

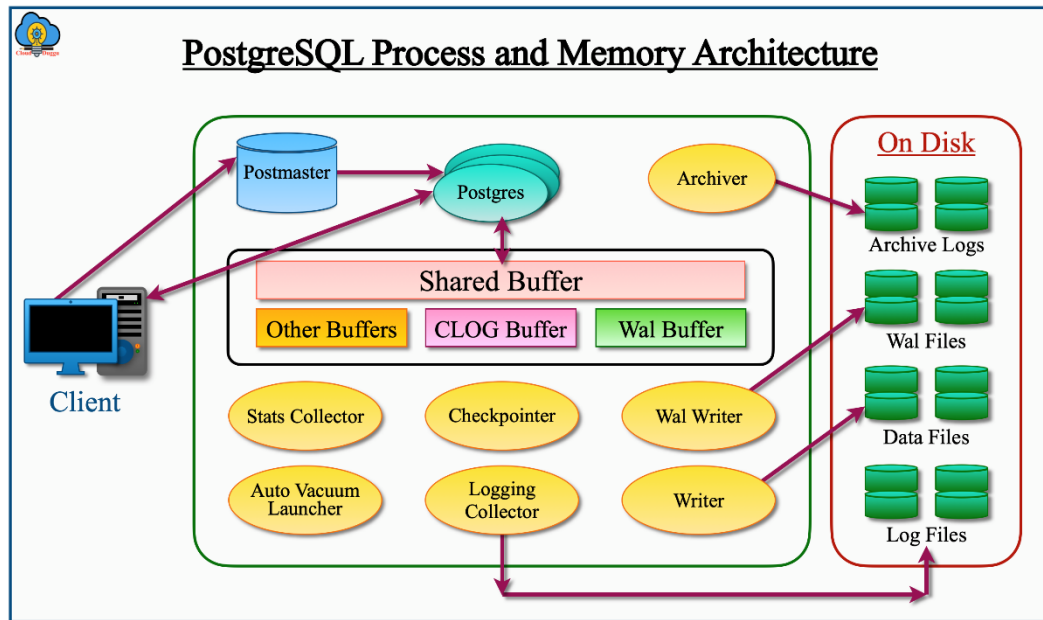


MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

DAY – 2

TOPIC: POSTGRES ARCHITECTURE



MIND MAP

Postgres Internal Architecture

- **Core Concepts & Design Choices**
 - Relational Database System
 - MVCC (multi – version concurrency control)
 - Append Model
 - New Tuple per Update
 - Last Tuple = Active
 - Process-based architecture
 - Not Thresh (Stability Reason)
 - Overhead (Virtual Memory Space)
 - Each process has own page table mapping
 - TLB thrashing (Translation Lookaside buffer)
 - Many process visible (top/Task Manager)
- **Main Processes**
 - Postmaster Process
 - Parent process for everything
 - Listener (port 1532 default)
 - Listens on specific address family/interface
 - Client connect to listener
 - Forks new backend process per client



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- Backed Processes
 - One per client connection
 - Capped by max_connections (default 100)
 - Forked from postmaster (Copy-On-Write optimization)
 - Handles query execution (local or outsource to BG Workers)
 - Work Memory (in backend / background workers)
- **Memory Components**
 - Shared Memory (Shared Buffers / buffer Pool)
 - Stores pages (table data, indexes)
 - Stores WAL records
 - Accessed by all processes
 - Allocated using mmap
 - Required mutexes / semaphores for concurrency
 - Default 128 MB (very low)
- **Background Workers**
 - Pool of processes
 - Outsource parallel plan execution from backend process
 - Scalable and predictable
 - Limited by max_worker_processes
 - Forked from Postmaster
- **Auxiliary Processes**
 - Background writer
 - Flushes dirty pages from shared memory to OS cache
 - Not for durability, but to free up shared buffers
 - OS writes to disk eventually
 - Check pointer
 - Flushes all WAL changes and dirty pages directly to disk
 - Uses O_Direct flag (bypass OS cache)
 - Creates checkpoint record (consistency point)
 - WAL has fixed size
 - Logger
 - Writer logs (standard output / error, csv) to disk
 - Autovacuum (Launcher & Worker)
 - Cleans up old/dead tuples (mvcc bloat)
 - Frees up pages for new data
 - Launcher spins up worker if needed
 - Limited by autovacuum_max_workers
 - WAL Processes
 - WAL Archiver
 - Backs up WAL entries for history. Recovery
 - Compressed changes (DNA of database)
 - Used for recovery and replication
 - WAL Receiver



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- Run on replicas, receives WAL records
- WAL Writer
 - Writes / Flushes WAL records to disk on commit
 - Critical for transaction during (even if pages are dirty)
- WAL senders
 - Master database pushes WAL changes to replicas
 - Configuration with max_wal_senders
- Startup Processes
 - First process to run (before portmaser)
 - Recovers database state after crash
 - Find last checkpoint
 - Redoes (applies) all WAL changes since last checkpoint to pages
 - Brings shared memory pages to consistent state before connections
 - WAL also called Redo Log

PostgreSQL's robust and highly scalable architecture is fundamentally built upon a client-server model, ensuring efficient data management and processing. This design involves a sophisticated interplay of several key components, primarily centered around its process architecture.

1. Process Architecture:

The core of PostgreSQL's operation lies in its carefully orchestrated process architecture, where different components are responsible for specific tasks, ensuring both data integrity and high performance.

Postmaster Daemon:

At the heart of every PostgreSQL instance is the **Postmaster Daemon**, a critical process that springs to life when the PostgreSQL server is initiated. Its role is akin to a central conductor, orchestrating the entire database system.

- **Connection Management:** The Postmaster's primary responsibility is to listen for incoming connection requests from client applications. When a client attempts to connect, the Postmaster receives the request, performs initial authentication, and then facilitates the establishment of a dedicated connection.
- **Recovery Operations:** In the event of an unexpected server shutdown or crash, the Postmaster is crucial for initiating recovery procedures. It ensures that the



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

database is restored to a consistent state, leveraging the Write-Ahead Log (WAL) to replay committed transactions and prevent data loss.

- **Shared Memory Initialization:** The Postmaster is responsible for initializing and managing the shared memory segments that are utilized by all PostgreSQL processes. This shared memory is vital for inter-process communication, caching data blocks, and storing system-wide information.
- **Background Process Management:** The Postmaster is also tasked with launching and overseeing the various background processes that perform essential maintenance and auxiliary tasks, ensuring the smooth operation of the database system.

Backend Processes (Postgres):

For every successful client connection, the Postmaster Daemon "forks" or spawns a dedicated **Backend Process**, often referred to simply as "Postgres." This design provides significant benefits for concurrency and isolation.

- **Client-Specific Handling:** Each backend process is exclusively responsible for handling the queries and interactions from its specific client connection. This isolation prevents contention between different clients and ensures that one client's operations do not negatively impact others.
- **Direct Communication:** After the initial connection and authentication phase handled by the Postmaster, the backend process communicates directly with its associated client. This direct link optimizes query execution and data transfer.
- **Query Execution:** The backend process parses and executes SQL queries received from the client, interacts with the data files, performs necessary computations, and returns the results to the client.
- **Transaction Management:** Each backend process manages the transactions initiated by its client, ensuring atomicity, consistency, isolation, and durability (ACID properties) for all database operations.

Background Processes:

Beyond the Postmaster and backend processes, PostgreSQL leverages a suite of specialized **Background Processes** that continuously perform crucial auxiliary and maintenance tasks. These processes run independently, contributing significantly to the database's reliability, performance, and recoverability.

- **Background Writer:** This process is responsible for periodically writing "dirty" data blocks (modified data that resides in the shared buffers) from memory to disk. By asynchronously writing these blocks, the Background Writer reduces the



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

burden on backend processes during COMMIT operations, improving overall transaction throughput.

- **Checkpoint:** The Checkpointer is a critical process for ensuring data consistency and facilitating fast crash recovery. Its primary role is to ensure that all modified pages in shared memory are written to disk at regular intervals (checkpoints). It also manages the lifecycle of Write-Ahead Log (WAL) files, signaling when they can be recycled. Checkpoints mark a point in time where all data up to that point is guaranteed to be on disk, minimizing the amount of WAL that needs to be replayed during recovery.
- **WAL Writer:** The **Write-Ahead Log (WAL)** is fundamental to PostgreSQL's durability and crash recovery. The WAL Writer is solely responsible for writing the contents of the WAL buffer (in-memory log records of changes) to the WAL files on disk. This ensures that every change is durably recorded before the actual data pages are written, guaranteeing that committed transactions are not lost even in the event of a system crash.
- **Autovacuum Daemon:** This essential process manages cleanup and optimization tasks automatically, preventing database "bloat" and improving query performance. It identifies and removes "dead" tuples (rows marked for deletion but not yet physically removed) from tables, reclaims disk space, and updates statistics (analyzes tables) to provide the query planner with accurate information for optimal query execution plans. Without the Autovacuum Daemon, database performance would degrade over time due to accumulating dead tuples and outdated statistics.
- **Logger:** The Logger process is responsible for handling and writing log messages generated by various server processes to the configured log files. These logs are invaluable for monitoring database activity, diagnosing issues, and auditing operations.
- **Archiver:** The Archiver process plays a vital role in PostgreSQL's backup and recovery strategy. It is responsible for copying completed WAL segments (files containing log records) to a designated archive storage location. This continuous archiving of WAL files is crucial for Point-in-Time Recovery (PITR), allowing administrators to restore the database to any specific moment in time.
- **Stats Collector:** The Stats Collector gathers and maintains statistical data about database activity, such as table and index usage, access patterns, and function execution counts. This statistical information is heavily utilized by the query planner to make intelligent decisions about how to execute queries efficiently, leading to optimized query performance.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Logical Replication Launcher:** Introduced in more recent versions, the Logical Replication Launcher manages the lifecycle of logical replication worker processes. Logical replication allows for selective replication of data changes (inserts, updates, deletes) to other databases, providing a flexible and powerful mechanism for data distribution and integration.
- **Background Workers:** Introduced in PostgreSQL 9.6, Background Workers provide a powerful extension mechanism, allowing developers to write custom C code that runs directly within the PostgreSQL server process space. These processes can execute queries in parallel, perform long-running tasks, or interact with external systems, significantly enhancing performance for certain operations and enabling advanced functionalities. This capability has paved the way for features like parallel query execution in later PostgreSQL versions.

PostgreSQL Architecture: A Detailed Overview

At its core, PostgreSQL operates through a series of interconnected processes and shared memory structures, all orchestrated to handle concurrent client requests and maintain database consistency.

1. Shared Memory

Shared memory is a critical component of PostgreSQL's architecture, serving as a high-speed caching mechanism to minimize disk I/O and facilitate efficient communication between various processes. The two most vital elements residing in shared memory are the **Shared Buffer** and the **WAL Buffer**.

a. Shared Buffer:

The **Shared Buffer** is the largest and most crucial part of shared memory. Its primary purpose is to **minimize disk I/O** by caching frequently accessed data blocks in memory. To achieve this, it adheres to several key principles:

- **Fast Access to Large Buffers:** The shared buffer pool is designed to provide very fast access to data blocks, even when dealing with extremely large datasets (tens to hundreds of gigabytes). This speed is paramount for quick query responses.
- **Minimizing Contention:** When numerous users concurrently access the database, the shared buffer must minimize contention. PostgreSQL employs



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

sophisticated locking mechanisms and concurrency control protocols to ensure that multiple processes can access and modify different parts of the buffer pool simultaneously without corrupting data.

- **Maximizing Resident Time for Hot Blocks:** Frequently used data blocks (often referred to as "hot" blocks) are kept in the buffer for as long as possible. This is achieved through algorithms that prioritize keeping actively used blocks in memory, reducing the need to fetch them repeatedly from disk.

b. WAL Buffer (Write-Ahead Log Buffer):

The **WAL Buffer** is a temporary in-memory staging area for database changes. Before any modification is written to the actual data files on disk, a record of that change is first written to the WAL buffer. This adherence to the **Write-Ahead Logging (WAL)** principle is fundamental to PostgreSQL's durability and recovery capabilities.

- **Durability:** The WAL buffer ensures that every committed transaction is recorded to disk (via the WAL files) before the actual data pages are updated. This guarantees that even in the event of a system crash, committed transactions can be recovered by replaying the WAL files.
- **Backup and Recovery:** From a backup and recovery perspective, the WAL buffers and the subsequent WAL files are immensely important. They enable Point-in-Time Recovery (PITR), allowing administrators to restore the database to any specific point in time by restoring a base backup and then replaying the necessary WAL segments.
- **Flush Mechanism:** The contents of the WAL buffer are periodically flushed (written) to the durable WAL files on disk at predetermined points, such as transaction commits, full WAL buffers, or a configurable timeout.

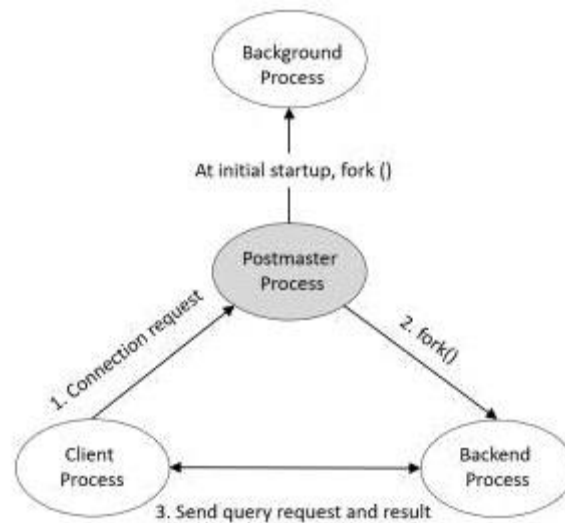
2. PostgreSQL Process Types

PostgreSQL's operational framework is built upon four distinct process types, each with specific responsibilities:

- **Postmaster (Daemon) Process:** The overarching parent process.
- **Background Process:** A set of specialized processes performing maintenance and auxiliary tasks.
- **Backend Process:** Dedicated processes handling individual client connections.
- **Client Process:** The external application connecting to the database.



3. Postmaster Process



The **Postmaster Process** (often referred to as the Postmaster Daemon) is the initial process launched when PostgreSQL starts. It acts as the central coordinator and supervisor for the entire database cluster.

- **Initialization:** Upon startup, the Postmaster performs crucial initialization tasks, including:
 - **Recovery:** If the server was not shut down cleanly, the Postmaster initiates crash recovery, utilizing WAL files to restore the database to a consistent state.
 - **Shared Memory Initialization:** It initializes and allocates the shared memory segments (like Shared Buffer and WAL Buffer) that will be used by all other PostgreSQL processes.
 - **Background Process Launch:** It launches and monitors all necessary background processes required for PostgreSQL's operation.
- **Connection Management:** The Postmaster continuously listens for incoming connection requests from client processes. When a request arrives, it handles the initial authentication and then **forks a new Backend Process** specifically for that client connection.
- **Parent Process:** As shown in the "Process relationship diagram" (Figure 1-2), the Postmaster process is the **parent process** of all other PostgreSQL processes (background and backend processes). This hierarchical structure allows the Postmaster to monitor and manage its child processes. If a child process crashes, the Postmaster can detect it and take appropriate action.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

```
$ pstree -p 1125
postgres(1125) /usr/local/pgsql/bin/postgres -D /usr/local/pgsql/data
├─ postgres(1249) postgres: logger process
├─ postgres(1478) postgres: checkpoint process
├─ postgres(1479) postgres: writer process
├─ postgres(1480) postgres: wal writer process
├─ postgres(1481) postgres: autovacuum launcher process
├─ postgres(1482) postgres: archiver process
└─ postgres(1483) postgres: stats collector process
```

4. Background Processes

PostgreSQL relies on a suite of **Background Processes** that run continuously to perform essential maintenance, logging, and statistical tasks. These processes are vital for the database's health, performance, and recoverability. Table 1-1 details their specific roles:

Process	Role
Logger	Writes error messages, warnings, and informational messages from various server processes to the configured log files. Essential for monitoring and troubleshooting.
Checkpoint	Ensures data consistency and facilitates crash recovery. When a checkpoint occurs, it writes all "dirty" buffers (modified data in shared memory) to disk and updates the control file.
Writer	(Often referred to as Background Writer) Periodically writes "dirty" buffers from shared memory to disk, reducing the I/O burden on backend processes during transactions.
WAL Writer	Writes the contents of the WAL buffer to the WAL files on disk, ensuring that all committed transactions are durably recorded before data pages are updated.
Autovacuum launcher	When the autovacuum feature is enabled (which is highly recommended), this process forks autovacuum worker processes as needed. It's responsible for initiating vacuum operations on bloated tables.
Archiver	When in archive_mode (for Point-in-Time Recovery), this process copies completed WAL files to a specified archive directory, enabling continuous archiving of transaction logs.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Process	Role
Stats Collector	Gathers and maintains various DBMS usage statistics, such as active session information (pg_stat_activity) and detailed table usage statistics (pg_stat_all_tables). This data is crucial for the query planner's optimization decisions.

5. Backend Process

A **Backend Process** is a dedicated child process forked by the Postmaster for each client connection. Its sole purpose is to serve the requests of that specific client.

- **Maximum Connections:** The maximum number of concurrent backend processes is controlled by the max_connections parameter (default: 100). This parameter effectively limits the number of simultaneous client connections.
- **Query Execution:** The backend process receives SQL queries from its associated client, parses them, plans their execution, retrieves or modifies data, and then transmits the results back to the client.
- **Local Memory:** For efficient query execution, each backend process allocates its own **local memory** structures. These are distinct from the shared memory and are only accessible by that specific backend process. Key parameters associated with local memory include:
 - **work_mem:** Defines the maximum memory space that can be used by internal sort operations, bitmap operations, hash joins, and merge joins before spilling data to temporary disk files. The default setting is 4 MB per operation, per backend. Increasing this can significantly improve performance for complex queries.
 - **maintenance_work_mem:** Specifies the maximum memory to be used by maintenance operations like VACUUM, CREATE INDEX, and ALTER TABLE ADD FOREIGN KEY. The default is 64 MB. A larger value can speed up these operations.
 - **temp_buffers:** Sets the maximum memory used for temporary tables and temporary disk files created during query execution. The default is 8 MB.

6. Client Process

The **Client Process** refers to the external application or tool (e.g., psql command-line tool, a custom application using libpq, or a GUI client) that initiates a connection to the PostgreSQL server.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Connection Request:** The client process sends a connection request to the Postmaster process.
- **Dedicated Backend:** Upon successful connection and authentication, a dedicated Backend Process is assigned to serve this client. All subsequent communication and query execution occur directly between the client process and its assigned backend process.

PostgreSQL's architecture, beyond its process and memory management, is deeply rooted in how it structures and stores data on disk. Understanding its database structure, tablespaces, and file organization is crucial for effective administration, performance tuning, and backup strategies.

Database Structure: A Deep Dive into PostgreSQL's On-Disk Organization

PostgreSQL organizes data hierarchically, starting from a "database cluster" down to individual table files, all managed with a focus on efficiency and integrity.

1. Items Related to the Database

PostgreSQL's fundamental organizational unit for user data is the **database**. A collection of these databases within a single PostgreSQL server instance is referred to as a **database cluster**.

- **Initial Databases upon initdb():** When you initialize a PostgreSQL data directory using the `initdb()` command, three foundational databases are automatically created:
 - **template0:** This is the pristine, immutable template database. It's guaranteed to be in its initial state, containing only the essential system catalog tables. Its `datallowconn` column is typically set to `f` (false), meaning no regular connections are allowed to it, preventing accidental modifications. This immutability makes it ideal for creating new databases that require a completely clean, default setup.
 - **template1:** This is the primary template database from which all user-created databases are cloned by default. Unlike `template0`, `template1` *can* be modified by users. This means you can add objects (tables, functions, extensions, etc.) to `template1`, and these objects will automatically be present in any new database created from `template1`. This is useful for establishing a standard set of objects across multiple user databases.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

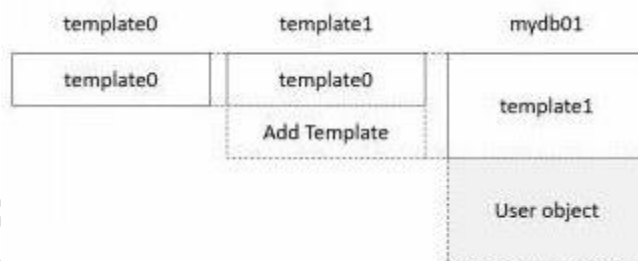
BY: Mukesh Chaurasia

- **postgres:** This is the default user database. If a client connects to the PostgreSQL server without specifying a database name, they will be connected to the postgres database. It is initially created as a clone of template1.
- **User Database Creation:** When a user creates a new database (e.g., CREATE DATABASE my_new_db;), PostgreSQL performs a **cloning operation** from template1 (unless template0 is explicitly specified). This means the new database starts with all the system catalogs and any user-defined objects that existed in template1 at the time of cloning. This cloning mechanism is efficient as it essentially creates a copy of the underlying file structure.

```
-- Create a T1 table in the template database.
template1=# create table t1 (c1 integer);

-- Create the mydb01 database.
postgres=# create database mydb01;

-- After connecting to mydb01 database, check if T1 table
exists. mydb01=# \d t1
      Table "public.t1"
  Column | Type   | Modifiers
-----+-----+-----
 c1      | integer |
```



- **Physical Database Location:** Each database within a PostgreSQL cluster is physically stored as a subdirectory under the \$PGDATA/base directory. The name of each subdirectory corresponds to the **Object Identifier (OID)** of that specific database. For example, if a database has an OID of 16384, its data files will reside in \$PGDATA/base/16384/. This OID-based naming scheme provides a unique identifier for each database's physical storage.

2. Items Related to Tablespaces

Tablespaces in PostgreSQL provide a powerful mechanism to control the physical storage location of database objects (tables, indexes, etc.). This allows administrators



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

to spread data across different storage devices for performance, capacity, or organizational reasons.

- **Initial Tablespaces upon initdb():** Immediately after initdb(), two default tablespaces are created:
 - **pg_default:** This is the default tablespace for all regular user tables and indexes unless explicitly specified otherwise during object creation. Its physical location is implicitly within the \$PGDATA/base directory structure. Specifically, for each database, there's a subdirectory inside \$PGDATA/base (named by the database's OID), and the pg_default tablespace for that database resides within that OID directory.
 - **pg_global:** This tablespace is reserved for storing objects that are managed at the **database cluster level** (i.e., they are globally accessible and consistent across all databases within the cluster). Examples include system catalog tables like pg_database or pg_authid (user roles). The physical location of pg_global is \$PGDATA/global.
- **Tablespace Usage and Subdirectories:**
 - A single tablespace can indeed be utilized by **multiple databases**. To maintain organization and prevent name collisions, when multiple databases use the same tablespace, a **database-specific subdirectory** is created within that tablespace's physical directory. This subdirectory is named after the OID of the database using it.
 - For pg_default, as mentioned, its logical location is \$PGDATA/base, and then within that, base/<database_OID> holds the actual table files.
 - For pg_global, its physical location is \$PGDATA/global.

```
postgres=# select oid, * from pg_tablespace;
```

oid	spcname	spcowner	spcACL	spcoptions
1663	pg_default	10		
1664	pg_global	10		

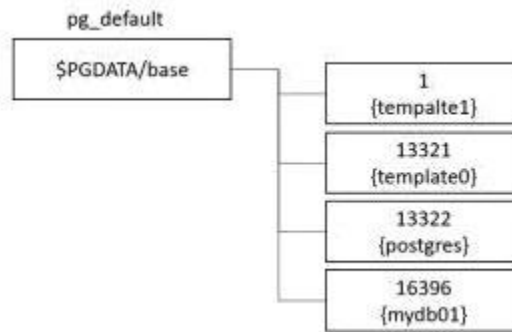
```
[postgres@pgserver ~]$ ls -l $PGDATA/base
```

```
drwx----- .2 postgres postgres 8192 Nov 4 20:06 1
drwx----- .2 postgres postgres 8192 Nov 4 20:02 13321
drwx----- .2 postgres postgres 8192 Nov 4 20:02 13322
drwx----- .2 postgres postgres 8192 Nov 4 20:06 16396
```

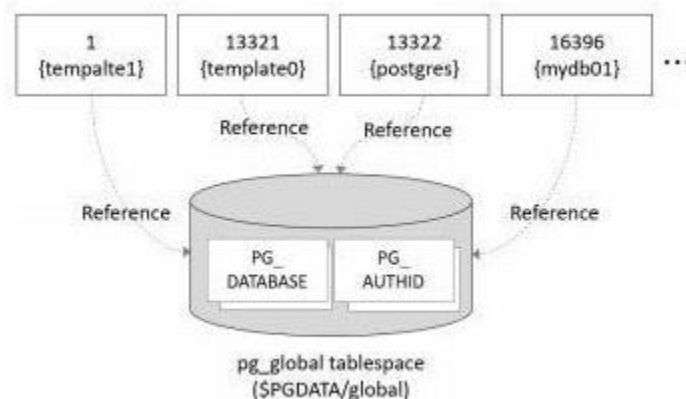


MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia



- **Creating User Tablespaces:** When a user creates a new tablespace (e.g., `CREATE TABLESPACE my_ts LOCATION '/data/my_tablespace';`), PostgreSQL does not directly store data in the specified directory. Instead, it creates a **symbolic link** in the `$PGDATA/pg_tblspc` directory that points to the user-defined physical location. The symbolic link's name in `$PGDATA/pg_tblspc` is the OID of the newly created tablespace. This allows PostgreSQL to manage tablespaces from a central location while providing flexibility in their physical placement.



3. Items Related to the Table (File Organization)

PostgreSQL stores table data and associated metadata in a specific file structure on disk. For each user table, there are typically three associated files:

- **Table Data File:** This file stores the actual row data of the table. Its filename is the **Object Identifier (OID)** of the table. For example, if a table has an OID of 12345, its data will be in a file named 12345.
- **Free Space Map (FSM) File:** This file manages the free space within the table's blocks. It keeps track of how much free space is available in each block, which helps in quickly locating blocks with enough space for new rows during inserts. Its filename is `OID_fsm`.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Visibility Map (VM) File:** This file is used to manage the visibility of data blocks within the table, particularly in the context of PostgreSQL's Multi-Version Concurrency Control (MVCC). It helps in quickly identifying blocks where all tuples are visible to all active transactions, allowing VACUUM operations to skip scanning those blocks. Its filename is `OID_vm`.
- **Index Files:** Unlike tables, indexes typically only have two associated files:
 - **Index Data File:** Stores the actual index structure (e.g., B-tree nodes). Its filename is the OID of the index.
 - **Free Space Map (FSM) File:** Manages free space within the index blocks, similar to table FSM files. Its filename is `OID_fsm`. Indexes do not have a `_vm` file because MVCC visibility is managed at the table level, not the index level.

4. Other Things to Remember... (`relfilenode` and `pg_relation_filepath`)

The relationship between an object's OID and its physical file name can be dynamic.

- **OID vs. relfilenode:** While the initial filename of a table or index is its OID, PostgreSQL uses an internal property called `relfilenode` to refer to the actual underlying file.
- **Rewrite Operations and File Name Changes:** When certain **rewrite operations** are performed on an object, its `relfilenode` value changes, and consequently, the physical filename also changes to the new `relfilenode` value. This is done to avoid contention and simplify cleanup of old versions of the data. Common operations that cause a `relfilenode` change include:
 - **TRUNCATE:** Removes all rows from a table and typically resets its `relfilenode`.
 - **CLUSTER:** Rewrites a table sorted by an index.
 - **VACUUM FULL:** Rewrites a table to reclaim disk space.
 - **REINDEX:** Rebuilds an index.
 - **ALTER TABLE** operations that fundamentally change the table's structure (e.g., adding a column with a default value that requires a table rewrite).
- **`pg_relation_filepath()` Function:** To easily determine the current physical location and filename of any table or index, you can use the `pg_relation_filepath('<object_name>')` function. This is extremely useful for troubleshooting and understanding the on-disk layout.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Running Tests and Verification

```
select oid, datname, datistemplate, datallowconn
from pg_database order by 1;
oid | datname | datistemplate | datallowconn
-----+-----+-----+-----
1 | template1 | t | t
13321 | template0 | t | f
13322 | postgres | f | t
```

The text provides practical examples to verify the concepts discussed:

- **pg_database View:** Querying the pg_database system catalog view after initdb() confirms the creation of template0, template1, and postgres databases. The datistemplate column clearly identifies template0 and template1 as template databases. The datallowconn column for template0 being f (false) highlights its immutable nature.
- **template0 vs. template1 Purpose:** The distinction between template0 (pristine, immutable) and template1 (customizable for user-defined templates) is crucial for database administrators.
- **postgres Database as Default:** The postgres database serves as the default connection target if no other database is specified.
- **Database OID as Directory Name:** The observation that databases are located under \$PGDATA/base with subdirectory names corresponding to their OIDs is a key physical layout detail.
- **User Database Creation (Cloning template1):** The example demonstrating the creation of a user table T1 in template1 and then verifying its existence in a newly created mydb01 database (cloned from template1) provides concrete proof of the cloning mechanism.
- **pg_tablespace View:** Querying pg_tablespace confirms the existence of pg_default and pg_global tablespaces.
- **pg_default Physical Location:** Figure 1-4 visually represents how pg_default is logically within \$PGDATA/base, with specific database OID subdirectories containing the actual data files.
- **pg_global Tablespace Role:** Figure 1-5 illustrates that pg_global stores cluster-level metadata, accessible uniformly across all databases, and is physically located at \$PGDATA/global.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

```
[postgres@pgserver data]$ ls -l $PGDATA/base
drwx----- . 2 postgres postgres 8192 Nov 4 19:34 1
drwx----- . 2 postgres postgres 8192 Nov 4 19:34 13321
drwx----- . 2 postgres postgres 8192 Nov 4 19:34 13322
```

Creating and Managing User Tablespaces in PostgreSQL

The ability to create and manage user-defined tablespaces allows administrators to precisely control where tables, indexes, and other database objects reside on the file system. This section delves into the details of creating tablespaces, their physical manifestation, and the considerations for changing their location.

Creating a User Tablespace: The CREATE TABLESPACE Command

To create a new tablespace in PostgreSQL, you use the CREATE TABLESPACE SQL command. This command requires two essential pieces of information: a unique name for the tablespace and the absolute path to the directory on the file system where the tablespace data will be stored.

Example:

```
postgres=# CREATE TABLESPACE myts01 LOCATION '/data01';
```

In this example:

- myts01 is the chosen name for the new tablespace. This name will be used in subsequent CREATE TABLE or CREATE INDEX statements to specify where an object should be stored.
- '/data01' is the absolute path to the directory on the server's file system that will serve as the root for this tablespace. **Crucially, this directory must exist and be empty, and the PostgreSQL operating system user must have appropriate read and write permissions to it.**

Once executed, you can verify the creation of the tablespace by querying the pg_tablespace system catalog view:

```
postgres=# SELECT spcname, spcllocation FROM pg_tablespace WHERE spcname = 'myts01';
```

This query will show myts01 as the tablespace name and /data01 as its specified location.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Physical Manifestation: Symbolic Links in \$PGDATA/pg_tblspc

While you specify an absolute path like /data01 when creating a tablespace, PostgreSQL doesn't directly store all tablespace metadata there. Instead, it maintains a structured approach within its \$PGDATA directory.

When CREATE TABLESPACE is executed, PostgreSQL creates a **symbolic link** (a shortcut or alias) within a special directory inside its data cluster: \$PGDATA/pg_tblspc/.

- The name of this symbolic link is the **Object Identifier (OID)** of the newly created tablespace (myts01).
- This symbolic link points to the physical directory you specified (e.g., /data01).

Why symbolic links? This design choice offers several advantages:

- **Centralized Management:** All tablespaces, regardless of their physical location, are logically referenced and managed from a single, well-known location within \$PGDATA.
- **Flexibility:** It allows tablespace directories to be located on different file systems, separate disk arrays, or even network-attached storage, without requiring PostgreSQL itself to be aware of the underlying storage complexities.
- **Simplified Cleanup:** When a tablespace is dropped, PostgreSQL can simply remove the symbolic link and the corresponding data directories.

Using the New Tablespace: Storing Database Objects

Once myts01 is created, you can instruct PostgreSQL to store new tables or indexes within it. This is achieved by using the TABLESPACE clause in your CREATE TABLE or CREATE INDEX statements.

Example: Creating a table in myts01 within the postgres database:

-- Connect to the 'postgres' database

```
postgres=# CREATE TABLE my_table_in_myts01 (id INT PRIMARY KEY, name TEXT)
TABLESPACE myts01;
```

Example: Creating a table in myts01 within the mydb01 database:

First, you would connect to mydb01:

```
postgres=# \c mydb01
```

You are now connected to database "mydb01" as user "postgres".



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

mydb01=# CREATE TABLE another_table_in_myts01 (code VARCHAR(10), description TEXT) TABLESPACE myts01;

Physical Storage within the Tablespace Directory (/data01)

After creating tables within myts01 (from both postgres and mydb01 databases), if you inspect the /data01 directory, you will observe a structured layout:

- 1. Database-Specific Subdirectories:** For each database that utilizes myts01, PostgreSQL will create a subdirectory within /data01. The name of each subdirectory will be the **OID of the database** that created objects in this tablespace.
 - For example, you might see /data01/12345/ (where 12345 is the OID of the postgres database) and /data01/67890/ (where 67890 is the OID of the mydb01 database). This mechanism ensures isolation and organization when multiple databases share the same tablespace.
- 2. Table Data Files:** Inside these database-specific subdirectories, you will find the actual data files for the tables created within that tablespace. As discussed previously, these files are named after the **relfilenode** (which often initially corresponds to the OID) of the table, along with their associated _fsm (Free Space Map) and _vm (Visibility Map) files.

Example Directory Structure (simplified):

/data01/

```
|— <OID_of_postgres_db>/
|   |— PG_9.x_20YYMMDD1/ (PostgreSQL version specific subdirectory)
|   |   |— <relfilenode_of_my_table_in_myts01>
|   |   |   |— <relfilenode_of_my_table_in_myts01>_fsm
|   |   |   |— <relfilenode_of_my_table_in_myts01>_vm
|   |   |   |— ...
|— <OID_of_mydb01_db>/
|   |— PG_9.x_20YYMMDD1/
|   |   |— <relfilenode_of_another_table_in_myts01>
|   |   |   |— <relfilenode_of_another_table_in_myts01>_fsm
|   |   |   |— <relfilenode_of_another_table_in_myts01>_vm
```



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

The PG_9.x_20YYMMDD1 (or similar) subdirectory within the database OID directory is a version-specific subdirectory, ensuring compatibility and allowing for upgrades without direct file system conflicts.

How to Change a Tablespace Location

```
-- Shutdown PostgreSQL.
/usr/local/pgsql/bin/pg_ctl -D /usr/local/pgsql/data stop

-- Copy the myts01 tablespace to the new file system.
[postgres@pgserver data01]$ cp -rp /data01/PG* /data02

-- Check the contents of the pg_tblspc
directory. [postgres@pgserver pg_tblspc]$ ls -l
lrwxrwxrwx. 1 postgres postgres 7 Nov 8 15:52 24576 -> /data01

-- Delete the symbolic link.
[postgres@pgserver pg_tblspc]$ rm 24576

-- Create a new symbolic link.
[postgres@pgserver pg_tblspc]$ ln -s /data02 24576

-- Confirm the contents
[postgres@pgserver pg_tblspc]$ ls -l
lrwxrwxrwx. 1 postgres postgres 7 Nov 8 15:53 24576 -> /data02

-- Startup PostgreSQL.
/usr/local/pgsql/bin/pg_ctl -D /usr/local/pgsql/data start
```

Changing the physical location of an existing tablespace is not a direct ALTER TABLESPACE ... SET LOCATION command. PostgreSQL does not natively support an online move of a tablespace's underlying directory due to the complexities of ensuring data consistency and handling active connections.

If the file system where a tablespace's directory is located becomes full, or if you need to relocate data for performance or archival reasons, the standard procedure involves these steps:

1. **Shutdown PostgreSQL:** This is crucial to ensure no ongoing writes or open file handles to the tablespace.
2. **Move the Directory:** Physically move the entire tablespace directory from its old location to the new desired location using operating system commands (e.g., mv /data01 /new_data01).



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

3. **Update the Symbolic Link:** After moving the directory, you *must* update the symbolic link in \$PGDATA/pg_tblspc/ to point to the new location.
 - First, remove the old symbolic link: `rm $PGDATA/pg_tblspc/<tablespace_OID>`
 - Then, create a new symbolic link pointing to the new location: `ln -s /new_data01 $PGDATA/pg_tblspc/<tablespace_OID>`
 - You can find the <tablespace_OID> by querying `SELECT oid FROM pg_tablespace WHERE spcname = 'myts01'`; while the database is running or by inspecting the \$PGDATA/pg_tblspc directory before the move.
4. **Restart PostgreSQL:** Once the symbolic link is updated, restart the PostgreSQL server. The server will now recognize the new location of the tablespace.

Important Note on Volume Managers:

The text correctly highlights that using a **volume manager** (like LVM on Linux or Storage Spaces on Windows) is a superior solution for managing file system capacity. Volume managers allow you to:

- **Extend Logical Volumes:** Add physical disks to an existing logical volume without downtime or data migration.
- **Stripe Data:** Distribute data across multiple physical disks for improved I/O performance.
- **Snapshotting:** Create point-in-time snapshots for backups.

By using a volume manager, you can make the underlying file system of a tablespace appear much larger or more resilient, eliminating the need to manually move tablespace directories.

Tablespaces and Partitioned Tables: A Powerful Combination

Tablespaces are particularly powerful when combined with **partitioned tables**. PostgreSQL's declarative partitioning allows you to divide a large table into smaller, more manageable child tables (partitions).

- **Flexible Capacity Management:** By assigning different partitions to different tablespaces, you can:
 - **Distribute I/O Load:** Place frequently accessed partitions on faster storage (e.g., SSDs) and less frequently accessed, historical data on slower, cheaper storage (e.g., HDDs).

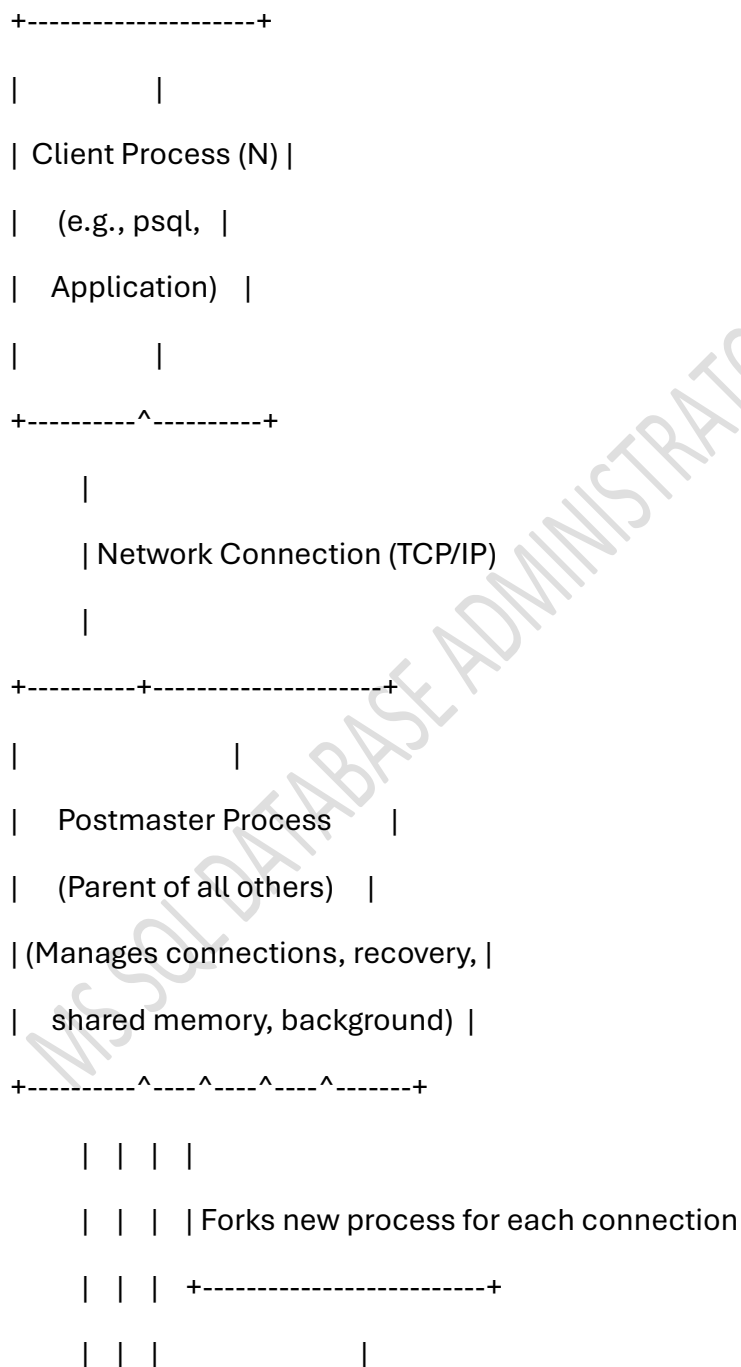


MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Manage Archiving:** Easily detach old partitions by simply changing their tablespace or even dropping the tablespace after archiving the data.
- **Cope with File System Capacity:** If one file system filling up, you can create new partitions on a new tablespace located on a different, less utilized file system, avoiding downtime for the entire table.

Diagram: PostgreSQL Process Relationship





MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia





MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

+-----+

All Backend & Background Processes access Shared Memory

Figure 1-2. Process relationship diagram

PostgreSQL's architecture, particularly its implementation of Multi-Version Concurrency Control (MVCC), necessitates a unique maintenance operation called **VACUUM**. Understanding VACUUM is fundamental to managing a healthy and performant PostgreSQL database.

What is VACUUM in PostgreSQL?

VACUUM is a crucial maintenance utility in PostgreSQL that primarily serves four vital functions:

1. **Gathering Table and Index Statistics:** VACUUM (specifically VACUUM ANALYZE or the ANALYZE command itself) collects statistical information about the distribution of data within tables and indexes. This information is then used by the PostgreSQL query planner to create efficient execution plans for SQL queries. Accurate statistics help the planner choose the best access paths (e.g., using an index vs. a sequential scan, or choosing the optimal join order).
2. **Reorganize the Table (Optional/Specific Forms):** While VACUUM generally reclaims space, it's important to differentiate. A standard VACUUM (without FULL) does *not* reorganize the table by rewriting it to a new, smaller disk footprint. It marks space as free for reuse. However, VACUUM FULL *does* reorganize the table by completely rewriting it to a new file, reclaiming all dead space and compacting the table on disk. This is a much more invasive and blocking operation.
3. **Clean Up Tables and Index Dead Blocks:** This is perhaps the most critical function directly tied to PostgreSQL's MVCC implementation. When a row is UPDATED or DELETED in PostgreSQL, the old version of the row is not immediately removed from the data page. Instead, it's marked as "dead" or "obsolete." VACUUM scans tables and indexes to identify these dead tuples and physically remove them, reclaiming the space they occupied for future use within the same table. This prevents "table bloat" – where a table occupies more physical space than necessary due to accumulated dead tuples.
4. **Frozen by Record XID to Prevent XID Wraparound:** This is another crucial function directly related to MVCC and transaction IDs (XIDs). PostgreSQL assigns



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

a unique XID to every transaction. These XIDs are 32-bit integers and, eventually, will "wrap around" (return to zero). To prevent data visibility issues that could arise from this wraparound (where very old transactions might appear to be in the future), VACUUM "freezes" old transaction IDs. This means it marks rows as "frozen" to indicate that they are permanently visible to all current and future transactions, effectively making their XID irrelevant for visibility checks. Without regular VACUUM and freezing, the database could become unusable due to **XID wraparound**.

Functions #1 and #2 (gathering statistics and certain forms of reorganization) are common requirements for most DBMSs to maintain performance. However, functions #3 (cleaning up dead blocks) and #4 (preventing XID wraparound) are **specifically necessary because of PostgreSQL's MVCC implementation**.

Differences in the MVCC Model: PostgreSQL vs. Oracle

The approach to MVCC is a fundamental architectural difference between PostgreSQL and systems like Oracle, which directly impacts the need for VACUUM.

Item	ORACLE	PostgreSQL
MVCC Model	UNDO Segments	In-place Versioning / Record within Block
Implementation	Stores <i>previous</i> versions of data in separate UNDO segments/tablespaces.	Stores <i>new</i> versions of data in the same data block (or a new block if space is insufficient), marking old versions as dead in place.
Shared Pool	Exists (Shared SQL Area, Data Dictionary Cache, Result Cache, etc.)	Does Not Exist (in the Oracle sense)
SQL Parsing	SQL parsing and execution plans are cached in the Shared Pool, accessible by all sessions.	SQL information is shared at the process level ; if the same SQL is executed multiple times <i>within one backend process</i> , it's hard-parsed only once.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Detailed Differences in the MVCC Model:

The core principle of MVCC is to enable high concurrency by ensuring "read operations do not block write operations and write operations should not block read operations."

Both Oracle and PostgreSQL achieve this, but their underlying mechanisms differ significantly:

Oracle's MVCC (UNDO Segments):

- **Mechanism:** When a row is modified or deleted in Oracle, the original (pre-change) version of the data is copied to a dedicated **UNDO segment** (or rollback segment) before the change is applied to the data block.
- **Reads:** Readers access the current data block. If a row was modified *after* the reader's transaction began, Oracle reconstructs the correct version of the row by applying the information from the UNDO segment.
- **Cleanup:** UNDO segments are automatically managed. Once a transaction commits and no other active transactions need the old versions of data in the UNDO segment, that space is automatically reclaimed. This means Oracle does not require a periodic "vacuum" like PostgreSQL to clean up old row versions, as they are managed separately.

PostgreSQL's MVCC (In-place Versioning / Record within Block):

- **Mechanism:** When a row is UPDATED in PostgreSQL, a *new* version of the row is inserted into the table (either in the same block if there's space, or a new block), and the *old* version is simply marked as "dead" (or "obsolete") by updating its xmax (transaction ID that deleted or updated the row). When a row is DELETED, the existing row is marked as dead by setting its xmax. The old data is not immediately removed.
- **Reads:** Readers examine the xmin (transaction ID that inserted the row) and xmax of each row version to determine if it was visible at the start of their transaction. If an xmin is too new or an xmax indicates the row was deleted by a committed transaction, that version is ignored.
- **Cleanup:** Because old row versions are left "in place" and simply marked as dead, they continue to occupy disk space. This accumulated "dead space" is known as **table bloat**. To reclaim this space and make it available for new rows *within the same table*, a VACUUM operation is required. Without regular vacuuming, tables can grow excessively large, impacting performance and disk utilization. This also applies to indexes, where entries pointing to dead tuples also accumulate.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Shared Pool: A Key Distinction

Another architectural difference that often surprises users familiar with Oracle is the absence of a "Shared Pool" (in the Oracle sense) in PostgreSQL.

Oracle's Shared Pool:

- **Purpose:** In Oracle, the Shared Pool is a critical SGA (System Global Area) component that caches parsed SQL statements, execution plans, data dictionary information, and results.
- **Benefit:** When multiple client sessions execute the same SQL query, Oracle can reuse the already parsed statement and execution plan from the Shared Pool, saving CPU cycles and improving performance significantly. This global caching benefits all connected sessions.

PostgreSQL's Approach (Process-Level Caching):

- **Mechanism:** PostgreSQL does not have a central, server-wide shared cache for SQL parsing and execution plans that is accessible by all concurrent backend processes. Each **backend process** (which handles a single client connection) has its own memory space where it caches parsed SQL statements and execution plans.
- **Benefit:** If a single backend process executes the *same SQL query multiple times*, it will typically "hard-parse" the query only once (i.e., perform the full parsing, semantic analysis, and plan generation). Subsequent executions within that *same process* can reuse the cached plan, leading to efficient "soft parsing."
- **Limitation:** However, if a *different* backend process (i.e., another client connection) executes the exact same SQL query, it will perform its own hard parse and create its own cached plan within its local memory. There is no global sharing of these parsed plans across different backend processes in the same way Oracle's Shared Pool operates. This means that if you have many concurrent connections all running the same query, each connection might parse it independently.

This design choice in PostgreSQL simplifies its architecture and avoids certain contention points that can arise in large shared memory areas. While it means less global caching for SQL plans, PostgreSQL relies on other optimization techniques and its efficient process model to achieve high performance.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Understanding VACUUM with Scripts

Let's illustrate the concepts of dead tuples and VACUUM with a simple example.

Scenario: We'll create a table, insert data, update some rows (creating dead tuples), and then observe the effect of VACUUM.

1. Initial Setup and Observation

-- Connect to your PostgreSQL database (e.g., 'postgres')

```
\c postgres
```

-- Create a sample table

```
CREATE TABLE products (  
    product_id SERIAL PRIMARY KEY,  
    name VARCHAR(100),  
    price DECIMAL(10, 2),  
    last_updated TIMESTAMP DEFAULT NOW()  
);
```

-- Insert some initial data

```
INSERT INTO products (name, price) VALUES  
('Laptop', 1200.00),  
('Mouse', 25.00),  
('Keyboard', 75.00),  
('Monitor', 300.00);
```

-- Check current table size on disk

-- The function pg_relation_size() shows the actual size of the data file for the table.

-- The size will be small initially.

```
SELECT pg_size_pretty(pg_relation_size('products'));
```

2. Create Dead Tuples (Updates and Deletes)



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Now, let's create some dead tuples by updating and deleting rows.

-- Update a row: This creates a new version of the row and marks the old one as dead.

```
UPDATE products SET price = 1250.00, last_updated = NOW() WHERE product_id = 1;
```

-- Update another row multiple times: This will create multiple dead versions.

```
UPDATE products SET price = 26.00, last_updated = NOW() WHERE product_id = 2;
```

```
UPDATE products SET price = 27.00, last_updated = NOW() WHERE product_id = 2;
```

```
UPDATE products SET price = 28.00, last_updated = NOW() WHERE product_id = 2;
```

-- Delete a row: This marks the row as dead.

```
DELETE FROM products WHERE product_id = 3;
```

-- Insert a new row (this might reuse some space *if* a vacuum happened, but not yet)

```
INSERT INTO products (name, price) VALUES ('Webcam', 50.00);
```

-- Check table size again.

-- You will likely see that the table size has increased or stayed the same,

-- despite updates/deletes, because dead tuples still occupy space.

```
SELECT pg_size_pretty(pg_relation_size('products'));
```

Observation: The table size has likely increased or remained stable, even though logically, some rows were modified or deleted. This is because the old versions (dead tuples) are still on disk.

3. Run VACUUM and Observe Space Reclamation

A standard VACUUM will reclaim the space from dead tuples, making it available for reuse *within the table's existing file*. It does *not* immediately shrink the file.

-- Run a standard VACUUM

```
VACUUM products;
```




MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

-- Check table size again.

-- In many cases, `pg_relation_size()` might **not** decrease after a plain `VACUUM`.

-- This is because `VACUUM` only marks space as free **within** the file, it doesn't

-- truncate the file unless dead tuples are at the very end.

-- However, the space is now available for new inserts/updates.

```
SELECT pg_size_pretty(pg_relation_size('products'));
```

-- To truly see the effect of space reuse, let's insert many rows.

-- These new rows will preferentially use the space freed by `VACUUM`.

```
INSERT INTO products (name, price) SELECT 'Item ' || generate_series(1, 100), 10.00;
```

-- Now check size again.

-- If the inserts fit into the reclaimed space, the size might not grow as much as

-- it would have without the `VACUUM`.

```
SELECT pg_size_pretty(pg_relation_size('products'));
```

4. Run `VACUUM FULL` to Shrink the Table

`VACUUM FULL` rewrites the entire table, aggressively reclaiming all dead space and shrinking the physical file size. This operation is much more resource-intensive and blocks concurrent access.

-- This will block concurrent reads and writes to the 'products' table.

```
VACUUM FULL products;
```

-- Check table size again.

-- This time, you **should** see a significant decrease in size, as the table file

-- has been rewritten and compacted.

```
SELECT pg_size_pretty(pg_relation_size('products'));
```

-- Clean up



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

DROP TABLE products;

5. Illustrating XID Wraparound Prevention (Conceptual)

You can't easily demonstrate XID wraparound and freezing with a simple script, as it involves billions of transactions. However, understanding the system views and parameters related to it is key:

- **autovacuum_freeze_max_age:** This parameter (default 200 million transactions) defines the maximum transaction age (in XIDs) before a table *must* be vacuumed to prevent XID wraparound. The autovacuum daemon will automatically trigger a VACUUM FREEZE on tables approaching this threshold.
- **pg_class.relfrozensxid:** For each table, this column in pg_class stores the TransactionId beyond which all older transaction IDs in that table are "frozen" (marked as visible to all transactions). VACUUM operations update this value.
- **pg_database.datfrozensxid:** Similar to relfrozensxid, but for the entire database.

Monitoring relfrozensxid (conceptual script):

-- This query shows how close each table is to the autovacuum_freeze_max_age threshold.

-- The higher the "age", the closer it is to needing a freeze vacuum.

SELECT

relname,

age(relfrozensxid) AS xid_age,

pg_size_pretty(pg_relation_size(oid)) AS current_size

FROM pg_class

WHERE relkind = 'r' -- 'r' for regular tables

AND relnamespace = (SELECT oid FROM pg_namespace WHERE nspname = 'public') -- only public schema

ORDER BY xid_age DESC;

-- You would need to run this over a long period with many transactions

-- to see the xid_age increase and then decrease after a VACUUM.

Conclusion:



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

VACUUM is not just a cleanup utility in PostgreSQL; it's an intrinsic part of its MVCC architecture, ensuring data consistency, preventing bloat, and safeguarding against critical issues like XID wraparound. Regular and proper vacuuming (often managed automatically by the autovacuum daemon) is essential for maintaining the performance and stability of any PostgreSQL database.

PostgreSQL's security model is built upon a robust and flexible system of **roles and privileges**. This hierarchical structure allows database administrators (DBAs) to meticulously control who can access what database resources and what actions they can perform. This fine-grained control is paramount for maintaining data integrity, confidentiality, and overall system security.

PostgreSQL Roles and Privileges: A Detailed Explanation

The combination of roles and privileges in PostgreSQL provides a comprehensive framework for authentication (verifying identity) and authorization (determining access rights).

1. Roles: The Foundation of Identity and Grouping

In PostgreSQL, a **role** is a fundamental concept that encompasses both database users and groups of users. This unified approach simplifies management by allowing an administrator to create a single entity that can either log in to the database directly (like a user) or serve as a container for privileges that can be inherited by other roles (like a group).

Key Characteristics and Functions of Roles:

- **User Identity:** A role can represent an individual user who needs to connect to the database. Such roles are typically created with login capabilities (i.e., they have a password and can initiate sessions).
 - **Example:** `CREATE ROLE sales_analyst WITH LOGIN PASSWORD 'securepass';`
- **Group Management:** Roles can also act as groups. This is a powerful feature for simplifying privilege management. Instead of granting permissions to individual users, you can grant permissions to a group role, and then grant that group role to multiple individual user roles. Any user who is a member of the group role will inherit all the privileges granted to that group.
 - **Example:**



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

CREATE ROLE data_readers; -- This will be our group role

GRANT SELECT ON financial_data TO data_readers; -- Grant privilege to the group

GRANT data_readers TO sales_analyst; -- Grant the group role to a user role

Now, sales_analyst automatically inherits SELECT privilege on financial_data.

- **Attributes (Permissions within the Role):** Roles can be created with specific attributes that define their global capabilities within the database cluster. These attributes are inherent permissions granted directly to the role itself, not necessarily on specific database objects. Common role attributes include:
 - **LOGIN:** Allows the role to connect to the database server.
 - **SUPERUSER:** Grants all privileges within the database cluster, bypassing all permission checks (use with extreme caution!).
 - **CREATEDB:** Allows the role to create new databases.
 - **CREATEROLE:** Allows the role to create, modify, and drop other roles.
 - **INHERIT / NOINHERIT:** Controls whether a role inherits the privileges of roles it is a member of (default is INHERIT). NOINHERIT is useful for more granular control, where inherited privileges must be explicitly "SET ROLE" by the user.
 - **BYPASSRLS:** Bypasses Row-Level Security (RLS) policies.
 - **CONNECTION LIMIT n:** Specifies the maximum number of concurrent connections for a role.
 - **PASSWORD '...':** Sets the password for login roles.
 - **VALID UNTIL 'timestamp':** Sets an expiration date for the role's password.
- **Default Role:** The initial superuser role created during initdb is typically postgres (or the operating system user who ran initdb). This role has SUPERUSER privileges by default.

Benefits of Using Roles:

- **Simplified Administration:** Instead of managing permissions for dozens or hundreds of individual users, administrators can define a smaller set of group roles, grant privileges to these roles, and then simply add or remove users from these groups.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Consistency:** Ensures that all users performing a similar function (e.g., "data entry specialists") have the exact same set of permissions, reducing errors and inconsistencies.
- **Auditability:** Makes it clearer to audit who has what level of access by inspecting role memberships and role privileges.
- **Flexibility:** Allows for dynamic adjustment of permissions. If a user's responsibilities change, you can simply change their role memberships rather than revoking and re-granting individual privileges.

2. Privileges: Defining Actions on Database Objects

Privileges (also known as permissions) define the specific actions that a role is allowed to perform on a particular **database object**. These are the granular access controls that govern interactions with data, schema elements, and other database components.

- **The GRANT and REVOKE Commands:** Privileges are assigned using the GRANT command and removed using the REVOKE command.
- **Common Privileges and Their Associated Objects:**
 - **SELECT:**
 - **On Tables/Views:** Allows reading data from the specified table or view.
 - **On Sequences:** Allows using curval() and nextval() functions for sequences.
 - **On Functions/Procedures:** Allows executing the function/procedure.
 - **On Schemas:** Allows looking up objects within the schema.
 - **INSERT:**
 - **On Tables/Views:** Allows adding new rows to the table or view.
 - **UPDATE:**
 - **On Tables/Views:** Allows modifying existing rows in the table or view.
 - **On Sequences:** Allows using setval() function for sequences.
 - **DELETE:**
 - **On Tables/Views:** Allows removing rows from the table or view.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **TRUNCATE:**
 - **On Tables:** Allows emptying a table completely, which is faster than DELETE but cannot be rolled back (unless within a transaction).
- **REFERENCES:**
 - **On Tables:** Allows creating a foreign key constraint that references the table.
- **TRIGGER:**
 - **On Tables:** Allows creating triggers on the table.
- **CREATE:**
 - **On Databases:** Allows creating new schemas within the database.
 - **On Schemas:** Allows creating new objects (tables, views, functions, etc.) within the schema.
 - **On Tablespaces:** Allows creating tables, indexes, and databases in the tablespace.
- **CONNECT:**
 - **On Databases:** Allows a role to connect to the specified database.
- **USAGE:**
 - **On Schemas:** Allows accessing objects within the schema (e.g., SELECT * FROM schema_name.table_name). SELECT on a table in a schema also implicitly requires USAGE on the schema.
 - **On Sequences:** Allows using the sequence (implicitly includes SELECT and UPDATE on the sequence).
 - **On Functions/Procedures:** Allows calling the function/procedure.
 - **On Foreign Data Wrappers/Servers/User Mappings:** Allows using these foreign data objects.
- **ALL PRIVILEGES:**
 - Grants all available privileges on the specified object type.
- **Grant Option (WITH GRANT OPTION):** When a privilege is granted WITH GRANT OPTION, the recipient role not only gains the privilege but also gains the ability to



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

grant that same privilege to other roles. This creates a chain of delegable permissions.

- **Example:** GRANT SELECT ON orders TO manager WITH GRANT OPTION; manager can now GRANT SELECT ON orders TO sales_team;
- **Default Privileges (ALTER DEFAULT PRIVILEGES):** PostgreSQL allows you to define **default privileges** for objects that will be created in the future by a specific role or within a specific schema. This is incredibly useful for automating permissions for newly created objects.
 - **Example:**

```
ALTER DEFAULT PRIVILEGES FOR ROLE app_owner IN SCHEMA public  
GRANT SELECT, INSERT, UPDATE, DELETE ON TABLES TO app_user_role;
```

This statement means that any *new* tables created by app_owner in the public schema will automatically grant SELECT, INSERT, UPDATE, DELETE to app_user_role. Existing tables are not affected.

- **Ownership Privileges:** Every database object (tables, views, functions, etc.) has an **owner** (the role that created it). The owner implicitly has all privileges on that object, including the ability to GRANT and REVOKE privileges on it. Ownership can be transferred using ALTER TABLE ... OWNER TO

How Roles and Privileges Work Together: The Access Control Flow

1. **Authentication:** A client attempts to connect to a PostgreSQL server using a specific **login role** (e.g., sales_analyst).
2. **Connection Privilege:** PostgreSQL first checks if the sales_analyst role has the LOGIN attribute and the CONNECT privilege on the requested database.
3. **Authorization (Role-Based):** Once connected, the user's effective privileges are determined by:
 - Any attributes directly granted to their login role (e.g., SUPERUSER, CREATEDB).
 - The sum of all privileges granted to the roles they are members of (if INHERIT is enabled).
4. **Authorization (Object-Based):** When the user executes an SQL command (e.g., SELECT * FROM financial_data;), PostgreSQL checks if the effective privileges of the user's current role allow that specific action on the target object



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

(financial_data). This involves checking SELECT privilege on the table and USAGE privilege on the schema containing the table.

5. **Row-Level Security (RLS):** If Row-Level Security policies are active on the table, these policies are evaluated *after* the initial object privileges, to determine which specific rows the user can see or modify.

Example Scenario: Setting up a Secure Environment

Let's illustrate with a practical example:

-- 1. Create Login Roles for individuals

```
CREATE ROLE alice WITH LOGIN PASSWORD 'alice_pass';
```

```
CREATE ROLE bob WITH LOGIN PASSWORD 'bob_pass';
```

```
CREATE ROLE charlie WITH LOGIN PASSWORD 'charlie_pass';
```

-- 2. Create Group Roles to define functional access

```
CREATE ROLE data_entry_team;
```

```
CREATE ROLE reporting_team;
```

```
CREATE ROLE db_developers CREATEROLE CREATEDB; -- Devs can create roles and  
dbs
```

-- 3. Grant Group Roles to Login Roles

```
GRANT data_entry_team TO alice;
```

```
GRANT reporting_team TO bob;
```

```
GRANT db_developers TO charlie;
```

-- 4. Create a Database (charlie can do this due to CREATEDB)

```
-- Connect as charlie first: \c postgres charlie
```

```
-- Then:
```

```
CREATE DATABASE project_db OWNER charlie;
```

```
-- Connect to the new database
```



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

\c project_db

-- 5. Create a Schema for application data

```
CREATE SCHEMA app_data AUTHORIZATION charlie;
```

-- 6. Grant Schema USAGE for all relevant teams

```
GRANT USAGE ON SCHEMA app_data TO data_entry_team;
```

```
GRANT USAGE ON SCHEMA app_data TO reporting_team;
```

-- 7. Create a Table in the new schema (charlie as owner of schema can do this)

```
CREATE TABLE app_data.customer_orders (  
    order_id SERIAL PRIMARY KEY,  
    customer_name VARCHAR(100),  
    order_date DATE,  
    total_amount DECIMAL(10, 2)  
);
```

-- 8. Grant specific privileges to group roles on the table

```
GRANT SELECT, INSERT, UPDATE ON app_data.customer_orders TO data_entry_team;
```

```
GRANT SELECT ON app_data.customer_orders TO reporting_team;
```

-- 9. Test as Alice (data_entry_team member)

```
-- Connect as alice: \c project_db alice
```

```
-- Alice can insert:
```

```
INSERT INTO app_data.customer_orders (customer_name, order_date, total_amount)  
VALUES ('Customer A', '2024-07-30', 150.75);
```

```
-- Alice can select:
```



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

```
SELECT * FROM app_data.customer_orders;
```

```
-- Alice cannot delete:
```

```
-- DELETE FROM app_data.customer_orders WHERE order_id = 1; -- ERROR: permission
denied for table customer_orders
```

-- 10. Test as Bob (reporting_team member)

```
-- Connect as bob: \c project_db bob
```

```
-- Bob can select:
```

```
SELECT * FROM app_data.customer_orders;
```

```
-- Bob cannot insert:
```

```
-- INSERT INTO app_data.customer_orders (customer_name, order_date, total_amount)
```

```
-- VALUES ('Customer B', '2024-07-31', 200.00); -- ERROR: permission denied for table
customer_orders
```

-- 11. Clean up (as a superuser or the owner of project_db)

```
-- \c postgres
```

```
-- DROP DATABASE project_db;
```

```
-- DROP ROLE alice;
```

```
-- DROP ROLE bob;
```

```
-- DROP ROLE charlie;
```

```
-- DROP ROLE data_entry_team;
```

```
-- DROP ROLE reporting_team;
```

```
-- DROP ROLE db_developers;
```

PostgreSQL: Object Hierarchies, Replication, Load Balancing, and High Availability

PostgreSQL is a sophisticated relational database management system that employs a logical and physical hierarchy for organizing data, coupled with advanced features to ensure data availability, performance, and resilience.

1. PostgreSQL Object Hierarchies



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

PostgreSQL organizes its data and logical structures in a clear hierarchical fashion, making it intuitive for users and developers to manage database resources effectively. This hierarchy spans from the highest level of database clusters down to granular objects like columns and triggers.

- **Clusters:** At the top of the hierarchy is the **database cluster**. A cluster represents a single running PostgreSQL server instance, managing a collection of individual databases. All databases within a cluster share the same PostgreSQL server processes, shared memory, and configuration. The \$PGDATA directory is the physical manifestation of a database cluster.
- **Databases:** Within a database cluster, **databases** serve as the primary, isolated containers for data storage. Each database is an independent entity with its own set of schemas, tables, and other objects. Databases within the same cluster cannot directly access data from each other without specific mechanisms like dblink or foreign data wrappers. Standard databases like postgres, template0, and template1 are created by default during initdb.
- **Schemas:** **Schemas** act as namespaces within a database. They provide a logical grouping for database objects like tables, views, functions, and indexes, preventing naming conflicts and aiding in organization. For instance, public is the default schema. You can have sales.customers and hr.customers tables in the same database without conflict. Schemas are also critical for security, as privileges can be granted at the schema level.
- **Tables:** **Tables** are the fundamental structures for storing data records in a relational database. Each table consists of rows (records) and columns (attributes). Tables enforce data types, constraints (e.g., PRIMARY KEY, FOREIGN KEY, NOT NULL), and can be partitioned for managing large datasets.
- **Columns:** **Columns** represent individual attributes or fields within a table. Each column has a specific data type (e.g., INTEGER, VARCHAR, DATE) and can have associated constraints (e.g., UNIQUE, DEFAULT values).
- **Indexes:** **Indexes** are special lookup tables that a database search engine can use to speed up data retrieval operations. They are created on one or more columns of a table and provide quick access to rows based on the indexed values. PostgreSQL supports various index types, including B-tree (most common), hash, GiST, SP-GiST, GIN, and BRIN.
- **Views:** **Views** are virtual tables defined by a stored query. They do not store data themselves but present data from one or more underlying tables as if it were a single table. Views are commonly used for:



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Simplifying Complex Queries:** Encapsulating complex joins or calculations.
- **Security:** Restricting access to certain rows or columns of a table without granting direct access to the base table.
- **Data Abstraction:** Providing a consistent interface even if the underlying table structure changes.
- **Functions (Stored Procedures): Functions** (often referred to as stored procedures in other DBMSs) are user-defined code routines that encapsulate a series of SQL statements and procedural logic. They can accept input parameters, perform operations, and return a value (or a set of values). Functions are vital for:
 - **Code Reusability:** Writing complex logic once and calling it multiple times.
 - **Performance:** Reducing network round trips by executing logic directly on the database server.
 - **Security:** Granting execution privileges on functions without granting direct access to the underlying tables. PostgreSQL supports functions written in SQL, PL/pgSQL (its procedural language), C, and other procedural languages.
- **Triggers: Triggers** are a special type of function that automatically execute in response to specific events (e.g., INSERT, UPDATE, DELETE) on a particular table. They are used to enforce complex business rules, maintain data integrity, or audit changes. Triggers can be defined to fire BEFORE or AFTER the event.
- **Sequences: Sequences** are database objects that generate unique, incremental numeric values. They are commonly used to generate primary key values for tables, ensuring uniqueness without manual intervention.

This hierarchical structure streamlines data organization and manipulation, providing a clear and logical framework for managing complex database environments in PostgreSQL.

2. Replication, Load Balancing, and High Availability

For mission-critical applications, ensuring continuous operation, high performance, and data integrity is paramount. PostgreSQL offers robust features for **Replication**, enables **Load Balancing**, and supports various architectures for **High Availability (HA)** and **Disaster Recovery (DR)**.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

Replication

Replication is the process of creating and maintaining multiple copies of a database across different servers (nodes). This technique significantly improves data availability, fault tolerance, and read performance. PostgreSQL offers two primary types of built-in replication:

- **Physical Replication (Streaming Replication):**
 - **Mechanism:** This is PostgreSQL's most common and robust replication method. It involves replicating the entire database cluster at the binary level. The primary server continuously streams its **Write-Ahead Log (WAL)** records to one or more replica (standby) servers. The replicas then apply these WAL records to their own data files, essentially replaying the changes made on the primary.
 - **Benefits:**
 - **High Fidelity:** Replicates the entire data directory, including all databases, schemas, tables, indexes, and even configuration changes.
 - **Disaster Recovery:** Excellent for DR as it provides an exact copy of the primary.
 - **High Performance:** Very efficient for replicating large volumes of changes.
 - **Read Scale-Out:** Replica servers can be used for read-only queries, offloading work from the primary and scaling read performance.
 - **Limitations:** Replicates the entire cluster; you cannot selectively replicate individual tables or databases. Manual promotion of a replica to primary might be required in case of failover without a cluster manager.
- **Logical Replication:**
 - **Mechanism:** Introduced in PostgreSQL 10, logical replication focuses on replicating specific changes at a logical level, rather than block-level changes. It decodes the WAL into a stream of logical changes (e.g., INSERT, UPDATE, DELETE statements) that can then be applied to a subscriber database.
 - **Benefits:**



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Selective Replication:** Allows replicating specific tables, schemas, or even entire databases, providing granular control.
- **Cross-Version/Platform Compatibility:** Can replicate between different major PostgreSQL versions and even to other database systems (using foreign data wrappers or custom applications to consume the logical stream).
- **Schema Transformation:** The subscriber can have a different schema than the publisher (e.g., different column order, additional columns), as long as the replicated columns are compatible.
- **Zero-Downtime Upgrades:** Can be used for major version upgrades with minimal downtime.
- **Limitations:** More CPU overhead than physical replication due to WAL decoding. Requires primary keys on tables for updates/deletes.

Load Balancing

Load balancing in the context of PostgreSQL refers to the distribution of database queries and client connections across multiple servers or nodes to enhance performance, improve throughput, and ensure optimal resource utilization. While PostgreSQL itself doesn't have a built-in load balancer for client connections, it integrates seamlessly with external tools:

- **Connection Poolers:** These are common solutions for load balancing in PostgreSQL. They sit between client applications and the PostgreSQL server(s), managing a pool of database connections and efficiently distributing incoming queries.
 - **PgBouncer:** A lightweight connection pooler primarily focused on reducing the overhead of establishing new connections. It can manage a pool of connections and route requests to the appropriate primary or replica, depending on the configuration (e.g., routing read-only transactions to replicas).
 - **Pgpool-II:** A more feature-rich middleware that offers connection pooling, load balancing (especially for read queries to replicas), automatic failover, and query caching. It can be configured to send write queries to the primary and distribute read queries among available replicas.
- **DNS Round Robin / Application-Level Load Balancing:** Simpler forms of load balancing can involve configuring DNS records to cycle through multiple IP



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

addresses of PostgreSQL servers or implementing load balancing logic directly within the application code.

High Availability (HA)

High Availability (HA) is a critical requirement for mission-critical applications, ensuring continuous database operation even in the face of hardware failures, network issues, or other outages. PostgreSQL's replication capabilities are the foundation for HA. A common approach involves setting up a **primary-replica (master-standby) architecture**:

- **Primary-Replica Architecture:**
 - One server acts as the **primary** (or master), handling all write operations and typically read operations.
 - One or more **replica** (or standby) servers maintain synchronized copies of the primary's data through physical or logical replication. These replicas can also serve read-only queries.
 - In the event of a primary server failure, one of the replicas can be **promoted** to become the new primary, minimizing downtime.
- **Cluster Managers for Automated HA:** While manual failover is possible, automated cluster managers are essential for true HA environments. These tools monitor the health of the primary and replicas, automatically detect failures, initiate failover (promoting a replica), and reconfigure the cluster.
 - **Patroni:** A highly popular and recommended open-source solution for PostgreSQL HA. Patroni is a template for setting up a HA PostgreSQL cluster using distributed consensus tools (like etcd, Consul, or ZooKeeper). It handles:
 - **Health Monitoring:** Continuously checks the status of all PostgreSQL nodes.
 - **Automated Failover:** If the primary fails, Patroni automatically selects the healthiest replica and promotes it.
 - **Replication Management:** Sets up and manages physical streaming replication.
 - **Configuration Management:** Stores and distributes PostgreSQL configuration.
 - **Service Discovery:** Integrates with DNS or load balancers to redirect traffic after failover.



MS SQL DATABASE ADMINISTRATOR 30 DAYS COURSE

BY: Mukesh Chaurasia

- **Other HA Tools:** Other tools include Keepalived (for virtual IP management), Corosync/Pacemaker (general-purpose cluster resource managers), and proprietary solutions from cloud providers.
- **Recovery Point Objective (RPO) and Recovery Time Objective (RTO):**
 - **RPO (Recovery Point Objective):** Defines the maximum acceptable amount of data loss measured in time. With streaming replication, RPO can be very low (seconds or even zero data loss if using synchronous replication).
 - **RTO (Recovery Time Objective):** Defines the maximum acceptable downtime measured in time. With automated HA solutions like Patroni, RTO can be very low (minutes or even seconds for failover).