

ABSTRACT

Hadoop is a framework that supports data-intensive distributed applications. And is used to process data on thousands of computers. Two scalable distributed clustering methods are proposed: the Single Clustered File method and the Multiple Clustered Files method. Both methods were executed on one to eight machines in Hadoop.

This Project is dealing with implementation of bigdata using hadoop framework. we have provided the support of data hadoop version 1.2.1 and version 2.6.4. The project includes implementation of Docker Technology. My project provide two way to set up hadoop cluster

- **AUTOMATIC CONFIGURATION**
- **ON DEMAND CONFIGURATION**

AUTOMATIC CONFIGURATION

In this phase we configure hadoop cluster in such a way that it uses computer resources in optimized way. So it selects automatically a node which has more storage and ram used as a Datanode and TaskTracker and which has less ram and storage used as a Namenode and JobTracker among all the available node and we are also providing a simple and well-structured web interface to the user so that he can use the cluster in efficient way.

ON DEMAND CONFIGURATION

In this phase we ask the user to enter number of datanode and task tracker and it creates instantly. To implement this we use a tool called Docker. Docker containers wrap a piece of software in a complete filesystem that contains everything needed to run: code, runtime, system tools, system libraries – anything that can be installed on a server. This guarantees that the software will always run the same, regardless of its environment.

LIST OF CONTENT

S.NO	CHAPTER NO.	CONTENTS	PAGE NO.
1	CHAPTER 1	BIG DATA	
		1.1 What is big data	1
		1.2 Human generated and machine generated data	1
		1.3 Where does big data comes from	1
		1.4 Examples of big data in real world	2
		1.5 Challenges of big data	2
2	CHAPTER 2	HADOOP AND BIG DATA	
		2.1 How hadoop solve the big data problem	4
		2.2 Business case for hadoop	5
3	CHAPTER 3	HADOOP DISTRIBUTED FILE SYSTEM(HDFS)	
		3.1 Namenode	8
		3.2 Secondary Namenode	9
		3.3 Data node	10
		3.4 Other features	11
		3.5 Compare with NFS	11
		3.6 Limitation of HDFS	12
4	CHAPTER 4	MAP REDUCE	
		4.1 Job Tracker	15
		4.2 Task Tracker	15
		4.3 Node Manager	15
		4.4 Resource Manager	16
		4.5 Features	16
5	CHAPTER 5	HADOOP ENVIRONMENT SETUP	19
6		CONCLUSION	24