

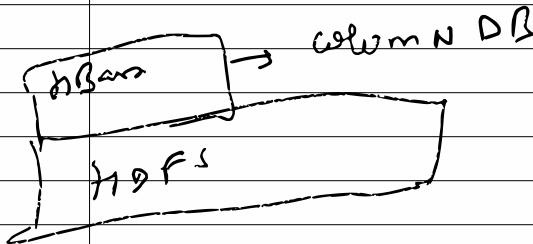
## Lecture 9 → HBase slides

Hbase → distributed column orientated data store built on top of HDFS

↳ Apache open source project

↳ storage for hadoop distributed computing

Logically organized into tables → rows and columns



Map Reduce or SQL → row based DB

### Hbase

Distributed system → scale to 100 or 1000s of nodes

↳ Fast record lookup  
record level insertion  
w/ updates

↳ by writing new version  
of values.

### hadoop

Distributed systems →  
scale to 100/1000 of nodes

→ good for batch processing  
Not good → record lookup  
operations

incremental addition  
of small batches

No random reads / writes → stick to HDFS

HBase run modes

↳ standalone  
Distributed

Out of box → Hbase runs in standalone mode.

→ configure hbase → edit files

↳ Hbase config directory

edit conf/hbase-env.sh → tell hbase which java to use.



configure

heapsize  
options for JVM

location for log files

JAVA\_HOME → points to root of java install.

1) standalone Mode

↳ default mode

Hbase does not use HDFS → uses local files

runs all Hbase daemons and a local Zookeeper  
all wrap in the same JVM.

Zookeeper binds to a well known port so clients  
may talk to Hbase.

### Local file system & durability

↳ using Hbase on local file system does not  
automatically guarantee durability.

↳ local FS will lose edits if files are not  
properly closed → very likely to happen  
when experimenting with a new download.

Need → run Hbase on HDFS to ensure all  
writes are preserved.

Running of local FS → fails off ground easily and  
makes it difficult to sync with system.

standalone  $\rightarrow$  runs all daemons and zookeeper on same VM

Running Hbase -

- Hbase shell  $\hookleftarrow$

$\rightarrow$  create table ~~new~~ named 'ts2' with a single column family 'cf'.

$\Rightarrow$  create 'ts2', 'cf' ;

view table

$\Rightarrow$  list 'ts2'  $\hookleftarrow$

create  $\rightarrow$  create 'table-name', 'column-family'  
view  $\rightarrow$  list

insert  $\rightarrow$  put 'table-name'; 'row1', 'cf:a'; 'value1'  
 $\quad$  put 'table-name'; 'row2', 'cf:b'; 'Value 2'

Verify the data by running scan

Scan 'table-name';  
for single row

Scan 'table-name'; 'row1';

delete table

↳ first disable

disable 'table-name'

then drop

drop 'table-name'

stop hbase

⇒ ./bin/stop-hbase.sh

2) distributed mode

↳ can be divided in

↳ pseudo-distributed → all daemons run on single node

↳ fully-distributed → daemons are spread across  
all the nodes in cluster

Pseudo-distributed mode can run against local  
file system or against an instance of HDFS

Fully-distributed mode can run only on HDFS

## Pseudo-distributed

↳ simply a fully distributed mode runs on a single host

Use this configuration testing and prototyping on Hbase

Don't use this mode for production nor for evaluating Hbase performance.

If want to run on HDFS rather than local FS, setup your HDFS

Ensure you have a working HDFS before proceeding.

### Configuration for Pseudo-distributed mode

i) Edit conf/hbase-site.xml

↳ tell hbase to run in (pseudo)-distributed mode rather than in default stand alone mode.

↳ set  $\rightarrow$  hbase.cluster.distributed  $\rightarrow$  true

$\Downarrow$   
default false

↳ with this Hbase starts up  
Hbase master process

Zookeeper server  
Region server

temp files will be saved in  $\rightarrow$  hbase-username

## Using HDFS

Setting up Hbase to write in local file system's temporary directory is ephemeral step.

↳ it will lose all the data once the system is shutdown.

For permanent setup we make use of an instance of HDFS; Hbase data will be written to the hadoop distributed file system rather than local file system's temp directory.

↳ changes → hbase-site.xml

```
<name> hbase.rootdir </name>
<value> hdfs://localhost:9000/hbase </value>
```

## Starting initial hbase cluster

start hbase → /sbin-hbase.sh

start fact up master → bin/local-master-backup.sh start 1.

We can start upto 10 backups.

Start region server → bin/local-regionserver.sh start 1

We can start total of 100 region servers.

## Fully distributed mode

for running hbase on fully distributed mode

set hbase.cluster.distributed to true -TURE

point hbase.rootdir to your namenode IP.

if namenode example is running on port 8020

<name> hdfs.rootdir /name </name>

<value> hdfs://namenodexample:8020/hbase </value>

## HDFS client configuration

add pointer to hadoop-conf-dir to the hbase-clientpath  
in hbase-env.sh

example of such hdfs client configuration :  
dfs.replication .

if we want to run with a replication factor of 5,  
hbase will create files with default of 3  
unless you do the above to make the configuration  
available to Hbase .

Master host → 16010

Region Server → 16020

## Hbase Data Model

↳ based on Google's Bigtable model  
↳ key - value pairs.

row key : ↓  
column family

timestamp + values.  
↳ for update or  
of values

## Hbase logical view

(row key  
implicit primary  
key in form  
terms) | part → is all byt[] in hbase

Dif type of rows have dif set of columns

table is sparse  
Different cells may have different values  
at different timestamp.

Each row has a key ✓

Each record is divided into column families →  
Each column family consists of one or more  
columns ↴

key

↳ byte array

- serves as the primary key for the table
- indexed for fast look up

column family

- Has a name : string
- contains one or more related columns

column

- belongs to one column family
- included inside the row
  - ↳ family Name : column Name.

Version Number

- Unique within each key
- By default → system's timestamp
- Datatype is long

Value (cell)

- Byte array

Notes:-

Hbase schema consists of several tables

Each table consists of a set of column families  
↳ columns are not part of schema.

Hbase has dynamic columns

↳ because column names are encoded inside the cells

↳ different cells have different columns.

The version number is user supplied

↳ does not have to be inserted in increasing order  
Version number can be unique within each key.

Tables can be very sparse  
Many cells are empty

Keys are indexed as primary key.

### Hbase physical model

Each column family is stored in a separate file  
↳ HTables.

keys and version number is replicated with each column family

Empty cell are not stored.

Hbase maintains multilevel index on values  
(key, column family, column name, timestamp)

Example

Row Key	Data
cutting	info: {height:'5', stick:'ai'} roles: {ASF: 297}

## info CF

key	column key	timestamp	cell value
0123	info:height	123456..	1

## roles CF

key	column key	timestamp	cell value
0123	roles:ASF	123456..	29A

## Column families

↳ different set of columns may have different properties and access patterns

Configurable by column family

- compression (none, gzip, lzo)
- version retention policies
- cache priority

CFs stored separately on disk: access one without waiting 10 ms the other.

## Tables regions

Each HTable (column family) is partitioned horizontally into regions

↳ common part to HDFS blocks

HDFS: blocks  $\Rightarrow$  HBase: HTable  $\Rightarrow$  regions

## Hbase Architecture:

3 Major components

↳ The Hbase Master  
↳ One master

↳ The HRegion Server

↳ Many Region Servers

↳ The Hbase Client

## Hbase components

### Region

- ↳ A subset of table rows
  - ↳ like horizontal range partitioning
- ↳ Automatically done

### Region Server (many slaves)

↳ manages data regions

↳ serves data for reads and writes (using - log)

### Master

↳ Responsible for coordinating the slaves

↳ Assign Regions, detects failures

↳ Admin functions.

HMaster

HRegion Server

↳ HRegion

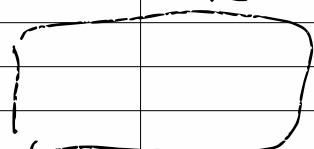
↓  
store file

Hfile

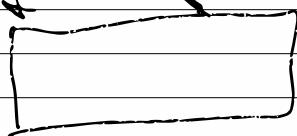
HMaster

HFile

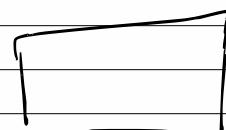
DFS  
client



Data Node



Data node



Hadoop

ZooKeeper

↳ Hbase depends on Zookeeper

→ By default hbase manages one Zookeeper instance → starts and stops

→ HMaster and HRegion Server register themselves with Zookeeper.

## Creating a Table

```
HBaseAdmin admin = new HBaseAdmin(config);
```

HColumnDescriptor [] column;

```
column = new HColumnDescriptor[2];
```

```
column[0] = new HColumnDescriptor("columnFamily1");  
column[1] = new HColumnDescriptor("columnFamily2");
```

HTableDescriptor desc;

```
desc = new HTableDescriptor(Bytes.toBytes("myTable"));
```

```
desc.addFamily(column[0]);
```

```
desc.addFamily(column[1]);
```

```
admin.createTable(desc);
```

Steps:

- 1) Create admin
- 2) Create column families
- 3) Create table and add column family to it
- 4) admin add the table into it.

Operations on Regions :- Get()

→ Given a key → return corresponding Record



for each value return the higher version.

⇒ To get current version

Get ger = new Get (Bytes. toBytes ("row1"));

Result r =htable. ger (ger)

byte [] b = r. gerValue (Bytes. toBytes ("c1")),  
Bytes. toBytes ("attr"));

⇒ we can control number of versions we want  
↳ following code will return last 3 versions  
of the row.

Get ger = new Get (Bytes. toBytes ("row1")),

get. setMaxVersions (3);

Result r = hashtable. ger (ger)

Current version as above

All versions by →

List <key Value> ka = r. gerColumns (Bytes. toBytes ("c1"),  
Bytes. toBytes ("attr"));

all versions

Operations On Regions: Scan()

↳ It allows iteration over multiple rows for specified attributes.

---

get()

↳ select value from table where  
key = "com.apache.www" and  
label = "anchor: apache.com"

Scan()

Scan select value from table  
where anchor = "cons1.com"

`put()`

- ↳ insert a new record with new key.
- ↳ insert a record for an existing key.

`Delete()`

Marking table cells as deleted

Multiple levels

- ↳ can mark an entire column family as deleted
- ↳ can mark all column family of a given row as deleted.

- All operations are logged by the RegionServer
- The log is flushed periodically

# Basic joins

- ↳ does not support joins

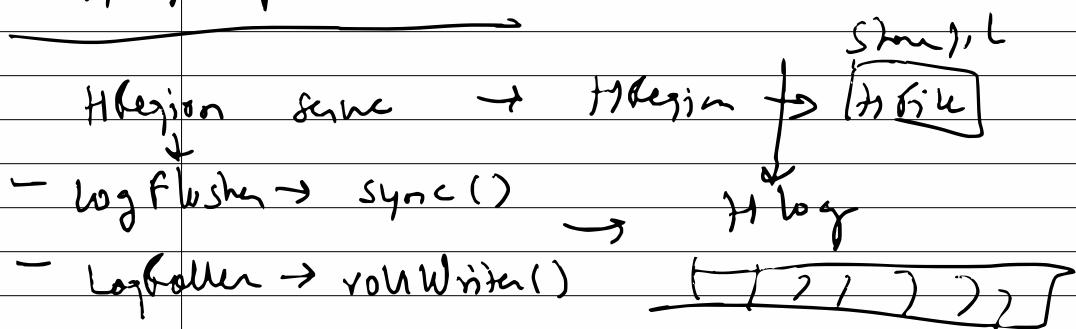
Can be done in application layer

- ↳ using `scan()` and `get()` operation.

## Altering a Table

- ↳ Disable the table before changing the schema
  - ↳ adding new column family
  - ↳ modifying existing column family.

## Logging operation



## Hbase Deployment

Master Node → NameNode  
Secondary NameNode  
HMaster  
Job Tracker  
ZooKeeper

Slave Nodes → RegionServer → 5+ slaves  
Data Nodes with HBase, HDFS,  
Task Tracker MR slave processes.

## HBase vs HDFS

### Plain HDFS / mR

Write pattern - Append only

Read patterns Full table scan, partition table scan

Hive key → Very good

Structured storage  
Distributed (JSON / CSV / log file)

Max data size → 30 + PB

### HBase

Random write, bulk incremental

Random read, small range scan or table scan

4-5x slower.

Sparse column family data model.

~ 1PB.

## HBase vs FDBms

### RDBms

Data layout Row-oriented

Transactions → multi-row ACID

Query lang SQL

Security Authentication / Authorization

indexes On arbitrary columns

Max data 7TB

size

Throughput  
limits

1000s queries / second

HBase

Column family oriented  
single row only  
get / put / scan / delete ex.

Work in progress

Row-key only

~ 2PB

Millions of queries /

second

## When to use HBase

Use if:-

- Need random read, random write, or both (but not neither)
- Need to do many thousands of operations per second on many TB of data
- Your access patterns are well-known and simple.

Don't use HBase If:-

- You only afford to your dataset and tend to need the whole thing
- Primarily do ad-hoc analysis (ill-defined access patterns)
- Your data easily fits on one beefy node.

## Case Study :- InMemory

Google lab published → BigTable



Index designed to manage very large datasets ("petabytes of data") in a cluster of data servers

BigTable supports key search, range search and high throughput file scans, and also provides a flexible storage for structured data.

HBase is open source clone of BigTable.

↳ Aimed at a large, distributed relational data base.

↳ Presents many aspects common to NoSQL systems :- distribution, fault tolerance, flexible modeling, absence of some features deemed essential in centralized DBms (e.g. concurrency), etc.

Key Value based storage → like map

↳ No table → No schema

can store value one another map  
multi map structure

Object remains un-updated by change.

An Hbase table is a multi-map structure

↳ instead of keeping one value for each property of an object, HBase allows the storage of several versions

↳ Each version is identified by a timestamp.

↳ Hbase simply replaces atomic values by a map where the key is the timestamp.

↳ helps to figure out the power and flexibility of the data representations

Now the document is built from two nested maps,

- 1) column
- 2) Time stamp

The document is globally viewed as a column map.

key ↴ map

Column key ↗ map → as many keys as there are time stamps.

The time stamp maps are completely independent from each other.

We can add as many timestamp keys as we want.

Same for column → can add as many columns as we need.

If a table has 2 families meta / dir  
then they are fixed

↳ Unlike column and timestamp maps,  
the keys in a family map is fixed.

→ We can not add new families to a table once it is created.

→ The family level therefore the equivalent of a relational schema,  
↳ Although content of a family value may be a quite complex structure

Row → tables

Row has value → key is chosen by user operation → get (key, value), put (key, value)

Table map is an sorted map.

↳ Rows are grouped in key values.

## 2) Apache Phoenix

(a) What is Phoenix?

Phoenix

↳ open source SQL skin for HBase

↳ Use standard JDBC APIs instead of the regular HBase client APIs. To create tables, insert data, query your HBase data

(b)

### 3) Apache Phoenix

A relational data base layer for Apache HBase

- Query Engine

- ↳ Transforms SQL queries into native HBase API calls

- ↳ pushes as much work as possible onto the cluster for parallel execution

- Metadatabase repository

- Typed access to data stored in HBase tables

- JDBC driver.

Apache HBase - A high performance horizontally scalable data store engine for big data, suitable as the store of record for mission critical data.

Not like a spreadsheet → a sparse, consistent, distributed, multi-dimensional, sorted map.

Scalability by dividing in regions.  
↓  
Region servers.

Phoenix puts the SQL back in NoSQL

↳ A complete relational system

↳ reintroduces the familiar declarative SQL interface to data (DDL, DML, etc)

- read only views on existing HBase data
- dynamic columns extend schema at runtime
- Schema is versioned - per row by HBase
  - allowing flashback queries using prior versions of meta data.

→ reintroduces typed data and query optimization possible with it

→ secondary indexes, query optimizations, statistics ...

→ integrates the scalable HBase data storage platform as just another JDBC driver.

Accessing HBase data with Phoenix can be substantially faster than direct HBase API use

Parallelizes query  
supports secondary indexing  
uses every trick based on - Hbase version, meta data, query, reverse scan, small scan, skip scan etc

Phoenix jar is installed on the HBase region server classpath.

The Phoenix JDBC driver is installed on the client.

The application speaks SQL to HBase-

Hbase Data model to the relational world.

Phoenix supports read only access to existing Hbase tables.

Empty cell will be inserted in each row in the native table to enforce primary key constraint

Note - serialized bytes in the table must match the expected phoenix type serialization.

g) Mahout is  
open source machine learning  
library from Apache  
↳ Algorithms fall under broad umbrella  
of machine learning

At the moment → recommendation  
engine, clustering and classification



scalable  
machine learning tool when the data  
is huge

↳ written in Java libraries

No UI, prepackaged server, or an installable  
framework of tools ~~and~~ to be adopted by  
developers.

Apache W�nes → Mahout

Much work is to convert mahout algorithms  
to work at the scale on top of Hadoop.

Many in experimental stage

Three core themes are evident

- 1) Recommender engines
- 2) Clustering
- 3) Classification

or currently has

↳ collaborative filtering

User & Item based recommenders

K-means, Fuzzy K-means clustering

Mean shift clustering

Discriminative clustering

Random forest, decision tree based classifier

High performance Java collections

All techniques work best when provided with large amounts of good input data.

Must not only work on large amounts of data but also must produce results very quickly

↳ scalability → makes key reason for being

Handout for large data management

↳ MapReduce → solving problems scale w/ hadoop

Requires Java b-

build around mapreduce → Higher level tools

↓  
data flow → lower level tools

Add matrix dependency in form.xml.

## Collaborative filtering

↳ recommendation based on user's relationship to item

↓  
No prior knowledge of properties of item needed.

Attribute → content based filtering.

Content base is good but ~~not~~ very specific

Like a content based recommendation system for books will not work for movie recommendation.

Example → User based Recommendation

Input:- large data → in form of preference

Strength of user's  $\leftarrow$  given a number  
preference to item

User Id, item id, preference

- ids are always a number in matrix-
- preference can be anything as long as larger value means stronger positive preference.

Data model  $\rightarrow$  stores and provides access to all the preferences, user and item data needed in the computation.

User Similarity  $\rightarrow$  provides notion of how similar two users are.

User Neighbourhood  $\rightarrow$  group of users that one

recommender  $\rightarrow$  pull all above together to recommend items to users.

User Test Data  $\rightarrow$  some training data set aside for testing.

difference  $\rightarrow$  estimate  $\Rightarrow$  real value

Root mean square of differences is used

$\downarrow$   
square root of the average of the squares of the difference between actual and estimated values.

$\downarrow$

Penalty penalises diff that are way off.

Precision  $\rightarrow$  portion of top recall from good recall

Recall  $\rightarrow$  portion of good recall that are in top recall.

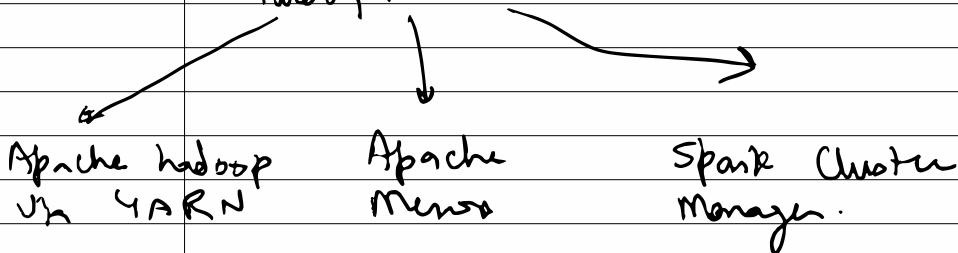
## Spark and Shark

Apache spark is a general purpose cluster in-memory computing system.

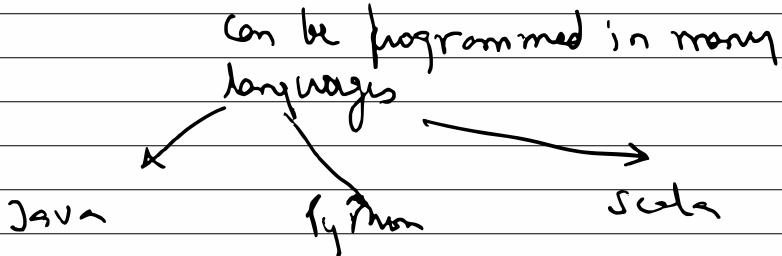
Provides high-level APIs in Java, Scala and Python, and an optimized engine that supports general execution graphs.

Provides various high level tools like Spark SQL for structured data processing, MLlib for machine learning and more.

Spark framework can be deployed through



Spark framework is polyglot



## Shark

↳ A fully apache hive compatible data warehousing system that can run 100x faster than Hive.

Spark → 100x faster than Hadoop for certain mapreduce application.

## Spark :-

Not a modified version of Hadoop

Separate, fast mr like engine

↳ in memory data storage for very fast iterative queries.

↳ General execution graph and powerful optimising

↳ 100x faster than Hadoop.

Compatible with → HDFS, HBase, Sequencefiles etc.

Shark → part of apache hive to run on spark.

Compatible with existing hive data, metastores, and queries (Hive QL, UDFs etc).

Similar speed up of upto 100x.

## Map Reduce Successful

↓ User wants

more complex multi stage application  
More interactive ad hoc queries

↓

Need faster data sharing across parallel jobs.

## Why spark ?

- Provides powerful caching and disk persistence capabilities
- Interactive data analysis
- faster batch
- iterative algorithms
- real time stream processing
- faster decision making.

## Data sharing in MapReduce :-

↳ slow due to replication, serialization and disk I/O

## Data sharing in spark

↳ 10 - 100x faster than network and disk.

## spark Programming Model

key idea  $\rightarrow$  RDD

$\hookrightarrow$  Resilient Distributed Datasets

- $\hookrightarrow$  distributed collections of objects that can be cached in memory across cluster nodes
- $\rightarrow$  manipulated through various parallel operators
- $\rightarrow$  automatically rebuilt on failure

Introdu

- $\hookrightarrow$  clean language integrated API in Scala
- $\hookrightarrow$  can be used interactively from Scala console

Example : log mining

scaled to 1TB in 5  $\rightarrow$  sec.

Fault tolerance

- $\hookrightarrow$  RDDs track the series of transformation used to build them to recompute lost data

```
merge = textFile(...).filter(~contains("error"))
    .map(~split(' ') (2))
```

Supervised operators

- $\hookrightarrow$  map, reduce, filter, count, groupBy, min, join, sort

## Other engine feature

- ↳ general graphs of operators
- ↳ hash-based reducers  
(faster than Hadoop's sort)
- ↳ controlled data partitioning to lower communication

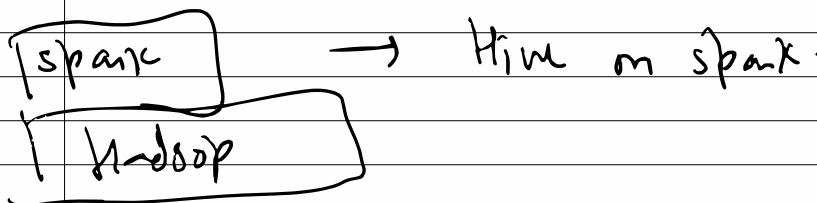
## Application

- ↳ in memory analytics & anomaly detection
- ↳ interactive queries on data streams
- ↳ exploratory log analysis
- ↳ traffic estimation w/ GPS data
- ↳ Twitter spam classification.

Conviva → 40x gain over hive

↓  
avoiding

- ↳ repeated reading, deserialization and filtering.



Caching hive records as java objects is inefficient due to high performance overhead

Instead, spark employs column-oriented

Storage using arrays of primitive types.

Benefit

- ↳ similarly compact size to serialized data → 5x faster access

Spark → interactive → multi stage apps

another → stream processing

- ↳ track and update state in memory as events arrive

- ↳ large scale reporting, click analysis, spam filtering etc.

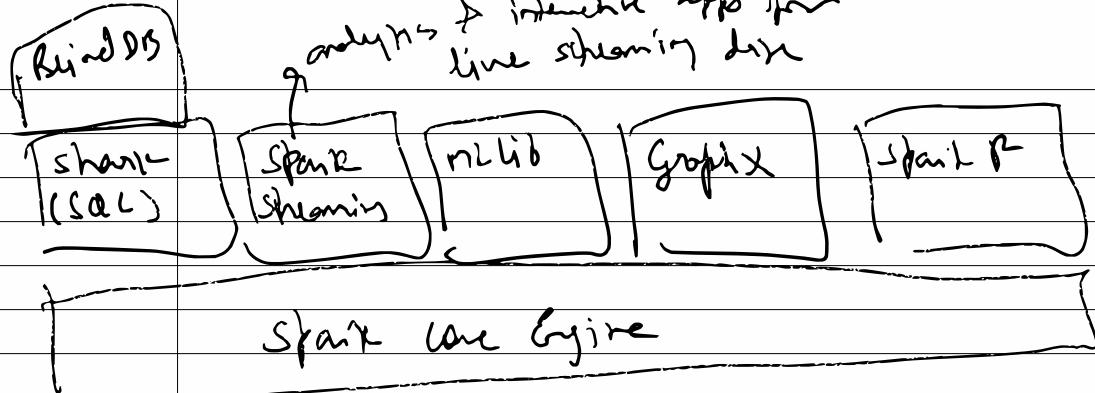
· flatmap() · map() · reduce by window

- Extend spark to perform streaming computation
- series of small ( $\sim 1\text{s}$ ) batch jobs, keeping state in memory as fault tolerance RPO.

↳ can process 44GB/second on 100 nodes at sub second latency

Spark & shark → speed up interactive and complex analytics on hadoop data

- ↳ easy to run locally, on EC2, on Mesos, YARN-



# PIG

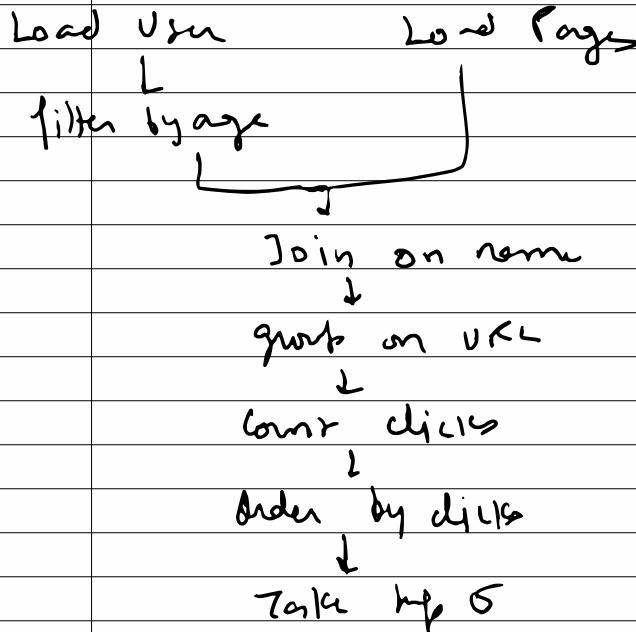
when

High level language  
Small learning curve  
Increases productivity

Insulates from complexity of Map Reduce  
↳ Job configuration library  
Mapper Reducer Optimization  
Data flow  
Job chains

Map Reduce Example

↳ Top 5 most visited by users 18-25



load and filter users by age

Users = LOAD 'users' AS (name, age)

users = FILTER users by age  $\geq 18$  and age  $\leq 28$

load pages

pages = LOAD 'pages' AS (users, url)

joined = JOIN users by name, pages by users

grouped = group joined by url

summed = FOR EACH grouped GENERATE  
group, COUNT (joined) AS clicks

sorted = Order summed by clicks desc

top 5 = limit sorted 5

Store top 5 into '/data/top5sids';

Hadoop is fine

↳ fewer lines of code

↳ less development time

↳ reasonably close to optimal performance

## Under the Hood

Pig LATIN programs

↳ Query Parser  
↓

Semantic checking  
↓

Logical Optimizer

↖

Logical to physical Translator

↓

Physical to m/R translator  
↓

Map Reduce Launcher

↓

Create a job jar to be submitted to hdfs cluster.

## How Pig differs from Map Reduce

Pig provides → join, filter, group by, order by, union etc.

MapReduce → group by operation directly (skipping reduce phase)

Order by → indirectly by Pig may it implements grouping.

filter and projection can be implemented trivially in the map phase

other operations (like join) are not provided and must be written by user.

Pig provides some complex, non trivial implementations of these standard data operations

Eg → records per key is not evenly distributed  
↳ some reducers get 10 times more data than others

Pig has join and order by operators that will handle this case and in some cases will rebalances the reducers.

Writing this in map reduce will be time consuming.

In MapR, the data processing inside the map and reduce phase is opaque to the system.

↳ MapR has no time to optimize or check user's code.

Pig can analyze pig latin script written by user and understand the dataflow the user is describing.

↳ Error checking and optimizing

Map Reduce does not have type system

↳ flexibility to user

↳ downside → no code check and error can occur during runtime.

---

### Pig Latin Scripts

```
users = LOAD '/hdfs/users' USING PigStorage('')  
as (name, age);
```

```
skinnyUsers = filter users by Age > 150;
```

```
skinnyUser = foreach skinnyUsers generate name, age;
```

```
Name = foreach Users generate name;
```

```
DistinctName = distinct names;
```

group then join in Pig Latin

txns = load 'transactions' as (customer, purchase)

grouped = group txns by customer

total = foreach grouped generate group, sum(txn.purchase);

load customer profile

↳ profile = load 'customers-profile' as (customer, zipcode)

answer = join total by group, profile by customer

Dump answer :

Using Userdefined functions

↳ write code in Java and compile in jar

Pig Latin Snipz

Register myUDF.jar

student = load 'students' using PigStorage(',') as (name);

upper name = foreach student generate myUDF. UPPERCASE(name)

Dump uppername.

## Running Pig

↳ Interactive mode

↳ Pig -n local → local mode

Pig -n mapreduce / map reduce mode  
or  
Pig

it will start command line tool

> grunt >>

## 2) Batch mode

write script of pig in pig script and  
save as id.pig (example)  
running :-  
local mode

pig -n local id.pig

## MapReduce Mode

pig id.pig

or

pig -n MapReduce id.pig.

multiple line comment

/ \* ... \* /

single line comment

--

## Pig Latin statements

All statements take input and produce output  
except LOAD and STORE which read data from  
and write data to the file system.

Must end with semicolon ( ; )

- ↳ A LOAD statement to read data from filesystem
- ↳ A series of transformation statements to process the data
- ↳ A dump statement to view results or a Show statement to save the results.
- DUMP and STORE is required to generate output
- Pig validates and execute LOAD, PIGLET only while running DUMP/STORE

loading part

load function

Filter → work with tuples, row of data  
for each → columns of data

group data in single relation

to group, inner join, outer join → join with how  
many relations

Union → merge

Split → partition the content.

Storing intermediate results in memory → /tmp  
It must be present.

Storage → store final result to the filesystem

Debugging

DUMP 1d display results

Desrib 1de descriptive schema

Explain 1e

Illustrate 1i step by step execution

1q → quit the grust shell.

## slide 2

### Apache Pig

Simple to understand data flow language used in the analysis of large data sets.

Pig scripts are automatically converted to mapreduce jobs by pig interpreter.

Simple to understand data flow language for analysis's familiar with scripting language

↳ fast iterative language with strong MapReduce compilation engine.

↳ Rich, multivalued, nested operations performed on large data sets.

Runs on user machine

↳ job submitted to cluster & executed on cluster

↳ No need to install anything extra on cluster.

Data type → int, long, double, character, byte array

Complex → bag, tuple, Map.

Pig Data Loader

↓

Pig storage

bin storage

binary storage

text loader

Pig Dump

ls = load → loader

gp = group by user → group → combiner

foreach gp generate group, comb(key)  
iterating  
skipping  
group

Storage.

Each statement defines a new dataset,  
possibly in terms of existing datasets.

Each dataset is immutable.

Datasets can be given aliases to use later.  
use dump

Multiple tuples → refined as bags.

foreach ... generate → another bag

## Grouping

groupin  $\rightarrow$  (key, bag)

group key

tuple for every record  
for that key.

group by makes an output bag containing tuples, containing more bags.

group = group key BY user;

in : Bag Of (user, query, time)

out : Bag Of (group, Bag Of (user, query, time), maxed by)

The grouping item is always name('group')

Thinking like Pig

↳ Data flow language  $\rightarrow$

HIVE  $\rightarrow$  like SQL  $\rightarrow$  not data flow language

declarative language

Pig eats anything



any data → relational, nested, semistructured,  
unstructured.



support complex data types like bag, tuples

that can be nested to form fairly  
sophisticated data structure.

## Hive

A distributed warehouse infrastructure build on the top of Hadoop for providing data summarization, querying and analysis

- ↳ RDBMS
- ↳ Structure
- ↳ Access to different storage

Key building principle:-

SQl is - familiar language  
Extensibility → Types, functions, forms,  
scripts  
Performance.

## Intuitive

- ↳ make the unstructured data look like tables regardless how it really looks or.
- ↳ SQL based query can be directed against these tables
- ↳ generate specify execution plan for this query.

## Hive

↳ A data warehousing system to store structured data on Hadoop file system

Provide an easy query these data by execution Hadoop map reduce plans.

## Application

↳ log processing  
Text mining  
Document indexing  
Customer facing (B)  
Predictive modeling, hypothesis testing.

Hive built on top of HDFS

↓  
stores data in HDFS

↓  
Compile SQL queries as MapReduce jobs  
and run the job in the cluster.

## Data Model

↳ Hive structures data into databases  
→ tables, columns, rows

## Tables

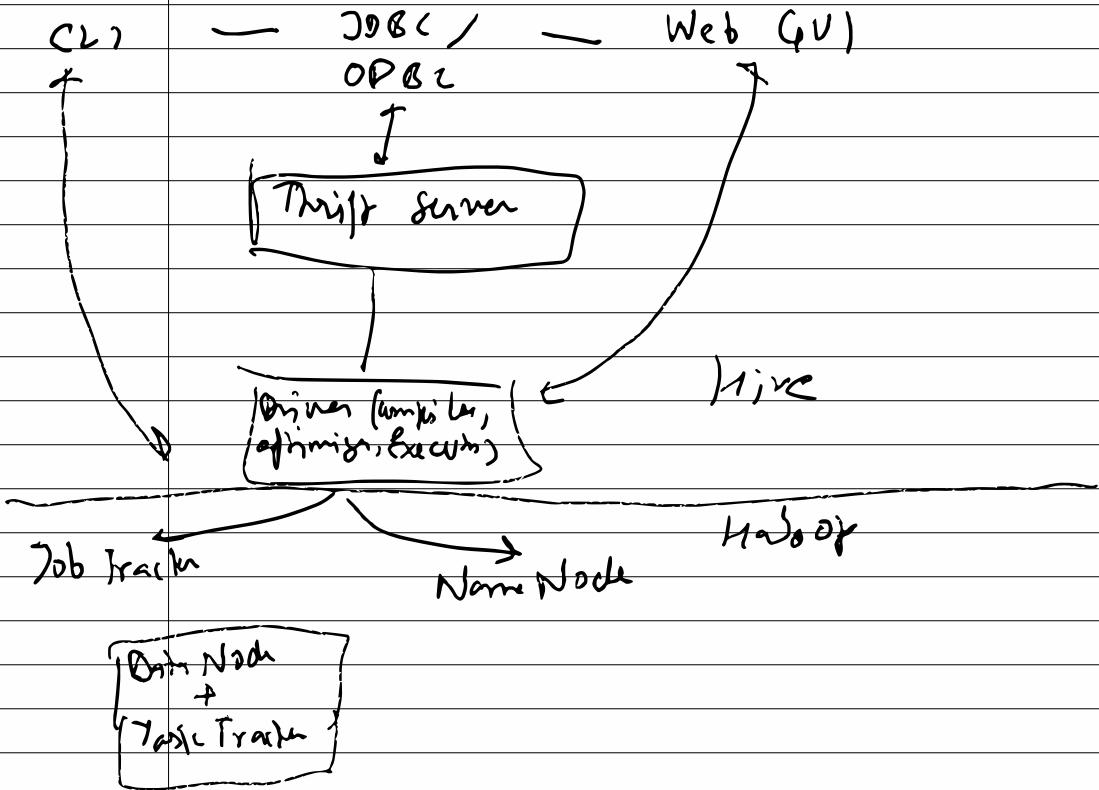
↳ Basic type columns (int, float, boolean)  
↳ complex type : list / map

## Partitions

bucket →

## Architecture

↓



# Architecture

1) External interfaces

↳ provides UI & CLI And API like JDBC & ODBC

2) Thrif2 server

↳ Exposes a very simple client API to execute HiveQL statements.

3) metastore

↳ is the system catalog  
All other components of Hive interact with the metastore

4) Driver

Manages the life cycle of a HiveQL statement during compilation, optimization and execution.

5) Compiler → translates into plan.

↳ translates statements into a plan which consists of a DAG of map-reduce jobs.

The driver submits the individual map-reduce jobs from the DAG to the Execution Engine in a topological order.

## 1) Metastore

↳ system catalog that has metadata about the tables stored in HDFS

## Database

↳ is the namespace for tables

## Table

↳ metadata for table contains list of columns and their types, owner, storage and SerDe information

## Partition

↳ Each partition can have its own column and SerDe and storage information.

## 2) Query Compiler

### Parser

↳ transforms a query string to a parse tree representation

### Semantic Analyzer

↳ transforms the parse tree to a block-based internal query representation

### Logical Plan Generator

↳ converts the internal query representation to a logical plan, which consists of a tree of logical operators.

## optimizer

↳ performs multiple passes over the logical plan and reanalyzes it in several ways.

## Physical Plan Generator

↳ converts logical plan into physical plan consisting of DAG of map-reduce jobs.

## Creating table

D Create table pokes (foo int, bar string) :-

2D create table pokes (foo int, bar string)  
partitioned by (ds string)



partition column called ds  
↳ virtual columns.

Default → tables → text input format  
delimited → ^A( + A + A )

browsing through tables

show tables ;

↳ lists all tables

show tables !. \* s ;

↳ show tables ending with s .

describe table-name ;

↳ shows the list of columns.

Alter table commands

Alter table emp ts RENAME TO 300 below;

Alter table pokes add columns (new\_col int)

Replace column replaces all the existing columns and only change the table's schema not the data .

Replace columns can also be used to drop columns from table's schema

Alter table invites replace column (too int);

↳ it drops the second column

Drop table table-name ; → drops the table .

## DML Operations

Loading data from flat files into Hive :

```
LOAD DATA LOCAL INPATH '/files/k1.txt'  
OVERWRITES INTO TABLE f01;
```



loads a file that contains two columns separated by ctrl-a into f01s table.

'LOCAL' signifies that the input file is on the local file system.

If LOCAL is omitted then it looks into HDFS.

```
LOAD DATA LOCAL INPATH '/(k).txt'
```

↑  
if local file /  
and if loading  
from hdfs  
OVERWRITES INTO TABLE  
f01s partition  
(ds = '2008-08-15');

If the file is in HDFS it is moved into hive unshredded name space.

The operation requires mounting so, it instantaneous.

## SQL operations

Select



select a.\* from invites a where

a.ds = '2008-08-15';

It shows Result; or connect.

INSERT command

To save in directory

```
INSERT OVERWRITE DIRECTORY '/tmp/hdfs_out'
SELECT a.* FROM invites a WHERE
a.ds = '2008-08-15'
```

writes in HDFS

OVERWRITE LOCAL DIRECTORY



It will write in local file system.

# Performance

## 1) Group by operation

Efficient execution plan based on:

Data skew:

- now evenly distributed data across number of physical nodes
- bottleneck vs load balance

Partial Aggregation:

- Group the data with the same group by value as soon as possible.
- In memory hash table for merger.
- Earlier than combine.

## 2) Join operation

Traditional map reduce join  
Early map side join

- ↳ very efficient for joining small table with large table.
- ↳ keep smaller table in memory first
- ↳ join with a chunk of larger table data each time
- ↳ spare complexity from time complexity.

### 3) SQL

- ↳ How to load data from file to make it look like a table

### 4) Legacy way

- ↳ Create the field object when necessary
- ↳ Reduce the overhead to create unnecessary objects in Hive.
- ↳ Java is expensive to create objects
- ↳ increases performance.

### Pros

- ↳ Easy way to process large scale data
- ↳ Support SQL-based queries
- ↳ provide more user defined interfaces to extend
- ↳ programmability
- ↳ efficient execution plans for performance
- ↳ interoperability with other database tools.

### Cons :-

- No easy way to append data
- files in HDFS are immutable

Hive component →

Shell interface : like the MySQL shell

Pipeline

↳ session handles, fetch, execution

Compiler

↳ Parse, plan, optimize.

Execution Engine

↳ DAG stage

↳ run map or reduce.

Data Units

↳ Data bases

↳ tables

↳ partitions

↳ buckets -

Type System

Primitive Types

- Integers
- Boolean
- Floating point numbers
- String.

Complex types

- structs
- Maps
- Array.

User defined functions

↳ Java code

Registering the class

↳ Create function my-lower as 'com.example';

Using the function

Sum(key).

→ select my-lower(hi)e, group by es

group by my-lower(hi)e )

## Oozie

Workflow scheduler for Hadoop



Designed at Yahoo! for their complex search engine workflows.



Now an open source Apache incubator project.

Allows user to create DAG for workflow



Can run parallelly / sequentially in hadoop.

→ Can run plain java classes, pig work flows and interact with HDFS

↳ Nice if you want to delete or move files before a job runs.

→ Major flexibility

↳ start, stop, suspend and re-run jobs.

→ Restart after failure.

## Features

- ↳ Java client API / command line Interface
- ↳ Web service API
- ↳ Run periodic jobs.
  - ↳ hourly, daily, weekly etc.
- ↳ Receive email when job is complete.

## How to make workflow

- ↳ make job
- ↳ make jar
- ↳ make workflow.xml with job configurations in it
  - ↓
    - input files
    - output files
    - input readers & writers
    - mappers and reducers
    - job specific arguments

Also need job.properties file. This file defines namenode, job tracker etc.

It also gives the location of the stored  
~~jobs~~ jars and other files.

Only really → run from command line.

Start Node → tells app where to start

<start to = "[Name - Node]" />

End Node → signals end of logic job

<end name = "[Name - Node]" />

Error Node - signals that an error occurred and a message describing the error should be printed out.

<error name = "[Node-name]" />

<message> "[A custom message"]

</message>

</error>

Logic fork join

Fork → Starts the parallel jobs

join → parallel jobs rejoin at this node. All forked jobs must be completed to continue the workflow.

Can also define decision nodes like switch statements.

for parameters in  $\$ \{ NAME NODE \}$

↓

skip nodes

from failed nodes.

can do either wr or bmr.

Sentiment Analysis

---