# Déployez un modèle dans le cloud

Data Science | Projet 8

Fırat Yasar
20/02/2022

# Sommaire

Présentation

# Présentation de la problématique

◉ **L'objectif :**

Notre start-up **Fruits** fait partie de l'**AgriTech de France** et cherche à proposer des solutions innovantes pour la récolte des fruits. Le but est de mettre 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.

◉ **La mission:**

Mettre en œuvre une première version du **moteur de classification** des images de fruits.

Construire dans un environnement **Big Data** une première chaîne de traitement des données qui comprendra le preprocessing et une étape de réduction de dimension.

◉ **La base de données :**

Pour cela, nous avons à notre disposition la base de données :

https://www.kaggle.com/moltean/fruits

# Découverte du jeu de données

```
1  display(data.withColumn('label', split(col('path'), '/').getItem(5)))
```

▸ (4) Spark Jobs

Table    Data Profile

| | path | modificationTime | length | content |
|---|---|---|---|---|
| 1 | dbfs:/mnt/mount_projet-8-bucket/data/Training/apple_hit_1/r0_116.jpg | 2022-02-13T10:42:01.000+0000 | 125373 | |
| 2 | dbfs:/mnt/mount_projet-8-bucket/data/Training/apple_hit_1/r0_114.jpg | 2022-02-13T10:42:01.000+0000 | 125088 | |
| 3 | dbfs:/mnt/mount_projet-8-bucket/data/Training/apple_hit_1/r0_108.jpg | 2022-02-13T10:42:01.000+0000 | 124905 | |
| | dbfs:/mnt/mount_projet-8-bucket/data/Training/apple_hit_1/r0_118.jpg | 2022-02-13T10:42:01.000+0000 | 124363 | |

Truncated results, showing first 36 rows.

☑ Show image preview ❔

**Data entraînement**
67 692 images

**Data test**
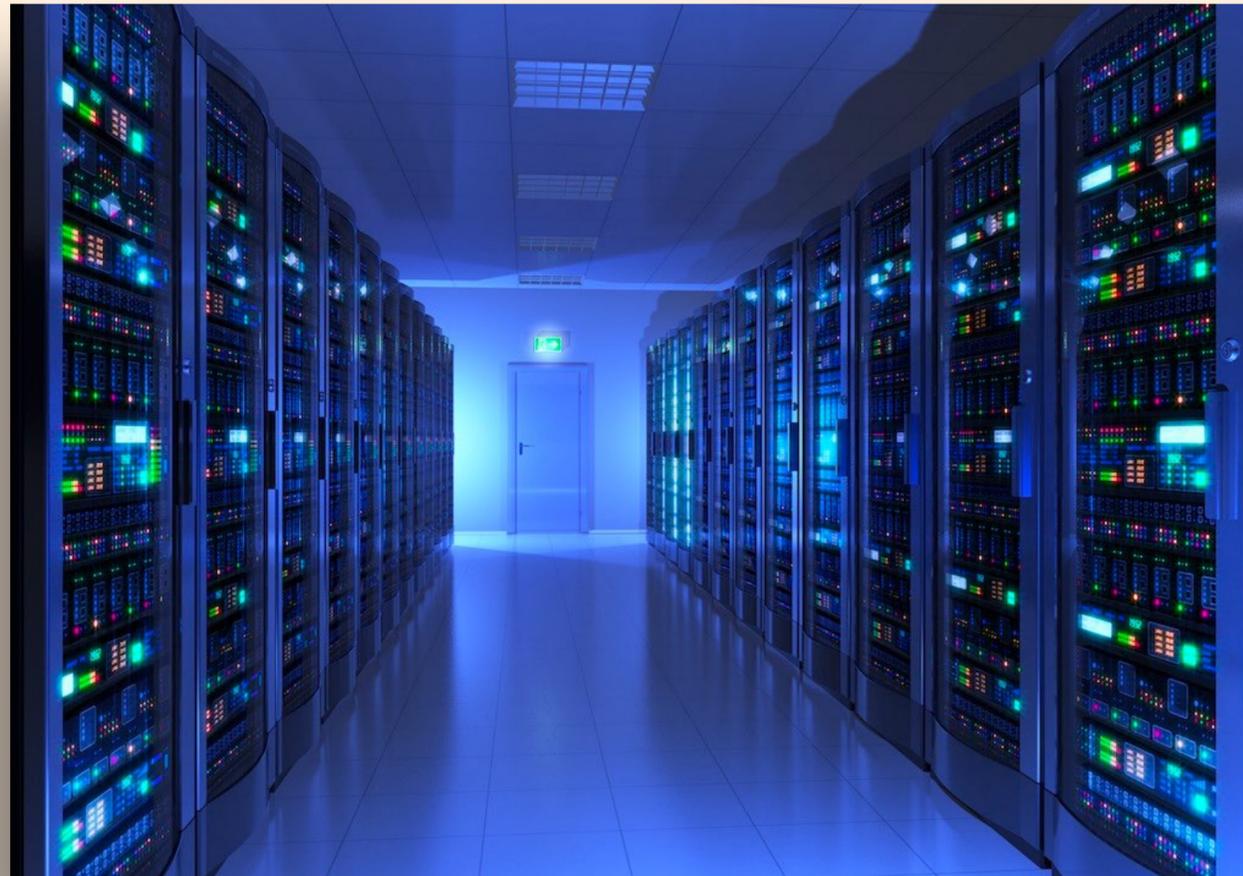22 688 images

une image = un fruit ou un légume

Taille d'une image : 100 x 100 pixels

Nombre de classes : 131
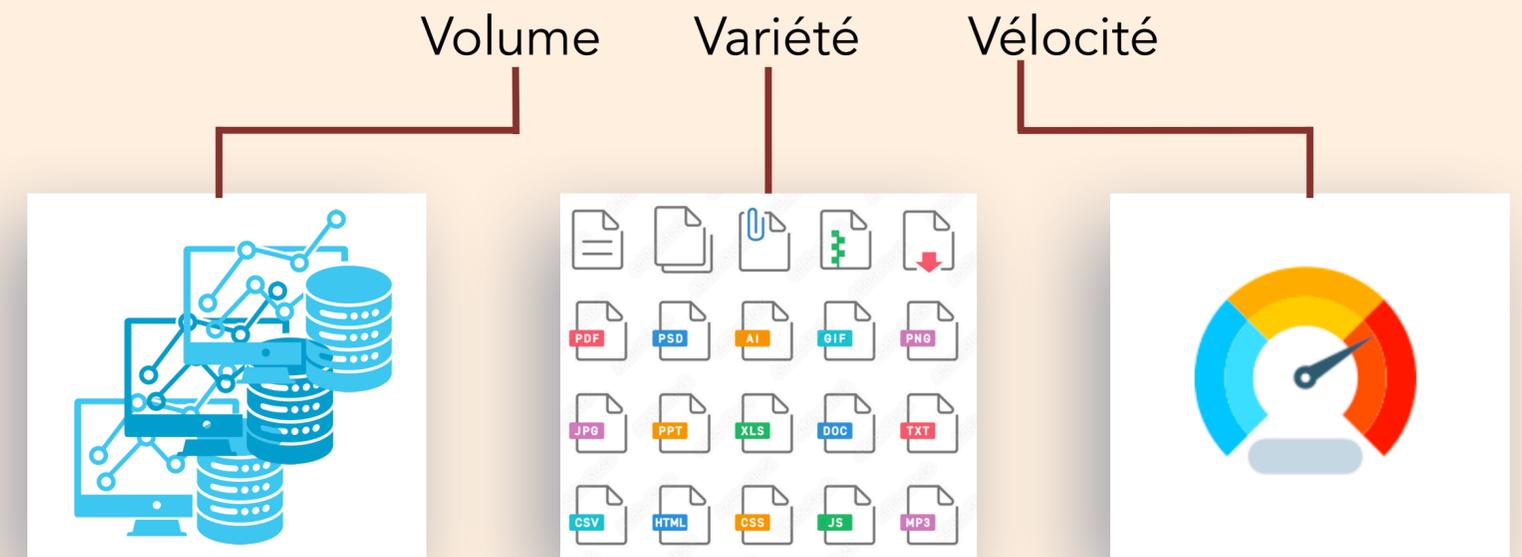
5

Architecture Big Data

# Qu'est-ce que le Big Data ?



◎ **Définition :**

Le **big data,** ou les données massives, désigne les ressources d'informations dont les caractéristiques en termes de **volume**, de **vélocité** et de **variété** imposent l'utilisation de technologies et de méthodes analytiques particulières pour générer de la valeur et qui dépassent en général les capacités d'une seule et unique machine et nécessitent des traitements parallélisés.
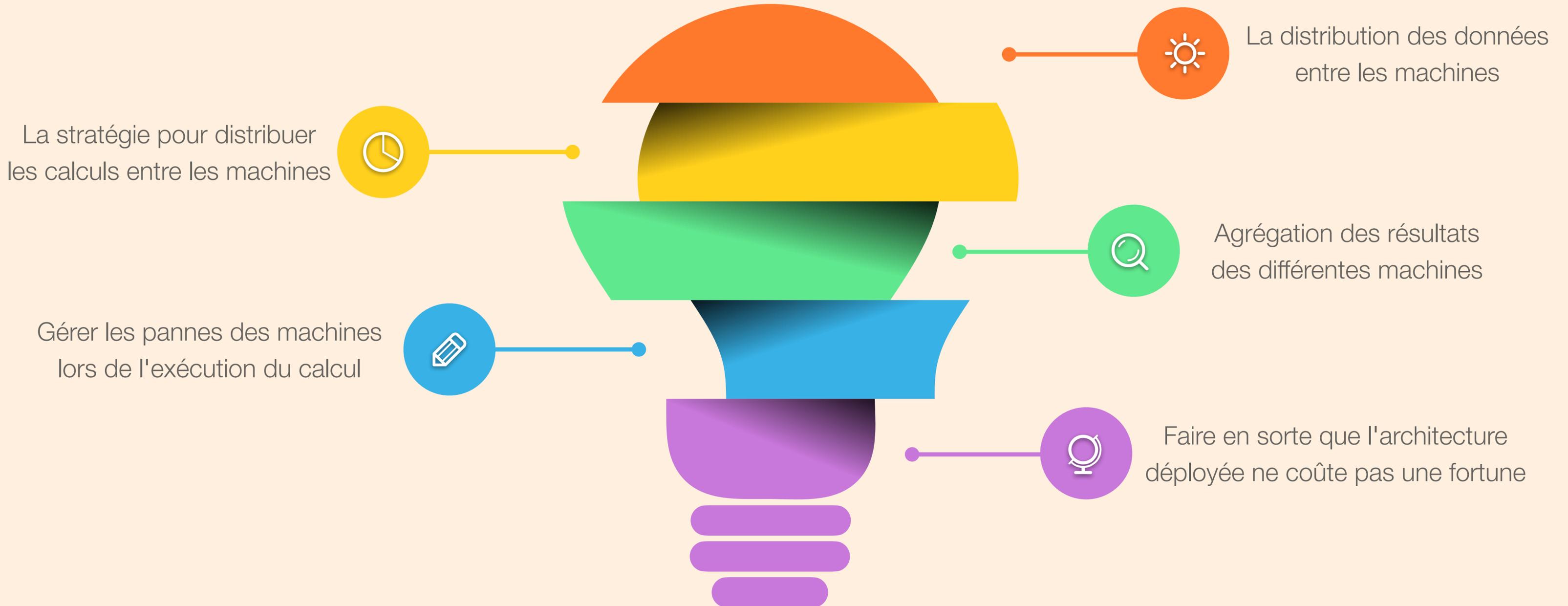
◎ **Les trois "V" de big data :**

Volume          Variété          Vélocité

# Traitement du Big Data

# Traitement du Big Data

La distribution des données entre les machines

La stratégie pour distribuer les calculs entre les machines

Agrégation des résultats des différentes machines

Gérer les pannes des machines lors de l'exécution du calcul

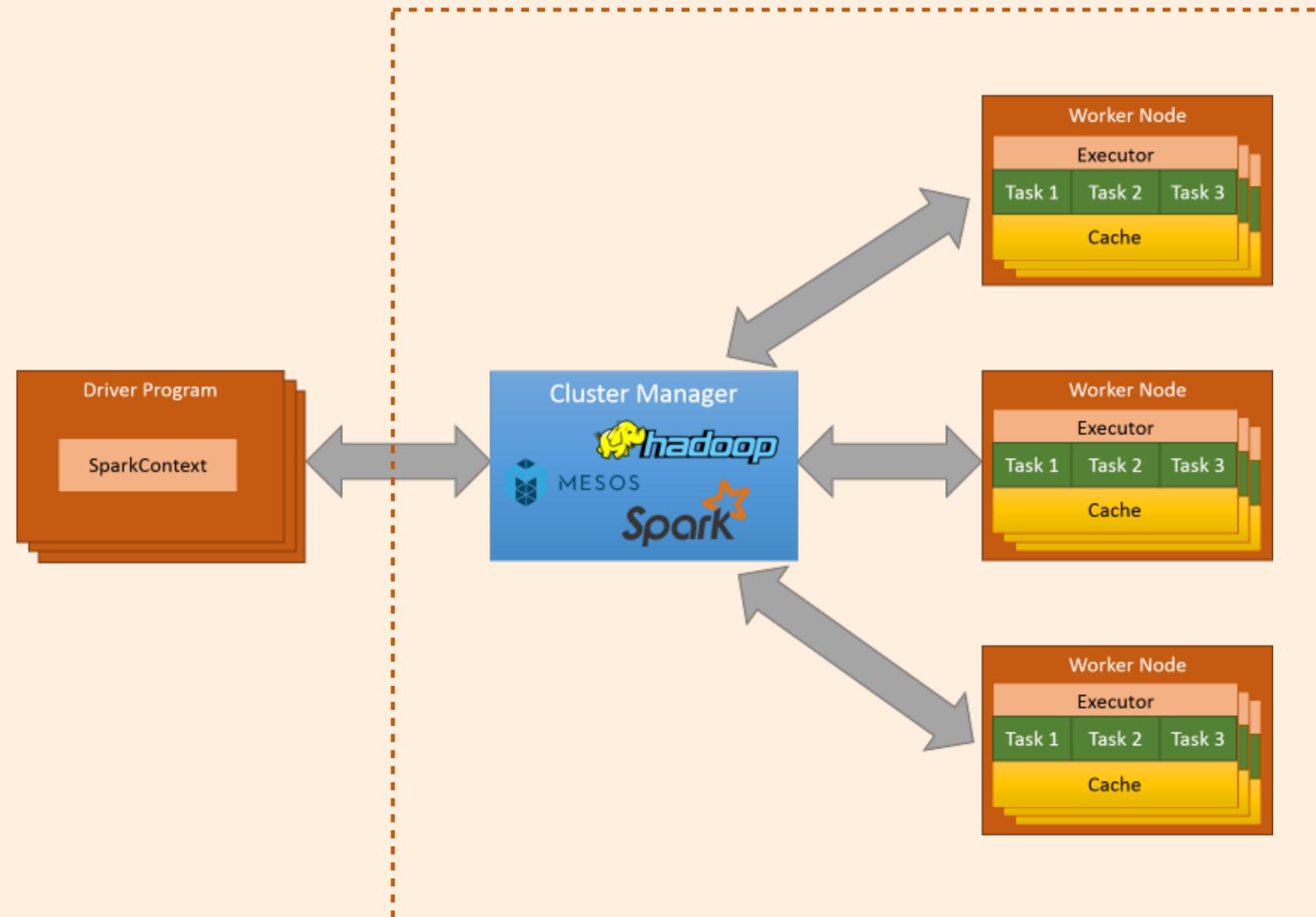Faire en sorte que l'architecture déployée ne coûte pas une fortune

# Automatisation du calcul distribué sur des données massives

◉ **Le calcul distribué :**

Les nœuds sur lesquels les calculs sont exécutés sont distants, autonomes et ne partagent pas de ressources ; la communication entre les nœuds s'effectue grâce à l'envoi de messages, au sein d'**un cluster.**

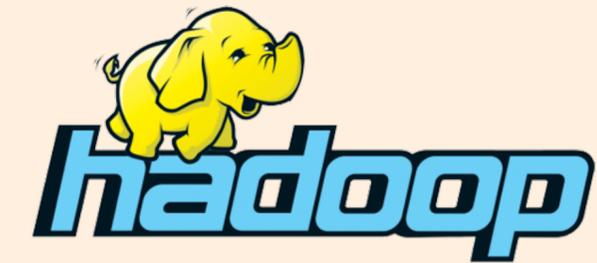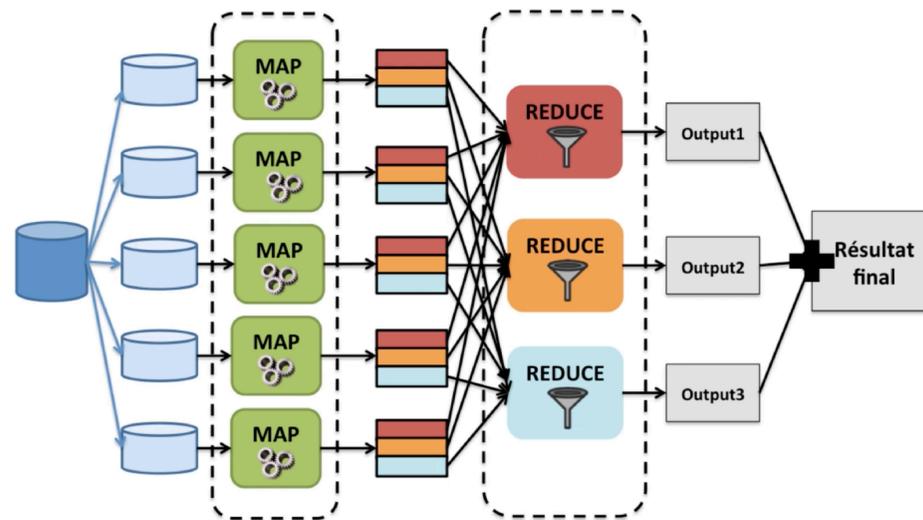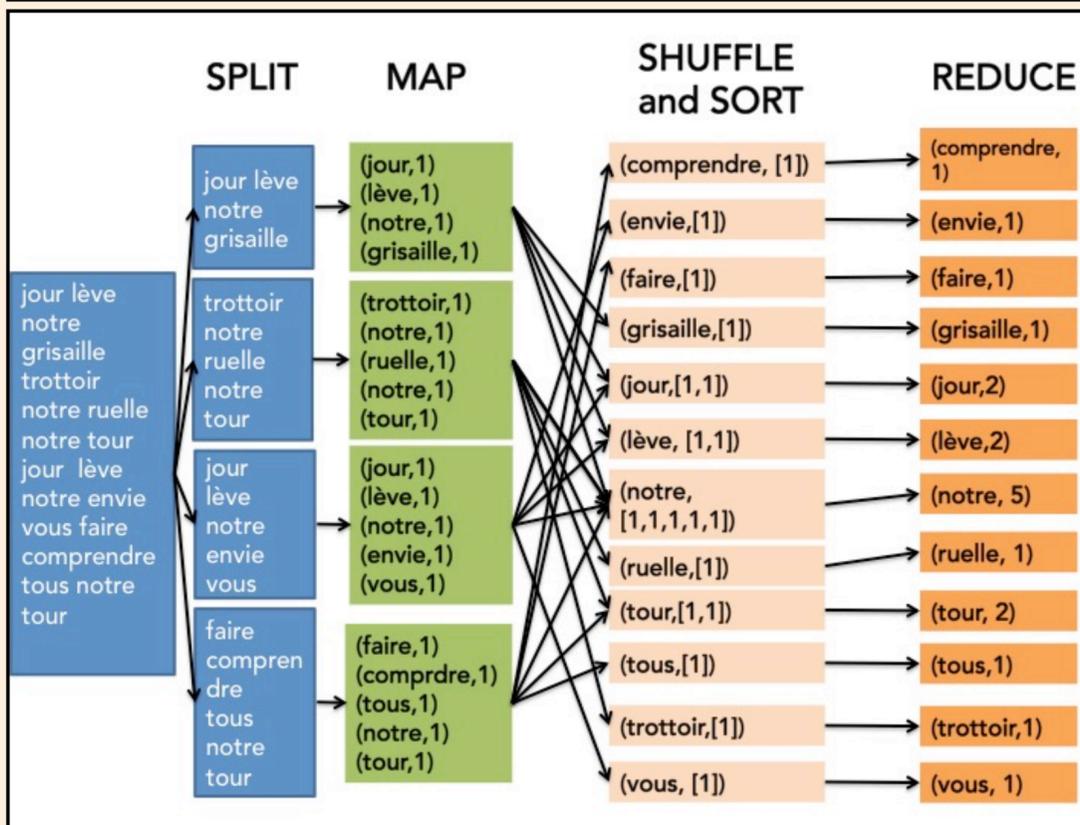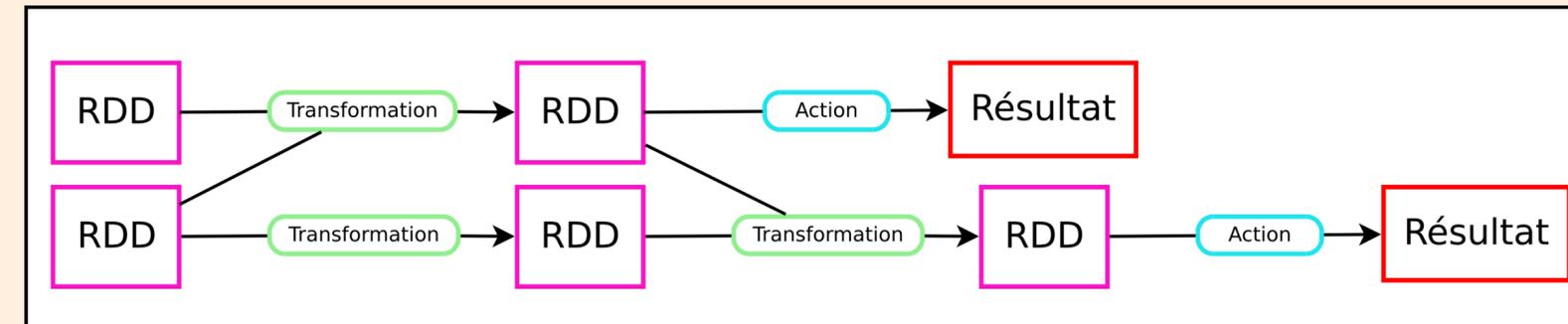**SparkContext** est le point d'accès à toutes les fonctionnalités de Spark.

Schéma d'exécution



◉ **Le fonctionnement général de MapReduce :**

1. L'ensemble des données à traiter est découpé en plusieurs lots ou sous-ensembles.

2. Dans une première étape, l'étape `MAP`, l'opération map, spécifiée pour notre problème, est appliquée à chaque lot. Cette opération transforme la paire `(clé, valeur)` représentant le lot en une liste de nouvelles paires (clé, valeur) constituant ainsi des résultats intermédiaires du traitement à effectuer sur les données complètes.

3. Avant d'être envoyés à l'étape REDUCE, les résultats intermédiaires sont regroupés et triés par clé. C'est l'étape de `SHUFFLE` and `SORT`.

4. Enfin, l'étape `REDUCE` consiste à appliquer l'opération reduce, spécifiée pour notre problème, à chaque clé. Elle agrège tous les résultats intermédiaires associés à une même clé et renvoie donc pour chaque clé une valeur unique.
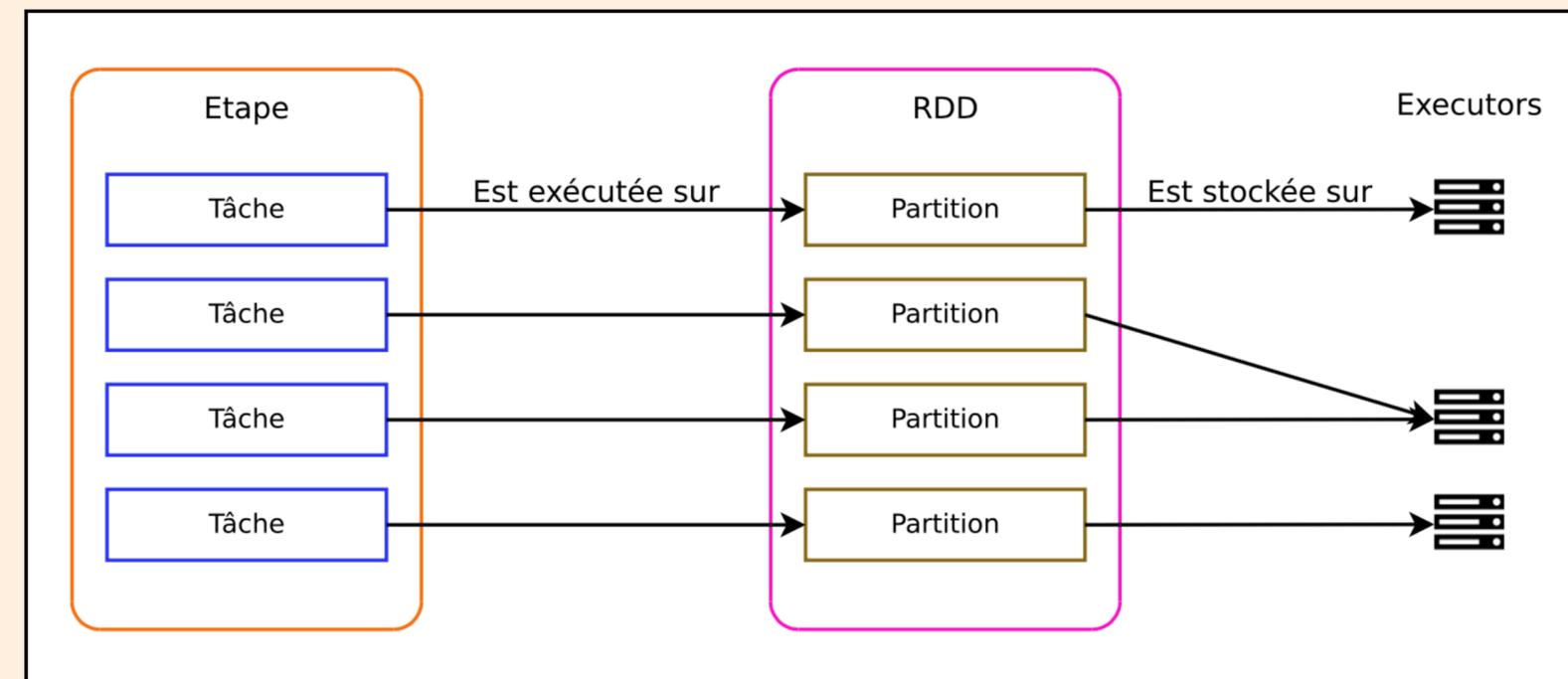
# Spark

## Le fonctionnement général de Spark :

- Utilisation des RDD (**R**esilient **D**istributed **D**ataset). RDD supporte deux type d'opérations : transformation et action.

  - Transformation : consiste à appliquer une fonction sur $n$ RDD et à retourner un nouveau RDD (map, filter, join etc.).
  - Action : consiste à appliquer une fonction sur $n$ RDD et à retourner une valeur (collect, show, count etc.).

- Les transformations évitent les calculs inutiles. Spark exécute les expressions uniquement lorsqu'elles sont nécessaires.

- Les RDD sont exécutés en mémoire de façon complètement tolérante aux pannes.

***

- Spark est basé sur la programmation fonctionnelle, **SCALA**
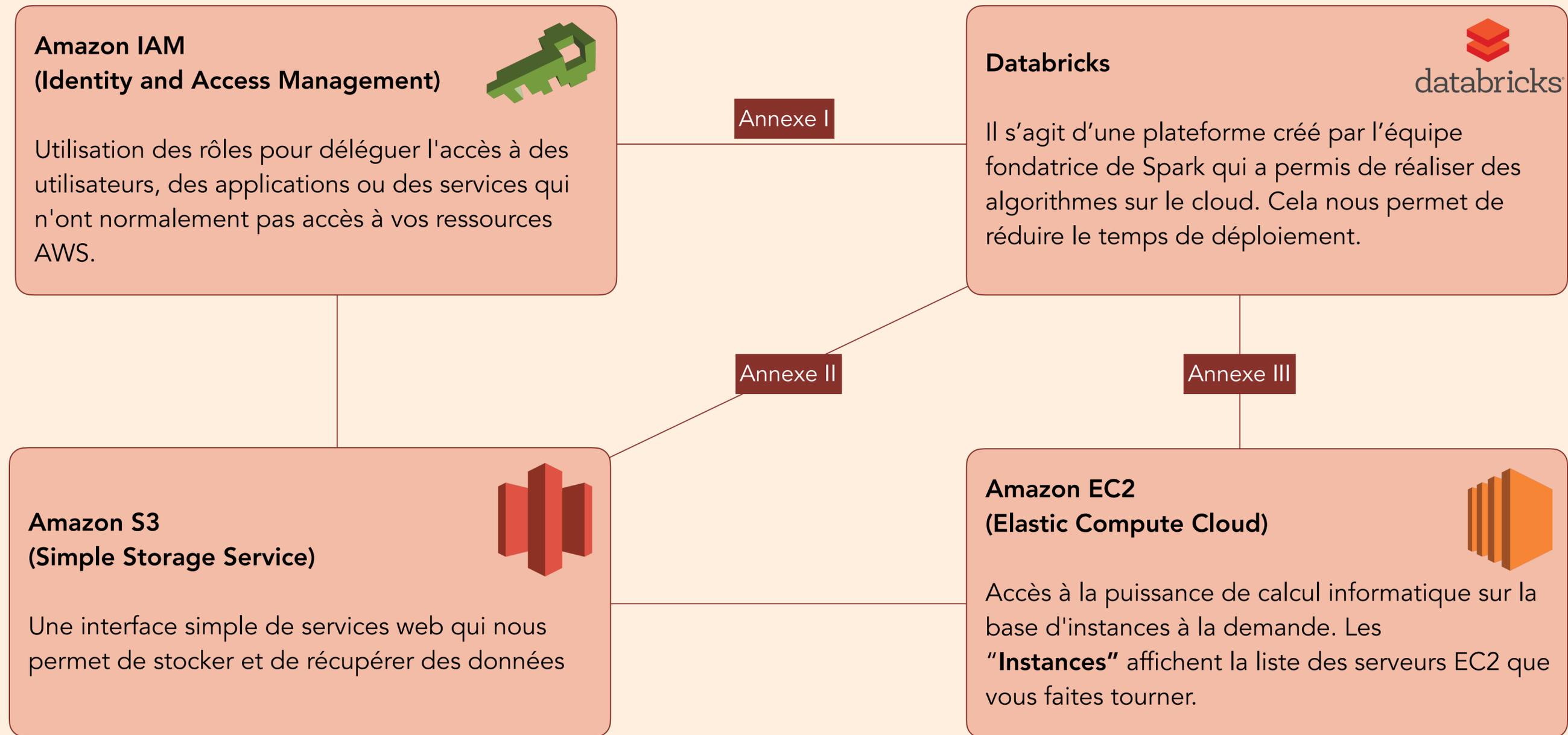
Amazon Web Services

13

# La console d'AWS

# Les services utilisés

## Amazon IAM
### (Identity and Access Management)

Utilisation des rôles pour déléguer l'accès à des utilisateurs, des applications ou des services qui n'ont normalement pas accès à vos ressources AWS.

Annexe I

## Databricks

Il s'agit d'une plateforme créé par l'équipe fondatrice de Spark qui a permis de réaliser des algorithmes sur le cloud. Cela nous permet de réduire le temps de déploiement.

Annexe II

Annexe III

## Amazon S3
### (Simple Storage Service)

Une interface simple de services web qui nous permet de stocker et de récupérer des données

## Amazon EC2
### (Elastic Compute Cloud)

Accès à la puissance de calcul informatique sur la base d'instances à la demande. Les **"Instances"** affichent la liste des serveurs EC2 que vous faites tourner.

# Création d'un cluster via Databricks



Create Cluster

New Cluster | Cancel | Create Cluster | DBU / hour: 4.11 - 12.33 | 2-8 Workers:64-256 GB Memory, 16-64 Cores / 1 Driver:32 GB Memory, 8 Cores

Free trial ends in **11** days. Continue with a pay-as-you-go subscription by providing your billing information.

UI | JSON

**Cluster name**
projet-8-cluster

**Cluster mode** ⍰
Standard

**Databricks runtime version** ⍰ — Learn more
Runtime: 9.1 LTS (Scala 2.12, Spark 3.1.2)

**Autopilot options**
☑ Enable autoscaling ⍰
☐ Enable autoscaling local storage ⍰
☑ Terminate after 120 minutes of inactivity ⍰

**Worker type** ⍰ — Min workers 2 — Max workers 8
m5d.2xlarge — 32 GB Memory, 8 Cores

**New** Configure separate pools for workers and drivers for flexibility. Learn more

**Driver type**
Same as worker — 32 GB Memory, 8 Cores

DBU / hour: 4.11 - 12.33 ⍰ — m5d.2xlarge

▸ Advanced options

Nom de cluster

Runtime version : 10.3 ML (includes Apache Spark 3.2.1, Scala 2.12)

Type d'instance EC2 : m5d.2xlarge
32 GB RAM
4 coeurs
2-8 noeuds

Type d'instance pour le driver

Tarification pour la configuration

# Le notebook

P8_Notebook  (Python)

Free trial ends in 5 days. Continue with a pay-as-you-go subscription by providing your billing information.

Schedule ∨    Share

projet-8-cluster | ∨    📄 File ▾    ✏ Edit ▾    🖼 View: Standard ▾    ▶ Run All    🧽 Clear ▾    ❓ Help

💬 Comments    🧪 Experiment    🕘 Revision history

## 1.1. Les bibliothèques nécessaires et initialisation Spark

Cmd 4

Python ▶▾ ∨ ─ ✕

```python
1   # Standard libraries
2   import numpy as np
3   import pandas as pd
4   import matplotlib.pyplot as plt
5   import seaborn as sns
6   from PIL import Image
7
8   # PySpark libraries
9   import pyspark
10  from pyspark.sql import SparkSession
11  from pyspark import SparkContext, SparkConf
12  from pyspark.sql.functions import split, col
13  import boto3
14
15  # MLlib, Spark's Machine Learning (ML) library
```

Command took 0.02 seconds -- by fyasar.fr@gmail.com at 18/02/2022, 02:23:28 on projet-8-cluster

**Associer le notebook avec le cluster configuré**

Cmd 5

```python
1   # Set up Hadoop Configurations for AWS_ACCESS_KEY & AWS_SECRET_KEY
2
3   ACCESS_KEY = "▮▮▮▮▮▮▮▮▮▮▮▮"
4
5   SECRET_KEY = "▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮"
6
7   ENCODED_SECRET_KEY = SECRET_KEY.replace("/", "%2F")
8
9   AWS_BUCKET_NAME = "projet-8-bucket"
10
11  MOUNT_NAME = "mount_projet-8-bucket"
12
13  dbutils.fs.mount(f"s3a://{ACCESS_KEY}:{ENCODED_SECRET_KEY}@{AWS_BUCKET_NAME}", f"/mnt/{MOUNT_NAME}")
14
15  print('AWS S3 bucket {} is successfully mounted.'.format(AWS_BUCKET_NAME))
```

⊞ java.rmi.RemoteException: java.lang.IllegalArgumentException: requirement failed: Directory already mounted: /mnt/mount_projet-8-bucket; nested exception is:

Command took 0.49 seconds -- by fyasar.fr@gmail.com at 18/02/2022, 02:23:28 on projet-8-cluster
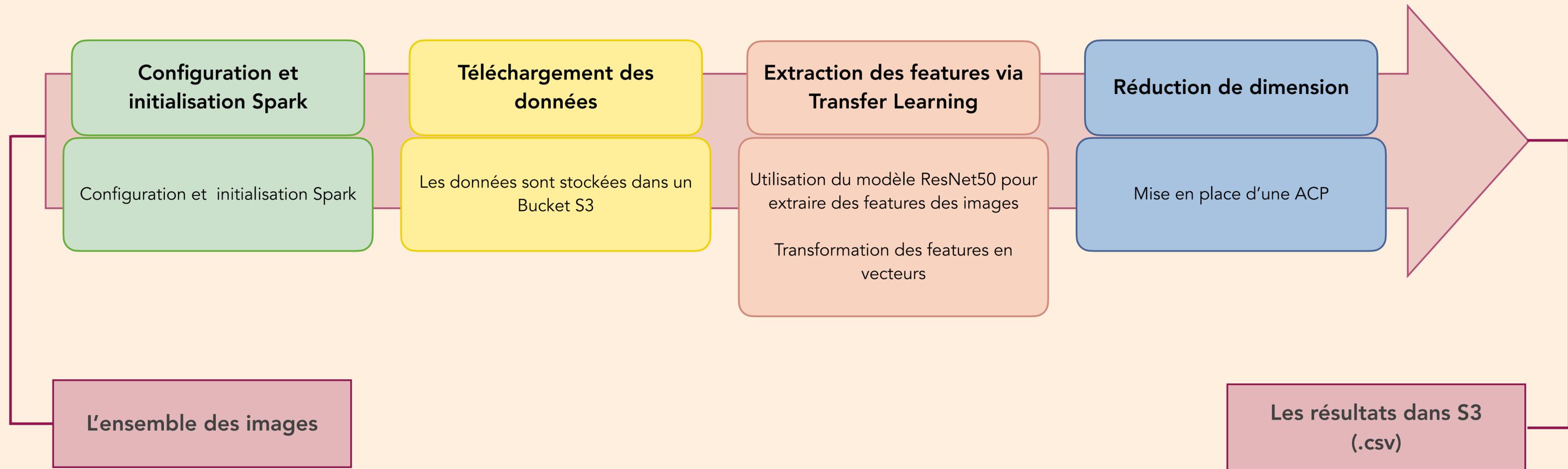
Cmd 6

## 1.2. Téléchargement du jeux de données

Cmd 7

La chaîne de traitement

# Le pipeline

| Configuration et initialisation Spark | Téléchargement des données | Extraction des features via Transfer Learning | Réduction de dimension |
|---|---|---|---|
| Configuration et initialisation Spark | Les données sont stockées dans un Bucket S3 | Utilisation du modèle ResNet50 pour extraire des features des images<br><br>Transformation des features en vecteurs | Mise en place d'une ACP |

**L'ensemble des images**

**Les résultats dans S3 (.csv)**

# Les résultats

Merci de votre attention

1. IAM Dashboard → Access management → Roles

2. Create role

23

**Add permissions**

Permissions policies (726)
Choose one or more policies to attach to your new role.

**Next**

Filter policies by property or policy name and press enter

| Policy name | Type | Description |
|---|---|---|
| AWSDirectConnectReadOnlyAccess | AWS m... | Provides read only access to AWS Direct Connect via the AWS Manage... |
| AmazonGlacierReadOnlyAccess | AWS m... | Provides read only access to Amazon Glacier via the AWS Managemen... |
| AWSMarketplaceFullAccess | AWS m... | Provides the ability to subscribe and unsubscribe to AWS Marketplace ... |
| AWSSSODirectoryAdministrator | AWS m... | Administrator access for SSO Directory |
| AWSIoT1ClickReadOnlyAccess | AWS m... | Provides read only access to AWS IoT 1-Click. |
| AutoScalingConsoleReadOnlyAccess | AWS m... | Provides read-only access to Auto Scaling via the AWS Management C... |
| AmazonDMSRedshiftS3Role | AWS m... | Provides access to manage S3 settings for Redshift endpoints for DMS. |
| AWSQuickSightListIAM | AWS m... | Allow QuickSight to list IAM entities |
| AWSHealthFullAccess | AWS m... | Allows full access to the AWS Health Apis and Notifications and the Per... |
| AlexaForBusinessGatewayExecution | AWS m... | Provide gateway execution access to AlexaForBusiness services |
| AmazonElasticTranscoder_ReadOnl... | AWS m... | Grants users read-only access to Elastic Transcoder and list access to r... |
| AmazonRDSFullAccess | AWS m... | Provides full access to Amazon RDS via the AWS Management Console. |
| SupportUser | AWS m... | This policy grants permissions to troubleshoot and resolve issues in an ... |
| AmazonEC2FullAccess | AWS m... | Provides full access to Amazon EC2 via the AWS Management Console. |
| SecretsManagerReadWrite | AWS m... | Provides read/write access to AWS Secrets Manager via the AWS Man... |
| AWSIoTThingsRegistration | AWS m... | This policy allows users to register things at bulk using AWS IoT StartTh... |
| AmazonDocDBReadOnlyAccess | AWS m... | Provides read-only access to Amazon DocumentDB with MongoDB co... |

---

**Name, review, and create**

**6**

**Role details**

Role name
Enter a meaningful name to identify this role.

Role-Projet-8

Maximum 128 characters. Use alphanumeric and '+=,.@-_' characters.

Description
Add a short explanation for this policy.

Maximum 1000 characters. Use alphanumeric and '+=,.@-_' characters.

Step 1: Select trusted entities     Edit

```
1  {
2    "Version": "2012-10-17",
3    "Statement": [
4      {
5        "Effect": "Allow",
6        "Action": "sts:AssumeRole",
7        "Principal": {
8          "AWS": "414351767826"
9        },
10       "Condition": {
```

---

Role Role-Projet-8 created     View role

**Roles (10)** Info

An IAM role is an identity you can create that has specific permissions with credentials that are valid for short durations. Roles can be assumed by entities that you trust.

Delete     Create role

**7**

| Role name | Trusted entities | Last activity |
|---|---|---|
| aws-elasticbeanstalk-ec2-role | AWS Service: ec2 | 2 days ago |
| aws-elasticbeanstalk-service-role | AWS Service: elasticbeanstalk | 2 days ago |
| AWSServiceRoleForAutoScaling | AWS Service: autoscaling (Service-Linked Rol | 2 days ago |
| AWSServiceRoleForEC2Spot | AWS Service: spot (Service-Linked Role) | - |
| AWSServiceRoleForElasticLoadBalancin | AWS Service: elasticloadbalancing (Service-Li | 2 days ago |
| AWSServiceRoleForSupport | AWS Service: support (Service-Linked Role) | - |
| AWSServiceRoleForTrustedAdvisor | AWS Service: trustedadvisor (Service-Linked | - |
| databricks-workspace-stack-CopyZipsF | AWS Service: lambda | 3 days ago |
| databricks-workspace-stack-functionRo | AWS Service: lambda | 3 days ago |
| Role-Projet-8 | Account: 414351767826 | - |

Databricks VPC

# Create policy

① ②

## Review policy

Before you create this policy, provide the required information and review this policy.

**10**

Name* `Policy-Projet-8`

Maximum 128 characters. Use alphanumeric and '+=,.@-_' characters.

Summary

This policy defines some actions, resources, or conditions that do not provide permissions. To grant access, policies must have an action that has an applicable resource or condition. For details, choose **Show remaining**. Learn more

🔍 Filter

| Service ▾ | Access level | Resource | Request condition |
|---|---|---|---|
| Allow (2 of 315 services) Show remaining 313 | | | |
| EC2 | **Full**: Tagging **Limited**: List, Write | All resources | None |
| IAM | **Limited**: Write | Path \| string like \| aws-service-role/spot.amazonaws.com, RoleName \| string like \| AWSServiceRoleForEC2Spot | iam:AWSServiceName \| string lil spot.amazonaws.com |

\* Required

Cancel   Previous   **Create policy**

---

**Identity and Access Management (IAM)** ✕

🔍 Search IAM

Dashboard

▾ **Access management**
User groups
Users
**Roles**
Policies
Identity providers
Account settings

▾ **Access reports**
Access analyzer
Archive rules
Analyzers
Settings
Credential report
Organization activity
Service control policies (SCPs)

IAM > Roles > Role-Projet-8

# Role-Projet-8

**11**

Delete

## Summary

Edit

| Creation date | ARN | Link to switch roles in console |
|---|---|---|
| February 11, 2022, 19:52 (UTC+01:00) | arn:aws:iam::744140778731:role/Role-Projet-8 | https://signin.aws.amazon.com/switchrole?roleName=Role-Projet-8&account=744140778731 |
| Last activity | Maximum session duration | |
| None | 1 hour | |

Permissions | Trust relationships | Tags | Access Advisor | Revoke sessions

### Permissions policies (1)

You can attach up to 10 managed policies.

🔄   Simulate   Remove

Add permissions ▾

🔍 Filter policies by property or policy name and press enter       ‹ 1 › ⚙

| ☐ | Policy name ⧉ | Type | Description |
|---|---|---|---|
| ☐ | ⊞ Policy-Projet-8 | Customer inline | |

### Permissions boundary - (not set)

Set a permissions boundary to control the maximum permissions this role can

## Cloud resources

Credential configuration    Storage configuration    Network configuration

### Credential configuration

For Databricks to launch clusters in your AWS account, you must create a cross-account IAM role that gives access to Databricks. Learn more

| Name | Role ARN | Created | |
|---|---|---|---|
| yasarigno-credentials | arn:aws:iam::744140778731:role/db-iam-role | last Tuesday at 6:27 PM | |

**12**

---

## Cloud resources

Credential configuration    Storage configuration    Network configuration

### Credential configuration

For Databricks to launch clusters in y...                                    ...atabricks. Learn more

### Add credential configuration                    ✕

For Databricks to launch clusters in your AWS account, you must create a cross-account IAM role that gives access to Databricks. To learn how, see Create a cross-account IAM role.

Once you have created the role, enter the role name and role ARN here. You are responsible for the AWS cost of clusters you create.

External ID ⑦

* Credential configuration name

CC-Projet-8

Human readable name to label your configuration

* Role ARN ⑦

arn:aws:iam::744140778731:role/Role-Projet-8

Cancel    Add

**13**

Annexe II
Databricks - AWS : Storage configuration

# Cloud resources

Credential configuration   Storage configuration   Network configuration

## Storage configuration

Databricks stores your account-wide assets, such as libraries, in an AWS S3 bucket that you must configure in your AWS account using a policy supplied by Databricks. Learn more

[Search]   Search

Add storage configuration

| Name | Bucket name | Created |
|------|-------------|---------|
| yasarigno-storage | db- | last Tuesday at 6:27 PM |

**1**

---

# Cloud resources

Credential configuration   Storage configuration   Network configuration

## Storage configuration

Databricks stores your account-wide... account using a policy supplied by Databricks. Learn more

### Add storage configuration

Databricks stores your account-wide assets, such as libraries, in an AWS S3 bucket that you must configure in your AWS account using a policy supplied by Databricks.

Enter the bucket name and click **Generate policy** to generate the necessary bucket policy to copy. Learn more.

\* Storage configuration name

cc-storage-projet-8

Human readable name to label your configuration

\* Bucket name

projet-8-bucket

Generate policy

Cancel   Add

**2**

---

## Amazon S3

Buckets
Access Points
Object Lambda Access Points
Multi-Region Access Points
Batch Operations
Access analyzer for S3

Block Public Access settings for this account

**Storage Lens**
Dashboards
AWS Organizations settings

Feature spotlight

AWS Marketplace for S3

Amazon S3

▶ **Account snapshot**   [View Storage Lens dashboard]
Storage lens provides visibility into storage usage and activity trends. Learn more

### Buckets (3) Info
Buckets are containers for data stored in S3. Learn more

[Copy ARN]  [Empty]  [Delete]  [Create bucket]

[Find buckets by name]

| Name | AWS Region | Access | Creation date |
|------|-----------|--------|---------------|
| databricks-workspace-stack-lambdazipsbucket-o62q1by5cxhe | US West (Oregon) us-west-2 | Objects can be public | February 8, 2022, 18:26:28 (UTC+01:00) |
| db-632143419cb5c21e3188dd00dac9277d-s3-root-bucket | US West (Oregon) us-west-2 | Bucket and objects not public | February 8, 2022, 18:26:28 (UTC+01:00) |
| elasticbeanstalk-eu-west-3-744140778731 | EU (Paris) eu-west-3 | Objects can be public | February 9, 2022, 00:29:10 (UTC+01:00) |

---

## Amazon S3

Buckets
Access Points
Object Lambda Access Points
Multi-Region Access Points
Batch Operations
Access analyzer for S3

Block Public Access settings for this account

**Storage Lens**
Dashboards
AWS Organizations settings

Feature spotlight

AWS Marketplace for S3

Amazon S3 > Create bucket

### Create bucket  Info
Buckets are containers for data stored in S3. Learn more

#### General configuration

Bucket name
projet-8-bucket

Bucket name must be unique and must not contain spaces or uppercase letters. See rules for bucket naming

AWS Region
EU (Frankfurt) eu-central-1

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

[Choose bucket]

#### Object Ownership  Info
Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

○ **ACLs disabled (recommended)**
All objects in this bucket are owned by this account. Access to this bucket and its objects is specified using only policies.

○ **ACLs enabled**
Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership

**3**

**Screenshot 1 (top-left): Amazon S3 Buckets**

Amazon S3

- Buckets
- Access Points
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- Access analyzer for S3

Block Public Access settings for this account

Storage Lens
- Dashboards
- AWS Organizations settings

Feature spotlight

AWS Marketplace for S3

✓ Successfully created bucket "projet-8-bucket"
To upload files and folders, or to configure additional bucket settings choose View details.

View details

Amazon S3

▶ Account snapshot                    View Storage Lens dashboard
Storage lens provides visibility into storage usage and activity trends. Learn more

Buckets (4) Info        Copy ARN   Empty   Delete   Create bucket
Buckets are containers for data stored in S3. Learn more

Find buckets by name                                    < 1 >

| Name | AWS Region | Access | Creation date |
|------|-----------|--------|---------------|
| databricks-workspace-stack-lambdazipsbucket-o62q1by5cxhe | US West (Oregon) us-west-2 | Objects can be public | February 8, 2022, 18:26:28 (UTC+01:00) |
| db-632143419cb5c21e3188dd00dac9277d-s3-root-bucket | US West (Oregon) us-west-2 | Bucket and objects not public | February 8, 2022, 18:26:28 (UTC+01:00) |
| elasticbeanstalk-eu-west-3-744140778731 | EU (Paris) eu-west-3 | Objects can be public | February 9, 2022, 00:29:10 (UTC+01:00) |
| projet-8-bucket | EU (Frankfurt) eu-central-1 | Bucket and objects not public | February 11, 2022, 20:22:56 (UTC+01:00) |

**Screenshot 2 (top-right): Edit bucket policy**

Amazon S3 > projet-8-bucket > Edit bucket policy

# Edit bucket policy Info

### Bucket policy
The bucket policy, written in JSON, provides access to the objects stored in the bucket. Bucket policies don't apply to objects owned by other accounts. Learn more

Policy examples     Policy generator

Bucket ARN
arn:aws:s3:::projet-8-bucket

**4**

Policy
```
1  {
2    "Version": "2012-10-17",
3    "Statement": [
4      {
5        "Sid": "Statement1",
6        "Principal": {},
7        "Effect": "Allow",
8        "Action": [],
9        "Resource": []
10     }
11   ]
12 }
```

Edit statement

Select a statement

Select an existing statement in the policy or add a new statement.

+ Add new statement

**Screenshot 3 (bottom): projet-8-bucket Objects**

Amazon S3 > projet-8-bucket

# projet-8-bucket Info

Objects | Properties | Permissions | Metrics | Management | Access Points

### Objects (0)
Objects are the fundamental entities stored in Amazon S3. You can use Amazon S3 inventory to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. Learn more

Copy S3 URI   Copy URL   Download   Open   Delete   Actions ▼

Create folder   Upload

Find objects by prefix                                  < 1 >

| Name | Type | Last modified | Size | Storage class |
|------|------|---------------|------|---------------|

No objects
You don't have any objects in this bucket.

Upload

**Add storage configuration** ×

Databricks stores your account-wide assets, such as libraries, in an AWS S3 bucket that you must configure in your AWS account using a policy supplied by Databricks. Learn more

Enter the bucket name and click **Generate policy** to generate the necessary bucket policy to copy. Learn more.

\* Storage configuration name

cc-storage-projet-8

Human readable name to label your configuration

\* Bucket name

projet-8-bucket

Generate policy

Cancel    Add

---

**Cloud resources**

Credential configuration    Storage configuration    Network configuration

**Storage configuration**

Databricks stores your account-wide assets, such as libraries, in an AWS S3 bucket that you must configure in your AWS account using a policy supplied by Databricks. Learn more

Add storage configuration

| Name | Bucket name | Created | |
|------|-------------|---------|--|
| cc-storage-projet-8 | projet-8-bucket | today at 8:27 PM | |
| yasarigno-storage | | last Tuesday at 6:27 PM | |

---

**Add storage configuration** ×

Human readable name to label your configuration

\* Bucket name

projet-8-bucket

Hide policy

Copy this policy and follow these instructions to complete the storage configuration process.

```
{
    "Version": "2012-10-17",
    "Statement": [
      {
        "Sid": "Grant Databricks Access",
        "Effect": "Allow",
        "Principal": {
          "AWS": "arn:aws:iam::414351767826:root"
        },
        "Action": [
          "s3:GetObject",
          "s3:GetObjectVersion",
          "s3:PutObject",
          "s3:DeleteObject",
          "s3:ListBucket",
          "s3:GetBucketLocation"
        ],
        "Resource": [
          "arn:aws:s3:::projet-8-bucket/*",
          "arn:aws:s3:::projet-8-bucket"
        ]
      }
    ]
}
```

Cancel    Add

**3**

## Workspaces

Search | Search

Create workspace ⌄

| Name | Status | Pricing tier | Region | Bucket name | Credential name | Created | |
|------|--------|-------------|--------|-------------|-----------------|---------|---|
| yasarigno | Running | Premium | us-west-2 | db-███ | yasarigno-credentials | last Tuesday at 6:27 PM | Open ⧉ ⋮ |

---

## Create workspace

Workspaces / Create workspace

**Creating workspace**

### Configurations

\* Workspace name

projet-8-workspace

Human readable name for your workspace

\* Subscription plan

Premium ⌄

\* Region

Frankfurt (eu-central-1) ⌄

\* Credential configuration

CC-Projet-8 ⌄

\* Storage configuration

cc-storage-projet-8 ⌄

Role ARN: arn:aws:iam::744140778731:role/Role-Projet-8

Bucket Name: projet-8-bucket

⌄ Advanced configurations

Save | Cancel

33

Annexe III
Databricks - AWS : Création d'un cluster et d'un notebook