**Why Use Columnar Storage (Amazon RedShift) Instead of Normal Indexing for a Data Warehouse?**

In **analytical workloads (OLAP)**, queries often involve scanning large amounts of data and aggregating values (e.g., SUM, AVG, COUNT). **Columnar storage is designed to optimize these operations**, whereas traditional indexing is more suited for transactional databases (OLTP).

**Key Advantages of Columnar Storage (Amazon Redshift) Over Normal Indexing**

**1. Faster Query Performance (Reduced I/O)**

- In a **row-based** system, scanning a table means reading **every row** in its entirety, even if you only need a few columns.

- **Columnar storage reads only the relevant columns**, drastically reducing I/O.

- Example: If a report only needs total  from a billion-row table, a **row-based** system would read **all columns** for every row, while a **columnar** system reads only total.

**2. Better Compression (Storage Efficiency)**

- **Columnar storage compresses better** because similar data types are stored together.

- Example: A column with status values (Active, Inactive) can use **dictionary encoding**, leading to much better compression than row-based storage.

- This saves **disk space** and **reduces query execution time**.

**3. Eliminates the Need for Multiple Indexes**

- **Row-based databases rely on indexes** to speed up queries. However:

  o Indexing takes **extra storage** and slows down **INSERTs/UPDATEs**.

  o Queries that use **different filter conditions** may require multiple indexes

- **Columnar storage naturally speeds up filtering**, so you **don't need as many indexes**.

**4. Optimized for Aggregations and Analytics**

- Analytical queries (e.g., **SUM, AVG, COUNT, GROUP BY**) process **entire columns**.

- In a **columnar database**, aggregations are performed **directly on compressed column data**, making them much faster.

- Traditional row-based databases require **scanning and filtering entire rows**, making them **slower for large datasets**.

**5. Parallel Processing & Query Optimization**

- Amazon Redshift uses **MPP (Massively Parallel Processing)**, where each node processes different **columns independently**.

- This allows Redshift to **execute queries faster** than a traditional row-based database.