# VLAIO TETRA Machine Learning Operations for Edge Condition Monitoring (MLOps4ECM)

# Part I – MLOps and Edge Intelligence in Practice: Industry Tech-Talks







# MLOps and Edge Intelligence in Practice

T

E

C

Н

T

Α

L

K

S

MLOps for video fire detection on surveillance cameras

Araani – Maarten Callens

Real-time processing on PLC's with AI

CTRL Engineering - Glenn De Loose

Managing an AI server: Kubeflow

Howest - Thomas Huyghebaert

Smarter data pipelines with Dagster

Leap Technologies – Wannes De Groote

How does Siemens enable industrial grade Al

Siemens - Niels Debusschere

Al Ops Capability Maturity Model

Superlinear – Pierre Gerardi

No-nonsense approach to MLOps

**Vintecc** – Gilles Ballegeer

# VLAIO TETRA Machine Learning Operations for Edge Condition Monitoring (MLOps4ECM)

Lunchbreak



# VLAIO TETRA Machine Learning Operations for Edge Condition Monitoring (MLOps4ECM)

Part II - Project results









- > 8 professors
- > 1 research manager
- 1 project office manager
- 6 post-docs
- > 40 junior researchers
- > 2 ATP test engineers

Mechanical Engineering

**Electrical Engineering** 

Computer Science

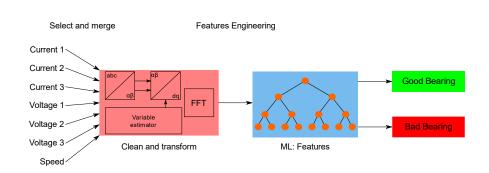






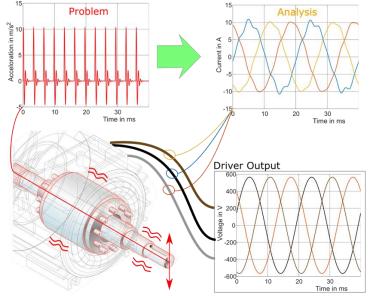
### The Ultimate Machine

- EMC test chambers
- Climate chambers
- Fuel cells
- Smart conditioning monitoring
- Edge computing







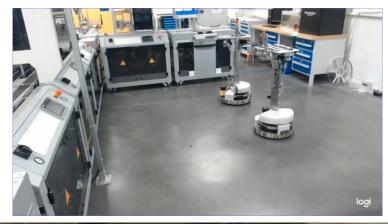






### The Ultimate Factory

- Data-driven process optimisation
- Autonomous systems
- Quality control











### **OPTIMISATION**

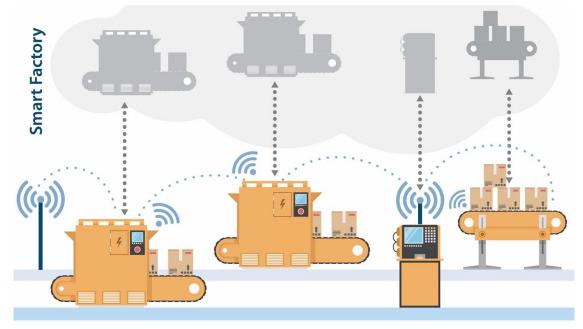
How can unwanted effects be mitigated?

Link between process parameters & conditions and product quality

### **MONITORING**

What is going wrong and what is the underlying reason?

- Abnormal machine behavior
- Product quality defects



Al for complex mechatronic systems and processes

### **CONTROL**

How can the parameter tuning be automated?

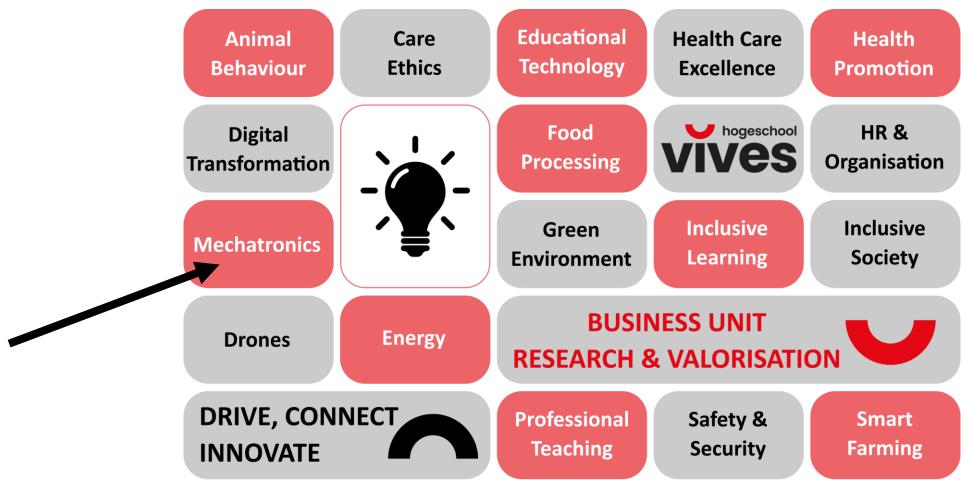
 Realise a stable system behavior

TRANSVERSAL ASPECTS: SAFETY, ROBUSTNESS, EFFICIENCY

### hogeschool VES



# Research at VIVES University of Applied Sciences







# Research Group Mechatronics – Team



Hanne Van den Berghe Research Manager

**Geert Furniere**Researcher

Teacher

**Emiel Wallays** 

PhD Researcher

**Noah Debaere** 

Researcher



Carl Grimmelprez

Jonas Lannoo Innovation Manager Senior Researcher

**Jeremy Lebon**Researcher
Teacher

Tom Van Gaever Researcher Teacher

**William Van Nieuwenhove** Researcher

Sille Van Landschoot

Researcher

Teacher

Alexander D'hoore

Researcher

Teacher

Philippe Van Overbeke

Researcher

Sam Lefebvre

Researcher

**Peter Vanbiervliet** 

Researcher

Teacher Philip Vanloofsvelt

Researcher

Teacher

Karolien Vandersickel

Researcher Teacher Catherine Middag

Researcher

**Tineke Verkest** 

Researcher

Researcher





## Research Group Mechatronics

Two complementary research tracks

### 1. Machinebuilding and product

- Sustainable product design
- Flexible and efficient production lines
- Al-driven autonomy

### 2. Internet of Things and data

- End-to-end prototyping
- Digitalisation for industry
- (I)IoT and Edge AI



### → Three labs in three locations







# Research Group Mechatronics

### Maaklab @ VIVES Campus Kortrijk







# Maaklab @ VIVES Campus Kortrijk

### Research & services

- Product design & materials
- Prototyping & proof-of-concept
- Custom parts & machine building
- Automation & 3D scanning

### Infrastructure

- Extensive machine park
  - CNC, Waterjet, Injection Molding
  - WAAM, EDM (spark erosion), 3D printing
- Scanning infrastructure
- Machine vision

















# Research Group Mechatronics

### IoT-lab @ VIVES Campus Brugge









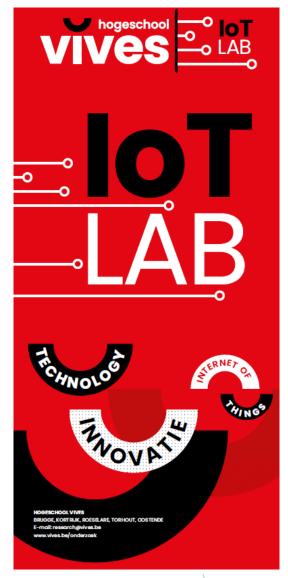
# IoT-lab @ VIVES Campus Brugge

### Development of IoT systems and architectures

- Custom hardware, firmware, software
- Connectivity & low-power
- Implementation of AI on IoT (edge/tinyML)

### Infrastructure

- Creating prototypes
- Test equipment
- Calculation & hosting servers









# New joint research lab

### Infinity lab @ Flanders Make Kortrijk 4th Floor

















# **Infinity lab**







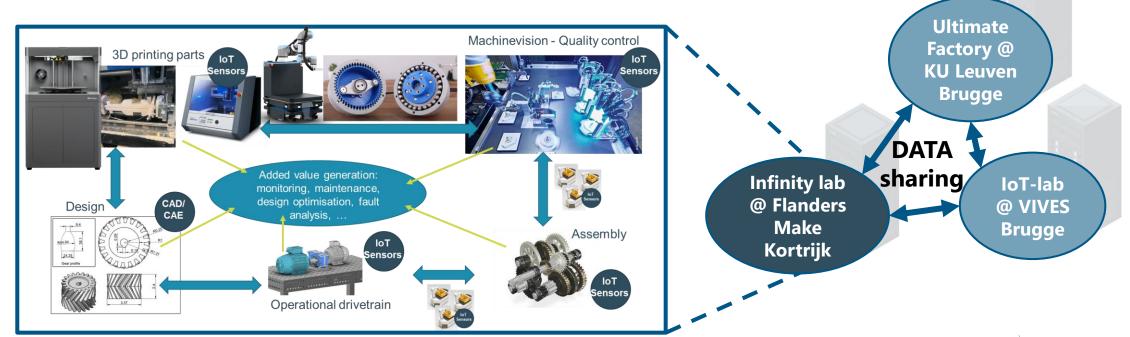


# **Infinity Lab**

• Infinity + Product Lifecycle Management + AI Robotics

> Straight-through digitalisation as motor for innovation





# VLAIO TETRA Machine Learning Operations for Edge Condition Monitoring (MLOps4ECM)

**Project Introduction** 







## MLOps4ECM - Motivation

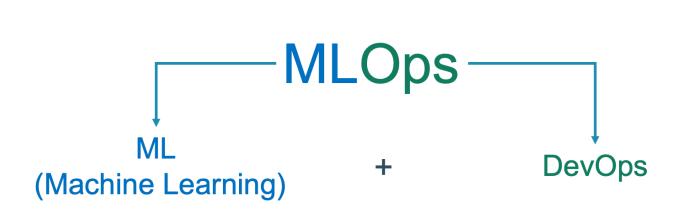
- Complex industrial processes require data-driven models to monitor and control operations
- Today, most AI developments remain limited to R&D or lab-scale test setups
- **Deployment in the field** demands continuous **monitoring** and feedback to ensure reliability
- Automation and model updates are essential to handle environmental changes, drift, and maintenance cycles
- Real-time decision-making close to the system = edge condition monitoring
  - → **MLOps** = cost reduction, consistent quality, increasing operational efficiency

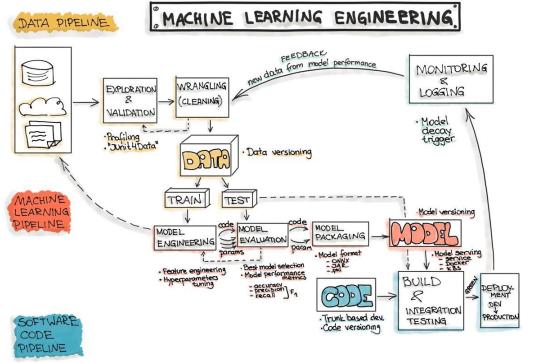




# MLOps for Edge Condition Monitoring

 MLOps = Paradigm used to create and maintain a machine learning model that will be deployed in production

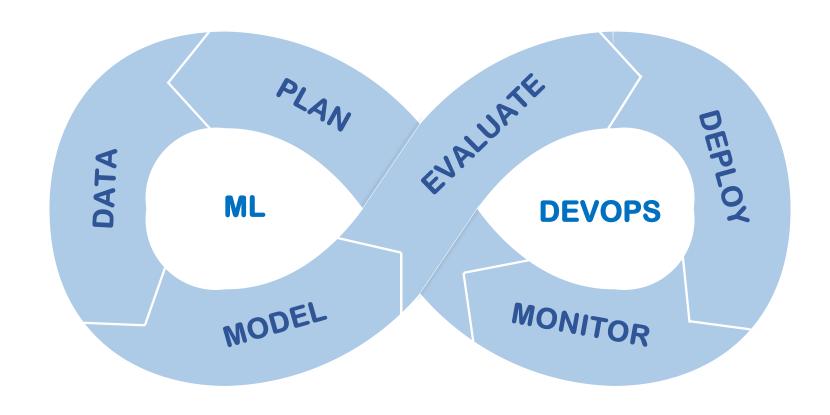








# MLOps lifecycle



# MLOps4Edge repertoire

### **Tools**

































### **Devices**



















# MLOps4ECM – Goals

Enabling companies to adopt MLOps practices for edge condition monitoring through the structured transfer of knowledge, tools, use cases and demonstrators in a pragmatic, industry-oriented manner.

Goal 1: Overview of available knowledge, needs, problems + definition of use-cases

Goal 2: Overview of DevOps and MLOps tools

Goal 3: Development and implementation of use-cases

Goal 4: Development and implementation of two demonstrators

Goal 5: Show the economic impact of using MLOps

Goal 6: Dissemination of the knowledge through seminars, workshops, education

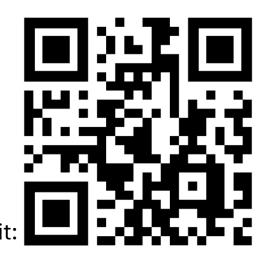






### Website

- General project info
- Workshops
- Guidelines
- Use-cases
- Demonstrators





### VLAIO TETRA MLOps4ECM

### MLOps for Edge Condition Monitoring

Binnen het kader van het VLAIO TETRA programma (TEchnologie TRAnsfer) voeren het IoT-lab van Hogeschool VIVES en de M-Group van de KU Leuven Campus Brugge gezamenlijk onderzoek uit om bedrijven praktische kennis, expertise en tools aan te reiken voor het operationaliseren en onderhouden van hun machine learning modellen. Dit stelt bedrijven in staat om de levenscyclus van hun machine learning modellen te bewaken en waar nodig te verfijnen.







### Projectsituering

Het bewaken van complexe industriële processen vraagt vaak naar datagedreven modellen om de gezondheidsstatus van de (sub-)systemen te bepalen. Tegenwoordig zijn er heel wat tools, handleidingen en opleidingen om die modellen te trainen en op te zetten. Het uitwerken van een datagedreven oplossing bestaat typisch uit dezelfde stappen: data verzamelen, verwerken, labelen, en mogelijke karakteristieken bepalen, om vervolgens het model te ontwerpen, trainen en optimaliseren, en implementeren. Ten slotte kan het resulterende model gebruikt worden in de bedrijfsomgeving, bv. om voorspellingen te maken. Het uitvoeren van deze stappen is tijdsintensief en verat specifieke kennis. Het resulterend model









### **Outline**

### Now

**Drift detection on edge devices** 

Lara Luys

Training, testing and deploying AI models on the edge

Sille Van Landschoot

Machine vision pipeline for efficient edge deployment

**Alexander D'hoore** 

Follow-up projects

16h00

Reception with demos



# Drift detection on edge devices

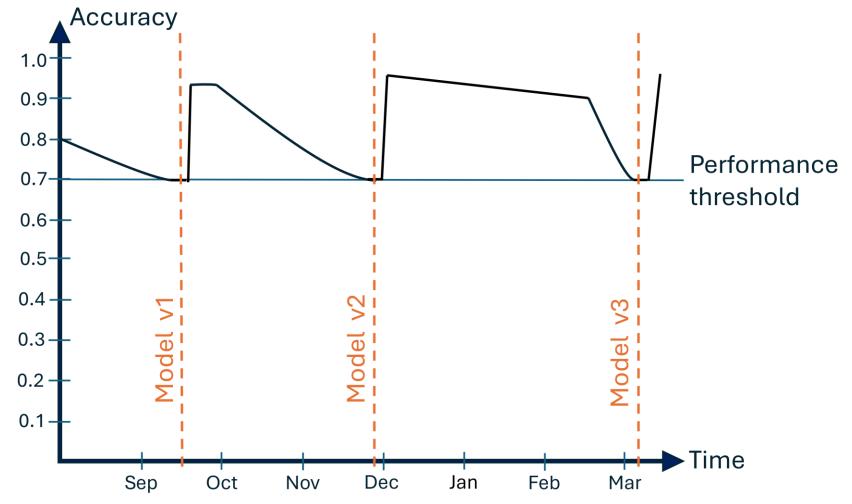
Lara Luys







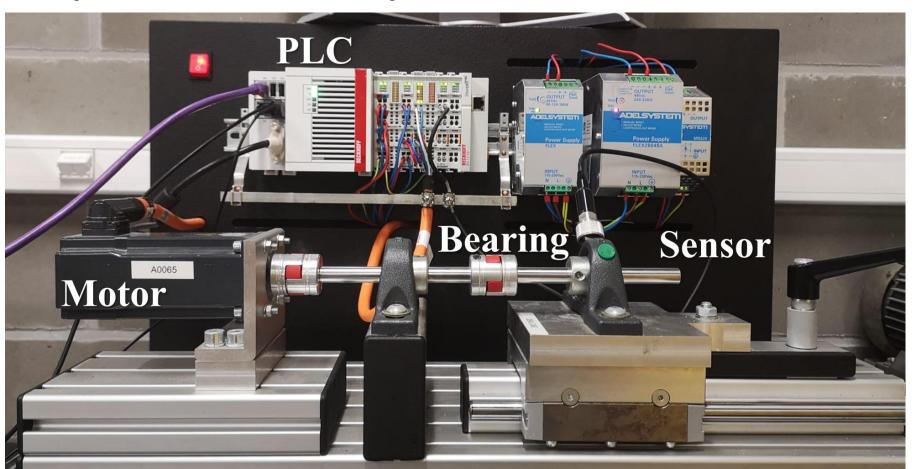
### What is drift?







# Time series drift detection demo (icw Beckhoff)











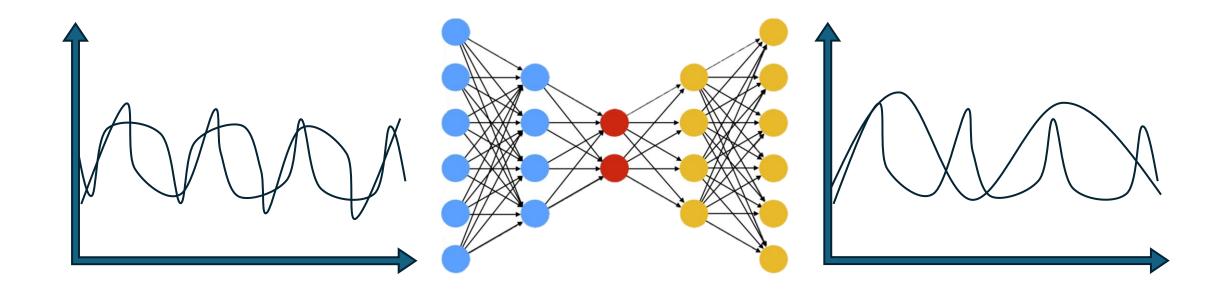
9/10/2025





# Drift detection demo (icw Beckhoff)

Goal: Drift detection of vibration data with autoencoder







# Magic Marble Machine demo (icw CTRL Engineering)

- Student projects
  - 1. Creation of the magic marble machine
- 2. Creation of machine vision models to compare timing for real-time applications.
- MLOps on PLC with vision







### Vandewiele use-case: The setup





- ML for fault detection in weaving machines
  - Post process
  - Realtime
- 2 datasets: Only faulty data
  - R&D dataset = training data
  - - + brakes
  - 3 possible drift moments
- Given labels fault detection: results of rulebased algorithms







# Vandewiele use-case: Research questions

1. Does the change of setting, sensors or brakes cause data drift?

2. Which type of edge device is needed for drift on the edge?







### Vandewiele use-case: Is there drift?

Try different types of drift detection with retraining



- Supervised (use labels)
  - DDM
  - ADWIN
  - EDDM
  - STEPD
- Unsupervised (No labels)
  - KS
  - D3
  - OCDD
  - SPLL
  - BNDM

Have hyperparameters to tune sensitivity of drift detector

Use multi-objective hyperparameter search with Optuna





# Vandewiele use-case: Multi-objective search

• 3 Metrics which compare results detectors to possible drift

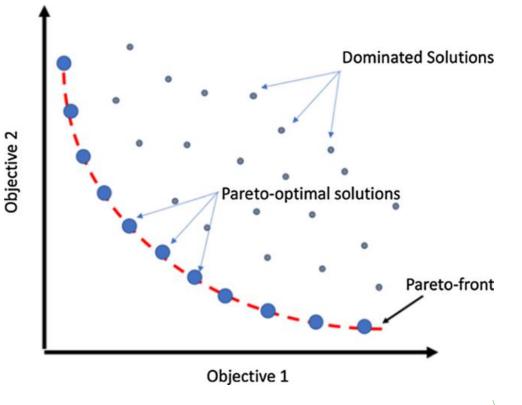
moments

Mean time to detection = MTD

Mean detection ratio = MDR

Number of false detections = FD

Result = Pareto front

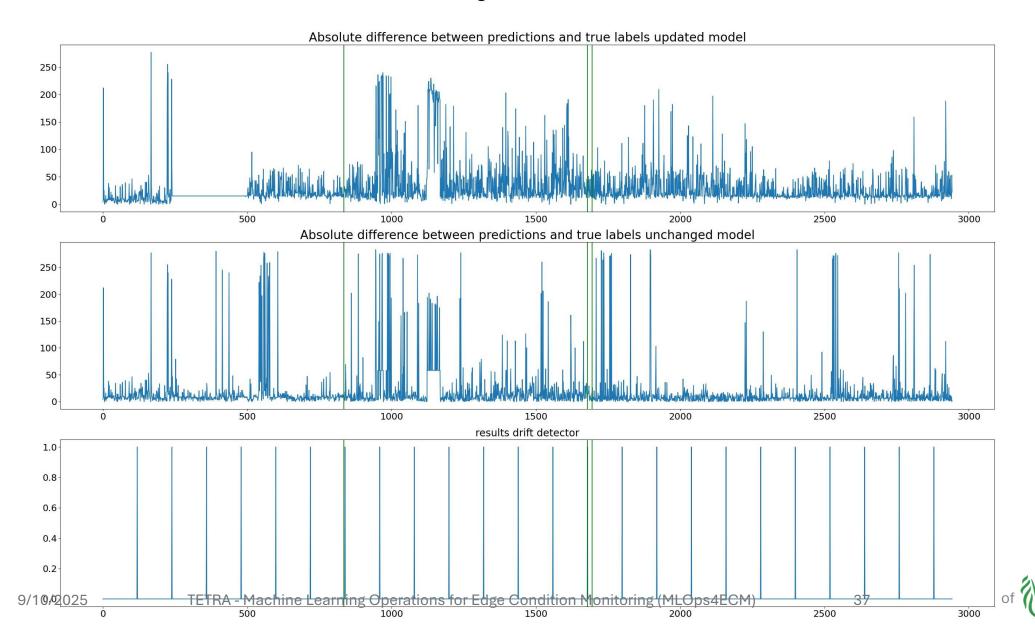


### Vandewiele use-case: Results





Monitoring results for bndm

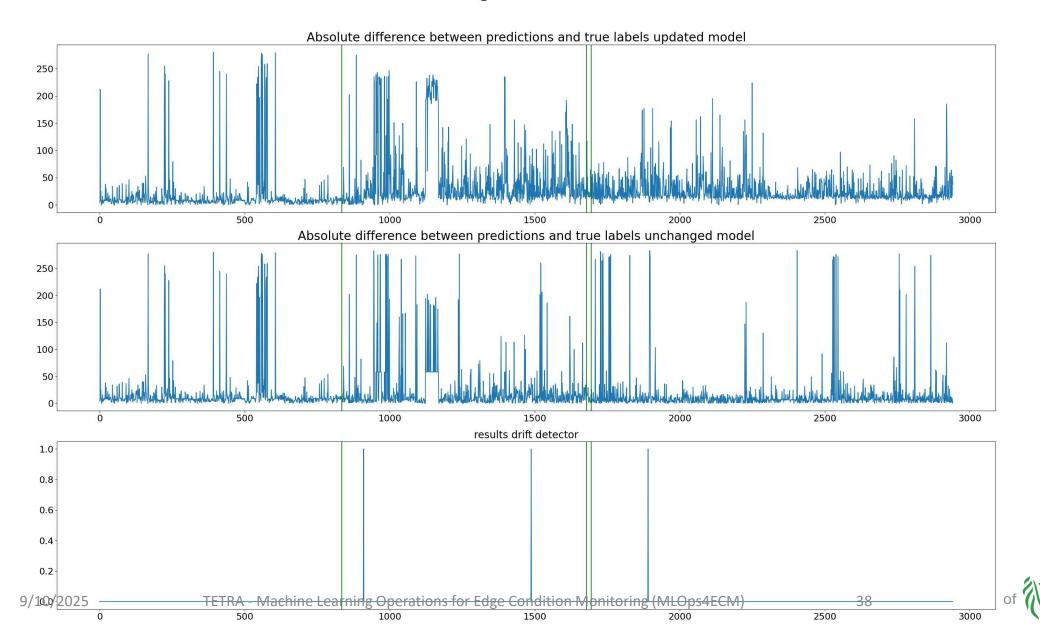


#### Vandewiele use-case: Results





Monitoring results for ocdd







# Vandewiele use-case: Conclusion question 1

- Given data is very noisy
  - > Faulty data is very inconsistent: Each failure is different
- How to improve?
  - → Retry with both faulty and working data this data will be more stable.
  - → Drift detectors = made for "streaming data"





# Vandewiele use-case: which edge device is needed for implementation?

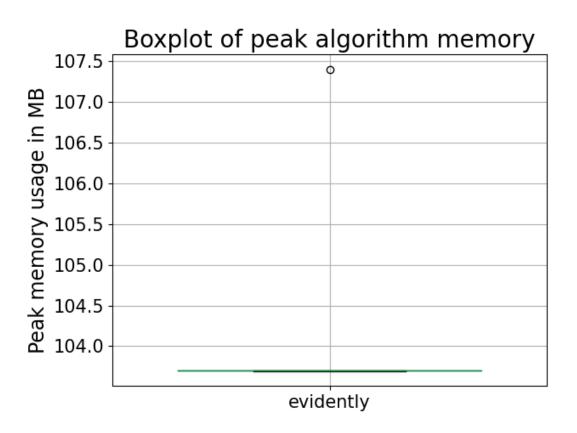
- Use of tool: Evidently
  - Advantage: Easy to implement + nice looking reports
  - Disadvantage: uses a lot of recourses, uses simple threshold techniques
- Drift detectors
  - Advantage: smaller resource usage (depends on detector)
  - Disadvantage: More difficult to implement

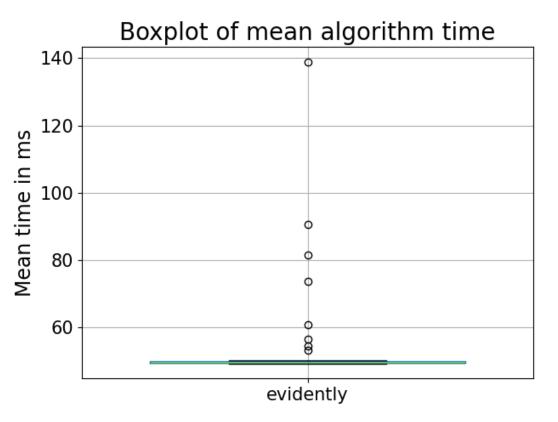






# Vandewiele use-case: Evidently resources





Includes model inference + drift detection

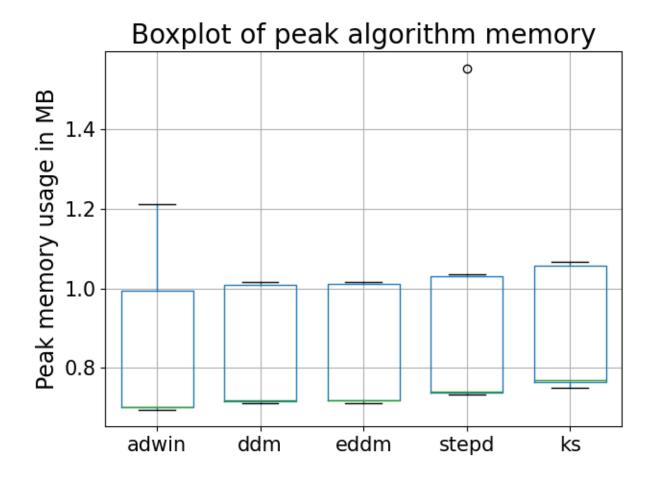
Includes only drift detection







# Vandewiele use-case: Drift detectors memory



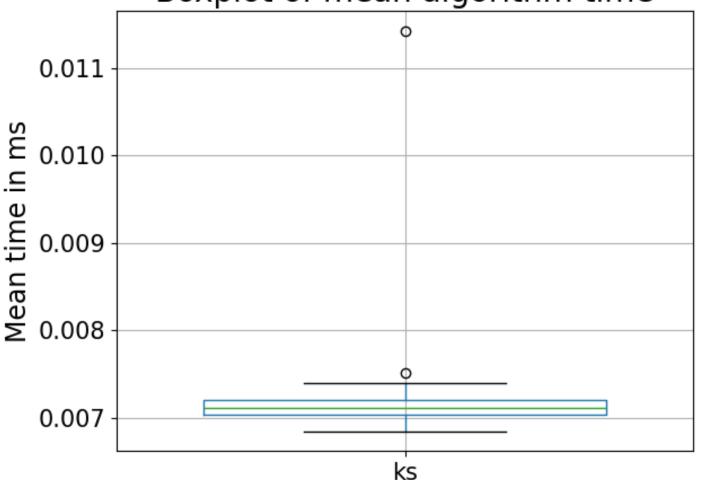
Includes model inference + drift detection





# Vandewiele use-case: Drift detectors timing

Boxplot of mean algorithm time



Includes only drift detection





# Vandewiele use-case: Type of edge device

- If Evidently: > 100MB
  - Microprocessors: Needs Linux or Windows to run Python
  - Linux boards, industrial PC's, PC-based PLC, etc.
  - E.g. Jetson, ARM Cortex A



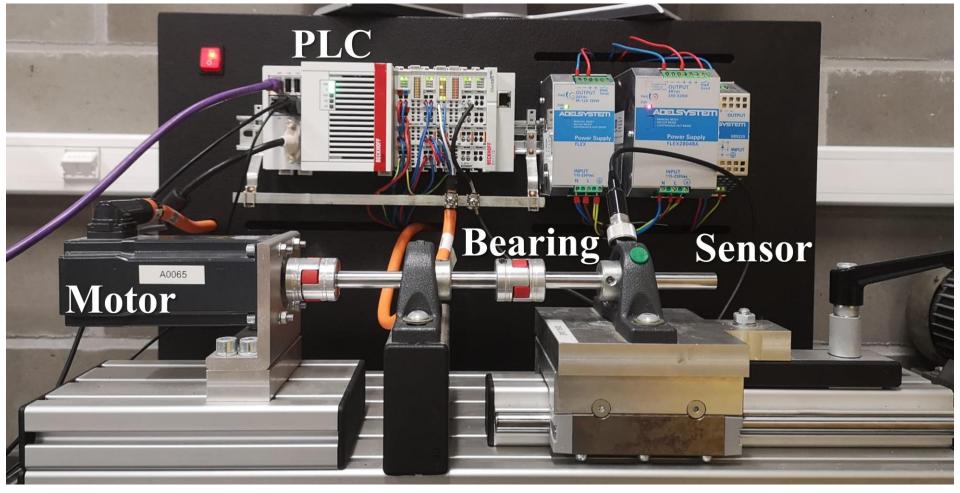
- If drift detectors: > 2-6 MB
  - Microcontrollers: Translation to C/C++ needed
  - E.g. ARM cortex M, ESP32







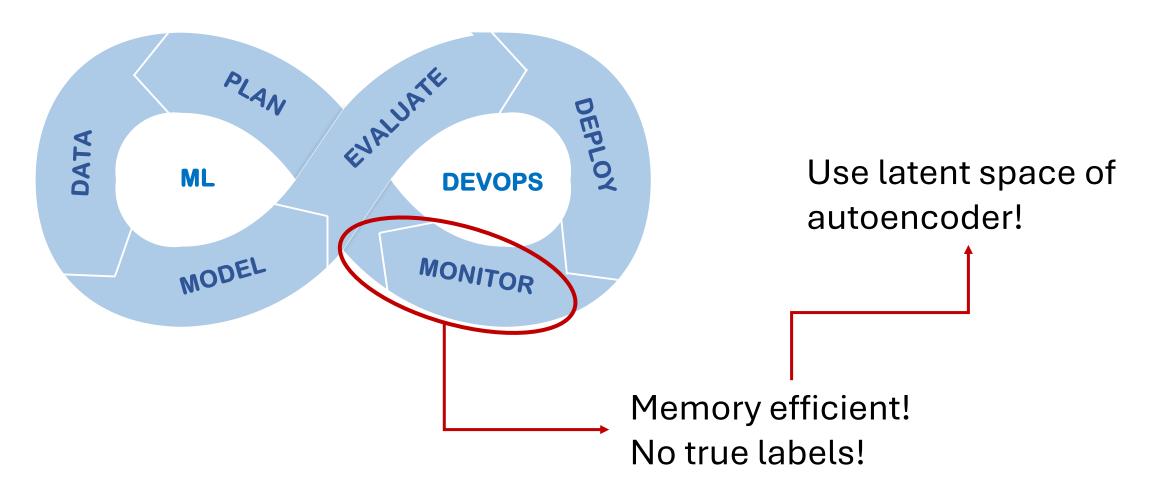
## Beckhoff use-case: The setup







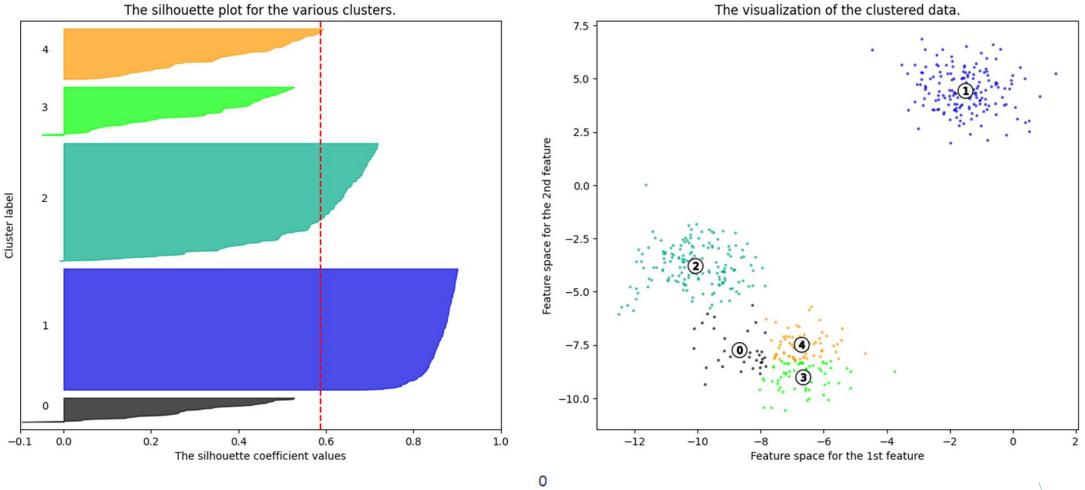
# Beckhoff use-case: Monitoring in production







# Beckhoff use-case: The latent space







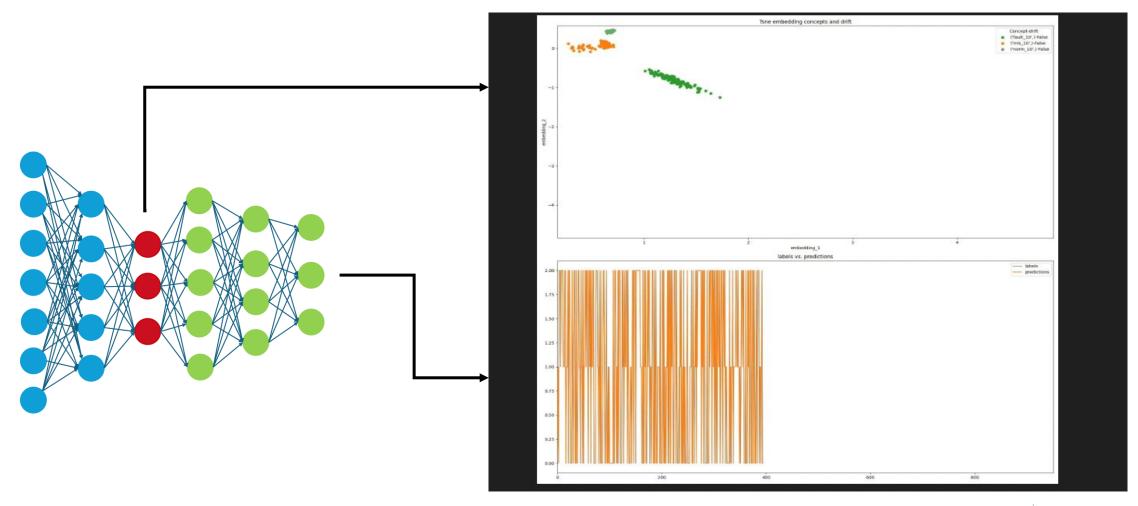
# Beckhoff use-case: Latent-Space Drift

etector Use Optuna to find A €stvith best Destentaingeted statemyessace based on silhouette score Concept clustering





#### Beckhoff use-case: The results

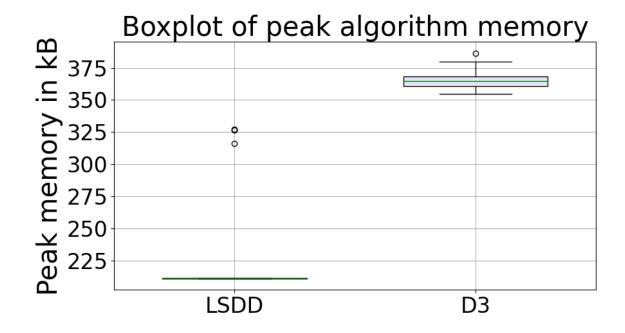


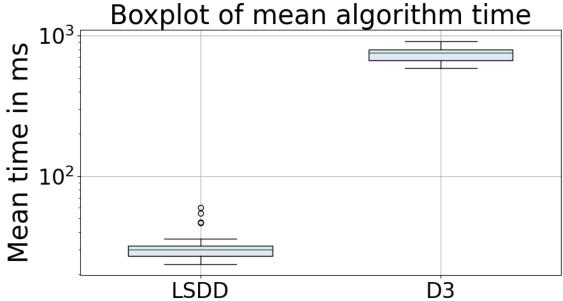




#### Beckhoff use-case: The results

Algorithm	Parameters	MTD	MDR	FD
LSDD	$N=13$ , $\gamma = 0.077$ , $k=73$	2.666	0	0
LSDD	$N=10, \gamma = 0.05, k=74$	2.666	0	0
D3	w=82, $\rho$ =0.544, threshold=0.6525, k=47	8.5	0.333	0



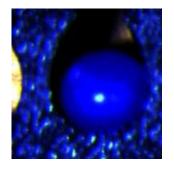




# CTRL-Engineering use-case: Does LSDD work on images?

- Artificially Brighten and darken images = 2 concepts:
  - Dark images
  - Bright images
  - Switch 9 times between 2 concepts





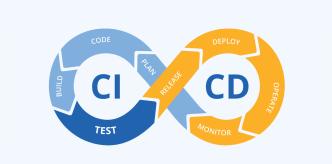
• MTD = 313.555, MDR = 0, FD = 0

# Training, testing and deploying Al models on the edge

Sille Van Landschoot

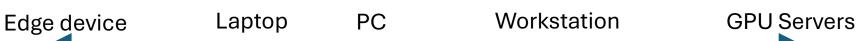


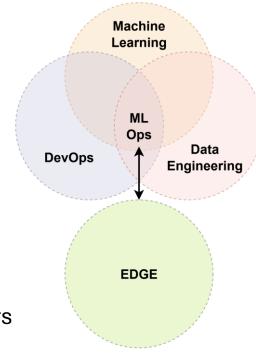






- MLOPS principles: creation of machine learning products
  - CI/CD automation, workflow orchestration and reproducibility;
  - versioning of data, model, and code;
  - collaboration;
  - continuous ML training and evaluation;
- Edge ML
  - Specialized hardware
  - Unique and custom approaches
  - Constrained resources









### Edge friction

There's no place like **127.0.0.1** 

- Developing for the edge is hard
  - We want to use the edge device for deployment and inference
  - But we want to use our machine for development and training
  - Even better is to use my machine and offload compute to dedicated servers
- Al projects
  - Large datasets might not fit the edge device or local develoment machine
  - Compute intensive workloads for training
  - Specialized hardware: Training requirements != inference requirements
- One size doesn't fit all:
   Different stages of the ML lifecycle have different requirements.







### Development vs. Production Mismatch

- Resource Requirements: demands powerful computational resources that are typically available in server environments.
- **Data Management Challenges:** Datasets, often exceed what can be *stored* on individual developer machines. This creates *bottlenecks*.
- Training Infrastructure Needs: Al training processes are computationally intensive and time-consuming, requiring dedicated infrastructure that can be scheduled and managed independently of development workstations.











## **Operational Constraints**

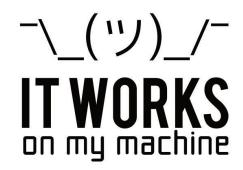
- Production System Isolation: Production edge devices must remain dedicated to inference operations and cannot be used for training or testing activities.
- **Development-Production Validation Gap:** Development teams need access to *edge-equivalent testing environments* to validate model performance and inference characteristics.
- Model Engineering Flexibility: The system must accommodate arbitrary model architecture choices and implementation decisions without constraining the development process.
- **Reproducibility and Tracking:** Dataset evolution and model versioning must be systematically tracked to ensure experiment *reproducibility*.



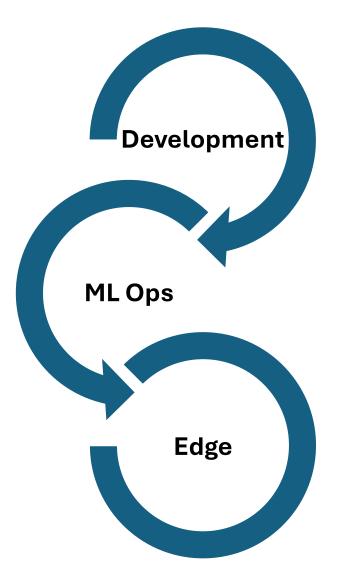


# Bridging the gap(s)

- The core challenge is bridging the gap between
- Al development: requiring powerful servers, large datasets, and iterative experimentation
- Production deployment: requiring lightweight, real-time inference on resource-constrained edge devices











# MLOps is bridging the gap

#### The MLOps solution

- Capture & Curate: Data and model versioning, reproducibility
- Automate: Orchestrate workflows for training, testing and development
- Package & Ship: Build reusable, self containing containers
- Observe & Learn: Traceability by tracking progress and results
- Decide: Maintain control over rollout and rollback







# Bridging the gap

- Capture & Curate:
  - Git
  - LakeFS
  - Minio
- Automate
  - Prefect
- Package & Ship
  - Docker containers
  - Docker Registry







• Open-source, self-hosted solutions



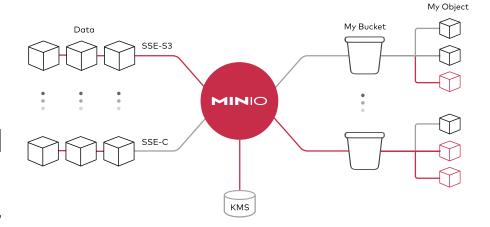


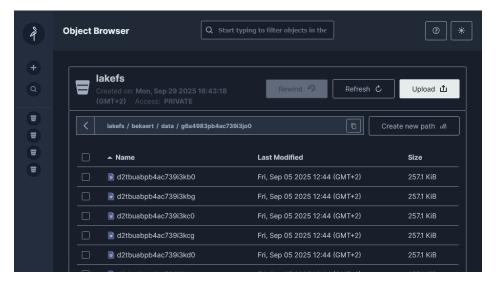


#### Minio



- S3-compatible object storage you run yourself—works with standard S3 SDKs/CLI and pre-signed URLs.
- **Keep data/models on-prem/edge** (latency data sovereignty).
- One **bucketed store** for datasets, models, artifacts, telemetry.
- Versioning + immutability → reproducibility and audits.
- Scales out; cheap per-TB versus managed cloud.







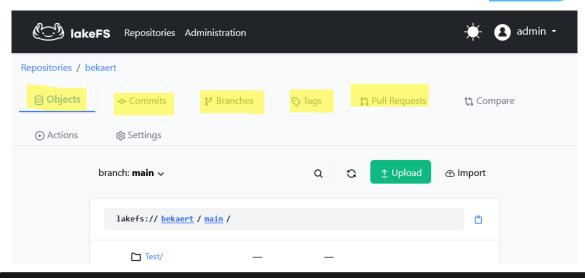
#### LakeFS



- Git-like version control for datasets
  - zero-copy branches/commits/merges
- Remote-first data management
- Features
  - Reproducible training: pin datasets/models by commit ID
  - **Isolated experiments**: branch fast, run pipelines, merge/promote when green
  - Safe rollbacks: revert a bad dataset/model instantly
  - Auditability: lineage from data → model
     → device release







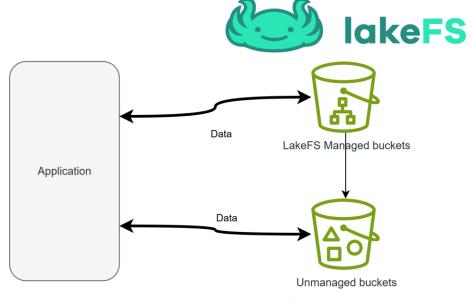




#### Minio + LakeFS

- Uses S3/Minio object storage as storage backend
- Open Source, Easy to self-host
- Containerized single command deployment using Docker, Docker Compose, Kubernetes,...
- Lightweight, perfect fit for development
- Access data over LakeFS or S3 API (open implementation)







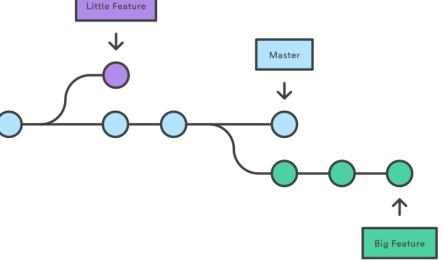


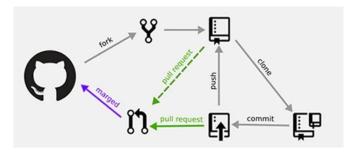


# Git Version Control



- Single source of truth for code, configs, pipelines, infrastructure as code
- GitOps triggers: commit/PR → CI runs tests/build/train → gated deploy
- Traceability: tag releases; link commits to data/model versions & device fleets
- Reproducibility: pin everything by commit/tag (env, params, artifacts)
- Rollbacks: revert commit/tag → auto-rollback pipeline, fast recovery





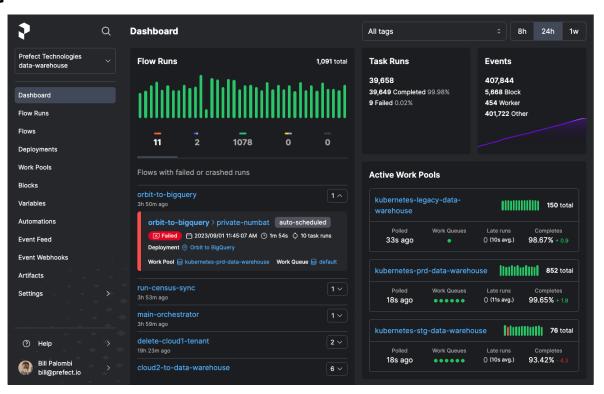




# Prefect



- "Modern workflow orchestration for data and ML engineers"
  - Open source
  - Manage, schedule, and monitor data workflows
  - Flexibility and scalability for various automation tasks
  - Useful in data engineering, MLOps, and machine learning pipelines
- Python-First and Code-Driven
- Integrations with ML libraries, data sources and platforms







# Prefect



#### Flow & Task Abstraction

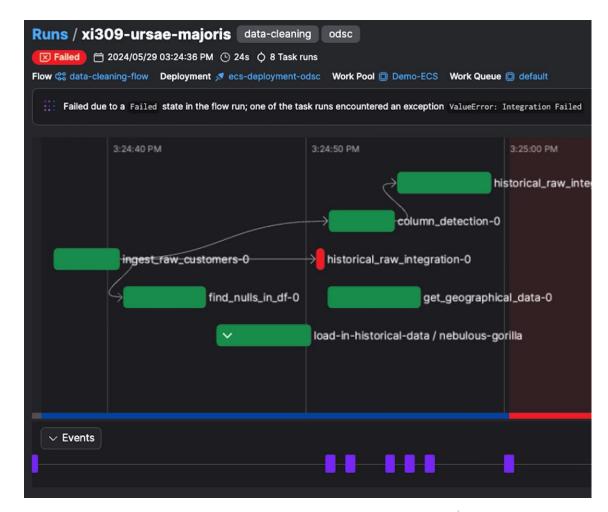
- Organize tasks into flows to structure workflows easily
- Task dependencies defined intuitively in Python

#### Resilience & Fault Tolerance

- Automatic retries, timeouts, and error handling
- Options to resume, skip, or rollback tasks upon failure

#### Observability & Real-Time Monitoring

- Track task status, view logs, and monitor performance
- Real-time alerts for task failures or anomalies







#### **Prefect Tasks and Flows**

- Tasks: Small, reusable functions representing individual steps in a workflow (e.g., data extraction, processing).
- **Flows**: Organized sequences of tasks forming a complete workflow or pipeline.
- Modular Design: Build, reuse, and manage tasks within flows for efficient development.
- **Dependency Management**: Control task order and relationships for accurate, predictable execution.

```
@task
def mbed_compile(target, compiler) -> None:
    run_mbed_container(f"mbed compile -m {target} -t {compiler}")

@task
def prepare_model(tflite_file: str) -> None:
    model_file_name = "model.h"
    shell_run_command(command=f"xxd -i {tflite_file_name} > {model_file_name}")
    replace_text = tflite_file_name.replace('/', '_').replace('.', '_')
    shell_run_command(command=f"sed -i 's/'{replace_text}'/g_model/g' {model_file_name}")
```

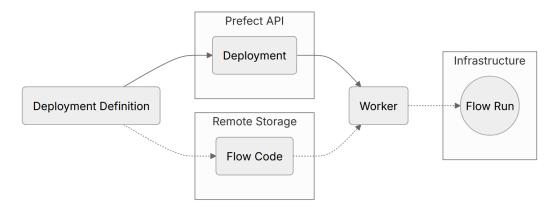


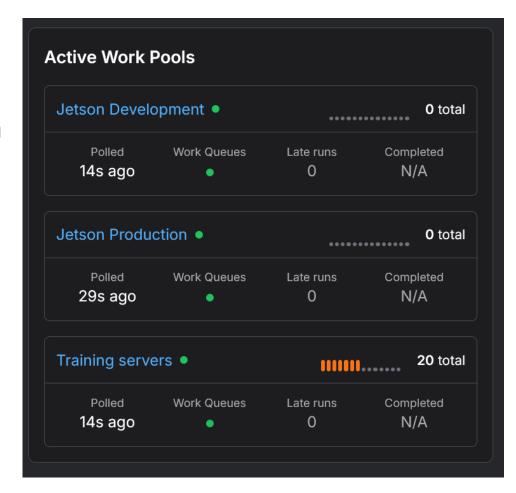




#### **Prefect Work Pools**

- Dynamic infrastructure provisioning
- Specialized pools for specific workloads with priorities and concurrency
- Different execution environments:
  - Docker
  - Process
  - Kubernetes
  - Azure, Google Cloud, AWS,...



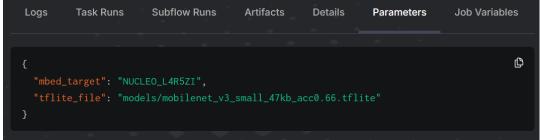




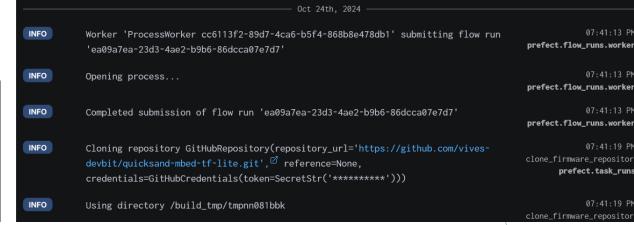


#### **Prefect Runs**

- Represents each instance of a workflow execution.
- Tracks status, logs, and outputs for real-time monitoring.
- Enables issue identification with success, failure, and retry info.











#### **Prefect Variables**

- **Dynamic configuration elements** used to streamline and adapt workflows.
- Centralized Configuration: Manage settings and constants across flows from a single place.
- Environment-Aware: Adjust workflow behavior based on environments (dev, staging, production).
- Reusable & Secure: Define once, reuse everywhere; securely stored to protect sensitive information.

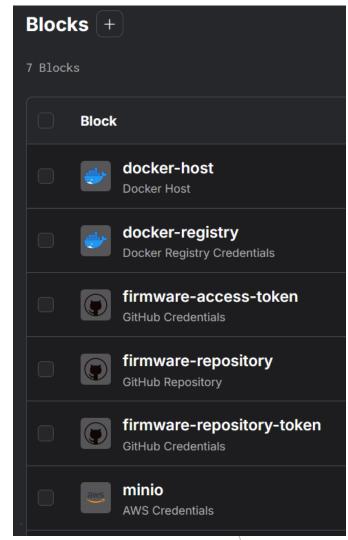






#### **Prefect Blocks**

- Reusable, modular building blocks that simplify the development, deployment, and management of workflows.
- **Pre-built integrations**: Connect to popular tools (databases, cloud providers, messaging services) without custom code.
- Customizable: Create and share custom blocks tailored to unique needs.
- Easy configuration: Define and reuse configurations across tasks to reduce setup time.







#### **Prefect Automations**

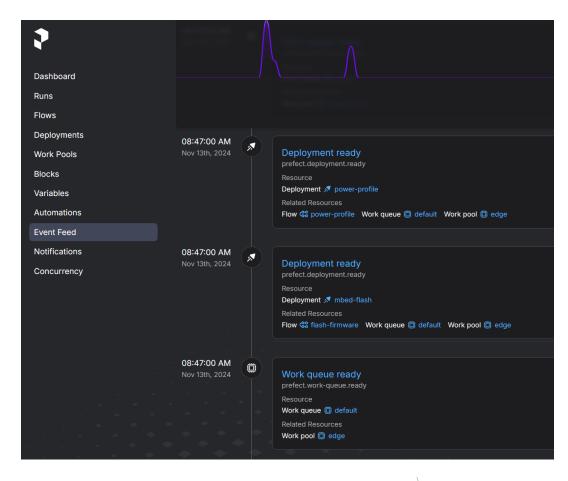
- Automated Actions: Trigger responses based on events (e.g., failures, successes).
- No-Code Setup: Easily configure alerts, retries, or escalations without custom code.
- Improves Reliability: Reduces manual intervention and ensures timely responses.
- Customizable: Tailor actions to specific workflows and business needs.





#### Prefect Event feed & Notifications

- Real-time updates on workflow events (starts, completions, failures).
- Centralized view for monitoring all activity.
- Quick identification of issues.
- Alerts for key events (failures, retries)
   via email, Slack, etc.
- Customizable **triggers** to focus on relevant updates.
- Keeps teams informed for proactive responses.

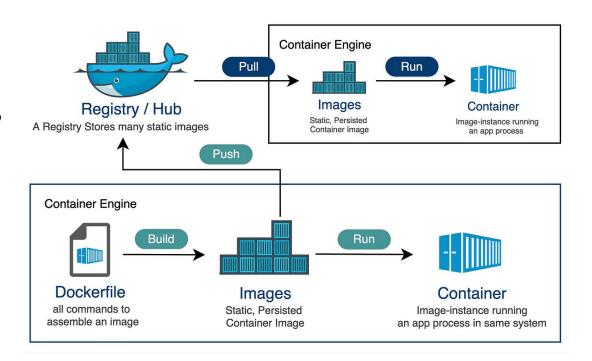






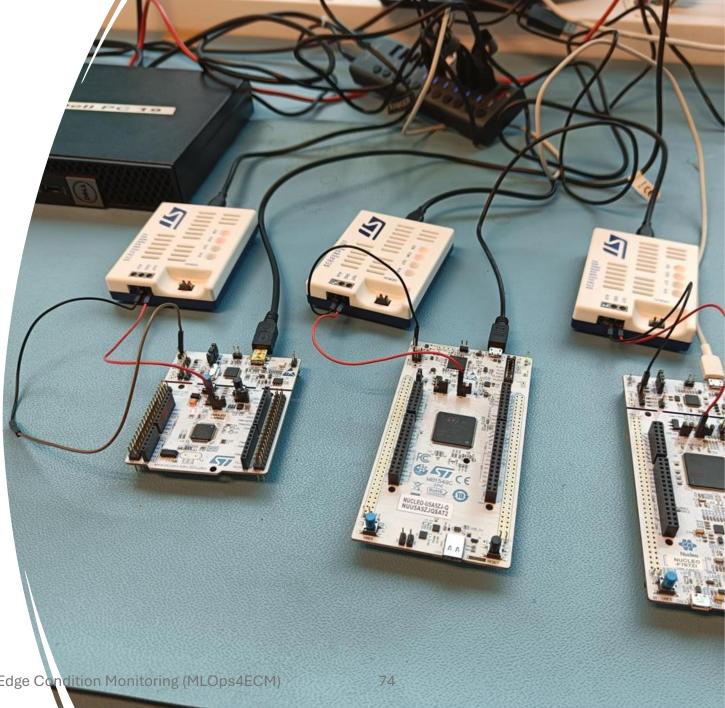
# Docker Image management

- Docker Registry for the management of Docker Images
- Efficient (re)use of Docker images, regardless of whether they need to be executed on the server or edge
- Centralized registry for easy management
- Easy local deployment for on-prem data sovereignty and low-latency pulls



#### Two use-cases

- Testing and validating 900+ Al models on a microcontroller system with energy consumption monitoring
- Pipeline for training, testing and deploying AI models on the Edge



# Pipeline for training, testing and deploying AI models on the Edge

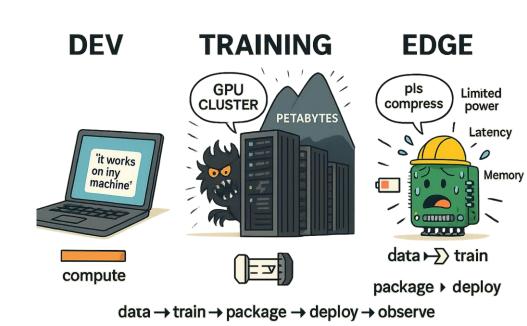






#### Machine vision on Jetson Orin Nano

- Large image dataset
  - Hard to manage and maintain
- Training on edge device is impractical
  - Best representation of production environment
  - Does not compare to a development environment
  - Not optimized for training, only inference
  - Constrained resources
    - Storage
    - Memory
    - CPU
    - Bandwidth







### Development on local machine

- Pros
  - Allows for fast development cycles and instant feedback
  - Use subsamples of the dataset
  - Offline work
  - Full control of OS, tools and drivers
- X Cons
  - Laptops with limited resources: CPU, GPU, RAM, Storage,...
  - Does not represent the production/edge device
  - Does not scale









#### Development on a server

- Pros
  - Virtually limitless in resources: CPU, GPU, RAM, Storage
  - Remote development, anywhere and anytime
  - Long running detached jobs
  - Data proximity for large dataset
- X Cons
  - Cost
  - Shared resources among multiple developers/projects
  - Still does not represent the edge system









## Development on the edge

- <a>Pros</a>
  - Represents the deployment environment
- X Cons
  - Very limited resources, No GPU,
  - Ineffective for training
  - No (full) support for training frameworks, tools and libraries
  - Often a different architecture

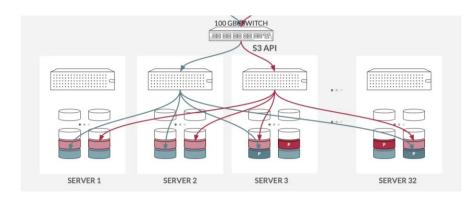




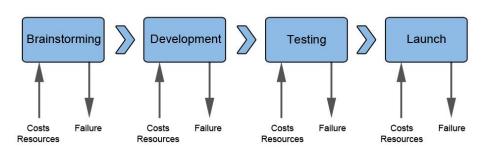


#### Development requirements

- Scalable data storage and versioning for large image datasets
- **GPU-accelerated training** infrastructure separate from developer machines
- Automated pipeline orchestration for training workflows
- Model validation and testing before deployment
- Reproducible experiments with dataset change tracking





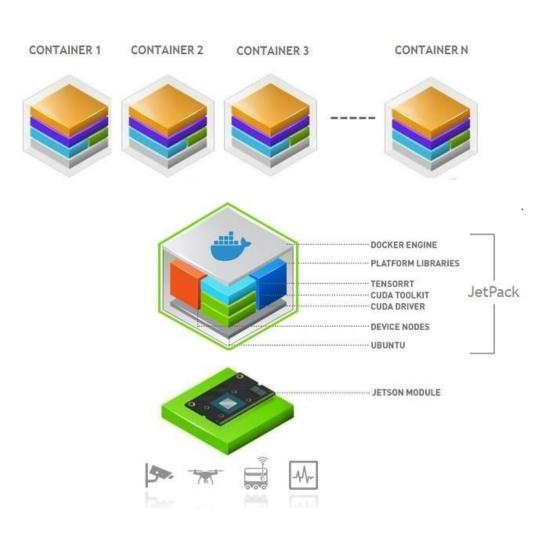






# Production requirements

- Real-time inference on NVIDIA Jetson edge devices
- Containerized deployment for consistent environments
- Automated model updates from development to edge
- Monitoring and feedback loops for continuous improvement
- Camera abstraction support via standardized protocols (RTSP or gstreamer)







## Integration Challenges

GPU/Data Servers Development Edge Device

- Resource optimization (server GPUs vs. edge efficiency)
- Seamless pipeline from development to production
- Version control for both code and data
- Automated testing and validation gates
- DevOps tooling for ML lifecycle management





#### Infrastructure

- MinIO: Object storage for models and artifacts
- Prefect: Workflow orchestration and scheduling
- Docker Registry: Container image distribution
- MediaMTX: RTSP streaming server for testing
- Grafana: Real-time monitoring and alerting
- Infrastructure should support the process of development and deployment

Infrastructure

Development

Deployment





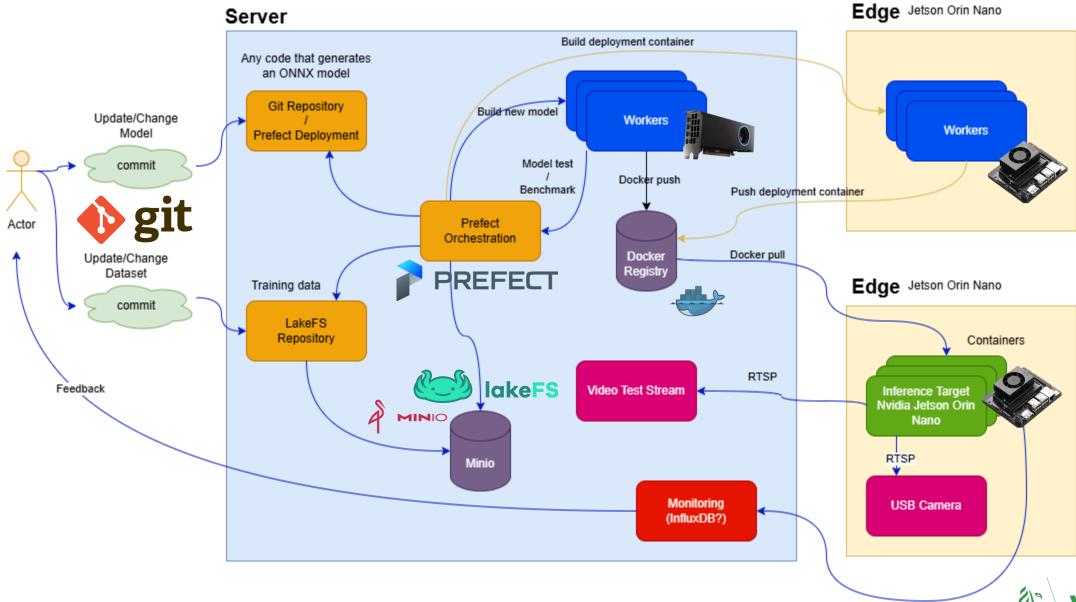


# **MLOps Workflow**

- **1. Code Commit** → Git repository triggers pipeline
- 2. Data Versioning → LakeFS manages dataset branches
- **3. Training Execution** → Prefect orchestrates containerized training
- **4. Model Validation** → Automated testing before deployment
- **5. Container Building** → Docker images with trained models
- **6. Edge Deployment** → Jetson devices pull and run new containers
- 7. Monitoring Loop → Performance feedback for continuous improvement





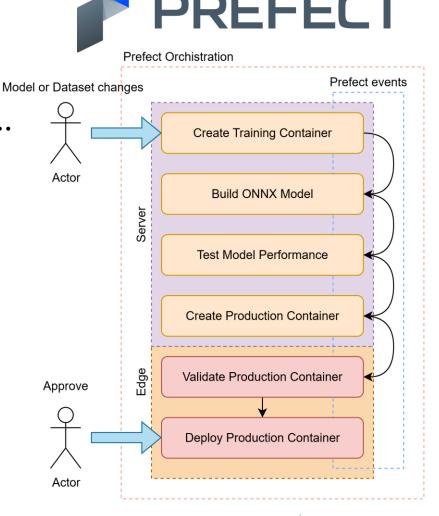






## **MLOps Workflow**

- Developer produces events
  - Commit code, Update Model, Change dataset, ...
- Prefect Deployments run flows and task
- Tasks generate events to trigger other deployments
- Flows are assigned to different work-pools
- Wen everything is fine, developer triggers deployment







#### What Worked Well

#### Architectural Decisions:

- - Clean Separation: Keeping ML core logic independent from orchestration proved invaluable for development speed and testing
- - Container Strategy: Docker containers provide consistency across development, testing, and production environments
- - Standard protocols: RTSP for cameras and ONNX for models ensure compatibility and interoperability

#### Development Experience:

- Dry Run Mode: Game-changer for rapid iteration without waiting for full training cycles
- - Data Versioning: LakeFS integration provides crucial experiment reproducibility
- Automated Testing: Built-in performance evaluation catches issues before deployment







## **Key Lessons Learned**

#### MLOps Pipeline Design:

- - Infrastructure Independence: Core ML functions should work standalone for faster development
- - **Testing Strategy**: Multiple levels of testing (structural, performance, integration) catch different types of issues
- - **Monitoring is Critical**: Real-time feedback from edge devices informs model improvements
- - **Version Everything**: Not just code, but data, models, and container images

#### Edge AI Deployment:

- Resource Optimization: ONNX format and proper model optimization crucial for Jetson performance
- - **Update Mechanisms**: Container-based deployment enables reliable model updates in production
- Camera Integration: RTSP provides standardized interface across different industrial camera systems





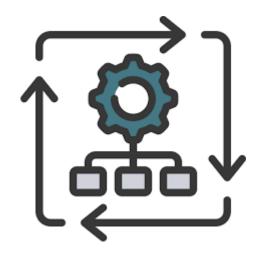


#### Conclusion

• This project successfully demonstrates a **complete MLOps pipeline** for industrial AI deployment, bridging the gap between research/development and production edge inference.

 This MLOps pipeline provides a reusable template for deploying Al models in industrial edge environments, combining the flexibility needed for Al development with the reliability required for production systems.

# Auto deployment & energy monitoring



Sille Van Landschoot







# Auto deployment & energy monitoring

- Energy monitoring of AI models on embedded devices
  - 150+ models
  - 6+ targets
- Multiple firmware implementations and revisions
- New targets and new models at any time
   #tests = #models x #targets x #version
- Solution: Automate!
  - How? Prefect?





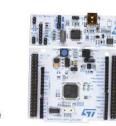




# Auto deployment & energy monitoring

- Version control: Git(Hub)
- Model storage: Minio, S3 compatible
- **Docker image repository**: Docker Registry
- Task orchestration: Prefect server
- Edge nodes: Prefect workers
- Targets: STM32 Nucleo development boards
- Energy monitoring: STLink-v3PWR













(32-pin MCU)

(64-pin MCU)

(144-pin MCU)





Prefect Server

Prefect Worker

Edge workload



#### **End-user**

Blackbox

#### Inputs

- Firmware
- Model
- Target

#### **Output**

Power report



Power monitor

STLink-v3PWR

Prefect Worker

Server workload

Model .....

Target .....

Server/Cloud

Target

**Embedded Target** 

Nucleo

Edge





### Targets

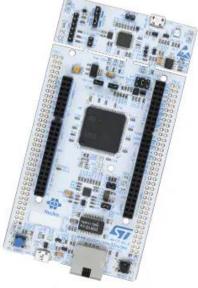
- Nucleo target development boards
  - Microcontrollers
  - Microcontroller architectures Cortex-M
  - Memory constraints
    - RAM
    - Flash
  - Peripherals
- USB Connections to edge computer
  - UART comunication



NUCLEO-32 (32-pin MCU)



NUCLEO-64 (64-pin MCU)



NUCLEO-144 (144-pin MCU)









#### Power monitor

- STLink v3PWR
- Managed over UART
  - Python script
  - Start/stop measurements
  - Set power and voltages
  - Reset targets
  - UART for target commands



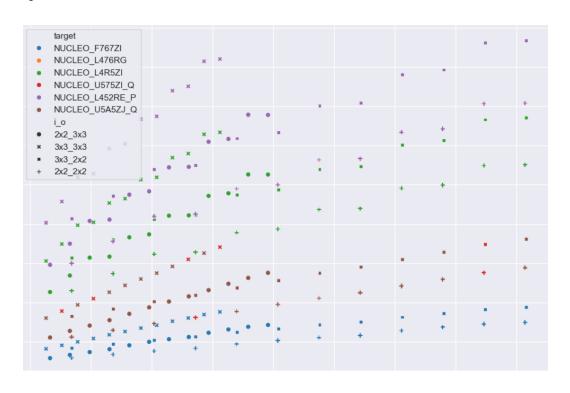


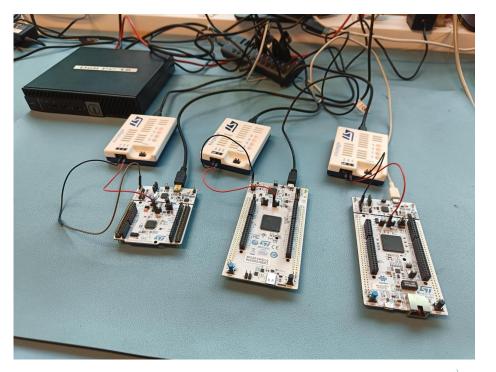




# Reports and logs

• Plain old **Python scripts or Jupyter Notebook** with matplotlib, pandas,... decorated with Prefect Tasks and Flows









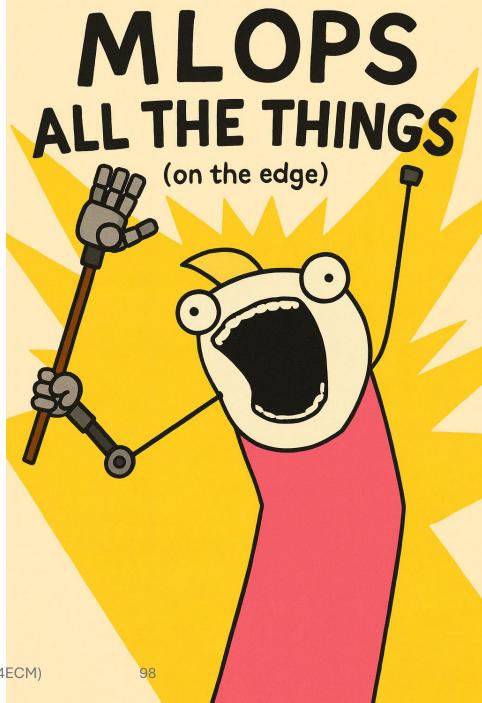
#### Lessons learned

- MLOps infrastructure requires some upfront effort but pays of in the long run
- Prefect provides workflow orchestration with minimal overhead
  - Simply add some annotations to your Python code
- Decoupling video streams from physical camera's allow for flexibility during development
  - Share camera streams across different systems and stages in the development cycle
- Implement 'dry-run' options to speed up MLOps development
- Open-source solutions exist to support the ML Ops workflows
- Combining lightweight tools can support development AND production processes



#### Conclusion

- Automation using Prefect is very flexible and easy for Al workflows
- Minio and LakeFS allow for large and remote datasets with version control
- Containers allow for flexible execution across different machines and targets
- Version control of code and models provides traceability and reproducibility and rollbacks



# Machine vision pipeline for efficient edge deployment

Alexander D'hoore







# Use-case: Leap Technologies

- **Goal**: Deploy an occlusion-aware **instance segmentation** model on an edge device. **Low-cost hardware** is preferred. What are the minimum hardware resources required?
- Tools used: PyTorch, Albumentations, Ray Data, Ray Tune, MinIO (S3), TensorBoard, Proxmox, ONNX Runtime, OpenVINO, ExecuTorch, RKNN Toolkit 2
- Edge device(s): Raxda Rock 5B Rockchip RK3588 SoC
- Result: Data and training pipeline, with model deployment to an NPU







### Use-case: Leap Technologies

- Where did this project start?
- Leap Technologies provided us with a model:

**ORCNN**: Learning to See the Invisible: End-to-End Trainable Amodal Instance Segmentation

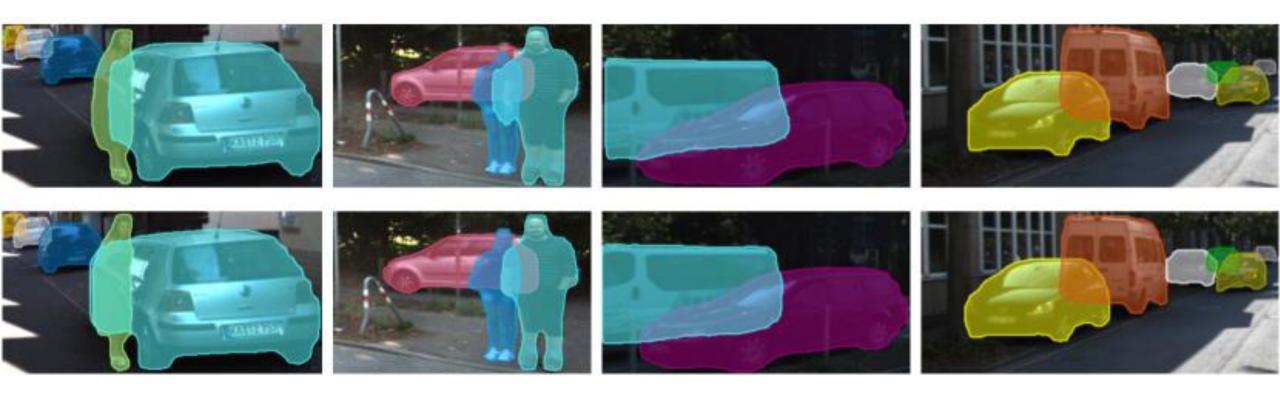
And a question: how can we run this efficiently on an edge device?







# Occlusion-aware instance segmentation







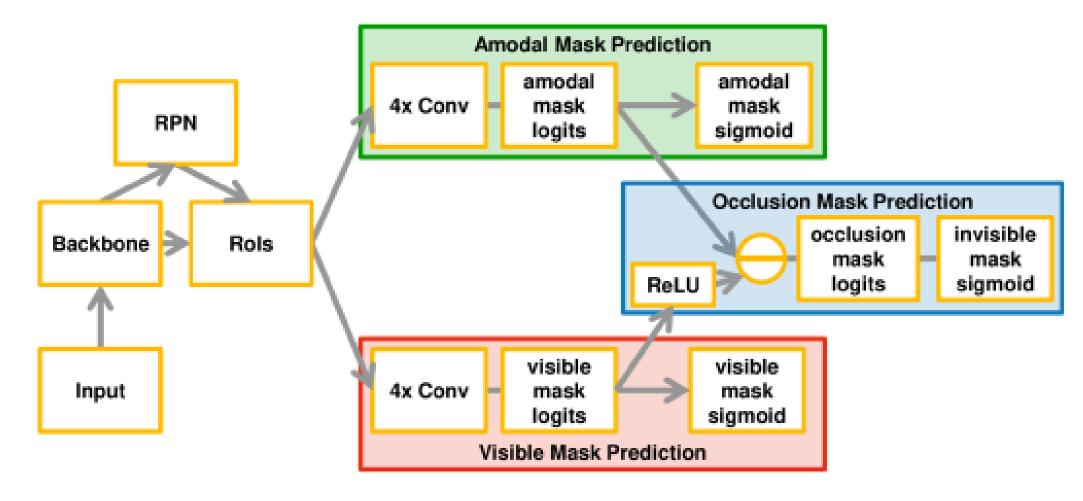
# Occlusion-aware instance segmentation

- Image classification: Assigns a single label to the entire image
- **Object detection**: Identifies objects with bounding boxes and labels for each
- Instance segmentation: Generates pixel-level segmentation masks for each detected object
- Occlusion-aware segmentation: Extends instance segmentation by predicting masks for both visible and occluded (invisible) object parts
- These techniques build upon one another, with each successive method being more computationally demanding, especially on edge devices.





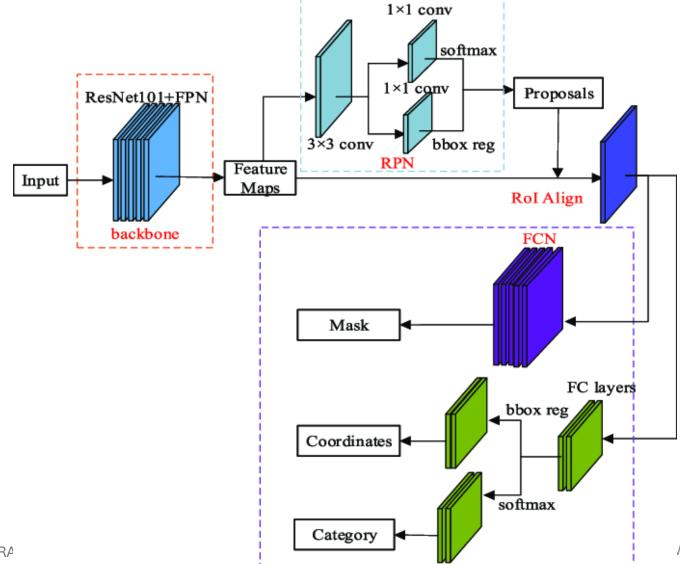
#### ORCNN = Mask-RCNN + extra head







#### Mask-RCNN architecture



Three branches



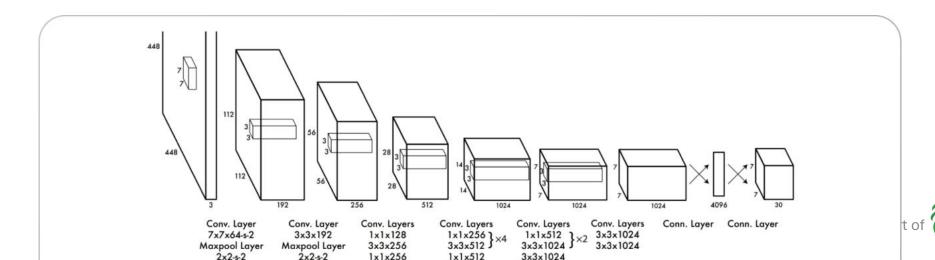


## Two-stage vs one-stage models

- Mask-RCNN: A classic two-stage segmentation model
  - Quite good accuracy, but a bit slow

9/10/2025

- Typical two-stage model: first backbone -> than boxes/masks
- Better for edge deployment: one-stage models
  - Popular: YOLO (You Only Look Once) family of models
  - Faster inference, but worse performance (in general)







# Note about licensing

- Open source comes in two flavours:
- Copy-left licenses (GPL, AGPL)
  - All code must remain open
  - Beware if you use in commercial
  - E.g. Linux, Ultralytics YOLO
- Permissive licenses (MIT, Apache)
  - Code free to use how you want
  - Often just need to mention authors
  - E.g. Pytorch, Torchvision

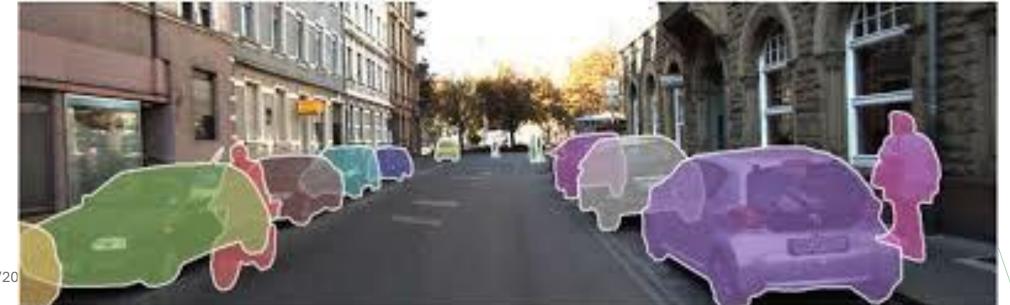






#### Occlusion-aware dataset?

- Some datasets exist: COCOA, D2SA, KINS
- Problem: Ground trust often hand-annotated
  - Humans must guess what's behind occluded areas
- => Resulting in **low-quality ground truth labels**









### Occlusion-aware dataset?

- Better idea:
- Generate a dataset by copy-pasting objects on top of each other
- This heavy data augmentation is a form of synthetic data

- Synthetic data => we know the ground trust labels
- Make sure to combine this with classic data augmentation
- => we use the Python library **Albumentations** for that







### Data augmentation: Albumentations







### Synthetic data generation

- To generate data we must:
  - Decode PNG/JPGs and segmentation masks (RLE)
  - Cut out objects, augment individual objects
  - Combine objects into scene, on top of a background
  - Encode to PNG/JPG again and store it
- => All of this is very **CPU-heavy**

How to go faster? Threading? Cluster?













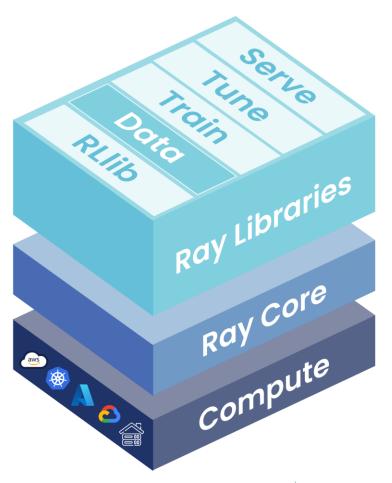
### Ray Data: Scalable Datasets for ML

• Solution: Ray Data cluster

 Handles the distribution of jobs across multiple worker nodes

- An alternative to Spark, SQL warehouses...
- But focused on Machine Learning

Ray Data => really likes Parquet files (in S3)







### Ray Data: Scalable Datasets for ML

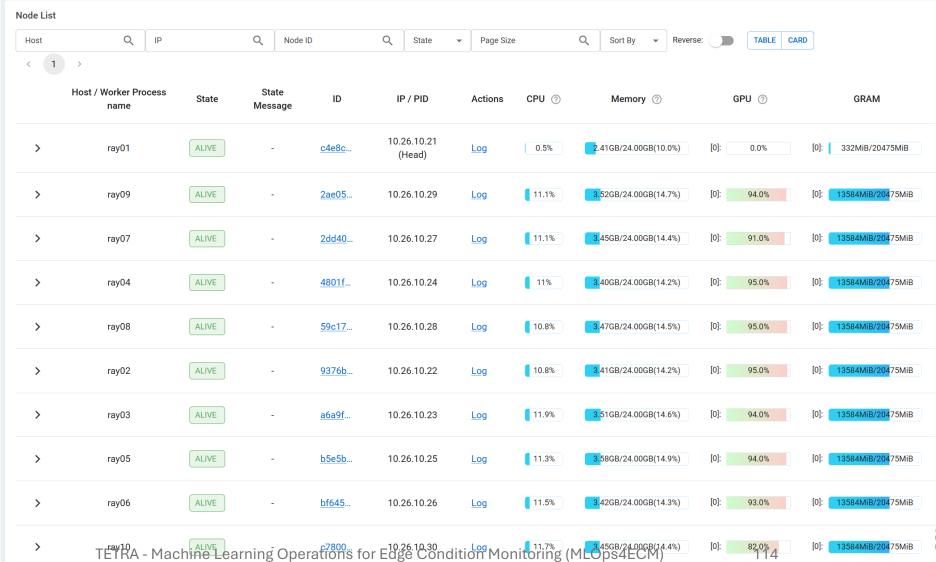
9/10/2025

```
from typing import Dict
import numpy as np
import torch
import ray
class TorchPredictor:
    def init (self):
        self.model = torch.nn.Identity()
        self.model.eval()
    def call (self, batch: Dict[str, np.ndarray]) -> Dict[str, np.ndarray]:
        inputs = torch.as tensor(batch["data"], dtype=torch.float32)
        with torch.inference_mode():
            batch["output"] = self.model(inputs).detach().numpy()
        return batch
ds = (
    ray.data.from_numpy(np.ones((32, 100)))
    .map batches(TorchPredictor, concurrency=2)
```





### Ray Data: Scalable Datasets for ML







### Where to store data? S3!

- Local disk => quite limited
- Network file system => possible, but not ideal

- Object storage in cloud => max scale, lowest cost
  - Amazon S3 => modern standard for object storage
- Object storage in cluster => local data, faster access
  - Sync to AWS, Azure, GCP ... for long term storage
- We use Ray Data + Minio S3 for data handling







### Fish Segmentation Dataset

We decided to apply these ideas to instance segmentation of **fish**, to detect **overlap** (singulation).

- Steps to take:
  - Use an existing fish dataset with masks (but without occlusions).
  - Generate a synthetic and augmented dataset using Ray Data
  - Train an edge Al model on that data







# Original dataset => little variation









# Compose into complex images





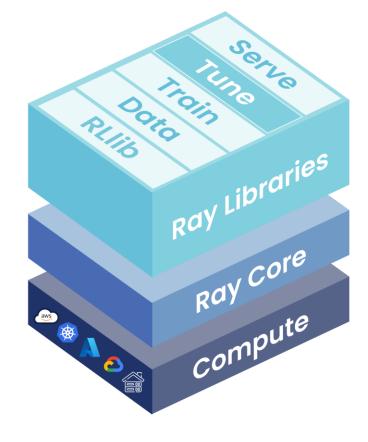


### Model training: Ray Train/Tune

How can we best train ML models?

- On a single laptop, workstation, server?
- What if we want to try many configurations?

- Solution:
- Ray Train: distributed model training (1 model, many servers)
- Ray Tune: hyper parameter search (many models, many servers)







# Model training: for the Edge

- We are going to train models for edge deployment
- These will not be huge models => fit on single GPU
- We choose: Ray Tune hyperparameter optimization
- Ray Tune will train many models for us,
   and finds the best combination of parameters
- Hyperparameters: LR, Regularization, layers...



👃 Ubuntu X 🍶 Ubuntu X 🕂 🗡

UVEN BRUGGE

(venv) root@ray01:~/leap-use-case/03-hyper-tune# python tune\_rfdetr.py
2025-10-08 05:14:57,386 INFO worker.py:1771 -- Connecting to existing Ray cluster at address: 10.26.10.21:6379...
2025-10-08 05:14:57,394 INFO worker.py:1942 -- Connected to Ray cluster. View the dashboard at http://10.26.10.21:8265

Configuration for experiment	rfdetr_lr_wd_sweep
Search algorithm	BasicVariantGenerator
Scheduler	AsyncHyperBandScheduler
Number of trials	9

View detailed results here: ray/rfdetr\_lr\_wd\_sweep

To visualize your results with TensorBoard, run: `tensorboard --logdir /tmp/ray/session\_2025-10-08\_05-02-33\_808306\_966/artifacts/2025-10-08\_05-14-57/rfdetr\_lr\_wd\_sweep/driver\_artifacts`

Trial status: 9 PENDING

Current time: 2025-10-08 05:14:57. Total running time: 0s

Logical resource usage: 0/144 CPUs, 0/9 GPUs (0.0/9.0 accelerator\_type:RTX)

Trial name	status	lr	weight_decay
train_rfdetr_tune_bb951_00000 train_rfdetr_tune_bb951_00001 train_rfdetr_tune_bb951_00002 train_rfdetr_tune_bb951_00003 train_rfdetr_tune_bb951_00004 train_rfdetr_tune_bb951_00006	PENDING PENDING PENDING PENDING PENDING PENDING PENDING	3e-05 0.0001 0.0003 3e-05 0.0001 0.0003 3e-05	3e-05 3e-05 3e-05 0.0001 0.0001 0.0001
train_rfdetr_tune_bb951_00007 train_rfdetr_tune_bb951_00008	PENDING PENDING	0.0001 0.0003	0.0003 0.0003

Trial train\_rfdetr\_tune\_bb951\_00000 started with configuration:







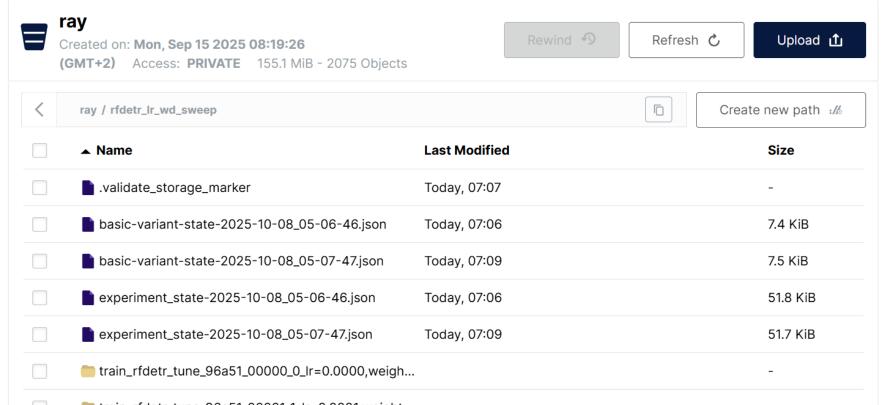
Trial name	status	lr	weight_decay
train_rfdetr_tune_abe18_00000	RUNNING	3e-05	3e-05
train_rfdetr_tune_abe18_00001	RUNNING	0.0001	3e-05
train_rfdetr_tune_abe18_00002	RUNNING	0.0003	3e-05
train_rfdetr_tune_abe18_00003	RUNNING	3e-05	0.0001
train_rfdetr_tune_abe18_00004	RUNNING	0.0001	0.0001
train_rfdetr_tune_abe18_00005	RUNNING	0.0003	0.0001
train_rfdetr_tune_abe18_00006	RUNNING	3e-05	0.0003
train_rfdetr_tune_abe18_00007	RUNNING	0.0001	0.0003
train_rfdetr_tune_abe18_00008	RUNNING	0.0003	0.0003





### Ray Tune: store artifacts in S3

- Training creates artifacts: models, metrics, predictions...
- We configure Ray Tune to store these in Minio S3

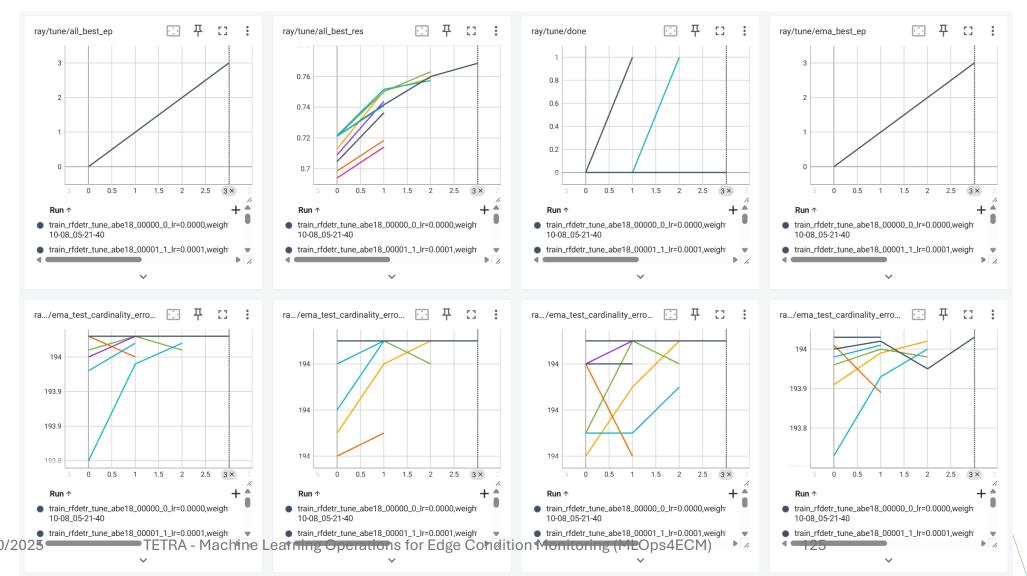






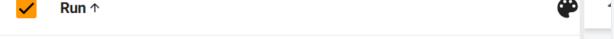


# Tensorboard: live training graphs





# Q Filter runs (regex)



- train\_rfdetr\_tune\_abe18\_00000\_0\_lr=0.0000,weight\_decay=0.0000\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00001\_1\_lr=0.0001,weight\_decay=0.0000\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00002\_2\_lr=0.0003,weight\_decay=0.0000\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00003\_3\_lr=0.0000,weight\_decay=0.0001\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00004\_4\_lr=0.0001,weight\_decay=0.0001\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00005\_5\_lr=0.0003,weight\_decay=0.0001\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00006\_6\_lr=0.0000,weight\_decay=0.0003\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00007\_7\_lr=0.0001,weight\_decay=0.0003\_2025-10-08\_05-21-40
- train\_rfdetr\_tune\_abe18\_00008\_8\_lr=0.0003,weight\_decay=0.0003\_2025-10-08\_05-21-40

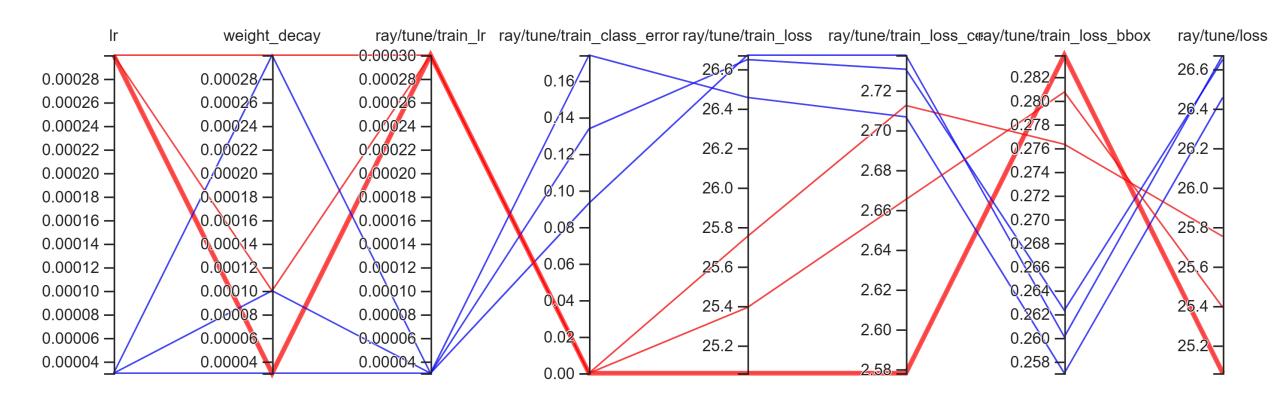
### ray/tune/all\_best\_res 0.76 0.74 0.72 0.7 0.5 1.5 2.5 0 Run ↑ train\_rfdetr\_tune\_abe18\_00000\_0\_lr=0.0000,weight 10-08\_05-21-40 train\_rfdetr\_tune\_abe18\_00001\_1\_lr=0.0001,weight

ra





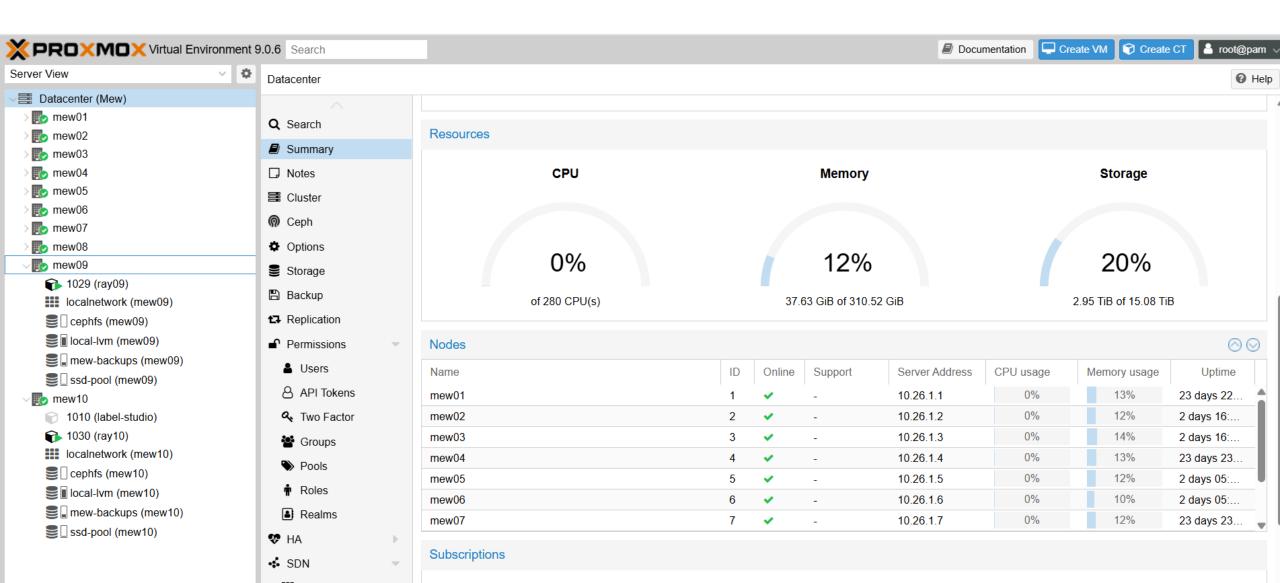
### Tensorboard + hyperparameter search







# How we manage servers? Proxmox







### Deployment: Which edge device?

- Most be quite powerful => instance segmentation is heavy
- But we want **cheap** hardware => so not a GPU server

Jetson? Linux ARM board?

- => Radxa Rock 5B, with Rockchip SoC
- Has an NPU for neural acceleration





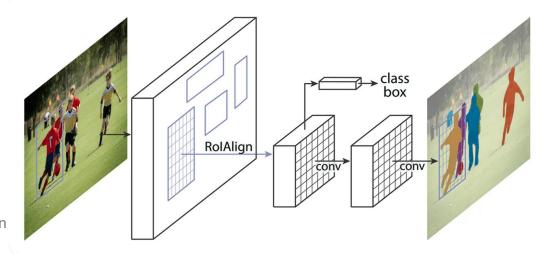


### First model: Mask-RCNN

- Torchvision Mask-RCNN (Apache)
- Pytorch version: 8 seconds per frame
- ONNX Runtime (CPU): Just as bad
- OpenVINO for ARM: much faster! 5s per frame
- OpenVINO pinned to big CPU core: 3s per frame
- Backbone (ResNet-50) on NPU:
   2.5 seconds per frame
- Results: CPU+NPU still quite slow











### Second model: RF-DETR

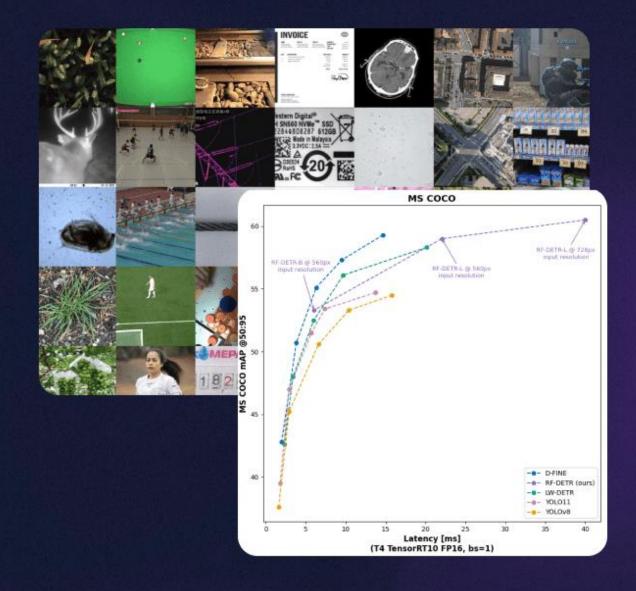
- Vision transformer (ViT) model
  - Two-stage model
  - DINOv2 backbone

 Faster and more accurate than Ultralytics YOLO11

Permissive license: Apache



RF-DETR:
A New SOTA,
Real-Time Object
Detection Model



roboflow





### **RF-DETR: Vision Transformer**

- Pure Pytorch model: 1300 ms per frame
- ONNX Runtime (CPU): 1250 ms per frame

- Does not run:
  - OpenVINO, ORT ACL, ORC XNNPACK, TVM, Executorch
- Transformers don't have great support on edge







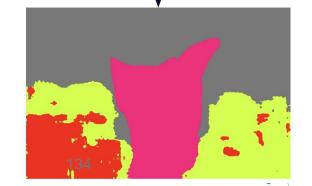
### RF-DETR backbone: DINOv2

- State-of-the-art vision encoder
- Trained using self-supervised learning

- ONNX Runtime (CPU): 950 ms per frame
- OpenVINO (CPU): 472 ms per frame
- Rockchip NPU: 367 ms per frame

• Surprise: runs on the NPU!









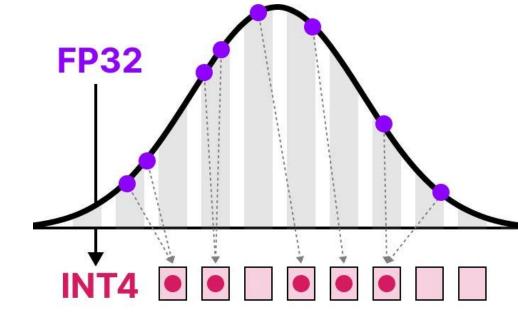


### RF-DETR: Pushing the limit

- Backbone NPU + Detection ONNX Runtime
  - => 360 ms + 350 ms = 710 ms per frame

- Lower resolution: 432x432 -> 384x384
  - => 260 ms + 260 ms = 520 ms per frame

- Dynamic quantization of detection model on ONNX Runtime
  - => 260 ms + 210 ms = 470 ms per frame







### Advantages of using the NPU

Rockchip NPU gives us a nice speed-up (higher FPS)

- NPU also has lower power usage
  - Rock 5B goes from 14W to 10W for higher FPS
  - NPU is more power efficient than the CPU when running the same neural network
- Now CPU is free to do other tasks

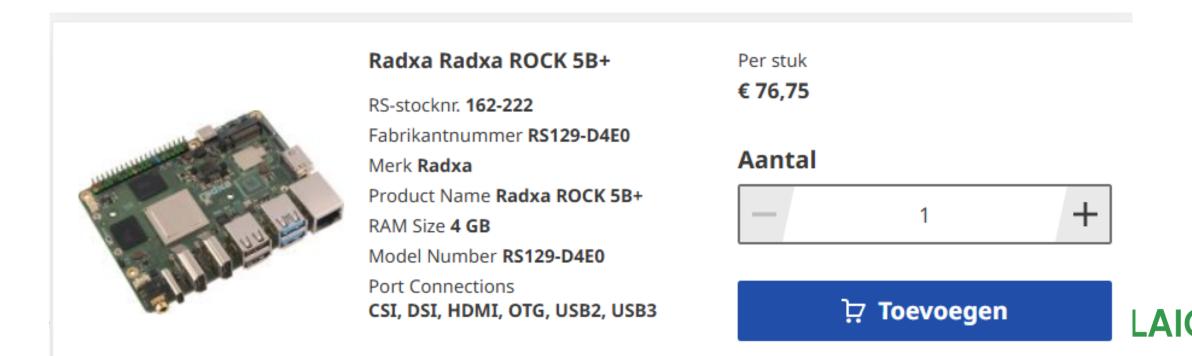






### Conclusion: SOTA on ARM board

- We have deployed a state-of-the-art model
- This runs at 2 FPS on a < 100 euro board</li>
- RF-DETR detection model runs at 5 FPS



### Demo: two models

- Top: detection @ 5 FPS
- Bottom: segmentation @ 2 FPS

- Bottom trained on synthetic data
  - Generated with Ray Data
  - Transformers are very data-hungry

- Bottom trained using Ray Tune
  - Transformers have many parameters



### Questions?

- Top: detection @ 5 FPS
- Bottom: segmentation @ 2 FPS

- Bottom trained on synthetic data
  - Generated with Ray Data
  - Transformers are very data-hungry

- Bottom trained using Ray Tune
  - Transformers have many parameters



# Follow-up projects





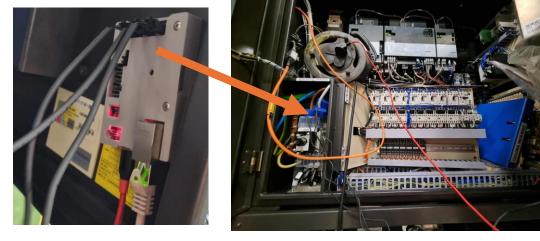


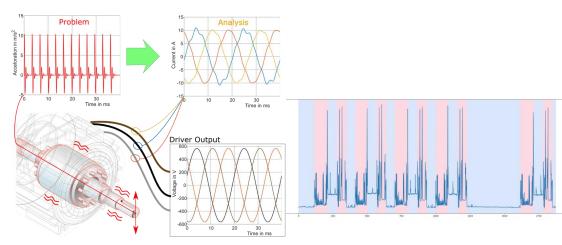
### C3 Retrokit

### **Optimising retrofits with RetroKit**

- Edge-based condition monitoring
- Modular machine learning approach
- Component- and system-level analysis

→ Call to action: Machine available?









### **TETRA Al4Automation**

### **Unlocking Intelligence for Smart Factory Automation**

- Goal: Supporting manufacturing companies in implementing Al within production environments through hands-on use cases, demonstrations, and guidance.
- Exploring the latest industrial AI tools (Siemens, Beckhoff, Advantech, Microsoft, NVIDIA) to identify practical, high-impactions.
- Call to action: Specific challenge? Let's talk!

### **Sectors:**

Food Processing - Polymer Processing - Machine building







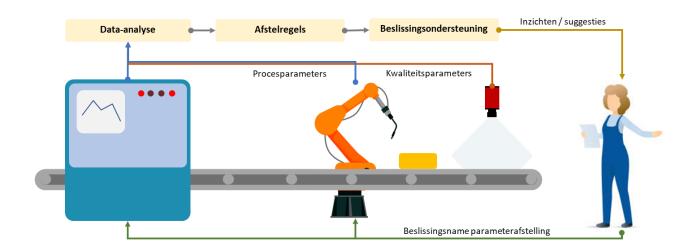
### CORNET-COOCK+ DeciQual





### Data-driven decision support for quality control and process optimization

- Determining the influence of process parameters on quality parameters
- Communicating insights to the operator
- Supporting the operator towards (optimal) process control and parameter tuning



### **Sectors:**

Food Processing – Polymer Processing – Machine building



# Do you recognize similar challenges in your own company context?

Get in touch via mathias.verbeke@kuleuven.be jonas.lannoo@vives.be







### Thanks to









































### Thanks to

### **Onderzoekers:**

Lara Luys

Alexander D'hoore

Sille Van Landschoot

Noah Debaere

(KU Leuven Brugge)

(Hogeschool VIVES Brugge)

(Hogeschool VIVES Brugge)

(Hogeschool VIVES Brugge)

### **Promotoren:**

Mathias Verbeke

Jonas Lannoo

(KU Leuven Brugge)

(Hogeschool VIVES Brugge)

# VLAIO TETRA Machine Learning Operations for Edge Condition Monitoring (MLOps4ECM)

**Reception with demos** 

