

Intro to Hadoop and MapReduce

~~more~~ more and more data is being generated

~~1/2 of world's data is generated every day~~

Big Data

"Data too big to be processed on a single machine"

~~created fast~~

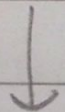
~~comes from different sources and with different formats~~

3Vs: volume, variety, velocity

size of
the data

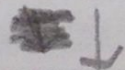
data coming
from a lot of
different sources
and formats

speed at which data
~~is generated~~
arrives, ready to be
processed



price to store
data has dropped
incredibly.
(storing data
reliably is more
expensive)

Some companies
throw away data
because it is too
expensive to store
it reliably, like logs



With databases
data must fit
in ~~pre~~ pre-defined
tables, but most
data is

unstructured or
semi-structured

store data in
original format
so you're not
throwing anything
away

↓
store data at a rate
of also TB/day

Predict what
the user wants
for example

On ~~hadoop~~ hadoop you
can store data in its
raw format, and manipulate
it and reformat it later

→ All data
is useful
at some point

~~What is Hadoop?~~ What Hadoop!

Hadoop consists of a way to store data, known as the Hadoop Distributed File System, and a way to process data with MapReduce.

Data is split and stored across the collection of machines, or cluster, when we want to process the data, we process it where it is actually stored.

Machines can be ~~added~~ ^{added} to the cluster as the amount of data grows.

The machines don't need to be high end.

Hadoop ecosystem!

~~to make it easier to use~~

to make Hadoop easier to use

- Hive ~~help~~ ^{help} querying data without knowing how to ~~code~~ ^{code} ~~mapreduce~~ ^{mapreduce} is by

SQL entered is then interpreted in mapreduce code and run in the cluster

- Pig, analyse data in a fairly simple scripting language rather than mapreduce, also here the code is turned into mapreduce

- Impala, queries data in SQL, but directly accesses it from HDFS, instead of running mapreduce jobs. It is optimized for low-latency queries (while hive is optimized for running long batch processing jobs)

- Sqoop, takes data from a relational database and puts it in HDFS

- Flume ingests data as generated by external systems, and puts it into the cluster
- HBase is a realtime database, built on top of HDFS
- Hue, graphical front end to the cluster
- Oozie, workflow ~~program~~ management tool
- Mahout, machine learning library

Making all these talk to one another and work well can be tricky, to make installing and maintaining a cluster like this easier, Cloudera, has put ~~together~~ together a distribution ~~called~~ of Hadoop called CDH

It takes all the ecosystem project and Hadoop, and packages them together so that installation ~~is~~ is easy and they compatible together ~~makes sure~~

HDFS

When a file is loaded into HDFS, ~~it~~ is splitted into chunks, which we call blocks.

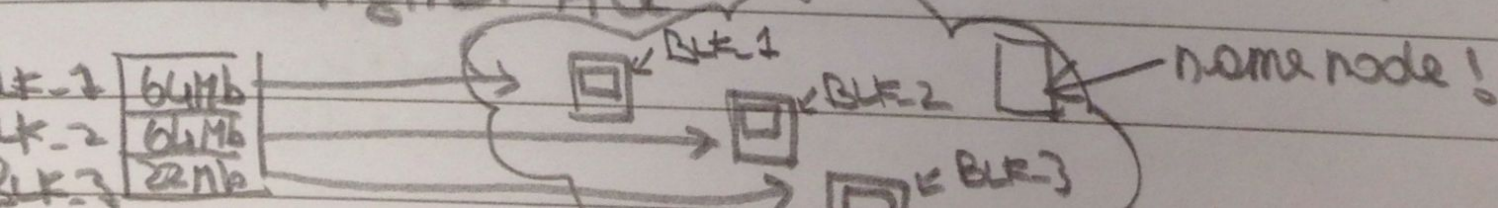
Each block is ~~64~~ 64MB by default, and ~~has~~ has a unique name (BLK_ some large number).

As the file is uploaded to HDFS, each block will get stored on one node in the cluster.

A daemon is running on each of the machines in the cluster ~~each machine is called data node~~

~~is~~ A separate machine called datanode runs a daemon called NameNode, it contains

the Metadatas that tells which blocks make up the original file



- Flume ingests data as generated by external systems, and puts it into the cluster
- HBase is a realtime database, built on top of HDFS
- Hue, graphical front end to the cluster
- Oozie, workflow ~~management~~ management Tool
- Mahout, machine learning library

Making all these talk to one another and work well can be tricky, to make installing and maintaining a cluster like this easier, Cloudera, has put ~~together~~ together a distribution ~~called~~ of Hadoop called CDH

It takes all the ecosystem project and Hadoop, and packages them together so that installation ~~is~~ is easy and they compatible together makes sense

HDFS

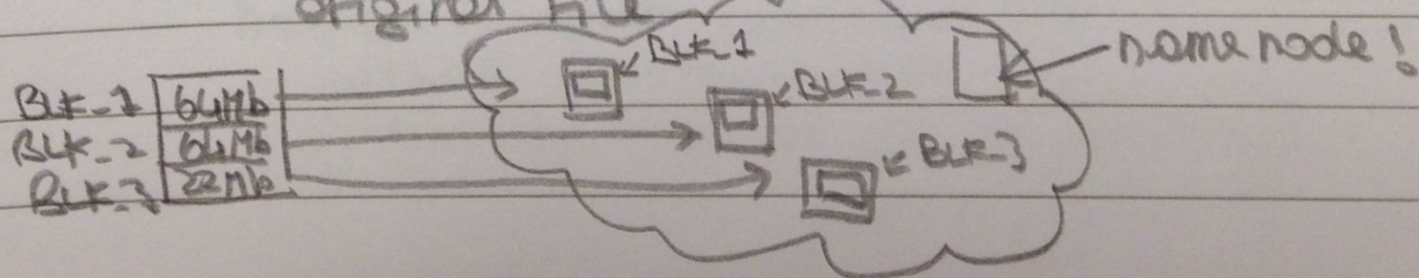
When a file is loaded into HDFS, ~~it~~ is splitted into chunks, which we call blocks.

Each block is ~~64~~ 64MB by default, and ~~has~~ has a unique name (BLK_ some large number).

As the file is uploaded to HDFS, each block will get stored on one node in the cluster.

A daemon is running on each of the machines in the cluster ~~each machine is called data node~~ ~~is~~ A separate machine called datanodes

runs a Daemon called NameNode, it contains the Metadata that says which blocks make up the original file



Data redundancy!

* If one node fails, the file has missing data.
For that reason, Hadoop replicates each block three times as it is stored in HDFS (picks three nodes at random and puts one copy of the block on each of the three)

~~The~~ If a node fails, the namenode sees that a block is under-replicated and ~~and~~ it will arrange to have it re-replicated so that there are three copies

~~If the namenode fails, it is~~

~~Namenode is configured to store the metadata not only on its own local disk but also somewhere on the network file system~~

If the namenode fails, the entire cluster is inaccessible, if the metadata is lost, the entire cluster's data is lost

To avoid that, most clusters have configured two namenodes.

- The active namenode works as before
- The standby can be configured to take over if the active fails

All of the commands that interact with the Hadoop File System, start with `hadoop fs`

`hadoop fs -ls` ← list of what's in home directory, in the Hadoop cluster

`hadoop fs -put purchases.txt`
upload purchases.txt file
destination filename can be specified

hadoop ~~fs~~ FS -tail purchases.txt (-cat for the entire content)
shows the last few lines of the file

other hadoop -FS commands closely mirror UNIX Filesystem commands

hadoop FS -mv purchases.txt purch.txt

hadoop ~~fs~~ FS -rm purch.txt

hadoop FS -mkdir myinput

hadoop FS -put purchases.txt myinput

=

Map Reduce

processing a large file from top to bottom
can take a long time

MapReduce is designed to process data
in parallel. ~~the~~ The file is broken into
chunks and processed in parallel

~~Each chunk is assigned to a mapper that
will at the same time look over a
different fraction of the data.~~

~~Reducers collect results that belong to the
same element but was computed in different
mappers.~~

= ~~Mapper: each deal with a small~~

Mappers are programs, each deal with a
small amount of data and work in parallel

~~Mapper~~ Their output is called the intermediate
records

Hadoop deals with all data in the form of key value
pairs. That's what intermediate records are

Shuffle and sort, shuffle is the movement of int records
from mappers to reducers, the ~~sort~~ sort is the part
that the reducers will organise the int. records.

Each reducer works on ~~one~~ one set of records at a time, it gets a key and a list of all the values that are processed in some way, then it writes out the final result

you can have the final result sorted by adding an extra step that merges the result

~~The MapReduce Framework~~

Daemons of MapReduce

↓
piece of code running all the time ~~on the~~ ~~machine~~

A daemon runs on each machine

When running a MapReduce job, you submit the job to the job tracker, which splits the work into mappers and reducers that will run on other cluster nodes.

Running map and reduce tasks is handled by a daemon called the task tracker, which runs on each datanode

Since the task tracker runs on the same machine as the datanodes, the Hadoop Framework will have the map tasks work directly on the pieces of data that ~~are~~ are stored on that machine

By default Hadoop will use a HDFS block as the input split for each mapper ~~and~~ and it'll make sure that a mapper works on data on the same machine.

If a block needs processing, the task tracker on the machine where the block is will likely be the one chosen to process it. But that's not always possible as TaskTrackers could be busy

... in which case a different node will be chosen to process the block and it will be streamed over the network

=
to submit the job

hadoop jar /usr/lib/hadoop-~~hadoop~~.../streaming/...

-mapper m.scale -reducer r.scale ~~File~~

-File m.scale -File r.scale -input dir/in

output dir/out
 \nwarrow the output directories must not already exist

=
hadoop produces _SUCCESS, _log5, part-0000

I can get the result by doing \downarrow or more ~~or~~ depending on
hadoop fs -get part#0000 localFile.txt reducers

A simplified command for the one above is

hadoop m.scale r.scale dir/in dir/out

=
~~the processing~~