

## 04 Row Vs Columnar

### How the data is stored on Disk?

a	b	c
a1	b1	c1
a2	b2	c2
a3	b3	c3

#### Row-Based Layout

a1	b1	c1	a2	b2	c2	a3	b3	c3
----	----	----	----	----	----	----	----	----

#### Column-Based Layout

a1	a2	a3	b1	b2	b3	c1	c2	c3
----	----	----	----	----	----	----	----	----

### Row VS Column oriented data

Aspect	Row-oriented	Column-oriented
<b>Data Storage</b>	Stores data together by row	Stores data together by column
<b>Read Efficiency</b>	Efficient for reading entire records	Efficient for reading specific columns
<b>Write Efficiency</b>	Generally faster for writing new records	May be slower for writing, as it needs to update multiple column files
<b>Compression</b>	Less efficient compression	More efficient compression, especially for columns with low cardinality
<b>Query Performance</b>	Better for queries that access entire rows	Better for analytical queries that access a subset of columns
<b>OLTP vs OLAP</b>	Suited for OLTP (Online Transaction Processing)	Suited for OLAP (Online Analytical Processing)
<b>Data Retrieval</b>	Retrieves all fields of a record, even if not needed	Can retrieve only the required columns, reducing I/O
<b>Update Operations</b>	Efficient for updating entire records	Less efficient for updating individual records
<b>Aggregations</b>	Less efficient for column-wise aggregations	Very efficient for column-wise aggregations
<b>Storage Overhead</b>	Generally less storage overhead	May have more storage overhead due to metadata
<b>Schema Changes</b>	Easier to add new columns	Can be more complex to add new columns
<b>Data Skipping</b>	Limited data skipping capabilities	Efficient data skipping using column statistics
<b>File Formats</b>	Examples: CSV, JSON	Examples: Parquet, ORC
<b>Memory Usage</b>	May use more memory when only specific columns are needed	Efficient memory usage when working with specific columns

Aspect	Row-oriented	Column-oriented
<b>Vectorized Processing</b>	Less suitable for vectorized processing	Well-suited for vectorized processing
<b>Encoding Schemes</b>	Limited encoding options	Supports various encoding schemes (dictionary, run-length, etc.)
<b>Data Integrity</b>	Easier to maintain referential integrity	Can be more challenging to maintain referential integrity
<b>Use Cases</b>	Transactional systems, real-time record lookups	Data warehouses, analytical workloads, BI tools
<b>Partitioning</b>	Typically partitioned by row ranges	Can be partitioned both by row ranges and column values
<b>Cache Efficiency</b>	Less cache-efficient for analytical queries	More cache-efficient for analytical queries
<b>Scan Performance</b>	Slower for full table scans on specific columns	Faster for full table scans on specific columns