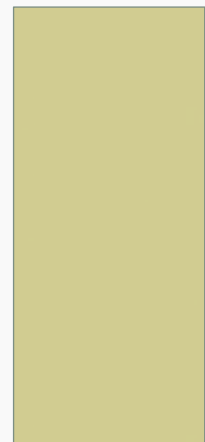


# PROCESSING BIG DATA WITH HADOOP

AHMET DEMIRELLI



# OUTLINE

- Why is Hadoop?
- What is Hadoop?
- When is Hadoop?
- Where is Hadoop?
- How is Hadoop?

# WHY HADOOP?

- Databases can process GBs of data
- TBs is too much for a RDBMS
- Solution
  - Faster processors & More Memory ??
  - Distributed Systems
    - Synchronisation ?
    - Bandwidth ?
    - Node Failure ?
    - Job fails ?

# WHAT IS HADOOP?

- <http://hadoop.apache.org>
- Based on papers published by Google
- GFS - <http://research.google.com/archive/gfs.html>
- MapReduce - <http://research.google.com/archive/mapreduce.html>
- Implemented by Apache Software Foundation
- A Platform provides both
  - Distributed storage (HDFS)
  - Distributed computation (MapReduce)

# WHAT IS HADOOP?

- Application developed in high level Programming Language (Java,Python,Scala,Pig, Hive..etc)
- Very little communication between nodes
- Data is distributed and replicated in advance (HDFS)
- Data is processed on the stored location
- Hadoop scalable and fault tolerant

# WHAT IS HADOOP?

- Scalable
  - Adding new nodes increases the system capability proportionally
- Fault Tolerance (If one of the nodes goes down)
  - No loss of data (replication)
  - No loss of tasks (failed task will be reassigned to another node)
  - Recovered nodes will rejoin to the cluster



# WHAT IS HADOOP?



Hadoop Ecosystem



Hadoop Core Components  
(Hadoop 1.0)

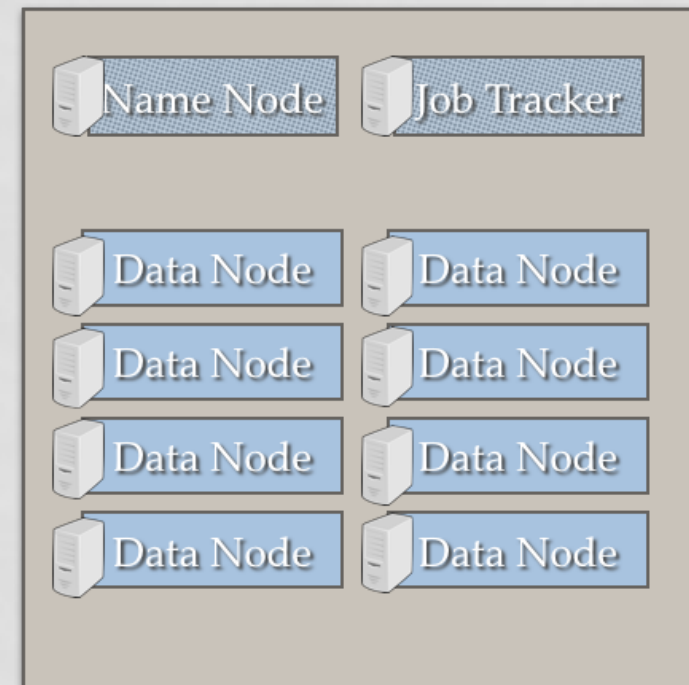
# WHAT IS HADOOP?

- HDFS (Hadoop Distributed File System)
  - Stores data on the cluster
  - A File system written in Java
  - Runs on top of native file system (like ext3,ext4,xfs)
  - Designed for dealing with massive amount of data
  - Keeps files Read-Only (for performance)
  - Performs best with large files (100MB or more)
  - Each block is replicated on 3 different nodes by default
  - Information about the files and locks (Metadata) is kept on the Name Node
- MapReduce
  - Processes data on the cluster
  - Manages Jobs

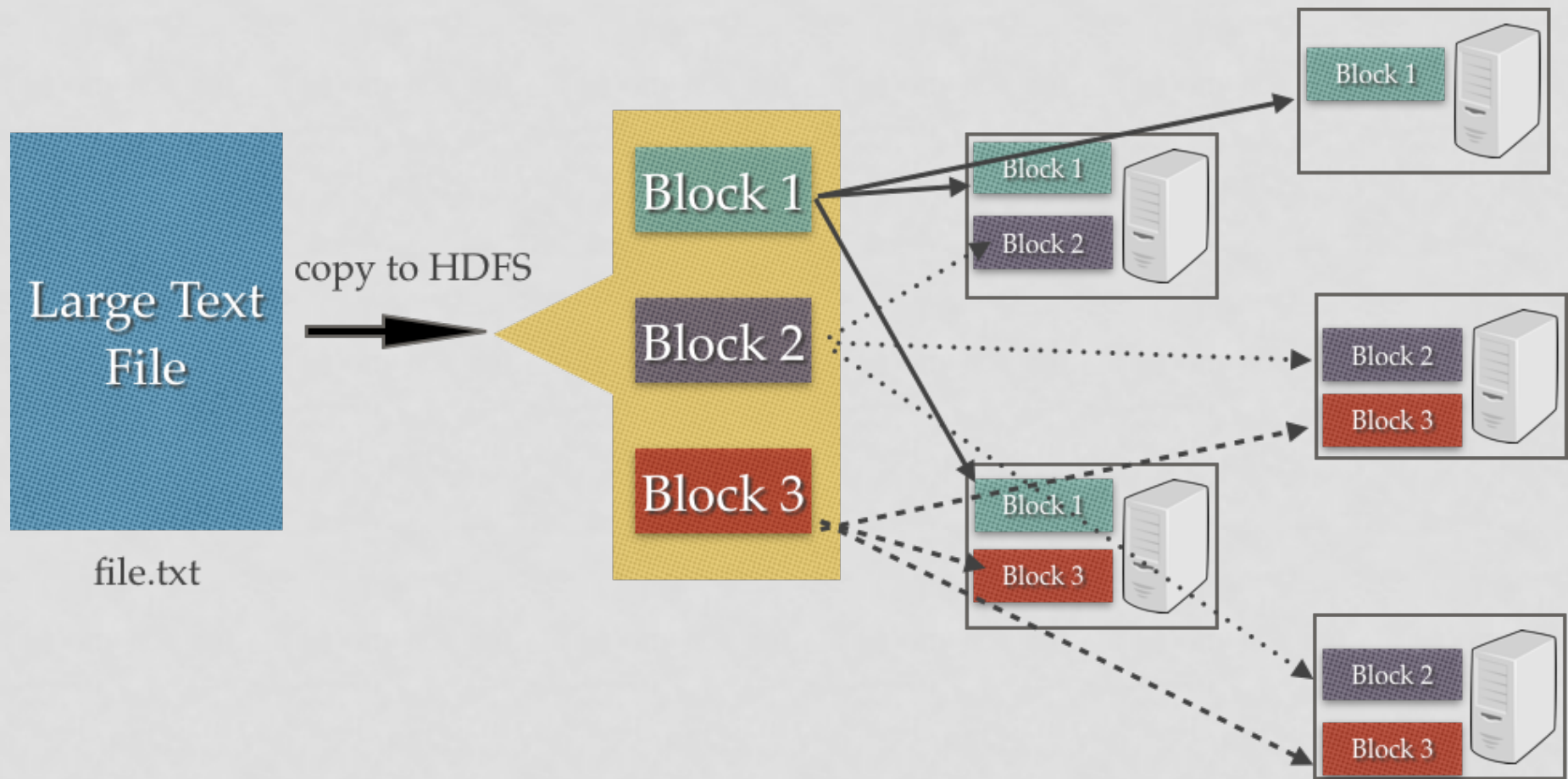


# WHAT IS HADOOP?

- Hadoop Cluster : Group of computers working together store and process the big data
- Hadoop Cluster Nodes ( Hadoop 1.0)
  - master nodes
    - cluster has two master nodes
      - Name Node (manages HDFS )
      - Job Tracker ( manages MapReduce)
  - slave (worker) nodes
    - HDFS
    - MapReduce



# WHAT IS HADOOP?



# WHEN IS HADOOP?

- You have massive data ( larger than GBs )
- Or operations that takes time on regular data
- And most importantly you have a parallelizable algorithm
- Or a better motivation
  - To get a good grade from this course 😊

# WHERE IS HADOOP?

- You can install from scratch (Bad idea)
- You can use a distribution (Easy Install)
  - Cloudera
  - Hortonworks
  - MapR
- You can download quick-start-vm from (No install)
  - Cloudera
  - Hortonworks

# WHERE IS HADOOP?

- Google Cloud
  - <http://cloud.google.com>
- Amazon EMR
  - <https://aws.amazon.com/elasticmapreduce>
- Microsoft Azure
  - <https://azure.microsoft.com>
- Ibm Bluemix
  - <https://console.ng.bluemix.net/catalog>
- Oracle Cloud
  - [https://cloud.oracle.com/en\\_US/bigdata?tabID=1448073786401](https://cloud.oracle.com/en_US/bigdata?tabID=1448073786401)

# HOW IS HADOOP?

- Demos
  - MapRaduce (Java)
  - Hive
  - Pig
  - Spark (Python)

# MAP REDUCE

- Framework for processing data parallel on HDFS
- Written in Java
- Map -> Shuffle -> Reduce steps

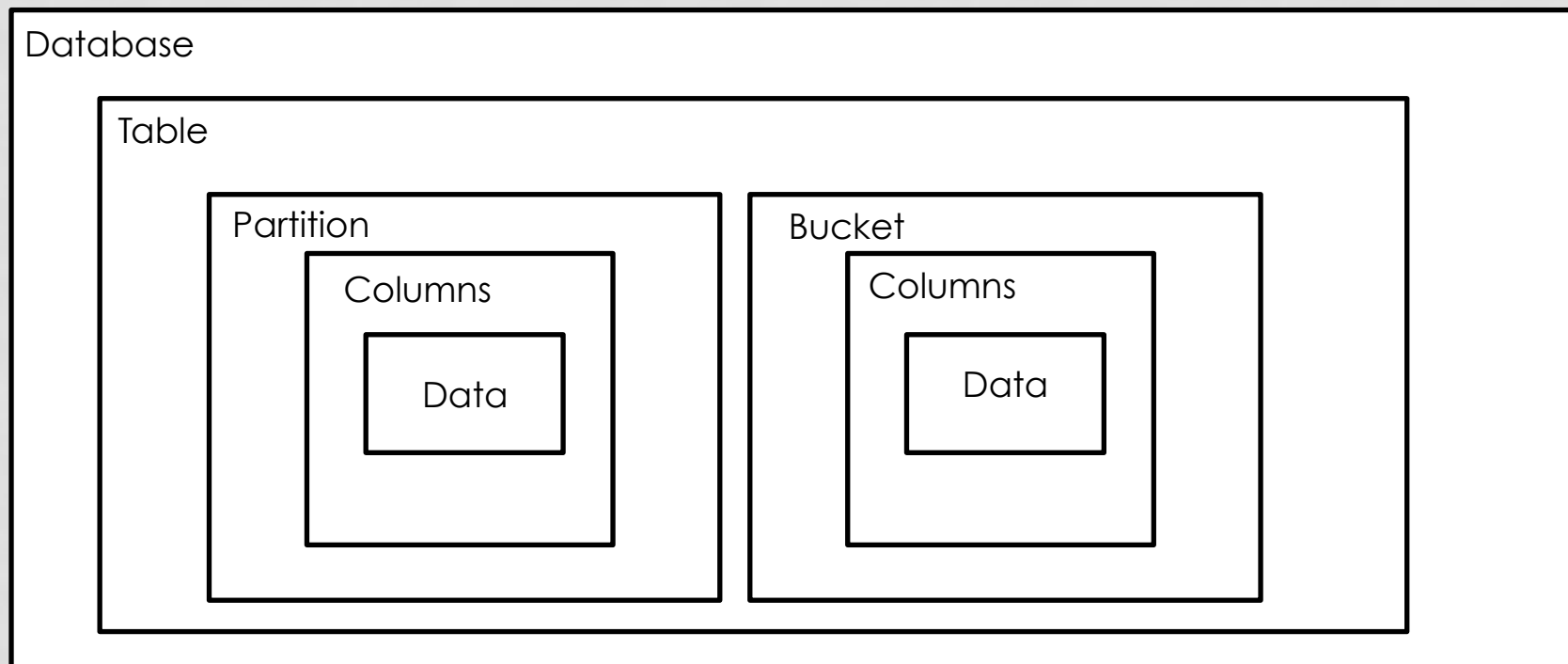


# HIVE

- <http://hive.apache.org>
- Data warehouse system for Hadoop
- ETL operations
- Easy data summarization and analysis of large volumes of data.
- Hive QL ( based on SQL )
- Query execution in MapReduce
- Compiled in to MapReduce jobs ( limitations )
- Runs with structured data (like csv )

# HIVE DATA UNITS

- Databases , Tables , Partitions , Buckets (or Clusters)



# HIVE PRIMITIVE TYPES

- Integers
  - TINYINT - 1 byte integer
  - SMALLINT - 2 byte integer
  - INT - 4 byte integer
  - BIGINT - 8 byte integer
- Boolean type
  - BOOLEAN - TRUE/FALSE
- Floating point numbers
  - FLOAT - single precision
  - DOUBLE - Double precision
- String type
  - STRING - sequence of characters in a specified character set

# HIVE

## HIVE QL

- Create tables, views and indexes
- Select ,where , group by, order by, joins
- No complex, correlated sub queries
- Not for small amounts of data
- No union, intersect, distinct queries

# HIVE LIKE SQL

## Hive SQL Datatypes

INT

TINYINT/SMALLINT/BIGINT

BOOLEAN

FLOAT

DOUBLE

STRING

TIMESTAMP

BINARY

ARRAY, MAP, STRUCT, UNION

DECIMAL

CHAR

VARCHAR

DATE

## Hive SQL Semantics

SELECT, LOAD, INSERT from query

Expressions in WHERE and HAVING

GROUP BY, ORDER BY, SORT BY

Sub-queries in FROM clause

GROUP BY, ORDER BY

CLUSTER BY, DISTRIBUTE BY

ROLLUP and CUBE

UNION

LEFT, RIGHT and FULL INNER/OUTER JOIN

CROSS JOIN, LEFT SEMI JOIN

Windowing functions (OVER, RANK, etc.)

INTERSECT, EXCEPT, UNION DISTINCT

Sub-queries in WHERE  
(IN/NOT IN, EXISTS/NOT EXISTS)

Sub-queries in HAVING

	Hive 0.10
	Hive 0.11
	Future

# HIVE LIKE SQL

## Metadata

### Function

Selecting a database

Listing databases

Listing tables in a database

Describing the format of a table

Creating a database

Dropping a database

### MySQL

```
USE database;
```

```
SHOW DATABASES;
```

```
SHOW TABLES;
```

```
DESCRIBE table;
```

```
CREATE DATABASE db_name;
```

```
DROP DATABASE db_name;
```

### Hive

```
USE database;
```

```
SHOW DATABASES;
```

```
SHOW TABLES;
```

```
DESCRIBE (FORMATTED|EXTENDED) table;
```

```
CREATE DATABASE db_name;
```

```
DROP DATABASE db_name (CASCADE);
```

# HIVE LIKE SQL

## Function

Retrieving Information (General)

Retrieving All Values

Retrieving Some Values

Retrieving With Multiple Criteria

Retrieving Specific Columns

Retrieving Unique Output

Sorting

Sorting Reverse

Counting Rows

Grouping With Counting

Maximum Value

Selecting from multiple tables (Join same table using alias w/"AS")

## MySQL

```
SELECT from_columns FROM table WHERE  
conditions;
```

```
SELECT * FROM table;
```

```
SELECT * FROM table WHERE rec_name =  
"value";
```

```
SELECT * FROM TABLE WHERE rec1 =  
"value1" AND rec2 = "value2";
```

```
SELECT column_name FROM table;
```

```
SELECT DISTINCT column_name FROM  
table;
```

```
SELECT col1, col2 FROM table ORDER BY  
col2;
```

```
SELECT col1, col2 FROM table ORDER BY  
col2 DESC;
```

```
SELECT COUNT(*) FROM table;
```

```
SELECT owner, COUNT(*) FROM table  
GROUP BY owner;
```

```
SELECT MAX(col_name) AS label FROM  
table;
```

```
SELECT pet.name, comment FROM pet,  
event WHERE pet.name = event.name;
```

## Hive

```
SELECT from_columns FROM table WHERE  
conditions;
```

```
SELECT * FROM table;
```

```
SELECT * FROM table WHERE rec_name =  
"value";
```

```
SELECT * FROM TABLE WHERE rec1 =  
"value1" AND rec2 = "value2";
```

```
SELECT column_name FROM table;
```

```
SELECT DISTINCT column_name FROM  
table;
```

```
SELECT col1, col2 FROM table ORDER BY  
col2;
```

```
SELECT col1, col2 FROM table ORDER BY  
col2 DESC;
```

```
SELECT COUNT(*) FROM table;
```

```
SELECT owner, COUNT(*) FROM table  
GROUP BY owner;
```

```
SELECT MAX(col_name) AS label FROM  
table;
```

```
SELECT pet.name, comment FROM pet JOIN  
event ON (pet.name = event.name)
```



# HIVE

## RUN HIVE QUERIES

### Command Line

Function	Hive
Run Query	<code>hive -e 'select a.col from tab1 a'</code>
Run Query Silent Mode	<code>hive -S -e 'select a.col from tab1 a'</code>
Set Hive Config Variables	<code>hive -e 'select a.col from tab1 a' -hiveconf hive.root.logger=DEBUG,console</code>
Use Initialization Script	<code>hive -i initialize.sql</code>
Run Non-Interactive Script	<code>hive -f script.sql</code>

# HIVE

## SAMPLE DATASET

- 144 MB
- 22,000,000 ratings and 580,000 tag applications applied to 33,000 movies by 240,000 users. Last updated 1/2016
- The data are contained in four files,
  - movies.csv (movieId,title,genres)
  - ratings.csv (userId,movieId,rating,timestamp)
  - tags.csv (userId,movieId,tag,timestamp)
  - links.csv (movieId,imdbId,tmdbId)
    - movieId is an identifier for movies used by <https://movielens.org>. E.g., the movie Toy Story has the link <https://movielens.org/movies/1>.
    - imdbId is an identifier for movies used by <http://www.imdb.com>. E.g., the movie Toy Story has the link <http://www.imdb.com/title/tt0114709/>.
    - tmdbId is an identifier for movies used by <https://www.themoviedb.org>. E.g., the movie Toy Story has the link <https://www.themoviedb.org/movie/862>.

# HIVE

## SAMPLE DATASET

- Upload data to your Google Cloud Bucket
- Copy data from you bucket to master instance
  - `gsutil cp gs://<your-bucket>/yourfilename.csv .`
- Create Table
  - `hive> CREATE TABLE ratings ( userid STRING, movieid STRING, rating DOUBLE, unixtime STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY ' ' STORED AS TEXTFILE;`
- Load data
  - `hive> LOAD DATA LOCAL INPATH ' /usr/ahmet/ratings.csv' OVERWRITE INTO TABLE ratings;`
- Query data
  - `hive> SELECT COUNT(*) FROM ratings;`

# HIVE ACCESS

- Hive Shell
- Java, Python using JDBC Driver
- .NET ODBC Driver
- Apache Thrift Client

# PIG

- <http://pig.apache.org>
- Pig Engine parses, optimizes and runs scripts as MapReduce jobs
- Pig-Latin high level scripting language
- ETL operations
- No need for Java, Scala or Python knowledge
- Is able to store data at any point during a pipe line

# PIG COMMANDS

- Load data
  - `Movies = LOAD '/user/hadoop/movies.txt' USING PigStorage(',') as (id,name,year,rating,duration);`
- Load data with types
  - `Movies = LOAD '/user/hadoop/movies.txt' USING PigStorage(',') as (id:int,name:chararray,year:int,rating:float,duration:int);`
- Check format of data
  - `Describe Movies;`
- Filter
  - `filtered= FILTER Movies BY rating>3.5;`

# PIG COMMANDS

- Select part of data and assign to another variable
  - selected= **foreach** filtered **generate** year, rating, name;
  - Dump selected
- Store processed data to HDFS
  - STORE selected INTO '/user/hadoop/processed.csv' USING PigStorage (',');



# APACHE SPARK

- Word count example