

Big Data Questions

General Data Engineering Questions

Data warehouse vs database

Warehouse: use aggregate functions for data integrity as a use case

Database: manipulation and aid in visualization

Difference b/w a python Pandas DF vs Spark DF

A Spark DataFrame and a Pandas DataFrame are both data structures used for processing and analyzing data, but they have some fundamental differences.

1. Distributed computing vs. in-memory computing: Spark DataFrames are distributed across a cluster of machines and can handle much larger datasets than Pandas DataFrames, which are stored in memory on a single machine. This makes Spark DataFrames well-suited for big data processing, while Pandas DataFrames are better suited for smaller datasets.
2. Language: Spark DataFrames are used in Spark, which is primarily implemented in Scala and Java, and also supports Python and R. Pandas DataFrames are used in Python.
3. Speed: Spark DataFrames are optimized for parallel processing and can handle large datasets with high speed, while Pandas DataFrames are optimized for in-memory computing and may slow down when handling large datasets.
4. API: Spark DataFrames have a more restricted API than Pandas DataFrames, with fewer operations available for data manipulation. However, Spark DataFrames are designed to work well with distributed computing and can handle complex processing tasks.
5. Memory management: Spark DataFrames manage memory differently from Pandas DataFrames. Spark uses a memory hierarchy and caching mechanism to optimize data processing across the cluster, while Pandas uses the computer's main memory to store the DataFrame.

In summary, Spark DataFrames are optimized for distributed computing, can handle larger datasets, and have a more restricted API than Pandas DataFrames. Pandas DataFrames are optimized for in-memory computing and provide a wider range of operations for data manipulation. The choice of which one to use depends on the specific use case and the size of the dataset.

Structured vs Unstructured Data

Structured data and unstructured data are two different types of data that differ in their format and level of organization.

Structured data refers to data that is organized in a structured format, typically in a tabular format with well-defined columns and rows. This data is highly organized and follows a specific schema or data model. Examples of structured data include data stored in relational databases,

spreadsheets, and data represented in XML or JSON formats. Structured data is easy to search, analyze, and process with the help of tools such as SQL or data analysis software.

On the other hand, unstructured data refers to data that does not have a specific structure or format. It is often text-heavy and does not follow a consistent pattern or schema. Examples of unstructured data include emails, social media posts, images, videos, and audio files.

Unstructured data is often more difficult to process and analyze than structured data, as it requires advanced NLP (natural language processing) and machine learning techniques to extract meaning and insights from the data.

While structured data is highly organized and easy to work with, it is limited in the type of information it can represent. Unstructured data, while more difficult to work with, often contains valuable information that cannot be captured by structured data alone. As a result, a combination of structured and unstructured data is often used in data analysis and decision-making processes.

star schema vs snowflake

Star schema and snowflake schema are two commonly used data modeling techniques in data warehousing.

Star schema is a type of dimensional data model that represents data in a star-shaped structure, with a central fact table and one or more dimension tables connected to it. The fact table contains the measures or numerical data that are being analyzed, such as sales revenue or customer orders, while the dimension tables provide descriptive information about the data, such as product information, customer information, or location information. In a star schema, the dimension tables are denormalized, meaning that they contain redundant data to improve query performance.

Snowflake schema, on the other hand, is a variation of the star schema that normalizes the dimension tables to reduce redundancy and improve data consistency. In a snowflake schema, each dimension table is split into multiple related tables, with each table representing a specific level or hierarchy of the dimension. This results in a more complex and normalized structure, with more tables and relationships than a star schema.

The choice between star schema and snowflake schema depends on the specific needs of the data warehouse and the query performance requirements. Star schema is often preferred when query performance is a priority, as it allows for faster and simpler queries with fewer joins. Snowflake schema is preferred when data consistency and normalization are more important, as it allows for better data management and fewer data anomalies.

In summary, star schema and snowflake schema are both effective data modeling techniques for data warehousing, with different strengths and weaknesses. The choice between them depends on the specific requirements of the data warehouse and the data analysis needs.

- SQL
 - select, filter, order
 - joins
 - grouping
 - de-duplication (!!)
- Coding
 - can you create a class in some language: can you code at all, is the code you write clean and clear
 - simple algorithmic puzzles
 - data wrangling Python or a similar language

SSIS

ETL lifecycle (i.e. benefits of streaming over batch/vice versa, or what a staging area is for)

SQL Specific

Can you explain the differences between inner join, left join, right join, and full outer join?

Inner join returns only the matching rows between two tables, left join returns all rows from the left table and matching rows from the right table, right join returns all rows from the right table and matching rows from the left table, and full outer join returns all rows from both tables.

Can you write a SQL query to find the second-highest salary from a table?

```
SELECT MAX(salary)
FROM employees
WHERE salary < (SELECT MAX(salary) FROM employees)
```

Can you write a SQL query to find the nth-highest salary from a table?

```
SELECT salary
FROM employees
ORDER BY salary DESC
LIMIT 1 OFFSET (N - 1)
```

Can you explain the differences between a clustered index and a non-clustered index?

A clustered index determines the physical order of data in a table, while a non-clustered index creates a separate structure that points to the physical data in the table.

Can you write a SQL query to find all the employees who have joined in the last 30 days?

```
SELECT *  
  
FROM employees  
  
WHERE join_date >= DATEADD(day, -30, GETDATE())
```

Can you write a SQL query to calculate the moving average of a column in a table?

```
SELECT value, AVG(value) OVER (ORDER BY date ROWS BETWEEN 2 PRECEDING AND CURRENT  
ROW) as moving_avg  
  
FROM my_table
```

Can you explain the differences between a primary key and a foreign key?

A primary key uniquely identifies each row in a table, while a foreign key references a primary key in another table.

Can you write a SQL query to find all the employees who have a manager, along with the name of their manager?

```
SELECT e.employee_name, m.employee_name AS manager_name  
  
FROM employees e  
  
JOIN employees m ON e.manager_id = m.employee_id
```

Can you explain the difference between a view and a table?

SQL views and tables are two different types of database objects that serve different purposes.

A table is a fundamental database object that stores data in a structured way, with rows representing individual records and columns representing different attributes or fields of the data. Tables are typically used to store and manage large volumes of data, and they can be queried using SQL to retrieve, update, or delete data.

A view, on the other hand, is a virtual table that does not store any data of its own, but instead derives its data from one or more underlying tables. Views are created using SQL queries that define the data to be included in the view, and they can be used to provide a simplified or customized view of the data that is more convenient or meaningful for specific purposes. Views can be used to hide complex queries or calculations, simplify data access, or restrict access to sensitive data.

The main difference between a view and a table is that a view does not store data directly, but rather provides a dynamic, up-to-date representation of the data in one or more underlying tables. This means that changes made to the underlying tables will be reflected in the view automatically, without the need to update the view explicitly. In contrast, changes made to a table will only affect that table and will not be reflected in any views that reference it.

Another difference is that views can be used to restrict access to certain columns or rows of the data, providing a security mechanism to protect sensitive information. Views can also be used to join multiple tables together and present the data as if it were coming from a single table, providing a convenient way to query data that is spread across multiple tables.

In summary, while both views and tables are important database objects that are used to manage and analyze data, they serve different purposes. Tables are used to store and manage large volumes of data, while views are used to provide a simplified or customized view of the data that is more convenient or meaningful for specific purposes.

Can you write a SQL query to find the top N products by revenue, where N is a parameter passed to the query?

```
SELECT product_name, SUM(revenue) AS total_revenue
FROM sales
GROUP BY product_name
ORDER BY total_revenue DESC
LIMIT N
```

What is the difference between DELETE and TRUNCATE?

DELETE and TRUNCATE are both SQL commands used to remove data from a table, but they differ in the way they operate.

DELETE is a DML (Data Manipulation Language) statement used to remove individual rows from a table. DELETE can be used with or without a WHERE clause to specify the conditions under which rows should be deleted. When a DELETE statement is executed, the specified rows are removed from the table, but the table structure and metadata remain unchanged. Deleted rows can be rolled back using the ROLLBACK command, provided that a transaction is in progress.

TRUNCATE, on the other hand, is a DDL (Data Definition Language) statement used to remove all rows from a table. TRUNCATE removes all data from the table, but it does not delete the table structure or metadata. Unlike DELETE, TRUNCATE cannot be rolled back, as it does not log

individual row deletions. TRUNCATE also resets any identity columns or sequences to their initial value.

In terms of performance, TRUNCATE is generally faster than DELETE, especially for large tables, as it does not need to log individual row deletions. TRUNCATE also releases any storage space used by the table back to the database, whereas DELETE leaves the space allocated for future use.

In summary, DELETE and TRUNCATE are both SQL commands used to remove data from a table, but they differ in their operation and impact. DELETE is used to remove individual rows from a table and can be rolled back, while TRUNCATE removes all rows from a table and cannot be rolled back. TRUNCATE is generally faster than DELETE, especially for large tables, and it releases any storage space used by the table back to the database.

What's the difference between WHERE and HAVING?

WHERE and HAVING are both SQL clauses used to filter data, but they differ in the way they operate and the types of data they can filter.

WHERE is used to filter rows based on specific conditions that are applied to individual rows. WHERE is typically used with SELECT, UPDATE, and DELETE statements to restrict the rows that are affected by the statement. The conditions in a WHERE clause can reference columns from the table being queried or other tables joined in the query. The WHERE clause is evaluated for each row in the table, and only rows that satisfy the conditions are returned in the query result.

HAVING, on the other hand, is used to filter groups of rows based on aggregate functions, such as COUNT, SUM, AVG, MIN, or MAX. HAVING is typically used with GROUP BY statements to restrict the groups that are returned by the statement. The conditions in a HAVING clause reference the results of the aggregate functions and are applied to the groups that are created by the GROUP BY clause. The HAVING clause is evaluated after the GROUP BY clause, and only groups that satisfy the conditions are returned in the query result.

In summary, WHERE is used to filter individual rows based on specific conditions, while HAVING is used to filter groups of rows based on aggregate functions. WHERE is evaluated before any grouping is done, while HAVING is evaluated after the grouping is done. WHERE can filter any data in a table, while HAVING can only filter aggregated data.

What is an index in SQL? When would you use an index?

An index in SQL is a database object that is used to improve the performance of data retrieval and manipulation operations, such as SELECT, UPDATE, DELETE, and JOIN statements. An index is created on one or more columns of a table and contains a sorted copy of the column values, along with a pointer to the corresponding rows in the table. This allows the database to quickly

locate and access the rows that match a given search condition, without having to scan the entire table.

An index is useful when there are large amounts of data in a table and queries are frequently performed to search or filter the data based on specific columns or combinations of columns. Without an index, the database must perform a full table scan to locate the desired rows, which can be slow and resource-intensive. With an index, the database can use the index to quickly locate the relevant rows and retrieve them efficiently.

Here are some scenarios where an index can be useful:

- a. Searching large tables: When searching for specific data in a large table, an index can help to quickly locate the rows that match the search condition, reducing the amount of time and resources required for the query.
- b. Joining tables: When joining tables together, an index can be created on the columns used for the join condition to speed up the process of matching rows between the tables.
- c. Sorting data: When sorting data based on specific columns, an index can help to avoid the need for a full table scan and speed up the sorting process.
- d. Constraints: When enforcing unique or primary key constraints, an index can be created on the relevant columns to ensure that the constraint is enforced efficiently.

It is important to note that creating too many indexes or creating indexes on columns that are rarely used in queries can have a negative impact on performance, as it can increase the amount of time and resources required for data manipulation operations. Therefore, it is important to carefully consider the data access patterns and query requirements when deciding which columns to index.

What are aggregate functions in SQL?

Aggregate functions in SQL are built-in functions that perform calculations on a set of values and return a single value as the result. Aggregate functions are used to perform calculations on groups of rows, such as finding the sum, average, minimum, maximum, or count of values in a column.

Here are some commonly used aggregate functions in SQL:

- e. SUM(): Calculates the sum of all values in a column.
- f. AVG(): Calculates the average of all values in a column.
- g. MIN(): Finds the minimum value in a column.
- h. MAX(): Finds the maximum value in a column.
- i. COUNT(): Counts the number of rows in a table or the number of values in a column.

What SQL commands are utilized in ETL?

ETL (Extract, Transform, Load) is a common process used in data warehousing and business intelligence to move data from source systems to a data warehouse or data mart. SQL commands are often used in ETL to perform various tasks, such as extracting data from source systems, transforming the data into the desired format, and loading the data into the target system.

Here are some SQL commands that are commonly used in ETL:

6. **SELECT:** Used to extract data from source systems, such as databases, files, or APIs.
7. **INSERT:** Used to load data into the target system, such as a data warehouse or data mart.
8. **UPDATE:** Used to modify data in the source system or target system as part of the transformation process.
9. **DELETE:** Used to remove data from the source system or target system as part of the transformation process.
10. **JOIN:** Used to combine data from multiple tables or sources during the transformation process.
11. **GROUP BY:** Used to group data by one or more columns during the transformation process.
12. **WHERE:** Used to filter data based on specific conditions during the extraction or transformation process.
13. **CASE:** Used to perform conditional transformations on data, such as converting values from one format to another.
14. **CAST/CONVERT:** Used to convert data types during the transformation process, such as converting string values to numeric values.
15. **MERGE/UPSERT:** Used to perform updates and inserts on target systems based on the data being loaded.

These SQL commands are often used in combination with other tools and technologies, such as ETL software, data integration platforms, and scripting languages, to create complex and efficient ETL processes.

Does JOIN order affect SQL query performance?

Yes, the order in which tables are joined can affect the performance of an SQL query. The order of joins determines the sequence in which the database engine retrieves and filters data from the tables. A poorly chosen join order can result in longer query execution times and increased resource consumption, while an optimal join order can improve query performance and reduce resource usage.

The database engine uses a query optimizer to determine the best join order for a given query. The optimizer considers various factors, such as table sizes, indexes, data distribution, and query complexity, to generate an execution plan that minimizes the cost of retrieving and filtering data from the tables.

However, in some cases, the optimizer may not be able to find the optimal join order, or the query may be too complex to be optimized efficiently. In such cases, specifying a join order manually can help to

improve query performance. By explicitly specifying the join order, the query can avoid unnecessary data scans and joins, and minimize the amount of data that needs to be processed.

It is important to note that the impact of join order on query performance can vary depending on the specific database engine, data volumes, and query complexity. Therefore, it is important to analyze query performance and experiment with different join orders to find the most efficient solution for a given scenario.

How do you change a column name by writing a query in SQL?

```
ALTER TABLE table_name
```

```
ALTER COLUMN old_column_name new_column_name data_type;
```

How do you handle duplicate data in SQL?

1) delete the duplicate data

2) add unique clause to the table

```
ALTER TABLE customers ADD CONSTRAINT uc_email UNIQUE (email);
```

<https://www.interviewquery.com/questions/subscription-overlap>

Python

What is the difference between "is" and "=="?

"==" compares the values of two objects, while "is" compares the memory addresses or identities of two objects. The choice of which operator to use depends on the specific use case and whether you need to compare values or identities.

What is a decorator?

decorator is a special type of function that can be used to modify or extend the behavior of another function, class, or object. Decorators allow you to add functionality to existing code without having to modify the code itself.

How would you perform web scraping in Python?

Web scraping is the process of extracting data from websites using automated tools. In Python, web scraping can be performed using several popular libraries, including BeautifulSoup, Scrapy, and Requests.

Are lookups faster with dictionaries or lists in Python?

In Python, lookups are generally faster with dictionaries than with lists. This is because dictionaries use hash tables to store key-value pairs, which allows for constant-time lookups, regardless of the size of the dictionary. Lists, on the other hand, are implemented as arrays and use sequential indexing, which means that the time required for a lookup increases linearly with the size of the list.

How familiar are you with TensorFlow? Keras? OpenCV? SciPy?

What is the difference between a list and a tuple?

Mutable vs immutable

What is data smoothing and how do you do it?

Data smoothing is a technique used in data analysis and signal processing to remove noise or other irregularities from data. The goal of data smoothing is to obtain a clearer and more accurate representation of the underlying patterns in the data.

What is a cache database? And why would you use one?

A cache database is a type of database that is designed to store frequently accessed data in a fast and efficient manner. The purpose of a cache database is to reduce the response time of frequently used data, and improve the performance of applications that rely on that data.

Cache databases work by storing a copy of the data in memory, rather than on disk. This allows for faster access times, as data can be retrieved directly from memory rather than having to be read from disk. Additionally, cache databases often use advanced algorithms and data structures to optimize performance, such as caching frequently accessed data in a priority queue, or using a least-recently-used (LRU) eviction policy to remove data that hasn't been accessed recently.

There are several reasons why you might want to use a cache database, including:

- a. Improved performance: Cache databases can significantly reduce the response time of frequently accessed data, which can improve the performance of applications that rely on that data.

- b. Scalability: Cache databases can help improve the scalability of applications by reducing the load on the primary database, allowing it to handle more requests.
- c. Availability: Cache databases can improve the availability of data by providing a fast and efficient way to access frequently used data, even in the event of a primary database failure.
- d. Cost savings: By reducing the load on the primary database, cache databases can help reduce the need for expensive hardware upgrades or additional database licenses.

Some common examples of cache databases include Redis, Memcached, and Apache Ignite. These databases are often used in conjunction with other databases, such as MySQL or PostgreSQL, to provide fast and efficient access to frequently used data.

In summary, a cache database is a type of database that stores frequently accessed data in memory, rather than on disk, in order to improve the performance of applications that rely on that data. By reducing response times, improving scalability and availability, and potentially reducing costs, cache databases can provide significant benefits for applications that require fast and efficient data access.

What are some primitive data structures in Python? What are some user-defined data structures?

In Python, there are several built-in primitive data structures, including:

- e. Integers: used to represent whole numbers, such as 1, 2, and 3.
- f. Floats: used to represent decimal numbers, such as 3.14 and 2.71828.
- g. Booleans: used to represent True or False values.
- h. Strings: used to represent text, such as "hello world" or "Python is awesome!".
- i. None: used to represent the absence of a value.

In addition to these primitive data structures, Python also supports several user-defined data structures, including:

- j. Lists: ordered collections of objects, which can contain objects of different types. Lists are defined using square brackets, such as [1, 2, 3].
- k. Tuples: ordered collections of objects, similar to lists, but are immutable (cannot be modified after creation). Tuples are defined using parentheses, such as (1, 2, 3).
- l. Sets: unordered collections of unique objects, used to perform mathematical set operations, such as union, intersection, and difference. Sets are defined using curly braces, such as {1, 2, 3}.
- m. Dictionaries: unordered collections of key-value pairs, used to map keys to values. Dictionaries are defined using curly braces, with each key-value pair separated by a colon, such as {'apple': 1, 'banana': 2, 'orange': 3}.

In addition to these built-in data structures, Python also supports user-defined classes, which can be used to create custom data structures with specific properties and methods. These classes can be used to define more complex data structures, such as trees, graphs, and queues.

Database Design

What are the features of a physical data model?

A physical data model is a representation of a database's physical structure, including tables, columns, indexes, and relationships. The physical data model is used to design the actual database schema and is used by database management systems to store and retrieve data.

The key features of a physical data model include:

- n. **Tables:** Physical data models include a detailed representation of each table in the database, including the table name, column names, and data types. The physical data model also specifies the relationships between tables, such as primary key-foreign key relationships.
- o. **Indexes:** Physical data models include information about the indexes that will be used to optimize data retrieval. Indexes can be created on one or more columns in a table and are used to speed up queries by allowing the database to quickly locate the required data.
- p. **Constraints:** Physical data models define the constraints that will be applied to the database, such as primary key constraints, unique constraints, and foreign key constraints. These constraints help to ensure data integrity and prevent data inconsistencies.
- q. **Data types:** Physical data models specify the data types that will be used for each column in the database. This includes data types for numbers, strings, dates, and other types of data.
- r. **Performance considerations:** Physical data models take into account performance considerations, such as disk space usage, query performance, and database maintenance requirements.
- s. **Storage considerations:** Physical data models specify how data will be stored on disk, including the file structures and storage mechanisms used by the database management system.

In summary, the features of a physical data model include tables, indexes, constraints, data types, performance considerations, and storage considerations. These features help to ensure that the database is designed and optimized for efficient data storage, retrieval, and maintenance.

What database relationships do you know?

In database design, there are three primary types of relationships between tables:

- t. One-to-One (1:1) relationship: In a one-to-one relationship, each record in one table corresponds to exactly one record in another table, and vice versa. This type of relationship is often used when one piece of data is dependent on another, but only in a one-to-one manner.
- u. One-to-Many (1:N) relationship: In a one-to-many relationship, each record in one table can correspond to many records in another table, but each record in the second table can correspond to only one record in the first table. This type of relationship is often used when one table contains related data that needs to be separated out for normalization purposes.
- v. Many-to-Many (N:M) relationship: In a many-to-many relationship, each record in one table can correspond to many records in another table, and vice versa. This type of relationship requires a third table, known as a junction table, to link the two tables together. The junction table contains foreign keys that reference the primary keys of the two related tables.

It's important to note that these relationships are not mutually exclusive, and a database can contain all three types of relationships. Additionally, there are variations and combinations of these relationships, such as recursive relationships and self-referencing relationships, which allow for more complex data modeling. Overall, understanding and correctly implementing database relationships is essential to building efficient and maintainable databases.

How would you handle data loss during a migration?

Handling data loss during a migration is a critical aspect of ensuring a successful migration. Here are some steps you can take to handle data loss during a migration:

- w. Backup the data: Before starting the migration process, ensure that you have a backup of all the data that is being migrated. This backup will act as a fail-safe in case of any data loss during the migration.
- x. Test the migration process: Prior to the actual migration, test the migration process on a small subset of data to identify any potential issues that could lead to data loss.
- y. Monitor the migration process: During the migration process, closely monitor the migration progress to identify any potential data loss issues. Use logs and alerts to stay informed of any issues that arise.
- z. Rollback plan: Have a rollback plan in place to revert the migration in case of data loss. This plan should include a timeline for reverting the migration, as well as procedures for ensuring data integrity during the rollback process.
- aa. Data recovery plan: In case of data loss, have a data recovery plan in place. This plan should include steps for recovering lost data, as well as procedures for testing the recovered data to ensure its integrity.

- bb. Post-migration validation: After the migration is complete, validate the migrated data to ensure that all data has been successfully migrated and that there is no data loss.

By following these steps, you can minimize the risk of data loss during a migration and ensure a successful migration process.

What are the three types of data models?

There are three main types of data models:

- cc. Conceptual data model: This type of data model describes the high-level relationships between entities in a business or organization. It focuses on identifying the major entities, attributes, and relationships between them, and does not contain specific details about how the data will be implemented in a physical database. Conceptual data models are used to communicate and document the business requirements for a database project.
- dd. Logical data model: This type of data model provides a detailed representation of the data requirements of a business or organization, and shows the relationships between entities, attributes, and data types. Logical data models are typically used by data analysts and database designers to design and implement a database schema. Logical data models are independent of any specific database management system, and can be used to design databases for various platforms.
- ee. Physical data model: This type of data model represents how data is physically stored in a database management system. It includes details such as the database schema, indexes, data types, and storage parameters. Physical data models are used by database administrators and developers to implement the logical data model in a specific database management system. They provide a blueprint for building and maintaining the database, and help to ensure that the database is efficient and performs well.

In summary, conceptual data models provide a high-level view of the data requirements, logical data models provide a detailed view of the data relationships and attributes, and physical data models provide a detailed view of the database schema and storage requirements. Together, these data models help to ensure that the database meets the business requirements, is well-designed, and is efficiently implemented.

What is normalization? Denormalization?

Normalization and denormalization are two related concepts in database design that describe different approaches to organizing data in tables.

Normalization is the process of organizing data in a database so that it is structured and optimized for efficient storage and retrieval. The goal of normalization is to minimize data redundancy and ensure data integrity by breaking down tables into smaller, more specialized tables with clearly defined relationships. Normalization involves applying a set of rules called normal forms to ensure that each table has a single, primary key and that each column in the table depends only on the primary key.

There are several levels of normalization, with each level building on the previous level to further refine the structure of the database. The most common normal forms are first normal form (1NF), second normal form (2NF), and third normal form (3NF).

Denormalization, on the other hand, is the process of intentionally adding redundancy to a database schema to improve query performance or simplify data access. Denormalization involves breaking some of the rules of normalization by adding redundant columns or duplicating data across tables. The goal of denormalization is to optimize performance by reducing the number of joins required to retrieve data from the database.

Denormalization can be useful in certain scenarios, such as when dealing with large, complex data sets or when querying data frequently in a read-heavy application. However, it can also increase data redundancy and complexity, and can make it more difficult to maintain data integrity over time.

Overall, normalization and denormalization are two different approaches to organizing data in a database, each with its own advantages and disadvantages. Normalization is typically used to ensure data integrity and minimize redundancy, while denormalization is used to optimize query performance or simplify data access. The choice between normalization and denormalization depends on the specific needs and requirements of the application and the data being stored.

What is the difference between 1NF, 2NF, 3NF?

Normalization is a process of organizing data in a database to minimize data redundancy and ensure data integrity. The process involves breaking down tables into smaller, more specialized tables with clearly defined relationships. There are several levels of normalization, each building on the previous level to further refine the structure of the database.

Here are the three most common normal forms and their defining characteristics:

ff. First Normal Form (1NF):

- A table is in 1NF if it contains no repeating groups or arrays of data.
- Each row must be unique, identified by a primary key.
- Each column must contain atomic values, meaning they cannot be further divided into smaller components (NO LIST).
- Example of a table that is not in 1NF: a table that stores multiple phone numbers for each customer in a single column.

gg. Second Normal Form (2NF):

- A table is in 2NF if it is in 1NF and each non-key column is fully dependent on the primary key (NO PARTIAL DEPENDENCIES, seen with composite keys).

- This means that each column in the table must depend on the entire primary key, not just a part of it.
- Example of a table that is not in 2NF: a table that stores order details with order ID, product ID, and product price in a single table.
- hh. Third Normal Form (3NF):
- A table is in 3NF if it is in 2NF and there are no transitive dependencies.
- A transitive dependency occurs when a non-key column is dependent on another non-key column.
- All non-key columns must depend only on the primary key or other non-key columns in the table.
- Example of a table that is not in 3NF: a table that stores employee details with department name and location, where the department name determines the location.

Normalization helps to ensure that data is organized efficiently and accurately, making it easier to maintain and query. However, achieving higher levels of normalization can sometimes result in a more complex database structure, which can lead to longer query times and reduced performance. The decision of which normal form to use ultimately depends on the specific needs and requirements of the application.

What are some things to avoid when building a data model?

When building a data model, there are several things to avoid to ensure that the model is efficient, scalable, and easy to maintain. Here are some common pitfalls to avoid when building a data model:

- ii. Over-normalization: Normalization is an important aspect of data modeling, but over-normalizing can lead to a complex, difficult-to-maintain data model. Avoid breaking tables down into too many small, specialized tables that require many joins to retrieve data.
- jj. Under-normalization: On the other hand, under-normalizing can lead to data redundancy and inconsistency. Avoid duplicating data unnecessarily across tables, which can lead to data inconsistencies and update anomalies.
- kk. Lack of consistency: Be consistent in naming conventions, data types, and other aspects of the data model to ensure that the model is easy to understand and maintain. Inconsistencies can make it difficult to query data and can lead to errors and inconsistencies over time.
- ll. Ignoring data types and constraints: Be sure to use appropriate data types and constraints to ensure data integrity and accuracy. For example, use a numeric data type for numeric data, and use constraints such as NOT NULL and UNIQUE to ensure that data is valid and consistent.
- mm. Poor performance: Design the data model with query performance in mind. Avoid using too many joins, and consider using denormalization or other performance optimization techniques to ensure that queries are fast and efficient.

- nn. Lack of scalability: Design the data model with future growth in mind. Consider the potential for increased data volume, new data sources, and changes to the business model, and design the data model to accommodate these changes.
- oo. Lack of documentation: Document the data model thoroughly to ensure that it is easy to understand and maintain. Include information such as table and column descriptions, relationships, data types, and constraints.

By avoiding these common pitfalls, you can ensure that your data model is efficient, scalable, and easy to maintain over time.

Why are NoSQL databases more useful than relational databases?

NoSQL (Not Only SQL) databases and relational databases serve different purposes and have different strengths and weaknesses. NoSQL databases are often better suited for certain types of data and applications than relational databases. Here are some reasons why NoSQL databases may be more useful than relational databases in some scenarios:

- pp. Scalability: NoSQL databases are designed to scale horizontally across multiple servers, making them ideal for handling large, distributed data sets. Relational databases are more difficult to scale horizontally, and may require complex sharding or replication techniques to achieve similar levels of scalability.
- qq. Flexibility: NoSQL databases are often schemaless or have flexible schema structures, which allows for greater flexibility in handling unstructured or semi-structured data. Relational databases have rigid, predefined schemas that can make it difficult to handle data with varying structures or changing requirements.
- rr. Performance: NoSQL databases can be optimized for specific use cases and workloads, such as high-volume, low-latency data processing. Relational databases may struggle to handle certain types of queries or data-intensive workloads, especially at scale.
- ss. Cost: NoSQL databases can be more cost-effective than relational databases for certain types of workloads, especially when it comes to scaling horizontally. NoSQL databases often use commodity hardware and open-source software, which can be less expensive than the proprietary software and hardware required for large-scale relational databases.
- tt. Availability: NoSQL databases are often designed for high availability and fault tolerance, making them ideal for mission-critical applications. Relational databases may require more complex architectures to achieve similar levels of availability and durability.

It's important to note that NoSQL databases are not a one-size-fits-all solution, and may not be appropriate for all types of data or applications. Relational databases still have their place in many applications, especially when dealing with structured data and complex relationships. The choice of database technology ultimately depends on the specific needs and requirements of the application.

Data Engineering Case Study

How would you design a relational database of customer data?

Designing a relational database of customer data involves identifying the entities, attributes, and relationships involved in the data, and creating a schema that reflects these relationships. Here's an example of how you might design a relational database of customer data:

Entities:

- Customers
- Orders
- Products

Attributes:

- Customers: customer_id, first_name, last_name, email, phone_number, address, city, state, zip_code
- Orders: order_id, customer_id, order_date, order_total
- Products: product_id, product_name, product_description, price

Relationships:

- Customers can place many orders, but each order belongs to only one customer.
- Orders can contain many products, and each product can appear in many orders.

How would this design process change for customer data? What factors would you need to consider in Step 1?

With what database design patterns do you have the most experience

Your task is working on building a notification system for a Reddit-style app. How would the backend and data model look?

Building a notification system for a Reddit-style app requires a backend architecture and data model that can efficiently process and manage large volumes of user-generated content and notifications. Here's an example of how you might design the backend and data model for such a system:

Backend Architecture:

- The backend architecture would likely be built on a microservices architecture, with different services handling different aspects of the notification system.
- One service might handle the creation and management of notifications, while another service might handle the delivery of notifications to users.
- The system might also include a message queue or event streaming platform to manage the flow of notifications between services.

Data Model:

- The data model would include tables to store user information, post information, and notification information.
- The user table would store user information such as usernames, email addresses, and notification preferences.
- The post table would store information about each post, such as the post ID, the user who created the post, and the content of the post.
- The notification table would store information about each notification, such as the notification ID, the user who should receive the notification, the type of notification, and the post or comment associated with the notification.
- The system might also include additional tables to store information about comments, upvotes, and other user-generated content.

Here's an example of how you might create the notification table:

You have two ETL jobs that feed into a single production table each day. What problems might this cause?

If two ETL jobs feed into a single production table each day, there are several potential problems that can arise, including:

- uu. Data conflicts: If the two ETL jobs are updating the same rows in the production table, conflicts can arise if the updates are not synchronized. For example, if one ETL job updates a row to reflect new information, while the other ETL job updates the same row to reflect different information, it can lead to inconsistencies and data errors.
- vv. Overwriting data: If the two ETL jobs are not properly coordinated, it can result in data being overwritten or deleted unintentionally. For example, if one ETL job deletes a row while the other ETL job is updating the same row, it can result in data loss and inconsistencies.
- ww. Performance issues: If the two ETL jobs are running simultaneously or overlapping, it can lead to performance issues, including slow queries, database locks, and resource contention.
- xx. Incomplete data: If one of the ETL jobs fails or encounters an error, it can result in incomplete data being loaded into the production table, leading to data inconsistencies and errors.

To address these issues, it's important to coordinate the ETL jobs to ensure that they are properly synchronized and that there are no conflicts or inconsistencies in the data. This can be achieved through various methods, including using a synchronization framework or implementing transactional processing to ensure that updates are atomic and consistent. Additionally, monitoring and logging tools can be used to track the status of the ETL jobs and identify any errors or issues that may arise.

What's the difference between ETL and ELT?

ETL (Extract, Transform, Load) and ELT (Extract, Load, Transform) are both data integration processes used to move data from source systems to a target system. However, they differ in the order in which they perform the transformation step:

- yy. ETL: In ETL, data is first extracted from source systems and stored in a staging area. Then, the data is transformed and manipulated into the desired format, and finally loaded into the target system. ETL requires a separate data integration tool or middleware to transform and manipulate the data.
- zz. ELT: In ELT, data is extracted from source systems and directly loaded into the target system. Once the data is loaded, it is transformed and manipulated using the native processing capabilities of the target system. ELT leverages the processing power of the target system and eliminates the need for a separate data integration tool.

The primary difference between ETL and ELT is the location where data transformation takes place. In ETL, data transformation occurs in a separate system or middleware, while in ELT, data transformation occurs in the target system. ELT can be faster and more efficient than ETL when the target system has robust processing capabilities, as it reduces the need for additional infrastructure and eliminates the latency associated with moving data to a separate staging area.

However, ETL may be a better choice when dealing with complex transformations or large volumes of data, as it provides greater control over the transformation process and can be more easily scaled to handle high volumes of data. Ultimately, the choice between ETL and ELT depends on the specific needs and requirements of the organization, as well as the capabilities of the target system.

What is an initial load in ETL? What about full load?

An initial load and a full load are both types of data loads in ETL (Extract, Transform, Load) processes.

- aaa. Initial Load: An initial load is the first time that data is loaded into a target system. It typically involves moving a large amount of data from source systems to a

target system, and may require additional processing and data validation. The goal of the initial load is to ensure that the target system contains all the necessary data for ongoing operations.

- bbb. Full Load: A full load is a complete refresh of data in the target system. It involves replacing all existing data in the target system with new data from the source systems. This type of load is typically performed periodically to ensure that the data in the target system is up-to-date and accurate. Full loads can be resource-intensive and time-consuming, especially when dealing with large volumes of data.

Both initial loads and full loads are important steps in the ETL process. The initial load establishes the baseline data in the target system, while the full load ensures that the target system is updated with the latest data. After the initial and full loads are complete, ongoing incremental loads can be performed to update the target system with new or changed data from the source systems.

Incremental loads involve only moving new or changed data from the source systems to the target system since the last load, reducing the amount of data that needs to be processed and improving the efficiency of the ETL process. The choice of load type (initial, full, or incremental) depends on the specific needs and requirements of the organization, as well as the size and complexity of the data being loaded.

With what ETL tools are you most familiar?

What are partitions? Why might you increase the number of partitions?

Partitions in a database refer to a way of dividing a large table into smaller, more manageable pieces based on a partition key. Each partition contains a subset of the data and can be stored on different physical devices or servers. Partitions can improve query performance, reduce data processing time, and increase overall scalability of the database.

Here are some reasons why you might increase the number of partitions:

- ccc. Improved performance: By dividing the data into smaller subsets, queries can be processed faster and more efficiently. Queries that only require data from one or a few partitions can be executed more quickly than queries that require data from the entire table.
- ddd. Reduced resource consumption: By processing data in smaller batches, the overall resource consumption of the database can be reduced. This can lead to faster processing times and reduced costs for the organization.
- eee. Increased scalability: As the size of the database grows, adding partitions can help ensure that the database remains scalable and can handle increasing volumes of data. By distributing data across multiple servers, partitions can help ensure that the database can handle more queries and users.

- fff. Improved availability: By storing partitions on different servers or devices, partitions can help ensure that the data is available even in the event of a hardware failure or other outage. If one partition is unavailable, other partitions can continue to function, ensuring that the data remains accessible to users.

In summary, partitions are a way of dividing large tables into smaller, more manageable pieces based on a partition key. Increasing the number of partitions can improve performance, reduce resource consumption, increase scalability, and improve availability of the database. The optimal number of partitions depends on the specific needs and requirements of the organization and the nature of the data being stored.

What are database snapshots? What's their importance?

Database snapshots are a read-only copy of a database that captures the state of the database at a specific point in time. Snapshots can be used for a variety of purposes, including backup and recovery, reporting, and data analysis.

Here are some key benefits of database snapshots:

- ggg. Faster backup and recovery: Snapshots can be used to create a backup of the database without the need to shut down the database. This can help reduce downtime and improve recovery time in the event of a failure or data loss.
- hhh. Improved data analysis: Snapshots provide a static view of the database at a specific point in time, allowing for accurate data analysis and reporting. This can be especially useful for historical analysis or compliance reporting.
- iii. Improved database performance: Snapshots can be used to offload reporting or other read-only queries from the production database, improving overall performance and reducing the risk of impacting production operations.
- jjj. Simplified testing and development: Snapshots can be used to create a copy of the database for testing and development purposes, without impacting the production database. This can help reduce the risk of introducing errors or issues into the production environment.

Database snapshots are important because they provide a way to capture the state of the database at a specific point in time, allowing for accurate backup and recovery, improved data analysis, improved database performance, and simplified testing and development. By providing a read-only copy of the database, snapshots also help improve overall database security and reduce the risk of data loss or corruption.

What are views in ETL? What is used to build them?

Views in ETL (Extract, Transform, Load) are virtual tables that are used to present data from one or more source tables in a specific format. Views are created using a SELECT statement, and they allow users to query the data in a more structured and meaningful way.

Here are some benefits of using views in ETL:

- kkk. Simplify data access: Views can simplify data access by providing a pre-defined query that users can use to access data. This can be especially useful when dealing with complex queries or large amounts of data.
- III. Enhance security: Views can be used to control access to sensitive data by limiting the columns or rows that users can access. This can help improve overall database security and reduce the risk of data breaches or unauthorized access.
- mmm. Improve performance: Views can improve performance by reducing the amount of data that needs to be processed or by pre-aggregating data in a way that is optimized for a specific query.
- nnn. Simplify data modeling: Views can be used to simplify data modeling by presenting data in a way that is easier to understand and manipulate. This can help reduce the complexity of the ETL process and make it easier to maintain over time.

Views can be built using a SQL query that retrieves data from one or more source tables. The query can include filters, joins, and other transformations to present the data in a specific format. Once the view is created, users can query the view as if it were a regular table, using SELECT statements to retrieve data.

Overall, views are an important tool in ETL that can simplify data access, enhance security, improve performance, and simplify data modeling. By presenting data in a meaningful way, views can help make the ETL process more efficient and effective.

What could be potential bottlenecks in the ETL process?

There are several potential bottlenecks that can arise in the ETL (Extract, Transform, Load) process, including:

- ooo. Data volume: The amount of data being processed can be a bottleneck, especially if the data volume is large. This can impact performance, increase processing time, and strain system resources.
- ppp. Data quality: Poor data quality, such as missing or incomplete data, can lead to errors and processing delays. This can also impact downstream applications that rely on the data.
- qqq. Complexity of transformations: Complex data transformations can be a bottleneck, especially if they require significant processing or multiple steps. This can increase processing time and resource consumption, impacting overall system performance.

- rrr. System resource limitations: System resource limitations, such as insufficient disk space, memory, or processing power, can cause the ETL process to slow down or fail.
- sss. Data source availability: The availability of data sources can also be a bottleneck, especially if the data sources are slow or unreliable. This can impact the overall processing time and delay the delivery of the final data output.
- ttt. Network latency: Network latency can impact the speed at which data is transferred between systems, leading to processing delays and slower overall performance.
- uuu. Concurrency: Multiple ETL processes running concurrently can compete for resources and impact overall system performance. This can be especially challenging if the ETL processes are working with the same data or data sources.

To mitigate these potential bottlenecks, it's important to optimize the ETL process by identifying and addressing any areas of weakness. This can include improving data quality, simplifying data transformations, optimizing system resources, and coordinating concurrent ETL processes to avoid resource contention. By optimizing the ETL process, organizations can improve data processing speed, accuracy, and overall performance.

How would you triage an ETL failure?

Triage is the process of identifying and prioritizing issues in order to resolve them in a timely manner. When an ETL (Extract, Transform, Load) process fails, the following steps can be taken to triage the failure:

- vvv. Identify the error message: The first step is to identify the error message or error code associated with the ETL failure. This can help identify the specific issue that caused the failure.
- www. Check logs and system metrics: Check logs and system metrics to gain insight into the overall health of the system and identify any other issues that may be contributing to the failure.
- xxx. Determine the scope of the issue: Determine the scope of the issue by identifying which processes, tables, or databases are affected. This can help prioritize the issue and determine the appropriate next steps.
- yyy. Assess the impact: Assess the impact of the ETL failure on downstream applications or processes. This can help determine the level of urgency for resolving the issue.
- zzz. Attempt to rerun the process: Attempt to rerun the ETL process to determine if the issue is persistent or intermittent.
- aaaa. Engage the appropriate team members: Engage the appropriate team members, such as developers or database administrators, to help diagnose and resolve the issue.

- bbbb. Prioritize and resolve the issue: Prioritize the issue based on its impact and level of urgency, and work to resolve it as quickly as possible.

By following these steps, organizations can quickly triage ETL failures and take appropriate action to resolve them. This can help ensure that the ETL process is running smoothly and that data is being processed accurately and efficiently.

Describe how to use an operational data store.

An operational data store (ODS) is a database that is designed to support real-time operational reporting and analysis. Here are some steps to using an ODS:

- cccc. Identify the data sources: Identify the data sources that will be used to populate the ODS. These can include transactional systems, customer databases, marketing systems, and other sources of operational data.
- dddd. Design the data model: Design the data model for the ODS, taking into account the data sources and the reporting and analysis requirements. The data model should be designed to support real-time reporting and analysis, with data organized in a way that is optimized for queries and analysis.
- eeee. Extract and transform the data: Extract the data from the source systems and transform it as necessary to meet the requirements of the ODS. This may include data cleansing, data validation, and other transformations to ensure that the data is accurate and consistent.
- ffff. Load the data: Load the transformed data into the ODS, ensuring that it is structured in a way that supports real-time reporting and analysis.
- gggg. Implement real-time reporting and analysis: Implement reporting and analysis tools that allow users to access and analyze the data in real-time. This can include dashboards, ad-hoc reporting tools, and other tools that allow users to explore and analyze the data in meaningful ways.
- hhhh. Monitor and maintain the ODS: Monitor the ODS to ensure that it is performing optimally and that data is being processed accurately and efficiently. This may include monitoring system performance, identifying and resolving issues as they arise, and implementing ongoing maintenance and upgrades as needed.

By following these steps, organizations can use an operational data store to support real-time reporting and analysis, enabling them to make data-driven decisions and respond quickly to changing business conditions. An ODS can provide valuable insights into operational data, helping organizations improve efficiency, optimize processes, and drive business growth.

Big Data Questions

What is Hadoop?

Open-source framework which is used for data manipulation running on clusters.

Components of Hadoop?

Hadoop common: all libraries and utilities that are used in a Hadoop application

HDFS

HDFS stands for Hadoop Distributed File System. It is a distributed file system that is designed to store and manage large volumes of data across multiple servers or nodes in a Hadoop cluster.

HDFS is a key component of the Apache Hadoop ecosystem and provides high-performance access to data that is stored in Hadoop. It is highly fault-tolerant and scalable, making it well-suited for handling big data workloads.

In HDFS, data is stored in blocks and distributed across the cluster. Each block is replicated multiple times for redundancy and fault tolerance. HDFS also provides a framework for processing and analyzing data stored in the system, such as MapReduce and other distributed computing tools.

Overall, HDFS is an essential component of the Hadoop ecosystem and enables efficient and reliable storage and processing of large data sets

YARN

YARN stands for Yet Another Resource Negotiator. It is a cluster management technology in the Hadoop ecosystem that allows multiple data processing engines such as Apache Hadoop MapReduce, Apache Spark, Apache Flink, and others to run on the same cluster. YARN acts as a resource manager and job scheduler, providing a unified platform for managing distributed computing resources across the Hadoop cluster.

YARN separates the resource management and job scheduling functions from the data processing engine, which allows multiple applications to run simultaneously on the same Hadoop cluster. It dynamically allocates resources to applications based on their requirements, ensuring that each application gets the necessary resources it needs to run efficiently. YARN also provides a framework for monitoring and managing the cluster resources, making it easier to manage the overall Hadoop cluster.

Overall, YARN is an important component of the Hadoop ecosystem, allowing multiple data processing engines to run on the same cluster and providing efficient resource management and scheduling capabilities for distributed computing.

MapReduce

MapReduce is a programming model and software framework used for processing large volumes of data in a distributed computing environment. It is a key component of the Apache Hadoop ecosystem and is commonly used for processing big data workloads.

The MapReduce programming model involves two main functions: Map and Reduce. The Map function processes input data and produces a set of key-value pairs, while the Reduce function takes the output of the Map function and produces a set of final output values. Both Map and Reduce functions can be written in any programming language that supports the Hadoop API.

In the MapReduce framework, the data is split into small chunks and distributed across a cluster of nodes for processing. Each node processes its portion of the data independently and in parallel. The results from each node are then combined to produce the final output.

MapReduce is a scalable and fault-tolerant framework, making it ideal for processing large volumes of data in parallel across a cluster of nodes. It is commonly used for data-intensive applications such as batch processing, log processing, and data analysis.

Name Node

In Hadoop Distributed File System (HDFS), the NameNode is a critical component that manages the file system namespace and regulates access to files by clients. It stores information about the location of each file in the Hadoop cluster, as well as metadata such as permissions, ownership, and replication factor.

The NameNode is responsible for processing client requests to create, delete, or modify files and directories in the file system. When a client wants to read or write data to a file, it first contacts the NameNode to obtain information about the location of the file blocks and the DataNodes that store them. The NameNode then communicates with the relevant DataNodes to initiate data transfer between the client and the DataNodes.

The NameNode is a single point of failure in the HDFS architecture. Therefore, Hadoop provides mechanisms such as backup and recovery tools to ensure that the NameNode can be restored quickly in the event of a failure. Additionally, newer versions of Hadoop provide an active-standby NameNode configuration to provide high availability for the HDFS NameNode service.

Streaming

Streaming is a data processing technique in which data is processed as a continuous stream of events, rather than being processed in batches. Streaming is commonly used for real-time data processing and analysis, such as in the case of data generated by social media, IoT devices, or financial transactions.

In a streaming data processing system, data is processed incrementally as it arrives, rather than being stored and processed later in batches. This allows for real-time analysis and decision-making based on the data.

Apache Kafka, Apache Flink, and Apache Spark Streaming are examples of popular open-source streaming data processing frameworks. These frameworks provide libraries and tools for processing and analyzing data streams in real-time.

One of the key advantages of streaming is that it allows for faster processing of data and quicker response times for real-time applications. Streaming data processing systems also enable near real-time data processing, allowing businesses to make decisions based on the most up-to-date information. However, streaming data processing systems can be more complex to implement and manage than batch processing systems.

4 v's of data

The "4 Vs of Data" is a popular concept used to describe the characteristics of big data. The 4 Vs are:

- iiii. Volume: Volume refers to the amount of data that is being generated and collected. With the explosion of digital data, there is an ever-increasing volume of data that needs to be processed and analyzed.
- jjjj. Velocity: Velocity refers to the speed at which data is generated and collected. With the rise of real-time data sources such as social media and IoT devices, data is being generated at an unprecedented velocity, requiring real-time or near-real-time processing and analysis.
- kkkk. Variety: Variety refers to the different types of data that are being generated, such as structured, unstructured, and semi-structured data. Big data often includes a mix of different types of data, which presents challenges for processing and analysis.
- IIII. Veracity: Veracity refers to the accuracy and reliability of the data. Big data can include a significant amount of noise, errors, and inconsistencies, which can impact the accuracy and reliability of analysis results. Ensuring the veracity of data is a critical challenge for big data processing and analysis.

Overall, the 4 Vs of Data highlight the unique challenges and opportunities presented by big data, and understanding these characteristics is important for effectively processing and analyzing large volumes of data.

Block and block scanner

In the Hadoop Distributed File System (HDFS), data is stored in the form of blocks. A block is a fixed-size segment of a file that is used to distribute data across the Hadoop cluster. The default block size in HDFS is 128MB, but it can be configured to a different size based on the needs of the application.

The block scanner is a service in HDFS that is responsible for verifying the integrity of data stored in Hadoop. It checks the health of each block in the file system by scanning the block contents and detecting any corruption or errors. If any errors are found, the block scanner takes corrective actions, such as copying the data to a new location or replicating the data to ensure data availability and integrity.

The block scanner runs in the background on each HDFS data node, continuously checking the health of data blocks. The block scanner works together with the NameNode to ensure that the data stored in HDFS is consistent and reliable.

Overall, the block scanner is an important component of the Hadoop ecosystem, ensuring data availability and reliability by verifying the health of data blocks in HDFS.

How does the block scanner deal with corrupted

When the block scanner detects a corrupted block in HDFS, it takes corrective actions to ensure the availability and integrity of the data. There are several ways that the block scanner can deal with corrupted files, including:

- mmmm. Marking the block as corrupt: The block scanner can mark the corrupted block as corrupt and report it to the NameNode. The NameNode then ensures that the corrupted block is not used for further processing, such as replication or recovery.
- nnnn. Re-replicating the block: If the corrupted block is part of a replicated set, the block scanner can re-replicate the block to ensure data availability. The block scanner identifies the corrupted replica and replaces it with a new, healthy replica.
- oooo. Copying the data to a new location: In some cases, the block scanner may be able to copy the data from the corrupted block to a new location. This ensures that the data is available for processing while also removing the corrupted block from the system.
- pppp. Removing the corrupted block: In extreme cases where the block cannot be repaired or replaced, the block scanner can remove the corrupted block from the system. This ensures that the corrupted data does not affect the availability or integrity of other data in the system.

Overall, the block scanner takes corrective actions based on the severity and type of corruption detected in the HDFS blocks. These actions ensure that the data in Hadoop is consistent and reliable, even in the face of corrupted files.

How communicate with data node

In Hadoop Distributed File System (HDFS), the NameNode communicates with the DataNodes to manage the storage and retrieval of data. The communication between the NameNode and DataNodes happens through a protocol called the Data Transfer Protocol (DTP).

The NameNode is responsible for managing the metadata of the file system, including the location of each file and the blocks that make up each file. When a client requests access to a file, the NameNode provides the location of the file blocks and the DataNodes that store them.

To retrieve or store data, the client then communicates directly with the DataNodes. The DataNodes are responsible for reading and writing data blocks to disk, and they send block reports to the NameNode to inform it of the state of the blocks they are storing.

The Data Transfer Protocol is used by the NameNode and the DataNodes to communicate and transfer data blocks. When a client needs to read data from a file, it sends a request to the NameNode, which responds with the locations of the blocks that make up the file. The client then establishes a connection with the DataNodes and requests the blocks directly from them.

Similarly, when a client needs to write data to a file, it sends a request to the NameNode, which determines the DataNodes that should store the new blocks. The client then establishes a connection with the DataNodes and sends the new blocks directly to them.

Overall, the communication between the NameNode and DataNodes in HDFS is critical for the proper functioning of the Hadoop ecosystem, ensuring that data is stored and retrieved efficiently and reliably.

HIVE partitions

Hive is a data warehousing tool that is built on top of Hadoop. It provides a SQL-like interface for querying and analyzing data that is stored in Hadoop Distributed File System (HDFS). Hive partitions are a way of organizing data within Hive tables to improve query performance and manageability.

Partitions in Hive are created by dividing the data in a table into smaller, more manageable chunks based on one or more columns. For example, if a table contains sales data, it could be partitioned by date, so that all data for a given day is stored in the same partition. Partitioning allows Hive to skip over irrelevant data during query processing, which can improve query performance.

Partitions can be created when a table is first created, or they can be added later. To create a partitioned table in Hive, the partition columns are specified as part of the table schema. Data can then be loaded into the table using the Hive load command or an external tool such as Apache Sqoop.

Hive supports both static and dynamic partitioning. Static partitioning involves specifying the partition column values explicitly when inserting data into the table. Dynamic partitioning involves automatically creating partitions based on the values of the partition columns.

Overall, Hive partitions are a powerful feature that allows users to improve the performance and manageability of their data in Hadoop. By organizing data into smaller, more manageable chunks, Hive partitions can make it easier to analyze and query large data sets stored in HDFS.

Skewing of Data

In data processing, data skew refers to an imbalance in the distribution of data among partitions or nodes in a distributed system. This can cause some nodes or partitions to be overloaded with data while others are underutilized, which can lead to performance issues and slow processing times.

Data skew can occur due to various reasons, such as uneven data distribution, skewed data values, or poor partitioning strategy. Data skew can be particularly challenging in big data environments, where large volumes of data are being processed and analyzed in a distributed computing environment.

To address data skew, there are several techniques that can be employed, such as:

- qqqq. Sampling: Sampling can be used to gain insight into the distribution of data and identify any skew. This can help in identifying and isolating skewed data and adjusting the processing strategy accordingly.
- rrrr. Repartitioning: Repartitioning involves redistributing data across partitions or nodes to balance the data distribution. This can be achieved by using a more effective partitioning strategy or by repartitioning the data based on the frequency of data access.
- ssss. Skewed joins: Skewed joins involve handling skewed data during join operations. This can be achieved by using techniques such as bucketing or using a distributed hash join.
- tttt. Dynamic partitioning: Dynamic partitioning involves automatically adjusting the partitioning strategy based on the data distribution. This can help in handling dynamic data and prevent skew from occurring.

Overall, addressing data skew is an important consideration in big data processing and analysis. By identifying and addressing skew in the data, performance can be improved, and processing times can be reduced.

Salting

Buckets

Avro vs Parquet

Parallel Distributed computation vs normal computing