# Big Data. Scalability and models

# The value of data -- the need for computing



In science since long ago -- generalized in last years

# Scalability



How is it achieved? At what cost? What do we gain?

What hardware architectures?  What programming models?

# Scaling computing



Infraestructura de cómputo

HPC

Algoritmos/Aplicaciones

Datos

Big Data

Recurso humano

# map reduce

```
2012-01-01 09:08 BOG Libros 88.56 Discover
2012-01-01 09:09 BGA Libros 337.71 Efectivo
2012-01-01 09:52 BGA Libros 62.41 Discover
2012-01-01 10:08 MED Musica 93.37 Visa
2012-01-01 10:22 BGA Musica 369.94 MasterCard
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover
```

# map reduce

```
2012-01-01 09:08 BOG Libros 88.56 Discover
2012-01-01 09:09 BGA Libros 337.71 Efectivo
2012-01-01 09:52 BGA Libros 62.41 Discover
2012-01-01 10:08 MED Musica 93.37 Visa
2012-01-01 10:22 BGA Musica 369.94 MasterCard
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover
```

1

```
BOG 88.56
BGA 337.71
BGA 62.41
MED 93.37
BGA 369.94
MED 119.12
BOG 296.76
```

# map reduce

```
2012-01-01 09:08 BOG Libros 88.56 Discover
2012-01-01 09:09 BGA Libros 337.71 Efectivo
2012-01-01 09:52 BGA Libros 62.41 Discover
2012-01-01 10:08 MED Musica 93.37 Visa
2012-01-01 10:22 BGA Musica 369.94 MasterCard
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover
```

## 1

```
BOG 88.56
BGA 337.71
BGA 62.41
MED 93.37
BGA 369.94
MED 119.12
BOG 296.76
```

## 2

```
BOG 88.56 296.76

BGA 337.71 62.41 369.94

MED 93.37 119.12
```

# map reduce

```
2012-01-01 09:08 BOG Libros 88.56 Discover
2012-01-01 09:09 BGA Libros 337.71 Efectivo
2012-01-01 09:52 BGA Libros 62.41 Discover
2012-01-01 10:08 MED Musica 93.37 Visa
2012-01-01 10:22 BGA Musica 369.94 MasterCard
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover
```

**1**

```
BOG 88.56
BGA 337.71
BGA 62.41
MED 93.37
BGA 369.94
MED 119.12
BOG 296.76
```

**2**

```
BOG 88.56 296.76

BGA 337.71 62.41 369.94

MED 93.37 119.12
```

**3**

```
BOG 385.32

BGA 770.06

MED 212.49
```

# map reduce

```
2012-01-01 09:08 BOG Libros 88.56 Discover
2012-01-01 09:09 BGA Libros 337.71 Efectivo
2012-01-01 09:52 BGA Libros 62.41 Discover
2012-01-01 10:08 MED Musica 93.37 Visa
2012-01-01 10:22 BGA Musica 369.94 MasterCard
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover
```

## 1 map            2 shuffle            3 reduce

```
BOG 88.56
BGA 337.71
BGA 62.41          BOG 88.56 296.76              BOG 385.32
MED 93.37
BGA 369.94         BGA 337.71 62.41 369.94       BGA 770.06
MED 119.12
BOG 296.76         MED 93.37 119.12              MED 212.49
```

map (k, v)                              reduce (k, [$v_1$, ...])

# map reduce

```
2012-01-01 09:08 BOG Libros 88.56 Discover      host A
2012-01-01 09:09 BGA Libros 337.71 Efectivo
2012-01-01 09:52 BGA Libros 62.41 Discover
2012-01-01 10:08 MED Musica 93.37 Visa          host B
2012-01-01 10:22 BGA Musica 369.94 MasterCard
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover     host C
```

## 1 map            ## 2 shuffle            ## 3 reduce

```
BOG 88.56
BGA 337.71
BGA 62.41         BOG 88.56 296.76           BOG 385.32
MED 93.37
BGA 369.94        BGA 337.71 62.41 369.94    BGA 770.06
MED 119.12
BOG 296.76        MED 93.37 119.12           MED 212.49
```

map (k, v)                              reduce (k, [$v_1$, ...])

# map reduce

```
2012-01-01 09:08 BOG Libros  88.56 Discover
2012-01-01 09:09 BGA Libros 337.71 Efectivo
```
host A
```
2012-01-01 09:52 BGA Libros  62.41 Discover
2012-01-01 10:08 MED Musica  93.37 Visa
2012-01-01 10:22 BGA Musica 369.94 MasterCard
```
host B
```
2012-01-01 10:58 MED Musica 119.12 Efectivo
2012-01-01 11:36 BOG Musica 296.76 Discover
```
host C

## 1 map

```
BOG  88.56
BGA  337
BGA  6
MED
    9.94
    119.12
BOG 296.76
```

programmer

map (k, v)

## 2 shuffle

```
BOG  88.56
BGA           62.41 369.94
MED   3.37 119.12
```

framework

## 3 reduce

```
BOG  38
    .06
    212.49
```

programmer
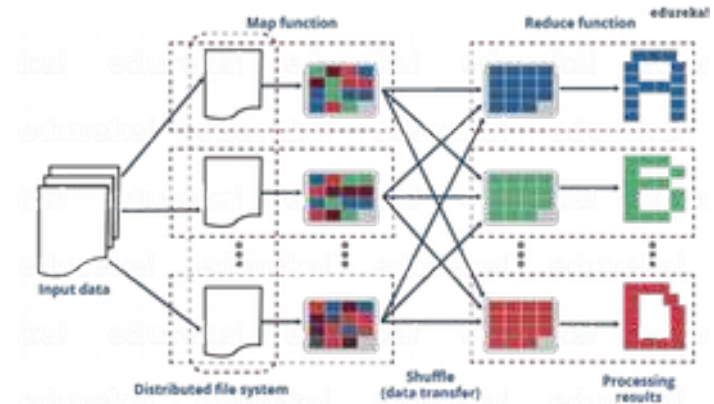
reduce (k, [$v_1$, …])

# map reduce



DATA TRANSFER ONLY IN SHUFFLE
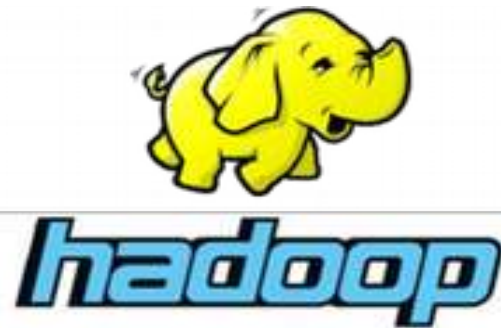
# map reduce



Data ALREADY exists in nodes

MR Programmer → forget about paralellism

Framework developer → optimize coordination and comms

RESTRICTED PROGRAMMING MODEL

# How to crunch 1PB

Lots of disk spinning all time
Redundancy, disks die
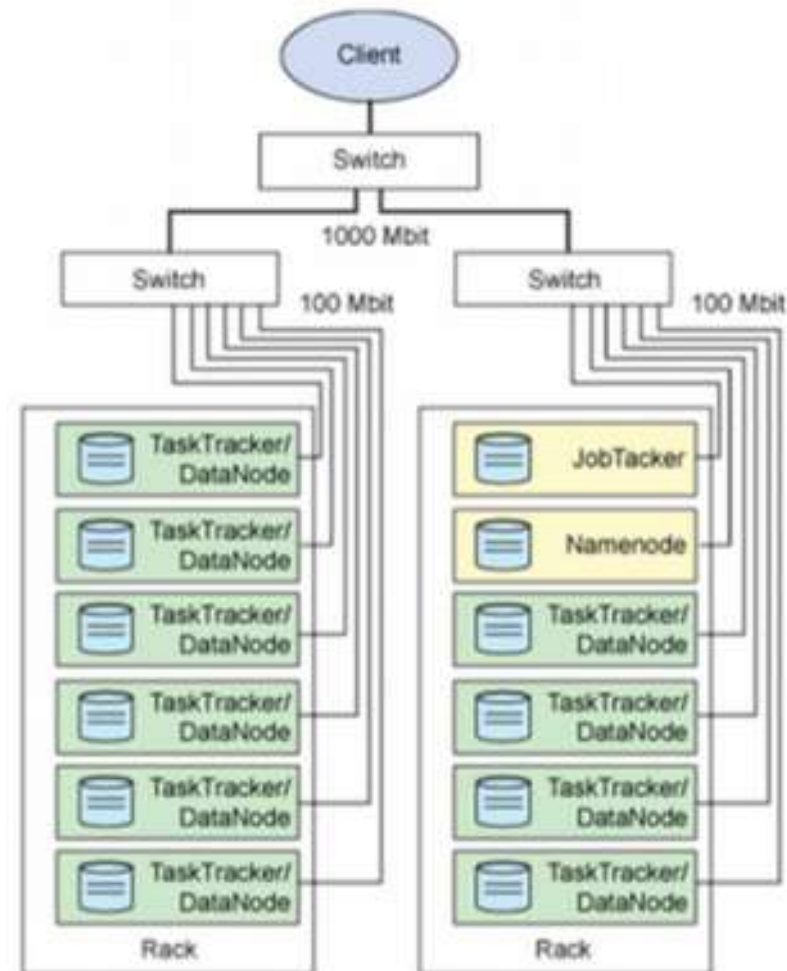Lots of CPUs, working all time
Retry, since errors happen

# Design qualities

Scalable – many servers
Reliable – redundant storage
Fault-tolerant – auto retry, self-healing

# Computing to Data

data goes to computing

computing goes to data

# noSQL

Expressivity SQL vs. Scalability



CA Category
RDBMS

CP Category
BigTable
HBase
MongoDB
Redis

Consistency

CA    CP

Availability    AP    Partition Tolerance

AP Category
Dynamo
Voldemort
Cassandra
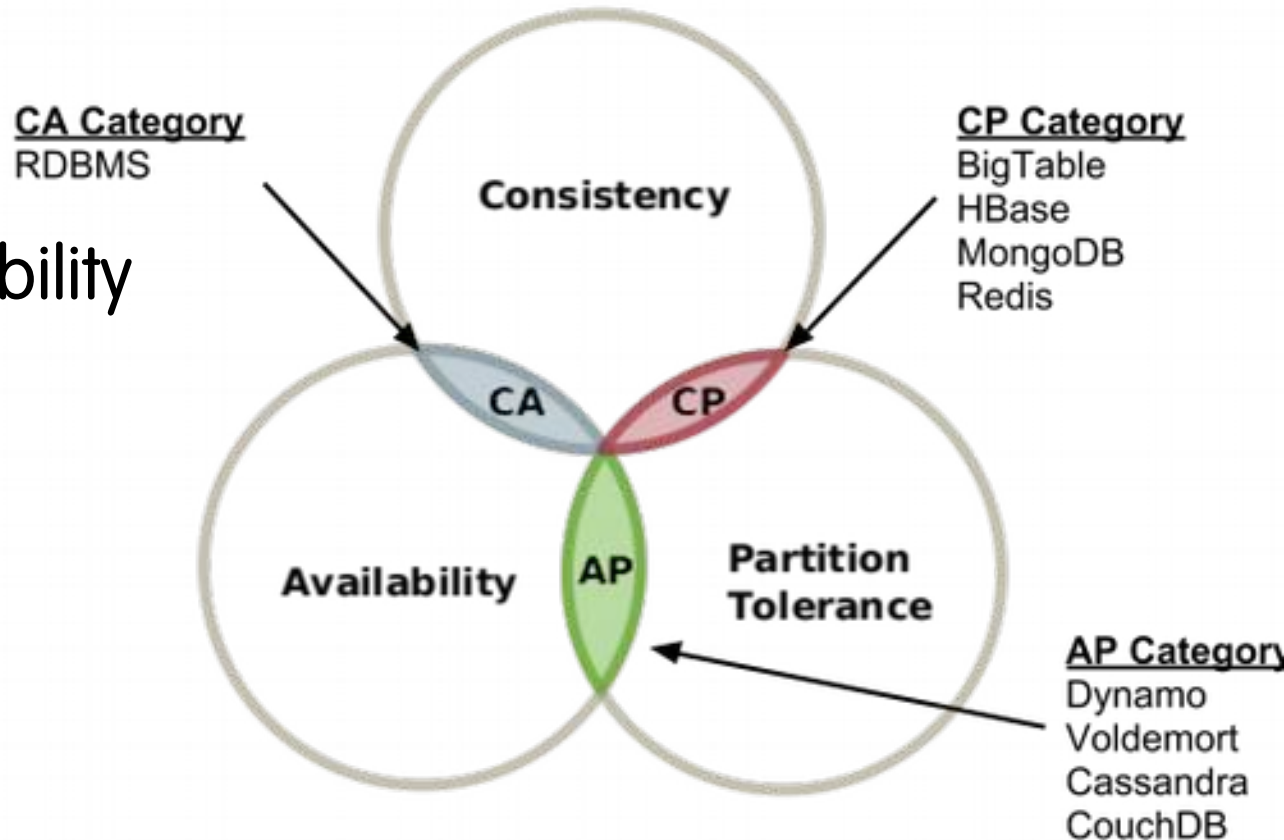CouchDB

Simpler data model (key, values)

Simpler operations

    Scan/access per key, basic transactions (check&put)

    No joins, no SQL language

Simple failover and scale up

Big table, Hbase, DynamoDB, Azure, Cassandra, etc.

# Why Big Data

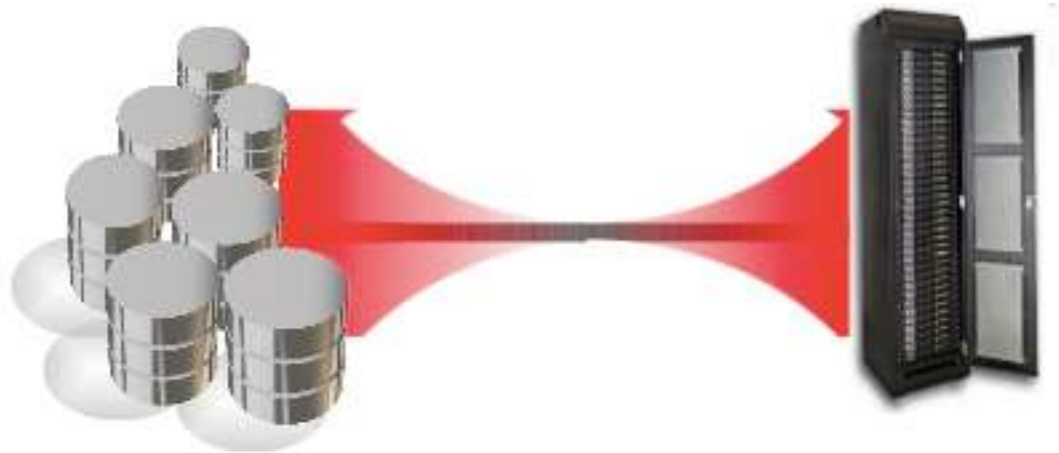Data growing faster than computation speeds

Storage and network bottlenecks

Facebook daily logs: 60TB

1000 genomes project: 200TB

Cost of 1TB of disk: USD 30
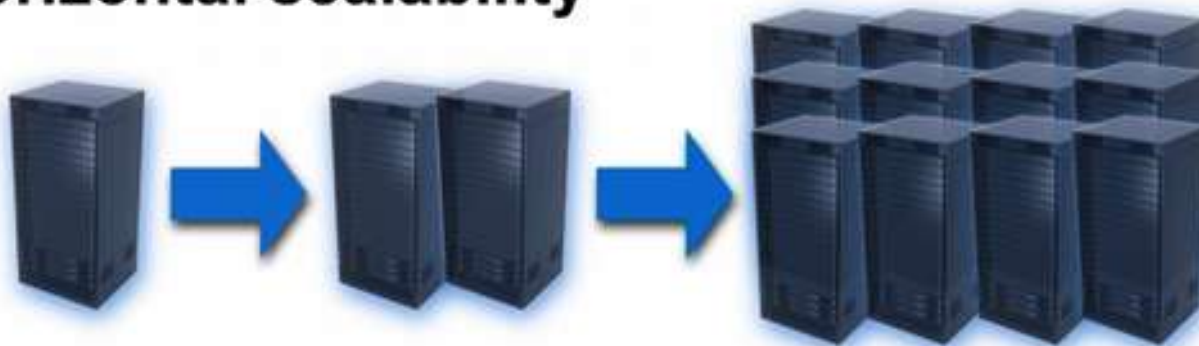
Time to read 1TB from disk: 3 hrs (100 MB/s)

# Scalability in Big Data

**Vertical "scalability"**



scale up
traditional DBs
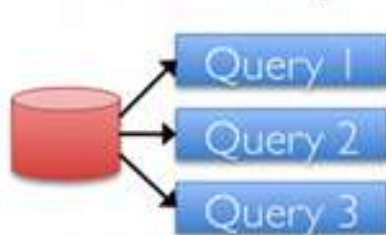
**Horizontal scalability**



scale out
noSQL

seek "triviality" for appropriate
sw+hw architectures

recent technologies (virtualization, etc.) tend to favor the cost of scaling out
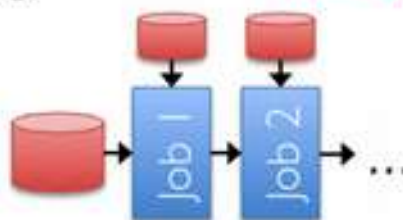
**Apache Spark Motivation**

- Using Map Reduce for complex jobs, interactive queries and online processing involves *lots of disk I/O*
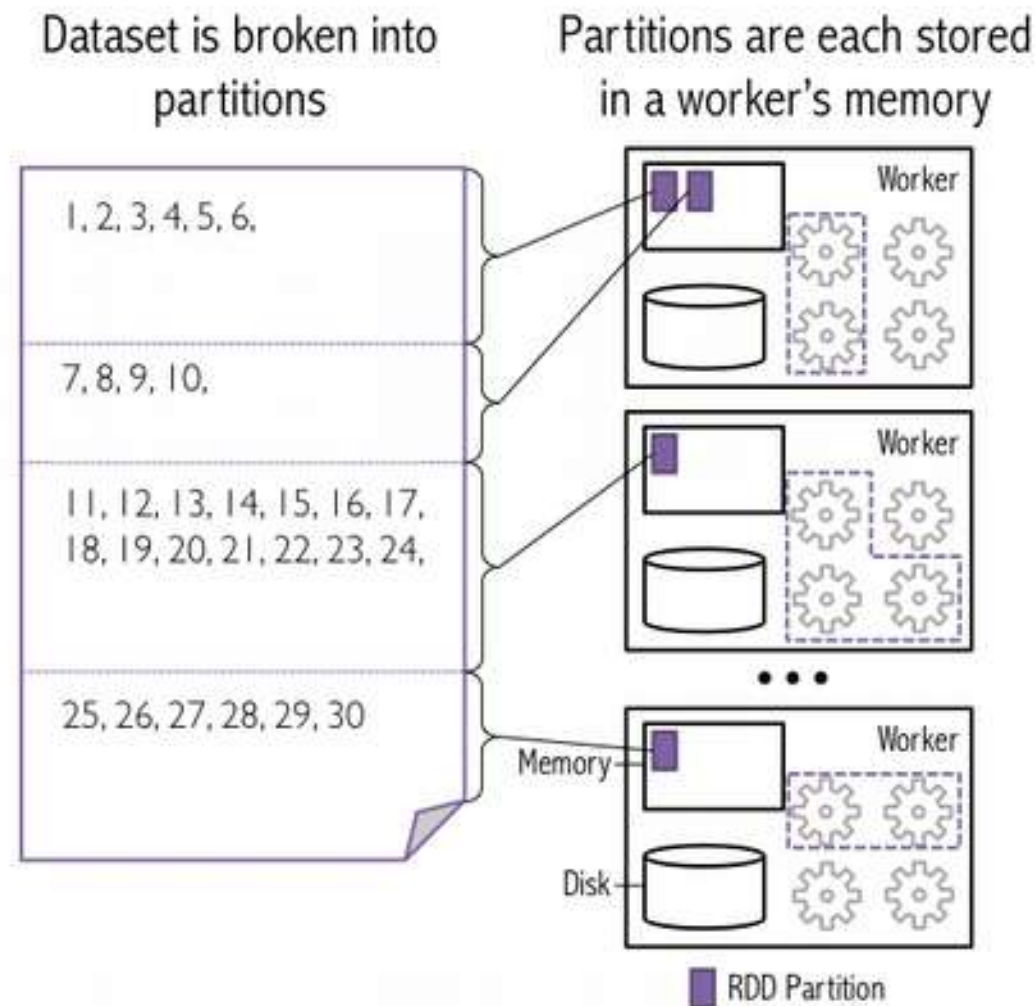
Interactive mining      Stream processing

Also, iterative jobs

Disk I/O is very slow

# Spark computing model

Resilient Distributed Datasets … **ON MEMORY**



Dataset is broken into partitions

Partitions are each stored in a worker's memory

| 1, 2, 3, 4, 5, 6, |
| 7, 8, 9, 10, |
| 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, |
| 25, 26, 27, 28, 29, 30 |

Worker

Worker

Memory

Disk

Worker

RDD Partition

# Spark computing model

A program is a set of transformations on RDDs
**(to/from the distributed memory)**

Here is an operation, run it on all the data.
 - Don't care where or even run it twice!!!

Large set of distributed primitives
 - *M/R, groupby, etc.*

Still computing goes to data!!!

# Big Data

Focused on data
Scalability for the masses
Tradeoff to scale
More cloud oriented

**Coarse grained parallelism**
independent tasks, localized synch
machine based partitioning

**Targets average performance**

how much data is big data?

# HPC

More science driven
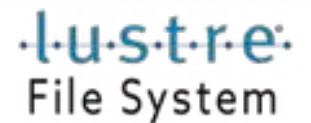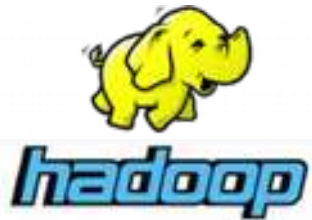Closer to hardware
Cutting edge algorithms
Better defined problems

**Fine grained parallelism**
intercommunicating synchronized tasks
also CPU/GPU based partitioning

**Targets peak performance**

let's get to the guts of your code

# Big Data and HPC

# Big data challenges for HPC

What kind of Big Data problems can be addressed with traditional HPC resources?

Big Data clusters (Spark/noSQL) managed very differently from job scheduling based clusters.

What is the "customer" base? Final users? Programmers?

What is the cost of using Big Data/HPC solutions?
$$$, people, opportunity?

Technological / Non functional requirements (security, streaming data, data delivery SLAs, etc.)

# Big data approaches for HPC

consider container based management (i.e. Openstack).

consider Big Data models (spark) for parallelizing scientific software.

HPC community is strong in algorithmics →
programming/deployment models for Big Data

complementarities: CPU/GPU software running on Spark clusters