**Projet 3a CS - Présentation ML Ops**

1. Schéma global d'une stack de ML Ops

https://mllops.vercel.app/



2. Présentation succincte de chaque technologie

EXPERIMENT TRACKING

Mlflow

EXPERIMENTATION

Apache Zeppelin

Jupyter

DATA VERSIONING

Apache Zeppelin

Jupyter

CODE VERSIONING

Git

GitLab

PIPELINE ORCHESTRATION

ZenML

Dagster

MODEL REGISTRY

Mlflow

MODEL SERVING

BentoML

NVIDIA Triton

ARTIFACT TRACKING

Mlflow

MODEL MONITORING

Prometheus

Grafana

Evidently

React Flow

| | | | |
|---|---|---|---|
| DS | | | (Notebook) / App python : <br> - Data preprocessing <br> - Model definition <br> - Model training <br> - Model testing / evaluation |
| DE | ML Eng / ML Ops | | Model packaging (Bento ML, Kserve, triton) <br> Model registry / versioning (mlflow, bentoml x yatai) <br> Model deployment  (Bento ML, triton) <br> Model Serving (Bento ML, triton) <br> Data versioning (DVC) <br> Experiment / Artifact tracking de prod / live (mlflow, W&B, Neptune) <br> -> équipe DS sur le sujet <br> Pipeline orchestration (ré-entrainement, validation, de déploiement) (ZenML x airflow, kubeflow) <br> Model monitoring (prometheus, grafana, loki) |
| | Runtime engine / pipeline | | Kafka, Spark, Flink, Temporal |

2 grandes familles de pipeline :
- Training / fine-tuning (exemple : fineweb)
- Inference (exemple : fineweb)

3.  Présentation de frameworks de ML Ops (ZenML, Flyte, Ray, ...)

Zen ML
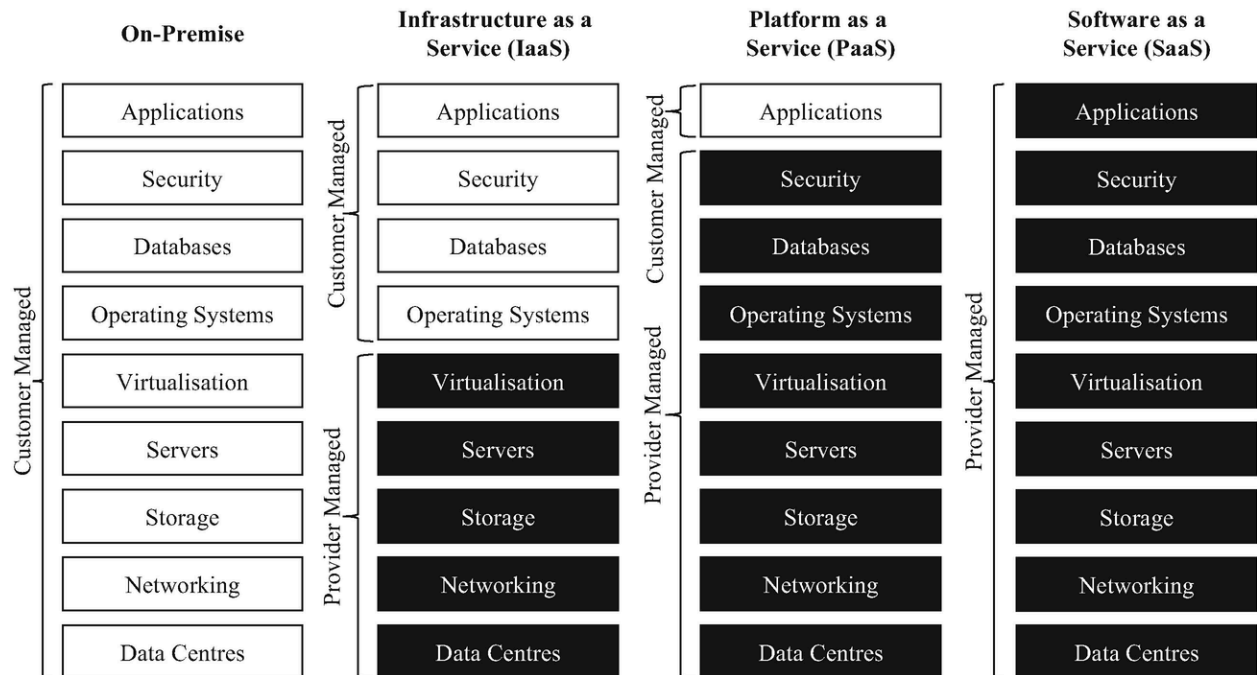https://www.zenml.io/
https://docs.zenml.io/
https://docs.zenml.io/user-guide/starter-guide
https://docs.zenml.io/how-to/build-pipelines
https://docs.zenml.io/stack-components/component-guide

Flyte
https://flyte.org/

Ray
https://ray.io

4. Présentation des solutions Cloud (GCP Vertex AI, AWS Sagemaker, Azure ML)

| On-Premise | Infrastructure as a Service (IaaS) | Platform as a Service (PaaS) | Software as a Service (SaaS) |
|---|---|---|---|
| Applications | Applications | Applications | Applications |
| Security | Security | Security | Security |
| Databases | Databases | Databases | Databases |
| Operating Systems | Operating Systems | Operating Systems | Operating Systems |
| Virtualisation | Virtualisation | Virtualisation | Virtualisation |
| Servers | Servers | Servers | Servers |
| Storage | Storage | Storage | Storage |
| Networking | Networking | Networking | Networking |
| Data Centres | Data Centres | Data Centres | Data Centres |

[Source wikipedia](#)

- Ready to use task-specific solutions (almost SaaS)
  - Vision
  - LLM
  - Document parsing
  - Etc…
- Ready to deploy models
  - Azure model catalog
  - [GCP model garden](#)
  - AWS Sagemaker Marketplace
- MLOps bricks
  - Data catalog
  - Pre-built environment registries (VM, training, inference)
  - Training
  - Inference RT / batch
  - Monitoring
  - Model registry