# Short Notes

| | | |
|---|---|---|
| ⊚ Created by | 🟩 c | chess pog |
| 🕐 Created time | @April 2, 2024 10:53 PM | |
| ☰ Tags | | |

Functional Programming

Cassandara

Elastic Search

MongoDB

Structured, Semi-structure and Unstructured Data

CAP Theorem

Master-Slave Architecture

Client-Server architecture

Hadoop Map reduce

Application of Big Data Analytics

Hbase Architecture

Analyzers in Lucene

Vertical and Horizontal Scalibility

## Amazon Cloud / Big data in cloud

Amazon Web Services (AWS) provides a variety of cloud computing services, including storage, computing power, databases, analytics, and more. One of

the services offered by AWS is Amazon Elastic MapReduce (EMR), which is particularly relevant in the context of Hadoop.

Hadoop is an open-source framework for distributed storage and processing of large datasets across clusters of computers. It consists of two main components: Hadoop Distributed File System (HDFS) for storage and MapReduce for processing.

Amazon EMR simplifies the process of running Hadoop and other big data frameworks on AWS infrastructure. It allows users to quickly and easily provision resizable clusters of virtual servers (called instances) on AWS, pre-configured with Hadoop and other relevant tools like Apache Spark, HBase, Presto, and Flink.

Here's how Amazon EMR works in the context of Hadoop:

1. **Cluster Provisioning**: With Amazon EMR, you can provision a Hadoop cluster with a few clicks or through API calls. You specify the number and type of instances you need for your cluster, and AWS takes care of provisioning and configuring them.

2. **Integration with AWS Services**: Amazon EMR integrates seamlessly with other AWS services. For example, you can use Amazon S3 for storing input and output data, Amazon DynamoDB for input data, and Amazon RDS for storing metadata.

3. **Scalability**: One of the key benefits of using Amazon EMR is its scalability. You can easily scale your cluster up or down based on your processing needs. This elasticity allows you to handle varying workloads without worrying about managing the underlying infrastructure.

4. **Managed Environment**: Amazon EMR manages the Hadoop cluster environment, including installation, configuration, and tuning. This frees you from the operational overhead of managing a Hadoop cluster and allows you to focus on analyzing your data.

5. **Cost-Effective**: AWS offers a pay-as-you-go pricing model for Amazon EMR, where you only pay for the resources you use. This can be more cost-effective than maintaining and managing your own on-premises Hadoop cluster.

6. **Security and Compliance**: Amazon EMR provides various security features, including encryption of data at rest and in transit, integration with AWS

Identity and Access Management (IAM) for access control, and integration with AWS Key Management Service (KMS) for managing encryption keys.

In summary, Amazon EMR simplifies the process of running Hadoop and other big data frameworks on AWS infrastructure, providing scalability, cost-effectiveness, and managed services, which are essential for processing large datasets efficiently.

## Characteristics of NoSQL database

NoSQL databases, which stand for "Not Only SQL," are a diverse set of database management systems that are designed to handle large volumes of data that are not easily organized into the tabular schema of traditional relational databases. These databases offer various features and characteristics that make them suitable for specific use cases. Some common characteristics of NoSQL databases include:

1. **Schema-less or Flexible Schema**: Unlike relational databases, NoSQL databases typically do not enforce a rigid schema. This means that the structure of the data can be dynamic, allowing for easier adaptation to changing data requirements.

2. **Horizontal Scalability**: NoSQL databases are often designed to scale horizontally across multiple servers or nodes. This allows them to handle large volumes of data and high throughput by distributing the load across multiple machines.

3. **High Availability**: Many NoSQL databases are built with fault tolerance and high availability in mind. They often support features like data replication and automatic failover to ensure that data remains accessible even in the event of hardware failures or network issues.

4. **Partition Tolerance**: NoSQL databases are often designed to be partition tolerant, meaning they can continue to operate even if a portion of the network fails or becomes unavailable. This is crucial for systems operating at scale where network partitions are common.

5. **Support for Unstructured and Semi-Structured Data**: NoSQL databases are well-suited for storing unstructured or semi-structured data types such as JSON, XML, key-value pairs, graph data, or columnar data. This

flexibility allows developers to store and query data in a way that best fits their application requirements.

6. **Designed for Specific Use Cases**: NoSQL databases are often optimized for specific use cases such as real-time analytics, content management, social networking, IoT applications, and more. Different types of NoSQL databases (document-oriented, key-value, columnar, graph, etc.) are tailored to excel in different scenarios.

7. **Eventual Consistency**: Some NoSQL databases prioritize availability and partition tolerance over strict consistency. They may provide eventual consistency guarantees, meaning that updates to the data will eventually propagate to all replicas but may not be immediately consistent across the entire system.

8. **Open Source and Community Support**: Many NoSQL databases are open source and have vibrant developer communities contributing to their development and support. This can provide access to a wealth of resources, documentation, and community-driven plugins and extensions.

# IOT

The Internet of Things (IoT) refers to the network of physical objects or "things" embedded with sensors, software, and other technologies to connect and exchange data with other devices and systems over the internet. These objects can range from everyday household items to sophisticated industrial tools and equipment.

Here are some key components and concepts of IoT:

1. **Devices and Sensors**: IoT devices are equipped with various sensors such as temperature sensors, motion sensors, GPS, cameras, etc., to collect data from the environment or the object itself. These devices can be anything from smart thermostats and wearables to industrial machinery and vehicles.

2. **Connectivity**: IoT devices use different communication protocols such as Wi-Fi, Bluetooth, Zigbee, RFID, cellular networks (2G, 3G, 4G, and now 5G), and others to connect and transmit data to other devices or servers over the internet.

3. **Data Processing and Analytics**: The data collected by IoT devices is often processed and analyzed either locally on the device or in the cloud.

Advanced analytics techniques like machine learning and artificial intelligence can be applied to derive insights from this data, enabling predictive maintenance, optimization of operations, and other applications.

4. **Cloud Computing**: Cloud platforms play a crucial role in IoT by providing scalable storage and computing resources for processing the massive amounts of data generated by IoT devices. Cloud services also facilitate remote device management, software updates, and integration with other systems.

5. **Security and Privacy**: Security is a significant concern in IoT due to the interconnected nature of devices and the sensitive data they collect. Measures such as encryption, authentication, access control, and regular software updates are essential to protect IoT systems from cyber threats and unauthorized access.

6. **Applications**: IoT finds applications in various industries and domains, including smart homes, healthcare, agriculture, transportation, manufacturing, energy management, and environmental monitoring. Examples include smart thermostats that adjust temperature based on occupancy, wearable fitness trackers, connected cars with real-time traffic data, and smart grids for efficient energy distribution.

7. **Challenges and Considerations**: Despite its potential benefits, IoT also presents challenges such as interoperability issues, data privacy concerns, security vulnerabilities, and scalability issues. Additionally, managing and analyzing the vast amounts of data generated by IoT devices can be complex and resource-intensive.

## DFS

Distributed File System (DFS) is a file system with data stored on multiple servers and accessed and managed as a single coherent system. It allows users to access files and folders from multiple locations as if they were stored on a single device.

Here are some key concepts and characteristics of Distributed File Systems:

1. **Scalability:** DFS offers scalability by distributing file storage across multiple servers. As the storage requirements grow, additional servers can be added to the DFS infrastructure, allowing for seamless expansion without affecting system performance or accessibility.

2. **Redundancy and Fault Tolerance**: Many DFS implementations incorporate redundancy and fault tolerance mechanisms to ensure high availability and data reliability. This may involve replicating files across multiple servers or using techniques like RAID (Redundant Array of Independent Disks) to protect against data loss in the event of hardware failures.

3. **File Access Transparency**: One of the essential features of DFS is its ability to provide transparent access to files and folders across the network. Users can access files using familiar file paths or network drive mappings without needing to know the physical location of the data.

4. **Concurrency Control**: DFS must support concurrent access to files by multiple users or processes while maintaining data consistency and integrity. This typically involves implementing mechanisms such as file locking to prevent conflicts and ensure that changes made by one user do not interfere with others.

5. **Performance Optimization**: DFS implementations may include features to optimize performance, such as caching frequently accessed data locally to reduce network latency and improve responsiveness. Additionally, load balancing techniques may be employed to evenly distribute file access requests across multiple servers and prevent bottlenecks.

6. **Interoperability**: DFS systems should be compatible with various operating systems and network protocols to support heterogeneous environments. This enables users to access files using different devices and platforms seamlessly.

Some well-known examples of Distributed File Systems include:

- **NFS (Network File System)**: A distributed file system protocol developed by Sun Microsystems, widely used in Unix and Linux environments.

- **HDFS (Hadoop Distributed File System)**: A distributed file system used by the Apache Hadoop framework for storing and processing large datasets across clusters of computers.

# Hive

Hive is an open-source data warehouse infrastructure built on top of Hadoop. It provides a mechanism to project structure onto the data in Hadoop and to query that data using a SQL-like language called HiveQL. Hive enables data

summarization, query, and analysis of large datasets stored in Hadoop's distributed file system (HDFS). It was developed by Facebook and later contributed to the Apache Software Foundation. Hive offers features such as partitioning, indexing, and user-defined functions (UDFs), making it suitable for data warehousing tasks and analytics in big data environments. With its familiar SQL interface, Hive lowers the barrier to entry for users who are already familiar with SQL-based querying.

## Oozie

Oozie is a workflow scheduler system used to manage Hadoop jobs. It allows users to define workflows to orchestrate multiple Hadoop jobs and processes, including Hive queries, Pig scripts, MapReduce jobs, and more. Oozie workflows are defined using XML, where actions, dependencies, and control flow are specified. This enables users to create complex data processing pipelines with dependencies and triggers between various tasks. Oozie provides a web-based interface for monitoring and managing workflows, making it easier to track the progress of jobs and troubleshoot issues. It is widely used in Hadoop ecosystems for automating and managing data processing workflows, making it an essential tool for big data processing and analytics pipelines.

## Distributed Searching

Distributed searching refers to the process of searching for information's across multiple sources or nodes that are distributed across a network. Unlike traditional centralized searching, where a single search index or database is queried, distributed searching involves querying multiple sources simultaneously to retrieve relevant information. This approach is commonly employed in large-scale systems such as peer-to-peer networks, distributed databases, and distributed file systems.

In distributed searching, the search query is typically distributed to multiple nodes, each of which holds a portion of the overall dataset. These nodes independently process the query and return relevant results to the requester. The results are then aggregated and presented to the user, often with a ranking based on relevance.

One of the key advantages of distributed searching is scalability. By distributing the search load across multiple nodes, distributed searching can handle large

volumes of data and high query rates more efficiently than centralized searching. Additionally, distributed searching can improve fault tolerance and reliability since the failure of one node does not necessarily disrupt the entire search process.

However, distributed searching also poses several challenges. Ensuring consistency and synchronization across distributed nodes, managing communication overhead, and maintaining a comprehensive index of distributed data can be complex tasks. Furthermore, achieving optimal load balancing and efficient query routing are critical for maximizing performance in distributed searching systems.

JSON (JavaScript Object Notation) and XML (eXtensible Markup Language) are both popular formats used for storing and exchanging data. While they serve similar purposes, they have distinct differences in their structures and applications.

## JSON

- JSON is a lightweight data interchange format primarily based on JavaScript syntax, but it's language-independent.

- It is widely used for transmitting data between a server and a web application, making it popular in web development.

- JSON data is represented as key-value pairs and nested structures, making it easy to understand and work with.

- It's human-readable and easy to parse, which contributes to its popularity.

- JSON is commonly used for APIs (Application Programming Interfaces) and configuration files.

XML:

- XML is a markup language designed for storing and transporting data.

- It's more verbose compared to JSON due to its use of tags to define data structures.

- XML is highly extensible, allowing users to define their own markup tags and structure data according to their needs.

- It's commonly used in situations where data needs to be structured in a hierarchical manner, such as in documents or configurations.

- XML is widely used in industries like finance, healthcare, and publishing, where data integrity and structured representation are crucial.

Comparison:

- JSON is generally considered more lightweight and easier to parse than XML due to its simpler syntax.

- XML offers more flexibility and extensibility, making it suitable for complex data structures and documents.

- JSON is often preferred in web development and APIs due to its simplicity and compatibility with JavaScript, while XML is preferred in industries where complex data structures and standards compliance are important.

JSON:

```json
jsonCopy code
{
  "name": "John Doe",
  "age": 30,
  "email": "john@example.com",
  "address": {
    "street": "123 Main St",
    "city": "Anytown",
    "state": "CA",
    "postal_code": "12345"
  }
}
```

XML:

```xml
xmlCopy code
<contact>
  <name>John Doe</name>
  <age>30</age>
  <email>john@example.com</email>
  <address>
```

```
        <street>123 Main St</street>
        <city>Anytown</city>
        <state>CA</state>
        <postal_code>12345</postal_code>
    </address>
  </contact>
```

In the JSON example, data is represented as key-value pairs, and the structure is nested using curly braces `{}`. It's straightforward and easy to read.

In the XML example, data is enclosed within tags, creating a hierarchical structure. Each element contains a start tag, content, and an end tag. While more verbose, XML allows for greater flexibility and extensibility with custom tags and attributes.

## Zookeeper

Apache ZooKeeper is a centralized open-source service for maintaining configuration information, providing distributed synchronization, and providing group services in large distributed systems. It's often used in big data ecosystems like Apache Hadoop to coordinate and manage distributed systems.

In the context of Apache Hadoop, ZooKeeper plays a crucial role in providing coordination and synchronization among the various components of the Hadoop ecosystem. Here's how ZooKeeper is typically used in Hadoop:

1. **Configuration Management**: ZooKeeper helps in managing and distributing configuration details across the Hadoop cluster. This ensures that all nodes in the cluster are aware of the current configuration settings, which is essential for the proper functioning of the distributed system.

2. **Leader Election**: ZooKeeper facilitates leader election in distributed systems. For example, in Apache Hadoop's HDFS (Hadoop Distributed File System), ZooKeeper helps in selecting a leader among the NameNode instances to ensure high availability and fault tolerance.

3. **Synchronization**: ZooKeeper provides distributed synchronization primitives like locks, barriers, and queues, which are essential for coordinating distributed tasks and ensuring consistency among different components of the Hadoop ecosystem.

4. **Dynamic Membership**: ZooKeeper allows nodes to dynamically join or leave the cluster. This dynamic membership management ensures that the cluster can adapt to changes in its size or configuration without service disruption.

5. **Health Monitoring**: ZooKeeper can be used for monitoring the health of nodes in the cluster. By registering ephemeral nodes, each node can indicate its status, and failure detection mechanisms can be implemented using ZooKeeper's watch mechanism.

Overall, ZooKeeper serves as a crucial component in the Hadoop ecosystem, providing the necessary coordination and management capabilities required for building and running large-scale distributed systems effectively. It ensures that distributed applications built on top of Hadoop can operate reliably and efficiently in a clustered environment.

**Apache Sqoop**: Sqoop (SQL-to-Hadoop) is a tool designed for efficiently transferring bulk data between Apache Hadoop and structured datastores such as relational databases (e.g., MySQL, PostgreSQL, Oracle). It allows users to import data from relational databases into Hadoop's HDFS (Hadoop Distributed File System) for analysis using tools like MapReduce, Hive, or Spark, and to export data from Hadoop back into relational databases. Sqoop is particularly useful when dealing with large volumes of data stored in traditional relational databases and when you want to leverage Hadoop's processing power for analysis.

**Apache Flume**: Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of streaming data (such as log files, events, etc.) from various sources to a centralized data store like Hadoop's HDFS or HBase. Flume operates as a distributed system where data flows from sources (e.g., web servers, sensors) through a series of channels and sinks. Sources produce data, channels act as pipelines for the data flow, and sinks are responsible for writing the data to the destination. Flume is commonly used for log collection and aggregation in real-time analytics scenarios.

## Apache Pig

Apache Pig is a high-level platform for analyzing large datasets in Apache Hadoop. It provides a simple and powerful scripting language called Pig Latin, which abstracts the complexities of MapReduce programming. Pig Latin allows users to express their data transformation and analysis tasks concisely, enabling rapid development of data processing pipelines. Pig automatically optimizes and executes these tasks into MapReduce jobs, making it suitable for processing diverse data types and structures efficiently. With its ease of use and scalability, Apache Pig is widely used in big data environments for tasks such as ETL (Extract, Transform, Load), data cleaning, and exploratory data analysis.

**Combiner Functions:**

In the context of MapReduce programming, a combiner function is a function that operates on the output of the mapper phase before passing it to the reducer phase. The primary purpose of the combiner function is to perform a local aggregation of intermediate key-value pairs produced by the mappers.

Combiner functions are optional components in MapReduce programs, and their use depends on the specific requirements of the computation. They are particularly useful when the mapper output contains a large number of intermediate key-value pairs that can be aggregated before sending them over the network to the reducers. By performing local aggregation, combiners reduce the amount of data shuffled between mapper and reducer nodes, which can significantly improve the overall performance of the MapReduce job.

However, it's important to note that combiners are not always applicable or effective. They must satisfy the property of being associative and commutative since they may be applied multiple times and in any order during the MapReduce process. Additionally, combiners should not alter the final output of the MapReduce job; they are used solely for optimization purposes.

**Fault-Tolerant Systems:**

Fault tolerance is a critical aspect of distributed systems, especially in the context of Big Data processing frameworks like Hadoop and systems like Google File System (GFS). A fault-tolerant system is designed to continue operating properly in the event of one or more component failures within the system.

In the context of Big Data systems, fault tolerance is particularly important due to the scale and complexity of the data being processed. These systems typically run on clusters of commodity hardware, where component failures are relatively common. Therefore, mechanisms for fault tolerance are essential to ensure uninterrupted operation and data integrity.

Some common techniques used to achieve fault tolerance in Big Data systems include:

1. **Replication:** Data replication involves storing multiple copies of data across different nodes in the cluster. If one node fails, the data can still be accessed from other replicas, ensuring availability and reliability.

2. **Redundancy:** Redundancy involves having backup components or systems that can take over in case of failure. For example, in Hadoop, NameNode High Availability (HA) ensures that if the primary NameNode fails, a standby NameNode can take over seamlessly.

3. **Checkpointing:** Checkpointing involves periodically saving the state of a computation or data processing job. In case of failure, the system can restart from the last checkpoint rather than from the beginning, reducing the impact of failures on job execution time.

4. **Monitoring and Self-healing:** Fault-tolerant systems often include mechanisms for monitoring the health of system components and automatically taking corrective actions in case of failures. This may involve restarting failed processes, reallocating resources, or triggering failover to backup components.

By implementing these and other fault tolerance techniques, Big Data systems can ensure continuous operation and data integrity despite the inherent challenges of distributed computing environments.