



# Sommaire:



- 1 Contexte et problématique
- 2 Jeu de donnée
- 3 Big Data
- 4 Architecture sur le cloud
- 5 Création de l'Environnement Big data
- 6 Chaine de traitement SparkUI
- 7 Démonstration sur le cloud
- 8 Conclusion



# 1-Contexte et problématique

La start-up AgriTech, nommée "Fruits!", cherche à proposer des solutions innovantes pour la récolte des fruits, tout en préservant la biodiversité de ces derniers.

L'objectif est de permettre des traitements spécifiques pour chaque espèce de fruits en développant des robots cueilleurs intelligents.

- Se faire connaître en mettant à disposition du grand public une application mobile qui permettrait aux utilisateurs de prendre en photo un fruit et d'obtenir des informations sur ce fruit.
- Cette application permettre de:
  - □ sensibiliser le grand public à la biodiversité des fruits et de mettre en place une première version du moteur de classification des images de fruits.
  - ☐ construire une première version de l'architecture Big Data nécessaire.

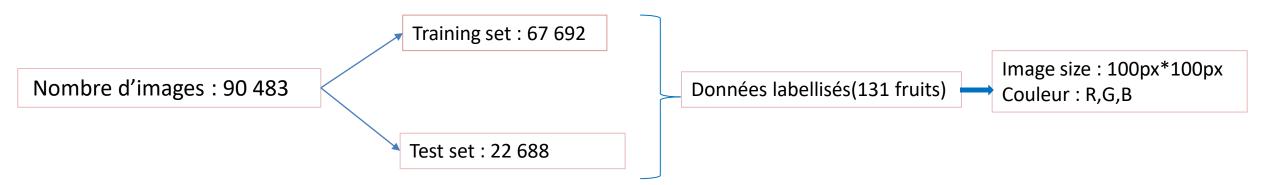




## 2- Jeu de données



Jeu de données **Fruits-360** : un ensemble de données d'images de format jpg contenant des fruits en plusieurs variété et plusieurs espèces.



Multi-fruits: 103 images.





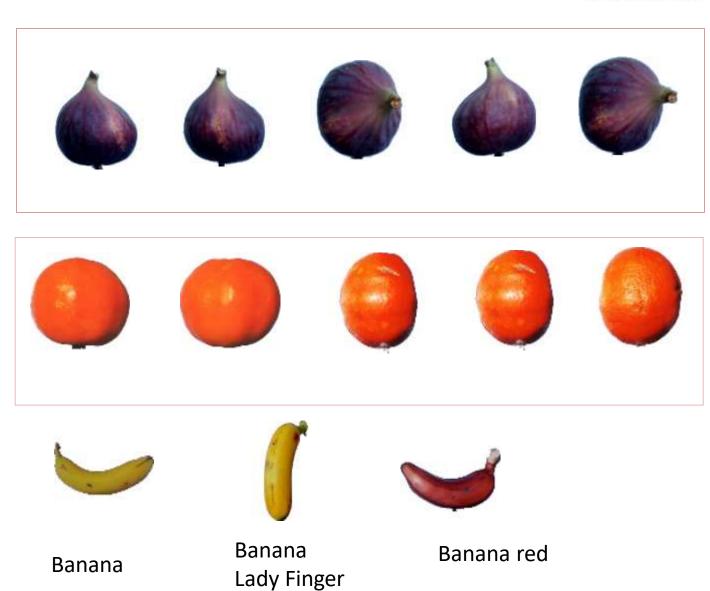


# 2-Jeu de données:



- Photographiés à plusieurs angles.
- Fond d'images éliminé en blanc.
- Plusieurs variétés de fruits .





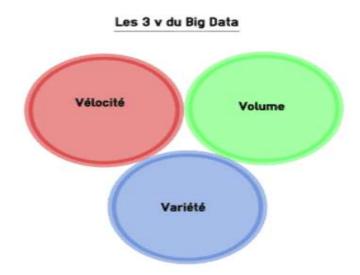


### 3- BIG DATA:

Le terme "Big Data" fait référence à des ensembles de données volumineux, complexes et variés qui dépassent les capacités des méthodes traditionnelles de gestion et d'analyse de données.

Le concept de Big Data est souvent caractérisé par les trois V suivants :

- Volume: La taille gigantesque de données.
- Variété: La diversité de types de données.
- Vélocité: La rapidité à laquelle de nouvelles données sont générées.



# 4-les différentes briques d'architecture choisies sur le cloud:

#### 4.1- Prestataire cloud choisi:

**AWS** (Amazon Web Services) : une plateforme de services cloud proposée par Amazon. Il offre une large gamme de services, y compris le calcul, le stockage, la base de données, l'analyse, l'apprentissage automatique, la sécurité, la gestion des identités, le réseau, le déploiement d'applications



#### 4.2-IAAS ou PAAS?

Solution IAAS:(Infrastructure as a Service)

- Les serveurs sont vierges.
- Installation de tous les outils pour notre script (Spark, Java, Python, Jupyter...) par nous-même.
- Inconvénient: Chronophage.

Solution PAAS: (Platform as a Service)

- Les Framework importants sont déjà installés.
- Les patchs de sécurité sont automatiquement mis à jour
- Faciliter de clonage ses clusters.

> Pour ce projet, la solution PaaS est choisie.

## 4-les différentes briques d'architecture choisies sur le cloud:

#### EC2: Elastic Compute Cloud

Des serveurs virtuels dans le cloud.

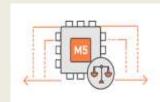


L'instance m5.large est un type d'instance de aws.

•Taille de l'Instance : x.Large.

•vCPU : 4 unités de calcul virtuelles.

•Mémoire: Environ 16Go de RAM.



S3 (Simple Storage Service).

 permet de stocker des objets, tels que des fichiers et des données



IAM (Identity and Access Management) est le service de gestion des identités et des accès d'Amazon Web Services (AWS).



### Spark:

Spark : est un framework <u>open-source</u>. conçu pour fournir une plateforme rapide et unifiée pour l'analyse de données à grande échelle.

- Traitement en Mémoire
- ➤ APIs Polyglottes
- > Traitement de Données Distribué
- Flexibilité dans le Stockage des Données

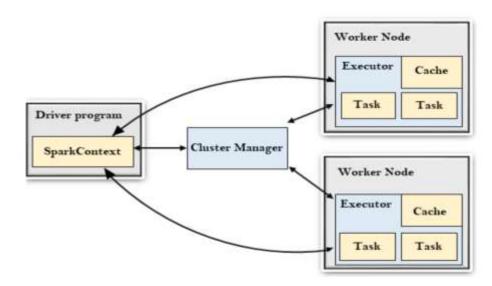
### Spark architecture:

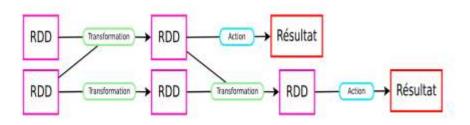
- Les instructions sont envoyées au driver lorsqu'on interroge Spark grâce à PySpark.
- ➤ La communication avec le driver s'effectue grâce au SparkContext.
- Le gestionnaire de cluster instancie les différents workers.
- Les workers instancient un exécuteur.
- Les exécuteurs sont chargés d'exécuter les différentes tâches de calcul.

### Spark RDD (Resilient Distributed Datasets):

#### **Tolérance aux pannes:**

- transformations et les actions réalisées sur les RDD permettent de construire un graphe acyclique orienté (DAG).
- Les connexions entre les nœuds sont soit des transformations, soit des actions.
- Le graphe est dit *acyclique* car aucun RDD ne permet de se transformer en luimême via une série d'actions.
- Un RDD peut être régénéré à partir de ses RDD parents.





graphe acyclique orienté (DAG: "directed acyclic graph«)

### 5-Création de l'Environnement Big data (IAM,S3, EC2):

# 1-) Amazon IAM (Identity and Access Management) - Autorisations Utilisateur

- permet de gérer l'accès aux ressources AWS de manière sécurisée
- Autorisations tout type de manipulation dans les buckets S3
   2-Amazon IAM Configuration AWS:

```
(database) LHadjemi@HPE2109P026:-$ aws configure

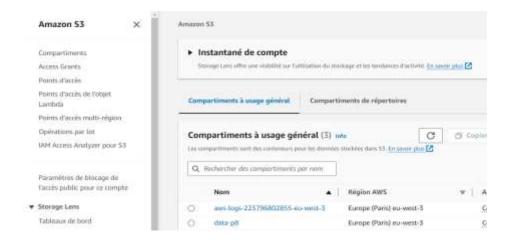
AWS Access Key ID [***
AWS Secret Access Key

Default region name [eu-west-3]:
Default output format [json]:
(database) LHadjemi@HPE2109P026:-$
```



Paire de clés: Permet de sécuriser l'accès à la connexion Utilisateur eu-west-3: Serveurs à Paris -> Respecte le RGPD RGPD: Règlement Général sur la Protection des Données Nous oblige à utiliser des serveurs en Europe uniquement

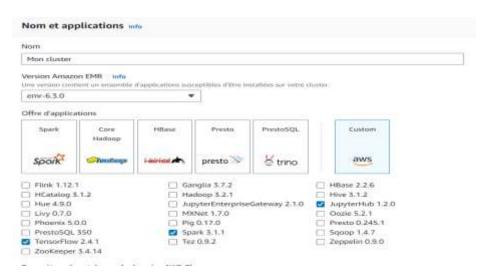
#### 3-Amazon S3:



#### 5-Amazon EMR Boostrapping:

- wheel : Accélère l'installation des packages Python.
- pillow : Bibliothèque de traitement d'image en Python.
- pyarrow : Lit les fichiers parquets et les convertit en Dataframe.
- boto3 + s3fs + fsspec : Permettent d'interagir avec Amazon S3.
- Pandas : Pour les differentes manipulations des dataframes.

#### 4-Amazon EMR:



#### 6-Amazon EMR choix d'instances:

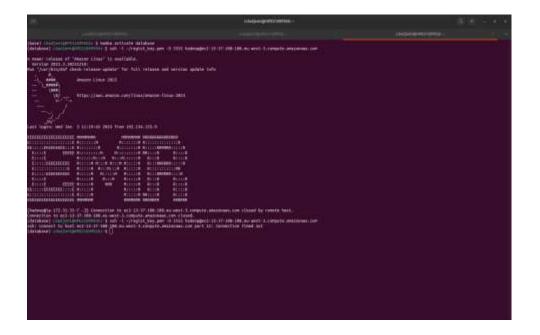
#### Groupes d'instances uniformes

#### Primaire Choisir un type d'instance EC2 m5.xlarge Actions \ 4 vCore 16 GIB mémoire EBS uniquement stockage Prix à la demande : 0.224 USD par instance/heure Prix Spot le plus bas : 0.068 USD (eu-west-3h) Utiliser la haute disponibilité Lancez des clusters hautement disponibles et plus résilients avec trois nœuds primaires sur des instances à la demande. configuration s'applique pendant toute la durée de vie de votre cluster. En savoir plus [2] ▶ Configuration de nœud - facultatif Unité principale Choisir un type d'instance EC2 m5.xlarge Actions \* 4 vCore 16 GiB mémoire EBS uniquement stockage Prix à la demande : 0.224 USD par irestance/heure Prix Spot le plus bas : 0.068 USD (eu-west-3b)

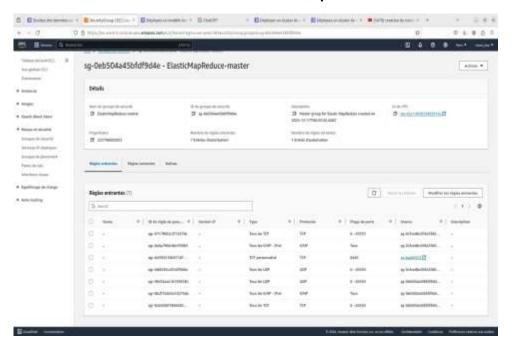
#### 7- Création des tunnel SSH:

Pour pouvoir s'éxecuter et ouvrir le JupyterHub, la création des clés ssh est indispensable.

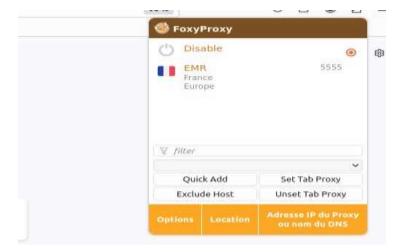
- On établit la connexion au bash une fois le clusters créer.



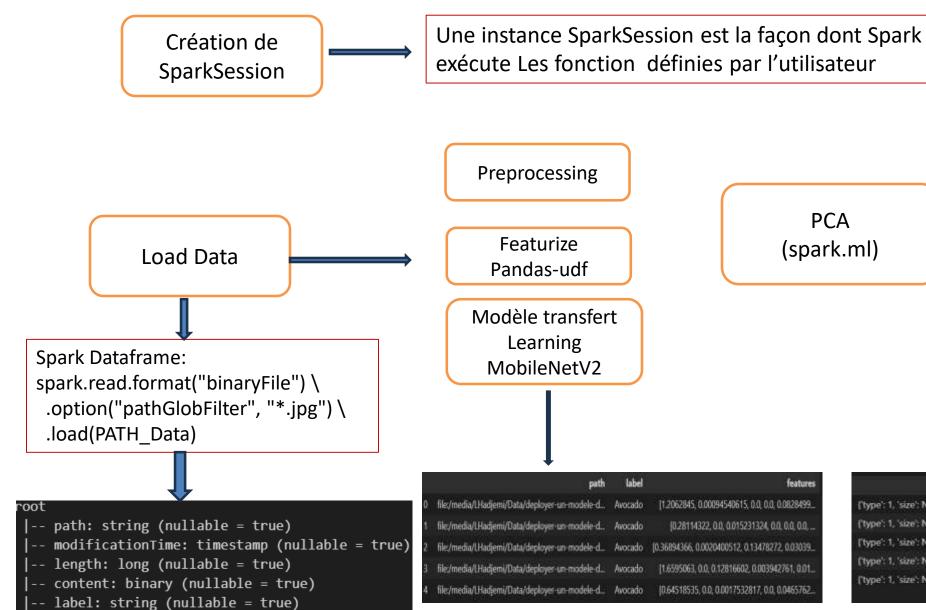
- Autorisation les tunnel ssh port 5555 sur aws.



- Configuration de l'extension Foxyproxy pour le navigateur.



### 6-Chaine de traitement Spark en local:



**PCA** (spark.ml) Ecriture et puis lecture vers pd.DataFrame au format paquet, engine pyarrow

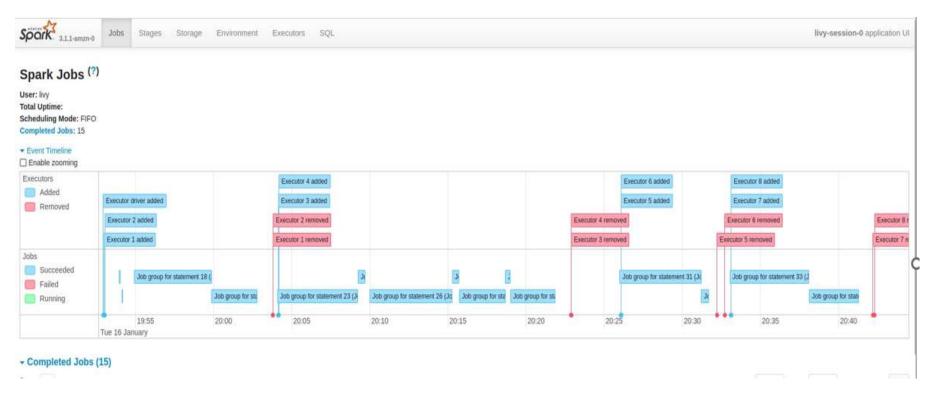
pca_vectors	features_array
('type': 1, 'size': None, 'indices': None, 'va	[-16.95106, -4.3774457, -11.170761, -10.679921
('type': 1, 'size': None, 'Indices': None, 'va	[-17.011036, -13.817756, -10.5200405, 2.280178_
('type': 1, 'size': None, 'indices': None, 'va	[-15.204871, 0.99675465, -8.867973, -0.4300575
('type': 1, 'size': None, 'indices': None, 'va	[-15.122531, 2.5644705, -9.907209, -1.5638206,
('type': 1, 'size': None, 'indices': None, 'va	[-1.7547172, -2.3506665, -12.355202, -4.484647_

### 6-Chaine de traitement SparkUI en local:



- 1-Activation des machine de clusters (création de driver)
- 2-Chargement des données images.
- 3-La conversion des images au format binaires.
- 4-Application du model mobilenetv2
- 5-Données au format paquet.
- 6-lecture de tableau pandas
- 6-chargement de PCA
- 7- Application de PCA sur les features.
- 8-chargement des données au format paquet

### 7-Chaine de traitement SparkUI sur le cloud :



### 7-1-Comparaison du temps d'exécution:

Nombres d'instances	Temps d'exécution
1 primaire (m5.xlarge) 2 unité principale (m5.xlarge)	3039.4540 Seconde
1 primaire (m5.xlarge) 4 unité principale (m5.xlarge)	1508,6541 Seconde

### 7-Chaine de traitement SparkUI sur le cloud :

#### 7-2- Résultats finaux:

```
data final.show(5, True)
   FloatProgress(value=0.0, bar style='info', description='Progress:', layout=Layout(height='25px', width='50%'),...
                             path
                                                  label
                                                                        cnn vectors
                                                                                                       pca vectors!
                                                                                                                                 features array
     s3://data-p8/Test...|Pineapple Mini|[0.0,3.6077606678...|[-6.0440096337977...|[-6.0440097, 2.89...
                                          Watermelon | [0.16857005655765... | [0.32359665142388... | [0.32359666, 4.18...
     s3://data-p8/Test...
     s3://data-p8/Test...|Pineapple Mini|[0.0,4.9579777717...|[-5.7261558602160...|[-5.7261558, 3.45...
     s3://data-p8/Test...|Pineapple Mini|[0.0,4.8031039237...|[-4.9360270642147...|[-4.936027, 2.071...
                                          Watermelon|[0.02986907027661...|[-3.0326533321854...|[-3.0326533, 3.54...
     s3://data-p8/Test...
   only showing top 5 rows
                                                                           Amazon 53
                                                                                                              Propriétés
                                                                           Compartiments
                                                                           Access Grants
                                                                                                       Objets (21) into
                                                                           Points d'accès
                                                                                                            Copier TURESS
                                                                                                                             Copier TURL
                                                                                                                                                                                         Créer un dossier.
                                                                           Points d'accès de l'objet.
                                                                           Lambda
7-3- Données sur S3:
                                                                           Points d'accès multi-région
                                                                                                       Les objets sont les entités fondamentales stockées dans Amazon S3. Vous pouvez utiliser l'inventaire Amazon S3 🔀 gour obterior une liste de trois les objets de votre compartiment. Pour que d'autres
                                                                                                       personnes pursonnt accider à vos objets, vous devez leur accorder explicitement des autorisations. En savoir plus
                                                                           Opérations par lot
                                                                           IAM Access Analyzer pour $3
                                                                                                        Q. Rechercher des objets en fonction du préfit
                                                                                                                                                                                              (1) 0
                                                                                                                                                 Dernière modification P Taille
                                                                                                                                                                                     ♥ Classe de stockage
                                                                           Paramétres de blocage de
                                                                           l'accès public pour ce compte
                                                                                                                                                  16 Jan 2024 09:41:17 PM
                                                                                                           SUCCESS.
                                                                                                                                                                                        Standard
                                                                         ▼ Storage Lens
                                                                                                            part-00000-fe6ca43b-
                                                                           Tableaux de bord
                                                                                                            71a8-42ab-8475-
                                                                                                                                                  16 Jan 2024 09:38:51 PM
                                                                                                                                                                                 6.3 Mo
                                                                                                                                                                                        Standard
```

Groupes Storage Lens

Fonctionnalité spot

Paramètres AWS Organizations

parquet

parquet

16 Jan 2024 09:38:58 PM

6.3 Mo

Standard

2de5be16c5de-£000 anappy parquet.

2de5be16c5de-

P part-00001-fe6ca43b-71a8-42ab-6475-

# 7-5-Démonstration sur le cloud aws

Démonstration sur le cloud aws

# **Conclusion:**

- Prise en main de Pyspark et découverte de plusieurs approche comme format paquet.
- Nombreuses possibilités techniques, choix d'instances, et choix de version adapté a EMR.
- AWS (IAM + S3 + EMR) nous propose un service big data complet.
- SparkUI nous permet de visualiser les détails des calculs distribués