



# **PROJET 8 :**

## **CALCULS DISTRIBUÉS ET DÉPLOIEMENT D'UN MODÈLE DANS LE CLOUD**

# PLAN DE LA PRÉSENTATION

1) Problématiques métier

2) Présentation des données

3) Solutions retenues

4) Traitement des images avec PySpark en local

5) Traitement des images avec PySpark dans le cloud (AWS EMR)

6) Conclusion

7) Perspectives

# PROBLÉMATIQUES DE LA START-UP "FRUITS!"

## Contexte :

- start-up qui veut mettre en place un moteur de classification d'images de fruits.
- construction d'une première version de son architecture Big Data.
- utilisation de données disponibles sur Kaggle comme point de départ pour la POC.

## Mission :

- développer dans un environnement Big Data une première chaîne de traitement des données : preprocessing et réduction de dimension.

## Contraintes :

- anticiper une augmentation très rapide du volume de données après la livraison de ce projet.
- développer des scripts en PySpark.

## Suggestions :

- utiliser le cloud AWS pour profiter d'une architecture Big Data.



Fruits!

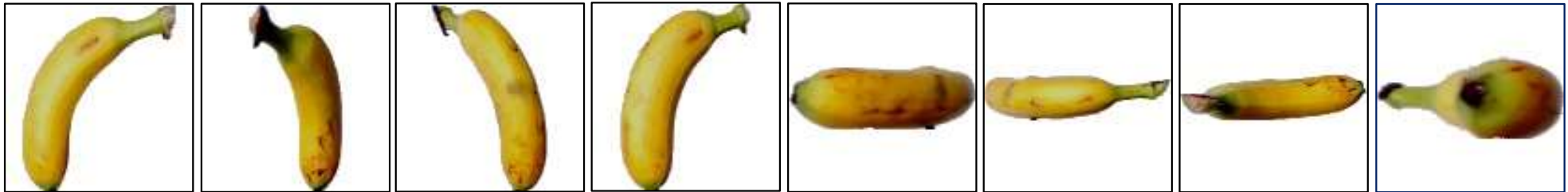
# PRÉSENTATION DES DONNÉES

**Disponibles sur Kaggle :** <https://www.kaggle.com/moltean/fruits>

**~90000 images :**

- 131 fruits et légumes différents,
- photos en studio avec fond blanc,
- sous différents angles,
- dimensions : 100 x 100 pixels x 3 (RGB),
- format JPEG,

**Exemple :**



**Données utilisées pour la POC de notre projet :**

- images de 10 fruits (10 dossiers) différents, soit ~5000 images.
- accès via un bucket S3.

# SOLUTIONS RETENUES

## Approaches envisagées pour générer des descripteurs :

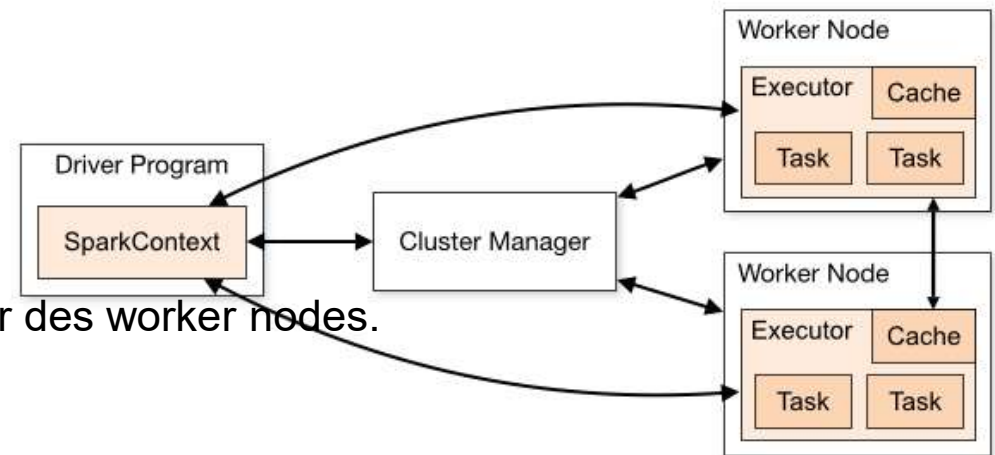
1- Techniques classiques de vision par ordinateur pour générer des descripteurs (SIFT, ORB, SURF) puis recherche de n clusters pour générer des features à n dimensions pour chaque image.

2- Transfer learning avec un réseau de neurones préentraîné.

→ Préférence pour l'approche 2 (qualité des features, vitesse de calcul, simplicité). **Choix de VGG16.**

## Apache Spark pour l'architecture de calcul Big Data :

- Tâches de calculs :
  - distribuées entre executors qui se trouvent sur des worker nodes.
  - réalisables en parallèle.
  - tolérance aux pannes.
- Driver program : le programme principal. Fera appel à des exécuteurs.
- Executor : process qui exécute les tâches de calcul distribué demandées par le driver.
- Cluster manager : gère les ressources du cluster, notamment les interactions entre driver et executor.
- PySpark : API Python pour Spark (langage natif : Scala).



# EXTRACTION DE DESCRIPTEURS PAR VGG-16

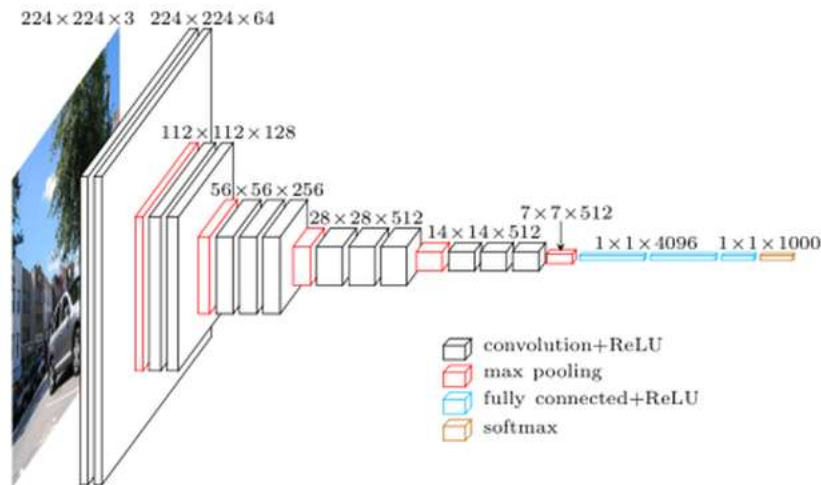
## VGG-16 :

Réseau de neurones convolutif à 16 couches pour la classification d'images.

Réseau pré-entraîné sur 1.3 millions d'images.

Possibilité de l'utiliser en transfer learning.

Dernière couche : classifieur softmax qui prend un vecteur de dimension 4096 en entrée.



## Génération du descripteur (4096 dimensions) d'une image :

- 1) Preprocessing de l'image (en particulier redimensionnement au format 224\*224\*3 pixels)
- 2) Calcul du descripteur via le modèle préentraîné.

# TRAITEMENT DES IMAGES AVEC PYSPARK

Procédé mis en œuvre, pour chaque image:



Map/Reduce sur toutes les images :

```
94 # broadcasting
95 bc_model_vgg16 = sc.broadcast(model_vgg16)
96
97 # job spark
98 resultat = sc.parallelize(megabatch_img) \
99     .map(lambda img: VGG16_extracteur('S3', img, bc_model_vgg16.value)) \
100     .collect()
```

Export des données collectées vers un .csv :

	dim_0	dim_1	dim_2	...	dim_4093	dim_4094	dim_4095	path
0	0	0	0	...	0.258413106	1.812143803	0	Apple Braeburn/0_100.jpg
1	0	0	0	...	0	1.141268373	0	Apple Braeburn/100_100.jpg
2	0	0	0	...	0.117657334	0.927160382	0	Apple Braeburn/101_100.jpg
3	0	0	0	...	0.293886185	0.592740119	0	Apple Braeburn/102_100.jpg
...	...	...	...	...	...	...	...	...

# TRAITEMENT DES IMAGES SUR UN PC LINUX EN LOCAL AVEC PYSPARK

## Résultats :

- fonctionnement de Spark en `--deploy-mode local` différent de celui observé en `--deploy-mode cluster`
- contrairement au mode **cluster (ressources en conteneur)**, en local les CPUs, la RAM et la Java VM sont partagés entre tous les exécuteurs.
- **résultats non significatifs par rapport à un cluster dans le cloud** avec des images stockées à proximité.
- en local, en se limitant à un seul exécuteur → **les temps d'I/O (S3)** sont supérieurs aux temps de calcul.
- en local, en augmentant le nb d'exécuteurs → performances boostées par partage des CPUs entre tâches Spark (pas possible en cluster).
- utilisation de **Spark session** pour transformer les données → lent et fréquentes erreurs de la Java VM.

## Conclusion :

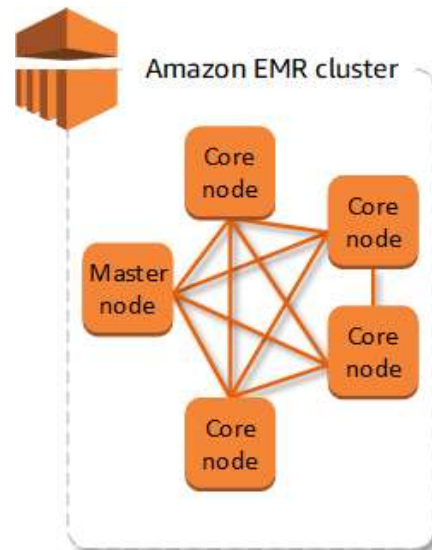
- mode local utile pour **mettre au point le script PySpark**.
- **Spark session non retenue** pour le traitement dans le cloud.
- **tuning** des paramètres Spark **non extrapolable** au cluster.



# TRAITEMENT DES IMAGES DANS LE CLOUD (AWS EMR) AVEC PYSPARK

## Présentation d'EMR (Elastic MapReduce) :

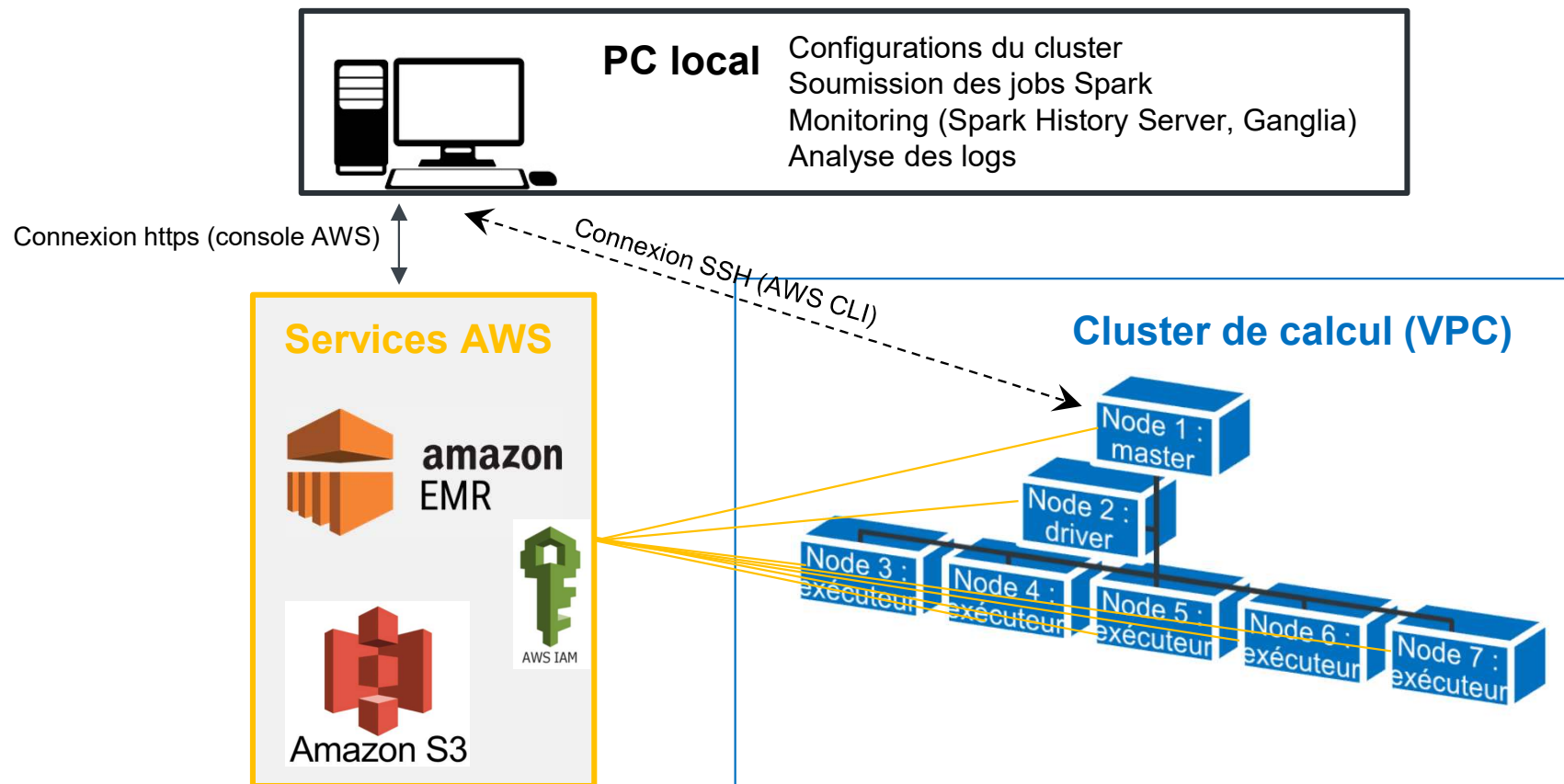
- plateforme de services pour le calcul distribué.
- simplifie l'utilisation des infrastructures de données massives, telles que Hadoop et Spark.
- évolutif à travers des instances EC2 dans le cloud AWS.



- cluster (composant principal) = ensemble de nodes (serveurs EC2), chacun ayant un rôle précis :
  - master node : effectue le suivi du statut des tâches et surveille l'état du cluster.
  - core node : exécutent les tâches dans un système de fichiers distribués (HDFS).
  - task node : facultatif, exécutent des tâches et ne stockent pas les données dans HDFS.

# TRAITEMENT DES IMAGES DANS LE CLOUD (AWS EMR) AVEC PYSPARK

**Cluster retenu :** 7 instances EC2 m5.xlarge (compromis entre : RAM / vCPUs / coût / limitations de mon compte AWS)



# TRAITEMENT DES IMAGES DANS LE CLOUD (AWS EMR) AVEC PYSPARK

## Configuration finale du cluster :

- EMR release 6.5.0.
- 7 instances EC2 m5.xlarge (7\*4 vCPUs, 7\*16 Go RAM, 7\*64Go SSD). 1 master et 6 cores.
- Bootstraps : « pip install » boto3, pandas==1.2.5, tensorflow==2.7.0, Pillow
- Softwares pré-chargés : Hadoop v3.2.1, Spark 3.1.2, Ganglia 3.7.2.

## Principaux écueils rencontrés avec EMR :

- Configuration initiale d'un cluster (certaines configurations par défaut proposées par AWS plantent !).
- Disponibilité des instances m5.xlarge (eu-west-1 > eu-west-3).

## Principaux debuggages pour faire passer le code du local vers un cluster EMR :

- Affecter 10Go au conteneur du driver (au lieu des 2Go par défaut) pour lancer VGG16.
- Versionning des modules Python.
- Sérialisation du modèle VGG16 (*TypeError: can't pickle weakref objects*) : utiliser tensorflow >= 2.7.0

## TRAITEMENT DES IMAGES DANS LE CLOUD (AWS EMR) AVEC PYSPARK

### Métriques du job Spark sur mini-clusters :

Worker nodes → Exécuteurs	vCPUs & RAM totale, par exécuteur	Nb partitions du RDD	Nb images	Durée job (s)	Commentaire
1 → 1	4 vCPUs & 16 Go	4 (défaut)	64	30	17s d'overhead au démarrage du job, puis 13s pour 64 images
			128	43	
		8 (forcé)	128	43	Pas d'impact temps du nb de partition
2 → 2		8 (défaut)	64	24	17s d'overhead au démarrage du job, puis 7s pour 64 images
			128	31	
2 → 4	2 vCPUs & 8 Go		128	32 à 52	« 6 tasks failed » en 7 jobs

→ Prévoir **16Go par exécuteur** pour éviter des « tasks failed »

## TRAITEMENT DES IMAGES DANS LE CLOUD (AWS EMR) AVEC PYSPARK

**Job final (montée en échelle sur 4987 images) :**

Worker nodes → Exécuteurs	vCPUs & RAM totale, par exécuteur	Nb images	Nb partitions du RDD	Durée job (s)	Commentaire
5 → 5	4 vCPUs & 16 Go	4987	20 (défaut)	224	2 répétitions (0 task failed)

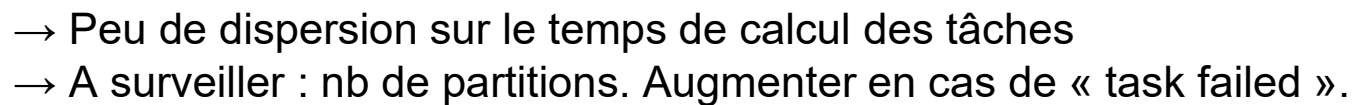
→ On reste **constant sur la vitesse de calcul par vCPU d'exécuteur** :

- 4100 images / (h.vCPU) sur le job du slide précédent (7s pour 64 images avec 8 vCPUs)
- 4400 images / (h.vCPU) sur le job de 4987 images en défalquant l'overhead de démarrage.

→ Augmentation du temps de calcul cohérente avec **variation linéaire** :  
**f(nb d'images, 1/(nb de vCPUs exécuteur))**

→ Taille du .csv de sortie du job : 161Mo.

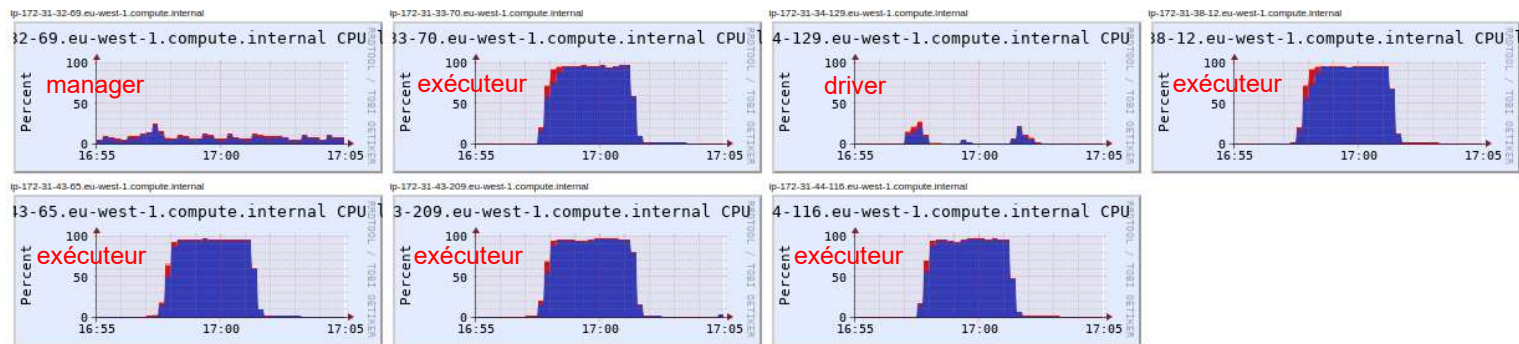
## Job final (4987 images) - Spark History Server - event timeline :



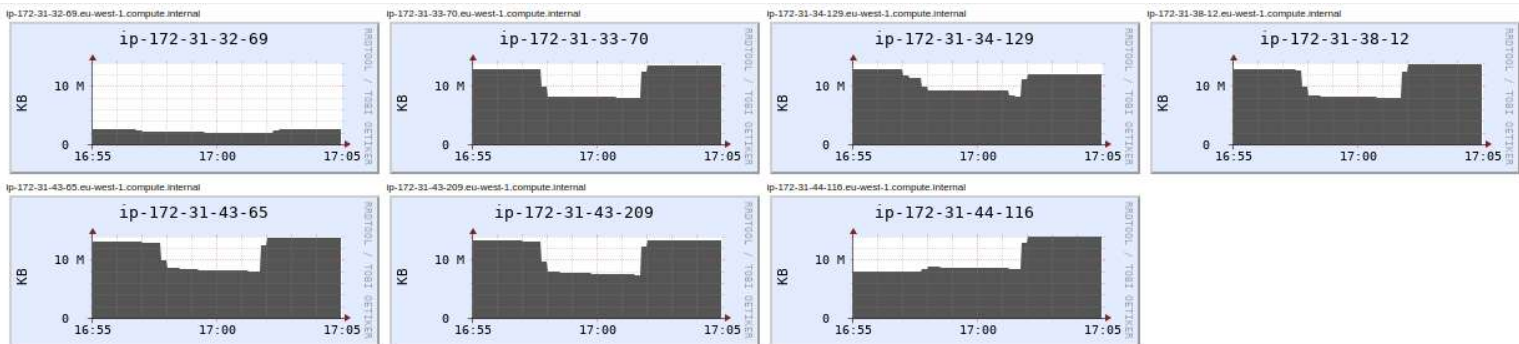
# TRAITEMENT DES IMAGES DANS LE CLOUD (AWS EMR) AVEC PYSPARK

## Job final (4987 images) - monitoring Ganglia :

Charge des CPUs :



Charge des RAMs :



- Possibilité de fusionner le manager et le driver sur un seul node.
- CPUs des exécuteurs au max durant les jobs (peu d'idle dû aux I/O) → envisager p2.
- RAM des exécuteurs : bien dimensionnée.
- RAS sur le monitoring des I/O réseau.

# CONCLUSION

## **POC** obtenue :

- sur ~5000 images,
- avec un algorithme d'extraction basé sur VGG-16,
- dans le cloud AWS EMR,
- avec le framework Spark.

**Augmentation du temps de calcul linéaire** lors de la montée en échelle  $f(\text{image}, 1/(\text{vCPUs}))$

## **Vitesse de calcul par vCPU d'exécuteur sur instance m5.xlarge :**

- ~4400 images / (h.vCPU)
- 1.000.000 images pour 11€ (0.2€/h pour m5.xlarge).

## **Sur un plan personnel :**

- acquisition de nouvelles compétences : Spark, PySpark, Boto3, AWS (en particulier S3, EMR).
- amélioration des compétences : Linux (tunnel SSH), Java VM.



# PERSPECTIVES

## **Accélérer les calculs :**

- 1) par le hardware : instances de calcul mieux adaptées telles que EC2 p2 (GPU).
- 2) par le code : solutions envisagées à l'instanciation de boto3 à chaque image :
  - `spark.read.format('Image')`,
  - pré-télécharger les images depuis S3 vers le VPC (task node de stockage).

## **Fiabiliser les calculs :**

augmenter le nb de partitions du RDD.

## **Points durs pour une montée en échelle :**

taille du fichier .csv de sortie trop gros si unique → générer un fichier .csv par fruit.

## **Tester images plus grande que 100\*100 pixels :**

code et cluster conçus pour des images de taille plus grande.



# SCRIPTS FINAUX

## Création du cluster depuis AWS CLI :

```
aws emr create-cluster --applications Name=Hadoop Name=Hive Name=Spark Name=Ganglia --ec2-attributes '{"KeyName":"cle_irelande","InstanceProfile":"EMR_EC2_DefaultRole","SubnetId":"subnet-00bbdbcf537bf57b5","EmrManagedSlaveSecurityGroup":"sg-08a05f6eec6c610ac","EmrManagedMasterSecurityGroup":"sg-09d0c339df01417fd"}' --release-label emr-6.5.0 --log-uri 's3n://aws-logs-327946743066-eu-west-1/elasticmapreduce/' --instance-groups '[{"InstanceCount":6,"EbsConfiguration":{"EbsBlockDeviceConfigs":[{"VolumeSpecification":{"SizeInGB":32,"VolumeType":"gp2"},"VolumesPerInstance":2}]},"InstanceGroupType":"CORE","InstanceType":"m5.xlarge","Name":"Core - 2"},{"InstanceCount":1,"EbsConfiguration":{"EbsBlockDeviceConfigs":[{"VolumeSpecification":{"SizeInGB":32,"VolumeType":"gp2"},"VolumesPerInstance":2}]},"InstanceGroupType":"MASTER","InstanceType":"m5.xlarge","Name":"Master - 1"}]' --auto-scaling-role EMR_AutoScaling_DefaultRole --bootstrap-actions '[{"Path":"s3://projet8-oc/scripts-pyspark/emr_bootstrap_v5.sh","Name":"Custom action"}]' --ebs-root-volume-size 100 --service-role EMR_DefaultRole --enable-debugging --auto-termination-policy '{"IdleTimeout":14400}' --name 'Cluster pour Projet 8' --scale-down-behavior TERMINATE_AT_TASK_COMPLETION --region eu-west-1
```

## Script bash du bootstrap :

```
#!/bin/bash
sudo python3 -m pip install boto3 pandas==1.2.5 tensorflow==2.7.0 Pillow==8.4.0
```

## Script de lancement du job PySpark depuis la CLI du node master :

```
spark-submit --deploy-mode cluster --driver-memory 10g s3://projet8-oc/scripts-pyspark/extraction_features_v5.1.py
```

# CODE DU DRIVER

```
1 # choix de la zone selon disponibilité :
2 BUCKET_NAME = "projet8-oc" # eu-west-1
3
4 # imports modules
5 import os
6 from tensorflow.keras.preprocessing.image import load_img
7 from tensorflow.keras.preprocessing.image import img_to_array
8 from tensorflow.keras.applications.vgg16 import preprocess_input
9 from tensorflow.keras.applications.vgg16 import VGG16
10 from tensorflow.keras.models import Model
11 from pyspark import SparkContext
12 import boto3
13 from datetime import datetime
14 from time import time
15 import pandas as pd
16
17 def timestamp():
18     return datetime.now().strftime("%Y-%m-%d--%H-%M-%S")
19
20 def logger_s3(s3_bucket, logfile, log, log_to_shell=True):
21     """Appends log to logfile. Uploads each update to S3 bucket.
22     Args :
23         - s3_bucket (Bucket).
24         - logfile (string) .
25         - log (string).
26         - log_to_shell (bool) : si True, imprime également dans le std.out.
27     Returns : /
28     """
29     time_log = timestamp() + " : " + log + "\n"
30     if log_to_shell : print(""*100 + "\n", time_log)
31     with open("./logs/" + logfile, 'a') as f:
32         f.write(time_log)
33     s3_bucket.upload_file("./logs/" + logfile, 'logs/' + logfile)
34
35 # fonction d'extraction de feature
36 def VGG16_extracteur_spark(path, nom_image, model):
37     """Transforme un fichier image en un vecteur de dimension 4096.
38     Args :
39         - path : chemin vers les images (chemin local ou 'S3').
40         - nom_image : chemin d'accès à l'image (exemple : 'Apple Braeburn/r_173_100.jpg')
41         - model : model_vgg16 ou bc_model_vgg16.
42     Returns :
43         - liste de dimension 4097 (4096 dimensions de VGG16 + nom_image).
44     """
45     # create path to the image stored locally
46     if path == 'S3':
47         # il faut une instance Bucket différente sur le driver et les exécuteurs
48         s3_bucket_vgg16 = boto3.resource('s3').Bucket(BUCKET_NAME)
49         s3_bucket_vgg16.download_file('input/' + nom_image, '/tmp/img')
50         path_nom_image = '/tmp/img'
51     else:
52         path_nom_image = os.path.join(path, nom_image)
53     # Load the image for keras processing
54     image = load_img(path_nom_image, target_size=(224, 224))
55     # convert the image pixels to a numpy array
56     image = img_to_array(image)
57     # reshape data for the model
58     image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
59     # prepare the image for the VGG model
60     image = preprocess_input(image)
61     # get extracted features
62     im_features = model.predict(image)
63     # convert to list and add nom_image
64     im_features = im_features[0].tolist()
65     im_features.append(nom_image)
66     return im_features
```

```
71 # start Logging
72 os.makedirs(os.path.join(os.getcwd(), 'logs/'), exist_ok=True)
73 logfile = "Log for job started " + timestamp() + ".txt"
74 logger_s3(s3_bucket, logfile, "logging starts")
75
76 # Liste des fruits dans le répertoire "input/" : megabatch_img
77 all_fruits = s3_bucket.objects.filter(Prefix="input")
78 megabatch_img = []
79 for obj in all_fruits:
80     megabatch_img.append(obj.key.lstrip("input/"))
81 del megabatch_img[0] # suppression d'un objet non pertinent
82
83 # Load model VGG16
84 model_vgg16 = VGG16()
85 # remove the output layer
86 model_vgg16 = Model(inputs=model_vgg16.inputs, outputs=model_vgg16.layers[-2].output)
87
88 # Spark context : Le SparkConf est supprimé car géré via le spark-submit
89 sc = SparkContext()
90
91 # Logger des paramètres du spark context Lancé
92 logger_s3(s3_bucket, logfile, "Configuration de Spark : " + str(sc.getConf().getAll()))
93
94 # broadcasting
95 bc_model_vgg16 = sc.broadcast(model_vgg16)
96
97 # Logging
98 logger_s3(s3_bucket, logfile, "spark job starts")
99 t0 = time()
100
101 # job spark
102 resultat = sc.parallelize(megabatch_img) \
103     .map(lambda img: VGG16_extracteur_spark('S3', img, bc_model_vgg16.value)) \
104     .collect()
105
106 # Logging
107 logger_s3(s3_bucket, logfile, f"spark job has ended (duration of {round(time() - t0, 1)}s) - output to S3 starts")
108
109 # output
110 df_output = pd.DataFrame(resultat)
111 df_output.columns = [f'dim_{i}' for i in range(4096)] + ['path']
112 file_name = f'df_output_{len(df_output)}_fruits.csv'
113 local_path = os.path.join(os.getcwd(), file_name)
114 df_output.to_csv(path_or_buf=local_path)
115 s3_bucket.upload_file(local_path, 'output/' + file_name)
116
117 # end Logging
118 logger_s3(s3_bucket, logfile, "output to S3 has ended")
```