# ~~Serving & training optimisation~~
# ML Pipeline

## Sprint 5 - Week 9

*INFO 9023 - Machine Learning Systems Design*
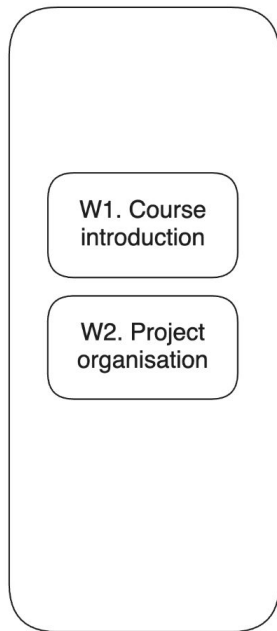
Thomas Vrancken (t.vrancken@uliege.be)
Matthias Pirlet (matthias.pirlet@uliege.be)

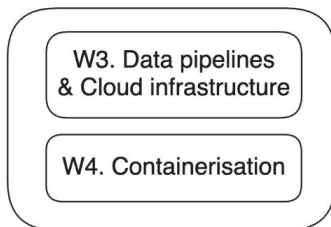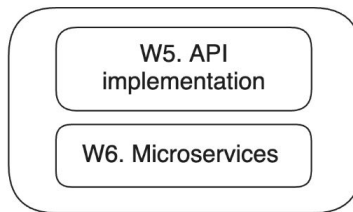*2025 Spring*

# Status on our overall course roadmap



Sprint 1:
Project organisation

- W1. Course introduction
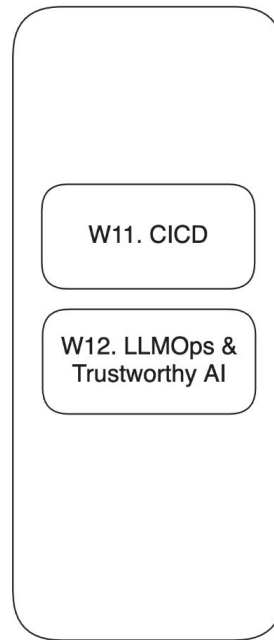- W2. Project organisation

Sprint 2:
Cloud & containerisation

- W3. Data pipelines & Cloud infrastructure
- W4. Containerisation

Sprint 3:
API implementation

- W5. API implementation
- W6. Microservices

Sprint 4: Model serving & optimisation

- W7. Model serving
- W8. Serving & training optimisation

Sprint 5: Pipeline & monitoring

- W9. Model pipeline
- W10. Monitoring & dashboarding

Order swapped

Sprint 6:
CICD

- W11. CICD
- W12. LLMOps & Trustworthy AI

# Agenda

## What will we talk about today

**Lecture**

1. ML model pipeline
2. ML platforms & orchestrators

**Directed Work**

3. Vertex Pipeline

LIÈGE
université

# ML model pipeline

# Why do we need ML pipelines?



The general idea is to not treat **ML workflows** as a one-off, but to treat them in a **reliable** and **reproducible** way.

# Why do we need ML pipelines?

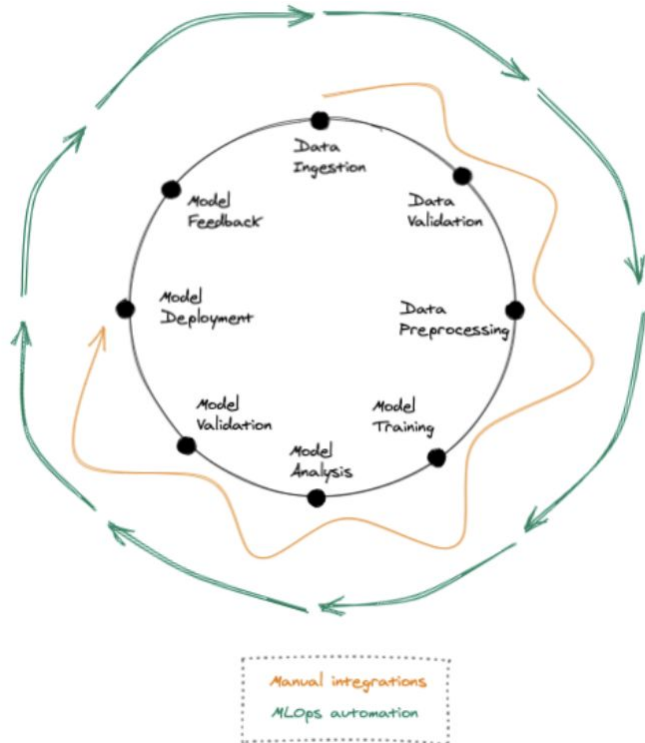The Spotify experience

## The Winding Road to Better Machine Learning Infrastructure Through Tensorflow Extended and Kubeflow

Posted on December 13, 2019 by josh baer and samuelngahane

"As we built these new Machine Learning systems, we started to hit a point where **our engineers spent more of their time maintaining data and backend systems in support of the ML-specific code** than iterating on the model itself. We realized we needed to standardize best practices and build tooling to bridge the gaps between data, backend, and ML: we needed a Machine Learning platform. "

# Why do we need ML pipelines?

Spotify accomplished an astonishing 7x for running machine learning pipelines



Spotify's Kubeflow Pipelines Runs

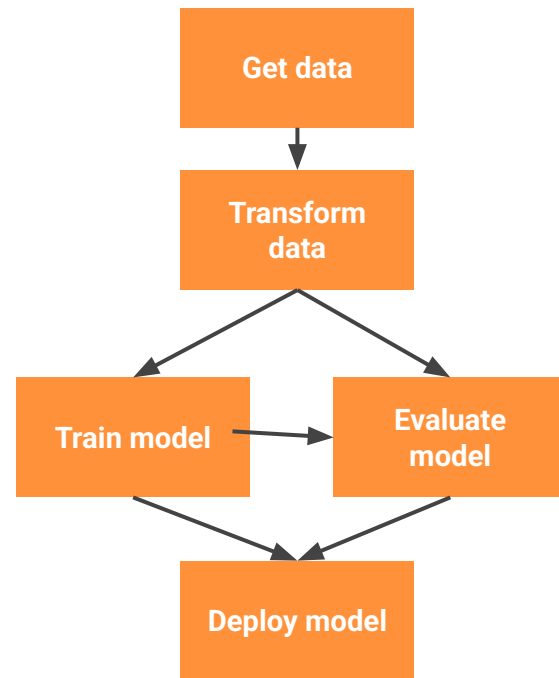# Represent your entire ML workflow as a DAG

Think about your workflow in terms of Directed Acyclic Graph (**DAGs**):

- What needs to be done sequentially vs in parallel
- Does the step process something itself or call an external service?
- Process first, implementation second

# Represent your entire ML workflow as a DAG



- Kubeflow pipelines
- Sagemaker pipelines
- Vertex pipelines
- Valohai
- MLFlow

# Definition of a ML Pipeline

"A machine learning pipeline is a series of interconnected data processing and modeling steps designed to automate, standardize and streamline the process of building, training, evaluating and deploying machine learning models."



https://www.ibm.com/topics/machine-learning-pipeline

# ML (pipeline) platforms

# ML (pipeline) platforms

# What are components of an ML pipeline?

**Components** are the building blocks of an ML pipeline. They are the nodes of the ML pipeline DAG.

They contain:

- Metadata
- Interface (input/output specifications)
- Command, Code & Environment.

Each component will be made of a script ran in a separated **Docker image**.

https://learn.microsoft.com/en-us/azure/machine-learning/concept-component?view=azureml-api-2

# Example ML Pipeline

Predict house prices with GCP
Vertex AI

# How to pass data between components?

Data staging

# How to pass data between components?

## Data staging

Data needs to be **staged** between components, stored into a Cloud service at the end of a component so it can be taken up by the next component.



Vertex pipelines integrate with Google Cloud Storage out of the box.

Google Cloud Storage

# How to pass data between components?

## Data staging

```python
from kfp.v2 import dsl, compiler
from kfp.v2.dsl import component, Input, Output, Dataset

@component
def generate_data(output_data: Output[Dataset]):
    import pandas as pd
    df = pd.DataFrame({'numbers': [1, 2, 3, 4, 5]})
    df.to_csv(output_data.path, index=False)

@component
def process_data(input_data: Input[Dataset], output_data: Output[Dataset]):
    import pandas as pd
    df = pd.read_csv(input_data.path)
    df['squared'] = df['numbers'] ** 2
    df.to_csv(output_data.path, index=False)
```

*Creates the following GCS objects*

```
gs://your-bucket/pipeline-root/runs/{run-id}/generate_data/output_data/
gs://your-bucket/pipeline-root/runs/{run-id}/process_data/output_data/
```

# Same example but using BigQuery for staging

```python
from kfp.v2 import dsl, compiler
from kfp.v2.dsl import component, Input, Output, Dataset, Model

BQ_DATASET = "your_bq_dataset"
BQ_SOURCE_TABLE = f"{PROJECT_ID}.{BQ_DATASET}.source_table"
BQ_TARGET_TABLE = f"{PROJECT_ID}.{BQ_DATASET}.processed_table"

@component
def create_bigquery_table():
    """Creates a BigQuery table and inserts sample data."""
    from google.cloud import bigquery

    client = bigquery.Client()
    schema = [
        bigquery.SchemaField("id", "INTEGER"),
        bigquery.SchemaField("value", "FLOAT"),
    ]

    table = bigquery.Table(BQ_SOURCE_TABLE, schema=schema)
    client.create_table(table, exists_ok=True)

    rows = [{"id": i, "value": float(i)} for i in range(1, 6)]

@component
def process_bigquery_data(bq_source: str, bq_target: str):
    """Reads from BigQuery, processes the data, and writes back to BigQuery."""
    from google.cloud import bigquery

    client = bigquery.Client()
    query = f"""
    SELECT id, value, value * value AS squared_value
    FROM `{bq_source}`
    """
    df = client.query(query).to_dataframe()

    job_config = bigquery.LoadJobConfig(write_disposition="WRITE_TRUNCATE")
    client.load_table_from_dataframe(df, bq_target, job_config=job_config)
```

Vertex lets you easily visualise artifacts - such as evaluation metrics

# Example ML Pipeline

## Predict house prices with GCP Vertex AI

Based on [Kubeflow Pipelines](#) (KFP)

```
1   # IMPORT REQUIRED LIBRARIES
2   from kfp.v2 import dsl
3   from kfp.v2.dsl import (Artifact,
4                           Dataset,
5                           Input,
6                           Model,
7                           Output,
8                           Metrics,
9                           Markdown,
10                          HTML,
11                          component,
12                          OutputPath,
13                          InputPath)
14  from kfp.v2 import compiler
15  from google.cloud.aiplatform import pipeline_jobs
```

# Example ML Pipeline

## Predict house prices with GCP Vertex AI

Define **components**.

Can be separate **Docker containers** or **python functions**.

*Can create a pipeline in one single notebook*

```python
1   @component(
2       base_image=BASE_IMAGE,
3       output_component_file="get_data.yaml"
4   )
5
6   def get_houseprice_data(
7       filepath: str,
8       dataset_train: Output[Dataset],
9   ):
10
11      import pandas as pd
12
13      df_train = pd.read_csv(filepath + '/train.csv')
14
15      df_train.to_csv(dataset_train.path, index=False)
```

*Data Ingestion component*

LIÈGE université

# Example ML Pipeline

## Predict house prices with GCP Vertex AI

```python
@component(
    base_image=BASE_IMAGE,
    output_component_file="preprocessing.yaml"
)

def preprocess_houseprice_data(
    train_df: Input[Dataset],
    dataset_train_preprocessed: Output[Dataset],
):

    import pandas as pd
    from src.data_preprocessing.preprocessing import data_preprocessing_pipeline

    train_df = pd.read_csv(train_df.path)

    # data_preprocessing_pipeline creates a copy of the df, removes id col, converts to
    # subtracts YearSold from temporal features and cosine transforms cyclic features.
    train_df_preprocessed = data_preprocessing_pipeline(train_df)

    train_df_preprocessed.to_csv(dataset_train_preprocessed.path, index=False)
```

*Data Preprocessing component*

```python
@component(
    base_image=BASE_IMAGE,
    output_component_file="train_test_split.yaml",
)

def train_test_split(dataset_in: Input[Dataset],
                     dataset_train: Output[Dataset],
                     dataset_test: Output[Dataset],
                     test_size: float = 0.2):

    import pandas as pd
    from sklearn.model_selection import train_test_split

    df = pd.read_csv(dataset_in.path)
    df_train, df_test = train_test_split(df, test_size=test_size, random_state=42)

    df_train.to_csv(dataset_train.path, index=False)
    df_test.to_csv(dataset_test.path, index=False)
```

*Train test split component*

# Example ML Pipeline

## Predict house prices with GCP Vertex AI

...

```
14    import pandas as pd
15    import pickle
16    import shap
17    from src.modelling.train import HousePriceModel
18    from src.utils.utils import get_image_data
19
20    TARGET = 'SalePrice'
21
22    # Read train and test data
23    train_data = pd.read_csv(dataset_train.path)
24    test_data = pd.read_csv(dataset_test.path)
25
26    # Instantiate the model class
27    house_price_model = HousePriceModel(test_data.copy(),    #we perform hyperparameter t
28                                        target=TARGET,
29                                        n_kfold_splits=3,
30                                        n_trials=100,
31                                        random_state=42)
32
33    # Create X_train and y_train
34    X_train = train_data.drop(TARGET, axis=1)
35    y_train = train_data[TARGET]
```

...

*Model training component*

```
1     @component(
2         base_image=BASE_IMAGE,
3         output_component_file="model_evaluation.yaml"
4     )
5     def evaluate_houseprice(
6         houseprice_model: Input[Model],
7         metrics_baseline: Output[Metrics],
8         metrics_train: Output[Metrics],
9         metrics_test: Output[Metrics]):
10
11        import pickle
12
13        file_name = houseprice_model.path
14        with open(file_name, 'rb') as file:
15            model_data = pickle.load(file)
16
17        scores = model_data["scores_dict"]
18
19        def log_metrics(scores, metric):
20            for metric_name, val in scores.items():
21                metric.log_metric(metric_name, float(val))
22
23        log_metrics(scores["baseline_scores"], metrics_baseline)
24        log_metrics(scores["train_scores"], metrics_train)
25        log_metrics(scores["test_scores"], metrics_test)
```

*Model evaluation component*

# Example ML Pipeline

## Predict house prices with GCP Vertex AI

```python
17    from google.cloud import aiplatform as vertex_ai
18    from pathlib import Path
19
20    # Checks existing Vertex AI Enpoint or creates Endpoint if it is not exist.
21    def create_endpoint ():
22        endpoints = vertex_ai.Endpoint.list(
23        filter='display_name="{}"'.format(model_endpoint),
24        order_by='create_time desc',
25        project=gcp_project,
26        location=gcp_region,
27        )
28        if len(endpoints) > 0:
29            endpoint = endpoints[0] # most recently created
30        else:
31            endpoint = vertex_ai.Endpoint.create(
32                display_name=model_endpoint,
33                project=gcp_project,
34                location=gcp_region
35            )
36        return endpoint
37
```
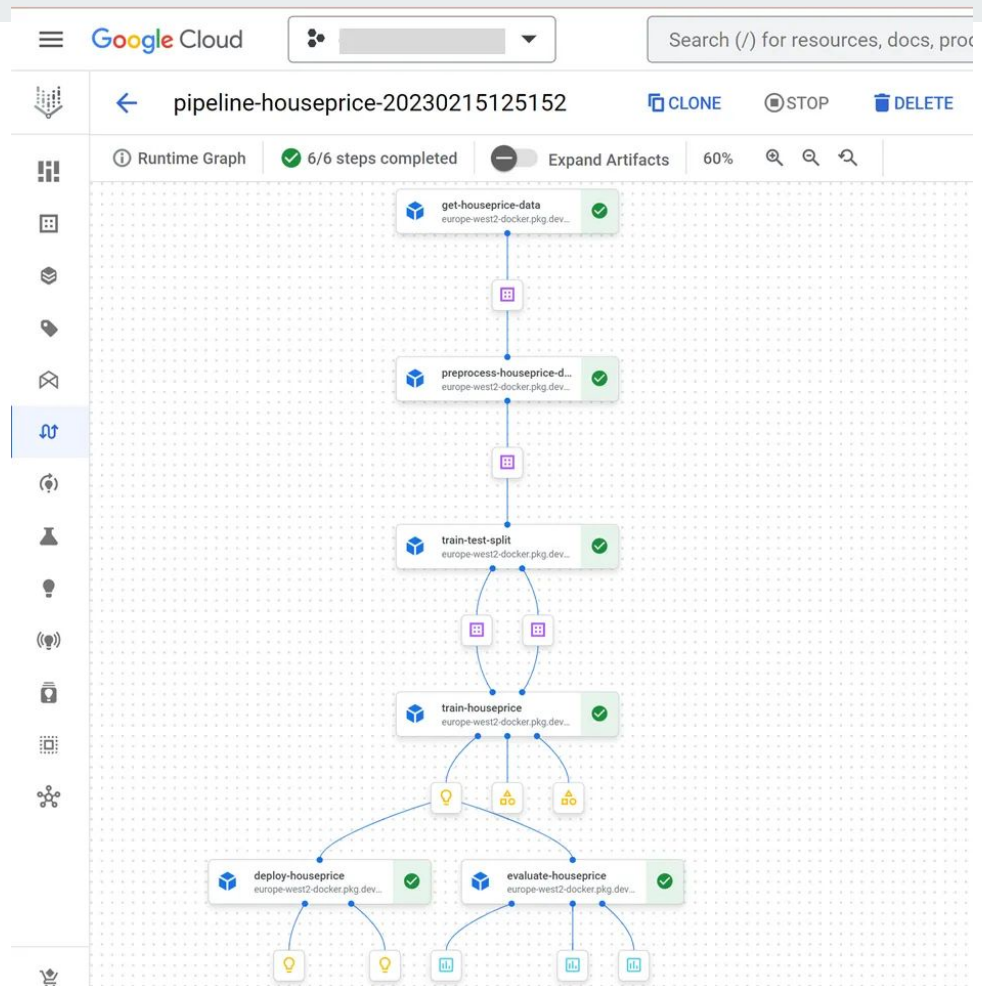
```python
40        # Uploads trained model to Vertex AI Model Registry or creates new model version in
41    def upload_model ():
42        listed_model = vertex_ai.Model.list(
43        filter='display_name="{}"'.format(display_name),
44        project=gcp_project,
45        location=gcp_region,
46        )
47        if len(listed_model) > 0:
48            model_version = listed_model[0] # most recently created
49            model_upload = vertex_ai.Model.upload(
50                    display_name=display_name,
51                    parent_model=model_version.resource_name,
52                    artifact_uri=str(Path(model.path).parent),
53                    serving_container_image_uri=serving_container_image_uri,
54                    location=gcp_region,
55                    serving_container_predict_route="/predict",
56                    serving_container_health_route="/health"
57                )
```
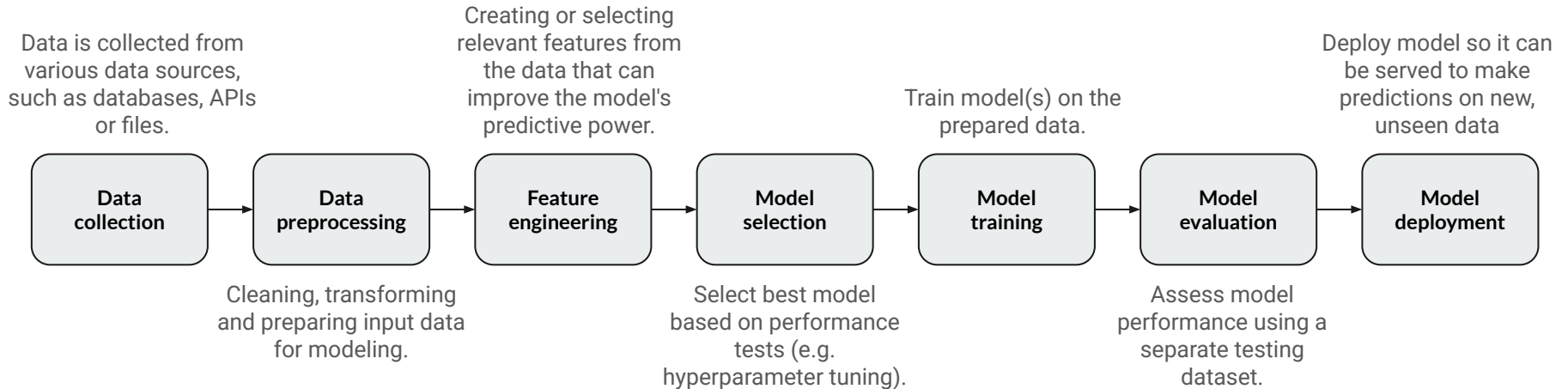
*Model deployment component*

# Example ML Pipeline

Predict house prices with GCP Vertex AI

# Typical components of an ML pipeline

Data is collected from various data sources, such as databases, APIs or files.

Creating or selecting relevant features from the data that can improve the model's predictive power.

Train model(s) on the prepared data.

Deploy model so it can be served to make predictions on new, unseen data

| Data collection | → | Data preprocessing | → | Feature engineering | → | Model selection | → | Model training | → | Model evaluation | → | Model deployment |

Cleaning, transforming and preparing input data for modeling.

Select best model based on performance tests (e.g. hyperparameter tuning).

Assess model performance using a separate testing dataset.

https://www.ibm.com/topics/machine-learning-pipeline

# Avoiding the endless POC loop…

- Let's change the model architecture
- Let's try with another model
- Let's gather more labeled data
- Ah wait, it was the wrong data, here is the new data, can we get results tomorrow?
- What were the accuracies again with the other model?
- Oops not sure what notebook was executed to get that model
- Could we try out the model to get feedback?
- …

*The key to a **successful POC outcome** is to **plan for production** during the proof of concept.*

# Benefits of an ML Pipeline

- **Modularization**: Breaks ML workflows into reusable components, making development faster and easier to manage.

- **Reproducibility**: Ensures experiments and results can be consistently recreated by tracking configurations and data.

- **Efficiency**: Automates repetitive tasks like data preprocessing and model training, saving time and effort.

- **Scalability**: Seamlessly handles large datasets and complex models by leveraging cloud infrastructure.

- **Experimentation**: Enables quick testing of different models and hyperparameters without manual reconfiguration.

- **Deployment**: Streamlines moving models from development to production with automated workflows.

- **Collaboration**: Facilitates teamwork by structuring workflows and sharing components across teams.

- **Version control and documentation**: Tracks changes in datasets, models, and configurations, ensuring transparency and traceability.

# Triggers

**Timed**

- **Schedule**: Pipeline runs repetitively in relation to the creation time of the Recurring Run. Set the unit (minutes, hours, etc.) and the scalar that goes with it

**Event based**

- **Manual**: Launch the pipeline manually
- **Triggered**: As part of another service (e.g. if users upload a new batch of training in a specific GCS bucket)

```
#  ┌──────────── minute (0–59)
#  │ ┌────────── hour (0–23)
#  │ │ ┌──────── day of the month (1–31)
#  │ │ │ ┌────── month (1–12)
#  │ │ │ │ ┌──── day of the week (0–6) (Sunday to Saturday;
#  │ │ │ │ │                          7 is also Sunday on some systems)
#  │ │ │ │ │
#  │ │ │ │ │
# * * * * * <command to execute>
```

# Microservices vs ML Pipeline

| | Microservices | ML Pipeline |
|---|---|---|
| **What is it?** | One full application made of different services that are calling each other to serve the overall purpose of the application. | Ran a single workflow made of components that run sequentially. |
| **Is made of…** | Services (~= APIs) | Components. Single process that is ran once (e.g. data preparation, model training, evaluation and deployment). |
| **Utilisation** | The microservices stay up. | One time run (triggered/scheduled). |
| **Purpose** | All over software engineering. | ML. |

LIÈGE université

# Example ML Pipeline (Kubeflow pipeline KFP)

# Example ML Pipeline (Sagemaker)

# Example ML Pipeline (Azure ML)

# ML Platforms

# ML (pipeline) platforms

# Sagemaker

Amazon SageMaker is an automated platform and comprehensive suite of tools that simplifies the development, training and deployment of machine learning (ML) models. It reduces the complexity of model development by providing a web-based interface for creating ML pipelines and pre-built algorithms.

Key features:
1. Notebooks
2. Training Job
3. Model registry
4. Prediction
5. Pipelines

# Sagemaker: Notebooks

Machine learning (ML) compute instance running the Jupyterlab. Use Jupyter notebooks in your notebook instance to prepare and process data, write code to train models, deploy models to SageMaker hosting, and test or validate your models.

**Advantages**:
- No use of data locally
- Select desired (GPUs)
- Better memory
- Manage access
- Collaborate as a team

*Equivalent*

Vertex Workbench

Azure Notebooks

https://docs.aws.amazon.com/sagemaker/latest/dg/gs-setup-working-env.html

# Sagemaker: Training Job

Managed way of launching a model training

1. Either be fully specify the ML algorithm, hyperparameters and input data location.
2. Or customize training job by selecting specific Sagemaker instance types and adding software libraries.

**Advantages**:
- Access more compute (GPUs)
- Access more memory
- Consistency
- Automation

https://docs.aws.amazon.com/sagemaker/latest/dg/how-it-works-training.html
https://docs.aws.amazon.com/sagemaker/latest/dg/debugger-supported-frameworks.html

**SageMaker AI-supported frameworks and algorithms**

TensorFlow

PyTorch

MXNet

XGBoost

SageMaker AI generic estimator

*Equivalent*

Vertex Training

Azure ML Training

LIÈGE université

# Sagemaker: Model registry

SageMaker Model Registry is a service for packaging model artifacts with deployment information. That information includes your deployment code, what type of container to use and what type of instance to deploy. The Model Registry decreases time to deployment, as you already have the necessary information about how the model should be deployed.

**Advantages**:
- Consistency
- Automation
- Scalability
- Time to deployment

*Equivalent*

Vertex AI Model Registry

Azure ML Model Registry

LIÈGE université

# Sagemaker: Prediction

After you've trained and deployed your model in SageMaker, you can use it to generate predictions based on new data. There are two main ways to do this.

- Endpoint deployment: Deploy your model to an endpoint, allowing users and applications to send API requests to get model predictions in real time.
- Batch predictions: Generate predictions on large amounts of data without needing an immediate response..

**Advantages**:
- Consistency
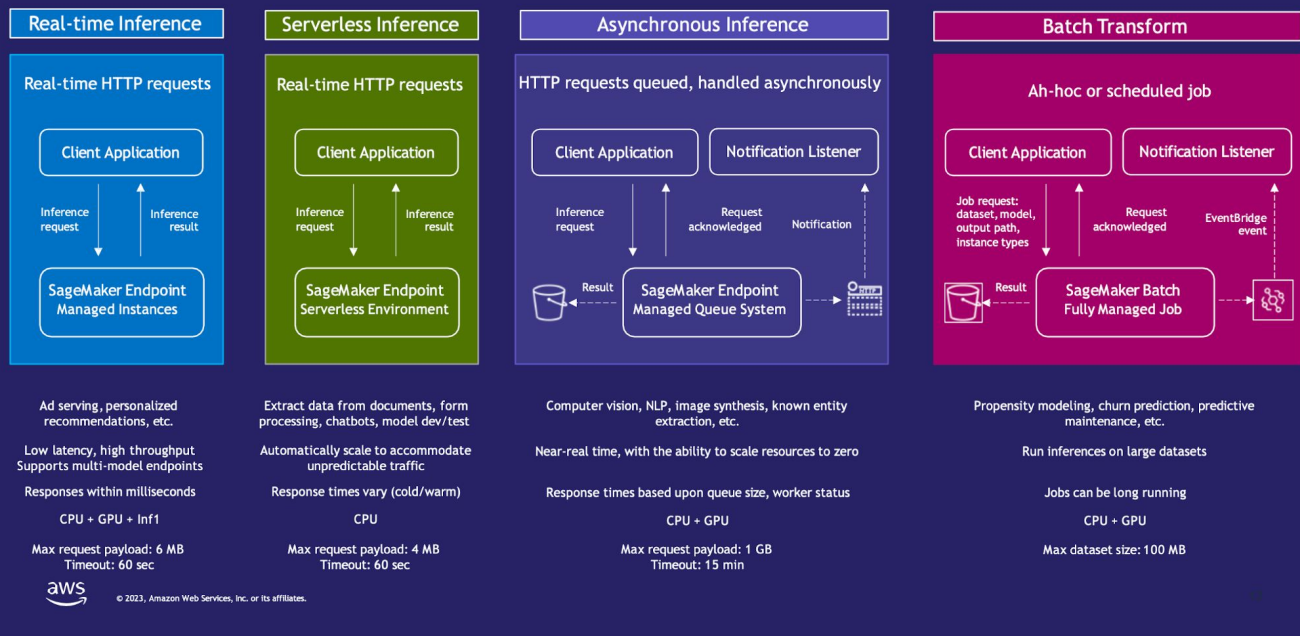- Automation
- Scalability
- Time to deployment

*Equivalent*

Vertex AI Prediction

Azure ML Prediction

https://docs.aws.amazon.com/sagemaker/latest/dg/how-it-works-training.html

# Sagemaker: Prediction

# Sagemaker: Pipeline

SageMaker Pipelines is a tool for building, deploying and managing end-to-end ML workflows. Users can create automated workflows that cover data preparation, model training and deployment — all from a single interface. Because Pipeline logs everything, you can easily track and re-create models.

**Advantages**:
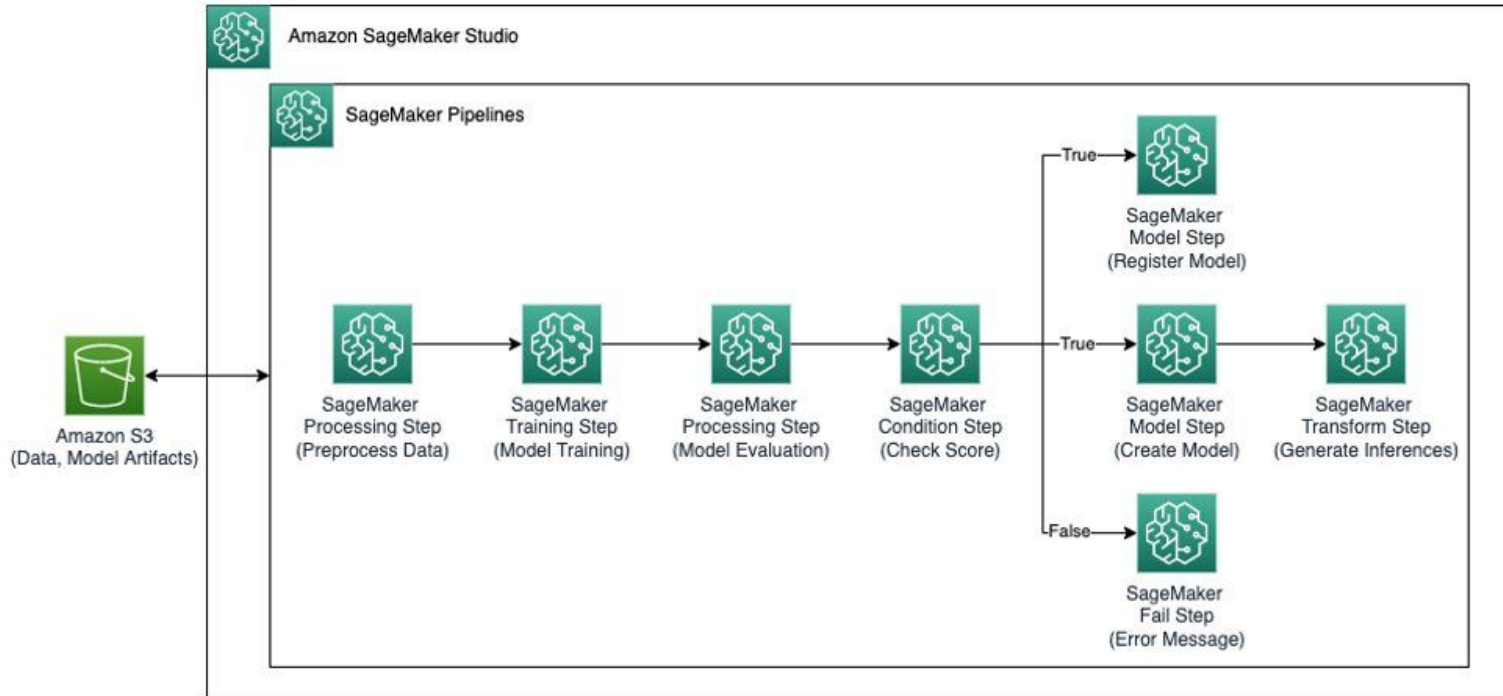- Resource consumption
- Consistency
- Automation
- Scalability
- Time to deployment

*Equivalent*

Vertex AI Pipeline

Azure ML Pipeline

LIÈGE université

# Sagemaker: Pipeline (example)

# Directed work:
# Vertex Pipeline

# Wrap-up

# Lecture summary

(split over next week)

| Topic | Concepts | To know for… | |
| --- | --- | --- | --- |
| | | **Project** | **Exam** |
| ML model pipeline | <ul><li>What it is</li><li>Standard steps of ML pipelines</li></ul> | | Yes |
| ML platforms & orchestrators | <ul><li>Sagemaker offerings (workbench, training, registry. Predictions & pipelines)</li></ul> | Possibly (Vertex equivalent) | |
| Vertex Pipeline | <ul><li>Build a vertex pipeline</li></ul> | Possibly | |

# Project objective for sprint 5

| # | Week | Work package | Requirement |
|---|------|-------------|-------------|
| 5.1 | W09 | Run your model training as a job in the Cloud. You can implement this in different ways:<br>• Containerise your training script and run it on a VM in the Cloud (e.g. on EC2 or on Cloud Run, example, )<br>• Use a managed service such as Vertex Training or Sagemaker Training<br><br>**Attention**: This can incur Cloud **costs**. Make sure to use a platform where you have credits and not burn through them. You can ask for support from the teaching staff in that regard. | Optional |
| 5.2 | W10 | Build a simple user interface or dashboard to show your results and deploy it on the Cloud. | Optional |

# That's it for today!


MADE IT THROUGH THE LECTURE
SEE YOU NEXT WEEK !
imgflip.com