Hadoop –
Spark
Overview

J. Allemandou

Generalities

Hadoop Spark

Demo

# Apache Hadoop & Spark – What is it ?

Joseph Allemandou

JOALTECH

17 / 05 / 2018

Joseph Allemandou

1 Generalities on High Performance Computing (HPC)

2 Apache Hadoop and Spark – A Glimpse
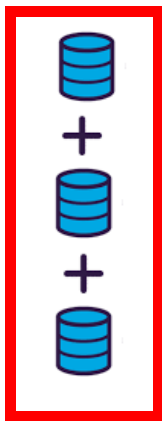
3 Demonstration

**Scale Up**

**Scale out**

**Scale Up**

**Scale out**

Things to consider when doing parallel computing:

- Partitioning (tasks, data)
- Communications
- Synchronization
- Data dependencies
- Load balancing
- Granularity
- IO

Livermore Computing Center - Tutorial

- big data
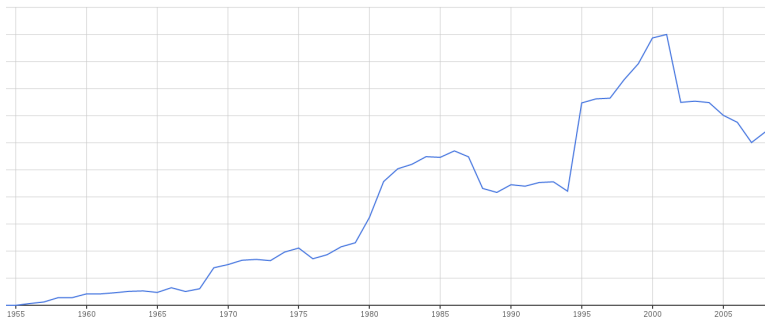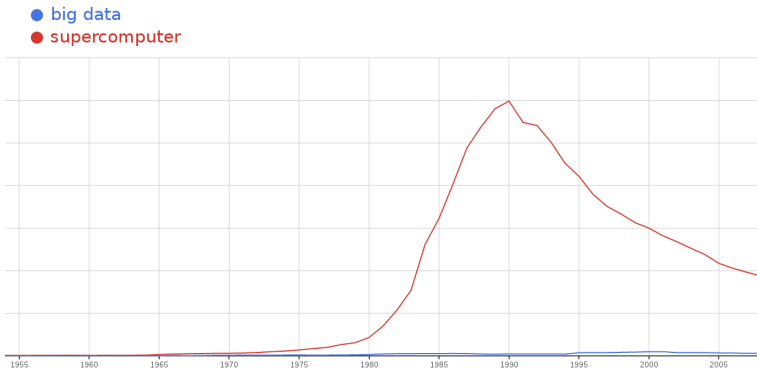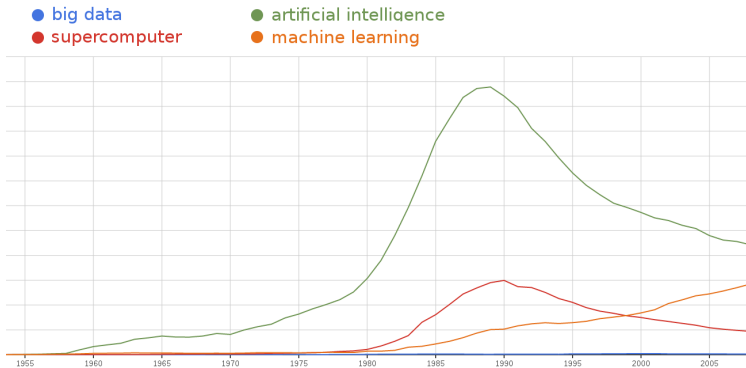


Figure: Google Ngram Viewer

Hadoop –
Spark
Overview

J. Allemandou

Generalities

Hadoop Spark

Demo



Figure: Google Ngram Viewer

Figure: Google Ngram Viewer

Figure: Google Trends

Hadoop –
Spark
Overview

J. Allemandou

Generalities
Hadoop Spark
Demo

**Supercomputer**

Dedicated hardware



Message Passing Interface

**Big Data**

Commodity hardware

**Ifremer**

Hadoop –
Spark
Overview

J. Allemandou

Generalities

Hadoop Spark

Demo

- C / C++ / Fortran / Python

- Low-level API - Send / receive messages

- a lot to do manually
  - split the data
  - assign tasks to workers
  - handle synchronisation
  - handle errors

Hadoop –
Spark
Overview

J. Allemandou

Generalities
Hadoop Spark
Demo

- Java / Scala / Python / R

- High-level API - dataflows

- Less performant than MPI (see this scientific paper)

- Easier to code than MPI (a lot!)

- High-availability oriented

1 Generalities on High Performance Computing (HPC)

2 Apache Hadoop and Spark – A Glimpse

3 Demonstration

- Distributed storage and computation platform (Java, open source)
  - HDFS – <u>H</u>adoop <u>D</u>istributed <u>F</u>ile <u>S</u>ystem
  - YARN – <u>Y</u>et <u>A</u>nother <u>R</u>esource <u>N</u>egotiator

Ifremer

- Distributed storage and computation platform (Java, open source)
  - HDFS – <u>H</u>adoop <u>D</u>istributed <u>F</u>ile <u>S</u>ystem
  - YARN – <u>Y</u>et <u>A</u>nother <u>R</u>esource <u>N</u>egotiator

- Built for scalability
  - Up to thousands of nodes (see this tweet )
  - Can be easily extended (heterogeneous hardware)
  - Is resilient to errors (to some extent)

**lfremer**

- Distributed storage and computation platform (Java, open source)
  - HDFS – Hadoop Distributed File System
  - YARN – Yet Another Resource Negotiator

- Built for scalability
  - Up to thousands of nodes (see  this tweet )
  - Can be easily extended (heterogeneous hardware)
  - Is resilient to errors (to some extent)

- Working with multiple computation engines
  - Hadoop MapReduce (the good old one)
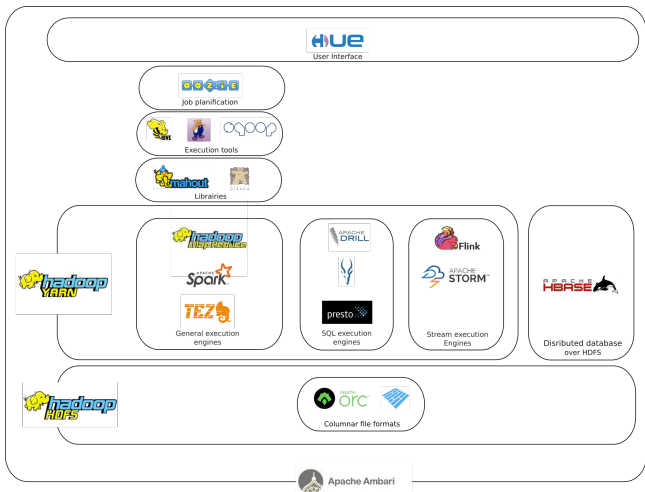  - Apache Spark (the new kid on the block)

**Ifremer**

- HDFS
    - Files are split in blocks – Partitioning !
    - Blocks are replicated – Failure tolerance !
    - Master - Slave architecture (NameNode & DataNodes)

- HDFS
  - Files are split in blocks – Partitioning !
  - Blocks are replicated – Failure tolerance !
  - Master - Slave architecture (NameNode & DataNodes)

- YARN
  - Execution containers with defined resources
  - Scheduler manages resource sharing
  - Master - Slave architecture (ResourceManager & NodeManagers)

- HDFS
    - Files are split in blocks – Partitioning !
    - Blocks are replicated – Failure tolerance !
    - Master - Slave architecture (NameNode & DataNodes)

- YARN
    - Execution containers with defined resources
    - Scheduler manages resource sharing
    - Master - Slave architecture (ResourceManager & NodeManagers)

- Global
    - Execution engines exploit *data locality*

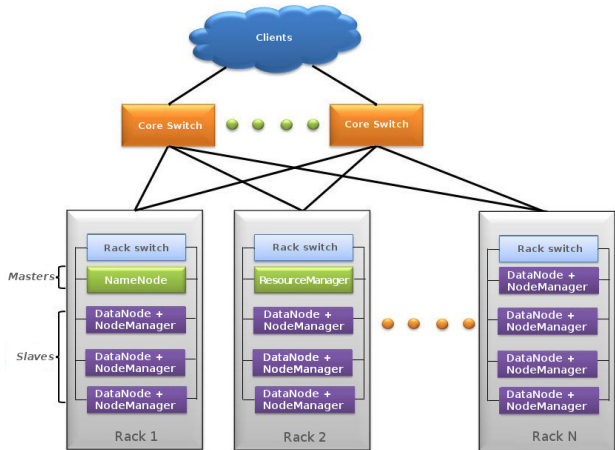# Apache Hadoop - The Frame (HDFS + YARN)
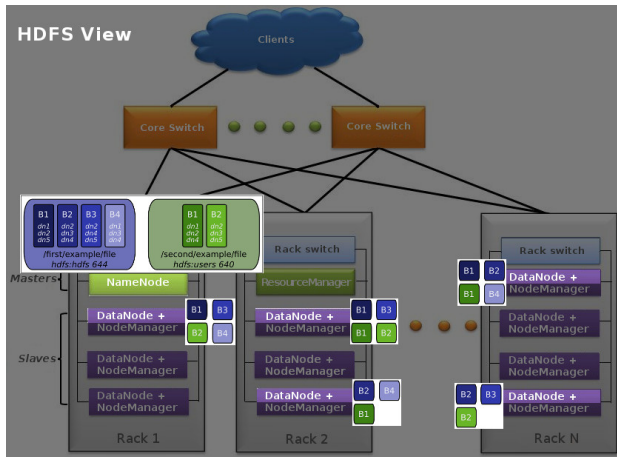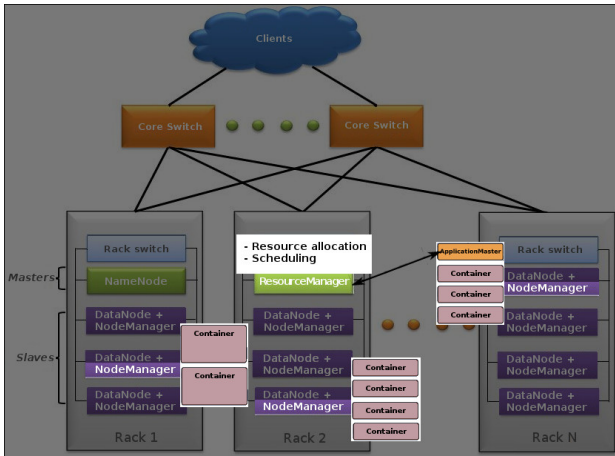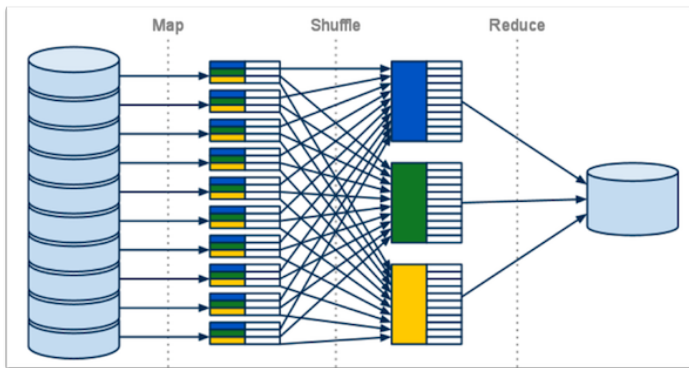
Hadoop –
Spark
Overview

J. Allemandou

Generalities

Hadoop Spark
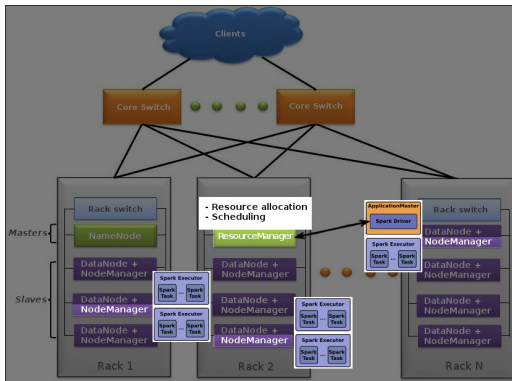
Demo

- Programming data-flows instead of single map-reduce steps
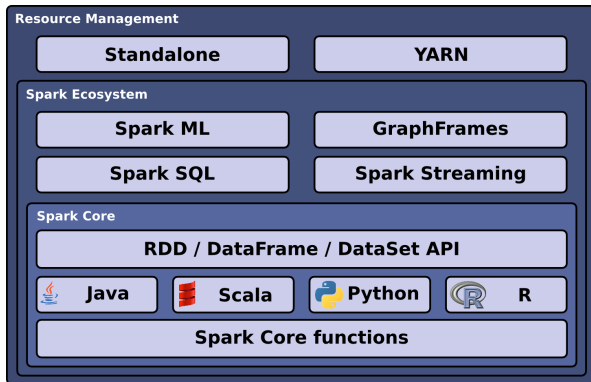  - A lot less code for complex flows
  - Data lineage allows for error recovery

- Programming data-flows instead of single map-reduce steps
  - A lot less code for complex flows
  - Data lineage allows for error recovery

- Decoupling of tasks and containers – Spark executors run multiple tasks

  - Less executor management overhead
  - Executors can reuse RAM - Caching!

# Spark - A big ecosysem as well

**Resource Management**

| Standalone | YARN |
|---|---|

**Spark Ecosystem**

| Spark ML | GraphFrames |
|---|---|
| Spark SQL | Spark Streaming |

**Spark Core**

RDD / DataFrame / DataSet API

| Java | Scala | Python | R |
|---|---|---|---|

Spark Core functions

Hadoop –
Spark
Overview

J. Allemandou

Generalities

Hadoop Spark

Demo

Mediawiki History Reconstruction Job

# Plan

1. Generalities on High Performance Computing (HPC)

2. Apache Hadoop and Spark – A Glimpse

3. Demonstration

Notebook with Spark on Docker

Hadoop –
Spark
Overview

J. Allemandou

Generalities

Hadoop Spark

Demo