



WSAI Montréal 2023

Deploying AI models at scale

Building production-grade AI / ML products.

Victor Pereboom - CTO @ UbiOps

Tanja Crijns - Software Engineer @ Gradyent

The screenshot displays the UbiOps web interface. On the left is a sidebar with navigation options: Project dashboard, Deployments, Pipelines (highlighted), Storage, Request schedules, Imports & Exports, Logging, Audit events, Monitoring, and Permissions. Below this is a 'My subscription' section showing credit usage (0 / 7500 Units) and a reset date (14-01-2023). The main content area shows the details for a pipeline named 'v1 default'. It includes metadata such as 'Created' (21-11-2022 15:05), 'Edited' (22-11-2022 13:24), and 'Endpoint URL' (https://api.staging.ubiops.com/v2.1/proje...). It also lists 'Request retention mode' (Full) and 'Request retention time' (1 month). A 'Notification groups' section shows 'Failed requests' and 'Finished requests' with 'No group yet'. A 'Description' field is empty. At the bottom, a visual workflow diagram shows a 'video-splitter' node connected to two parallel processing nodes, which then merge before a final 'frame-processor' node.

UbiOps

anouk-staging > test-project-anouk > Pipelines > video-processing > v1 > Details

v1 *default* EDIT LOGS DELETE

General Requests Metrics

Created 21-11-2022 15:05

Edited 22-11-2022 13:24

Last request -

Endpoint URL [?] https://api.staging.ubiops.com/v2.1/proje...

Request retention mode [?] Full (Metadata + request in- and output)

Request retention time [?] 1 month (2419200)

Notification groups

Failed requests	No group yet	EDIT
Finished requests	No group yet	EDIT

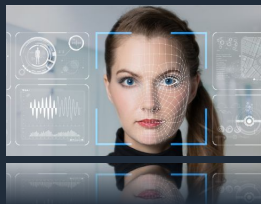
Description

No description available

video-splitter

frame-processor

／ There is a huge growth in AI services. Many businesses want to turn **ML & AI models into** products.

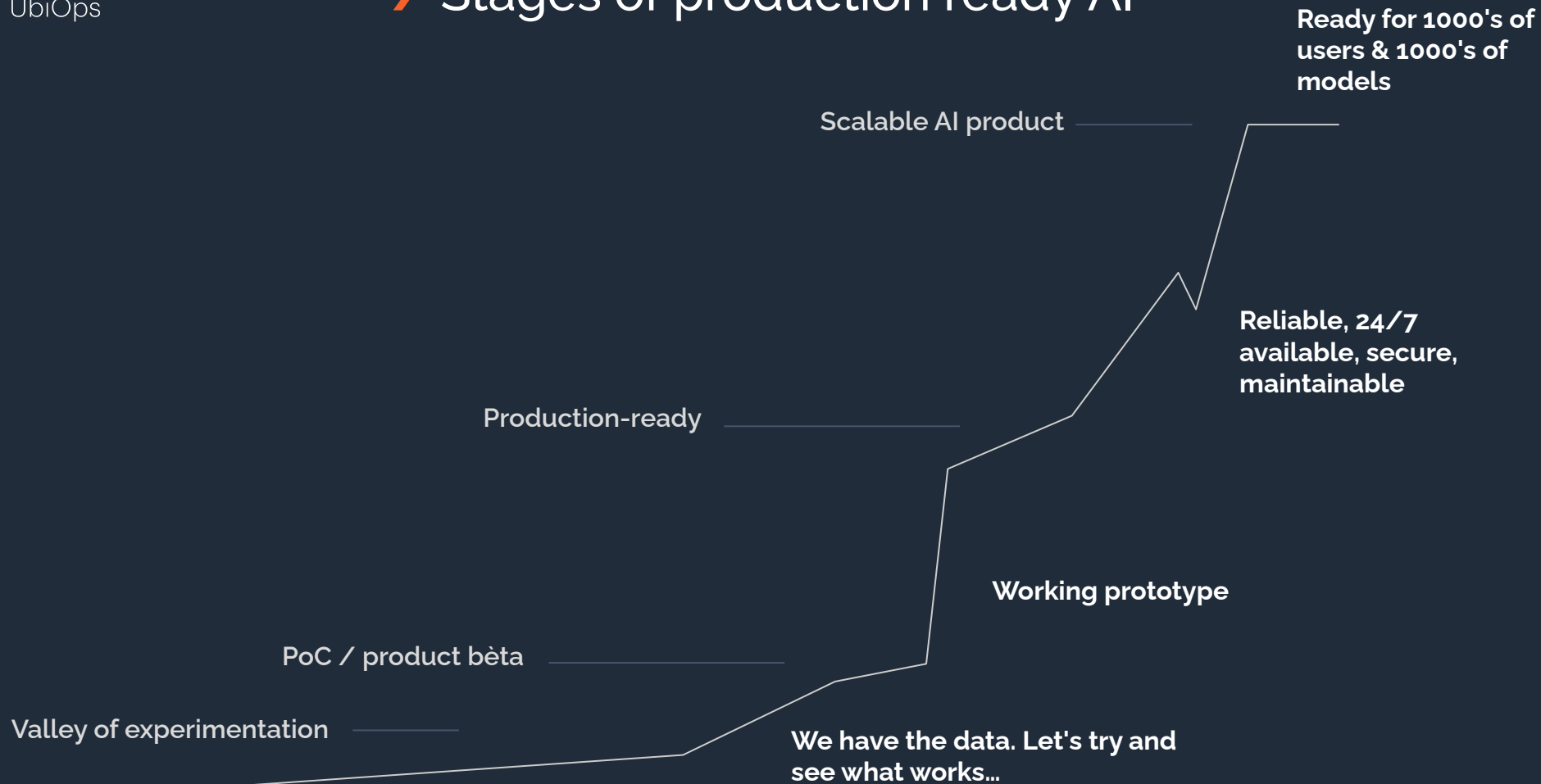


/ Productionizing AI & ML: The next big thing?

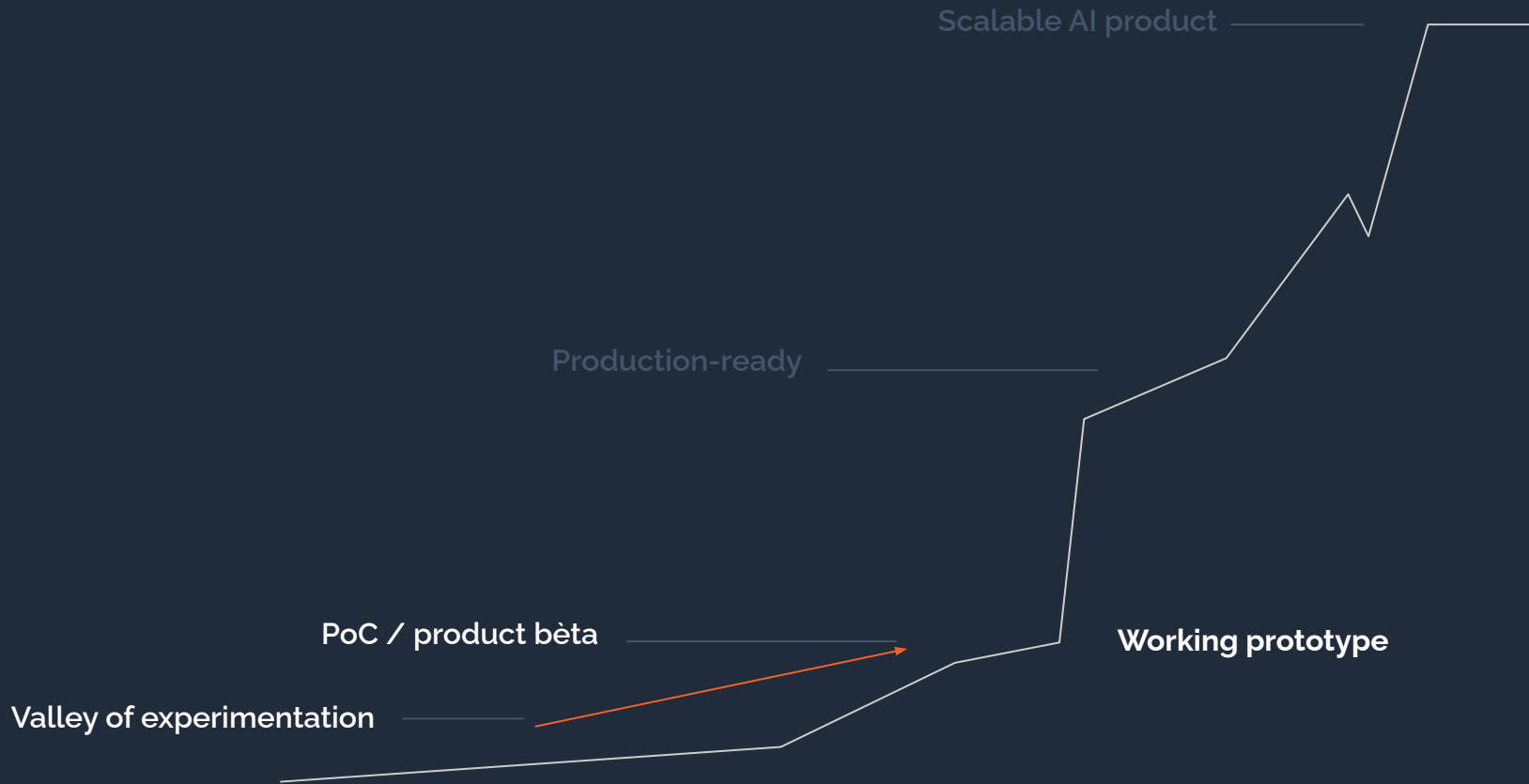
✓ **Inference** is becoming more important than **training**.

We're entering a new era of AI applications

/ Stages of production ready AI



/ Stages of production ready AI

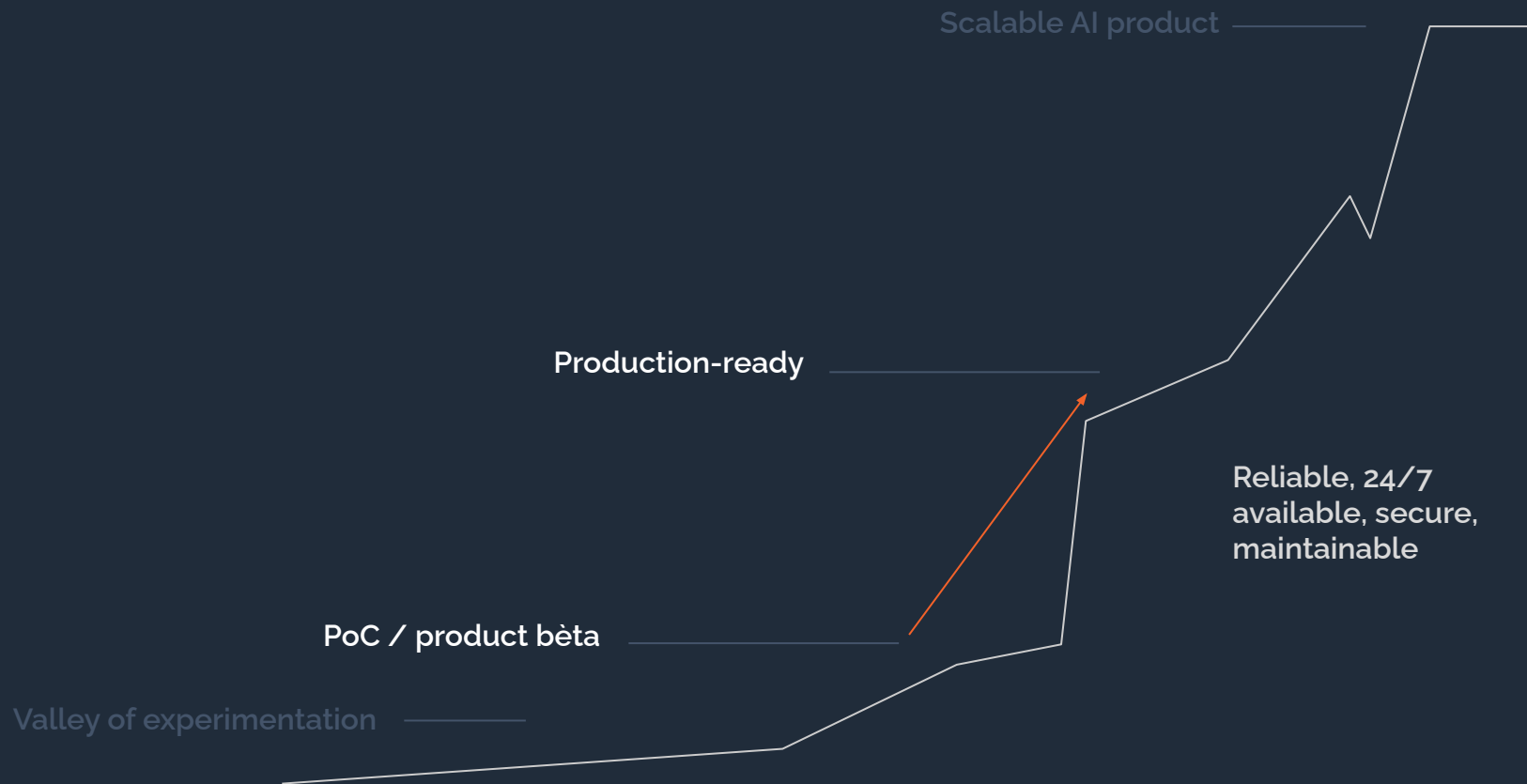


✓ **Foundation models** will change how we approach AI system development. But there are also **thousands** of business cases that need a more traditional data science approach.

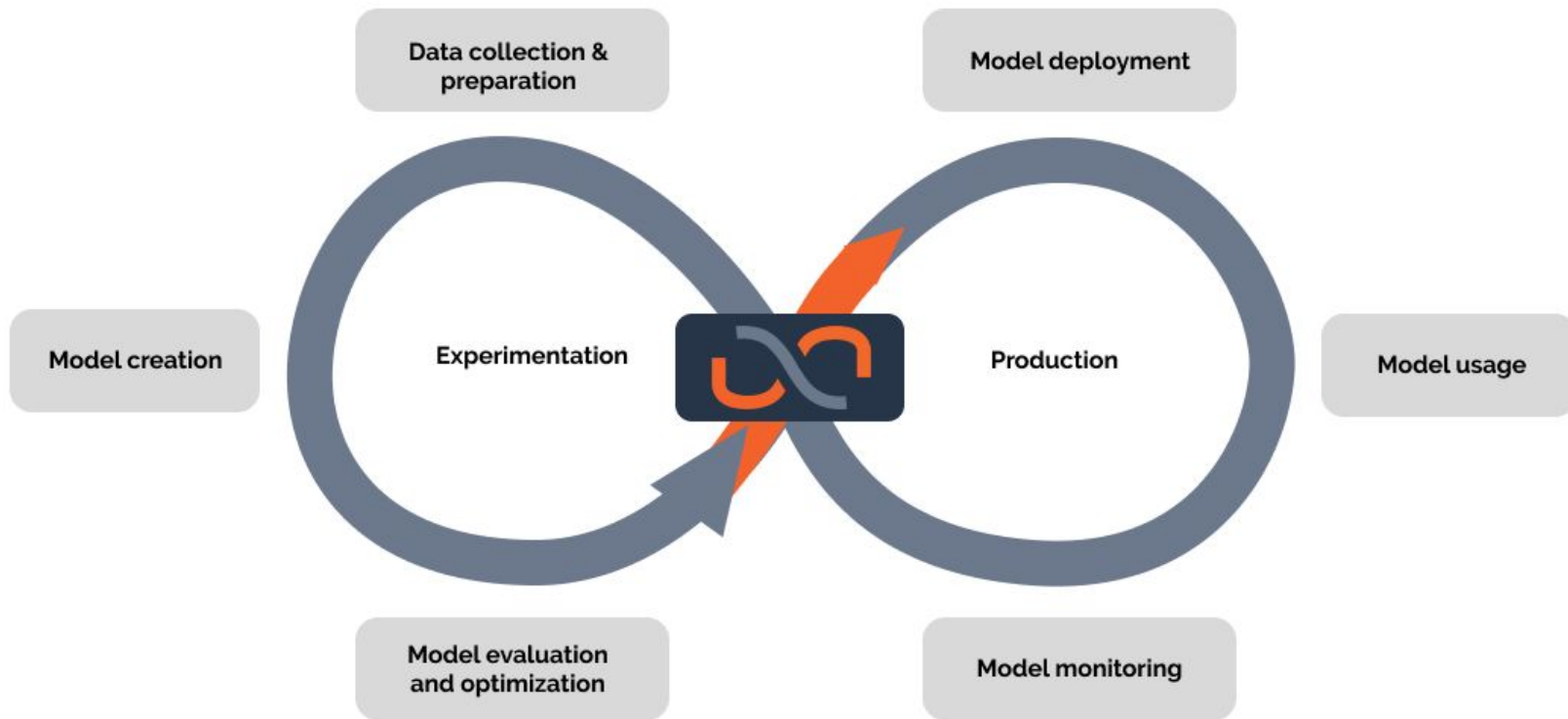
"Not everyone is Google"

A large, solid orange geometric shape, resembling a stylized arrow or a corner, located in the bottom right corner of the slide.

/ From PoC to a model in production



/ What does **productionization** even mean?

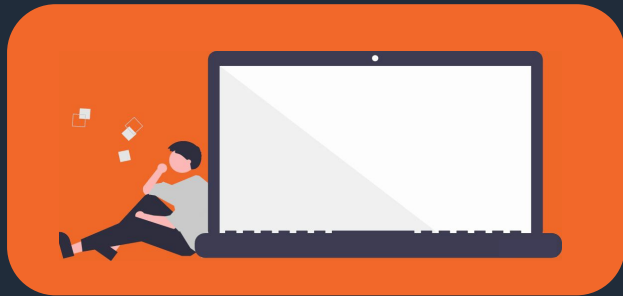


/ In the end, with AI in production, teams suddenly need to care about **reliability**, **uptime** and **security** of their solutions

/ What does **productionization** even mean?

New, non AI / ML related things to deal with

- Reproducibility (*"It worked on my machine this morning ..."*)
- Maintainability (*"Was it 'churn_v5_thisone.py' or 'churn_v5_final.py'...?"*)
- Robustness (*"Let's hope they don't kill our spot instance during the demo..."*)
- Security & compliance (*"Do we have someone named Cody in our team??"*)

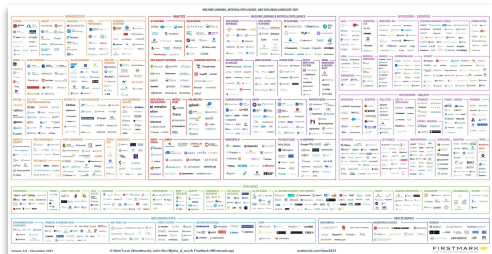


/ The **AI model** stays the same, the complexity of the solution explodes

AI Model



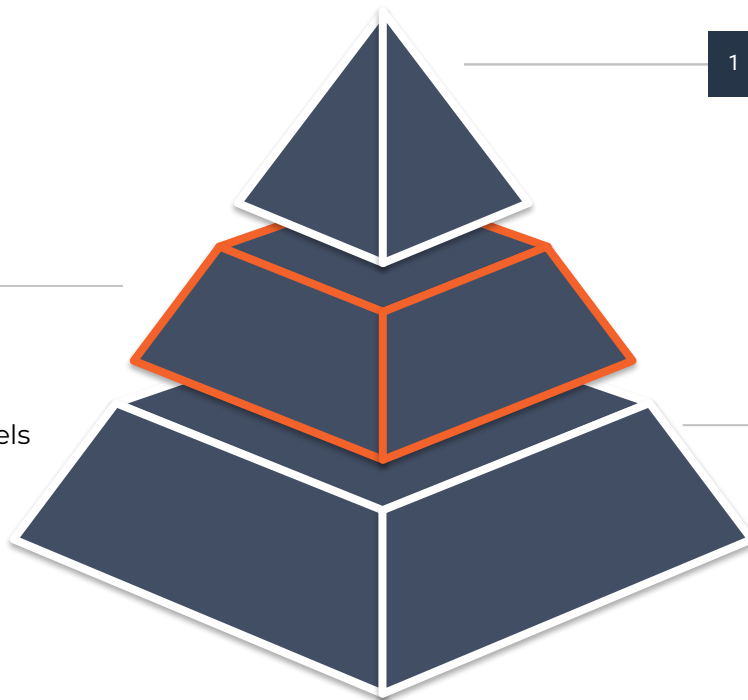
/ Solving the model deployment challenge



MLOps/ModelOps

2

Tools, technologies, and practices to deploy, monitor, and govern ML/AI algorithms and other analytical models in production-grade applications.



1

Data Science

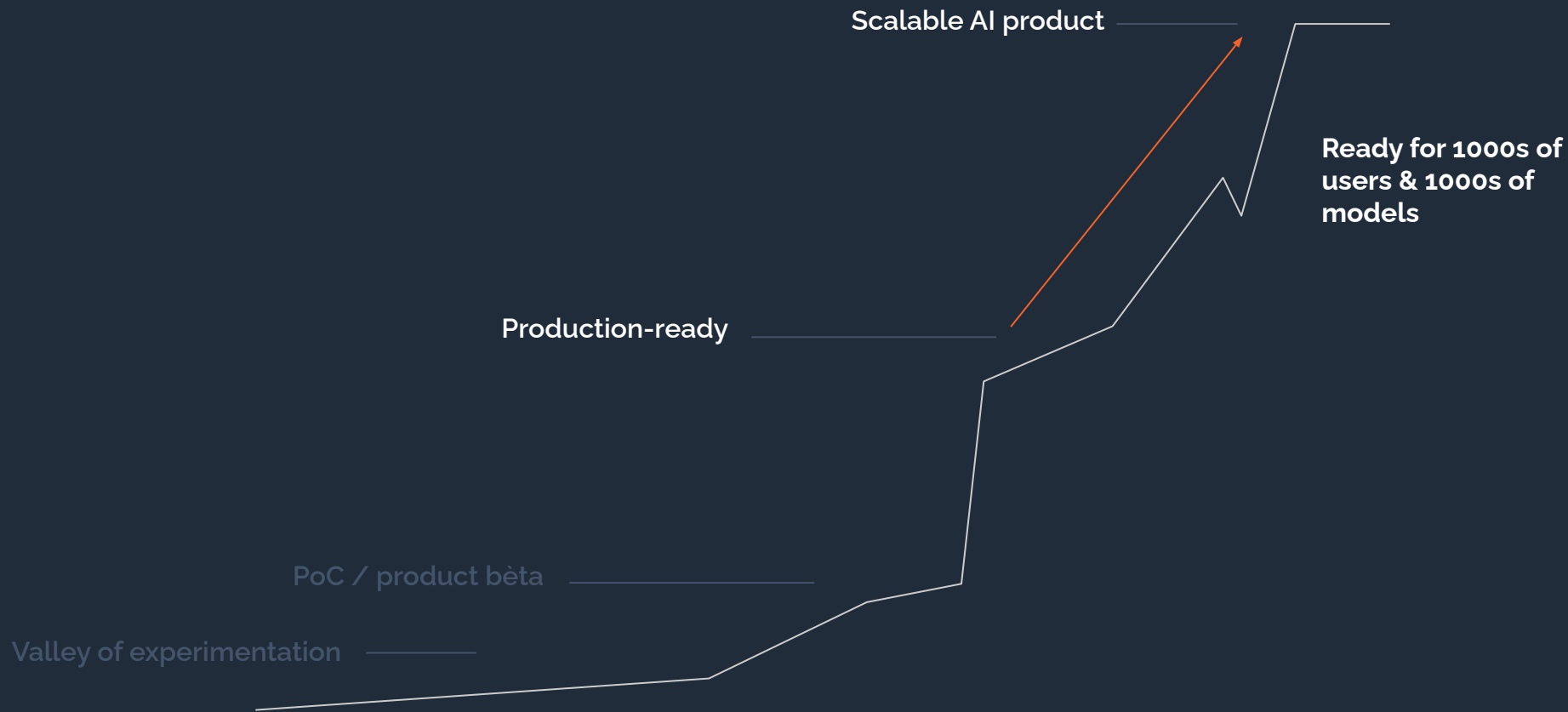


3

Infrastructure



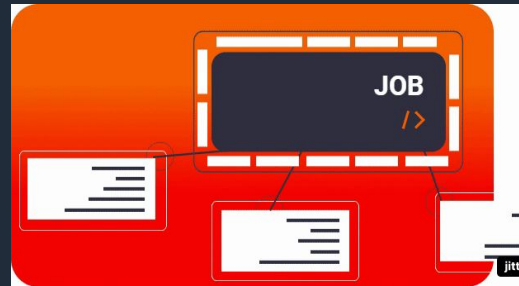
/ Scaling it up (and down)



/ Scaling up

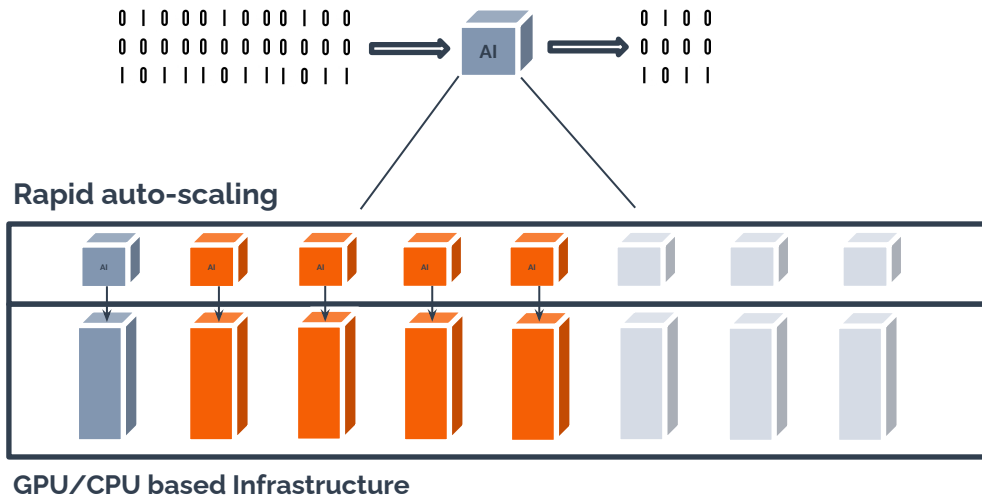
More non AI / ML related things to deal with

- Automatic scaling (*"We need to be ready for all those new users"*)
- Cost efficiency (*"What's wrong? Ah, he saw the cloud bill..."*)
- Resource availability (*"What do you mean, they ran out of GPUs ?!"*)
- Hybrid cloud (*"This data can't move to the cloud...sorry"*)



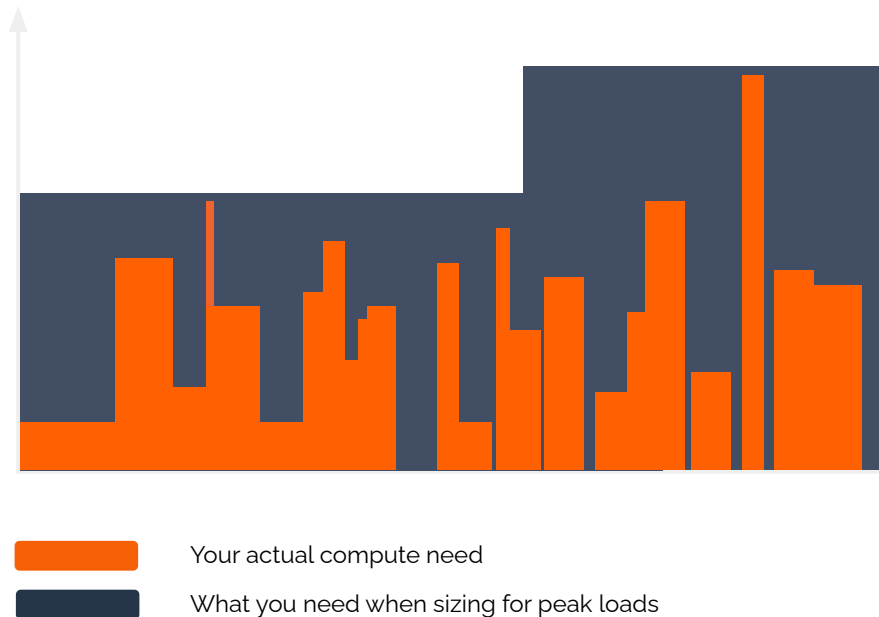
/ Scaling up...

Making sure the AI service scales with demand. Ensuring it can deal with spiky usage without the whole solution going down



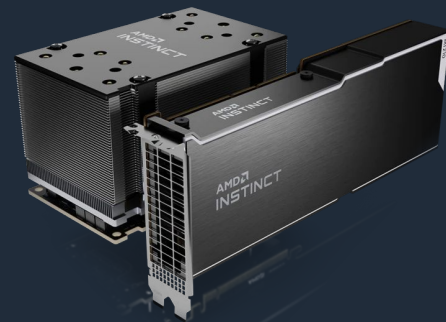
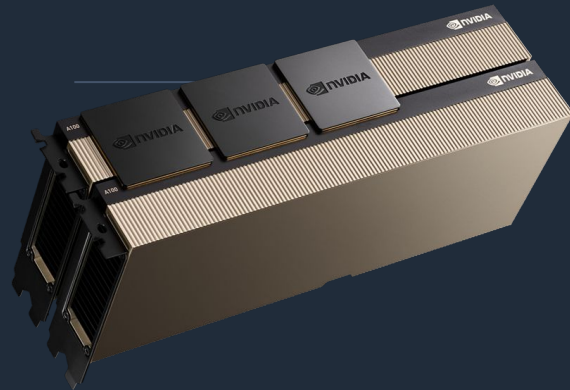
/ ...and scaling back down

Scaling up and down (back to zero)
rapidly ensures you follow your compute
demand curve as close as possible,
optimizing resource usage



Challenges with GPUs

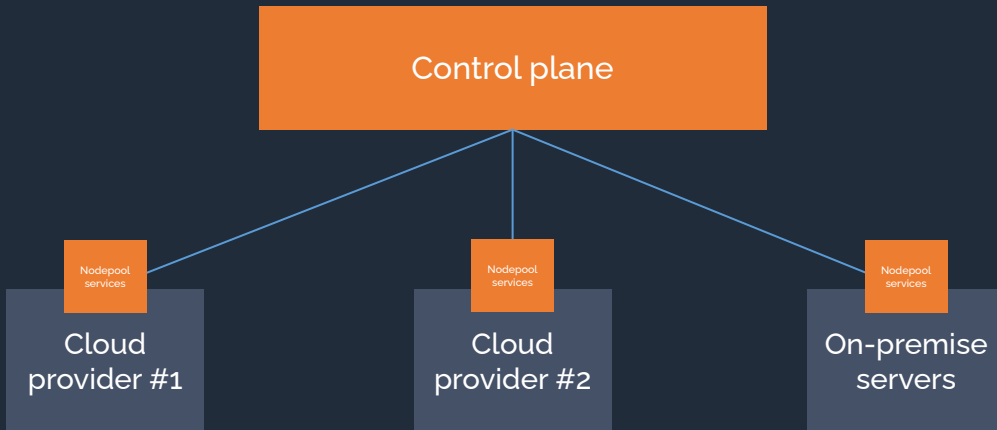
- Necessary for acceleration of AI models (and training)
- They're expensive
- Sharing a GPU between containers is not easy
- Issues with GPU availability, driven by supply chain issues and AI models that get heavier and heavier



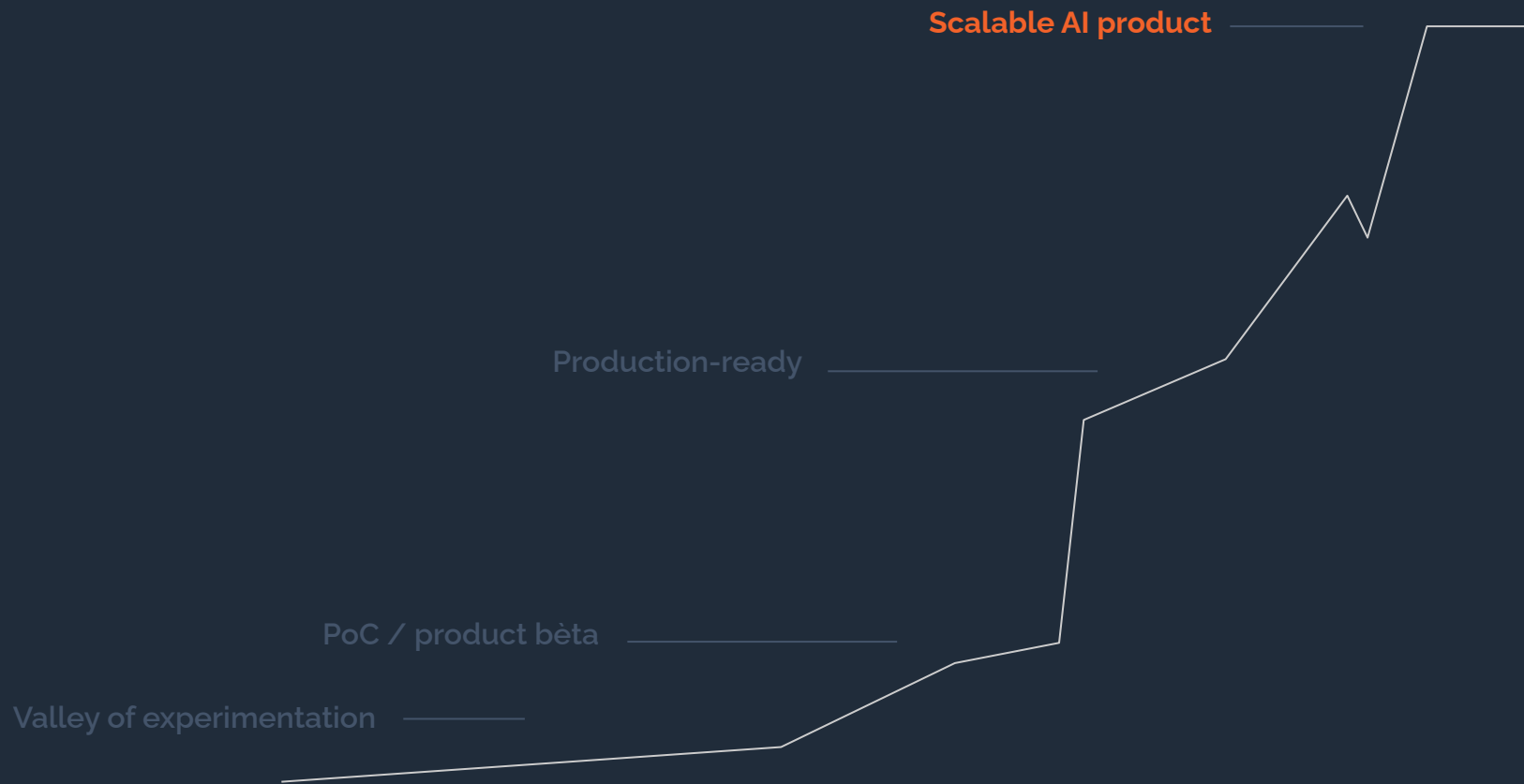
/ Run across environments: Hybrid & multi-cloud

Enabling to **lift and shift** workloads dynamically across different environments.

- Higher availability of resources (and GPUs)
- Bring compute to where the data is.
- Solve compliance and security blockers



/ Reaching the summit



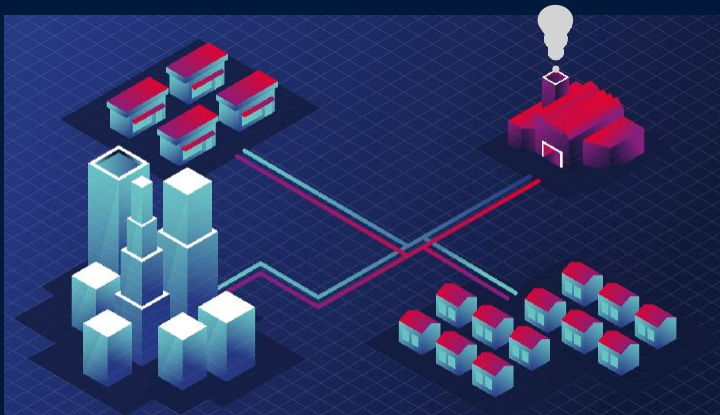


Usecase Gradyent Digital Twin for heating grids

Tanja Crijns



District and industrial heating networks need to transform



High temperature & losses
Single / few heat sources
Stand alone

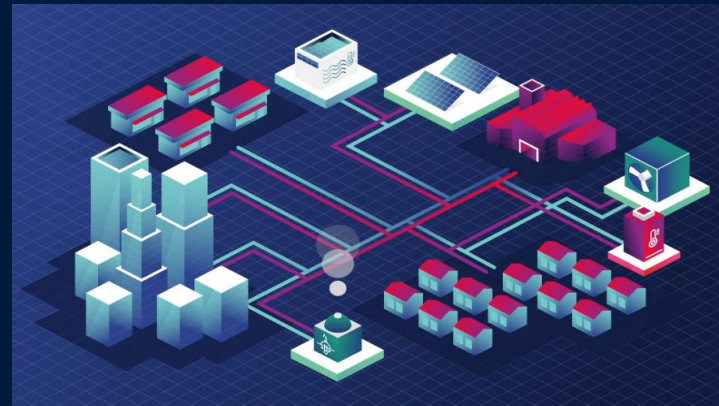
Optimize



Decarbonize



Grow



Lower temperatures
Multiple, different types of heat sources
Integrated with industry, electricity & cooling

The Gradyent **Digital Twin** brings a new level of optimization & design by covering the **entire system** in real time



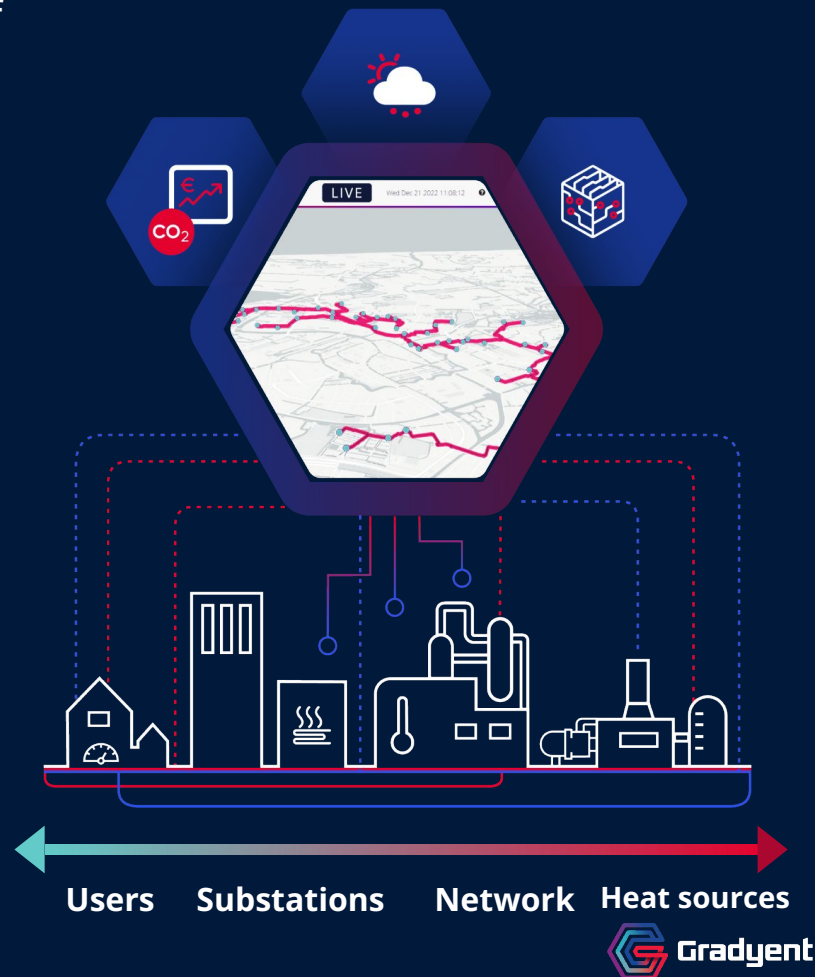
Real time operational optimization

By optimizing source scheduling, temperatures and pressure, debottlenecking, demand response, issue resolution

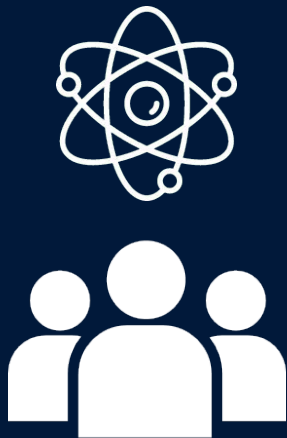


Design & simulation of future proof networks

Enhancing decision making around low carbon sources, system changes and extensions



UbiOps provides an
abstraction layer
between our engineers
and the cloud



anouk-staging > test-project-anouk > Pipelines > video-processing

v1 **default**  EDIT  LOGS  DELETE

General Requests Metrics

Created 21-11-2022 15:05

Edited 22-11-2022 13:2

Last request -

Endpoint URL  https://api.sta

Request retention mode  Full (Metadat

Request retention time  1 month (24

Not groups and requ

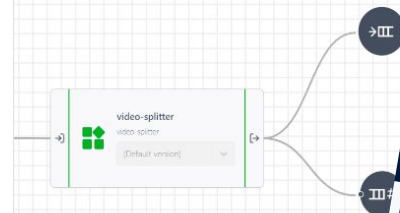
ed re

Description

No description available

Units

11-2023

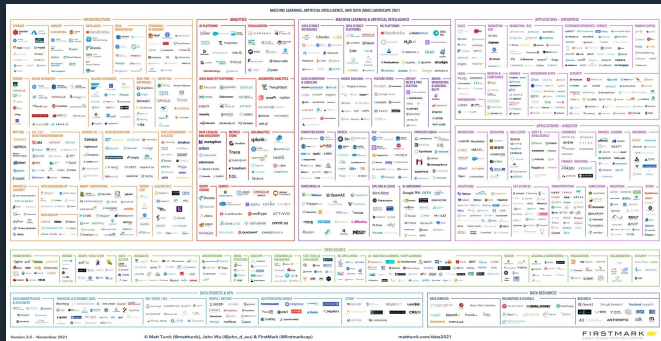




Gradyent - Five key benefits

- Ease of use
- Fast paced development
- Production grade stability
- Security and compliance
- Clear insight in compute usage

Fast track to productized AI



Scalable AI product

Production-ready

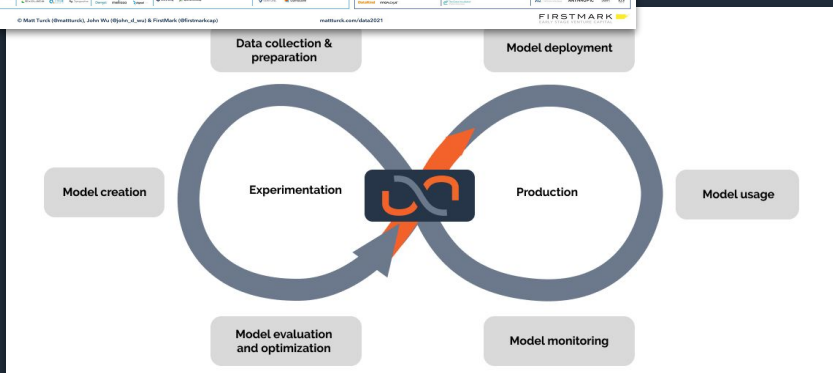
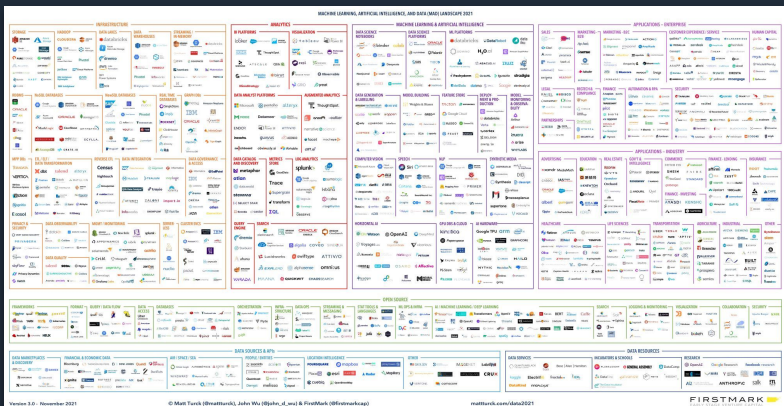
PoC / product beta

Valley of experimentation

- Short time to market
- Lower costs
- Increased reliability and maintainability of a solution

/ The future for AI products

- ❑ A **new reality** where implementation of AI has a significant impact on competitiveness, market share and profitability.
- ❑ The MLOps ecosystem is expanding rapidly
- ❑ Teams will start picking their own preferred stack of tools
- ❑ Inference will become increasingly important
- ❑ Foundation models will drive product creation, but also hardware and accelerator needs

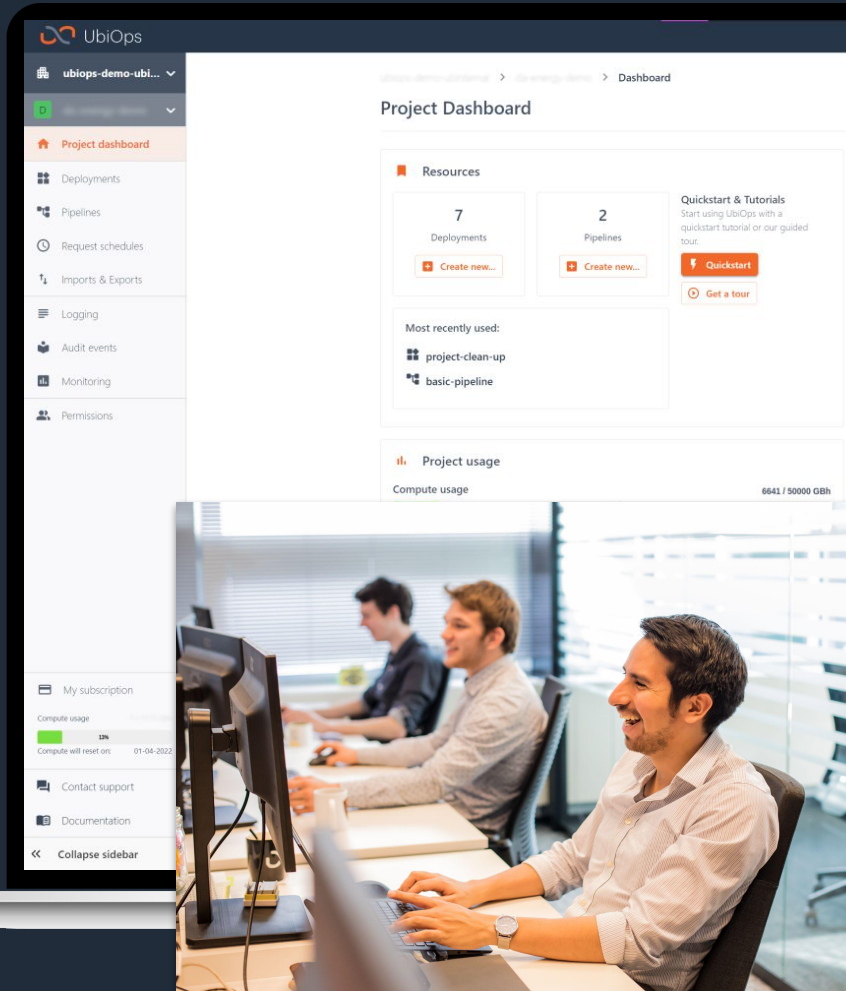




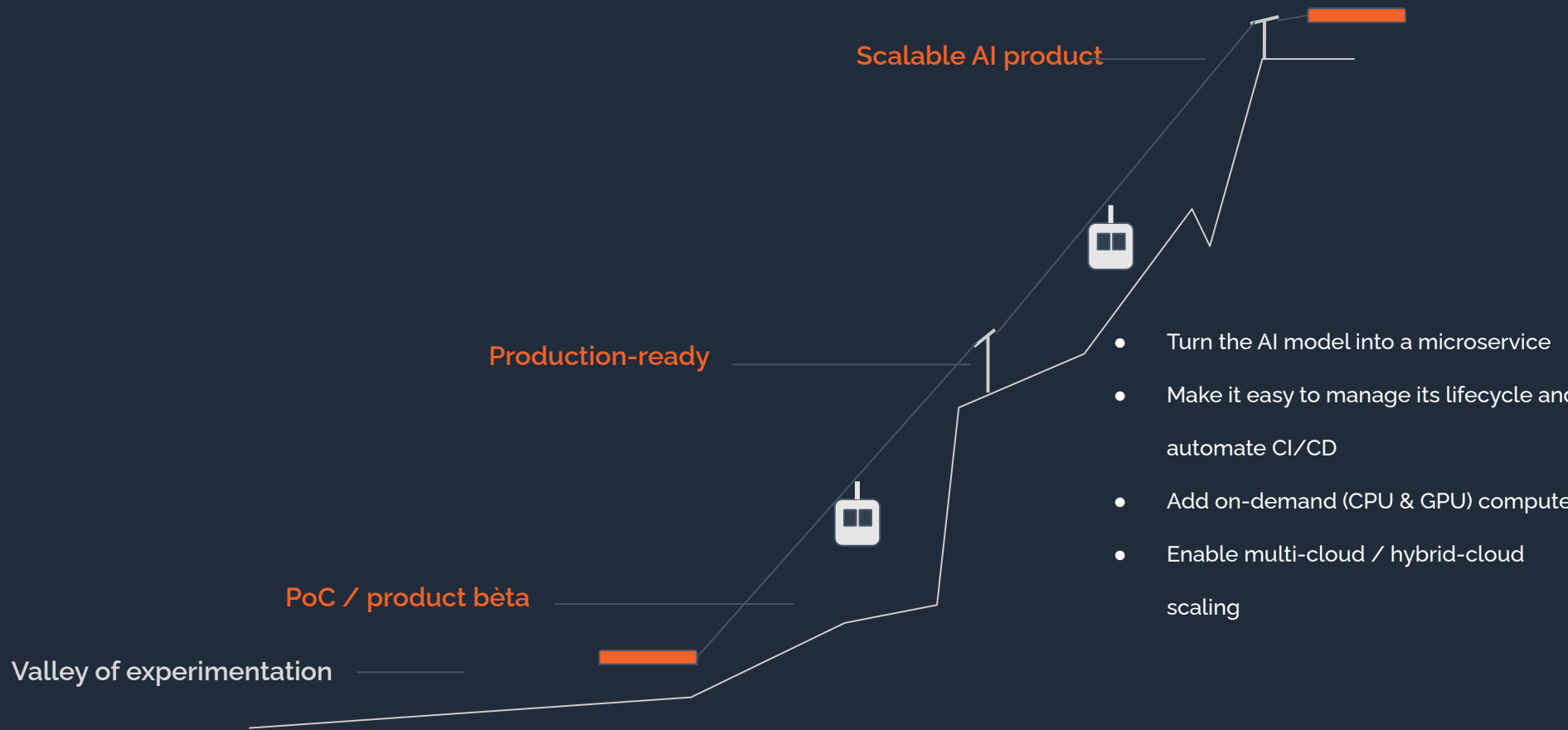
UbiOps

Build production-grade AI / ML products.

Our mission is to help teams of all sizes solve the technical challenges around building and running the next generation of AI products and services



/ Fast-track to production ready AI





WSAI Montréal 2023

Thank you!

Victor Pereboom - CTO @ UbiOps
victor@ubiops.com | visit ubiops.com

Tanja Crijns - Software Engineer @ Gradyent
tanja@gradyent.ai

Meet us at booth **A70**

