

M2 SIAD Data Science  
2021 - 2022

---

# **La Data Science Orientée Production**

## **1. Cloud Computing**

---

Corentin Vasseur

## SOMMAIRE

1

### Introduction

Qu'est ce que le Cloud Computing ?  
3 grands acteurs  
Les fondamentaux  
Comment l'utiliser ?

2

### Le Stockage

Différents types de stockage  
Accès CLI  
Accès SDK

3

### Le Calcul

Machines Virtuelles  
Les coûts



### Machine Learning

Machine learning system  
ML Pipeline  
Tools chain  
API ML

4

### Platform AI

SageMaker (AWS)  
AI Platform (GCP)  
Azure ML (Microsoft Azure)

5

### Workflow ML

Full Data Engineering Pipeline  
Ex. GCP  
Ex. AWS

6

# 1

---

## Introduction : les fondamentaux

---

[RETOUR SOMMAIRE](#)

## 1. Introduction : les fondamentaux

### Qu'est-ce que le cloud computing ?



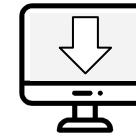
On-demand  
self-service



Broad network  
access



Ressource  
pooling



Rapid elasticity



Pay for what you  
use

# 1. Introduction : les fondamentaux

## 3 grands acteurs



**AWS**

+ 165 services

2006



**GCP**

+ 90 services

2008



**Microsoft Azure**

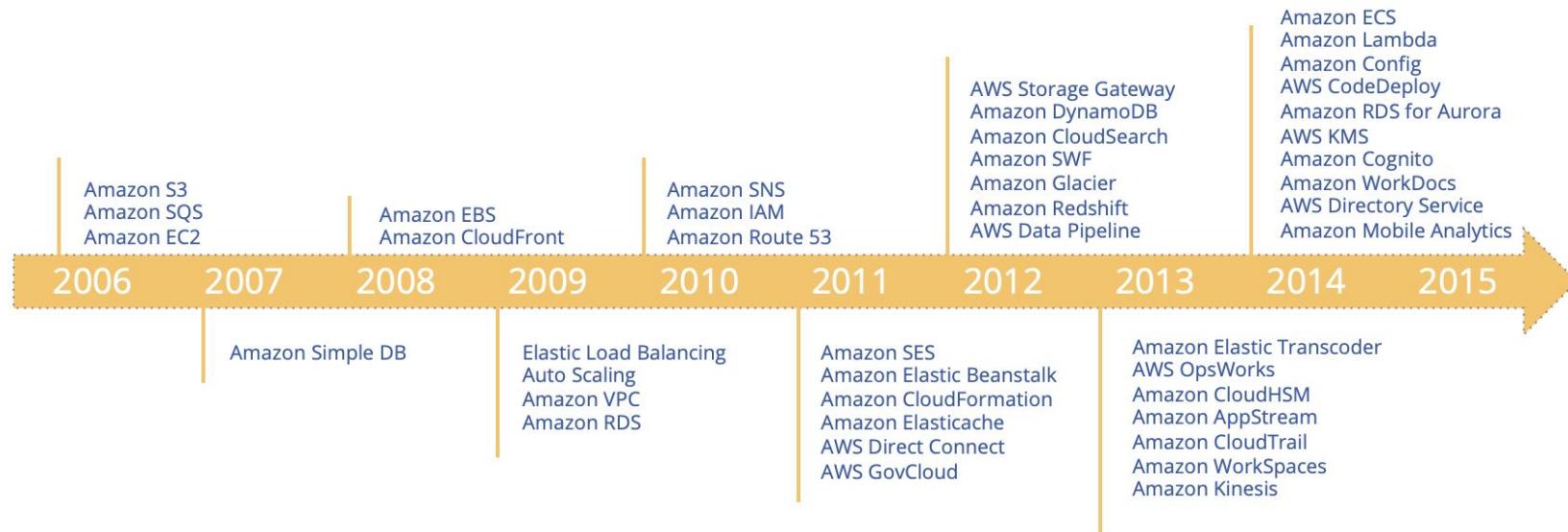
+ 600 services

2010



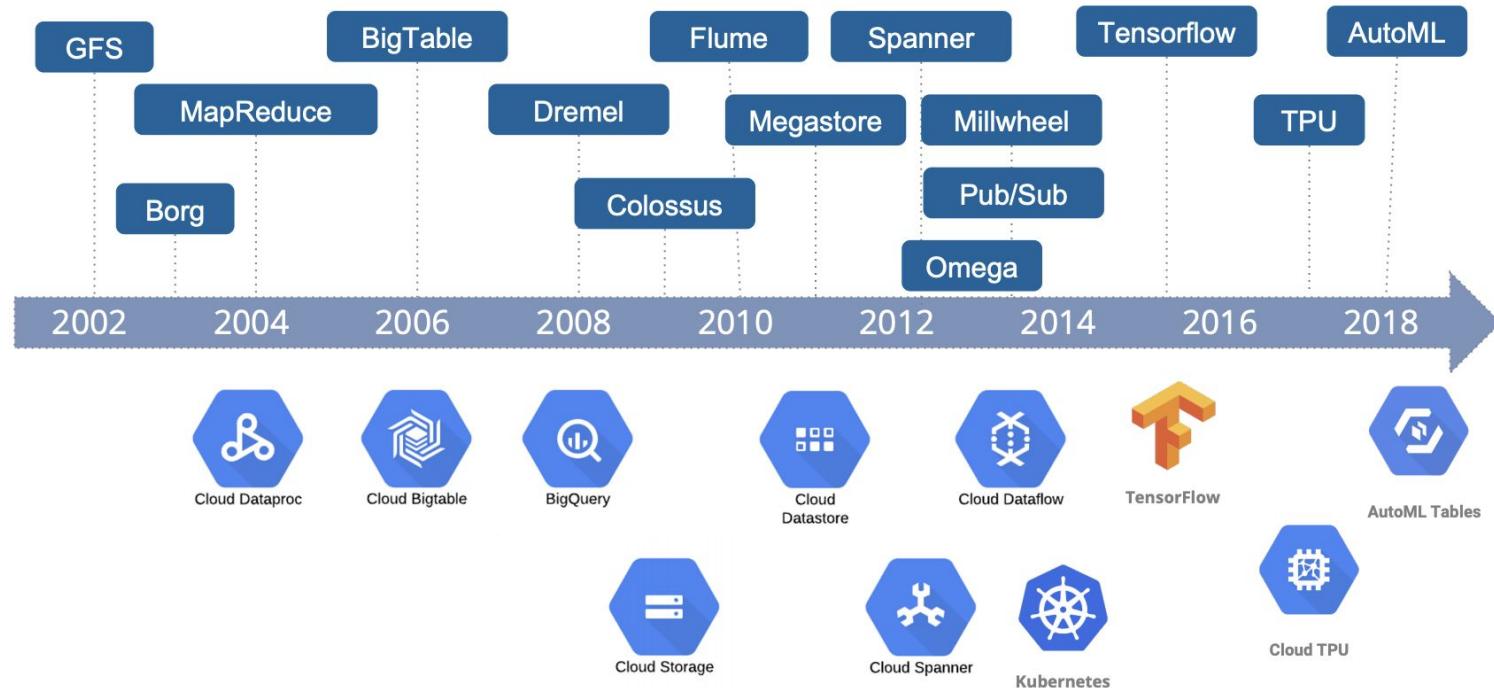
# 1. Introduction : les fondamentaux

**AWS**



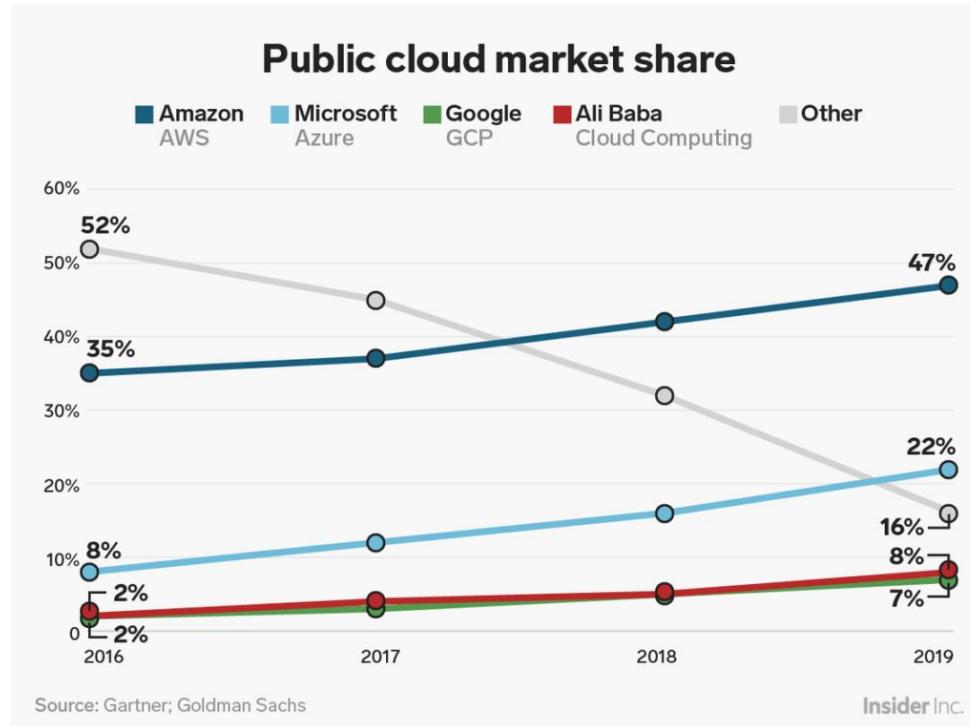
# 1. Introduction : les fondamentaux

**GCP**



## 1. Introduction : les fondamentaux

### Comment ces acteurs se partagent-ils le marché ?



# 1. Introduction : les fondamentaux

## Les solutions



Physical



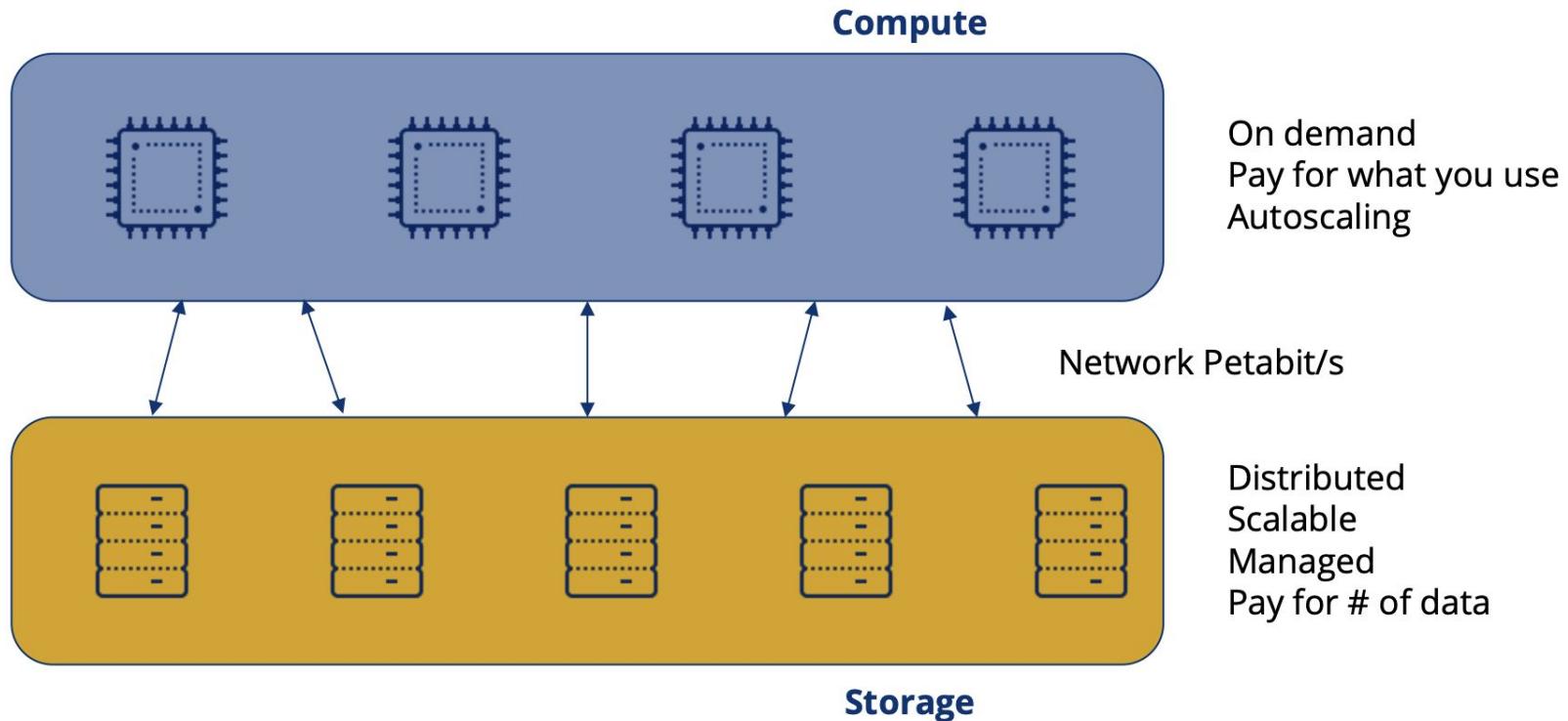
Virtualized



Serverless

## 1. Introduction : les fondamentaux

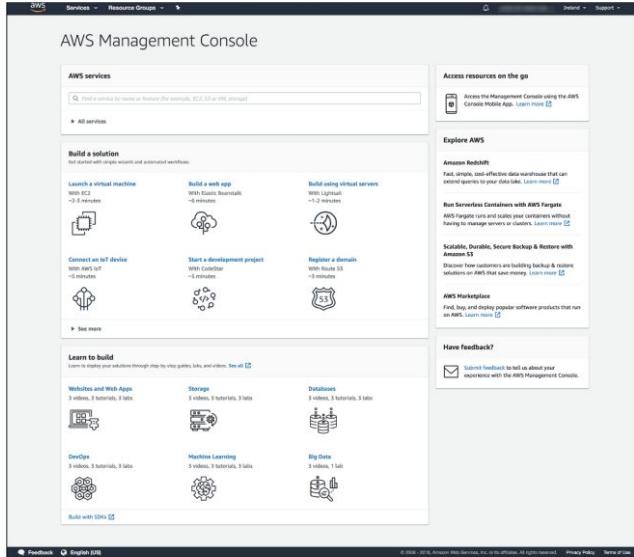
### Le stockage et le calcul



# 1. Introduction : les fondamentaux

## Comment utiliser le cloud ?

UI



CLI

```
altiok ~ $ gcloud alpha interactive
Welcome to the gcloud interactive shell environment.

Tips:
  * start by typing "gcloud " to get auto-suggestions
  * run gcloud alpha interactive --update-cli-trees to enable
    autocomplete for gsutil and kubectl
  * run gcloud alpha interactive --help for more info

Run $ gcloud feedback to report bugs or request new features.

$ g
  gcloud
    gsutil
```

API



In Production - Infrastructure as Code : Terraform, Cloudformation, Deployment Manager, ARM

## 1. Introduction : les fondamentaux

### **Le SDK (Software Development Kit)**

#### **AWS**

Boto 3

```
pip install boto3  
https://boto3.amazonaws.com/v1/documentation/api/latest/index.html
```

#### **GCP**

```
pip install --upgrade google-cloud-XXX  
https://github.com/googleapis/google-cloud-python#google-cloud-python-client
```

#### **Azure :**

```
pip install azure-XXX  
https://pypi.org/project/azure/
```

**2**

---

# **Le Stockage**

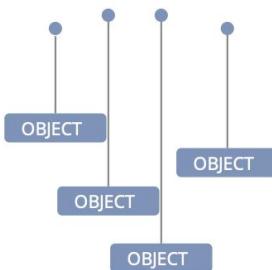
---

[RETOUR SOMMAIRE](#)

## 2. Le Stockage

# Object Storage vs File Storage vs Block Storage

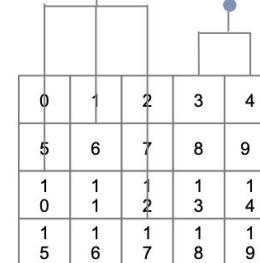
**Object Storage**



Client app (API)  
Internet  
Images, PDFs, Video  
Custom Metadata



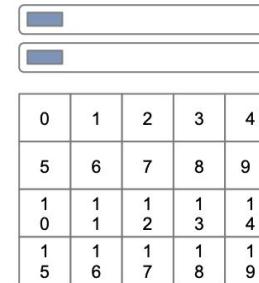
**File Storage**



User, Client via OS  
LAN / 10Gb  
user data, web content  
Fixed FS Attributes



**Block Storage**



Operating System  
Dedicated Fibre / 100Gb  
OS, Database, Transactional  
Fixed Sys Attributes



## 2. Le Stockage

### Cloud storage

- Stockage d'objets sécurisé, durable, hautement évolutif et rentable
- Les objets peuvent être structurés ou non structuré
- Facile à utiliser avec une interface de services Web qui rend facile le stockage et le téléchargement des objets
- Il est possible de stocker de grande quantité de données de n'importe où sur le Web
- On parle souvent de 'Data Lake'



Amazon Simple  
Storage Service (S3)



Cloud Storage



Storage blob



Microsoft Azure

## 2. Le Stockage

### Le stockage à plusieurs niveaux

- Stockage d'objets distribué (**Distributed object storage**) ou magasins de valeurs-clés redondants dans lesquels vous pouvez stocker des données objets.
- “**Cool Storage**” ou services de stockage conçus pour stocker des sauvegardes de données. Accès peu fréquent
- Stockage d'archivage (**Cold (archival) storage**) ou services de stockage conçus pour stocker des données d'archivage à des fins de conformité ou à des fins d'analyse.

#### AWS S3

Standard  
IA  
One Zone IA  
Intelligent Tiering  
Glacier  
Glacier Deep Archive

#### Google Cloud Storage

Standard  
Nearline  
Coldline  
Archive

#### Azure Blob Storage

with GRS or RA-GRS  
with ZRS  
with LRS  
Cool tier  
Archive Tier

## 2. Le Stockage

### Accès avec CLI

<https://docs.aws.amazon.com/cli/latest/reference/s3/index.html>

<https://aws.amazon.com/blogs/developer/leveraging-the-s3-and-s3api-commands/>

```
aws s3 mb s3://mybucket
```

```
aws s3 cp test.txt s3://mybucket/test2.txt
```

```
aws s3 cp myDir s3://mybucket/ --recursive --exclude "*.jpg"
```

<https://cloud.google.com/storage/docs/gsutil>

```
gsutil mb -l asia gs://some-bucket
```

```
gsutil cp *.txt gs://my-bucket
```

```
gsutil cp -r dir gs://my-bucket
```

## 2. Le Stockage

# Accès avec le SDK

### AWS

Boto 3

```
pip install boto3  
https://boto3.amazonaws.com/v1/documentation/api/latest/index.html
```

### GCP

google-cloud-storage

```
pip install --upgrade google-cloud-storage  
https://pypi.org/project/google-cloud-storage/
```

### Azure :

```
pip install azure-storage-blob  
https://github.com/Azure/azure-sdk-for-python
```

**3**

---

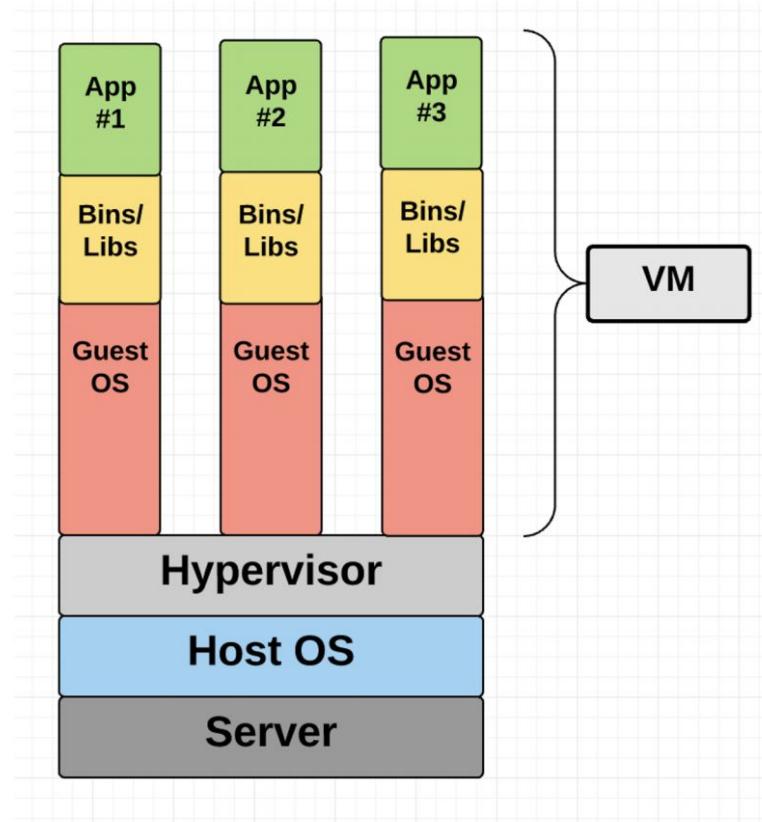
# **Le calcul**

---

[RETOUR SOMMAIRE](#)

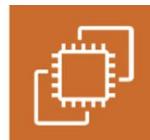
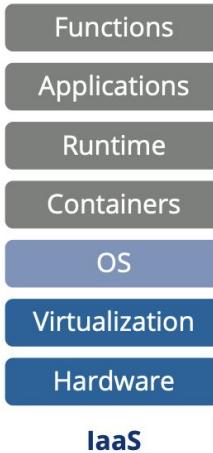
### 3. Le Calcul

## Les machines virtuelles



### 3. Le Calcul

## Les machines virtuelles



Amazon EC2



Compute  
Engine



Virtual machine

Secure, resizable compute capacity in the cloud

- Scalable and high-performance virtual machines
- Create Linux and Windows virtual machines (VMs) in seconds.

### 3. Le Calcul

## AWS : EC2 for ML

EMR supports GPU instances types

Some EC2 instances types targeted at ML tasks

- compute optimized
- accelerated computing (GPUs)

ml.\* instances only for SageMaker

Appropriate instance for deep learning :

- P3: 8 Tesla V100 GPU's
- P2: 16 K80 GPU's
- G3: 4 M60 GPU's (all Nvidia chips)

Conda-based Deep Learning AMI's

[AWS Deep Learning AMIs](#)

[Amazon Elastic Inference](#)

Librairies

- Tensorflow
- Keras
- MXNet
- Gluon
- Pytorch
- ...

GPU acceleration:

- CUDA
- cuDNN
- NCCL
- NVidia Driver

### 3. Le Calcul

## GCP : Compute Engine for ML

GPUs are independent resources

- can be added to some GCE instances (depends on regions/zones too)

Deep Learning VMs : [Deep Learning VM](#)

Cloud TPU : [NEXT TPU](#)

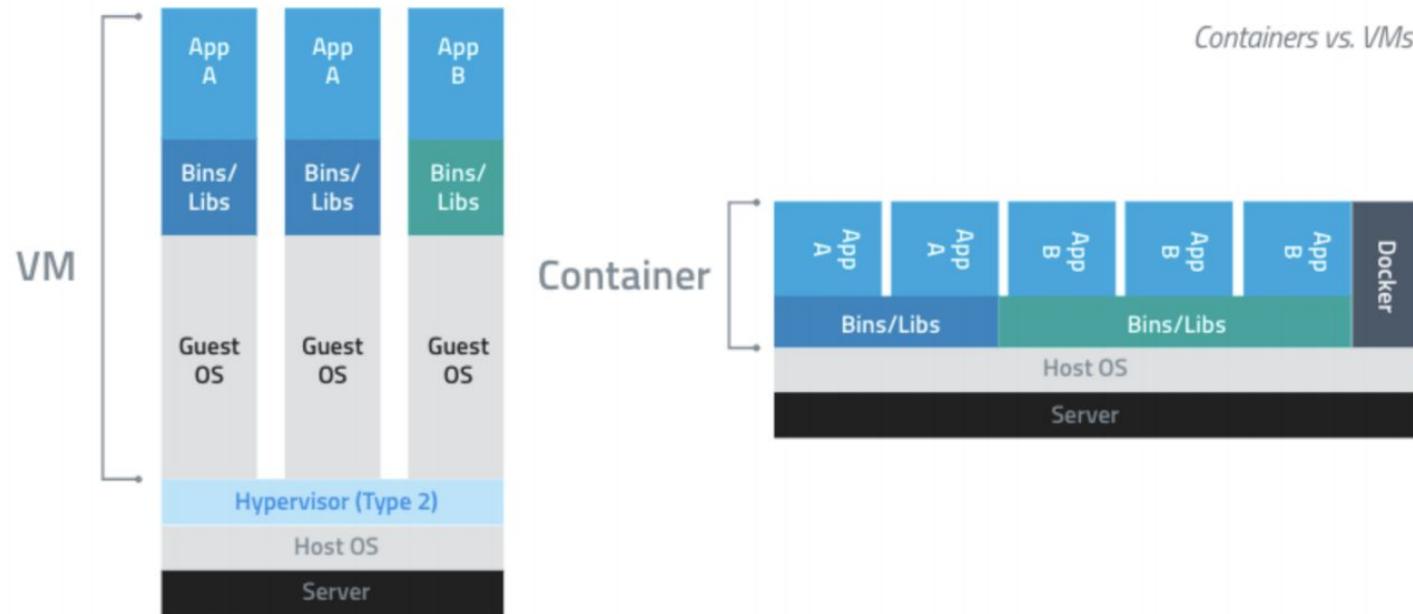
### 3. Le Calcul

## Les coûts

On demand	Reserved	Interruptible	Sole Tenant
On-demand Pay as you go	Reserved Sustained use discounts Committed use discounts Sole-tenant nodes	Spot Preemptible instances	Dedicated Sole-tenant nodes
			

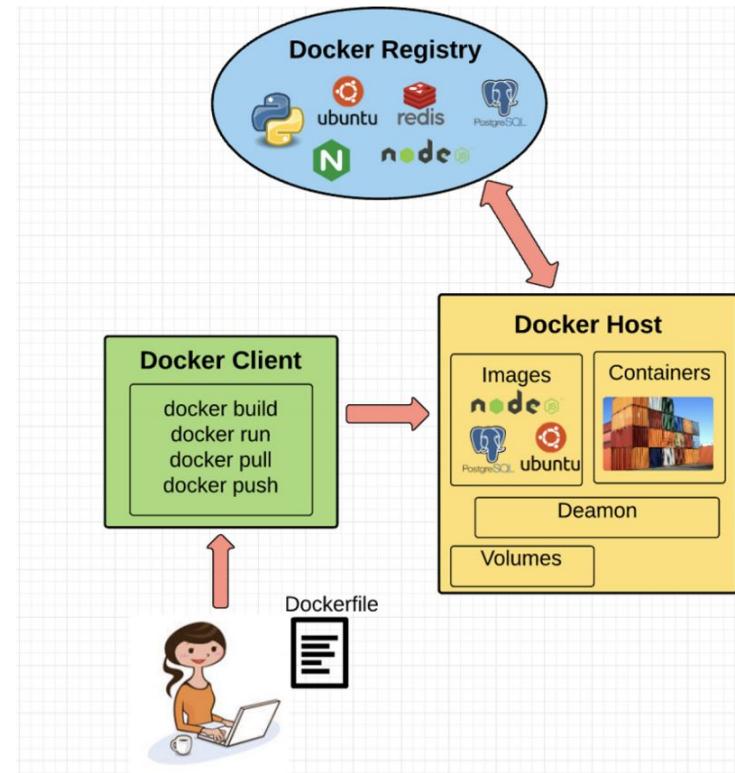
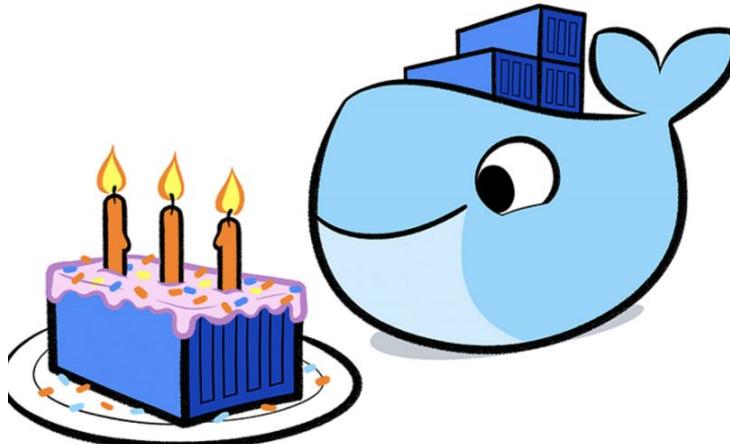
### 3. Le Calcul

## Les containers



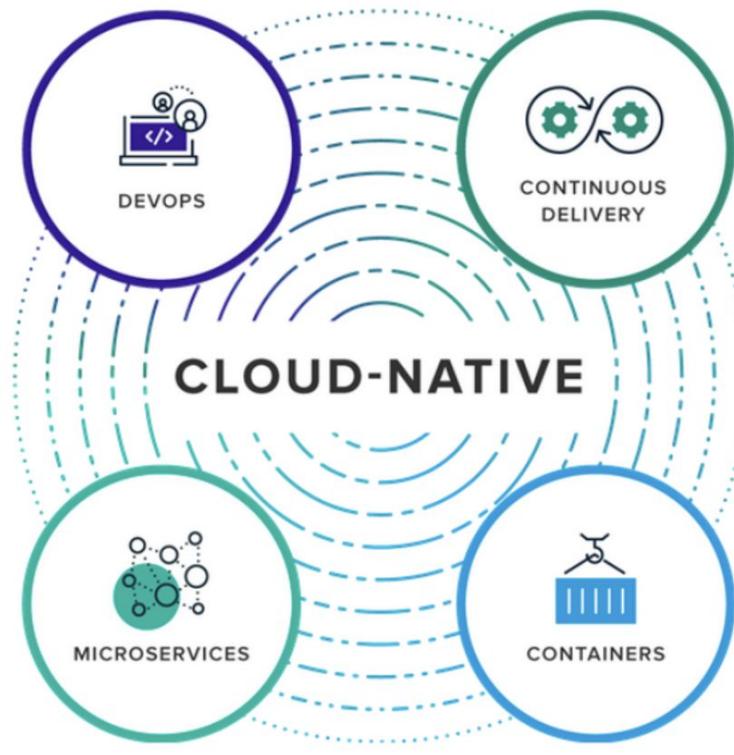
### 3. Le Calcul

## Les containers : docker



### 3. Le Calcul

## Solutions Cloud Native



### 3. Le Calcul

## Docker python web app

app.py

```
from flask import Flask
app = Flask(__name__)

@app.route("/")
def hello():
    return "Hello World!\n"

@app.route("/version")
def version():
    return "Helloworld 1.0\n"

if __name__ == "__main__":
    app.run(host='0.0.0.0')
```

requirements.txt

```
Flask==0.12
uwsgi==2.0.15
```

Dockerfile

```
FROM ubuntu:18.10
RUN apt-get update -y && \
    apt-get install -y python3-pip python3-dev
COPY requirements.txt /app/requirements.txt
WORKDIR /app
RUN pip3 install -r requirements.txt
COPY . /app
ENTRYPOINT ["python3", "app.py"]
```

4

# Machine Learning Engineering

[RETOUR SOMMAIRE](#)

#### 4. Machine Learning Engineering

## Machine Learning System

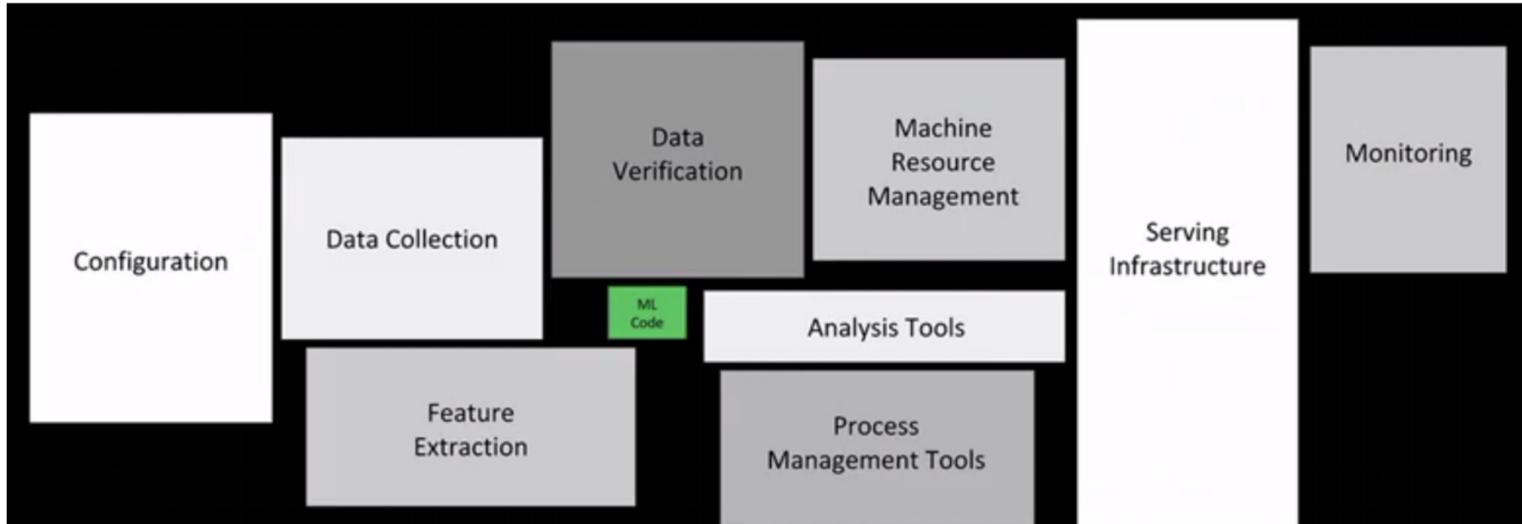
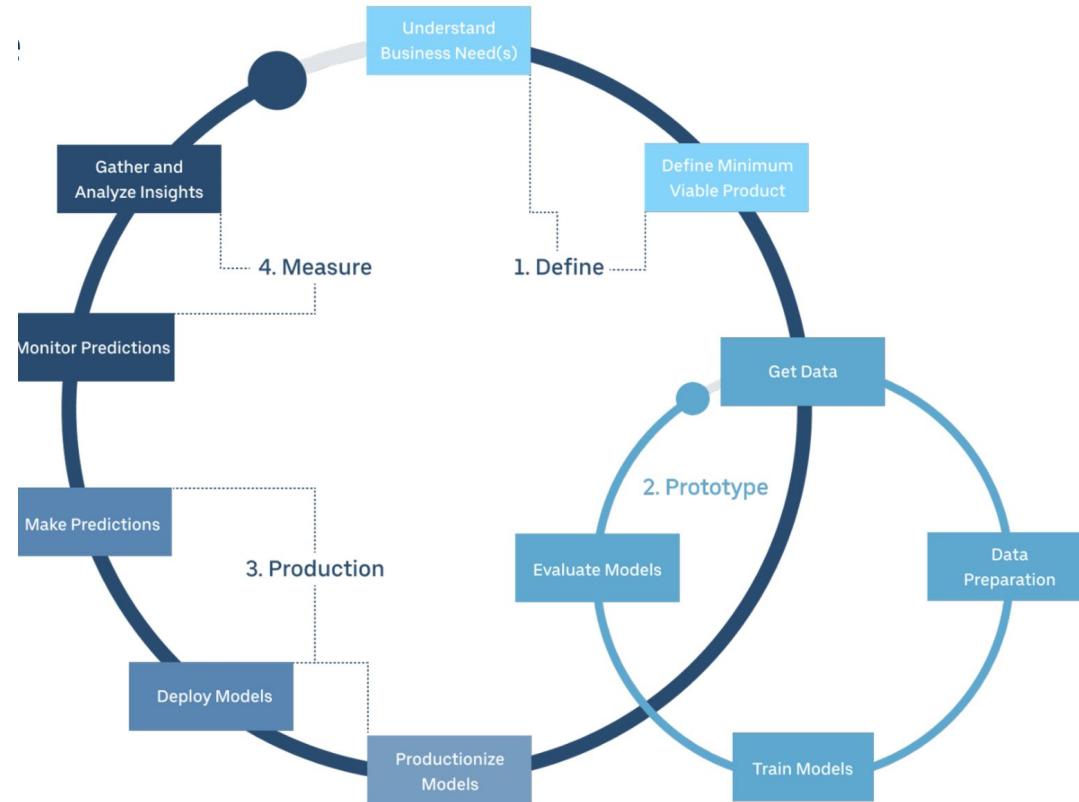


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small green box in the middle. The required surrounding infrastructure is vast and complex.

Source : Hidden Technical Debt in Machine Learning Systems, Google NIPS 2015

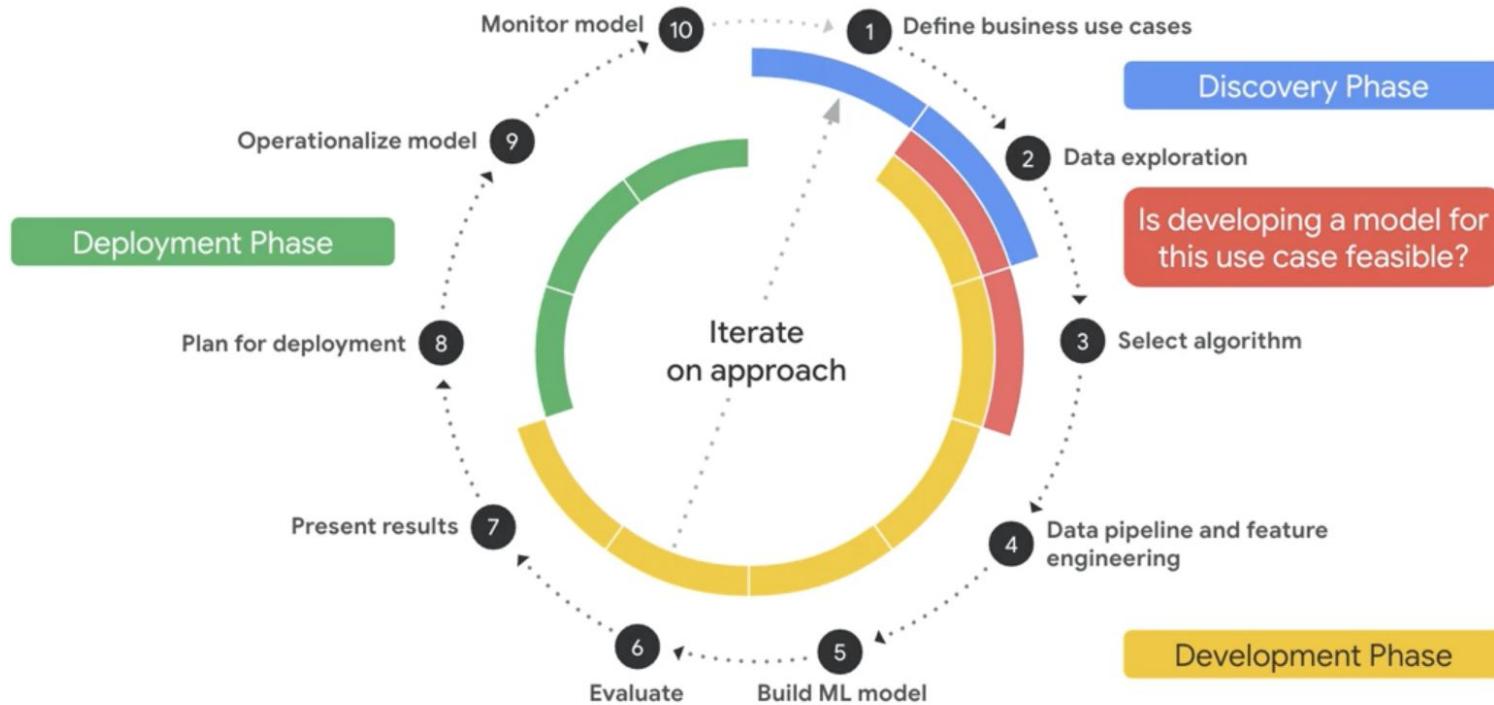
## 4. Machine Learning Engineering

### ML Pipeline



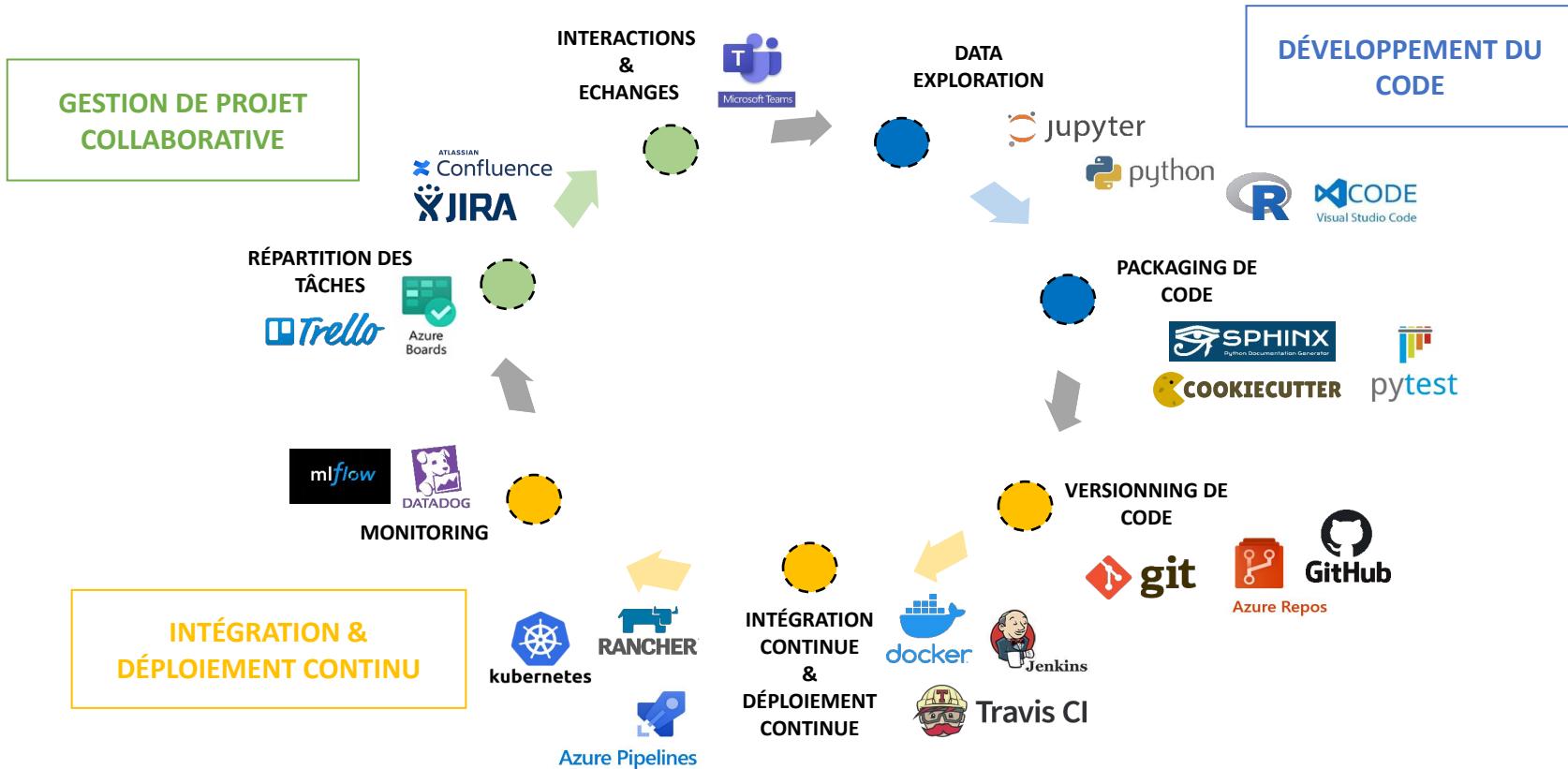
## 4. Machine Learning Engineering

### ML Pipeline



## 4. Machine Learning Engineering

### Tools chain



## 4. Machine Learning Engineering

### ML options on the cloud

#### Cloud Compute



Cloud TPU



TensorFlow on AWS



AWS Deep Learning AMIs



Apache MXNet on AWS



GPUs



Kubernetes

#### As A Service



AI Platform



Amazon SageMaker

#### Built-ins algorithms

Linear learner  
Wide and deep  
XGBoost  
Image classification  
Object detection  
k-means  
k-nn  
LDA  
PCA  
...



BigQuery ML

#### Auto ML



Cloud AutoML

Vision  
Video  
Natural Language  
Translation  
Data Tables

#### Pretrained Model as a Service



Cloud Vision API



Cloud Speech-to-Text



Cloud Video Intelligence API



Cloud Text-to-Speech



Cloud Natural Language API



Cloud Translation API



Recommendations AI

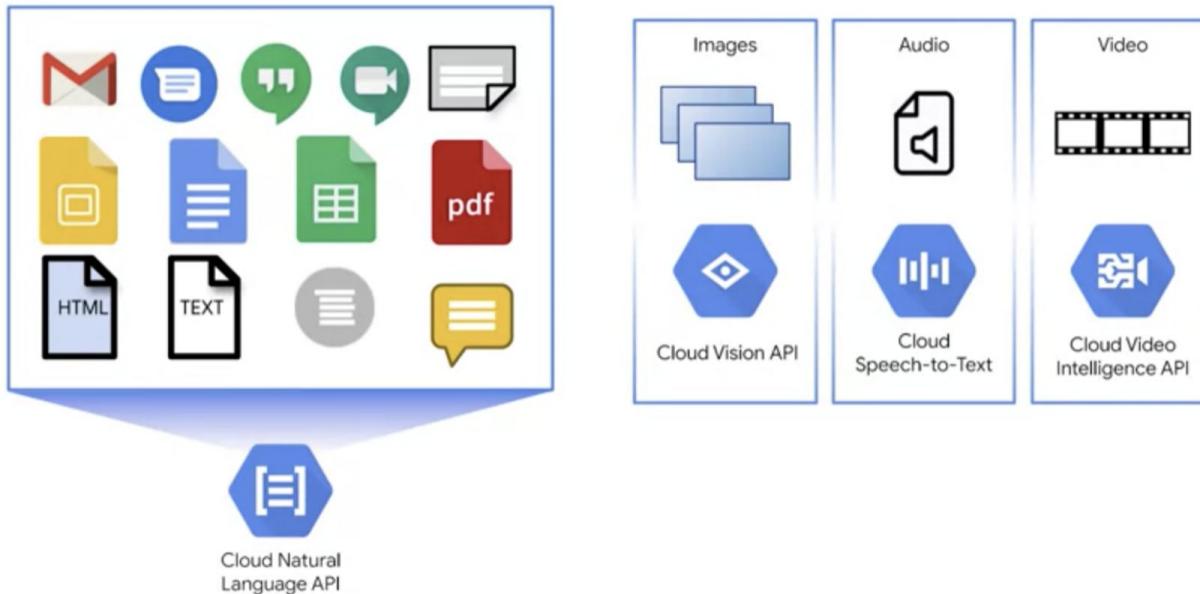
More control for advanced users

Bring your own data

Ready to go

## 4. Machine Learning Engineering

### API ML (GCP)



## 4. Machine Learning Engineering

### **API ML (AWS)**



Amazon Comprehend



Amazon Forecast



Amazon Personalize



Amazon Textract



Amazon Translate



Amazon Lex



Amazon Polly



Amazon Rekognition

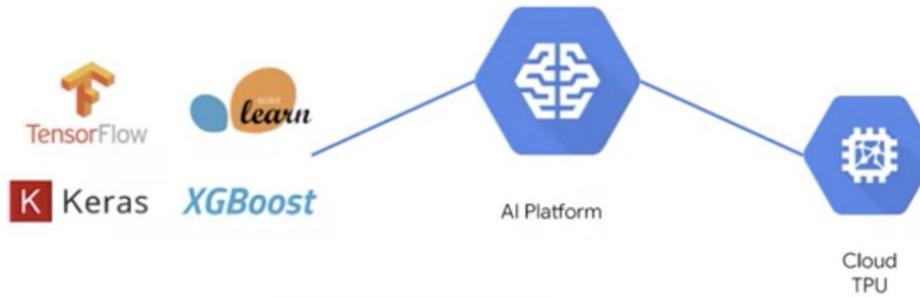


Amazon Transcribe

## 4. Machine Learning Engineering

### ML Platform : AI Platform (GCP)

Cloud AI Platform is a fully managed service for custom machine learning models



- Scales to production
- Batching and distribution of model training
- Performs transformations on input data
- Hyper-parameter tuning
- Host and autoscale predictions
- Serverless - self-tuning - manages overhead

## 4. Machine Learning Engineering

### ML Platform : SageMaker (AWS)



Amazon SageMaker

#### Build

- Preprocessing
- Ground Truth
- Notebooks

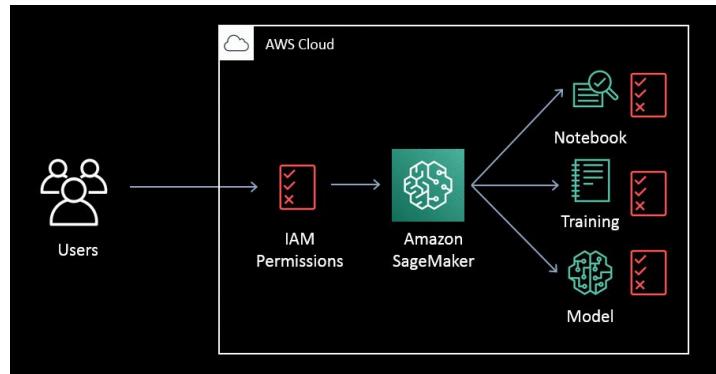
#### Train

- Built-in algorithms
- Hyperparameter tuning
- Notebooks
- Infrastructure

#### Deploy

- Realtime
- Batch
- Infrastructure
- Notebooks
- Neo (Edge)

Control :  
Console (UI), SDK, Notebooks



## 4. Machine Learning Engineering

### ML Platform : Azure ML



Azure Machine Learning

[Home](#)

Author

[Automated ML](#)

[Designer](#)

[Notebooks](#)

Assets

[Datasets](#)

[Experiments](#)

[Models](#)

[Endpoints](#)

Manage

[Compute](#)

[Datastores](#)

[Notebook VMs](#)

Welcome!

[Create new](#)

**Automated ML**  
Automatically train and tune a model using a target metric.  
[Start now](#)

**Designer**  
Drag-n-drop interface from prepping data to deploying models.  
[Start now](#)

**Notebooks**  
Code with Python SDK and run sample experiments.  
[Start now](#)

**My recent resources**

Runs			
Run Number	Experiment	Status Updated Time	Status
1	Sample_1_-_Regression...	9/27/2019, 1:38:37 PM	Completed
1474	category-based-prope...	9/18/2019, 4:37:10 PM	Completed
1475	category-based-prope...	9/18/2019, 3:49:21 PM	Completed

5

---

## Ex. Workflow

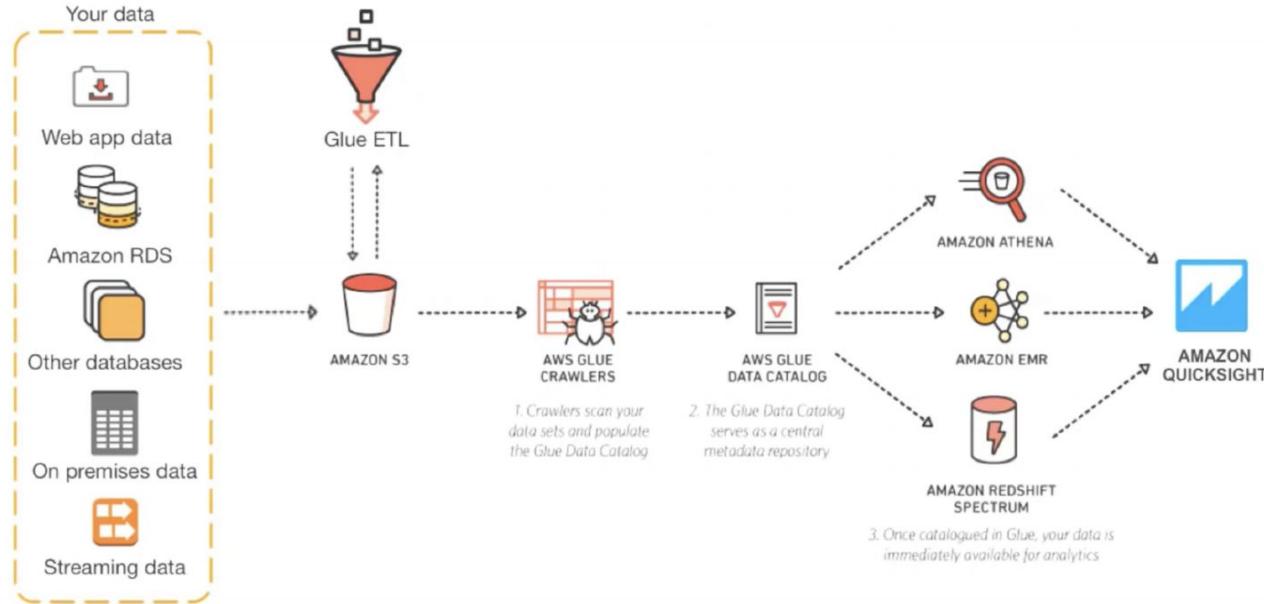
---

[RETOUR SOMMAIRE](#)

## 6. Workflow

# Full Data Engineering Pipeline : AWS

## Data Lake on S3

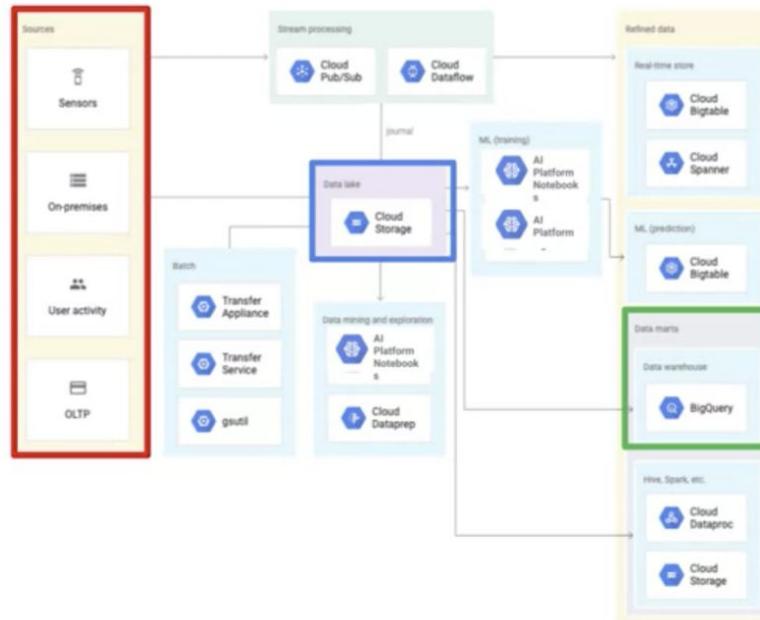


## 6. Workflow

# Full Data Engineering Pipeline : GCP

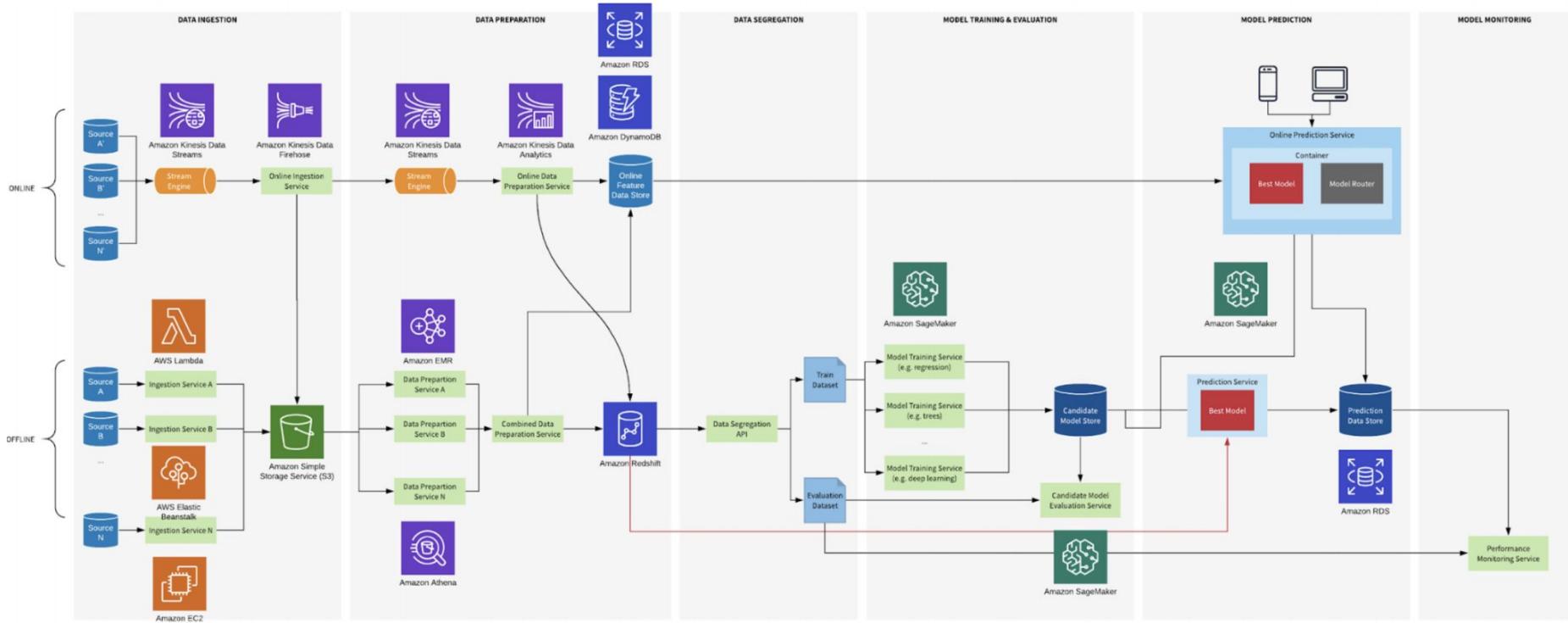
## Example Architecture

1. **Data sources**
2. **Data Lake**
3. **Data Pipelines**
4. **Data Warehouse**
5. Used for ML and analytics workloads



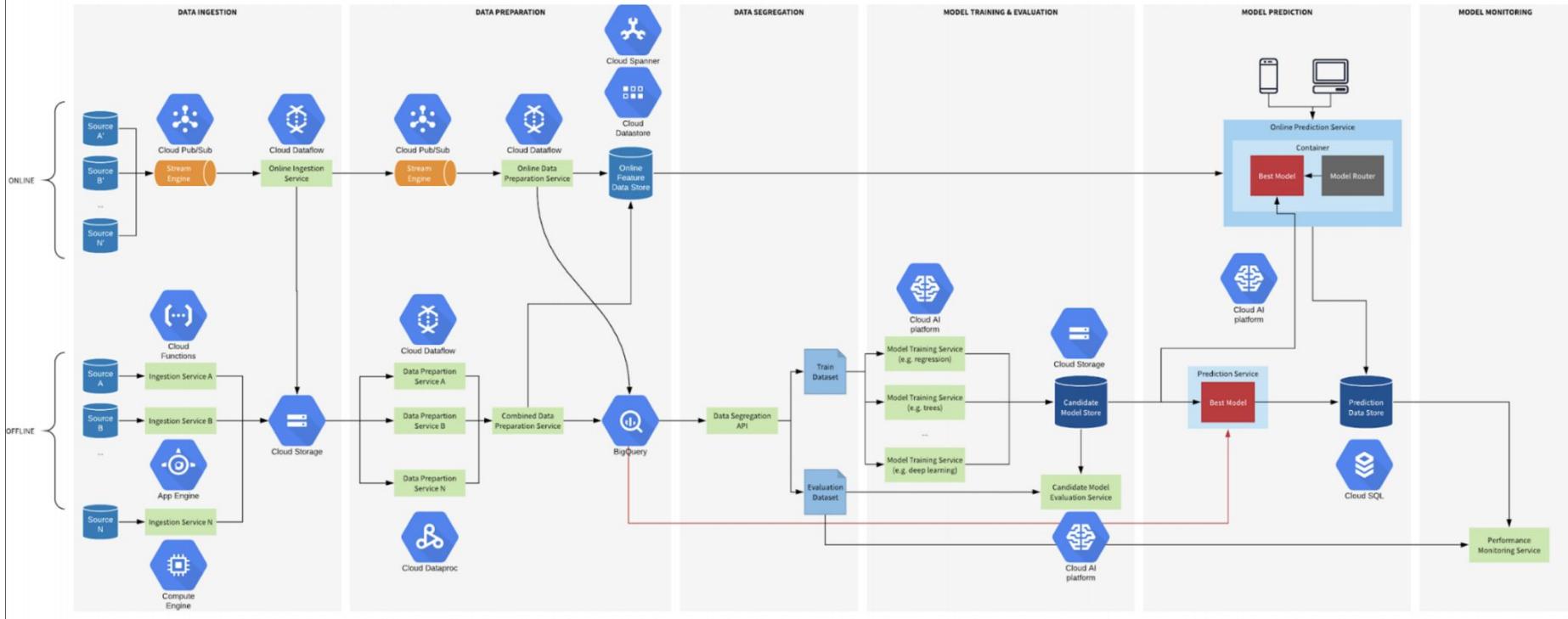
## 2. Le Stockage

### Ex. full ML Pipeline (AWS)



## 6. Workflow

### Ex. full ML Pipeline (GCP)





UNE QUESTION ?

**corentin.vasseur@decathlon.com**

Machine Learning Engineer - MLOps