

Creating Directory Structure in HDFS

Creating the directory structure in HDFS is the first step to organizing data in Hadoop. It involves planning the hierarchy of directories for efficient data storage and access. The structure should align with the organization's data requirements and access patterns.

KA by Kakumanu Kalyani

Overview

Setting Up Hadoop Cluster

Started:	<div>1</div> <div>Tue Apr 03 14:20:51 -0400 2018</div>
Version:	<div>Resource Allocation</div> <div>3.0.1, r496dc57cc2e4f4da117f7a8e3840aaeac0c1d2d0</div>
Compiled:	<div>Allocate and configure hardware resources for the Hadoop cluster.</div> <div>Fri Mar 16 19:00:00 -0400 2018 by lei from branch-3.0.1</div>
Cluster ID:	<div>2</div> <div>CID</div>
Block Pool ID:	<div>Cluster Configuration</div> <div>BP</div>
	<div>Install and set up Hadoop components across the cluster nodes.</div>
	<div>3</div> <div>Network Configuration</div>
	<div>Optimize network settings for seamless communication within the cluster.</div>

Configuring Hadoop Environment

Hadoop Configuration Files

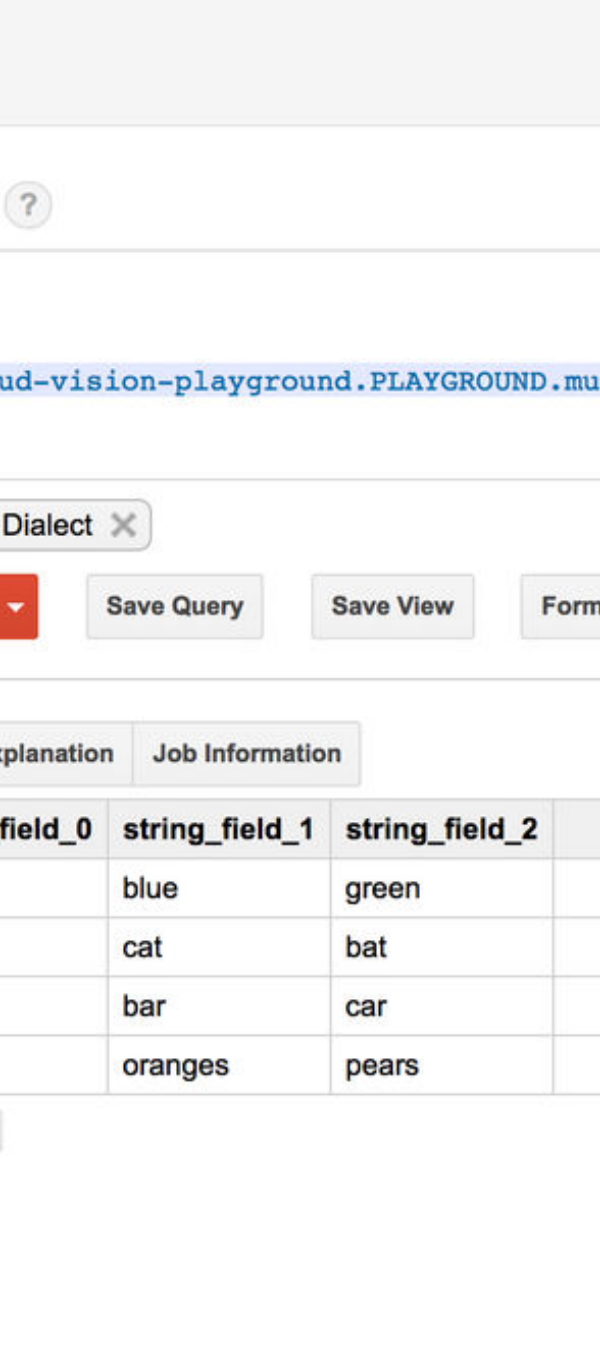
Modify core-site.xml and hdfs-site.xml to define Hadoop cluster configurations.

Memory Management

Adjust memory settings for efficient execution of Hadoop components.

Security Setup

Implement security measures, including user authentication and authorization.



Loading Data into HDFS

1 Data Ingestion Methods

Choose between direct data load, file upload, or distributed copy methodologies.

2 Data Replication Strategy

Determine the replication factor for fault tolerance and data redundancy.

3 Integration with ETL Tools

Integrate with ETL tools for seamless data extraction, transformation, and loading.

Running MapReduce Jobs



Map Phase

Divide and conquer data processing using the map phase.



Reduce Phase

Aggregate and summarize map outputs during the reduce phase.



Job Execution

Monitor and manage MapReduce jobs for optimal execution.

Monitoring Hadoop Cluster

Cluster Health

Monitor node status, resource utilization, and overall cluster health.

Alerting & Notifications

Set up alerts and notifications for critical cluster events and anomalies.

Performance Optimization

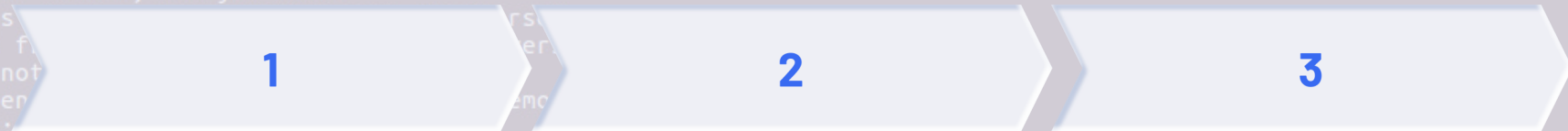
Analyze and optimize cluster performance based on resource usage and job execution.



```
naveen@naveen-ubuntu:~$ su nvnhadoop
Password:
To run a command as administrator (user "root"), use "sudo <command>".
See "man sudo_root" for details.

nvnhadoop@naveen-ubuntu:/home/naveen$ cd
nvnhadoop@naveen-ubuntu:~$ hdfs dfs -safemode get
No command 'hdfs' found, did you mean:
  Command 'hdfsls' from package 'hdf4-tools' (universe)
  Command 'hfs' from package 'hfsutils-tcltk' (universe)
hdfs: command not found
nvnhadoop@naveen-ubuntu:~$ sudo hdfs dfs -ls /
[sudo] password for nvnhadoop:
sudo: hdfs: command not found
nvnhadoop@naveen-ubuntu:~$ hdfs dfsadmin -safemode get
No command 'hdfs' found, did you mean:
  Command 'hdfsls' from package 'hdf4-tools' (universe)
  Command 'hfs' from package 'hfsutils-tcltk' (universe)
hdfs: command not found
nvnhadoop@naveen-ubuntu:~$ bin/hdfs dfs -ls /
bash: bin/hdfs: No such file or directory
nvnhadoop@naveen-ubuntu:~$ sudo ln -s /usr/local/hadoop/bin /usr/local/bin
[sudo] password for nvnhadoop:
nvnhadoop@naveen-ubuntu:~$ hdfs dfs -ls /
No command 'hdfs' found, did you mean:
  Command 'hfs' from package 'hfsutils-tcltk' (universe)
  Command 'hdfsls' from package 'hdf4-tools' (universe)
hdfs: command not found
nvnhadoop@naveen-ubuntu:~$ bin/hdfs dfs -ls /
bash: bin/hdfs: No such file or directory
nvnhadoop@naveen-ubuntu:~$
```

Troubleshooting Hadoop Issues



Issue Identification

Thoroughly diagnose and identify the root cause of Hadoop issues.

Troubleshooting Steps

Follow systematic troubleshooting steps to resolve identified issues.

Log Analysis

Review logs and error messages to trace and resolve problems.

Best Practices for Hadoop Deployment

1

Performance Tuning

Fine-tune configurations for optimal performance and resource utilization.

2

Security Optimization

Implement encryption, authentication, and role-based access controls for data security.

3

Backup & Recovery

Establish robust backup and recovery mechanisms for data protection and disaster recovery.