity: Waiting transactions use heavyweight locks on XIDs to form a queue (lock info not centralized).

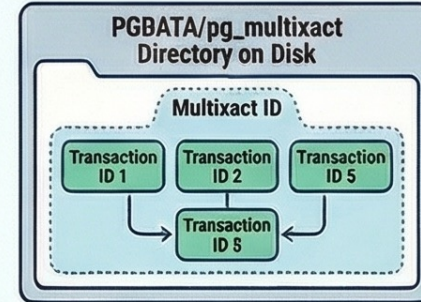


3. Handling Shared Locks with Multixactions

How can multiple transactions lock one row?

XID	xmax
	Multixact ID

The 'xmax_is_multi' Hint Bit: Special flag in row header indicates 'xmax' holds a Multixact ID.

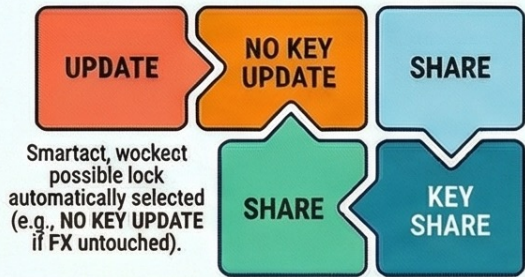


Multixact Details are Stored on Disk: Information on grouped transactions is not in the row itself.

2. The Four Modes of Row-Level Locks

EXCLUSIVE MODES

SHARED MODES



Smartest, weakest possible lock automatically selected (e.g., NO KEY UPDATE if FK untouched).

Granular control over concurrency.

LOCK CONFLICT MATRIX (X = Conflict)

	Key Share	Share	No Key Update	Update
Key Share			X	X
Share			X	X
No Key Update	X	X		X
Update	X	X	X	

Simple SELECT Queries Never Lock Rows: To acquire a row-level lock during a read, you must explicitly use a SELECT ... FOR [MODE] command.

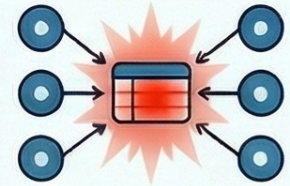
4. The Lock Queue and Its Pitfalls



Preventing Resource Starvation: Heavyweight tuple lock creates orderly queue, prevents indefinitely waiting race conditions.



The Waiting Queue Can Collapse: Waiting transactions in READ COMMITTED may abandon orderly queue and race to acquire lock after primary release.



QUOTES ON HOTSPOTS: "Updating the same table rows from several processes at once is a bad idea." - High contention on "hotspot" rows becomes severe performance bottleneck.

5. Strategies for Managing Lock Waits



'NONAZT': Don't Wait, Fail Fast. Error out immediately if locked; handle programmatically.

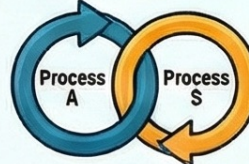


'SKIP LOCKED': Process What's Available. Ignore locked rows, proceed to next available (ideal for parallel job queues).



'lock_timeout': Set a Time Limit on Waiting. Abort statement after specified duration, prevents indefinite stuck operations.

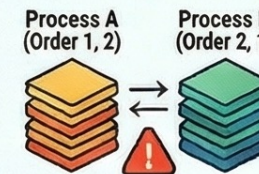
6. Deadlocks: The Vicious Cycle



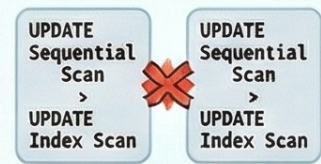
What is a Deadlock? Two or more transactions in a circular dependency, each waiting for a resource held by another.



Automatic Deadlock Detection & Resolution: Proactively checks for circular wait-for graph if wait exceeds "deadlock_timeout" (default 1s). Terminates one transaction to release locks.



Common Cause: Inconsistent Lock Order. Different application processes locking same resources in different orders. Solution: Always lock resources in a consistent, predetermined order.



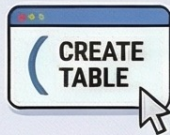
The Hidden Deadlock Trap: Even two "WRATE" statements can deadlock if query plans cause opposite locking orders.

A Guide to Specialized Locking in PostgreSQL



NON-OBJECT LOCKS

Locks resources that are not standard tables or rows. Used to protect system catalog objects like tablespaces, schemas, roles, and data types during transactions.



AccessShareLock on public schema, prevents being dropped

- 1 Identified by a trio of values: database (database OID) (system catalog OID), and objid (object OID within the catalog).



RELATION EXTENSION LOCKS

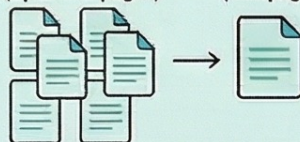
Protects a table or index when it needs to grow. Ensures only one process can add a new physical page to a relation's file at a time, preventing data corruption.

- 1 Released immediately after use. Unlike most locks, released right after the page is added, not at the end of the transaction.

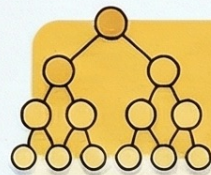
- ✓ Does not cause deadlocks. Not included in the deadlock detection graph.

- ✓ Released immediately after use. Unlike most locks, released right after the page is added, not at the end of the transaction.

HEAP FILES (Up to 512 pages) → B-TREE INDEXES (One page)



Heap files extend faster than B-tree indexes.



GIN INDEX

PAGE LOCKS

A specialized lock used exclusively for GIN indexes. Manages concurrent insertions of composite values.

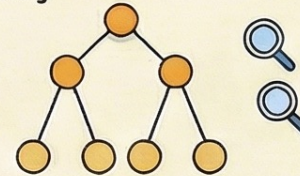
- 1 Enables high-performance "fast update" indexing.



UNSORTED PENDING LIST



METAPAGE LOCK (temporary)



MAIN GIN STRUCTURE

- ✓ Does not block normal index reads. The exclusive lock on the metapage during the batch move does not interfere with processes using the index for searching.
- ✓ Like extension locks, they are released immediately and don't cause deadlocks.



ADVISORY LOCKS

A cooperative, application-defined locking mechanism.

Explicitly acquired and released by developers to manage application-level logic.



APPLICATION PROCESS



CUSTOM LOCK ID (e.g., via hashtext)



EXTERNAL RESOURCE / API

- 1 Can persist across transactions. By default, advisory locks are session-level, remaining held until explicitly released or the session ends.
- 3 Managed with a dedicated set of functions: `pg_advisory_lock`, `unlock`, `try`, `share`, `zact`.

PREDICATE LOCKS

A "lock" that doesn't actually lock anything. Used by the **SERIALIZABLE** isolation level to track data dependencies between transactions.

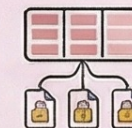
- 1 The core of Serializable Snapshot Isolation (SSI). It tracks "read-write dependencies." If a transaction reads a row that another concurrent transaction then modifies, a predicate lock records this dependency.

SEQUENTIAL SCAN

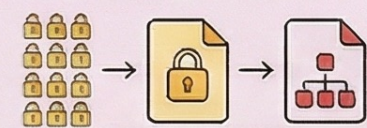


Index Scans are more granular than Sequential Scans.

INDEX SCAN



- 1 Prevents "phantom reads" and "write skew" anomalies. When a transaction commits, PostgreSQL checks the dependency graph. If a dangerous pattern (potential anomaly) is found, the transaction is aborted to preserve true serializability.



Features automatic lock escalation to conserve memory.

PostgreSQL

A Deep Dive into PostgreSQL's Memory Locks

The Locking Toolkit: Spinlock vs. LWLock

Understanding the Core Mechanisms



Spinlock: The Sprinter



Type: Exclusive lock using atomic CPU commands for extremely short durations (a few CPU cycles).



How It Handles Contention: "Busy-Waiting". Processes repeatedly check in a tight loop ("spin") until free. Efficient only for rare, brief contention.



Limitations: Exclusive mode only. No built-in deadlock detection or monitoring instrumentation.



LWLock (Lightweight Lock): The Marathoner



Type: Acquired for longer periods to manage data structures, sometimes spanning I/O operations.



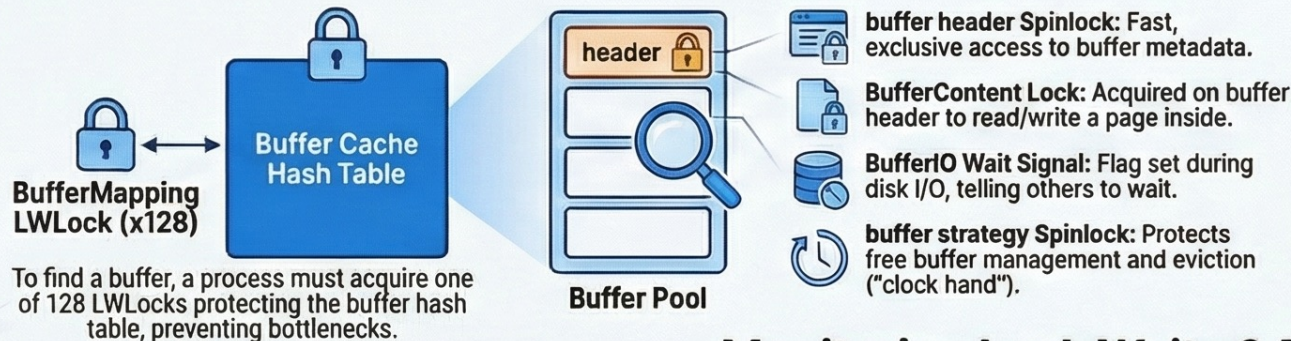
How It Handles Contention: Shared & Exclusive Access. Supports multiple simultaneous readers and single writer, offering flexibility.



Limitations: LWLocks lack deadlock detection but include instrumentation for observability. No queue; access is relatively random.

Locks in Action: Critical Shared Memory Structures

The Buffer Cache



The WAL Buffer



Monitoring Lock Waits & Performance

Turning Abstract Concepts into Actionable Insights

Why Monitor Waits?



Locks can cause waits, creating performance bottlenecks. Tracking "wait events" is crucial.



Step 1: Real-Time Check with `pg_stat_activity`. Live snapshot of current `'wait_event_type'` and `'wait_event'` for active processes.



Step 2: Historical Analysis Since `'pg_stat_activity'` is not cumulative, an extension like `'pg_wait_sampling'` is needed to collect wait statistics over time for a complete picture.

Common Wait Event Categories & Analysis

- LWLock**
- Lock (heavy)**
- BufferPin**
- ID (disk)**
- Client (network)**
- IPC (inter-process)**

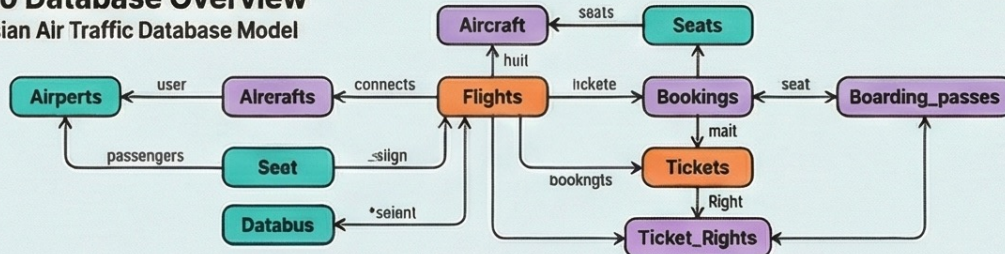


Sample Analysis: `'pg_wait_sampling_profile'` might reveal frequent `'IO'` events like `'WALTyoc'` and `'WALWrite'`, indicating a disk I/O bottleneck.

The Life of a PostgreSQL Query: From Text to Result

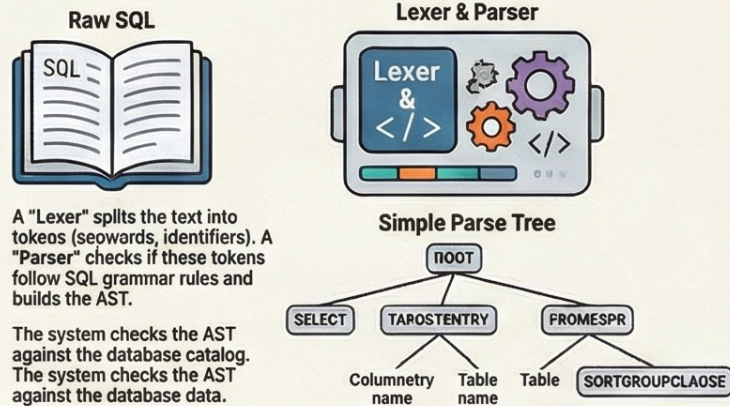
Demo Database Overview

A Russian Air Traffic Database Model



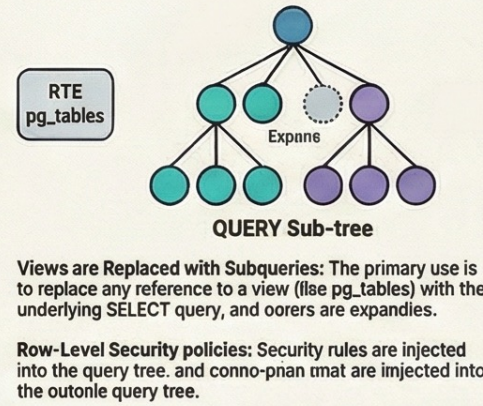
Step 1: Parsing

The server reads the raw SQL text to understand its structure.



Step 2: Transformation (Rewrite)

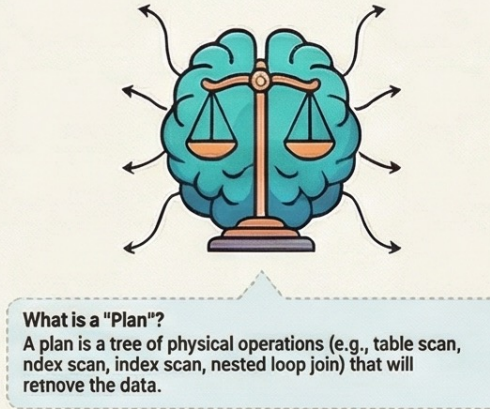
The initial parse tree is modified and expanded.



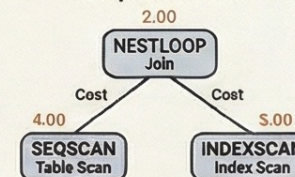
Step 3: Planning (Optimization)

The Planner finds the most efficient execution path.

SQL is declarative... The planner, a cost-based optimizer, evaluates many potential plans to find the "cheapest" one.



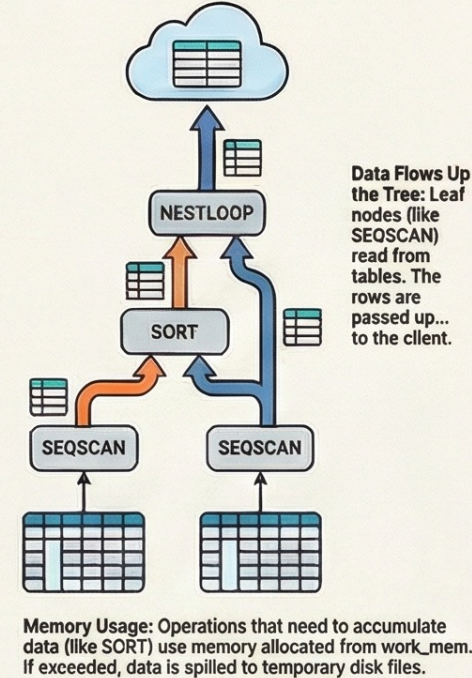
Sample Plan Tree



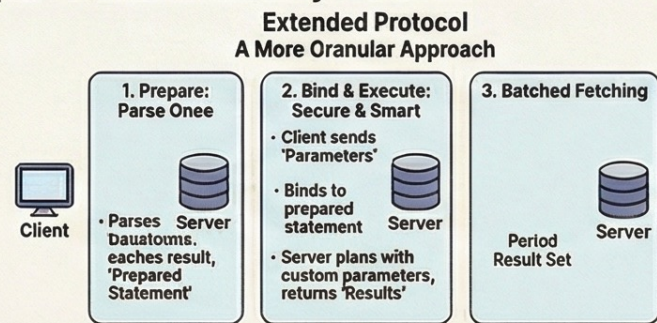
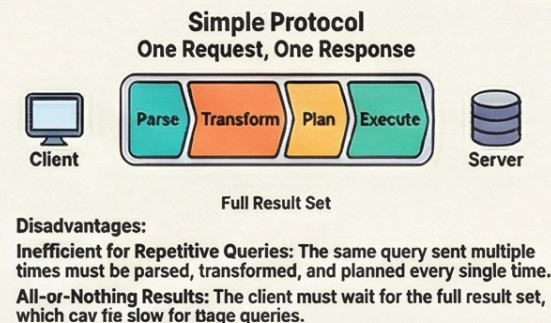
Step 4: Execution

The Executor runs the chosen plan.

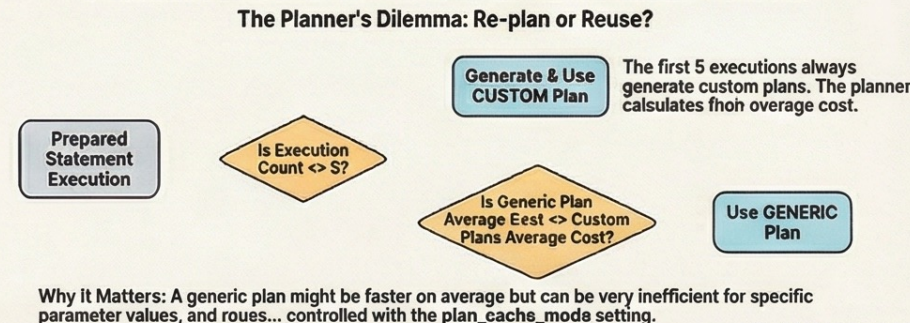
The executor works its way through the plan tree, with each parent node "pulling" rows from its child nodes.



Simple vs. Extended Query Protocols

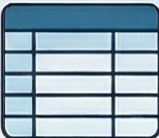


Custom vs. Generic Plans in the Extended Protocol



The PostgreSQL Planner's Secret Weapon: A Guide to Query Statistics

From basic table counts to advanced multivariate analysis, here's how PostgreSQL estimates costs and chooses efficient execution plans.




Basic counts for every table.

PostgreSQL stores high-level statistics for each table in the `pg_class` system catalog, providing the planner with a starting point.

How are they collected?

Statistics are gathered during manual or automatic **ANALYZE** operations, as well as other maintenance tasks like **VACUUM**, **CLUSTER**, and **CREATE INDEX**.




KEY FINDING


The planner is smarter than its data.

If actual file size is larger than `relpages` suggests, the planner adjusts its row estimate upwards.


KEY METRICS: `reltuples` & `relpages`



reltuples (Row Count)
Number of rows; default estimate for queries without filters.




relpages (Table Size in 8KB Pages)
Total table size in pages.




STATISTICS COLLECTION

EXAMPLE




Before ANALYZE
Planner assumes default (e.g., 410 rows, `reltuples` = -1).



After ANALYZE
Estimate becomes accurate based on gathered data.

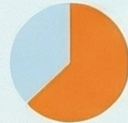
DEEPER DIVE: COLUMN-LEVEL STATISTICS



Understanding the data inside each column.


Stored in `pg_statistic`, these stats describe data distribution, crucial for estimating **WHERE** clause selectivity.

`null_frac`: The frequency of NULLs



Fraction of NULL entries, used to estimate **IS NULL** and **IS NOT NULL** conditions.

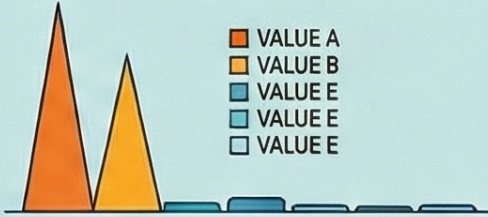
`n_distinct`: The number of unique values



For uniform data, planner divides total rows by `n_distinct` for equality checks.

VISUAL CONCEPT

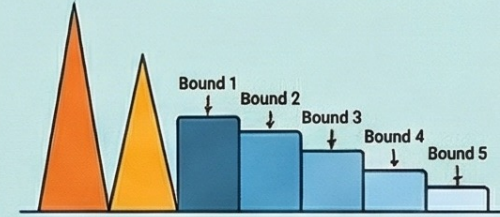
For non-uniform data: Most Common Values (MCV)



Stores exact frequencies of popular values for highly accurate estimates on specific value filters.

VISUAL CONCEPT


For many unique values: Histograms



Groups remaining data into buckets with roughly equal items, storing only boundaries for range queries (`<`, `+`).

ADVANCED TOOLS: EXPRESSION & MULTIVARIATE STATISTICS

PROBLEM



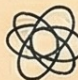
`WHERE extract(month FROM date_col) = 1`

Planner can't guess function results, uses inaccurate fixed guess (e.g., 0.5%).

SOLUTION

Planner can't guess function results, uses inaccurate fixed: (e.g., 0.5%).

PROBLEM



`flight_ny` — `departure_sirport`

SOLUTION

Columns are not always independent. Planner assumes independence, leading to severe underestimates for correlated columns.

Solution 1: Extended Statistics on Expressions

Use `CREATE STATISTICS` to collect detailed stats (`n_distinct`, MCV) on expression results.

Solution 2: Indexing an Expression


Creating an index on an expression automatically generates and stores statistics.

SOLUTION

Solution: Multivariate Statistics
Use `CREATE STATISTICS` to analyze relationships, including Functional Dependencies, `M-Distinct` Combinations, and Multivariate MCV Lists.


OTHER KEY STATISTICS FOR FINE-TUNING

`avg_width`: Average value size



Average size of data in a column, critical for memory estimates for sorting/hashing.

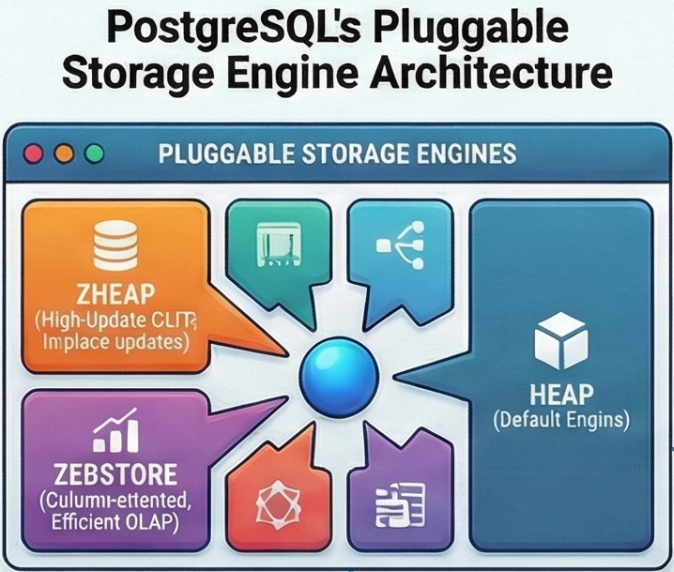
`correlation`: Physical vs. Logical Order



Highly Sorted
Measures correlation between disk storage and logical value order. High correlation (+1 or -1) significantly speeds up index scans.

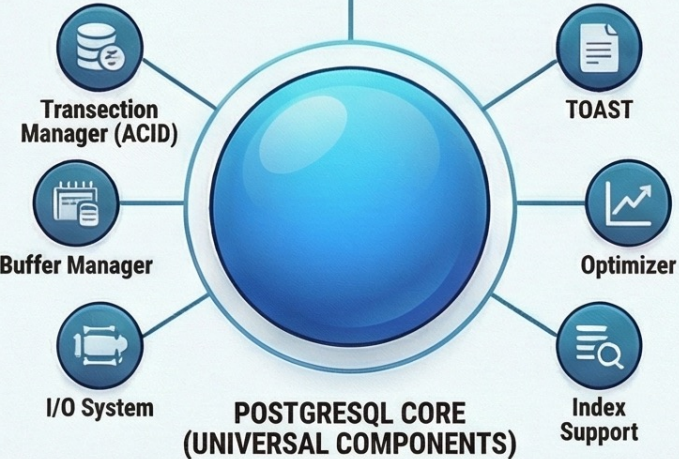
NotebookLM

Anatomy of a PostgreSQL Table Scan: From Sequential to Parallel



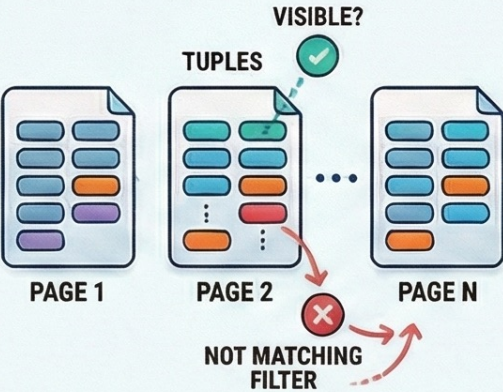
Defined by the Storage Engine:
Tuple Format & Data Structure,
Table Scan Implementation,
INSERT/DELETES/UPDATE Logic,
Visibility Rules, Backup/Analyze
Procedures.

CORE / STORAGE ENGINE
Provided by PostgreSQL Core / Defined by the Storage Engine

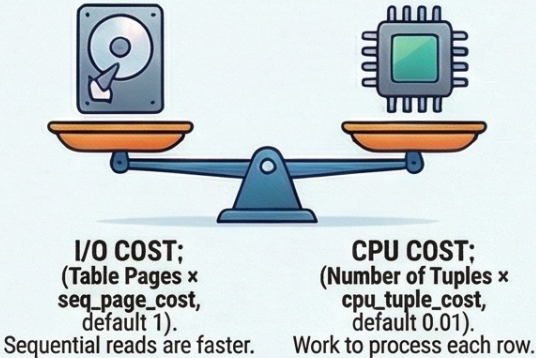


The Sequential Scan (Seq Scan)

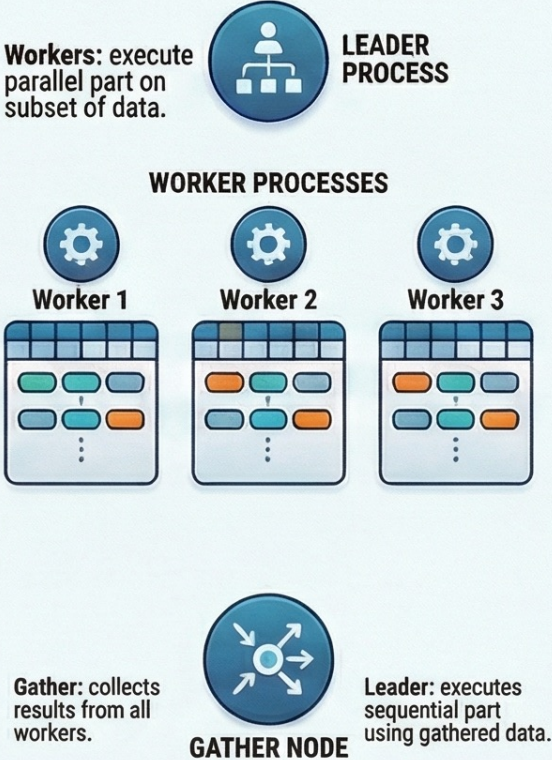
The Foundational Access Method.
Reads the entire table file, page by page, checking each tuple for visibility and filtering out those that don't match the query.



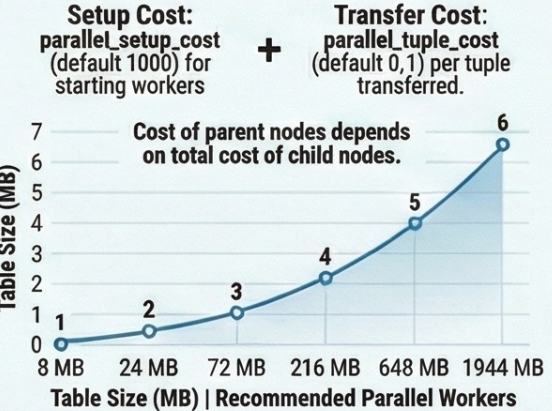
COST ESTIMATION FORMULA:
Total Cost = I/O Cost + CPU Cost



The Parallel Sequential Scan



Cost of Parallelism



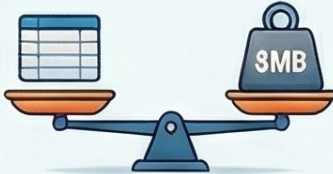
Rules and Limits of Parallelism

Parallelism is Not Automatic
Conditions and parameters must be met.

When Parallel Plans are NOT Used

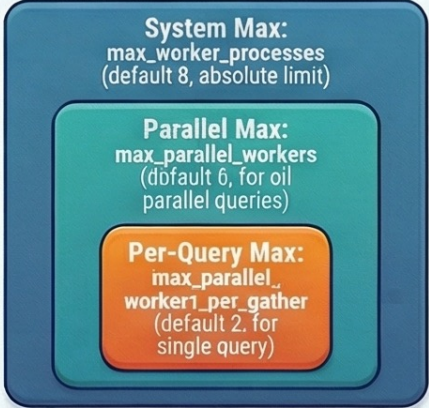
- ✓ **Data Modification** (UPDATE, DELETE, SELECT FOR UPDATE)
- ✗ **Pausable Queries** (Cursors, PL/pgSQL FOR loops)
- ✗ **Unsafe Functions** (marked as PARALLEL UNSAFE)
- ✗ **Restricted Plan Nodes** (CTE Scan, SubPlan, InitPlan run sequentially)

Minimum Table Size



min_parallel_table_scan_size
Parallel scan only considered if table size exceeds min_parallel_table_scan_size.

Worker Process Limits



ANATOMY OF A POSTGRESQL INDEX: THE EXTENSIBLE ENGINE

SECTION 1: THE CORE ARCHITECTURE OF POSTGRESQL INDEXING

WHAT IS A POSTGRESQL INDEX?

A database object used to speed up data retrieval. It links indexed values (keys) to data rows (tuples) using a Tuple ID (TID), avoiding full table scans.



B-TREE

Suited for range queries, equality. Default.



HASH

Efficient for simple equality lookups.



GiST

Generalized Search Tree. Extensible for complex types (geometry, text search).



GIN

Generalized inverted index, ideal for full text search, arrays, JSONB.



SP-GiST

Space Partitioned GiST. Good for non-balanced datasets.



BRIN

Block Range Index. Very compact for large, sequential datasets.

INDEXING ENGINE

The common engine that coordinates all index operations: retrieving TIDs, checking data visibility, and re-checking conditions.

OPERATOR CLASSES

(e.g., int4_ops, text_ops, gist_int4_ops)

DATA TYPES

(e.g., integer, text, point)

int4_ops

integer

text_ops

text

gist_int4_ops

boolean

bool_ops

ACCESS METHODS

(e.g., btree, gist)

SECTION 2: CUSTOMIZING INDEXES WITH OPERATOR CLASSES

WHAT IS AN OPERATOR CLASS?

A set of operators (e.g., "<", "<=", ">") and functions that teaches an access method how to handle a specific data type. A single data type can have multiple operator classes.

WHAT IS AN OPERATOR FAMILY?

A collection of related operator classes that handle similar data types (e.g., integer, ops family for aaalint, integer, bigint).

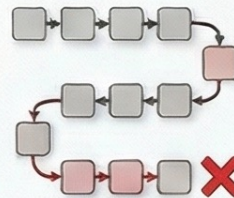
CASE STUDY: SEARCHING TEXT WITH 'LIKE'

PROBLEM:

Standard 'text_ops'

table	name
	'Elena'
	'Elena'
	'Elena'
	'Elena'
	'Elena'

...WHERE name LIKE 'ELENA%'



SEQUENTIAL SCAN

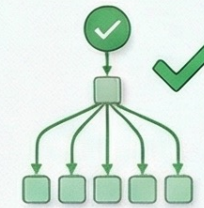
Falls for LZKE queries with now C outlows, leading to slow Sequential Scan.

SOLUTION:

'text_pattern_ops'

table	name
	'Elena'
	'Elena'
	'Elena'
	'Elena'
	'Elena'

...WHERE name LIKE 'ELENA%' using text_pattern_ops



BITMAP INDEX SCAN

Explicitly use text_pattern_ops for LZKE queries, enabling a much faster Bitmap Index Scan.

SECTION 3. INDEX CAPABILITIES: THE 3 LEVELS OF PROPERTIES

1. ACCESS METHOD PROPERTIES

Define core capabilities of the entire index type (e.g., all B-tree indexes).



can_order: Can return data in sorted order.



can_unique: Can enforce UNIQUE & PRIMARY KEY constraints.



can_multi_col: Can be null on more than one column.



can_exclude: Can support EXCLUDE constraints for advanced checks.



can_include: Can store extra key columns (enforcing mass).

2. INDEX-LEVEL PROPERTIES

Apply to a specific, existing index on a table.



clusterable: Table rows can be reordered to match index (CLUSTER).



index_scan: Can return row IDs one by one.



bitmap_scan: Can return a bitmap of all matching row IDs at once.



backward_scan: Can scan the index in reverse order.

3. COLUMN-LEVEL PROPERTIES

Define behaviors for a specific column within an index.



asc/desc: Specifies column storage order.



orderable: Can actify an ORDER BY clause.



returnable: Value can be read directly from the index (index only scan).



search_array: Efficiently search for multiple values (M).



search_nulls: Efficiently search for IS NULL or IS NOT NULL.

A Visual Guide to PostgreSQL Index Scans

HOW A STANDARD INDEX SCAN WORKS



1. Find matching Tuple IDs (TIDs) in index.
2. Fetch corresponding data row (tuple) from heap.

ANATOMY OF AN EXPLAIN PLAN

```
Index Scan using bookings_pkey on bookings
Index Cond: (book_ref = '9AEB6'::bpcher)
Filter: (total_amount = 46500.00)
```

Checked using
index

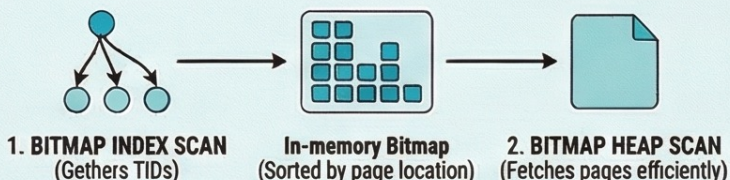
Re-checked after
fetching data

THE COST OF AN INDEX SCAN

Total cost = Index page access +
Heap page access & tuple
processing



BITMAP SCANS: THE BEST OF BOTH WORLDS

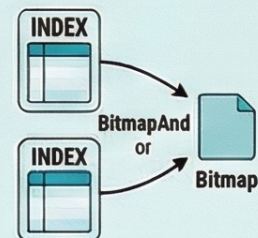


A Smarter Two-Phase Approach: Collects all TIDs, sorts them by page location, then fetches pages in physical order (each page read once).

THE SOLUTION FOR LOW CORRELATION

Turns slow, random I/O into predictable, ordered I/O.

COMBINING MULTIPLE INDEXES



THE ROLE OF WORK_MEM

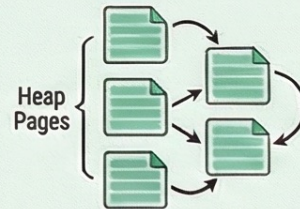
Bitmap built in memory (limited by work_mem): If space runs out, becomes 'lcasy', requiring re-checks.

THE DECISIVE FACTOR: DATA CORRELATION

What is Correlation? Relationship between physical row order and logical index order.

HIGH CORRELATION: THE BEST CASE (Efficient Sequential Reads)

Rows physically next to each other;
efficient, sequential reads.

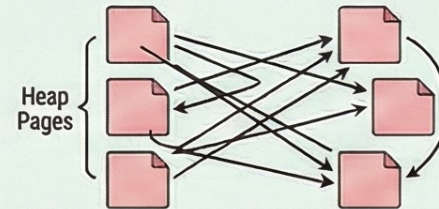


High Correlation

Near 1 or -1

LOW CORRELATION: THE WORST CASE (Inefficient Random Reads)

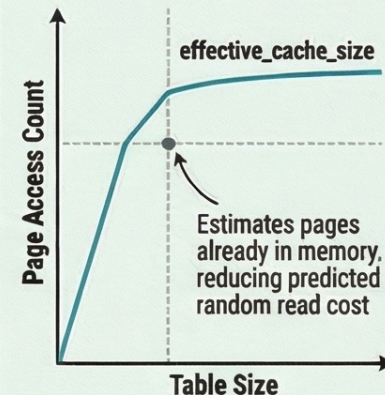
Rows scattered randomly; many
slow, random I/O operations.



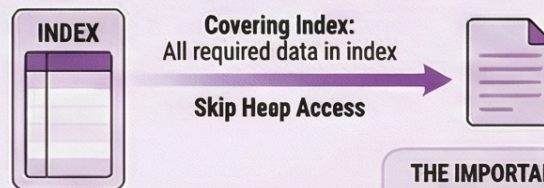
Low Correlation

Near 0

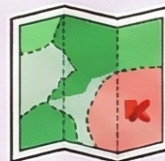
CACHING AS A PERFORMANCE BUFFER



INDEX-ONLY SCANS: SKIPPING THE HEAP



If all required columns are in the index, heap access is avoided for significant performance boost.



THE VISIBILITY MAP CHECK
PostgreSQL checks Visibility Map to ensure rows are visible to current transaction.

THE IMPORTANCE OF VACUUM

Before VACUUM:
Heap Fetches 192,109
(Planner thinks heap visit needed)

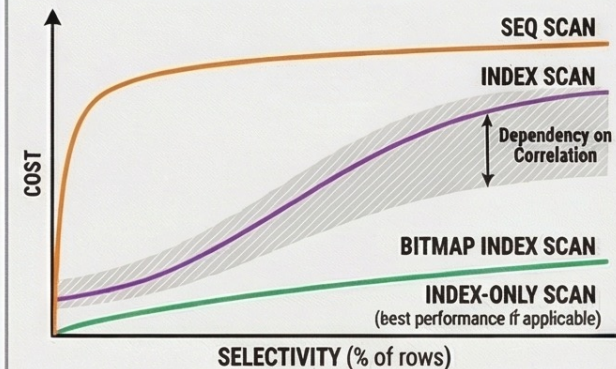
After VACUUM:
Heap Fetches 0
(Visibility Map updated, fewer heap visits)

CREATING COVERING INDEXES WITH INCLUDE

Use **INCLUDE** clause to add extra columns to index without being part of the key.

A COMPARATIVE SUMMARY

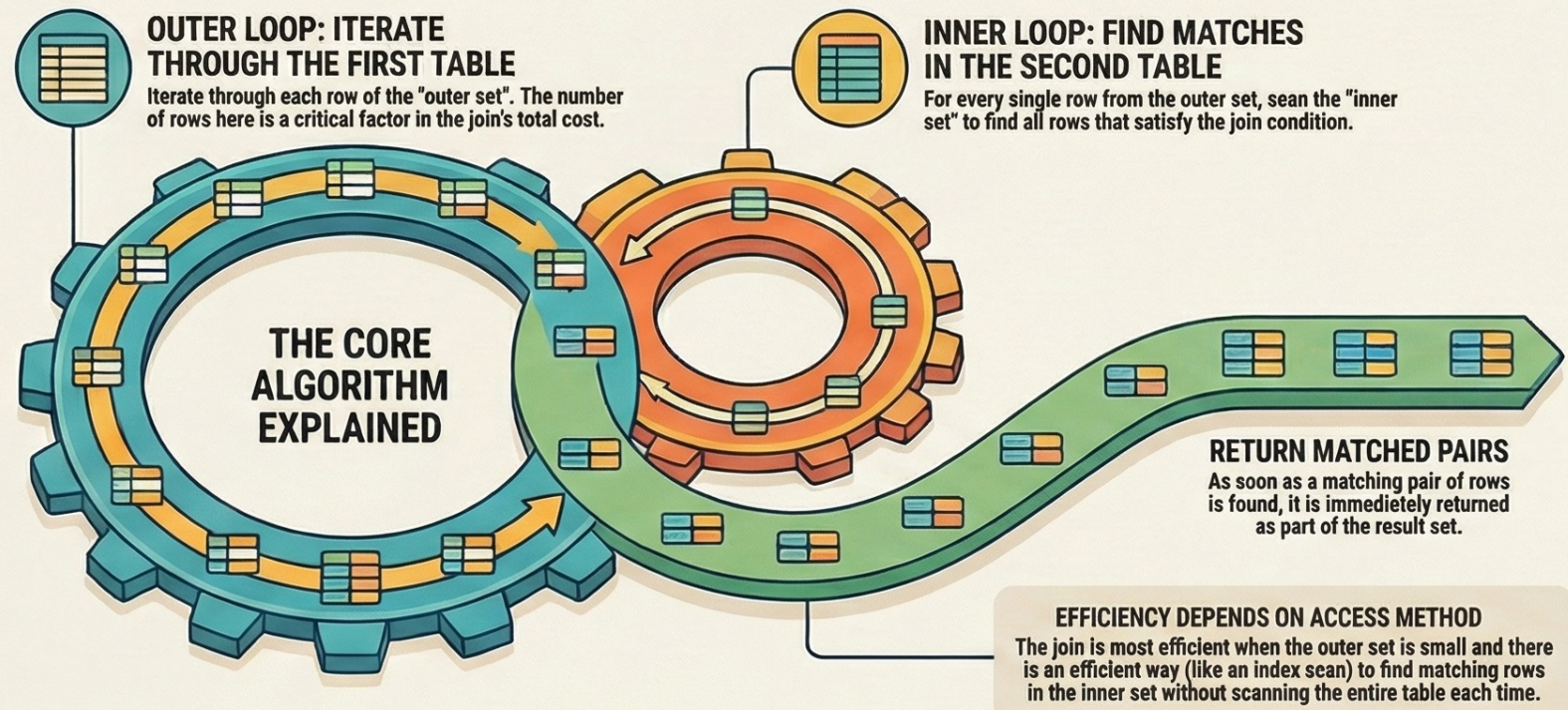
CHOOSING THE RIGHT SCAN FOR THE JOB



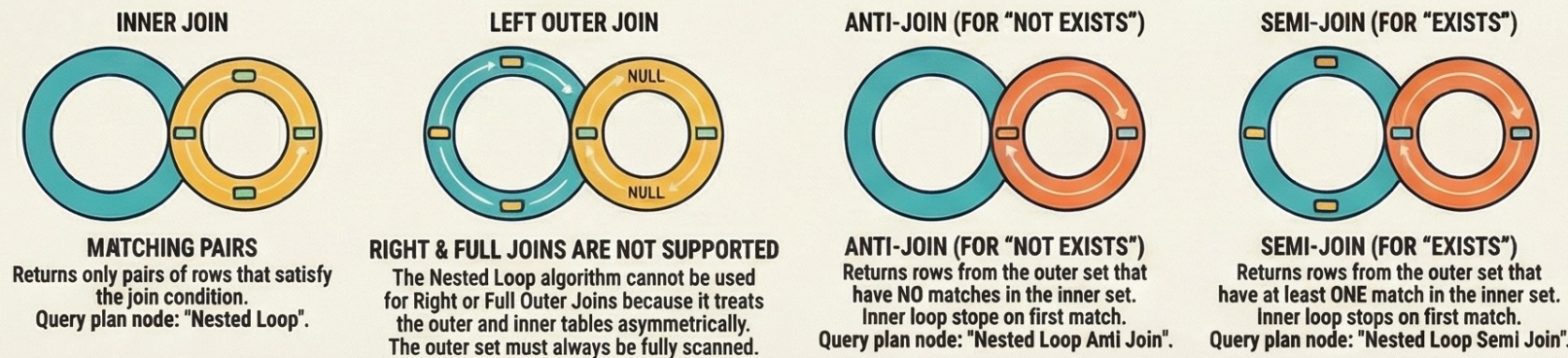
WHEN TO EXPECT EACH SCAN TYPE

Sequential Scan	Large percentage of table read.
Index Scan	Fetching very few rows (high selectivity); degrades with low correlation.
Index-Only Scan	Extremely fast, requires covering index and up-to-date visibility map.
Bitmap Scan	Strong in low correlation: for queries too selective for seq scan but not selective enough for index scan.

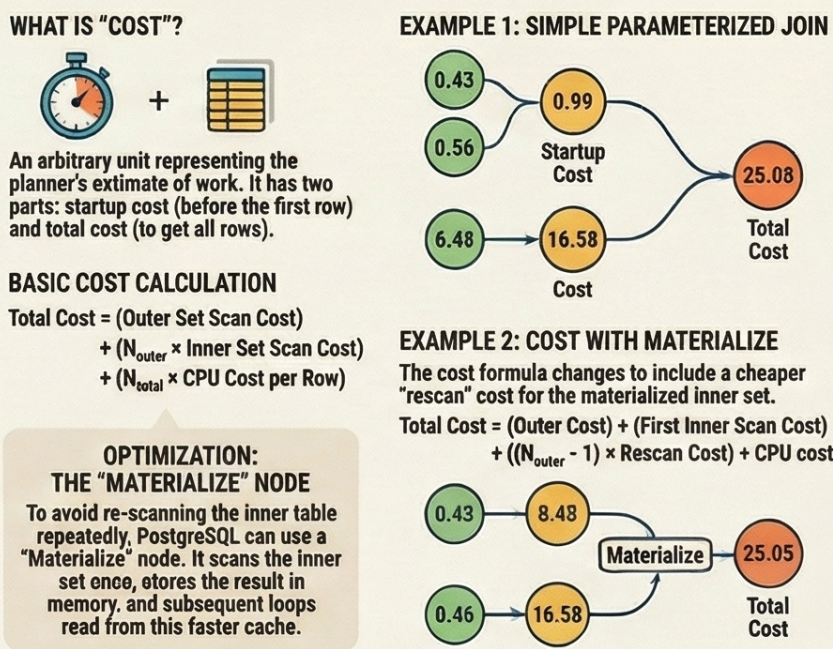
A Visual Guide to PostgreSQL's Nested Loop Join



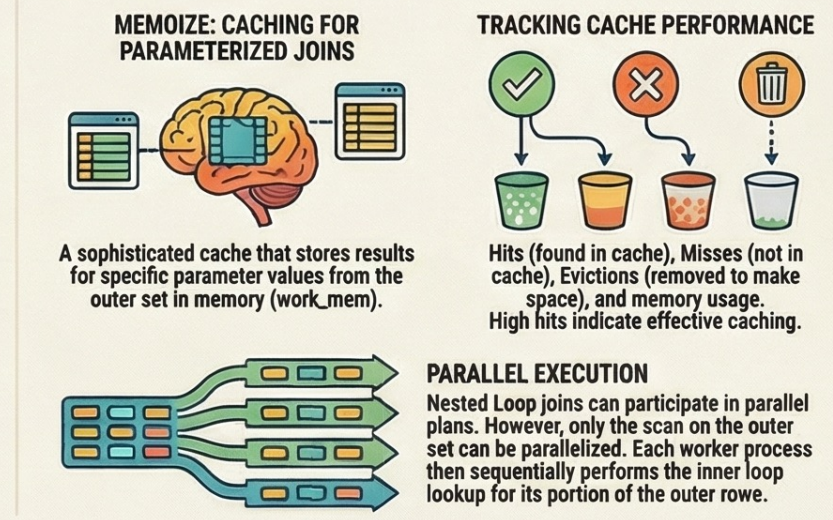
HOW NESTED LOOPS HANDLE DIFFERENT JOIN TYPES



UNDERSTANDING COST ESTIMATION



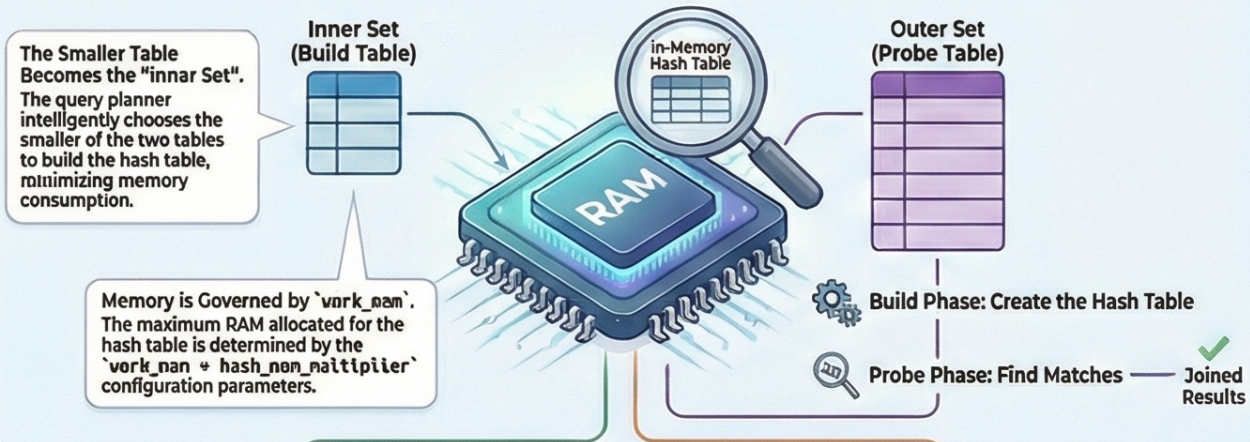
ADVANCED OPTIMIZATIONS



How Database Hash Joins Work: From Single Pass to Parallel Processing

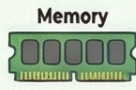
The Anatomy of a Hash Join

A hash join uses an in-memory hash table for fast lookups.



1-Pass Hash Join: The Ideal Scenario

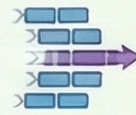
1 Build Phase: Create the Hash Table
The database scans the entire inner set, calculates a hash value from the join key for each row, and stores the rows in the in-memory hash table.



2 Probe Phase: Find Matches
The database scans the outer set row-by-row. For each row, it calculates the hash value of its join key and checks the hash table for matching entries.



3 Return Phase: Output Results
Once a match is found, the combined row is returned. This process continues until the entire outer set has been scanned.



The join is fastest when it completes in a single pass (Batches: 1). This indicates the entire hash table fit within the allocated `work_mem`, avoiding any slow disk I/O.

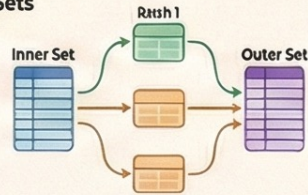


2-Pass Hash Join: When Data Overflows Memory

Problem: The hash table is too large for RAM. If the inner set is larger than the available `work_mem`, the database must spill data to disk.

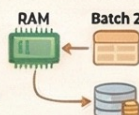
1 Step 1: Partition Both Sets

Rows from both inner and outer sets are divided into multiple smaller "batches" using a hash function. The first batch of the inner set is loaded into a hash table in RAM, while all other batches (from hash tables) are written to temporary files on disk.



2 Step 2: Process Batches Sequentially

After the first batch is processed, the memory is cleared. Then, for each subsequent pair of batch files, the inner batch is loaded into the hash table, and the outer batch is probed against it.

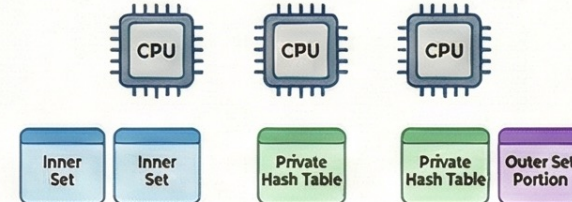


2-pass joins are less efficient due to heavy disk I/O. The `EXPLAIN ANALYZE BUFFERS` output will show significant "temp read" and "temp written" values, indicating performance overhead.



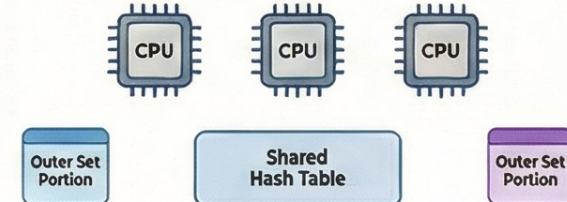
Parallel Hash Joins: Scaling with CPU Cores

Standard Parallel Join: Each worker builds its own hash table.



All parallel workers read the entire inner set to build identical, private hash tables. They then each process a different portion of the outer set. Total memory usage is multiplied by the number of workers.

Shared Parallel Join: Workers collaborate on a single, shared hash table.



All workers build one large hash table in shared memory. This pools their work, increasing the likelihood of a 1-pass join for very large inner sets.

Parallel 3-Pass Join to Handle Massive Data. If even the combined memory is insufficient, a complex 3-pass parallel algorithm partitions data to disk, and workers process batches independently using smaller hash tables in shared memory.

Optimization and Best Practices



Avoid `SELECT *` to reduce hash table size. Only select the columns you need. Forcing columns in the hash table means less memory usage, making a 1-pass join more likely.



Example: Memory Savings
A query using `SELECT *` might require 145MB for its hash table, while the same join selecting only one column might require just 113MB.



Keep Table Statistics Up-to-Date
Outdated statistics can cause the planner to underestimate data size, leading to an initial memory allocation that's too small and forcing a costly dynamic resize of batches during execution.

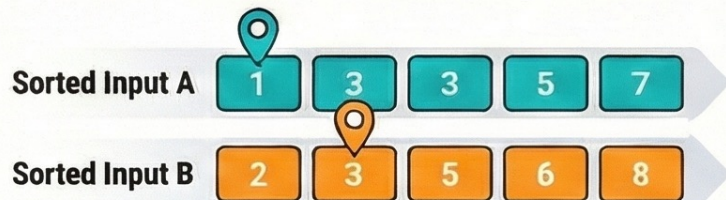


The Planner is Smart About Join Direction. Even if you write a `LEFT JOIN`, the planner may execute it as a `RIGHT JOIN` if it knows the smaller table can be used to build the hash table, ensuring optimal performance.

Inside the Database Engine: A Guide to Sorting & Merging

The Merge Join Explained

Merge Joins process two pre-sorted datasets.



How It Works: The Two-Pointer Scan



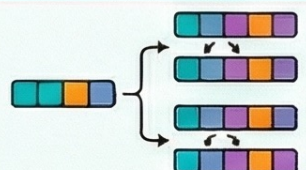
Limited to Equi-Joins:
Only supports equality (=) operators.

Usable in Parallel Queries:
Outer dataset can be scanned in parallel; inner scanned entirely by each worker.

The Sorter's Toolkit: Database Sorting Algorithms

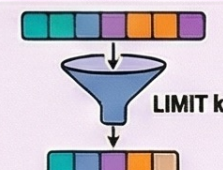
A diverse set of methods for different data sizes and query needs.

Quicksort (In-Memory)



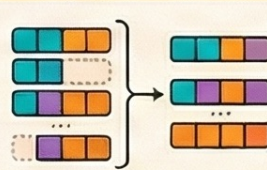
Entire dataset fits in memory (work_mem).
Complexity: $O(n \log_2 n)$.

Top-N Heapsort



Efficiently finds top 'k' items.
Complexity: $O(n \log_2 k)$.

Incremental Sort

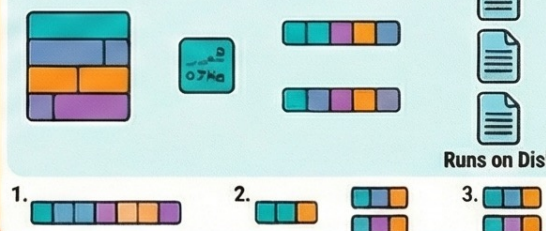


Optimizes partially sorted data (e.g., from index scan), reducing memory use.

External Merge Sort (For Large Data)

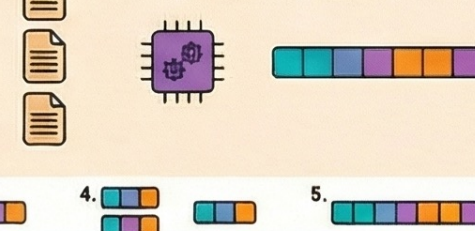
Phase 1: Sort & Write

Reads chunks, sorts in memory, writes sorted runs to temp files.



Phase 2: Merge

Merges sorted runs into final sorted result set. May require multiple passes.



The Join Showdown: A 3-Way Comparison

Nested Loop Join

Best for: Small datasets, OLTP, critical first-row latency.

- ✓ No startup cost
- ✓ Early exit with index
- ✓ Supports ALL join conditions
- ✗ Very slow with large datasets ($O(N*M)$ complexity).

Hash Join

Best for: Large datasets, analytical OLAP queries.

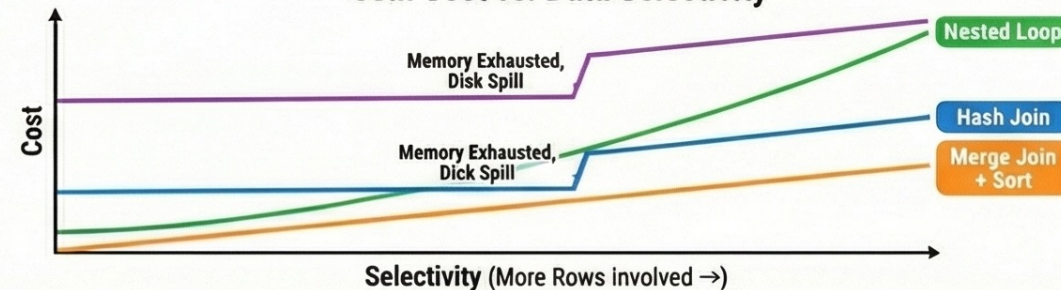
- ✓ Linear complexity ($O(N+M)$)
- ✓ Very efficient if hash table fits in memory
- ✗ High startup cost (builds entire hash table)
- ✗ Only works for equi-joins.

Merge Join

Best for: Versatile for OLTP & OLAP if data is pre-sorted.

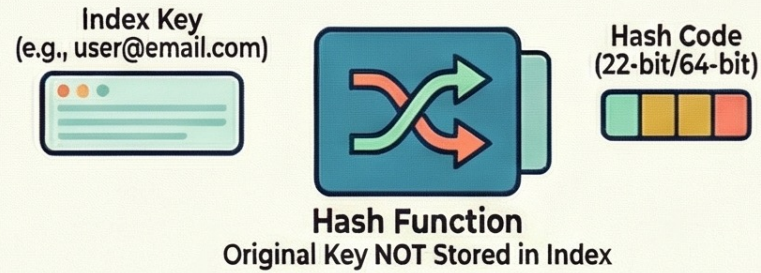
- ✓ Low memory usage
- ✓ No startup delay
- ✓ Linear complexity
- ✗ Requires sorted data. Explicit sort ($O(n \log_2 n)$) often makes Hash Join better.

Join Cost vs. Data Selectivity

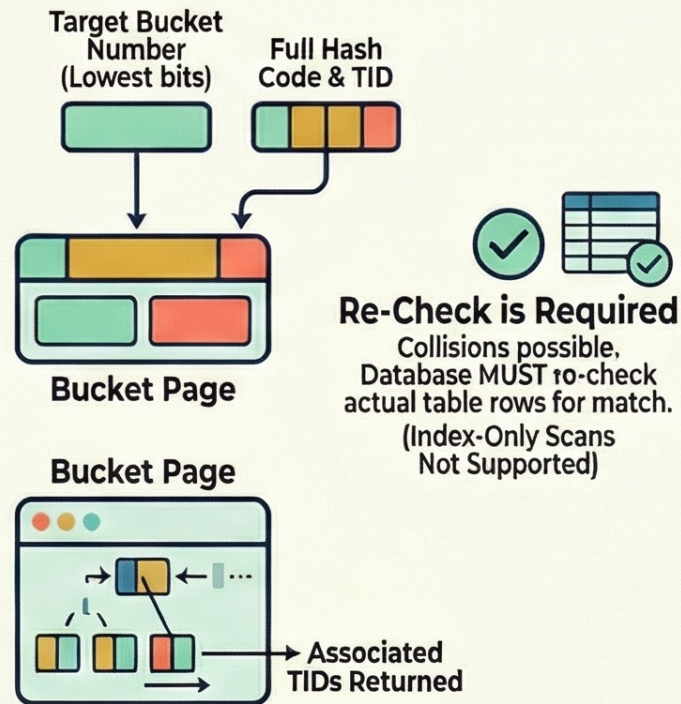
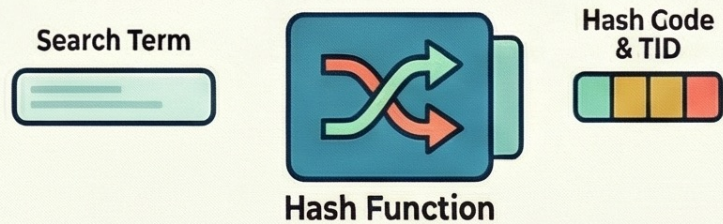


A Deep Dive into PostgreSQL Hash Indexes

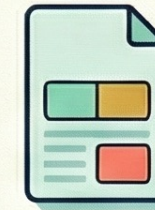
1. Insertion: Storing a New Entry



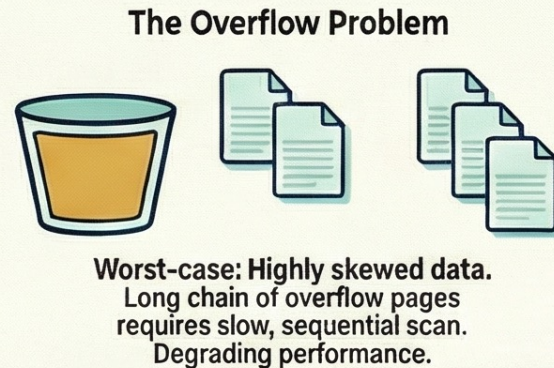
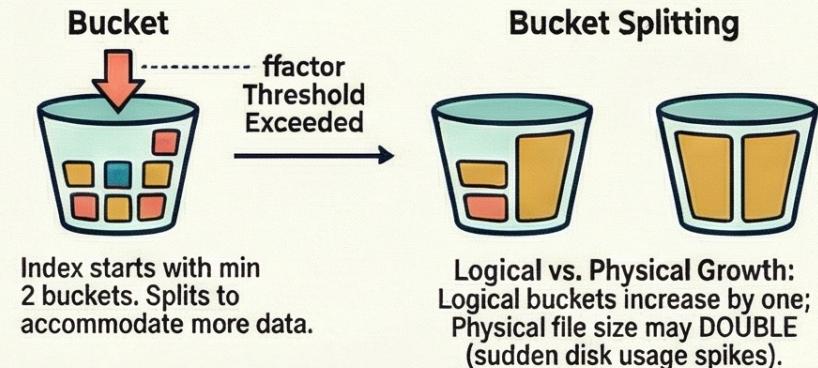
2. Search: Finding a Record



On-Disk Anatomy: The Four Page Types



Dynamic Growth & Performance Pitfalls



Properties & Limitations

Supported Operations

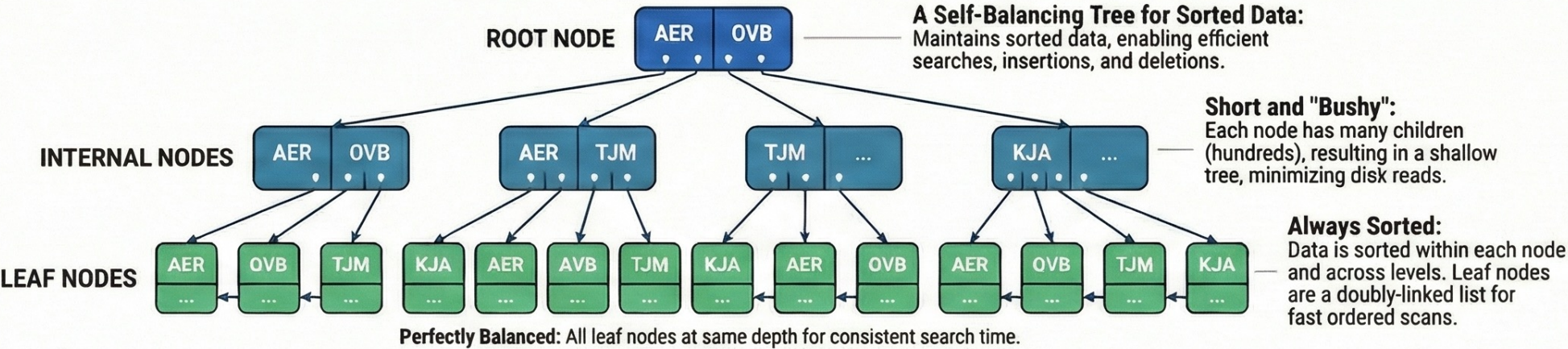
- ✓ Equality Search ('=')
(Primary Use Case)
- ✓ Exclusion Constraints
(Unique-like constraints)
- ✓ Bitmap & Index Scans

Unsupported Operations

- ✗ Ordering ('<', '>')
(No order preservation)
- ✗ Unique Constraints
(Use exclusion instead)
- ✗ Multi-Column Indexes
(Single-column only)
- ✗ Index-Only Scans
(Heap fetch always required)
- ✗ NULLs
(Equality undefined)

A Visual Guide to B-Trees in PostgreSQL

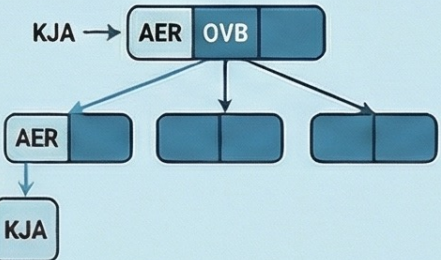
ANATOMY OF A B-TREE



HOW B-TREES HANDLE DATA

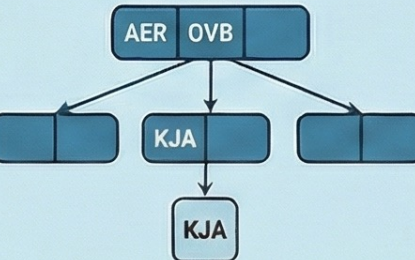
1. Searching: Start at the Root

Compare search value to keys to decide child branch (AER ≤ KJA < OVB).



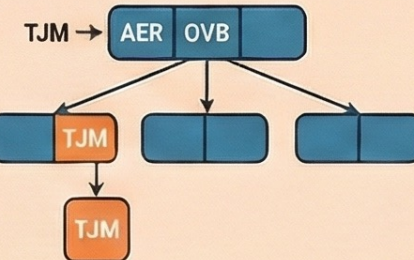
2. Searching: Descend the Tree

Repeatedly narrow search path until a leaf node is reached, fast due to shallow depth.



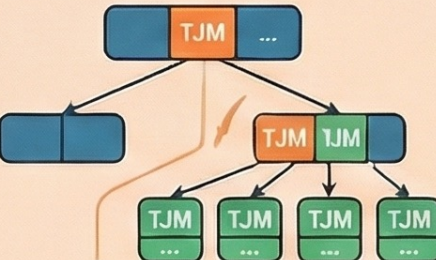
3. Insertion: Find the Spot

Traverse tree to find correct leaf node to maintain sorted order.



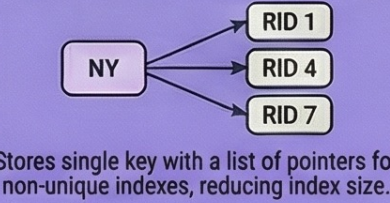
4. Insertion: Split if Full

If target leaf is full, split it and add a reference to the parent. Split can propagate to the root, growing the tree.

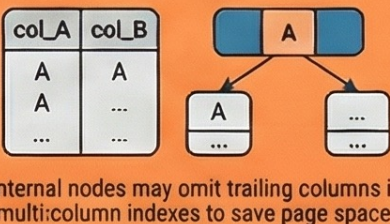


ADVANCED FEATURES IN POSTGRESQL

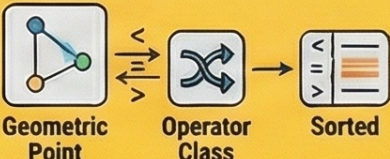
Deduplication for Efficiency



Suffix Truncation



Custom Sorting with Operator Classes



Multi-Column Indexing Order Matters

An index on (col_A, col_B) can be used efficiently for searches on col_A or on (col_A and col_B), but not for searches on col_B alone. Column order is critical for performance.

B-TREE PROPERTIES



Supports Ordering & Uniqueness
Only access method in PostgreSQL to enforce data uniqueness and return sorted data.



Full Scan Capabilities
Supports Index Scans, Bitmap Scans, and Backward Scans (for DESC queries) via leaf node linked list.



Search Flexibility
Supports searching for NULL values and retrieval directly from index (index-only scans).

A Visual Guide to PostgreSQL's GiST Indexes

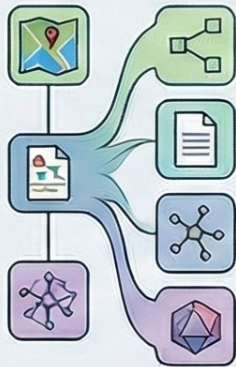
What is a GiST Index?

Standard B-Tree



Optimized for
Ordinated Data
(Numbers, Text)

GiST Index



Balanced, Tree-Structured
Access Method for Complex Data.
Adaptable via Operator Classes.



Operator Class
Defines Core Indexing Logic
for Specific Data Types.



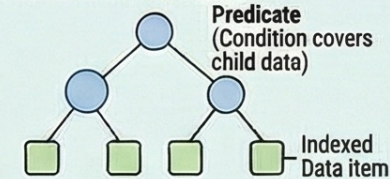
Handles Page Layout,
Locking, WAL
Automatically.

GiST (Generalized Search Tree) Framework: Flexible & Extensible

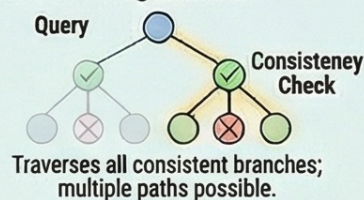


How GiST Works: The Core Mechanics

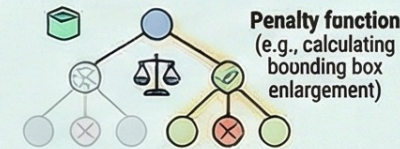
1. Hierarchical Structure



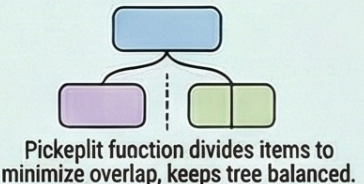
2. Searching the Tree



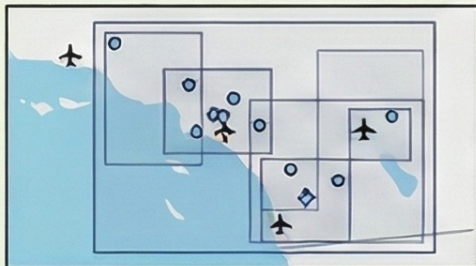
3. Inserting a New Value



4. Splitting a Full Page

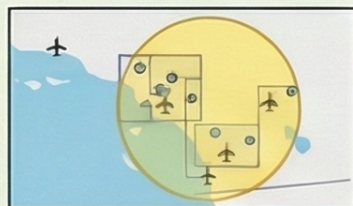


Use Case 1: R-Tree for Spatial Data (e.g., Airport Coordinates)

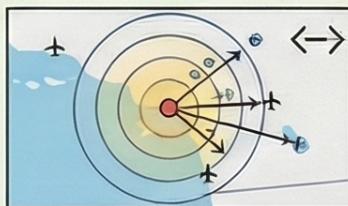


Groups points/
polygons into
Minimum Bounding
Rectangles (MBRs).

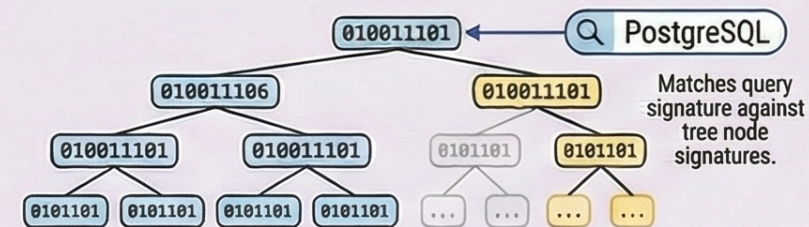
Query: "Contained Within" Search



Query: "k-Nearest Neighbor" (k-NN) Search



Use Case 2: RD-Tree for Full-Text Search



Accuracy vs. Size Trade-off: False positives possible; requires re-check against actual table data. Larger signature reduces false positives, increases index size.

GiST Properties & Other Use Cases



Supports Multi-column
Indexes, Exclusion
Constraints, included
Columns.



Does **NOT** support
Unique Constraints
or Native Ordering.

Other Supported Data Types:



Range
Types



Network
Addresses
(inet)



Integer
Arrays
(intarray)



Key-Value
Stores
(hstore)



Tree-like
Data
(ltree)

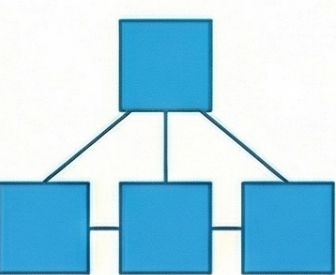


btree_gist
(Combines
B-Tree data)

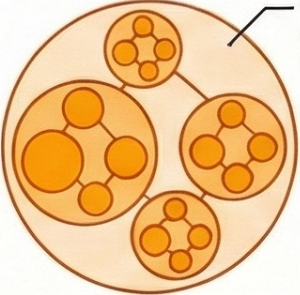
SP-GiST Explained: A Visual Guide to Space-Partitioned Indexing

A framework for creating unbalanced, space-partitioned trees optimized for specialized data types like spatial points and text strings.

What is an SP-GiST Index?



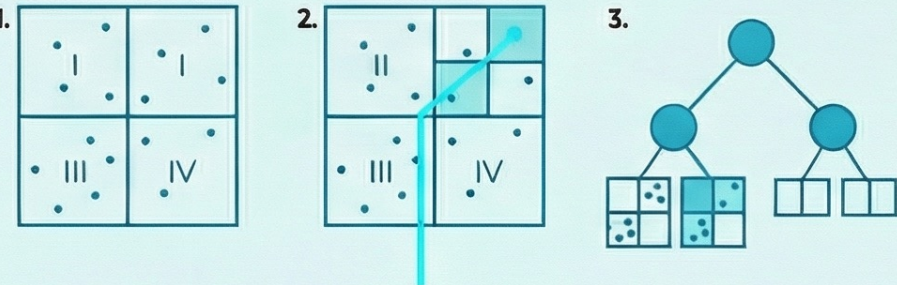
❌ **B-tree/GiST Nodes:** One Node per Page, Balanced



✅ **SP-GiST:** Packed, Small Nodes on Single Page, Unbalanced Structure (Branch depths vary)

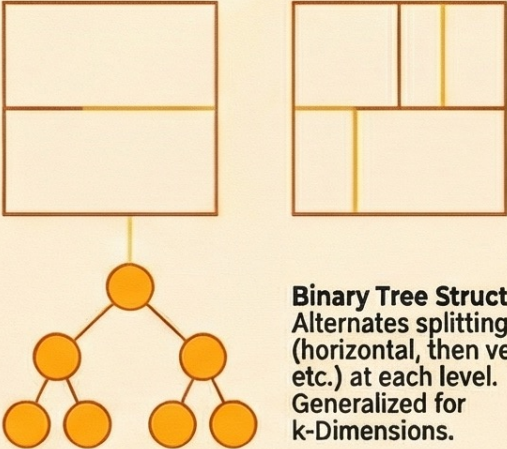
Internal nodes contain a shared condition (prefix); Leaf nodesets & TIDs.

Use Case 1: Quad-tree for 2D Points



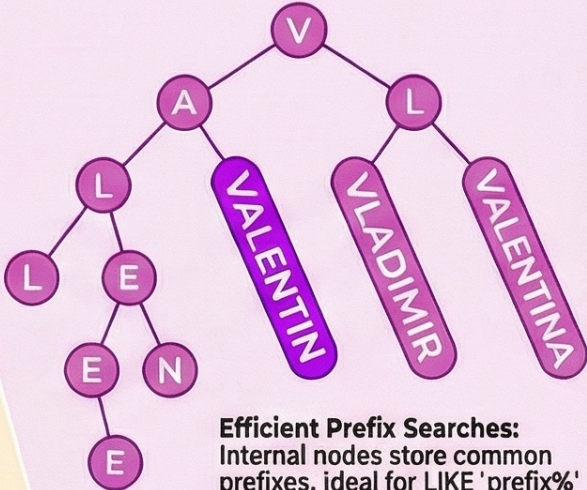
Searching: Prunes irrelevant quadrants using a "consistency function" to find points above (3,7). New points added to corresponding quadrants, triggering "picksplit" when full.

Use Case 2: k-D Tree for Points



Binary Tree Structure: Alternates splitting axis (horizontal, then vertical, etc.) at each level. Generalized for k-Dimensions.

Use Case 3: Radix Tree for Strings



Efficient Prefix Searches: Internal nodes store common prefixes. Ideal for LIKE 'prefix%' or starts_with(). More compact than B-trees, only stores necessary parts of strings.

Key SP-GiST Properties: Supported vs. Not Supported

Feature	Supported by SP-GiST	Description
Exclusion Constraints	Yes (✅)	Can enforce constraints like "no overlapping ranges".
Covering Indexes (INCLUDE)	Yes (✅)	Can include non-key columns for index-only scans.
Bitmap Scans	Yes (✅)	Efficiently combines results of multiple conditions.
NULL Value Indexing	Yes (✅)	NULLs are supported in a separate tree structure.
Ordering Results	No (❌)	Cannot return results in sorted order directly from index.
Unique Constraints	No (❌)	Cannot be used to enforce uniqueness.
Multi-Column Indexes	No (❌)	An SP-GiST index can only be created on a single column.
Clustering	No (❌)	Cannot be used for the CLUSTER command.

Unlocking PostgreSQL: A Deep Dive into GIN Indexes

What is a GIN Index?



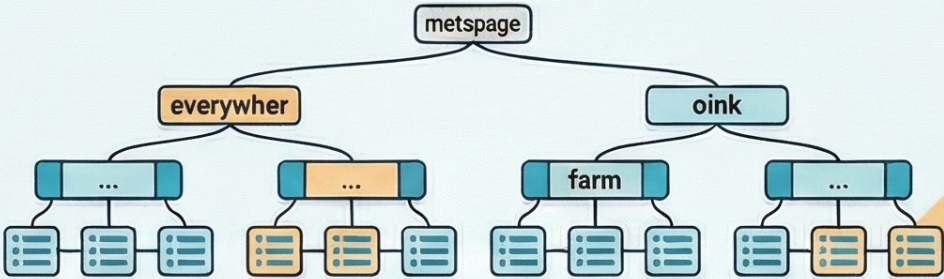
GIN stands for Generalized Inverted Index.

It's designed for composite data types composed of separate elements, like words in a document or items in an array.

It indexes the elements, not the whole value.

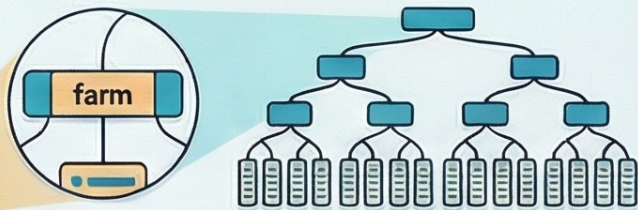
GIN maps each individual element to all the table rows (TIDs) that contain it, similar to a book's index.

The Core Structure is a B-Tree of Elements.



The main data structure is a B-tree where keys are the indexed elements (e.g., words). Each leaf entry in this tree points to a list of TIDs where that element appears.

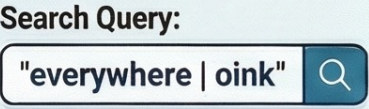
Lists of row IDs are called "Posting Lists".



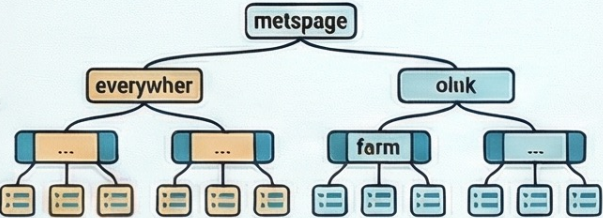
If a posting list for an element becomes too long (i.e., the element is very common), it's stored in a separate B-tree called a "posting tree" for efficiency.

GIN in Action: A Full-Text Search Example

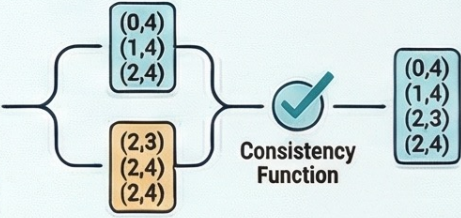
1 Extract Keys from the Query.



2 Find TIDs for Each Key.



3 Merge and Check Results.

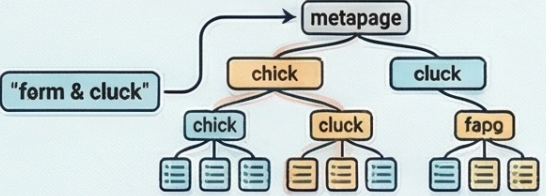


Consistency Check Process for Query: "everywhere | oink"

TID	Contains "everywher"	Contains "oink"	Consistency Function (result)
(0,4)	✓	—	✓
(1,4)	✓	—	✓
(2,3)	—	✓	✓
(2,4)	✓	✓	✓

Performance & Trade-offs

Search is optimized by prioritizing rare terms.



For a query like 'farm & cluck', GIN first finds the few documents with the rare word ('cluck') and only then checks if those specific documents also contain the common word ('farm'), saving significant work.

Querying for a rare term is drastically faster.

Searching for 'wrote' (231,173 docs) took 243ms, while searching for 'wrote & tattoo' (1 doc) took only 8 ms, almost as fast as searching for 'tattoo' alone (2.2me).

GIN updates can be slow.

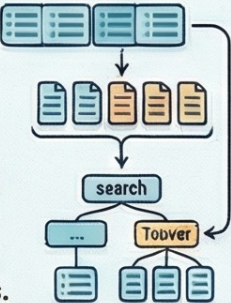


Indexing a single document can require many changes across the index tree, as each word in the document is a separate entry.

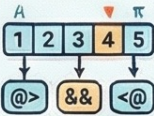
"Fast Update" speeds up writes by delaying them.

During a search, PostgreSQL must scan both the main index tree and the separate, unsorted pending list, potentially reducing read performance until the list is merged.

The pending list makes writes faster but can slow down reads.

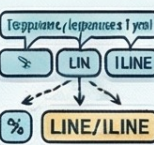


Key Use Cases



Indexing Arrays

Speeds up queries that check for element containment, overlap, or if an array is contained by another.



Fuzzy String Search with Trigrams

The pg_trgm extension allows GIN to index three character segments of text, enabling very text similarity searches and pattern matching.



Indexing JSONB Documents

GIN offers two operator classes for JSONB: jsonb_ops (default) indexes every key and value, while jsonb_path_ops indexes the full path to each value, which is often more efficient.

GIN Index Properties & Limitations

Strengths

- ✓ Supports multi-column indexes.
- ✓ Returns results via BITMAP SCAN.
- ✓ Perfect for checking the existence of elements within composite types.

Weaknesses

- ✗ Cannot enforce UNIQUE constraints.
- ✗ Does not support ordering, so it can't be used to avoid sorting steps.
- ✗ Index only scans are not possible.

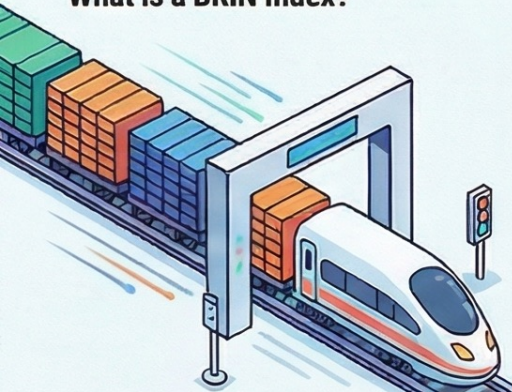
Inefficient with LIMIT clauses.

GIN always builds a full bitmap of all matching rows before fetching any, making it inefficient to look just for the first few results.

Alternative: RUM Index. The RUM extension is based on GIN but adds features GIN lacks, like storing positional information (for phrase searching) and supporting ordering, at the cost of larger index size and slower writes.

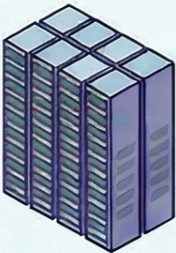
PostgreSQL BRIN Indexes: Small Footprint, Big Performance for Massive Tables

What is a BRIN Index?



A "Coarse" Index for Filtering, Not Finding

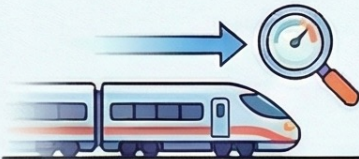
Instead of locating specific rows, BRIN quickly eliminates large chunks of a table that don't match a query's criteria.



Optimized for Massive Tables

Designed for multi-terabyte tables where index size is a primary concern, prioritizing a small footprint over pinpoint accuracy.

An Accelerator for Sequential Scans

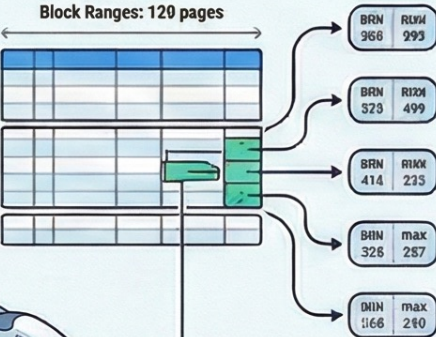


It can be thought of as a smart way to speed up full table scans or as an alternative to table partitioning.

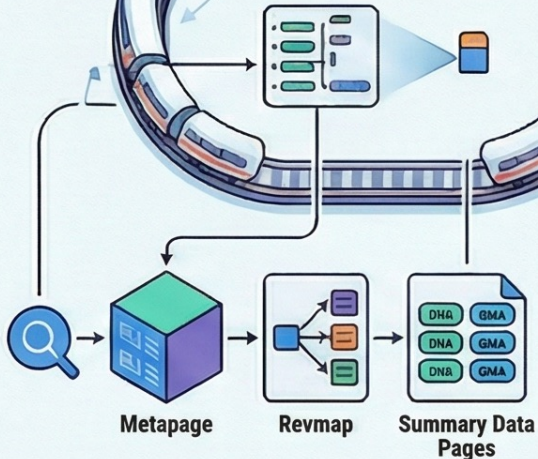
The Anatomy of a BRIN Index

Summarizing Data in Ranges

The index stores a summary for each range, not pointers to individual rows.



A query interacts with three main components: Metapage, Revmap, and Summary Data.



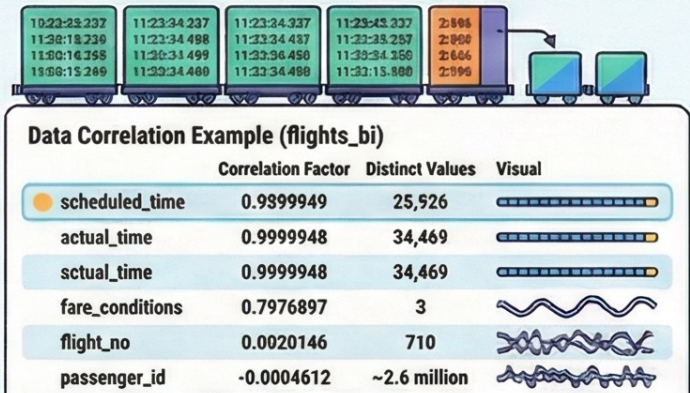
When BRIN Shines: The Power of Correlation

High Correlation is the Key to Success

BRIN is most effective when the physical order of rows on files closely matches the logical order of the indexed values.

Ideal for Append-Only Data

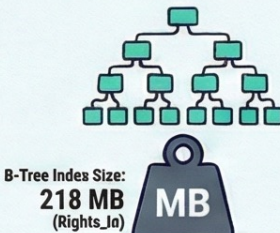
Perfect for tables like event logs or financial transactions where new data is inserted chronologically and old data is rarely updated.



BRIN vs. B-Tree: Size vs. Precision

A Massive Difference in Size

For a 4GB table, a BRIN index can be over 1,600 times smaller than a comparable B-Tree index.



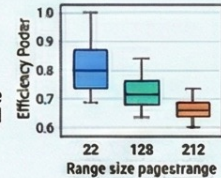
The Trade-Off

B-Tree offers higher precision for finding individual rows, but its also can be a prohibitive cost on very large tables.

Flavors of BRIN: Operator Classes for Every Need

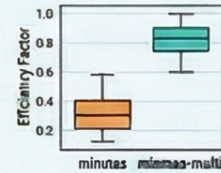
minmax

The Standard for Ordered Data Stores the minimum and maximum value for each range. Best for highly correlated data like timestamps or acrial numbers.



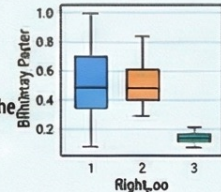
minmax-multi

The Solution for Updated Data Solves the problem of updates breaking correlation by storing multiple value subranges per summary, isolating outliers. Restores efficiency at the cost of a larger index.



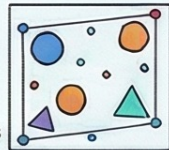
bloom

For Uncorrelated Data with Localized Values. Uses a bloom filter to check for the presence of values within a range. Ideal for columns like product codes that lack natural ordering but appear clustered. Only supports equality (<=) checks.



inclusion

For Geometric and Range Types Stores a "bounding box" that contains all values in a range. Useful for spatial data (point, lines) but may require planner hints as correlation statistics are often unavailable.



Summary of BRIN Properties

What BRIN CAN Do

- ✓ Multi-column Indexes
- ✓ Bitmap Scans
- ✓ NULL value searches

What BRIN CANNOT Do

- ✗ Guarantee Uniqueness
- ✗ Provide Ordered Results (ORDER BY)
- ✗ Be used for Index-Only Scans
- ✗ Support Exclusion Constraints

1. WAL



Write-Ahead Log
Durable & Fast

2. COMMIT

Guaranteed Durabiity
No Data Loss

3. CHECKPOINT

Smooth I/O.
Stable Performance

4. BGWRITER & CHECKPOINTER

Efficient Writes



Simpl & Robust
Eable Recovery

5. LSN



Simple & Robust.
Easy Recovery

THANKS

