



COLLÈGE
ROSEMONT

Valorisation de données

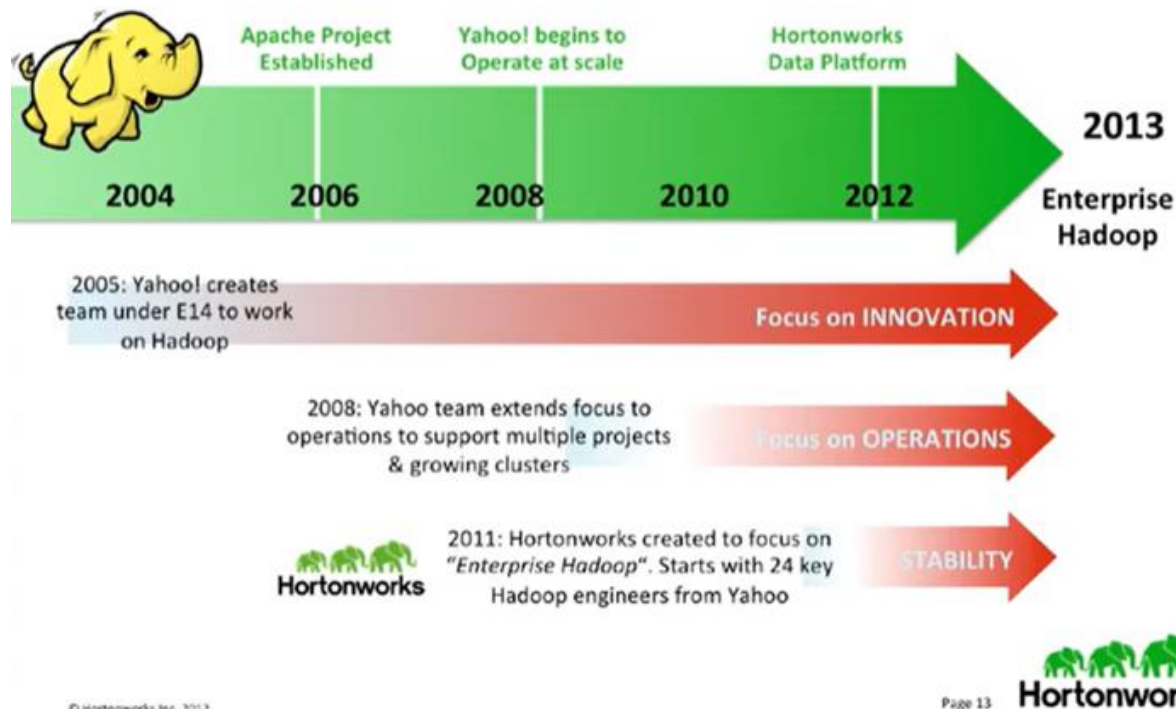
Par : Abderrazak Sahraoui

Sommaire Plateformes de Big Data

- Plateforme HortonWorks Data Platform (HDP)
- Écosystème Hadoop
- Hadoop vs Bases de données relationnelles
- Architecture de plateforme d'entreprise
- Installation HDP
- Définitions : Ambari, Hcatalog, Hive, Pig, Sparks, Sqoop, Flume, Hbase, Oozie, ...
- Autres Plateformes Cloud :
 - Microsoft Azure – HDInsight
 - Amazone - AWS

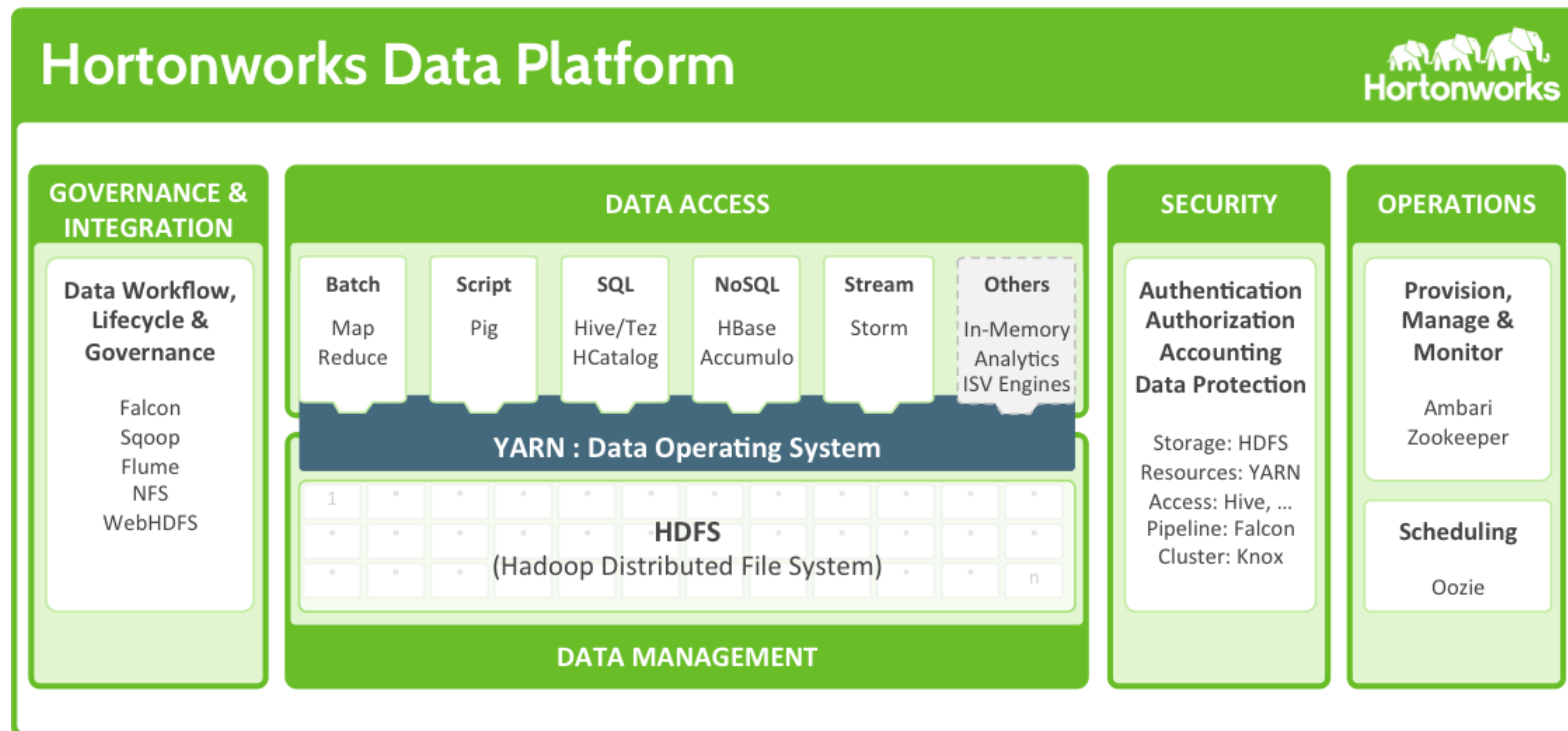
Hortonworks Data Platform (HDP)

A Brief History of Apache Hadoop



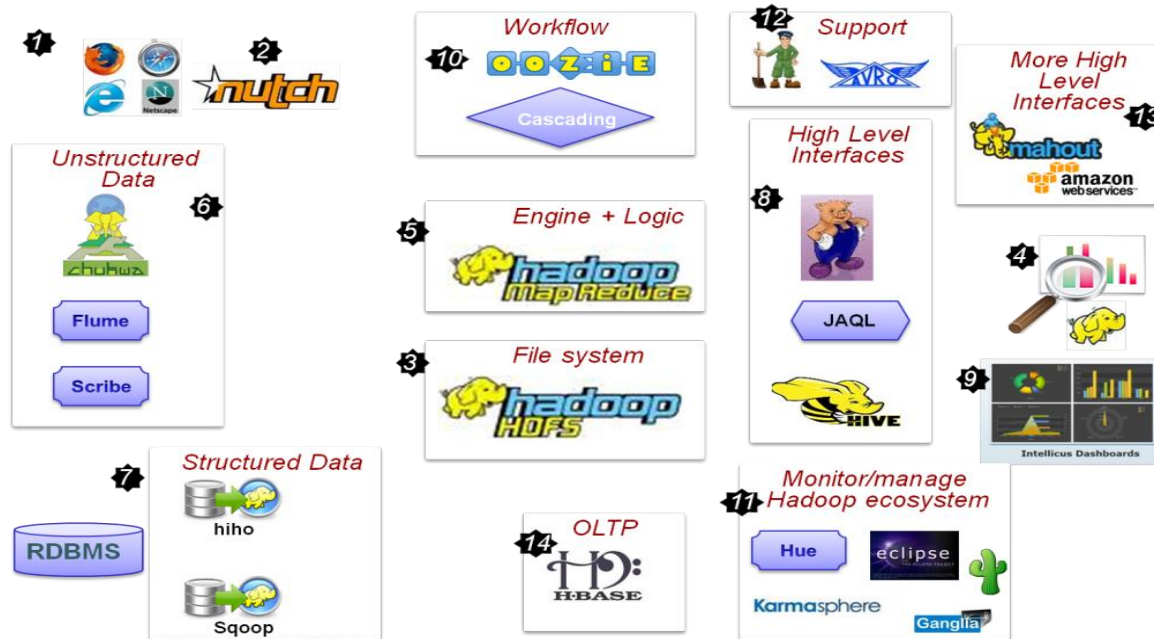
- Lancement d'une plateforme Hadoop pour entreprise open source en 2012 par Hortonworks

(HDP) Hortonworks Data Platform

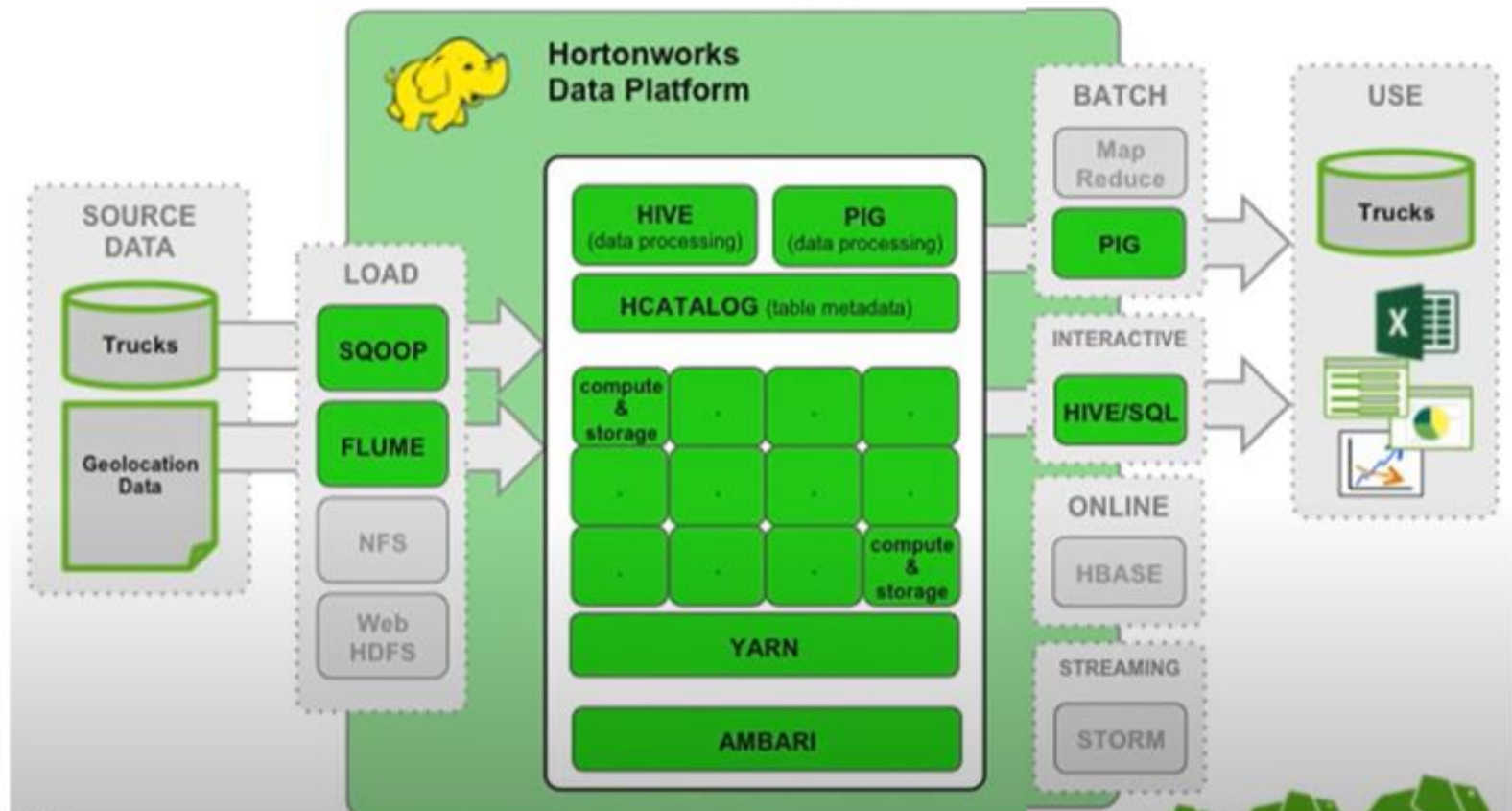


Une autre vision de l'écosystème de Hadoop

Hadoop Ecosystem Map

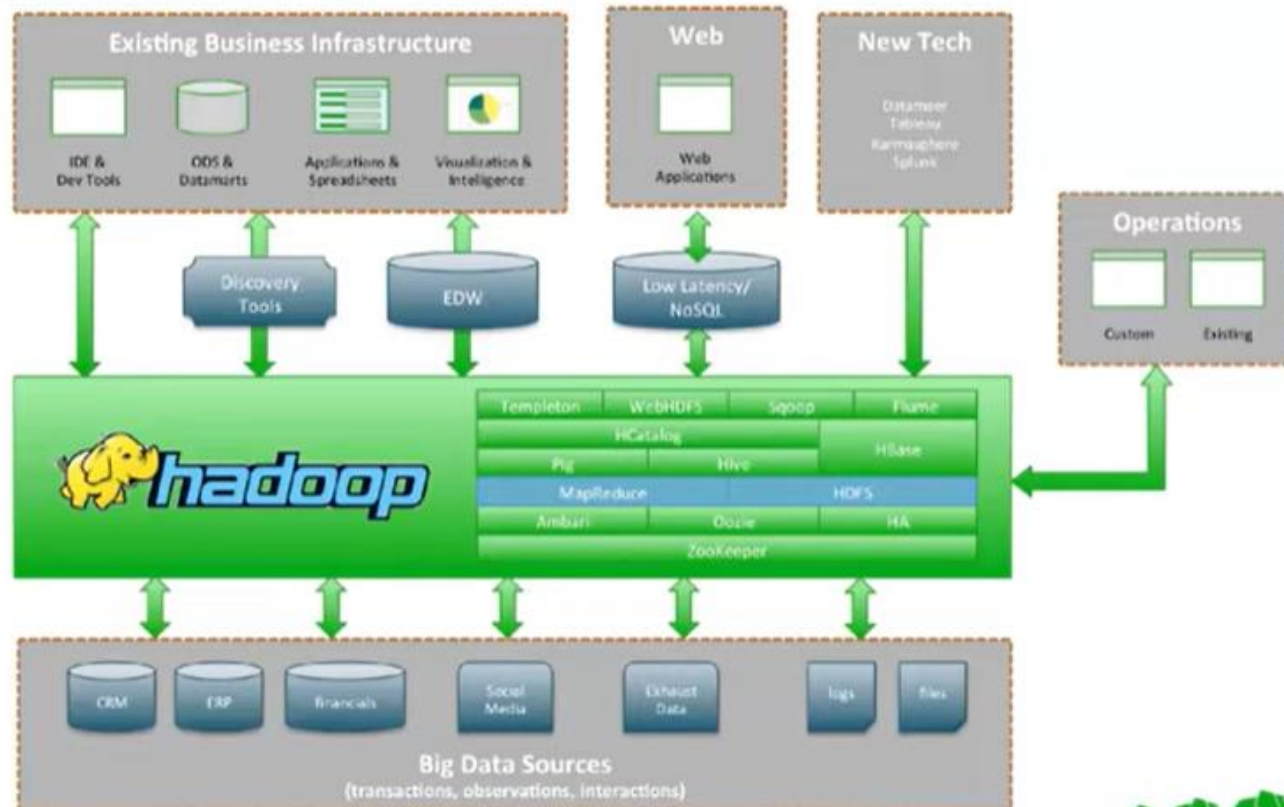


Ecosystème d'une plateforme Hadoop



Architecture d'entreprise avec Hadoop

Hadoop in Enterprise Data Architectures



Installation HDP

- Télécharger et installer Oracle Virtual Box
 - <https://www.virtualbox.org/>
- Exécuter Virtual Box
- Télécharger et installer Hortonworks HDP pour Virtual Box
 - <https://www.cloudera.com/downloads/hortonworks-sandbox.html>
- Lancer HDP sur Virtual Box
- Ouvrir dans un navigateur web la page : 127.0.0.1:8080

https://www.virtualbox.org/



The screenshot shows the VirtualBox website homepage. At the top, there's a navigation bar with various links. The main header features the VirtualBox logo and the text "Welcome to VirtualBox.org!". Below this, a paragraph describes VirtualBox as a powerful x86 and AMD64/Intel64 virtualization product. A large blue banner in the center promotes "Download VirtualBox 7.0". To the left, a sidebar lists links like "About", "Screenshots", "Downloads", "Documentation", "End-user docs", "Technical docs", "Contribute", and "Community". On the right, a "News Flash" section lists recent releases, including VirtualBox 7.0.6, 7.0.4, 7.0.2, 6.1.40, 6.1.0, and 6.1.38. At the bottom, there's a footer with the Oracle logo and links for "Contact", "Privacy policy", and "Terms of Use".

VirtualBox
Welcome to VirtualBox.org!

VirtualBox is a powerful x86 and AMD64/Intel64 virtualization product for enterprise as well as home use. Not only is VirtualBox an extremely feature rich, high performance product for enterprise customers, it is also the only professional solution that is freely available as Open Source Software under the terms of the GNU General Public License (GPL) version 3. See "About VirtualBox" for an introduction.

Presently, VirtualBox runs on Windows, Linux, macOS, and Solaris hosts and supports a large number of guest operating systems including but not limited to Windows (NT 4.0, 2000, XP, Server 2003, Vista, Windows 7, Windows 8, Windows 10), DOS/Windows 3.x, Linux (2.4, 2.6, 3.x and 4.x), Solaris and OpenSolaris, OS/2, and OpenBSD.

VirtualBox is being actively developed with frequent releases and has an ever growing list of features, supported guest operating systems and platforms it runs on. VirtualBox is a community effort backed by a dedicated company: everyone is encouraged to contribute while Oracle ensures the product always meets professional quality criteria.

Download VirtualBox 7.0

Hot picks:

- Pre-built virtual machines for developers at [Oracle Tech Network](#)
- Hyperbox** Open-source Virtual Infrastructure Manager [project site](#)
- phpVirtualBox** AJAX web interface [project site](#)

News Flash

- New January 17th, 2023 VirtualBox 7.0.6 released!**
Oracle today released a 7.0 maintenance release which improves stability and fixes regressions. See the [Changelog](#) for details.
- New January 17th, 2023 VirtualBox 6.1.42 released!**
Oracle today released a 6.1 maintenance release which improves stability and fixes regressions. See the [Changelog](#) for details.
- New November 18th, 2022 VirtualBox 7.0.4 released!**
Oracle today released a 7.0 maintenance release which improves stability and fixes regressions. See the [Changelog](#) for details.
- New October 20th, 2022 VirtualBox 7.0.2 released!**
Oracle today released a 7.0 maintenance release which improves stability and fixes regressions. See the [Changelog](#) for details.
- New October 11th, 2022 VirtualBox 6.1.40 released!**
Oracle today released a 6.1 maintenance release which improves stability and fixes regressions. See the [Changelog](#) for details.
- New October 10th, 2022 VirtualBox 7.0.0 released!**
Oracle today released a significant new version of Oracle VM VirtualBox, its high performance, cross-platform virtualization software. [Changelog](#) for details.
- New September 2nd, 2022 VirtualBox 6.1.38 released!**
Oracle today released a 6.1

ORACLE
Contact - Privacy policy - Terms of Use

<https://www.virtualbox.org/wiki/Downloads>

Get Started with Hortonworks Sandbox

Hortonworks Sandbox can help you get started learning, developing, testing and trying out new features on HDP and DataFlow.

Hortonworks HDP

The HDP Sandbox makes it easy to get started with Apache Hadoop, Apache Spark, Apache Hive, Apache HBase, Druid and Data Analytics Studio (DAS).

[Download Now](#)

Cloudera DataFlow (Ambari)

The Cloudera DataFlow (Ambari)—formerly known as Hortonworks DataFlow—Sandbox makes it easy to get started with Apache NiFi, Apache Kafka, Apache Storm, and Streaming Analytics Manager (SAM).

[Download Now](#)

Machine Virtuelle

Oracle VM VirtualBox - Gestionnaire de machines

Fichier

Machine

Aide

Outils

Nouvelle

Ajouter

Configuration

Oublier

Afficher

Hortonworks Sandbox HD...

En fonction

vm

Éteinte

Général

Nom :

Hortonworks Sandbox HDP 2.6.5

Système d'exploitation :

Red Hat (64-bit)

System

Mémoire vive :

8192 Mo

Processeurs :

4

Ordre d'amorçage :

Disque dur, Optique

Accélération :

Pagination imbriquée, PAE/NX ,
Paravirtualisation KVM

Prévisualisation

```

216231 1676331 R03 grace-period kthreod stack dump:
216231 1776371 R03 rca_sched detected stall on CPU=task:
216231 1776371 30  7  0 ticks this dP, idle=0x00000000, softirq=00000000,
Ap=4
216231 1780241 R03 rca_sched self-detected stall on CPU
216231 1780241 (detected by 2, t=2333376, jiffies: g=2548832, c=2548832, q=8)
216231 1780241 30  7  0 ticks this dP, idle=0x00000000, softirq=00000000,
Ap=4
216231 1780841 R03
216231 1780841 rca_sched kthreod started for 2335152 jiffies g=2548832 c=2548832
R03 R03 GP soft_FPU(C1) rotate=0x00000000
216231 1780841 R03 grace-period kthreod stack dump:
216231 1780841 R03 rca_sched kthreod started for 2335152 jiffies g=2548832 c=2548832
R03 R03 GP soft_FPU(C1) rotate=0x00000000
216231 1780841 R03 rca_sched self-detected stall on CPU=task:
216231 1780841 R03 grace-period kthreod stack dump:
216231 1780841 R03
216231 1780841 All QoS ones, last rca_sched kthreod activity 2335154 (4513244)
216231 1780841 jiffies: (111) next dP=2, next rca_sched R03
216231 1780841 rca_sched kthreod started for 2335176 jiffies g=2548832 c=2548832
R03 R03 GP soft_FPU(C1) rotate=0x00000000
216231 1780841 R03 grace-period kthreod stack dump:
216231 1780841 R03
216231 1780841 system(111) system(111) service watchdog timeout (limit 3min)

```

Affichage

Mémoire vidéo :

8 Mo

Contrôleur graphique :

VBoxVGA

Port serveur bureau distant :

5905

Enregistrement :

Désactivé

Stockage

Contrôleur :

IDE Controller

Maître primaire IDE :

Hortonworks Sandbox HDP 2.6.5-disk001.vdi (Normal, 117,19 Gio)

Audio

Pilote hôte :

Windows Audio Session

Contrôleur :

ICH AC97

Réseau

Interface 1 :

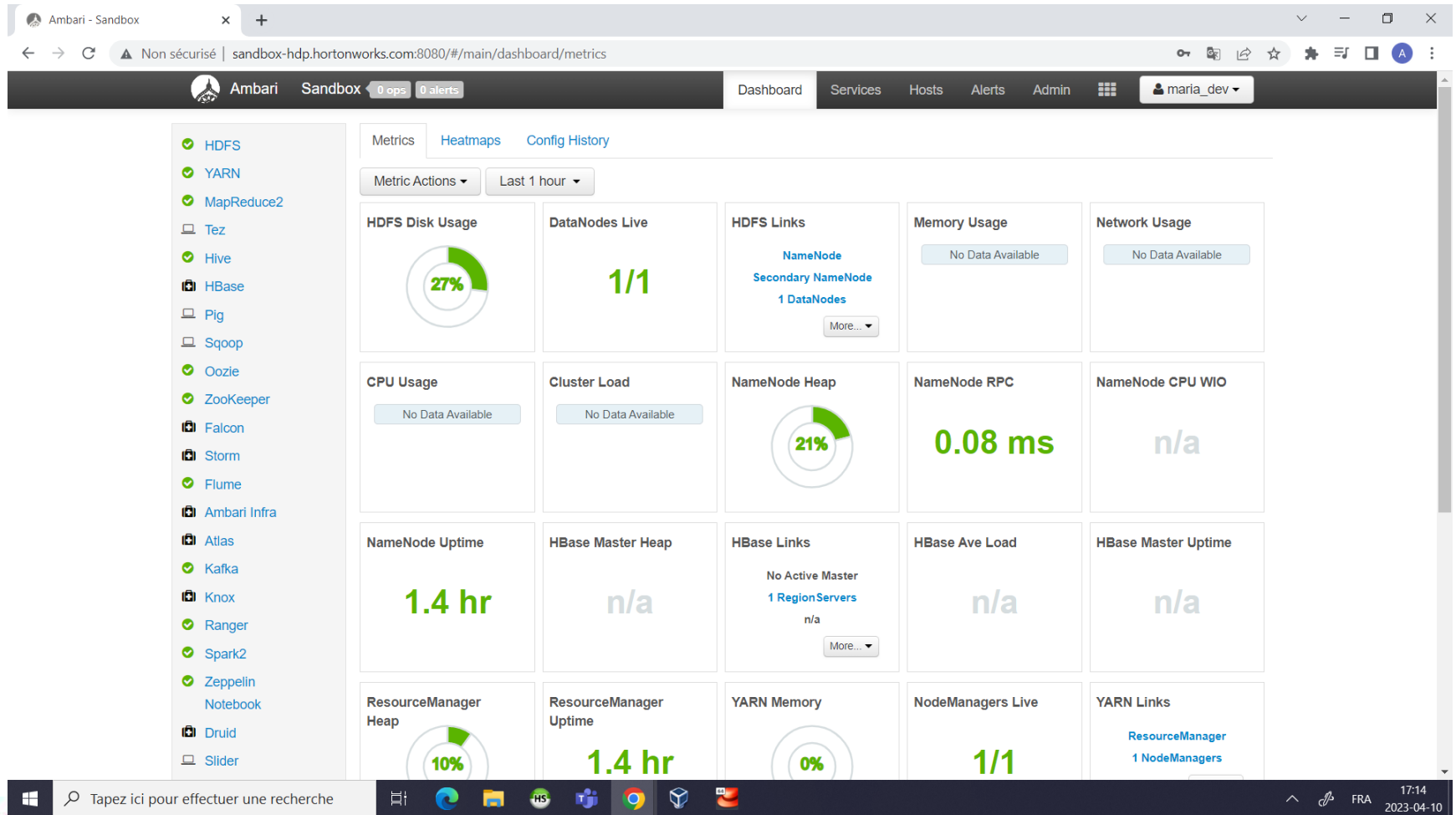
Intel PRO/1000 MT Desktop (NAT)

USB

127.0.0.1:8080

- Le tableau de bord Ambari s'exécute sur le port :8080 . Par exemple, <http://sandbox-hdp.hortonworks.com:8080>

Administration HDP par Ambari



Installation Putty

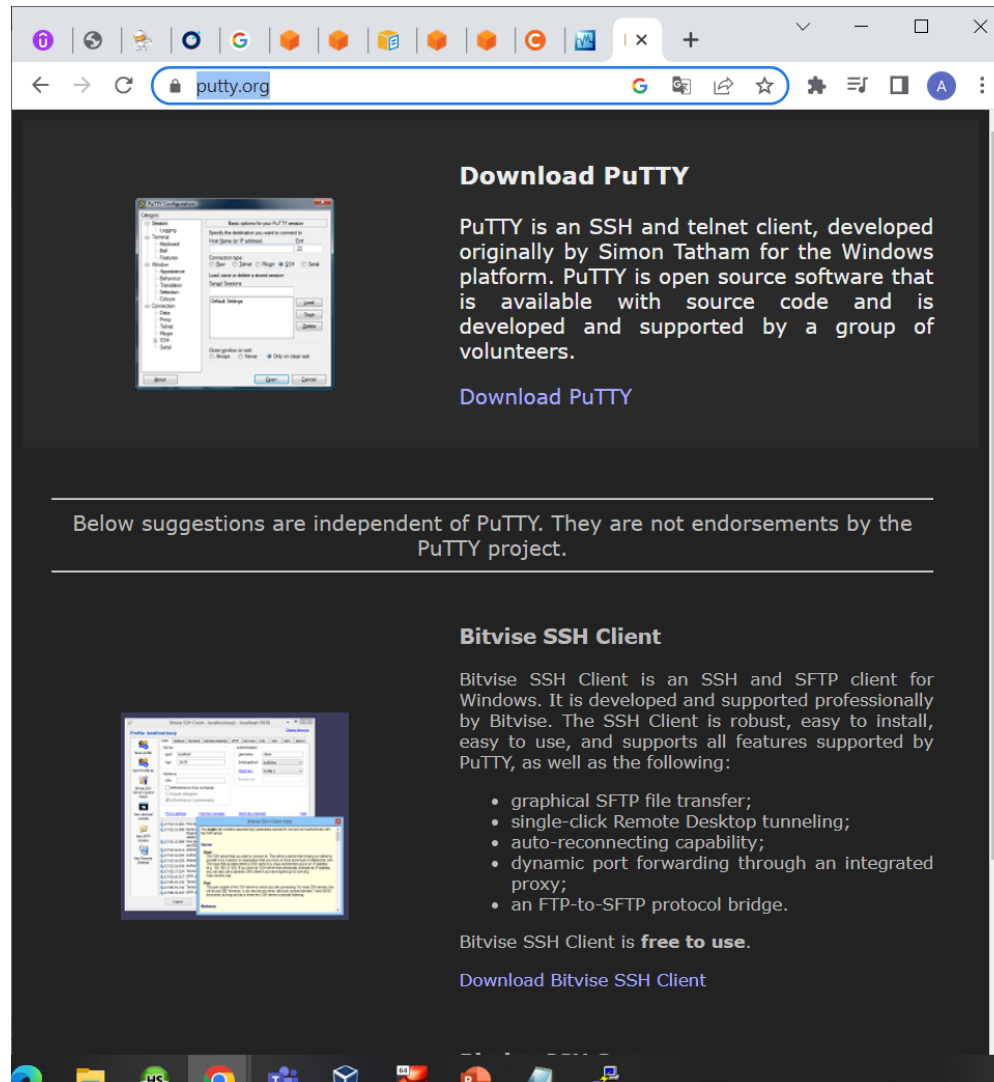
- Télécharger et installer Putty
 - <https://www.putty.org/>
- Exécuter Putty
- Configurer la ligne HDP

Host <u>N</u> ame (or IP address)	<u>P</u> ort
maria_dev@127.0.0.1	2222

Saved Sessions	
hdp	
Default Settings	
hdp	<input type="button" value="Load"/>
	<input type="button" value="Save"/>

- Cliquer sur Open

<https://www.putty.org/>



Download PuTTY

PuTTY is an SSH and telnet client, developed originally by Simon Tatham for the Windows platform. PuTTY is open source software that is available with source code and is developed and supported by a group of volunteers.

[Download PuTTY](#)

Below suggestions are independent of PuTTY. They are not endorsements by the PuTTY project.

Bitvise SSH Client

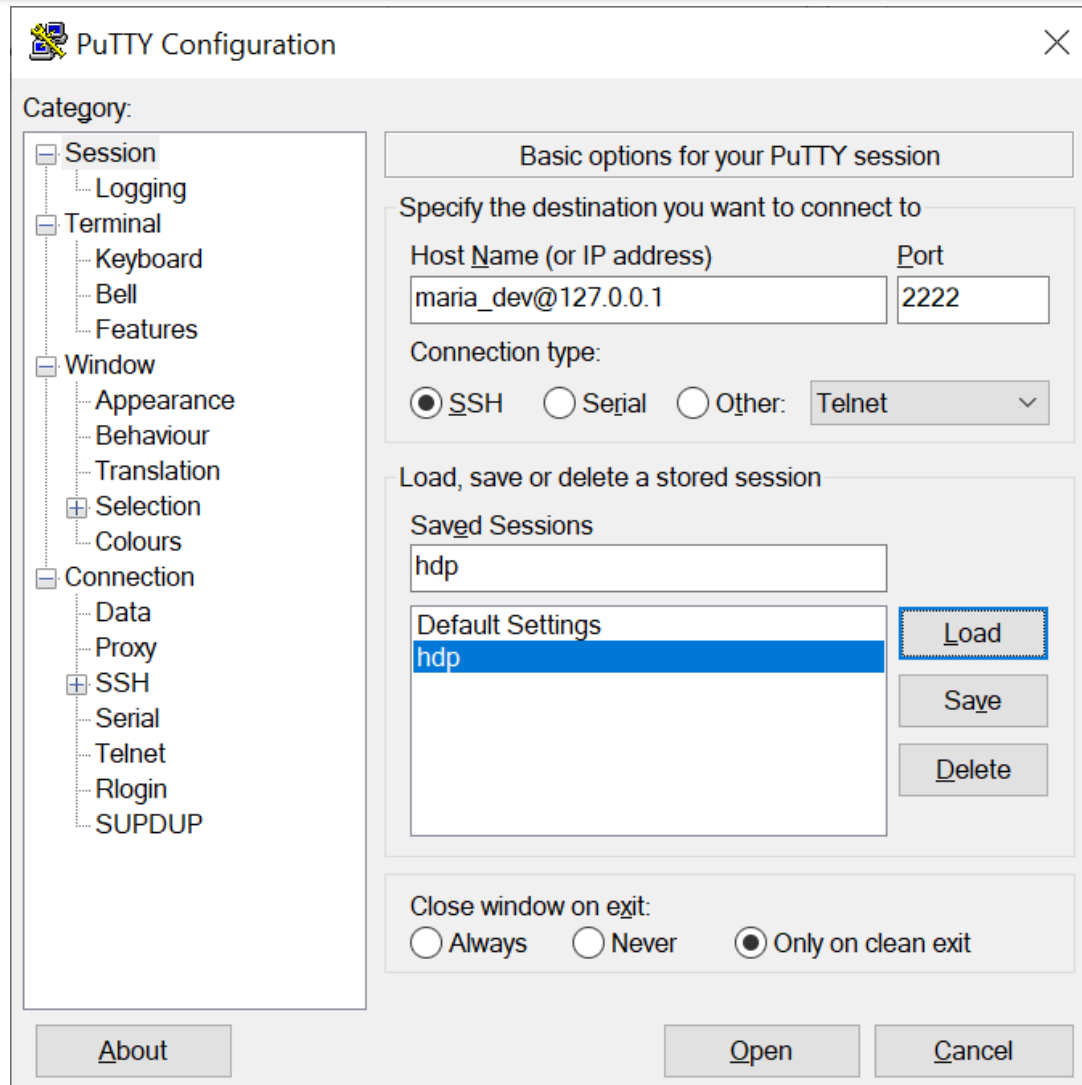
Bitvise SSH Client is an SSH and SFTP client for Windows. It is developed and supported professionally by Bitvise. The SSH Client is robust, easy to install, easy to use, and supports all features supported by PuTTY, as well as the following:

- graphical SFTP file transfer;
- single-click Remote Desktop tunneling;
- auto-reconnecting capability;
- dynamic port forwarding through an integrated proxy;
- an FTP-to-SFTP protocol bridge.

Bitvise SSH Client is **free to use**.

[Download Bitvise SSH Client](#)

Commande en ligne



The image shows the PuTTY Configuration dialog box. On the left is a tree view under 'Category:' with sub-items: Session, Logging, Terminal, Keyboard, Bell, Features, Window, Appearance, Behaviour, Translation, Selection, Colours, Connection, Data, Proxy, SSH, Serial, Telnet, Rlogin, and SUPDUP. The 'SSH' category is expanded. The main area is titled 'Basic options for your PuTTY session'. It contains a section 'Specify the destination you want to connect to' with a 'Host Name (or IP address)' field containing 'maria_dev@127.0.0.1' and a 'Port' field containing '2222'. Below this is a 'Connection type:' section with radio buttons for 'SSH' (selected), 'Serial', and 'Other:', followed by a dropdown menu showing 'Telnet'. Another section 'Load, save or delete a stored session' contains a 'Saved Sessions' list with 'hdp' and 'Default Settings' (highlighted in blue). To the right of this list are 'Load', 'Save', and 'Delete' buttons. At the bottom is a 'Close window on exit:' section with radio buttons for 'Always', 'Never', and 'Only on clean exit' (selected). At the very bottom are 'About', 'Open', and 'Cancel' buttons.

PuTTY Configuration

Category:

- Session
 - Logging
- Terminal
 - Keyboard
 - Bell
 - Features
- Window
 - Appearance
 - Behaviour
 - Translation
 - Selection
 - Colours
- Connection
 - Data
 - Proxy
 - SSH
 - Serial
 - Telnet
 - Rlogin
 - SUPDUP

Basic options for your PuTTY session

Specify the destination you want to connect to

Host Name (or IP address) Port

maria_dev@127.0.0.1 2222

Connection type:

☒ SSH ☐ Serial ☐ Other: Telnet

Load, save or delete a stored session

Saved Sessions

hdp

Default Settings

hdp

Load

Save

Delete

Close window on exit:

☐ Always ☐ Never ☒ Only on clean exit

About Open Cancel

Exécution de Putty

```

maria_dev@sandbox-hdp:~
Using username "maria_dev".
maria_dev@127.0.0.1's password:
[maria_dev@sandbox-hdp ~]$ hadoop fs -ls
Found 2 items
drwxr-xr-x  - maria_dev hdfs          0 2023-04-07 21:37 .Trash
drwxr-xr-x  - maria_dev hdfs          0 2023-03-10 00:11 hive
[maria_dev@sandbox-hdp ~]$ ls
[maria_dev@sandbox-hdp ~]$ wget http://media.sundog-soft.com/hadoop/ml-100k/u.data
--2023-04-10 20:35:48--  http://media.sundog-soft.com/hadoop/ml-100k/u.data
Resolving media.sundog-soft.com (media.sundog-soft.com)... 52.217.196.153, 52.217.192.41, 54.
231.228.81, ...
Connecting to media.sundog-soft.com (media.sundog-soft.com)|52.217.196.153|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 2079229 (2.0M) [application/octet-stream]
Saving to: 'u.data'

100%[=====>] 2,079,229    5.03MB/s   in 0.4s

2023-04-10 20:35:48 (5.03 MB/s) - 'u.data' saved [2079229/2079229]

[maria_dev@sandbox-hdp ~]$
```

Gestion des données sur HDP

- Data Management – Store and process vast quantities of data in a storage layer that scales linearly.
- Hadoop Distributed File System (HDFS) is the core technology for the efficient scale out storage layer, and is designed to run across low-cost commodity hardware
- Apache Hadoop YARN is the pre-requisite for Enterprise Hadoop as it provides the resource management and pluggable architecture for enabling a wide variety of data access methods to operate on data stored in Hadoop with predictable performance and service levels.

Accès aux données sur HDP

- Data Access – Interact with your data in a wide variety of ways – from batch to real-time.
- Apache Hive is the most widely adopted data access technology
- Apache Pig provides scripting capabilities, Apache Storm offers real-time processing, Apache HBase offers columnar NoSQL storage and Apache Accumulo offers cell-level access control.

Accès aux données sur HDP

- [Apache Hive](#) – Built on the MapReduce framework, Hive is a data warehouse that enables easy data summarization and ad-hoc queries via an SQL-like interface for large datasets stored in HDFS.
- [Apache Pig](#) – A platform for processing and analyzing large data sets. Pig consists of a high-level language (Pig Latin) for expressing data analysis programs paired with the MapReduce framework for processing these programs.
- [MapReduce](#) – MapReduce is a framework for writing applications that process large amounts of structured and unstructured data in parallel across a cluster of thousands of machines, in a reliable and fault-tolerant manner.
- [Apache Spark](#) – Spark is ideal for in-memory data processing. It allows data scientists to implement fast, iterative algorithms for advanced analytics such as clustering and classification of datasets.

Accès aux données sur HDP

- [Apache Storm](#) – Storm is a distributed real-time computation system for processing fast, large streams of data adding reliable real-time data processing capabilities to Apache Hadoop 2.x
- [Apache HBase](#) – A column-oriented NoSQL data storage system that provides random real-time read/write access to big data for user applications.
- [Apache Tez](#) – Tez generalizes the MapReduce paradigm to a more powerful framework for executing a complex DAG (directed acyclic graph) of tasks for near real-time big data processing.
- [Apache Kafka](#) – Kafka is a fast and scalable publish-subscribe messaging system that is often used in place of traditional message brokers because of its higher throughput, replication, and fault tolerance.
- [Apache HCatalog](#) – A table and metadata management service that provides a centralized way for data processing systems to understand the structure and location of the data stored within Apache Hadoop.

Accès aux données sur HDP

- [Apache Slider](#) – A framework for deployment of long-running data access applications in Hadoop. Slider leverages YARN's resource management capabilities to deploy those applications, to manage their lifecycles and scale them up or down.
- [Apache Solr](#) – Solr is the open source platform for searches of data stored in Hadoop. Solr enables powerful full-text search and near real-time indexing on many of the world's largest Internet sites.
- [Apache Mahout](#) – Mahout provides scalable machine learning algorithms for Hadoop which aids with data science for clustering, classification and batch based collaborative filtering.
- [Apache Accumulo](#) – Accumulo is a high performance data storage and retrieval system with cell-level access control. It is a scalable implementation of Google's Big Table design that works on top of Apache Hadoop and Apache ZooKeeper.

Gouvernance et Intégration de données

- Data Governance and Integration – Quickly and easily load data, and manage according to policy. Workflow Manager provides workflows for data governance, while Apache Flume and Sqoop enable easy data ingestion, as do the NFS and WebHDFS interfaces to HDFS.

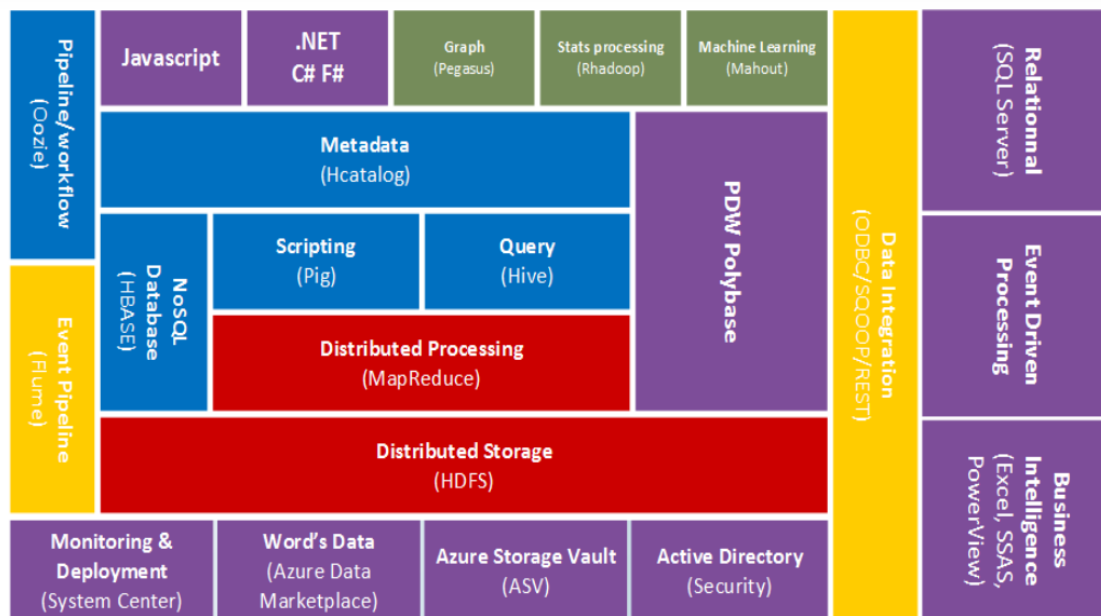
- Workflow Management – Workflow Manager allows you to easily create and schedule workflows and monitor workflow jobs. It is based on the Apache Oozie workflow engine that allows users to connect and automate the execution of big data processing tasks into a defined workflow.
- Apache Flume – Flume allows you to efficiently aggregate and move large amounts of log data from many different sources to Hadoop.
- Apache Sqoop – Sqoop is a tool that speeds and eases movement of data in and out of Hadoop. It provides a reliable parallel load for various, popular enterprise data sources.

Gouvernance et Intégration de données

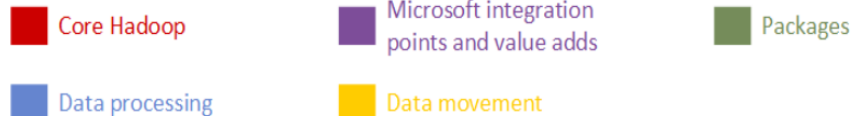
- Security – Address requirements of Authentication, Authorization, Accounting and Data Protection. Security is provided at every layer of the Hadoop stack from HDFS and YARN to Hive and the other Data Access components on up through the entire perimeter of the cluster via Apache Knox.
- Apache Knox – The Knox Gateway (“Knox”) provides a single point of authentication and access for Apache Hadoop services in a cluster. The goal of the project is to simplify Hadoop security for users who access the cluster data and execute jobs, and for operators who control access to the cluster.
- Apache Ranger – Apache Ranger delivers a comprehensive approach to security for a Hadoop cluster. It provides central security policy administration across the core enterprise security requirements of authorization, accounting and data protection.

- Operations – Provision, manage, monitor and operate Hadoop clusters at scale.
 - [Apache Ambari](#) – An open source installation lifecycle management, administration and monitoring system for Apache Hadoop clusters.
 - [Apache Oozie](#) – Oozie Java Web application used to schedule Apache Hadoop jobs. Oozie combines multiple jobs sequentially into one logical unit of work.
 - [Apache ZooKeeper](#) – A highly available system for coordinating distributed processes. Distributed applications use ZooKeeper to store and mediate updates to important configuration information.

Plateforme Microsoft HDInsight



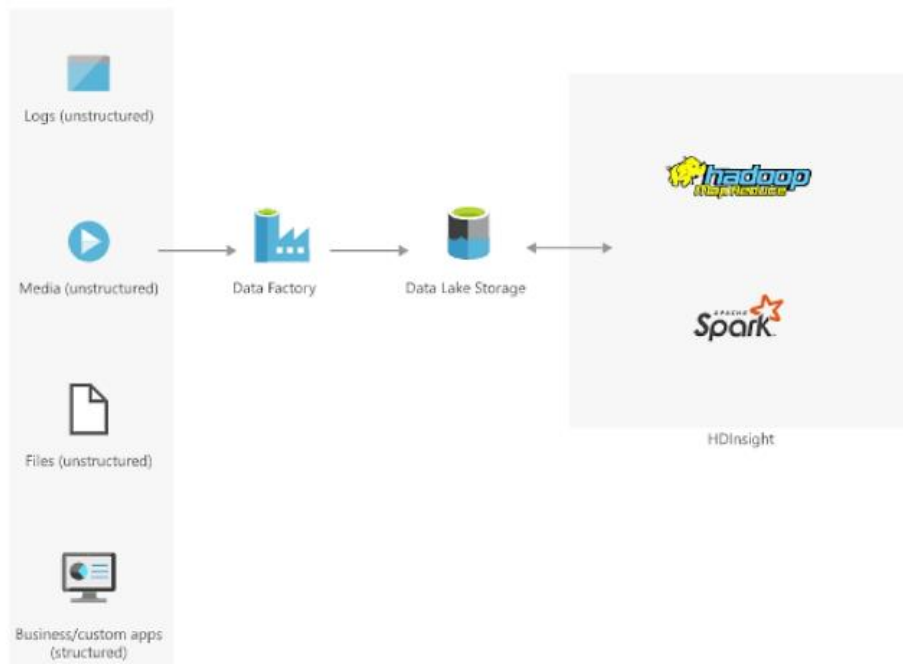
Legend



- HDInsight est basée sur la plateforme Hortonworks
- Les briques violettes dans l'écosystème HDInsight sont les composants ajoutés par Microsoft au produit Hortonworks.

<https://blog.octo.com/hdinsight-le-big-data-selon-microsoft/>

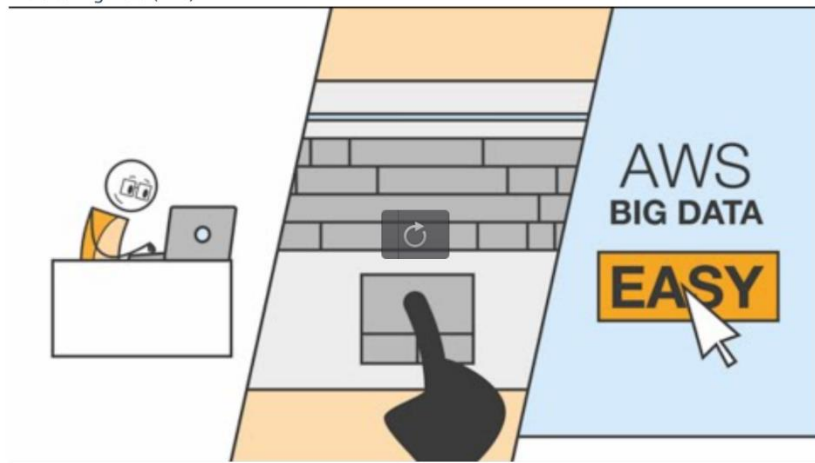
Big data sur Azure



<https://azure.microsoft.com/fr-ca/products/hdinsight>

- Lac de données grâce à l'intégration des solutions et services de stockage de données Azure.
- Mise à l'échelle automatique en fonction de la charge ou de la planification.
- Tableaux de bord pour surveiller l'intégralité de lac de données à l'aide des

Big data sur Amazon Web Services and Cloud (AWS)



<https://aws.amazon.com/fr/big-data/what-is-big-data/>

- AWS propose des services dédiés : mouvement des données, stockage des données, lacs de données, analyse des données du big data, analyse des journaux, analyse de streaming, informatique décisionnelle (BI), Apprentissage (machine learning ML).
- Offre clusters et écosystèmes de big data composés de Hadoop, Hive, Pig, Sparks, etc...
- Permet à travers les services EC2 de monter sa propre solution big data.

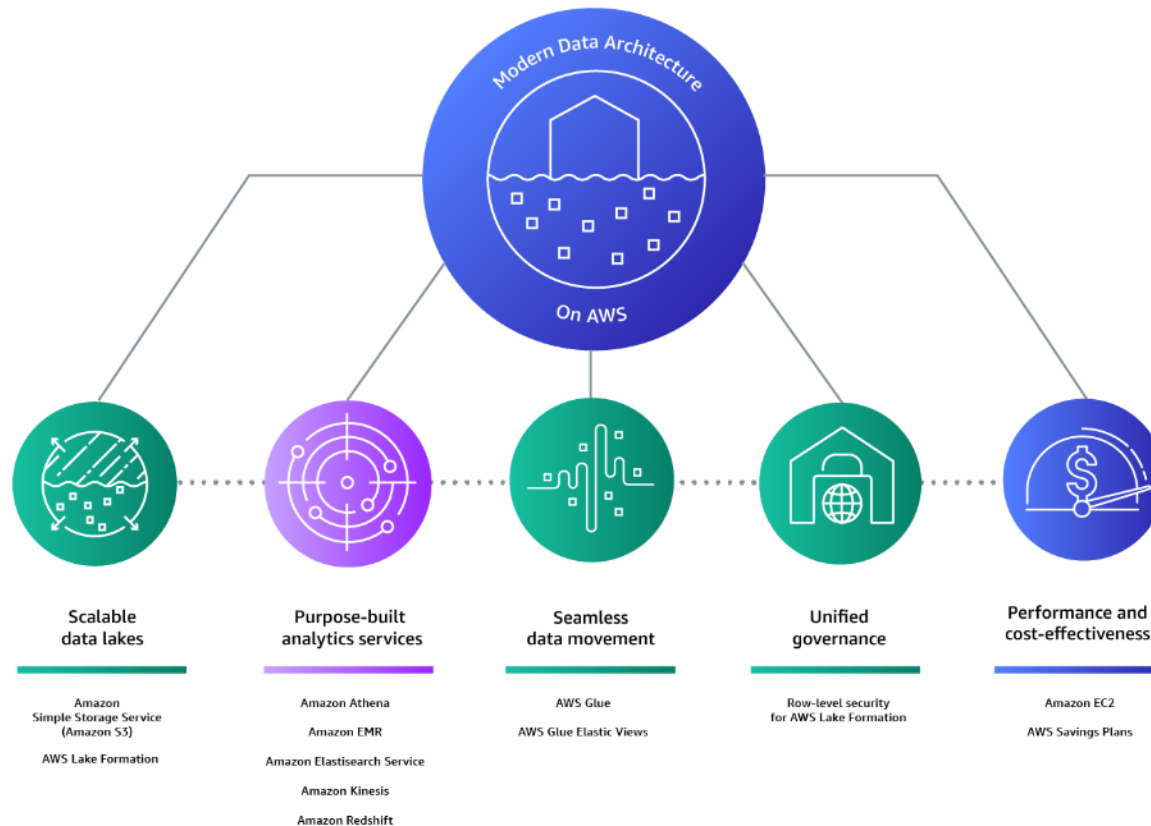
Architecture sur AWS



- Les organisations prennent leurs données stockées dans des silos et les déplacent dans un seul endroit pour les analyser et les utiliser pour des tâches de ML (machine learning). Pour réaliser cette opération de façon optimale, elles ont besoin d'utiliser une architecture de données moderne leur permettant de déplacer facilement des données entre des lacs et des magasins de données spécialisés.

<https://aws.amazon.com/fr/big-data/what-is-big-data/>

Architecture sur AWS



- Intégrer un lac de données, un entrepôt de données et des magasins spécialisés, afin d'unifier la gouvernance et de faciliter le mouvement des données.

<https://aws.amazon.com/fr/big-data/what-is-big-data/>

Partenaires sur AWS



Cloudera

L'exécution de Cloudera Enterprise sur AWS fournit aux utilisateurs informatiques et professionnels une plateforme de gestion des données qui peut servir de base au traitement et à l'analyse modernes des données.

[En savoir plus »](#)



Informatica Cloud

Informatica Cloud offre une intégration optimisée aux services de données AWS avec une connectivité native à plus de 100 applications.

[En savoir plus »](#)



Dataguise

Dataguise est le leader de l'exécution métier sécurisée, fournissant des solutions de sécurité centrées sur les données qui détectent et protègent les données sensibles d'une entreprise, peu importe où elles se trouvent ou qui a besoin de les exploiter.

[En savoir plus »](#)



Alluxio Data Orchestration

Alluxio Data Orchestration permet aux clients de mieux exploiter les principaux services AWS, tels que EMR et S3 pour les charges de travail d'analytique et d'IA.

[En savoir plus »](#)

Amazon EC2

What is Amazon EC2?

Amazon Elastic Compute Cloud (Amazon EC2) provides scalable computing capacity in the Amazon Web Services (AWS) Cloud. Using Amazon EC2 eliminates your need to invest in hardware up front, so you can develop and deploy applications faster. You can use Amazon EC2 to launch as many or as few virtual servers as you need, configure security and networking, and manage storage. Amazon EC2 enables you to scale up or down to handle changes in requirements or spikes in popularity, reducing your need to forecast traffic.

For more information about cloud computing, see [What is cloud computing?](#)

<https://docs.aws.amazon.com/AWSEC2/latest/WindowsGuide/concepts.html>

- Monter sa propre solution big data avec EC2