

Advance Database management

Week-2 (class-2)

Continuation from first

ppt

- **Bigdata** - is the main focus of this course.

Week-3 will be more

Hadoop.

Highly scalable is mostly

consist of

commodity hardware

SQL -

two Application

OLAP, OLTP

Big data Application

OLTP - Strictly follows ACID properties

but for Big data we usually use OLAP as ACID is not that important rather performance is important

Big Data Technologies

- Goals
 - Support of non-structural or semi-structural data
 - Highly Scalable

- Two major technologies
 - Hadoop ecosystem
 - NoSQL DBMS

- Support of
 - Web-based Real-time Application
 - OLAP
 - Online Streaming Processing

Historical data - This data is processed in batches

Online streaming process -
is Basically real-time OLAP

Hadoop Ecosystem

Hadoop Ecosystem - History -

- Hadoop Distributed File System/MapReduce
 - Initially developed at Google for a project named Nutch (2003) to crawl and index millions of web pages
 - The storage engine and processing parts which Nutch packaged as the core of **Hadoop** product at Yahoo
 - First Apache open-source release in 2006
- Pig
 - Initially developed at Yahoo
 - Component built on top of MapReduce
 - Provides applications with SQL-like interface to access data in HDFS
 - First Apache open-source release in 2008
- Hive
 - Initially developed at Facebook
 - Component built on top of MapReduce
 - Provides applications with SQL-like interface, DDL, JDBC/ODBC interface to access data in HDFS
 - First Apache open-source release in 2010

Map Reduce.

lies

Hadoop developed at google then it was open sourced

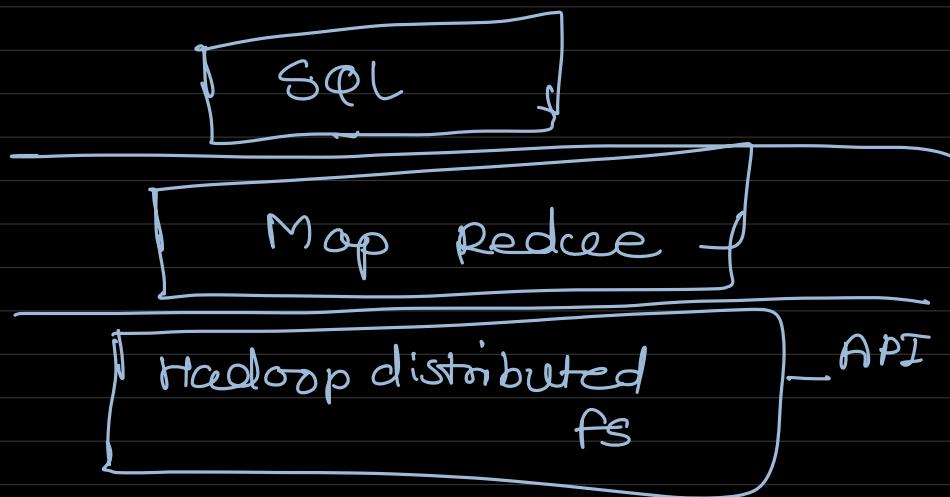
Hadoop distributed file system (HDFS) = JAVA pro

In case you don't wanna use

java then we can use SQL
using pig & Hive.

(YARN - Yet another resource negotiator)

Generic distributed processing distribution.



With YARN - you can build Application on top
that's another use of YARN with negleated
Map reduce which speed up the process

so YARN is key

Descriptive analysis
pig, Hive, (Map reduce)

OLAP

Spark (YARN)

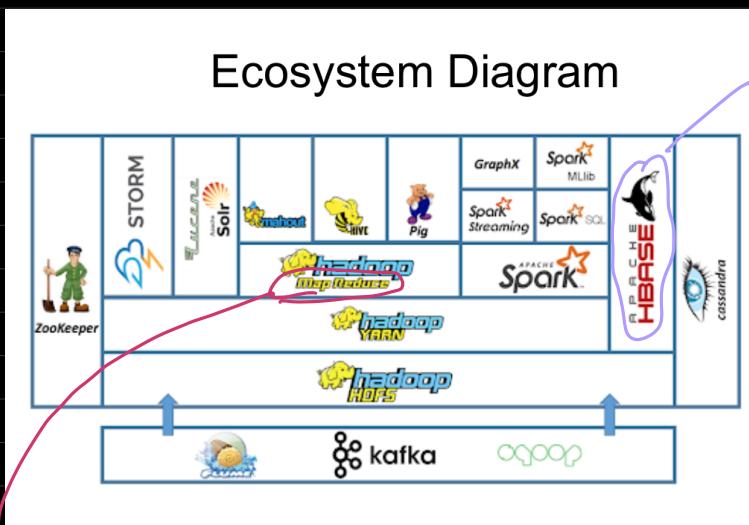
proactive analytics
Mahout

Unified Engine
it's product it can do
descriptive analysis, predictive analysis
& also, online streaming process

Online streaming processing



One of the focus of class
Map reduce



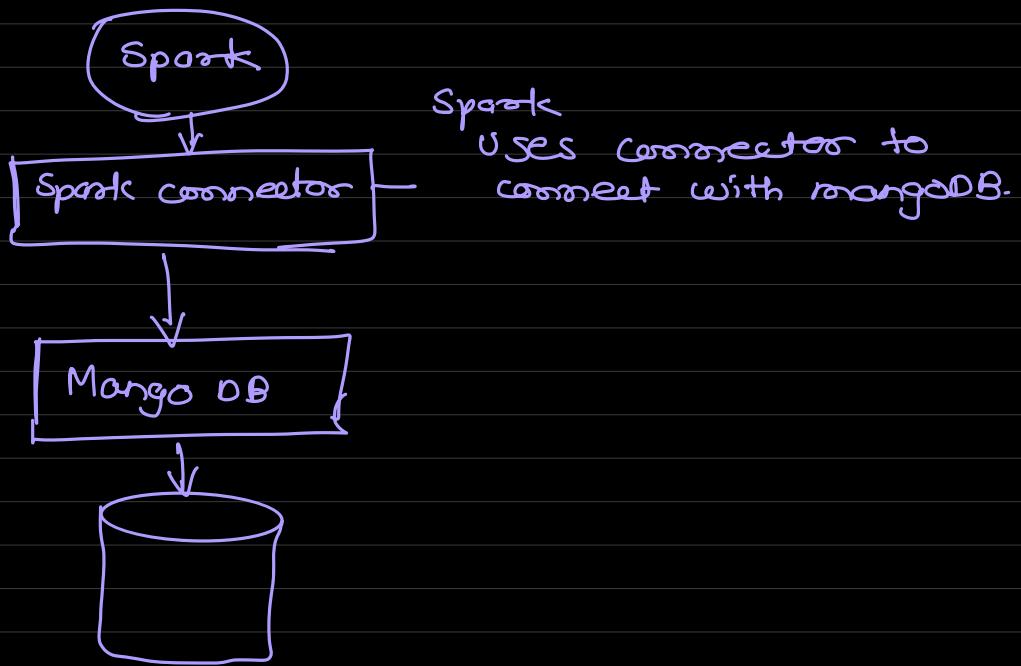
It's itself is
DBMS

According to poor. map reduce should be below YARN

What's the diff betⁿ DBMS & file system

logical data view
(doesn't care where file is stored)
High level view to manage data

physical location



New Big Data Concepts

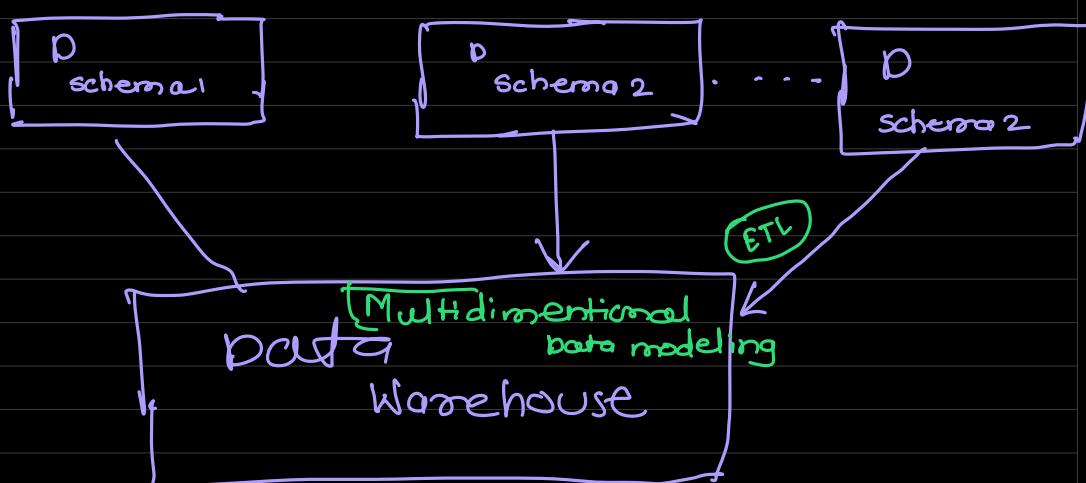
On-line streaming processing

- Real-time OLAP

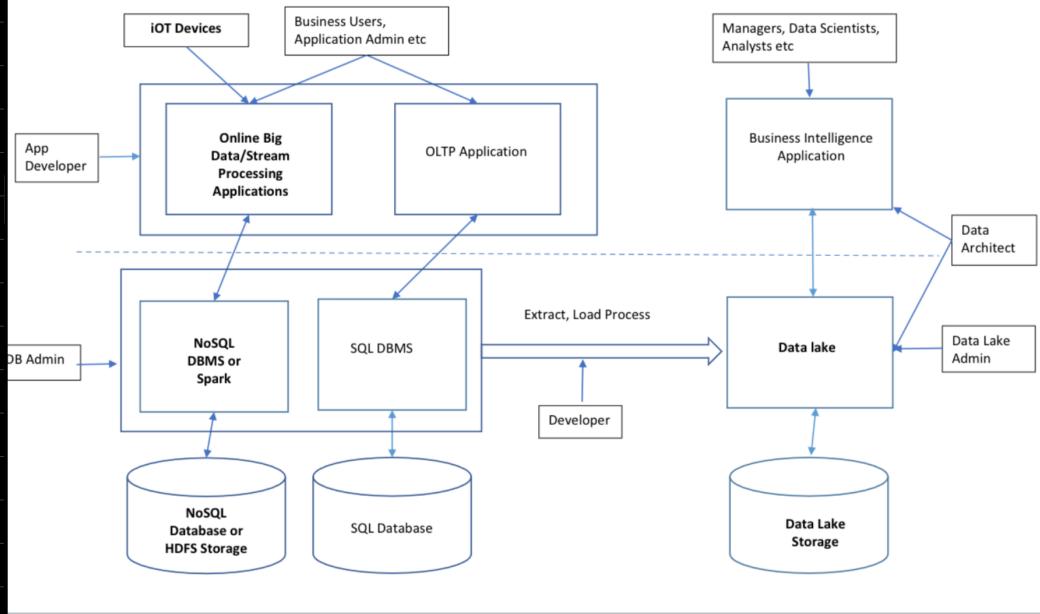
Data Lake (vs Data Warehouse)

- Data warehouse store data in a multidimensional model
- Data lake store data as it is

needs
thing



IT Operational Diagram for Big Data era



final ~~data~~ - project | Data lake
NoSQL
Big data

Course Roadmap

- Distributed Data Processing
- Hadoop Ecosystem
 - Hadoop Distributed File System
 - YARN
 - MapReduce Programming Model
 - Spark
 - HBase
- NoSQL DBMS
 - Data Sharding
 - Data Availability
 - Data Consistency
 - Data Backup & Recovery

— — — X — —
chapters - 3

Data processing

1. Multiple threading
Execution time = $T_p \times N$ Improving Execution time

with multithreading ex. 20

$$= \frac{T_p \times N}{20}$$

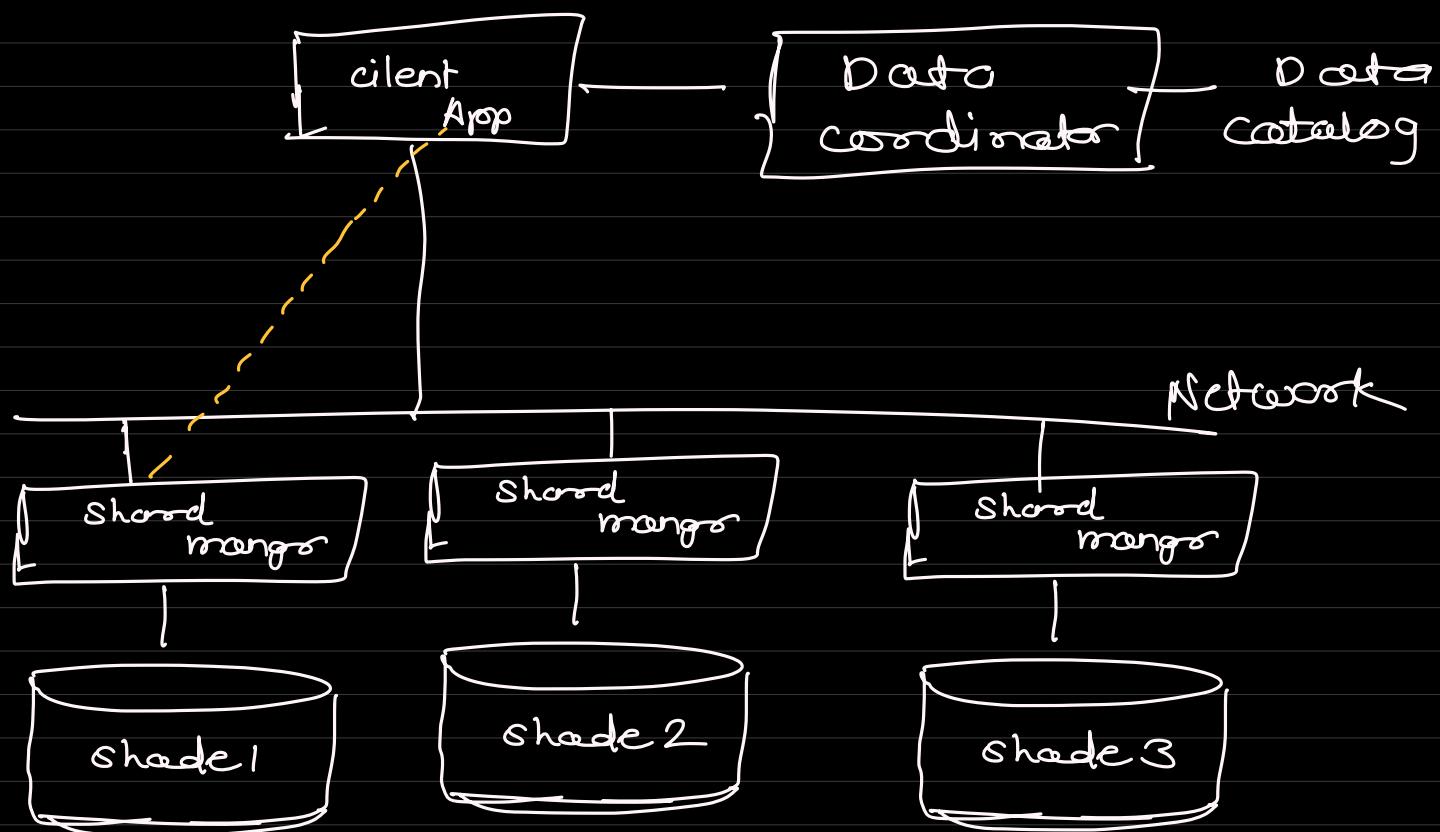
so we significantly reduce Execution time.

2. Multiple process multiple processors

3. Horizontal scaling is more cost effective than vertical scaling

▲ Data sharding — Search more implement parallelism into this send data into blocks to diff machines

more group = more parallelism



Distributed Algorithms
ex.

- Hadoop Mapreduce
- distributed query processing - send query to diff systems
- distributed ML algo. - spark already implemented it (i think)

framework for executing distributed algo

- Generic - Hadoop, Mesos

Distributed job processing

①

Task scheduling aka. DAG =

Distributed job should know following

- where job is located
- DAG
- Resource requirements to run each task

IPC - interprocess communication

= distance from data.
Data locality? = How close the data is to your task

medium How much cost the overhead have to get to the data
if Overhead is higher then locality is low

- ↓
2. if data is on diff machine data overhead is Higher.
3. if task is on another machine & data on another machine
then that's the highest locality.

→ worst

Data Availability.

- Usage of redundancy & proper failover will achieve high availability.

Disk array - Backup disk

- so when using shards to maintain the high availability we keep backup of that data to another disk

Map reduce Application
Spark Application

Arch & implementation of Hadoop Ecosystem.

final project

2 Member team.