

Rs - 120/-

HADOOP

CLASS ROOM NOTES



Flat No: 212, 2nd Floor, Annapurna Block,
Aditya Enclave, Ameerpet, Hyderabad, AP.

Ph No: 040 6462 6789, 0998 570 6789
E-mail: info@kellytechno.com, www.kellytechno.com.

Hadoop

By Mr. GOPAL KRISHNA |

- * Inspired by Google Bigtable and mapreduce papers circa 2004.
- * Created by Doug Cutting.
- * Originally built to support distribution for search engine.
- * Named after a shuffled elephant.

Big Data :-

Big Data is about the growing challenge that organizations face as they deal with large and fast-growing sources of data or information that also present a complex range of analysis and use problem. These can include:

- * Having a computing infrastructure that can ingest, validate, and analyze high volumes (size and/or rate) of data.
- * Assessing mixed data (structured and unstructured) from multiple sources.
- * Dealing with unpredictable content with no apparent schema or structure.
- * Enabling real-time near-real-time collection, real-time analysis and answers.

BigData technologies describe a new general of technologies and architectures, designed to extract value from very large volumes of a wide variety of data, by enabling high-velocity capture, discovery, and/or analysis.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 988 570 6739

The data comes from everywhere: sensors used to gather climate information, posts to social media sites, digital pictures and videos, purchase transaction records, and cell phone GPS signals to name a few. This data is Big Data.

Three vs of BigData :-

Volume describes the amount of data generated by organizations or individuals.

Variety describes structured and unstructured data, such as text, sensor data, audio, video, click streams, log files and more.

Velocity describes the frequency at which data is generated, captured and shared.

BigData - Hadoop :-

Hadoop is the Apache open source software framework for working with BigData. It was derived from Google technology and put to practice by yahoo and others. But, Bigdata is too varied and complex for a strict one-size-fits-all solution.

while Hadoop has surely captured the greatest name, recognition, it is just one of three classes of technologies well suited to storing and managing BigData.

Hadoop is software framework, which means it includes no. of components that were specifically designed to solve large-scale distributed data storage, analysis and retrieval tasks.

Lead of BigData :-

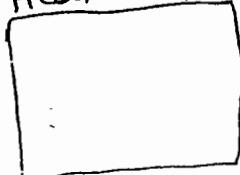
* Google process 20 PB a day

* Wayback machine has 3 PB + 100 TB/month

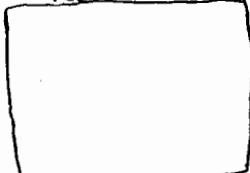
* Facebook has 2.5 PB of user data + 15 TB/day

* eBay has 6.5 PB of user data + 50 TB/day

Healthcare



Retail



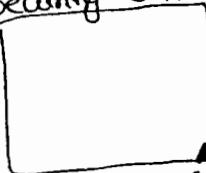
Finance

By Mr. GOPAL KRISHNA

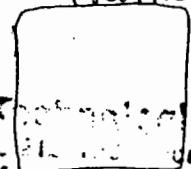
Log Analysis



Security Surveillance



Traffic Control

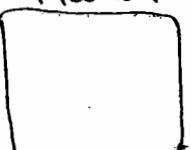
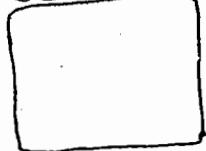


Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Endave,
Amarpet, Hyderabad-500 016.
Ph: 040-2462 6789, 098 970 5789

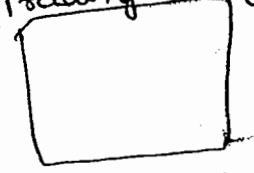
Telecom



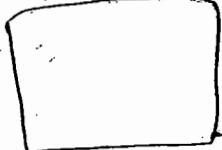
Search Quality



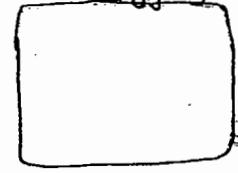
Trading Analytics



Fraud & Risk



Energy & utilities



Hadoop for the research users a quick timeline of how things have progressed :-

2004:- Initial versions of what is now HDFS (Hadoop distributed file system) and mapReduce implemented by DougCutting & Mike Cafarella.

Dec 2005:- Nutch ported to the new framework. Hadoop runs reliably on 20 nodes.

Jan 2006:- Doug Cutting joins yahoo!

Feb 2006:- Apache Hadoop project off mapReduce and support the standalone development of mapReduce and HDFS. Adoption of Hadoop by yahoo! Gridteam.

April 2006:- sort benchmark (10GB/node) run on 188 nodes in 47.0 hours [better than April benchmark]

May 2006:- yahoo! setup a hadoop research cluster - 300 nodes run on 500 nodes in 42 hours [better set of machines].

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Endave,
Amarpet, Hyderabad-500 016.
Ph: 040-2462 6789, 098 970 5789

2006 :- Research cluster reaches 600 nodes.

In 2006 :- sort benchmark run on 20 nodes in 1.8 hrs
100 nodes in 3.3 hrs
500 nodes in 5.2 hrs
900 nodes in 7.8 hrs.

Jan 2007 :- Research cluster reaches 900 nodes.

Apr 2007 :- Research clusters - two clusters of 1000 nodes.

Apr 2008 :- Won the 1-TeraByte sort benchmark in 209 seconds on 900 nodes.

Oct 2008 :- Loading 10 TeraBytes of data per day onto the research clusters.

March 2009 :- 17 clusters with a total of 24,000 nodes.

April 2009 :- won the minute sort by sorting 500GB in 59 seconds (on 1400 nodes) and the 100-terabyte sort in 173 minutes [on 8,400 nodes].

Example for a Statistical Analysis of ONE MINUTE internet search :-

- Google Receives over 2,000,000 search queries
- Facebook receive 34,722 "Likes".
- consumers spend \$ 272,070 on web shopping.
- Apple receives 47,000 Apps downloads.
- 370,000 minutes of calls on Skype;
- 98,000 posts on tweets,
- 20,000 posts on Tumblr,
- 13,000 hours of music streaming on Pandora,
- 12,000 new ads on Craigslist,
- 6,600 pictures uploaded to flickr
- 1,500 new blog posts,
- 600 new YouTube videos.

Hadoop's History :-

By Mr. GOPAL KRISHNA

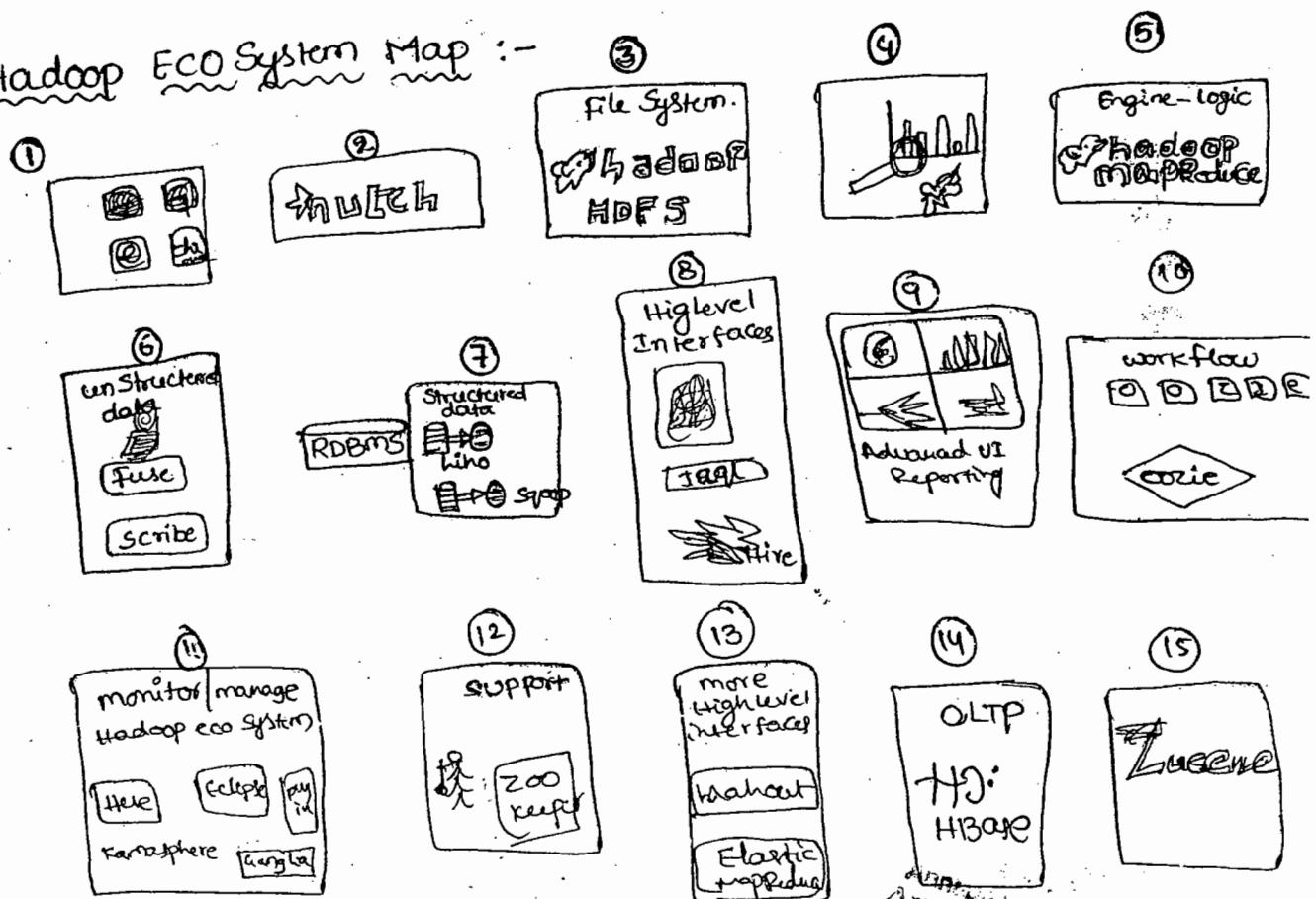
3

→ Hadoop was created by Doug Cutting & Michael J. Cafarella. Doug, who was working at yahoo at the time, named it after his son's toy elephant. It was originally developed to support distribution for the Nutch search engine project.

→ Hadoop is based on work done by Google in the late 1990s/early 2000s.

specifically, on papers describing the Google File System (GFS) published in 2003, and the MapReduce paper published in 2004.

Hadoop ECO System Map :-



① Large data on the web.

② Nutch built to crawl this web data.

③ Large volume of data had to stored - HDFS introduced.

④ How to use this data? Report

⑤ MapReduce Framework built for coding & running analytics

- ⑥ unstructured data - weblogs, click streams, Apache logs, Server logs - fuse, webdav, chukwa, flume and scribe.
- ⑦ sqoop and hilo for loading data into HDFS - RDBMS data.
- ⑧ High level interfaces required over low level map reduces programming - Hive, Pig, Jaql
- ⑨ BI tools with advanced UI Reporting.
- ⑩ workflow tools over Map-Reduce processes and High level languages - Oozie.
- ⑪ monitor & manage hadoop, run jobs/hive, view HDFS - high level view - Hue, Karmasphere, eclipse plugin, cacti, ganglia.
- ⑫ Support frameworks - Avro (Serialization), Zookeeper (Coordination).
- ⑬ More High level interfaces/uses - Mahout, Elastic MapReduce
- ⑭ OLTP also possible in HBase.
- ⑮ Lucene is a text search engine library written in Java.

Different Ecosystems of Hadoop :-

- * Hadoop is best known for mapReduce and it's Distributed file System (HDFS, renamed from NDFS)
- [note:- NDFS is also used for projects that fall under the umbrella of infrastructure for distributed computing and large scale data processing]

HDFS :- if you want 4000+ computers to work on your data, then you'd better spread your data across 4000+ computers. HDFS does this for you. HDFS has a few moving parts. The Datanodes store your data, and the NameNode keeps track of where stuff is stored. There are other pieces, but you have enough to get started.

MapReduce :- This is the programming model for Hadoop. There are two phases, not surprisingly called map and Reduce. To impress your friends tell them there is sort & shuffle between the map and Reduce phase. The job tracker manages the 4000+ components of your mapReduce job. The task trackers take orders from the job tracker. If you like Java then code in Java. If you like SQL or other non-Java languages you are still in luck, you can use a utility called Hadoop Streaming.

By Mr. GOPAL KRISHNA

Hadoop Streaming :- A utility to enable MapReduce code in any language: C, perl, python, C++, Bash etc. The examples include a python mapper and an AWK reducer.

Hive and Hue :- If you like SQL, you will be delighted to hear that you can write SQL and Hive convert it to a MapReduce job. No, you don't get a full ANSI-SQL environment, but you do get 4000 nodes and multi-peta-byte scalability. Hue gives you a browser based graphical interface to do your Hive work.

Pig :- A higher-level programming environment to do MapReduce coding. The pig language is called Pig Latin. You may find the naming conventions somewhat unconventional, but you get incredible price-performance and high availability.

Sqoop :- provides bi-directional data transfer between Hadoop and your favorite relational database.

Oozie :- Manages Hadoop workflow. This does not replace your scheduler or BPM tooling, but it does provide if-then-else branching and control Hadoop jobs.

HBase:- A super-scalable key-value store. It works very much like a persistent hash-map (for python fans think dictionary). It is not a relational database despite the name HBase.

Flume:- A real time loader for streaming your data into Hadoop. It stores data in HDFS and HBase. You will want to get started with Flume, which improves on the original Flume.

Mahout:- machine learning for Hadoop. Used for predictive analytics and other advanced analysis.

Fuse:- makes the HDFS system look like a regular filesystem so you can use ls, rm, cd, and others on HDFS data.

ZooKeeper:- used to manage synchronization for the cluster. You won't be working much with Zookeeper, but it is working hard for you. If you think you need to write a program that uses Zookeeper you are either very, very, smart and could be a committee for an Apache project, or you are about to have a very bad day.

How Big is 'BigData'?

With time, data volume is growing exponentially. Earlier we used to talk about Megabytes or Gigabytes. But time has arrived when we talk about data volume in terms of terabytes, petabytes and also zettabytes! Global data volume was around 1.82 ZB in 2011 and is expected to be 7.9 ZB in 2015. It is also known that the global information doubles in every two years!

How analysis of BigData is useful for organizations? 5

Effective analysis of Big Data provides a lot of business advantage as organizations will learn which areas to focus on and which areas are less important. Big Data analysis provides some early key indicators that can prevent the company from a huge loss or help in grasping a great opportunity with open hands! A precise analysis of Big Data helps in decision making. For instance, now a days people rely so much on Facebook and Twitter or service.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 998 570 6789

HDFS and mapReduce features in Hadoop versions:

Feature	1.0	0.22	2.0
Secure Authentication	yes	no	yes
old configuration Names	yes	deprecated	deprecated
New Configuration Names	no	yes	yes
old mapReduce API	yes	yes	yes
New MapReduce API	yes [with some missing libraries]	yes	yes
MapReduce 1 runtime [classic]	yes	yes	no
MapReduce 2 runtime [YARN]	no	no	yes
HDFS Federation	no	no	yes
HDFS high availability	no	no	yes

what is Hadoop?

- Hadoop is analysis for Bigdata.
- Hadoop is an open source framework for creating distributed applications that process huge amount of data.

one definition of huge

Ex:- 25,000 machines

more than 10 clusters

3 PB of data (compressed, unreplicated)

700+ users

10,000+ jobs/week.

another definition

- Hadoop is a open source, distributed, Batch processing and fault-tolerance system which is capable of storing huge amount of data [TB, PB, Zettabytes etc] along with processing on the same amount of data.
- Hadoop easy to use parallel programming model.

- Hadoop framework consists two main core components:

① HDFS ② mapReduce.

it is responsible
for storing massive
amount of data
on the cluster.

it is responsible
for processing
massive amount
of data on
the cluster.

- HDFS & mapReduce capabilities are it's Kernel for Hadoop.

HDFS (Hadoop Distributed File System)

6

Hadoop modes of Installation :-

① stand alone mode

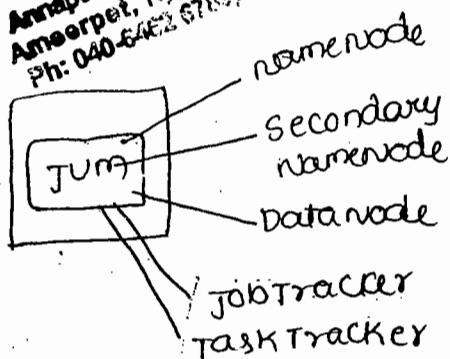
② pseudo distributed mode [labwork]

③ Fully Distributed mode [Realtime]

① stand alone mode :< Development Test Debug..

- single machine
- no daemons are running
- Everything runs in single JUM.
- standard OS storage.
- Good for development and test with small data, but will not catch all errors.

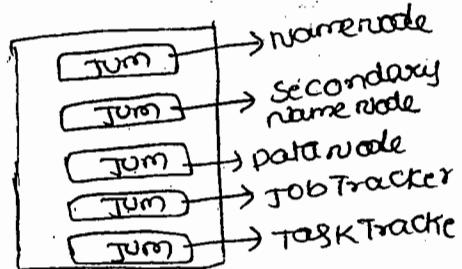
Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6462 6719, 988 570 6789



By Mr. GOPAL KRISHNA

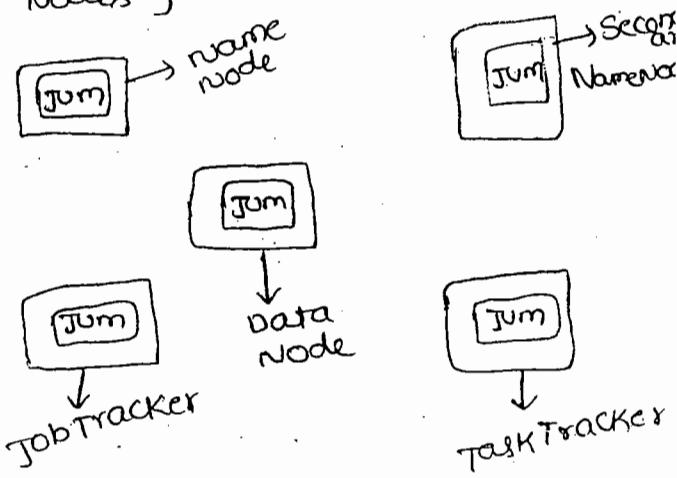
② pseudo distributed mode :- [multiple JUM's in single node]

- single machine but cluster is simulated.
- Daemons run in separate JUM.
- separate JUM's in single node.



③ Fully Distributed mode :- [All components run in separate nodes]

- Run Hadoop on cluster of machines
- Daemons run
- production environment
- Good for staging & production.



what is distributed file system :-

- System that permanently store data.
- Support concurrency, distribution, replication access to file and remote servers.
- divided into logical units (files, shards, chunks, blocks).
- DFS's are working w/w based approach because of DFS's are more complex than Regular disk file systems.

Ex:- file system tolerate node failure without suffering data loss.

Hadoop Distributed File System :-

- Hadoop is a distributed file system and it uses to store bulk amounts of data like terabytes (or) even petabytes.
- HDFS support high throughput mechanism for accessing this large amount information.
- In HDFS files are stored in sequential redundant manner over the multiple machines and this guaranteed the following ones.
 - ① Durability to failure
 - ② High availability to very parallel applications.

Example: NFS (Network file system)

- NFS gives remote access to a single logical volume stored on a single machine.
- NFS server can visible a portion of its local files system to external clients and also the client can mount this remote file system directly into their own Linux file system, and interact with it as though it were part of the local drive.

Advantage of NFS

- It is transparency [that is clients do not need to be particularly aware that they are working on files stored remotely]

Advantages of HDFS :-

- ① HDFS store large amount of information.
- ② HDFS is simple and robust coherency model.
- ③ That is it should store reliability.
- ④ HDFS is scalable and fast access to this information and it also possible to serve large number of clients by simply adding more machines to the cluster.
- ⑤ HDFS should integrate well with Hadoop mapReduce, allowing data to be read and computed upon locally when possible.
- ⑥ HDFS providing streaming read performance.
- ⑦ Data will be written to the HDFS once and then read several times.
- ⑧ The overhead of caching is helps to read data from HDFS source.
- ⑨ Fault-tolerance by detecting faults and applying quick, automatic recovery.
- ⑩ processing logic close to the data, rather than the data close to the processing logic.
- ⑪ portability across heterogeneous commodity hardware and operating Systems.
- ⑫ Economy by distributing data and processing across clusters of commodity personal computers.
- ⑬ Efficiency by distributing data and logic to process it in parallel on nodes where data is located.
- ⑭ Reliability by automatically maintaining multiple copies of data and automatically redeploying processing logic in the event of failures.

By Mr. GOPAL KRISHNA

HDFSTechologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aanya Enclave,
Ameerpet, Hyderabad - 500 046
P: +91 98490 67777, M: +91 98490 67777

DisAdvantages of HDFS :-

In distributed file system, it is limited in its power. The files in an NFS volume all reside on a single machine. This will create some problems.

- ① It does not give any reliability guarantees if that machine goes down.

Ex:- By replacing the files to other machine.

- ② All the clients must go to this machine to retrieve their data. This can overload the server if a large no. of clients must be handled.
- ③ clients need to copy the data to their local machines before they can operate on it.

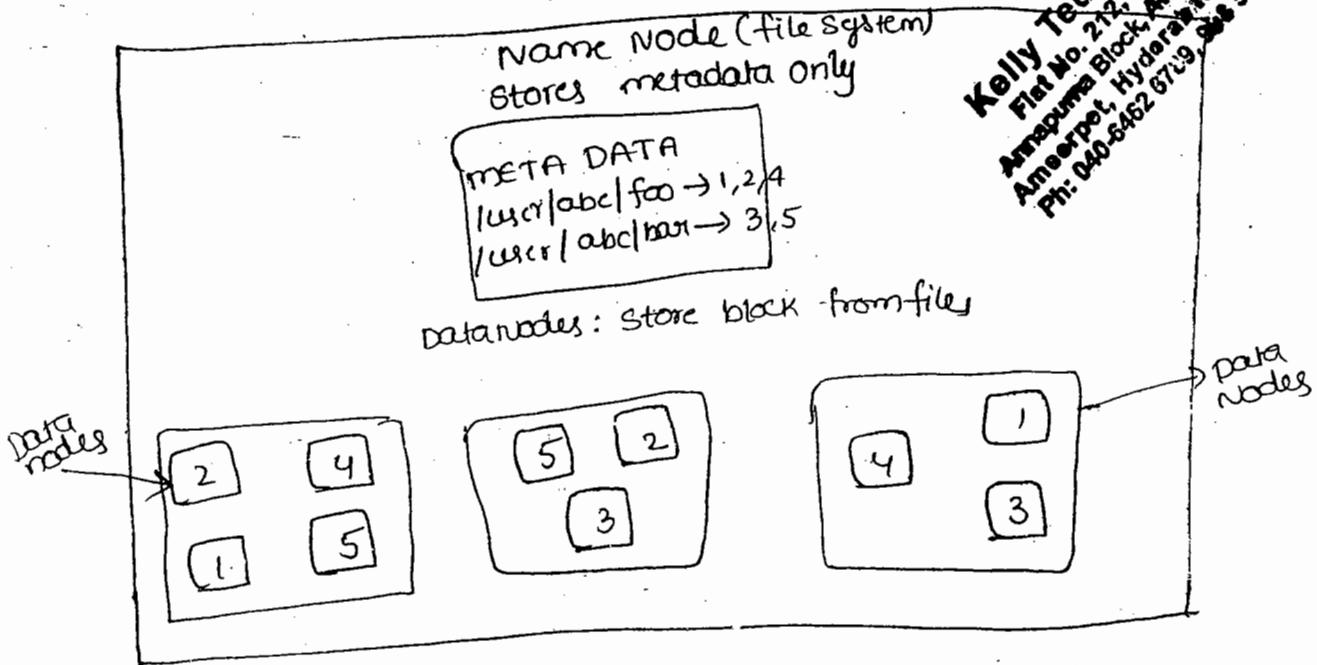
Goals of HDFS:-

- ① very large distributed file System :- 10K nodes, 100 million files, 10 PB.
- ② Assumes commodity hardware : files are replicated to handle hardware failure. Detect failures and recover from them.
- ③ optimized for batch processing :- Data locations exposed so that computations can move to where data resides. It provides very high aggregate bandwidth.

HDFS is a block structured file System:-

Block :- Block is the minimum unit of data that is said in HDFS, which is typically 64 MB by default. However we can increase the blocksize in the multiples of 64 MB.

- Each file is broken into blocks of a fixed size of and these blocks are stored across a cluster of one or more machines with data storage capacity.
- Individual machines in the cluster are called the data nodes.
- A file can be made of several blocks and not necessarily stored on the same machine.
- The target machine chose each block randomly on a block-by-block basis. **By Mr. GOPAL KRISHNA**
- So access permission to a file may need the cooperation of multiple machines and it supports file size for larger than a single machine DFS.
- Individual files sometimes need large space than a single hard drive ~~you~~ could hold. If several machines must be involved in the serving of a file then a file could be rendered unavailable by the loss of any one of those machines. HDFS controls this problem by replicating each block across a no. of machines (by default 3).



- In the above figure the variances represents multiple files with replication factor of 2 and the NameNode maps the filenames onto the block ids.
- In block structured file systems commonly use a block size on the order of 4 or 8 KB.
- The default block size in HDFS is 64 KB. This permits HDFS to decrease the amount of metadata storage required per file.
- In HDFS block structured file system, all the information's are handled by single machine called NameNode.
- The NameNode stores all the metadata for the file system. All the informations like tracks file names, permissions, and the locations of each block of each files etc can be stored in the main memory of the NameNode machine and it permits fast access to the metadata.
- To allow open a file the client first contacts the NameNode and access a list of locations for the block that comprise the files and these locations identify the DataNodes which hold each block.
- Clients then read file data directly from the DataNode servers in parallel. So the NameNode is not directly involved in the bulk data transfer keeping its overhead to a minimum.
- NameNode information should be preserved even if the NameNode machine fails.
- NameNode failure is more severe for the cluster than DataNode failure.
- When individual data nodes may crash and the entire cluster will continue to operate, the loss of the NameNode will render the cluster inaccessible until it is manually restored.

Features of HDFS :-

9

HDFS is file system designed for storing have some key characteristics. They are

- ① Support for very large files.
- ② commodity Hardware
- ③ streaming data access
- ④ High-latency data access
- ⑤ lots of small files
- ⑥ multiple writers, arbitrary file modifications
- ⑦ moving computation is than moving data.

By Mr. GOPAL KRISHNA

① Support for very large files :-

files that are ^{support} hundreds of megabytes, Giga By Tera Bytes in size.

Hadoop clusters are running today that store PB, ZB of data. ^{set} of nodes.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 998 570 6789.

② commodity Hardware :-

- HDFS requires a commodity hardware [The H/W which is been available for most of the vendors and Hadoop does not require high configuration hardware, expensive SW to be part of it's basic installation.]

- HDFS is always working without a noticeable interruption to the user. [In the face of failures]

- For the commodity hardware chance of node failure across the cluster is high, at least for large clusters.

3) Streaming data access :-

- HDFS is the most efficient data processing pattern. That pattern is write once/read many times.
- HDFS will follow streaming data access [Sequential flow]

Ex:-

file data
1
2
3
:
1,00,000

→ search for 6564 record.

→ Search come from 1 to 8564 Sequentially
(no random, index access)

Ex:- update



if we can update the 500 record. It is not possible here. Take into local, update the file.

- sequential flow means, when we can save the file on HDFS file, we can read the data (or) write the data in the sequential fashion only.
- HDFS having the feature of read many times and write once.

write once/
read many
times

Ex:-

update the
file

Take to local

A.txt
Initialize

local
linux
environment

put it
back by
deleting A.txt

HDFS environment

once we can store a file on top of HDFS file, we can not override that file, instant we can take the file into local environment, if we modify the content and we can put the same file on top of the HDFS file. By deleting whole HDFS file.

④ High - latency data access :-

10

- Generally applications that require low - latency access data.
- As we are using in HDFS, used to very large amount of data. Because to take more time - this may be at the expense of latency.

Eg:- Hadoop (Sequential Access)

hive >
↓
high latency
↓
12 min
↓
↓
↓

Here HDFS is used to large amount of data But to take more time.

In the tens of milliseconds with HDFS.

ROBMS (Randomly Access)

SQL > ↓
↓ 5 min
↓
low latency

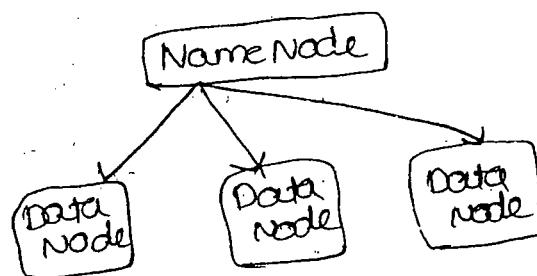
Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarpet, Hyderabad-500 016.
Ph: 040-6462 6789, 988 570 6789

range, will not work well

By Mr. GOPAL KRISHNA

⑤ lots of small files :-

- NameNode is responsible for maintaining the metadata information of the hadoop file system.
- hadoop file system have no. of files in the amout of memory on the NameNode.



⑥ multiple writers, arbitrary file modifications:

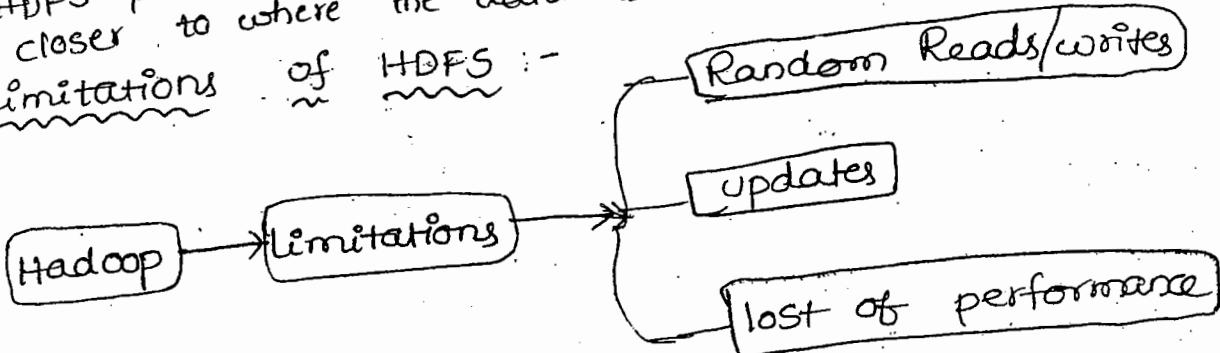
- files in HDFS may be written to by single writer.
- writers are always made at the end of the file.
- There is no support for multiple writers (or) modifications at arbitrary offset in the file.
[These might be supported in the feature, but they are likely to be relatively]

⑦ moving computation is easy than moving data:-

In Traditional distributed system data from SAN (Shared Access Network) device will be copied to.

- A computation requested by an application is much more efficient if it is executed near the data it operates on. This is especially true when the size of the data set is huge.
- The assumption is that it is often better to migrate the computation closer to where the data is located rather than moving the data to where the application is running.
- HDFS provides interfaces for application to move themselves closer to where the data is located.

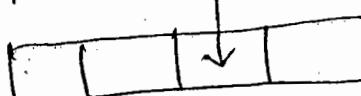
Limitations of HDFS :-



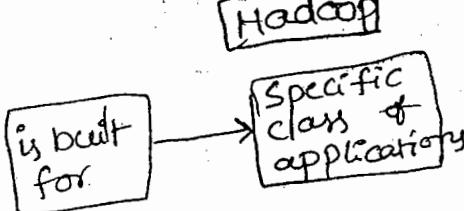
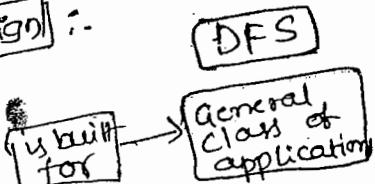
Sequential Read



Random Read



Design :-



Xcell Technologies
Amritapet, Hyderabad - 500 016
Ph: 040-6462 6789, 08 570 6789

Hadoop Architecture

11

→ Hadoop Architecture will be classified into 5 daemons.

① NameNode and it's functionality.

② DataNode and it's functionality.

③ JobTracker and it's functionality.

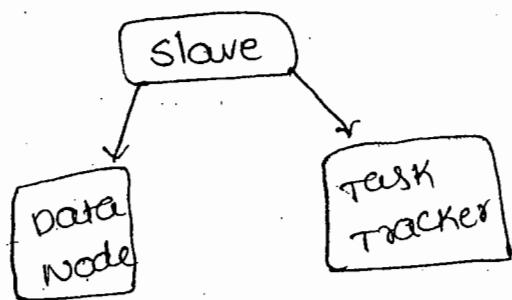
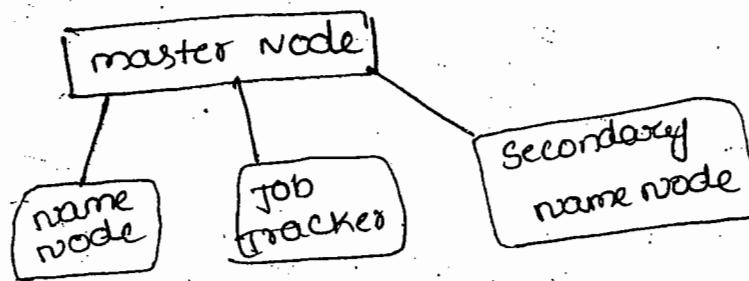
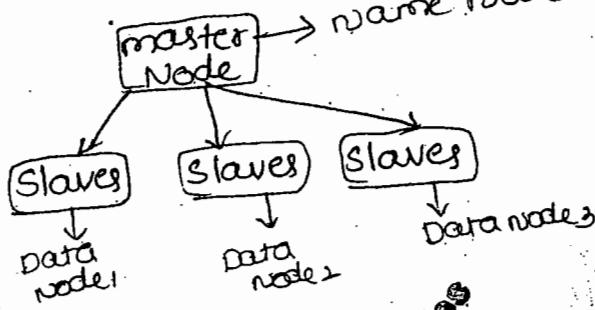
④ TaskTracker and it's functionality.

⑤ Secondary NameNode and it's functionality.

By Mr. GOPAL KRISHNA

→ Hadoop architecture follows master-slave architecture.

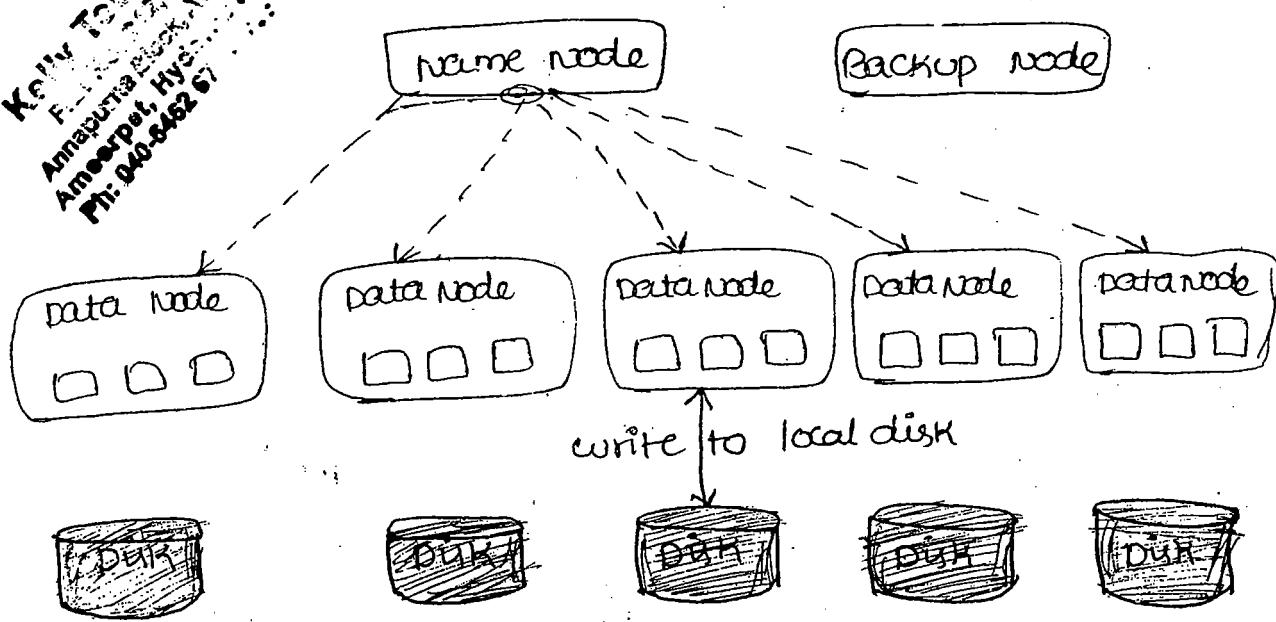
Data
only
raw data



Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-4462 6789, 998 570 6789

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-4462 6789, 998 570 6789

HDFS Architecture



NameNode :-

- The master node in Hadoop architecture it is treated as name node.
- NameNode is responsible for maintaining the metadata information of the Hadoop file System.
 - process the data
 - data about data.
- NameNode maintains the file System. The file System have metadata for all the files & directories. This information is stored persistantly on the local disk.
 - permanently.
- NameNode maintain the file system Namespace.
- The NameNode executes the Namespace file, The file operations like opening, closing & renaming files & directories.
- The NameNode will only update two important permanent files of Hadoop file System called the
 - fs namespace Image
 - edit log
- Namespace access to files by clients.

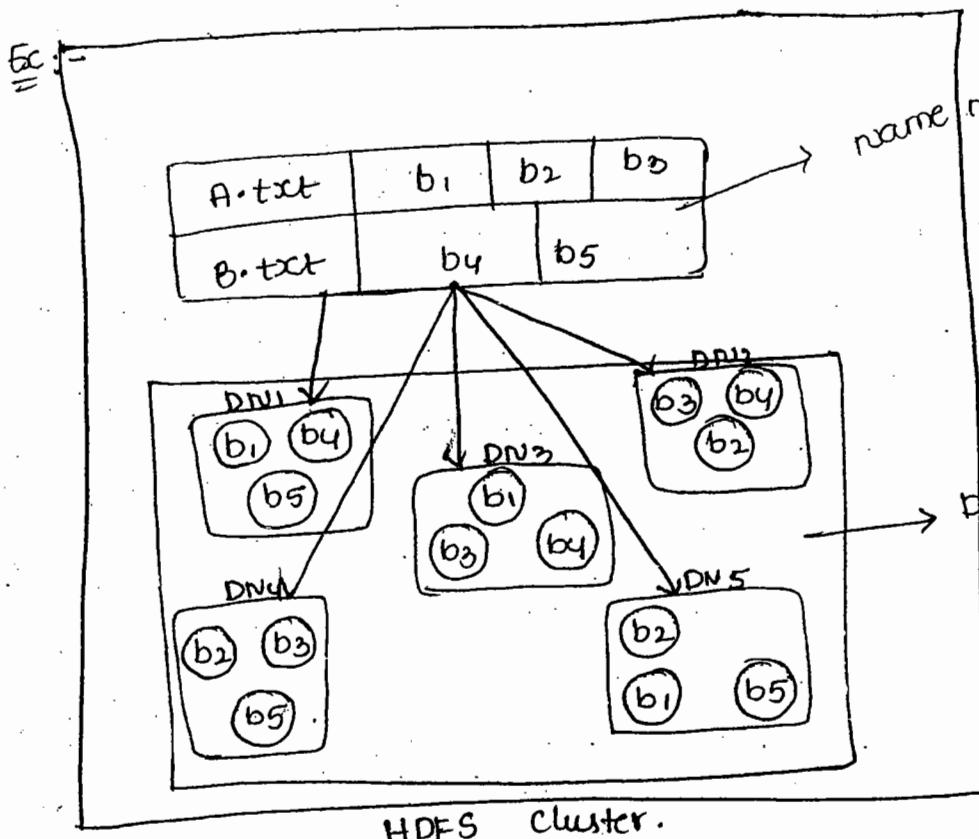
physical locations
of each & every
block of hadoop
cluster

By Combining which
blocks so and so
file will generated

12

meta data
Information

By Mr. GOPAL KRISHNA



Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave
Ameerpet, Hyderabad - 500 083
Ph: 040-5462 5789, 9880 6765

Flat No. 1A
Annapurna Block, Aditya Enclave
Ameerpet, Hyderabad - 500 083
Ph: 040-5462 5789, 9880 6765

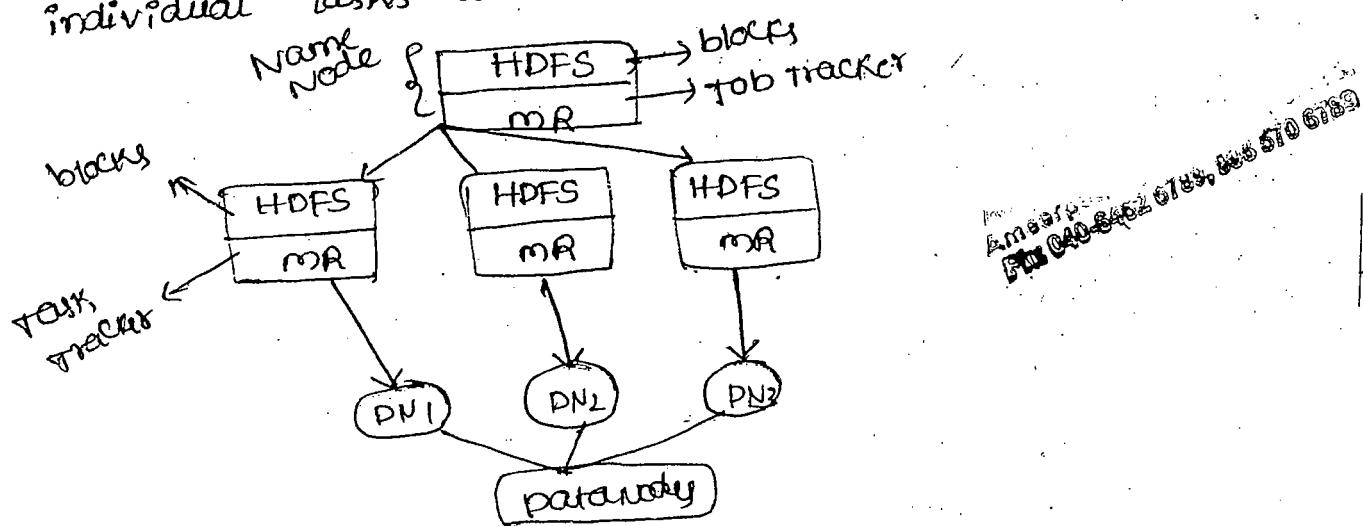
- A client access the file System on behalf of user by communicating with the Namenode.
- The client access file system to a portable operating System Interface (POSIX). So, the user code does not to know about the NameNode.
- Assign Blocks to datanodes.
- Keeps track of live nodes (through heartbeats)
- Limitates re-replication in case of data node loss.
- Block metadata is held in memory.
 - will run out of memory when too many files exist.
- Is a single point of failure in the System
 - some solutions exist

② Data Node :-

- Data Node is a place hold of the data i.e. actual data in datanodes only in the form of HDFS blocks. [By default each block size is 64MB]
not the metadata.
- Data Nodes are store and retrieve blocks, reporting to Namenode.
- A client access the file system on behalf of the user by communicating with the Data Nodes.
- uses underlying regular FS for storage (e.g.: ext3)
 - takes care of distribution of blocks across disks
 - Don't use RAID
 - more disks means more IO throughput.
- Does not know about the rest of the cluster (Shared nothing)

③ Job Tracker :-

- Job Tracker is one of the SW daemons for Hadoop Architecture.
- Job Tracker is responsible for scheduling & Rescheduling the tasks in the form of mapReduce jobs.
- Job Tracker is also getting the acknowledgement [response] back from the Task Tracker.
- Generally JobTracker will decide on top of the NameNode.
- JobTracker manages the mapreduce jobs, distributes individual tasks to machines running the Task Tracker.



④ Task Tracker :-

13

- Task Tracker is responsible for instantiating & monitoring individual map & reduce works.
- Task Tracker is also known as Slaves daemon for Hadoop architecture.
- Task Tracker is primarily responsible for executing the tasks assign by the JobTracker in the form of MRJOBS.
- Generally Tasktracker will resides on top of the data nodes.

By Mr. GOPAL KRISHNA

Note :-

- ① JobTracker & Task Tracker are the two important daemons of Hadoop architecture, which are completely responsible for the processing of the data by the means of mapReduce programming.

⑤ Secondary NameNode :-

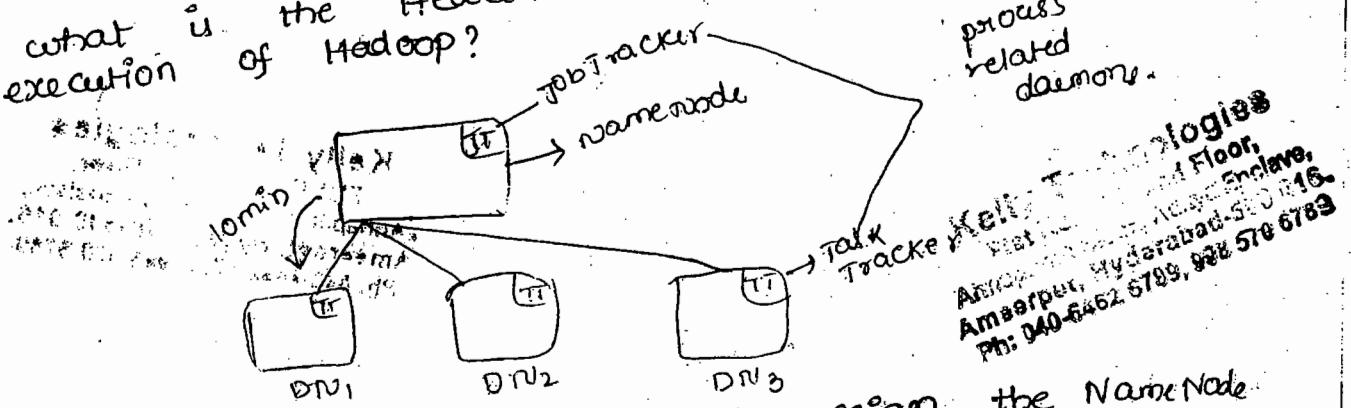
- Secondary NameNode is deprecated.
- It performs periodic checkpoints of the namespace and helps keep the size of file containing log of HDFS modifications within certain limits at the Namenode.
- It is replaced by Checkpoint node.
- Secondary NameNode will acts as separate physical file.
- whenever the primary node is down in Hadoop architecture, the Secondary NameNode will come into picture.
- we can never ever call a Secondary NameNode as the replica our primary Node.
- Secondary NameNode is responsible for only read the fsImageSpace & editlog.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-2462 6789, 098 570 0789

- > what is the checkpoint mechanism?
 - Through the checkpoint mechanism only Hadoop cluster will make sure that all the metadata information of Hadoop file system will get updated in the two persistence files
 - * fsnameSpaceImage
 - * editlog
- The checkpoint mechanism can be configured either an hourly, weekly, daily.

The checkpointing module periodically creates checkpoints of the namespace.

- what is the HeartBeat mechanism & speculative execution of Hadoop?

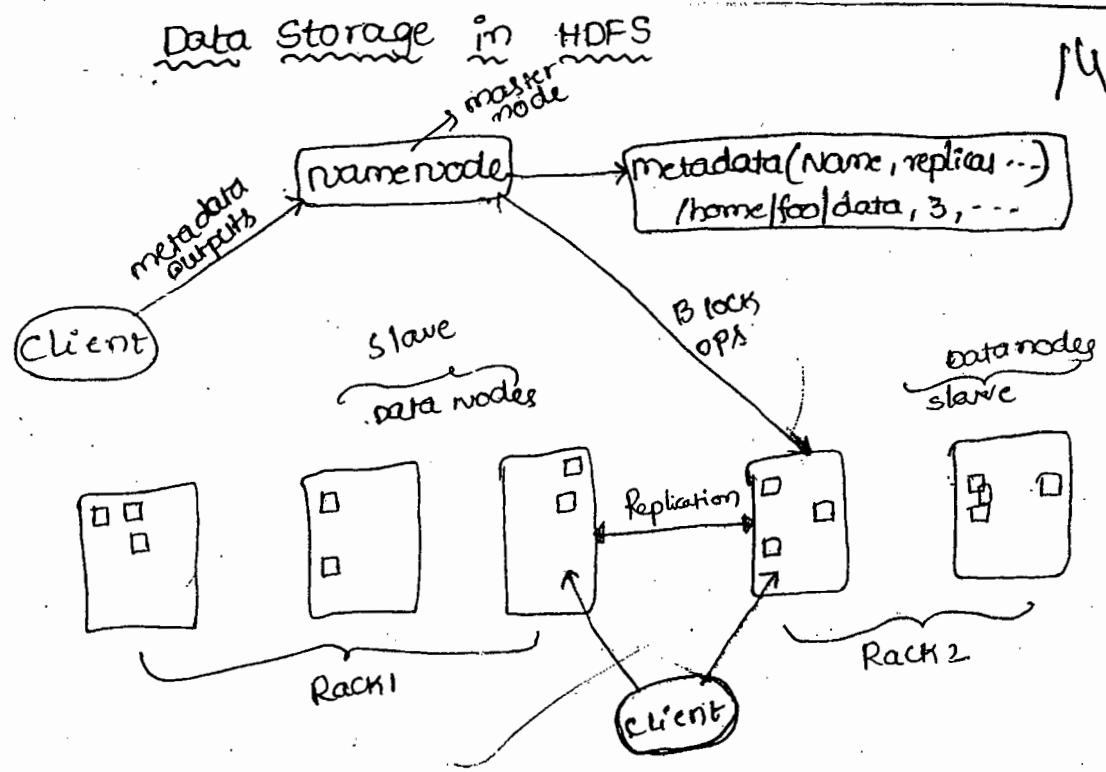


when a particular jobtracker is assigned the NameNode will expects a kind of acknowledgement from the datanodes [task tracker] and regular interval bases. This is known as HeartBeat mechanism

if a datanode fails to response within 10 min time the namenode will mark that node either it is dead, running slow (or) not functionality slow.

Immediately the jobtracker will assign the same task to some other Idle nodes with hadoop cluster. This is known as speculative execution of Hadoop.

By the means of speculative execution of hadoop only the end user will not feel any delay will setting the response back from Hadoop processing.



- HDFS Architecture is a master-slave architecture.
- In HDFS cluster consists of master node is a NameNode.
- NameNode manages the file system name space regulates access to by clients.
- In HDFS, file System Name Space allows user data to be stored in files. Internally file is split into blocks. Blocks are stored into set of DataNodes.
- NameNode executes the file system namespace operations like opening, closing & renaming files & directories. It also determines the mapping of blocks to DataNodes.
- The DataNodes are responsible for serving read & write requests from the file system clients.
- The DataNodes also perform block creation, deletion and replication upon instruction from the NameNode.

Introduction about blocks :-

- HDFS is designed to support very large datasets.
- HDFS supports write once/read many times semantics on files.
- In HDFS data is split into blocks and distributed across multiple nodes in the cluster.
- Each block is typically 64MB (or) 128 MB in size.
- Each block is replicated multiple times. By default replication factor of block is 3 times. Replicas are stored on different data nodes.
- HDFS utilizes the local file system to store each HDFS block as a separate file.
- HDFS blocksize can not be compared with the traditional file system blocksize.

Data Replication:-

- An application can specify the no. of replicas of a file.
- The replication factor can be specified at file creation time and can be changed later.
- The replication factor are configurable per hdfs-site.xml.
- Default replication factor of file - 3

$$\text{no. of copies of a file} = \text{Replication factor of that file}$$

[hdfs - site . xml]

```
<name>dfs.replication</name>
<value>3</value>
```

variable specifies how many times each HDFS block should be replicated.

Kelly Technologies
Flat No. 212, 2nd Floor,
Ampara Block, Adyar Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6452 6719, 984570 6789

- The placement of the replicas is critical to HDFS reliability and performance.
- Optimizing replica placement distinguishes HDFS from other distributed file systems.
- Rack-aware replica placement
goal: improve reliability, availability and N/W bandwidth utilization.
- Research topic.
- In many racks, communication between racks are through switches.
- Network bandwidth between machines on the same rack is greater than those in different racks.
- NameNode determines the rackid for each datanode.
- Replicas are typically placed on unique racks
 - simple but non-optimal
 - writes are expensive
 - replication factor is 3
 - Another research topic?
- Replicas are placed one on a node in a local rack, one on a different node in the local rack, one on a node in a different rack.
- $\frac{1}{3}$ of the replica on a node, $\frac{2}{3}$ on a rack and $\frac{1}{3}$ distributed evenly across remaining racks.

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Atmapuram Block, Ameerpet,
 Hyderabad, 500 016.
 Ph: 040 462 6789, 098 570 6789

By Mr. GOPAL KRISHNA

Replica selection :-

- Replica selection for READ operation: HDFS tries to minimize consumption and latency.
- if there is a replica on the Reader node then that is preferred.
- HDFS cluster may span multiple data centers: replica center is preferred over the remote one.

→ HDFS by interacting with the interfaces have two best approaches. They are

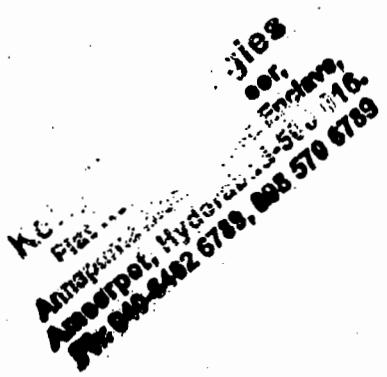
① Commandline Interface.

② JavaBased Approach.

i) Commandline Interface :-

- Commandline Interface is one of the simplest and many developers, the most familiar.
- Commandline Interface is the interactive Shell.

ii) Java Based Approach :-



The persistence of file System metadata :-

1/6

- In the HDFS, namenode is responsible for the metadata. metadata have file System name space & edit log.
- namenode will only update the **fsImage**
 - ① fsImage
 - ② Edit log.

Kelly Technologies
Flat No. 212, 2nd Floor
Annapurna Block, Antara Enclave
Ameerpet, Hyderabad - 500 018.
Ph: 040-6452 6729, 984570 6739

fsImage :-

- the fsimage file is a persistent checkpoint of the file system metadata.
- However, it is not updated for every file system write operation, since writing out the fsimage file, which can grow to be gigabytes in size, would be very slow.
- this does not compromise resilience.

By Mr. GOPAL KRISHNA

when Namenode fails :-

- if the name node fails, then the latest state of it's metadata can be reconstructed by loading the fsimage from disk into memory, then applying each of the operations in the edit log.
- In fact, this is precisely what the namenode does when it starts up [Learn about safe mode]

Edit Log :-

- when a client system performs a write operation (such as creating or moving a file) it is first recorded in Edit log.
- The namenode also has in memory representation of the file system metadata, which it updates after the editlog has been modified.
- Edit log to permanently record every change that occurs to file system metadata.

when a client request a data:-

- The in-memory metadata is used to serve read requests.
- The editlog is flushed and synced after every write before a success code is returned to the client.
- For Namenodes that write to multiple directories, the write must be flushed and synced to every copy before returning successfully.
- This ensures that no operation is lost due to machine failure.

Ex:- * creating the new file in HDFS.

- * NameNode to insert a record into the editlog.
- * Similarly, changing the replication factor of a file causes a new record to be inserted into the editlog.
- * The NameNode uses a file in its local host OS file system to store the editlog.

Note:

- HDFS namespace is stored by the NameNode. The NameNode uses a transaction log called the editlog.
- HDFS namespace including the mapping of blocks to files and file system properties is stored in a file called FSImage.
- FSImage is stored as file in the NameNode's local file system too.
- The NameNode keeps an image of the entire file system namespace and file block map in memory.

Safemode :- whenever cluster is starting up in Hadoop certain things will be done by the NameNode. Some thing

① loading all system related configuration files.

② check for the satisfactory replication for the data.

③ All system related dependent files.

- while doing these above all operations. The NameNode is HDFS puts on read-only-mode [HDFS can not be reached at that moment] this stage is known as safemode ON
- After doing all these stuff automatically safemode will come out of safemode ON → OFF which means that HDFS will be accessible mode.
- But sometimes safemode will not be turned into OFF mode. At that point of time Hadoop administrator to run the below command.

By Mr. GOPAL KRISHNA

hadoop dfsadmin -safemode LEAVE ↵

→ safemode is in OFF

hadoop dfsadmin -safemode ENTER ↵

→ safemode is in ON.

hadoop dfsadmin -safemode get ↵

→ safemode is OFF

hadoop dfsadmin -report ↵

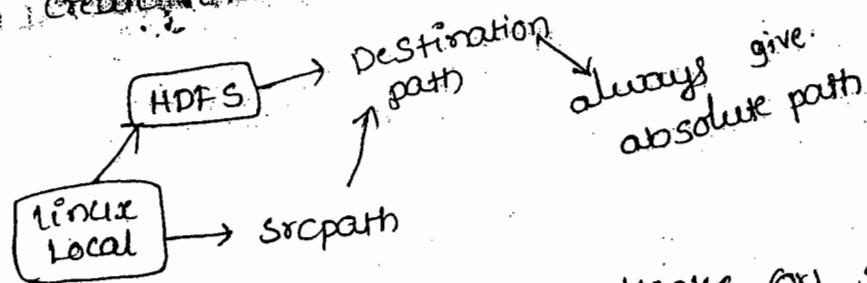
→ looks like the GUI mode ↵
http://127.0.0.1:50070/ ↵

↓
cluster health report command.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6452 6779, 986 572 6739

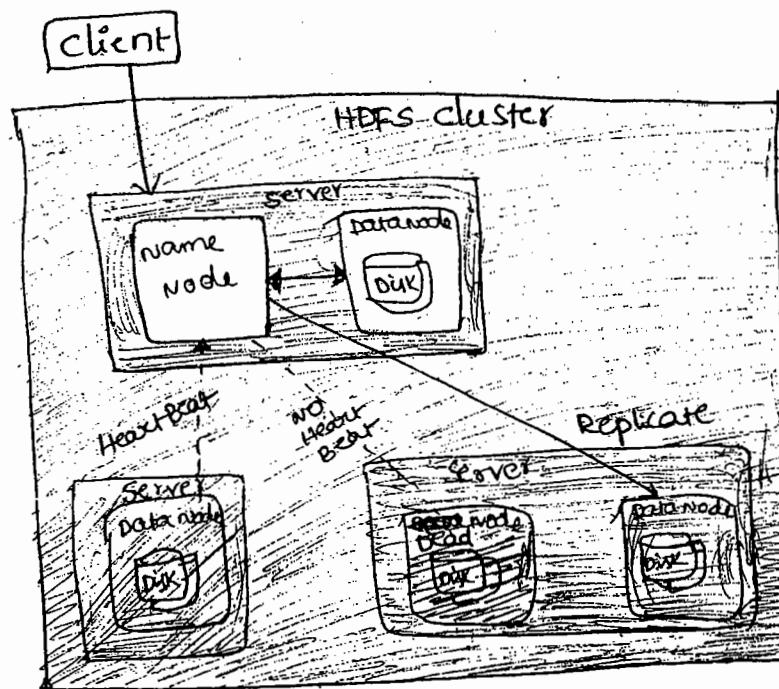
Hadoop file System

- hadoop fs is indicating the compiler to interact with Linux local environment to HDFS environment.
- hadoop fs is not support the -vi command. Because HDFS is writeonce.
- hadoop fs is support -touchz command.
- hadoop fs looks only HDFS directories & file but not for localdirectories & files.
- In HDFS file overloading is not possible.
- if the one of the file moving one environment to another environment there have only default privileges.
- we can not create file on top of HDFS.
- we can't create the file on local directory.



- hadoop fs does not support hardlinks or softlinks.
- hadoop fs does not implement user quotas.
- Error information is sent to stderr & output is sent to stdout.

Kelly Technologies
Flat No 1002, 10th floor,
Rajiv Gandhi Infotech Park,
Banjara Hills, Hyderabad-500 036.
Andhra Pradesh, India
Ph: 040-3452 6789, 988 670 5789

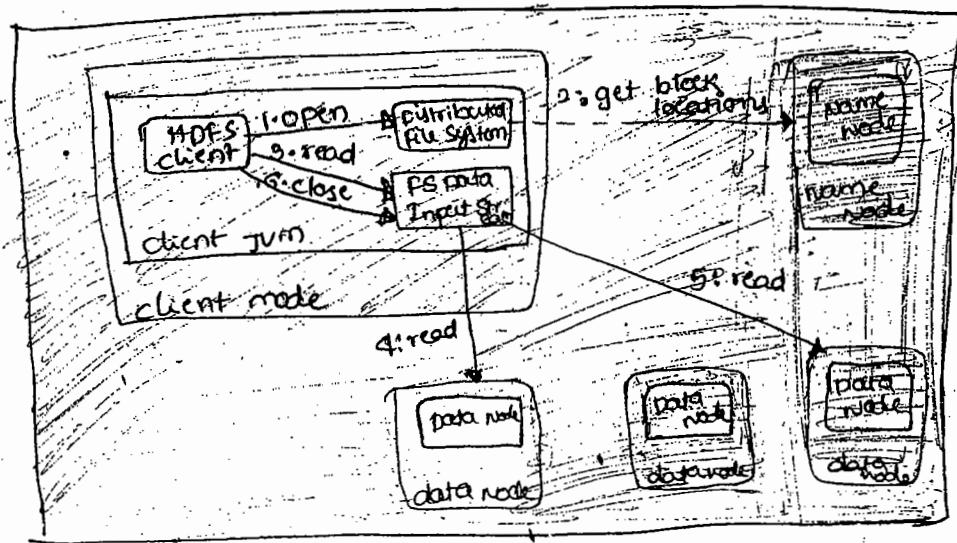


By Mr. GOPAL KRISHNA

Kelly Technologies
Flat No. 212, 2nd Floor,
Amarpura Block, JntuHd, Hyderabad-500 016,
Pmk Opposite 6789, 988 570 0789

HDFS Debugging Steps :-

Client Reading Data from HDFS



- The client opens the file it wishes to read by calling `open()` on the `File System` object, which for HDFS is an instance of DFS.
- DFS calls the namenode, using RPC, to determine the locations of the blocks for the first few blocks in the file. For each block, the namenode

returns the addresses of the datanodes that have a copy that block and the datanodes are stored according to their proximity to the client.

- The DFS returns a FSDataInputStream (an input stream that supports file seeks) to the client for it to read data from.

The client then calls read() on the stream.

- DFS Stream connects to the first (closest) datanode for the first block in the file.

Data is streamed from the datanode back to the client, which calls read() repeatedly on the stream.

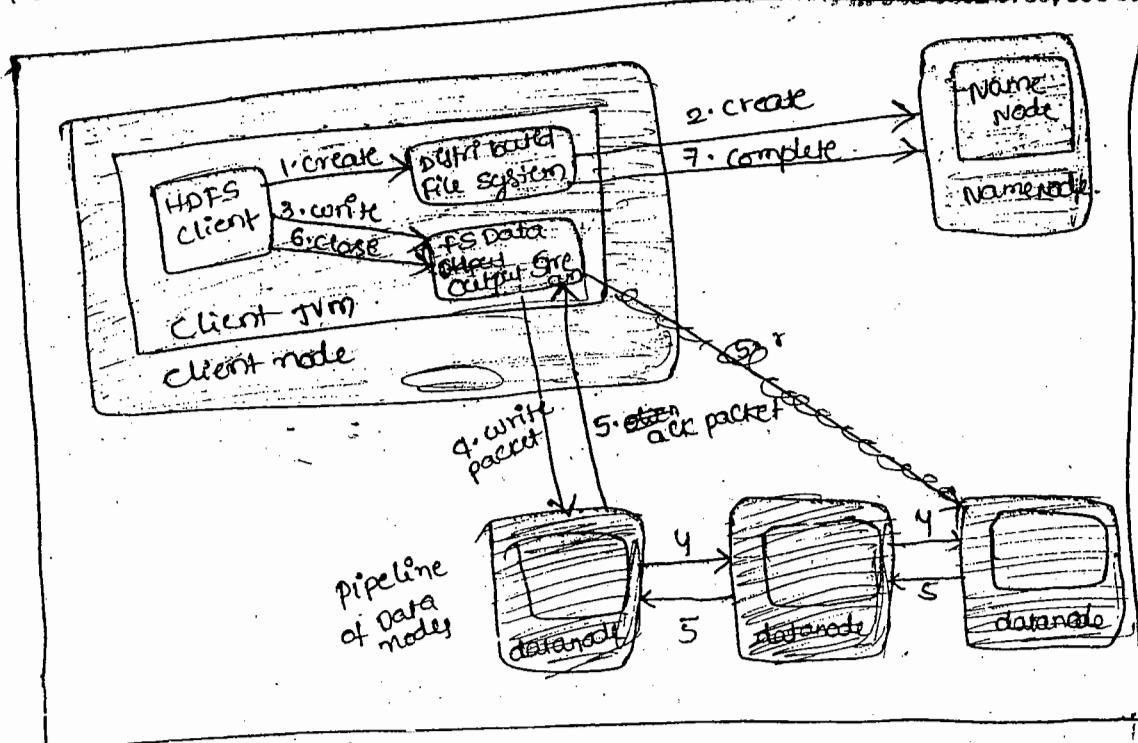
when the end of the block is reached, DFSInputStream will close the connection to the datanode, then find the best datanode for the next block.

when the client has finished reading, it calls close() on the FSDataInputStream.

Client Writing Data to HDFS

Kelly Technologies

Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-2462 6789, 098 510 6789



- The client creates the file by calling `create()` on DFS.
- DFS makes an RPC call to the namenode to create a new file in the filesystem's namespace, with no blocks associated with it.
- The namenode performs various checks to make sure the file doesn't already exist, and then the DFS returns a `FSDataOutputStream` for the client to start writing data to.
- As the client writes data, DFS output Stream splits it into packets, which it writes to an internal queue, called the data queue.

By Mr. GOPAL KRISHNA

- The data queue is consumed by the DataStreamer whose responsibility it is to ask the namenode to allocate new blocks by picking a list of suitable datanodes to store the replicas.
- The list of datanodes forms a pipeline.
- The DataStreamer streams the packets to the first datanode in the pipeline, which stores the packet and forwards it to the second datanode in the pipeline. Similarly, the second datanode stores the packet and forwards it to the third (and last) data node in the pipeline.

DFSOutput Stream also maintains an internal queue of packets that are waiting to be acknowledged by datanodes, called the ack queue.

A packet is removed from the ackqueue only when it has been acknowledged by all the datanodes in the pipeline.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 018
Ph: 040 6462 6789, 998 570 6789

what if a datanode fails while it is writing to it.

- ② The pipeline is closed.
- ③ Any packets in the ackqueue are added to the front of the data queue so that datanodes that are downstream from the failed node will not miss any packets.
- ④ The partial block on the failed datanode will be deleted if the failed datanode recovers later on the failed data node is removed from the pipeline and the remainder of the block's data is written to the two good datanodes in the pipeline.
- ⑤ The namenode notices that the block is under-replicated, and it arranges for a further replica to be created on another node.

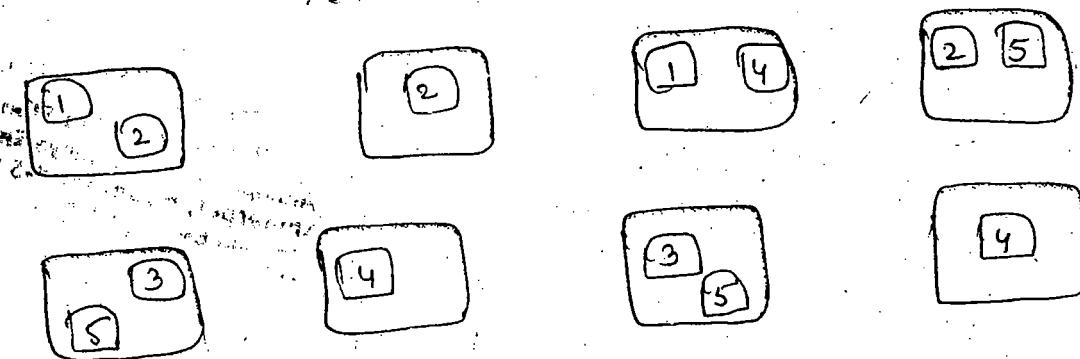
Replica placement

Block Replication

Kelly Technologies
Flat No. 272, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 986 570 6789

```
namenode(filename, numReplicas, blockSize --  
/users/root/data/part-0, r:2, {1,3} --  
/users/root/data/part-1, r:3, {2,4,5} --
```

Data Nodes



Qn
• The placement of replicas is critical to HDFS reliability and performance.

Ans

- The purpose of rack-aware replica placement policy is to improve 3 different ways. They are
 - ① Readability
 - ② availability
 - ③ N/w bandwidth utilization.

By Mr. GOPAL KRISHNA

- communication between two nodes in different racks has to go through switches. In most cases, network bandwidth between machines in the same rack is greater than bandwidth between machines in different racks.

- HDFS placement policy is to put one replica on one node in the local rack, another on a node in a different rack, and last on a remote rack.

Kelly Technologies
Plot No. 212, 2nd Main,
Amarpetta Block, Adyar
Amarpetta, Hyderabad - 500 026
Ph: 040-6462 6789, 043 570 6720

Replica Section :-

- To minimize global bandwidth consumption and read latency.
- HDFS tries to satisfy a read request from a replica that is closer to the reader.
- if there exists a replica on the same rack as the reader node, then that replica is preferred to satisfy the read request.
- if any HDFS cluster spans multiple data centers, then a replica that is resident in the local data center is preferred over any remote replica.

Replication pipelining :-

- Data is pipelined from datanode to next datanode.

For example:-

- When a client is writing data to an HDFS file, it's data is first written to a localfile.
- Suppose HDFS file has a replication factor is 3.
- The client retrieves a list of datanodes from the NameNode.
- The list of datanodes will host a replica of that block.
- The client then flushes the data blocks to the first datanode.
- The first datanode starts receiving the data in small portions (4 KB) writes each portion to its local repository and transfers the portion to the second datanode in the list.
- The second datanode starts receiving each portion of the block, writes each portion to its repository and then flushes that portion to the 3rd datanode.
- Finally, the datanode writes the data to its local repository. One datanode can be receiving one in the pipeline and at the same time forwarding data to the next one in the pipeline.

Communication protocols :-

- All HDFS communication protocols are based on top of the TCP/IP protocol.
- A client establishes a connection to a configurable TCP port on the NameNode machine. It talks the client protocol with the NameNode.
- The datanodes talk to the NameNode using the DataNode protocol.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad, Telangana 500 016.
Ph: 040-5462 0100, 040 570 6700

- A Remote Procedure call (RPC) abstraction wraps both the client protocol and the datanode protocol.
- By design, NameNode never initiates any RPCs. Instead, it only responds to RPC requests issued by Datanodes (or) clients. 21

By Mr. GOPAL KRISHNA

Robustness:-

- The primary objective of HDFS is to store data reliably even in the presence of failures. The three common types of failures are
 - ① NameNode failures
 - ② Data Node failures
 - ③ N/w partitions.

Kelly Technologies
Flat No. 212, 2nd Floor,
Nampalli Block, Addya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6462 6719, 55-570 6789

Data Disk failure, Heartbeats & Re-Replication:-

- Each data node sends a Heart-beat message to the NameNode periodically.
- A new partition can cause a subset of Datanodes to lose connectivity with the NameNode.
- The NameNode detects this condition by the absence of a Heartbeat message.
- The NameNode marks Datanodes without recent Heartbeats as dead Datanode if not available to HDFS any more.
- Datanode death may cause the replication factor of some blocks to fall below their specified value.
- The NameNode constantly tracks which blocks need to be replicated and initiates replication whenever necessary.
- The necessity for re-replication may arise due to many reasons. (1) Datanode may become unavailable, a replica may become corrupted, a harddisk on a datanode may fail, or the replication factor of a file may be increased.

cluster Rebalancing :-

- The HDFS architecture is compatible with data rebalance schemes.
- A scheme might automatically move data from one datanode to another datanode.
- if the freespace on a datanode falls below a certain threshold.
- if in the event of sudden high demand for a particular file, a scheme might dynamically create additional replicas and rebalance other data in the cluster. These types of data rebalancing schemes are not yet implemented.

Data Integrity :-

- Data integrity is possible that a block of data fetched from a datanode arrives corrupted. This corruption can occur because of faults in storage device, N/w faults, or buggy SW.
- The HDFS client SW implements checksum checking on the contents of HDFS files.
- When a client creates an HDFS file, it computes a checksum of each block of the file and stores these checksums in a separate hidden file in the same HDFS namespace.
- When a client retrieves file contents it verifies that the data it received from each datanode matches the checksum stored in the associated checksum file.
- If not, then the client can opt to retrieve that block from another datanode that has a replica of that block.

Kelly T
Sathya
Anupama B
Rohit
Abhishek
Ph. 044 42267712, 094 570 670

metadata disk failure :-

22

- FSImage & Edit log are central data structures of HDFS. A corruption of these files can cause the HDFS instance to be non-functional.
- For this reason, the NameNode can be configured to support maintaining multiple copies of FSImage & EditLog.
- The NameNode will update two important persistence files of HDFS.

① FSImage ② EditLog.

to get updated synchronously.

By Mr. GOPAL KRISHNA

- Synchronous updating of multiple copies of the FSImage & EditLog may degrade the rate of name space transactions per second that a NameNode can support.
- However, this degradation is acceptable because even though HDFS applications are very data intensive in nature, they are not metadata intensive.
- When a NameNode restarts, it selects the latest consistent FSImage & EditLog to use.
- The NameNode machine is a single point of failure for an HDFS cluster.
- If the NameNode machine fails, manual intervention is necessary. Currently automatic restart and failover of the NameNode software to another machine is not supported.

Snapshots :-

- Snapshots support storing a copy of data at a particular instant of time.
- one usage of the snapshot feature may be to roll back a corrupted HDFS instance to a previously known good point in time.
- HDFS does not currently support snapshots but will in a future release.

Data Organization

① Data Blocks :-

- HDFS is designed to support very large files.
- HDFS supports write once & read many times semantics on files.
- A typical blocksize used by HDFS is 64MB. Thus an HDFS file is chopped up into 64MB chunks, and if possible, each chunk will reside on a different data node.

② Staging :-

- A client request to create file does not reach the NameNode immediately.
- In fact, initially the HDFS client caches the file data into a temporary localfile.
- Application writers are transparently redirected to this temporary local file.
- When the localfile accumulates data worth over one HDFS blocksize, the client contacts the NameNode.
- The NameNode inserts the file name into the file system hierarchy and allocates a data block for it.

Kelly Technologies
First Floor, 2nd Main,
1st Cross, Jayanagar 12th Block,
Bengaluru - 560055
Ameerpet, Hyderabad - 500016.
Ph: 040-64526789, 098 570 6789

- The NameNode responds to the client request with the identify of the Datanode and the destination data block.
- then the client flushes the block of data from the local temporary file to the specified Datanode.
- when a file is closed, the remaining un-flushed data in the temporary localfile is transferred to the Datanode.
- The client then tells the NameNode that the file is closed. At this point, the NameNode commits the file creation operation into a persistent store
- if the NameNode dies before the file is closed, the file is lost.

By Mr. GOPAL KRISHNA

The above approach has been adopted after careful consideration of target applications that run on HDFS.

- These applications need streaming writes to files. if a client writes to a remote file directly without any client side buffering, the N/w speed and the congestion in the N/w impacts throughput considerably.
- This approach is not without precedent. Earlier DFS's [ex: AFS] have used client side caching to improve performance]

[POSIX → requirement has been relaxed to achieve higher performance of data uploads].

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-5462 6722, 570 6789

Accessibility

- HDFS can be accessed from applications in many different ways.
- HDFS provides a JAVA API for applications to use.
- HTTP Browser can also be used to browse the files of an HDFS instance.
- work is in progress to expose HDFS through the webDAV protocol.

FS Shell :-

- HDFS allows user data to be organized in the form of files and directories. It provides a command line interface called FS Shell which user interact with the data in HDFS.
- The syntax of this command set is similar to other shells (e.g. bash, csh) that users are already familiar with.
- FS shell is targeted for applications need a scripting language to interact with the stored data.

DFS Admin :-

- The DFS Admin command set is used for administering an HDFS cluster. These are commands that are used only by an HDFS Administrator.

Browser Interface :-

- A typical HDFS install configures a webserver to expose the HDFS namespace through a configurable TCP port. This allows a user to navigate the HDFS namespace and view the contents of its files using a web browser.

Space Reclamation [File deletes & undeletes] :-

24

- when a file is deleted by a user or an application it is not immediately removed from HDFS.
- HDFS first renames it to a file in the /trash directory.
- The file can be restored quickly as long as it remains in /trash.

By Mr. GOPAL KRISHNA

- A file remains in /trash for a configurable amount of time. After the expiry of its life in /trash, the namenode deletes the file from the HDFS namespace.
- The deletion file causes the block Kelly Technologies with the Flat No. 212, 2nd Floor, Annapurna Block, Aditya Enclave, Antespet, Hyderabad-500 016. Ph: 040-6452 6789, 989 570 6789 to be freed.

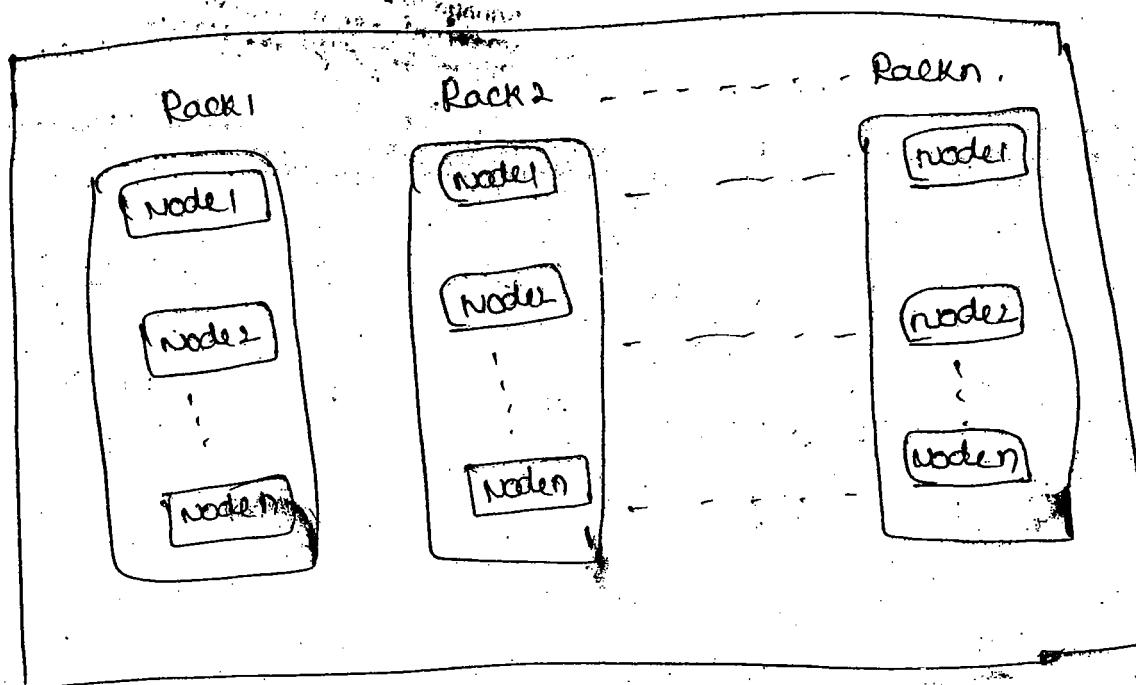
[Note:- There could be an appreciable time delay b/w, the time a file is deleted by a user & the time of the corresponding increase in free space in HDFS.]

- A user can undelete a file after deleting it as long as it remains in the /trash directory. If a user wants to undelete a file that he/she has deleted, he/she can navigate the /trash directory and retrieve the file. The /trash directory contains only the latest copy of the files that were deleted.
- The /trash directory is just like any other directory with one special feature: HDFS applies specified policies to automatically delete files from this directory.
- The current default policy is to delete files from /trash that are more than 6 hours old.
- In the future, this policy will be configurable through a well defined interface.

Decrease Replication factor :-

- when the replication factor of files is reduced, the NameNode selects excess replicas that can be deleted.
- The next heartbeat transfers this information to the DataNode.
- The DataNode then removes the corresponding blocks and the corresponding free space appears in the cluster.
- once again, there might be time delay between the completion of the Set Replication API call and the appearance of free space in the cluster.

Hadoop cluster :- (set of machines)



Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016
Ph: 040-6462 6789, 995 570 6789

Linux Commands

25

File/Directory Basics

① mkdir → making the directory

Ex:- mkdir anitha ↵

② cd → change the directory.

Ex:- cd anitha ↵

/anitha #

By Mr. GOPAL KRISHNA

③ clear → clear the commands.

④ pwd → present working directory.

⑤ date → display the date.

⑥ who am i → display the user name.

⑦ rm → Delete files.

⑧ ls → List files

⑨ cp → copy files

Ex:- cp sample.txt
↓
src path

Sample.txt
↓
destination path:

(or) → CP -r dir1 dir2
→ CP -r dir1/* dir2
→ CP -r dir1/* .txt dir2

⑩ mv → renaming the files

Ex:- mv sample.txt input.txt
↓ ↓
src path renaming path.

⑪ ln → link files

Ex:- ln sample.txt ↵
ln: creating hardLink .| sample.txt

⑫ rmdir → remove the directory

Ex:- rmdir anitha ↵
delete the directory.

Kelly's
Flat No. 1001
Annamma Block A-77 Flats,
Amarpet, Hyderabad-500 014
Ph: 040-2452 6739, 984570 6739

File viewing

① cat → view files

Ex: cat input.txt ↴
This is a Hadoop Note Book.

② less → Page through files.

Ex: less input.txt ↴
This is a Hadoop Note Book.
input.txt (END)

③ head → view file beginning [1st 10 records displayed]

Ex: head input.txt.

④ tail → view file ending last 10 records display.

Ex: tail input.txt.

⑤ nl → number lines

nl input.txt ↴

⑥ od → view binary data.

od input.txt ↴

⑦ xxd → view hexa binary data.

xxd input.txt

⑧ gv → view Post Script / PDF files.

⑨ xdvi → view Tex dvi files.

File Creation & Editing

① emacs → Text editor

⑥ abiword → Edit word doc documents.

② pico → Text editor

⑦ gnumeric → Edit excel documents,

③ vi → Text editor

④ umask → set default file protection.

⑤ soffice → edit word/excel
power point docs.

File properties

26

- ① stat → Display file attributes.

Ex:- stat input.txt ↵

File:	- - -	Device:	- - -
Size:	- - -	Accesses:	- - -

- ② wc → Count bytes/words/lines

Ex:- wc input.txt ↵

1 1 18

- ③ du → measure disk usage.

Ex:- du ↵

20.

- ④ file → Identify file types.

Ex:- file input.txt ↵

- ⑤ touch → change file stamps.

Ex:- touch input.txt ↵

- ⑥ chown → change file owner.

Ex:- chown 777 input.txt ↵

- ⑦ chgrp → change file group.

Ex:- chgrp 43 input.txt ↵

11 ↵
rw-r--r-- 1 777 utmp 18

- ⑧ chmod → change file protections.

Ex:- chmod 777 input.txt ↵

11
rwx rwx rwx - - -

- ⑨ chattr → change advanced file attributes.

- ⑩ lsattr → list advanced file attributes.

By Mr. GOPAL KRISHNA

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Adyar Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 5719, 986 570 6789

File Location

① find → locate files.

Ex:- find <

· |sample1.txt
· |input1.txt

② S locate → locate files via index.

③ which → locate commands

④ whereis → locate standard files.

Ex:- whereis <srcpath>

File Text manipulation

① grep → Search text for matching lines.

grep hadoop sample.txt

② cut → Extract columns

③ paste → append columns

④ tr → translate characters

⑤ sort → sort lines

⑥ uniq → locate identical lines

⑦ tee → copy stdin to a file and to stdout simultaneously.

file compression

① gzip → compress files (gnu zip)

② compress → compress files (unix)

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 018.
Ph: 040-6462 6789, 985 570 6789

③ bzip2 → compress file (Bzip2)

27

④ zip → compress files (windows zip)

File Comparison

① diff → compare files line by line.

Ex: diff input.txt sample.txt

② comm → compare sorted files.

Ex: comm input.txt sample.txt

③ cmp → compare files byte by byte.

cmp input.txt sample.txt

④ md5sum → Compute check Sums.

disks and file systems

df → show free disk space

mount → make a disk accessible

fsck → check a disk for errors

sync → flush disk caches.

printing

lpr → print files

lpq → view print queue

lprm → remove print jobs.

Spelling Operations

look → lookup Spelling

aspell → check spelling interactively

spell → check spelling

in Batch.

Kelly Technologies
Flat No. 212, 2nd Floor,
Amaravati Block, Kotha Enclave,
Amberpet, Hyderabad-500 078.
Ph: 040-252 6789, 989-570 6789

Backups and Remote Storage

mt → control a tape device

dump → Backup a disk

restore → Restore a dump

tar → Read/write tape archives.

cd record → Burn a CD

rsync → mirror a set of files.

Audio & Video

grip → play CDs & ripmp35

xmms → play Audio files.

Processes

- ① ps → List all processes
- ② w → List user's processes
- ③ uptime → View the system load
- ④ top → Monitor processes
- ⑤ xload → Monitor system load
- ⑥ free → Display free memory
- ⑦ kill → Terminate process
- ⑧ nice → Set process priority
- ⑨ renice → Change process priorities

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Annapurna Block, Aditya Enclave,
 Ameerpet, Hyderabad - 500 016.
 Ph: 040-6462 370 6789

~~Jobs~~ Scheduling Jobs

- ① sleep → Wait for sometime
- ② watch → Run programs at set intervals
- ③ at → Schedule a job.
- ④ crontab → Schedule repeated jobs.

Hosts

- ① uname → Print system information
- ② hostname → Print the system's hostname
- ③ ifconfig → Set/display n/w information
- ④ host → Lookup DNS
- ⑤ whois → Look up domain registrants
- ⑥ ping → Check if host is reachable
- ⑦ traceroute → View n/w path to host

Networking

28

- ① ssh → securely log into remote hosts
- ② telnet → log into remote hosts.
- ③ scp → securely copy files between hosts.
- ④ sftp → " "
- ⑤ ftp → copy files b/w hosts.
- ⑥ evolution → GUI email client
- ⑦ mutt → Textbased email client
- ⑧ mail → minimal email client
- ⑨ mozilla → web browser
- ⑩ lynx → Text only web browser
- ⑪ wget → Retrieve webpages to disk.
- ⑫ srln → Read usenet news.
- ⑬ gaim → Instant messaging / IRC
- ⑭ talk → Linux / unix chat.
- ⑮ write → send msgs to terminal
- ⑯ mesg → prohibit talk / write.

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Annapurna Block, Attitya Enclave,
 Ameerpet, Hyderabad-500 016.
 Ph: 040-6452 6789, 988 570 6789

By Mr. GOPAL KRISHNA

HDFS Blocks :-

- Blocks are traditionally either 64MB or 128 MB.
- Default is 64MB.
- The motivation is to minimize the cost of seeks as compared to transfer rate.
 - Time to transfer' > Time to seek.
- For example, let's say
 - seek time = 10ms
 - transfer rate = 100 MB/s
- To achieve seek time of 1% transfer rate.
 - Blocksize will need to be = 100MB

→ what is the difference between FS Shell & DFS Shell

FS Shell

distributed file system

DFS shell

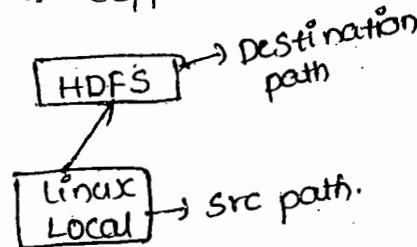
- ① FS relates to a generic file system which can point to any file systems like local, HDFS - etc.
- ② FS shell is invoked by bin/hadoop fs. All the FS shell commands take path URI as arguments.
- ③ The URI format is Scheme://authority/path.
The scheme & authority are optional.
~~If authority is not specified, the default scheme specified in the configuration is used.~~
- ④ For the HDFS the scheme is hdfs, and for local file system the scheme is file.
- ⑤ An HDFS file or directory such as /parent/child can be specified as hdfs://namenode:port/parent/child or simply as (/parent/child)
- ⑥ most of the commands in FS shell behave like corresponding UNIX commands.
- ① DFS is very specific to HDFS.
- ② The HDFS shell is invoked by bin/hadoop dfs. All the HDFS shell commands take path URI's as arguments.
- ③ The URI format is scheme://authority/path.
The scheme & authority are optional.
if not specified, the default scheme specified in the configuration is used.
- ④ For the HDFS the scheme is hdfs and for local file system the scheme is file.
- ⑤ An HDFS file or directory such as /parent/child can be specified as hdfs://namenode:port/parent/child or simply as (/parent/child)
- ⑥ most of the commands in HDFS shell behave like corresponding UNIX commands.

Shell

Hadoop file System (hadoop fs) :-

29

- hadoop fs indicating the compiler to interact with linux local environment to HDFS environment.
- hadoop fs is not support the -vi command. Because HDFS is writeOnce.
- hadoop fs is support for -touchz command.



Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6452 6789, 998 570 6789.

- hadoop fs looks the only HDFS directory but not local directory.
- we can not create a file on top of HDFS.
- we can create the file on local.
- we can not update a file on top of HDFS. we can update in local, after that file is part into the HDFS.

By Mr. GOPAL KRISHNA

Ques

- Hadoop fs & ssh does not support hardlinks (or) softlinks.

Ques

- hadoop fs does not implement user quotas.

Ques

- Error information is sent to stderr & output is sent to stdout.

Display detailed help for a command

hadoop fs -help <command-name>

HDFS Shell Commands

:S shell :-

① cat → content display (or) View content

the

② chgrp

②4) test

②7) count

③ chmod

②5) text

④ chown

②6) touchz..

⑤ copyFromLocal

User commands

⑥ copyToLocal

① archive

⑦ cp

② distcp

⑧ du

③ fs

⑨ duu

④ fsckr

⑩ expunge

⑤ fetchdat

⑪ get

⑥ jar

⑫ getmerge

⑦ job

⑬ ls

⑧ pipes

⑭ mr

⑨ queue

⑮ mkdir

⑩ version

⑯ moveFromLocal

⑪ CLASS NAME

⑰ mv

⑫ classpath

⑲ put

Administration Commands

⑳ rm

① balancer

㉑ rmdir

② daemon log

㉒ setrep

③ datanode

㉓ stat

④ dfsadmin

tail

⑤ Jobtracker

⑥ namenode

⑦ secondary namenode

⑧ taskTracker

⑨ safemode

Kelly Technologies
Flat No. 212
Annapurna Society
Ameerpet, Hyderabad
Ph: 040-64222449, 040-24227889

Shell commands:

① **mkdir** → make the directory.

30

Ex:- hadoop fs -mkdir /anitha ↴
It is just created on
hadoop fs -mkdir anitha ↴
It is created on default hdfs
path [/user/root]

② **help** → display help for all commands

Ex:- hadoop fs -help ↴

③ **ls** → list the file display stats, for a directory displays
immediate children.

Ex:- hadoop fs -ls ↴ if path is not specified.

hadoop fs -ls <directory name>
< path >

By Mr. GOPAL KRISHNA

④ **lsr** → Recursively list the contents that match the
specified file pattern.

Ex:- hadoop fs -lsr <path>

⑤ **du** → show the amount of space, if we are using
the files in HDFS

Ex:- hadoop fs -du <path> <file>

Note:- Add -h option to display in human-readable
format instead of bytes.

Ex:- hadoop fs -du -h /somedir ↴

206.3K /somedir

=

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-2452 6719, 295 570 6789

⑥ `du` → show the amount of space, in bytes, used by the files that match the specified file pattern.

Ex: `hadoop fs -du` `st`
`hadoop fs -du` `<path>` file (or) directory.

⑦ `df` → show the capacity, free and used space of the file system.

Note: if the file system has multiple partitions, and no path to a particular partition is specified then the status of the partitions will be shown.

`hadoop fs -df <path>`

⑧ `(cp -r)` copy files that match pattern `<src>` to `<destination>`

Note: when copying multiple files, the `<dest>` must be a directory.

Ex: `hadoop fs -cp anitha/input.txt` /anitha/
`<src path>` `<dest path>`
must be
HDFS paths.

⑨ `mv` → move files that match the specified file pattern `<src>` to `<dest>`

Note: when moving multiple files, the destination must be a directory.

`hadoop fs -mv <srcpath>` `<dest path>`

must be
HDFS paths

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad - 500 018,
Ph: 040-6462 6789, 989 570 6789

(10) `rm` → Delete the files and particular non-empty directories.

Ex: `hadoop fs -rm <path>`

(11) `rmr` → remove, recursively, ^{directories}
Ex: `hadoop fs -rmr <path>` ↴ directory.

31
Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6709, 988 570 6789

(12) Count → Count the no. of directories, files and bytes under the paths that match the specified file pattern.

Note: The output columns are Dir-Count, File-Count, Content-Size, file-name.

Ex: `hadoop fs -count <path>` ↴ directory.

Note: The output columns with -q are QUOTA, QUOTA, REMAINING-QUOTA, SPACE-QUOTA, REMAINING-SPACE-QUOTA, DIR-COUNT, FILE-COUNT, CONTENT-SIZE, FILE-NAME.

Ex: `hadoop fs -count -q <path>` ↴ directory.

(13) `put` → copy single src (or) multiple src's from local file system to the destination file system. Also reads input from stdin and writes to destination file system.

Ex: `hadoop fs -put <src path> <dest path>`
 ↑
 copy from local local file ↓
 HDFS dir

multiple file from Local to HDFS

By Mr. GOPAL KRISHNA

Ex: `hadoop fs -put *.* <dest path>`
 ↓
 HDFS path.

only txt file from Local to HDFS

Ex: `hadoop fs -put *.txt <dest path>`
 ↓
 local path.
 ↓
 HDFS path

14) get → copy files to Local file System.

syn.: hadoop fs -get ^{or} <src> <dest>
 copyToLocal HDFS path local path

multiple files from HDFS to Local

sign: hadoop fs -get /anitha/*.* $\frac{\text{<src-path>}}{\downarrow \text{local path}}$ (dest path)

⑯ expunge → empty the trash

sgn: hadoop fs -expunge ↵

⑯ move from Local → move
copy single src (or) multiple src's from
local file system to destination file system

59. hadoop fs -movefromlocal <local src> <dest>

(17) `getmerge` → takes a source directory and destination file as input and concatenates file in src into the destination local file.

syn: hadoop fs -getmerge <src path> <local destination path>
↓
HDFS path → local path.

Note :- addnl can be set to enable adding a newline character at the end of each file.

character at the
e.g.: `hadoop fs -getmerge <srcpath> ~<localpath> <addn>`
Compulsory
Hadoop
directory

(18) `text` → takes a source file and outputs the file in text format.
The allowed formats are zip & Text Read Input Stream.

Syn: `hadoop fs -text <src path>`

37

(19) `touchz` → creates a file of zero length (or) size.

Syn: `hadoop fs -touchz <path>`

(20) `test` →
-e check to see if the file exists.
returns 0 if true.
-z check to see if the file is zero length
returns 0 if true.
-d check return 1 if path is directory
else return 0.

Syn:
`hadoop fs -test -ezd <path name>`
`hadoop fs -test -e <path>`
`hadoop fs -test -z <path>`
`hadoop fs -test -d <path>`

Kelly Technologies
Plot No. 212, Zeta Park,
Ameerpet, Hyderabad-500 014
Ph: 040 462 6789 988 510 014

(21) `stat` → Returns the stat information on the path.

Syn: `hadoop fs -stat <path name>` local & HDFS.

(22) `Tail` → Displays last kilobyte of the file to std out
-f option can be used as in unix.

Syn: `hadoop fs -tail <path name>`
↓ HDFS path.

By Mr. GOPAL KRISHNA

(23) `setrep` → changes the replication factor of a file.
-R option is for recursively increasing the replication factor of files within a directory.

Syn: `hadoop fs -setrep -w 3 -R <HDFS path name>`

-w 3 -R

↓

/user/root/dir

↓
<HDFS path name>

(24) `chgrp` → change group association of files. with -R make the change recursively through the directory structure.

The user must be the owner of files (or) else a super-user.

Syn:- `hadoop fs -chgrp -R [Group URI]`

`hadoop fs -chgrp -R [group] [directory name]` (or) `hadoop fs -chgrp [filename]`

(25) `chmod` → change the permissions of files. with -R make the change recursively through the directory structure. The user must be the owner of files or else a Super-user.

Syn:- `hadoop fs -chmod 744 -R [filename]`

permissions

- u - user who owns the file
- g - group that owns the file
- o - others & - all
- r - read the file, w - write or edit the file
- x - execute or run the file as a program

Numeric permissions

chmod can also be attributed by using numeric permissions

Ex:- `file -rwx-r--0` → owner group every one only owner

Ex:- `hadoop fs -chmod 777 /tmp` → path

400	→ read by owner
040	→ read by group
004	→ read by anybody (other)
200	→ write by owner
020	→ write by group
002	→ write by anybody
100	→ execute by owner
010	→ execute by group
001	→ execute by anybody

(26) `chown` → change the owner of files. with -R make the change recursively through the directory structure. The user must be the owner of the file (or) else a Super-user.

Ex:- `hadoop fs -chown chope file.txt` → file name

Give permission → owner to user chope for the file file.txt.

`chown -R chope work` → (path dir name)

Given chown permissions to chope for all files in the work directory.

User Commands

33

① archive:-

Hadoop stores the small files inefficiently such as each file get stored in a block & Namenode has to keep the metadata information in memory. So with this reason most of the namenode memory will get eat up by this small files only which results in a wastage of memory.

To avoid the same problem we use hadoop archives (or) har files (.har is the extension for all the archive files).

when creating archive directory the input is converted to mapreduce jobs, so we can call hadoop archives as a input for our mapreduce programming.

By Mr. GOPAL KRISHNA

* Hadoop archives are special format archives.

* Hadoop archive maps to a file system directory.

* Hadoop archive always has a *.har extension.

* Hadoop archive directory contains metadata [in the form of -index and -master-index] and -data.

contain the files
that are part of archive
and location with in
the part files.

Kelly Tech Solutions
Plot No. 212, 2nd Floor,
Ananya Block, Ananya Enclave,
Hyderabad-500 015.
Ph. 040 452 5700, 900 570 5700

Syn:- create the archive file:-

hadoop archive -archiveName name -P <src> <dest>

E:- hadoop archive -archiveName foo.har -P \$Anitha /myarch

hadoop fs -ls /myarchive ↴

hadoop fs -lsr /myarchive/foo.har ↴

hadoop fs -cat /myarchive/foo.har/part-0 ↴

② distcp: Distributed copy.

- The distcp command is tool used for large inter and intra cluster copying.
- we have multiple hadoop clusters are running. we want to transfer several terabytes of data one cluster to another.
Ex: hadoop clusters are loaded with terabytes of data.
It will take forever to transfer terabytes of data from one cluster to another.
- distributed (or) parallel copying of data can be good solution for this & that is what distcp does.
- distcp runs mapreduce job to transfer your data from one cluster to another.

syn: `hadoop distcp <srcpath> <dest.path>`

`hadoop distcp hdfs://nn1:8020/foo/bar \ hdfs://nn2:8020/bar/foo`

Katty Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6729, 988 570 0789

multiple src list :-

`hadoop distcp hdfs://nn1:8020/foo/bar \ hdfs://nn1:8020/foo/b \ hdfs://nn2:8020/bar/foo`

From file using :-

`hadoop distcp -f hdfs://nn1:8020/srclist \ hdfs://nn2:8020/bar/foo`

③ `fs → Runs a generic file system user client`

syn: `hadoop fs [Generic-options] [command-options]`

command options can be found at hadoop shell guide.

④ jar → Runs a jar file. users can bundle their mapreduce code in a jar file and execute it using this command.

The Streaming jobs are run this command.

↓
Hadoop streaming is a utility that comes with the hadoop distribution.

syn:- hadoop jar <jar name>
Ex:- hadoop jar <wordcount.jar>

⑤ job → This command to interact with mapreduce jobs

hadoop job [generic options]

⑥ pipes →

By Mr. GOPAL KRISHNA

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

c

MapReduce

By Mr. GOPAL KRISHNA

34A

Introduction:-

- mapReduce published in 2004 by google.
- Hadoop can run mapreduce programs written in various languages like java, Ruby, python and c++.
- mapReduce is a parallel programming model for processing the huge amount of data.
- mapReduce making the Structured data and out of some unstructured data ..etc.
- mapReduce provides automatic parallelization & distribution, Fault-tolerance, I/O scheduling, monitoring and status.

MapReduce Overview:-

- Applications processing data on Hadoop using the mapReduce paradigm.
- A mapReduce job usually splits the input data-set into independent chunks, which are processed by the map tasks in a completely parallel manner. The framework sorts the outputs of the maps, which are then input to the reduce tasks. Typically both the input and the output of the job are stored in a file system. The framework takes care of scheduling tasks, monitoring them and re-executes the failed tasks.
- mapReduce applications specify the input/output locations and supply map and reduce functions via implementation of appropriate Hadoop interfaces such as Mapper and Reducer. These, and other parameters comprise the job configuration. The Hadoop job client then submits the job (job executable, etc) and configuration to the JobTracker, which then assumes the responsibility of distributing the software / configuration to the

Kelly Technologies
Flat No. 212, 2nd Floor, Aditya Enclave,
Amaravati Block, Hyderabad-500 016.
Ph: 040-34526777, 9885706709

Written

slaves, scheduling tasks and monitoring them, providing status and diagnostic information to the job-client.

- The map/reduce framework operates exclusively on <key,value> pairs - that is, the framework views the input to the job as a set of <key,value> pairs and produces a set of <key,value> pairs as the output of the job, conceivably of different types.

what is mapReduce:

functional style
programming

- map/reduce
 - programming model from LISP and other functional language.
- many problems can be phrased this way.
- easy to distribute across nodes.
- nice retry/failure semantics.
- sort/merge based distributed computing.
- The underlying system takes care of the partitioning of the input data, scheduling the program's execution across several machines, handling machine failures, and managing required inter-machine communication.
- computational processing occurs on both
 - unstructured data : file system
 - structured data : database
- Tried and tested in production
- many implementation options.
- fault-tolerant, reliable, and supports thousands of nodes and petabytes of data.

Kelly T
Flat No.
Annapurna
Ameerpet
Ph: 040-23110000
Mobile:
9840000166
35706789

Motivation for MapReduce (why) :-

35

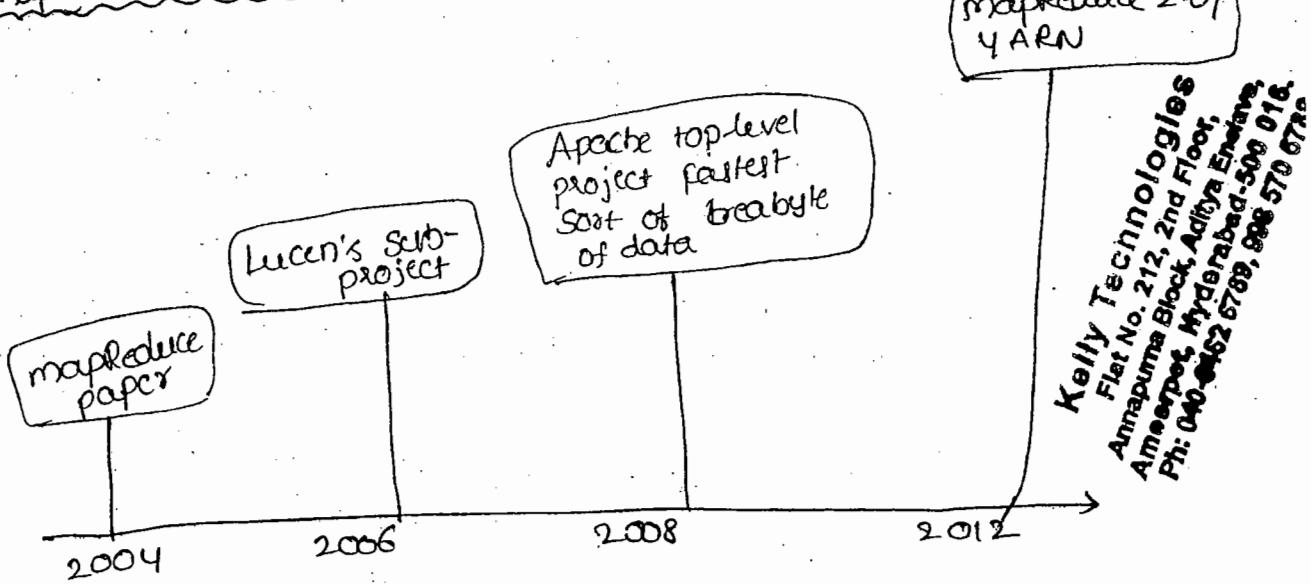
By Mr. GOPAL KRISHNA

- Large scale data processing.
- Want to use 1000's of CPUs.
- But don't want hassle of managing things.
- MapReduce architecture provides
 - automatic parallelization & distribution
 - fault-tolerance
 - I/O scheduling
 - monitoring & status updates.

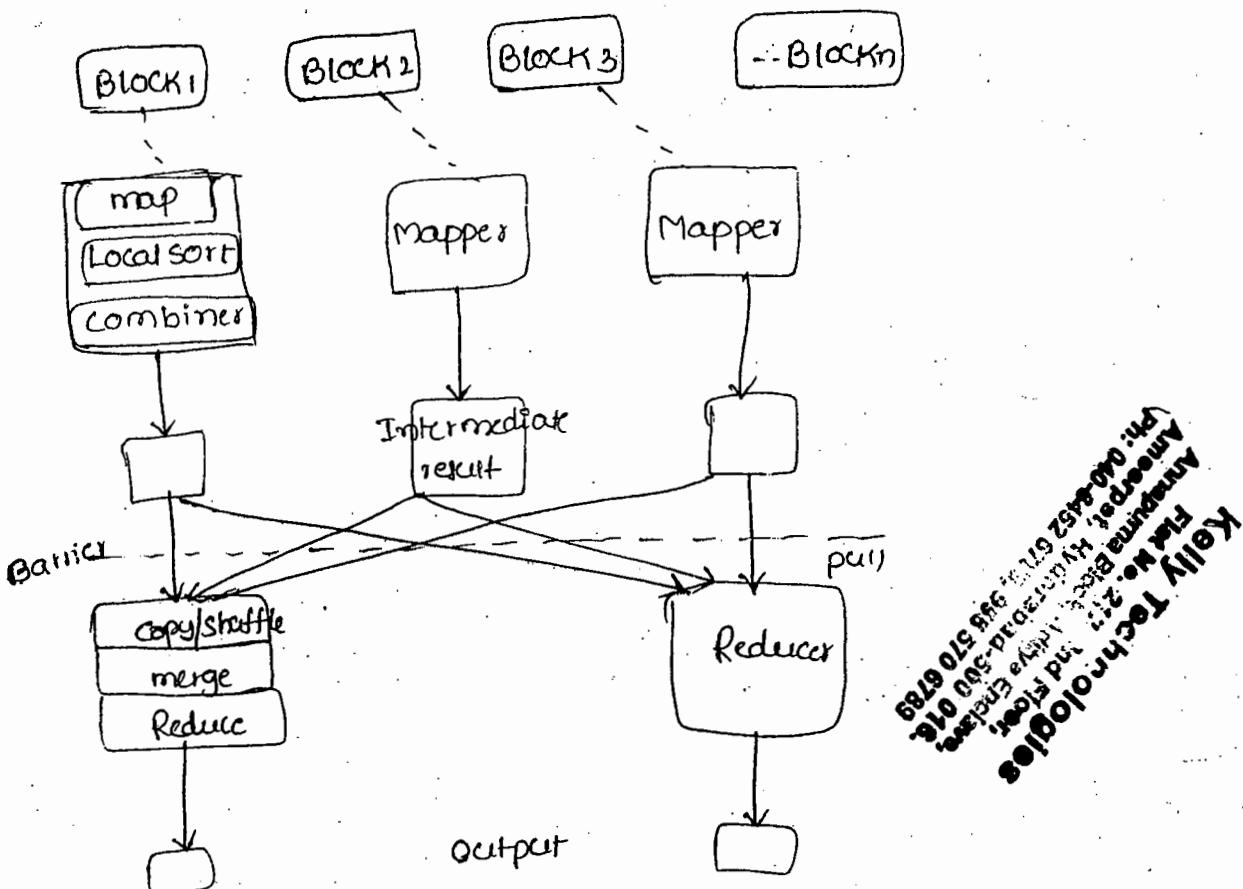
MapReduce Limitations :-

- Can not control the order in which the maps or reductions are run.
- For maximum parallelism, you need maps and Reduces to not depend on data generated in the same MapReduce job (i.e. stateless).
- A database with an index will always be faster than a MapReduce job on unindexed data.
- Reduce operations do not take place until all maps are complete (or have failed then been skipped).
- General assumption that the output of Reduce is smaller than the input to map; large data source used to generate smaller final values.

MapReduce History :-



mapReduce Architecture:



- mapReduce Architecture is a master-slave architecture.
 - mapReduce is a programming model only. we can expect parallel processing in Hadoop which means that if an input file coming into Hadoop file as part of processing of the same file.
 - Input will get divided into multiple chunks and each & every chunk will be processed in different nodes.
 - The whole process of mapReduce will be controlled by
 - ① Job Tracker
 - ② Task Tracker
 - Job Tracker will assign the tasks in the form of the available Job Trackers.
 - mapReduce jobs to all the
 - Task trackers has to act up on the assigned tasks, and they have send the progress information to the Job tracker on Timely basis.

- In case of failure (w.r.t tasks) \rightarrow Job Trackers will reassign the task to some other still idle available Tasktrackers.

36

\rightarrow It consists of two phases: map and Reduce.

By Mr. GOPAL KRISHNA

MapReduce : the Mapper :-

- when the map function starts producing output, it is not simply written to disk. The process is more involved, and takes advantage of buffering writes in memory and doing some presorting for efficiency reasons.
- Each map task has a circular memory buffer that it writes the output to. The buffer is 100 MB by default, it can change based on the size. When the contents of the buffer reaches a certain threshold size 80% a background thread will start to spill the contents to disk during this time, the map will block until the spill is complete.

Before it writes to disk, data is partitions corresponding to the reducers. Each partition, sort by key, and if there is a combiner function, it is run on the output of the sort. So there is less data to write to local disk and transfer to the Reducer.

The mapper reads data in the form of key/value pairs and it outputs zero or more key/value pairs.

mapReduce flow \rightarrow The Mapper

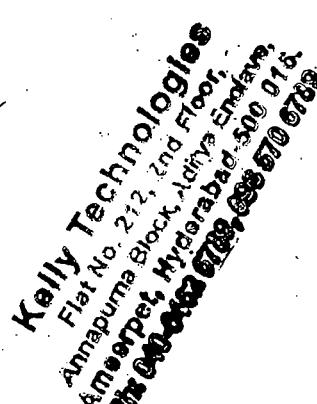
Each of the portions (Record Reader

mapper

partitioner

Reducer

... etc)



mapReduce : the Reducer

- The mapoutput file is sitting on the local disk of machine. The reduce task needs the map output for its particular partition from several tasks across the cluster.
- The map tasks may finish at different times, so the reduce task starts copying their outputs as soon as each completes. This is known as the copy phase of the reduce task, by default 5 threads can copy, this we can change by property.
- when all the map outputs have been copied, the reduce task moves into the sort phase. For example, if there were 50 map outputs, and the merge factor was 10, then there would be 5 rounds. Each round would merge 10 files into one, so at the end there would be five intermediate results files. final round that merges these 5 files into a single sorted file and move to the reduce phase. The output of this phase is written directly to the output file System, typically HDFS.

In the above diagram, After the map phase is over, all the intermediate values for a given intermediate key are combined together into a list. This list is given to a Reducer:

Kelly Technologies
Flat No. 212, 2nd floor,
Endave,
Annapurna Block
Amarpet, Hyderabad
Ph: 040-6452 6789

- There may be a single Reducer, multiple Reducers.
- This is specified as part of the job configuration.
- All values associated with a particular intermediate Key are guaranteed to go the same Reducer.
- The intermediate key, and their value lists, are passed to the Reducer in sorted key order.
- This step is known as the 'Shuffle & Sort'.

- The Reducer outputs zero or more final key/value pairs.

These are written to HDFS

- In practice, the Reducer usually emits a single key/value pair for each input key.

The mapreduce flow: Reducers → outputs

Each of the portions (Reducer, RecordWriter, output file).

MapReduce Different programming model :-

- ① Different phases of mapReduce algorithm.
- ② Different data types in mapReduce.

Different phases of mapReduce algorithm :-

mapReduce works by breaking the process into 3 phases.

① mapper phase

② sort & shuffle phase [logical layer]

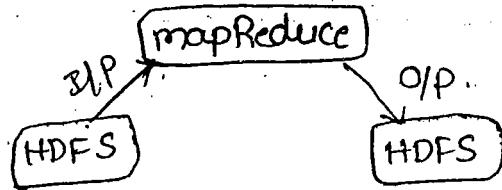
③ Reduce phase.

- In mapReduce each phase as input and output functions.

↓
Mapper functions

Key-value pairs as
↓
Reducer functions.

Kelly Tech, No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ph: 040-6462 6719, 988 570 6789



- mapReduce will expect the input in the form of (key, value) pairs from HDFS layer only. Once it is done with the processing it will produce the o/p again on top of HDFS in the form of (key, value) pair.

phase	input	output
Mapper	(K, v)	(K, v)
Sort & shuffle	(K, v)	(K, list(v))
Reducer	(K, list(v))	(K, v)

Example for mapReduce :-

Input.txt

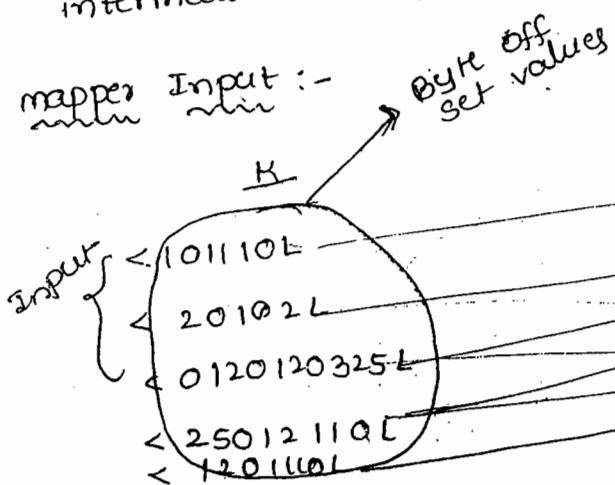
Hadoop is a Bigdata analysis storage & processing for the data. Hadoop is having good market now. Good demand will give good openings.

Explanation :-

1) Map phase :- map phase will expect input & output forms.

- map phase will work on input of key, value pairs. It is inherently parallel. Each list element processed independently.
- completely written by the system in the mapper phase.

- In the map phase key value is in the form of byte off set values.
- A list of data elements are provided to mapper function called the map() which transforms input data to an intermediate output data element.

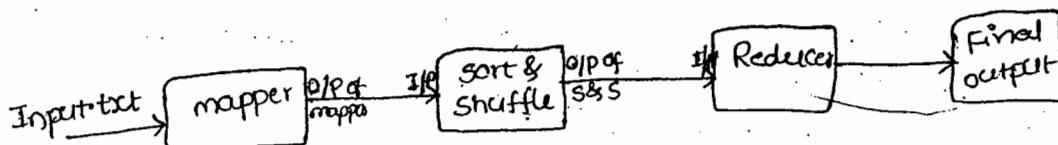


By Mr. GOPAL KRISHNA

✓
hadoop is Bigdata Analysis

Storage & processing for the data. Hadoop is having good market now. Good demand will give good openings.

Kelly Tech Solutions
Flat No. 212, 2nd Floor,
Annapurna Block, Hyd-521008
Ph: 040-66626739, 9845057500



Mapper Output :-

- Mapper output will be displayed in the form of (k,v) pairs.
- In this mapper output value is always 1.

K	V
hadoop	1
is	1
a	1
bigdata	1
:	
opening	1

	Input	Output
mapper	k,v	k,v
sort & shuffle	k,v	k, list(v)
Reducer	k, list(v)	(k,v)

Sort & Shuffle :-

- In this phase take the sort & shuffle phase input.

Sort :- Sort is used to list the inputs in sorted order.

Shuffle :- mapReduce makes the guarantee that the input to every reducer is sorted by key. It is known as the shuffle.

sort & shuffle Input

K	V
hadoop	1
is	1
a	1
Bigdata	1
:	1
openings	1

sort & shuffle output

K	list(values)
hadoop	1, 1
is	1, 1
Bigdata	1
:	1
openings	1

Reducer :- is inherently

sequential. (unless processing multiple lists)

- A Reducer function receives an iterator input values from an output list. It combines these values together returning a single output value.

Reducer Input

K	list(values)
hadoop	1, 1
is	1, 1
Bigdata	1
:	1
openings	1

Reducer output

K	V
hadoop	2
is	2
Bigdata	1
:	1
openings	1

Different data types in mapReduce

By Mr. GOPAL KRISHNA

39

normal language	mapReduce
Int	IntWritable
float	FloatWritable
double	DoubleWritable
long	LongWritable
String	Text
boolean	BooleanWritable

How to write basic mapReduce algorithm :-

- In any mapReduce program with irrespective of Business logic will divided into 3 phases.

- ① Driver code → configuration level details
- ② mapper code → Business logic
- ③ Reducer code → final output.

Driver code :-

- configuration level details with respect to job , JAR creation etc.
- mapper, Reducer class level details.
- final output key , value , data type details.
- Input and output HDFS paths.

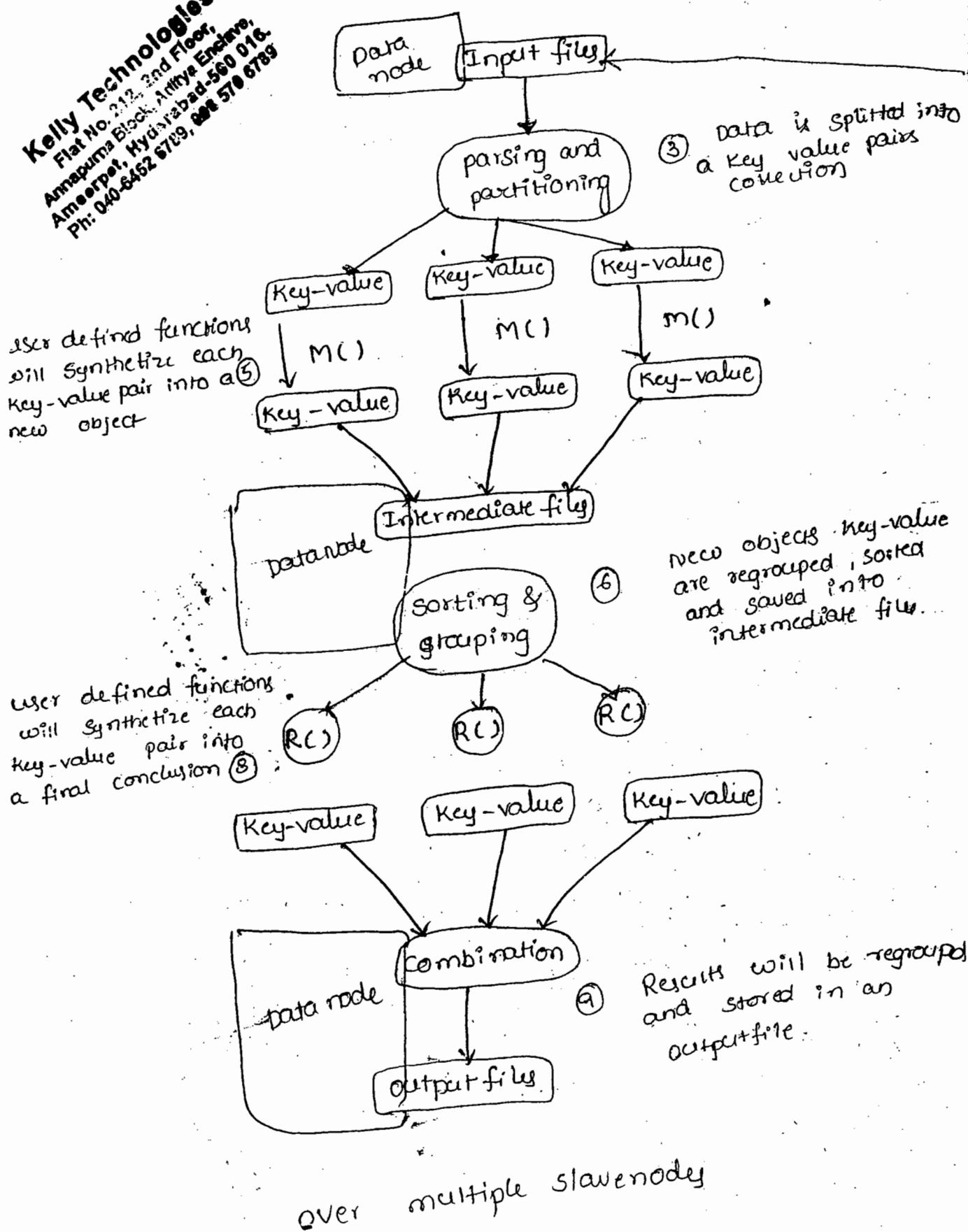
[Mapper - Reducer]
One → One

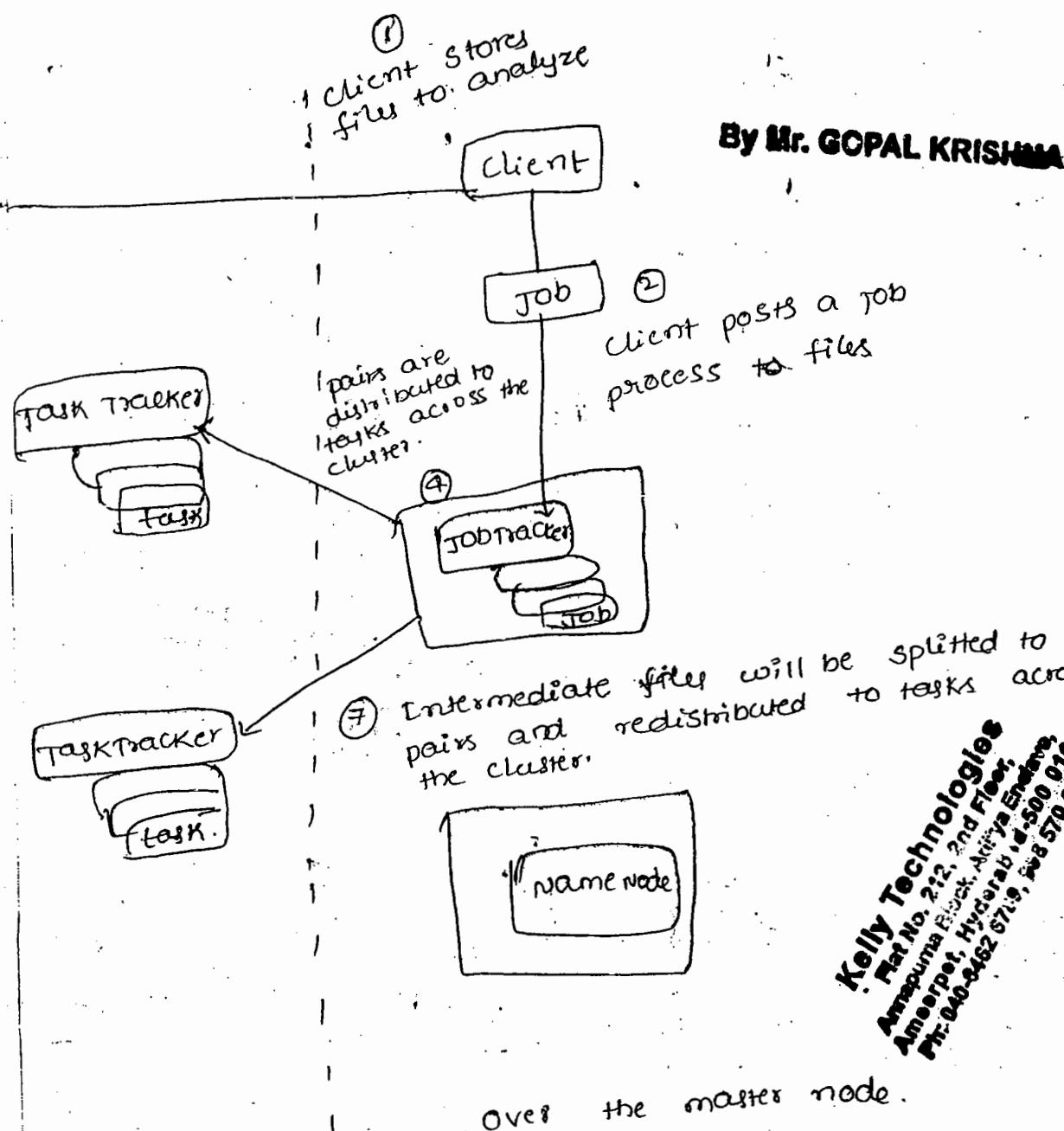
many → one
 mapper → Reducer
 mapper
 mapper

Kelly Technologies
 Flat No. 212, 2nd Floor
 Annapurna Block, Ameerpet, Hyderabad - 500 011
 Phone: 040 46626789, 9885061833

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Annapurna Block, Amity Enclave,
 Ameerpet, Hyderabad-500 016.
 Ph: 040-6462 6729, 040 570 6789

Hadoop mapReduce life cycle

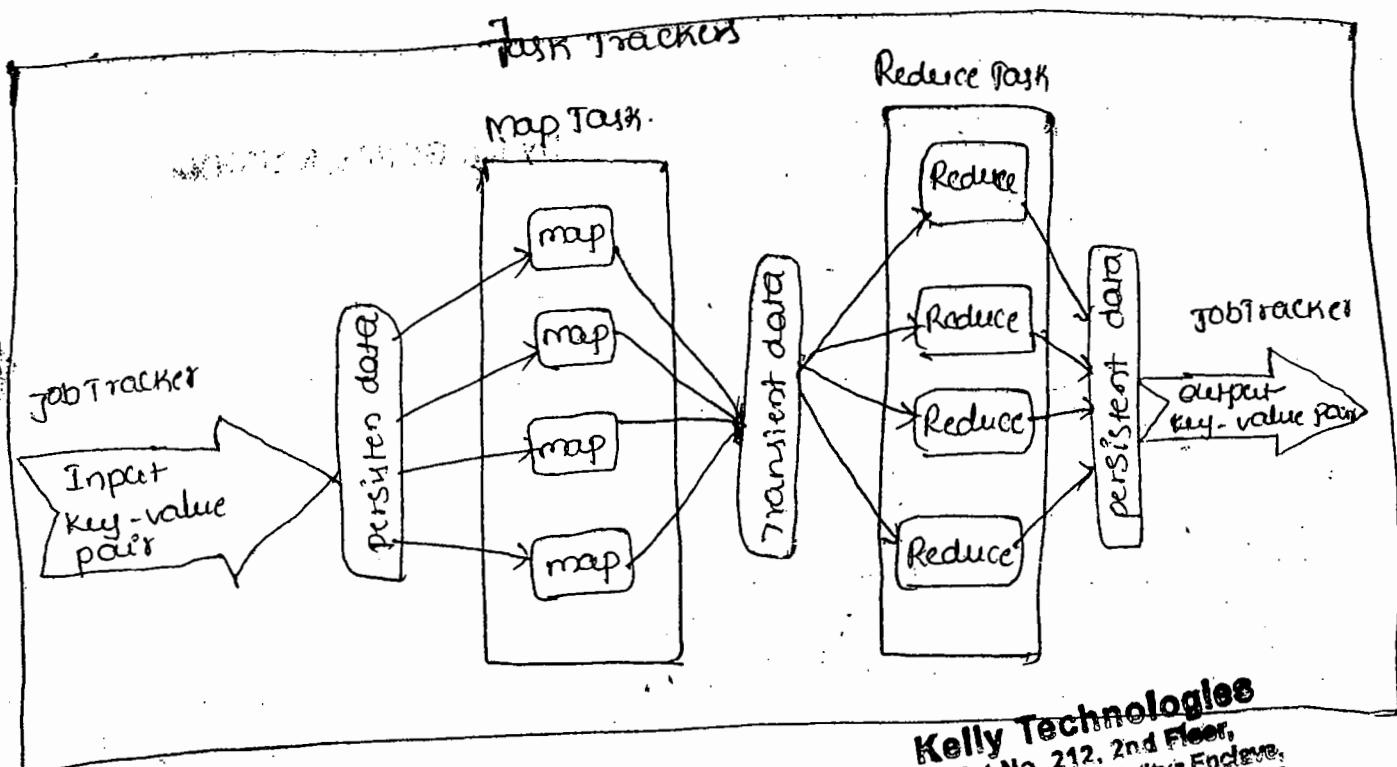




Kelly Technologies
 Plot No. 212, 2nd Floor,
 Amegama Layout, Ameerpet, Hyderabad - 500 018
 Tel: 040-2462 6719, 248 570 6789

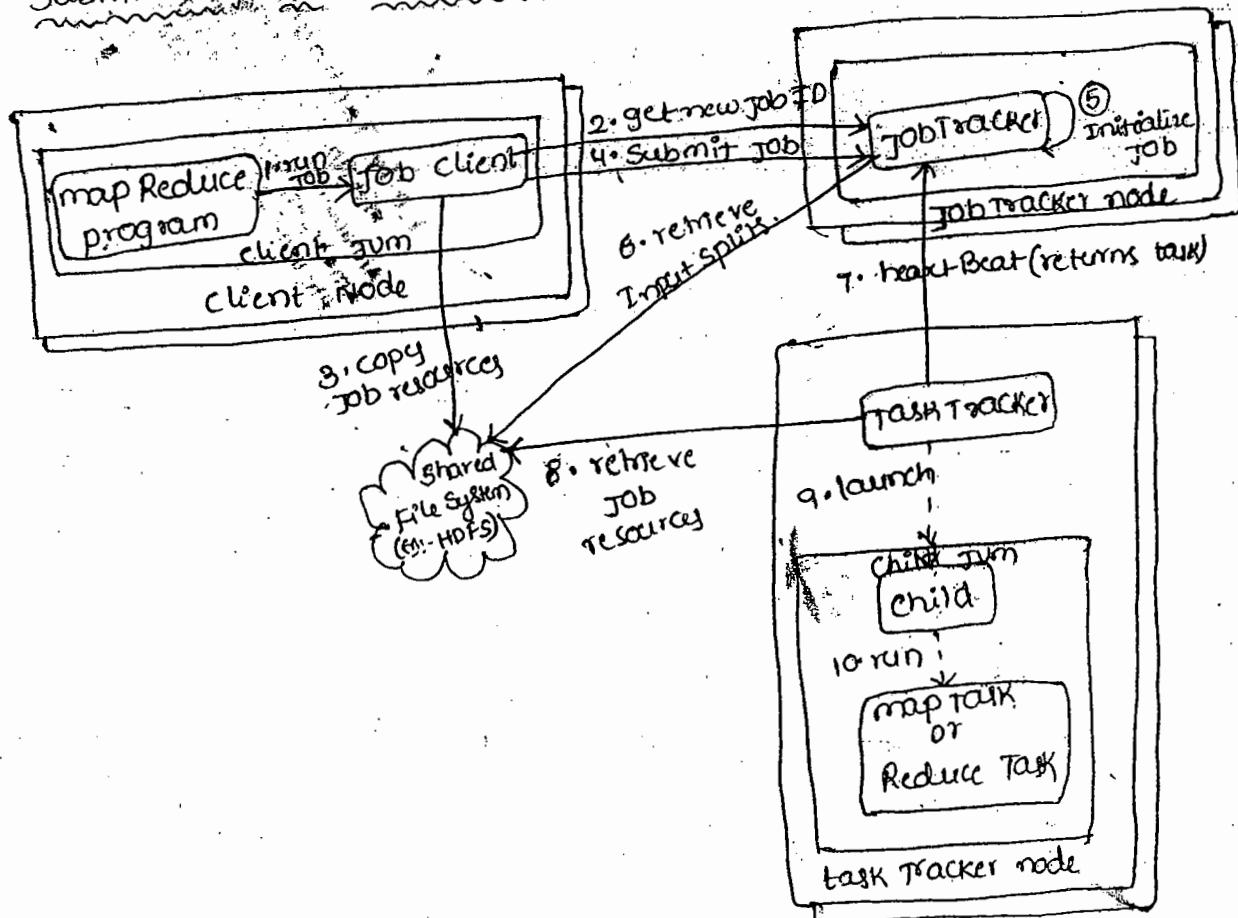
-
- Record Reader is a class. Each input format provides its own Record Reader implementation.
- LineRecord Reader :- Reads a line from a text file.
- Key Value Record Reader :- used by Key value Text Input Format.

map Reduce execution framework



Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6662 6773, 563 570 6783

Submitting a mapreduce job :-



• whenever client submits the input file to the mapreduce in the form of jobconf object (.JAR file). Jobtracker will receive the file to the available tasktrackers in the hadoop cluster.

• By default every tasktracker supports two map tasks and two reduce tasks. So it will accept the tasks by jobtracker and it will perform the maptasks first and sends the progress information back to the jobtracker by the means of heartBeat.

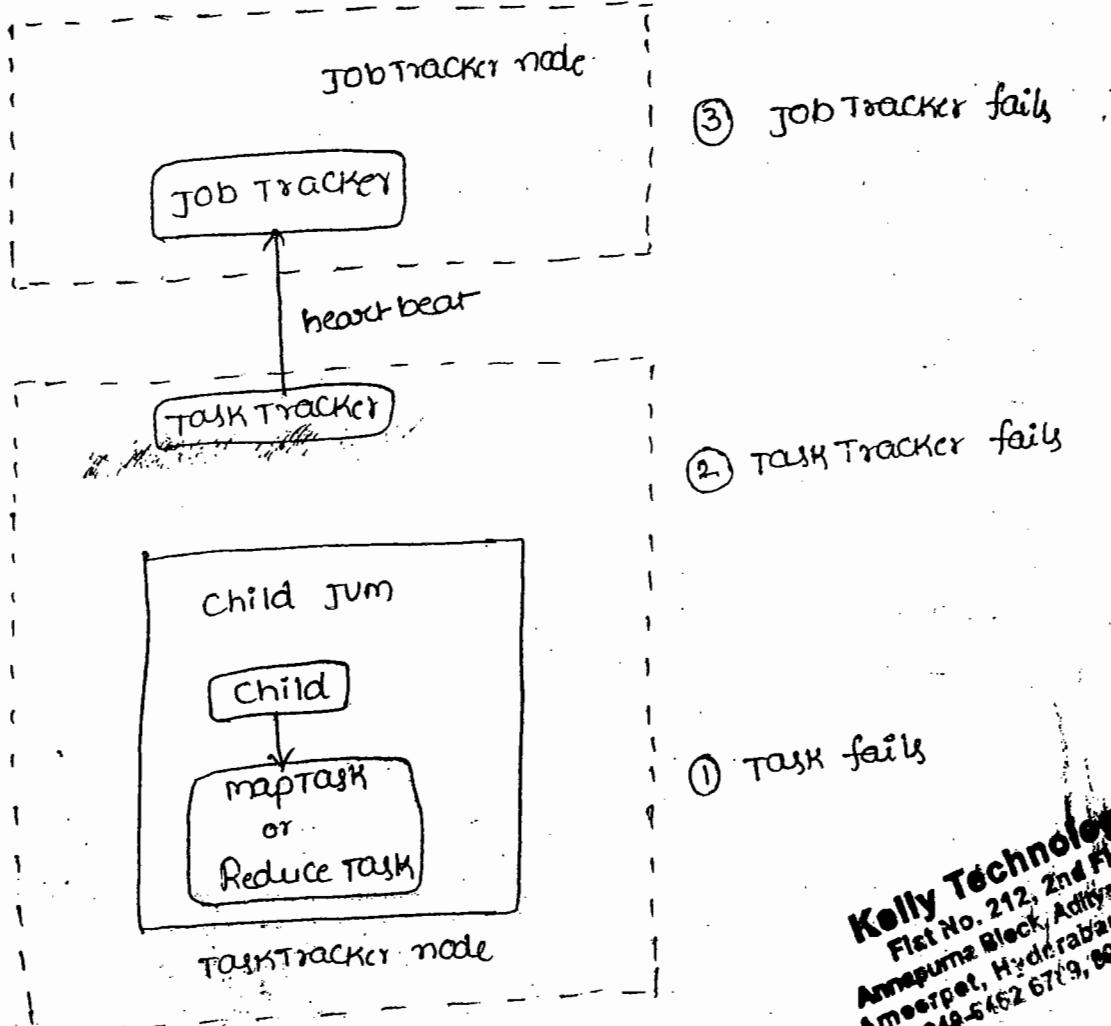
By Mr. GOPAL KRISHNA

• Based up on the progress information sent by the tasktracker to the jobtracker, jobtracker will initiate the next phase called Reducer (provided if the tasktracker has completed map phase (100%)).

• Once jobtracker receives all the Reducer outputs from all the tasktrackers to whom it has assigned the task then it is called job complete.

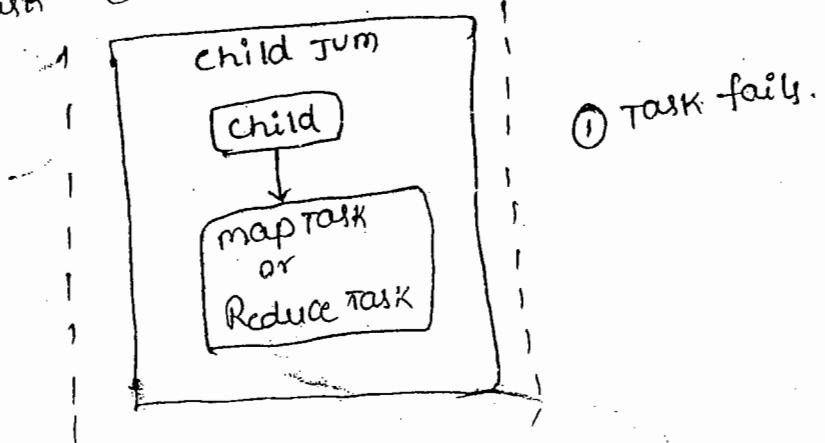


Fault - Tolerance



① Task Failure:

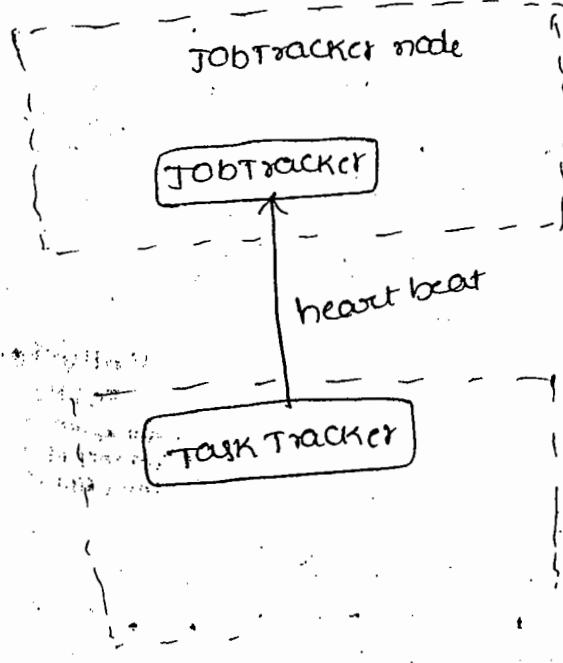
- if a child task fails, the child JUM reports to the Tasktracker before it exists. Attempt is marked failed, freeing up slot for another task.
- if the child task hangs, it is killed. Jobtracker reschedules the task on another machine.
- if task continues to fail, job is failed.



Kelly Technologies
 Flat No. 212, 2nd Floor,
 Nimmajpura Block, Aditya Enclave,
 Amravati, Hyderabad-500 016.
 Ph: 040-5462 6779, 988 578 6783

② Task Tracker failure:

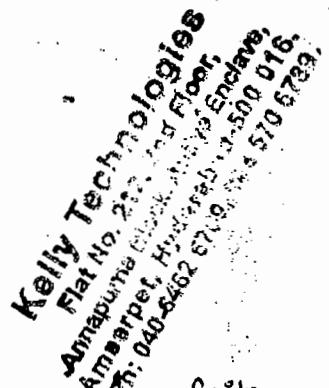
- Job Tracker receives no heartbeat.
- Removes TaskTracker from pool of Tasktrackers to schedule task on.



→ if at all jobtracker does not receive heartbeat report from anyone of the Tasktrackers it considers the particular Tasktracker is down and it assigns the same task to other idle replicated node.

② Task Tracker fails.

By Mr. GORAL KRISHNA



③ jobtracker failure

- Single point of failure. Job fails.



③ JobTracker fails.

→ if at all genuinely there is a flaw in the logic written in both mapper & reducer if at all we are getting corrupted/bad records in all these cases task will fail.

In this case jobtracker will try to execute the same task 4 times by default.

In the 4 attempts the same task is been failing then jobtracker will mark the entire job as failed (as we can say job is complete when all the tasks are completed successfully).

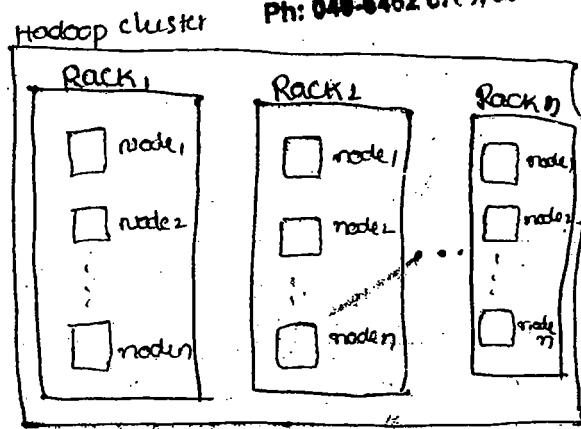
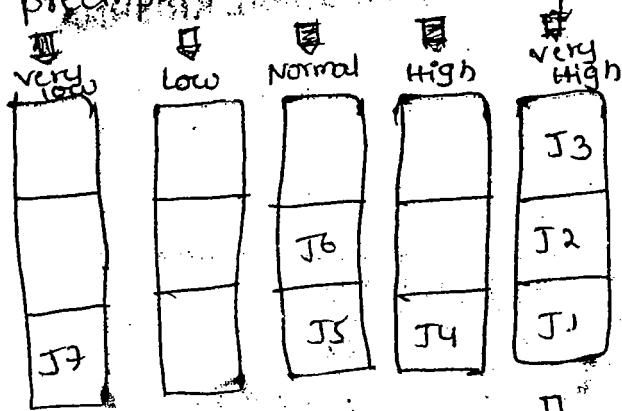
Scheduling :-

- ① FIFO scheduler (with priorities)
- ② Fair scheduler
- ③ Capacity scheduler.

FIFO scheduler (with priorities):-

- Not suitable for shared production-level cluster
- Each job uses the whole cluster, so jobs wait their turn.
- can set priorities for the jobs in the queue (5 queues with priorities)

- preemption is not supported.

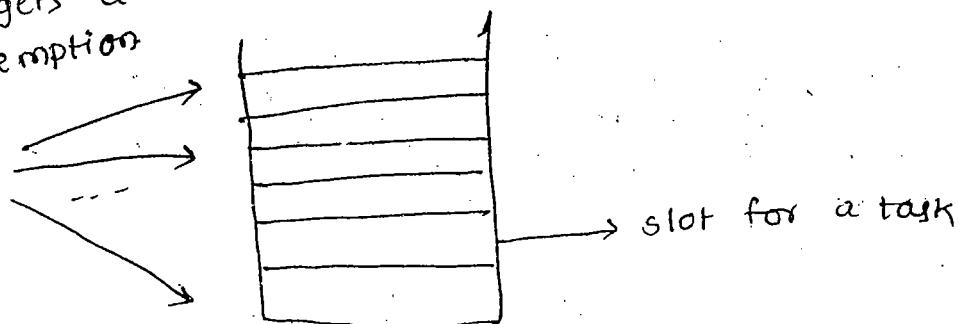


Fair scheduler :-

- jobs are assigned to pools (1 pool per user by default)
- A user job pool is a number of slots assigned for tasks for that user.

- Each pool gets the same no. of tasks slots by default.
- Every user gets a fair share of the cluster capacity overtime.
- supports preemption

many submits
many jobs



long jobs can monopolize cluster
short jobs can hold back and have no guarantees on executing

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarapuram, Hyd - 500 016.
Ph: 040-6462 6777 / 570 0789

- multiple users can run jobs on the cluster at the same time.

W3

By Mr. GOPAL KRISHNA

- Example:

- mary, John, peter submit jobs that demand 80, 80, and 120 tasks respectively.
- say the cluster has a limit to allocate 60 tasks at most.
- Default behaviour: Distribute task fairly among 3 users (each get 20).

maximum amount of tasks that can be allocated in this cluster : 60



mary
Demand: 80



John
Demand: 80



peter
Demand: 120

- minimum share can be set for a pool

- In the previous example, say mary has a minimum share of 40.

- mary would be allocated 40, then the rest is distributed evenly to other pools.



mary
Demand: 80
minimum
Share: 40



John
Demand: 80



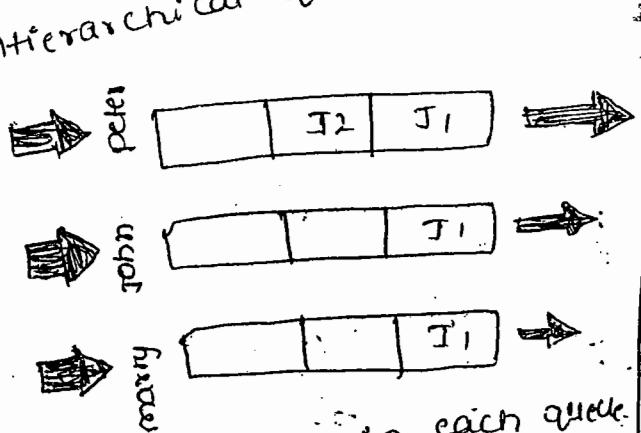
peter
Demand: 120

Kelly Technologies
Flat No. 212, 2nd Floor,
Amapuram Block, Aman Estate,
Ameerpet, Hyderabad 500 016
Ph: 040-4462 570 0789

③ Capacity Scheduler:-

- Similar goal to Fair scheduler: distribute jobs fairly among users.
- works with queues instead of pools.
- A user thinks he has the cluster to himself with FIFO scheduling, but actually is sharing resources with other users.
- Hierarchical queues (mimic an organization)

Hadoop cluster:



- FIFO scheduling in each queue.
- Supports priority.

Task execution - Improving Performance:-

Speculative execution:-

- Job execution is time sensitive
- Hadoop detects slow running tasks and launches another, equivalent task as backup.
- The output from the first of these tasks to finish is used.

Task Jars Reuse:-

- Tasks run in their own JARs for isolation.
- Starting up a JAR is relatively expensive if jobs are short.
- In these cases when jobs have many short tasks, improves performance by configuring JAR reuse.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarpet, Hyderabad 500 016.
Ph: 040 6662 6712, 98495 57899

Joining datasets in MapReduce jobs :-

44

- ① map side joins
- ② Reduce Side Joins.

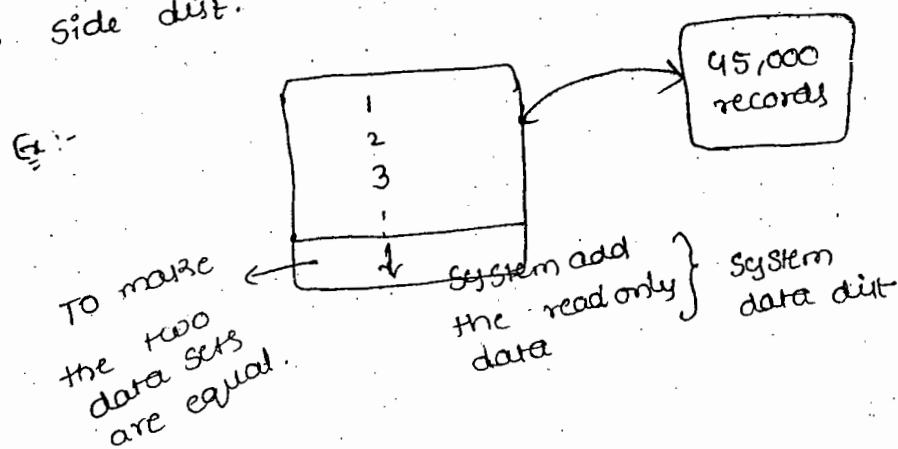
By Mr. GOPAL KRISHNA

mapside joins :-

- map-side join is done in the map phase and done in memory.
- the most common problems with map-side joins are out of memory exceptions on slave nodes.
- map-side join is faster because join operation is done in memory.
- Replicate a relatively smaller input source to the cluster. put the replicated data set into a local hash table.
- join a relatively larger input source with each local hash table.
- mapper do mapper-side join.

Disadvantage :-

- All the data should be sorted order.
- side dist.

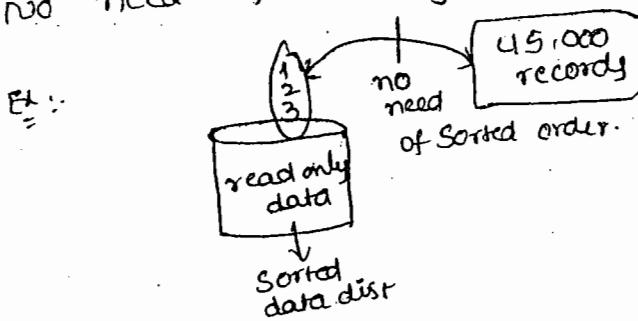


Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 014
Ph: 040 4462 6789, 988 570 6789

Reduce-side join :-

Reduce side join is a technique for merging data from different sources based on a specific key. There are no memory restrictions.

no need of sorting [overridden is not here]



Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040 6452 6789, 988 578 6789

Distributed Cache :-

- Distributed cache distributes files efficiently.
- read only
- Distributed cache is a facility provided by the map/reduce framework to cache files (text, archives, jars and so on) needed by applications.

The framework will copy the necessary files to the slave node before any tasks for the job are executed on that node.

It is efficiency stems from the fact that files are only copied once per job and the ability to cache archives which are un-archived on the slaves. It can also be used as a rudimentary s/w distribution mechanism for use in the map and/or reduce tasks. It can be used to distribute both jars and native libraries and they can be put on the classpath or native library path for the map/reduce tasks.

- The distributed cache is designed to distribute a small no. of medium size artifacts, ranging from a few mbs to few tens of mb's

- one drawback of the current implementation of the distributed cache is that there is no way to specify map or reduce specific artifacts. us

Counters :-

Counters are the useful channel for gathering statistics about the job.

- for the quality control OR application level statistics
- For problem diagnosis.

By Mr. GOPAL KRISHNA

Built-in counters :-

- Hadoop maintains some built-in counters for every job, which report various metrics for our job.

Eg:- expected amount of INPUT consumed.

expected amount of OUTPUT produced.

Some - built in counters:

- ① map input Records :- number of input records consumed by all the maps in the job. Incremented every time a record is read from InputSplit (through RecordReader) before passing to map() method of mapper.

- ② map output Records :- number of output records produced by all the maps in the job. Incremented every time a collect() method is called on context object, like wise,

- ③ Reduce Input Records

④ Redu

- ④ Reduce output records.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Annya Enclave,
Ameerpet, Hyderabad-500 018.
Ph: 040-5762 6719, 944 570 6789

- Counters are maintained by the task with which they are associated, and periodically sent to the task tracker and then to jobtracker so they all can be globally aggregated.
- The built-in job counters are actually maintained by the jobtracker, so they do not need to be set across the N/w unlike the all other counters, including the user defined ones.

Compression - MapReduce

Kelly T.
 Flat No. 1001,
 Annapurna Block, Jaya Enclave,
 Amravati, Hyderabad-500 016.
 Ph: 040-4452 0789, 098 570 6789

Compression In MapReduce :-

- compression reduces no. of bytes written to / read from HDFS.
- compression effectively improves the efficiency of network bandwidth and disk space.
- this saves the amount of data being transferred between map nodes to reduce nodes.

LZO:-

LZO is a compression/decompression library.

compression format : LZO

Hadoop Compression Codec:

com.hadoop.compression.lzo.Lzopcodec

Note:- "Codec" is the implementation of a compression-decompression algorithm. In Hadoop, a "codec" is represented by an implementation of the "compression codec" interface.

Why using compression:-

① Reduce storage requirements

② Speed up data transfers (across the N/w or from disks)

LZO Key characteristics:-

By Mr. GOPAL KRISHNA
If 6

- very fast decompression.
- Requires an additional buffer during the compression (size of 8 Kb or 64 Kb depends on the compression level).
- It does not require the additional buffer during the decompression other than the source and destination buffers... that's why fast decompression is possible with LZO.
- Allows the user to adjust the balance between compression ratio and compression speed, without affecting the speed of decompression.

Hadoop provides the below compression codecs:

- com.hadoop.compression.DefaultCodec
- com.hadoop.compression.lzoCodec
- com.hadoop.compression.SnappyCodec

Configurations Required for LZO compression - In mapred-site.xml:

By default compression is not enabled in the MapReduce (i.e. value is false). To achieve the compression we have to edit the below property of mapred-site.xml:

```
<property>
  <name>mapred.output.compress</name>
  <value>false</value>
</property>
```

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad, 500 016
Ph: 040 670 670 670 670 670

Note :- To enable the compression, value should be "true".

which compression codec to be used while compressing job output. By default "DefaultCodec" will be used. Inorder to use other than default (Lzo or Snappy), we have to replace their corresponding codecs.

Like Below:

```
<property>
  <name> mapred.output.compression.codec </name>
  <value> org.apache.hadoop.io.compress.DefaultCodec </value>
</property>
```

Note: 1. in place of DefaultCodec, give LzoCodec, SnappyCodec
2. we can also specify which compression codec to be used while compressing the map outputs.

Snappy Compression :-

Snappy is also compression/decompression library. It does not aim for maximum compression, or compatibility with other compression library.

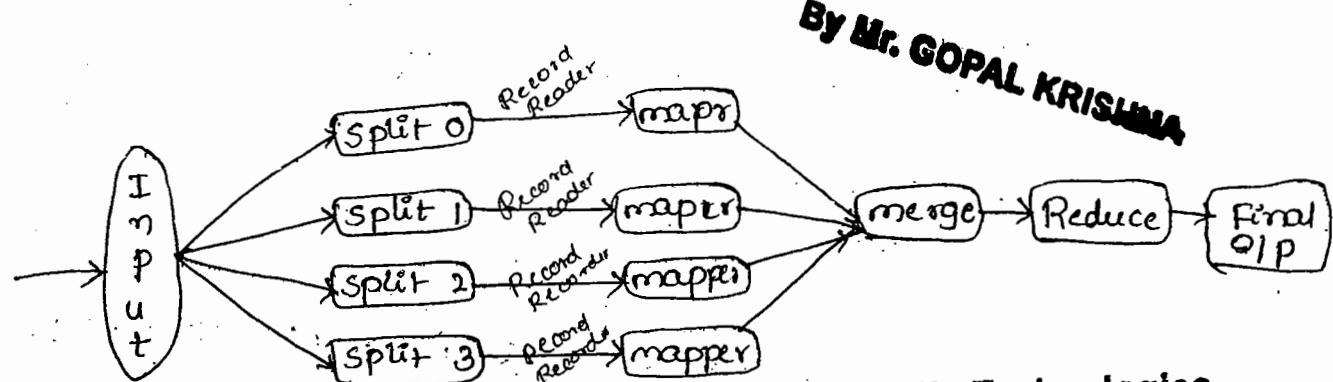
Snappy aims for very high speeds and reasonable compression.

Kelly Technologies
Flat No. 211, 3rd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 988 570 6789

Input Splits :-

17

- if mapReduce jobs can process the entire input in single shot there will not be any concept of parallelism with the same reason hadoop divides the input to mapReduce jobs into number of fixed size chunks known as Input splits (or) splits.



mapped · splits · max · size
mapped · split · min · size

Part - m - ooo → typical output of mapper
part - r - ooo → typical output of Reducer

Kelly Technologies

Flat No. 212, 2nd Floor,
Annapurna Block, Adivya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040 6462 6789, 928 570 6789

- mapReduce will not take the input as it is which the client gives instead; it divides the input data into multiple chunks which is known InputSplits. (or) splits.
- split size should always be equal to or greater than blocksize.

Note:- General practise would be blocksize should be equal to splitsize.

By the means of split concept mapReduce achieves the parallelism in Hadoop.

• if our split size is less than block size we will have so many smaller size splits and thereafter multiple mappers will be created on each and every split which will result into a poor performance.

• mapReduce job is a unit of work which client is expecting. This mapReduce jobs can be driven by two daemons like:

① JobTracker [which will co-ordinate the task & Scheduling & Rescheduling the task]

② TaskTracker [which is exactly responsible for execution of the task on the datanodes]

Creating Input and output formats in mapReduce jobs :-

① Input Formats

② Output formats.

Input Formats:-

- specification for reading data.
- defines how to break the inputs.
- provides RecordReader objects to read the files.
- Read data and converts to pairs.
- It is possible to write custom Input Formats.

File Input Format:- It is base implementation class for all the file formats.

• Base implementation for all the classes.

1) Text Input Format:-

- default format

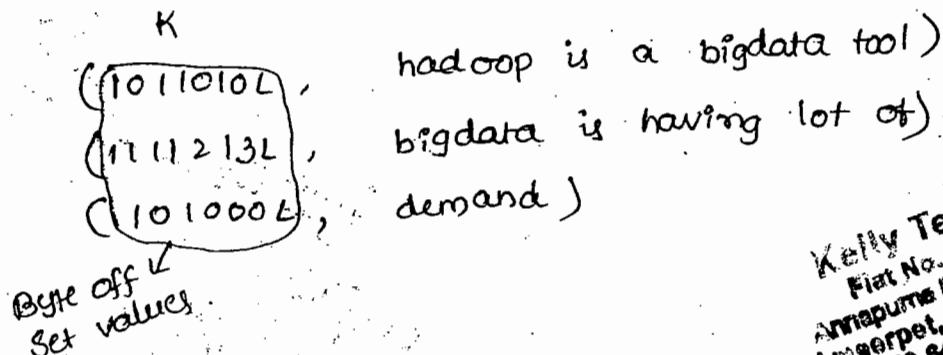
Here Key → System Byte Offset values

value → each and every line treated as value.

Kelly Technologies
Flat No. 212, 2nd Floor,
Parashuram Block, Antyodaya Enclave,
Ameerpet, Hyderabad, 500 016.
Ph: 040-5462 6716, 988 570 6730

- It is the default file format of mapreduce which will be used when the incoming data is in the form of "Text".
- Each and every line of the code is the record here and each & every record will be separated by a new line character '\n'.
- Text Input file format, Generally
 - Key → Byte offset values
 - value → The entire text of the record.

By Mr. GOPAL KRISHNA



hadoop is a bigdata tool)
bigdata is having lot of demand)

Kelly Technologies
Fiat No. 212, 2nd Floor,
Annapurna Block, Adivya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 5789, 988 570 6789

- plain Text Input
- split → single HDFS block (can be configured)
- Record → single line of text; line feed or carriage return used to locate end of line.

key → LongWritable - position in the file

value → Text - line of text.

=

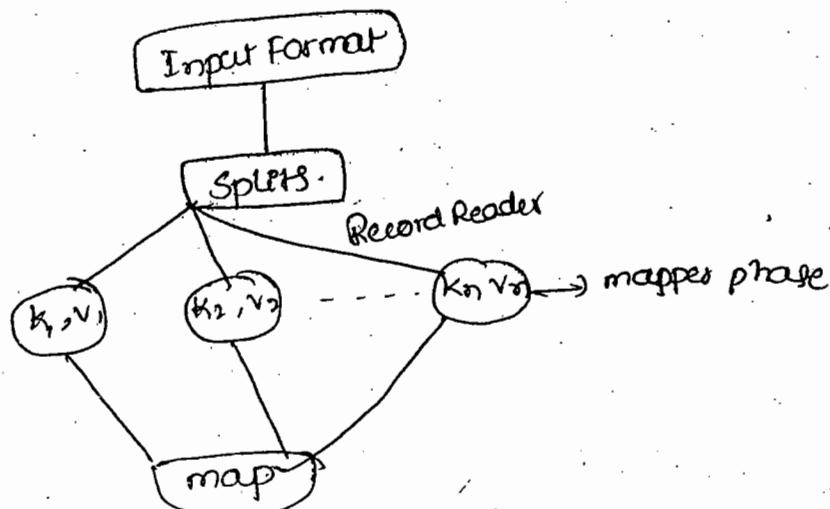
② Key-value Text Input format:- Each line as key-value (tab delimited)

- Key value Text Input Format will be used in mapReduce programming whenever we are getting the input data in the form of (K, V) .
- By default (K, V) TextInputFormat is '\t' 't' is the default delimiter (However we can change the same in code).
- As we have the specific key in this input formats the Byteoffset values will not get generated.

Ex:-
K
hadoop - t - bigdata.
Bigdata - t - Analysis.
Android - t - mobileApps.

③ N-line Input format :-

- Same as Text Input format but splits equal to configured N-lines.
- Number of the maximum records in each split and every split.



- Record Reader is a class which is primarily responsible for the (K, V) pairs given in the Input Splits.

• whenever we are providing the data multiple splits into mapper functions, if we have variable no. of records and each and every split then we will not have proper control on when exactly a particular split will be completed.

• with the same reason if we want to put a fixed no. of records each & every split then we can go ahead with NLineInputFormat.

Split → NLines: configured via mapred.line.input.format or NLineInputFormat.setNumLinesPerSplit(job, 100);

Record → single line of text

Key → LongWritable - position in the file

value → Text - line of text.

By Mr. GOPAL KRISHNA

④ Table Input Format:-

- converts data in HTable to format consumable to mapReduce.
- mapper must accept proper key/values.

split → Rows in on HBase Region (provided scan may narrow down in the result)

Record → Row, returned columns are controlled by a provided scan.

Key → Immutable Bytes writable.

value → Result (HBase class)

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad - 500 018.
Ph: 040 2462 6739, 988 570 0739

Sequence file Input Format:-

- Hadoop specific binary representation.
- Special type of file to store key-value pairs.
- Store keys and values as byte arrays.
- uses length encoded bytes as format.
- often used as input or output format for MR jobs.
- Has built in compression on values.

Output Formats:- specification for writing data.

How the result $\langle \text{key}, \text{value} \rangle$ pairs are written into files.

- possible to write custom output formats.
- Hadoop default output formats.
 - Text output Format
 $\langle \text{key} \text{ TAB } \text{value} \rangle$
 - Sequence File Output Format.
 - Hadoop specific binary representation.

validate output specification for that job.
you may have seen annoying messages that output directory already exists.

- creates implementation of Record writer.
- Responsible for actually writing data.

- creates Implementation of output committer.
 - Setup and clean-up job's and task's artifacts (ex:- directories)
 - commit or discard tasks output.

Kelly Tech Solutions
Flat No. 211, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarpet, Hyderabad-500 018
P.M. 040 262 6739, 988 570 6739

① Text output format :-

- outputs plain text.
- saves key-value pairs separated by tab.
- configured via mapreduce.output.textoutputformat property
- set output path

TextOutputFormat.setOutputPath(job, new Path("mypath"));

② Table output Format :-

- saves data into HTable.

- Reducer output key is ignored.
- Reducer output value must be objects.

Data localization in mapReduce :-

Mapper will be stored in localfile.

- Every mapper will be completed by tasktracker, once the mapper phase has been completed by tasktracker, the mapper o/p will be stored in the localfile system (but not on HDFS). The reason might be all the mapper output from respective tasktracker put on the NIC in order to consume by the Reducer phase.

- In this whole process due to some N/W errors or between related errors, if we loss the mapper o/p, we have to redo the entire process of mapper which is a time consuming process and performance overhead.

- To avoid the same, the mapper o/p will be stored in the localfile system. This is called data localization.

Note :- Data localization is only for the mapper but not for reducer, The reason might be the o/p whatever we get from reducer is the final output.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph. 040 492 07000, 070 070 070 070

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph. 040 492 07000, 070 070 070 070

By Mr. GOPAL KRISHNA

Combiner :-

- To optimize the N/W & bandwidth limitations combiner concept will be used in mapReduce programming.
- Combiner will act as local Reducer (or) mini Reducer as whatever the data consumed by Reducer phase, the same copy will reside in the combiner.
- Hadoop does not provide any guarantee on combiner's execution.
- Hadoop may call combiner function zero, one or many times for a particular map output record.

Note :-

- while calling reduce method, Hadoop does not provide any guarantee on values in stored order corresponding to a particular key, to achieve this we used the Secondary sorting.
- The combiner function does not replace the reduce function.
- Combiner is specify the combiner function.

Partitioners :-

- partitioner allows you to distribute how outputs from the map storage are sent to the reducers. Basically it partitions the Keyspace.
- partitioner controls the partitioning the keys of the intermediate map-outputs. The key [subset of the key] is used to derive the partition.

The total no. of partitions is same as the no. of reduce tasks for the job.

Kelly Technologies
Flat No. 212, 2nd floor
Annapurna Block, Adyar
Amarpet, Hyderabad - 500 010.
Ph: 040 4992 5789, 988 570 6789

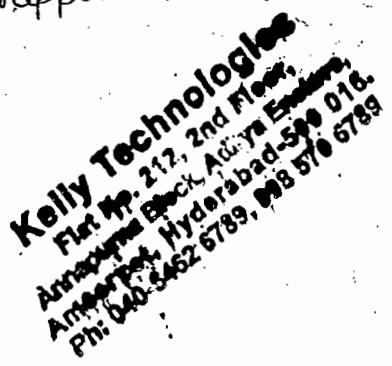
- partitioner run on the same machine after mapper had completed it's execution, by consuming mapper o/p. so entire mapper output(record) is sent to partitioner and partitioner forms r (no. of reduce task) groups from the mapper output.

- By default hadoop framework Hashbased partitioner. This particular partitions the keyspace by using the hashCode.
- The following is logic Hashpartitioner executes to determine a reducer for a particular 'key'.

$$(\text{key.hashCode}() \& \text{Integer.MAX_VALUE}) \% \text{num Reduce Tasks.}$$

How to write custom partitioner :- By Mr. GOPAL KRISHNA

- To have Hadoop use a custom partitioner you will have to do minimum the following three.
 - Create a new class that extends partitioner class.
 - Override method get partition.
 - In the wrapper that runs the mapReduce, either
 - add the custom partitioner to the job programmatically using method set partitioner class or
 - * add the custom partitioner to the job as a config file (if your wrapper reads from config file or oozie).



Melly Technologies
Flat No. 213, 2nd Floor,
Annakurni Block, 14th Main
Amarapuri, Hyderabad - 500 073
Ph: 243-8862 8789, 988-570 572

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;
```

```
public class WordCount {
```

```
    public static class TokenizerMapper extends Mapper  
<Object, Text, Text, IntWritable> {
```

```
        private final static IntWritable one = new IntWritable(1);
```

```
        private Text word = new Text();
```

```
        public void map(Object key, Text value, Context context)
```

```
            throws IOException, InterruptedException {
```

```
            StringTokenizer itr = new StringTokenizer(value.toString());
```

```
            while (itr.hasMoreTokens()) {
```

```
                word.set(itr.nextToken());
```

```
                context.write(word, one);
```

```
}
```

```
}
```

map
red
ct

```

public static class IntSumReducer extends Reducer
    <Text, IntWritable, Text, IntWritable> {
    private IntWritable result = new IntWritable();
    public void reduce(Text key, Iterable<IntWritable> values,
        Context context) throws IOException,
        InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        result.set(sum);
        context.write(key, result);
    }
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "word count");
        job.setJarByClass(WordCountNew.class);
        job.setMapperClass(TokenizerMapper.class);
        job.setCombinerClass(IntSumReducer.class);
        job.setReducerClass(IntSumReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}

```

52

*Kelly Technologies
Flat No. 212, 2nd Floor,
Anapurna Block, 10th Main,
Ameerpet, Hyderabad - 500 016
Ph: 040 452 6789, 988 578 6789*

Configuration is a base class from Hadoop (mapred-site.xml)

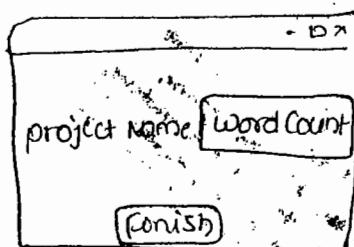
By Mr. GOPAL KRISHNA

- Job is a particular class. It is derived from base class i.e. Configuration.
- We can create the job compulsory conf object will be created.

How to create program in NetBeans (or) MyEclipse IDE :-

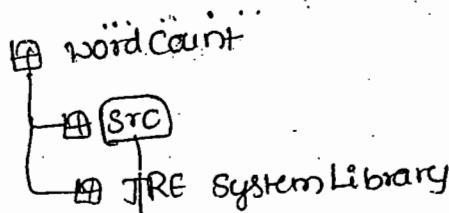
Step 1 :-

File → New → Java project → click OK.

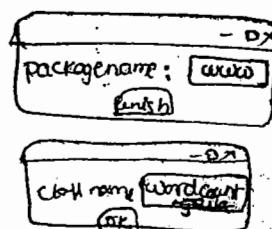


Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarpet, Hyderabad 500 018.
Ph: 040-6462 6789, 988 570 6789

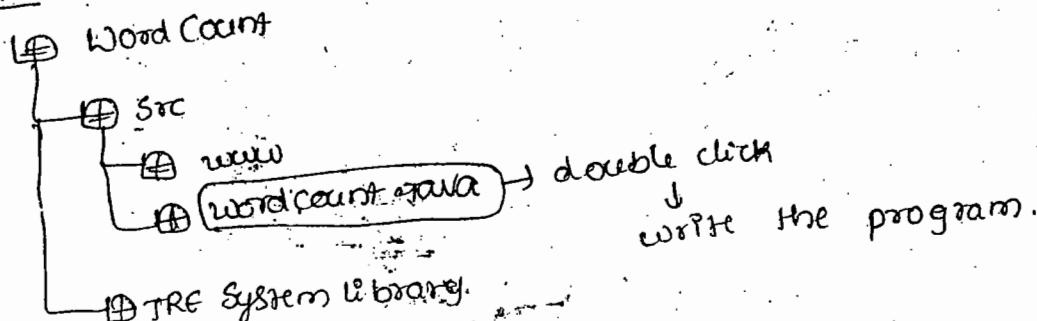
Step 2 :-



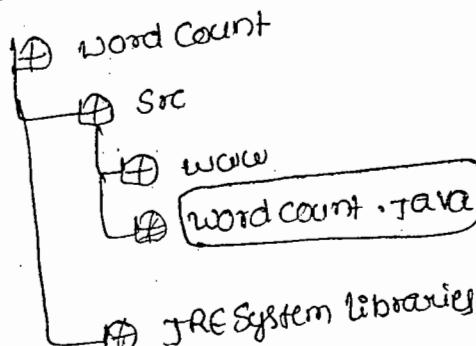
Right click → New → package →
→ class →



Step 3 :-



Step 4 :- How to create add jar file.



right click → Build path
↓
configure Buildpath
↓
libraries
↓
Add external jar files
OK

Steps : How to export file jarfile to Linux.

53

④ wordcount

④ SRC

④ ww

④ wordCount.java

④ JRE System Library.

Right click → export → Java

JAR file

double click

Browse

→ click

give the

JAR name Save

Step 6 : Open the vmware.

ifconfig

192.168.2.35.131

Step 7 : Start → Rightclick → open windows explorer



//192.168.2.35.131 ↵

By Mr. GOPAL KRISHNA



downloads.

open → In this Downloads copy paste

the jar file.

Step 8 :-

Downloads # ls ↵

wordcount.jar

Step 9 : # nano Input.txt ↵

hadoop fs -mkdir anitha ↵

hadoop fs -put Input.txt anitha ↵

Kelly Technologies
Flat No. 2/2 2nd Floor,
Annapurna Bhk X 1/17 Emanya,
Ameerpet, Hyderabad - 500 018.
Ph: 940-64526789, 938-570-6789

Step 11 : # hadoop jar <Runnable jar name> class name <without extns>
<HDFS Inputpath with file> <HDFS output path file>

hadoop jar wordcount.jar wordcount anitha/Input.txt
anitha/output ↵

hadoop fs -ls anitha/output ↵ hadoop fs -cat anitha/output

	Traditional RDBMS	MapReduce
Data Size	GB - TBs	TBs - PBs
Access	Interactive & Batch	Batch
Updates	Read/write many times	Write once, read many times
Structure	static schema	Dynamic Schema
Integrity	high	Low
Scaling	non linear	Linear

MapReduce optimizations:-

- overlapping maps, shuffle and sort
- mapper locality.
 - * schedule mappers close to the data.
- combiner
 - * mappers may generate duplicate keys
 - * side effect free reducer can be run on mapper node.
 - * minimizes data size before transfer
 - * Reducer still is run

Speculative execution to help with load-balancing.

- * Some nodes may be slower.
- * Run Duplicate task on another node, take first answer as correct and abandon other duplicate tasks.
- * only done as we start getting toward the end of the tasks.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Adya Enclave,
Ameerpet, Hyderabad-500 010.
Ph: 040-462 6789, 986 570 0789

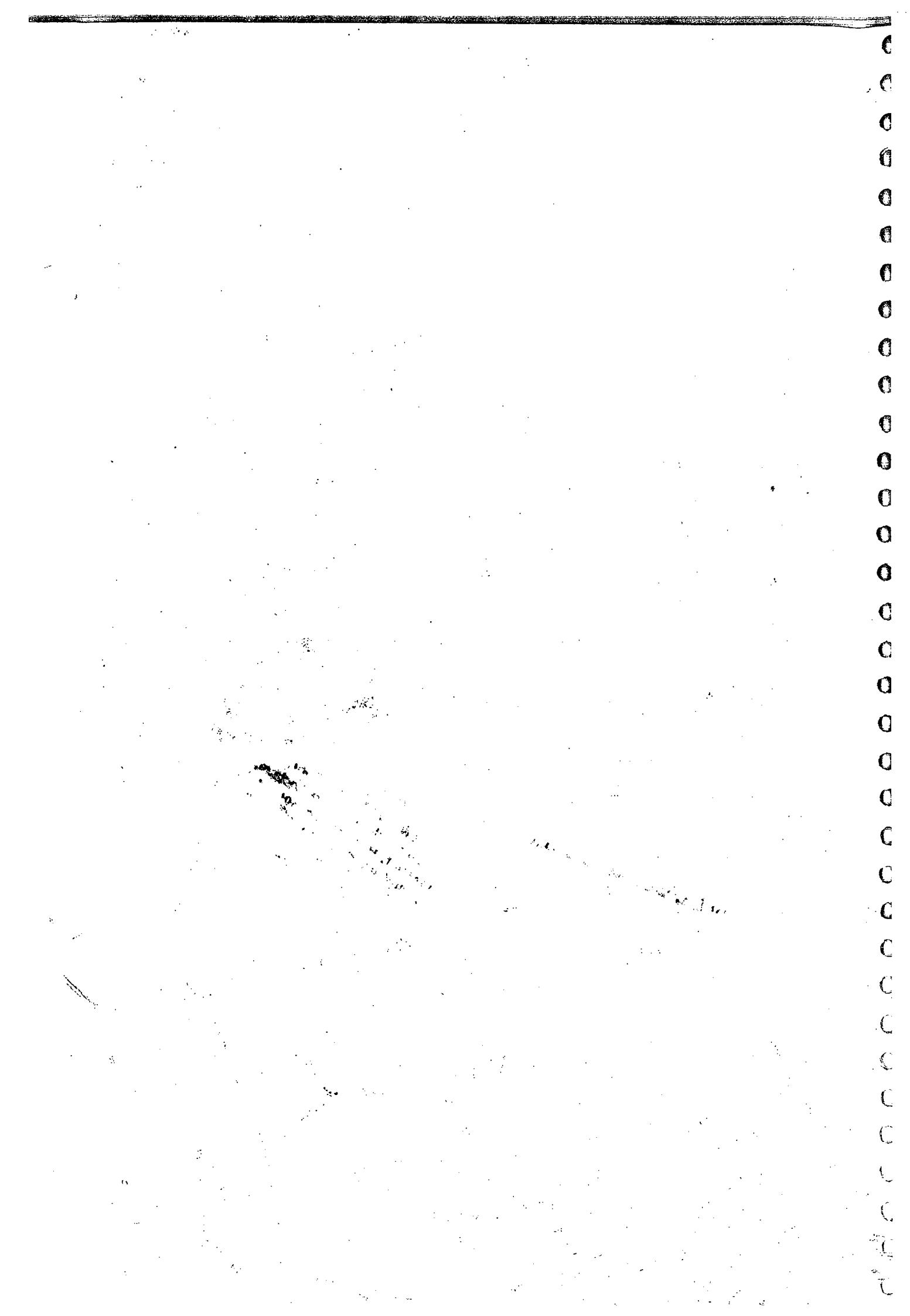
Serialization:

54

- Transforms structured objects into a byte stream.
 - for transmission over the network. Hadoop uses RFT
 - for persistent storage on disks.
- Hadoop uses its own serialization format, `Writable`
 - comparison of types is crucial (shuffle and sort phase)
Hadoop provides a custom `RawComparator`, which avoids deserialization.
 - `Custom Writable` for having full control on the binary representation of data.
 - Also external frameworks are allowed: enter Avro.
- Fixed length or variable-length encoding?
 - fixed length: when the distribution of values is uniform.
 - variable length: when the distribution of values is not uniform.

~~By Mr. Durai Nadesan~~

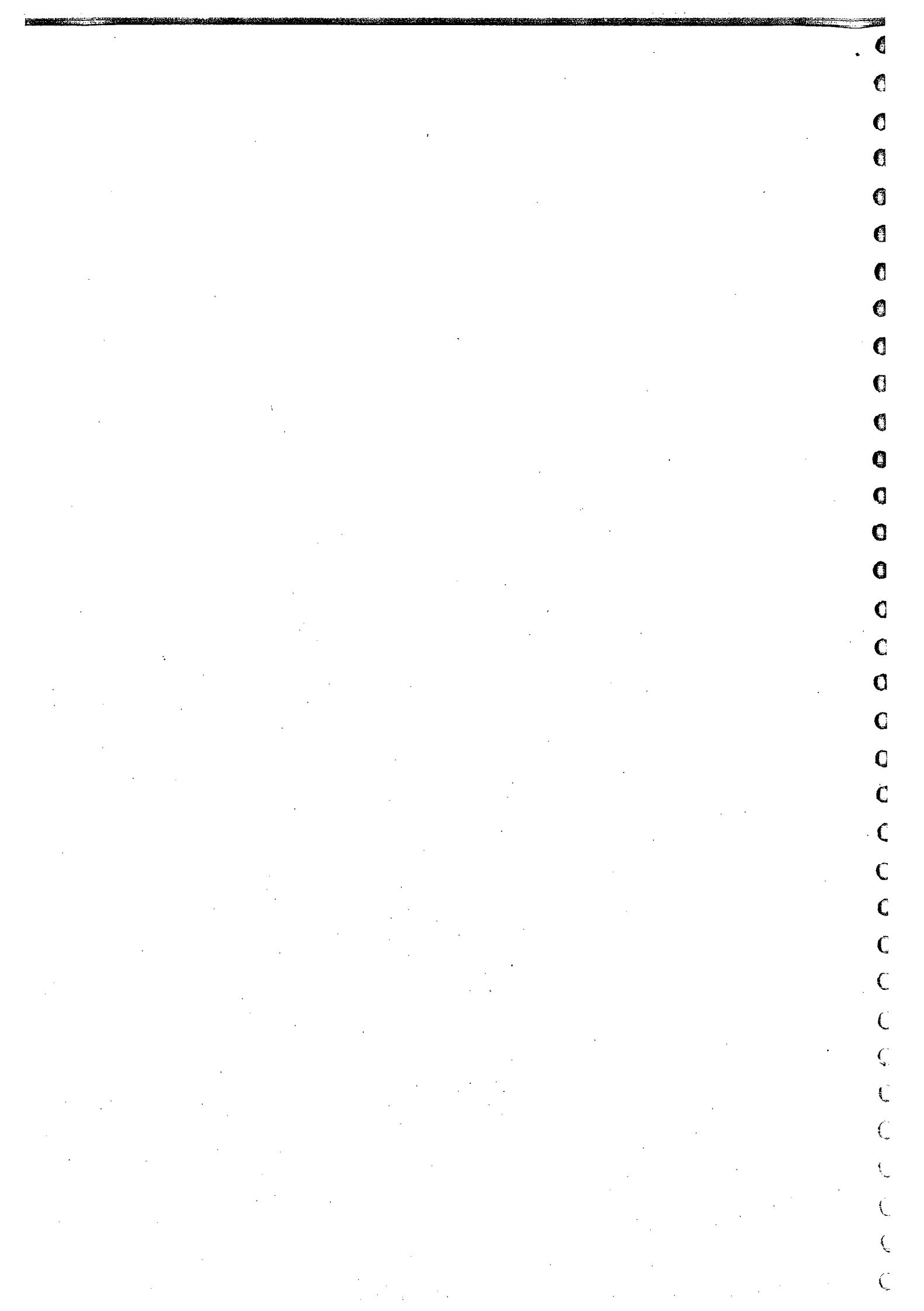
Kelly Technologies
Flat No. 212, 2nd Floor,
Arimparama Block, 11th Main Road,
Aeroport, Hyderabad - 500 021
Ph: 040-44526788, 080-5705100



[Pig]

55

- Apache Pig is one of the components of Hadoop which is an abstract layer (high level procedural language) on top of Hadoop and mapreduce platform.
- Pig is also known as data flow language (or) Transformation language.
- pig was initially introduced at yahoo laboratories at very first time inorder to address their adhoc request.
- pig language was officially opted by Apache in the year 2006 from there onwards ~~by Mr. Gopal Krishna~~ used to pig language called Apache Pig.
- The language to express pig language transformation is known as "PigLatin".
- As we are dealing with multiple transformation and dataflow will be through the transformations. we can call pig is a transformation language (or) Dataflow language.
- pig is a simple language that externally executes statements.
- A statement is an operator that takes input (such as bag represents a bag as its output.)
- pig latin are not case-sensitive. UDF names are case sensitive.
Ex: COUNT → count
By: COUNT → COUNT



pig vs SQL

5b.

pig

SQL

- | | |
|-------------------------------------|--|
| ① pig is procedural (How) | ① SQL is declarative. |
| ② Nested relational data model | ② Flat relational data model |
| ③ Schema is optional | ③ schema is required |
| ④ scan - centric analytic workloads | ④ OLTP + OLAA workloads |
| ⑤ Limited query optimization | ⑤ significant opportunity for query optimization |

Advantages of using pig:-

- ① pig can be treated as a ~~higher level language~~ By M. GOPAL KRISHNA
 - increases programmer productivity
 - decreases duplication of effort.
 - opens the MR programming system to more users
 - reduces hadoop complexity.
- ② pig insulates against hadoop version upgrade
 - Hadoop configuration tuning.
 - job configuration tuning.

Pig Features :-

- Data-flow Language ; - a sequence of steps where each step specifies only a single high-level data transformation.

The step-by-step method of creating a program in Pig is much cleaner and simpler to use than the single block method of SQL. It is easier to keep track of what your variables are, and where you are in the process of analyzing your data.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad 500016.
Ph: 040-64626716, 64626723.

- user Defined function (UDF)

- Debugging Environment

- Nested data model

what is pig :-

- pig is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs.
- pig generates and compiles a map/Reduce programs on the internally.

why Pig :-

- Ease of programming :- It is trivial to achieve parallel execution of simple, "embarrassingly parallel" data analysis tasks.
complex tasks comprised of multiple interrelated data transformations are explicitly encoded as data flow sequences, making them easy to write, understand and maintain.

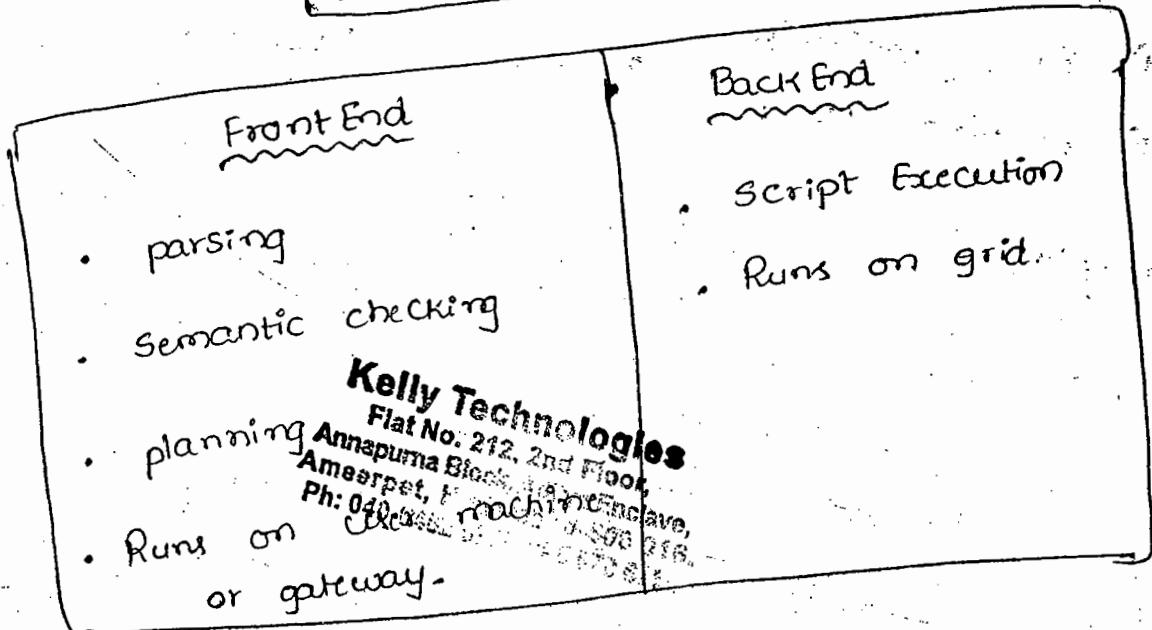
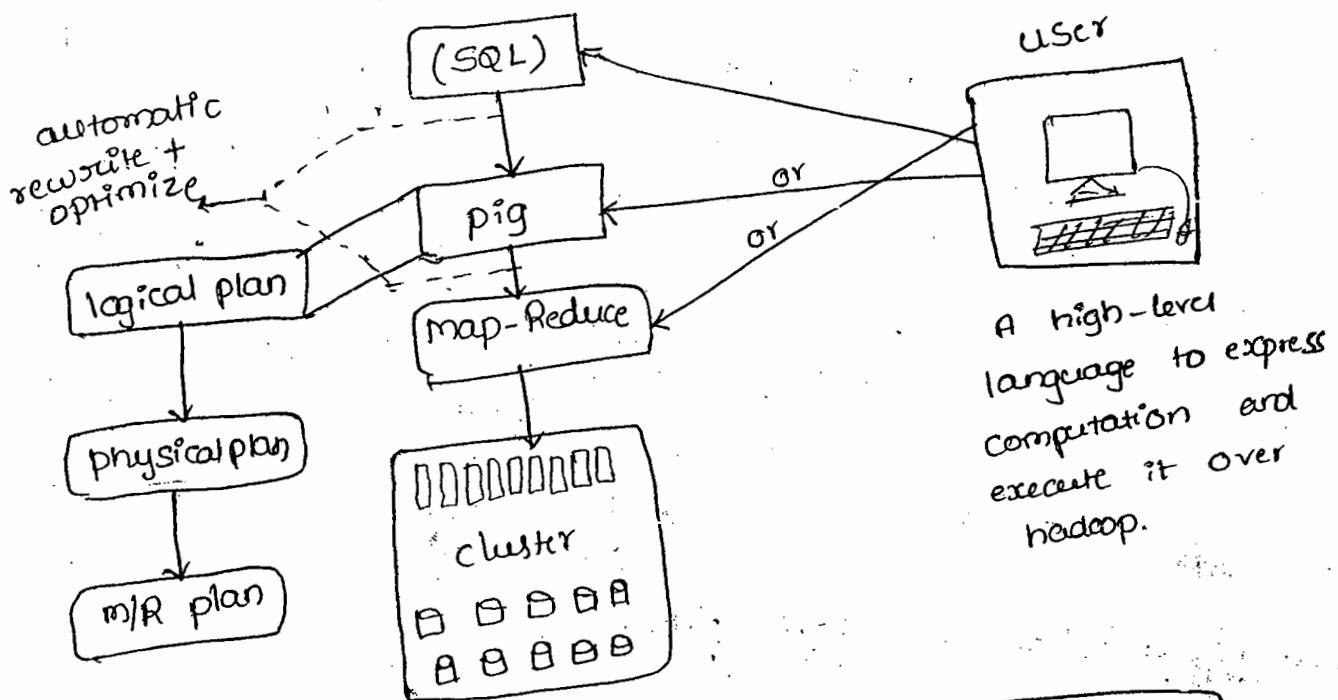
Pig takes care of ...

- Schema and type checking.
- Translating into efficient physical Dataflow.
• (i.e sequence of one or more MapReduce jobs)
• Exploiting data reduction opportunities.
• (e.g. early partial aggregation via a combiner).

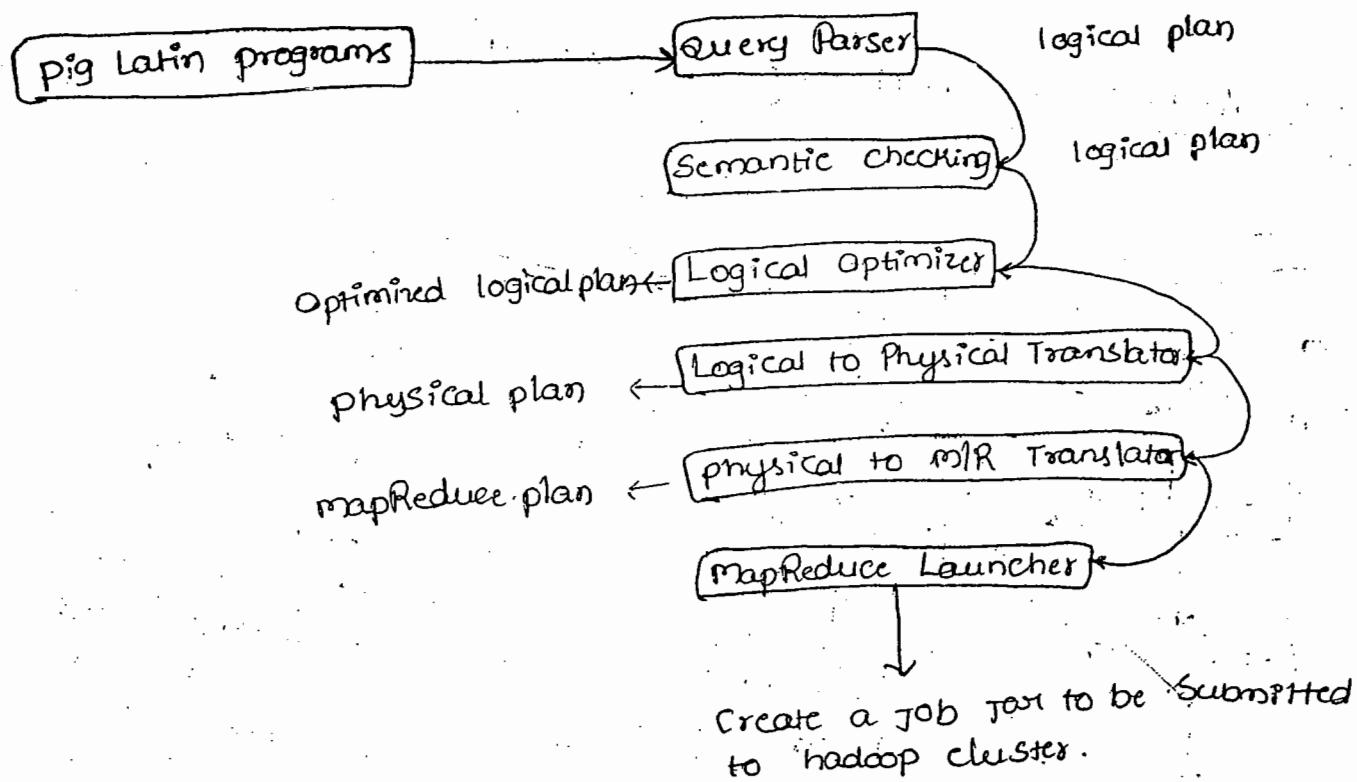
- Executing the system-level dataflow.
(i.e running the mapreduce jobs)
- Tracking progress, errors, etc.

57

Pig Architecture (mapReduce as Backend)



Front End Architecture :-



During optimization of logical plan

- ① Early projections
- ② pushing limit
- ③ Streaming Load and Store.

Logical plan :-

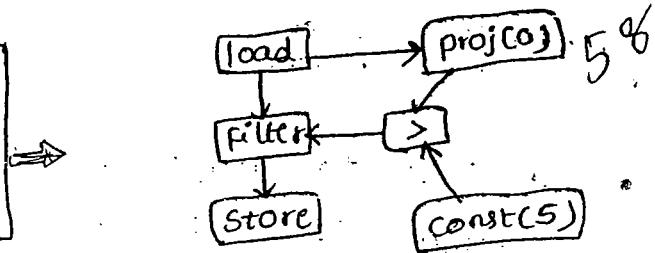
- Directed Acyclic Graph
- logical operators as node
- dataflow as edges
- Logical operators
 - one per pig statement
 - type checking with schema

Eg:-

```

a = load 'my file';
b = filter a by $0 > 5;
store b into 'myfiledir';

```



physical plan :-

- Pig supports two back-ends.

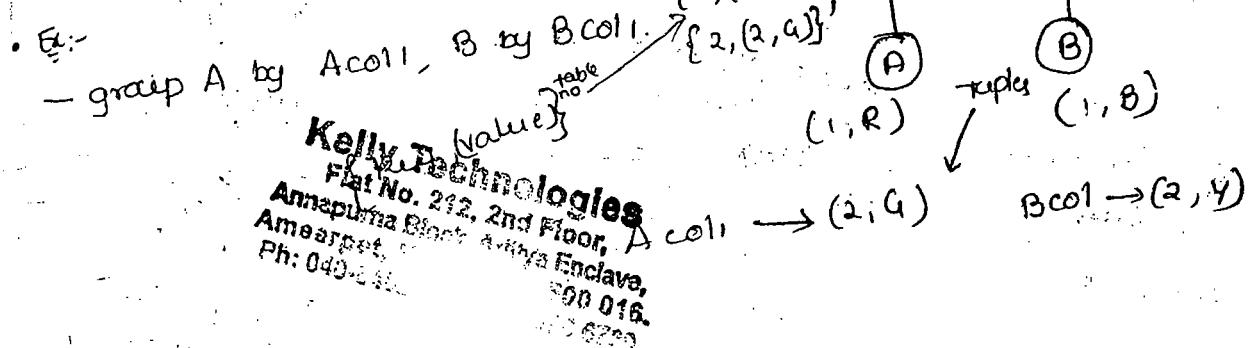
- Local
- Hadoop mapReduce

logical to physical translation

- 1:1 correspondence with most logical operators.
- Except cross, Distinct, Group, a co-group and Order.

Eg:- logical to physical plan for group operator :-

- Logical operator for co-group/group is converted to 3 physical operators.
- Local Rearrange (LR)
- Global Rearrange (GR)
- Package (PKG)



Kelly Technologies
Flat No. 212, 2nd Floor, A col1 → (2, Y)
Annapurna Block, Kalya Enclave,
Amarpet, 500 016.
Ph: 040 2441 6829

mapReduce plan :-

- physical to MapReduce (m/R) plan conversion happens through the MR Compiler.
 - converts a physical plan into a DAG of m/R operators.
- Boundaries for m/R include cogroup/group, distinct, cross, order by, limit (in some cases).
- push all subsequent operators between co group to next co group into reduce.
- order by is implemented as 2 m/R jobs.
- JobControl compiler then uses the m/R plan to construct a JobControl object.
- job.jar is created and submitted to Hadoop Job Control.

Back End Architecture :-

- pig operators are placed in a pipeline and each record is then pulled through.

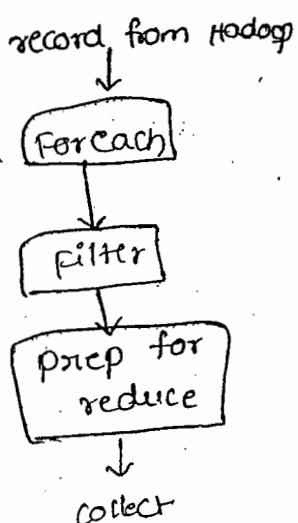
A = load ...

B = foreach ...

C = filter ...

D = group ...

prep phase will have foreach & filter



Lazy Execution :-

- nothing really executed, until you request output.
- only when the store, dump, explain, describe, illustrate command is encountered does the MR execution takes place either on your cluster or front end.

Advantages :-

- In-memory pipelining.
- filter re-ordering across multiple commands.

parallelism :-

- split wise parallelism on map-side operators.
- By default, 1 reducer
- PARALLEL keyword.
- group, cogroup, cross, join, distinct, order

Running modes of execution

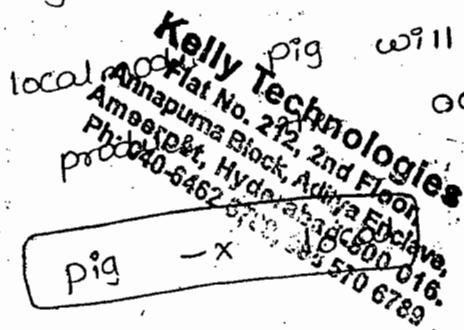
① Local mode

② mapreduce mode

Local mode :-

TO run pig in local mode, you need single machine. (All files are installed and run using your local host & file system)

In the local mode &



Syntax :-

pig -x

- local file system is used.
- local host and useful for prototyping & debugging.

② Hadoop(MapReduce) Mode :- TO run pig in Hadoop (mapreduce) mode, you need access to a Hadoop cluster and HDFS installation. mapreduce mode is the default mode.

In mapreduce mode, pig will expect the input from HDFS and produce the output from on top of HDFS.

Syntax:- `pig (or) pig -x mapred`

• Take advantage of a Hadoop cluster and HDFS.

Running Pig :-

① Interactive mode (Grunt shell)

② Batch mode (script mode)

③ Java mode (embed pig statements in java)
• Keep pig.jar in the class path

① Grunt Shell :- Grunt shell is an interactive based shell. which means where exactly we will get the output than and their itself. whether it is success or fail.

`syn. grunt >`

this mode
is default
mode.

`pig -x local <
grant>`

local mode
of execution of
pig

`> pig <
grant >`

mapreduce mode
of execution of
pig.

② script mode :- In script mode all the commands [with transformations] will be included in a single file called .pig file.

Ex:- sample.pig ↴

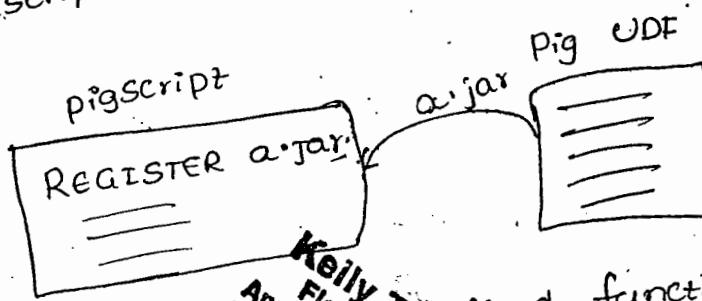
- And this execute in .pig file, the whole file will be executed in a single shot and we will get the result.

local mode	pig -x local sample.pig
mapreduce mode	pig sample.pig

③ Embedded mode :-

- if certain things are directly possible through the transformation user defined functions (Pig we can also provision of writing the some UDF code with customizations) and we will COPAL KRISHNA in our pigscript.

- In Embedded mode of pig executions with in the our pigscript we will call the javameade.



Incorporate

REGISTER → Registers

DEFINE → creates an script (or)

IMPORT → import macros defined in separate file into a script

Kelly Technologies
Flag No. 212, 2nd Floor,
Amanapuri, Hyderabad
Ph: 040-64444444

Defined functions :-

With the pig runtime

a macro, UDF, streaming alias command specification.

Data types of Pig :- pig data types can be divided into two types.

- ① scalar types → contain a single value
- ② complex types → contain other types.

- ① scalar types :-

Normal language	Pig
int	int
String	chararray
float	float
double	double
long	long
boolean	boolean
byte	bytearray
tuple	tuple
bag	bag
map<object, object>	map

- ② complex types :-

Custom data types in pig :-

- ① atom :- atom is the smallest unit in apache pig which we can denote with a single character.

Each and every character is a atom.

Ex:- empid ename esal
 100 xyz 10000

② tuple :- collection of all these fields (or) ordered set of fields.

• Field can be simple (or) Complex type

• Nested relational model.

Ex:-	empid	empname	esal
	100	xyz	10000

Field :- field is nothing but collection of atoms.

Ex:-	empid	empname	esal
	100	xyz	10000

Collection of all these tuples is known as

③ Bag :-

Bag.

Ex:-

empid	empname	esal
100	xyz	10000
101	abc	4000
102	cde	6000

Mr. GOPAL KANISHKA

④

map :- Set of (key, value) pairs. → keys must be unique

get to be
→ character

can be scalar (or)
complex datatype.

Type

① Data Atom
an atomic value

Example

'alice'

② Tuple
a record

(alice, lakers)

Kelly Technologies
Flat No. 212, 2nd Floor,
Ameerpet, Hyderabad
Ph: 040 6716 8899

③ Data Bag
Collection supporting scan

{(alice, lakers),
(alice, kings)}

④ Data Map
Collection with key
lookup

[(lakers, {lakers}),
(kings, {kings})]
'age' # 22

→ maps must be enclosed by square bracket.
→ Each key-value pair must be separated by #.

→ Key must be a character.

Pig Storage :-

- Loads or stores relations using field-delimited text format.
- Each line is broken into fields using a configurable field delimiter (defaults to a tab character) to be stored in the tuple's fields. It is the default storage when none is specified.

Bin Storage :-

- Loads or stores relations from or to binary files. An internal pig format is used that uses hadoop writable objects.
- It stores arbitrarily nested data.
 - used for loading intermediate results that were previously stored using it.

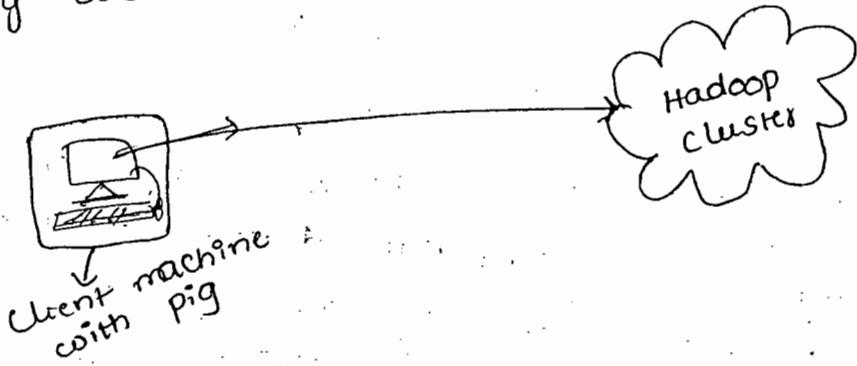
Text Loader :-

- loads relations from a plain-text format. Each line corresponds to a tuple whose single field is the line of text.

where does pig live?

62

- pig is installed on user machine.
- No need to install anything on the Hadoop cluster.
- pig and Hadoop versions must be compatible.
- pig submits and executes jobs to the Hadoop cluster.



what is pig used for?

ICM UW → discover knowledge from large collections
of documents (academic papers)

LinkedIn :- Ex: discovering "people you may know"

Mendeley :- Ex: finding trending keyword

Nokia(Ovi Maps) :- exploring unstructured data sets from logs,
database dumps and data feeds.

Real web :- computing statistics on banner views, clicks and
user behaviour on target websites after click.

- weblog processing
- Data processing for web search platforms.
- Adhoc queries across heterogeneous data sets.
- Rapid prototyping algorithms for processing large data sets
- ETL
- machine learning
- Text processing (Semi, Structured)
- data analysis
- processing many data sources

Kenny Technologies
Flat No. 212, 2nd Floor,
Aneguduma Block, Aonisa Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-64626116, 98455706789

→ what is the difference between mapReduce & pig.

mapreduce

- ① mapReduce expects the programming language skills for writing the business logic.

pig

- ① pig there is no much of programming skills. As we are writing whole logic will making use of pig transformation (or) operators.

- ② if we can do any change in the mr program, we need to certain problems we can change the process entire.

- ① compiling the program
 - ② executing the program
 - ③ packing up the "
 - ④ deploying the same (cluster environment)
- ③ As a General saying of Hadoop mr program write 200 lines of mr code.

- ⑤ mapReduce requires multiple stages, leading to long development life cycles.
Rapid prototyping, Increased productivity.

- ② In the pig, we can completely dealing with simple scripting we can avoid other transaction process.

- 5% of the mr code
- 5% of the mr development time
- Readable & reusable.
- Increases programmer productivity
- 25% of the mr execution time.

- ③ In pig we can that type of mr program, we can write 10 lines of code.

- ④ pig provides the log analysis
 - ad hoc queries across various large data sets.
 - analytics.

Note:- Even though we can do processing, internally convert to a MR job.
every statement

Pig Transformations (or) Pig Operators :-

13

<u>Statement</u>	<u>Description</u>
① load	→ read data from the file System
② store	→ write data to the file system
③ dump	→ write output to stdout
④ foreach	→ Apply expression to each record and generate one or more records.
⑤ filter	→ Apply predicate to each record and return records where false.
⑥ group/cogroup	→ collect records with the same key from one more inputs.
⑦ join	→ joins two or more inputs based on a key.
⑧ order	→ sort records based on a key.
⑨ distinct	→ Remove duplicate records.
⑩ union	→ merge datasets.
⑪ limit	→ limit the number of records.
⑫ split	→ split data into more sets, based on filter conditions.
⑬ cross	→ creates the cross product of two (or) more relations.

Pig Diagnostic Operators :-

<u>Statement</u>
① describe
② dump
③ explain
④ illustrate

<u>Description</u>
→ Returns the schema of the relation.
→ Returns results to the screen.
→ Amrita Bhawan Block A, 2nd Floor, 1st Phase, HSR Layout, Bangalore, India - 560 016.
→ Displays a plan.
→ Displays a step-by-step execution of a sequence of statements.

Relational operators in pig

- ④ Filter, foreach, Group, orderBy, sample, distinct, limit, join(inner), join(outer).

Built-in functions in pig :-

- ① EVAL ③ string ⑤ Tuple, Bag, Map Specific functions
- ② MATH ④ date/time

Evaluation(EVAL) functions in Pig :-

- ① AVG : - calculates the average value of entries in a bag.
- ② CONCAT : - concatenates two byte arrays or two character arrays together.
- ③ COUNT : - calculates the no. of entries in a bag.
- ④ COUNT-STAR : - calculates the no. of elements in a bag.
- ⑤ DIFF : - compares two fields in a tuple.
- ⑥ IsEmpty : - checks if a bag or map is empty.
- ⑦ MAX : - computes the Max function to maximum of the numeric values or chararrays in a single-column bag.
- ⑧ MIN : - MIN function to compute the minimum of a set of numeric values or char-arrays in a single column bag.
- ⑨ SIZE : - computes the no. of elements based on any pig data type.
- ⑩ SUM : - compute the sum of a set of numeric values in a single - column bag.
- ⑪ TOKENIZE : - TOKENIZE function to split a string of words (all words in a single tuple) into bag of words (each word in a single tuple)

math functions in Pig :-

b6

- ① ABS :- Returns the absolute value of an expression.
note :- expression result is int, long, float or double.
- ② ACOS :- Returns the arc cosine of an expression.
expression whose result is type double.
 ACOS(expression)
- ③ ASIN :- Returns the arc sine of an expression.
- ④ ATAN :- Returns the arc tangent of an expression.
- ⑤ CBRT :- Returns the cube root of an expression.
- ⑥ CEIL :- Returns the value of an expression rounded up to the nearest Integer.
- ⑦ COS :- Returns the trigonometric cosine of an expression.
- ⑧ COSH :- Returns the hyperbolic cosine of an expression.
- ⑨ EXP :- Returns Euler's number e raised to power of x.
- ⑩ FLOOR :- Returns the value of an expression rounded down to the nearest Integer.
- ⑪ LOG :- Returns the natural algorithm (base e) of an expression.
- ⑫ LOG10 :- Returns the base 10 algorithm of an expression.
- ⑬ RANDOM :- Returns a pseudo random number.
- ⑭ ROUND :- Returns the value of an expression rounded to an integer.
- ⑮ SIN :- Returns the sine of an expression.

Kelli Technologies
Plot No. 212, 2nd Floor,
Annapurna Block, Dabba Enclave,
Ameerpet, Hyderabad-500 076
Ph: 040-5462 6719, 069 570 6789

- (15) SINH :- Returns the hyperbolic sine of an expression.
- (16) SQRT :- Returns the positive square root of an expression.
- (17) TAN :- Returns the trigonometric tangent of an angle.
- (18) TANH :- Returns the hyperbolic tangent of an expression.

Pig commands :-

- (1) Load :- read the data from file System.

`grant > A = load 'input.txt' using PigStorage();`

- (2) dump :- display the output itself.

`dump A;`

- (3) Store :- write data to the file system.

`Store 'B' into 'output';`

- (4) Foreach :- display specific column.

`B = foreach A generate id;`

- (5) order :- display the records ascending order.

`> B = order A by name;`

- (6) desc :- display the records descending order.

`B = order A by name desc;`

⑦ Filter :- Apply predicate to each record and remove records where false.

B = filter A by \$1 == 'Gopal' (or) name == 'Gopal'

⑧ limit :- display the records from begining.

B = limit A 2;

⑨ Distinct :- Removing all duplicate values

B = Distinct A;

⑩ Group By / coGroup By :- Group the data in single Relation
is called the Group By.

Group the data in two (or) more relations is
called the coGroup By.

B = group A by id;

C = coGroup A by id;

⑪ split :- splits a relation into two (or) more relations

split A into B if id == 100,
C if name == 'Anitha';

⑫ Join :- Joins two or more relations.

C = Join A by B - by id;

intersection

KELLY TECHNOLOGIES
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-5462 6709, 088 570 6789

(13) Left outer join :- Joins two (or more) inputs based on a key is LEFT.

$C = \text{join} A \text{ by id LEFT, } B \text{ by id;}$

(14) Right outer join :- Joins two (or more) inputs based on a key is Right.

$C = \text{join} A \text{ by id RIGHT, } B \text{ by id;}$

(15) Cross :- Creates the cross product of two (or more) relations.

$C = \text{cross } A, B;$

(16) Describe :- prints relation's schema.

$\text{describe } A;$

(17) Illustrate :- Displays a step-by-step execution of a sequence of statements.

$\text{illustrate } A;$

(18) Explain :- prints the logical and physical plans.

$\text{explain } A;$

(19) union :- Combines two or more relations into one.

$A = \text{union } A, B;$

(20) FLATTEN :- unpack (or) unnest all the things.

$B = \text{group } A \text{ by id;}$

$c = \text{foreach } B \text{ generate group, FLATTEN}(A);$

2) TOKENIZE :- Tokenize function Split a string of words
B = foreach A generate TOKENIZE(record);

Aggregate functions :-
To make the aggregate functions compulsory - make the group.

① COUNT :- calculate the no. of entries in a bag.
B = order A by id;

C = group B by age;

D = foreach C generate group, COUNT(B);

② SUM :- calculate the sum of a set of numeric values
in a single-column bag.

E = foreach C generate group, SUM(B.sal);

③ MAX :- calculate the maximum value or char-arrays in a single column bag.

F = foreach C generate group, MAX(B.sal);

④ MIN :- compute the minimum of a set of numeric values or char-arrays in a single column bag.

G = foreach C generate group, MIN(B.sal);

⑤ AVG :- It calculates the average value of entries in a bag.

H = foreach C generate group, AVG(B.age);

- ⑥ CONCAT :- Concatenates two byte arrays (or) two character arrays together.

I = foreach c generate group, CONCAT(B.name)

UDF → user defined functions.

- ① UDF's are required to define custom processing.
- ② UDF's can be written in the Java, python, JavaScript and ruby and permit pig to support custom processing.
- ③ Support for writing UDF's in python, javascript and Ruby still evolving.
- ④ UDF's provide an opportunity to extend pig into your application domain.

How to write java UDF:-

- UDF's can be developed by extending EvalFunc class and overriding exec method.

Embedded mode of pig:-

① putting one thing to another thing is called embedded mode.

UDF String Convert to Uppercase Letters ?

b7

```
Package myudf;
import java.io.IOException;
import org.apache.pig.EvalFunc;
import org.apache.pig.data.Tuple;
import org.apache.pig.impl.util.WrappedIOException;
public class UPPER Extends EvalFunc<String>
{
    public String exec(Tuple input) throws IOException
    {
        if (input == null || input.size() == 0)
            return null;
        try
        {
            String str = (String) input.get(0);
            return str.toUpperCase();
        }
        catch (Exception e)
        {
            throw wrappedIOException.wrap("caught exception process
                input row", e);
        }
    }
}
```

Execution :- if config @

192.168.144.129

→ nano embed script

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Ashiya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6462 3739, Mob: 570 6789

Downloads

myudf.jar

REGISTER

myudfs.jar; → Jar name.

A = load 'input.txt' using PigStorage('t') as (id:int,
name:chararray, marks:int);

B = foreach A generate myudfs.UPPER(name);
 ↓ package name ↓ class name.

dump B;

Run!

Pig -x local embedscript pig
input.txt

Hive

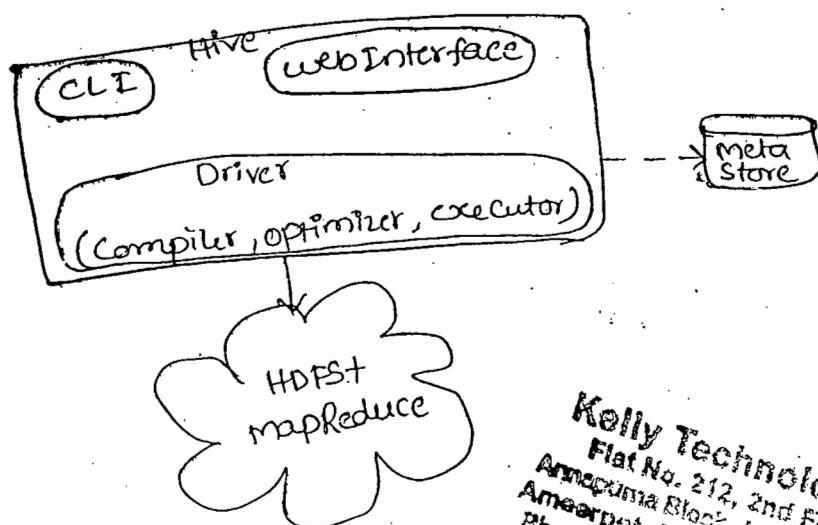
b6

- ① Hive is one of the components of the Hadoop. Hive was initially introduced at facebook in order to fulfill requirements with respect to ETL tools.

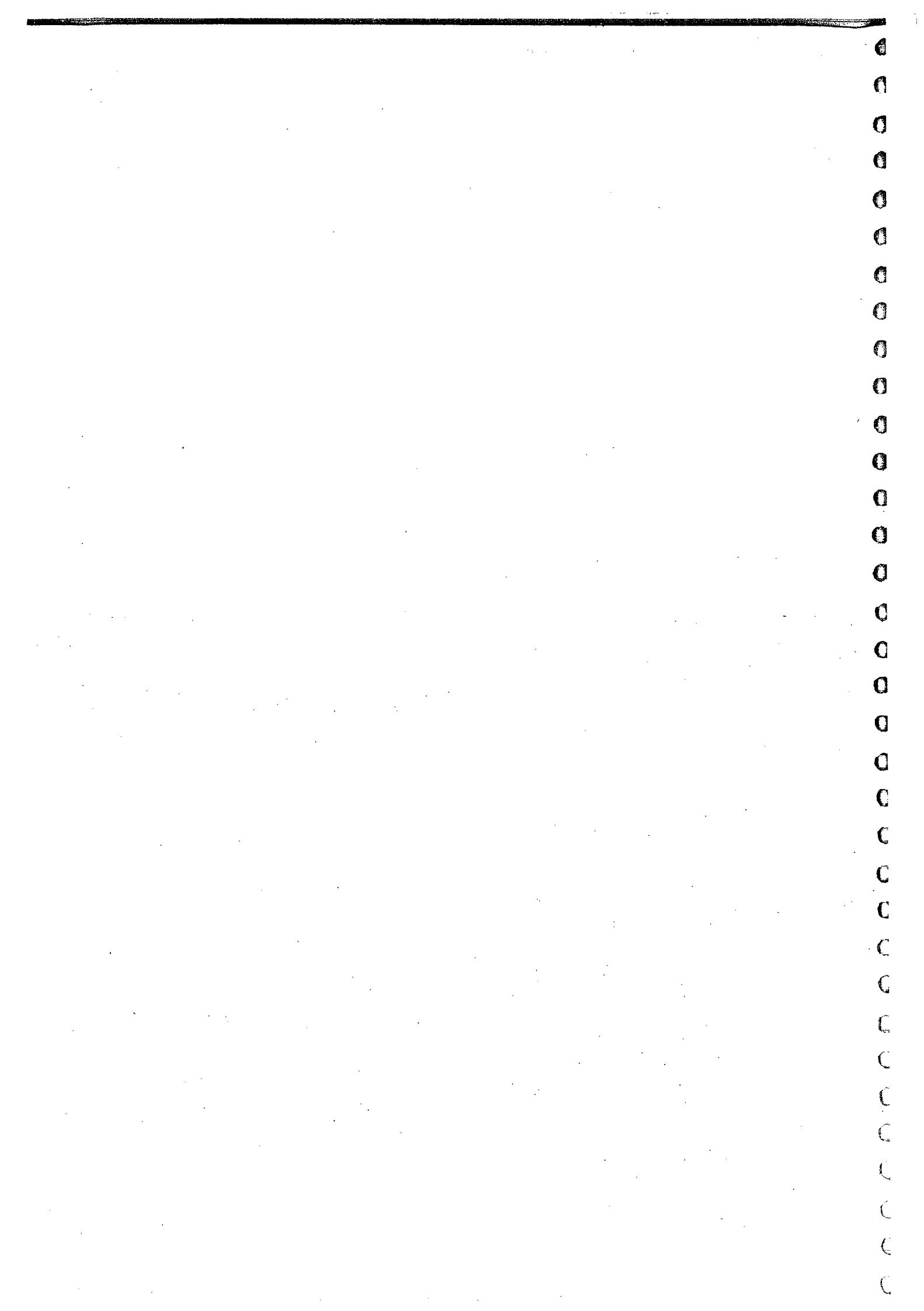
Extract Transformation Load.

Hive is a dataware house system for hadoop that facilitates easy data summarization, adhoc queries, and the analysis of large datasets stored in Hadoop compatible file systems.

Hive provides a mechanism to project structure onto this data and query the mechanisms data using a SQL like language called HiveQL. At the same time this language also allows traditional map/reduce programmers to plug in their custom mappers and reducers when it is inconvenient or ~~inconvenient~~ to express this logic in HiveQL.



Kelly Technologies
Flat No. 212, 2nd Floor,
Amarapura Block, Ashoknagar Enclave,
Amarpet, Hyderabad - 500 016.
Ph: 040-64677777, 040-6788



Hive is an abstraction on top of mapreduce. It allows users to query data in the Hadoop cluster without knowing Java or mapreduce. It uses HiveQL language, very similar to SQL.

About Hive:-

- Hive was originally developed at Facebook.
- provides a very SQL-like language.
- can be used by people who know SQL.
- Enabling Hive requires almost no extra work by the system administrator.
- Hive "layers" table definitions on top of data in HDFS.
- Hive tables are stored in Hive's "warehouse" directory: HDFS, By default /user/hive/warehouse.
- Tables are stored in subdirectories of the warehouse directory.
- Actual data is stored in ~~by~~ ⁱⁿ ~~Hive~~ - control character - de delimited text, or sequence files.

Hive is data warehousing tool on top of Hadoop. It is same as

SQL:-

- SQL like queries
- show, tables, describe, drop table
- create table, alter table
- select, insert

Hive Limitations:-

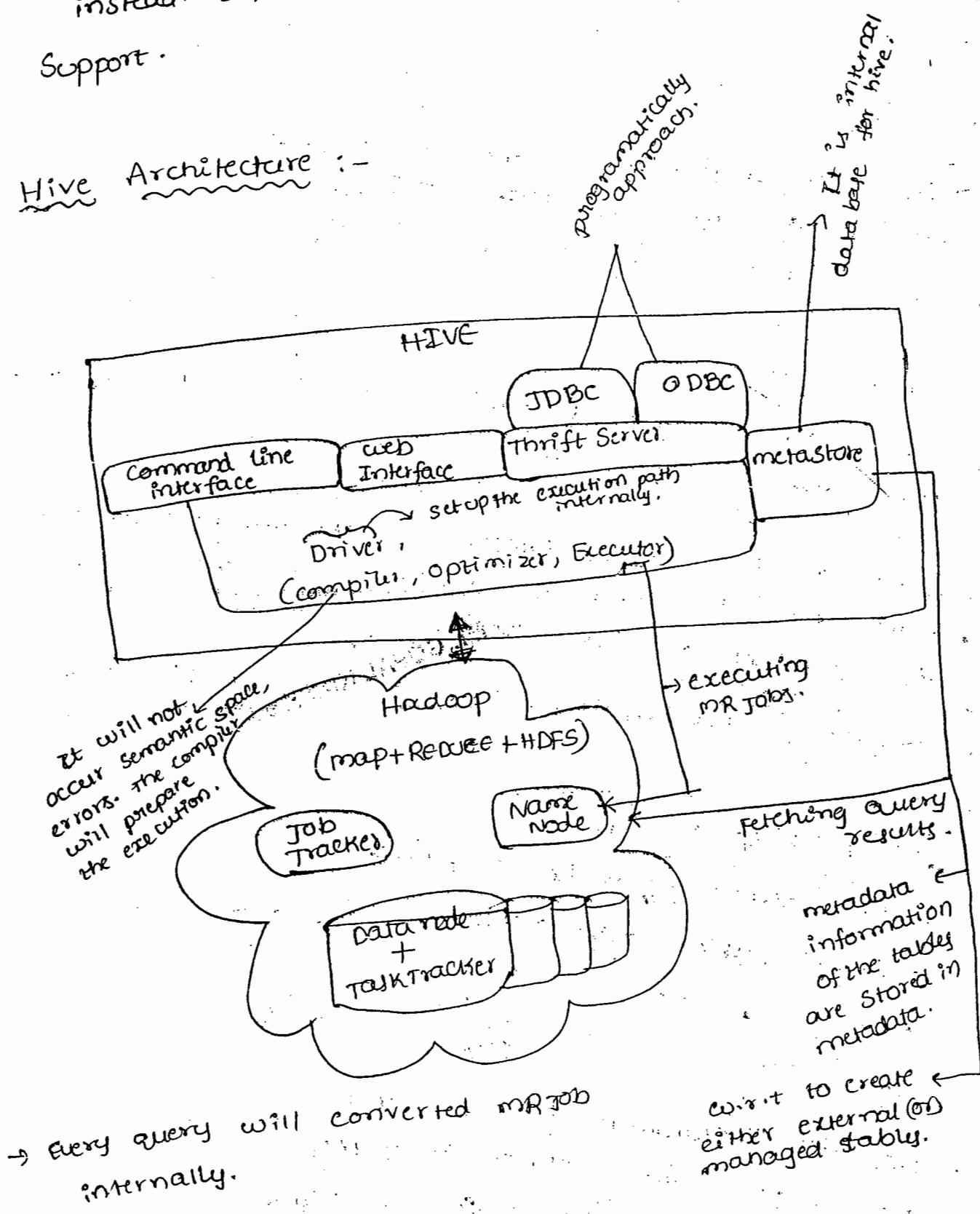
- Not all "standard" SQL is supported.
- No support for UPDATE or DELETE.
- No support for INSERTING single rows

Flat No. 212, 2nd Floor,
Amarpet, Hyderabad-500 016
Ph: 040-6462 6789, 998 570 6789

- Relatively limited no. of built-in functions.
- No data types for date or time - use the STRING datatype instead. In new version date or time datatype will

Support.

Hive Architecture :-



- metastore :- stores system catalog.
- Driver :- manages life cycle of HiveQL query as it moves thru 'HIVE'; also manages session handle and session statistics.
- query compiler :- compiles HiveQL into a directed acyclic graph of map/reduce tasks.
- Execution engine :- the component executes the tasks in proper dependency order; interacts with hadoop.
- Hive Server :- provides thrift interface and JDBC/ODBC for integrating other applications.
- client components :- CLI, webInterface, JDBC/odbc interface
- extensibility interface include SerDe, user defined function and user defined aggregate function.

By Mr. GOPAL KRISHNA

Hive CLI

- DDL (Data Definition Language)
 - create table / drop table / rename table.
 - alter table add column.
- Browsing
 - show tables
 - describe table
 - cat table
- Loading data
- Querying

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Building,
Ameerpet, Hyderabad,
Telangana, India
Pin: 410 013

web UI for Hive :-

- metastore UI:

- Browse and navigate all tables in the system.
- Comment on each table and each column.
- Also captures data dependencies.

why hive ?

- Need a multi petabyte warehouse.
- Files are insufficient data abstractions.
- Need tables, schemas, partitions, indices.
- SQL is highly popular.
- Need for an open data format
 - RDBMS have a closed data format
 - flexible schema
- Hive is a Hadoop subproject!

Hive used for types of applications :-

- Summarization
 - Ex:- Daily/weekly aggregations of impression/click counts.
 - complex measure of user engagement.
- adhoc analysis
 - Ex:- how many group admins broken down by state/country.
- Data mining (Assembling Training Data)
 - Ex:- user engagement as a function of user attributes.

- Spans Detection:-

- anomalous patterns for site integrity.
- Application API usage patterns.

- Ad optimizations

- Document indexing

- Customer facing business intelligence (e.g. Google analytics)
- predictive modeling, hypothesis testing.

Hive Components :-

- Shell :- allows interactive queries like mysql shell connected to database - also supports web and JDBC clients.
- Driver :- session handles, fetch, execute.
- Compiler :- parse, plan, optimize
 - when Driver invokes compiler with HiveQL, the compiler converts string into a plan.

- plan can be

- metadata inference operation for ~~LOAD~~ statement.
- HDFS operation for LOAD statement.
- For insert / queries consists of DAG (Directed Acyclic Graph) of mapReduce (mr) jobs.

- parser transform query into a parse tree representation

- semantic analyzer based internal query representation - ~~retrive~~ **Kenny Technologies**

Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-64626789, 9685706789. verification of the

input table from metastore, verify the
column names, expand select * and does type
including implicit type conversions.

- checking

- physical plan Generator converts logical plan into physical plan consisting of DAG of map-reduce jobs.
- logical plan Generator Converts internal query representation to a logical plan consists of a tree of logical operators.
- optimizer perform multiple passes over logical plan and resulted in several ways.
 - Combine multiple joins which share the join key into a single multi-way JOIN \rightarrow a single mapReduce job.
 - Prune columns early and pushes predicates closer to the table scan operator to minimize data transfer.
 - Prunes unneeded partitions by query.
 - For sampling query - Prunes unneeded bucket.

• Execution Engine :- DAG of stages(MR, HDFS or metadata).

• MetaStore :- MetaStore is the central repository of Hive metadata.

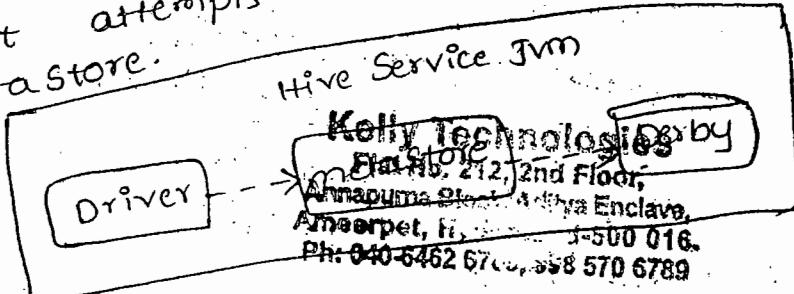
The metastore is divided into two parts.

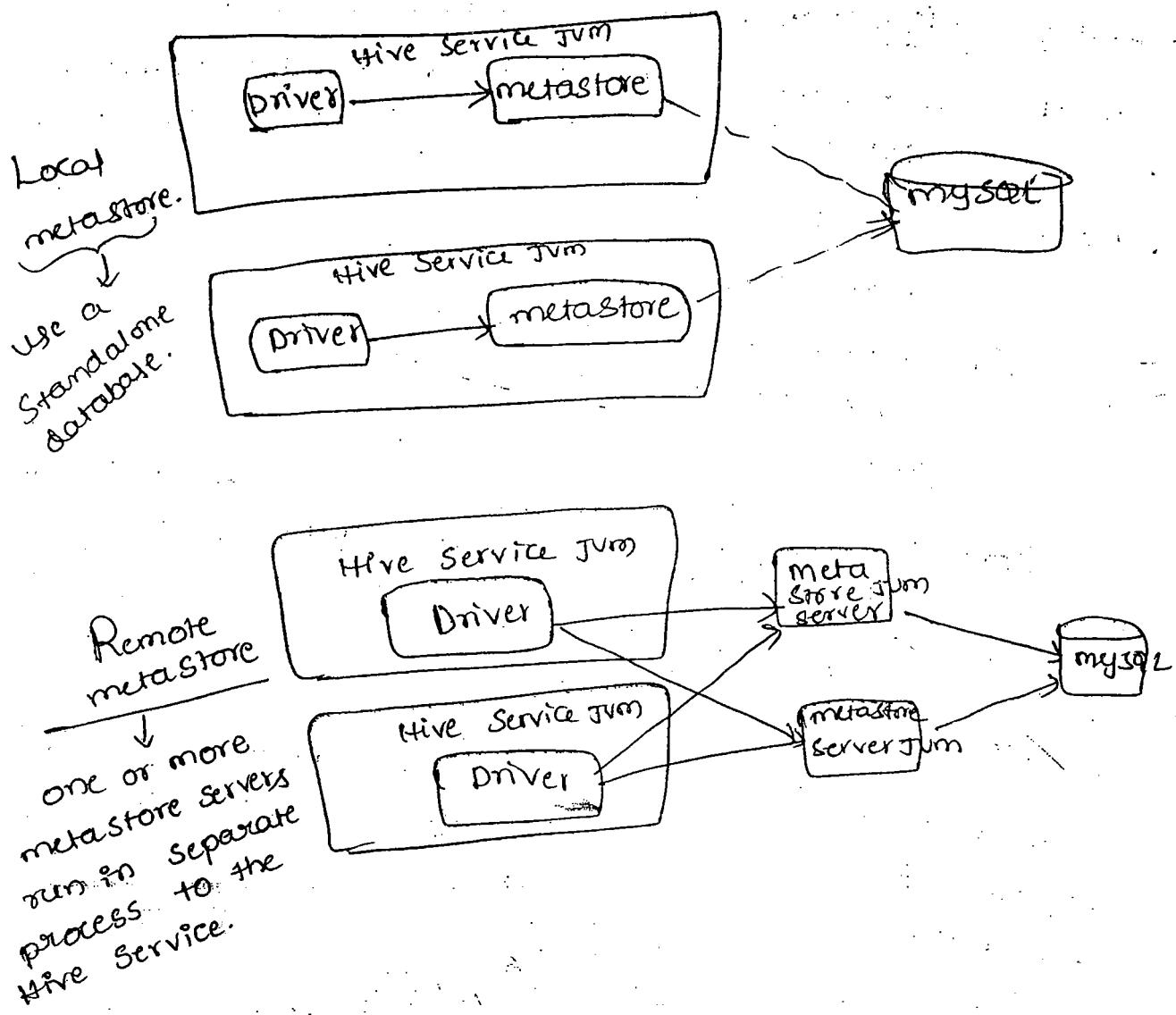
- ① Service
- ② backing store for the data.

By default, metastore service runs in the same Jvm as the Hive Service and it contains an embedded Derby database instance backed by the local disk. This is called the Embedded metastore configuration.

metastore configuration :-

- all the metadata information of the hive table will reside in internal database of hive So called metastore.
- metastore will provide below information.
 - (1) schema information of the table
 - (2) column information of the table
 - (3) Table locations
 - (4) partition keys - etc.
- metastore is a keyrole in hive architecture. metastore is not running, hive process the data is not possible without existing of metastore concepts it is not possible process the data in hive.
- By default hive comes with derby database as its metastore. however we can very well configure other database(mysql, thrift Server - etc) as the metastore of hive.
- the metastore configuration file is "hive-site.xml"
- Failed to start database "metastore-db" :-
 - when it attempts to open a connection to the metastore.





Hive Database :-

① Data Model

① Tables :-

- analogous to tables in relational database
- each table has a corresponding HDFS Dir
- Data is serialized and stored in files within Dir
- supports external tables on data stored in HDFS, NFS or local directory.
- typed columns (int, float, string, date, boolean) - also list, map (for JSON-like data).

② partitions

① partitions :-

- analogous to dense indexes on partition column.
- nested subdirectories in HDFS for each combination of partition column values.
- allows users to efficiently retrieve rows.
- A table can have one or more partitions (1-64) which determine the distribution of data within subdirectories of table directory.

Ex:- Table T under /whl/T and is partitioned
column on ds + chry.

for ds = 20090101

chry = US.

Then data is stored within dir.

/whl/T/ds = 20090101/chry = US.

By Mr. GOPAL KRISHNA

② Buckets :-

- split data based on hash function of some column of a table - mainly for parallelism.
- Data in each partition may in turn be divided into buckets based on the value of a hash function of some column of a table.

Example :-

- Bucket
- HDFS file for user hash bucket 20
- HDFS file for user hash bucket 21

Katty Technologies
Plot No. 212, 2nd Floor,
Amarpura Block, Kalyanpur,
Ph. 022 2466 5555, 022 2466 5556

Antimata Block, Kalyanpur,
Amarpura, Ph. 022 2466 5555, 022 2466 5556

500 676

500 676

500 676

500 676

500 676

500 676

500 676

500 676

Hive physical layout :-

- warehouse directory in HDFS
Ex:- /home/hive/warehouse.
- Tables stored in subdirectories of warehouse - Partitions, buckets from subdirectories of tables.
- Actual data stored in flat files - Control char-delimited text, or sequence files - with custom SerDe, can use arbitrary format.

Hive datatype :-

- ① Support primitive column types:
 - Integers: TINYINT, SMALLINT, INT, BIGINT
 - Boolean: BOOLEAN
 - Floating point numbers: FLOAT, DOUBLE
 - String : STRING
 - Date : ~~DATE~~ String
 - Nested collections :- (complex type)
 - Maps : M ['group']
 - Arrays : ['a', 'b', 'c'], A[1] returns 'b'.
- ② user defined data types:
 - user can also define their own type programmatically.
 - Structures with attributes of any type.
 - Structs : {a INT, b INT}

structured data :-

① Type System

- primitive types
- Recursively build up using composition/maps/lists.

- ① Object Inspector interface for user-defined data types:-
- TO recursively list schema
- TO recursively access fields within a row object.

② Generic (De)Serialization Interface SerDe :- (serialization/deserialization)

③ Serialization families implement interfaces:-

- Thrift DDL base SerDe
- Delimited text based SerDe
- you can write your own SerDe (JSON, XML ---)

Serialization / Deserialization

• Generic (De) serialization Interface SerDe.

• use LazySerde

• flexible Interface to translate unstructured data into structured data.

• Designed to read data separated by different delimited characters.

• These SerDes are located *By Mr. GOPAL KRISHNA* in `hive-contrib.jar`.

Hive File Formats

Hive stores different file formats.

- Hive lets users store different file formats.
- Helps in performance improvement.

SQL Example

Create table dept (key INT, value STRING) STORED AS SequenceFileInputFormat

INPUTFORMAT 'org.apache.hadoop.mapred.SequenceFileInputFormat'

OUTPUTFORMAT 'org.apache.hadoop.mapred.SequenceFileOutputFormat'

HiveQL :-

- Support SQL-like query language called HiveQL for select, join, aggregate, union all and subquery in the from clause.
- Support DDL statements such as CREATE table with serialization format, partitioning and bucketing columns.
- Command to load data from external sources and INSERT into HIVE tables.
- Do not support UPDATE and DELETE.
- Support multitable INSERT
- Support user-defined column transformation (UDF) and aggregation (UDAF) function written in Java.

Hive vs pig

Hive

- ① Language is SQL-like
- ② Schema: Table definitions that are stored in a metastore
- ③ Hive programmatical access is JDBC, ODBC
- ④ In the hive have partitions
- ⑤ Server is optional (Thrift)
- ⑥ UDFs (Java)
- ⑦ Custom Serializer/Deserializers

pig

- ① Language is pig latin.
- ② A Schema is optionally defined at runtime.
- ③ pig access: is pigServer
- ④ There's no partitions
- ⑤ no server
- ⑥ UDFs (Java)
- ⑦ Custom Serializer/Deserializer

- ⑧ DFS direct access at run time
 - ⑨ join/order/SET is possible
 - ⑩ Shell command interface is possible
 - ⑪ Streaming is supported
 - ⑫ webInterface is possible
- The entire data of hive will be residing in two types of tables.

- ⑧ DFS direct access at default.
- ⑨ join/order/SET is possible

- ⑩ Shell command interface is possible

- ⑪ Streaming is supported

- ⑫ There is no web interface.

→ The entire data of hive will be residing in two types of tables.

① managed tables

② External tables

By Mr. Gopal

① managed tables :- managed tables ~~are~~ the hive default tables i.e. these tables will be completely managed by hive warehouse only.

All the managed tables will reside under the subfold of the default hive directory is known as

/user/hive/warehouse **Kelly Technologies**

Flat No. 212, 2nd Floor,

Annapurna Block, Aditya Enclave

Ameerpet, Hyderabad - 500016

Ph: 040-64621166

Mobile: 9840629999

② External tables :- As the name suggests, these tables are hive external warehouse, i.e. whatever the external tables that we are creating will reside in an external path will be hive warehouse directory.

Syn: create external table <Table-name>

(-----) Keyword - it is treated as external table.

(-----) otherwise if it is not used that keyword it is treated as managed table.

) Location <<HDFS path>

→ what is the difference between Managed Tables & External tables.

Managed Table

- ① point to existing data directories in default hive warehouse.
- ② can not create partitions.
- ③ data is assumed to be in hive compatible format.
- ④ Dropping managed table, drops both metadata and data.
- ⑤ managed tables are stored in /user/hive/warehouse,
hive default path.

∴ Create table empid int,
name string;

hive>
create table emp(cmpid int,
ename string,
esal double);

rowformat delimited fields
terminated by '\t'. stored
as textfile ↪

External Table

- ① point to existing data directories in HDFS
- ② can create table and partitions
- ③ data is assumed to be in hive compatible format.
- ④ Dropping external table, drops the only metadata.
metadata → no. of columns,
types of columns,
terminators.
- ⑤ external tables are stored in external path of the
hive warehouse directory.

Eg: hive > create

create ~~table~~ External table emp
(empid int,
ename string,
esal double) location

'/suitablelocation';

It will create by
default by hive.

Hive Commands :-

76

To launch the Hive shell, start a terminal and run `bin/hive`.

Note: example is the table name for all query.

`hive >`

Hive : Creating Tables

- ① Create table : example(`id INT, name STRING`) Row form
delimited fields terminated by '`\t`' stored as textfile

→ `describe` example;

→ `show tables;`

Hive : loading Data into Hive :-

Data is loaded into Hive with ~~By Mr. GOPAL KARISHMA~~ LOAD DATA INPATH

Statement - Assumes that the data is already in HDFS.

`hive > load data inpath 'file-txtdata.txt' into table example;`

If the data is on the localfile system, use `load data local inpath`

automatically loads it into HDFS.

`hive > load data local inpath "file-txtdata.txt" into table example;`

Hive : Select queries :-

Hive supports most fam

`hive > select * from example`

`hive > select * from example where id > 100 ORDER BY name`

ASC limit 10;

JD

Kelly Technologies

Flat No. 212, 2nd Floor,

Aeroplane Block, Ameerpet, H

Pb: 040-6465706789

Enclave, 500 016.

joining Tables :-

```
SELECT e.name, e.dep., s.id FROM example e JOIN sample s  
ON (e.dep = s.dep) WHERE e.id >= 20;
```

Creating user-defined functions :-

Insert over write table U-data-new

Select transform(userid, movieid, rating, unixtime) using

'`movieLensMapper.py`' AS (`userid`, `movieid`, `rating`,

'python weekday - map[
weekday] From u-data'

weekday) from U-data.

join Query : sample

Select A # from example A join(Select id, max(start-time) as

Start time from example B.

where start-time < 25 group by id) MAXSR
MAXSR-STAC

ON A.id = MAXsp.id and A.start-time = MAXsp.start-time;

using NOT IN / IN hive query :-

Select * FROM example where id).

Select * FROM example

After tables in Hive :-

After tables in HIVE
use table example rename to new example;

hive > alter table example ADD COLUMNS (doj string, emplike int);
hive > alter table example add column empname string;

partition :-

partitions :-
we > create table parttab(logid int, logerror string)
PARTITIONED BY (logdate string, columns string) row
delimited fields terminated by '\t' zones

partitions :-
hive > create table parttab(logid int, logerror string)
PARTITIONED BY (logdate string, columns string) row
delimited fields terminated by '\t' using
textfile.

partitions :-
hive > create table parttab(logid int, logerror string)
PARTITIONED BY (logdate string, columns string) ROW
delimited fields terminated by '\t' stored
as textfile.

→ hive supports 3 file formats.

- ① Textfile
- ② Sequence file
- ③ Record character file...

11

→ CTAS :- (create table as select)

hive > create table CTASTAB as select id, name, sal from
new managed tab;

~~hiveso~~

→ JOIN :- create table merge as select n.id, e.name,
e.sal from new managed tab n JOIN empmanaged e ON
(n.id = e.id);

→ LEFT OUTER JOIN :-

hive > create table merge as select n.id, e.name, e.sal from
newmanagedtab n LEFT OUTER JOIN empnewmanagedtab
e on (n.id = e.id).

By Mr. GOPAL KRISHNA

→ partitions :- schema information is same.

These columns not at all be input file.

partition column information provided at run time or

Distribute by :-

① Distribute by controls how output will be divided

among Reducers.

Flat No. 212, 2nd Floor,

Annapurna Block, Alpha Enclave

Amitabh Bhattacharya

DISTRIBUTE BY ensures that records of ename will go to the same reducer and then to SORT the same data the way we want (ascending order of empld).

③ Distribute by works similar to GroupBy in the sense that how it controls the reducers to receive the rows for processing.

Note :- Hive requires Distribute by clause comes Before the Sort by if we are using both in a query.

Hive > select * from disttab;

> Select empid, ename, esal from disttab DISTRIBUTE BY ename SORT BY ename asc, esal asc;

CLUSTER BY :-

① CLUSTER BY clause is a short-hand way of expressing the DISTRIBUTE BY with SORT BY.

Ex:- Hive > select empid, ename, esal from disttab CLUSTER BY ename;

ORDER BY :-

① The ORDER BY clause is familiar from other SQL dialects.

② It performs the 'total' ordering of the query result set' i.e "all the data will be passed to a single reducer" which will take longer time if we have larger datasets.

③ Hive requires the clause for ORDER BY as it will take longer time (if hive property hive.mapred.mode is set to 'strict' which is by default 'nonstrict').

SORT BY :-

① where as SORT BY ("Hive's alternative to ORDER BY"), will orders the data only within each reducer and performing the 'local' ordering of the data! where each REDUCER'S output will be SORTED.

② Better performance is traded for total ordering.

Ex:- Select empid, ename, esal from testudf order by empid asc, esal desc;

Ex:- Select empid, ename, esal from testudf sort by empid asc, esal desc;

Indexes in HIVE:-

① > create index parttab_index on table parttab(logid) as 'org.apache.hadoop.hive.al.index.compact.CompactIndexHandler' with DEFERRED REBUILD;

② create table tt (i int, jint);

create index * on Table tt(i) As 'org.apache.hadoop.hive.al.index.compact.CompactIndexHandler' with DEFERRED REBUILD;

→ UDF

com.hive.gopal;
import org.apache.hadoop.hive.al.exec.UDAF;
import org.apache.hadoop.hive.al.exec.UDAF.Evaluator;
import org.apache.hadoop.io.IntWritable;
public class Maximum extends UDAF{
 public static class MaximumInt UDAFEvaluator implements

private IntWritable result;

public void init(){
 result = null;

}
public boolean iterate(IntWritable value){
 if (value == null)
 return true;

}

By Mr. GOPAL KRISHNA

Kally Technologies
Flat No. 210, 2nd Floor,
Annapurna Apartments,
Amarapuri, Hyderabad,

```
if (result == null) {  
    result = new IntWritable (value.get());  
}  
else {  
    result.set (Math.max(result.get(), value.get()));  
}
```

Execution :-

```
hive > add jar maximum.jar  
> create temporary function maximum as  
    com.hive.gopal.Maximum;
```

sqoop create Hive-table :-

```
# sqoop create emptable --connect jdbc:mysql://localhost/  
Gopal --table student --emp-tables GopalStudent;
```

sqoop hive-import :-

```
# sqoop import --connect jdbc:mysql://localhost/Gopal  
--table student -hive-import;
```

Kelly Technologies

Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 988 570 6789. HBase

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789.

79

what is HBase?

HBase is a database: the Hadoop database. It is indexed by rowKey, columnKey and timestamp. HBase stores structured and semistructured data naturally so you can load it with tweets and parsed log files and a catalog of all your products right along with their customer reviews. It can store unstructured data too, as long as it's not too large.

HBase is designed to run on a cluster of computers instead of single computer. The cluster can be built using commodity hardware. HBase scales horizontally as you add more machines to the cluster. Each node in the cluster provides a bit of storage, a bit of cache, and a bit of computation as well. This makes HBase incredibly and forgiving. No node is unique, if one of those machines breaks down, you simply replace it with another. This adds upto a powerful, scalable approach to data that, until now, has not been commonly available to more mortals.

BY MR. GOPAL KRISHNA
Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 988 570 6789.

How can HBase help?

- ① Replicate dimension tables from transactional databases with low latency and without sharding.
(Fact data can stay in Hive since it is append only)
- ② only move changed rows
- "full scrape" is too slow does not scale as data keeps growing.

- Hive by itself is not good at row-level operations.
- ③ Integrate into Hive's map/reduce query execution plans for parallel distributed processing.
- ④ multiversioning for snapshot consistency?

why HBase ?

- ① we own infrastructure, no usage limits.
 - ② Data model
 - Semistructured data in HBase (easily handles multiple types in same table)
 - Time-Series ordered.
 - scaling is built-in (just add more servers)
 - But extra indexing is DIY.
 - ③ very active developer community.
 - ④ Established, mature project (in relative terms!)
 - ⑤ matched our own toolset (Java/Linux based)
- HBase: (HDFS + random read/write)

- ⑥ clone of Google's Bigtable
 - Distributed (automatic partitioning)
 - column-oriented.
 - Semi-structured (columns can be added just by inserting)
 - Built-in versioning.

- ⑦ not an RDBMS

- NO joins

- NO SQL

- Data usually not normalized

- Transactions & built-in Secondary indexes available (as contrib) but immature.

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Annapurna Block, Indya Enclave,
 Ameerpet, Hyderabad - 500 016.
 Ph: 040-6461 570 6789

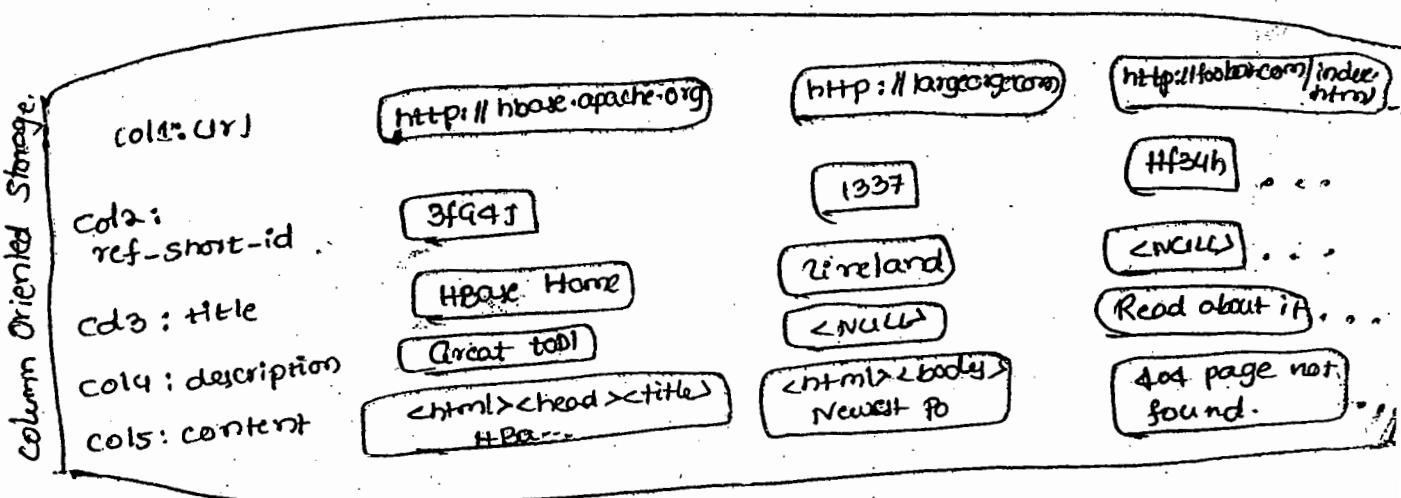
- need to think differently now about how you structure data.
 - denormalize your data where necessary.
 - structure data & row keys around common access.

HBase Data Storage:-

* HBase is column-oriented data storage.

* column-oriented means

- The reason to store values on a per column basis instead is based on the assumption.
- that for specific queries, not all of the values are needed.
- Reduced I/O



* column oriented databases save their data grouped by columns. Subsequent column values are stored contiguously on disk. This differs from the usual row-oriented approach of traditional databases, which store entire rows contiguously.

- The reason to store values on a per column basis instead is based on the assumption that, for specific queries, not all of the values are needed. This is often the case in analytical databases in particular, and therefore they are good candidates for this different storage schema.
- Reduced I/O is one of the primary reasons for this new layout, but it offers additional advantages playing into the same category: since the values of all one column are often very similar in nature or even vary only slightly between logical rows, they are often much better suited for compression than the heterogeneous record structure; most compression algorithms only look at a finite window. Specialized algorithms - for example, delta and/or prefix compression - selected based.

Row-oriented Storage

Row#	URL	HTTP Status	Content Length	Content Type	Content
Row1:	1 http://hibate.apache.org	304	1	HTTP/1.1	<html><head><title>Home</title>
Row2:	2 http://largeorg.com	1337	1	text/html	<html><body>Newest post
Row3:	3 http://foobaz.com/index.html	1434	1	text/html	<html><head>Read about it</head><body>404 page not found.

HBase Data model :-

- Tables are sorted by Row
- Table schema only define it's column families
- Each family consists of any no. of columns.
- Each column consists of any no. of versions.
- Columns only exist when inserted, Null are free.
- Columns within a family are sorted and sorted together.
- Everything except table names are byte[]
- (Row, family: Column, Timestamp) → value.

Rowkey	columnKey	Timestamp	cell
com.cnn.www	anchor:cnnst.com	T9	CNN
com.cnn.www	Anchor:mylook.ca	T8	CNN.com
com.cnn.www	Contents:	T6	<html>...
com.cnn.www	Contents:	T5	<html>...
com.cnn.www	Contents:	T3	<html>..

By Mr. GOPAL KRISHNA

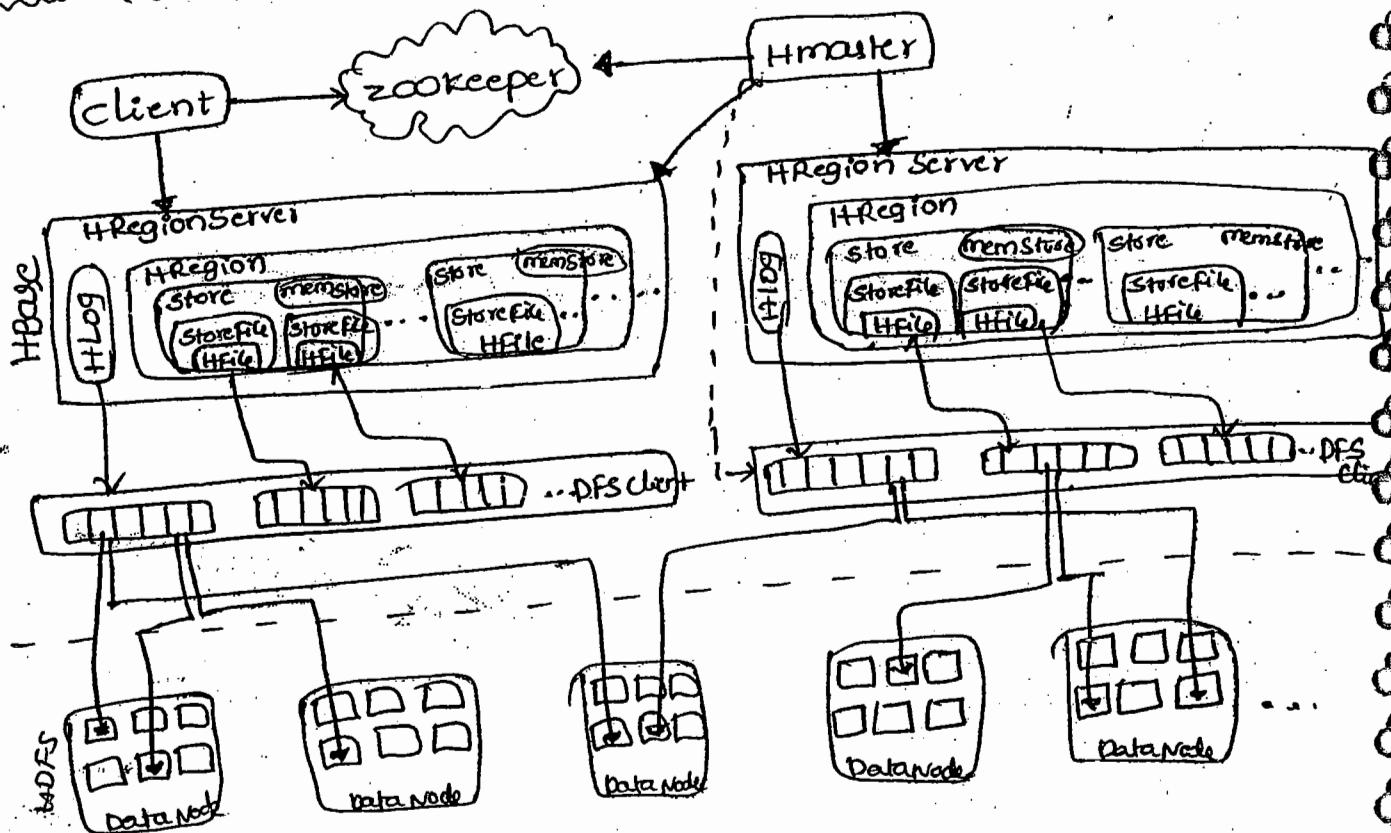
Column Families

- Different sets of columns may have different properties and access patterns.
- Configurable by column family.
 - compression (none, gzip, LZO)
 - Version retention policies
 - Cache priority.
- CF's stored separately on disk. access one without waiting to on the other.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, City Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6462 07
570 6789

- Row's columns are grouped into families.
- column family members identified by a common 'printable' prefix.
- column family should be predefined.
 - but column family members can be added dynamically
 - member name can be bytes.

HBase Architecture :-



Components

- HBase Master
- HRegion Server
- Client

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6452 570 6789

① HMaster :-

- master server is responsible for monitoring all RegionServer instances in the cluster, and is the interface for all metadata changes, it runs on the server which hosts namenode.
- master controls critical functions such as RegionServer failover and completing region splits. so while the cluster can still run for a time without the Master, the master should be restarted as soon as possible.
- master does not handle write any request.
(not a DB master)
- Does not handle location finding requests.
- not involved in the read/write path.
- Even master(s) is (are) down, cluster can respond to write/read request.
- Generally does very little most of the time.

master is stateless

all the data & state info stored in HDFS & co-processor.

Master is not SPOF!

② HRegionServer :-

Table : HBase Table

Region : Regions for the table.

store : store per column family for each region for the table.

MemStore : MemStore for each store for region for the table.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6452 6789, 998 570 6789

By Mr. GOPAL KRISHNA

Store file : Store files for each state for each region for the table.

Block : Blocks within a store file within a store for each region for the table.

Hlog used for recovering.

- Send heartbeat (loadinfo) to master
- write requests handle
- read request handle.
- flush
- compaction
- Region splits. (mange)

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-5452 6789, 998 570 6789

Regions :-

- ① horizontal partitioning
- ② Every region has a subset of the table's rows.
- ③ Region identified as
[- table, first row(+), last row(-)]
- ④ Table starts on a single region
- ⑤ split into two equal sized regions as the original region grows bigger and so on..

ZooKeeper :-

- master election and Server availability.
- cluster Management
- assignment transaction state Management.
 - contacts zookeeper to bootstrap connection
- client to the HBase cluster.
- Region key ranges, region Server address.
- Guarantees consistency of data across clients.

HBase Client

- The HBase client is responsible for finding RegionServers that are serving the particular row range of interest.
- It does this by querying the .META. and -Root- catalog tables in ZooKeeper.
- After locating the required region(s), the client directly contacts the RegionServer serving that region.

HBase Tables and Regions:-

- HBase table is made up of roughly equal sized regions.
- Each region may live on different node and is made up of several HDFS files and blocks, each of which is replicated by Hadoop.
- Region is specified by its startkey and end key.

Accessing HBase

- Java API (thick client)
- REST/HTTP
- Apache Thrift (any language)
- Hive/pig for analytics

HBase API

- get(row)
- put(row, Map<column, value>)
- scan(keyrange, filter)

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Annapurna Block, Aditya Enclave,
 Ameerpet, Hyderabad-500 010.
 Ph: 040-6452 6789, 898 570 6739

By Mr. GOPAL KRISHNA

- increment (row, columns)
- ... (checkAndPut, delete - etc)
- MapReduce/Hive

Kelly Technologies
 Flat No. 212, 2nd Floor,
 Annapurna Block, Aditya Enclave,
 Ameerpet, Hyderabad-500 016.
 Ph: 040-6452 6789, 998 570 6789

HBase vs RDBMS

RDBMS	HBase
① row oriented.	① column-family oriented
② multi-row ACID	② single row only.
③ SQL	③ get/put/scan/etc *
④ Authentication/Authorization on arbitrary columns	④ work in progress.
⑤ TB's	⑤ Row-key only.
⑥ 1000s queries/second	⑥ ~1 PB
	⑦ millions of queries/ second

HBase vs NOSQL

- Favors consistency over availability (but availability is good in practice).
- Great Hadoop integration (very efficient bulk loads, mapReduce analysis)
- Ordered range partitions (not hash)
- Automatically shards | scales (just run on more servers)
- sparse column storage (not key-value)

HBase vs HDFS

<u>HDFS</u>	<u>HBASE</u>
append-only	Random write, bulk incremental.
full table scan, partition table scan	Random read, small range scans, or table scan. 4-5 x slower.
higher performance very good in hive	sparse-column-family data model.
structured storage (sequence file)	~ 1 PB.
max data size (30+PB)	

Client side Buffering

in HBase :-

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016,
Ph: 040 3462 6789, 988 570 6789

HBase:

HBase is an open source, non-relational, distributed database modeled after Google's. It runs on top of Hadoop distributed file system, providing big table like capabilities for Hadoop. Hbase provides a way to store large amounts of data in a fault-tolerant way. Hbase can be also termed as the hadoop database. we can perform operations like insertion into Hbase using put method. put operation - is like an Rpc that transfers data from client to server and back. when the no. of insertions are less, it is feasible but when we need to have like 1000 insertions per second, it's better that we go for buffering.

By Mr. GOPAL KRISHNA

Client Side Buffering:-

HBase API has a built-in client side buffer that can buffer all the puts and send it in a single RPC call. The following methods are used for this purpose.

```
void setAutoFlush(boolean autoFlush)  
    boolean isAutoFlush()
```

We need to set the 'autoFlush' to 'false' to enable the client side buffer as by default it's not enabled.

```
table.setAutoFlush(false)
```

Once this has been done, we can check the state of the 'autoFlush' flag using the 'isAutoFlush()' method. Now, all the put instances we created are stored in the memory in the client process. To force the data out of the buffer, we need to call another function provided by API.

```
void flush()  
void flushCommits() throws IOException
```

This method flushes the data from buffer and ships it to the servers. Client batches the puts and sends them to the appropriate region servers in a single RPC call. You also have the option of setting the bufferSize (client-side write buffer).

In this case, API will track the size of data buffered and the required heap size for each instance. So it tracks the total overhead of the data. When it goes beyond a certain limit, it will implicitly flush the data. We can configure the max client side bufferSize using the following.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarpet, Hyderabad-500 016.
Ph: 040-6452 6789, 988 570 6789

void setWriteBufferSize (long writeBufferSize) throws IOException;

you can retrieve the client side bufferSize using getWriteBufferSize() method in the API. the default size is 2MB.

the data is hence flushed in 2 cases:

(a) Implicit Flush:-

By Mr. GOPAL KRISHNA

It's triggered whenever put() or setWriteBufferSize() is called. This will cause the API to compare the current bufferSize with the configured maximum bufferSize. if it exceeds the limit, then it triggers the flush. if the client buffer is disabled then each invocation of the put() will cause a differ buffer flush. close() also triggers buffer flush.

(b) Explicit Flush:-

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Atliya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 98 570 6789

the flush commits() is called explicitly to clear the buffer and sent the data to the Server. Here's an example of client side buffering,

```
HTable table = new HTable(config, "example-table");
table.setAutoflush(false); //setting autoflush to false to
                           //enable diff buffering.
Put put1 = new Put(Bytes.toBytes("row-1"));
put1.add(Bytes.toBytes("column-fam1"), Bytes.toBytes("column-
                           equal"));
put1.addColumn("value-1");
table.put(put1);
```

```

put put2 = new put(Bytes.toBytes("row-2"), put1.add(Bytes.toBytes("column-fam1"),
Bytes.toBytes("column-qual2"),
Bytes.toBytes("value-2")));
table.put(put2);

put put3 = new put(Bytes.toBytes("row-3"));
put1.add(Bytes.toBytes("column-fam3"),
Bytes.toBytes("column-qual3"),
Bytes.toBytes("value-3")));
table.put(put3);

```

Get get = new Get(Bytes.toBytes("row-1"));

Result res1 = table.get(get);

Result res1 = table.get(get); // an explicit flush is performed here.

- In this example, we are creating puts one by one and inserting into the table.
- When you need to access the write buffer content, you can use

`ArrayList<put> getWriteBuffer();`

- This will return the internal list of the put instances that have been buffered so far by invoking table.
- In the previous example, the put instances were being added one by one. We can also create a list of put

`void put(List<put> puts) throws IOException`

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Amarpet, Hyderabad-500 016.
Ph: 040-6452 6789, 943 570 6789

HBase shell mechanism :-

hbase > create 'test', 'cf'

hbase > put 'test', 'row1', 'cf:a', 'value1'

hbase > scan 'test'

Row

Column + Cell

row1

column = cf:a, timestamp = 1332311427315,
value = value1.

Java mechanism :-

```
import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.hbase.client.HTable;
import org.apache.hadoop.hbase.client.Result;
import org.apache.hadoop.hbase.client.ResultScanner;
import org.apache.hadoop.hbase.client.Scan;
import org.apache.hadoop.hbase.util.Bytes;

public class HBaseClient {
    public static void main(String[] args) throws IOException {
        Configuration config = HBaseConfiguration.create();
        config.set("hbase.zookeeper.quorum", "hadoop1");
        HTable testTable = new HTable(config, "test");
        for(int i=0; i<100; i++) {
            byte[] family = Bytes.toBytes("cf");
            byte[] qual = Bytes.toBytes("a");
            Scan scan = new Scan();
            scan.addColumn(family, qual);
            ResultScanner rs = testTable.getScanner(scan);
            for( Result r = rs.next(); r!=null; r=rs.next()) {
                byte[] valueObj = r.getValue(family, qual);
                String value = new String(valueObj);
                System.out.println(value);
            }
        }
        testTable.close();
    }
}
```

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-6462 6789, 968 570 6789

By Mr. GOPAL KRISHNA

Dependent JAR files :-

- commons-configuration-1.6.jar
- commons-lang-2.5.jar
- commons-logging-1.1.1.jar
- hadoop-core-1.0.0.jar
- hbase-0.92.1.jar
- log4j-1.2.16.jar
- slf4j-api-1.5.8.jar
- slf4j-log4j12-1.5.8.jar
- zookeeper-3.4.3.jar.

→ HBase shell commands are mainly categorized into 6 parts :-

- ① General HBase Shell Commands
- ② Tables Management Commands
- ③ Data Manipulation Commands
- ④ HBase Surgery tools
- ⑤ Cluster replication tools
- ⑥ Security tools.

D) General HBase Shell commands :-

status :- Show cluster status. Can be 'summary', 'simple', or 'detailed'. The default is 'summary'.

hbase > status
→ status 'simple'
→ status 'summary'
→ status 'detailed'.

version :- output this HBase version.

hbase > version.

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-3462 6789, 988 570 6789

who am i :- Show the current hbase user.

hbase > whoami

② Tables management Commands:-

minmay@ubuntu:~\$ su root ↵
password: ↵

root@ubuntu:~/home/gopal# jps ↵

1100 DataNode
1745 QuorumPeerMain
1702 TaskTracker
1443 NameNode
1311 JobTracker
2251 JPS
1590 Secondary NameNode

How to start HBase :-

kill 1745

cd /usr/lib/hbase/bin/ ↵

root@ubuntu:/usr/lib/hbase/bin# ./start-hbase.sh ↵

starting master, logging to

starting master, logging to

root-master-ubuntu.out ↵

root@ubuntu:/usr/lib/hbase/bin# jps ↵

1100 DataNode
2353 Master
1702 TaskTracker
1443 NameNode
2407 JPS
1311 JobTracker
1590 Secondary NameNode

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-4462 6777, 43 570 6789

By Mr. GOPAL KRISHNA

HBase Shell Commands:-

root@ubuntu: /usr/lib/hbase/bin# hbase shell ↵

hbase > version ↵

hbase > whoami ↵

hbase > status; ↵

~~hbase~~ ↵

To see the list of tables in HBase:-

hbase > list ↵

How to create a table in HBase:-

hbase > create 'Kelly', 'Hadoop'
Table name column family name.

hbase > create 'HBaseTab', 'KellyHadoop', 'GopalHadoop', 'StuHadoop'

Describe Command:-

① hbase > describe 'Kelly'

② hbase > describe 'HBaseTab'.

③ hbase > Select * from table --> equivalent.

④ hbase > scan 'Kelly'

⑤ How to enable and disable the tables

hbase > enable 'Kelly' ↵

hbase > disable 'Kelly' ↵

Altering tables:-

hbase > alter 'Kelly', NAME => 'NEWHADOOP', VERSIONS = 6

hbase > alter 'Kelly', NAME = 'NEWHADOOP', VERSIONS = 6

hbase > describe 'Kelly' ↵

Kelly Technologies
Flat No. 212, 2nd Floor,
Amarpura Block, Aditya Enclave,
Ameerpet, Hyderabad, 500 016.
Ph: 040 6482 6719, 998 570 6789

How should we know whether a table is enabled or disabled?

hbase > is-enabled 'Kelly'

> is-enabled 'HBaseTab'

IS table exists or not?

hbase > exists 'Kelly'

> exists 'HBaseTab'

> exists 'HBaseTable'

hbase > create 'table1', 'ColFam1'

> put 'table1', 'row1', 'ColFam1: colName1', 'firstValue'

> put 'table1', 'row1', 'ColFam1: colName2', 'firstValue'

> put 'table1', 'row2', 'ColFam1: colName3', 'firstValueFirst'

> put 'table1', 'row3', 'ColFam1: colName4', 'secondValueFirst'

hbase > scan 'table1'

> scan 'table1'

How to Insert (Load) the data in table having multiple Column Families (CFs) :-

Column Families (CFs) :-

hbase > create 'sampletable', 'colfam1', 'colfam2', 'colfam3'.

hbase > create 'sampletable', 'rowkey1', 'ColFam1: colName1', 'column-

value'

hbase > put 'sampletable', 'rowkey1', 'ColFam2: colName1', '

Column-1 value in ColumnFamily2'

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Iya Enclave,
Ameerpet, Hyderabad - 500 016.
Ph: 040-6452 6111, 570 6789

By Mr. GOPAL KRISHNA

hbase > Scan 'sampletable'

How to get the values from the table :-

hbase > get 'sampletable', 'rowkey3'

Count:-

hbase > count 'sampletable'

Kelly Technologies
Flat No. 212, 2nd Floor,
Annapurna Block, Aditya Enclave,
Ameerpet, Hyderabad-500 016.
Ph: 040-3452 6789, 998 570 6789

Delete:-

hbase > delete 'sampletable', 'rowkey2', 'colFamily:columnName'

that have a copy of that block and the datanodes are sorted according to their proximity to the client

- c) The DFS returns a FSDataInputStream (an input stream that supports file seeks) to the client for it to read data from.
- d) The client then calls read() on the stream
- e) DFSStream connects to the first (closest) datanode for the first block in the file.
- f) Data is streamed from the datanode back to the client, which calls read() repeatedly on the stream
- g) When the end of the block is reached, DFSInputStream will close the connection to the datanode, then find the best datanode for the next block
- h) When the client has finished reading, it calls close() on the FSDataInputStream

2) Writing Data Flow:

- a) The client creates the file by calling create() on DFS
- b) DFS makes an RPC call to the namenode to create a new file in the filesystem's namespace, with no blocks associated with it
- c) The namenode performs various checks to make sure the file doesn't already exist, and then The DFS returns aFSDataOutputStream for the client to start writing data to
- d) As the client writes data, DFSOutputStream splits it into packets, which it writes to an internal queue, called the data queue.
- e) The data queue is consumed by the DataStreamer, whose responsibility it is to ask the namenode to allocate new blocks by picking a list of suitable datanodes to store the replicas. The list of datanodes forms a pipeline
- f) The DataStreamer streams the packets to the first datanode in the pipeline, which stores the packet and forwards it to the second datanode in the pipeline. Similarly, the second datanode stores the packet and forwards it to the third (and last) datanode in the pipeline

DFSOutputStream also maintains an internal queue of packets that are waiting to be acknowledged by datanodes, called the ack queue.

A packet is removed from the ackqueue only when it has been acknowledged by all the datanodes in the pipeline

What if a datanode fails while data is being written to it.

- a) The pipeline is closed
- b) Any packets in the ackqueue are added to the front of the data queue so that datanodes that are downstream from the failed node will not miss any packets.
- c) The partial block on the failed datanode will be deleted if the failed datanode recovers later on the failed datanode is removed from the pipeline and the remainder of the block's data is written to the two good datanodes in the pipeline.
- d) The namenode notices that the block is under-replicated, and it arranges for a further replica to be created on another node

HDFS Administration

As a administrator you must know all the basic configurations, admin commands, upgrading and health check related commands.

Namenode directory structure:

After formatting namenode creates the following directory structure:

```
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/name$ ll
total 20
drwxr-xr-x 5 kelly1 kelly1 4096 2012-09-19 20:23 .
drwxr-xr-x 5 kelly1 kelly1 4096 2012-09-13 00:00 ..
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-19 20:23 current/
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-12 23:53 edits/
-rw-r--r-- 1 kelly1 kelly1 0 2012-09-19 20:23 in_use.lock
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-18 23:19 previous.checkpoint/
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/name/current$ ll
total 24
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-19 20:28 .
drwxr-xr-x 5 kelly1 kelly1 4096 2012-09-19 20:28 ..
-rw-r--r-- 1 kelly1 kelly1 4 2012-09-19 20:28 edits
-rw-r--r-- 1 kelly1 kelly1 3198 2012-09-19 20:28 fsimage
-rw-r--r-- 1 kelly1 kelly1 8 2012-09-19 20:28 fstime
-rw-r--r-- 1 kelly1 kelly1 100 2012-09-19 20:28 VERSION
```

edilog-->filesystem metadata(for every fs operation)

fsimage(periodically updated)

When namenode goes down, the latest state of its metadata can be reconstructed by loading the fsimage from disk into memory, then applying each of the operations in the edit log. This is what happens when namenode restarts.....called as SafeMode

Safe Mode:

When the namenode starts, the first thing it does is load its image file(fsimage) into memory and apply the edits from the edit log(edits). Once it has reconstructed a consistent in-memory image of the filesystem metadata, it creates a new fsimage file and an empty edit log. Only at this point does the namenode start listening for RPC and HTTP requests. However, the namenode is running in safe mode, which means that it offers only a read-only view of the filesystem to clients.

The below properties are related to hadoopsafemode:

- 1) dfs.replication.min
- 2) dfs.safemode.threshold.pct
- 3) dfs.safemode.extension

Commands for entering and leaving from safemode:

1) To check whether namenode is in safemode or not:

% hadoopdfsadmin -safemode get

Safe mode is ON

The front page of the HDFS web UI provides another indication of whether the namenode is in safe mode.

2) When you want to wait for the namenode to exit safe mode before carrying out a command, particularly in scripts.

```
% hadoopdfsadmin -safemode wait  
3) To enter safe mode, use the following command:  
% hadoopdfsadmin -safemode enter  
Safe mode is ON  
[This is useful while upgrading the cluster]  
4) To make the namenode leave safe mode:  
% hadoopdfsadmin -safemode leave
```

Safe mode is OFF
fsimage contains a serialized form of all the directory and file inodes in the filesystem. Each inode is internal representation of a file or directory's metadata, and contains such information as the file's replication level, modification and access times, access permissions, block size, and blocks a file made up of. For directories, the modification time, permissions, and quota metadata is stored.

The fsimage file does not record the datanodes on which the blocks are stored. Instead the namenode keeps this mapping in memory, which it constructs by asking the datanode for their block lists when they join cluster and periodically afterward to ensure the namenode's block mapping is up-to-date.

The edits file would grow without bound. Though this state of affairs would have no impact on the system while the namenode is running, if the namenode were restarted, it would take a long time to apply each of the operations in its edit log. During this time the filesystem would be offline, which is generally undesirable.

The solution is to run the secondary namenode, whose purpose is to produce check points of the primary's in memory filesystem metadata (This is applicable only the versions <= 0.21.0). The checkpointing process proceeds as follows.

- i) The secondary asks the primary to full its edits file, so new edits go to a new file.
- ii) The secondary retrieves fsimage and edits from the primary (using HTTP GET).
- iii) The secondary loads fsimage into memory, applies each operation from edits, then creates a new consolidated fsimage file.
- iv) The secondary sends the new fsimage back to the primary (using HTTP POST).
- v) The primary replaces the old fsimage with the new one from the secondary, and the old edits file with the new one it started in step 1. It also updates the fstime file to record the time that the checkpoint was taken.

At the end of the process, the primary has an up-to-date fsimage file, and a shorter edits file (it is not necessarily empty, as it may have received some edits while the checkpoint was being taken). It is possible for an administrator to run this process manually while the namenode is in safe mode, using the "hadoopdfsadmin -saveNamespace" command.

The secondary has similar memory requirements to the primary (since it loads the fsimage into memory), which is the reason that the secondary needs a dedicated machine on large clusters.

The schedule for checkpointing is controlled by two configuration parameters. The secondary namenode checkpoints every hour (fs.checkpoint.period in seconds) or sooner if the edit log has reached 64MB (fs.checkpoint.size in bytes), which it checks every five minutes.



Secondary namenode directory structure:

```
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/namesecondary$ ls
in_use.lock
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/namesecondary$ cd current/
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/namesecondary/current$ ls
 edits  fsimage  fstime  VERSION
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/namesecondary/current$ ll .
total 24
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-18 23:19 .
drwxr-xr-x 4 kelly1 kelly1 4096 2012-09-19 20:23 ..
-rw-r--r-- 1 kelly1 kelly1 4 2012-09-18 23:19 edits
-rw-r--r-- 1 kelly1 kelly1 2864 2012-09-18 23:19 fsimage
-rw-r--r-- 1 kelly1 kelly1 8 2012-09-18 23:19 fstime
-rw-r--r-- 1 kelly1 kelly1 100 2012-09-18 23:19 VERSION
```

The layout of secondary's current directory is identical to the namenode's. This is by design, since in the event of total namenode failure, it allows recovery from a secondary namenode. This can be achieved either by copying the relevant storage directory to a new namenode, or, if the secondary is taking over as the new primary namenode, by using the `-importCheckpoint` option when starting the namenode daemon. The `-importCheckpoint` option will load the namenode metadata from the latest checkpoint in the directory defined by the `fs.checkpoint.dir` property, but if there is no metadata in the `dfs.name.dir` directory, so there is no risk of overwriting previous metadata.

Datanode directory structure:

As datanodes creates their storage directories automatically on startup, so no need to do explicitly formatted.

```
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-19 20:24 current/
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-12 23:59 detach/
-rw-r--r-- 1 kelly1 kelly1 0 2012-09-19 20:23 in_use.lock
-rw-r--r-- 1 kelly1 kelly1 157 2012-09-12 23:59 storage
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-19 20:23 tmp/
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/data$ cd current/
```

```

kelly1@kelly1-LIFEBOOK-LH530:~/dfs/data$ cd current/
kelly1@kelly1-LIFEBOOK-LH530:~/dfs/data/current$ ll
total 210072
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-19 20:24 .
drwxr-xr-x 2 kelly1 kelly1 4096 2012-09-19 20:23 ..
-rw-r--r-- 1 kelly1 kelly1 25551 2012-09-15 18:32 blk_15665771561366023_1903.meta
-rw-r--r-- 1 kelly1 kelly1 450 2012-09-15 18:29 blk_-1627393122322682795
-rw-r--r-- 1 kelly1 kelly1 11 2012-09-15 18:29 blk_-1627393122322682795_1005.meta
-rw-r--r-- 1 kelly1 kelly1 20413 2012-09-15 18:34 blk_1961868600997551736
-rw-r--r-- 1 kelly1 kelly1 167 2012-09-15 18:34 blk_1961868600997551736_1012.meta
-rw-r--r-- 1 kelly1 kelly1 67108364 2012-09-14 06:22 blk_-2805769607607470537
-rw-r--r-- 1 kelly1 kelly1 524295 2012-09-14 06:22 blk_-2805769607607470537_1003.meta
-rw-r--r-- 1 kelly1 kelly1 16880 2012-09-15 18:40 blk_-2947544845638000018
-rw-r--r-- 1 kelly1 kelly1 139 2012-09-15 18:40 blk_-2947544845638000018_1024.meta
-rw-r--r-- 1 kelly1 kelly1 67108364 2012-09-15 18:32 blk_3451198358097522603
-rw-r--r-- 1 kelly1 kelly1 524295 2012-09-15 18:32 blk_3451198358097522603_1006.meta
-rw-r--r-- 1 kelly1 kelly1 20755 2012-09-15 18:39 blk_4015418824033689825
-rw-r--r-- 1 kelly1 kelly1 171 2012-09-15 18:39 blk_4015418824033689825_1021.meta
-rw-r--r-- 1 kelly1 kelly1 67108364 2012-09-15 18:32 blk_-417367468448558263
-rw-r--r-- 1 kelly1 kelly1 524295 2012-09-15 18:32 blk_-417367468448558263_1006.meta
-rw-r--r-- 1 kelly1 kelly1 14786560 2012-09-14 06:22 blk_-5461212830652165877
-rw-r--r-- 1 kelly1 kelly1 115527 2012-09-14 06:22 blk_-5461212830652165877_1003.meta
-rw-r--r-- 1 kelly1 kelly1 10028 2012-09-15 18:39 blk_6147703354942118379
-rw-r--r-- 1 kelly1 kelly1 87 2012-09-15 18:39 blk_6147703354942118379_1023.meta
-rw-r--r-- 1 kelly1 kelly1 16718 2012-09-15 18:34 blk_-8102873251210258478
-rw-r--r-- 1 kelly1 kelly1 139 2012-09-15 18:34 blk_-8102873251210258478_1015.meta
-rw-r--r-- 1 kelly1 kelly1 688 2012-09-15 18:34 blk_-8252010316073749959
-rw-r--r-- 1 kelly1 kelly1 15 2012-09-15 18:34 blk_-8252010316073749959_1014.meta
-rw-r--r-- 1 kelly1 kelly1 36 2012-09-16 11:45 blk_8787259867977429757
-rw-r--r-- 1 kelly1 kelly1 11 2012-09-16 11:45 blk_8787259867977429757_1026.meta
-rw-r--r-- 1 kelly1 kelly1 4 2012-09-19 20:24 blk_-8860307820344948466
-rw-r--r-- 1 kelly1 kelly1 11 2012-09-19 20:24 blk_-8860307820344948466_1028.meta
-rw-r--r-- 1 kelly1 kelly1 2796 2012-09-19 20:23 dnccp_block_verification.log.curr
-rw-r--r-- 1 kelly1 kelly1 153 2012-09-19 20:23 VERSION

```

In VERSION, the namespaceID, cTime, and layoutVersion are all the same as the nameNode VERSION values. The namespaceID is retrieved from the namenode when the datanode first connects). The storageID is unique to the datanode(it is the same across all storage directories) and is used by the namenode to uniquely identify the datanode. The storageType identifies this directory as a datanode storage directory.

The files in the datanode's current storage directory are two types:

1) blk_<id>; It just consists of the raw bytes of a portion of the file being stored.

2) blk_<id>.metadata: It is made up of a header with version and type information, followed by a series of checksums for sections of the block.

Datanode creates a new subdirectory every time the number of blocks in a directory reaches 64(It can be configured using dfs.datanode.numblocks). By taking this measure, it ensures that there is a manageable number of files per directory, which avoids the problems that most operating systems encounter when there are a large number of files in a single directory.

Some important tools for administrator:

1) **hadoopdfsadmin:** To find information about the state of HDFS and to perform administration operations on HDFS

```
hadoop dfsadmin -help [cmd]
[ -status ]
[ -safeMode enter | leave | get | wait ]
[ -saveNamespace ]
[ -refreshNodes ]
[ -finalizeUpgrade ]
[ -upgradeProgress status | details | force ]
[ -metasave filename ]
[ -refreshServiceAcl ]
[ -refreshUserToGroupsMappings ]
[ -refreshSuperUserGroupsConfiguration ]
[ -setQuota <quota> <dirname>...<dirname> ]
[ -clrQuota <dirname>...<dirname> ]
[ -setSpaceQuota <quota> <dirname>...<dirname> ]
[ -clrSpaceQuota <dirname>...<dirname> ]
[ -setBalancerBandwidth <bandwidth in bytes per second> ]
[ -help [cmd]]
```

2) **hadoopfsck:** To provide an fsck utility for checking the health of the files in HDFS.

```
kelly1@kelly1-LIFEBOOK-LH530:~/Work/hadoop/bins ./hadoop fsck
Usage: DFSck <path> [-move | -delete | -openforwrite] [-files [-blocks [-locations [-racks]]]]
      <path> start checking from this path
      -move move corrupted files to /lost+found
      -delete delete corrupted files
      -files print out files being checked
      -openforwrite print out files opened for write
      -blocks print out block report
      -locations print out locations for every block
      -racks print out network topology for data-node locations
```

Datanode block scanner:

All the datanodes runs the block scanner, which periodically verifies all the blocks stored on the datanode. This allows bad blocks to be detected and fixed before they are read by clients.

It maintains

- 1) a list of blocks to verify
- 2) it scans them one by one for checksum errors.
- 3) block scanner report can be verify by visiting <http://datanode:50075/blockScannerReport>.

Upgrading the Cluster:

Steps for upgrading the cluster:

- 1) Make sure that any previous upgrade is finalized before proceeding with another upgrade.
- 2) Shut down MapReduce and kill any orphaned task processes on the tasktrackers.
- 3) Shut down HDFS and backup the namenode directories.
- 4) Install new versions of Hadoop HDFS and Mapreduce on the cluster and on clients.
- 5) Start HDFS with the -upgrade option.
- 6) Wait until the upgrade is complete.

7) Perform some sanity checks on HDFS.

8) Start MapReduce.

9) Roll back or finalize the upgrade

We can define two environment variables as HADOOP_OLD and HADOOP_NEW.

Start the upgrade:

```
% $HADOOP_NEW/bin/start-dfs.sh -upgrade
```

[This causes the namenode to upgrade its metadata, placing the previous version in a new directory called previous:

Wait until the upgrade is complete:

```
% $HADOOP_NEW/bin/hadoopfsadmin -upgradeProgress status
```

Do sanity check using fsck

If there is any problem in new version just rollback:

```
% $HADOOP_NEW/bin/stop-dfs.sh
```

```
% $HADOOP_OLD/bin/start-dfs.sh -rollback
```

Finalize the upgrade:

```
% $HADOOP_NEW/bin/hadoopfsadmin -finalizeUpgrade
```

```
% $HADOOP_NEW/bin/hadoopfsadmin -upgradeProgress status
```

How to kill the running mapreduce job?

```
kelly1@kelly1-LIFEBOOK-LH530:~/work/hadoop/bin$ ./hadoop job
Usage: JobClient <command> <args>
      [-submit <job-file>]
      [-status <job-id>]
      [-counter <job-id> <group-name> <counter-name>]
      [-kill <job-ids>]
      [-set-priority <job-id> <priority>]. Valid values for priorities are: VERY_HIGH HIGH NORMAL LOW VERY_LOW
      [-events <job-id> <from-event-> <# of events>]
      [-history <jobOutputDir>]
      [-list [all]]
      [-list-active-trackers]
      [-list-blacklisted-trackers]
      [-list-attempt-ids <job-id> <task-type> <task-state>]

      [-kill-task <task-id>]
      [-fail-task <task-id>]
```

```
kelly1@kelly1-LIFEBOOK-LH530:~/work/hadoop/bin$ ./hadoop job
Usage: JobClient <command> <args>
      [-submit <job-file>]
      [-status <job-id>]
      [-counter <job-id> <group-name> <counter-name>]
      [-kill <job-ids>]
      [-set-priority <job-id> <priority>]. Valid values for priorities are: VERY_HIGH HIGH NORMAL LOW VERY_LOW
      [-events <job-id> <from-event-> <# of events>]
      [-history <jobOutputDir>]
      [-list [all]]
      [-list-active-trackers]
      [-list-blacklisted-trackers]
      [-list-attempt-ids <job-id> <task-type> <task-state>]

      [-kill-task <task-id>]
      [-fail-task <task-id>]
```

MapReduce

1. MapReduce Introduction
2. Combiners
3. Partitioners
4. Writables
5. Writable Comparables
6. Distributed Cache
7. Input Formatters
8. Out Formatters
9. Schedulers
 - a. FIFO Scheduler
 - b. FAIR Scheduler
 - c. Capacity Scheduler
10. Compression
 - a. LZO
 - b. Snappy
11. Hadoop cluster up
 - a. Prerequisites
 - b. Pseudo Cluster Mode
 - c. Cluster Mode

MapReduce Introduction

MapReduce is a programming model for processing BigData. It processes data in parallel. To take advantage of the parallel processing that Hadoop provides, express your business logic as a Map Reduce Job.

MapReduce works by breaking the processing into two phases: the map phase and the reduce phase. Each phase has key-value pairs as input and output. Input key-value pair of map phase is determined by the type of Input Formatter being used, whereas rest of the key-value pair types may be chosen by the programmer. The programmer also specifies two functions: the map function and the reduce function.

The map function is represented by extending the Mapper Abstract class. Mapper class is a generic type, with four formal type of parameters that specify the input key, input value, output key and output value types of map function.

The reduce function is similarly defined by extending the Reducer Abstract class. Reducer class is also a generic type, with four formal type of parameters that specify the input key, input value, output key and output value types of map function.

Reducer aggregates the map outputs. So one needs to remember that map output key type, map output value type should match with reduce input key type and reduce input value type.

Let us write a MapReduce application for counting words in a text file.

Example-1 shows the implementation of map function for WordCount example.

```
public static class WordCountMapper  
extends Mapper<LongWritable, Text, Text, IntWritable>{  
  
    private final static IntWritable one = new IntWritable(1);  
    private Text word = new Text();  
  
    public void map (LongWritable key, Text value, Context context  
                    ) throws IOException, InterruptedException {  
        StringTokenizer itr = new StringTokenizer(value.toString());  
        while (itr.hasMoreTokens()) {  
            word.set(itr.nextToken());  
            context.write(word, one);  
        }  
    }  
}
```

Example-1

WordCountMapper extends Mapper<KEYIN, VALUEIN, KEYOUT, VALUEOUT> with specific types like LongWritable, Text, Text, IntWritable.

Hadoop uses its own data types rather than using Java data types. These data types are optimized for network serialization. These data types are found in the org.apache.hadoop.io package. Here LongWritable and Text corresponds to long and String respectively in Java.

For the above example, input key is line offset in the file and value is the entire line. The third parameter is a context object, which allows the input and output from the task. It is only supplied to the Mapper and Reducer.

For the above example, as we are counting number of times each word occurred, take the value, which is the entire line and convert it into String using toString() method in Text class. Split the line into words using StringTokenizer class. Iterate over all the words in the line and output word and its count 1 (the word has already occurred one time, that's why 1). For

outputting word and its count, use the context object, as discussed above.
Context has write method to write the map output.

Map outputs are stored in the Local FileSystem, where mapper has run. Also Map outputs are stored in the sorted order, sorted on Key. If the Job is Map only
Job, Map outputs are not sorted and directly stored in the configured File System in core-site.xml
(In most cases it is HDFS).

Let us consider the following input which has 4 Lines for the WordCount example

Hadoop is used for BigData Analysis.

Hadoop has two components namely HDFS and MapReduce.

HDFS is a Distributed Data Storage Framework

MapReduce is a Distributed Computing Framework.

input-1

The output of the Example-1 is

a 1
a 1
components 1
<complete this>

Example-2 shows the implementation of reduce function for WordCount example.

```
public static class WordCountReducer  
    extends Reducer<Text, IntWritable, Text, IntWritable> {  
    private IntWritable result = new IntWritable();  
  
    public void reduce(Text key, Iterable<IntWritable> values,  
                      Context context  
                      ) throws IOException, InterruptedException {  
        int sum = 0;  
        for (IntWritable val : values) {  
            sum += val.get();  
        }  
        result.set(sum);  
        context.write(key, result);  
    }  
}
```

Example-2

WordCountReducer extends Reducer<KEYIN, VALUEIN, KEYOUT, VALUEOUT> with specific types like Text, IntWritable, Text, IntWritable.

Note that WordCountReducer input key type is Text and value type is IntWritable, matches with the output key type and value type of WordCountMapper.

Reducer first copies all the map outputs and merges them. This process is called Shuffle and Sorting. During this phase, all the values corresponding to each key will be aggregated. reduce method is called for each key with all the values aggregated for that key.

For the above map output, Reducer first copies map output and as there is only one map output, merge is trivial. reduce method is called for each key and its aggregated values. So reduce method in WordCountReducer get the following input:

A 1,1
Distributed 1,1
<Complete this>

In the reduce method, iterate over all the values and counting how many them. Finally writing the reduce method output, which is word and its count using write method provided by Context class.

Finally we need to provide a driver class, which submits the Job to the HADOOP cluster with our own mapper and reduce implementation.

```
public class WordCount {  
  
    public static void main(String[] args) throws Exception {  
        Configuration conf = new Configuration();  
  
        if (otherArgs.length != 2) {  
            System.err.println("Usage: WordCount <in><out>");  
            System.exit(2);  
        }  
        Job job = new Job(conf, "kelly word count");  
        job.setJarByClass(WordCount.class);  
        job.setMapperClass(WordCountMapper.class);  
        job.setReducerClass(WordCountReducer.class);  
        job.setOutputKeyClass(Text.class);  
        job.setOutputValueClass(IntWritable.class);  
        FileInputFormat.addInputPath(job, new Path(otherArgs[0]));  
        FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));  
        System.exit(job.waitForCompletion(true) ? 0 : 1);  
    }  
}
```

Example-3

In Example-3, created a Job object with configuration object and name of the job being the arguments.

Configuration is an object through which Hadoop reads all the configuration properties mentioned in core-site.xml, hdfs-site.xml and mapred-site.xml.

By passing the conf object while job object creation, job object avails all the properties from the afore mentioned files during its life cycle (like map or reduce phase). Also a user can define their own property with a value (value type could be any data type) using setter methods provided in the Configuration object like conf.setInt (use the appropriate API based on the value type) method. Hence Configuration is very useful if you need to pass a small piece of metadata to your tasks. To retrieve the values in the task using context object and get configuration object using context.getConfiguration() and then get the required metadata.

Job object has been specified with Mapper, Reducer classes, output key type, output value type. Input path is mentioned using the static method of FileInputFormat class's addInputPath having job instance and input directory as arguments. Similarly output directory has been mentioned using static method provided by FileOutputFormat class. Finally submit the job to cluster using waitForCompletion method and wait for job completion. Once job completes successfully, find the output of the job in the output directory.

How to submit the job to cluster

1. Go to HADOOP_HOME directory
2. Run the following command to submit the job

bin/hadoop jar <path of jar file><package-name>.WordCount /input /output.

Suppose name of the jar file is jobs.jar and let us say this jar file is in /home/hadoop directory and package-name is kelly.training.hadoop then above command becomes

bin/hadoop jar /home/hadoop/jobs.jar kelly.training.hadoop.WordCount /input /output

Exercises

Write a MapReduce Job to display all the lines which has given word. (Simulating UNIX 'grep' functionality).

Ex: Suppose given word is : BigData.

For the input-1, MapReduce Job output should contain only the following line, since this line contains BigData word.

Output : Hadoop is used for BigData Analysis.

2. Write a MapReduce Job to display all the lines by replacing the given word and another input word. (Simulating UNIX 'sed' functionality).

Ex: Suppose given word is : *BigData* and another word is : *LargeData*.

For the input-1, MapReduce job output should contain only the following line, since this line contains BigData word. However BigData should be replace with LargeData, so the output should be

Output : Hadoop is used for LargeData Analysis.

3. Write a MapReduce Job to find the average of frequency of each word.

The format of the Text file is as follows

hadoop	10
hbase	15
hive	30
hadoop	47
solr	33
katta	68
hbase	26
hive	22

input-2

For input-2, the output should be

hadoop	28.5
hbase	20.5
hive	26
solr	33
katta	68.

Combiners

Combiners are used to reduce the amount of the data being transferred over the network. This is to optimize the usage of the network bandwidth, which is limited in most of the cases.

Combiners are treated as local reducers. They run by consuming the mapper output and run on the same machine where mapper had run earlier. Hadoop does not provide any guarantee on combiner's execution. Hadoop may call combiner function zero, one or many times for a particular map output record.

Imagine the first mapper produced the output

hadoop	1
--------	---

```
hadoop    1
hadoop    1
hadoop    1
hadoop    1
```

and second mapper produced

```
hadoop    1
hadoop    1
hadoop    1
```

We could use a combiner function that, just like reduce function, counts the word frequency.
The reduce would then be called with:

```
hadoop    <5,3>
```

Note: While calling reduce method, Hadoop does not provide any guarantee on values in sorted order corresponding to a particular key. To achieve this we used Secondary Sorting.

The combiner function does not replace the reduce function. Combiner is also implemented by extending the Reduce Abstract class and overriding the reduce method.

Specifying a combiner function

For the WordCount example WordCountReducer can also be used as Combiner and the same can be set using setCombinerClass API from Job class as highlighted below.

```
public class WordCount {
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        if (otherArgs.length != 2) {
            System.err.println("Usage: wordcount <in><out>");
            System.exit(2);
        }
        Job job = new Job(conf, "word count");
        job.setJarByClass(WordCount.class);
        job.setMapperClass(WordCountMapper.class);
        job.setCombinerClass(WordCountReducer.class);
        job.setReducerClass(WordCountReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
    }
}
```

```
FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}
```

Exercises

4. Can we use a combiner class for problem 3? Provide your explanation.

5. Write an MR job to find the highest repeated words in a file.

For the input-1, "is" repeated 3 times so the output should be
is 3

If there are more than one such word, print all those words.

Partitioners

Partitioner allows you to distribute how outputs from the map stage are sent to the reducers.
Basically it partitions the key space.

Partitioner controls the partitioning the keys of the intermediate map-outputs. The key (or subset of the key) is used to derive the partition.

The total number of partitions is the same as the number of reduce tasks for the job. Hence Partitioner controls which of the r reduce tasks the intermediate key (and hence record) is sent for reduction.

Partitioner run on the same machine after Mapper had completed its execution, by consuming Mapper output. So entire Mapper output (record) is sent to Partitioner and Partitioner forms r (number of reduce tasks) groups from the mapper output.

By default Hadoop framework Hash based partitioner. This Partitioner partitions the key space by using the hashCode.

The following logic HashPartitioner executes to determine a reducer for a particular Key 'key'.

~~How to Write Custom Partitioner~~

Hadoop provides Partitioner abstract class with a single method which can be extended to write a custom partitioner.

```
public abstract class Partitioner<KEY, VALUE> {
    public abstract int getPartition(KEY key, VALUE value, int numPartitions);
```

}

As discussed earlier KEY and VALUE type must match with MapOutput key and MapOutput value.

Example

Let us take the WordCount Example and suppose say requirement is all the words which starts with "a" go to one reducer and all the words which starts with "b" go to another reducer etc...

So the number of reducer required is 26.

Let us write a Partitioner for that

```
public class WordCountPartitioner extends Partitioner<Text, IntWritable> {  
  
    public int getPartition(Text key, IntWritable value, int numPartitions) {  
        return (key.toString().charAt(0) - 'a') % numPartitions;  
    }  
}
```

WordCountMapper's output key is Text and value type is IntWritable, so WordCountPartitioner Key type also Text and Value type is IntWritable.

In the getPartition method, take the first character of the word and finding its position relative to 'a'.

Suppose say firstChar is 'c'. So the ascii code of c is 99 and ascii code of a is 97.

so 'c' - 'a' gives 2, which is relative to 'a'.

Hence all words which starts with 'c' goes to 2nd reducer.

Specifying a Partitioner function

Partitioner class can be set using the setPartitionerClass method from the job class. The same has been highlighted below. Also specify number of reduce tasks as 26.

```
public class WordCount {  
  
    public static void main(String[] args) throws Exception {  
        Configuration conf = new Configuration();
```

```

if (otherArgs.length != 2) {
    System.err.println("Usage: wordcount <in><out>");
    System.exit(2);
}
Job job = new Job(conf, "word count");
job.setJarByClass(WordCount.class);
job.setMapperClass(WordCountMapper.class);
job.setCombinerClass(WordCountReducer.class);
job.setPartitionerClass(WordCountPartitioner.class);
job.setNumReduceTasks(26);
job.setReducerClass(WordCountReducer.class);
job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);
FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

By default number of reduce task be 1.

Caution: It is always recommended to provide implementation of hashCode when you write your own composite key.

In the WordCountPartitioner, observe that getPartition() method has been implemented without considering hashCode of the key. But it may not possible to write getPartition() method without considering hashCode of the key always.

Writable

Writable is an interface, which Hadoop uses for serialization and de-serialization purpose.
Any key or value type in Hadoop Map-Reduce framework must implement this interface.

This interface is a simple, efficient, serialization protocol, based on DataInput and DataOutput.

To implement your own Writable class, you must implement the following methods.

void write(DataOutput out) throws IOException;

Serialize of the fields of this object (state of the object) to out (DataOutput binary stream) using write implementation.

void readFields(DataInput in) throws IOException;

De-serialize the fields of this object (state of the object) from in (DataInput binary stream) using readFields implementation.

WritableComparable

WritableComparable is also an interface which extends both Writable interface and java.lang.Comparable interface.

Any Key in Hadoop Map-Reduce framework must implement this interface. During the sorting phase keys are compared with one another using implementation provided for the following method

```
public int compareTo(T o);
```

The Following are some of mostly used WritableComparable classes provided by Hadoop framework

BooleanWritable	-- boolean
ByteWritable	-- byte
IntWritable	-- int
FloatWritable	-- float
LongWritable	-- long
DoubleWritable	-- double
Text	-- String

Please go through the org.apache.hadoop.io API documentation for more Writable and WritableComparable classes

Distributed Cache

Distributed cache is a mechanism for copying large and read only files and archives to the task nodes in time for the tasks use them when they run. To save network bandwidth, files or archives are normally copied to any particular node once per job.

How to specify Distributed Cache entries?

Hadoop provides 2 options to achieve the same

a) API :

Job class have provided various APIs to add the files/archives into Distributed Cache.

The following are APIs from Job class

`addCacheArchive(URI uri)` : Add any archive to be cached using URI object.
`addCacheFile(URI uri)` : Add any read only file (like lookup) to be cached using URI object.
`addFileToClassPath(Path file)` : If a file is required to be in Child's classpath, use this API. It also adds the file to cache as well.
Ex: Suppose say your job depends on third party jars, you can use this API to place all those jar in child's classpath.
`addArchiveToClassPath(Path archive)` : If an archive is required to be in Child's classpath, use this API. It also adds the archive to cache as well.

To get the cache entries in either map/ reduce use the following APIs.

Use the context object in map/reduce/setup method in Mapper/Reducer class.

`getCacheFiles()` : To get the cache files.
`getCacheArchives()` : To get the cache archives.

b) Using CLI.

Rather than using the API, while running the job, we can use "-libjars" and "-files" options to add the entries to Distributed Cache.

-libjars : comma separated jar files to include in the classpath.
-files : comma separated files to be copied to the map reduce cluster
-archives : comma separated archives to be unarchived on the compute machines.
Go through the following example which illustrates how to use Distributed Cache .
PhoneRetrievalJob is to retrieve the phone numbers of the users. Input to the job contains list of users. A lookup file contains user and phone number. If the input user is present in the lookup file, we emit his phonenumber otherwise we just ignore the user.
package kelly.hadoop.training.jobs.dc;

```
import java.net.URI;  
  
import org.apache.hadoop.conf.Configuration;  
import org.apache.hadoop.filecache.DistributedCache;  
import org.apache.hadoop.fs.Path;  
import org.apache.hadoop.io.LongWritable;  
import org.apache.hadoop.io.Text;  
import org.apache.hadoop.mapreduce.Job;  
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
import org.apache.hadoop.mapreduce.lib.input.KeyValueTextInputFormat;  
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;  
import org.apache.hadoop.util.Tool;  
import org.apache.hadoop.util.ToolRunner;
```

```
public class PhoneRetrievalJob implements Tool {  
  
    private Configuration conf;  
  
    @Override  
    public int run(String[] args) throws Exception {  
        Job phoneJob = new Job(getConf());  
        phoneJob.setJobName("Kelly Phone Job");  
        phoneJob.setJarByClass(this.getClass());  
  
        phoneJob.setNumReduceTasks(0);  
  
        phoneJob.setMapperClass(PhoneRetrievalMapper.class);  
        phoneJob.setOutputKeyClass(Text.class);  
        phoneJob.setOutputValueClass(LongWritable.class);  
  
        phoneJob.setInputFormatClass(KeyValueTextInputFormat.class);  
  
        Path lookupPath = new Path(args[2]);  
        URI[] uris = new URI[1];  
        uris[0] = lookupPath.toUri();  
  
        DistributedCache.setCacheFiles(uris, phoneJob.getConfiguration());  
  
        FileInputFormat.setInputPaths(phoneJob, new Path(args[0]));  
        FileOutputFormat.setOutputPath(phoneJob, new Path(args[1]));  
  
        return phoneJob.waitForCompletion(true) == true ? 0 : -1;  
    }  
  
    @Override  
    public Configuration getConf() {  
        return conf;  
    }  
  
    @Override  
    public void setConf(Configuration conf) {  
        this.conf = conf;  
    }  
  
    public static void main(String[] args) throws Exception {  
        ToolRunner.run(new Configuration(), new PhoneRetrievalJob(), args);  
    }  
}
```

To run the above job, user need to mention the location of the lookup file in the File System. If you are run this job in local mode, the lookup file should present in Local File System. If you are running this job in Cluster mode, copy the lookup file into HDFS, and then give the path of HDFS while running the job.

The same procedure should be followed for all the jobs, when your job uses files in the Distributed cache.

Highlighted text contains, how to set the file in the Distributed cache to job instance.

In the mapper, we get the file from the Distributed Cache, open the lookup file and if the input record (user) present in the lookup file we emit user and his phone number otherwise we just ignore the record. Usage of the Distributed Cache in **PhoneRetrievalMapper**, has been highlighted in bold.

```
package kelly.hadoop.training.jobs.dc;
```

```
import java.io.BufferedReader;
import java.io.File;
import java.io.InputStream;
import java.io.InputStreamReader;

import org.apache.hadoop.filecache.DistributedCache;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class PhoneRetrievalMapper extends Mapper<Text, Text, Text, Text> {

    private BufferedReader reader;
    private Path lookUp;
    private Text phone = new Text();

    @Override
    protected void setup(Context context) throws java.io.IOException,
        InterruptedException {
        Path[] cacheFiles = DistributedCache.getLocalCacheFiles(context.getConfiguration());
        lookUp = cacheFiles[0];
    }

    @Override
    protected void map(Text key, Text value, Context context)
        throws java.io.IOException, InterruptedException {
        File file = new File(lookUp.toString());
        reader = new BufferedReader(new InputStreamReader(new FileInputStream(
            file)));
        String line;
```

```
while ((line = reader.readLine()) != null) {  
    if (line.contains(key.toString())) {  
        String phoneNum = line.split("\t")[1];  
        phone.set(phoneNum);  
        context.write(key, phone);  
    }  
}  
reader.close();  
};  
}
```

InputFormatters

Inputformatters describe the input specification of the job. The responsibility of the Input Formatters

- 1) Validate the input specification of the job
- 2) Split up the input file(s) into logical Input Split. Each input split is then assigned to an individual Mapper.
- 3) Provides a Record Reader to glean input record from logical Input Split for processing by Mapper.

Input split is computed using the following formulae

Max (MinSplitSize,Min (MaxSplitSize, HDFS Blk))

Where split size max is configured using mapred.max.split.size property and min split size is configured using mapred.min.split.size property.

Each record is splitted into Key-Value pairs. These key-value pairs will act as Map input key and Map input value. The type of input formatter to use is depends on the nature of the data.

There are various types of input formatters provided by the framework.

TextInputFormatter

This IF works on the text files. TIF recognizes each line (by default delimited by \n) in the file as one record. For TIF, key will be line offset in the file (LongWritable) and value will be entire line (Text). By default Hadoop used TIF.

Ex:

I am new to Hadoop.

Hadoop is great software.

As there are two lines in the files, so there are two records.

For the first record: key: 0(line offset) value: I am new to Hadoop

For the second record: key: 21(line offset) value: Hadoop is great software.

KeyValueTextInputFormatter

This IF works on the text files; however, it works on the text files which are in the key value format. KVTIF recognizes each line (by default delimited by \n) in the file as one record. Each line is divided into key and value parts by a separator byte. By default key and values in the file are delimited by \t (tab). But it can be overridden. If no such a byte exists, the key will be the entire line and value will be empty.

Ex:

Andhra-Pradesh	Hyderabad
Madhya Pradesh	Bhopal
Karnataka	Bangalore.

From the above input file, it is obvious that there are 3 records (as there are 3 lines in the file).

Also observe that above file is in Key Value format (Keys and Values are delimited by tab).

For the first record: key: Andhra Pradesh value: Hyderabad

For the second record: key: Madhya Pradesh value: Bhopal

So here both key and value types are Text.

Let us illustrate the same using an example

package kelly.hadoop.training.jobs.inputformatter;

```

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.KeyValueTextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

public class WordPairCountJob implements Tool {

    private Configuration conf;

    @Override
    public int run(String[] args) throws Exception {
        Job wordPairCountJob = new Job(getConf());
        wordPairCountJob.setJobName("Kelly WordPair Count");
        wordPairCountJob.setJarByClass(this.getClass());

        wordPairCountJob.setMapperClass(WordPairCountMapper.class);
        wordPairCountJob.setMapOutputKeyClass(WordPairCountKey.class);
        wordPairCountJob.setMapOutputValueClass(LongWritable.class);

        wordPairCountJob.setReducerClass(WordPairCountReducer.class);
        wordPairCountJob.setOutputKeyClass(WordPairCountKey.class);
        wordPairCountJob.setOutputValueClass(LongWritable.class);
    }
}

```

```
wordPairCountJob.setInputFormatClass(KeyValueTextInputFormat.class);

FileInputFormat.setInputPaths(wordPairCountJob, new Path(args[0]));
FileOutputFormat.setOutputPath(wordPairCountJob, new Path(args[1]));

return wordPairCountJob.waitForCompletion(true) == true ? 0 : -1;
}

@Override
public Configuration getConf() {
    return conf;
}

@Override
public void setConf(Configuration conf) {
    this.conf = conf;
}

public static void main(String[] args) throws Exception {
    ToolRunner.run(new Configuration(), new WordPairCountJob(), args);
}
}

package kelly.hadoop.training.jobs.inputformatter;

import java.io.DataInput;
import java.io.DataOutput;
import java.io.IOException;

import org.apache.hadoop.io.WritableComparable;

public class WordPairCountKey implements WritableComparable<WordPairCountKey> {

    private String word1;
    private String word2;

    public String getWord1() {
        return word1;
    }

    public void setWord1(String word1) {
        this.word1 = word1;
    }
}
```

```
public String getWord2() {
    return word2;
}

public void setWord2(String word2) {
    this.word2 = word2;
}

@Override
public void readFields(DataInput in) throws IOException {
    word1 = in.readUTF();
    word2 = in.readUTF();
}

@Override
public void write(DataOutput out) throws IOException {
    out.writeUTF(word1);
    out.writeUTF(word2);
}

@Override
public int compareTo(WordPairCountKey o) {
    int diff = word1.compareTo(o.word1);
    if (diff == 0) {
        diff = word2.compareTo(o.word2);
    }
    return diff;
}

@Override
public int hashCode() {
    return word1.hashCode() + 31 * word2.hashCode();
}

@Override
public String toString() {
    return "[word1=" + word1 + ", word2=" + word2 + "]";
}

}

package kelly.hadoop.training.jobs.inputformatter;

import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
```

```
import org.apache.hadoop.mapreduce.Mapper;

public class WordPairCountMapper extends
    Mapper<Text, Text, WordPairCountKey, LongWritable> {

    private final static LongWritable one = new LongWritable(1);

    @Override
    protected void map(Text key, Text value, Context context)
        throws java.io.IOException, InterruptedException {
        WordPairCountKey wpKey = new WordPairCountKey();
        wpKey.setWord1(key.toString());
        wpKey.setWord2(value.toString());
        context.write(wpKey, one);
    }

}

package kelly.hadoop.training.jobs.inputformatter;

import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.mapreduce.Reducer;

public class WordPairCountReducer extends
    Reducer<WordPairCountKey, LongWritable, WordPairCountKey, LongWritable> {

    @Override
    protected void reduce(WordPairCountKey key,
        java.lang.Iterable<LongWritable> values, Context context)
        throws java.io.IOException, InterruptedException {
        long sum = 0;
        while(values.iterator().hasNext()) {
            sum += values.iterator().next().get();
        }
        context.write(key, new LongWritable(sum));
    }

}
```

Sequence File Input Formatter

This IF works only on special binary files called Sequence file. Sequence files are specific to Hadoop. Sequence files are also in key value format, however, unlike KVTIF, keys and values can be any type.

Sequence files can be generated as the output of other MapReduce tasks and are an efficient intermediate representation for data that is passing from one MapReduce job to another.

This IF is used to process binary data like images, video and audio files. Also it can also be used to process small files. While processing small files, we can use file name as key and entire file as value.

We can apply compression to values while storing them in Sequence Files.

```
package kelly.hadoop.training.jobs.inputformatter;

import java.io.IOException;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.FSDataInputStream;
import org.apache.hadoop.fs.FileStatus;
import org.apache.hadoop.fs.FileSystem;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.BytesWritable;
import org.apache.hadoop.io.SequenceFile;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.SequenceFile.CompressionType;
import org.apache.hadoop.io.SequenceFile.Writer;

public class SequenceFileGenerator {

    private static final int BUFFER_SIZE = 2 * 1024 * 1024;

    public static void main(String[] args) throws IOException {
        Writer writer = null;
        Configuration conf = new Configuration();
        Path inputDir = new Path(args[0]);
        Path seqFile = new Path(args[1]);

        try {
            writer = SequenceFile.createWriter(seqFile.getFileSystem(conf),
                conf, seqFile, Text.class, BytesWritable.class,
                CompressionType.BLOCK);
            FileSystem fs = inputDir.getFileSystem(conf);
            BytesWritable content = new BytesWritable();
            Text fileName = new Text();
            byte[] data = new byte[BUFFER_SIZE];
            for (FileStatus file : fs.listStatus(inputDir)) {
                FSDataInputStream stream = fs.open(file.getPath());
                stream.read(data);
                content.set(data, 0, (int) file.getLen());
                fileName.set(file.getPath().getName());
                writer.append(fileName, content);
                stream.close();
            }
        } catch (IOException e) {
            e.printStackTrace();
        } finally {
            if (writer != null) {
                writer.close();
            }
        }
    }
}
```

```
        }
        writer.close();
    } finally {
        if (writer != null) {
            writer.close();
        }
    }
}

package kelly.hadoop.training.jobs.inputformatter;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.SequenceFileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

public class WordCountPerFileJob implements Tool {

    private Configuration conf;

    @Override
    public int run(String[] args) throws Exception {
        Job wordCountPerFileJob = new Job(getConf());
        wordCountPerFileJob.setJobName("Kelly Word Count Per File");
        wordCountPerFileJob.setJarByClass(this.getClass());

        wordCountPerFileJob.setMapperClass(WordCountPerFileMapper.class);
        wordCountPerFileJob.setMapOutputKeyClass(WordCountPerFileKey.class);
        wordCountPerFileJob.setMapOutputValueClass(IntWritable.class);

        wordCountPerFileJob.setReducerClass(WordCountPerFileReducer.class);
        wordCountPerFileJob.setOutputKeyClass(WordCountPerFileKey.class);
        wordCountPerFileJob.setOutputValueClass(IntWritable.class);

        wordCountPerFileJob.setInputFormatClass(SequenceFileInputFormat.class);

        FileInputFormat.setInputPaths(wordCountPerFileJob, new Path(args[0]));
        FileOutputFormat.setOutputPath(wordCountPerFileJob, new Path(args[1]));
    }
}
```

```
        return wordCountPerFileJob.waitForCompletion(true) == true ? 0 : -1;
    }

    @Override
    public Configuration getConf() {
        return conf;
    }

    @Override
    public void setConf(Configuration conf) {
        this.conf = conf;
    }

    public static void main(String[] args) throws Exception {
        ToolRunner.run(new Configuration(), new WordCountPerFileJob(), args);
    }
}
```

```
package kelly.hadoop.training.jobs.inputformatter;

import java.io.DataInput;
import java.io.DataOutput;
import java.io.IOException;

import org.apache.hadoop.io.WritableComparable;

public class WordCountPerFileKey implements
    WritableComparable<WordCountPerFileKey> {

    private String fileName;
    private String word;

    public String getFileName() {
        return fileName;
    }

    public void setFileName(String fileName) {
        this.fileName = fileName;
    }

    public String getWord() {
        return word;
    }
}
```

```
public void setWord(String word) {  
    this.word = word;  
}  
  
@Override  
public void readFields(DataInput in) throws IOException {  
    fileName = in.readUTF();  
    word = in.readUTF();  
}  
  
@Override  
public void write(DataOutput out) throws IOException {  
    out.writeUTF(fileName);  
    out.writeUTF(word);  
}  
  
@Override  
public int compareTo(WordCountPerFileKey o) {  
    int diff = fileName.compareTo(o.fileName);  
    if (diff == 0) {  
        diff = word.compareTo(o.word);  
    }  
    return diff;  
}  
  
@Override  
public String toString() {  
    return "[fileName=" + fileName + ", word=" + word + "]";  
}  
}  
  
package kelly.hadooptraining.jobs.inputformatter;  
  
import java.util.StringTokenizer;  
  
import org.apache.hadoop.io.BytesWritable;  
import org.apache.hadoop.io.IntWritable;  
import org.apache.hadoop.io.Text;  
import org.apache.hadoop.mapreduce.Mapper;  
  
public class WordCountPerFileMapper extends  
    Mapper<Text, BytesWritable, WordCountPerFileKey, IntWritable> {
```

```
private final static IntWritable one = new IntWritable(1);

@Override
protected void map(Text key, BytesWritable value, Context context)
    throws java.io.IOException, InterruptedException {
    String fileName = key.toString();
    String content = new String(value.getBytes());
    WordCountPerFileKey wcKey = new WordCountPerFileKey();
    wcKey.setFileName(fileName);
    StringTokenizer strTock = new StringTokenizer(content, " ");
    while (strTock.hasMoreTokens()) {
        wcKey.setWord(strTock.nextToken());
        context.write(wcKey, one);
    }
}

package kelly.hadoop.training.jobs.inputformatter;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.mapreduce.Reducer;

public class WordCountPerFileReducer
    extends
    Reducer<WordCountPerFileKey, IntWritable, WordCountPerFileKey, IntWritable> {

    @Override
    protected void reduce(WordCountPerFileKey key,
        java.lang.Iterable<IntWritable> values, Context context)
        throws java.io.IOException, InterruptedException {
        int count = 0;
        for (IntWritable intW : values) {
            count += intW.get();
        }
        context.write(key, new IntWritable(count));
    }
}
```

OutFormatters

The responsibility of the output formatters is to write the output of jobs to underlying storage(either HDFS or LocalFS or DataBase etc...). The output of the job is written using

help of Record Writers. So it means final output is a set of records. The type of output depends on the type of OF being used.

Also Outputformatters verify the whether the output directory already exists or not.

Mapreduce frame work provides the following Output Formatters

TextOutputFormatter

If we don't specify any output formatter, Hadoop by default considers this output formatter. This OF writes the output to plain text files.

SequenceFileOutputFormatter

This OF writes output to Sequence Files

Schedulers

The functionality of the Scheduler is to schedule the tasks(either Map Task or Reduce Task) to TaskTrackers. There are 3 Schedulers available in Hadoop.

- 1) Fifo Scheduler
- 2) Fair Scheduler
- 3) Capacity Task Scheduler

1) Fifo Scheduler

This is the default scheduler, Hadoop uses to schedule the tasks. This Scheduler simply schedules the tasks of a job in order of submission. Each job uses the whole cluster, so jobs have to wait their turn. Users can assign priority to a job while submission to the Cluster through mapred.job.priority property. Priorities could be one among VERY_HIGH, HIGH, NORMAL, LOW, VERY_LOW. Default priority of a job is NORMAL.

When the job scheduler is choosing the next job to run, it selects one with the highest priority. However, with the FIFO scheduler, priorities do not support preemption, so a high-priority job can still be blocked by a long running low priority job that started before the high-priority job scheduled.

2) Fair Scheduler

Fair scheduling is a method of assigning resources to jobs such that all jobs get, on average, an equal share of resources over time.

When there is a single job running, that job uses the entire cluster. When other jobs are submitted, tasks slots that free up are assigned to the new jobs, so that each job gets roughly the same amount of CPU time. Unlike Fifo scheduler, this lets short jobs finish in reasonable time while not starving long jobs. It is also an easy way to share a cluster between multiple of users.

Fair sharing can also work with job priorities - the priorities are used as weights to determine the fraction of total compute time that each job gets.

The fair scheduler organizes jobs into pools, and divides resources fairly between these pools.

By default, there is a separate pool for each user, so that each user gets an equal share of the cluster.

Within each pool, jobs can be scheduled using either fair sharing or first-in-first-out (FIFO) scheduling.

In addition to providing fair sharing, the Fair Scheduler allows assigning guaranteed *minimum shares* to pools, which is useful for ensuring that certain users, groups or production applications always get sufficient resources. When a pool contains jobs, it gets at least its minimum share, but when the pool does not need its full guaranteed share, the excess is split between other pools.

If a pool's minimum share is not met for some period of time, the scheduler optionally supports *preemption* of jobs in other pools. The pool will be allowed to kill tasks from other pools to make room to run. A pool can also be allowed to preempt tasks if it is below half of its fair share for a configurable timeout. When choosing tasks to kill, the fair scheduler picks the most-recently-launched tasks from over-allocated jobs, to minimize wasted computation. Preemption does not cause the preempted jobs to fail, because Hadoop jobs tolerate losing tasks; it only makes them take longer to finish.

The Fair Scheduler can limit the number of concurrent running jobs per user and per pool. This can be useful when a user must submit hundreds of jobs at once, or for ensuring that intermediate data does not fill up disk space on a cluster when too many concurrent jobs are running. Setting job limits causes jobs submitted beyond the limit to wait until some of the user/pool's earlier jobs finish. Jobs to run from each user/pool are chosen in order of priority and then submit time.

The Fair Scheduler can limit the number of concurrent running tasks per pool. This can be useful when jobs have a dependency on an external service like a database or web service that could be overloaded if too many map or reduce tasks are run at once.

Installation

set the following property in the Hadoop config file *HADOOP_CONF_DIR/mapred-site.xml* to have Hadoop use the fair scheduler:

```
<property>
  <name>mapred.jobtracker.taskScheduler</name>
  <value>org.apache.hadoop.mapred.FairScheduler</value>
</property>
```

Once you restart the cluster, you can check that the fair scheduler is running by going to <http://<jobtracker URL>/scheduler> on the JobTracker's web UI.

Configuration

The Fair Scheduler contains configuration in two places -- algorithm parameters are set in *HADOOP_CONF_DIR/mapred-site.xml*, while a separate XML file called the *allocation file*, located by default in *HADOOP_CONF_DIR/fair-scheduler.xml*, is used to configure pools, minimum shares, running job limits and preemption timeouts. The allocation file is reloaded periodically at runtime, allowing you to change pool settings without restarting your Hadoop cluster.

Important setting for *mapred-site.xml*

`mapred.fairscheduler.preemption` : Boolean property for enabling preemption. Default: false.

`mapred.fairscheduler.pool` : Specify the pool that a job belongs in.

`mapred.fairscheduler.sizebasedweight` : Take into account job sizes in calculating their weights for fair sharing. By default, weights are only based on job priorities. Setting this flag to true will make them based on the size of the job (number of tasks needed) as well, though not linearly. This lets larger jobs get larger fair shares while still providing enough of a share to small jobs to let them finish fast. Boolean value, default: false.

Important settings for fair-scheduler.xml

The allocation file configures minimum shares, running job limits, weights and preemption timeouts for each pool.

- `pool` elements, which configure each pool. These may contain the following sub-elements:
 - `minMaps` and `minReduces`, to set the pool's minimum share of task slots.
 - `maxMaps` and `maxReduces`, to set the pool's maximum concurrent task slots.
 - `schedulingMode`, the pool's internal scheduling mode, which can be *fair* for fair sharing or *fifo* for first-in-first-out.
 - `maxRunningJobs`, to limit the number of jobs from the pool to run at once (defaults to infinite).
 - `weight`, to share the cluster non-proportionally with other pools. For example, a pool with weight 2.0 will get a 2x higher share than other pools. The default weight is 1.0.
 - `minSharePreemptionTimeout`, the number of seconds the pool will wait before killing other pools' tasks if it is below its minimum share (defaults to infinite).
 - `minSharePreemptionTimeout`, the number of seconds the pool will wait before killing other pools' tasks if it is below its minimum share (defaults to infinite).
 - `user` elements, which may contain a `maxRunningJobs` element to limit jobs.
 - `poolMaxJobsDefault`, which sets the default running job limit for any pools whose limit is not specified.
 - `userMaxJobsDefault`, which sets the default running job limit for any users whose limit is not specified.
 - `fairSharePreemptionTimeout`, which sets the preemption timeout used when jobs are below half their fair share.

An example allocation file is given below:

```
<?xml version="1.0"?>
<allocations>
    <pool name="kelly_pool">
        <minMaps>5</minMaps>
        <minReduces>5</minReduces>
        <maxMaps>25</maxMaps>
        <maxReduces>25</maxReduces>
        <minSharePreemptionTimeout>300</minSharePreemptionTimeout>
    </pool>
    <user name="kelly_user">
```

```
<maxRunningJobs>6</maxRunningJobs>
</use!><use!>
<useMaxJobsDefault>3</useMaxJobsDefault>
<useSharePreemptionTimeout>600</useSharePreemptionTimeout>
</use!>
```

This example creates a pool `kelly_pool` with a guarantee of 5 map slots and 5 reduce slots. The pool also has a minimum share preemption timeout of 300 seconds (5 minutes), meaning that if it does not get its guaranteed share within this time, it is allowed to kill tasks from other pools to achieve its share. The pool has a cap of 25 map and 25 reduce slots, which means that once 25 tasks are running, no more will be scheduled even if the pool's fair share is higher. The example also limits the number of running jobs per user to 3, except for `kelly_user`, who can run 6 jobs concurrently. Finally, the example sets a fair share preemption timeout of 600 seconds (10 minutes). If a job is below half its fair share for 10 minutes, it will be allowed to kill tasks from other jobs to achieve its share.

Administration

The fair scheduler provides support for administration at runtime through two mechanisms:

1. It is possible to modify minimum shares, limits, weights, preemption timeouts and pool scheduling modes at runtime by editing the allocation file. The scheduler will reload this file 10-15 seconds after it sees that it was modified.
2. Current jobs, pools, and fair shares can be examined through the JobTracker's web interface, at <http://<JobTracker URL>/scheduler>. On this interface, it is also possible to modify jobs' priorities or move jobs from one pool to another and see the effects on the fair shares (this requires JavaScript).

The following fields can be seen for each job on the web interface:

- *Submitted* - Date and time job was submitted.
- *JobID, User, Name* - Job identifiers as on the standard web UI.
- *Pool* - Current pool of job. Select another value to move job to another pool.
- *Priority* - Current priority. Select another value to change the job's priority
- *Maps/Reduces Finished*: Number of tasks finished / total tasks.
- *Maps/Reduces Running*: Tasks currently running.
- *Map/Reduce Fair Share*: The average number of task slots that this job should have at any given time according to fair sharing. The actual number of tasks will go up and down depending on how much compute time the job has had, but on average it will get its fair share amount.

In addition, it is possible to view an "advanced" version of the web UI by going to <http://<JobTracker URL>/scheduler?advanced>. This view shows two more columns:

- *Maps/Reduce Weight:* Weight of the job in the fair sharing calculations. This depends on priority and potentially also on job size and job age if the *sizebasedweight* and *NewJobWeightBooster* are enabled.

3. CapacityScheduler

The CapacityScheduler is designed to run Hadoop Map-Reduce as a shared, multi-tenant cluster in an operator-friendly manner while maximizing the throughput and the utilization of the cluster while running Map-Reduce applications.

The CapacityScheduler is designed to allow sharing a large cluster while giving each organization a minimum capacity guarantee. The central idea is that the available resources in the Hadoop Map-Reduce cluster are partitioned among multiple organizations who collectively fund the cluster based on computing needs. There is an added benefit that an organization can access any excess capacity not being used by others. This provides elasticity for the organizations in a cost-effective manner.

Sharing clusters across organizations necessitates strong support for multi-tenancy since each organization must be guaranteed capacity and safe-guards to ensure the shared cluster is impervious to single rogue job or user. The CapacityScheduler provides a stringent set of limits to ensure that a single job or user or queue cannot consume disproportionate amount of resources in the cluster. Also, the JobTracker of the cluster, in particular, is a precious resource and the CapacityScheduler provides limits on initialized/pending tasks and jobs from a single user and queue to ensure fairness and stability of the cluster.

The primary abstraction provided by the CapacityScheduler is the concept of *queues*. These queues are typically setup by administrators to reflect the economics of the shared cluster.

Features

The CapacityScheduler supports the following features:

- **Capacity Guarantees** - Support for multiple queues, where a job is submitted to a queue. Queues are allocated a fraction of the capacity of the grid in the sense that a certain capacity of resources will be at their disposal. All jobs submitted to a queue will have access to the capacity allocated to the queue. Administrators can configure soft limits and optional hard limits on the capacity allocated to each queue.
- **Security** - Each queue has strict ACLs which controls which users can submit jobs to individual queues. Also, there are safe-guards to ensure that users cannot view and/or modify jobs from other users if so desired. Also, per-queue and system administrator roles are supported.
- **Elasticity** - Free resources can be allocated to any queue beyond its capacity. When there is demand for these resources from queues running below capacity at a future point in time, as tasks scheduled on these resources complete, they will be assigned to jobs on queues running below the capacity. This ensures that resources are available in a

predictable and elastic manner to queues, thus preventing artificial silos of resources in the cluster which helps utilization.

- **Multi-tenancy -** Comprehensive set of limits are provided to prevent a single job, user and queue from monopolizing resources of the queue or the cluster as a whole to ensure that the system, particularly the JobTracker, isn't overwhelmed by too many tasks or jobs.
- **Operability -** The queue definitions and properties can be changed, at runtime, by administrators in a secure manner to minimize disruption to users. Also, a console is provided for users and administrators to view current allocation of resources to various queues in the system.
- **Resource-based Scheduling -** Support for resource-intensive jobs, wherein a job can optionally specify higher resource-requirements than the default, thereby accommodating applications with differing resource requirements. Currently, memory is the the resource requirement supported.
- **Job Priorities -** Queues optionally support job priorities (disabled by default). Within a queue, jobs with higher priority will have access to the queue's resources before jobs with lower priority. However, once a job is running, it will not be preempted for a higher priority job, *preemption* is on the roadmap is currently not supported.

Installation

set the following property in the Hadoop config file *HADOOP_CONF_DIR/mapred-site.xml* to have Hadoop use the fair scheduler:

```
<property>
  <name>mapred.jobtracker.taskScheduler</name>
  <value>org.apache.hadoop.mapred.CapacityTaskScheduler</value>
</property>
```

Once you restart the cluster, you can check that the fair scheduler is running by going to <http://<jobtracker URL>/scheduler> on the JobTracker's web UI.

Setting up queues

You can define multiple queues to which users can submit jobs with the CapacityScheduler. To define multiple queues, you should use the *mapred.queue.names* property *mapred-site.xml*. The CapacityScheduler can be configured with several properties for each queue that control the behavior of the Scheduler. This configuration is in the *capacity-scheduler.xml*.

Important settings for capacity-scheduler.xml

mapred.capacity-scheduler.queue.<queue-name>.capacity : Percentage of the number of slots in the cluster that are made to be available for jobs in this queue. The sum of capacities for all queues should be less than or equal 100.

mapred.capacity-scheduler.queue.<queue-name>.maximum-capacity : maximum-capacity defines a limit beyond which a queue cannot use the capacity of the cluster. This provides a means to limit how much excess capacity a queue can use. By default, there is no limit. The

maximum-capacity of a queue can only be greater than or equal to its minimum capacity. Default value of -1 implies a queue can use complete capacity of the cluster. This property could be to curtail certain jobs which are long running in nature from occupying more than a certain percentage of the cluster, which in the absence of pre-emption, could lead to capacity guarantees of other queues being affected. One important thing to note is that maximum-capacity is a percentage, so based on the cluster's capacity it would change. So if large no of nodes or racks get added to the cluster, maximum Capacity in absolute terms would increase accordingly.

mapred.capacity-scheduler.queue.<queue-name>.minimum-user-limit-percent : Each queue enforces a limit on the percentage of resources allocated to a user at any given time, if there is competition for them. This user limit can vary between a minimum and maximum value. The former depends on the number of users who have submitted jobs, and the latter is set to this property value. For example, suppose the value of this property is 25. If two users have submitted jobs to a queue, no single user can use more than 50% of the queue resources. If a third user submits a job, no single user can use more than 33% of the queue resources. With 4 or more users, no user can use more than 25% of the queue's resources. A value of 100 implies no user limits are imposed.

mapred.capacity-scheduler.queue.<queue-name>.supports-priority : If true, priorities of jobs will be taken into account in scheduling decisions.

Capacity scheduler lazily initializes the jobs before they are scheduled, for reducing the memory footprint on jobtracker. Following are the parameters, by which you can control the initialization of jobs per-queue.

mapred.capacity-scheduler.maximum-system-jobs : Maximum number of jobs in the system which can be initialized, concurrently, by the CapacityScheduler. Individual queue limits on initialized jobs are directly proportional to their queue capacities.

mapred.capacity-scheduler.queue.<queue-name>.maximum-initialized-active-tasks : The maximum number of tasks, across all jobs in the queue, which can be initialized concurrently. Once the queue's jobs exceed this limit they will be queued on disk.

mapred.capacity-scheduler.queue.<queue-name>.maximum-initialized-active-tasks-per-user : The maximum number of tasks per-user, across all the of the user's jobs in the queue, which can be initialized concurrently. Once the user's jobs exceed this limit they will be queued on disk.

mapred.capacity-scheduler.queue.<queue-name>.init-accept-jobs-factor : The multiple of (maximum-system-jobs * queue-capacity) used to determine the number of jobs which are accepted by the scheduler. The default value is 10. If number of jobs submitted to the queue exceeds this limit, job submission are rejected.

Example

Flat No. 212, 2nd Floor, Annapurna Block, Aditya Enclave, Ameerpet, Hyd.
E-mail: info@kellytechno.com www.kellytechno.com 040-6462 6789; 0998 570 6789.



```
<?xml version="1.0"?>

<configuration>

    <!-- system limit, across all queues -->

    <property>
        <name>mapred.capacity-scheduler.maximum-system-jobs</name>
        <value>3000</value>
        <description>Maximum number of jobs in the system which can be initialized, concurrently, by the CapacityScheduler.</description>
    </property>

    <!-- queue: kellyA -->
    <property>
        <name>mapred.capacity-scheduler.queue.kellyA.capacity</name>
        <value>80</value>
    </property>
    <property>
        <name>mapred.capacity-scheduler.queue.kellyA.supports-priority</name>
        <value>false</value>
    </property>
    <property>
        <name>mapred.capacity-scheduler.queue.kellyA.minimum-user-limit-percent</name>
        <value>20</value>
    </property>

    <!-- queue: kellyB -->
    <property>
        <name>mapred.capacity-scheduler.queue.kellyB.capacity</name>
        <value>20</value>
    </property>
    <property>
        <name>mapred.capacity-scheduler.queue.kellyB.supports-priority</name>
        <value>false</value>
    </property>
    <property>
        <name>mapred.capacity-scheduler.queue.kellyB.minimum-user-limit-percent</name>
        <value>20</value>
    </property>
    <property>
        <name>mapred.capacity-scheduler.queue.kellyB.user-limit-factor</name>
        <value>1</value>
    </property>
```

```
<property>
  <name>mapred.capacity-scheduler.queue.kellyB.maximum-initialized-active-
  tasks</name>
  <value>200000</value>
</property>
<property>
  <name>mapred.capacity-scheduler.queue.kellyB.maximum-initialized-active-tasks-per-
  user</name>
  <value>100000</value>
</property>
```

Compression

Compression reduces the number of bytes written to/read from HDFS.

Compression effectively improves the efficiency of network bandwidth and disk space. This saves the amount of data being transferred between map nodes to reduce nodes.

LZO

LZO is a compression/decompression library. It has following characteristics

- Very fast decompression
- Requires an additional buffer during compression (size of 8KB or 64KB, depends on compression level)
- Requires no additional memory for decompression other than source and destination buffers.
- Allows the user to adjust the balance between compression ratio and compression speed, without affecting the speed of decompression.

Snappy

Snappy is a compression/decompression library. It does not aim for maximum compression, or compatibility with any other compression library; instead, it aims for very high speeds and reasonable compression.

For instance, compared to the fastest mode of zlib, Snappy is an order of magnitude faster for most inputs, but the resulting compressed files are anywhere from 20% to 100% bigger. Typical compression ratios are about 1.5-1.7x for plain text, about 2-4x for HTML, and of course 1.0x for JPEGs, PNGs and other already-compressed data. Similar numbers for zlib in its fastest mode are 2.6-2.8x, 3-7x and 1.0x, respectively.

CompressionCodecs are implementation compression algorithms. To use either LZO or Snappy compression codecs, LZO or snappy libraries must be installed on all the nodes in the cluster.

Hadoop provides Compression Codecs for Gzip, LZO and Snappy and are

`org.apache.hadoop.io.compress.DefaultCodec`
`org.apache.hadoop.io.compress.LzoCodec`
`org.apache.hadoop.io.compress.SnappyCodec`

By Default Compression is not enabled in Mapreduce.Edit the following properties in mapred-site.xml to achieve compression

If job outputs are to be compressed set value as true.

```
<property>
<name>mapred.output.compress</name>
<value>false</value>
<description>Should the job outputs be compressed?
</description>
</property>
```

Which compression codec to be used while compressing job output. By default DefaultCodec uses. To use lzo or snappy replace with their corresponding codecs.

```
<property>
<name>mapred.output.compression.codec</name>
<value>org.apache.hadoop.io.compress.DefaultCodec</value>
<description>If the job outputs are compressed, how should they be compressed?
</description>
</property>
```

If map outputs are to be compressed set the value as true

```
<property>
<name>mapred.compress.map.output</name>
<value>false</value>
<description>Should the outputs of the maps be compressed before being
sent across the network. Uses SequenceFile compression.
</description>
</property>
```

Specify which compression codec to be used while compressing map output

```
<property>
<name>mapred.map.output.compression.codec</name>
<value>org.apache.hadoop.io.compress.DefaultCodec</value>
<description>If the map outputs are compressed, how should they be
compressed?
</description>
</property>
```

By Default Hadoop comes with the following CompressionCodecs. If you want to use any other codec apart from these, mention the same in **core-site.xml**

```
<property>
<name>io.compression.codecs</name>
<value>org.apache.hadoop.io.compress.DefaultCodec,
      org.apache.hadoop.io.compress.GzipCodec,
      org.apache.hadoop.io.compress.BZip2Codec
</value>
<description>A list of the compression codec classes that can be used for
      compression/decompression.</description>
</property>
```

Setting Up a Hadoop Cluster

The following are prerequisites for installing Hadoop.
Note step1 is mandatory, but recommended.

1. Create a user 'hadoop' in your machine. (If not present)

To test whether 'hadoop' user present in your machine, enter the following command at the console

```
$ su hadoop
```

If there is no 'hadoop' user in your machine, it gives

Unknown id: hadoop

If it is already present just ignore the rest of the point. Otherwise carry on

Enter the following commands at console to create hadoop user.

```
$ sudo adduser hadoop
Adding user `hadoop' ...
Adding new group `hadoop' (1001) ...
Adding new user `hadoop' (1001) with group `hadoop' ...
Creating home directory /home/hadoop ...
Copying files from /etc/skel ...
Enter new UNIX password: <give hadoop as password>
Retype new UNIX password: <just enter hadoop again>
passwd: password updated successfully
Changing the user information for hadoop
Enter the new value, or press ENTER for the default
  Full Name []: <Just Press enter>
  Room Number []: <Just Press enter>
  Work Phone []: <Just Press enter>
  Home Phone []: <Just Press enter>
  Other []: <Just Press enter>
Is the information correct? [Y/n] Y
```

2. Install Java (Java 6 or later is required to run Hadoop.)

Check whether java is already installed or not with the following command

\$java

If console comes back with

Usage: java [-options] class [args...]

(to execute a class)

It means java is installed in your machine.

Otherwise install java with the following command (considering the UNIX flavor is Ubuntu)

\$ sudo apt-get install openjdk-6-jdk

The latest stable Sun JDK is preferred option.

If you install java with above command JAVA_HOME would be /usr/lib/jvm/java-6-openjdk".

3. Install openssh-server and openssh-client with the following commands

\$ sudo apt-get install openssh-server

\$ sudo apt-get install openssh-client

4. Go to hadoop user with the following command

\$ su hadoop

Password:hadoop <we have given this earlier>

5. Download Hadoop from the Apache Hadoop releases page (<http://hadoop.apache.org/core/releases.html>) and unpack the contents of the distribution in a directory such as "/home/hadoop/".

All the following commands has to be executed while you are at hadoop user.

6. Untar the hadoop distribution in /home/hadoop with the following command

\$ sudo tar xzf <Path of downloaded directory>/hadoop-1.0.3.tar.gz

Finally you will be able find /home/hadoop/hadoop-1.0.3 directory. Here onwards we refer this directory as \$HADOOP_HOME.

7. Change the owner of the Hadoop files to be the hadoop user and group.

\$sudo chown -R hadoop:hadoop hadoop-1.0.3

8. SSH configuration (For passwordless)

\$ ssh-keygen -t rsa

Just press on enter, you are prompted to enter any thing. Finally you find the following

Generating public/private rsa key pair.

Enter file in which to save the key (/home/hadoop/.ssh/id_rsa):

Created directory '/home/hadoop/.ssh'.

Enter passphrase (empty for no passphrase):

Enter same passphrase again:

Your identification has been saved in /home/hadoop/.ssh/id_rsa.

Your public key has been saved in /home/hadoop/.ssh/id_rsa.pub.

The key fingerprint is:

43:02:31:17:b8:59:a9:a2:38:d0:79:3f:12:76:29:5a hadoop@kelly-desktop

The key's randomart image is:

+--[RSA 2048]----+

```
| +000 |  
| .+o |  
| ..=.. |  
| .+ E oo |  
| o. *= S |  
| +...o . |  
| . .. |  
| | |  
+-----+
```

9. Go to .ssh directory. You will find 2 files in that directory namely

id_rsa id_rsa.pub

10. Enter the following command to place the public key in authorized_keys

\$ cat id_rsa.pub >> authorized_keys

Pseudo Distributed Mode

11. Go to HADOOP_HOME/conf directory.

12. Find the host name of your machine with the following command

\$ hostname

hadoop-desktop

In my machine "hadoop-desktop" is my host name.

13. Edit the core-site.xml with the following properties.

```
<configuration>  
  <property>  
    <name>default.name</name>  
    <value><hostname of your machine>:9000</value>  
  </property>  
</configuration>
```

Description

fs.default.name : The property `fs.default.name` is a HDFS filesystem URI, whose host is the namenode's hostname or IP address, and port is the port that the namenode will listen on for RPCs. If no port is specified, the default of 8020 is used.

14. Edit the `hdfs-site.xml` with the following properties.

```
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
  <property>
    <name>dfs.data.dir</name>
    <value>/home/hadoop/dfs/name</value>
  </property>
  <property>
    <name>dfs.name.dir</name>
    <value>/home/hadoop/dfs/data</value>
  </property>
</configuration>
```

Description

dfs.replication : Each block of the file is stored in number of servers specified for this property.

If the value of this property is 3, then each block of the file is stored in 3 DataNodes.

dfs.data.dir : This property specifies a list of the directories for a datanode to store its blocks. A datanode round robin writes between its storage directories.

dfs.name.dir: This property specifies a list of the directories where the namenode stores persistent filesystem metadata (the edit log, and filesystem image). A copy of each of the metadata file is stored in each directory for redundancy.

15. Edit the `mapred-site.xml` with the following properties.

```
<configuration>
  <property>
    <name>mapred.job.tracker</name>
    <value><hostname of your machine>:9001</value>
  </property>
  <property>
    <name>mapred.local.dir</name>
    <value>/home/hadoop/mapred/local</value>
  </property>
  <property>
    <name>mapred.system.dir</name>
    <value>/hadoop/mapred/system</value>
  </property>
</configuration>
```

```
</property>  
</configuration>
```

Description

mapred.job.tracker : Specify the hostname or IP address and port that the job tracker will listen on, Default port is 8021.

mapred.local.dir : This property specifies the list of directories separated by commas. During MapReduce job, intermediate data and working files are written to temporary local files. Since this data includes the potentially very large output of map tasks, you need to ensure the mapred.local.dir property which controls the location of local temporary storage.

mapred.system.dir : MapReduce uses a distributed file system to share files (such as the job JAR file) with the tasktrackers that run the MapReduce tasks. This property is used to specify a directory where these files can be stored.

16. Edit the masters file by replacing localhost with host name of your machine.

17. Edit the slaves file by replacing localhost with host name of your machine.

18. Edit the hadoop-env.sh by uncommenting the below line and replacing the JAVA_HOME of your machine (# is comment in shell script)

```
export JAVA_HOME=/usr/lib/jvm/java-6-openjdk
```

19. Go to HADOOP_HOME/bin directory

20. Format the Namenode with the following command.

```
./hadoop namenode format
```

Caution: Format the namenode only when are setting up the cluster first time.

If you format the namenode, when there is any data in HDFS, you will loose all the data in HDFS. So don't format the namenode. If you stop the Hadoop cluster and while restarting the cluster you will find all the data in HDFS.

21. Start the Hadoop cluster with the following command.

```
./start-all.sh
```

The above command starts Namenode, DataNode, Secondary Namenode, JobTracker and TaskTracker.

22. To confirm the above, just type the following command in the CLI.

```
$ jps
```

The above command lists all the java processes started by the user.

23. There are separate scripts to start/stop HDFS, MapReduce and are

start-dfs.sh -- To start HDFS only

stop-dfs.sh -- To stop HDFS only

start-mapred.sh -- To start MapReduce only

stop-mapred.sh -- To stop MapReduce only
stop-all.sh -- To stop entire Hadoop Cluster.

Cluster Mode

24. Perform the steps 1-9 on every machine, which we call node here onwards.
25. Designate one machine as master and others as slaves.
26. Assuming hostname of the master be master and slaves be slave1, slave2, slave3 etc.
27. Copy the id_rsa.pub (public key) to all the slave nodes. Execute the following command on master to copy the id_rsa.pub of the master machine.

```
$ ssh-copy-id hadoop@slave1  
$ ssh-copy-id hadoop@slave2  
$ ssh-copy-id hadoop@slave3
```

(Assuming 3 node cluster)

The above command copies the id_rsa.pub to authorized_keys of slave1, slave2, slave3.

The following steps has to be performed on all machines (both master and slaves)

28. Edit the masters file by replacing localhost with master on all the machines (both master and slaves)

29. Edit the slaves file by listing all slaves in the cluster (mention hostname of all the slaves)

Example: slaves file entries would be

```
slave1  
slave2  
slave3
```

30. Perform steps 11, 13, 14 and 15. Replace 1 by 3 in dfs.replication property in hdfs-site.xml.

31. Perform steps from 18 to 22.

References

Hadoop: The Definitive Guide by Tom White, First Edition

<http://code.google.com/p/snappy/>

<http://en.wikipedia.org/wiki/Lempel%E2%80%93Ziv%E2%80%93Oberhumer>

http://hadoop.apache.org/docs/r1.0.3/fair_scheduler.html

http://hadoop.apache.org/docs/r1.0.3/capacity_scheduler.html

MRV2 Source code set-up and debugging

- 1) To build mrv2 we need to install protbuff and maven.

Once you install both of them set the following

```
root@kelly1:/usr/local/apache-maven# export M2_HOME=/usr/local/apache-maven/apache-maven-3.0.3
```

```
root@kelly1:/usr/local/apache-maven# export M2=$M2_HOME/bin
```

```
root@kelly1:/usr/local/apache-maven# export MAVEN_OPTS="-Xms256m -Xmx512m"
```

```
root@kelly1:/usr/local/apache-maven# export PATH=$M2:$PATH
```



```
root@kelly1:/usr/local/apache-maven# export JAVA_HOME=/usr/lib/jvm/java-1.6.0-openjdk
root@kelly1:/usr/local/apache-maven# export PATH=$JAVA_HOME/bin:$PATH
root@kelly1:/usr/local/apache-maven# mvn --version
Apache Maven 3.0.3 (r1075438; 2011-02-28 23:01:09+0530)
Maven home: /usr/local/apache-maven/apache-maven-3.0.3
Java version: 1.6.0_23, vendor: Sun Microsystems Inc.
Java home: /usr/lib/jvm/java-6-openjdk/jre
Default locale: en_IN, platform encoding: UTF-8
OS name: "linux", version: "3.0.0-12-generic", arch: "amd64", family: "unix"

export LD_LIBRARY_PATH=/usr/local/lib
```

For protobuf installation you can also reffer following links:

<http://svn.apache.org/repos/asf/hadoop/common/trunk/hadoop-mapreduce-project/hadoop-yarn/README>

<http://svn.apache.org/repos/asf/hadoop/common/trunk/hadoop-mapreduce-project/INSTALL>

2) clone the hadoop

```
mkdir mr2
```

```
cd mr2
```

```
svn checkout http://svn.apache.org/repos/asf/hadoop/common/trunk/ hadoop-trunk
```

3) After cloning the trunk: mvn clean install -DskipTests mvn eclipse:eclipse mvn clean package -Pdist -Dtar -DskipTests

4) General settings for MR2

```
export JAVA_HOME=/usr/lib/jvm/java-1.6.0-openjdk
export PATH=$JAVA_HOME/bin:$PATH

export HADOOP_DEV_HOME=`pwd`/hadoop-trunk
export HADOOP_MAPRED_HOME=${HADOOP_DEV_HOME}
export HADOOP_COMMON_HOME=${HADOOP_DEV_HOME}
export HADOOP_HDFS_HOME=${HADOOP_DEV_HOME}
export YARN_HOME=${HADOOP_DEV_HOME}
export HADOOP_CONF_DIR=${HADOOP_DEV_HOME}/conf/
export YARN_CONF_DIR=~${HADOOP_DEV_HOME}/conf/
export LD_LIBRARY_PATH=/usr/local/lib
```

5) You can use the attched configurations for reference

6) If everything is proper then now you should be able start hadoop
hadoop-0.24.0-SNAPSHOT\$ bin/hadoop namenode -format

And start the HDFS services:

```
hadoop-0.24.0-SNAPSHOT$ sbin/hadoop-daemon.sh start namenode  
hadoop-0.24.0-SNAPSHOT$ sbin/hadoop-daemon.sh start datanode
```

And the new MR2 services

```
hadoop-0.24.0-SNAPSHOT$ bin/yarn-daemon.sh start resourcemanager  
hadoop-0.24.0-SNAPSHOT$ bin/yarn-daemon.sh start nodemanager  
hadoop-0.24.0-SNAPSHOT$ bin/yarn-daemon.sh start historyserver
```

Make sure all needed services are running:

```
hadoop-0.24.0-SNAPSHOT$ jps
```

```
Jps  
ResourceManager  
JobHistoryServer  
NodeManager  
DataNode  
NameNode
```

you can also reffer <http://www.cloudera.com/blog/2011/11/building-and-deploying-mr2/>

7) For contribution:

Please refer <http://wiki.apache.org/hadoop/HowToContribute>

For debugging:

As you are already built the mapreduce project, you can add debug conf's in yarn-env.sh and put debug points in the code and start analyzing it.

1) Copy yarn-env.sh(hadoop-dist/target/hadoop-0.24.0-SNAPSHOT/etc/hadoop/yarn-env.sh) into YARN_CONF_DIR

2) Add following debug conf's

```
YARN_RESOURCEMANAGER_OPTS="-Xdebug  
Xrunjdwp:transport=dt_socket,server=y,address=5000"  
YARN_NODEMANAGER_OPTS="-Xdebug  
Xrunjdwp:transport=dt_socket,server=y,address=5001"  
YARN_JOB_HISTORYSERVER_OPTS="-Xdebug -  
Xrunjdwp:transport=dt_socket,server=y,address=5002"
```

3) In eclipse, go to corresponding class file for example ResourceManager.java -> debug as ->
Remote Java Application ----> new ---> set corresponding port

Now it's ready to debug.



For testing your patch you can set-up separate testbed.

For contribution: <http://wiki.apache.org/hadoop/HowToContribute>

References

Hadoop definite guide

HBase definite Guide

The Hadoop Distributed FileSystem paper(<http://dl.acm.org/citation.cfm?id=1914427>)

Data-Intensive Text Processing with
MapReduce(<http://www.morganclaypool.com/doi/abs/10.2200/S00274ED1V01Y201006HLT007>)

Scheduling for real-time mobile
MapReducesystems(<http://dl.acm.org/citation.cfm?id=2002305&dl=ACM&coll=DL&CFID=113190881&CFTOKEN=92325790>)

Apache Pig

Introduction

- Pig is a high-level procedural language for querying large semi-structured data sets using Hadoop and MapReduce platform.

- Pig is a simple language that executes statements.
- A statement is an operation that takes input (such as bag, which represents a set of tuples) and emits another bag as its output.
- Pig follows the following format while executing a script
- Read the data from the file system.
- Perform a number of operations on the data.
- Store the final result back to file system.

Ques Example

Pig Code

```
lines = LOAD '/input/sample.txt';
hadoopLines = FILTER lines BY $0 MATCHES '/hadoop+';
STORE hadoopLines INTO '/output/cleanedLines';
```

Explanation

- Read the '/input/sample.txt' into a bag 'sample' that represents a collection of tuples.
- Filter each tuple (represented as \$0) with a regular expression, looking for the character sequence 'hadoop'.
- Store the 'hadoopLines' bag, which contains all of those tuples from '/input/sample.txt' that contain 'hadoop' into a new file called '/output/cleanedLines'

Installing Pig

- Download a recent stable release from Apache download Mirrors.
- Unpack the downloaded Pig installation.
- Export HADOOP_HOME as Hadoop Home Directory.
- export HADOOP_HOME=<path-to-hadoop>

- Go to Pig directory.
- Test the Pig installation with the simple command
- \$<Pig_HOME>/bin/pig -help.

Running Pig

Pig has two execution modes or exec types:

- Local Mode – To run Pig in local mode, you need access to a single machine; all files are installed and run using your local host and file system. Specify local mode using -x flag

\$ pig -x local

- Mapreduce Mode – To run Pig in mapreduce mode, you need access to a Hadoop cluster and HDFS installation. Mapreduce mode is the default mode.
- Pig commands can be executed using the following modes
- Interactive Mode
- Batch Mode

Interactive Mode & Batch Mode Execution

- To execute the pig statements in interactive mode, invoke the Grunt shell by typing the "pig" command.
- Enter the pig statements interactively at the grunt prompt.
- Running Grunt Example in interactive mode

```
grunt> lines = LOAD '/input/sample.txt';
grunt> hadoopLines = FILTER lines BY $0 MATCHES '^.*hadoop.*$';
2012-09-26 22:19:05,548 [main] WARN  org.apache.pig.PigServer - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
grunt> STORE hadoopLines INTO '/output/cleanedLines';
2012-09-26 22:19:30,128 [main] WARN  org.apache.pig.PigServer - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
2012-09-26 22:19:30,211 [main] WARN  org.apache.pig.PigServer - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
```

- Save the Pig latin statements in a file run it as below.

\$ bin/pig <file-name>.pig

Pig Operators

- LOAD – Loads data from the file system or other storage into a relation
- STORE – Saves a relation to the file system or other storage
- DUMP – Prints a relation to the console
- FILTER – Removes unwanted rows from a relation
- DISTINCT – Removes duplicate rows from a relation
- FOREACH...GENERATE – Adds or removes fields from a relation
- STREAM – Transforms a relation using an external program
- JOIN – Joins two or more relations
- COGROUP – Groups the data in two or more relations
- GROUP – Groups the data in a single relation
- CROSS – Creates the cross product of two or more relations
- ORDER – Sorts a relation by one or more fields
- LIMIT -- Limits the size of a relation to a maximum number of tuples
- UNION – Combines two or more relations into one
- SPLIT – Splits a relation into two or more relations.
-

Utility Functions

- AVG – Calculates the average value of entries in a bag.
- CONCAT – Concatenates two byte arrays or two character arrays together.
- COUNT – Calculates the number of entries in a bag.
- MAX – Calculates the maximum value of entries in a bag.
- MIN – Calculates the minimum value of entries in a bag.

- **SUM** – Calculates the sum of the values of entries in a bag.
- **TOKENIZE** – Tokenizes a character array into a bag of its constituent words.
- **PigStorage** – Loads or stores relations using a field-delimited text format. Each line is broken into fields using a configurable field delimiter (defaults to a tab character) to be stored in the tuple's fields. It is the default storage when none is specified.
- **BinStorage** – Loads or stores relations from or to binary files. An internal Pig format is used that uses Hadoop Writable objects.
- **TextLoader** – Loads relations from a plain-text format. Each line corresponds to a tuple whose single field is the line of text.
-

Example using PigStorage & foreach

```
$ pig -x local

grunt> passwdLines = load '/etc/passwd' using PigStorage':' As
(user:chararray,a:chararray,b:chararray,c:chararray,d:chararray,e:chararray,shell:chararray);

grunt> dump passwdLines;

grunt> userShell = foreach passwdLines generate user,shell;
grunt> dump userShell;
```

Example using GROUP BY, ORDER BY & AVG

Example to find the avg click for each URL.

```
$pig
grunt> urls = LOAD '/input/urlcount.txt' AS (url:chararray, count:int);
grunt> urlCount = GROUP urls BY url;
grunt> urlAvg = foreach urlCount generate group, AVG(urls.count) as urlC;
grunt> d = ORDER urlAvg BY urlC DESC;
```

```
grunt> STORE d INTO '/output/urlAvg' using PigStorage('#');
```

GROUP BY groups all the same urls together.

group is a reserved word in the pig to identify each group.

ORDER BY orders the urls by their average in the descending order. If we don't mention DESC, it will be ordered in the ascending order.

PARALLEL

- PARALLEL clause increases the parallelism of a job
- PARALLEL sets the number of reduce tasks for the MapReduce jobs generated by Pig. The default value is 1.
- PARALLEL only affects the number of reduce tasks. Map parallelism is determined by the input file, one map for each Input Split.
- If you don't specify PARALLEL, you still get the same map parallelism but only one reduce task.
- Specify the PARALLEL clause with any operator that starts a reduce phase, which includes COGROUP, CROSS, DISTINCT, GROUP, JOIN (inner), JOIN (outer), and ORDER.

PARALLEL - An Example

Let us consider the avg click for each URL pig script, with PARALLEL clause.

```
$ pig
grunt> urls = LOAD '/input/urlcount.txt' AS (url:chararray, count:int);
grunt> urlCount = GROUP urls BY url PARALLEL 3;
urlCount generate group, AVG(urls.count) as urlC;
grunt> d = ORDER urlAvg BY urlC DESC;
grunt> STORE d INTO '/output/urlAvg' using PigStorage('#');
```

For the above example 3 reducers will be spawned.

- UDFs are required to define custom processing.
- UDFs can be written in the Java, Python, JavaScript and Ruby and permit Pig to support custom processing.
- Support for writing UDFs in Python, JavaScript and Ruby still evolving.
- UDFs provide an opportunity to extend Pig into your application domain.
-

How to Write Java UDF

UDFs can be developed by extending EvalFunc class and overriding exec method.

Example : This UDF replaces a given string with another string.

```
package kelly.training.pig.udf;  
  
import java.io.IOException;  
  
import org.apache.hadoop.conf.Configuration;  
  
import org.apache.pig.EvalFunc;  
  
import org.apache.pig.data.Tuple;  
  
import org.apache.pig.map.util.UDFContext;  
  
public class Transform extends EvalFunc<String> {  
    public String exec(Tuple input) throws IOException {  
        if (input == null || input.size() == 0) {  
            return null;  
        }  
    }  
}
```

```
Configuration conf = UDFContext.getUDFContext().getJobConf();
```

```
String from = conf.get("replace.string");
```

```
if (from == null) {  
    throw new IOException("replace.string should not be null");  
}  
  
String to = conf.get("replace.by.string");  
  
if (to == null) {  
    throw new IOException("replace.by.string should not be null");  
}  
  
try {  
    String str = (String) input.get(0);  
    return str.replace(from, to);  
} catch (Exception e) {  
    throw new IOException("Caught exception processing input row", e);  
}
```

```
}
```

```
}
```

```
}
```

```
-- sed.pig
```

```
-- pig-training.jar contains custom udfs developed by a user. It has to be registered
```

```
before running the pig script.
```

```
REGISTER pig-training.jar;
```

```
A = LOAD '/input/sample.txt' AS (line: chararray);
```

```
B = FOREACH A GENERATE kelly.training.pig.udf.Transform(line);
```

```
STORE B INTO '/output/sed';
```

pig-training.jar contains the Transform UDF. These UDFs jar has to be registered.

While using custom UDF, mention UDF with package name.

How to Run the Sed Example using UDF

\$ pig -x local -P user.properties sed.pig

Or

\$ pig -P user.properties sed.pig

user.properties is properties file, which contains the following parameters

replace.string=

replace.by.string=

While running the sed.pig, mention the path of the properties file using -P option.

~~Kelly Technologies~~ sed.pig:

```
-- sed.pig
-- pig-training.jar contains custom udfs developed by a user. It has to be registered
-- before running the pig script.
REGISTER /home/kelly1/Work/Hadoop\ training/pig_training/pig-training.jar;
-- A = LOAD '/input/sample.txt' AS (line: chararray);
A = LOAD '/home/kelly1/Work/Hadoop\ training/pig_training/sample.txt' AS (line:
chararray);
B = FOREACH A GENERATE kelly.training.pig.udf.Transform(line);
-- STORE B INTO '/output/sed';
STORE B INTO '/home/kelly1/Work/sed';
```

~~Kelly Technologies~~ user.properties:

```
replace.string=hadoop
replace.by.string=Hadoop
```

~~Kelly Technologies~~ Transform.java:

```
package kelly.training.pig.udf;

import java.io.IOException;

import org.apache.hadoop.conf.Configuration;
import org.apache.pig.EvalFunc;
import org.apache.pig.data.Tuple;
import org.apache.pig.impl.util.UDFContext;

public class Transform extends EvalFunc<String> {
    public String exec(Tuple input) throws IOException {
        if (input == null || input.size() == 0) {
            return null;
        }
        Configuration conf = UDFContext.getUDFContext().getJobConf();
        String from = conf.get("replace.string");
        if (from == null) {
            throw new IOException("replace.string should not be null");
        }
        String to = conf.get("replace.by.string");
        if (to == null) {
            throw new IOException("replace.by.string should not be null");
        }
        try {
            String str = (String) input.get(0);
            return str.replace(from,to);
        } catch (Exception e) {
            throw new IOException("Caught exception processing input row", e);
        }
    }
}
```

sample.txt

hadoop is a framework to process big data.

hadoop has two components namely hdfs and mapreduce.

hdfs is a distributed data storage framework whereas mapreduce is distributed data processing framework.

Apache Pig

Introduction

- Pig is a high-level procedural language for querying large semi-structured data sets using Hadoop and MapReduce platform.
- Pig is a simple language that executes statements.
- A statement is an operation that takes input (such as bag, which represents a set of tuples) and emits another bag as its output.
- Pig follows the following format while executing a script
- Read the data from the file system.
- Perform a number of operations on the data.
- Store the final result back to file system.

Grep Example

Pig Code

```
lines = LOAD '/input/sample.txt';
hadoopLines = FILTER lines BY $0 MATCHES 'hadoop+.*';
STORE hadoopLines INTO '/output/cleanedLines';
```

Explanation

- Read the '/input/sample.txt' into a bag 'sample' that represents a collection of tuples.
- Filter each tuple (represented as \$0) with a regular expression, looking for the character sequence hadoop.
- Store the 'hadoopLines' bag, which contains all of those tuples from '/input/sample.txt' that contain 'hadoop' into a new file called '/output/cleanedLines'

Installing Pig

- Download a recent stable release from Apache download Mirrors.
- Unpack the downloaded Pig installation.
- Export HADOOP_HOME as Hadoop Home Directory.

- export HADOOP_HOME=<path-to-hadoop>
- Go to Pig directory.
- Test the Pig installation with the simple command
- \$<Pig_HOME>/bin/pig -help.

Running Pig

Pig has two execution modes or exec types:

- Local Mode – To run Pig in local mode, you need access to a single machine, all files are installed and run using your local host and file system. Specify local mode using -x flag

\$ pig -x local

- Mapreduce Mode – To run Pig in mapreduce mode, you need access to a Hadoop cluster and HDFS installation. Mapreduce mode is the default mode.
- Pig commands can be executed using the following modes
- Interactive Mode
- Batch Mode

Interactive Mode & Batch Mode Execution

- To execute the pig statements in interactive mode, invoke the Grunt shell by typing the "pig" command.
- Enter the pig statements interactively at the grunt prompt.

Running Grunt Example in interactive mode

```
grunt> lines = LOAD '/input/sample.txt';
grunt> hadoopLines = FILTER lines BY $0.MATCHES '^hadoop+.*';
2012-09-26 22:19:05,548 [main] WARN  org.apache.pig.PigServer - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
grunt> STORE hadoopLines INTO '/output/cleanedLines';
2012-09-26 22:19:30,128 [main] WARN  org.apache.pig.PigServer - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
2012-09-26 22:19:30,128 [main] WARN  org.apache.pig.PigServer - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
```

- Save the Pig latin statements in a file run it as below.

\$ bin/pig <file-name>.pig

Pig Latin Functions

- LOAD – Loads data from the file system or other storage into a relation
- STORE – Saves a relation to the file system or other storage
- DUMP – Prints a relation to the console
- FILTER – Removes unwanted rows from a relation
- DISTINCT – Removes duplicate rows from a relation
- FOREACH...GENERATE – Adds or removes fields from a relation
- STREAM – Transforms a relation using an external program
- JOIN – Joins two or more relations
- COGROUP – Groups the data in two or more relations
- GROUP – Groups the data in a single relation
- CROSS – Creates the cross product of two or more relations
- ORDER – Sorts a relation by one or more fields
- LIMIT -- Limits the size of a relation to a maximum number of tuples
- UNION – Combines two or more relations into one
- SPLIT – Splits a relation into two or more relations.
-

Built-in Functions

- AVG – Calculates the average value of entries in a bag.
- CONCAT – Concatenates two byte arrays or two character arrays together.
- COUNT – Calculates the number of entries in a bag.
- MAX – Calculates the maximum value of entries in a bag.

- MIN – Calculates the minimum value of entries in a bag.
- SUM – Calculates the sum of the values of entries in a bag.
- TOKENIZE – Tokenizes a character array into a bag of its constituent words.
- PigStorage – Loads or stores relations using a field-delimited text format. Each line is broken into fields using a configurable field delimiter (defaults to a tab character) to be stored in the tuple's fields. It is the default storage when none is specified.
- BinStorage – Loads or stores relations from or to binary files. An internal Pig format is used that uses Hadoop Writable objects.
- TextLoader – Loads relations from a plain-text format. Each line corresponds to a tuple whose single field is the line of text.
-

Example Using PigStorage & foreach

```
$ pig -x local
grunt> passwdLines = load '/etc/passwd' using PigStorage':' As
(user:chararray,a:chararray,b:chararray,c:chararray,d:chararray,e:chararray,shell:chararray);
grunt> dump passwdLines
grunt> userShell = foreach passwdLines generate user,shell;
grunt> dump userShell;
```

Example using GROUP BY, ORDER BY & AVG

Example to find the avg click for each URL.

```
$ pig
grunt> urls = LOAD '/input/urlcount.txt' AS (url:chararray, count:int);
grunt> urlCount = GROUP urls BY url;
grunt> urlAvg = foreach urlCount generate group, AVG(urls.count) as urlC;
```

```
grunt> d = ORDER urlAvg BY urlC DESC;
grunt> STORE d INTO '/output/urlAvg' using PigStorage('#');
```

GROUP BY groups all the same urls together.

group is a reserved word in the pig to identify each group.

ORDER BY orders the urls by their average in the descending order. If we don't mention **DESC** it will be ordered in the ascending order.

PARALLEL

- PARALLEL clause increases the parallelism of a job
- PARALLEL sets the number of reduce tasks for the MapReduce jobs generated by Pig. The default value is 1.
- PARALLEL only affects the number of reduce tasks. Map parallelism is determined by the input file, one map for each Input Split.
- If you don't specify PARALLEL, you still get the same map parallelism but only one reduce task.
- Specify the PARALLEL clause with any operator that starts a reduce phase, which includes COGROUP, CROSS, DISTINCT, GROUP, JOIN (inner), JOIN (outer), and ORDER.

PARALLEL - An Example

Let us consider the avg click for each URL pig script, with PARALLEL clause.

```
grunt> urls = LOAD '/input/urlcount.txt' AS (url:chararray, count:int);
grunt> urlCount = GROUP urls BY url PARALLEL 3;
      urlCount generate group, AVG(urls.count) as urlC;
      grunt> d = ORDER urlAvg BY urlC DESC;
      grunt> STORE d INTO '/output/urlAvg' using PigStorage('#');
```

For the above example 3 reducers will be spawned.

- ➔ UDFs are required to define custom processing.
- ➔ UDFs can be written in the Java, Python, JavaScript and Ruby and permit Pig to support custom processing.
- ➔ Support for writing UDFs in Python, JavaScript and Ruby still evolving.
- ➔ UDFs provide an opportunity to extend Pig into your application domain.
- ➔

UDFs can be developed by extending EvalFunc class and overriding exec method.

Example : This UDF replaces a given string with another string.

```
package kelly.training.pig.udf;
import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.pig.EvalFunc;
import org.apache.pig.data.Tuple;
import org.apache.pig.impl.eval.UDFContext;
public class Transform extends EvalFunc<String> {
    public String exec(Tuple input) throws IOException {
        if(input == null || input.size() == 0) {
            return null;
        }
        Configuration conf = UDFContext.getUDFContext().getJobConf();
```

Flat No. 212, 2nd Floor, Annapurna Block, Aditya Enclave, Ameerpet, Hyd.
E-mail: info@kellytechno.com www.kellytechno.com 040-6462 6789, 0998 570 6789.

```
String from = conf.get("replace.string");
if (from == null) {
    throw new IOException("replace.string should not be null");
}
String to = conf.get("replace.by.string");
if (to == null) {
    throw new IOException("replace.by.string should not be null");
}
try {
    String str = (String) input.get(0);
    return str.replace(from, to);
} catch (Exception e) {
    throw new IOException("Caught exception processing input row", e);
}
}
```

Pig UDF Example Using UDF

-- see pig
pig-training.jar contains custom udfs developed by a user. It has to be registered
-- before running the pig script.

```
REGISTER pig-training.jar;
A = LOAD '/input/sample.txt' AS (line: chararray);
B = FOREACH A GENERATE kelly.training.pig.udf.Transform(line);
```