

# Initiation au Big Data et traitement de données massives avec Apache Spark

# A propos de Moi

## **2017-2019:**

- Msc en Big Data, AIMS Mbour
- Msc en Algèbre appliquée à la Cryptographie, Université de Thiès

## **2019-2020:** Big Data Engineer @ MNS Consulting

- Installation d'un environnement Big Data et Mis en place d'une plateforme d'analyse et de prise de décision, se basant sur les données mobile.
- Mis en place de plateforme de visualisation des données provenant des capteurs IOT's.

## **2020-Present:** Big Data Engineer @ Atos Sénégal

- Aider les métiers sur le développement des KPI's.
- Mise en place de pipeline d'ingestion de données pour alimenter le Datalake
- Mettre à la disposition des data scientists, les données agrégées.



Ibrahima FALL  
Lead Big Data Engineer @ Atos Senegal,  
Email: [ibrahima.fall@atos.net](mailto:ibrahima.fall@atos.net)  
Alt-email: [iboudofall@gmail.com](mailto:iboudofall@gmail.com)

# Sommaire



1. Le Big Data et ses caractéristiques

2. Hadoop et son écosystème

3. Initiation à Spark

- Architecture
- Spark API's
- Spark operations
- Exécution Plan
- Demo

# Definitions



BIG DATA

## **What is Data ?**

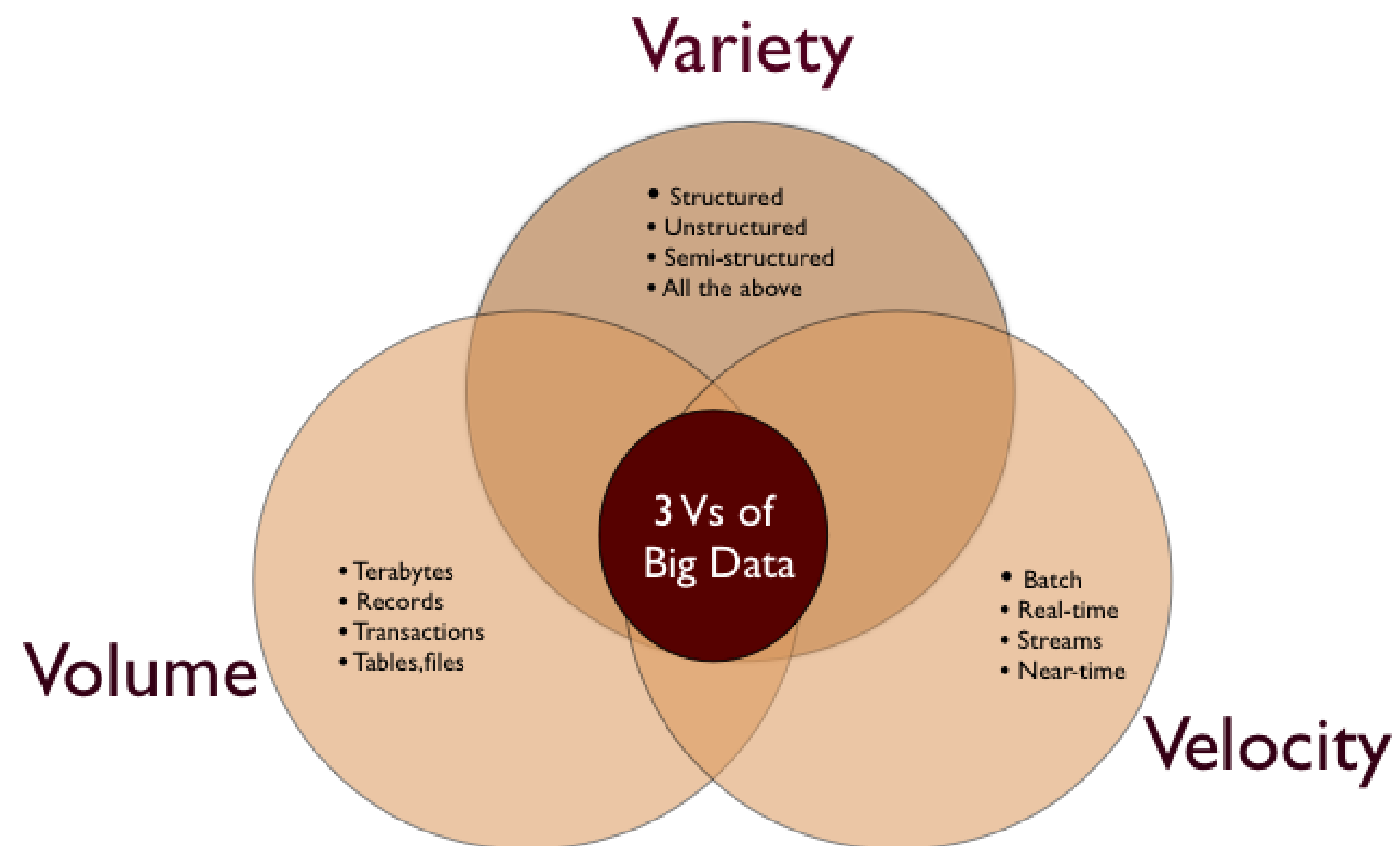
# Definitions



BIG DATA

## What is Big Data ?

Data with specific characteristics



Technology, Strategy, Technics





# Fast Growing

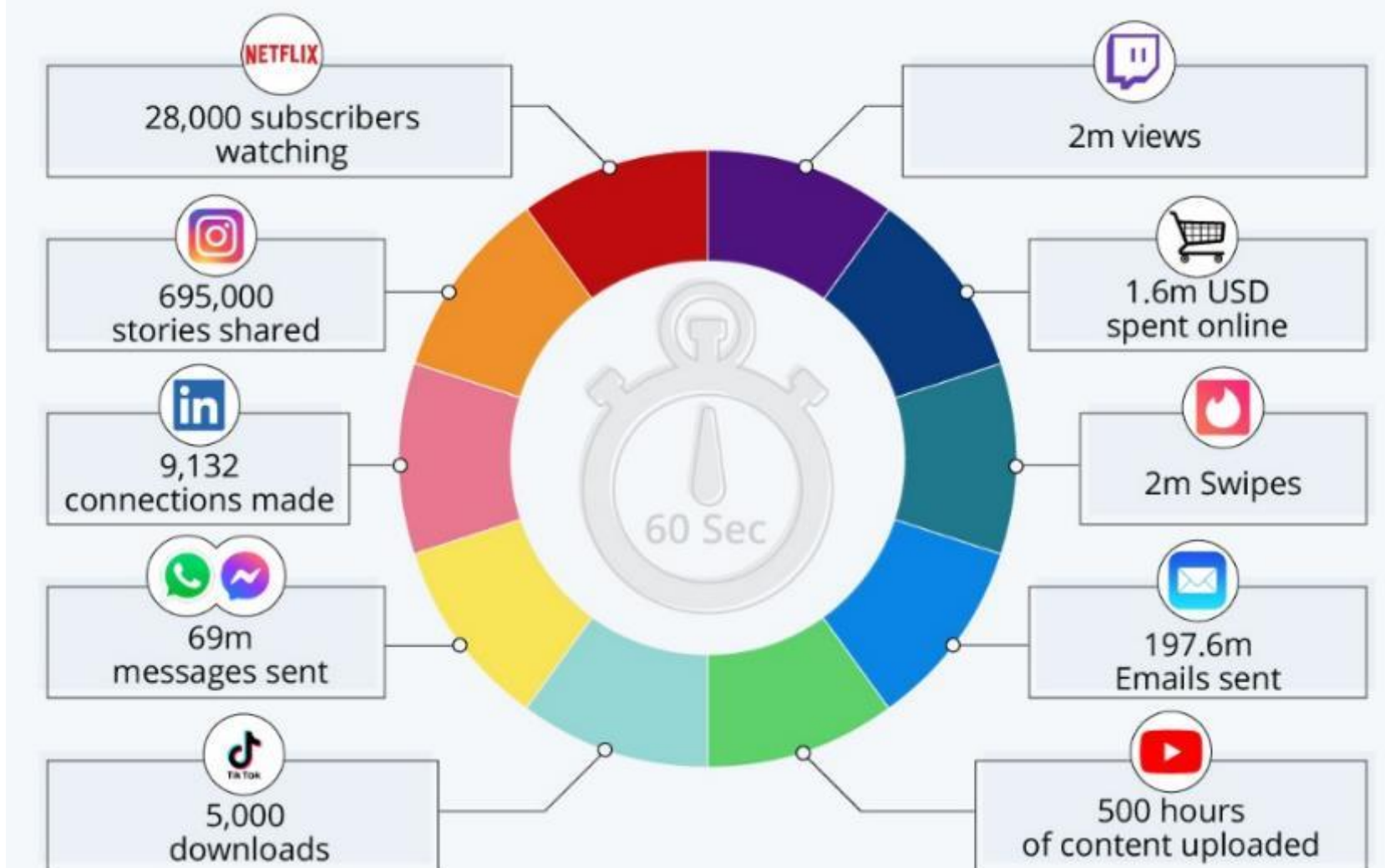


BIG DATA

- Facilité de collecter des données grâce à une automatisation des systèmes.
- Des capteurs omniprésents pour générer des données
  - Téléphone portable
  - Logs des logiciels
  - Caméras
  - Capteurs IOT
  - ...
- Le stockage devient accessible
- Prise de conscience que les données étaient très précieuses pour être supprimées

## A Minute on the Internet in 2021

Estimated amount of data created on the internet in one minute



Source: Lori Lewis via AllAccess

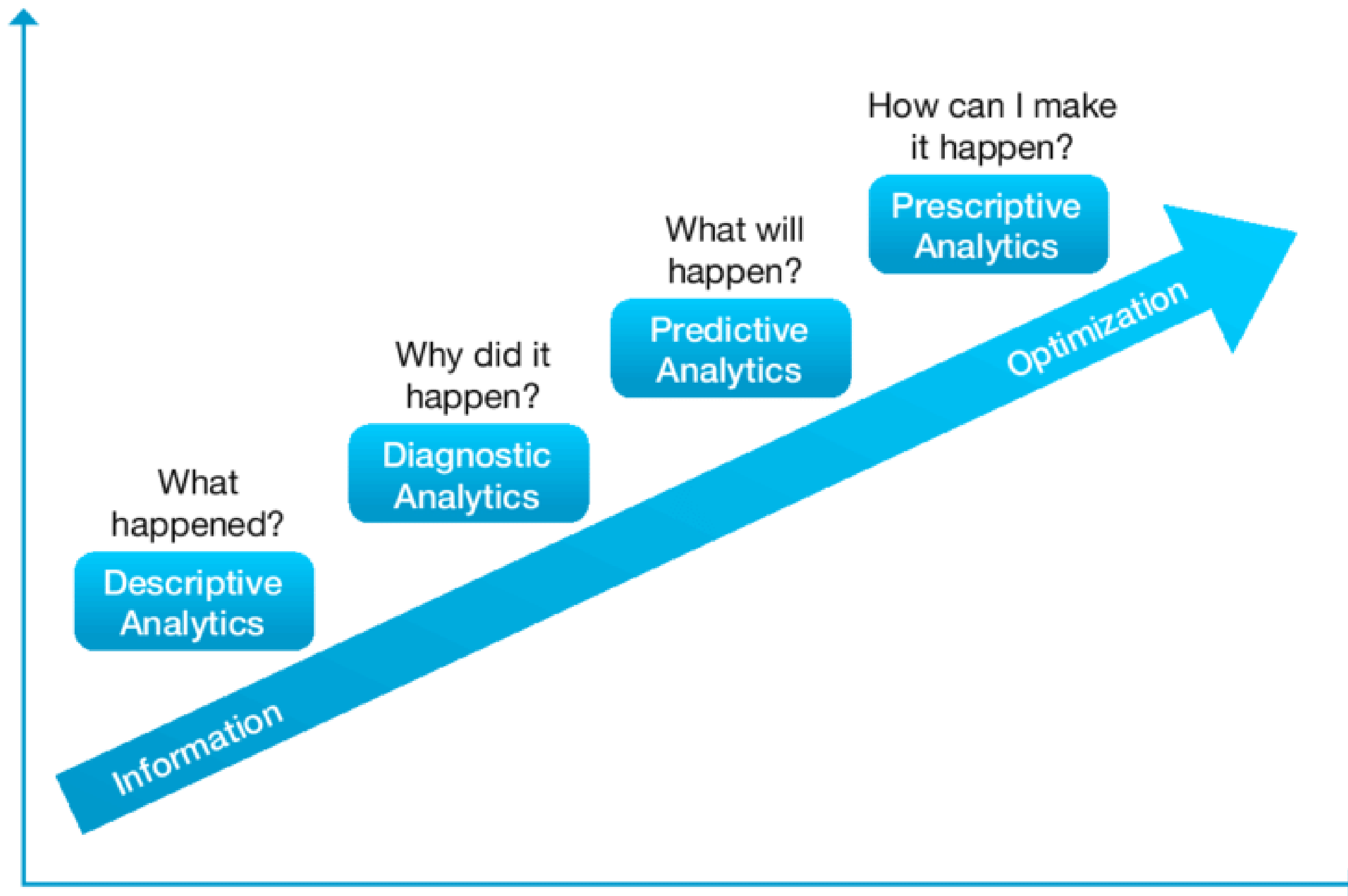


statista

# Type of analysis



BIG DATA



# Data Lifecycle



BIG DATA





# Big Data Jobs



## Data Scientist

La science des données est une méthode systématique dédiée à la connaissance, découverte via l'analyse des données

- En entreprise, optimiser les processus organisationnels pour Efficacité
- En sciences, analyser des données expérimentales/observationnelles pour
- dériver des résultats

## Skills:

- Computer Science
- Statistiques + Mathématiques
- Programmation: Python, R, SAS, JAVA, SQL
- Machine Learning
- Visualisation
- Connaissances du domaine

# Big Data Jobs



## Data Engineer

L'ingénierie des données est le domaine qui développe et fournit

- Des systèmes de gestion et d'analyse de données volumineuses
- Construire des plates-formes de données modulaires et évolutives pour les données scientifiques
- Déployer des solutions Big Data

Skills:

- Computer Science
- Bases de données
- Traitement de données en temps réel / données massives
- Programmation: Python, Scala, SQL, ...
- Comprendre les facteurs de performance et les limites des systèmes

# Big Data Jobs



## **Data Analyst**

Analyser les données pour les transformer en informations exploitables

- Définir la stratégie Data-Driven de l'entreprise
- Créer et maintenir les bases de données de l'entreprise
- Elaborer les critères de segmentation

Skills:

- Computer Science
- Statistiques
- Utilisation Excel
- Visualisation de données
- Programmation: Python, Scala, SAS, SQL, ...



# Data Engineer Vs Others



BIG DATA

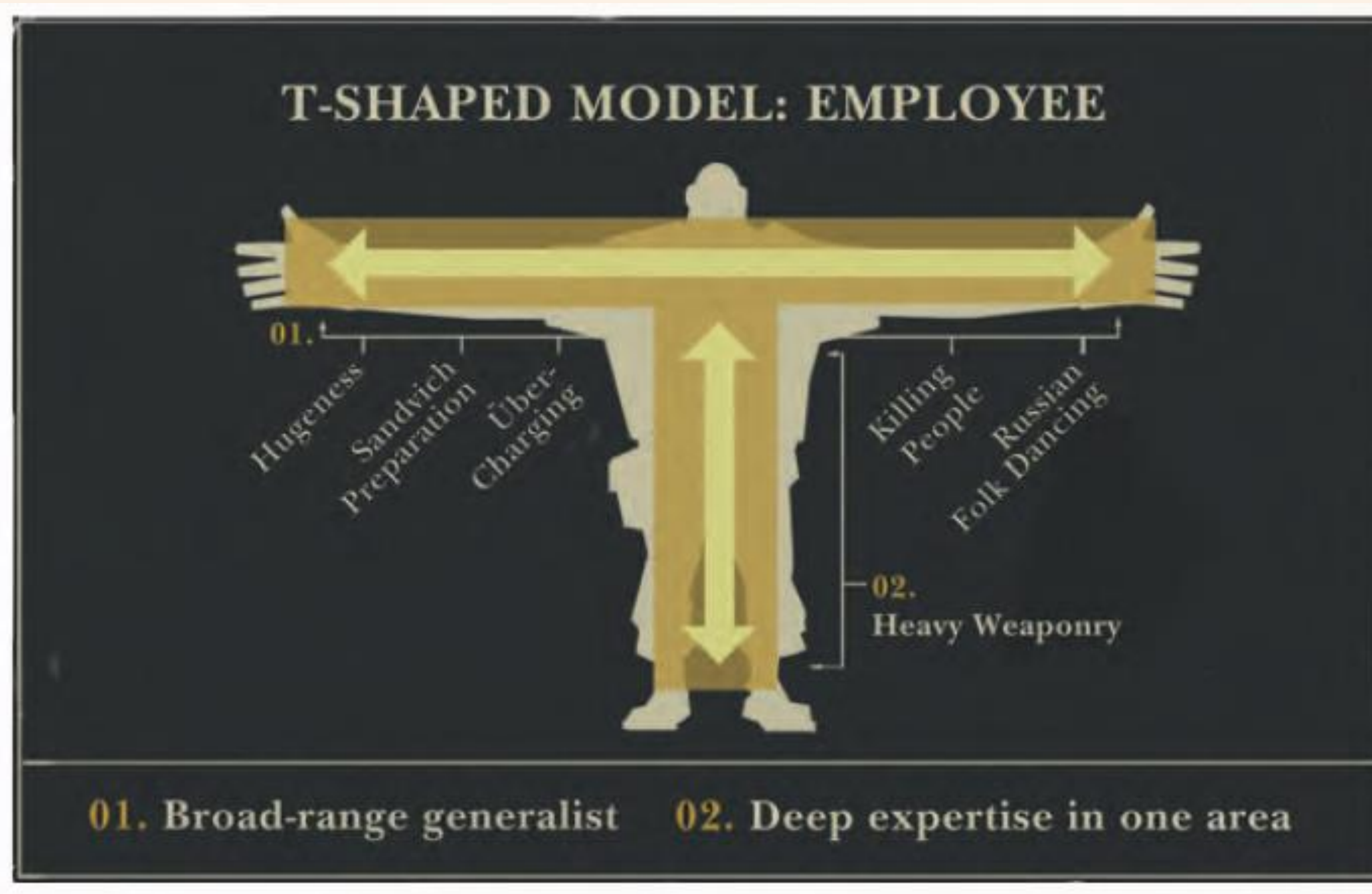




# Be T-shaped

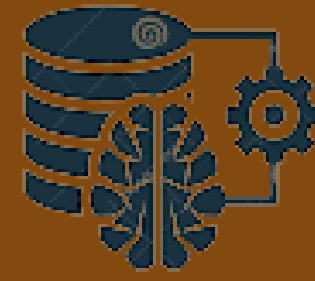


BIG DATA





# Use Cases



BIG DATA

## Smarter Healthcare



## Multi-channel



## Finance



## Log Analysis



## Homeland Security



## Traffic Control



## Telecom



## Search Quality



## Manufacturing



## Trading Analytics



## Fraud and Risk



## Retail: Churn, NBO





# Big data Challenges



BIG DATA



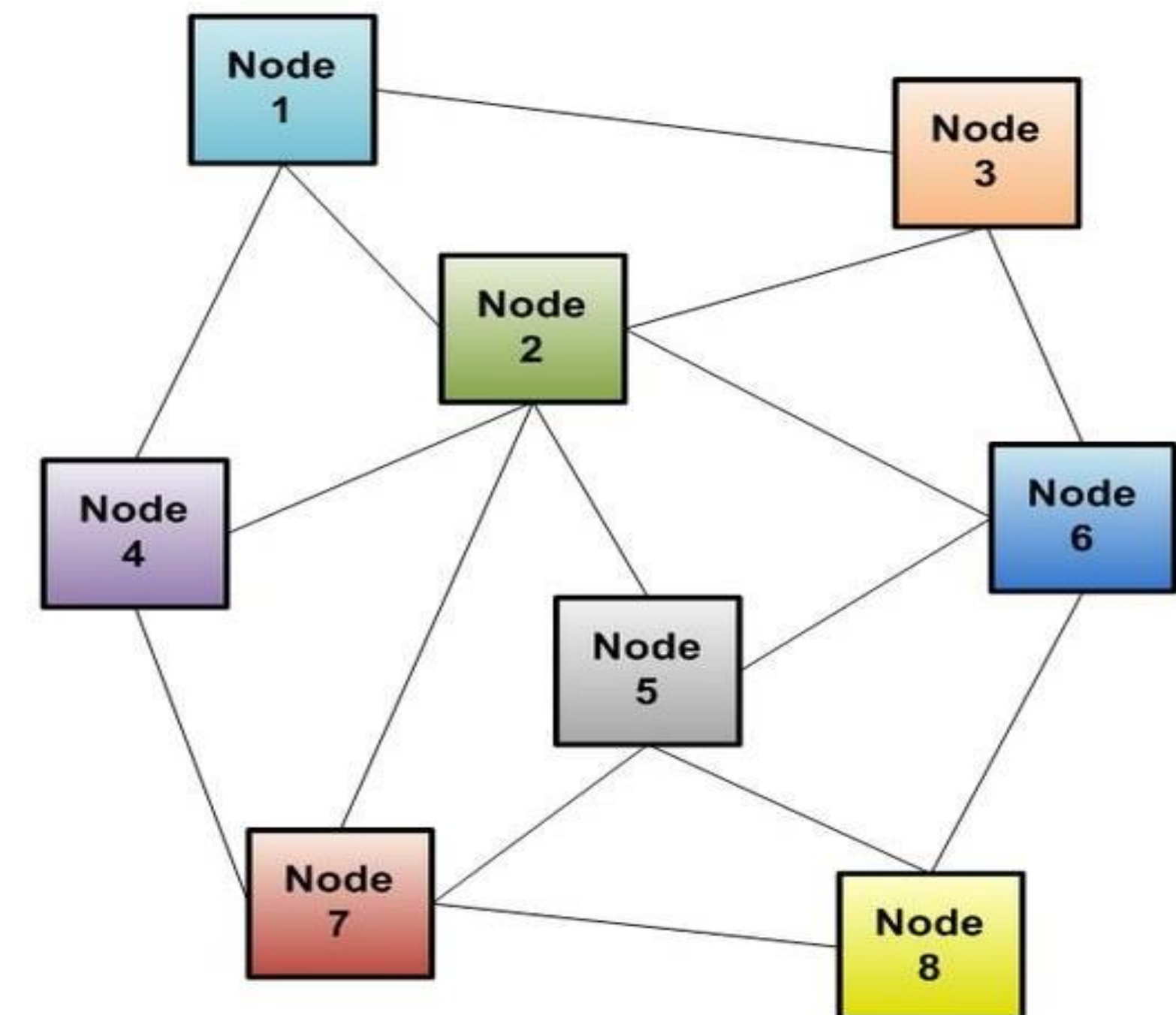
# Distributed System



Un système distribué est un environnement informatique dans lequel divers composants sont répartis sur plusieurs ordinateurs (ou autres appareils informatiques) sur un réseau. Ces appareils divisent le travail, coordonnent leurs efforts pour accomplir le travail plus efficacement.

## Caractéristiques:

- Scalabilité
- Concurrency
- Disponibilité / fault tolerance
- Transparence
- Hétérogénéité
- Réplication
- Géo-distribution





# What is Hadoop *Spark*



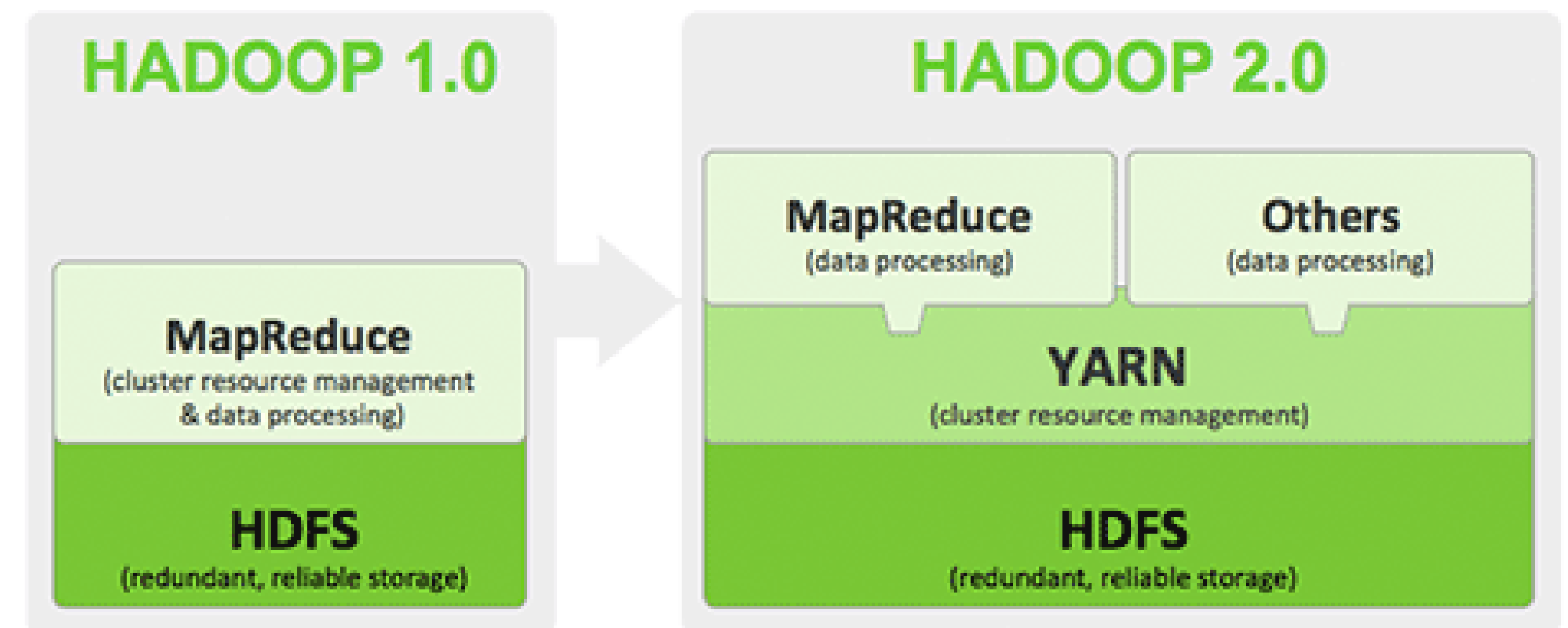
The Apache™ Hadoop® project develops open-source software for reliable, scalable, distributed computing.

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

➤ Hadoop est une plateforme Big Data open-source, utilisée pour stocker et traiter d'immenses volumes de données d'une manière distribuée.

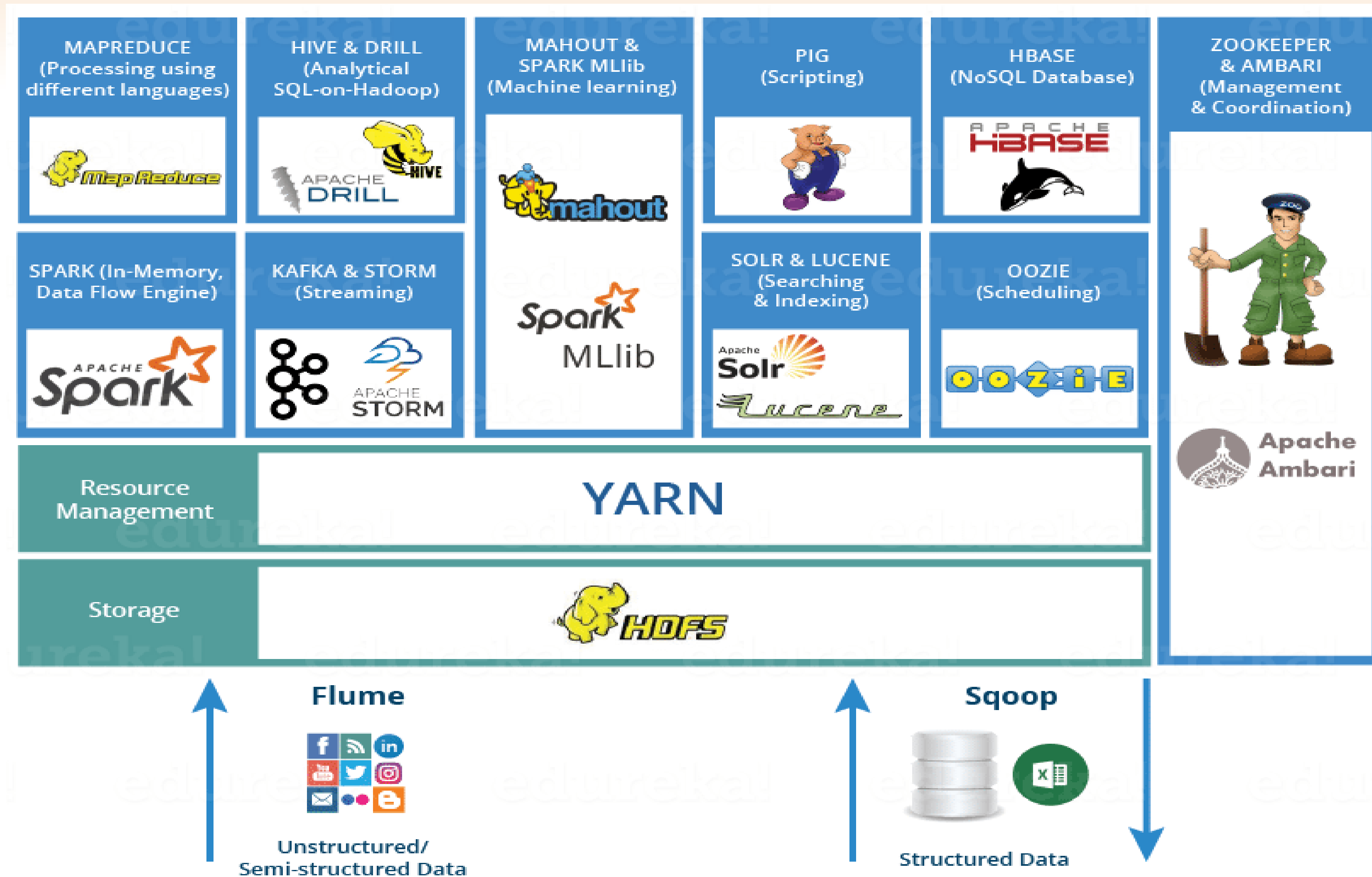
➤ Composants principaux:

- HDFS (Hadoop Distributed File System)
- YARN (Yet Another Resource Negotiator)
- MapReduce



# Hadoop Ecosystem

Spark





# What is Spark?



Apache Spark est un moteur de traitement de données ,

- in-memory
- open-source
- ultra-rapide
- tolérant aux pannes
- supporte plusieurs languages
- Moteur unifié
- lazy computations

pour le big data et le machine learning, avec des modules  
pour des traitements de données en temps reel, sql et calcul  
de graph.



Lightning-fast unified analytics engine

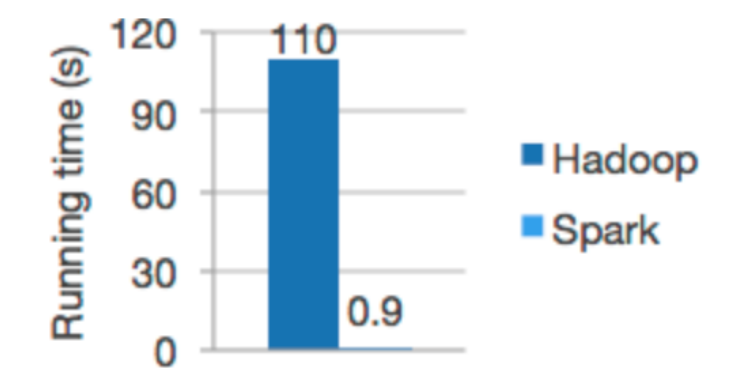
[Download](#)[Libraries ▾](#)[Documentation ▾](#)[Examples](#)[Community ▾](#)[Developers ▾](#)

**Apache Spark™** is a unified analytics engine for large-scale data processing.

## Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.



Logistic regression in Hadoop and Spark

## Ease of Use

Write applications quickly in Java, Scala, Python, R, and SQL.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python, R, and SQL shells.

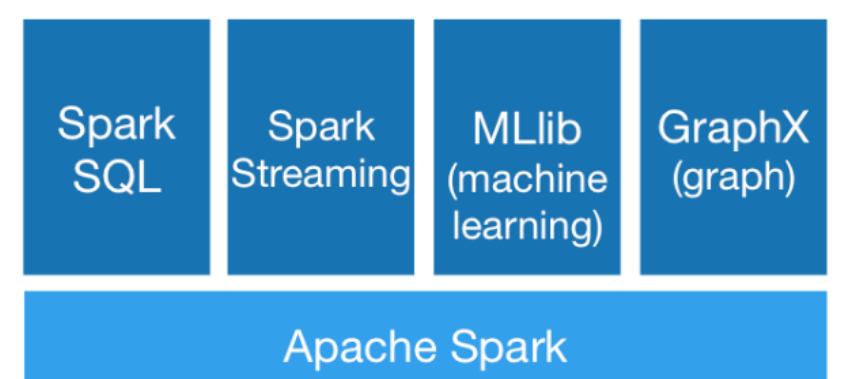
```
df = spark.read.json("logs.json")
df.where("age > 21")
  .select("name.first").show()
```

Spark's Python DataFrame API  
Read JSON files with automatic schema inference

## Generality

Combine SQL, streaming, and complex analytics.

Spark powers a stack of libraries including [SQL and DataFrames](#), [MLlib](#) for machine learning, [GraphX](#), and [Spark Streaming](#). You can combine these libraries seamlessly in the same application.



# Spark vs Hadoop ?



- En réalité, Spark est comparable à MapReduce et non à Hadoop.
- Spark peut se brancher à Hadoop et utiliser HDFS pour bénéficier du système de stockage de cette dernière.
- Spark est plus rapide que MapReduce (10 à 100 fois)
- Spark est écrit en Scala, MapReduce est écrit en Java
- En traitement, Spark supporte du **batch**/ **real-time**/ **graph** et MapReduce ne supporte que du **batch**
- MapReduce ne supporte pas du Caching, contrairement à Spark

# Why Spark ?

Spark



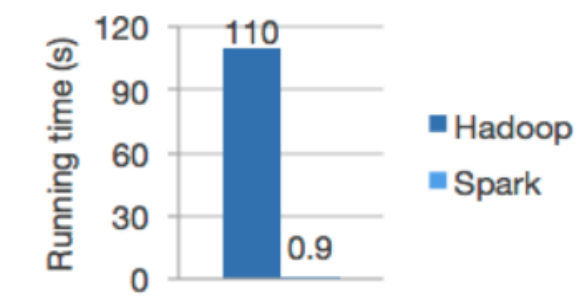
[Download](#) [Libraries](#) [Documentation](#) [Examples](#) [Community](#) [Developers](#)

**Apache Spark™** is a unified analytics engine for large-scale data processing.

## Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.



Logistic regression in Hadoop and Spark

Syntaxes très simples

```
df = spark.read.json("logs.json")
df.where("age > 21")
  .select("name.first").show()
```

Spark's Python DataFrame API  
Read JSON files with automatic schema inference

Moteur unifié

Rapidité

Facilite le développement avec les api Java, Scala, Python, R, SQL

## Ease of Use

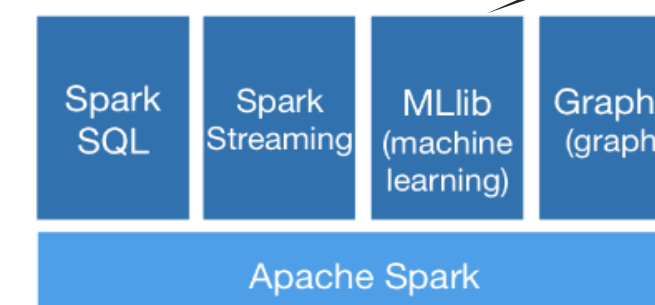
Write applications quickly in Java, Scala, Python, R, and SQL.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python, R, and SQL shells.

## Generality

Combine SQL, streaming, and complex analytics.

Spark powers a stack of libraries including [SQL and DataFrames](#), [MLlib](#) for machine learning, [GraphX](#), and [Spark Streaming](#). You can combine these libraries seamlessly in the same application.



## Runs Everywhere

Spark runs on Hadoop, Apache Mesos, Kubernetes, standalone, or in the cloud. It can access diverse data sources.

You can run Spark using its [standalone cluster mode](#), on [EC2](#), on [Hadoop YARN](#), on [Mesos](#), or on [Kubernetes](#). Access data in [HDFS](#), [Alluxio](#), [Apache Cassandra](#), [Apache HBase](#), [Apache Hive](#), and hundreds of other data sources.



Flexibilité

# Use of Spark



Suivant le projet, Spark peut être utilisé

- En mode single node
- En mode cluster, sur plusieurs machines
  - On Premise
  - On the cloud: Amazon Web Services(AWS), Microsoft Azure

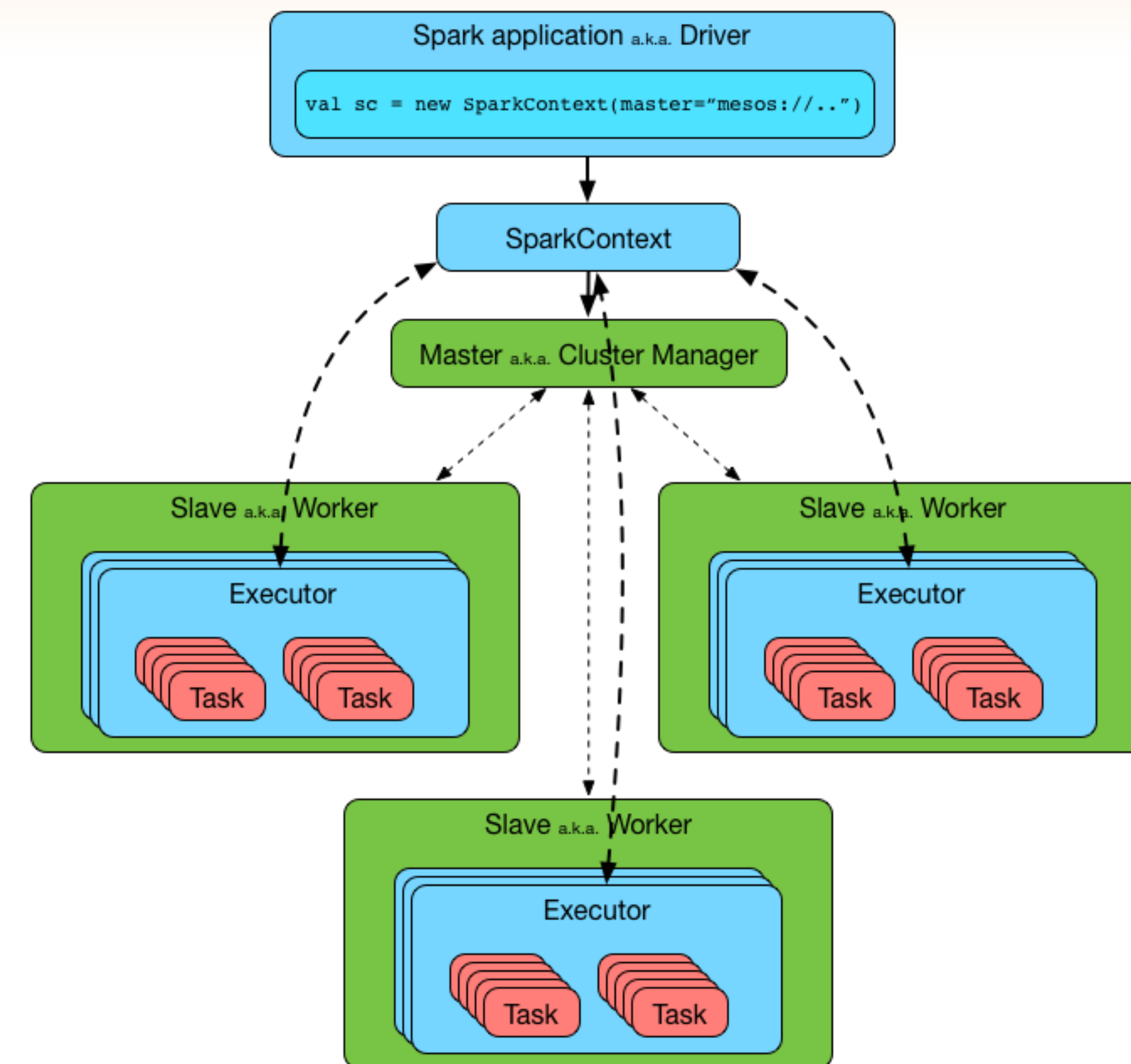


# Spark Architecture



- Driver : c'est le processus contenant le programme spark, il crée le DAG (Directed Acyclic Graph), planifie et coordonne l'exécution du programme
- SparkContext (SparkSession) : il représente la porte d'entrée du Spark Cluster avec des configurations prédéfinies (Master, nom de l'application, mémoire, nombre d'exécuteurs ...)
- Cluster Manager : négocie de la ressource et la rend disponible pour le driver, suit le statut des workers.
- Executors : lieu d'exécution des tasks du DAG

Master et le Cluster Manager peuvent être fusionnés comme en mode StandAlone ou séparés comme en mode YARN



Source: <https://sunbiaobiao.gitbooks.io/sa/content/spark-architecture.html>



# Spark API's



## Spark APIs: RDD, Dataset et DataFrame

### RDD

- RDD = Resilient Distributed Dataset (jeux de données distribués et résilients)
- RDDs sont les structures fondamentales de Spark
- RDDs sont tolérants aux pannes, immuables, partitionnés, distribués, supportent l'évaluation paresseuse (lazy evaluation)
- Compile-time Type Safe

### DataFrames

- Même structure que pandas dataframe et dataframes dans R
- Représente une abstraction des RDDs
- Supporte Catalyst Optimisation
- Erreurs syntaxiques détectées à la compilation
- Erreurs analytiques détectées au runtime

### Dataset

- Un union de RDD et DataFrame

# Spark Operations



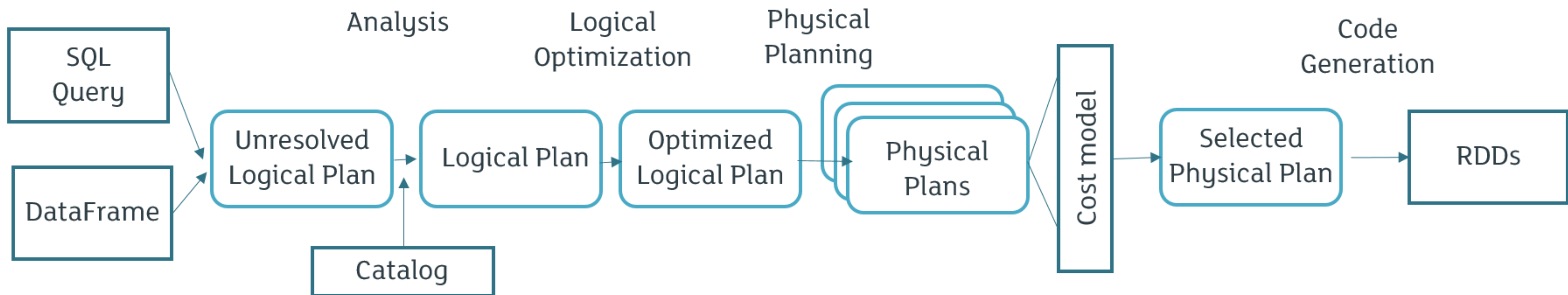
Spark dispose de deux types d'operations:

- Transformations: créer une nouvelle dataset à partir d'une autre existante
- Actions: retourne une valeur finale ou écrit le résultat dans un espace de stockage externe.

Les transformations sont dites "lazy": elles ne s'exécutent pas immédiatement: en effet, elles gardent juste l'ensemble des operations qui ont permis de construire la dataset résultante (lineage). Les transformations ne sont exécutées que lorsqu'une action est requise.

L'ensemble des operations faites sur un dataset forme un DAG (Directed Acyclic Graph) dont les sommets sont des RDD et les arêtes représentent les operations.

# Plan d'exécution *Spark*



**Enough Talk Let's Code**

