# Apache Hive : Basic Concepts & Interview Questions

**Deepa Vasanthkumar**

# Apache Hive: Transforming Big Data into Insights 🚀

Apache Hive, an open-source data warehouse infrastructure built on top of Hadoop, continues to revolutionize the way we analyze and process large volumes of data. 🌐 💡

## 🔍 What is Apache Hive?

Apache Hive provides a SQL-like interface to query and manage large datasets stored in distributed storage. It enables users to analyze data using familiar SQL syntax, making it accessible to data analysts and engineers with SQL proficiency.

🚀 Key Features:

**Scalability**: Hive seamlessly scales to handle petabytes of data across a cluster of commodity hardware.

**Performance Optimization**: With features like partitioning, indexing, and query optimization, Hive ensures efficient data processing.

**Integration with Hadoop Ecosystem**: Hive integrates with other Hadoop ecosystem components such as HDFS, HBase, and Spark, enabling seamless data workflows.

[Deepa Vasanthkumar – Medium](#)
[Deepa Vasanthkumar -| LinkedIn](#)

**Data Warehousing Capabilities**: It supports ACID transactions and complex data types, making it suitable for data warehousing applications.

**Extensibility**: Hive's architecture allows for the development of custom extensions and UDFs (User Defined Functions) to meet specific use case requirements.

Apache Hive plays a crucial role in the field of data engineering by providing a high-level abstraction and SQL-like interface for querying and analyzing large datasets stored in distributed environments, typically Hadoop-based. Here's how Hive contributes to the data engineering landscape:

**Data Processing and Analysis:** Hive allows data engineers to process and analyze massive volumes of structured and semi-structured data using familiar SQL queries. This simplifies the data processing tasks and makes it accessible to a broader audience, including data analysts and business users.

**Scalability**: Hive is designed to scale horizontally, allowing it to handle petabytes of data across a distributed cluster of commodity hardware. Data engineers can leverage Hive to process and analyze data at scale without worrying about infrastructure limitations.

**Integration with Hadoop Ecosystem**: Hive seamlessly integrates with other components of the Hadoop ecosystem, such as HDFS (Hadoop Distributed File System), HBase, Spark, and others. This integration enables data engineers to build end-to-end data pipelines that incorporate various data processing and storage technologies.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

**Data Warehousing Capabilities:** Hive provides features commonly found in traditional data warehouses, such as support for ACID (Atomicity, Consistency, Isolation, Durability) transactions, partitioning, indexing, and complex data types. This makes it suitable for building data warehousing solutions on top of Hadoop infrastructure.

**Performance Optimization:** While Hadoop MapReduce, the underlying execution engine for Hive, is known for batch processing and may have performance limitations, Hive optimizes query execution through techniques like query planning, query optimization, and runtime optimizations. Additionally, Hive supports vectorized query execution and cost-based optimization, improving performance for complex queries.

**Extensibility**: Hive's architecture allows for extensibility through custom user-defined functions (UDFs), user-defined aggregates (UDAFs), and user-defined types (UDTs). Data engineers can leverage these extensibility mechanisms to incorporate custom logic and processing into their Hive queries, tailored to specific use cases.

**Data Governance and Security:** Hive provides features for managing access control, data governance, and security, ensuring that sensitive data is protected and access to data is appropriately restricted based on user roles and permissions.

Apache Hive uses a SQL-like language called HiveQL (HQL) for querying and manipulating data stored in Hadoop Distributed File System (HDFS) or other compatible distributed storage systems.

 **Creating Tables:**

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

## Internal Table

```
CREATE TABLE table_name (

    column1 data_type,

    column2 data_type,

    ...

)

[PARTITIONED BY (partition_column data_type, ...)]

[CLUSTERED BY (bucket_column_name) INTO num_buckets BUCKETS]

[ROW FORMAT ...]

[STORED AS ...]

[TBLPROPERTIES (...)];
```

## External Table

```
CREATE EXTERNAL TABLE table_name (


    column1 data_type,


    column2 data_type,


    ...


)


[PARTITIONED BY (partition_column data_type, ...)]


[CLUSTERED BY (bucket_column_name) INTO num_buckets BUCKETS]


[ROW FORMAT ...]


[STORED AS ...]


[LOCATION 'hdfs://path/to/external/data']


[TBLPROPERTIES (...)];
```

- `PARTITIONED BY`: Specifies partition columns for the table.
- `CLUSTERED BY`: Specifies bucketing columns and the number of buckets.
- `ROW FORMAT`: Defines how data rows are formatted.
- `STORED AS`: Specifies the storage format of the table (e.g., TEXTFILE, ORC, PARQUET).
- `LOCATION`: Specifies the HDFS path where external table data resides.
- `TBLPROPERTIES`: Allows setting additional table properties like serde properties, table comments, etc.

**Loading Data into Tables:**

```
LOAD DATA [LOCAL] INPATH 'input_path' [OVERWRITE] INTO TABLE table_name
```

**Querying Data:**

```
SELECT column1, column2, ...


FROM table_name


[WHERE condition]
```

Deepa Vasanthkumar – Medium

Deepa Vasanthkumar -| LinkedIn

```
[GROUP BY column]


[ORDER BY column]
```

**JoiningTables:**

```
SELECT ...


FROM table1 t1


JOIN table2 t2 ON t1.column = t2.column
```

## Example

```
SELECT e.emp_name, d.dept_name FROM employee e JOIN department d ON
e.emp_dept = d.dept_id;
```

**Aggregating Data:**

```
SELECT emp_dept, COUNT() FROM employee GROUP BY emp_dept;
```

**Filtering Data:**

```
SELECT FROM employee WHERE emp_dept = 'IT' AND emp_id > 100;
```

[Deepa Vasanthkumar – Medium](#)
[Deepa Vasanthkumar -| LinkedIn](#)

**Subqueries:**

```
SELECT emp_name FROM (SELECT * FROM employee WHERE emp_dept = 'IT' ) AS
it_employees;
```

**Conditional Logic**

```
SELECT emp_name, CASE WHEN emp_dept = 'IT' THEN 'Information Technology'
WHEN emp_dept = 'HR' THEN 'Human Resources' ELSE 'Other' END AS department
FROM employee;
```

**Creating Views:**

```
CREATE VIEW it_employees_view AS SELECT emp_name FROM employee WHERE
emp_dept = 'IT';
```

# Hive Bucketing and Partitioning

Apache Hive offers powerful features like bucketing and partitioning, providing data engineers with efficient ways to organize and manage large datasets in distributed environments.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

🔍 **Partitioning**:

Partitioning involves dividing data into logical partitions based on one or more columns, typically date, region, or category. This allows Hive to store and process data more efficiently by organizing it into directories based on partition keys.

Benefits:

1. Improved Query Performance: Partitioning enables Hive to prune unnecessary data partitions during query execution, reducing the amount of data scanned and improving query performance.

2. Data Organization: Partitioning makes it easier to manage and query large datasets, especially when dealing with time-series or hierarchical data.

3. Optimized Data Loading: Partitioning can accelerate data loading processes by partitioning data into separate directories, allowing parallelism during data ingestion.

🗃️ **Bucketing**:

Bucketing involves dividing data into fixed-sized buckets based on the hash value of a chosen column. Each bucket contains a subset of data rows, and Hive ensures that rows with the same hash value are stored in the same bucket.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Benefits:

1. Data Skew Handling: Bucketing helps mitigate data skew issues by distributing data evenly across buckets, ensuring balanced query processing and resource utilization.

2. Join Optimization: Bucketed tables can accelerate join operations by reducing data shuffling and improving join locality.

3. Sampling and Sampling-Based Queries: Bucketing facilitates efficient sampling and sampling-based queries by providing a uniform distribution of data across buckets.

🌟 **Best Practices:**

- Choose Appropriate Partition and Bucket Columns: Select partition and bucket columns wisely based on query patterns, data distribution, and use cases.

# Hive with Cloud

**Amazon Web Services (AWS):**

  - Amazon EMR (Elastic MapReduce): Amazon EMR supports Hive out-of-the-box, allowing users to deploy Hive clusters on AWS infrastructure quickly. EMR provides managed Hadoop clusters, simplifying the setup and management of Hive deployments.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

- Amazon S3 (Simple Storage Service): Hive can directly interact with data stored in Amazon S3, allowing users to query and analyze data without needing to provision and manage HDFS clusters. This provides flexibility and cost-effectiveness in storing and processing data on AWS.

- AWS Glue: AWS Glue provides a fully managed extract, transform, and load (ETL) service that integrates with Hive metastore, allowing users to define, orchestrate, and monitor data pipelines for data preparation and transformation tasks.

## Microsoft Azure:

- Azure HDInsight: Azure HDInsight offers managed Hadoop clusters that support Apache Hive. Users can deploy Hive clusters on Azure and leverage Azure Blob Storage as the underlying storage for data processing. HDInsight provides integration with other Azure services for seamless data workflows.

- Azure Data Lake Storage: Hive can interact with data stored in Azure Data Lake Storage (ADLS), enabling users to run Hive queries on data stored in ADLS without needing to move or replicate data to other storage systems.

## Google Cloud Platform (GCP):

- Google Dataproc: Google Dataproc provides managed Apache Hadoop and Apache Spark clusters on GCP. Users can deploy Hive clusters on Dataproc and leverage Google Cloud Storage (GCS) as the underlying storage for data processing.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

- BigQuery: While not a direct integration with Hive, Google BigQuery offers a serverless, highly scalable data warehouse solution. Users can load data from Hive or other sources into BigQuery for analysis using SQL queries, providing an alternative to traditional Hive deployments on GCP.

Hive also supports generic JDBC and ODBC connectors, allowing users to connect to external data sources and cloud-based databases for data querying and analysis.

## Limitations of Hive

**High Latency:** Hive queries typically have higher latency compared to traditional relational databases due to the underlying MapReduce or Tez execution engine. This makes it less suitable for use cases requiring low-latency responses or real-time data processing.

**Batch Processing:** Hive is primarily designed for batch processing rather than interactive or real-time data processing. It may not be the best choice for use cases requiring real-time analytics or interactive querying.

**Limited SQL Support:** While Hive provides a SQL-like interface, it does not support the full SQL standard or advanced SQL features found in traditional relational databases. Certain SQL operations or functions may not be supported, leading to limitations in query expressiveness and functionality.

**Performance Overhead:** Hive introduces performance overhead compared to executing SQL queries directly on relational databases. The translation of SQL queries into

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

MapReduce or Tez jobs, along with data serialization and deserialization, can result in slower query execution times, especially for complex queries.

**Schema Evolution Challenges:** Hive's schema-on-read approach allows flexibility in handling semi-structured or evolving data schemas. However, it can also lead to challenges in schema management and data consistency, particularly when dealing with evolving data schemas or schema migrations.

**Lack of Indexing Support:** While Hive supports basic indexing mechanisms like partitioning and bucketing, it lacks advanced indexing capabilities found in traditional databases. This can impact query performance, especially for ad-hoc or complex queries that cannot leverage existing indexes.

**Limited Support for Transactions:** While Hive supports ACID transactions in certain scenarios (e.g., using ORC or Parquet file formats), it may not provide full support for transactions across all data operations and storage formats, leading to potential data consistency issues in certain use cases.

# Basic Interview Questions

**Q. What is the difference between Hive and traditional databases like MySQL or Oracle?** Hive is designed for querying and analyzing large datasets stored in Hadoop's distributed file system (HDFS), whereas traditional databases like MySQL or Oracle are relational databases primarily used for transactional processing.

**Q. What are the components of Hive?**

 Hive consists of three main components:

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Metastore: Stores metadata about Hive tables and partitions.

Query Processor: Translates HiveQL queries into MapReduce, Tez, or Spark jobs.

Execution Engine: Executes the query plans generated by the Query Processor.

**Q. What is HiveQL?**    Hive Query Language (HiveQL) is a query language similar to SQL. It provides a familiar SQL-like interface to interact with data stored in Hive.

**Q. What is a Hive metastore?**  The Hive metastore is a central repository that stores metadata for Hive tables, including their schema, location, and partition information. It acts as an interface between HiveQL and underlying storage.

**Q. What is the difference between an external table and a managed table in Hive?**

Managed Table: Hive manages the lifecycle of data for managed tables, including data deletion. Data is stored in a location controlled by Hive.

External Table: External tables are created over data that is not managed by Hive. Hive does not manage the data lifecycle, and data can reside anywhere in HDFS or other supported file systems.

**Q. What are Hive partitions?**

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Hive partitions allow you to organize data in a table based on one or more columns. Partitioning can significantly improve query performance by restricting the amount of data scanned during query execution.

**Q. What is dynamic partitioning in Hive?**    Dynamic partitioning in Hive allows you to automatically create partitions based on the values of specific columns during data insertion. It simplifies the process of managing partitions by eliminating the need for explicit partition creation.

**Q. Explain the difference between bucketing and partitioning in Hive.**

Partitioning: Partitioning divides data into logical units based on the values of one or more columns. It helps in organizing and querying data efficiently.

Bucketing: Bucketing, also known as clustering, divides data into fixed-sized buckets or files based on the hash value of a specified column. It is useful for improving query performance by reducing data skew and enabling efficient sampling.

**Q. How does Hive handle schema evolution?**

Hive supports schema evolution, allowing you to alter the structure of existing tables without losing data. You can add columns, change column data types, or rename columns using ALTER TABLE statements. Hive also supports backward compatibility for reading data with different schema versions.

**Q. Explain the different file formats supported by Hive.**

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Hive supports various file formats for data storage, including:

TextFile: Plain text files with each line representing a record.

SequenceFile: Binary file format optimized for storing key-value pairs.

ORC (Optimized Row Columnar): Columnar storage format providing improved query performance and compression.

Parquet: Columnar storage format similar to ORC, offering efficient compression and query performance.

**Q. What is ACID compliance, and does Hive support it?**  ACID (Atomicity, Consistency, Isolation, Durability) compliance ensures data integrity and consistency in database transactions. Hive supports ACID compliance starting from version 0.14 through the use of transactions and corresponding storage formats like ORC and ACID-enabled tables.

**Q. Explain the concept of SerDe in Hive.**

SerDe (Serializer/Deserializer) is a mechanism in Hive that facilitates the serialization and deserialization of data between Hive tables and Hadoop's internal representation. It allows Hive to work with various data formats by defining how data is serialized into storage and deserialized for querying.

Q**. What are the different join types supported by Hive?**

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Hive supports the following join types:

Inner Join: Returns only the rows that have matching values in both tables.

Left Outer Join: Returns all rows from the left table and matching rows from the right table.

Right Outer Join: Returns all rows from the right table and matching rows from the left table.

Full Outer Join: Returns all rows when there is a match in either table.

Cross Join: Returns the Cartesian product of two tables (every row from the first table is combined with every row from the second table).

**Q. Explain the difference between SORT BY and ORDER BY in Hive.**

ORDER BY: Orders the result set globally across all reducers, which requires shuffling and sorting all data before returning the final result. It is suitable for small result sets but can be inefficient for large datasets.

SORT BY: Sorts the output of each reducer independently before producing the final result. It avoids the need for shuffling all data to a single reducer, making it more efficient for large datasets but may not guarantee a globally sorted result.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

**Q. What is dynamic partition pruning, and how does it improve query performance?**

Dynamic partition pruning is a feature in Hive that optimizes query performance by selectively pruning partitions based on the query predicates at runtime. It analyzes the query filters and eliminates unnecessary partitions from being scanned, reducing the amount of data processed during query execution.

**Q. Explain the use cases for bucketing in Hive.**

Bucketing, or clustering, is beneficial in the following scenarios:

Improving query performance: Bucketing can reduce data skew and improve query performance by organizing related data into fixed-sized buckets, enabling efficient sampling and join operations.

Data skew mitigation: Bucketing helps distribute data evenly across buckets, reducing data skew and improving parallelism in query processing.

**Q. How does Hive handle skew join?** Hive handles skew join, which occurs when one or more keys have disproportionately large numbers of associated records, by performing the following optimizations:

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Map-side join: Hive tries to avoid skew join by performing a map-side join for small tables, where data can be replicated to all mappers.

Map join hint: Developers can use the MAPJOIN hint to force Hive to perform a map-side join for specific queries, which can mitigate the impact of skew join.

**Q. What is vectorization in Hive, and how does it improve query performance?**

Vectorization is a performance optimization technique in Hive that processes multiple rows at once using vectorized operations instead of row-by-row processing. It reduces CPU overhead and improves query performance by leveraging modern CPU architectures and instruction sets.

**Q. How does Hive integrate with external systems like Spark or HBase?**

Hive integrates with external systems like Spark or HBase through connectors or storage handlers. For example:

Hive on Spark: Allows Hive queries to run on Spark, leveraging Spark's in-memory processing capabilities for improved performance.

HBase integration: Hive provides HBase storage handler to access data stored in HBase tables, enabling seamless integration with Hive's querying capabilities.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

# Advanced Hive interview questions:

**Q. Explain Hive's cost-based optimizer (CBO) and its significance.**   The cost-based optimizer (CBO) in Hive is a query optimization framework that estimates the cost of executing different query plans and selects the most efficient one. It uses statistics such as column and table-level metadata, data distribution, and data size to estimate the cost of each query plan. CBO improves query performance by selecting optimal join strategies, join order, and access paths, resulting in faster query execution.

**Q. What are the benefits and limitations of using LLAP (Live Long and Process) with Hive?** Benefits of using LLAP with Hive:

   Improved query performance: LLAP maintains long-lived daemons that cache and process data in memory, reducing query latency.

   Interactive querying: LLAP enables sub-second query response times, making Hive suitable for interactive data exploration and ad-hoc analysis.

   Enhanced concurrency: LLAP supports concurrent query execution, allowing multiple users to execute queries simultaneously without contention.

   Limitations:

   Resource overhead: LLAP daemons consume memory and CPU resources, requiring careful resource management in multi-tenant environments.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

Configuration complexity: Configuring and tuning LLAP requires expertise to optimize resource utilization and performance.

**Q. Explain Hive transactions and ACID compliance in detail.**

Hive transactions provide Atomicity, Consistency, Isolation, and Durability (ACID) properties for managing concurrent data modifications in Hive tables. ACID compliance ensures that transactions are executed reliably and consistently, even in the presence of concurrent updates and failures. Hive transactions support the following operations:

Insert, update, delete: Hive supports transactional data manipulation operations on ACID-enabled tables.

Multi-statement transactions: Hive allows multiple data manipulation statements to be grouped into a single transaction.

Snapshot isolation: Hive transactions provide snapshot isolation to ensure that each transaction sees a consistent view of the data at the time of transaction start.

**Q. Explain how Hive supports vectorization and its impact on query performance.**
Vectorization in Hive is a query optimization technique that processes multiple rows at once using vectorized operations instead of row-by-row processing. It improves query performance by reducing CPU overhead and memory consumption associated with processing individual rows. Vectorized query execution in Hive operates on batches of rows, leveraging modern CPU architectures and instruction sets for enhanced parallelism and efficiency. As a result, vectorization can significantly accelerate query performance, especially for CPU-bound workloads involving complex analytical queries.

**Q. What are the considerations for optimizing Hive queries for performance?**
Optimizing Hive queries for performance involves several considerations, including:

  Data partitioning and bucketing: Organize data into partitions and buckets to optimize query performance and improve data locality.

  Predicate pushdown: Push filters and predicates down to storage-level operations to minimize data scanned during query execution.

  Join optimization: Choose appropriate join strategies (e.g., broadcast join, map join) based on the size and distribution of data.

  Parallelism: Tune parallelism settings (e.g., number of reducers, map tasks) to maximize resource utilization and minimize query execution time.

  Indexing: Use indexing for faster data retrieval, especially for point queries and range scans.

  Data skew handling: Address data skew issues through techniques such as skew join optimization, data skew detection, and redistribution.

**Q. Explain the use cases for UDFs (User-Defined Functions) and UDAFs (User-Defined Aggregate Functions) in Hive.**

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn

UDFs: User-Defined Functions in Hive allow developers to extend Hive's functionality by defining custom scalar functions. UDFs are useful for implementing complex transformations, data cleansing, and enrichment logic within Hive queries.

 UDAFs: User-Defined Aggregate Functions in Hive enable developers to define custom aggregate functions for summarizing and aggregating data. UDAFs are beneficial for calculating custom aggregations, such as median, percentile, or weighted averages, that are not supported by built-in aggregate functions.

**Q. Explain the role of Tez and LLAP in improving Hive query performance.**

 Tez: Apache Tez is an alternative execution engine for Hive that provides a more efficient and flexible way to execute Hive queries. Tez executes Hive queries as directed acyclic graphs (DAGs), allowing for dynamic task scheduling, data flow optimization, and resource sharing. Tez improves Hive query performance by reducing query latency, optimizing resource utilization, and enabling more complex query plans.

  LLAP (Live Long and Process): LLAP is a long-running daemon process introduced in Hive that provides caching, data processing, and query execution capabilities in memory. LLAP improves Hive query performance by maintaining cached data in memory, reducing data serialization overhead, and enabling interactive query processing with low-latency responses.

Deepa Vasanthkumar – Medium
Deepa Vasanthkumar -| LinkedIn