



Big Data Training 2015

Apache Hadoop Fundamentals and Architecture

April 2015



Gonzalo Zarza

Big Data Architect
PhD in High-Performance Computing

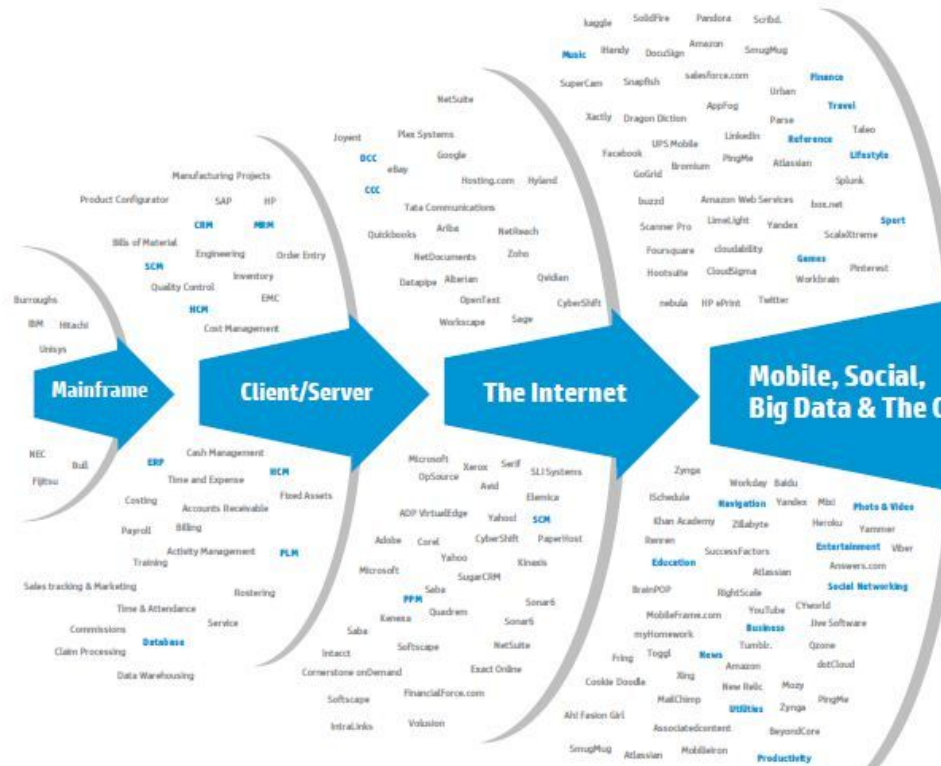
***Note:** this lecture has been prepared by Gonzalo Zarza, José Muguerza and Juan Gaviria*

Hadoop

We create innovative software products that
appeal to global audiences.



A new style of IT emerging



Every 60 seconds



98,000+ tweets



695,000 status updates



11 million instant messages



698,445 Google searches



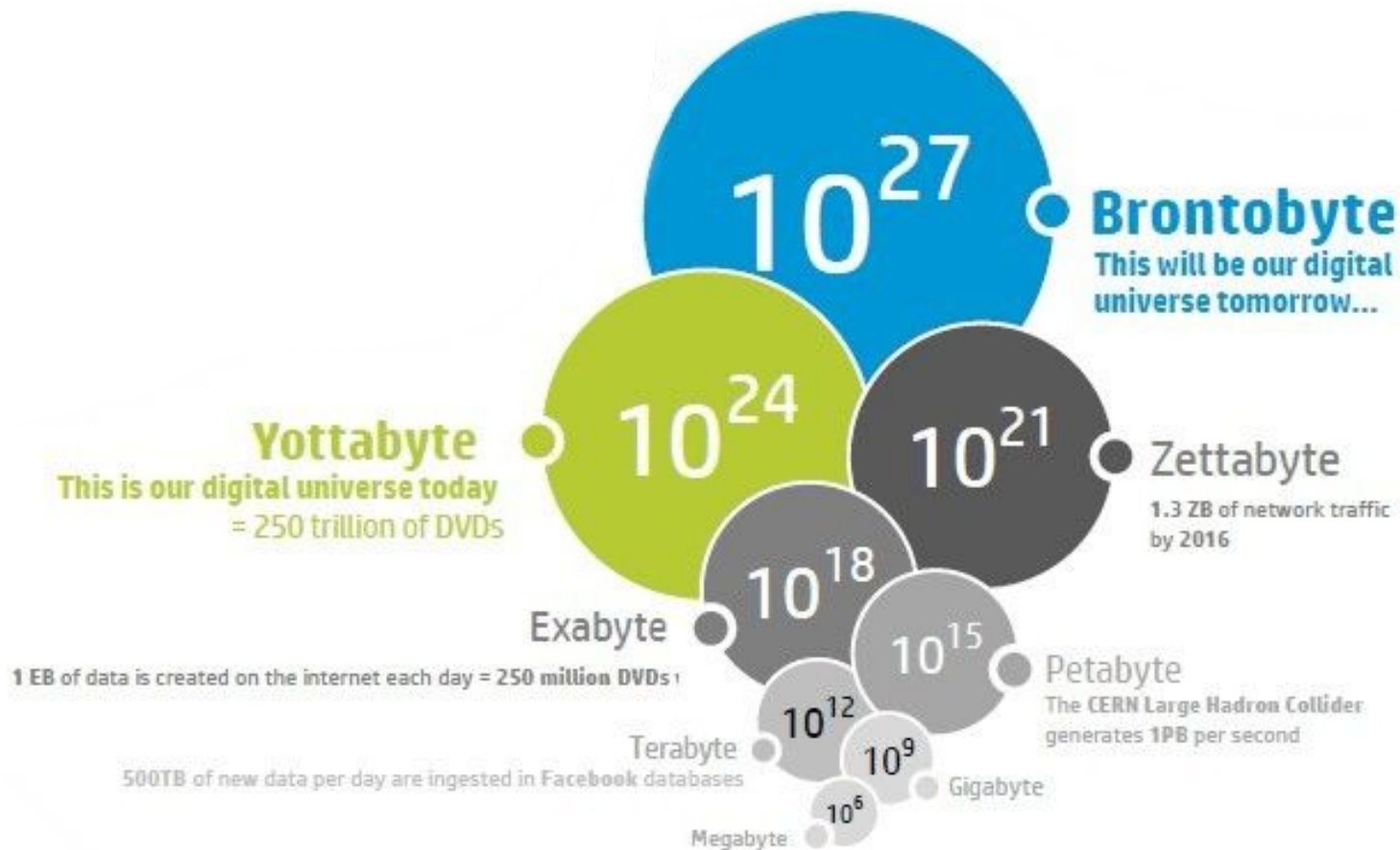
168 million+ emails sent



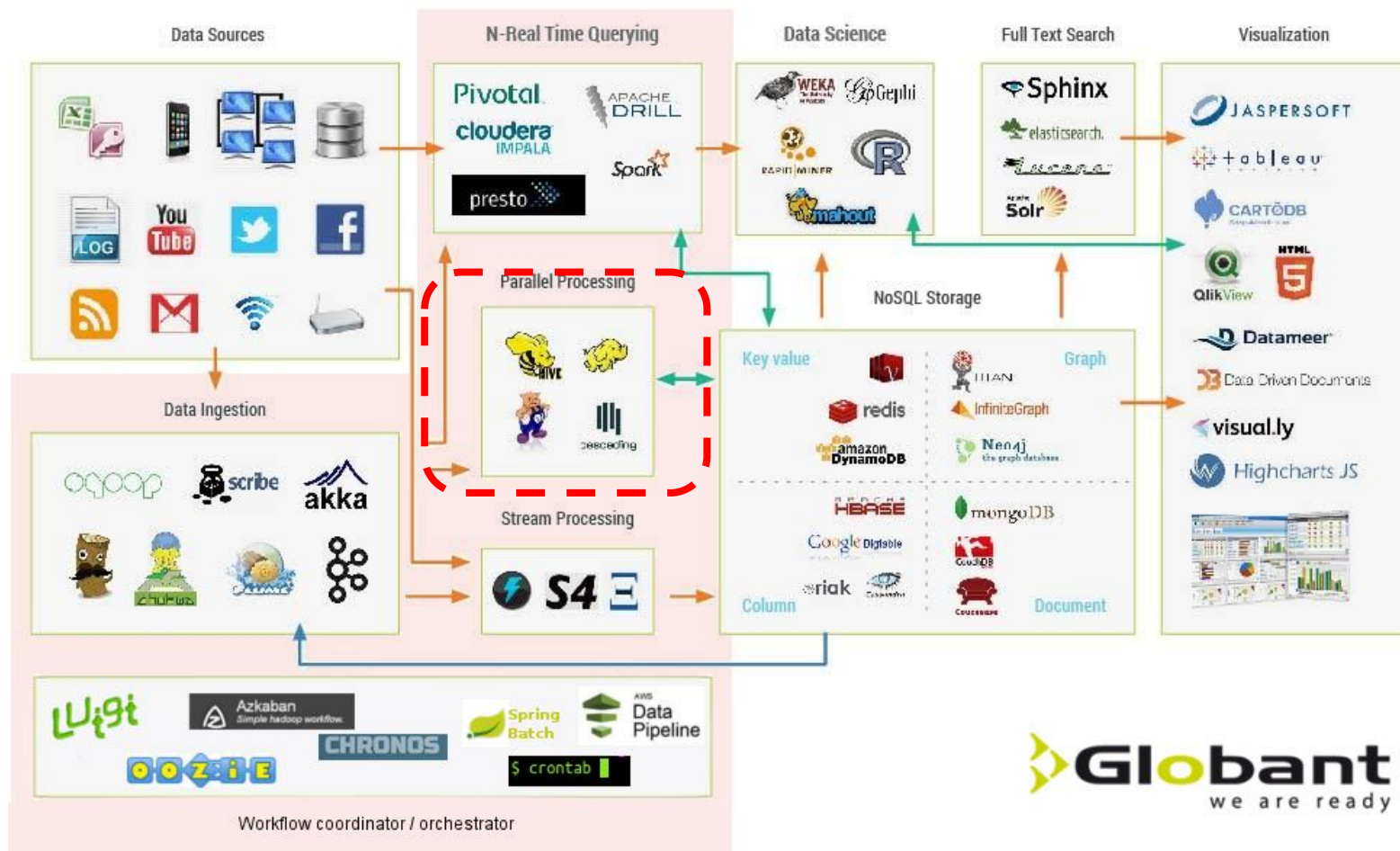
1,820TB of data created



217 new mobile web users



BIG DATA ECOSYSTEM



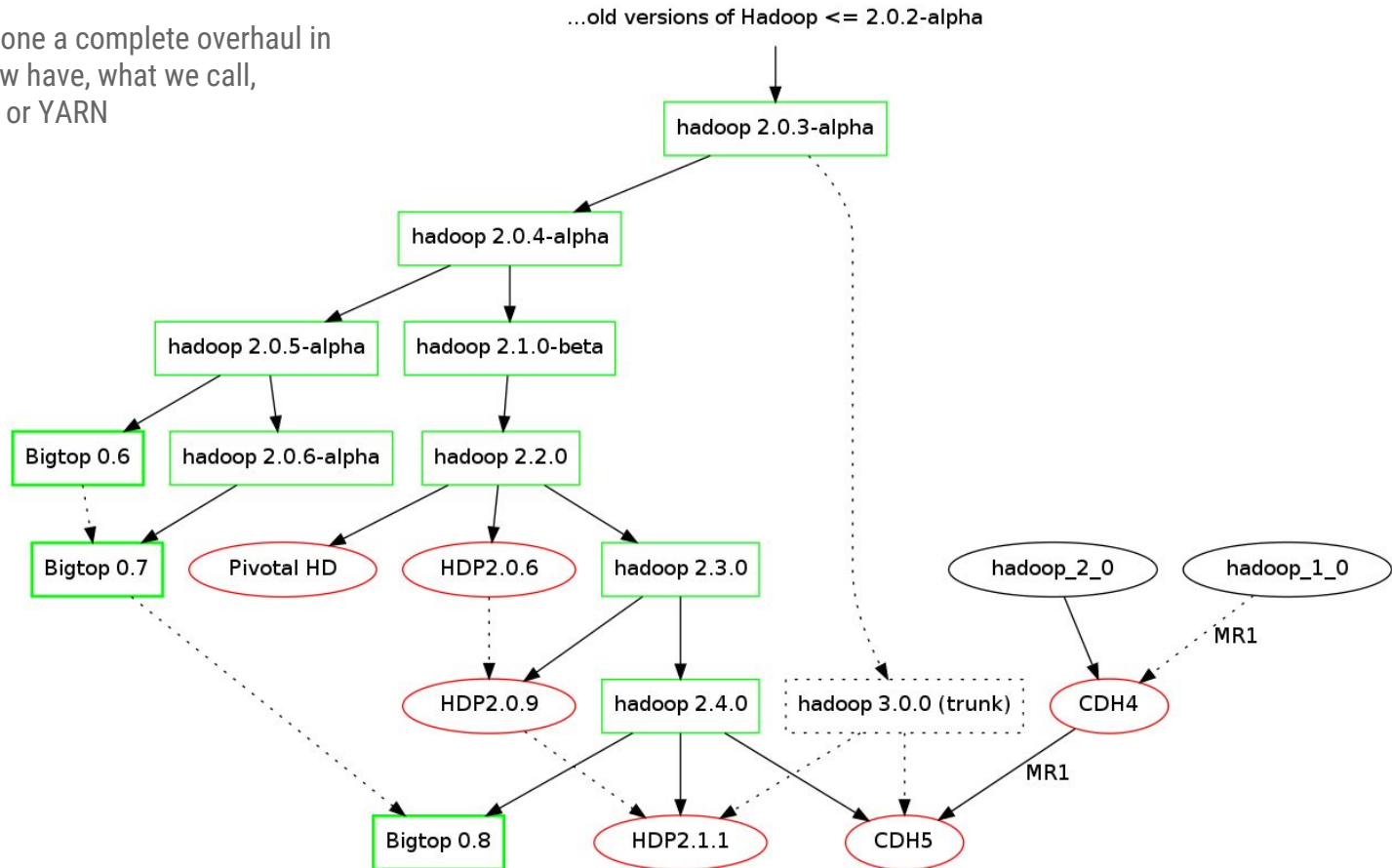


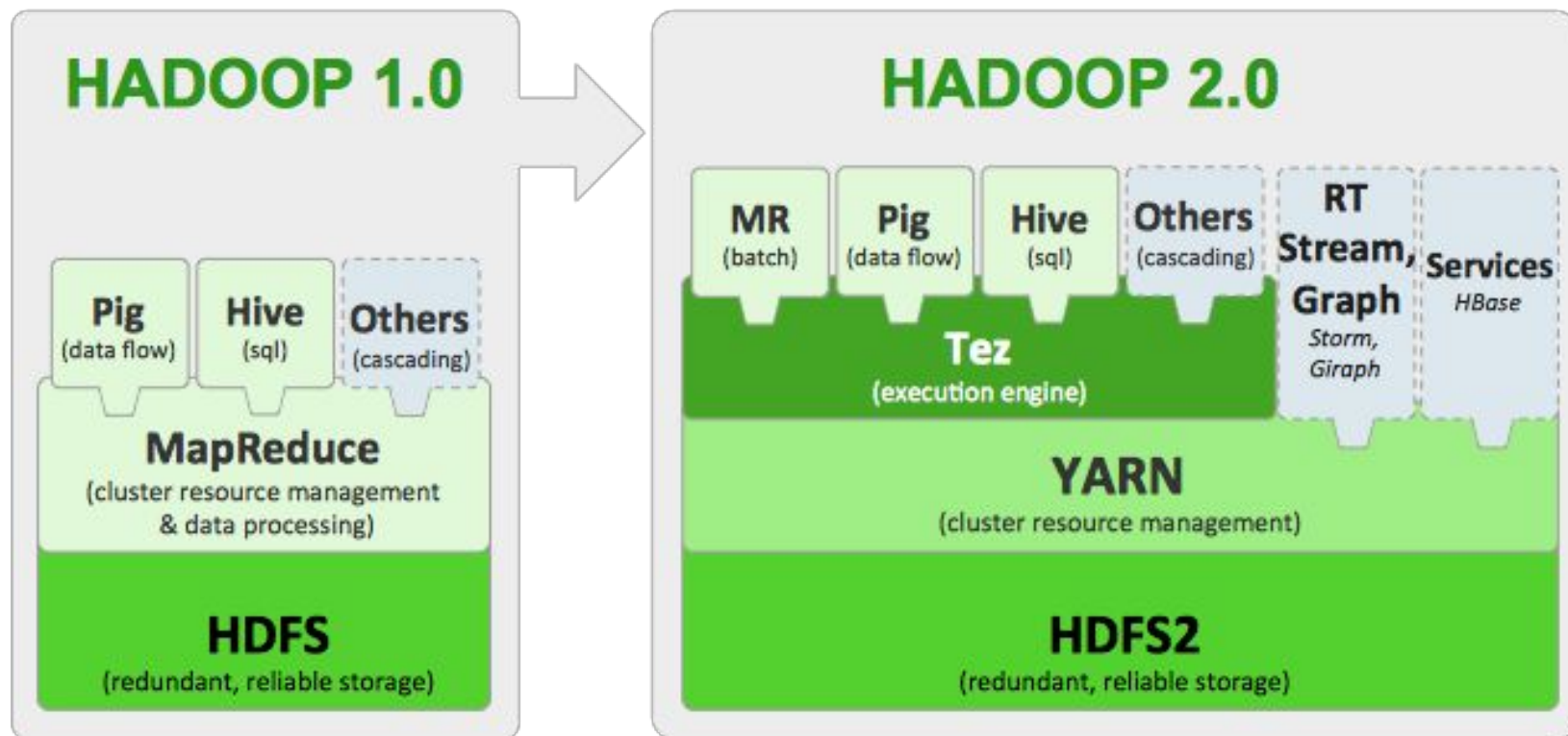
Description	Apache Hadoop es un framework de software que soporta aplicaciones distribuidas bajo una licencia libre . ¹ Permite a las aplicaciones trabajar con miles de nodos y petabytes de datos. Hadoop se inspiró en los papers de Google para MapReduce (2004) y Google File System (GFS) (2003) .
Last version	2.6.0 (Nov , 2014)
Maturity	Estable
Supported languages	Principalmente Java (también Python). APIs REST
Combined with other frameworks	Muchos, gran variedad de interacción....
Websites	https://hadoop.apache.org/ - Apache Hadoop YARN Book
Enterprise Support	Hortonworks , Cloudera , MapR , WanDisco , Intel, Altiscale
License	Apache License 2.0





MapReduce has undergone a complete overhaul in hadoop-0.23 and we now have, what we call, MapReduce 2.0 (MRv2) or YARN





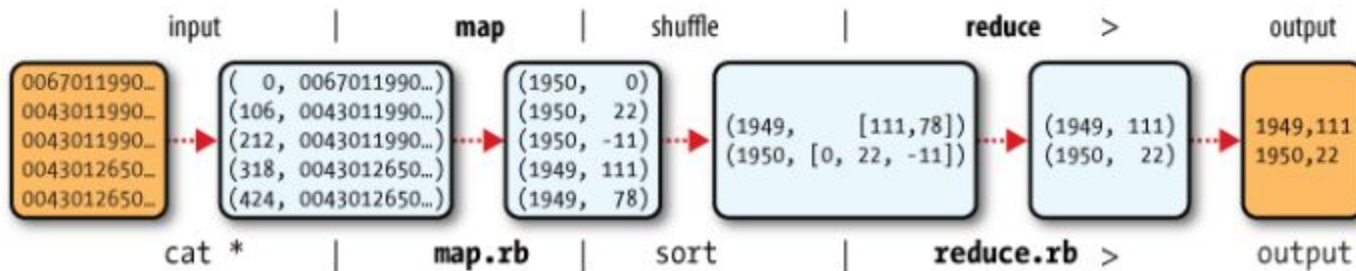
Hadoop MapReduce

We create innovative software products that
appeal to global audiences.



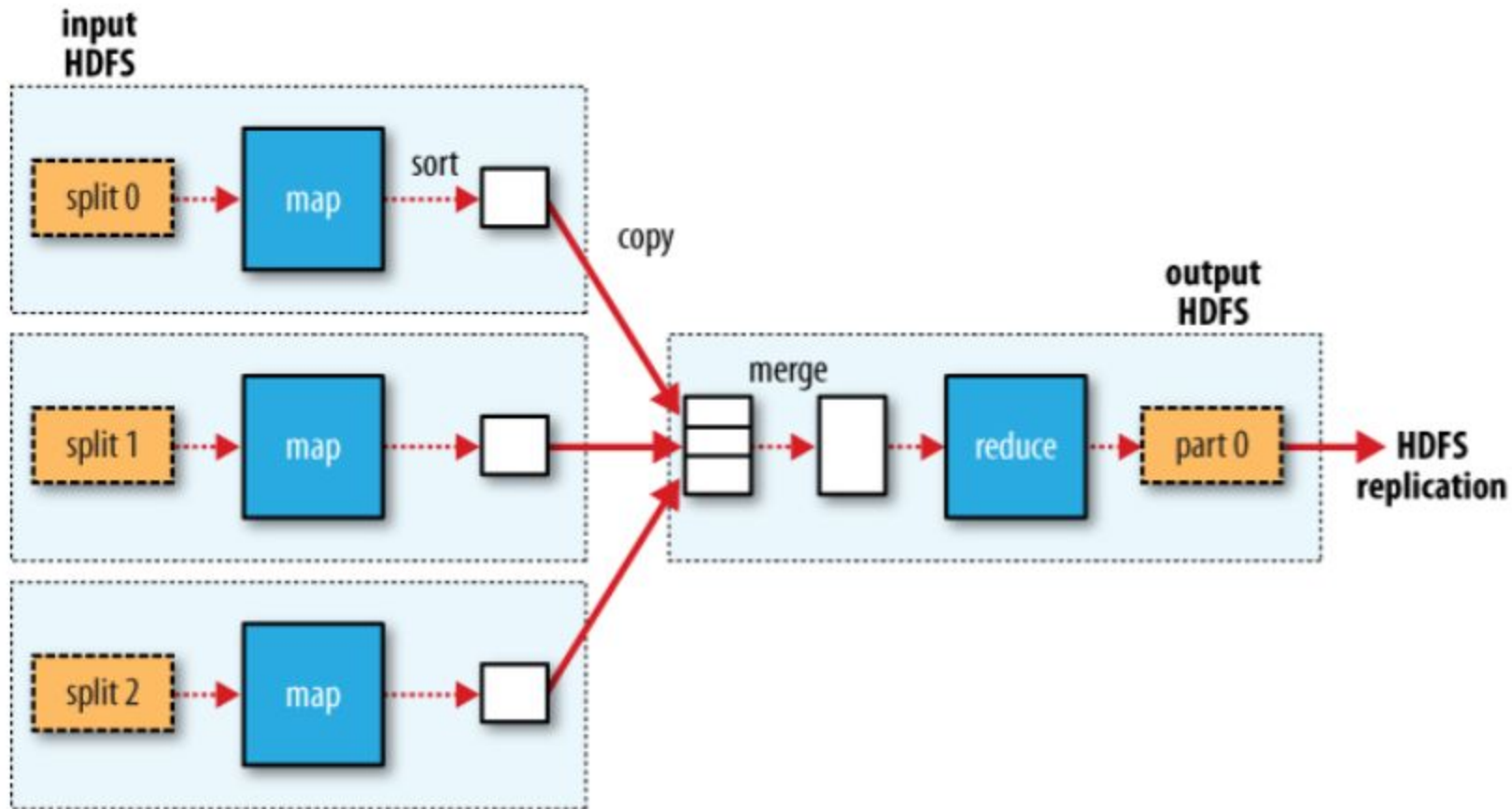
MapReduce is the heart of Hadoop. It is this programming paradigm that allows for massive scalability across hundreds or thousands of servers in a Hadoop cluster.

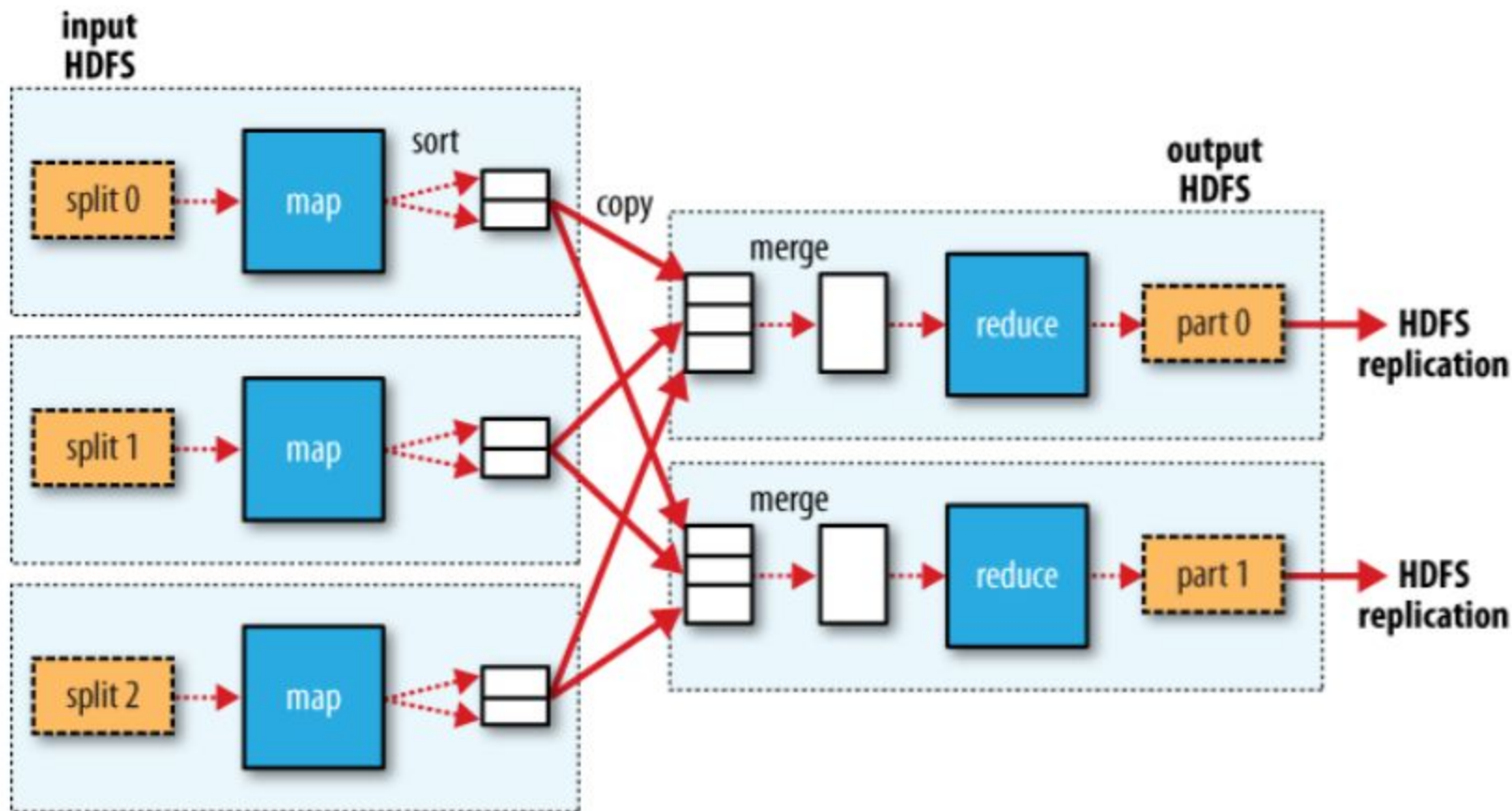
MapReduce works by breaking the processing into **two phases**: the **map** phase and the **reduce** phase. Each phase has **key-value pairs** as input and output, the types of which may be chosen by the programmer. The programmer also specifies two functions: the map function and the reduce function.

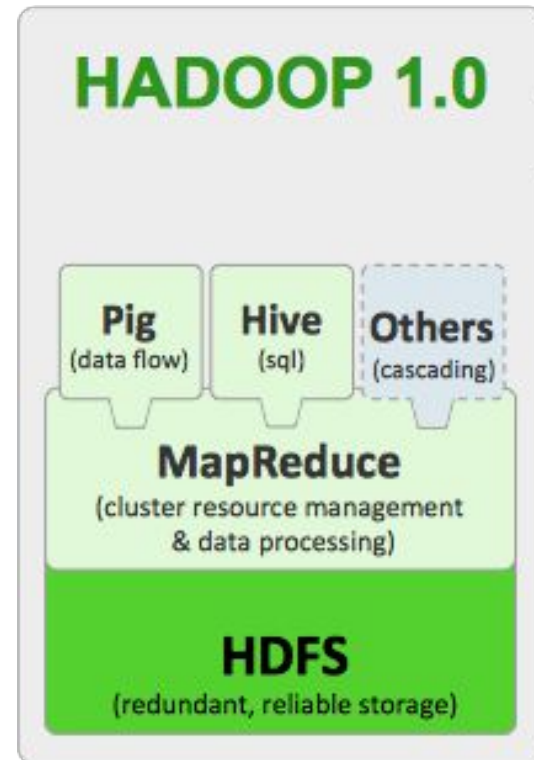
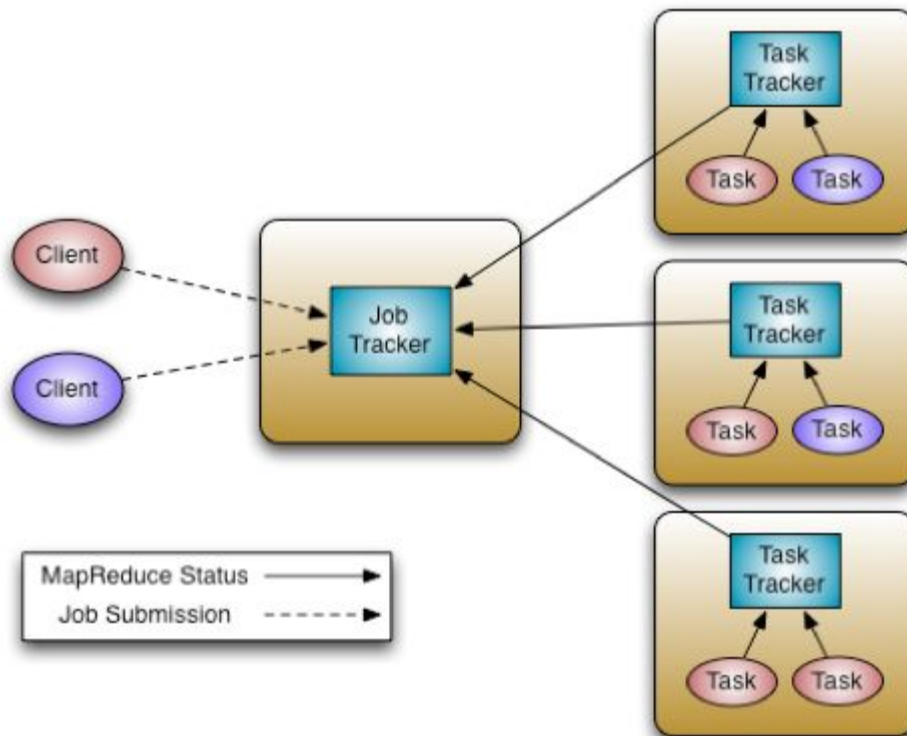


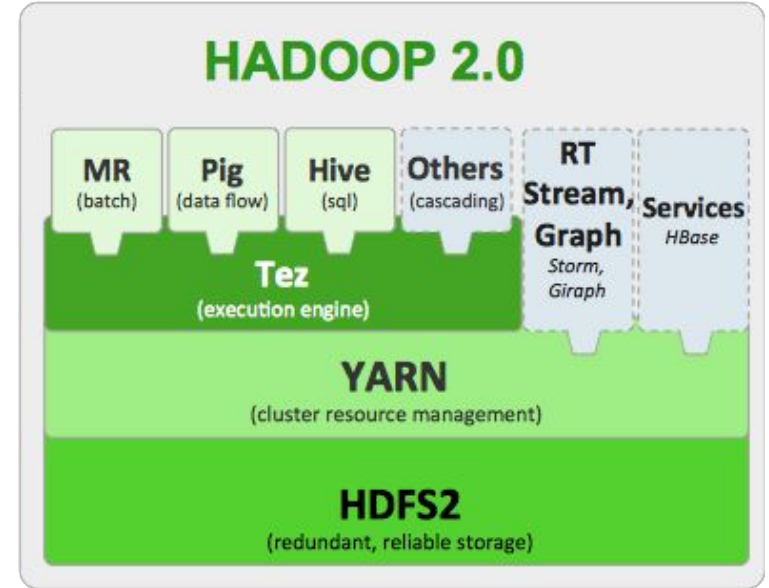
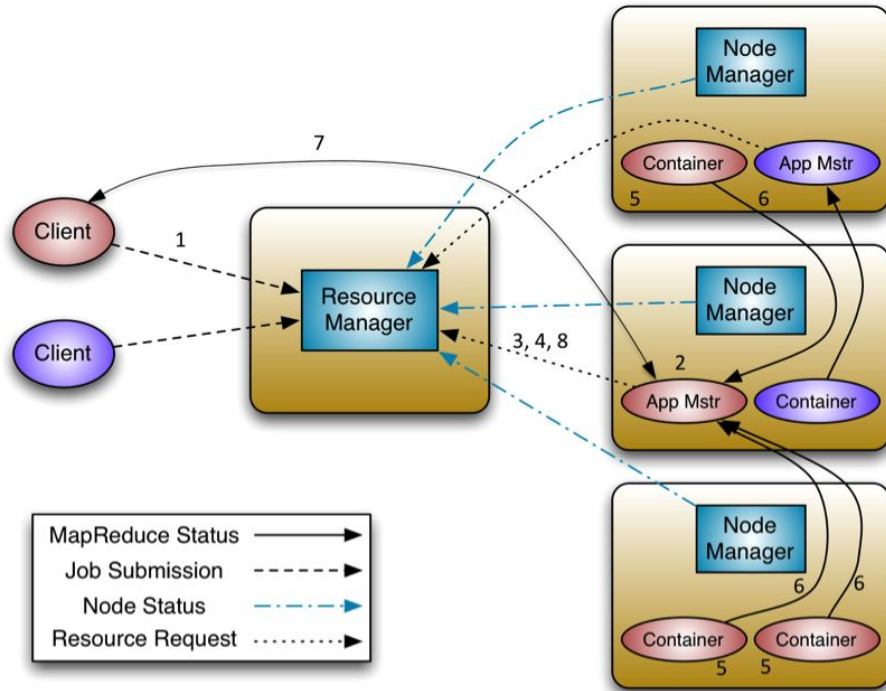
Note: A MapReduce job is a unit of work that the client wants to be performed; it consists of the input data, the MapReduce program, and configuration information. Hadoop runs the job by dividing it into tasks, of which there are two types: map tasks and reduce tasks.

Hadoop divides the input to a MapReduce job into fixed-size pieces called input splits, or just splits. Hadoop creates one map task for each split, which runs the user-defined map function for each record in the split.









Applications Run Natively IN Hadoop

Pig

Script

Hive

SQL

HBase

NoSQL

Accumulo

NoSQL

Storm

Stream

Solr

Search

Spark

In-Memory

Cascading

Java

OthersISV
Engines

YARN: Data Operating System

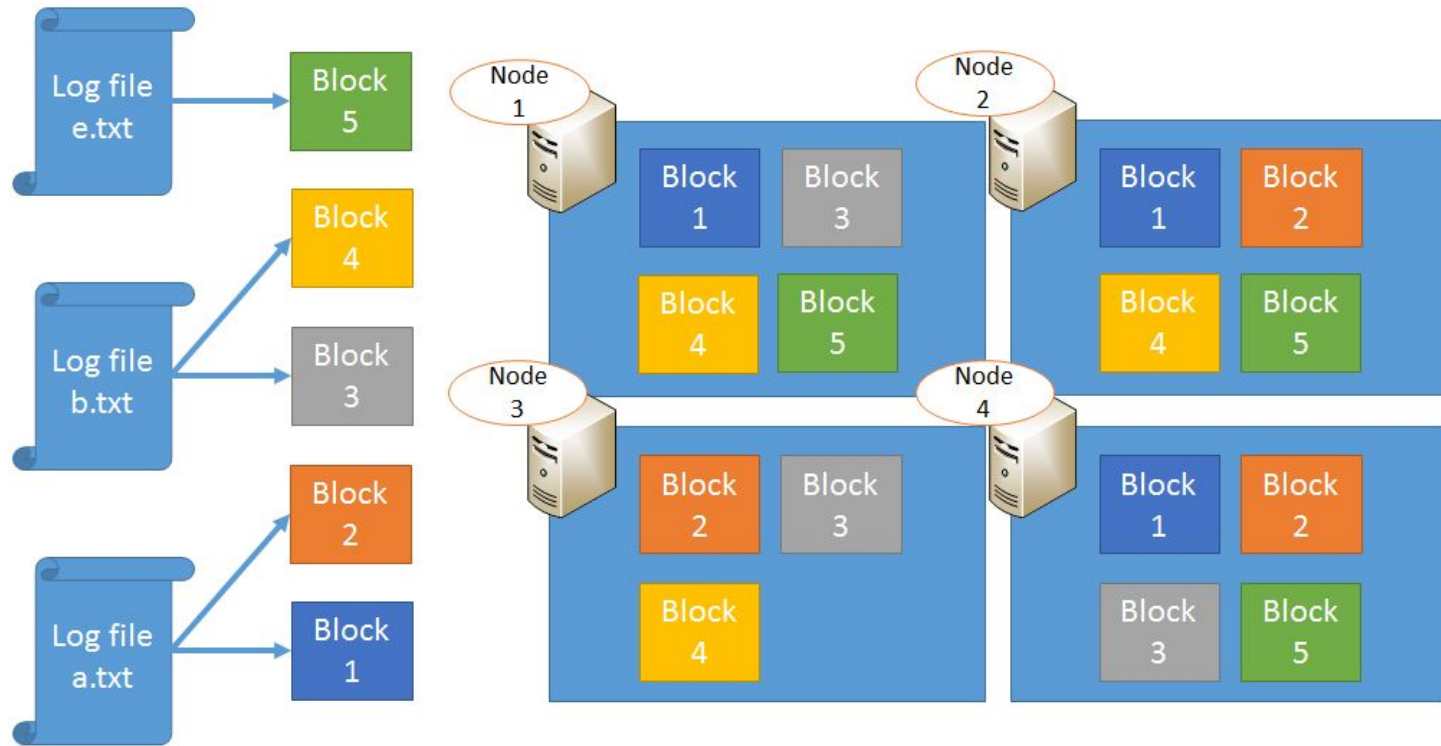
HDFS

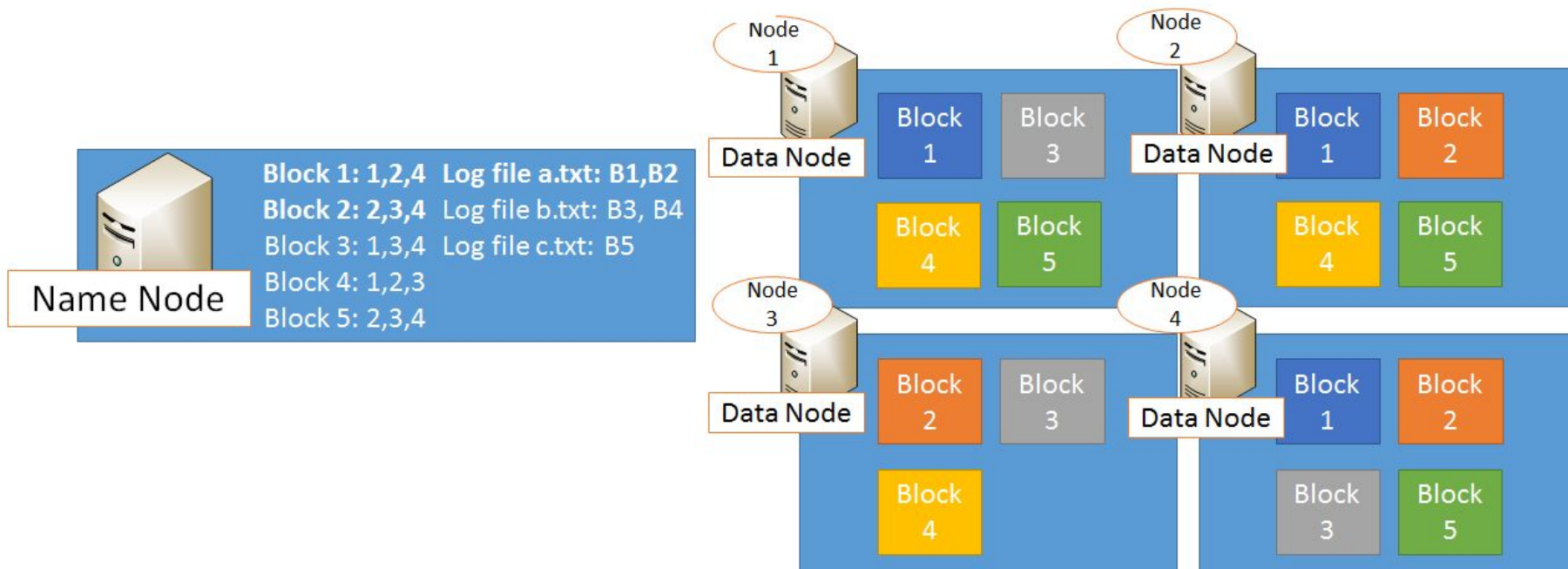
(Hadoop Distributed File System)

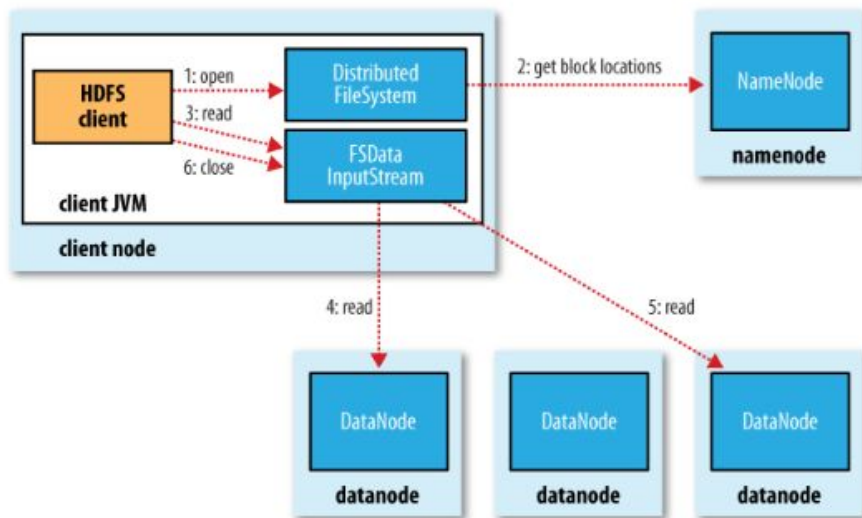
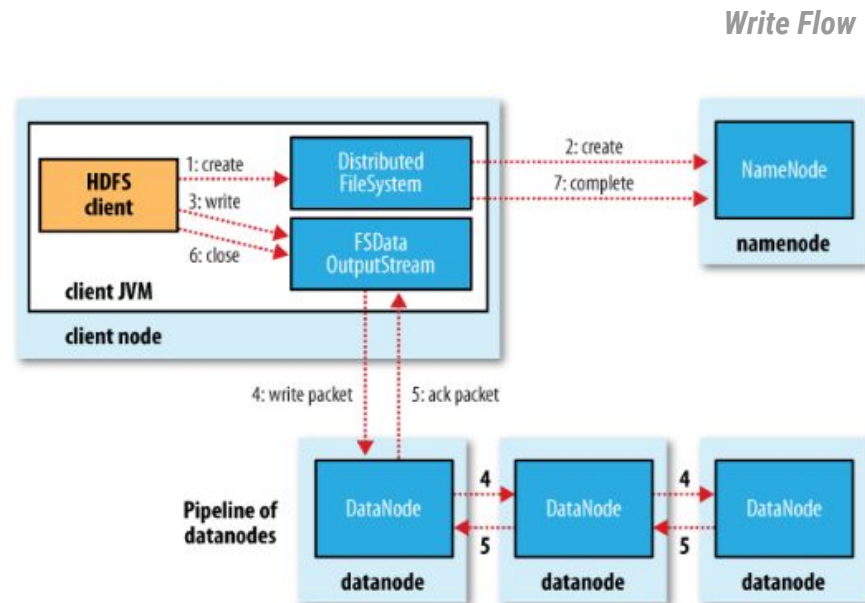
Hadoop HDFS

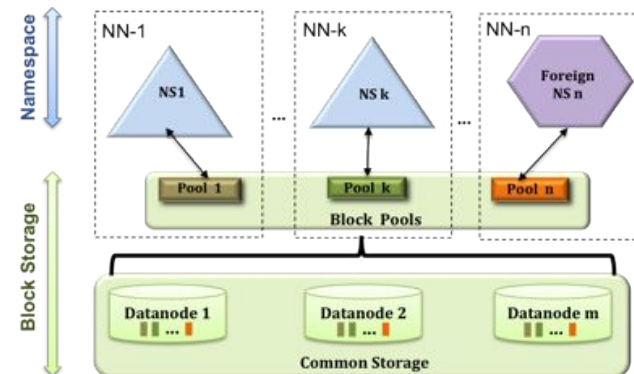
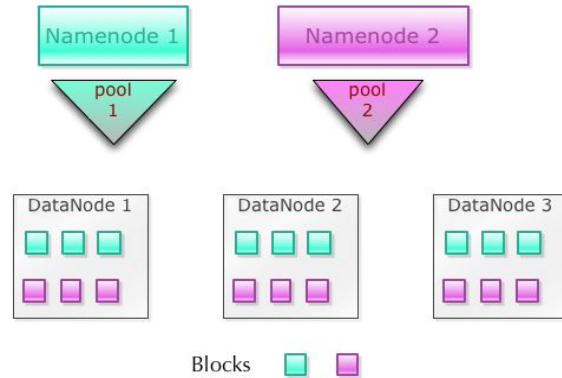
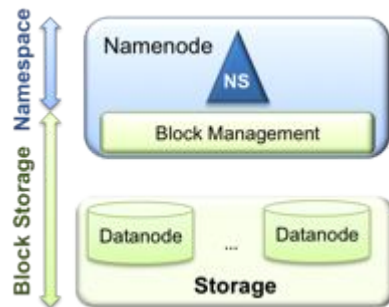
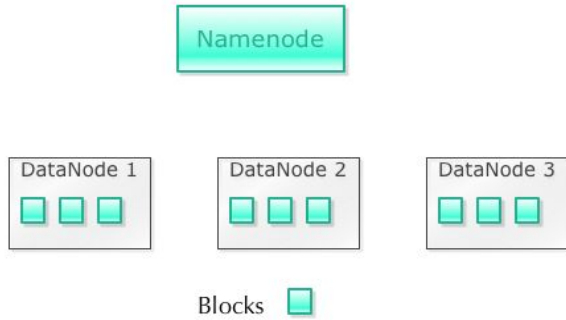
We create innovative software products that
appeal to global audiences.



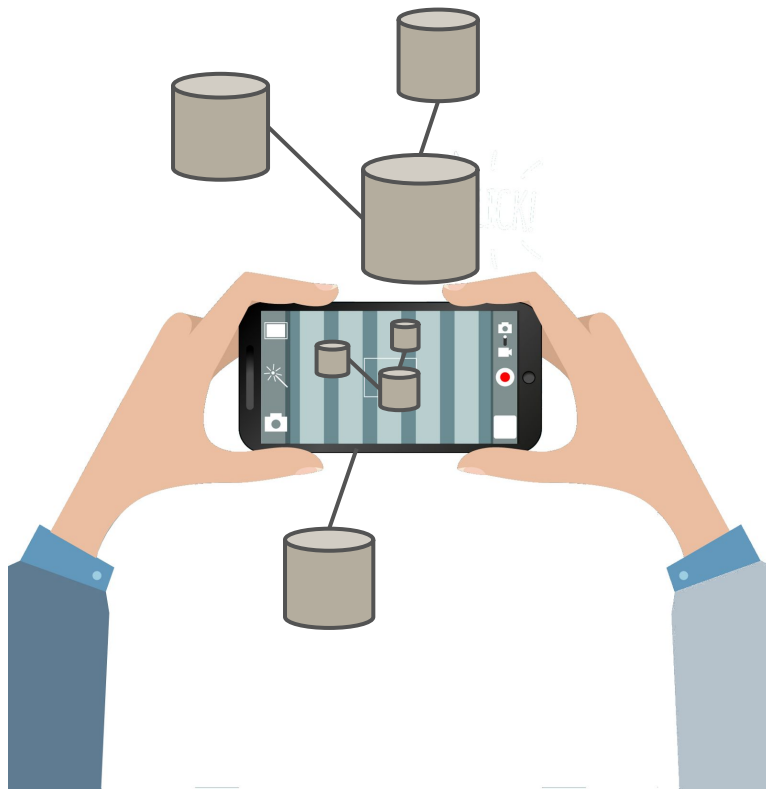




*Read Flow**Write Flow*



*HDFS Snapshots are read-only **point-in-time** copies of the file system.*





1. Low-latency Data Access

Applications that require real-time query, and low-latency access to data in tens of milliseconds will not work well with Hadoop.

Hadoop is not a substitute for a database. Database index records that will gains low-latency and fast response.

But if you really want to replace the database for real time needs, try HBase, which is a column-oriented database for random and real time read/write.

2. Structured Data

Hadoop is not fit for structured data with strong relationship. Hadoop works well for semi-structured and unstructured data. It stores data in files, doesn't index them like RDBMS. Therefore, each ad hoc query for Hadoop is processed by MapReduce job which will bring the latency cost.

3. When data isn't that big

How big the data is big enough for Hadoop? The answer is TB or PB. When your analytics data is only tens of GB, Hadoop is heavy. Don't follow the fashion and use Hadoop, just follow your requirements.

4. Too many small files

When there are too many small files, the NameNode will hit its memory limit where the block map and the metadata are hosted. And to handle the NameNode bottleneck, Hadoop introduces HDFS Federation.

5. MapReduce may not be the best choice

MapReduce is a simple programming model in parallel. But for MapReduce parallelism, you need to make sure each MR job and the data where the job runs on is independent from all the others. Every MR shouldn't have dependencies.

But if you want to do some data sharing during MR, you can do like this:

- Iteration: run multiple MR jobs, with the output of one being the input of the next MR.
- Shared state information. But don't share information in memory, since each MR job is run on single JVM.

References

We create innovative software products that
appeal to global audiences.





- **Hadoop: The Definitive Guide, 3rd Edition**. O'Reilly Media / Yahoo Press - [Online @ Globant's Big Data Training Site](#)
- **Guía de Instalación de Hadoop**. Globant Team - [Online @ Globant's Big Data Training Site](#)
- *To be continued...*

Questions?

We create innovative software products that
appeal to global audiences.



Thank You!