

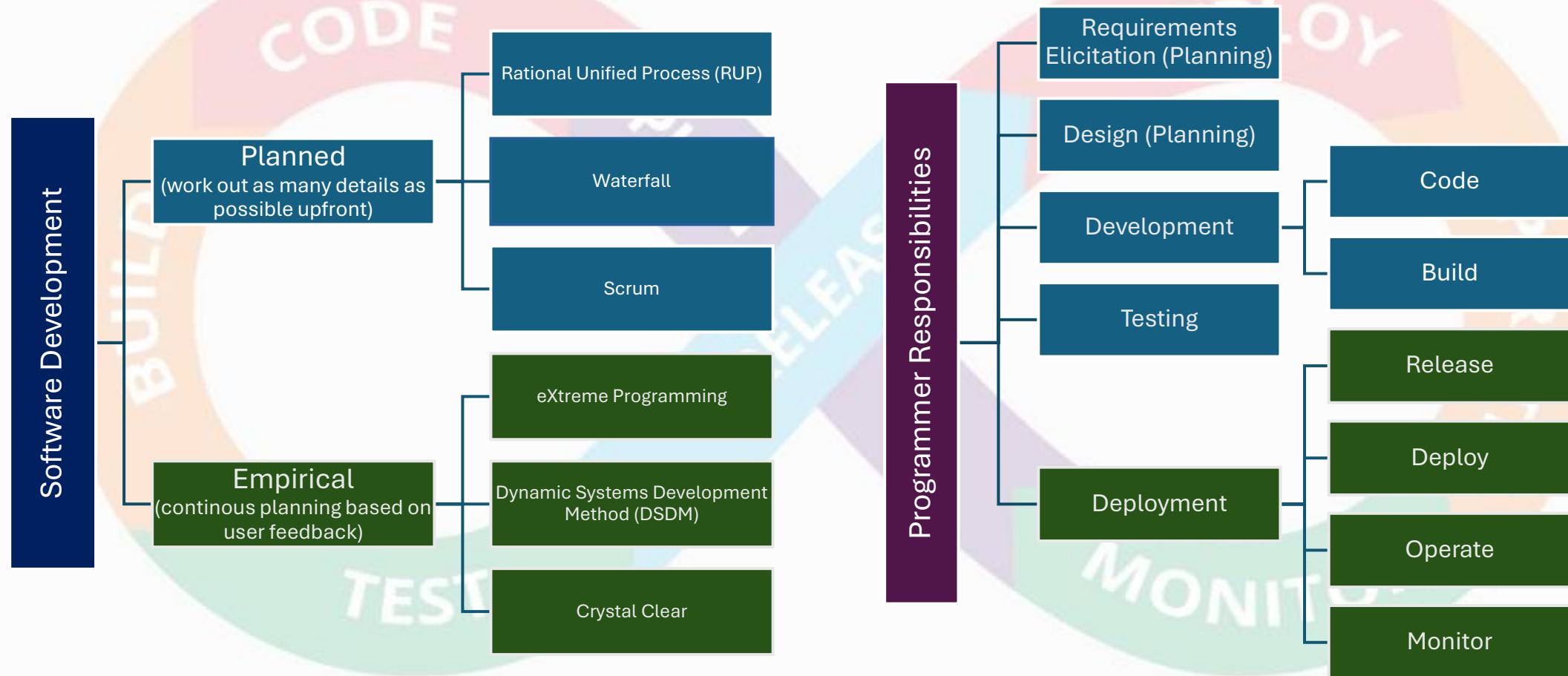
EduML

MLOps on Kubernetes

Andrei Olar

<andrei.olar@ubbcluj.ro>

DevOps¹



¹Kim, G., Humble, J., Debois, P., Willis, J., & Forsgren, N. (2021). *The DevOps handbook: How to create world-class agility, reliability, & security in technology organizations*. It Revolution.
Image courtesy of <https://mia-platform.eu/blog/devops/>

Machine Learning



Span¹

Narrow AI

General Purpose AI

AGI

ASI



Process^{2,3}

Level 0 – Manual (solitary development)

Level 1a – Automated Training (team collaboration)

Level 1b – Automated Serving (learning from evaluation)

Level 2 – Model monitoring (automated feedback loop)

¹Bostrom, N. (2014). *Superintelligence: Paths, dangers, strategies*. Oxford University Press.

²<https://mlops-for-all.github.io/en/docs/introduction/levels/>

³Sculley, D et. al, *Hidden Technical Debt in Machine Learning Systems*, Advances in Neural Information Processing Systems, 2015.

Image credits: OpenAI + Google

ML Ops



Goals

DevOps extension for ML professionals¹

Generic process for releasing ML artifacts²

Automate testing of ML artifacts²

Use established software development best practices with ML artifact²

¹<https://github.com/cdfoundation/sig-mlops/blob/main/roadmap/2020/MLOpsRoadmap2020.md#what-is-mlops>

²<https://ml-ops.org/content/mlops-principles>

ML Ops

Traits

Framework, language and method
agnostic

Reproducible

Automated

ML Ops

Stages

Data Preparation

Training / Experimentation

Deployment / Serving

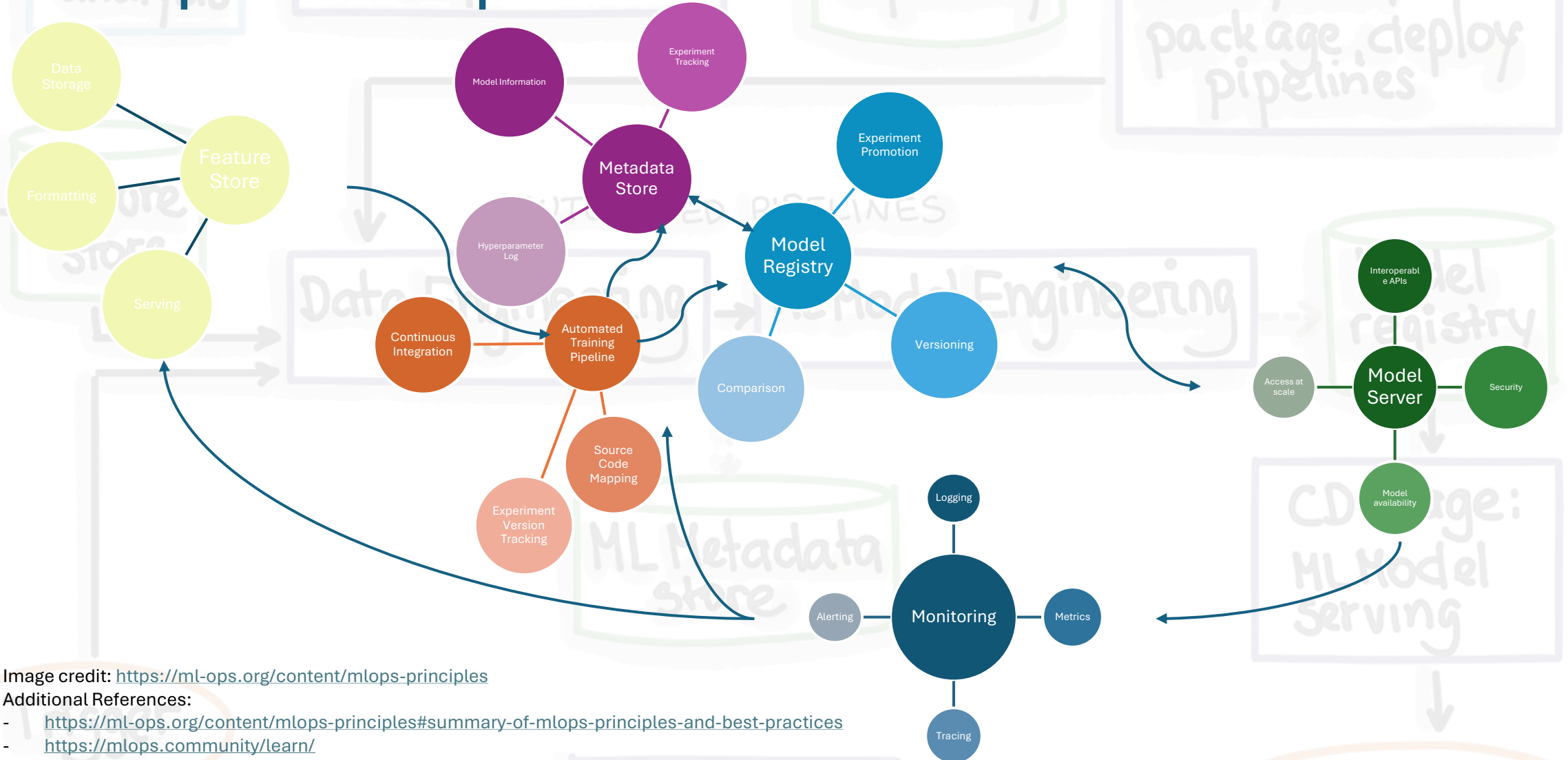
Monitoring

Image: <https://ubuntu.com/blog/what-is-mlops>

References:

- <https://blog.ml.cmu.edu/2020/08/31/3-baselines/>
- <https://cloud.google.com/blog/products/ai-machine-learning/key-requirements-for-an-mlops-foundation>

ML Ops - Components



ML Ops

Tooling

[MLFlow](#)



[Polyaxon](#)

[Kubeflow](#)

[MLRun](#)

Demo

The screenshot shows a web browser displaying the GitHub repository page for 'cs-ubb-eduml-app' by user 'koosie0507'. The repository is public and has 0 stars, 0 forks, and 3 watchers. The main content area shows a list of files and their commit history. The files listed are:

- `.github/workflows`: better image tagging (3 weeks ago)
- `copier-template`: update templated app readme (2 days ago)
- `data`: dog fooding: make this repo a sample for the template (last week)
- `helm`: dog fooding: make this repo a sample for the template (last week)
- `src/cs_ubb_eduml_app`: fix tuple order in `_process_wrapped_func_result` (last week)
- `.copier-answers.yml`: fix tuple order in `_process_wrapped_func_result` (last week)
- `.dockerignore`: create basic copier template (last week)
- `.gitignore`: ignore copier answers and mac finder markers (last week)
- `Dockerfile`: seems I have disabled job syncing entirely which is not w... (3 weeks ago)
- `LICENSE`: dog fooding: make this repo a sample for the template (last week)
- `MLProject`: dog fooding: make this repo a sample for the template (last week)
- `README.md`: use shell env var syntax to refer to config values in readme (2 days ago)
- `copier.yml`: dog fooding: make this repo a sample for the template (last week)
- `pyproject.toml`: dog fooding: make this repo a sample for the template (last week)

The repository also has a README file and is licensed under EUPL-1.2. The right sidebar shows the repository's activity, including a list of releases (2 tags) and packages (2 packages: `cs-ubb-mi-ops-test-1` and `cs-ubb-eduml-app`). The languages section shows the following distribution:

Language	Percentage
Python	46.0%
Smarty	19.4%
Jinja	33.5%
Dockerfile	1.1%

Specs

Kubernetes v1.31.1+k3s1

Cluster Nodes

- 3 control plane nodes: 8GB RAM, 4 core Intel [Xeon@2.2GHz](#)
- 2 GPU worker nodes: 32GB RAM, 16 core Intel [Xeon@2.3GHz](#)

Available GPUs

- One A100D-8C, CUDA Version: 12.2, 8GB RAM / GPU worker node
- Can't share or use NVLink (can't mine bitcoin either)

Recap

ML is complex both in terms of span and in terms of process

DevOps provides a solid framework to build upon automated ML processes

JupyterNotebooks – good Level 0 solution

MLFlow – good Level 1a solution

We have a nice, serial job running cluster



@Andrei Olar in [MIRPR-2024-2025](#)



@lauradiosan/MIRPR-UBB/2024-2025

@koosie0507/cs-ubb-eduml-app



andrei.olar [at] ubbcluj.ro