# From Relational to Hadoop
## Part 1: Introduction to Hadoop

Gwen Shapira , Cloudera
and
Danil Zburivsky, Pythian

# Tutorial Logistics

# Got VM?

- Grab a USB
- USB contains:
  - Cloudera QuickStart VM
  - Slides
  - Exercises and Solutions
  - Data

I talk about Hadoop.

You multi-task between:
- Listening to me talk about Hadoop
- Installing Virtualbox and Running the VM

PERCONA LIVE

# Using VirtualBox

1. Install and open VirtualBox on your computer
2. Under the menu "File", select "Import…"
3. Navigate to where you unpacked the .ovf file
4. and select it

5. You will find a "troubleshooting" file on the USB

# Agenda – First 90 minutes

- Why use Hadoop for ETL?
- How HDFS Works?
- How MapReduce Works?
- What other tools we'll use?
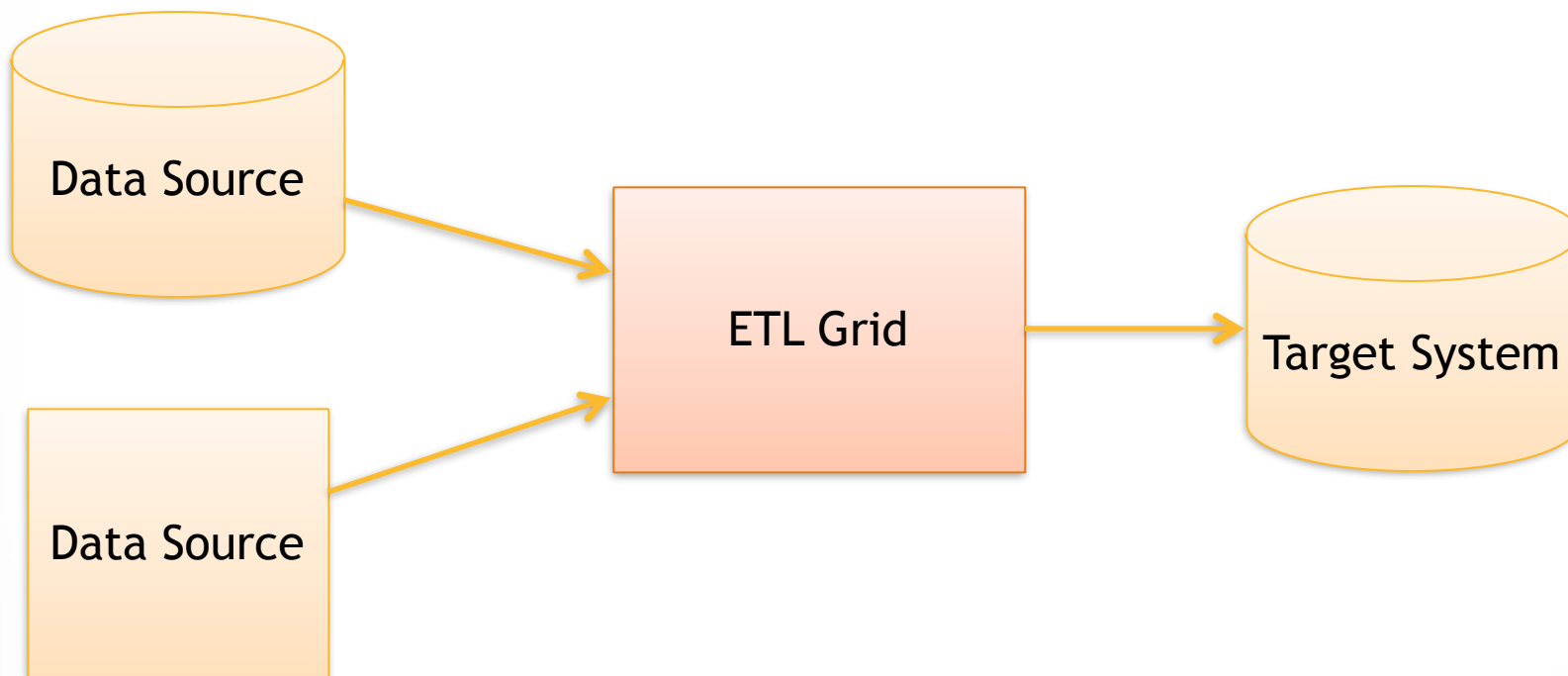
# Use Hadoop for ETL?

# ETL

- Extracting data from outside sources
- Transforming it to fit operational needs
- Loading it into the end target

- (Wikipedia: http://en.wikipedia.org/wiki/Extract,_transform,_load)

- Transform data before loading it to DB
  - XMLs, JSONs...
- Clean and standardize data
  - Units, names...
- Aggregate data, recommendations
- De-normalize for a DWH
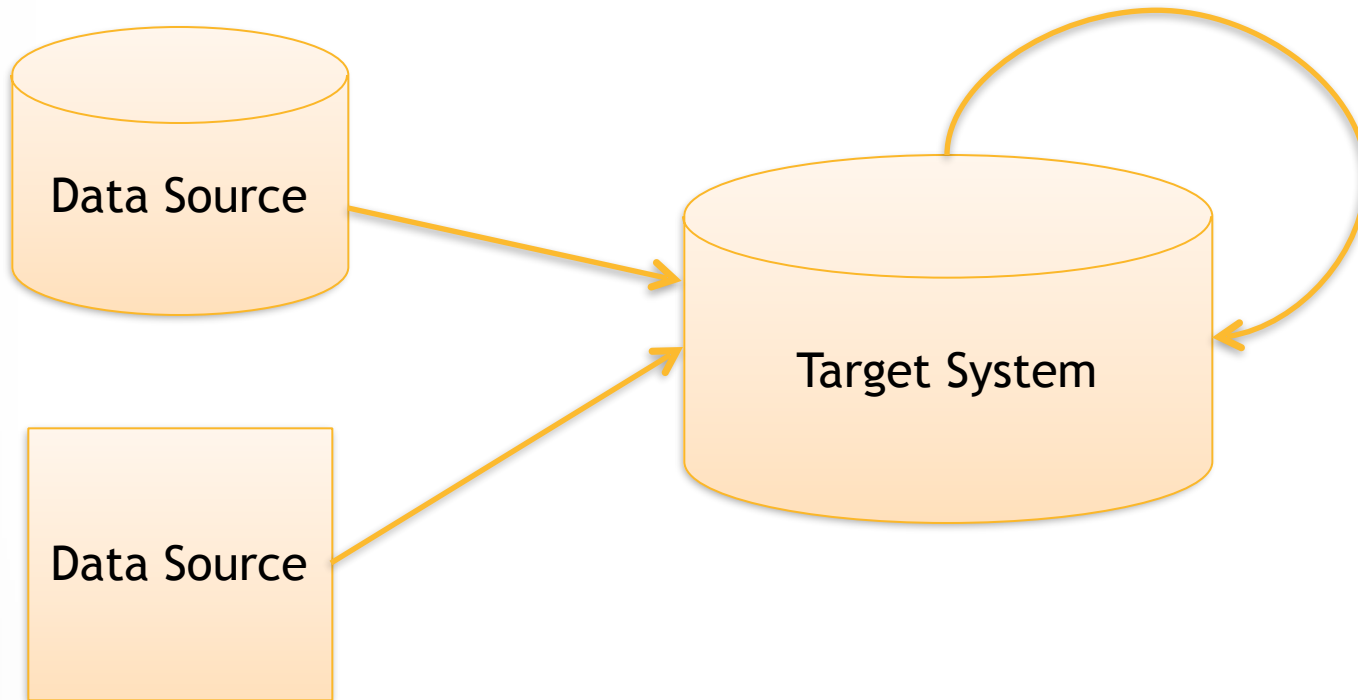
# Some people do it like this:

# One Problem

**8 hour processing pipeline can fail after 7.5 hours**

# OK, two

# Or like this:

# Few problems here too

- Unstructured data is a bigger challenge
- More contenders for CPU resources
  - Analysts
  - And even transactions
- Storage can get expensive
- Scaling DWH is a pain
- SQL is a bit limited

# Use Hadoop!

- Handles unstructured data
- Scales horizontally
  - Storage
  - CPU
- Flexible and powerful
- Lots of tools



PERCONA LIVE

# In this tutorial we'll show:

- Basics of using Hadoop
- Transform unstructured data
- Get RDBMS data to Hadoop
- Use Hadoop to join and aggregate
- Load the results back to an RDBMS

PERCONA LIVE

# HDFS:
Distributed, fault-tolerant file system

# Design Assumptions

- Failures are common
  - More scale == more failure
- Files are append-only
- Files are large (Gigabytes +)
- Access is large and sequential

# Quick Disk Primer

- Disk does a seek for each I/O operation
- Seeks are expensive (~10ms)
- Throughput / IOPS tradeoff
  - 100 MB/s and 10 IOPS
  - 10MB/s and 100 IOPS
- Big I/Os mean better throughput

# Quick Networking Primer

Core
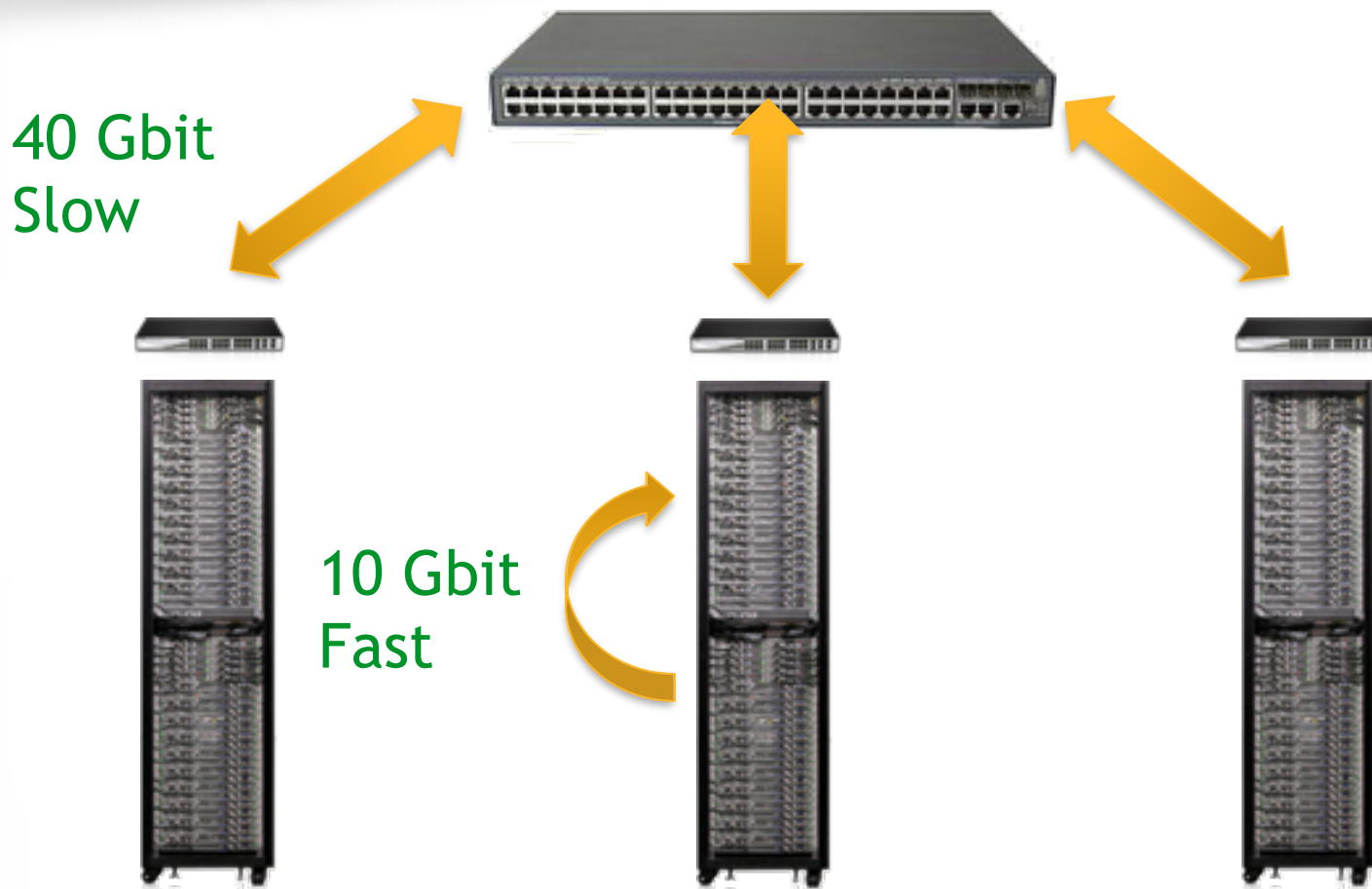switch

Top-of-rack switch

Rack

PERCONA
LIVE

# Quick Networking Primer



40 Gbit
Slow

10 Gbit
Fast

Metadata ⟶ NameNode

Paths,
filenames, file
sizes, block
locations, ...

Data blocks

DataNode   DataNode   DataNode   DataNode

NameNode

create("/tmp/myfile")

Client

Write to
[DN4,DN3,DN2]

[DN3,DN2
]

[DN2]

DN 1    DN 2    DN 3    DN 4

Rack 1    Rack 2

NameNode

open("/tmp/myfile","r")

Client

Read from
[DN4,DN3,DN2]

Data    Read

DN 1    DN 2

DN 3    DN 4

Rack 1    Rack 2

# Using HDFS

- HDFS – command line interface
- hdfs dfs –mkdir
- hdfs dfs –ls
- hdfs dfs –put
- hdfs fsck

# Practice Time

- Create HDFS directory:
  - /etl/earthquakes/landing
- Write earthquake data to directory
- Get first 10 lines of file
- How many blocks we have in the file?
- What is the replication factor?

# Solution

- ```
  sudo -u hdfs hdfs dfs -mkdir /etl
  sudo -u hdfs hdfs dfs chown
  cloudera:cloudera /etl
  hdfs dfs -mkdir /etl/earthquakes/landing
  ```

- ```
  hdfs dfs -put ~/datasets/earthquakes.json /etl/earthquakes/landing
  ```

- ```
  hdfs dfs -cat /etl/earthquakes/landing/earthquakes.json | head -10
  ```

- ```
  hdfs fsck /etl/earthquakes/landing/earthquakes.json
  ```

# MapReduce:

Programming and execution framework

# Just implement two functions

- Map:
  - Operate on every element
  - Filter, transform
- Reduce:
  - Combine and aggregate results

**Good news**

- You don't need to know MR
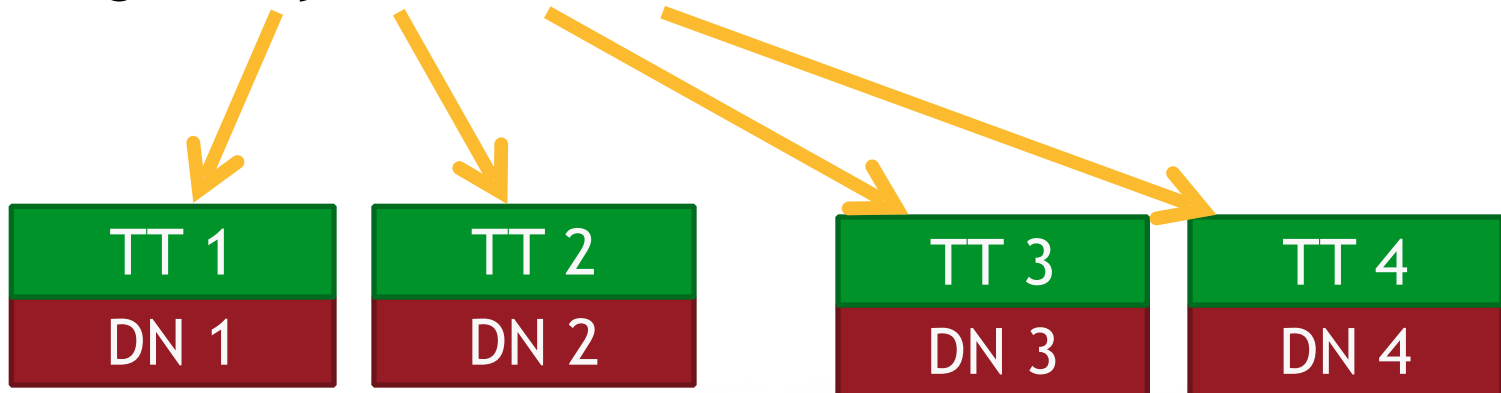- Many abstractions
- Alternative frameworks

**Bad news**

- You still need to know MR
- To understand how things work
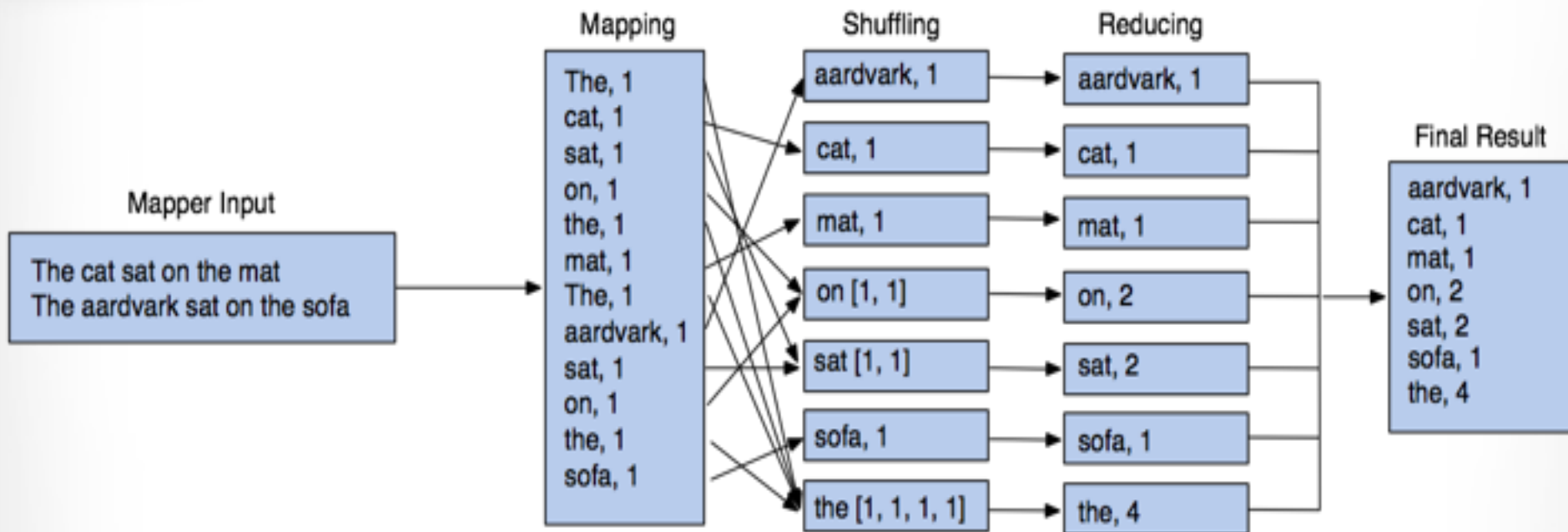- Still widely in use

# Map Reduce Architecture

NameNode | JobTracker

- Gateway for users
- Assigns tasks to TaskTrackers
- Tracks job status

TaskTrackers execute Map and Reduce tasks assigned by JT

TT 1
DN 1

TT 2
DN 2

TT 3
DN 3

TT 4
DN 4

# WordCount

# Running MapReduce

NameNode  JobTracker  ← `wordcount(<files>)`

[cat, 1] [the, 1] [dog, 1]  [sat, 1]

| M1 | M2 | M3 | M4 | | R1 | |
|----|----|----|----|---|----|---|
| TT 1 | | TT 2 | | | TT 3 | TT 4 |
| DN 1 | | DN 2 | | | DN 3 | DN 4 |

NameNode JobTracker ← `wordcount(<files>)`

[mat, 1][cat, 1] [bad, 1] [for, 1]

| M1 | M2 | M3 | M4 | | R1 |
|----|----|----|----|--|----|

| TT 1 | TT 2 | | TT 3 | TT 4 |
|------|------|--|------|------|
| DN 1 | DN 2 | | DN 3 | DN 4 |

# MapReduce

NameNode    JobTracker    ← wordcount(<files>)

```
[a,    5]
[cat, 2]
[dog, 1]
[the, 4]
[mat, 1]
```

M1  M2    M3  M4    R1

TT 1    TT 2    TT 3    TT 4

DN 1    DN 2    DN 3    DN 4

PERCONA LIVE

# Using MapReduce

- Submit a JAR

- Specify the class that contains the mapper and reducer

- hadoop jar <jar> <class> <parameters>

# Practice Time

- Load "Works of Shakespeare" or other file to HDFS
- Count words
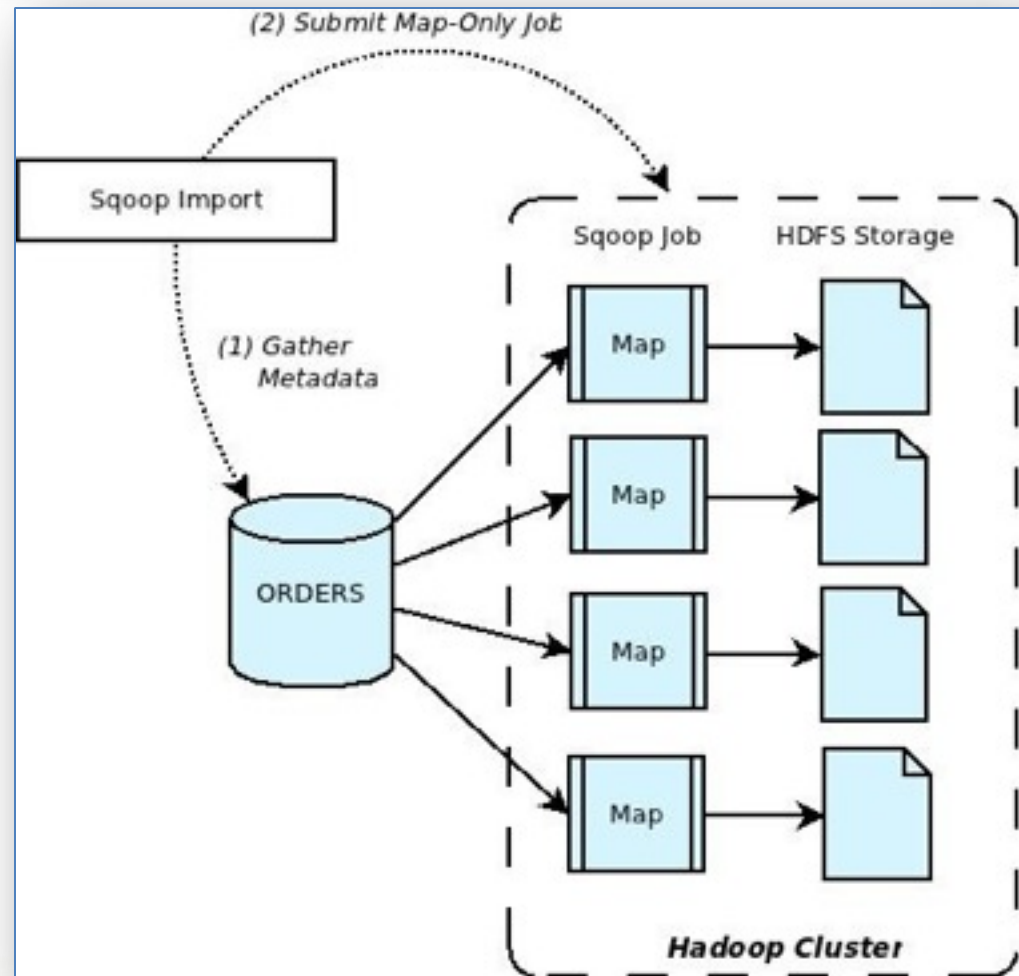- Generate 100M of random data.
- Run "terasort".

# Solution

- `hdfs dfs -put ~/pl_tutorial/datasets/ shakespeare.txt /user/cloudera`

- `hadoop jar /usr/lib/hadoop-0.20-mapreduce/ hadoop-examples.jar wordcount /user/cloudera/ shakespeare.txt /user/cloudera/cnt`

- `hadoop jar /usr/lib/hadoop-0.20-mapreduce/ hadoop-examples.jar teragen  1000000 /user/ cloudera/terasort-in`

- `hadoop jar /usr/lib/hadoop-0.20-mapreduce/ hadoop-examples.jar terasort /user/cloudera/ terasort-in /user/cloudera/terasort-out`
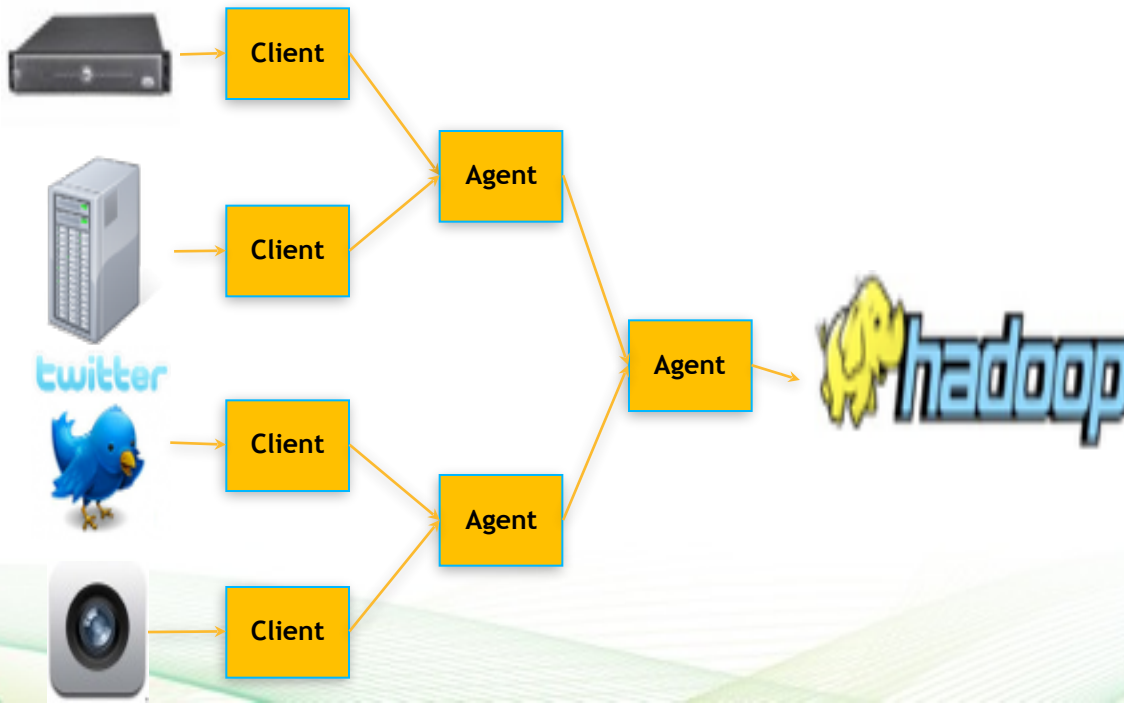
# Ecosystem Tools

# Sqoop

Transfers data
between Hadoop
and almost any
SQL database
with a JDBC driver

# Flume

Streaming data collection

And aggregation for:
JMS queues, HTTP APIs,  Log4J, Syslog, etc.

# Hive

Translate SQL to MapReduce
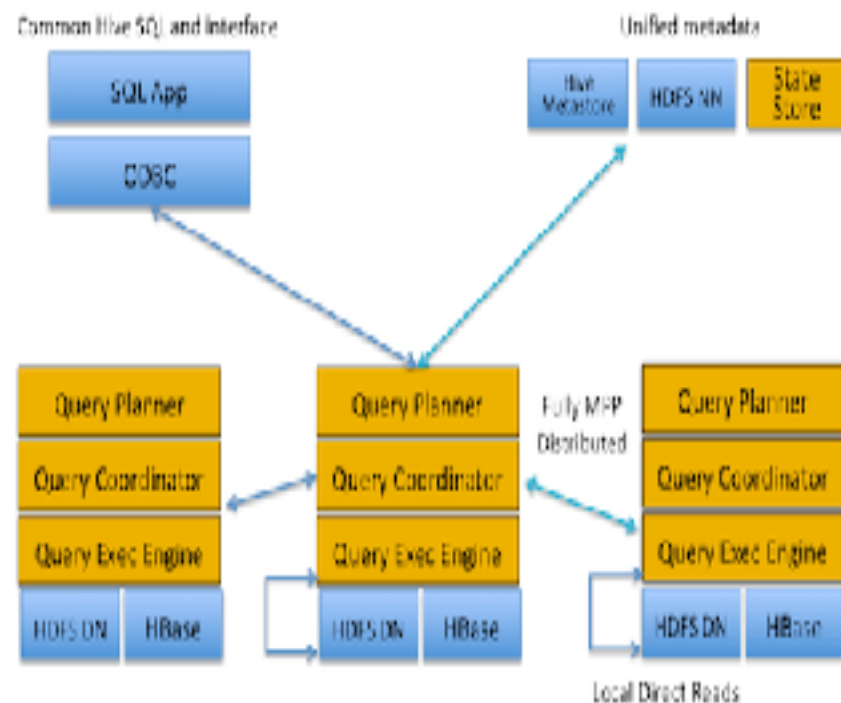
Select word,count(*) from shakespear
Group by word
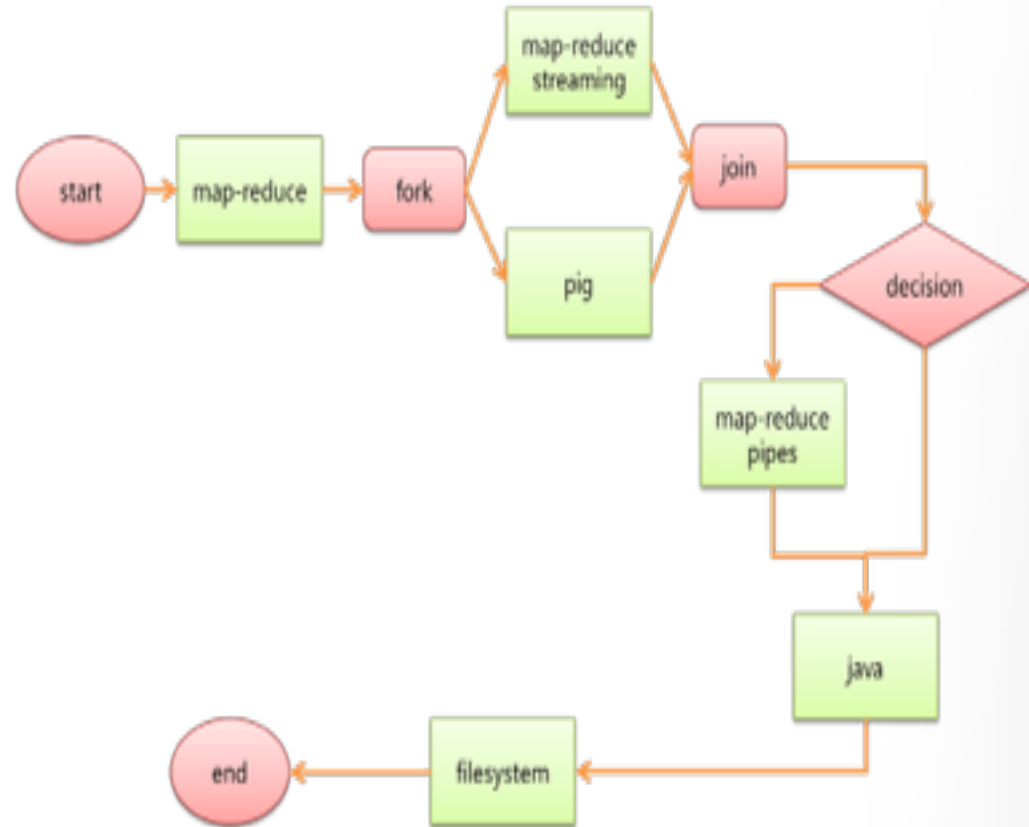
Modern MPP database
built on top of HDFS

Really fast! Written in C++
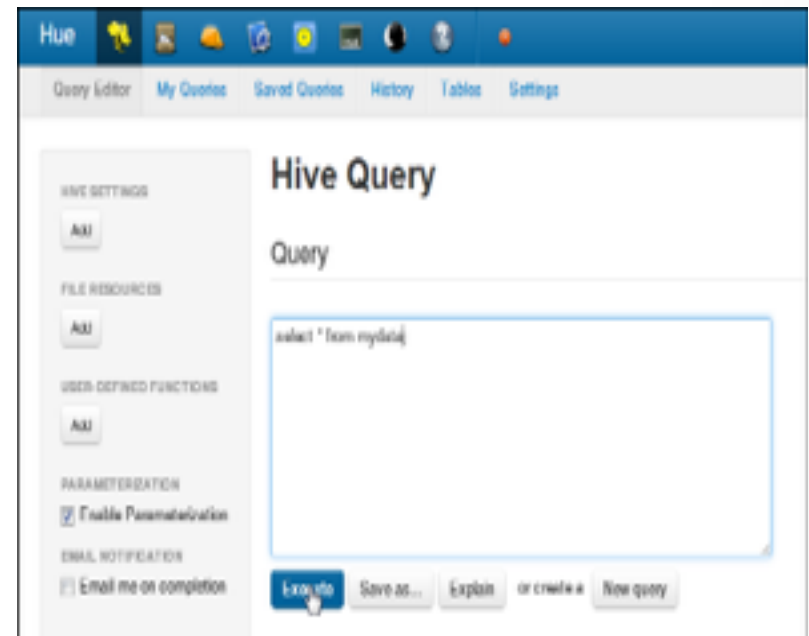
10-100x faster than Hive

# Oozie

A workflow engine and scheduler built specifically for large-scale job orchestration on a Hadoop cluster

# Hue

- Hue is an open source web-based application for making it easier to use Apache Hadoop.
- Hue features
  - File Browser for HDFS
  - Job Designer/Browser for MapReduce
  - Query editors
  - Oozie

# Practice Time

- Login to Hue
- Play around
- We will dive into Sqoop, Hive and Oozie in the next hour