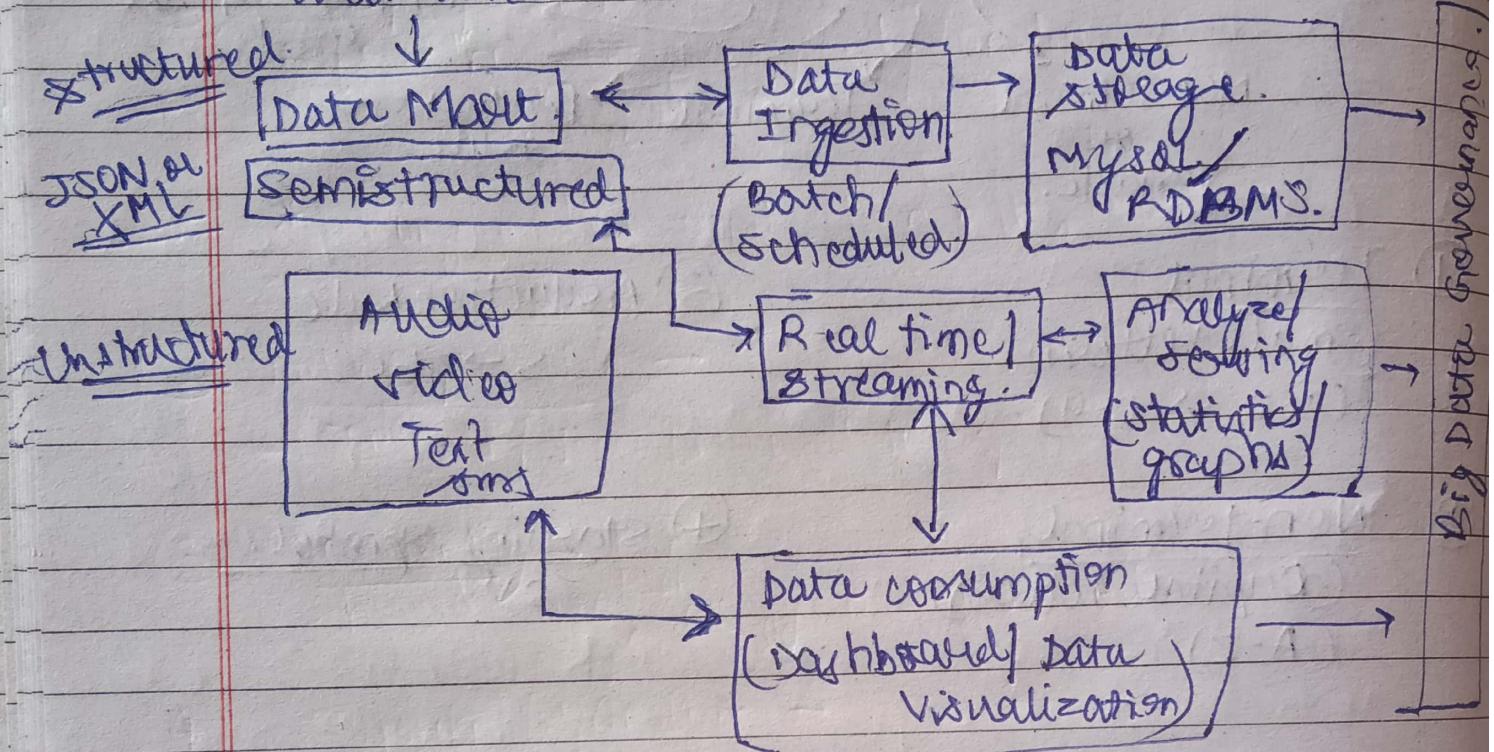


Architecture of Big Data

Data Sources.



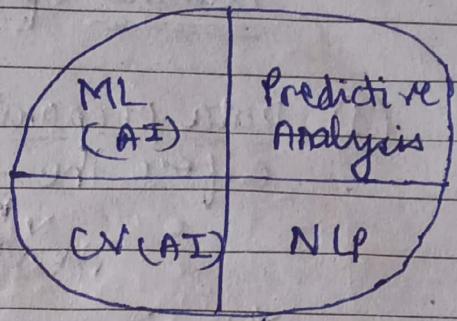
Data Flow.

- ① Capture data.
- ② Storing data.
- ③ Processing data / Data Ingestion
- ④ Data Visualization.

Fields in which Big Data is used :-

- ① Predictive Analysis → ~~Machine Learning~~
[Trends / graphs]

② Machine Learning. (Pattern matching / Recognition) (part of AI)



③ Natural Language Processing (NLP)
(Text as data)

④ Computer Vision (CV).
(Images as data)

Technologies used for Big Data:— (large dataset)

① Apache Hadoop (It is an open source framework)

② Spark. (by Apache)

↓ used.
in combination
with Hadoop

→ Open source processing
engine.

distributed
computing.
scalability
↓
Availability

③ Apache Flink. (Open source stream processing framework)
→ (User friendly api)
→ (live streaming of real-time data).

④ Presto (Open source SQL engine support #
interactive analysis on huge dataset)

⑤ Druid (open source analytical data storage.
→ designed for querying event-based data)
→ (log files)

Hadoop
MapReduce
Apache
HDFS

Hortonworks
Big Data Appliance
(Oracle)

Features of Big Data :-

① Data preparation →

Before the ~~data~~ is applied / constructed.

model
iterative

* During model construction but, data iteration before

② Data Exploration →

(Insights / Visuals used to explore data)

③ Scalability →

should use less energy and less cost.

④ Supports for Analytics →

Various types of analysis

⑤ Version control → For every version, checks are made.

⑥ Data Management →

⑦ Data Integration → collection of different types of data from diff. sources.

⑧ Data governance →

data should be accurate, usable, reliable, and easy to govern.

⑨ IDS → Intrusion Detection System.

⑩ Visualization → visual representation.

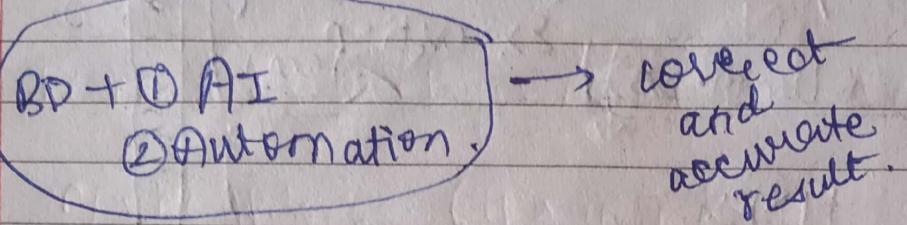
Date → verification of strategies adopted by the organisation and application of big data on it.

Audit and Analysis :-

- ① Reduction in operational cost.
- ② Improved decision making.
- ③ High customer retention.
- ④ High satisfaction rate.
- ⑤ Banking → Manufacturing.

* Prediction of Disease Sympt.

[candidate selection algorithm.]



* Knowledge base → Dataset with output.

⑤ Fraudulent cases

Big Data

Unit-2 (Hadoop)

- large scale batch data processing system.
- Distributed cluster system.
- parallel data processing

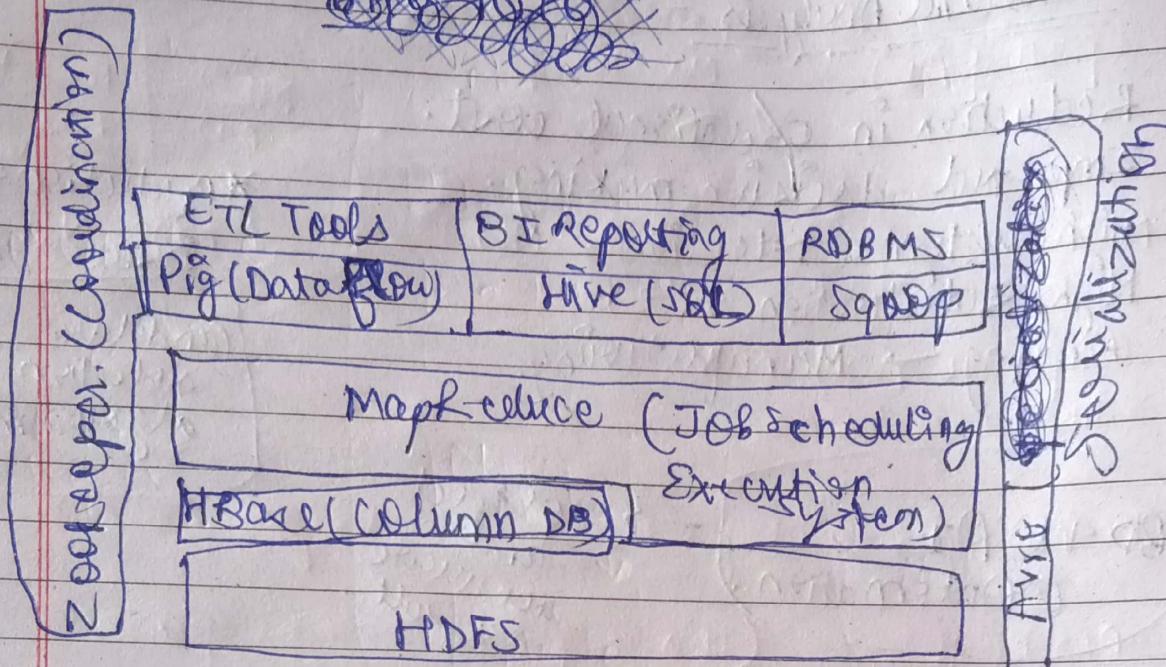
Features.

- ① Access to file systems -
 - common file
 - Hadoop package
 - contributor section
 - docs
- ② source code.

"Let us always meet each other with a smile, for the smile is the beginning of love." —Mother Teresa



Scanned with OKEN Scanner



Hadoop Architecture.

Spoon → CLI.

ETL → Extract, Transform, Load

BI → Business Intelligence.

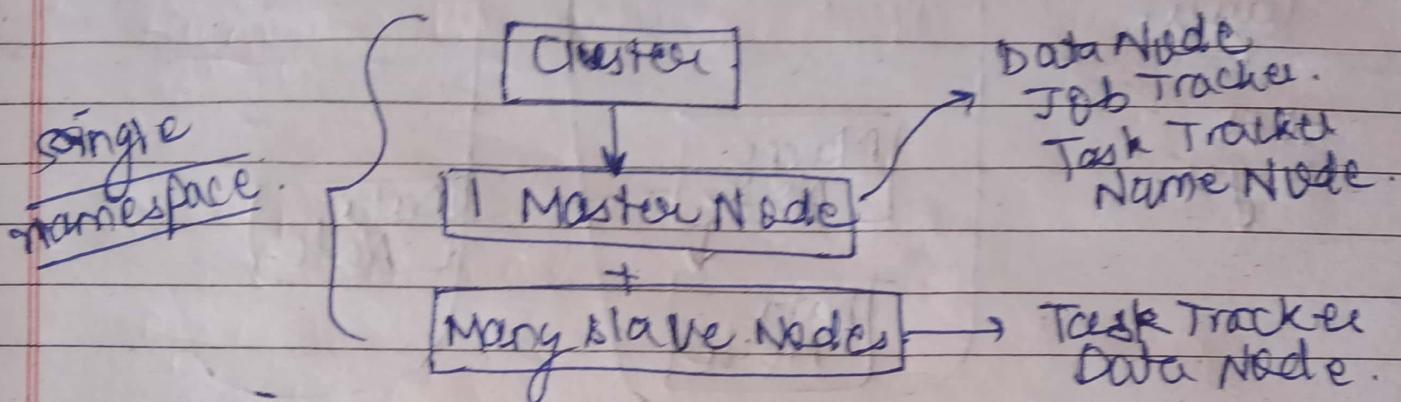
Core components of Hadoop.

HDFS MapReduce

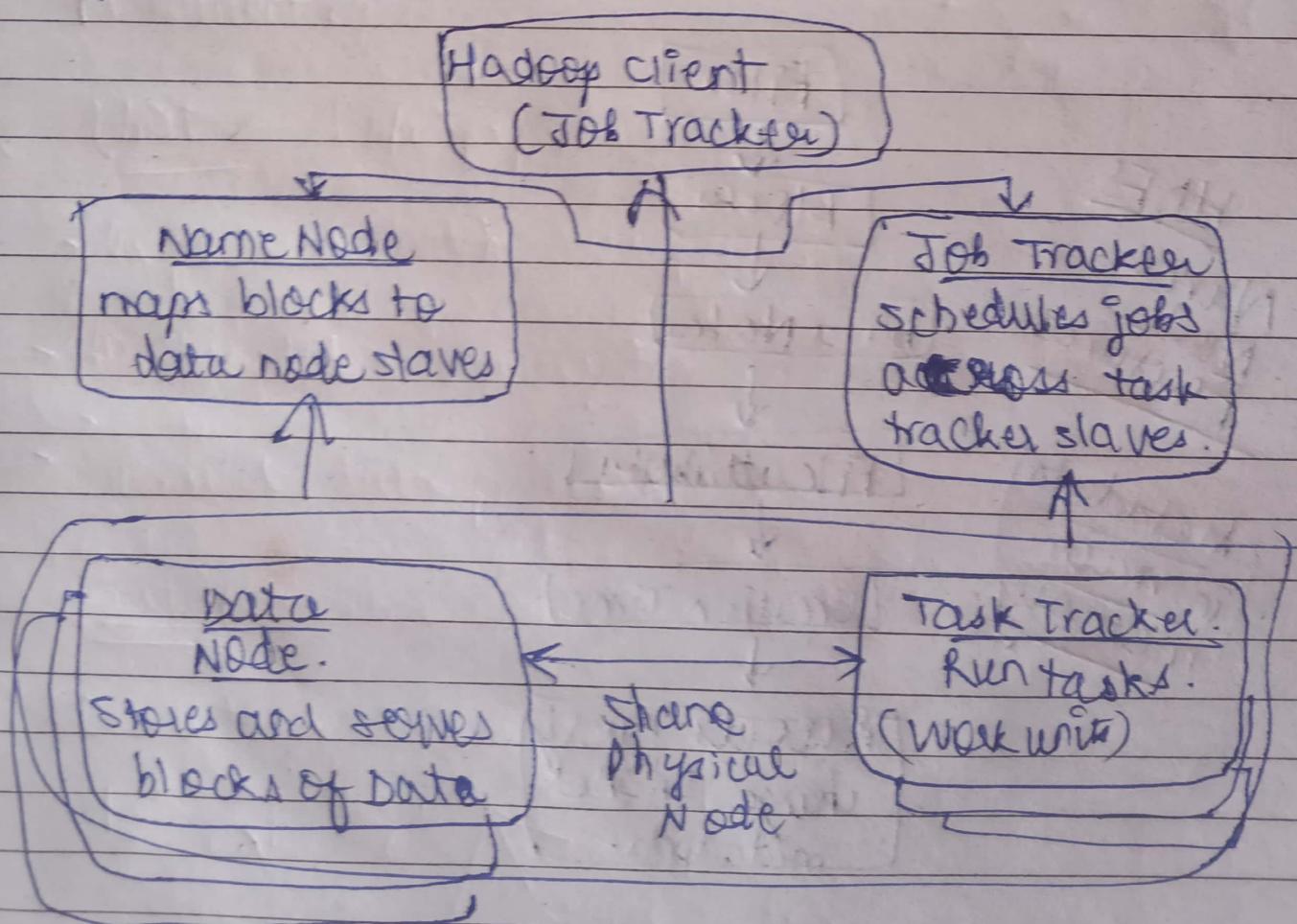
(Hadoop distributed file system).

HDFS

- Single namespace for the entire cluster.
- Traditional hierarchical file organization
- write-once read-many access model.
- Aware of network topology.



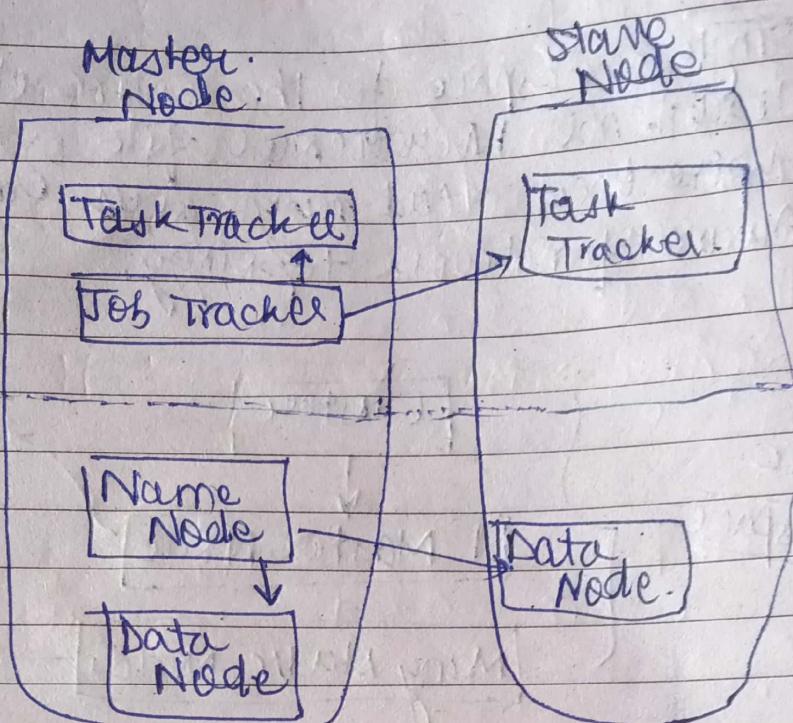
High Level Architecture.



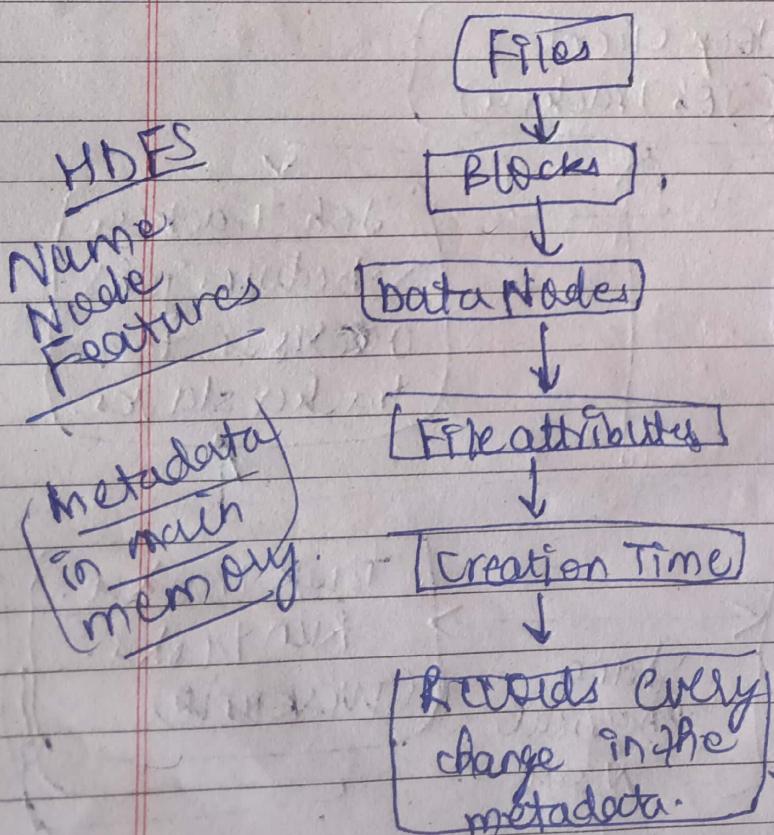
Date _____ / _____ / _____

MapReduce
layer.

HDFS
layer.



Hadoop Cluster

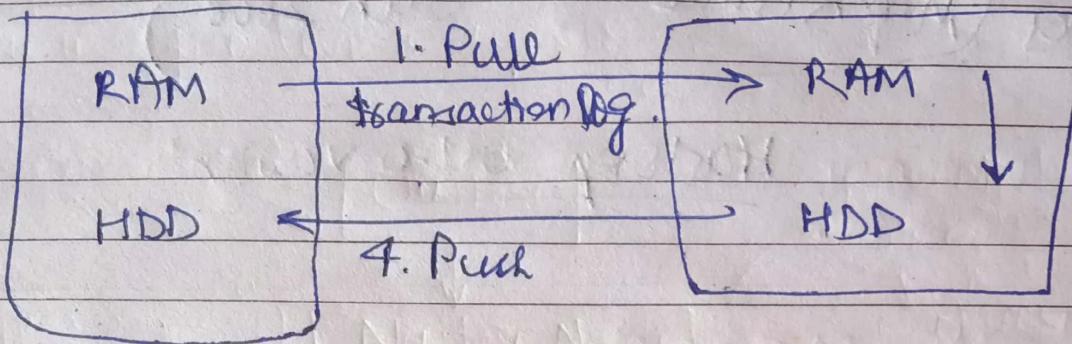


"Happiness is when what you think, what you say, and what you do are in harmony." —Mahatma Gandhi

HDFS - Name Node Architecture.

Primary Name Node

2. Merge changes
Secondary Name Node.



3. Stores to HDD

Data Nodes.

- Block stores data
- Periodic validation of checksums.
- Periodically sends a report of all ~~respective~~ existing blocks to the name nodes.

* MapReduce.

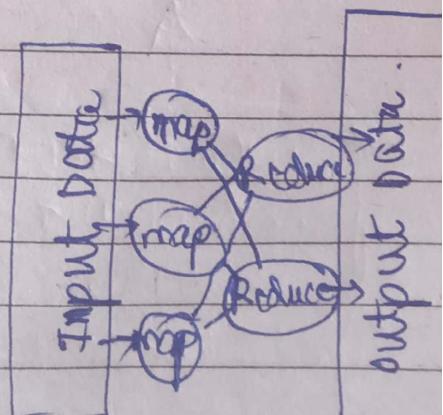
Job Tracker. Splitting into map and reduce tasks. Scheduling tasks on a cluster node.

based on
weights

split [kl, y1]

sort by kl

Merge. [kl, y1, y2, ys]



Task Tracker.

Runs MapReduce task periodically.

"Let your life lightly dance on the edges of time, like dew on the tip of a leaf." —Rabindranath Tagore

5000 lines

HDFS * Date

Apache Software Foundation.

Page No.: _____

Date _____

→ Doug cutting

→ Mike Cafarella.

Oct. 2003

6000 lines

MapReduce

→ Jan 2006

05)

Submit

BB

Hadoop 0.1.0 version.

Inp

"Happiness is when what you think, what you say, and what you do are in harmony." —Mahatma Gandhi

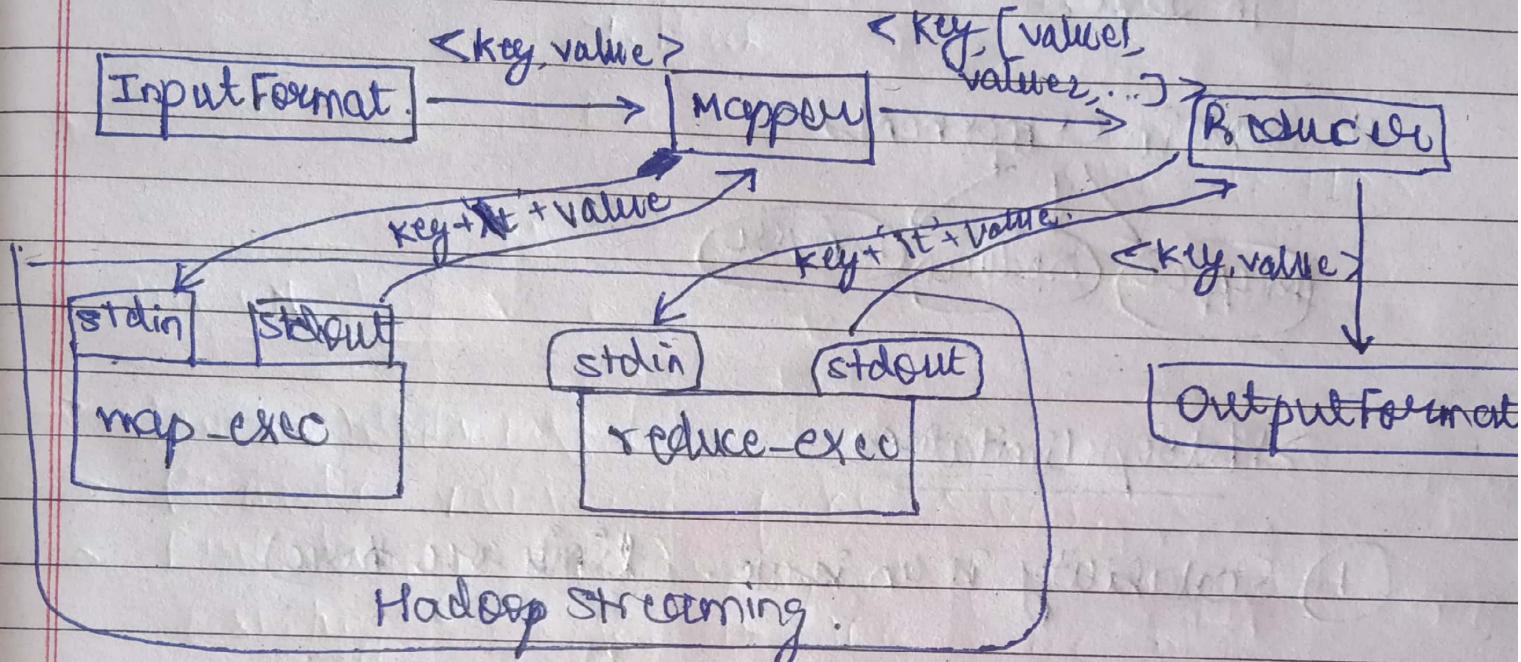


Scanned with OKEN Scanner

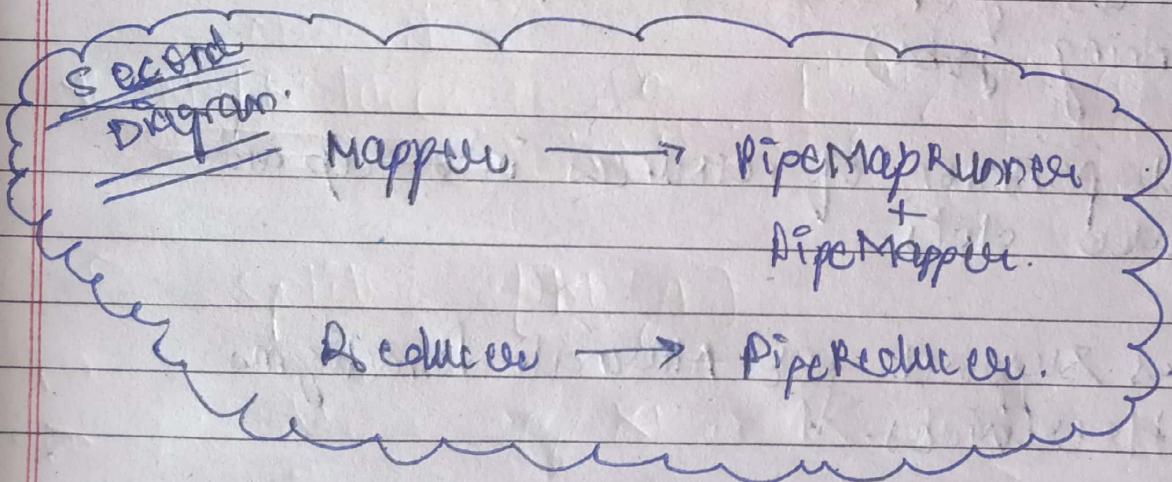
Q5) what is string? Explain any 5 string library functions used in Python with examples?

Submit → 6th september (Wednesday).

Big Data. (Hadoop pipes).



Hadoop Streaming.

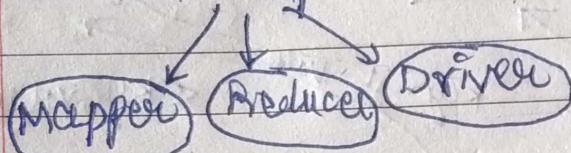


Date I/O Format

- inputFormat <javaClassName>
- JobConf: setInputFormat()
- outputFormat <javaClassName>

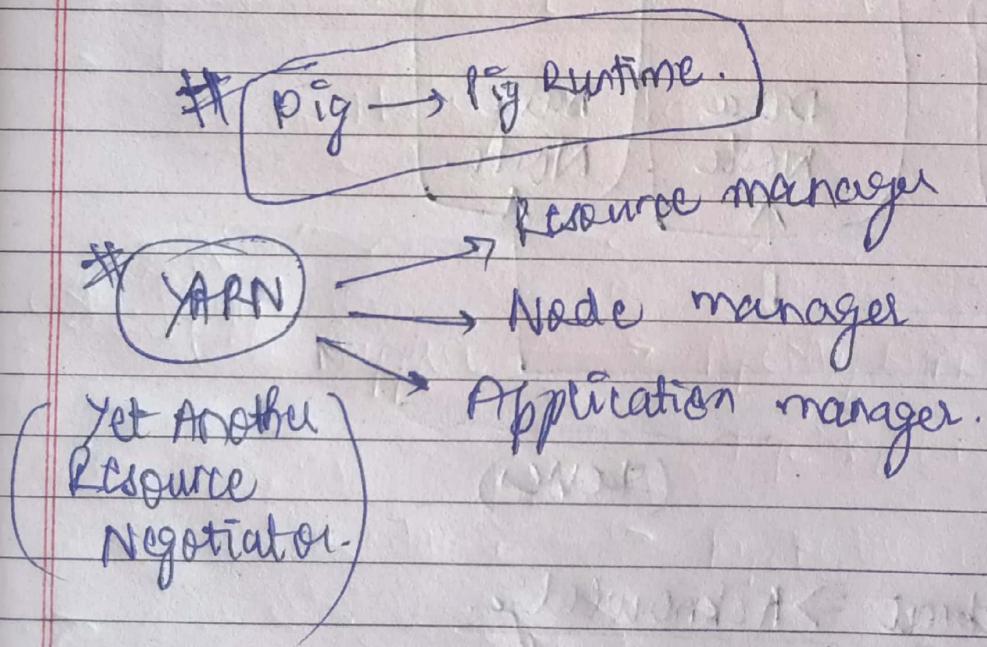
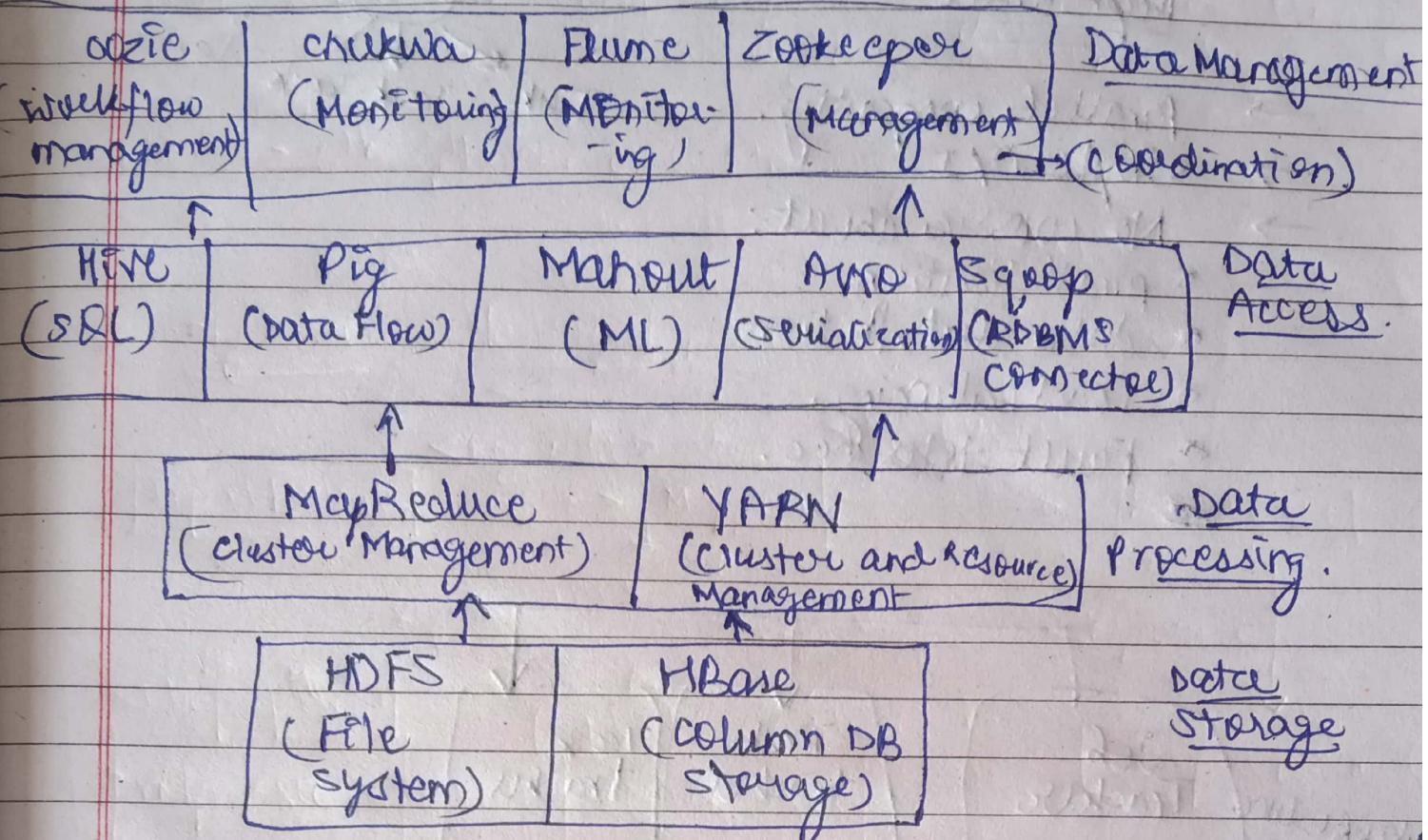
Written in Java but not necessary to code in Java.

3 main components



→ MapReduce Limitations

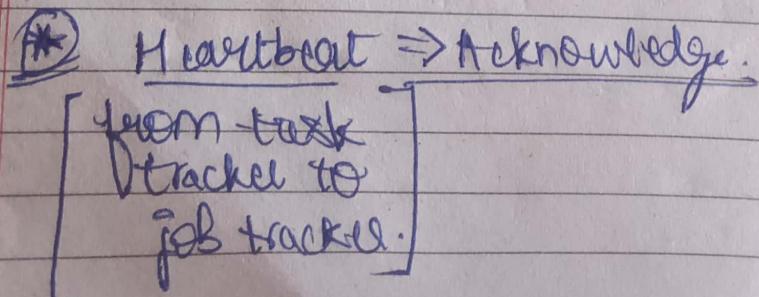
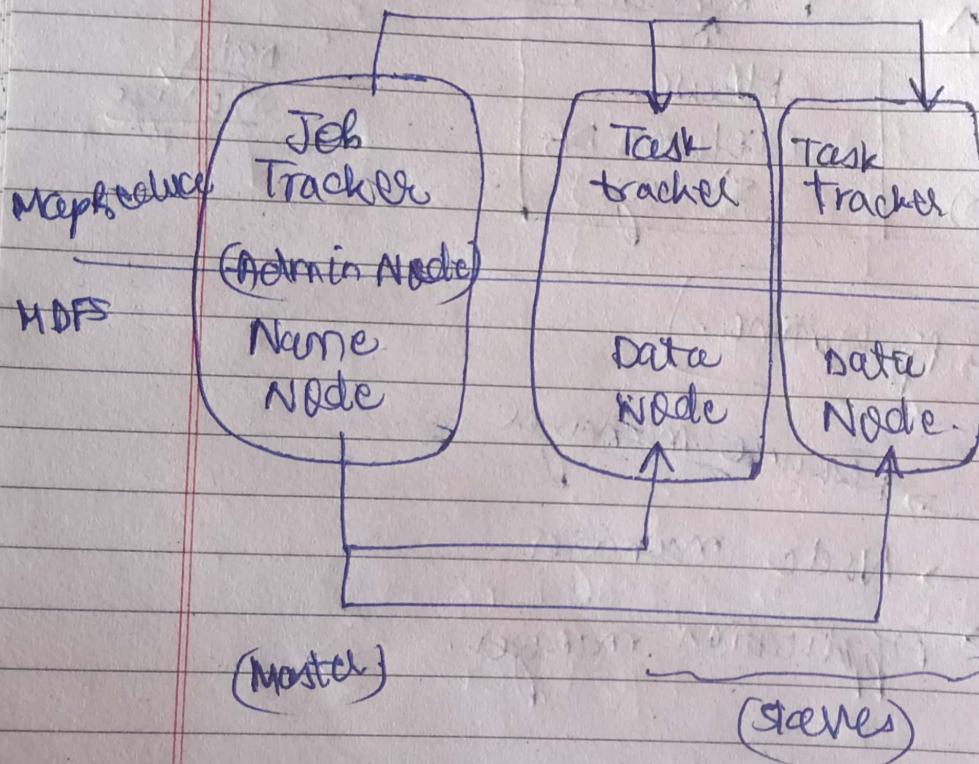
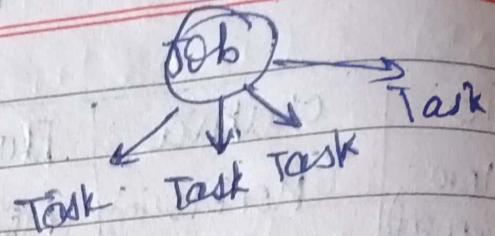
- ① Scalability is an issue. (limits core theme)
- ② Availability (Failure kills queued and running jobs)
- ③ Hard partition of resources into map and reduce slots.
- ④ lacks support for alternate paradigms and services.



Date _____

MapReduce.

- Parallel computing.
- Simple API
- No worries about:
 - parallelization
 - data distribution
 - Load balancing.
 - Fault tolerance.



Data is transformed from I/P files and fed into mappers.

