

Formation Traitement distribué pour les 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	<p>Les participants seront capables de :</p> <ol style="list-style-type: none"> 1. Stocker des données sur Hadoop 2. Implémenter en Python des algorithmes MapReduce et les exécuter sur une plateforme Hadoop 3. Charger des données sur Hive et les traiter à l'aide du langage HiveQL 4. Charger des données dans Spark et les traiter à l'aide du langage Spark SQL
Dates	08 au 10 Juin 2021
Durée	3 jours

PUBLIC CIBLE

- Développeurs
- Ingénieurs recherche et développement
- Experts en business intelligence
- Analystes de données

PRE-REQUIS

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

AGENDA ET CONTENU DU COURS

<p>Contenu Pédagogique</p>	<p>L'écosystème du Big Data</p> <ul style="list-style-type: none"> ▪ Les origines et caractéristiques du Big Data ▪ Les avancées technologiques ▪ L'informatique parallèle ▪ Les acteurs du Big Data <p>Le paradigme MapReduce</p> <ul style="list-style-type: none"> ▪ Historique ▪ Le pattern MapReduce ▪ Travaux pratiques : programmation MapReduce <p>Hadoop</p> <ul style="list-style-type: none"> ▪ Historique ▪ Socle technique ▪ HDFS ▪ Hadoop . YARN ▪ Ecosystème Hadoop ▪ Les distributions Hadoop ▪ Travaux pratiques : programmation Hadoop MapReduce <p>Apache Hive</p> <ul style="list-style-type: none"> ▪ Présentation ▪ Création d'une base de données sous Hive ▪ Commandes HQL ▪ Travaux pratiques : traitement de données avec HiveQL <p>Spark SQL</p> <ul style="list-style-type: none"> ▪ Apache Spark ▪ Les RDD ▪ Travaux pratiques : calcul distribué avec PySpark ▪ Spark SQL et DataFrames ▪ Travaux pratiques : traitement de données avec Spark SQL
---------------------------------------	---

MODES D'ANIMATION PEDAGOGIQUE

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

SUPPORTS PEDAGOGIQUES

- Support de cours ;
- Les Travaux Pratiques se feront en local et dans le Cloud de AWS via le service Amazon Elastic Map Reduce (EMR) ;
- Le langage de programmation utilisé sera Python.