

Date: / /

## Demonstrate Ambari's process to start and stop Hadoop services.

Apache Ambari is a web based management tool for provisioning, managing, and monitoring Apache Hadoop clusters.

Use a step-by-step guide on how to start and stop Hadoop services using Ambari.

- Access Ambari Web interface
- Navigate to Services Page
- Start or Stop a service
- Confirm the Action
- Monitor Service Status

## Show how to monitor health using Ambari Metrics.

Apache Ambari provides a comprehensive set of metrics to monitor the health and performance of your Hadoop cluster. Step-by-step guide on how to monitor cluster health using Ambari Metrics.

- Access Ambari web interface
- Navigate to Dashboard
- View Cluster Metrics
- Drill Down into Service Metrics
- View Host Metrics
- Set up Alerts
- Analyze Historical Data

Date: / /

## Write and explain a Pig Script for Student marks Analysis.

- Loading the Data
- Grouping by student
- Calculating Total marks
- Calculating Average Marks
- Filtering High Scores
- Output

Describe a real-world use case of IoT contributing to Big Data.

### Smart cities

Traffic management in smart cities, IoT sensors are embedded in traffic lights, roadways, and vehicles to collect real-time data on traffic flow, vehicle speed, congestion levels, and accidents.

The data is transmitted continuously to continuously to centralized Big Data systems, where it is processed and analyzed.

City planners and traffic control systems use this Big Data to optimize traffic signal timings, predict and reduce congestion, provide real-time updates to commuters via

Date: 11/1

Q1:

Improve emergency response time.

- ① Big Data ~~System~~ Pipeline using Hive and Pig

A Big Data pipeline using Hive and Pig enables efficient processing, querying and transformation of large datasets stored in Hadoop Distributed File system.

Here a pipeline:- Data Ingestion  
Data cleaning & Transformation  
Data warehousing & Querying  
Output Visualization

- ② Show how to monitor cluster health using Ambari Metrics.

To monitor cluster health using Apache Ambari follow using its metrics and dashboard tools.

Ambari provides real-time metrics and visualization for hadoop services, helping administrators detect and resolve performance issues.

Steps to monitor cluster Health using Ambari Metrics

- (1) Login to Ambari Web UI
- (2) Dashboard Overview

Date: / /

Service Health monitoring  
Knox Metrics view  
Ambari Metrics system (AMS)  
Set and manage Alert  
Use Keathab  
Check Log for Deep Analysis

Common metrics to watch

HDFS :- Live / Dead DataNodes, Disk usage,  
Block status

YARN :- Running Applications, containers,  
Resource usage.

hive :- Query Latency, connections

Hbase :- Region server status, Read / Write  
Latency

General :- CPU, Memory, Network, Disk usage

Ques Describe how Spark handles lazy  
evaluation.

Lazy evaluation means Spark does not  
execute transformations (like map, filter,  
etc) immediately when they are called.

Instead, it builds a logical execution  
plan and waits until an action (like to

Date: / /

(count(), collect(), saveAsTextFile()) is called.

How it works

Transformation

Action triggers execution

Benefits of Lazy Evaluation

Performance optimization

Failure Recovery

Efficient memory use

Show a use case for real-time processing using spark streaming

Real time Fraud Detection in Bank using spark streaming :-

A bank wants to monitor credit card transaction in real time to detect fraudulent activity. They need to analyze transaction data as it arrives to flag suspicious patterns like:

- Unusually high-value transaction
- multiple transaction from different locations within a short time
- Transactions from blacklisted IPs or merchants

## How Spark Streaming Helps

- ① Data ingestion
- ② Real Time Processing
- ③ Action

### Example Technologies used:-

- Apache Kafka  
For real time transaction data ingestion.
- Apache spark streaming  
For processing and Analysis.
- HBase / Cassandra  
To store flagged transactions.
- Dashboard  
For visualization and alert monitoring.

### Benefits

- Detect fraud within seconds instead of hours.
- Prevent financial losses and improve customer trust.
- Scalable and Fault tolerant real-time analytics.

(A) Write A command to list all files in an HDFS directory.

Ans: To list all files in an HDFS directory use the following command

hdfs dfs -ls <directory-path>

(B) Which SQL command is used to create a table in RDBMS?

M The CREATE TABLE command is used to create a new table in RDBMS.

The CREATE command is in DDL command.

(C) Which is the purpose of the hdfs-rm command?

M The rm command is used to remove objects such as files, directories, symbolic links and so on from the file systems like Unix.

d) Name any one DCL (Data Control Language) command!

M The GRANT and REVOKE command is used to assign specific privileges to a user.

(c) Which command is used to insert data into an RDBMS table?

L The INSERT command is used to insert data into a table in a RDBMS

(a) What is cloudera?

M Cloudera is a hybrid data platform designed for unmatched freedom to choose - any cloud, any analytics, any data. Cloudera delivers faster and easier data management and data analytics for data anywhere, with optimal performance, scalability, and security.

(b) What is the Role of Apache Ambari in Hadoop?

M Apache Ambari, an open-source software is a management platform that provides web-based user interfaces and APIs for monitoring, managing and provisioning Apache Hadoop clusters.

(c) Define mapreduce-

M mapreduce is a programming model that runs

Date: / /

on hadoop - a data analytical engine ~~application~~  
~~that~~ widely used for big data - and  
written applications that run in parallel  
to process large volumes of data  
stored on clusters.

d) What is partitioning in Hive?

M A partitioning strategy that is used  
to split a table into multiple files  
based on partition keys.

e) Define a Resilient Distributed Dataset (RDD)?

M A resilient distributed dataset (RDD) is  
Apache Spark's fault tolerant  
distributed collection of data elements  
that can be processed in parallel across  
a cluster of machines.

f. Discuss the role and usage of  
Apache Spark in the Hadoop  
ecosystem.

C Apache Spark is a powerful open source  
distributed processing system designed for  
big data workloads. It enhances the  
Hadoop ecosystem by providing faster and  
more versatile data processing capabilities  
than Hadoop's traditional MapReduce framework.

Date: 1/1

## Role in the Hadoop Ecosystem

Processing Engine

Resource management

Enhanced capabilities for fault and

availability across different clusters

## Use of Apache Spark

Batch Processing with 3 Part

Real-time Stream Processing

Interactive Queries

Machine Learning

Graph Processing

## Key Advantages

Speed, Flexibility, Unified Engine,

Each of the various components

Date: 1/1

## 5. Explain the evolution from traditional data processing to Big Data Processing

M The evolution from traditional data processing to Big data processing has been driven by the rapid growth of data volume, variety, and velocity.

Traditional data processing involved:-  
Structured data

Batch processing

Centralized Architecture

## Big Data Processing -

Distributed Architecture

Scalability

Flexibility

Real time processing

Some key technologies that have enabled the evolution to Big Data processing include:

Hadoop

Spark

NoSQL databases

Cloud Computing

The evolution to Big Data processing has provided several benefits:- Improved insight

Increased efficiency

Competitive advantage

Date: / /

Q6. Explain Ambari alerts and their role in cluster management.

- Ambari alerts are automated health checks that monitor the status and performance of various services and components in a Hadoop cluster.

### Type and customization

Ambari supports several types of alerts, including:

- Web, port, metric, aggregate, script, service, and security alerts.
- Alerts can be customized for check intervals, thresholds and notification methods.

### Role in cluster management

- Continuous monitoring of cluster components and services for health and availability.
- Immediate notification when thresholds are breached.
- Quick identification and troubleshooting of issues.
- Centralized management.

Date: / /

100. Create a step-by-step guide to installing and configuring Apache Spark.

1  
Step-by-step Guide:

Installing and Configuring Apache Spark Using Ambari

- (1) Download Apache spark
- (2) Install Apache spark
- (3) Configure Environment Variables
- (4) Configure spark
- (5) Start spark
- (6) Configure spark cluster
- (7) Verify spark cluster

List and describe Apache spark libraries.

Apache Spark offers several core and extension libraries designed to handle different data processing and analytical tasks. The main ones:

Core Spark Libraries

• Spark Core

The foundation of memory computing scheduling. Spark providing in and distributed task

Date: / /

Handles basic I/O task scheduling and fault recovery.

### Spark SQL

- Enables querying structured data using SQL on Data Frame API.
- Integrates with Hive, HDFS and other data sources for relational processing.

### Spark Streaming

- Process live data stream in real time.
- Integrates with sources like Kafka, Flume, and Kinesis.

### MLlib

- Scalable ml library.
- Includes algorithms for classification, regression, clustering and collaborative filtering.
- Compatibile with Java, Scala, Python, and R.

### GraphX

- graph processing library.
- Provides API for graph-parallel computation and analytics.