

## Formation Hadoop pour le traitement distribué des Big Data

### APERÇU / PLAN DU COURS

#### DESCRIPTION DU COURS

##### Brève description :

Les entreprises produisent d'énormes quantités de données chaque jour. Ces données sont stockées puis traitées et analysées pour en tirer de la valeur. Grâce aux plateformes de stockage et de traitements distribués de type Hadoop, il est devenu plus facile pour les ingénieurs de répondre aux problématiques du Big Data avec une grande efficacité et à un coût réduit.

Ce séminaire vous offre l'occasion de vous essayer au traitement distribué de données massives via la plateforme Hadoop et ses outils comme Hive et Spark.

Objectifs	Les participants seront capables de : <ol style="list-style-type: none"><li>1. Stocker des données sur Hadoop</li><li>2. Décrire les différentes étapes du modèle de programmation MapReduce</li><li>3. Charger des données sur Hive et les traiter à l'aide du langage HiveQL</li><li>4. Charger des données dans Spark et les traiter à l'aide du langage Spark SQL</li></ol>
Dates	<b>14 au 17 mars 2023</b>
Durée	<b>4 jours</b>

#### PUBLIC CIBLE

- Développeurs
- Experts en business intelligence
- Analystes de données

#### PRE-REQUIS

- Notions en administration système sous Linux.
- Notions en programmation orientée objet
- Connaissances de base du langage SQL

## AGENDA ET CONTENU DU COURS

<p><b>Contenu Pédagogique</b></p>	<p><b>L'écosystème du Big Data</b></p> <ul style="list-style-type: none"> <li>▪ Les origines et caractéristiques du Big Data</li> <li>▪ Les avancées technologiques</li> <li>▪ L'informatique parallèle</li> <li>▪ Aperçu des plateformes Big Data</li> <li>▪ Aperçu des services cloud de AWS</li> <li>▪ <b>Démonstration</b> : Requêtes SQL sur des données massives avec Amazon Athena</li> </ul> <p><b>Le paradigme MapReduce</b></p> <ul style="list-style-type: none"> <li>▪ Historique</li> <li>▪ Le pattern MapReduce</li> <li>▪ Cas d'utilisation de MapReduce</li> <li>▪ <b>Travaux pratiques</b> : programmation MapReduce</li> </ul> <p><b>Hadoop</b></p> <ul style="list-style-type: none"> <li>▪ Historique</li> <li>▪ Socle technique</li> <li>▪ Hadoop – YARN</li> <li>▪ Ecosystème Hadoop</li> <li>▪ <b>Travaux pratiques</b> : programmation Hadoop MapReduce</li> </ul> <p><b>Apache Hive</b></p> <ul style="list-style-type: none"> <li>▪ Présentation</li> <li>▪ Langage de définition de données</li> <li>▪ Interrogation des données</li> <li>▪ Gestion des données complexes</li> <li>▪ <b>Travaux pratiques</b> : analyse de données avec HiveQL</li> </ul> <p><b>Spark SQL</b></p> <ul style="list-style-type: none"> <li>▪ Apache Spark</li> <li>▪ Les RDD</li> <li>▪ Installation de Apache Spark en local</li> <li>▪ <b>Travaux pratiques</b> : calcul distribué avec PySpark</li> <li>▪ Spark SQL et DataFrames</li> <li>▪ <b>Travaux pratiques</b> : analyse de données avec Spark SQL</li> </ul>
-----------------------------------	--

## MODES D'ANIMATION PEDAGOGIQUE

- Exposés théoriques ;
- Questions-réponses ;
- Exercices d'application

**Formation continue**

- Démonstrations
- Travaux pratiques

**SUPPORTS PEDAGOGIQUES**

---

- Support de cours ;
- Les Travaux Pratiques : sur le poste de travail et dans le Cloud de AWS via le service Amazon Elastic Map Reduce (EMR) ;
- Le langage de programmation : Python.