# Kubernetes and MLOps for Scalable and Reproducible Generative AI Workflows

What will this
session be about?

# Who am I?

@AnnieTalvasto

CMO at VSHN

o CNCF Ambassador

o Azure MVP

o Kubernetes & CNCF meetup co-organizer

o Startup-coach

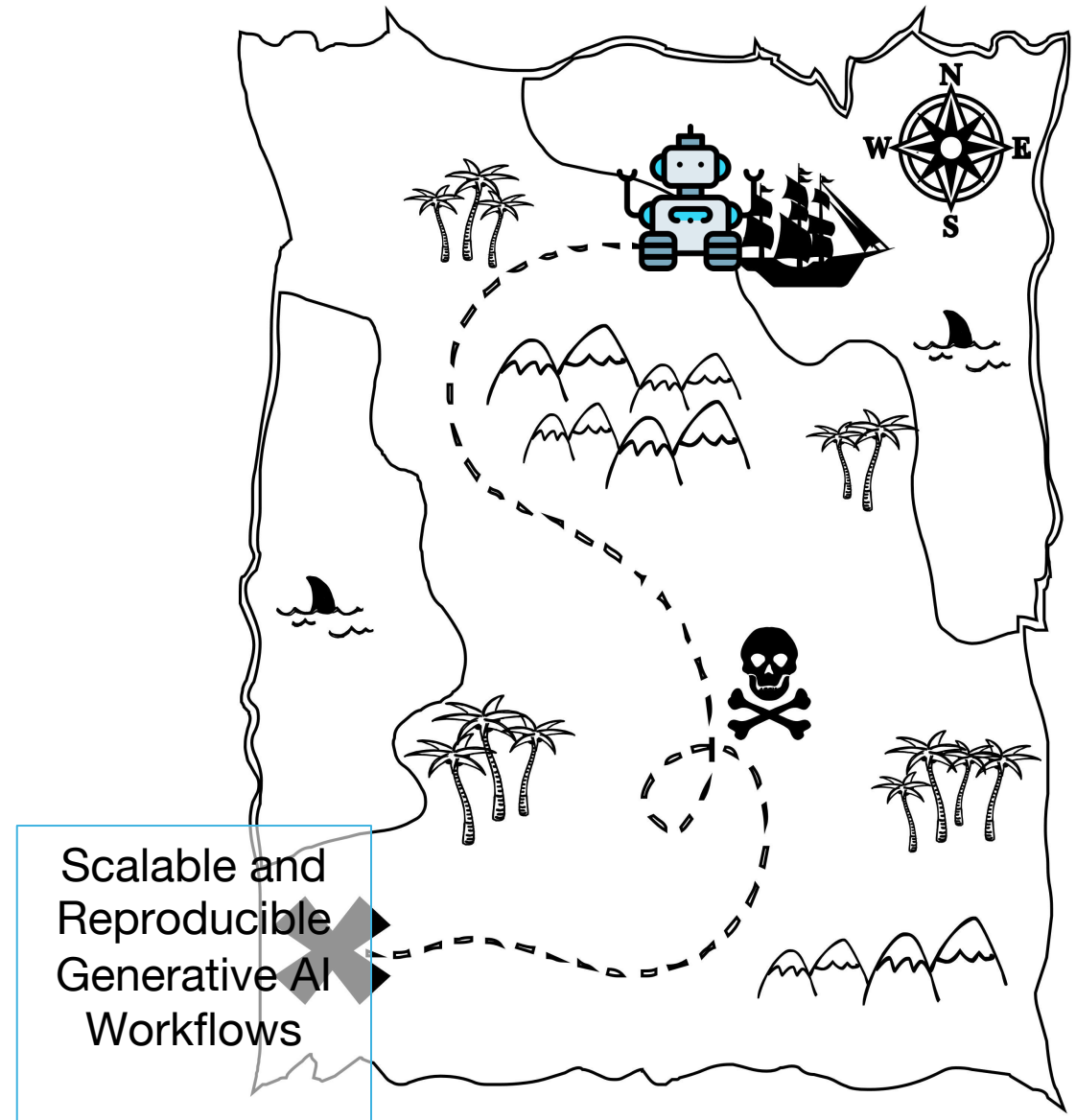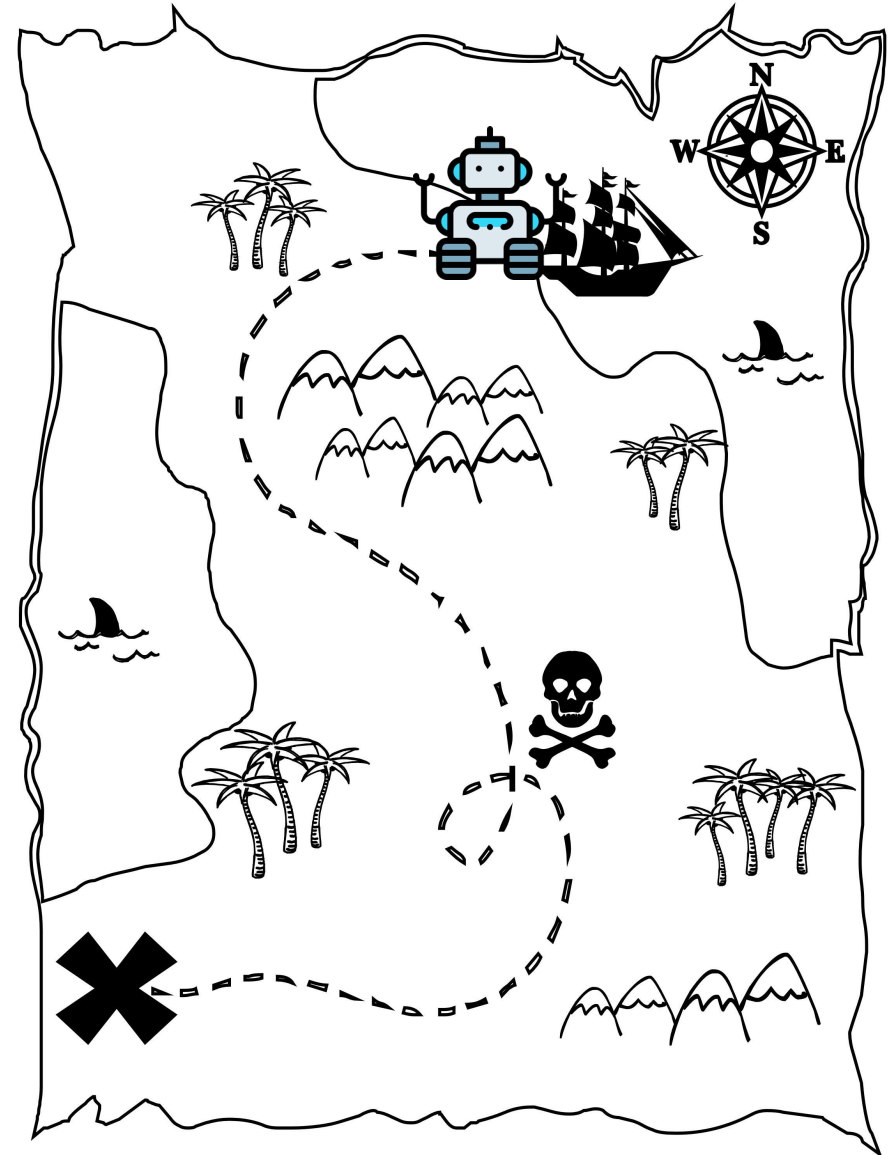o Created TechCraft Show – Tech interview show with Minecraft ⛏️💎🪨

# Agenda

o Introduction

o Definitions

o Why is AI/ML different

o Kubernetes & MLOps

o Kubeflow

o Best practices

o Considerations

o Wrap up & Resources

# Agenda

o Introduction

o Definitions

o Why is AI/ML different

o Kubernetes & MLOps

o Kubeflow

o Best practices

o Considerations

o Wrap up & Resources

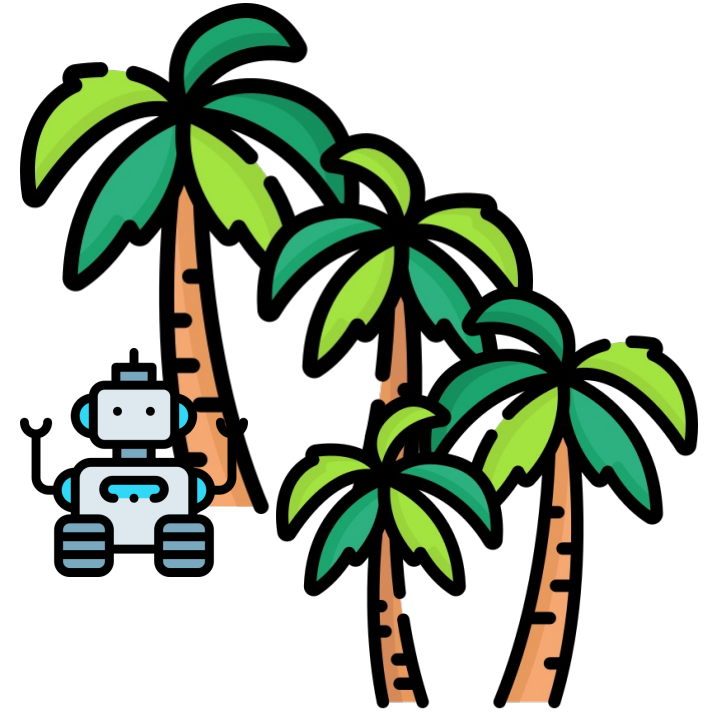Scalable and Reproducible Generative AI Workflows

# Agenda

o ~~Introduction~~

o Definitions

o Why is AI/ML different

o Kubernetes & MLOps

o Kubeflow

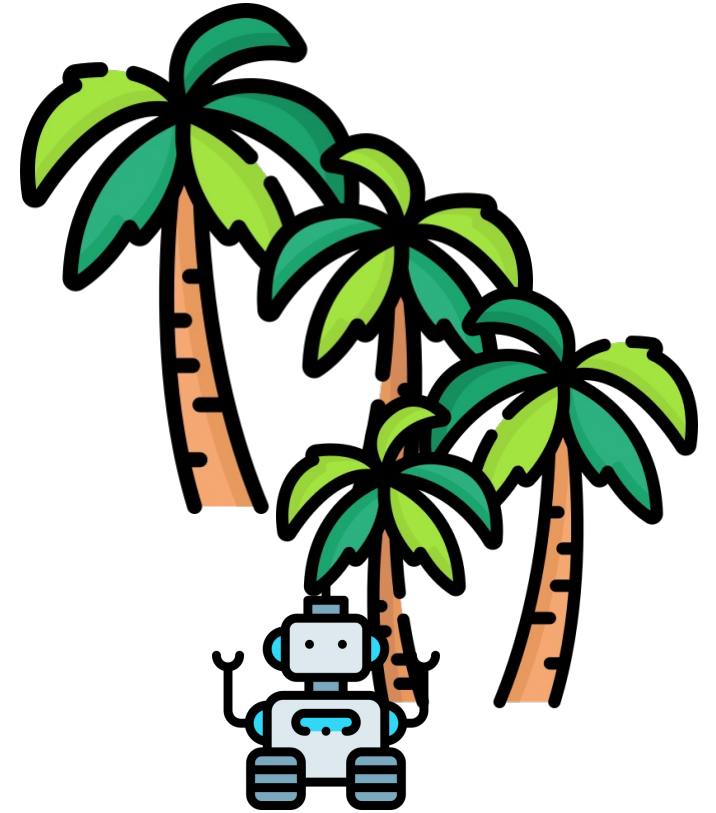o Best practices

o Considerations

o Wrap up & Resources

# AI

o Artificial Intelligence (AI) refers to the simulation of human intelligence in machines that are programmed to think and learn like humans.
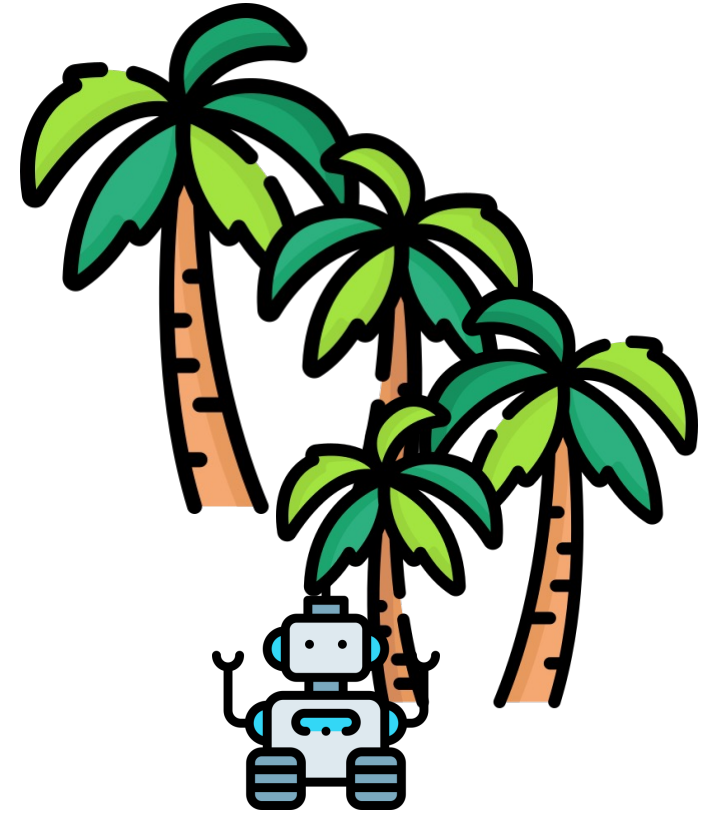
# ML

o Machine Learning (ML) is a subfield of artificial intelligence (AI) that focuses on the development of algorithms and statistical models that enable computer systems to improve their performance on a specific task through learning from data, without being explicitly programmed.
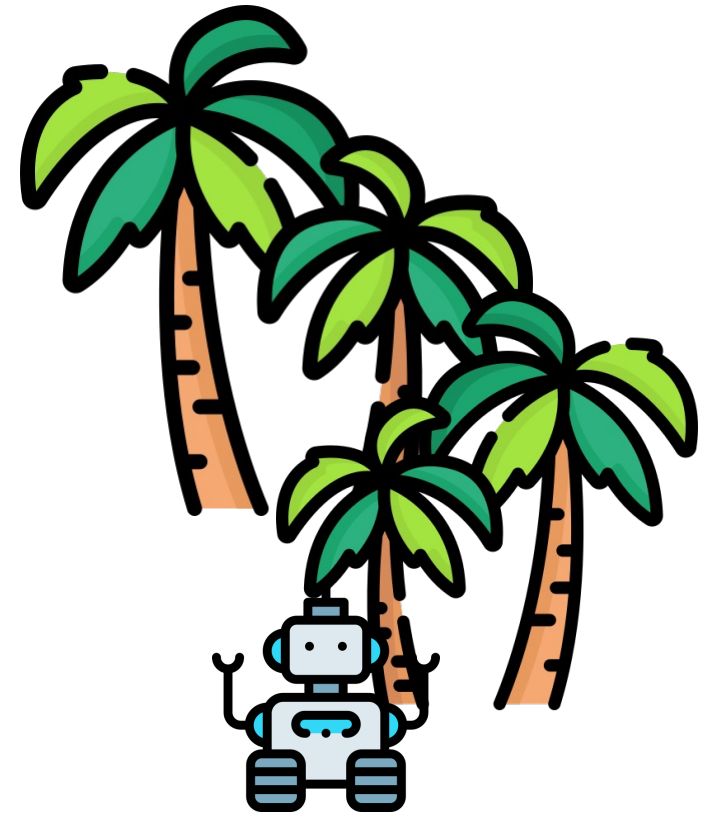
# Data science

o Data science is a multidisciplinary field that uses various techniques, algorithms, processes, and systems to extract insights and knowledge from data.
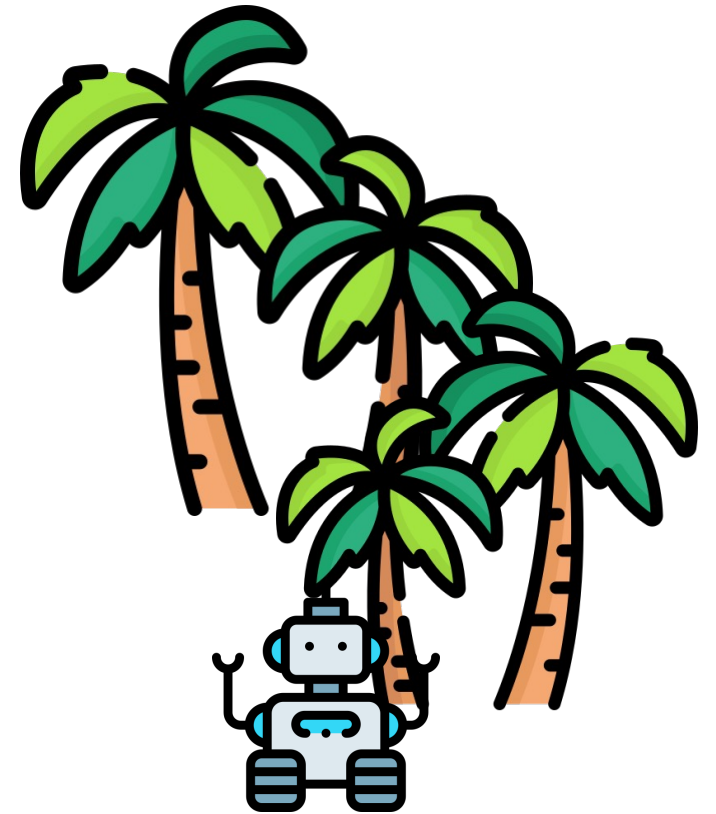
# DevOps

o DevOps, short for "Development" and "Operations," is a set of practices, principles, and cultural philosophies that aim to improve and streamline collaboration between software development (Dev) and IT operations (Ops) teams.
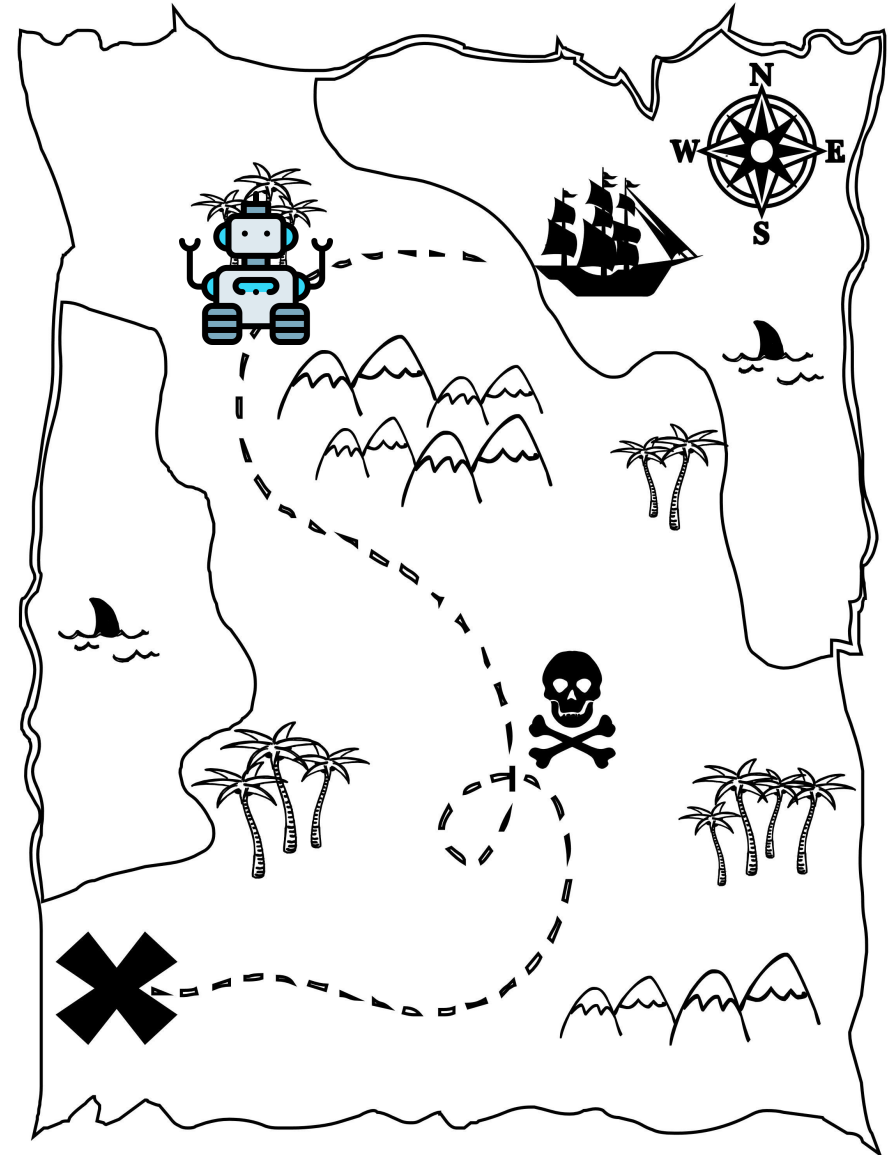
# MLOps

o MLOps, short for "Machine Learning Operations," is a set of practices, principles, and tools that combine machine learning (ML) with the practices of DevOps.
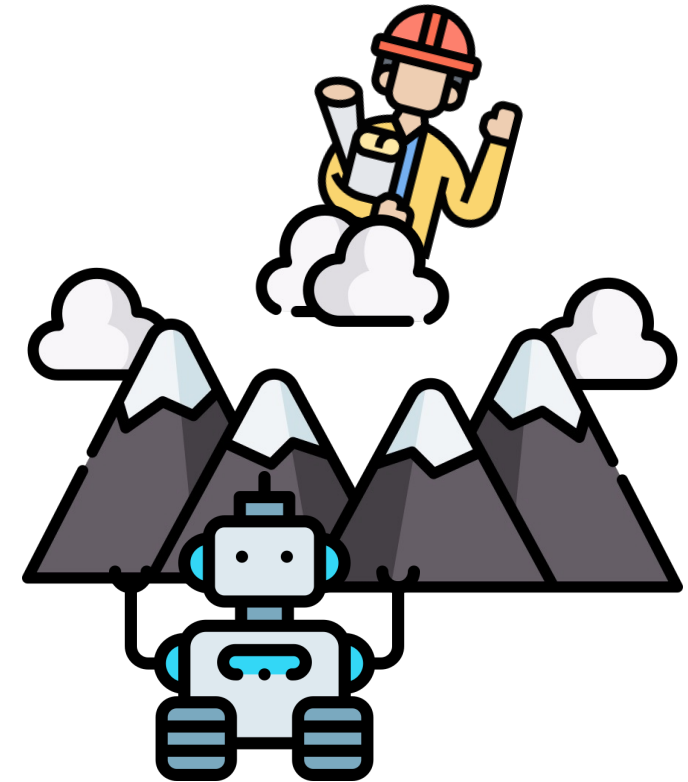
# Agenda

o ~~Introduction~~

o ~~Definitions~~

o Why is AI/ML different

o Kubernetes & MLOps

o Kubeflow

o Best practices

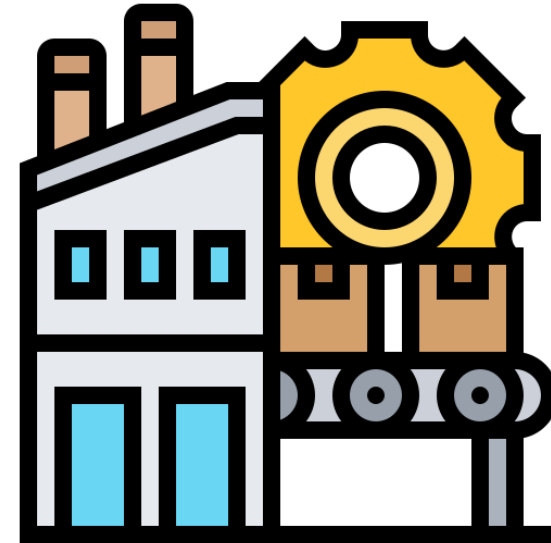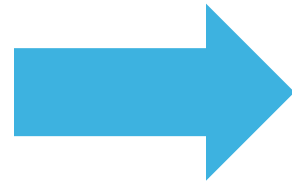o Considerations
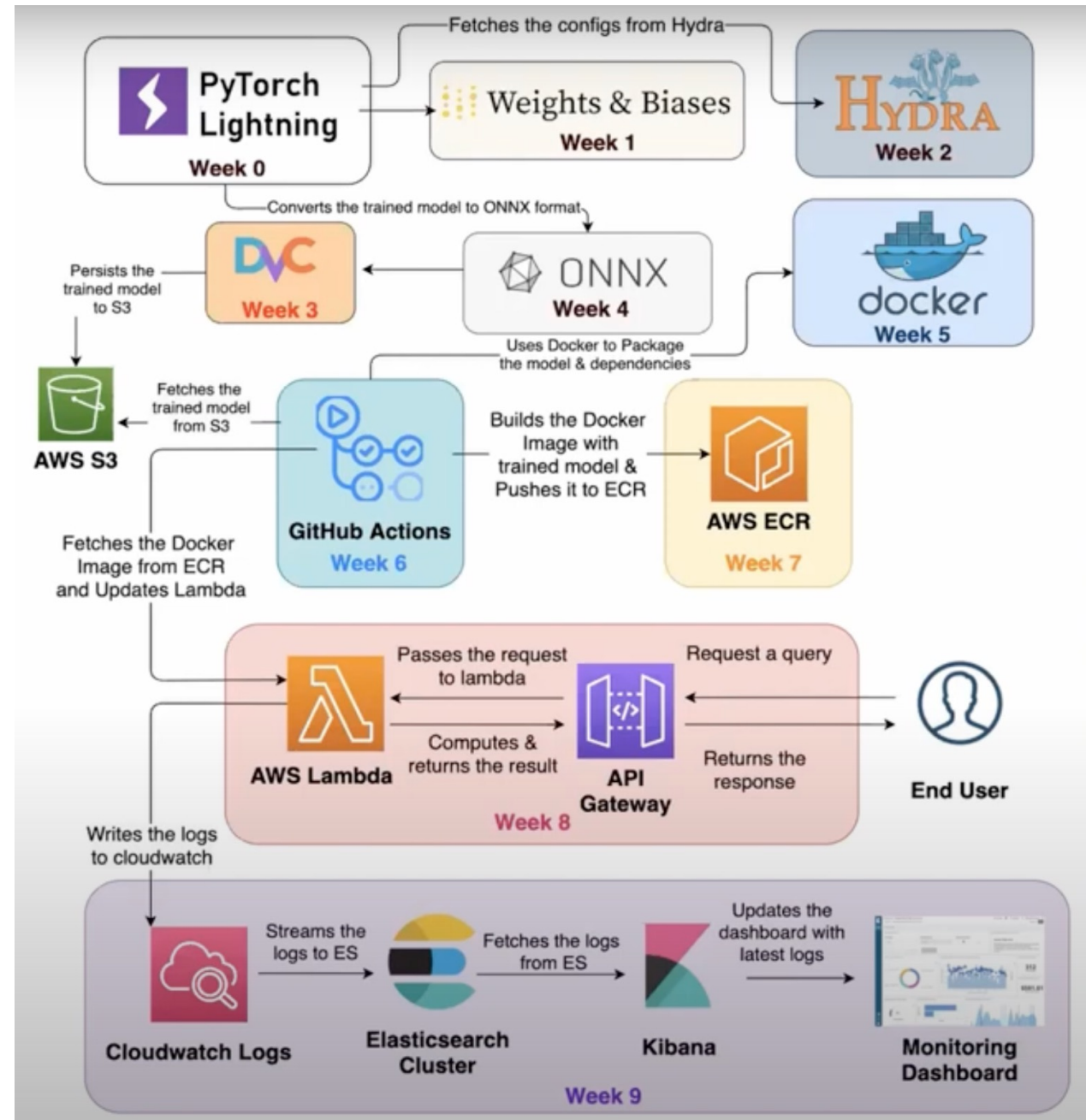
o Wrap up & Resources

# Differences with AI/ML

o Data-Centric Nature

o Model Complexity

o Iterative Development

o Scalability

o Monitoring and adaptation

o Testing and validation

o Explainabity and bias

o Rapid advancements

o Collaboration

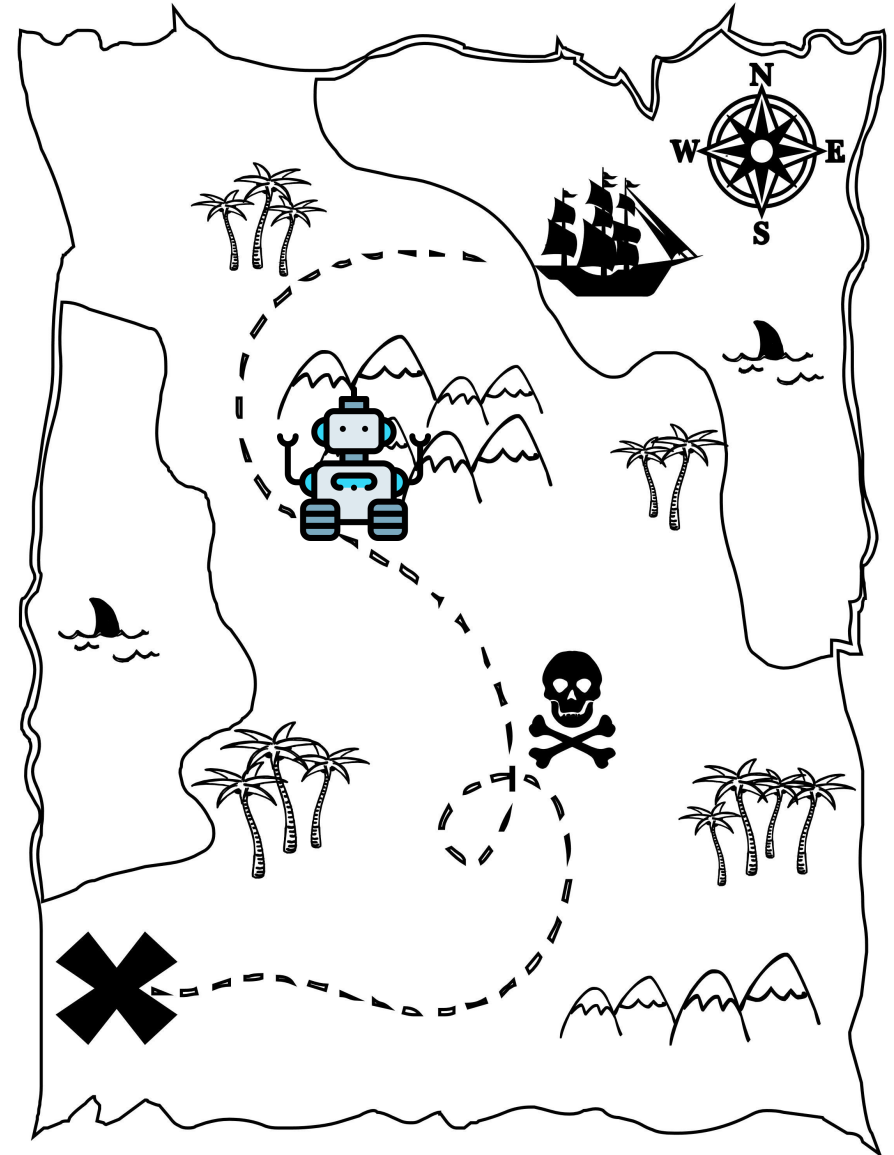o Deployment challenges

o Cost

# From Research to Production

Example basic MLOps flow by Raviraja Ganta

@AnnieTalvasto

# Agenda

# Kubernetes & MLOps

o Portability

o Customizability

o Performance

o Consistency

o Microservices
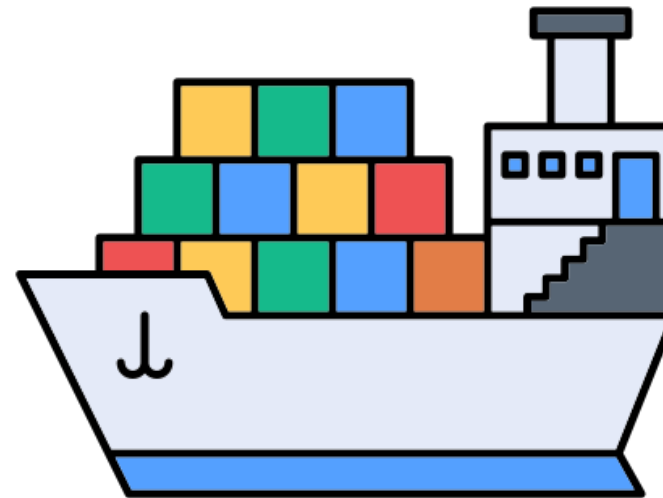
o Composability

# Agenda

- ~~Introduction~~
- ~~Definitions~~
- ~~Why is AI/ML different~~
- ~~Kubernetes & MLOps~~
- Kubeflow
- Best practices
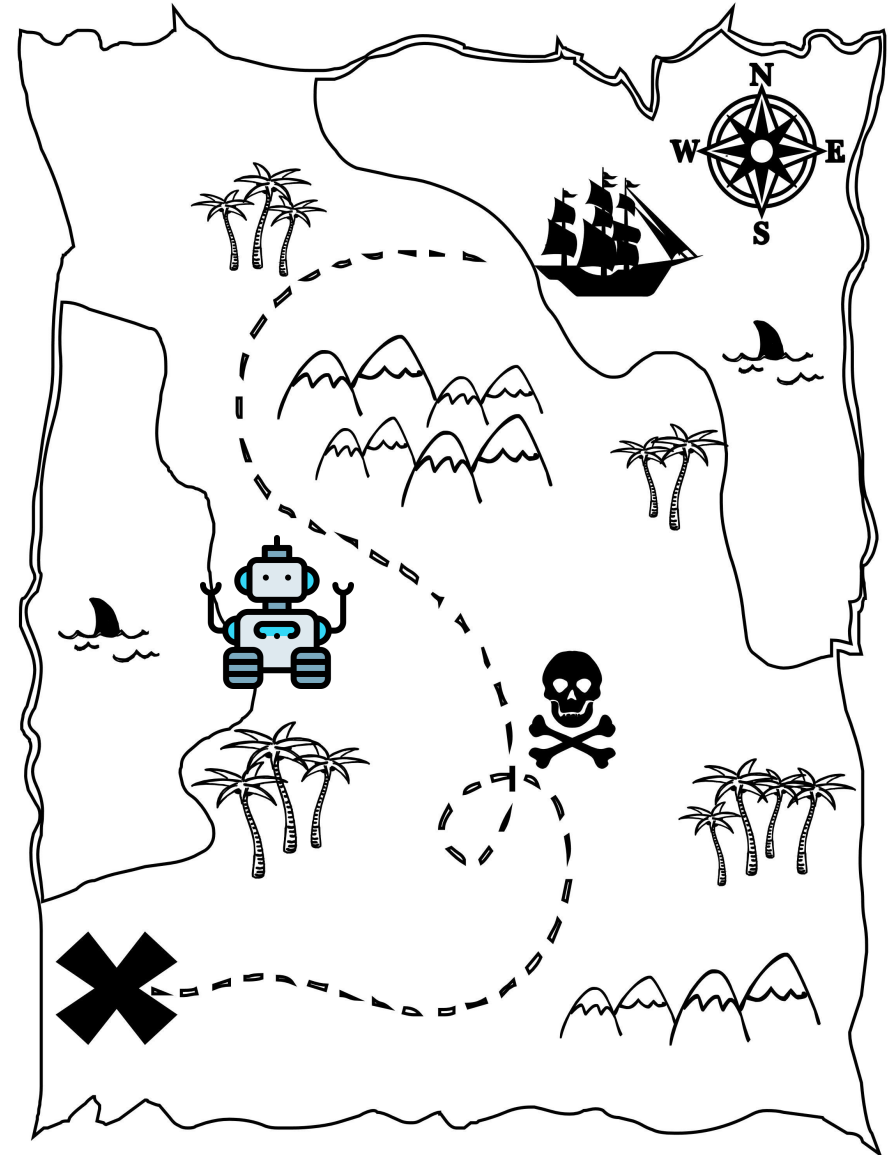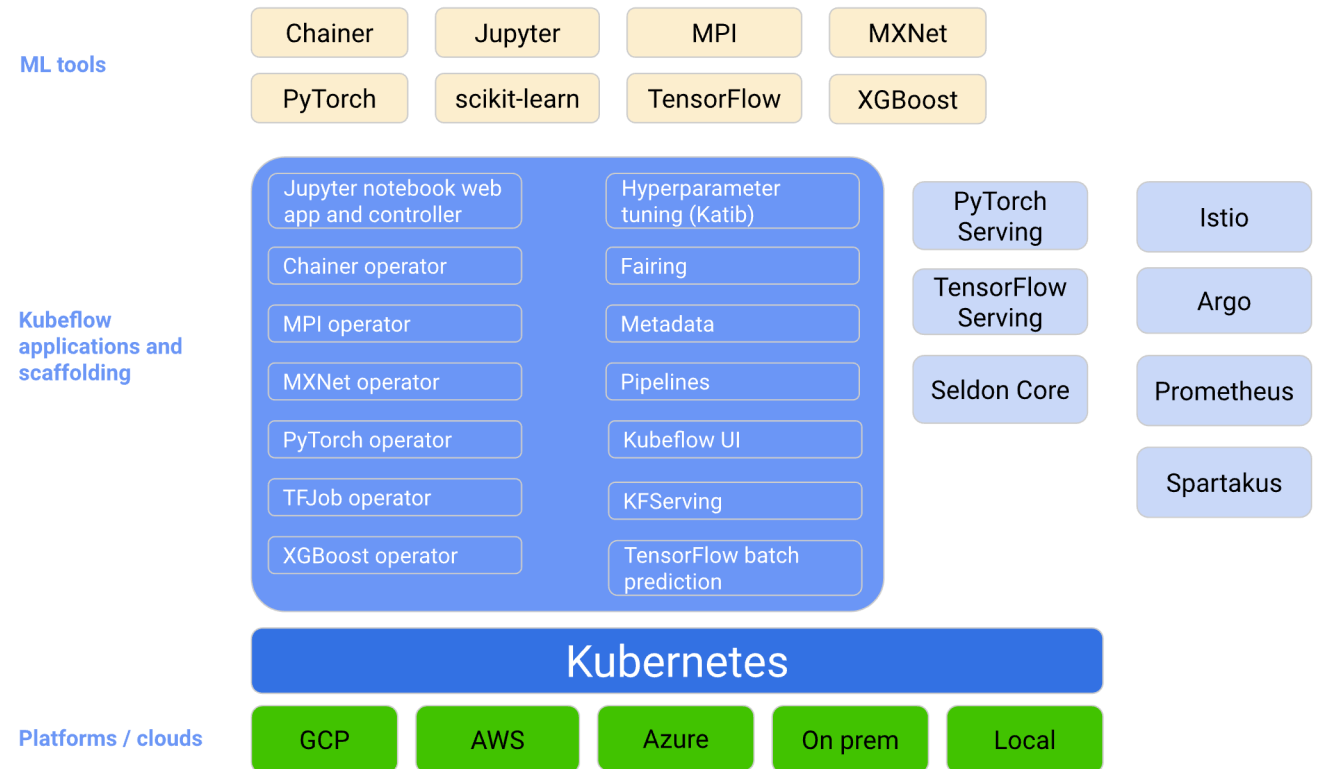- Considerations
- Wrap up & Resources
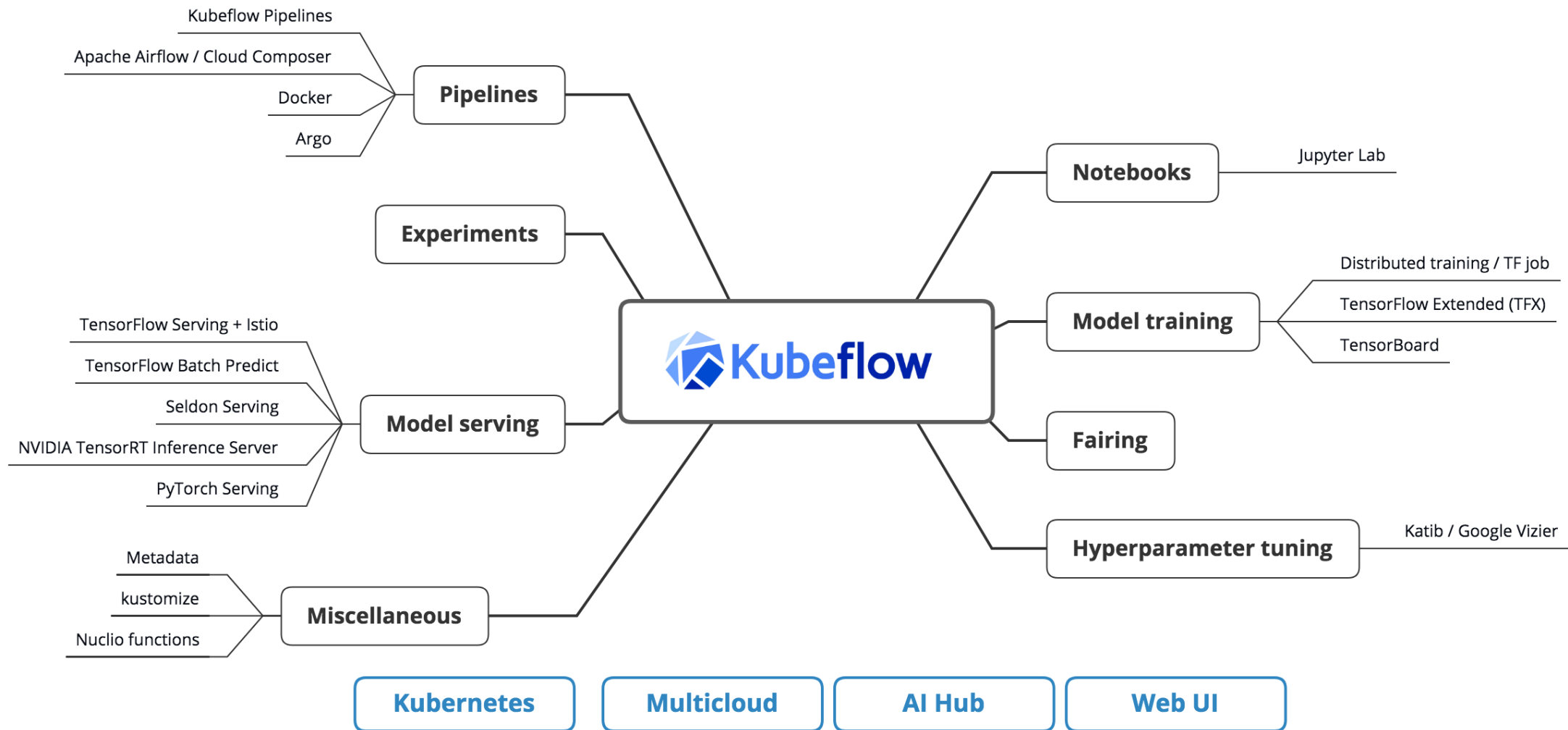
# Kubeflow

The machine learning
toolkit for Kubernetes.

# What is Kubeflow?

o End-to-End ML workflow

o Scalability and resource management

o Reproducibility and collaboration



**ML tools**

| Chainer | Jupyter | MPI | MXNet |
| PyTorch | scikit-learn | TensorFlow | XGBoost |

**Kubeflow applications and scaffolding**

| Jupyter notebook web app and controller | Hyperparameter tuning (Katib) |
| Chainer operator | Fairing |
| MPI operator | Metadata |
| MXNet operator | Pipelines |
| PyTorch operator | Kubeflow UI |
| TFJob operator | KFServing |
| XGBoost operator | TensorFlow batch prediction |

PyTorch Serving

TensorFlow Serving

Seldon Core

Istio

Argo

Prometheus

Spartakus

**Kubernetes**

**Platforms / clouds**

| GCP | AWS | Azure | On prem | Local |

Source: https://www.kubeflow.org/docs/started/architecture/
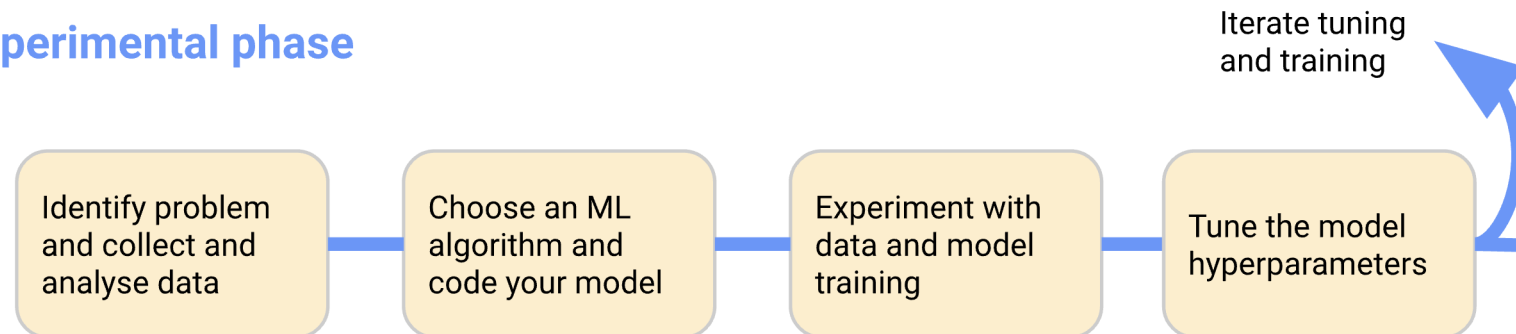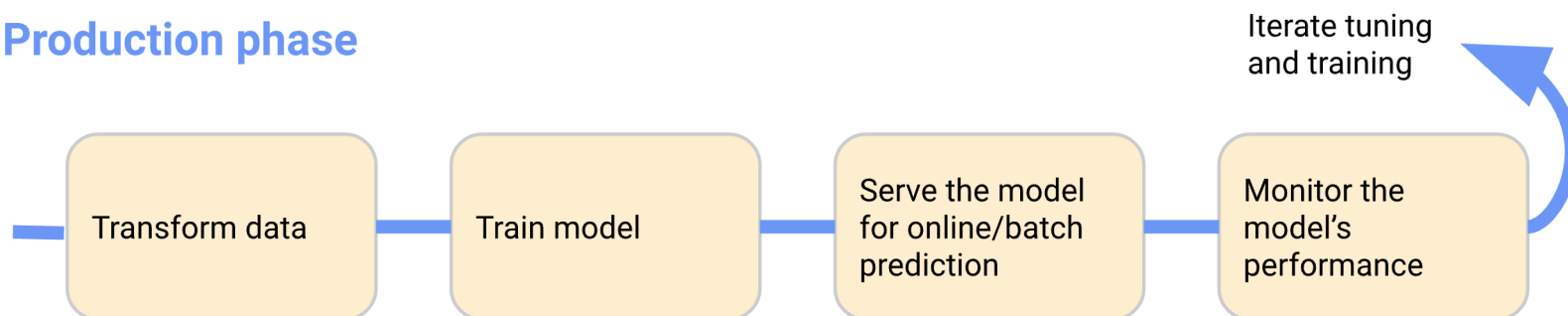
**@AnnieTalvasto**

Source: https://www.analyticsvidhya.com/blog/2023/01/kubeflow-streamlining-mlops-with-efficient-ml-workflow-management//

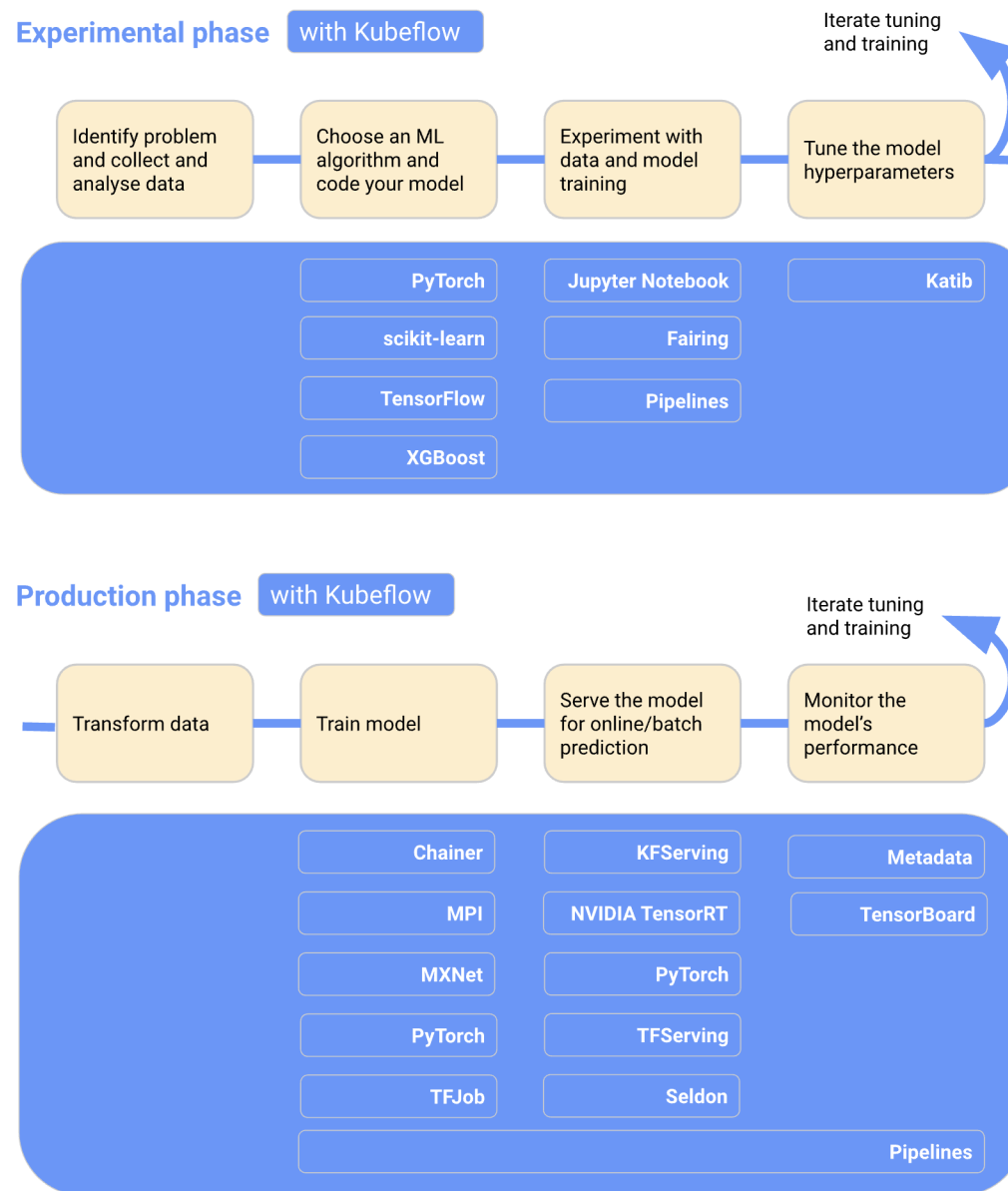@AnnieTalvasto

# Before Kubeflow

**Experimental phase**

| Identify problem and collect and analyse data | — | Choose an ML algorithm and code your model | — | Experiment with data and model training | — | Tune the model hyperparameters |

Iterate tuning and training

**Production phase**

| Transform data | — | Train model | — | Serve the model for online/batch prediction | — | Monitor the model's performance |

Iterate tuning and training

Source: https://www.kubeflow.org/docs/started/architecture/

# With Kubeflow



**Experimental phase** with Kubeflow

Identify problem and collect and analyse data → Choose an ML algorithm and code your model → Experiment with data and model training → Tune the model hyperparameters → Iterate tuning and training

- PyTorch
- scikit-learn
- TensorFlow
- XGBoost
- Jupyter Notebook
- Fairing
- Pipelines
- Katib

**Production phase** with Kubeflow

Transform data → Train model → Serve the model for online/batch prediction → Monitor the model's performance → Iterate tuning and training

- Chainer
- MPI
- MXNet
- PyTorch
- TFJob
- KFServing
- NVIDIA TensorRT
- PyTorch
- TFServing
- Seldon
- Metadata
- TensorBoard
- Pipelines

Source: https://www.kubeflow.org/docs/started/architecture/

# Example: Shell



Keynote: Machine Learning on Kubernetes Made Easy With Kubeflow - Masoud Mirmomeni & Jimmy Guerrero

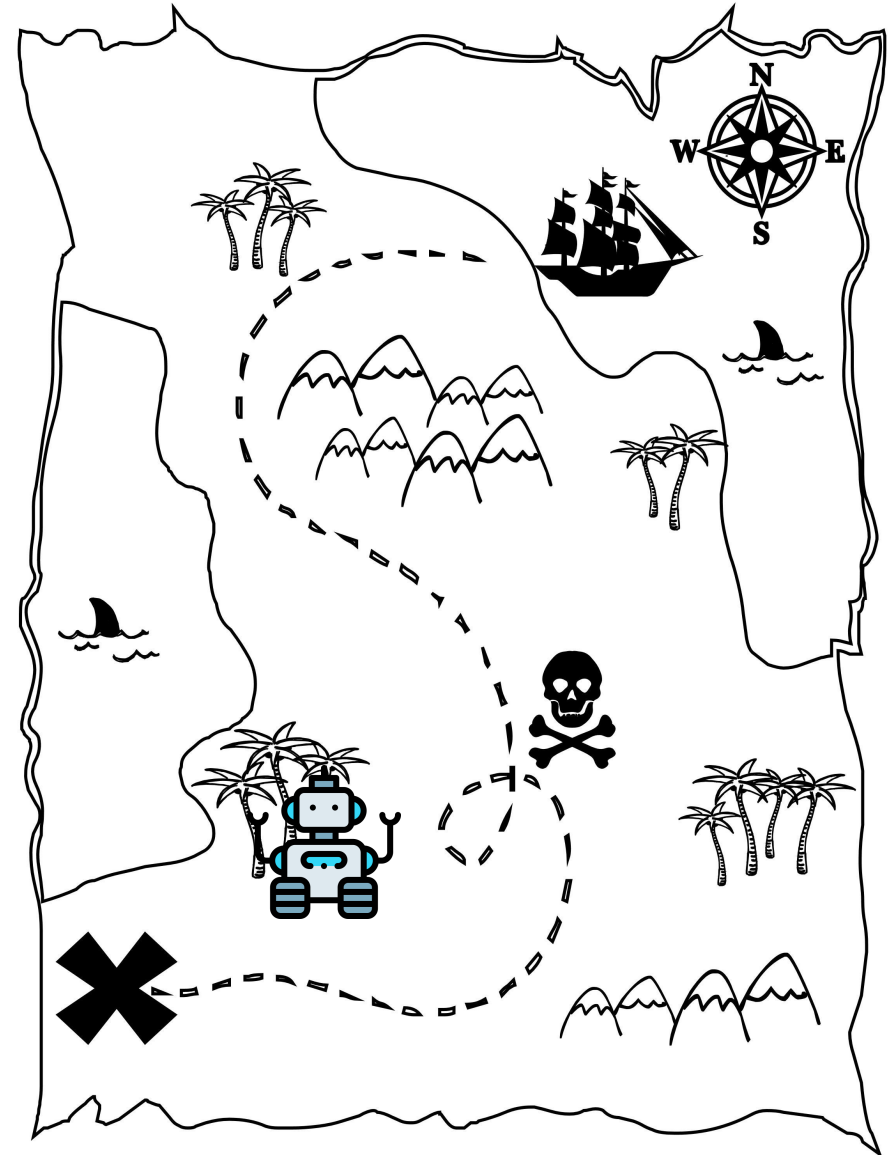CNCF [Cloud Native Computing Foundation]
105K subscribers

@AnnieTalvasto

# Agenda

o ~~Introduction~~
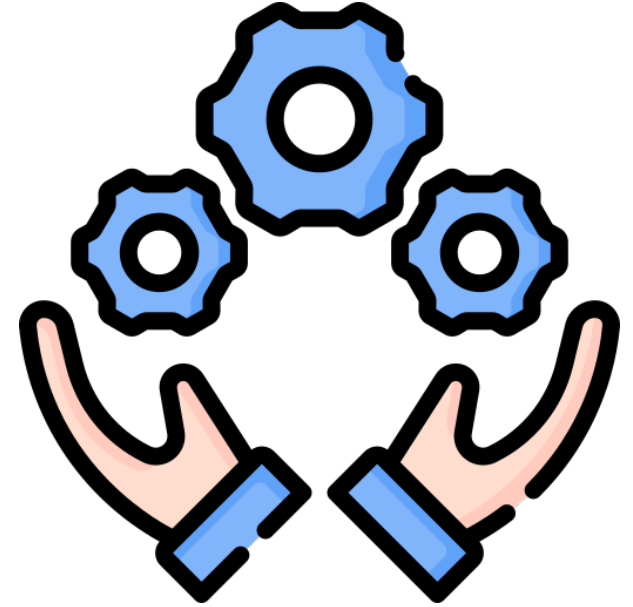
o ~~Definitions~~

o ~~Why is AI/ML different~~

o ~~Kubernetes & MLOps~~

o ~~Kubeflow~~

o Best practices

o Considerations

o Wrap up & Resources

# Best Practices

o It is always start with data

o Master all the necessary skills

o Know the tools you need and choose them wisely

o Understand the bottlenecks you face
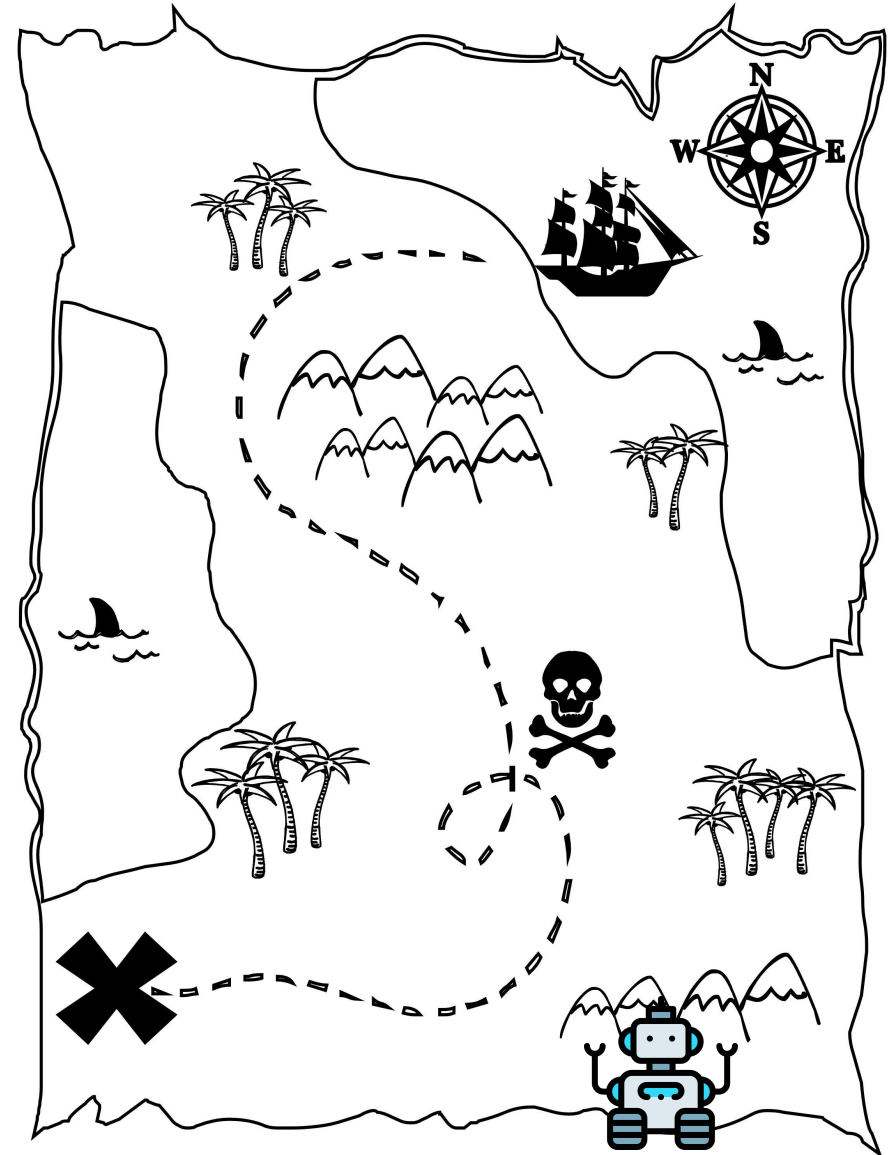
o MLOps ❤️ DevOps

o AutoML, Kaizen

# DevOps

o Automation

o CI/CD

o Infrastructure as Code

o Monitoring

o Culture, collaboration and comms

o Version Control

o Security

# Agenda

- ~~Introduction~~
- ~~Definitions~~
- ~~Why is AI/ML different~~
- ~~Kubernetes & MLOps~~
- ~~Kubeflow~~
- ~~Best practices~~
- Considerations
- Wrap up & Resources

# Ethical considerations

o Fairness & Bias
o Transparency
o Privacy
o Accountability
o Consent and user empowerment
o Security
o Human-Centric Design
o Data Governance
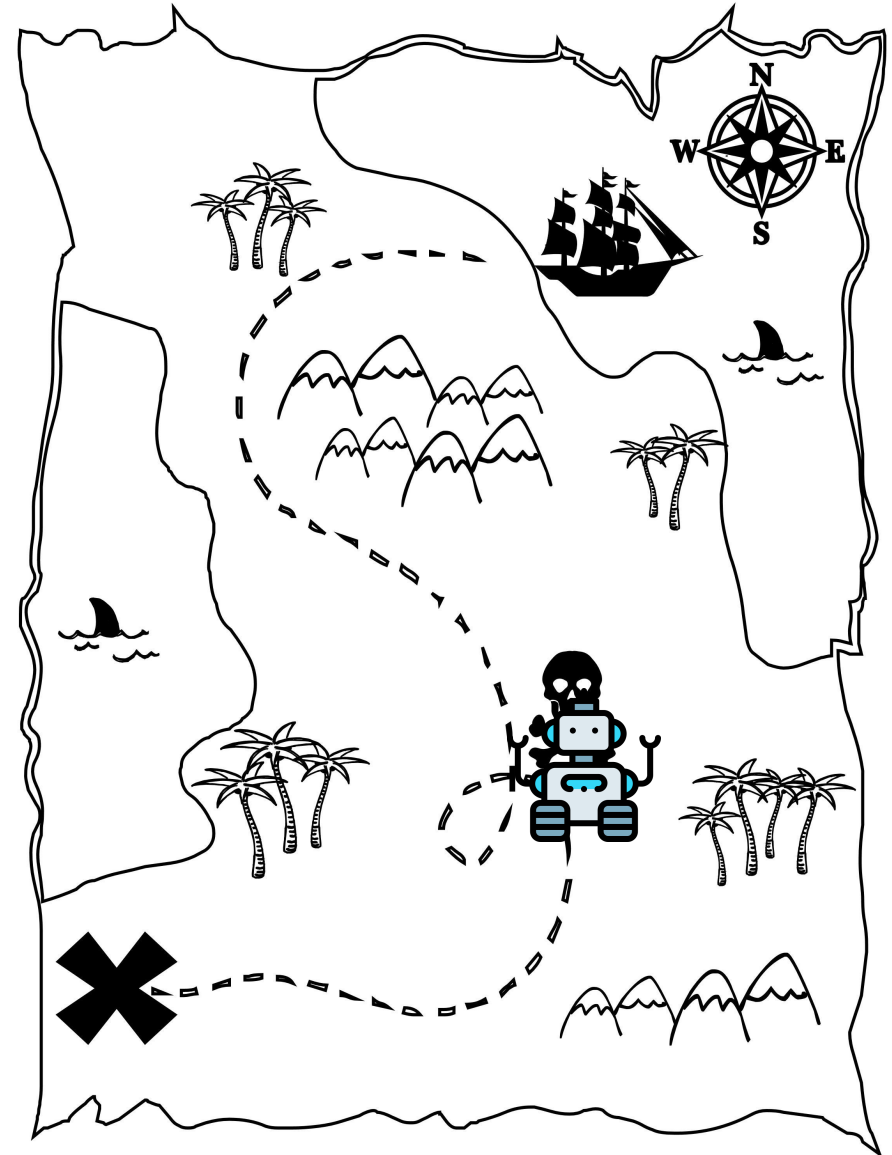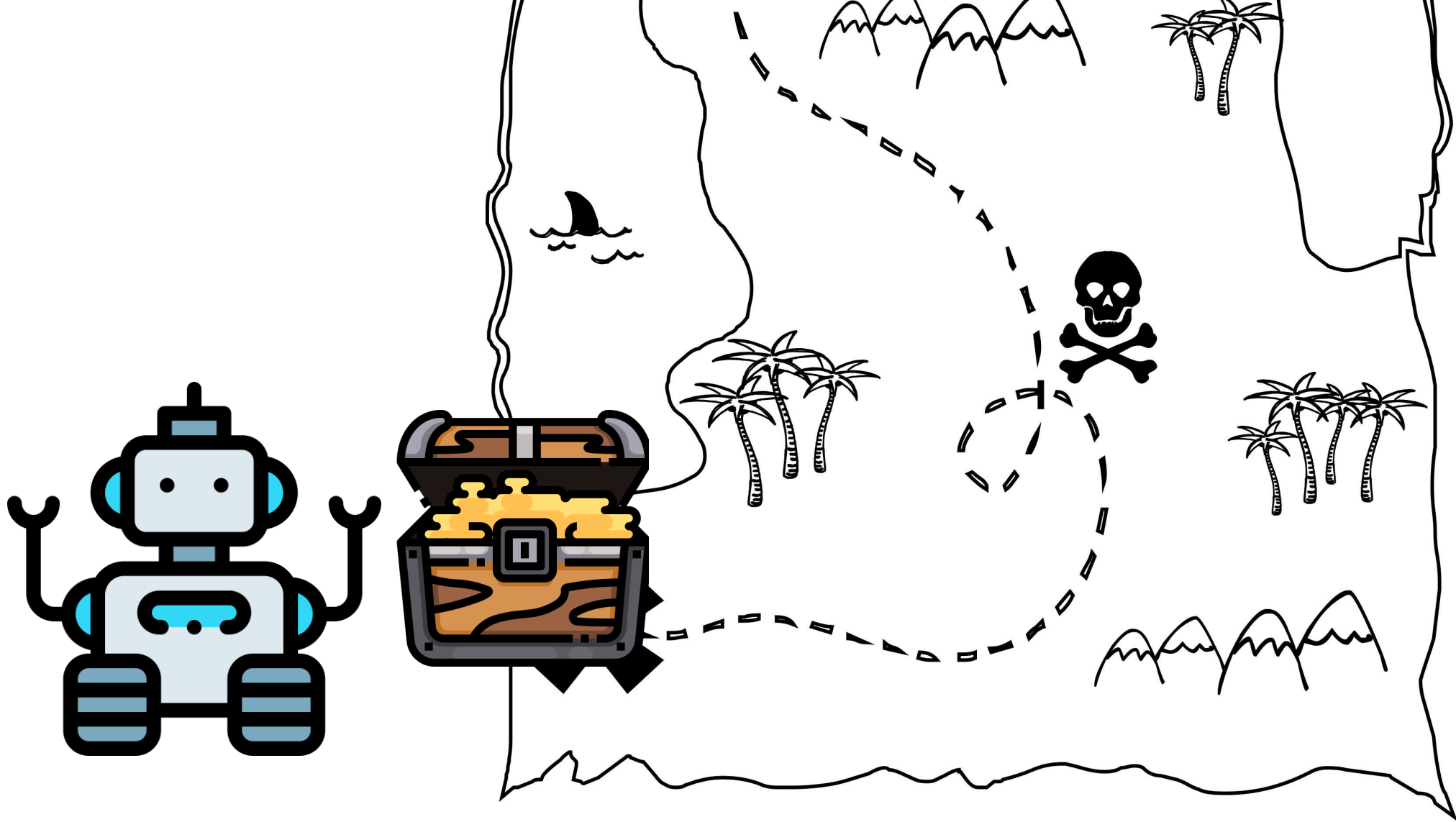o Long-term consequences
o Regulatory Compliance

# Learn more

o <mark>Links and slides: github.com/annietalvasto</mark>

o Practical MLOps by Noah Gift and Alfredo Deza (O'Reilly). Copyright 2021 Noah Gift and Alfredo Deza, 978-1-098-10301-9.

o Keynote: Machine Learning on Kubernetes Made Easy With Kubeflow - Masoud Mirmomeni & Jimmy Guerrero
  o https://www.youtube.com/watch?v=ick5hI5YI0k

o End to End MLOps Basics // Raviraja Ganta // MLOps Meetup #82
  o https://youtu.be/B1t_Vb2MkRw?si=90TzpNKa-veq-ifK

o Webinar: MLOps automation with Git Based CI/CD for ML
  o https://www.youtube.com/watch?v=VCUDo9umKEQ

o All Kubernetes AI day Sessions

o Welcome + Opening Remarks: The State of Production MLOps in the Cloud Native… - Alejandro Saucedo
  o https://www.youtube.com/watch?v=xymbp8RWaCQ&list=PLj6h78yzYM2M9oVaU3amsqL5RXUwcu

# Agenda

- Introduction
- Definitions
- Why is AI/ML different
- Kubernetes & MLOps
- Kubeflow
- Best practices
- Considerations
- Wrap up & Resources

@AnnieTalvasto

Thank you!