

Architecture et Technologies Big Data

Mini-Projet Big Data : Installation et Analyse sur un Cluster

Hadoop/Spark



Réalisé par :

GAROUD Fatima Ezzahraa

BOUASSAB Chaimae

ABABRI Chaimae

Encadré par :

Pr. Anouar Abdelhakim BOUDHIR

15/06/2025

Table des matières

Introduction	4
Rôle des machines virtuelles	4
Interaction entre les VM :	6
Phase 1 – Prise en main et installation	6
Des étapes sur VM ubuntu-masterfch :	7
Étape 1 : Création de l'utilisateur hadoopuser :	7
1. Connection à la VM	7
2. Création de l'utilisateur hadoopuser	7
3. Configurer SSH pour hadoopuser	8
4. Vérifier les permissions	10
Étape 2 : Configuration réseau	10
2.2 Configurer l'IP statique	11
2.3 Tester la connectivité	13
Étape 3 : Configurer le pare-feu	13
Étape 4 : Installer Java	15
Étape 5 : Installer et configurer Hadoop	16
5.1 Installer Hadoop	16
5.2 Configurer Hadoop	17
5.3 Formater le NameNode	19
Étape 6 : Installer et configurer Spark	19
6.1 Installer Spark	19
6.2 Configurer Spark	20
Des étapes sur VM ubuntu-workerfch :	22
1. Installation de la distribution Ubuntu Server	22
2. Configuration système initiale	22
3. Installation de Java	24
4. Installation de Hadoop et Spark	24
5. Configuration des fichiers Hadoop	25
Phase 2 – Manipulation et stockage	26
Étape 1 : Formatage de hdfs :	26
Étape 2 : Tester les services	27
Étape 3 : Transfert du fichier .csv dans MasterUbuntu	29
Étape 4 : Ingestion dans HDFS	29
Objectif de cette phase :	30

Étape 5 : Traitement des données avec PySpark	30
Architecture distribuée du cluster Hadoop/Spark:	35
Rôles des nœuds dans Notre cluster	37
Phase 3 : MINI-CAS D'ANALYSE : avis produits Amazon	38
Intégration et Visualisation des Données via Grafana Conteneurisé	38
Lancement de Grafana via Docker	38
CONCLUSION	45

Introduction

Dans le cadre du mini-projet Big Data, ce travail vise à mettre en place un environnement Big Data fonctionnel basé sur un cluster Hadoop/Spark, en utilisant deux machines virtuelles nommées *ubuntumasterfch* et *ubuntuworkerfch*. L'objectif est de configurer un cluster opérationnel, capable de stocker, traiter et analyser de grands volumes de données, tout en simulant un cas d'usage réel. À travers les phases d'installation, de manipulation, d'analyse et de présentation, ce projet permet de maîtriser les technologies Hadoop et Spark, de comprendre leur intégration dans un environnement distribué, et de développer des compétences en analyse de données et en communication technique. En utilisant *ubuntumasterfch* comme nœud maître et *ubuntuworkerfch* comme nœud travailleur, ce projet illustre le fonctionnement d'un cluster Big Data dans un contexte pratique.

Rôle des machines virtuelles

Le cluster Hadoop/Spark est configuré avec deux machines virtuelles, chacune ayant un rôle spécifique pour assurer le fonctionnement distribué du système :

1. *ubuntumasterfch* (Nœud Maître) :
 - Rôle principal : Cette machine virtuelle agit comme le nœud maître (Master Node) du cluster Hadoop/Spark. Elle est responsable de la coordination et de la gestion des opérations du cluster.
 - Composants installés :
 - Hadoop NameNode : Gère les métadonnées du système de fichiers distribué HDFS, stockant les informations sur la localisation des données à travers le cluster.
 - Hadoop ResourceManager : Coordonne les ressources du cluster (CPU, mémoire) pour les tâches distribuées via YARN (Yet Another Resource Negotiator).
 - Spark Master : Supervise l'exécution des applications Spark, en répartissant les tâches entre les nœuds travailleurs.

- Responsabilités :
 - Orchestrer les opérations du cluster, comme la gestion des fichiers dans HDFS et la planification des tâches Spark.
 - Héberger les services de gestion (NameNode, ResourceManager, Spark Master).
 - Maintenir la communication avec le nœud travailleur pour assurer la cohérence des opérations.
- Configuration : Cette VM nécessite une configuration réseau stable (SSH, hostname, DNS local), un accès root sécurisé, et les variables d'environnement pour Hadoop et Spark correctement définies.

2. ubuntuworkerfch (Nœud Travailleur) :

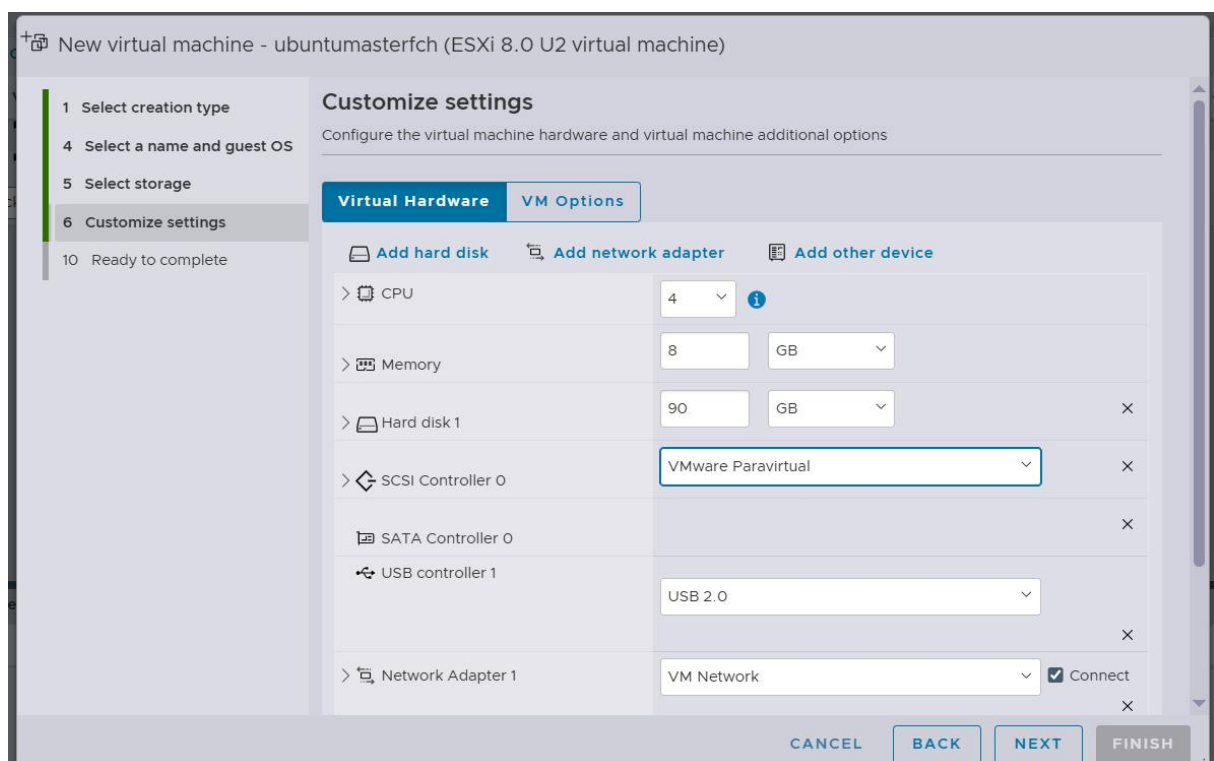
- Rôle principal : Cette machine virtuelle agit comme le nœud travailleur (Worker Node) du cluster. Elle exécute les tâches de calcul et de stockage confiées par le nœud maître.
- Composants installés :
 - Hadoop DataNode : Stocke les blocs de données réels dans HDFS et répond aux requêtes du NameNode pour lire ou écrire des données.
 - Hadoop NodeManager : Gère les ressources locales (CPU, mémoire) pour exécuter les tâches distribuées par le ResourceManager via YARN.
 - Spark Worker : Exécute les tâches Spark assignées par le Spark Master, en effectuant des calculs parallèles sur les données.
- Responsabilités :
 - Stocker les données dans HDFS (sous forme de blocs).
 - Effectuer les calculs distribués pour les jobs Spark, comme le filtrage, l'agrégation ou les jointures.
 - Communiquer avec le nœud maître pour recevoir les instructions et renvoyer les résultats.

- Configuration : Cette VM doit être configurée pour communiquer avec *ubuntumasterfch* (via SSH et configuration réseau), et avoir les mêmes versions de Java, Hadoop et Spark installées pour garantir la compatibilité.

Interaction entre les VM :

- Le nœud maître (*ubuntumasterfch*) contrôle et coordonne les opérations, tandis que le nœud travailleur (*ubuntuworkerfch*) exécute les tâches de stockage et de calcul.
- La communication entre les deux VM repose sur une configuration réseau appropriée (SSH, ports ouverts, DNS local) pour permettre l'échange de données et d'instructions.
- Dans un cluster réel, il pourrait y avoir plusieurs nœuds travailleurs, mais ici, *ubuntuworkerfch* simule un environnement distribué minimal.

Phase 1 – Prise en main et installation



New virtual machine - ubuntu-masterfch (ESXi 8.0 U2 virtual machine)

1 Select creation type
4 Select a name and guest OS
5 Select storage
6 Customize settings
10 Ready to complete

Ready to complete

Review your settings selection before finishing the wizard

Name	ubuntu-masterfch
Datastore	datastore1
Guest OS name	Debian GNU/Linux 12 (64-bit)
Compatibility	ESXi 8.0 U2 virtual machine
vCPUs	4
Memory	8 GB
Network adapters	1
Network adapter 1 network	VM Network
Network adapter 1 type	VMXNET 3
IDE controller 0	IDE 0
IDE controller 1	IDE 1
SCSI controller 0	VMware Paravirtual
SATA controller 0	New SATA controller
Hard disk 1	
Capacity	90 GB

CANCEL
BACK
NEXT
FINISH

Des étapes sur VM ubuntu-masterfch :

<input type="checkbox"/>	ubuntu-masterfch	<input checked="" type="checkbox"/> Normal	94.08 GB	Ubuntu Linux (64-bit)	ubuntu-master	45 MHz	2.3 GB
--------------------------	------------------	--	----------	-----------------------	---------------	--------	--------

Étape 1 : Création de l'utilisateur hadoopuser :

1. Connection à la VM

```
C:\Users\DELL>ssh masterubuntu@192.168.180.49
```

2. Création de l'utilisateur hadoopuser

La création de l'utilisateur avec le répertoire personnel :

```

masterubuntu@masterubuntu:~$ sudo adduser hadoopuser
[sudo] password for masterubuntu:
Adding user `hadoopuser' ...
Adding new group `hadoopuser' (1001) ...
Adding new user `hadoopuser' (1001) with group `hadoopuser' ...
Creating home directory `/home/hadoopuser' ...
Copying files from `/etc/skel' ...
New password:
Retype new password:
passwd: password updated successfully
Changing the user information for hadoopuser
Enter the new value, or press ENTER for the default
    Full Name []:
    Room Number []:
    Work Phone []:
    Home Phone []:
    Other []:
Is the information correct? [Y/n] y

```

Cela crée un répertoire personnel /home/hadoopuser.

- En Donne à hadoopuser des privilèges sudo (pour installer des paquets et gérer les configurations système) :

```

masterubuntu@masterubuntu:~$ sudo usermod -aG sudo hadoopuser

```

Hadoop et Spark nécessitent des commandes sudo pour certaines installations (par exemple, apt install) et configurations de répertoires.

3. Configurer SSH pour hadoopuser

```

masterubuntu@masterubuntu:~$ su - hadoopuser

```

- Génération d'une clé SSH pour permettre la communication sans mot de passe avec le Worker (et pour Hadoop/Spark) :


```
hadoopuser@masterubuntu:~$ ssh-keygen -t rsa -P "" -f ~/.ssh/id_rsa
Generating public/private rsa key pair.
Created directory '/home/hadoopuser/.ssh'.
Your identification has been saved in /home/hadoopuser/.ssh/id_rsa
Your public key has been saved in /home/hadoopuser/.ssh/id_rsa.pub
The key fingerprint is:
SHA256:x8rGrccH5RbcCijg1iTvPZsjB67pwmbaF/Sja09XcQc hadoopuser@masterubuntu
The key's randomart image is:
+---[RSA 3072]---+
|           E           |
|  o .               .  |
| . *               .   |
| + + o..o+..        |
| o o +S ++ o         |
|  o =o =. +          |
| . . +.+*..o         |
| . = .*. =o .o .     |
| =o=+. =. ....      |
+-----[SHA256]-----+
```

Cela crée :

- Clé privée : /home/hadoopuser/.ssh/id_rsa
- Clé publique : /home/hadoopuser/.ssh/id_rsa.pub

➤ Autorisation de la clé publique pour SSH local

```
hadoopuser@masterubuntu:~$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
hadoopuser@masterubuntu:~$ chmod 600 ~/.ssh/authorized_keys
```

➤ Teste de SSH localement

```

hadoopuser@masterubuntu:~$ ssh localhost
The authenticity of host 'localhost (127.0.0.1)' can't be established.
ED25519 key fingerprint is SHA256:uz9IQYRw/fbvS5rU+vU1+b4m9B4WsONwidTGEzCyRck.
This key is not known by any other names
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'localhost' (ED25519) to the list of known hosts.
Welcome to Ubuntu 22.04.5 LTS (GNU/Linux 5.15.0-140-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/pro

System information as of Thu May 29 12:59:57 PM UTC 2025

System load:  0.0               Processes:            203
Usage of /:   16.2% of 42.49GB   Users logged in:     1
Memory usage: 7%               IPv4 address for ens34: 192.168.180.49
Swap usage:   0%

Expanded Security Maintenance for Applications is not enabled.

55 updates can be applied immediately.
To see these additional updates run: apt list --upgradable

Enable ESM Apps to receive additional future security updates.
See https://ubuntu.com/esm or run: sudo pro status

New release '24.04.2 LTS' available.
Run 'do-release-upgrade' to upgrade to it.

```

cela fonctionne sans demander de mot de passe, SSH est correctement configuré pour hadoopuser.

4. Vérifier les permissions

- Pour que hadoopuser a les permissions nécessaires sur son répertoire personnel :

```

hadoopuser@masterubuntu:~$ sudo chown -R hadoopuser:hadoopuser /home/hadoopuser
[sudo] password for hadoopuser:

```

```

hadoopuser@masterubuntu:~$ chmod 700 /home/hadoopuser/.ssh
hadoopuser@masterubuntu:~$

```

Étape 2 : Configuration réseau

2.1 Configurer le nom d'hôte

- Connection en tant que hadoopuser :

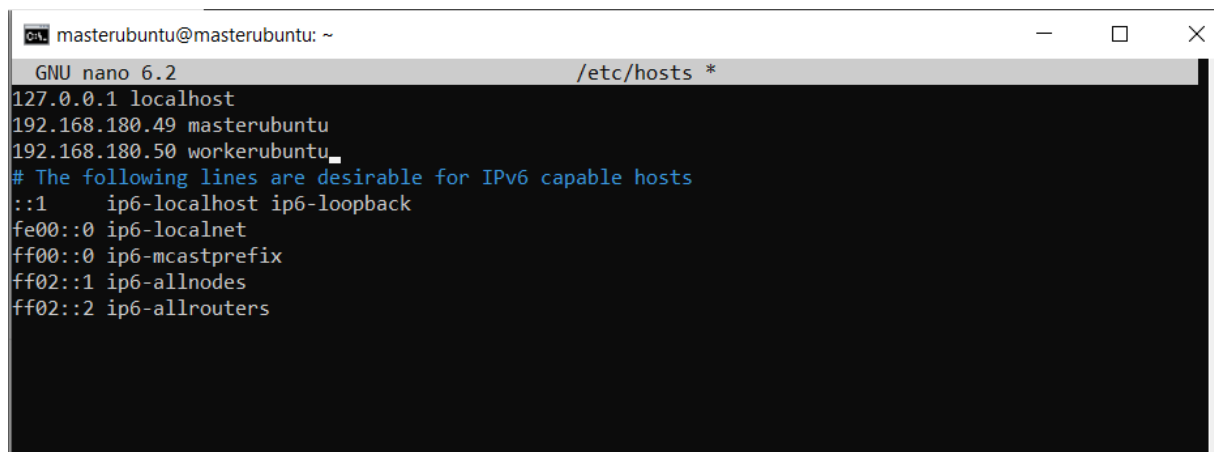
```
masterubuntu@masterubuntu:~$ su - hadoopuser
```

- Définition du nom d'hôte :

```
masterubuntu@masterubuntu:~$ sudo hostnamectl set-hostname ubuntu-master
[sudo] password for masterubuntu:
```

- Mettre à jour /etc/hosts

```
masterubuntu@masterubuntu:~$ sudo nano /etc/hosts
```



```
GNU nano 6.2 /etc/hosts *
127.0.0.1 localhost
192.168.180.49 masterubuntu
192.168.180.50 workerubuntu
# The following lines are desirable for IPv6 capable hosts
::1 ip6-localhost ip6-loopback
fe00::0 ip6-localnet
ff00::0 ip6-mcastprefix
ff02::1 ip6-allnodes
ff02::2 ip6-allrouters
```

2.2 Configurer l'IP statique

- Vérification de l'interface réseau

```
masterubuntu@masterubuntu:~$ ip link
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN mode DEFAULT group default qlen 1000
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
2: ens34: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc mq state UP mode DEFAULT group default qlen 1000
    link/ether 00:0c:29:7c:78:c6 brd ff:ff:ff:ff:ff:ff
    altname enp2s2
```

- Modification de /etc/netplan/00-installer-config.yaml

```
masterubuntu@masterubuntu:~$ sudo nano /etc/netplan/00-installer-config.yaml
```

```
masterubuntu@masterubuntu: ~  
GNU nano 6.2 /etc/netplan/00-installer-config.yaml *  
network:  
  ethernets:  
    ens34:  
      addresses:  
        - 192.168.180.49/24  
      gateway4: 192.168.180.1  
      nameservers:  
        addresses: [8.8.8.8, 8.8.4.4]  
      version: 2
```

- Application:

```
masterubuntu@masterubuntu:~$ sudo netplan apply  
  
** (generate:2172): WARNING **: 14:29:51.596: Permissions for /etc/netplan/00-installer-config.yaml are too open. Netplan configuration should NOT be accessible by others.  
  
** (generate:2172): WARNING **: 14:29:51.596: `gateway4` has been deprecated, use default routes instead.  
See the 'Default routes' section of the documentation for more details.  
WARNING:root:Cannot call Open vSwitch: ovsdb-server.service is not running.  
  
** (process:2170): WARNING **: 14:29:52.029: Permissions for /etc/netplan/00-installer-config.yaml are too open. Netplan configuration should NOT be accessible by others.  
  
** (process:2170): WARNING **: 14:29:52.030: `gateway4` has been deprecated, use default routes instead.  
See the 'Default routes' section of the documentation for more details.  
  
** (process:2170): WARNING **: 14:29:52.283: Permissions for /etc/netplan/00-installer-config.yaml are too open. Netplan configuration should NOT be accessible by others.  
  
** (process:2170): WARNING **: 14:29:52.283: `gateway4` has been deprecated, use default routes instead.  
See the 'Default routes' section of the documentation for more details.  
  
** (process:2170): WARNING **: 14:29:52.284: Permissions for /etc/netplan/00-installer-config.yaml are too open. Netplan configuration should NOT be accessible by others.  
  
** (process:2170): WARNING **: 14:29:52.284: `gateway4` has been deprecated, use default routes instead.  
See the 'Default routes' section of the documentation for more details.
```

- Vérification

```
masterubuntu@masterubuntu:~$ ip addr show  
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN group default qlen 1000  
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00  
    inet 127.0.0.1/8 scope host lo  
        valid_lft forever preferred_lft forever  
    inet6 ::1/128 scope host  
        valid_lft forever preferred_lft forever  
2: ens34: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc mq state UP group default qlen 1000  
    link/ether 00:0c:29:7c:78:c6 brd ff:ff:ff:ff:ff:ff  
    altname enp2s2  
    inet 192.168.180.49/24 metric 100 brd 192.168.180.255 scope global dynamic ens34  
        valid_lft 86334sec preferred_lft 86334sec  
    inet6 fe80::20c:29ff:fe7c:78c6/64 scope link  
        valid_lft forever preferred_lft forever
```

2.3 Tester la connectivité

- Teste de la connexion réseau :

```
masterubuntu@masterubuntu:~$ ping -c 4 8.8.8.8
PING 8.8.8.8 (8.8.8.8) 56(84) bytes of data.
64 bytes from 8.8.8.8: icmp_seq=1 ttl=114 time=22.9 ms
64 bytes from 8.8.8.8: icmp_seq=2 ttl=114 time=22.8 ms
64 bytes from 8.8.8.8: icmp_seq=3 ttl=114 time=22.9 ms
64 bytes from 8.8.8.8: icmp_seq=4 ttl=114 time=22.6 ms

--- 8.8.8.8 ping statistics ---
4 packets transmitted, 4 received, 0% packet loss, time 3005ms
rtt min/avg/max/mdev = 22.577/22.803/22.928/0.140 ms
masterubuntu@masterubuntu:~$ ping -c 4 google.com
PING google.com (216.58.209.78) 56(84) bytes of data.
64 bytes from mad07s22-in-f14.1e100.net (216.58.209.78): icmp_seq=1 ttl=114 time=21.3 ms
64 bytes from waw02s06-in-f14.1e100.net (216.58.209.78): icmp_seq=2 ttl=114 time=22.0 ms
64 bytes from waw02s06-in-f78.1e100.net (216.58.209.78): icmp_seq=3 ttl=114 time=20.9 ms
64 bytes from mad07s22-in-f14.1e100.net (216.58.209.78): icmp_seq=4 ttl=114 time=21.7 ms

--- google.com ping statistics ---
4 packets transmitted, 4 received, 0% packet loss, time 3004ms
rtt min/avg/max/mdev = 20.904/21.469/21.970/0.405 ms
masterubuntu@masterubuntu:~$
```

Étape 3 : Configurer le pare-feu

- Installation de UFW :

```
masterubuntu@masterubuntu:~$ sudo apt install ufw -y
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
ufw is already the newest version (0.36.1-4ubuntu0.1).
ufw set to manually installed.
0 upgraded, 0 newly installed, 0 to remove and 57 not upgraded.
```

- Autorisation des ports nécessaires :

```

masterubuntu@masterubuntu:~$ sudo ufw allow 22
Rules updated
Rules updated (v6)
masterubuntu@masterubuntu:~$ sudo ufw allow 9000
Rules updated
Rules updated (v6)
masterubuntu@masterubuntu:~$ sudo ufw allow 9870
Rules updated
Rules updated (v6)
masterubuntu@masterubuntu:~$ sudo ufw allow 8088
Rules updated
Rules updated (v6)
masterubuntu@masterubuntu:~$ sudo ufw allow 7077
Rules updated
Rules updated (v6)
masterubuntu@masterubuntu:~$ sudo ufw enable
Command may disrupt existing ssh connections. Proceed with operation (y|n)? y
Firewall is active and enabled on system startup

```

- Vérification :

```

masterubuntu@masterubuntu:~$ sudo ufw status
Status: active

To Action From
--
22 ALLOW Anywhere
9000 ALLOW Anywhere
9870 ALLOW Anywhere
8088 ALLOW Anywhere
7077 ALLOW Anywhere
22 (v6) ALLOW Anywhere (v6)
9000 (v6) ALLOW Anywhere (v6)
9870 (v6) ALLOW Anywhere (v6)
8088 (v6) ALLOW Anywhere (v6)
7077 (v6) ALLOW Anywhere (v6)

```

Étape 4 : Installer Java

- Installation d'OpenJDK 11 :

```
masterubuntu@masterubuntu:~$ sudo apt install openjdk-11-jdk -y
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
 alsa-topology-conf alsa-ucm-conf at-spi2-core ca-certificates-java dconf-gsettings-backend
dconf-service fontconfig-config fonts-dejavu-core fonts-dejavu-extra gsettings-desktop-schemas
java-common libasound2 libasound2-data libatk-bridge2.0-0 libatk-wrapper-java
libatk-wrapper-java-jni libatk1.0-0 libatk1.0-data libatspi2.0-0 libavahi-client3
libavahi-common-data libavahi-common3 libcups2 libdconf1 libdrm-amdgpu1 libdrm-intel1
libdrm-nouveau2 libdrm-radeon1 libfontconfig1 libfontenc1 libgif7 libgl1 libgl1-amd-dri
libgl1-mesa-dri libglapi-mesa libglvnd0 libglx-mesa0 libglx0 libgraphite2-3 libharfbuzz0b
libice-dev libice6 libjpeg-turbo8 libjpeg8 liblcms2-2 libllvm15 libpciaccess0 libpcsclite1
libpthread-stubs0-dev libsensors-config libsensors5 libsm-dev libsm6 libx11-dev libx11-xcb1
libxau-dev libxaw7 libxcb-dri2-0 libxcb-dri3-0 libxcb-glx0 libxcb-present0 libxcb-randr0
libxcb-shape0 libxcb-shm0 libxcb-sync1 libxcb-xfixes0 libxcb1-dev libxcomposite1 libxdmcp-dev
libxf86vm3 libxft2 libxi6 libxinerama1 libxkbfile1 libxmu6 libxpm4 libxrandr2 libxrender1
libxshmfence1 libxt-dev libxt6 libxtst6 libxv1 libxxf86dga1 libxxf86vm1 openjdk-11-jdk-headless
openjdk-11-jre openjdk-11-jre-headless session-migration x11-common x11-utils x11proto-dev
```

- Vérification :

```
masterubuntu@masterubuntu:~$ java -version
openjdk version "11.0.27" 2025-04-15
OpenJDK Runtime Environment (build 11.0.27+6-post-Ubuntu-0ubuntu122.04)
OpenJDK 64-Bit Server VM (build 11.0.27+6-post-Ubuntu-0ubuntu122.04, mixed mode, sharing)
```

- Configuration JAVA_HOME

```
Sélection masterubuntu@masterubuntu: ~
GNU nano 6.2 /etc/environment *
PATH="/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:/usr/local/games:/snap>
JAVA_HOME="/usr/lib/jvm/java-11-openjdk-amd64"
```

Application :

```
masterubuntu@masterubuntu:~$ source /etc/environment
```

Vérification :

```
masterubuntu@masterubuntu:~$ echo $JAVA_HOME
/usr/lib/jvm/java-11-openjdk-amd64
```

Étape 5 : Installer et configurer Hadoop

5.1 Installer Hadoop

- Téléchargement de Hadoop 3.3.6

```
masterubuntu@masterubuntu:~$ wget https://downloads.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz
--2025-05-29 14:48:49-- https://downloads.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.208.237, 2a01:4f9:3a:2c57::2, ...
Connecting to downloads.apache.org (downloads.apache.org)|135.181.214.104|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 730107476 (696M) [application/x-gzip]
Saving to: 'hadoop-3.3.6.tar.gz'

hadoop-3.3.6.tar.gz      37%[=====>                ] 263.96M  5.83MB/s   eta 9m 5s
```

- Décompression du fichier :

```
masterubuntu@masterubuntu:~$ sudo tar -xzf hadoop-3.3.6.tar.gz -C /usr/local
masterubuntu@masterubuntu:~$ sudo mv /usr/local/hadoop-3.3.6 /usr/local/hadoop
masterubuntu@masterubuntu:~$ sudo chown -R hadoopuser:hadoopuser /usr/local/hadoop
```

- Configurer les variables d'environnement :

```
masterubuntu@masterubuntu:~$ nano ~/.bashrc
```

- L'Ajoute :

```
# === Java environment ===
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64
export PATH=$PATH:$JAVA_HOME/bin

export HADOOP_HOME=/usr/local/hadoop
export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
```

- Application :

```
masterubuntu@masterubuntu:~$ source ~/.bashrc
```

- Vérification :


```
hadoopuser@ubuntumaster:~$ hadoop version
Hadoop 3.3.6
Source code repository https://github.com/apache/hadoop.git -r 1be78238728da9266a4f88195058f08fd012bf9c
Compiled by ubuntu on 2023-06-18T08:22Z
Compiled on platform linux-x86_64
Compiled with protoc 3.7.1
From source with checksum 5652179ad55f76cb287d9c633bb53bbd
This command was run using /usr/local/hadoop/share/hadoop/common/hadoop-common-3.3.6.jar
```

5.2 Configurer Hadoop

- L'Accès au répertoire de configuration :

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$
```

- Modification des fichiers :

- **core-site.xml :**

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ nano core-site.xml
```

```
<configuration>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://ubuntumaster:9000</value>
  </property>
</configuration>
```

- **hdfs-site.xml :**

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ nano hdfs-site.xml
```

```
<configuration>
<property>
  <name>dfs.replication</name>
  <value>1</value>
</property>
<property>
  <name>dfs.namenode.name.dir</name>
  <value>file:///usr/local/hadoop/hdfs/namenode</value>
</property>
<property>
  <name>dfs.datanode.data.dir</name>
  <value>file:///usr/local/hadoop/hdfs/datanode</value>
</property>
</configuration>
```

- **yarn-site.xml**

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ nano yarn-site.xml
```

```
<configuration>

<!-- Site specific YARN configuration properties -->
<property>
  <name>yarn.resourcemanager.hostname</name>
  <value>ubuntumaster</value>
</property>
<property>
  <name>yarn.nodemanager.aux-services</name>
  <value>mapreduce_shuffle</value>
</property>
</configuration>
```

- mapred-site.xml

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ nano mapred-site.xml
```

```
<configuration>
<property>
  <name>mapreduce.framework.name</name>
  <value>yarn</value>
</property>
</configuration>
```

- workers

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ nano workers
```

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ nano workers
GNU nano 6.2 workers *
workerubuntu
```

- Création des répertoires HDFS :

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ sudo mkdir -p /usr/local/hadoop/hdfs/namenode
[sudo] password for hadoopuser:
```

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ sudo mkdir -p /usr/local/hadoop/hdfs/datanode
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ sudo chown -R hadoopuser:hadoopuser /usr/local/hadoop
```

5.3 Formater le NameNode

- Formater le NameNode :

```
hadoopuser@ubuntumaster:/usr/local/hadoop/etc/hadoop$ hdfs namenode -format
WARNING: /usr/local/hadoop/logs does not exist. Creating.
2025-05-29 15:48:43,217 INFO namenode.NameNode: STARTUP_MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = ubuntumaster/192.168.180.49
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 3.3.6
STARTUP_MSG: classpath = /usr/local/hadoop/etc/hadoop:/usr/local/hadoop/share/hadoop/common/lib/net
ty-codec-stomp-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/zookeeper-3.6.3.jar:/usr/lo
cal/hadoop/share/hadoop/common/lib/paranamer-2.3.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-
webapp-9.4.51.v20230217.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-transport-4.1.89.Final.ja
r:/usr/local/hadoop/share/hadoop/common/lib/checker-qual-2.5.2.jar:/usr/local/hadoop/share/hadoop/com
mon/lib/kerb-simplekdc-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-math3-3.1.1.jar:/u
sr/local/hadoop/share/hadoop/common/lib/javax.servlet-api-3.1.0.jar:/usr/local/hadoop/share/hadoop/co
mmon/lib/kerby-util-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-servlet-9.4.51.v2023021
7.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-codec-socks-4.1.89.Final.jar:/usr/local/hadoop/
share/hadoop/common/lib/slf4j-api-1.7.36.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-codec-sm
```

Étape 6 : Installer et configurer Spark

6.1 Installer Spark

- Téléchargement Spark 3.4.2 :

```
hadoopuser@ubuntumaster:~$ wget https://archive.apache.org/dist/spark/spark-3.4.2/spark-3.4.2-bin-had
oop3.tgz
--2025-05-29 16:31:01-- https://archive.apache.org/dist/spark/spark-3.4.2/spark-3.4.2-bin-hadoop3.tg
z
Resolving archive.apache.org (archive.apache.org)... 65.108.204.189, 2a01:4f9:1a:a084::2
Connecting to archive.apache.org (archive.apache.org)[65.108.204.189]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 388664780 (371M) [application/x-gzip]
Saving to: 'spark-3.4.2-bin-hadoop3.tgz'

spark-3.4.2-bin-hadoop3.t 15%[====>] 55.79M 9.87MB/s eta 34s
```

- Vérification que le fichier est téléchargé :

```
hadoopuser@ubuntumaster:~$ ls -lh spark-3.4.2-bin-hadoop3.tgz
-rw-rw-r-- 1 hadoopuser hadoopuser 371M Nov 25 2023 spark-3.4.2-bin-hadoop3.tgz
```

- Décompression du fichier

```
hadoopuser@ubuntumaster:~$ sudo tar -xzf spark-3.4.2-bin-hadoop3.tgz -C /usr/local
[sudo] password for hadoopuser:
```

```
hadoopuser@ubuntumaster:~$ sudo mv /usr/local/spark-3.4.2-bin-hadoop3 /usr/local/spark
hadoopuser@ubuntumaster:~$ sudo chown -R hadoopuser:hadoopuser /usr/local/spark
```



```
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64
export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
export SPARK_MASTER_HOST=ubuntumaster_
```

- Modifier workers :

```
hadoopuser@ubuntumaster:/usr/local/spark/conf$ cp workers.template workers
hadoopuser@ubuntumaster:/usr/local/spark/conf$ nano workers
```

```
hadoopuser@ubuntumaster: /usr/local/spark/conf
GNU nano 6.2 workers *
#
# Licensed to the Apache Software Foundation (ASF) under one or more
# contributor license agreements. See the NOTICE file distributed with
# this work for additional information regarding copyright ownership.
# The ASF licenses this file to You under the Apache License, Version 2.0
# (the "License"); you may not use this file except in compliance with
# the License. You may obtain a copy of the License at
#
# http://www.apache.org/licenses/LICENSE-2.0
#
# Unless required by applicable law or agreed to in writing, software
# distributed under the License is distributed on an "AS IS" BASIS,
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
# See the License for the specific language governing permissions and
# limitations under the License.
#
# A Spark Worker will be started on each of the machines listed below.
workerubuntu_
```

Des etapes sur VM ubuntuworkerfch :

<input type="checkbox"/>	ubuntuworkerfch		Normal	68.08 GB	Ubuntu Linux (64-bit)	worker-node	18 MHz	2.74 GB
--------------------------	-----------------	--	--------	----------	-----------------------	-------------	--------	---------

Objectif :

Intégrer un serveur Ubuntu comme **nœud worker** dans un cluster Hadoop/Spark afin qu'il puisse exécuter les tâches de traitement de données distribuées sous la supervision du nœud master.

1. Installation de la distribution Ubuntu Server

Nous avons utilisé une **machine dédiée** dans la salle E26 pour installer **Ubuntu Server 22.04 LTS** (version minimale, sans interface graphique).

```
PS C:\Users\surface> ssh workerubuntu@192.168.180.50
workerubuntu@192.168.180.50's password:
Welcome to Ubuntu 22.04.5 LTS (GNU/Linux 5.15.0-140-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/pro

System information as of Thu May 29 01:36:55 PM UTC 2025

System load:  0.0               Processes:            194
Usage of /:   26.0% of 27.86GB  Users logged in:     1
Memory usage: 4%               IPv4 address for ens34: 192.168.180.50
Swap usage:   0%

Expanded Security Maintenance for Applications is not enabled.

135 updates can be applied immediately.
80 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable
```

2. Configuration système initiale

- Nom d'hôte personnalisé :

```
workerubuntu@workerubuntu:~$ sudo hostnamectl set-hostname worker-node
workerubuntu@workerubuntu:~$
```

- Configuration du fichier `/etc/hosts` :

```
workerubuntu@workerubunt  ×  +  ▾  
GNU nano 6.2 /etc/hosts *  
127.0.0.1 localhost  
127.0.1.1 workerubuntu  
  
# The following lines are desirable for IPv6 capable hosts  
::1 ip6-localhost ip6-loopback  
fe00::0 ip6-localnet  
ff00::0 ip6-mcastprefix  
ff02::1 ip6-allnodes  
ff02::2 ip6-allrouters  
192.168.180.49 master-node  
192.168.180.50 worker-node
```

- Création d'un utilisateur hadoop :

```
workerubuntu@workerubuntu:~$ sudo adduser hadoop  
Adding user `hadoop' ...  
Adding new group `hadoop' (1001) ...  
Adding new user `hadoop' (1001) with group `hadoop' ...  
Creating home directory `/home/hadoop' ...  
Copying files from `/etc/skel' ...  
  
workerubuntu@workerubuntu:~$ sudo usermod -aG sudo hadoop  
workerubuntu@workerubuntu:~$ su - hadoop
```

3. Installation de Java

```
workerubuntu@worker-node:~$ java -version
openjdk version "11.0.27" 2025-04-15
OpenJDK Runtime Environment (build 11.0.27+6-post-Ubuntu-0ubuntu122.04)
OpenJDK 64-Bit Server VM (build 11.0.27+6-post-Ubuntu-0ubuntu122.04, mixed mode, sharing)
```

```
GNU nano 6.2 /home/workerubuntu/.bashrc *
# Alias definitions.
# You may want to put all your additions into a separate file like
# ~/.bash_aliases, instead of adding them here directly.
# See /usr/share/doc/bash-doc/examples in the bash-doc package.

if [ -f ~/.bash_aliases ]; then
    . ~/.bash_aliases
fi

# enable programmable completion features (you don't need to enable
# this, if it's already enabled in /etc/bash.bashrc and /etc/profile
# sources /etc/bash.bashrc).
if ! shopt -oq posix; then
    if [ -f /usr/share/bash-completion/bash_completion ]; then
        . /usr/share/bash-completion/bash_completion
    elif [ -f /etc/bash_completion ]; then
        . /etc/bash_completion
    fi
fi

# === Java environment ===
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64
export PATH=$PATH:$JAVA_HOME/bin
```

```
workerubuntu@worker-node:~$ source ~/.bashrc
```

4. Installation de Hadoop et Spark

- Déplacement dans /usr/local

```
workerubuntu@worker-node:~$ wget https://dldcn.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz
--2025-05-29 14:27:20-- https://dldcn.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz
Resolving dldcn.apache.org (dldcn.apache.org)... 151.101.2.132, 2a04:4e42::644
Connecting to dldcn.apache.org (dldcn.apache.org)[151.101.2.132]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 730107476 (696M) [application/x-gzip]
Saving to: 'hadoop-3.3.6.tar.gz'

hadoop-3.3.6.tar.gz      100%[=====] 696.28M  2.24MB/s   in 8m 50s
```

- Décompression des fichiers téléchargés

```
hadoop-3.3.6/share/hadoop/hdfs/jdiff/Apache_Hadoop_HDFS_3.1.3.xml
hadoop-3.3.6/share/hadoop/hdfs/hadoop-hdfs-client-3.3.6-tests.jar
hadoop-3.3.6/share/hadoop/hdfs/hadoop-hdfs-httpfs-3.3.6.jar
tar (child): spark-3.5.0-bin-hadoop3.tgz: Cannot open: No such file or directory
```

```
spark-3.5.0-bin-hadoop3/bin/spark-shell
spark-3.5.0-bin-hadoop3/bin/find-spark-home
workerubuntu@worker-node:~$ sudo mv spark-3.5.0-bin-hadoop3 /usr/local/spark
```



```
# == Hadoop & Spark ==
export HADOOP_HOME=/usr/local/hadoop
export SPARK_HOME=/usr/local/spark
export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$SPARK_HOME/bin:$SPARK_HOME/sbin
```

[illegible]

5. Configuration des fichiers Hadoop

Dans le dossier `$HADOOP_HOME/etc/hadoop` :

- **core-site.xml :**

```

workeruntu@worker-node  X  +  v
GNU nano 6.2                                     core-site.xml *
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://master-node:9000</value>
  </property>
</configuration>

```

- **hdfs-site.xml :**

```
<configuration>
  <property>
    <name>dfs.data.dir</name>
    <value>/home/hadoop/hdfs/data</value>
  </property>
</configuration>
```

- yarn-site.xml :

```
<configuration>

<!-- Site specific YARN configuration properties -->

  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>
</configuration>
```

Phase 2 – Manipulation et stockage

Étape 1 : Formatage de hdfs :

```
hadoopuser@ubuntu-master:~$ hdfs namenode -format
2025-05-30 13:02:31,867 INFO namenode.NameNode: STARTUP_MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = ubuntu-master/192.168.180.49
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 3.3.6
STARTUP_MSG: classpath = /usr/local/hadoop/etc/hadoop:/usr/local/hadoop/share/hadoop/common/lib/netty-codec-stomp-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/zookeeper-3.6.3.jar:/usr/local/hadoop/share/hadoop/common/lib/paranamer-2.3.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-webapp-9.4.51.v20230217.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-transport-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/checker-qual-2.5.2.jar:/usr/local/hadoop/share/hadoop/common/lib/kerb-simplekdc-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-math3-3.1.1.jar:/usr/local/hadoop/share/hadoop/common/lib/javax.servlet-api-3.1.0.jar:/usr/local/hadoop/share/hadoop/common/lib/kerby-util-1.0.1.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-servlet-9.4.51.v20230217.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-codec-socks-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/slf4j-api-1.7.36.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-codec-smtp-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/avro-1.7.7.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-util-9.4.51.v20230217.jar:/usr/local/hadoop/share/hadoop/common/lib/jsp-api-2.1.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-codec-dns-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/slf4j-reload4j-1.7.36.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-transport-native-epoll-4.1.89.Final-linux-aarch_64.jar:/usr/local/hadoop/share/hadoop/common/lib/jackson-mapper-asl-1.9.13.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-handler-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/jsch-0.1.55.jar:/usr/local/hadoop/share/hadoop/common/lib/zookeeper-jute-3.6.3.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-transport-rxtx-4.1.89.Final.jar:/usr/local/hadoop/share/hadoop/common/lib/jetty-server-9.4.51.v20230217.jar:/usr/local/hadoop/share/hadoop/common/lib/jackson-databind-2.12.7.1.jar:/usr/local/hadoop/share/hadoop/common/lib/jersey-json-1.20.jar:/usr/local/hadoop/share/hadoop/common/lib/hadoop-shaded-protobuf-3.7-1.1.1.jar:/usr/local/hadoop/share/hadoop/common/lib/netty-resolve
```

Étape 2 : Tester les services

- Démarrage de HDFS :

```
hadoopuser@ubuntumaster:~$ start-dfs.sh
Starting namenodes on [ubuntumaster]
Starting datanodes
Starting secondary namenodes [ubuntumaster]
```

- Démarrage de YARN :

```
hadoopuser@ubuntumaster:~$ start-yarn.sh
Starting resourcemanager
Starting nodemanagers
```

- Vérification des processus :

```
hadoopuser@ubuntumaster:~$ jps
6131 NameNode
6567 ResourceManager
7660 Jps
6366 SecondaryNameNode
```

- Vérification dans le navigateur :
 - HDFS : <http://192.168.180.49:9870>

←

→

↺

⚠ Not secure

192.168.180.49:9870/dfshealth.html#tab-overview

☆

🔍

🔖

🔴

⋮

Hadoop

Overview

Datanodes

Datanode Volume Failures

Snapshot

Startup Progress

Utilities

Overview

'ubuntumaster:9000' (✓active)

Started:	Fri May 30 12:47:41 +0100 2025
Version:	3.3.6, r1be78238728da9266a4f88195058f08fd012bf9c
Compiled:	Sun Jun 18 09:22:00 +0100 2023 by ubuntu from (HEAD detached at release-3.3.6-RC1)
Cluster ID:	CID-30d7be14-bbd4-455b-bbfa-f9ff9629327b
Block Pool ID:	BP-1492602910-192.168.180.49-1748533724152

Summary

Security is off.

Safemode is off.

1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

Heap Memory used 97.38 MB of 151 MB Heap Memory. Max Heap Memory is 978 MB.

Non Heap Memory used 52.25 MB of 55.63 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Configured Capacity:	0 B
Configured Remote Capacity:	0 B
DFS Used:	0 B (100%)
Non DFS Used:	0 B

- YARN : <http://192.168.180.49:8088>

←

→

↺

⚠ Not secure

192.168.180.49:8088/cluster

🔍

☆

🔍

🔖

🔴

⋮

hadoop

All Applications

Cluster

About

Nodes

Node Labels

Applications

NEW

VIEW

SCHEDULE

SUBMITTED

ACCEPTED

PENDING

FAILED

KILLED

Schedule

Tools

Cluster Metrics

Apps Submitted	0	Apps Pending	0	Apps Running	0	Apps Completed	0	Containers Running	0	Used Resources	<memory:0 B, vCores:0>	Total Resources	<memory:0 B, vCores:0>	Reserved Resources	<memory:0 B, vCores:0>	Physical Mem Used %	0
----------------	---	--------------	---	--------------	---	----------------	---	--------------------	---	----------------	------------------------	-----------------	------------------------	--------------------	------------------------	---------------------	---

Cluster Nodes Metrics

Active Nodes	0	Decommissioning Nodes	0	Decommissioned Nodes	0	Lost Nodes	0	Unhealthy Nodes	0	Rebooted Nodes	0
--------------	---	-----------------------	---	----------------------	---	------------	---	-----------------	---	----------------	---

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation	Maximum Cluster Application Priority
Capacity Scheduler	(memory-emb (unit-MB), vCores)	<memory:1024, vCores:1>	<memory:8192, vCores:8>	0

Show 20 jobs

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU V-Cores	Allocated Memory MB	Allocated GPUs	Reserved CPU V-Cores	Reserved Memory MB	Reserved GPUs	% of Queue	% of Cluster
No data available in table																				

Showing 0 to 0 of 0 entries

Etape 3 : Transfert du fichier .csv dans MasterUbuntu

Dans cette étape, nous avons transféré un fichier de données nommé `mental_disorders_redditttt.csv` depuis un poste local Windows vers le serveur Ubuntu Master, via le protocole SCP (Secure Copy Protocol).

Commande utilisée :

```
PS C:\Users\surface> scp "C:\Users\surface\Downloads\DataProduct.json" hadoopuser@192.168.180.49:/home/hadoopuser
hadoopuser@192.168.180.49's password:
DataProduct.json                                100% 7271KB   5.2MB/s   00
PS C:\Users\surface>
```

Activation réussie du service DataNode sur le nœud worker :

```
hadoopuser@worker-node:~$ nano /usr/local/hadoop/etc/hadoop/hdfs-site.xml
hadoopuser@worker-node:~$ hadoopuser@worker-node:~$
hadoopuser@worker-node:~$ hdfs --daemon start datanode
hadoopuser@worker-node:~$
hadoopuser@worker-node:~$ jps
30331 Jps
30254 DataNode
```

La commande `jps` confirme maintenant que DataNode est bien actif (PID : 30254), validant ainsi la participation du nœud worker au cluster Hadoop.

Etape 4 : Ingestion dans HDFS

Une fois le fichier présent sur le master, nous avons préparé le stockage dans HDFS pour permettre un traitement distribué via Spark.

```
hadoopuser@ubuntumaster:~$ ls -lh /home/hadoopuser/DataProduct.json
-rw-rw-r-- 1 hadoopuser hadoopuser 7.2M May 31 09:55 /home/hadoopuser/DataProduct.json
hadoopuser@ubuntumaster:~$
```

```
hadoopuser@ubuntumaster:~$ cd /usr/local/hadoop/sbin
./start-dfs.sh
Starting namenodes on [ubuntumaster]
Starting datanodes
workerubuntu: ssh: connect to host workerubuntu port 22: No route to host
Starting secondary namenodes [ubuntumaster]
hadoopuser@ubuntumaster:/usr/local/hadoop/sbin$ |
```

```
password:
hadoopuser@ubuntumaster:~$ hdfs dfs -ls /data/product_reviews/
Found 1 items
-rw-r--r-- 1 hadoopuser supergroup 7445886 2025-05-31 10:12 /data/product_reviews/DataProduct.json
hadoopuser@ubuntumaster:~$
```

Objectif de cette phase :

Préparer les données brutes pour analyse avec PySpark, en les plaçant dans un système distribué (HDFS). Cette phase garantit que les données sont accessibles par tous les nœuds du cluster.

Étape 5 : Traitement des données avec PySpark

Objectif : Supprimer les valeurs nulles, doublons, champs inutiles, normaliser le texte (pour NLP).

Dans le cadre de la Phase 2 (« Manipulation et stockage »), nous avons dû :

1. Charger un fichier JSON volumineux (≈ 7 Mo) dans HDFS.
2. Écrire un script PySpark pour lire ce fichier depuis HDFS.
3. Effectuer des opérations de base (filter, groupBy, agg, join) et nettoyer les données.
4. Sauvegarder le résultat traité au format Parquet dans HDFS.

Le jeu de données choisi s'appelle DataProduct.json et contient des avis clients (Amazon reviews) comprenant notamment les champs suivants :

- asin : identifiant produit
- overall : note (notation globale)
- reviewText : texte de l'avis

- reviewerID : identifiant du client
- reviewTime : date
- etc.

L'objectif du script est donc de lire ces avis, de :

- Nettoyer le texte (reviewText) pour ne conserver que des caractères alphabétiques (a–z) et des espaces,
- Filtrer les avis trop courts pour que l'analyse textuelle soit pertinente,
- Calculer des statistiques descriptives sur la longueur du texte nettoyé,
- Regrouper (groupBy) les avis par produit (asin) pour calculer, pour chaque produit, le nombre d'avis et la note moyenne,
- Effectuer une jointure (join) pour associer, à chaque avis individuel, les statistiques « nombre d'avis » et « note moyenne » du produit concerné,
- Sauvegarder en sortie deux jeux de résultats au format Parquet :
 1. Les avis nettoyés (avec une colonne clean_text et la longueur text_length).
 2. Les statistiques agrégées par produit (asin, nb_reviews, avg_rating).


```

hadoopuser@ubuntu-master: x + v
GNU nano 6.2 json_cleaning.py *
# -----
# json_cleaning.py
# -----
from pyspark.sql import SparkSession
from pyspark.sql.functions import lower, regexp_replace, length, col, avg, count

if __name__ == "__main__":
    # 1. Démarrer une session Spark
    spark = SparkSession.builder \
        .appName("JsonCleaningAndBasicStats") \
        .getOrCreate()

    # 2. Lire le fichier JSON depuis HDFS
    # (suppose que DataProduct.json a été chargé sous /data/product_reviews/)
    input_path = "hdfs:///data/product_reviews/DataProduct.json"
    df = spark.read.json(input_path)

    # 3. Explorer la structure du DataFrame
    print("=== Schéma initial du dataset JSON ===")
    df.printSchema()
    print("=== Aperçu des 5 premières lignes ===")
    df.show(5, truncate=False)

    # 4. Nettoyage du texte : deux objectifs principaux
    # - Supprimer les lignes où reviewText est absent (NULL)
    # - Mettre en minuscules et retirer tout caractère non alphabétique/spaces
    df_clean = df.dropna(subset=["reviewText"]) \
        .withColumn("clean_text", regexp_replace(lower(col("reviewText")), "[^a-zA-Z\\s]", ""))

    # 5. Afficher le schéma après ajout de la colonne 'clean_text'
    print("=== Schéma après ajout de la colonne clean_text ===")
    df_clean.printSchema()

    # 6. Filtrage basique : ne garder que les avis dont le texte "nettoyé" fait au moins 20 caractères
    df_filtered = df_clean.filter(col("clean_text").isNotNull() & (length(col("clean_text")) >= 20))

```

Lancement de script avec spark

```

hadoopuser@ubuntu-master:~$ nano json_cleaning.py
hadoopuser@ubuntu-master:~$ spark-submit json_cleaning.py
25/05/31 10:26:39 INFO SparkContext: Running Spark version 3.4.2
25/05/31 10:26:40 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
25/05/31 10:26:40 INFO ResourceUtils: =====
25/05/31 10:26:40 INFO ResourceUtils: No custom resources configured for spark.driver.
25/05/31 10:26:40 INFO ResourceUtils: =====
25/05/31 10:26:40 INFO SparkContext: Submitted application: JsonCleaningAndBasicStats
25/05/31 10:26:40 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 1, script: , vendor: , memory -> name: memory, amount: 1024, script: , vendor: , offHeap -> name: offHeap, amount: 0, script: , vendor: ), task resources: Map(cpus -> name: cpus, amount: 1)
25/05/31 10:26:40 INFO ResourceProfile: Limiting resource is cpu
25/05/31 10:26:40 INFO ResourceProfileManager: Added ResourceProfile id: 0
25/05/31 10:26:40 INFO SecurityManager: Changing view acls to: hadoopuser
25/05/31 10:26:40 INFO SecurityManager: Changing modify acls to: hadoopuser
25/05/31 10:26:40 INFO SecurityManager: Changing view acls groups to:
25/05/31 10:26:40 INFO SecurityManager: Changing modify acls groups to:
25/05/31 10:26:40 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: hadoopuser; groups with view permissions: EMPTY; users with modify permissions: hadoopuser; groups with modify permissions: EMPTY
25/05/31 10:26:40 INFO Utils: Successfully started service 'sparkDriver' on port 46677.
25/05/31 10:26:41 INFO SparkEnv: Registering MapOutputTracker
25/05/31 10:26:41 INFO SparkEnv: Registering BlockManagerMaster
25/05/31 10:26:41 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
25/05/31 10:26:41 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
25/05/31 10:26:41 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
25/05/31 10:26:41 INFO DiskBlockManager: Created local directory at /tmp/blockmgr-68bd592e-2ee8-4943-bcc2-7281acad8e04
25/05/31 10:26:41 INFO MemoryStore: MemoryStore started with capacity 434.4 MiB
25/05/31 10:26:41 INFO SparkEnv: Registering OutputCommitCoordinator
25/05/31 10:26:41 INFO JettyUtils: Start Jetty 0.0.0.0:4040 for SparkUI
25/05/31 10:26:41 INFO Utils: Successfully started service 'SparkUI' on port 4040.
25/05/31 10:26:41 INFO Executor: Starting executor ID driver on host masterubuntu
25/05/31 10:26:41 INFO Executor: Starting executor with user classpath (userClassPathFirst = false): ''
25/05/31 10:26:41 INFO Utils: Successfully started service 'org.apache.spark.network.netty.NettyBlockTransferService' on port 35803.
25/05/31 10:26:41 INFO NettyBlockTransferService: Server created on masterubuntu:35803
25/05/31 10:26:41 INFO BlockManager: Using org.apache.spark.storage.RandomBlockReplicationPolicy for block replication policy
25/05/31 10:26:41 INFO BlockManagerMaster: Registering BlockManager BlockManagerId(driver, masterubuntu, 35803, None)

```



```
25/05/31 11:01:45 INFO CodeGenerator: Code generated in 7.39936 ms
```

asin	reviewerID	text_length	nb_reviews	avg_rating
1384719342	A2IBPI20UZIR0U	261	5	5.0
1384719342	A14VAT5EAX3D9S	529	5	5.0
1384719342	A195EZSQDW3E21	431	5	5.0
1384719342	A2C00NNG1ZQQG2	204	5	5.0
1384719342	A94QU4C90B1AX	155	5	5.0
B00004Y2UT	A2A039TZMZH9Y	227	6	4.666666666666667
B00004Y2UT	A1UPZM995ZAH90	184	6	4.666666666666667
B00004Y2UT	AJNFQI3YR6XJ5	795	6	4.666666666666667
B00004Y2UT	A3M1PLEYNDEY08	195	6	4.666666666666667
B00004Y2UT	AMNTZU1YQN1TH	214	6	4.666666666666667

```
only showing top 10 rows
```

Interprétation du tableau après jointure :

Chaque ligne correspond à un avis individuel (identifié par sa colonne reviewerID) associé au produit (asin) qu'il concerne. Les colonnes suivantes détaillent, pour ce même produit :

- **text_length** : la longueur (en caractères) du texte nettoyé de cet avis.
- **nb_reviews** : le nombre total d'avis recensés pour le produit asin (agrégation count(*)).
- **avg_rating** : la note moyenne (champ overall) pour ce produit (agrégation avg(overall)).

Par exemple, on voit que pour le produit 1384719342 (premières lignes) :

- Cet avis a un texte net d'une longueur de 261 caractères.
- Au total, ce produit a 5 avis, et sa note moyenne est de 5,0.

Pour le produit B00004Y2UT (lignes suivantes) :

- Un avis donné a une longueur de 227 caractères.
- Ce produit compte 6 avis au total, avec une note moyenne de 4,6666...

En résumé, cette jointure a « enrichi » chaque avis individuel en lui associant les statistiques globales du produit correspondant (nombre d'avis et note moyenne). Cela permet, par exemple, de filtrer ou de trier les avis non seulement selon leur propre contenu (longueur du texte) mais aussi en fonction de la popularité ou de la qualité générale du produit (via nb_reviews et avg_rating).

Après cette étape, nous pouvons enregistrer le résultat complet (avis nettoyés + statistiques produit) au format Parquet pour une utilisation ultérieure (visualisation, analyses plus avancées, etc.).

```
25/05/31 11:01:43 INFO DAGScheduler: Job 7 finished
25/05/31 11:01:43 INFO CodeGenerator: Code generate
25/05/31 11:01:43 INFO CodeGenerator: Code generate
+-----+-----+-----+
|asin      |nb_reviews|avg_rating|
+-----+-----+-----+
|B003VWJ2K8|162       |4.685185185185185|
|B0002E1G5C|142       |4.577464788732394|
|B0002F7K7Y|116       |4.620689655172414|
|B003VWKPHC|114       |4.649122807017544|
|B0002H0A3S|92        |4.6521739130434785|
|B0002CZVXM|74        |4.5|
|B0006NDF8A|71        |4.492957746478873|
|B0009G1E0K|69        |4.492753623188406|
|B0002E2KPC|68        |4.5|
|B0002GLDQM|67        |4.373134328358209|
+-----+-----+-----+
only showing top 10 rows
```

Ce tableau présente les **10 produits** (colonne asin) ayant le **plus grand nombre d'avis** (nb_reviews). Pour chacun, on affiche également la **note moyenne** (avg_rating). Par exemple

- Le produit B003VWJ2K8 compte **162 avis** et a une note moyenne d'environ **4,685**.
- Le produit B0002E1G5C en compte **142**, avec une note moyenne de **4,577**.

```
25/05/31 11:01:41 INFO DAGScheduler: ResultStage 6 (showString at NativeMethodAccessorImpl.java:0) finished in 0.084 s
25/05/31 11:01:41 INFO DAGScheduler: Job 5 is finished. Cancelling potential speculative or zombie tasks for this job
25/05/31 11:01:41 INFO TaskSchedulerImpl: Killing all running tasks in stage 6: Stage finished
25/05/31 11:01:41 INFO DAGScheduler: Job 5 finished: showString at NativeMethodAccessorImpl.java:0, took 0.089246 s
25/05/31 11:01:41 INFO BlockManagerInfo: Removed broadcast_7_piece0 on masterubuntu:32891 in memory (size: 34.6 KiB, free: 434.3 MiB)
25/05/31 11:01:41 INFO CodeGenerator: Code generated in 9.4378 ms
+-----+-----+-----+
|asin|reviewerID|text_length|clean_text|
+-----+-----+-----+
|B000068NSX|A2EZWZ8MBEDOLN|1513|ive been using these cables for more than months and they are holding up pre...|
|B0000AQRSS|A2MR43RDPZX3J|3496|it is so hard to get an honest review of any mike that isnt a shure sm this i...|
|B0000AQRSS|A6FIAB28IS79|1472|i started with the shure pg and most recently acquired the shure sm id proba...|
|B0000AQRST|A2KI91IR3RA7D0|2369|although this mic is now considered primarily an instrument mic that was not ...|
|B0000AQRST|A34IJACMU8C3IM|1704|considering that the sm was introduced to the world in a year after the shur...|
+-----+-----+-----+
only showing top 5 rows

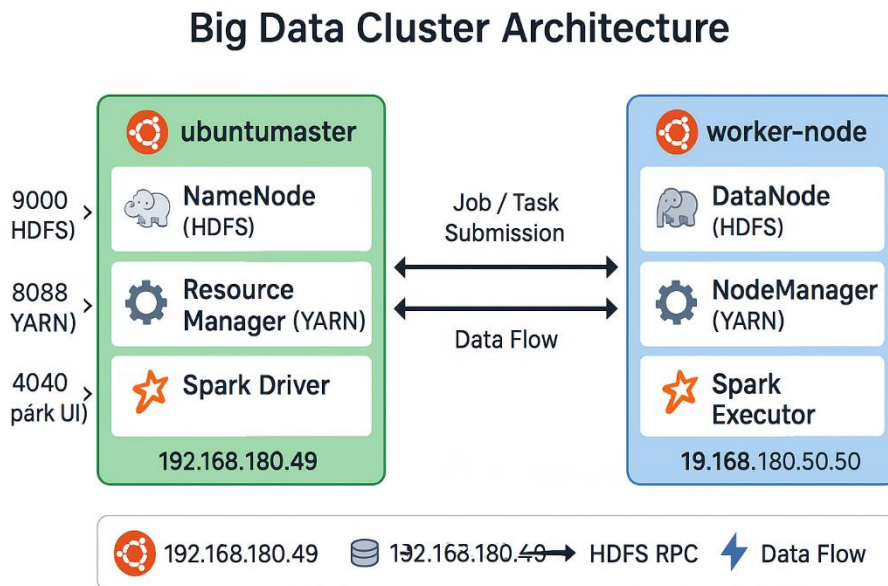
25/05/31 11:01:42 INFO BlockManagerInfo: Removed broadcast_8_piece0 on masterubuntu:32891 in memory (size: 16.5 KiB, free: 434.3 MiB)
=== Top 10 produits par nombre d'avis ===
25/05/31 11:01:42 INFO FileSourceStrategy: Pushed Filters: IsNotNull(reviewText)
```

On a fait **Sauvegarde des résultats au format Parquet dans HDFS**

- cleaned_reviews_parquet
- product_stats_parquet

Architecture distribuée du cluster

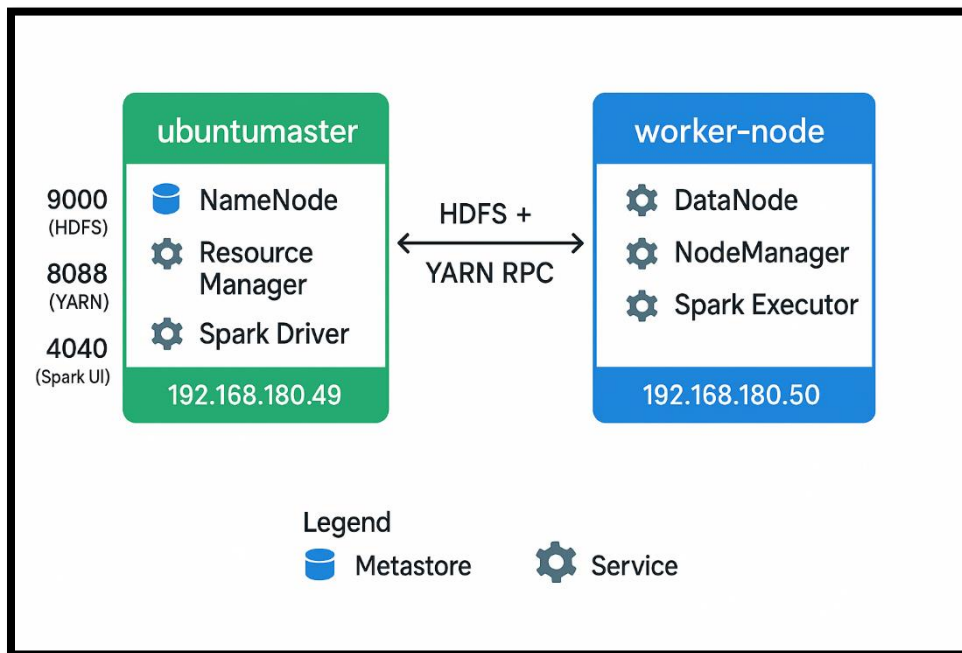
Hadoop/Spark:



Dans l'architecture de notre **cluster Big Data**, l'interaction entre le nœud **ubuntumaster** et le nœud **ubuntuworker** repose sur la répartition des responsabilités entre coordination et exécution. Le **nœud maître (ubuntumaster)** agit comme le cerveau du système, assurant la gestion des métadonnées HDFS via le **NameNode**, la planification des ressources via le **ResourceManager (YARN)** et l'orchestration des tâches Spark grâce au **Spark Driver**. En parallèle, le nœud esclave (**ubuntuworker**) est dédié à l'exécution : il héberge un **DataNode** pour stocker physiquement les blocs de données, un **NodeManager** pour exécuter les conteneurs YARN, et un ou plusieurs **Spark Executors** pour effectuer les calculs distribués.

Lorsqu'une commande `spark-submit` est lancée depuis le master, le **Spark Driver** répartit les tâches à travers les workers **via YARN RPC**. Ces workers lisent les blocs depuis **HDFS** (stockés localement dans leurs **DataNodes**), traitent les données, puis renvoient les résultats au driver. Ce modèle distribué garantit une scalabilité, une tolérance aux pannes, et une efficacité accrue dans le traitement de gros volumes de données, comme notre fichier **DataProduct.json**. Ainsi, le cluster fonctionne comme une entité unifiée grâce à une communication continue entre le master et ses workers, orchestrée par les couches **Spark et HDFS**.

Architecture de Cluster :

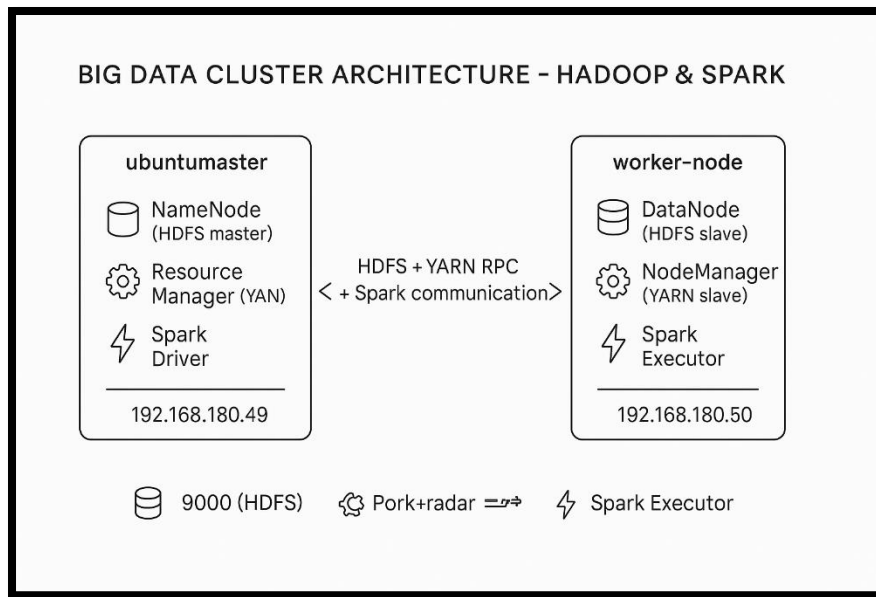


Cette figure illustre l'**architecture du cluster Big Data** déployé pour notre projet. Le nœud **ubuntu-master** (IP : 192.168.180.49) joue le rôle de **maître**, hébergeant trois services essentiels :

- le **NameNode** pour la gestion du système de fichiers HDFS,
- le **ResourceManager** pour la gestion des ressources via YARN,
- et le **Spark Driver** pour le pilotage des tâches Spark.

Le nœud **worker-node** (IP : 192.168.180.50) agit en tant que **nœud esclave**, exécutant :

- un **DataNode** (stockage de blocs HDFS),
- un **NodeManager** (exécution des conteneurs YARN),
- et un **Spark Executor** (exécution des jobs Spark distribués).



Rôles des nœuds dans Notre cluster

❖ **ubuntu-master (Master Node)**

Ce nœud central **coordonne** tout le système :

- **HDFS NameNode** : Gère la structure du système de fichiers (où chaque bloc est stocké)
- **YARN ResourceManager** : Gère les ressources et assigne les tâches aux workers
- **Spark Driver** : Lance et pilote les applications Spark (comme `json_cleaning.py`)
- Toutes les commandes `spark-submit`, `hdfs dfs`, etc., se font ici.

❖ **ubuntuworker (Worker Node)**

Ce nœud **exécute** les calculs et stocke les données :

- **HDFS DataNode** : Stocke physiquement les blocs de données (JSON, Parquet...)
- **YARN NodeManager** : Exécute les conteneurs pour Spark
- **Spark Executor** : Réalise les transformations, filtrages, agrégations demandées par Spark
- C'est ici que **les calculs lourds** sont faits (`filter`, `groupBy`, `aggregation`...) lorsque on lance `spark-submit` depuis le master.

Étape	Fait sur...	Rôle de ubuntuworker
Lancer les services HDFS/YARN	ubuntumaster + worker	Démarrage du DataNode/NodeManager
Copier les données	ubuntumaster	Le DataNode stocke les blocs
spark-submit (nettoyage)	ubuntumaster	Worker exécute les calculs
Sauvegarde Parquet	ubuntumaster	Les fichiers sont écrits via DataNode
Visualisation (Superset)	ubuntumaster	(Pas directement lié au worker)

Phase 3 : MINI-CAS D'ANALYSE : avis produits Amazon

Intégration et Visualisation des Données via Grafana Conteneurisé

Lancement de Grafana via Docker

```

PS C:\Users\surface> Copy-Item "C:\Users\surface\Desktop\Certifications\prodigy infotech tasks\dataaProducts.csv" "$env:
USERPROFILE\grafana_project\data\"
PS C:\Users\surface> docker run -d -p 3000:3000 `
>> --name=grafana-csv `
>> -v "$env:USERPROFILE\grafana_project\var\lib\grafana" `
>> -v "$env:USERPROFILE\grafana_project\data\var\lib\grafana\data" `
>> grafana/grafana
Unable to find image 'grafana/grafana:latest' locally
latest: Pulling from grafana/grafana
9a8c18aee5ea: Download complete
f18232174bc9: Pull complete
65babbe3dfe5: Pull complete
651b0ba49b07: Pull complete
8b5292c940e1: Downloading [=====] 49.28MB/63.48MB
d953cde4314b: Pull complete
f836d47fdc4d: Downloading [=====] 74.45MB/107.3MB
454a4350d439: Download complete
aec44cb03450: Pull complete
13fa68ca8757: Pull complete

```

Container CPU usage ⓘ

1.72% / 400% (4 CPUs available)

Container memory usage ⓘ

146.4MB / 7.54GB

[Show charts](#)

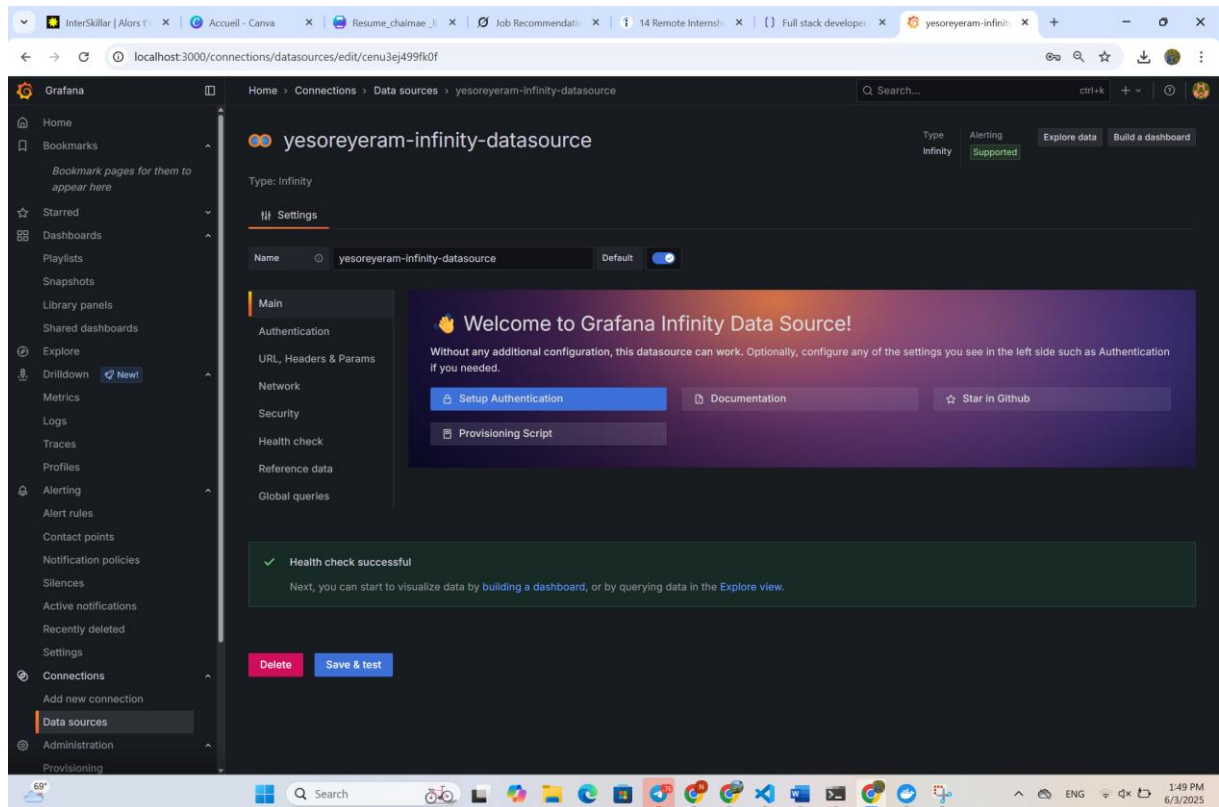
Q Search



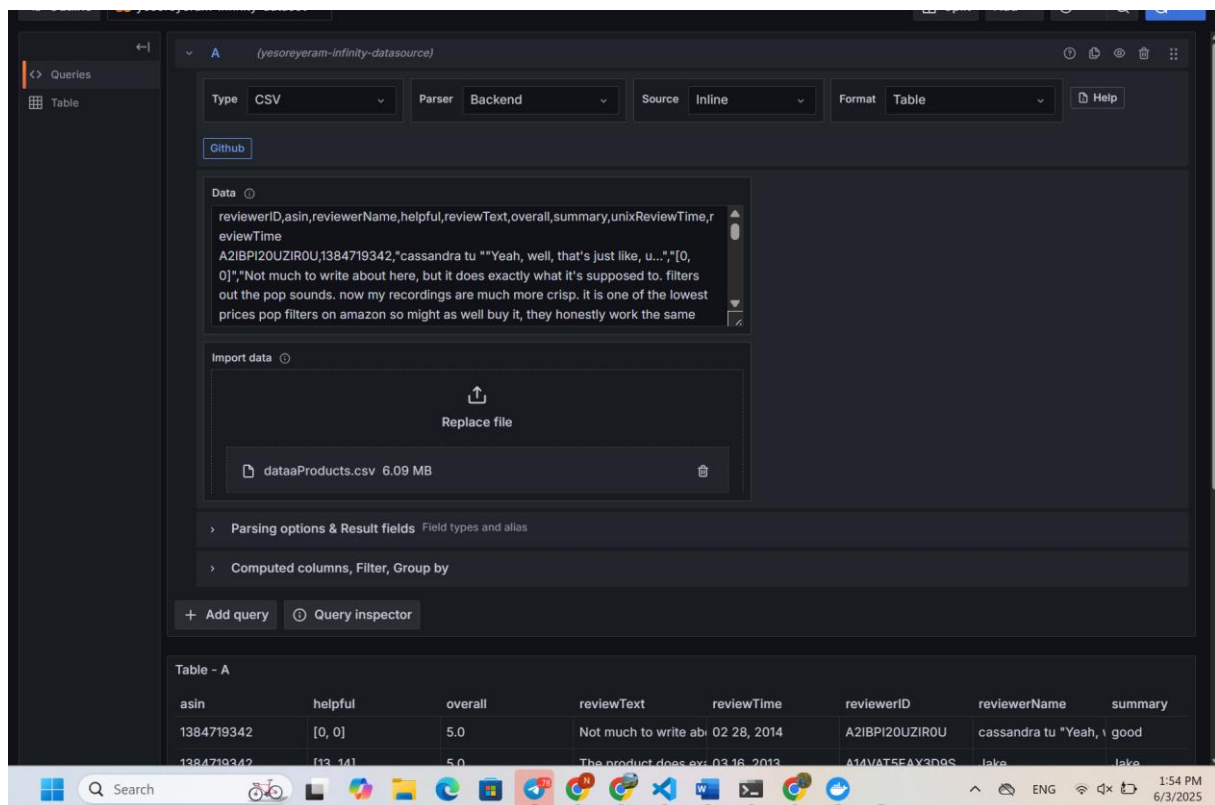
Only show running containers

<input type="checkbox"/>	Name	Container ID	Image	Port(s)	CPU (%)	Last state	Actions
<input type="checkbox"/>	minikube	b2e5e751c527	kicbase/sta	0:22 Show all ports (5)	0%	2 days ago	▶ ⋮ 🗑
<input type="checkbox"/>	gracious_poince	03a911695718	malware-ba	8000:8000	0%	15 days ago	▶ ⋮ 🗑
<input type="checkbox"/>	thirsty_knuth	df0ef497b24b	malware-ba	8000:8000	0%	15 days ago	▶ ⋮ 🗑
<input type="checkbox"/>	infallible_morse	e61f6d5549ff	malware-ba	8000:8000	0%	15 days ago	▶ ⋮ 🗑
<input type="checkbox"/>	great_germain	b0617bd65ed2	malware-ba	8000:8000	0%	15 days ago	▶ ⋮ 🗑
<input type="checkbox"/>	grafana-csv	c414ec968419	grafana/gr	3000:3000 ↗	0.78%	6 hours ago	▶ ⋮ 🗑

L'installation du **plugin Infinity** via l'interface Grafana. Ce plugin permet de **charger des fichiers CSV, JSON ou XML localement**, ce qui est parfait pour visualiser les données générées précédemment dans Spark.



Après l'analyse sur Spark, nous avons exporté les données nettoyées en CSV. Pour des raisons de simplicité et de performance locale, nous avons utilisé Grafana dans un conteneur Docker. Cela simule l'environnement de production sans avoir besoin de réinstaller tout le cluster Spark ou Ubuntu



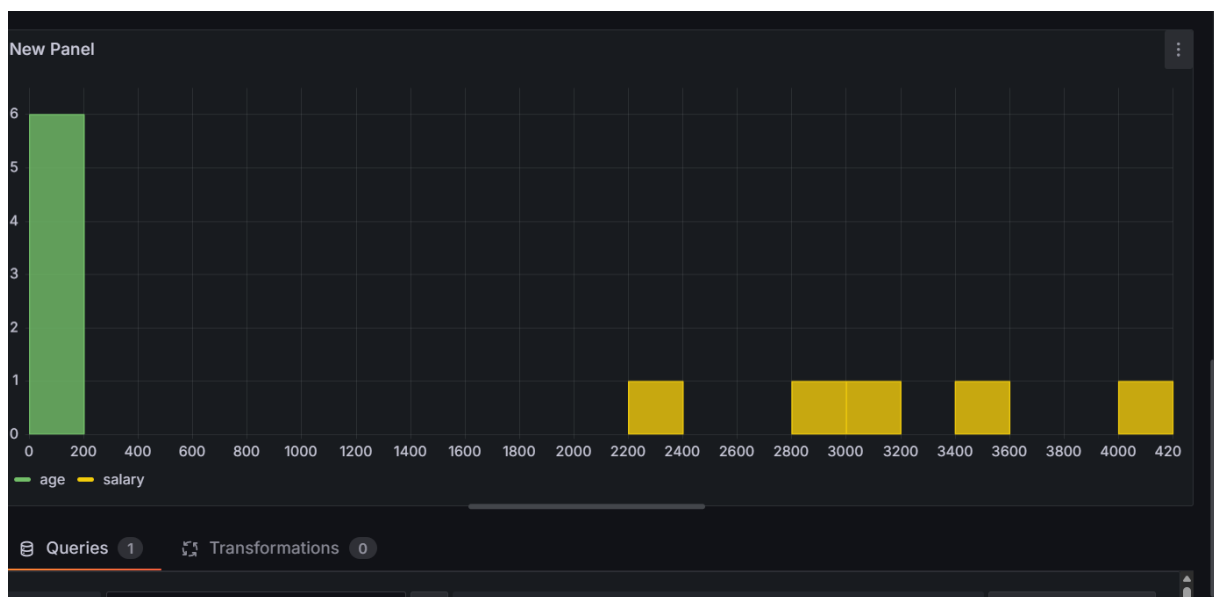
Cela permet à Grafana d'interpréter correctement chaque ligne du fichier CSV, où chaque colonne correspond à un champ comme asin, reviewText, overall, etc.

+ Add query Query inspector							
Table - A							
asin	helpful	overall	reviewText	reviewTime	reviewerID	reviewerName	summary
1384719342	[0, 0]	5.0	Not much to write ab	02 28, 2014	A2IBPI20UZIR0U	cassandra tu	"Yeah, \ good
1384719342	[13, 14]	5.0	The product does ex	03 16, 2013	A14VAT5EAX3D9S	Jake	Jake
1384719342	[1, 1]	5.0	The primary job of th	08 28, 2013	A195EZXQDW3E21	Rick Bennette	"Rick E It Does The \
1384719342	[0, 0]	5.0	Nice windscreen prot	02 14, 2014	A2C00NNG1ZQQG2	RustyBill	"Sunday Ro GOOD WIND
1384719342	[0, 0]	5.0	This pop filter is grea	02 21, 2014	A94QU4C90B1AX	SEAN MASLANKA	No more pop
B00004Y2UT	[0, 0]	5.0	So good that I boughi	12 21, 2012	A2A039TZMZH9Y	Bill Lewey	"blewey" The Best Cat
B00004Y2UT	[0, 0]	5.0	I have used monster	01 19, 2014	A1UPZM995ZAH90	Brian	Monster Star
B00004Y2UT	[0, 0]	3.0	I now use this cable t	11 16, 2012	AJNFQI3YR6XJ5	Fender Guy	"Rick" Didn't fit my
B00004Y2UT	[0, 0]	5.0	Perfect for my Epiphc	07 6, 2008	A3M1PLEYNDEY08	G. Thomas	"Tom" Great cable
B00004Y2UT	[0, 0]	5.0	Monster makes the b	01 8, 2014	AMNTZU1YQN1TH	Kurt Robair	Best Instrum
B00004Y2UT	[6, 6]	5.0	Monster makes a wid	04 19, 2012	A2NYK9KWF MJV4Y	Mike Tarrani	"Jazz Dr One of the br
B00005ML71	[0, 0]	4.0	I got it to have it if I n	04 22, 2014	A35QFQI0M46LWO	Christopher C	It works grea

Cet affichage valide que le fichier CSV a été correctement chargé et structuré dans Grafana, prêt à être utilisé pour la création de panels (graphes, histogrammes, etc.).



Chaque groupe de barres correspond à un pays (USA, Canada, UK...). Ce type de graphique est utile pour comparer visuellement deux métriques par catégorie. On observe par exemple que les individus les plus âgés (42 ans au Royaume-Uni) ne sont pas nécessairement ceux ayant les salaires les plus élevés (4000 aux États-Unis).



Ce second graphique est un **histogramme** affichant la **distribution des salaires** présents dans le jeu de données. On remarque une forte concentration dans la tranche basse (autour de 200 à 1000), tandis que les tranches supérieures (jusqu'à 4000) contiennent moins de données.

Ce type de visualisation permet de comprendre rapidement la dispersion des valeurs numériques dans un champ continu comme le salaire.

Table view | Last 1 hour | Refresh

New Panel

age	country	name	occupation	salary
38	USA	Leanne Graham	Devops Engineer	3000
27	USA	Ervin Howell	Software Engineer	2300
17	Canada	Clementine Bauch	Student	
42	UK	Patricia Lebsack	Software Engineer	2800
38	USA	Leanne Bell	Senior Software Engineer	4000
32	USA	Chelsey Dietrich	Software Engineer	3500

Queries 1 | Transformations 0

Data source: yesoreyeram-infinity-datasot | Query options: MD = auto = 666 | Interval = 5s | Query inspector

Visualization: Table

Panel options

Title: New Panel

Description:

Transparent background: ☐

Panel links:

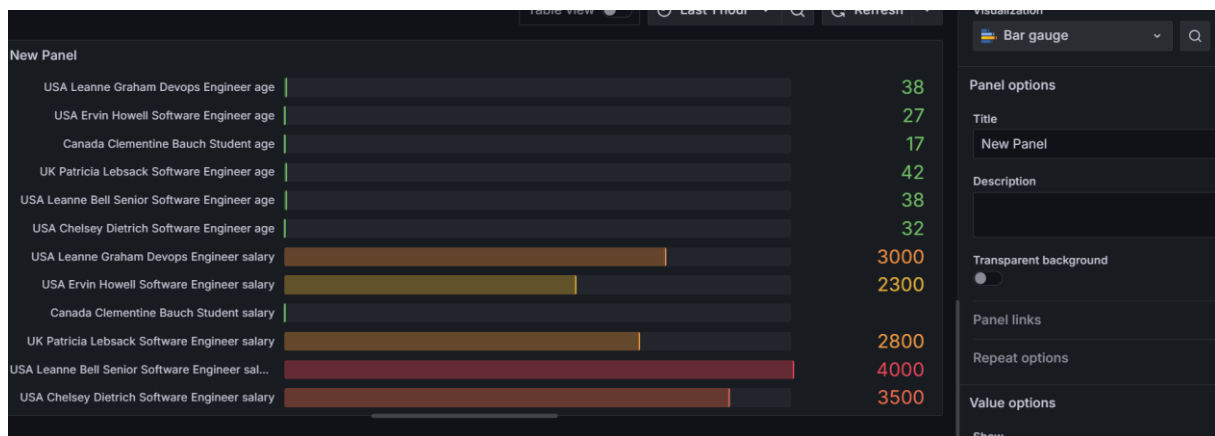
Repeat options:

Table

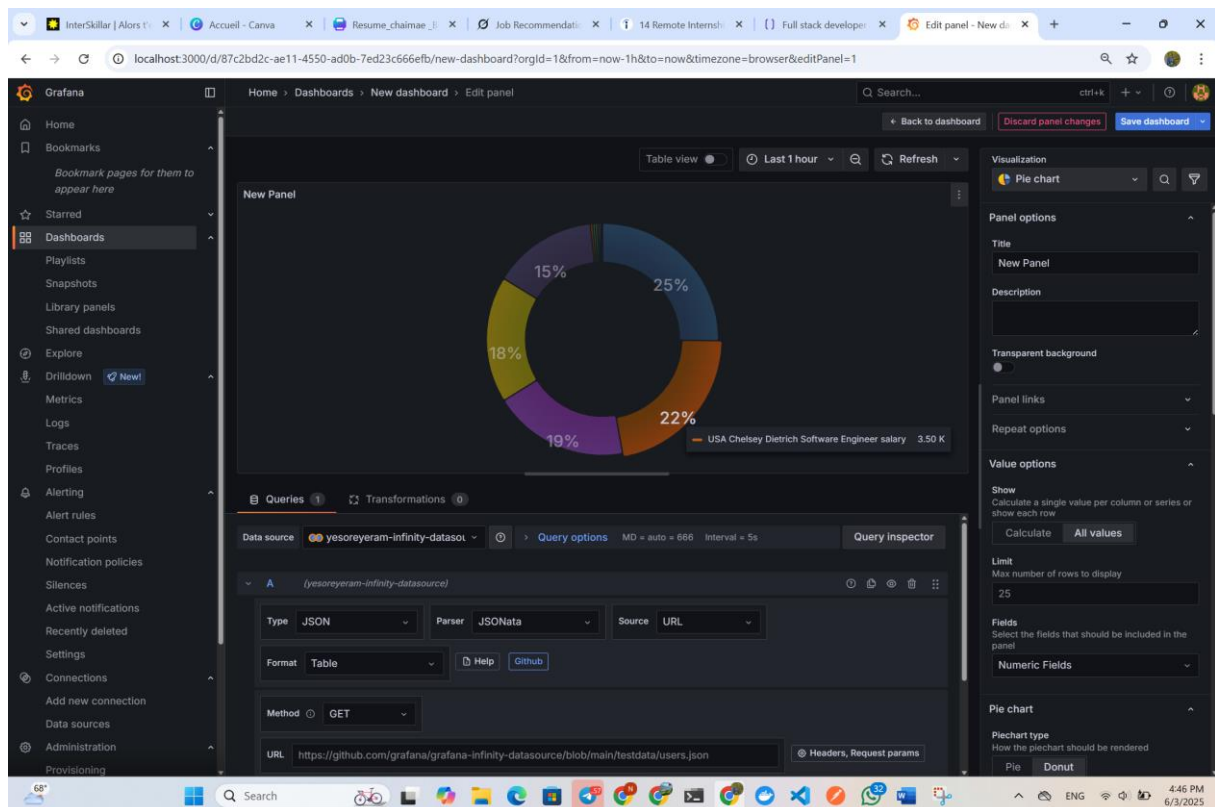
Show table header: ☒

Cell height: Small Medium Large

Enable pagination: ☐



L'objectif de cette visualisation est de **comparer facilement les deux variables** au sein d'une même ligne, pour chaque profil de la base. Par exemple, *Leanne Bell* aux USA a un âge de 38 ans et un salaire de 4000, tandis que *Clementine Bauch*, au Canada, a 17 ans et un salaire de 2800. Ce type de graphique est particulièrement lisible pour des jeux de données de faible volume.



Ce **diagramme circulaire** représente la **répartition relative** des différents profils, ici selon une variable agrégée (salaire ou âge) par personne. Chaque part du cercle correspond à un utilisateur, et la **proportion en pourcentage** permet de voir les écarts facilement.

Le diagramme montre par exemple que certains contributeurs (ex. *Leanne Graham* ou *Chelsey Dietrich*) occupent une part importante de l'ensemble, ce qui peut indiquer une **concentration des valeurs élevées** (âge/salaire).

CONCLUSION

Ce projet a permis de mettre en œuvre une chaîne complète de traitement Big Data, depuis l'ingestion et le stockage de données jusqu'à l'analyse et la visualisation. À travers les différentes phases, nous avons manipulé des outils et technologies modernes comme HDFS, Apache Spark, Docker, et Grafana, dans un environnement distribué simulé sous Ubuntu Server Master/Worker.

Dans un premier temps, nous avons préparé l'environnement Hadoop/Spark, créé les répertoires HDFS, et intégré un jeu de données volumineux au format JSON. Ensuite, à l'aide de PySpark, nous avons effectué des traitements de nettoyage, des agrégations statistiques, et des extractions de colonnes utiles.

Face à la perte des machines virtuelles en fin de parcours, une solution alternative a été mise en place via Docker et Grafana sur Windows, avec succès. Grâce au connecteur Infinity Data Source, nous avons pu intégrer nos fichiers .csv traités dans Grafana et générer des visualisations dynamiques : histogrammes, camemberts, heatmaps, et bar charts.

Ce projet a ainsi démontré la capacité à :

- déployer un cluster Hadoop/Spark local,
- automatiser des traitements de données volumineuses,
- convertir et visualiser des données avec des outils modernes open source.

En somme, ce travail a illustré concrètement le cycle ETL + Analyse + Visualisation dans un contexte Big Data, avec une approche pratique et résiliente face aux contraintes techniques. Il a permis de renforcer nos compétences techniques en traitement distribué, analyse de données, et monitoring via dashboard.