# Big Data

## Algorithmes PySpark pour traitement de donnée

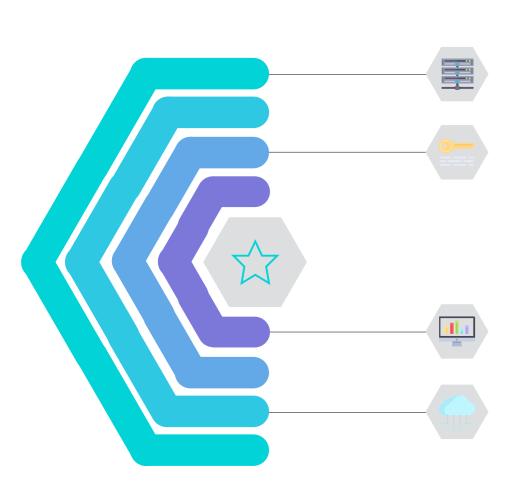
Nous créons environs 2,5 trillions d'octets de données par jour. En 2 ans nous avons créé 90% des données présentes dans le monde. Ce nombre de données provient de l'avancée de toutes les technologies et aussi de nos habitudes usuelles telles que l'envoie de messages, mails ou vidéos. C'est pourquoi il en devient difficile de travailler sur ce volume de données avec les outils classiques de gestion de base de données.

De plus le nombre de données à traiter ne cesse d'augmenter et étant lié à nos habitudes de vie, ne cesseront d'augmenter, c'est pourquoi la manipulation des données représente bel et bien un enjeu extrêmement important pour le futur.

# Le Big Data

Le Big Data désigne des ensembles de donnés devenant tellement volumineux qu'il en devient Vitesse impossible de les manipuler avec des outils classiques de gestion de base de données. Ces données sont caractérisés par ce qu'on appelle les "3V" Variété Volume

# L'écosystème de Spark



## **SPARK STREAMING**

Peut être utilisé pour traitement temps-réel des données en flux.

### SPARK SQL

Permet d'exposer les jeux de données Spark via API JDBC et d'exécuter des requêtes de type SQL en utilisant les outils BI et de visualisation traditionnels.

## **SPARK MLLIB**

Librairie de machine learning contenant tous les algorithmes et utilitaires d'apprentissage classiques, comme la classification, la régression, le clustering, le filtrage collaboratif, la réduction de dimensions, en plus des primitives d'optimisation sous-jacentes. **SPARK GRAPHX** 

La nouvelle API (en version alpha) pour les traitements de graphes et de parallélisation de

# Qu'est ce que Spark? Spark?



### LE DÉVELOPPEMENT

Développé par AMPLab, de l'Université UC Berkeley et passé sous forme de projet Apache en 2010

### **HADOOP**

Permet à des applications sur clusters Hadoop d'être exécutées jusqu'à 100 fois plus vite en mémoire et 10 fois plus vite sur disque



Framework de traitements

Big Data open source

réalisé pour effectuer des

analyses sophistiquées et

conçu pour la rapidité et

la facilité d'utilisation

**APACHE SPARK** 

# LE FRAMEWORK

Spark propose un framework complet et unifié pour subvenir aux besoins de traitements Big Data pour divers jeux de données, divers par leur nature aussi bien que par le type de source

## **APPLICATIONS**

Écrire rapidement des applications en Java, Scala ou Python et inclut plus de 80 opérateurs haut-niveau. Possible d'utiliser de façon interactive pour requêter les données depuis un shell. Spark supporte les requêtes SQL, le streaming de données et des fonctionnalités de machine learning et de traitements orientés graphe

# Les fonctionnalités de Spark

## LANGAGE DE SPARK

Spark est écrit en Scala et s'exécute sur la Machine Virtuelle Java (JVM)



AMÉLIORE MAPREDUCE Apporte des améliorations à MapReduce grâce à des étapes de shuffle moins coûteuses. Mais Spark n'en reste pas là et propose d'autres fonctions que Map et Reduce

# **FLEXIBLE**

Spark maintient les résultats intermédiaires en mémoire plutôt que sur disque. Mais le moteur d'exécution est conçu pour travailler aussi bien sur les deux.



## **CINQ LANGAGES**

Spark propose une interface de programmation de haut-niveau pour une meilleure productivité (API en Java, Scala et Python) et un shell interactif pour Scala et Python. Spark supporte également le Clojure et le R

## **EVALUATION PARESSEUSE**

Spark supporte les évaluations paresseuses ("lazy evaluation") des requêtes, ce qui aide à l'optimisation des étapes de traitement.

## **OPTIMISE LES GRAPHES**

L'optimisation de graphes d'opérateurs arbitraires

La méthode

Conclusion

JVM: La machine virtuelle Java est un appareil informatique virtuel qui exécute des programmes compilés sous forme de bytecode Java.

SQL: Langage informatique normalisé servant à exploiter des bases de données relationnelles. Machine Learning: Mise en place d'algorithmes en vue d'obtenir une analyse prédictive à partir de données dans un but précis.

MapReduce : Modèle de programmation massivement parallèle adapté au traitement de très grandes quantités de données.

Framework: Ensemble d'outils constituant les fondations d'un logiciel informatique ou d'applications web, et destiné autant à faciliter le travail qu'à augmenter la productivité du programmateur qui l'utilisera. Open Source: Permet de distribuer et d'utiliser gratuitement un logiciel, ainsi que de le modifier et de l'améliorer en donnant accès à son code source.

Shell: interface utilisateur d'un système d'exploitation destinée à lancer d'autres programmes et gérer leurs interactions. API: (Application Programming Interface) interface de programmation qui permet de se "brancher" sur une application pour échanger des données.

JDBC: (Java Database Connectivity) intergiciel qui permet à une application Java de manipuler plusieurs bases de données. Intergiciel : logiciel servant d'intermédiaire de communication entre plusieurs applications, généralement complexes ou distribuées sur UN réseau informatique.

Cluster: concentration géographique d'entreprises reliées ensemble, de fournisseur, et d'institutions associés dans un domaine particulier. Graphe: ensemble de points nommes noeuds, sommets ou cellules reliés par un segment fléché ou non nommé arrête