

# Course Introduction

Sprint 1 - Week 1

*INFO9023 - Machine Learning Systems Design*

Thomas Vrancken ([t.vrancken@uliege.be](mailto:t.vrancken@uliege.be))

Matthias Pirlet ([matthias.pirlet@uliege.be](mailto:matthias.pirlet@uliege.be))

2025 Spring

# Agenda

What will we talk about today.

## Lecture (2.5 hours)

1. Introduction to the staff
2. Introduction to ML Systems Designs & MLOps
3. Real-world use case
4. Development phases & challenges
5. Course organisation
6. Use case definition framework

# **Introduction to the staff**

# Introduction to the staff



Thomas Vrancken

(Instructor)

[t.vrancken@uliege.be](mailto:t.vrancken@uliege.be)



Matthias Pirlet

(Teaching assistant)

[matthias.pirlet@uliege.be](mailto:matthias.pirlet@uliege.be)

# Our experience & expertise make us leading AI specialists.

## EXCEPTIONAL TALENT & SKILLS



110+ experts spread over 3 different EU locations.



Known for technical expertise  
Loved for our business results



Talent magnet: 16 applications each day

## STATE-OF-THE-ART TECH KNOWLEDGE & ASSETS



6 Mio downloads per month of our open source packages



150+ clients, 300+ projects, 3 spin-offs  
17% of time in R&D, 250+ publications,  
15 awards, avg time to value 11.5 weeks



Security, Legal and Ethical AI experts

# We work with customers across industries and geographies.

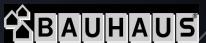
## Life Sciences & Healthcare



de volksbank

## Manufacturing & utilities

## CPG, Retail & Ecommerce



MediaMarktSaturn  
Retail Group



ML6



## Financial Services



AGENTSCHAP  
WEGEN & VERKEER

MARCH

Keypoint  
Outsourcing strategies

Fostplus

FUNKE  
MEDIEN  
GRUPPE

Booking.com

Google



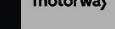
avrotros



SPRINGER NATURE

## Communication, Media & Technology

## Public & Professional Services



# ML6 - your partner in AI.

We accompany organisations through their entire AI journey

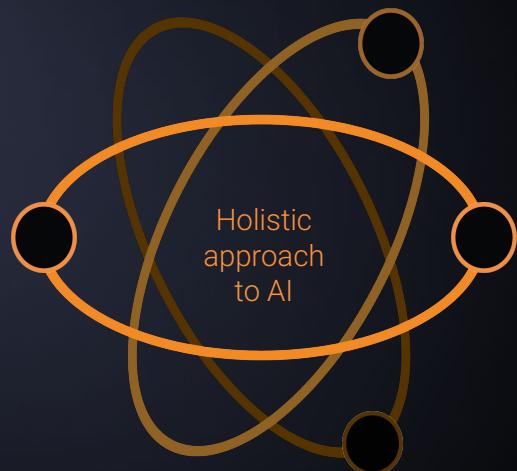
- Use case definition & assessment
- AI solution design, engineering and deployment
- Managed support and maintenance services
- Production-level scaling & evolution of AI solutions

We help remove barriers to technology adoption

- Security
- Ethics & Regulation
- Business case building
- Selecting the right tech stack
- Facilitating user adoption

We cover all AI domains

- Machine Vision
- NLP
- Structured Data
- Reinforcement Learning & Generative AI
- MLOps & Engineering best-practices



We engineer bespoke AI solutions

- Tailored to complex client needs
- Agile development & use of boiler plates where relevant
- Reliable, robust & maintainable solutions

We deliver end-to-end

- Data labelling
- Sourcing of internal and external data
- Hardware selection and/or integration (incl. IOT & edge devices)
- Front-end development

We are technology-agnostic

- Cloud agnostic: AWS, Azure, GCP
- Open source minded
- Tech radar for stack selection
- Hybrid cloud - on premise; and edge deployment

# We are recognized as leaders by the industry.

Don't just take our word for it



**1000**  
Europe's Fastest  
Growing Companies

#386 (EU) | #4 BE

**Data  
News**

Nominated in 2023,  
2022, 2021

**Trends  
GAZELLEN202**

Multiple nominations &  
one award win

**Deloitte.**  
Technology Fast 50



Nominated in 2022,  
2021, 2020, 2019,  
2018



Multiple nominations &  
award wins in 2022, 2015

**C**

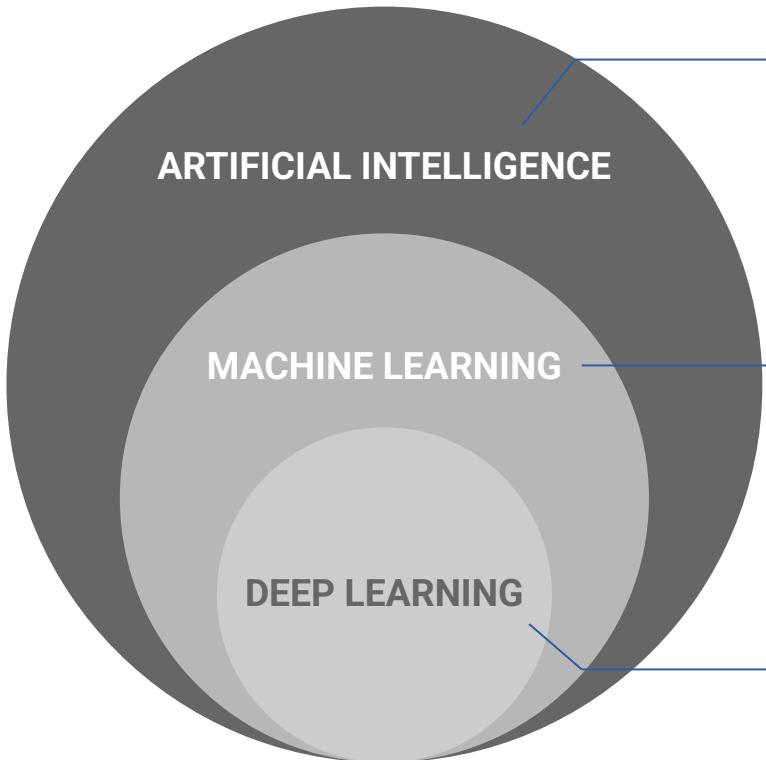
Multiple nominations &  
award win in 2023

**EY**

Scale-up of the year  
finalists in 2023

# **General introduction to ML Systems Design & MLOps**

# AI vs ML vs DL



## ARTIFICIAL INTELLIGENCE

Ability of a machine to perform cognitive functions we associate with human minds, such as perceiving, reasoning, learning, problem solving, and even creativity

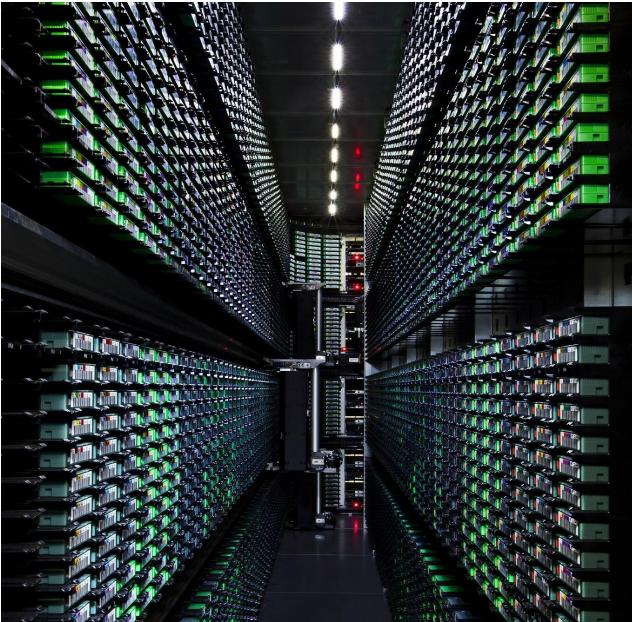
## MACHINE LEARNING

AI techniques that give machines the ability to learn from data without being explicitly programmed, i.e. to automatically improve through experience

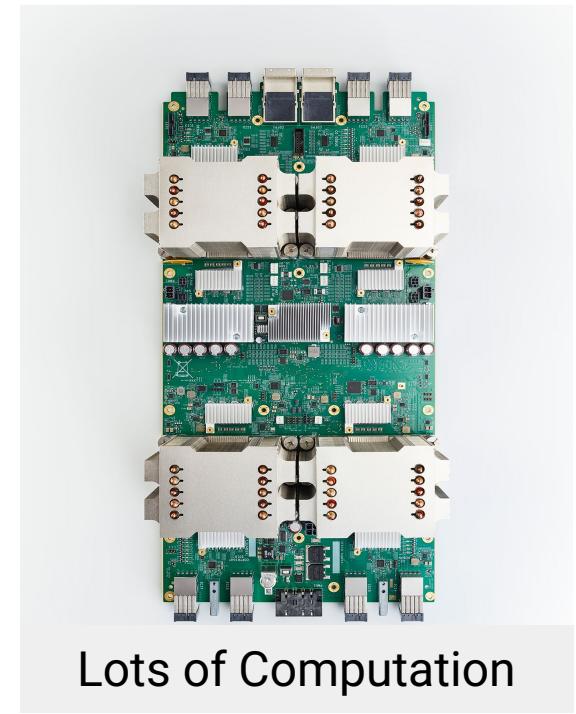
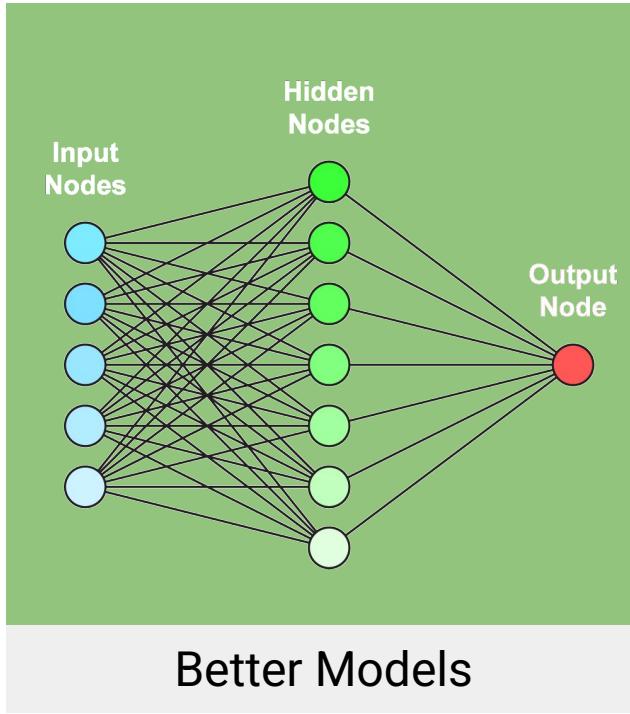
## DEEP LEARNING

Type of Machine Learning built upon the concept of interconnected layers known as “neurons” that form a neural network.

# Why now ?



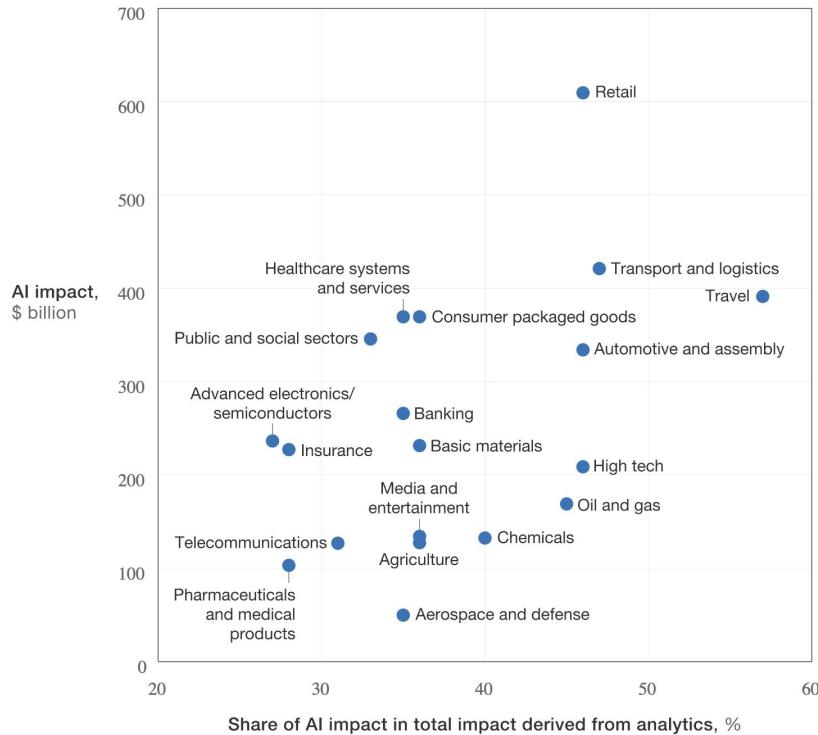
Large Datasets



Lots of Computation

# AI is bringing a revolution in many different industries

Artificial intelligence (AI) has the potential to create value across sectors.



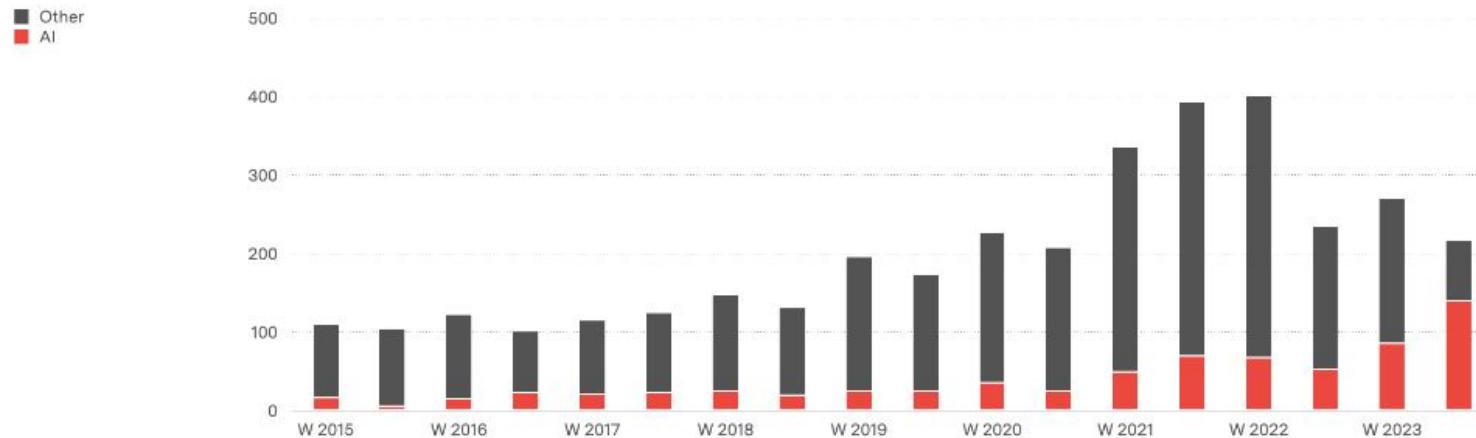
AI value creation by 2030

**13 trillion USD**

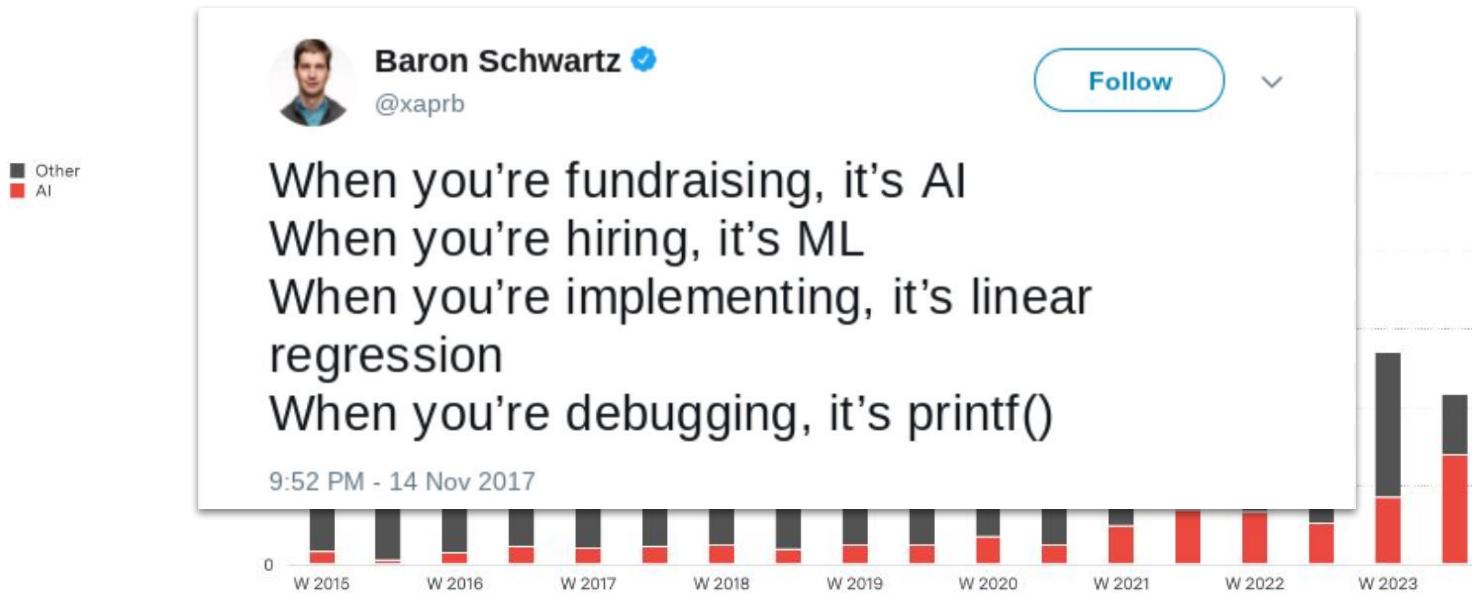
Most of it will be outside the consumer internet industry

# Investment in AI ventures is skyrocketing.

Y Combinator startups by field

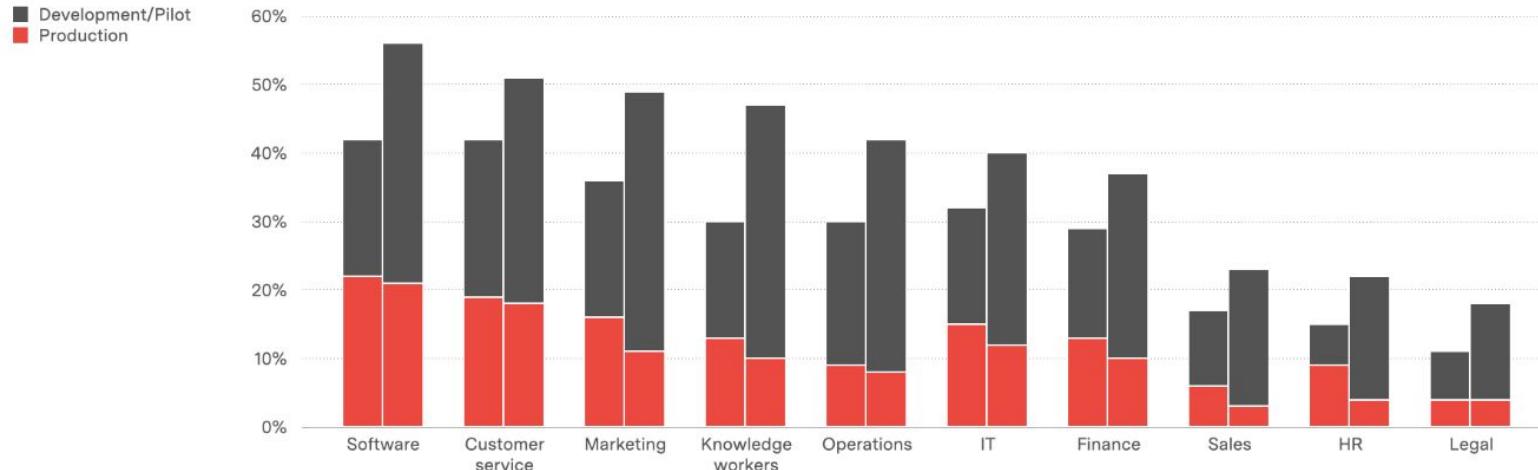


# Investment in AI ventures is skyrocketing.

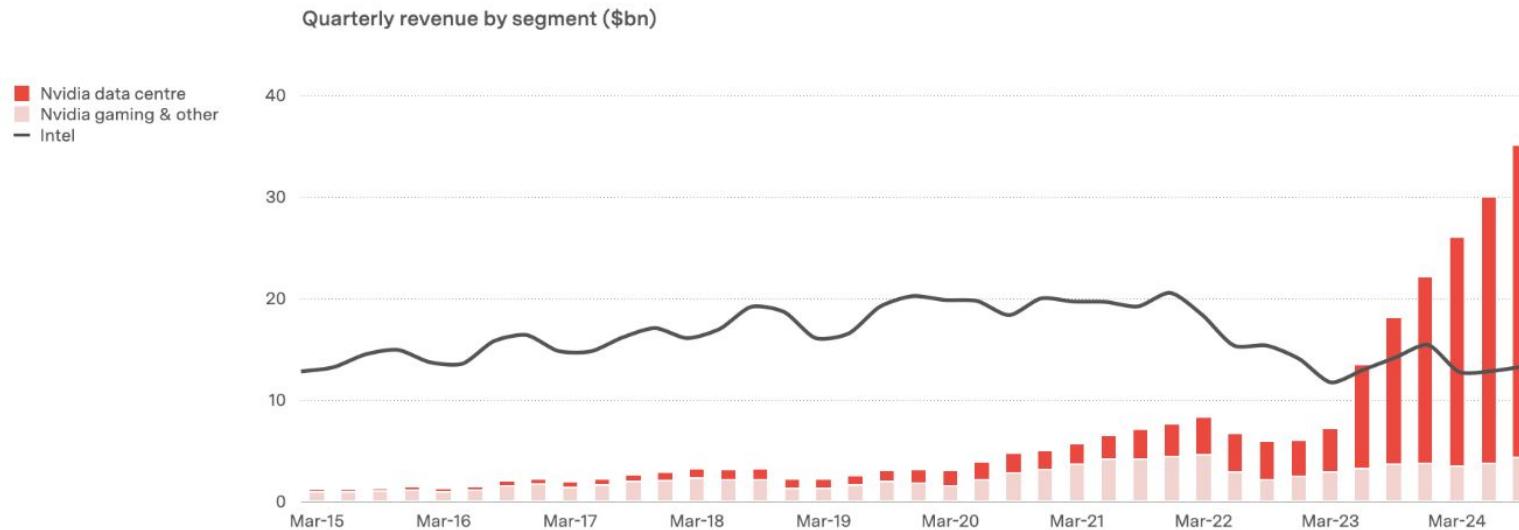


# Companies have a large interest for AI, but often struggle with getting their use cases to production.

Enterprise use case adoption rates for generative AI, October 2023 & February 2024



# Demand for AI ripples down to other segments, such as Cloud infrastructure.



Future innovations & impact with AI might not be in training better models but rather in applying them and making them more efficient.



# The “API model” might not last forever.

“Everyone in tech is giving someone else’s business model away for free”

Meta’s open source

Turn models into commodity infrastructure!



The screenshot shows a web browser window with the URL [about.fb.com](https://about.fb.com). The page title is "Meta". Below it, a post is displayed with the title "Open Source AI Is the Path Forward" and the author "By Mark Zuckerberg, Founder and CEO". The date "July 23, 2024" is also visible. The main content of the post discusses the historical development of Unix and how open source Linux became the industry standard, eventually surpassing closed source Unix versions.

In the early days of high-performance computing, the major tech companies of the day each invested heavily in developing their own closed source versions of Unix. It was hard to imagine at the time that any other approach could develop such advanced software. Eventually though, open source Linux gained popularity – initially because it allowed developers to modify its code however they wanted and was more affordable, and over time because it became more advanced, more secure, and had a broader ecosystem supporting more capabilities than any closed Unix. Today, Linux is the industry standard foundation for both cloud computing and the operating systems that run most mobile devices – and we all benefit from superior products because of it.



AI is everywhere!

# A few example of ML applications.

Facial  
recognition



Product  
recommendation



Email spam  
filtering



Autocomplete



Finance  
predictions



Healthcare  
imaging



Weather  
forecast



...

# Going beyond the notebook

You were tasked with implementing a model to predict electricity productions of a wind turbine farm.

What you did:

- **Export** key data for the last 10 years (electricity production, weather, ...)
- **Analyse** the data and build **features**
- **Train** and **optimise** a ML models
- Make and visualise **predictions**

# Going beyond the notebook

You were tasked with implementing a model to predict electricity productions of a wind turbine farm.

Now what?... How will you:

- Automatically do predictions every X minutes with new data
- Automatically retrain the model with new data
- Collaborate with a larger team
- Continuously integrate new changes to your application
- Not run your model locally on your own laptop
- Optimise your model so it does not consume too many resources
- Monitor your models
- ...



# Why do we need ML Systems Design?

Building a ML application means implementing much more than just your ML model.

INFO 9023 -  
Machine  
Learning Systems  
Design

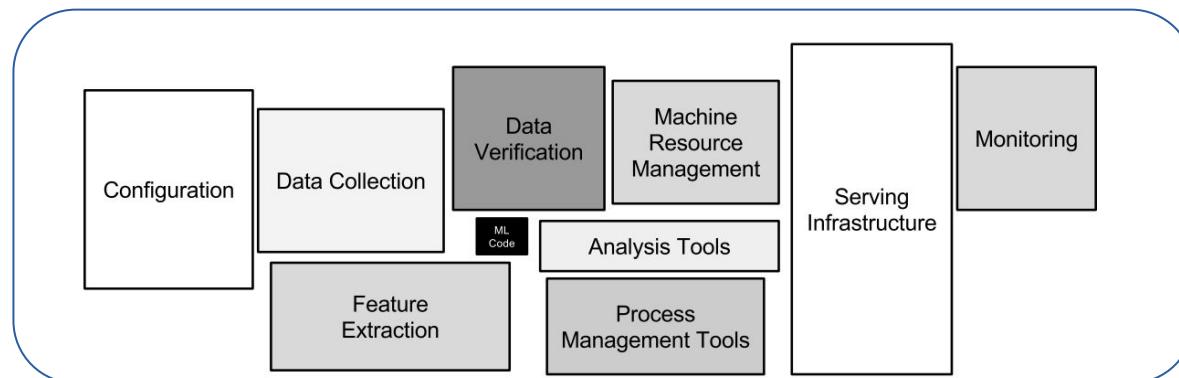


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

Sculley, D. et al. (2015). Hidden technical debt in machine learning systems.

[https://papers.nips.cc/paper\\_files/paper/2015/hash/86df7dcfd896fcfa2674f757a2463eba-Abstract.html](https://papers.nips.cc/paper_files/paper/2015/hash/86df7dcfd896fcfa2674f757a2463eba-Abstract.html)

# Important definitions

**ML Application:** The final solution or program powered by a Machine Learning model.

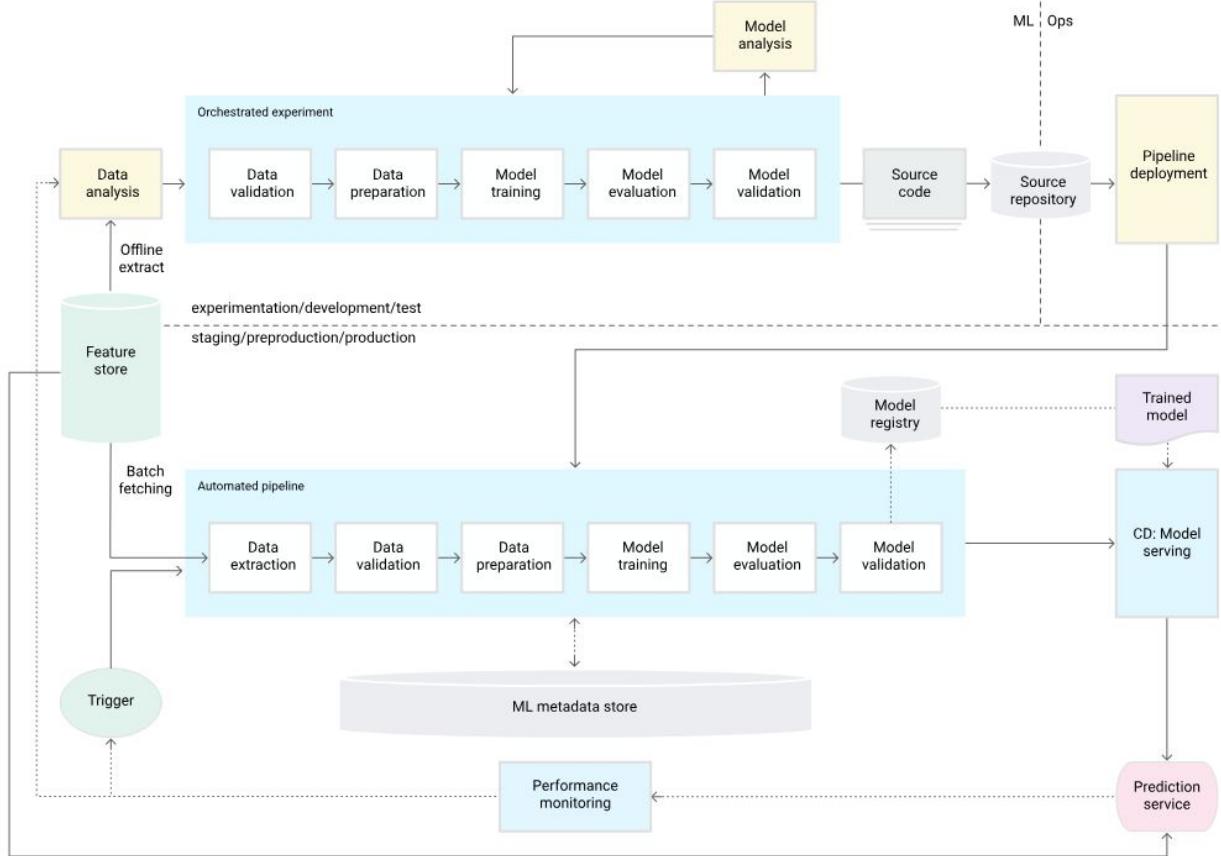
**ML System:** All the components responsible for the implementation and management of the data and models powering an ML application.

**ML Systems Design:** The act of designing the architecture and implementing an ML System.

**MLOps:** Set of practices that aim at implementing and maintaining ML systems in production reliably and efficiently.

# **Key concepts of ML Systems Design**

# Typical architecture of an ML system



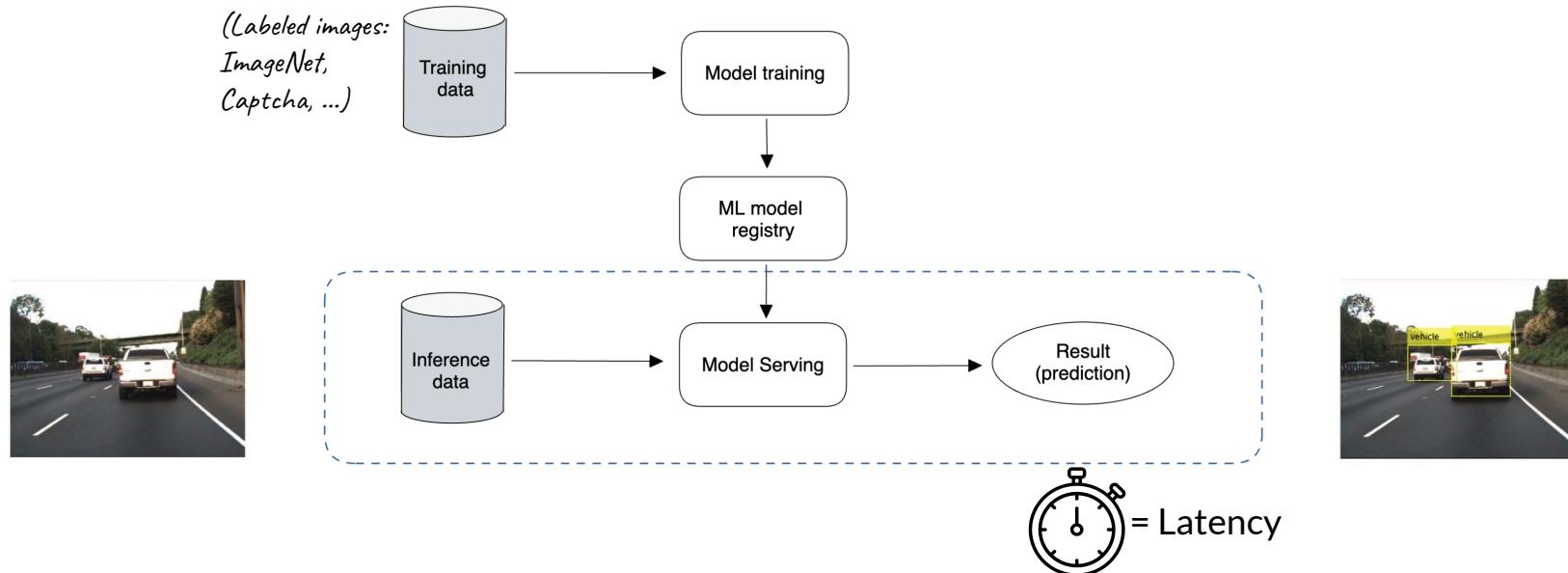
# Key concept: Data preparation

It all starts with data. How to go through all these steps efficiently and effectively.



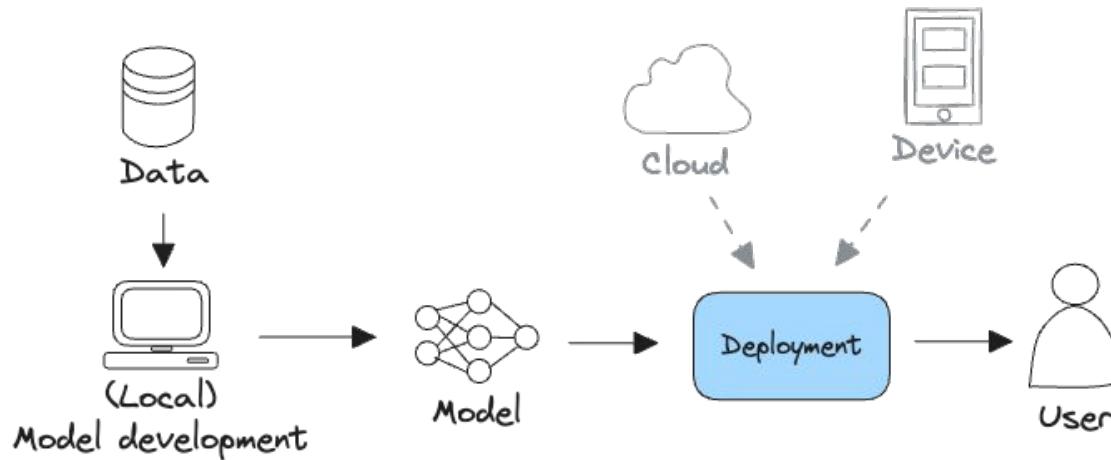
# Key concept: ML model serving

How to efficiently serve ML model to client.



# Key concept: ML model deployment

How to efficiently deploy your model for serving.



# Key concept: Containerisation

**Containers** encapsulate an application as a **single executable package** that contains all the information to **run it on any hardware**:

- Application code
- configuration files
- libraries
- dependencies

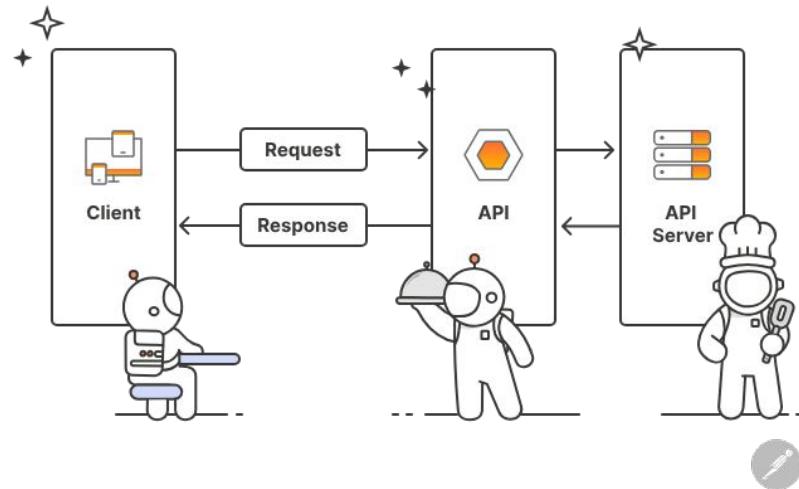
**Abstracts** the application from its **host operating system**.

Containers can be easily transported from a desktop computer to a virtual machine (VM) or from a Linux to a Windows operating system, and they will run consistently on virtualized infrastructures or on traditional “bare metal” servers, either on-premise or in the cloud.



# Key concept: APIs

Allow other services to call your model or application.



An **Application Programming Interface (API)** is a set of protocols that enable different software components to communicate and transfer data.

Developers use APIs to bridge the gaps between small, discrete chunks of code in order to create applications that are powerful, resilient, secure, and able to meet user needs.

# Key concept: Cloud infrastructure

Cloud infrastructure allow for data storage, compute allocation, training and deploying model, monitoring, ...

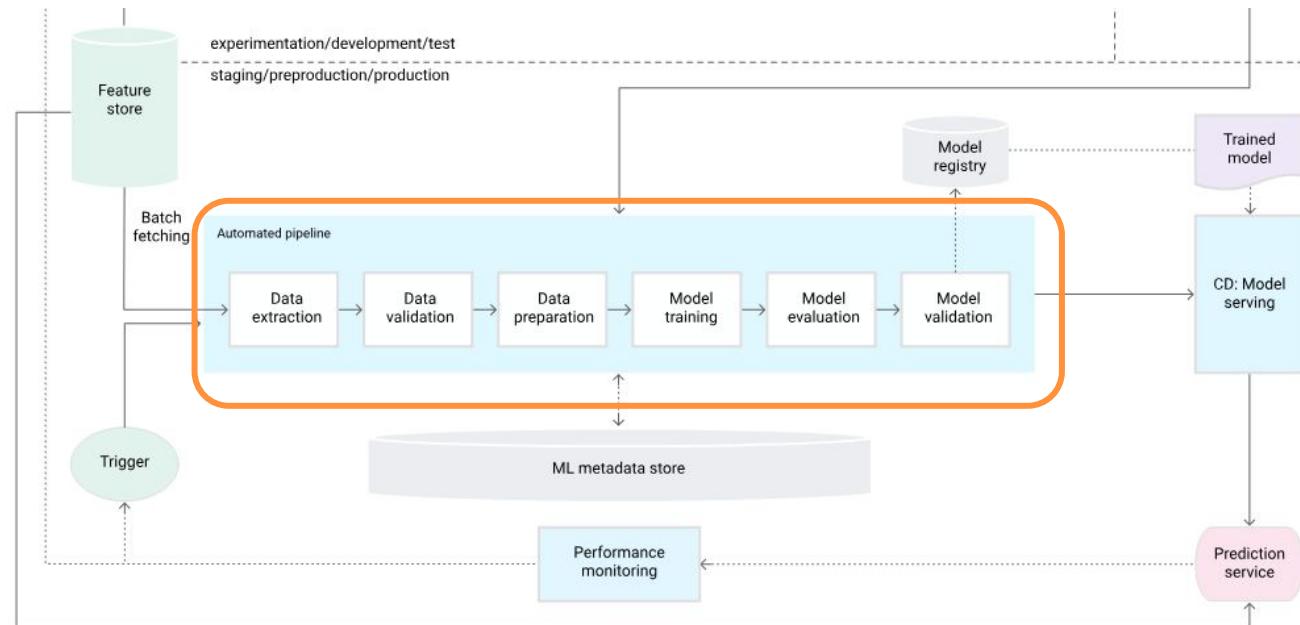


Google Cloud



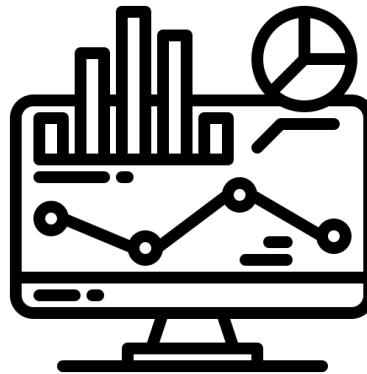
# Key concept: ML Pipeline

Orchestrates components to prepare data, train, evaluate and deploy ML models  
(among other things)



# Key concept: Monitoring

Ensuring that models in production are performing well.



**Resource level** (performance and usage of resources used by the model serving)

- How much is it being used by users?
- Are the CPU, RAM, network usage, and disk space as expected?
- What are the Cloud costs?
- Are requests being processed at the expected rate?
- What is the system uptime? Some maintenance contract depend on it.

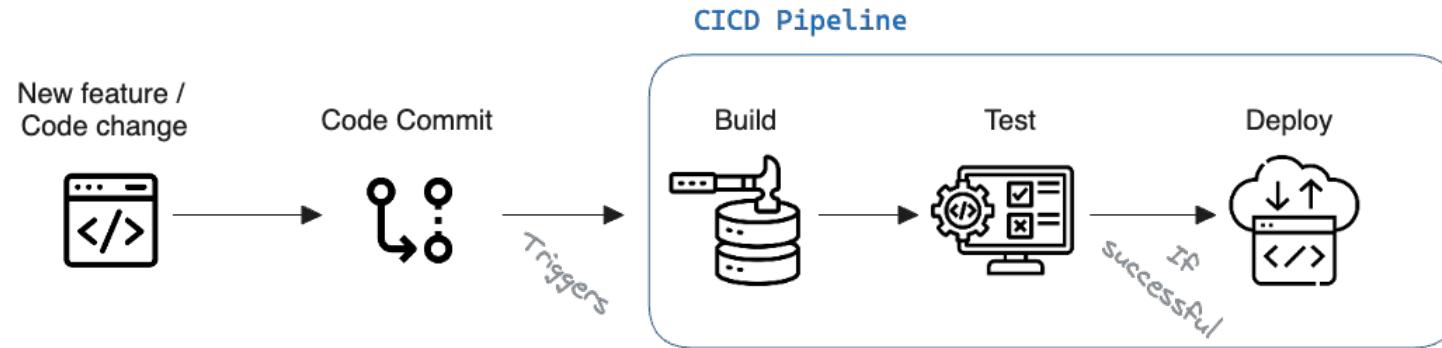
**Performance level** (performance/accuracy of the model over time)

- Is the model still doing accurate predictions with the new data coming in?
- Is the data distribution changing?
- Is the target variable changing?
- Are concepts around the model changing?

# Key concept: CICD

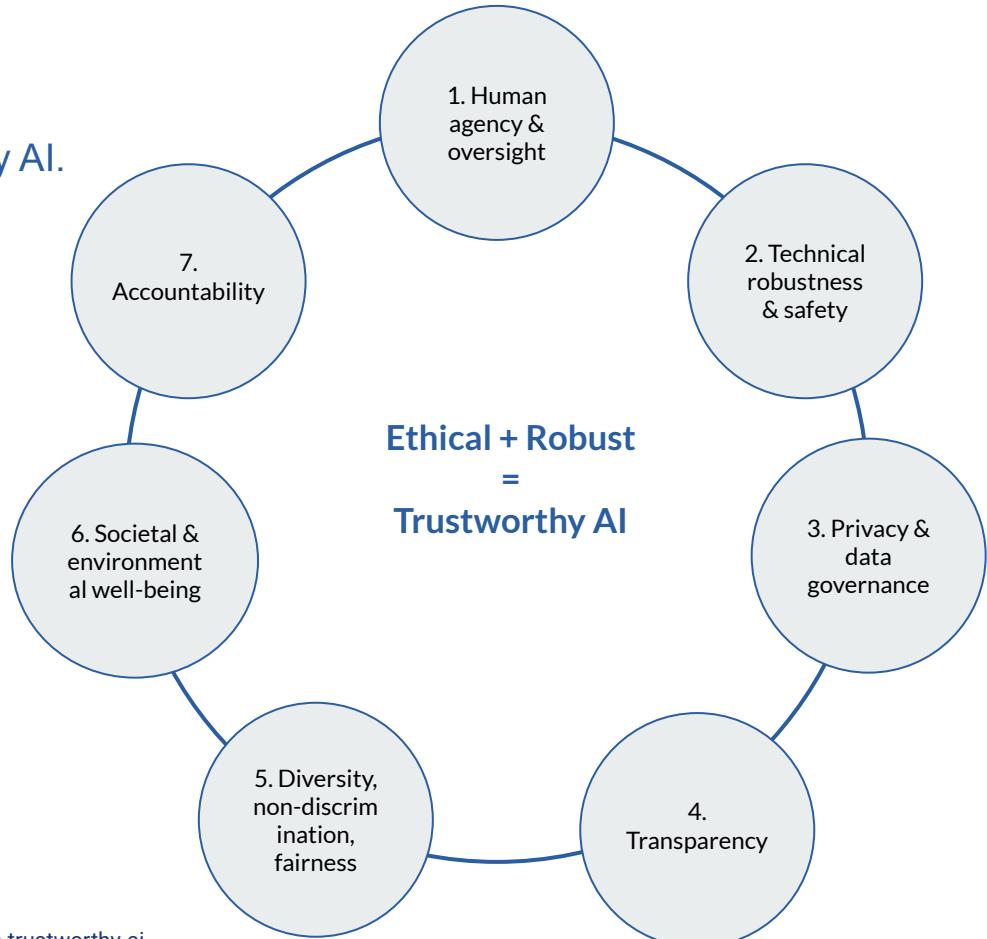
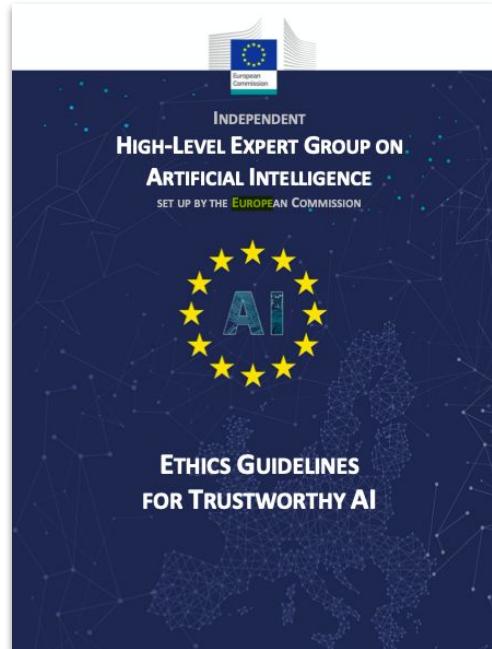
Allows you to continuously work on your application and efficiently deploy new changes to it.

Continuous Integration and Continuous Delivery.



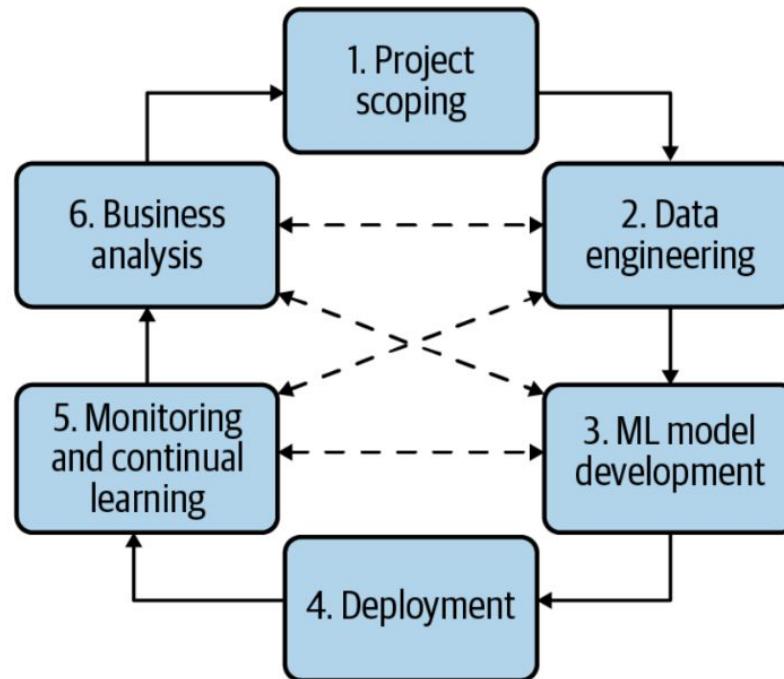
# Key concept: Ethical AI

Guidelines & legislation on building trustworthy AI.

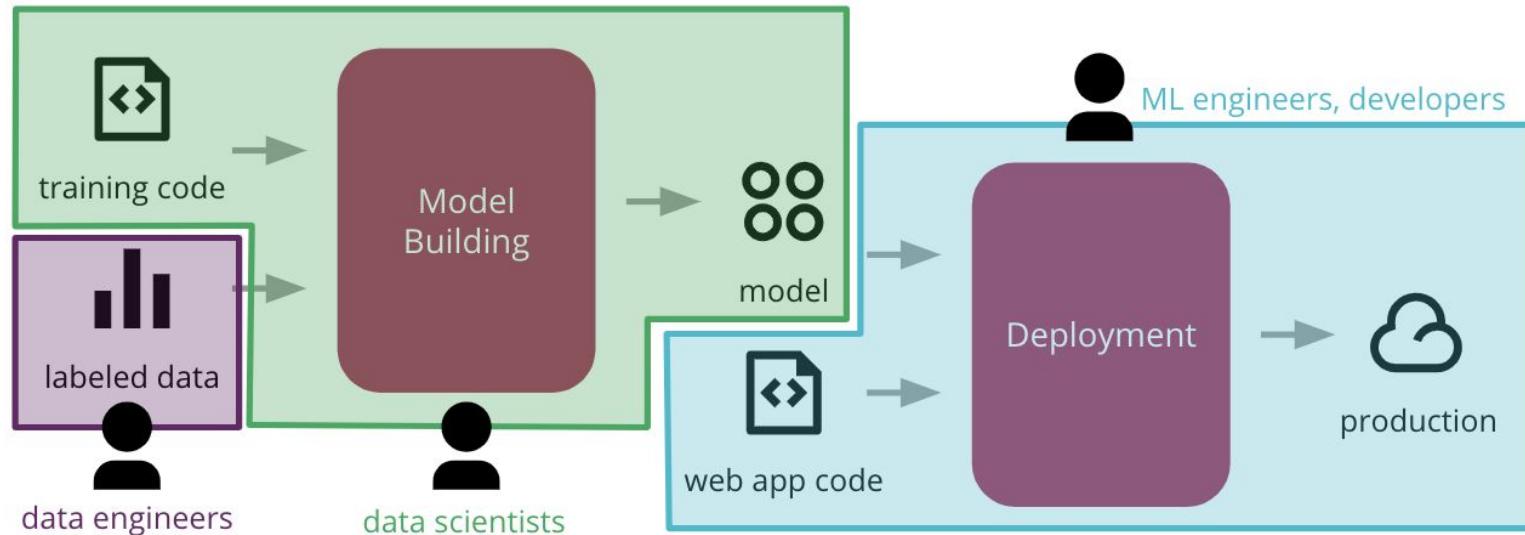


# **Roles & organisation of ML projects**

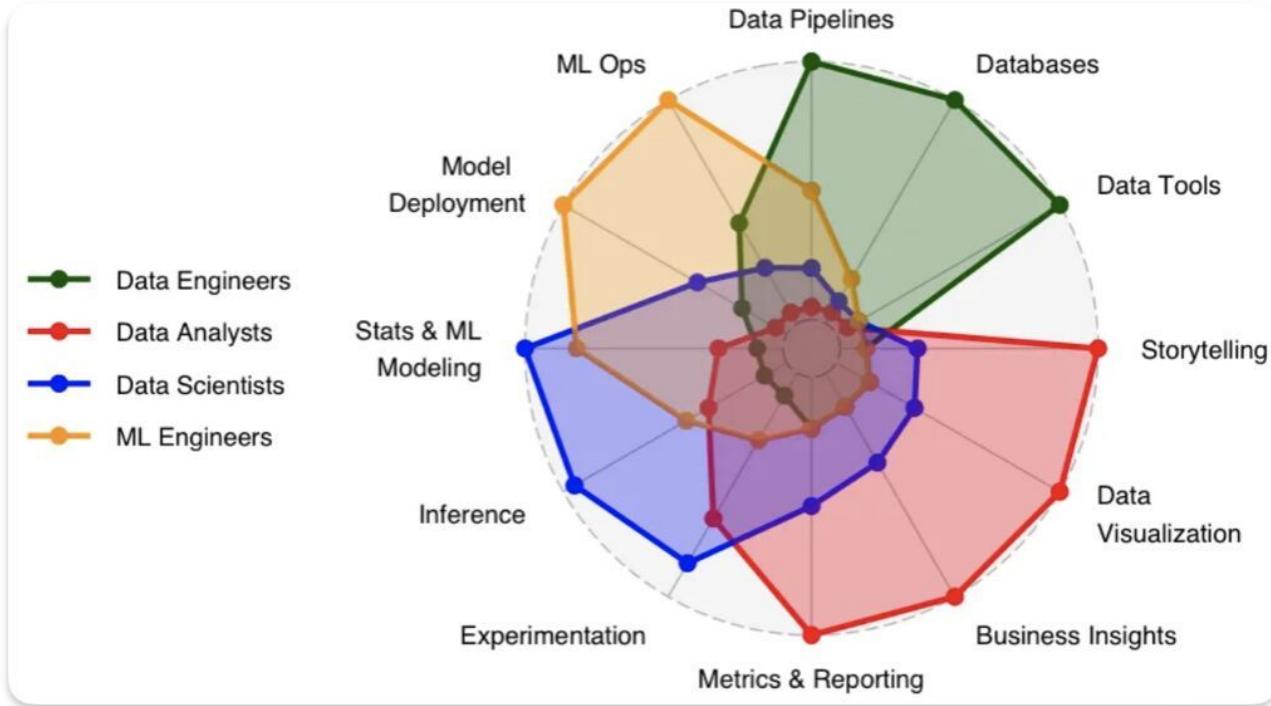
# Typical ML project lifecycle



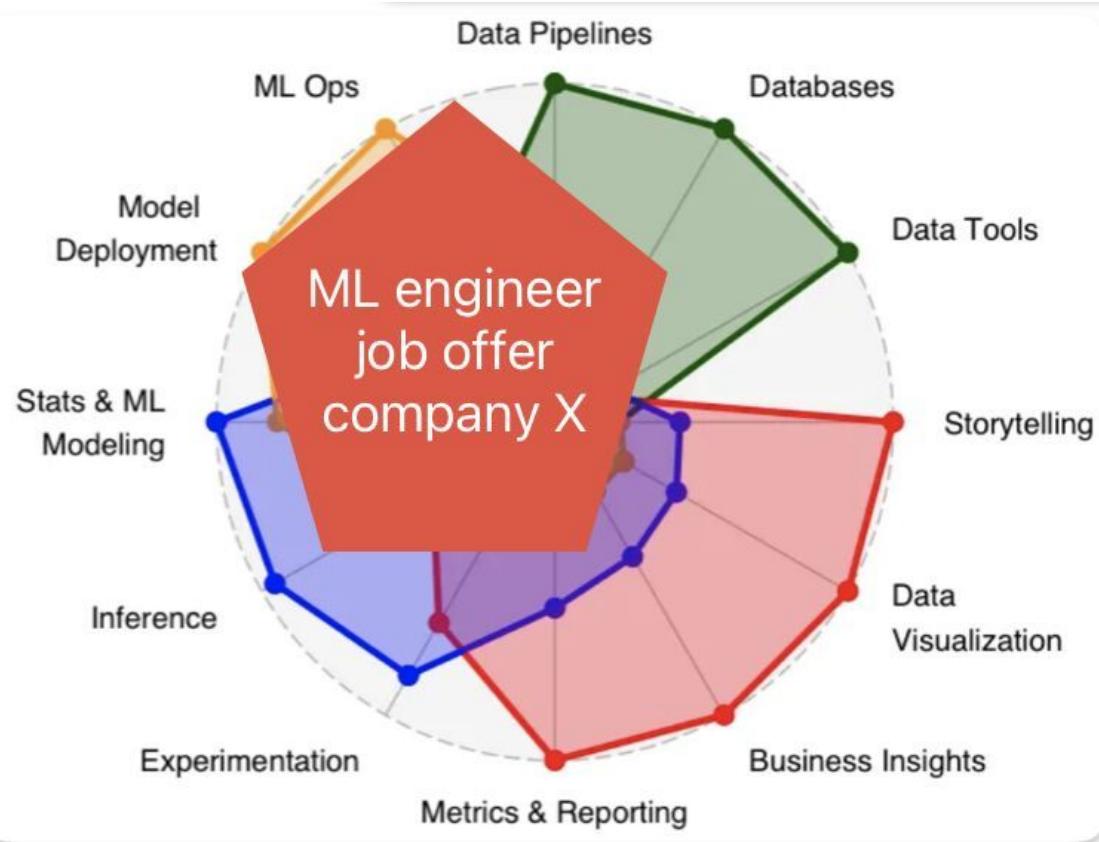
# Roles around a ML system implementation.



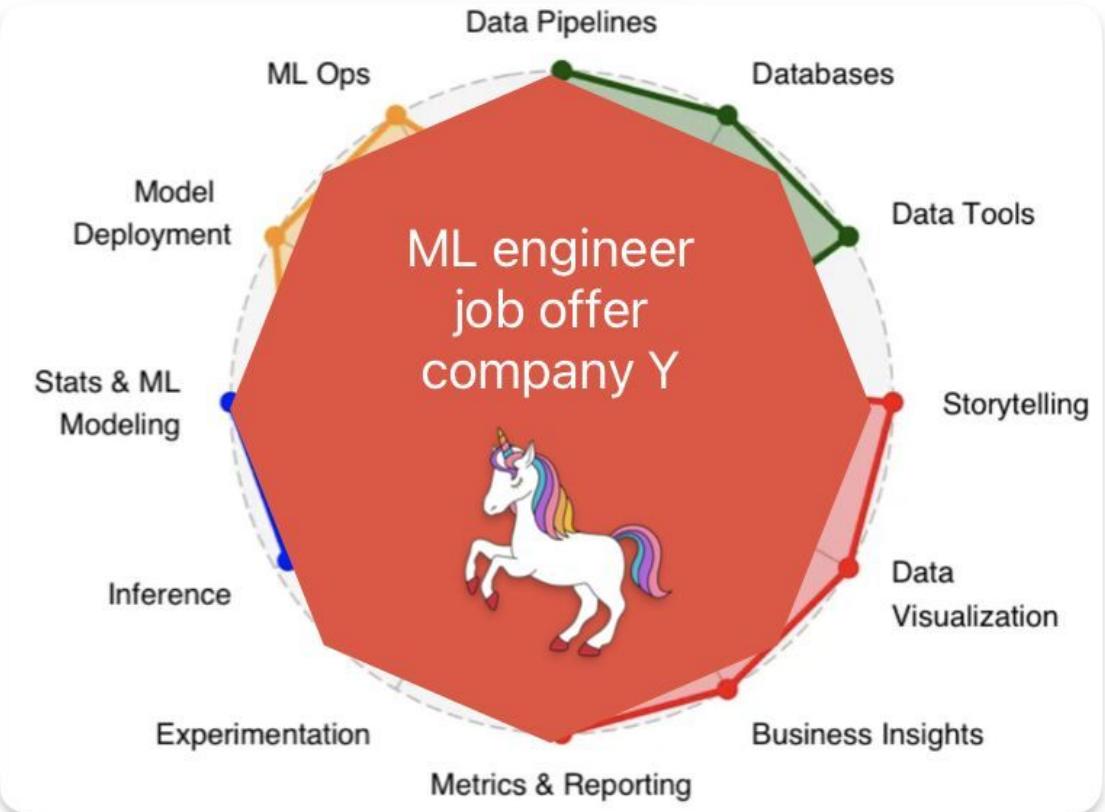
# Different set of skills per roles



In reality it's  
a bit blurry



In reality it's  
a bit blurry



# ML Engineering skills are in high demand

Chip Huyen @chipro · Oct 12, 2020  
Machine learning engineering is 10% machine learning and 90% engineering.  
88 608 7.6K

You Retweeted  
Elon Musk @elonmusk  
Replies to @chipro  
Yeah  
11:09 PM · Oct 12, 2020 · Twitter for iPhone  
93 Retweets 16 Quote Tweets 5,293 Likes



Andrej Karpathy · Following  
(Former) Director of AI at Tesla, Op...  
1yr • Edited • 3

I am hiring Deep Learning Engineers for the Tesla AI team. Strong software engineering is the primary requirement. Except for the scientist role, deep learning interest or knowledge is only a bonus (we will teach you). For the deep learning scientist role any domain outside of computer vision (e.g. speech, NLP, etc.) works great too.

# Teams can adopt different MLOps maturity levels

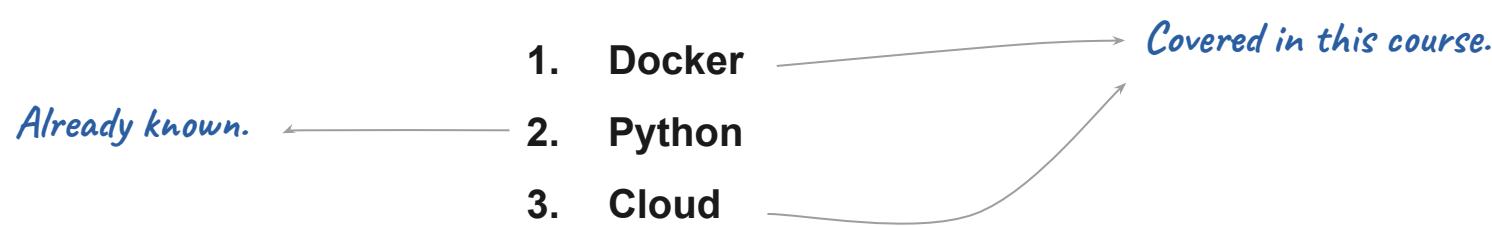


| Level                                       | Highlights  | Technology   |
|---|---|--|
| <b>Level 0<br/>No MLOps</b>                 | <ul style="list-style-type: none"><li>Difficult to manage full ML model lifecycle</li><li>Teams are disparate and releases are painful</li><li>"black boxes," little feedback during/post deployment</li></ul>            | <ul style="list-style-type: none"><li>Manual training, builds and deployments</li><li>Manual testing of model and application</li><li>No centralized tracking of model performance</li></ul> |
| <b>Level 1<br/>DevOps but<br/>no MLOps</b>  | <ul style="list-style-type: none"><li>Releases are less painful than No MLOps</li><li>Limited feedback on how well a model performs in production</li><li>Difficult to trace/reproduce results</li></ul>                  | <ul style="list-style-type: none"><li>Automated builds</li><li>Automated tests for application code</li></ul>  |
| <b>Level 2<br/>Automated<br/>Training</b>   | <ul style="list-style-type: none"><li>Training environment is fully managed and traceable</li><li>Easy to reproduce model</li><li>Releases are manual, but low friction</li></ul>   | <ul style="list-style-type: none"><li>Automated model training</li><li>Centralized tracking of model training performance</li><li>Model management</li></ul>                                 |
| <b>Level 3<br/>Automated<br/>Deployment</b> | <ul style="list-style-type: none"><li>Releases are low friction and automatic</li><li>Full traceability from deployment back to original data</li><li>Entire environment managed: dev &gt; test &gt; production</li></ul> | <ul style="list-style-type: none"><li>Integrated A/B testing of model performance</li><li>Automated tests for all code</li><li>Centralized tracking of model training performance</li></ul>  |
| <b>Level 4<br/>Full MLOps</b>               | <ul style="list-style-type: none"><li>Full system automated and easily monitored</li><li>Automated feedback collection and retraining</li><li>Close to zero-downtime</li></ul>  | <ul style="list-style-type: none"><li>Automated model training and testing</li><li>Verbose, centralized metrics from deployed model</li></ul>  |

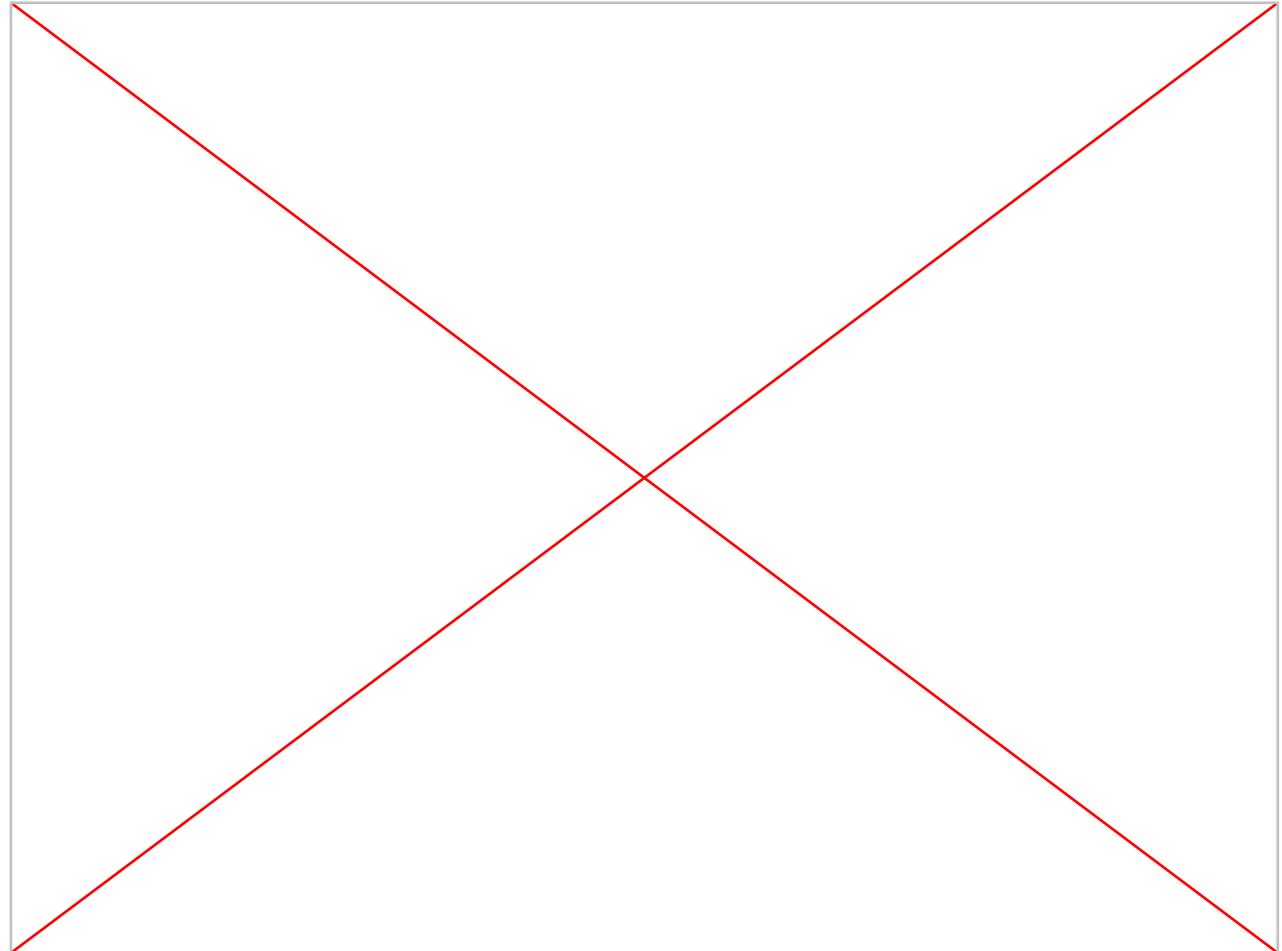
# Study on demanded skills for MLOps engineers.

Looking at 310 job offers on MLOps in Q4 2023.

Top 3 highest demanded skills:



Going from  
standard ML  
Engineer to  
MLOps master...



# **Real-life example of a MLOps organisation (or AI Platform)**

Linkedin case study

# Linkedin integrates many ML applications

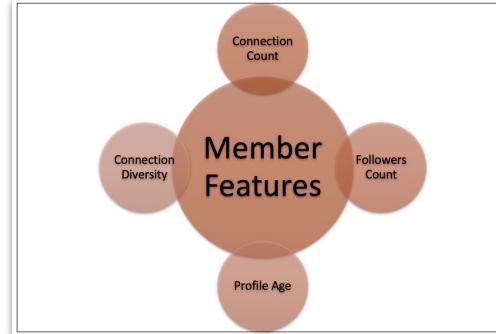
## Viral spam content detection

Detecting spam content...

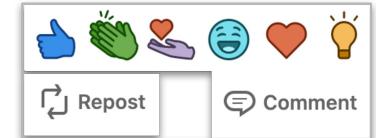


... Using boosted tree algorithm  
on the following features:

### Post features



### Member features



### Engagement features

# Linkedin integrates many ML applications

## Personalised LinkedIn News Feed

Select personalised content for users...



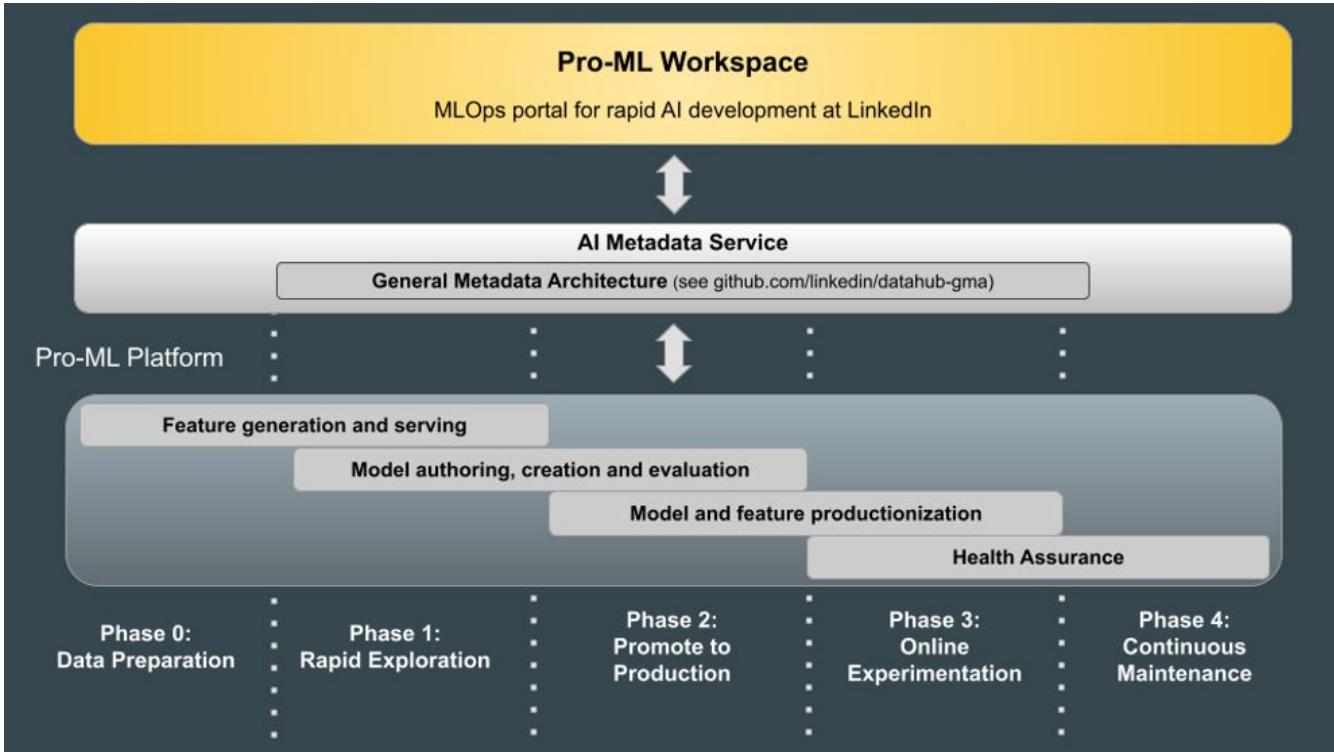
... Using boosted tree algorithm on the following features:

**Identity:** Who are you? Where do you work? What are your skills? Who are you connected with?

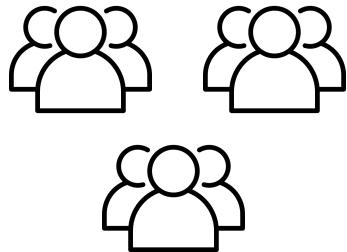
**Content:** How many times was the update viewed? How many times was it “liked”? What is the update about? How old is it? What language is it written in? What companies, people, or topics are mentioned in the update?

**Behavior:** What have you liked and shared in the past? Who do you interact with most frequently? Where do you spend the most time in your news feed?

# Linkedin's Productivity Machine Learning (Pro-ML) platform.



*Teams of  
data scientists*



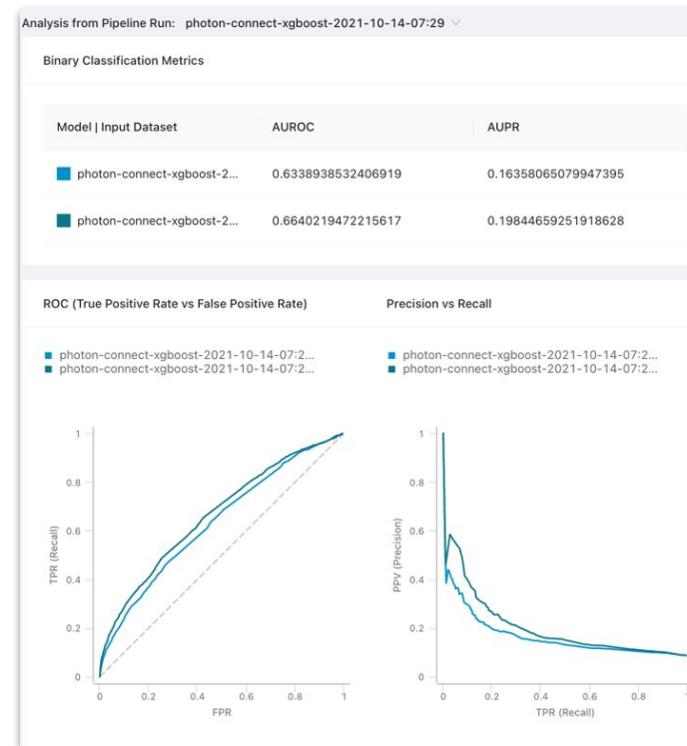
# Linkedin Pro-ML platform.

## Step: Model authoring, creation, and evaluation

Model tracking and experimentation platform

Similar to **MLFlow** or **Weights & Biases** (which we will cover in this course).

The screenshot shows the LinkedIn Pro-ML platform's 'Training' section. On the left is a sidebar with icons for Pipeline Runs, Models, and Projects. The main area has tabs for Pipeline Runs, Models (selected), and Projects. Under Models, there are tabs for My Models (60), All (9999+), and ... . A search bar at the top right allows searching by model name. Below the tabs is a table with columns: Model, Date & Time (UTC), Project, Step Type, Component, Location, and Model status. Two rows are visible: one for 'photon-connect-xgboost-2021-10-14-07-29-xgbconst\_training\_autotune' and another for 'photon-connect-xgboost-2021-10-14-07-29-quasar-servingconfig-replacer'. Both rows show the date and time of the run, the project name 'photon-connect-v2-demo-elong', the step type (Model Training or Model Rewrite), the component (XGBoostTrainer or QuasarServingConfigReplacer), and a status of 'Unpublishable'.



# Linkedin Pro-ML platform.

## Step: Model productionisation

The screenshot shows the LinkedIn Pro-ML Platform's interface. On the left is a dark sidebar with icons for Pro-ML Workspace, Training, Publishing (which is selected and highlighted in blue), Monitoring, Search, and Help. The main area has a header "Publishing" with dropdown menus for "Models in progress", "Published models", and "Model groups". Below this is a table with columns: Model Name, Publish Name, Version, Model Group, Date & Time Created, Created By, Status, and Actions. There are two rows in the table:

| Model Name    | Publish Name  | Version | Model Group        | Date & Time Created    | Created By | Status       | Actions |
|---------------|---------------|---------|--------------------|------------------------|------------|--------------|---------|
| tg_mre-demo   | zetastg11     | 0.0.1   | test-model-group-3 | 05/06/2020<br>18:17:30 | [redacted] | ● Publishing | (edit)  |
| kabootarModel | test-approval | 0.0.1   | test-model-group   | 10/29/2020<br>21:32:35 | [redacted] | ● Publishing | (edit)  |

Workflows to publish or deprecate models.

# Linkedin Pro-ML platform.

Step: “Health insurance”  
(aka monitoring)



# **Use case deep dives**

Real-estate valuation  
assistant

# Context & Problem Statement

...heard of Fednot?



## Fednot

- = Koninklijke Federatie van het Belgisch Notariaat
- = Fédération Royale du Notariat belge
- = Royal Federation of the Belgian Notaryship

Fednot supports the notary studies with juridical advice, office management, IT solutions, trainings, and information for the general public.

# Valuation assistant.

N Val

e-notariaat.acc.credoc.be/valuation\_v1/estimation/result

FEDNOT | Waarderingsassistent immo

Thomas UYTTENHOVE TEST ETUDE 12 NL

TERUG | RESULTAAT

Dataset van het pand

Adres  
10 Sportstraat, 9000 - Gent  
Percelennummer  
4480910810/00F006

Waardering

STUUR UW FEEDBACK

Resultaat van de waardering ⓘ

Indicatieve prijs € 397 000

| Prijsbereik           | Aantal huizen | Percentage (%) |
|-----------------------|---------------|----------------|
| < € 100.000           | 10            | ~3%            |
| € 100.000 - € 150.000 | 15            | ~5%            |
| € 150.000 - € 200.000 | 20            | ~7%            |
| € 200.000 - € 250.000 | 30            | ~10%           |
| € 250.000 - € 300.000 | 40            | ~13%           |
| € 300.000 - € 350.000 | 50            | ~16%           |
| € 350.000 - € 400.000 | 60            | ~20%           |
| € 400.000 - € 450.000 | 70            | ~23%           |
| € 450.000 - € 500.000 | 80            | ~27%           |
| € 500.000 - € 550.000 | 70            | ~23%           |
| € 550.000 - € 600.000 | 60            | ~20%           |
| € 600.000 - € 650.000 | 50            | ~17%           |
| > € 650.000           | 40            | ~13%           |

Indicatieve prijs en distributie van de geïndexeerde verkoopprijs van 319 huizen binnen een straal van 1 km.

HOE HEEFT HET MODEL DEZE INDICATIEVE PRIJS BEREIKT?

Map details: The map shows the location of the property at 10 Sportstraat. Surrounding landmarks include AZ Jan Palfijn, Zwembad GUSE, Delhaize Watersport, McDonald's, KASK & Conservatorium, Muziekcentrum De Blijfje, and various streets like Neermeerskaai, Gordunakai, and sportstraat. Numbered pins (1-10) mark specific locations along the roads.

Services 1.8.0 UI 8.9.0

Keyboard shortcuts | Map data ©2022 Terms of Use Report a map error

© 2024 ML6. All rights reserved. ML6 Public Information

# How the ML model conceptually works

| Known values that the model will use as input to make predictions |             |          |           |                         |        | What the model needs to predict |
|---|-------------|----------|-----------|-------------------------|--------|---------------------------------|
| Feature variables   |             |          |           |                         |        | Target variables                |
| ID  | Size (sqft) | Bedrooms | Bathrooms | Distance to City Center | Garage | House Price (k\$)               |
| 1   | 2200        | 3        | 2         | 5                       | Yes    | 300                             |
| 2   | 1800        | 4        | 2         | 3                       | No     | 275                             |
| 3   | 1400        | 2        | 1         | 10                      | Yes    | 200                             |
| ...   | ...         | ...      | ...       | ...                     | ...    | ...                             |
| 80 000  | 3000        | 5        | 4         | 4                       | Yes    | 400                             |
| 80 001  | 1600        | 3        | 2         | 12                      | No     | ? (Test)                        |
| ...   | ...         | ...      | ...       | ...                     | ...    | ...                             |
| 100 000   | 2100        | 4        | 2         | 6                       | Yes    | ? (Test)                        |

**Single house** ←

**Train set**

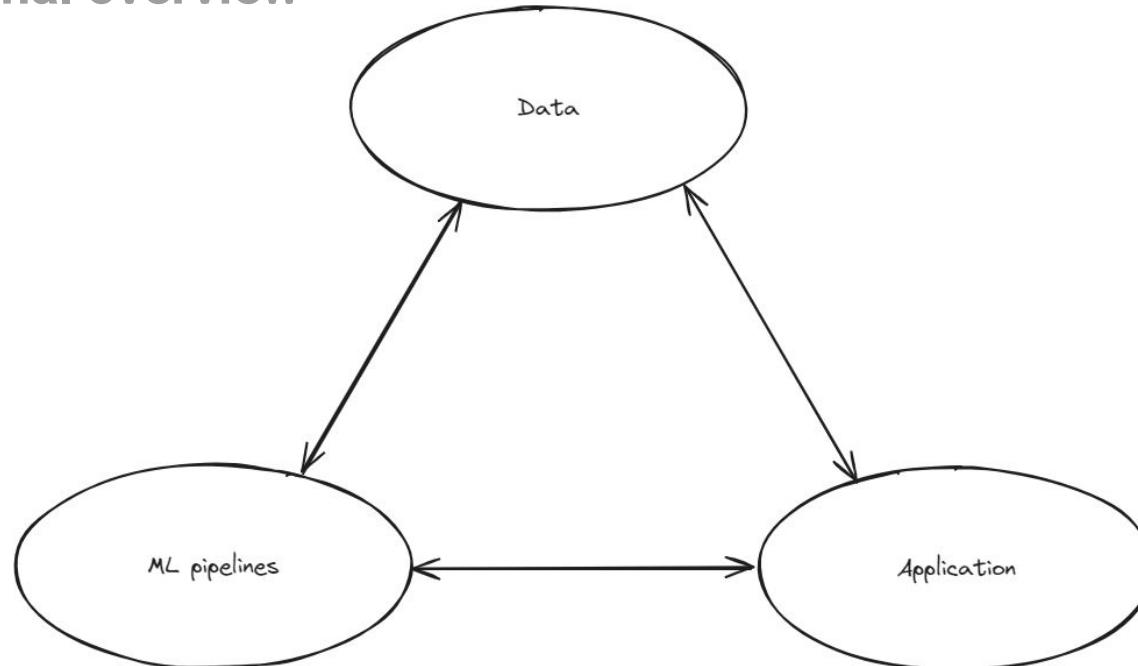
**Test set**

The ML model will see many **observations** (houses) defined as a set of **features** (information, variables). From it the model will learn patterns and what impacts **target variable** (house prices).

If given new observations, the model can **predict** the target variable based on the input features.

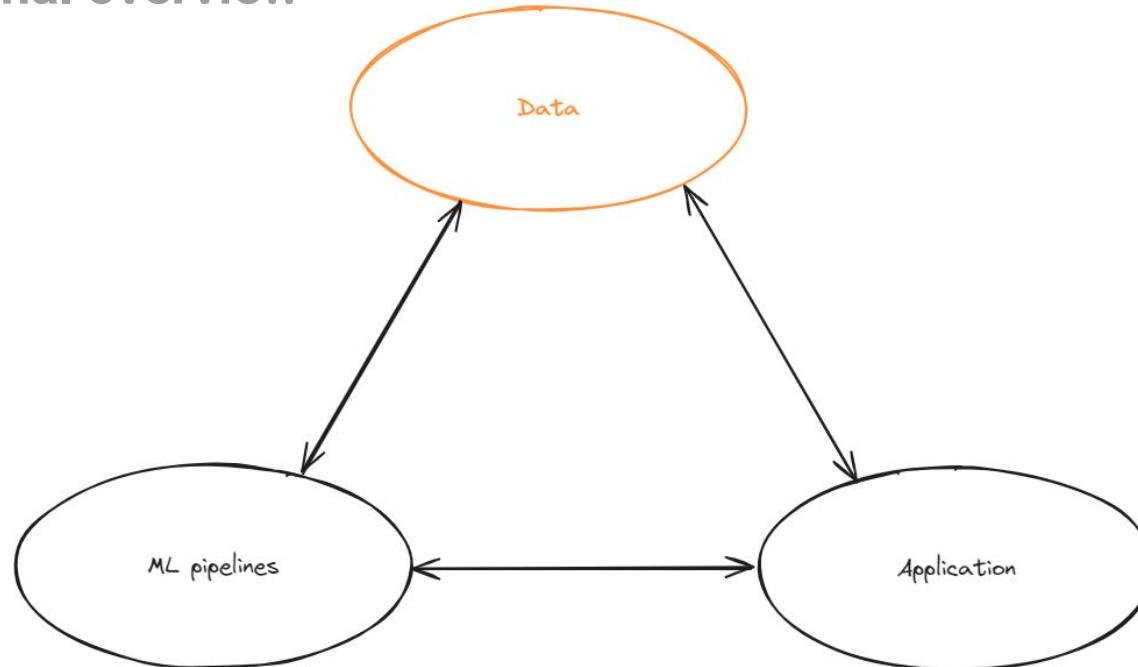
# Solution architecture.

## A functional overview



# Solution architecture.

## A functional overview



# Valuation features.

## Legal real-estate data

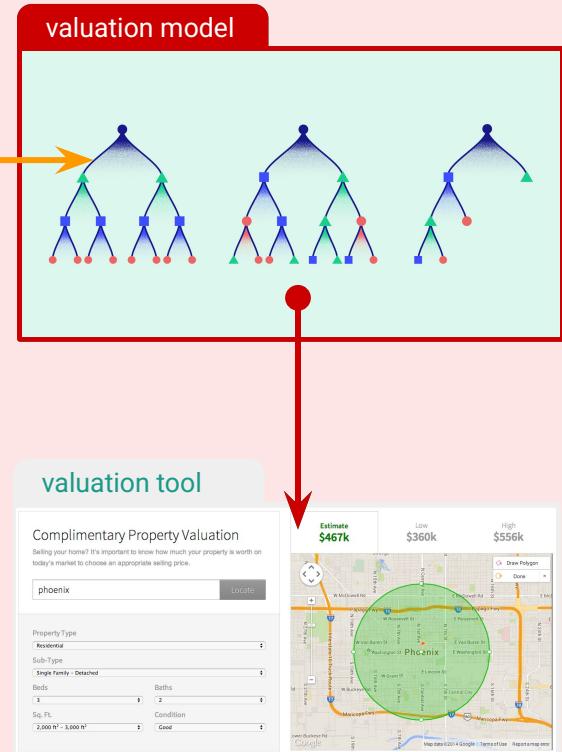
processed document

legal platforms



Legal ('AI') features

Legal features alone do not capture sufficient information to accurately predict real estate prices...

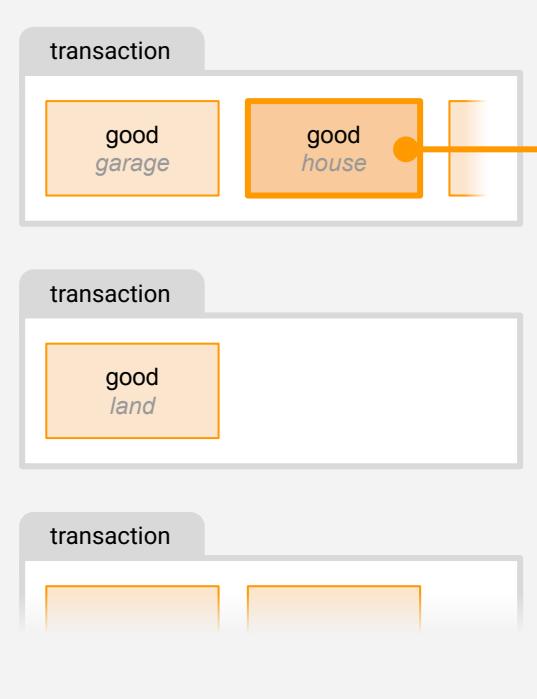


# Valuation features.

## Open real-estate data

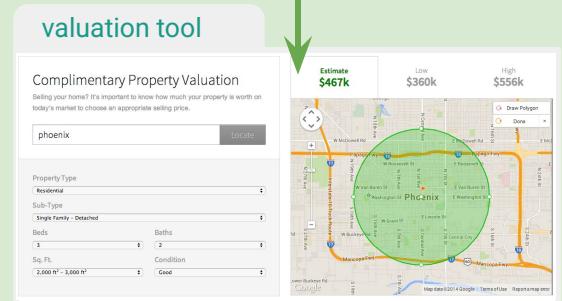
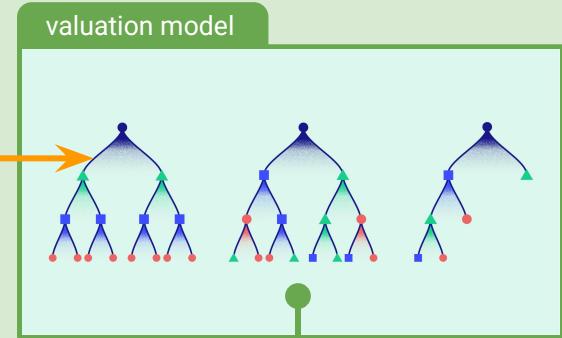
processed document

legal platforms



Legal ('AI') features  
+  
Leverage open data to expand limited feature set

- Various types of data sources provided by the government or community
- (Mostly) freely accessible
- Opportunities for complex, more informative features!



# Valuation features.

## Open real-estate data

### ■ Cadastral information

- Parcel area, street width, ratio, and orientation;
- Building area, type (“open”, “half-open”, “closed”), facade width, and orientation



# Valuation features.

## Open real-estate data

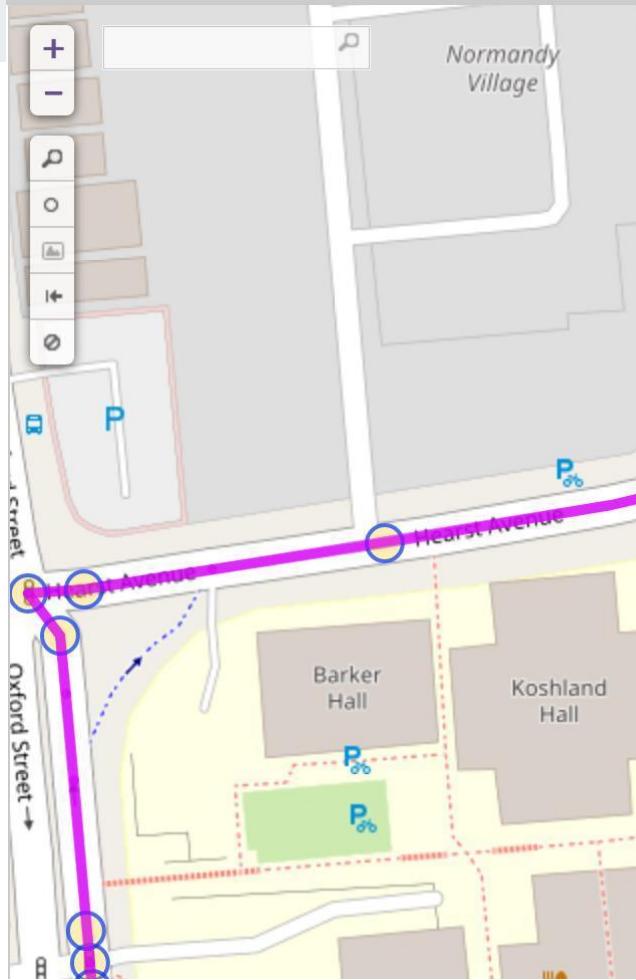
- Cadastral information
- Height information
  - Building height, volume, number of stories.



# Valuation features.

## Open real-estate data

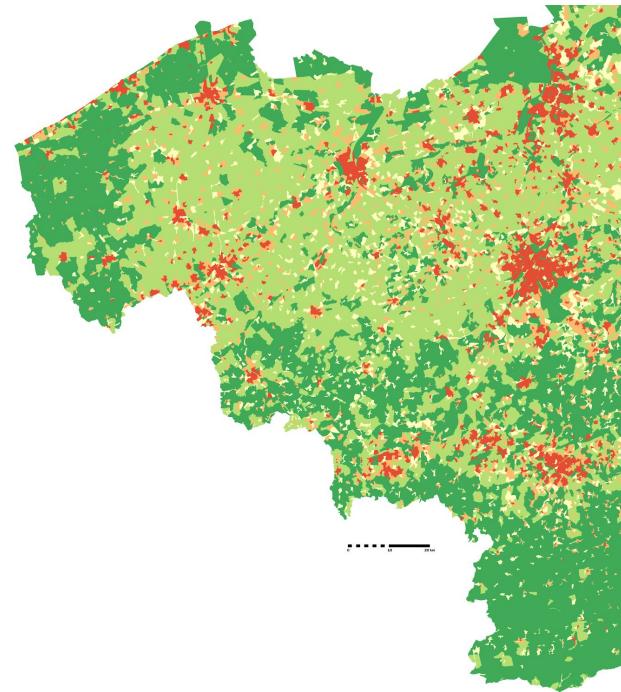
- Cadastral information
- Height information
- Location information
  - Distance to major cities and to nearest city center, highway (entry), primary road, railway, station, bus stop, etc.



# Valuation features.

## Open real-estate data

- Cadastral information
- Height information
- Location information
- Local socio-economic and demographic statistics
  - Municipality population size, tax percentage, prosperity index, avg. income;
  - Statistical sector cadastral income percentiles



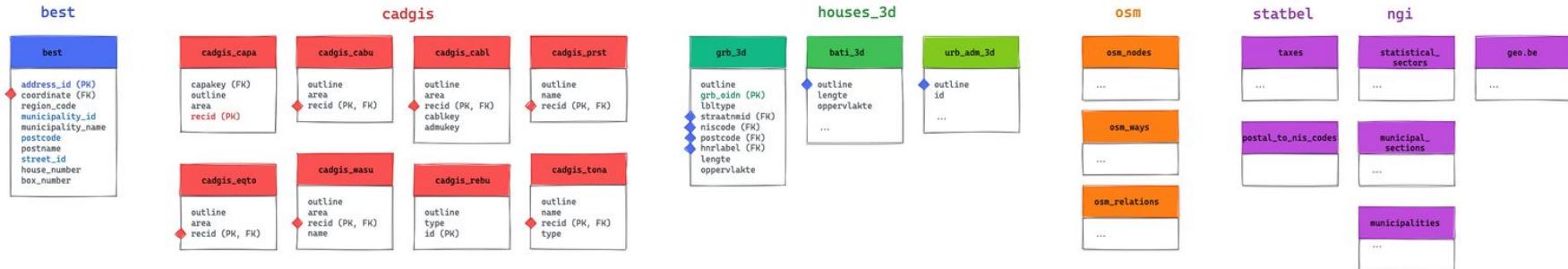
2011 Density (people / km<sup>2</sup>)

|              |
|--------------|
| 0 - 50.4     |
| 50.4 - 392   |
| 392 - 1080   |
| 1080 - 2120  |
| 2120 - 45700 |

Source: statbel.fgov.be

# Open Data

## Sources



**Update Frequencies:**

Weekly

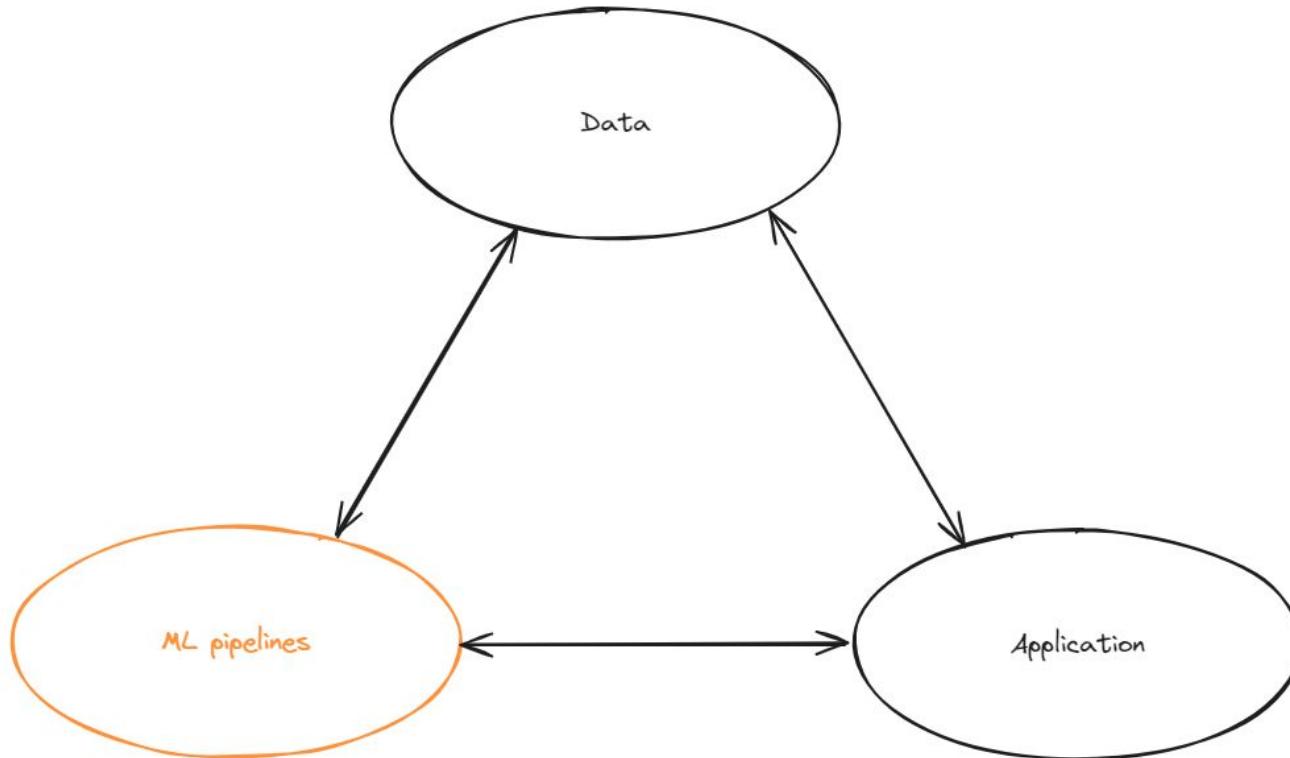
Yearly

Wallonia - Not (2016)  
Flanders - Not (2016)  
Brussels - Yearly

Weekly

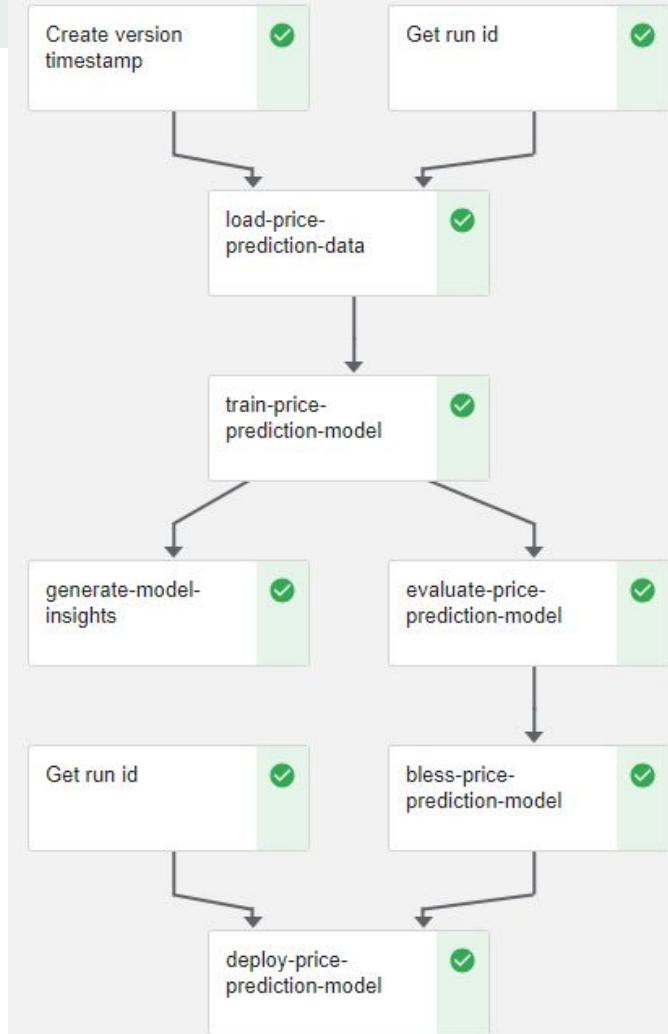
Yearly

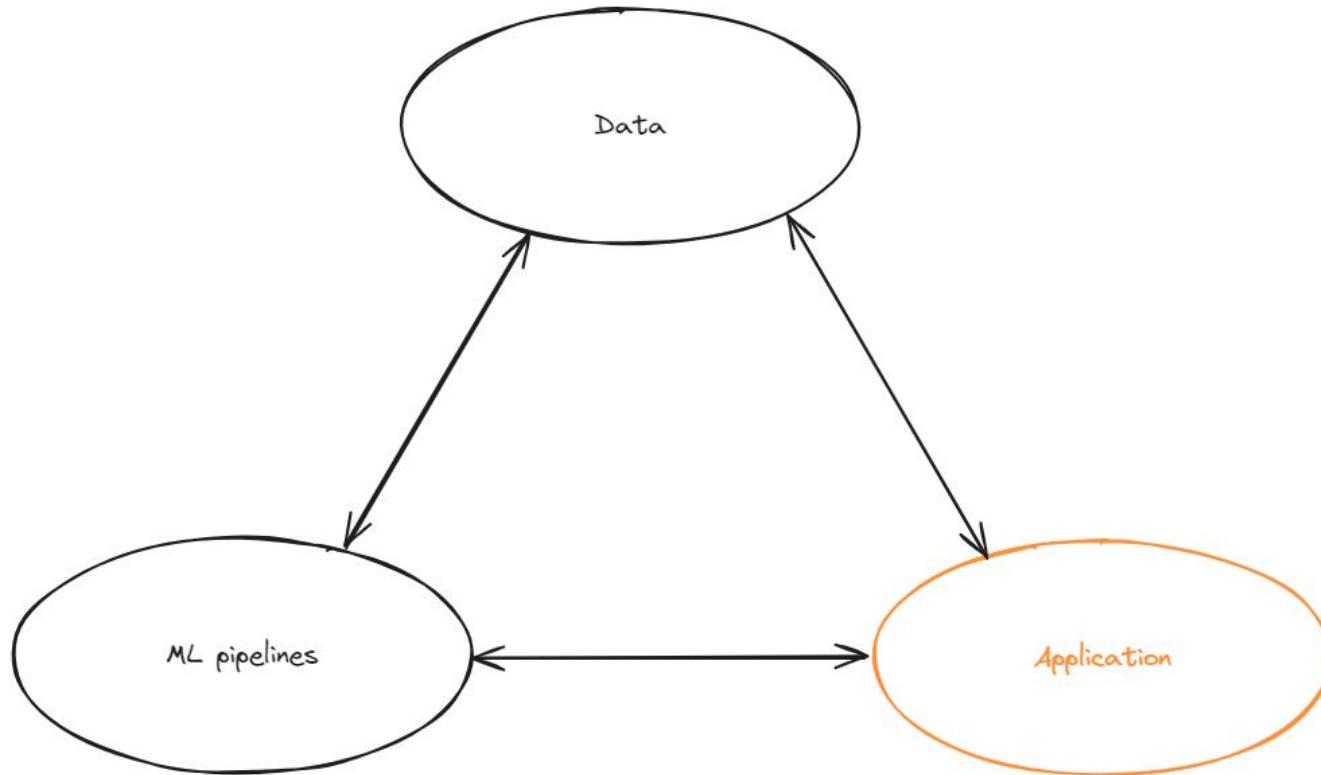
Yearly



# Automated pipeline to train and deploy new price prediction models

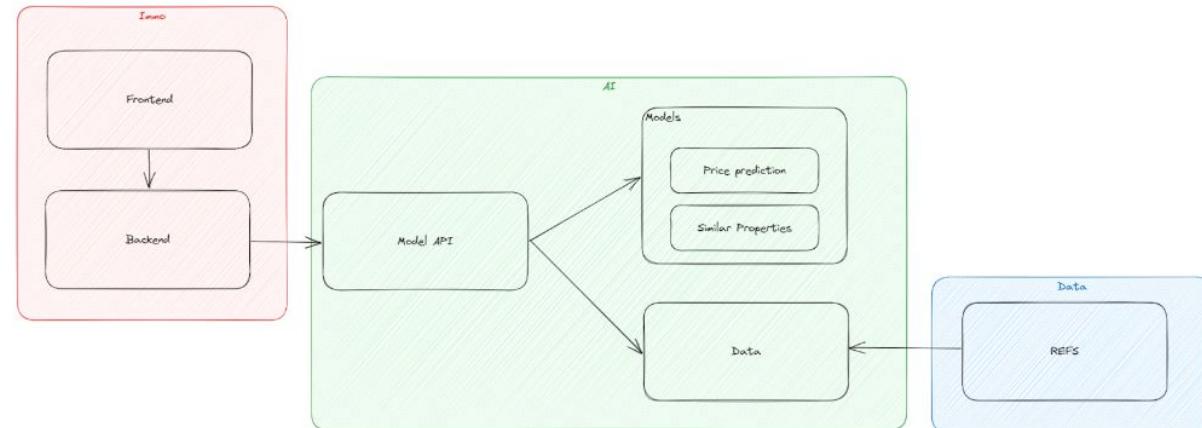
- Allows you to implement a **ML pipeline** made of different components, usually ran sequentially.
- Each component can be a **Docker image**
- Hosted on a **kubernetes cluster** (set of node machines for running containerized applications). Can be on the **Cloud**.
- **Benefits**
  - Modularized
  - Reproducible
  - Efficient
  - Scalable
  - Deployments
  - Collaboration
  - Version control and documentation





# Model API

- **Cloud Run**
  - Hosting the front-end
  - Hosting APIs to connect components
- **GKE**
  - Hosting the models
- **Bigquery**
  - Hosting the data
- **Storage**
  - Hosting artifacts



# **Project phases & challenges**

# Build different stages of your solution

## Proof of Concept

Use easily available data to show that your model or solution can work.  
Low efforts.  
Prove the feasibility and value.  
Iterate fast.

## Minimum Viable Product

Just enough features for a small set of users to start using it.  
Gather feedback and make sure that it is designed in an optimal way.

## Productionisation / scaling

Build the infrastructure to finally deploy your solution and let users use it.  
Gradual roll-out to more and more users in more and more markets.  
Deploy better models, attract more users, go to new markets, maintain the solution, ...

## Maintenance

Keep the solution up and running.  
Monitor resources and performance.  
Update packages and dependencies (software around solution change).  
Security and up-time.

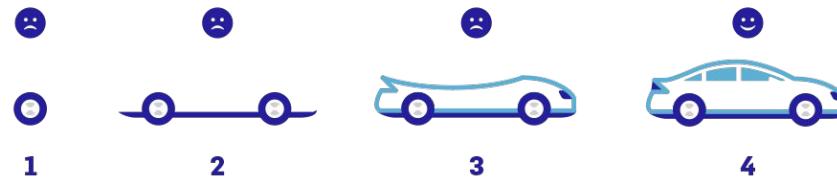
| POC     | MVP      | Productionisation / scaling | Maintenance           | ... |
|---------|----------|-----------------------------|-----------------------|-----|
| 2 weeks | 2 months | 6 months                    | As long as it's up... |     |

# Build different stages of your solution

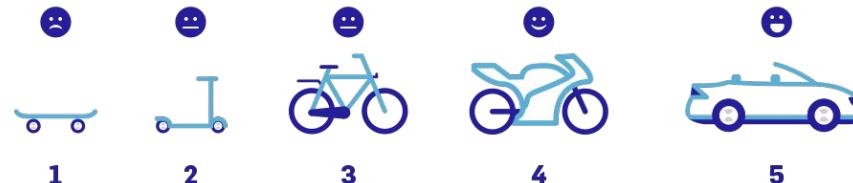


# At each stage, your product should be usable

**NOT LIKE THIS!**



**LIKE THIS!**



# Data science projects are challenging to bring to production

Breaking the myth

**Forbes**

*“87% of data science projects never make it into production...”*

**VentureBeat**

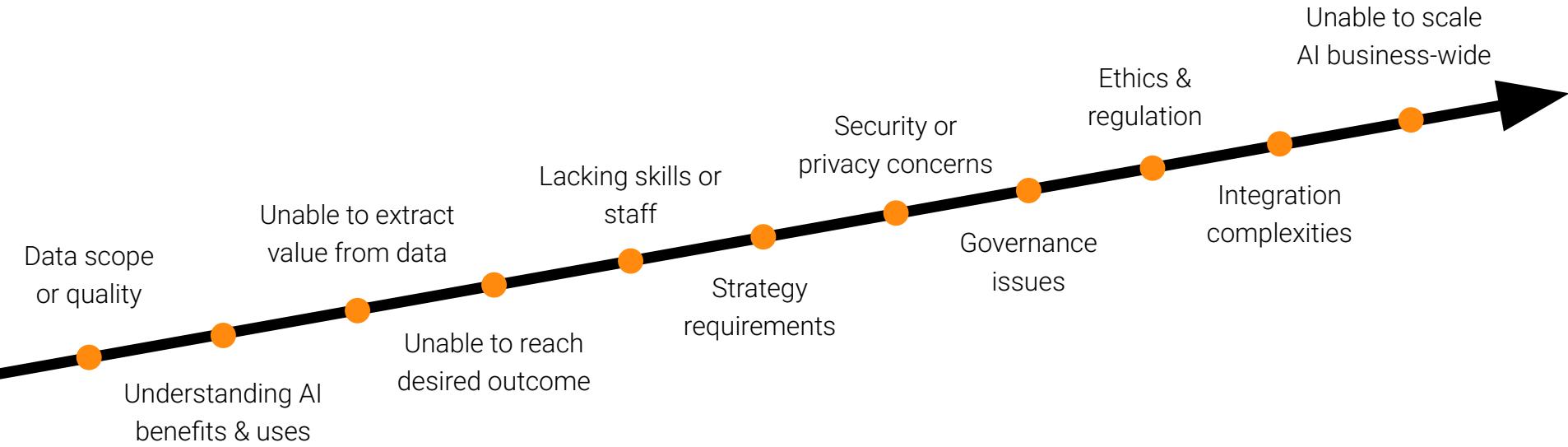


*(Might not be a factual number...)*

*But data science project are still challenging to actually roll-out to the real world!*

# AI Journey Challenges.

While AI is an enabler for strategic priorities, it doesn't come without its challenges.

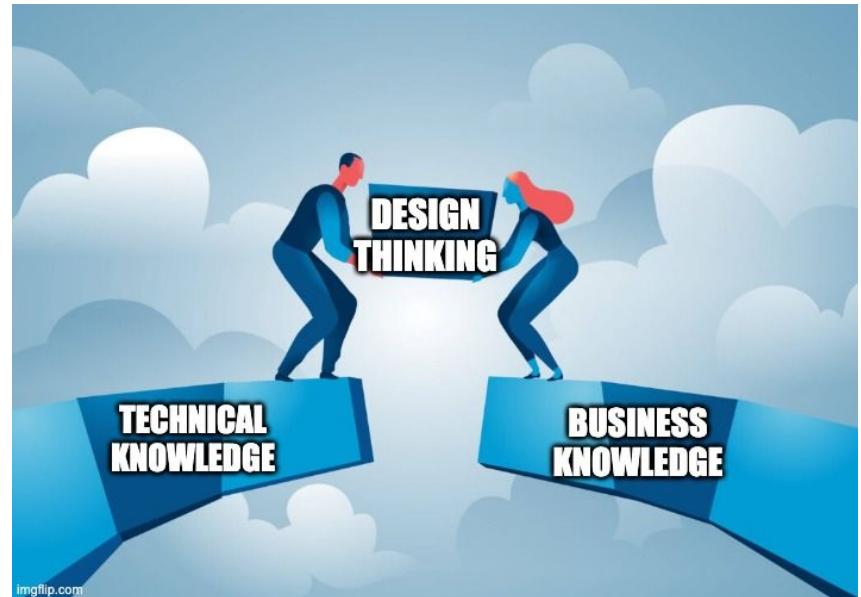


# **Project definition framework**

# Process to define new use cases.

## How to get started?

- **New ideas** do not come spontaneously
- Proactively organise **workshops** to identify how ML can create value in an organisation.
- Use **design thinking** techniques.
- Make sure to have the **right people around the table** (decision makers, stakeholders, users and (ofc) engineers).
- Spend enough time in it - **starting in the right direction** is key.



# Other concept: Design thinking

Same ideas, different framework  
(coming from front-end engineering)

## 1. Empathize

Engage in qualitative research methods such as interviews and workshops to deeply understand the users, their needs, and their pain points.

## 2. Define

Clearly articulate the user's needs and challenges based on the insights gathered during the empathize phase. Map out the user's interaction with the solution.

## 6. Implement

Once the design is finalized, begin the development process using appropriate technologies and frameworks.



## 3. Ideate

Engage in collaborative sessions to generate a wide range of solutions and ideas.

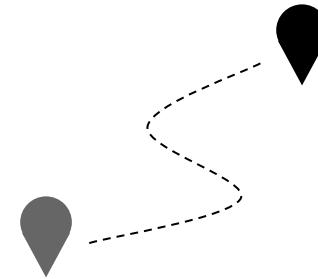
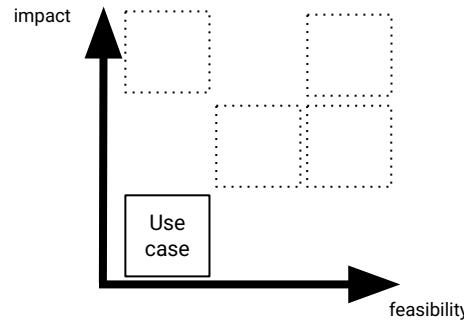
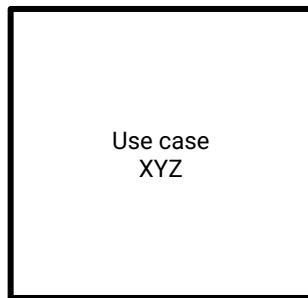
## 4. Prototype

Build a **mock** application to validate whether it fits your users needs.

## 5. Test

Monitor user interactions and gather data to measure the application's success. Maintain an ongoing feedback loop with users to continually refine and improve the application.

# Framework to define an AI use case.



**1** Identify AI opportunities

**2** Evaluate and refine selected use cases and their feasibility

**3** Prioritize top use cases to kickstart AI

**4** Define the roadmap towards this AI use case

# Identify opportunities

- Ideate and map user process
  - Identification of **business opportunities**
  - Identification of **challenges**
  - **Opportunities:** where can AI help?
- Cluster opportunities
- Name AI use cases



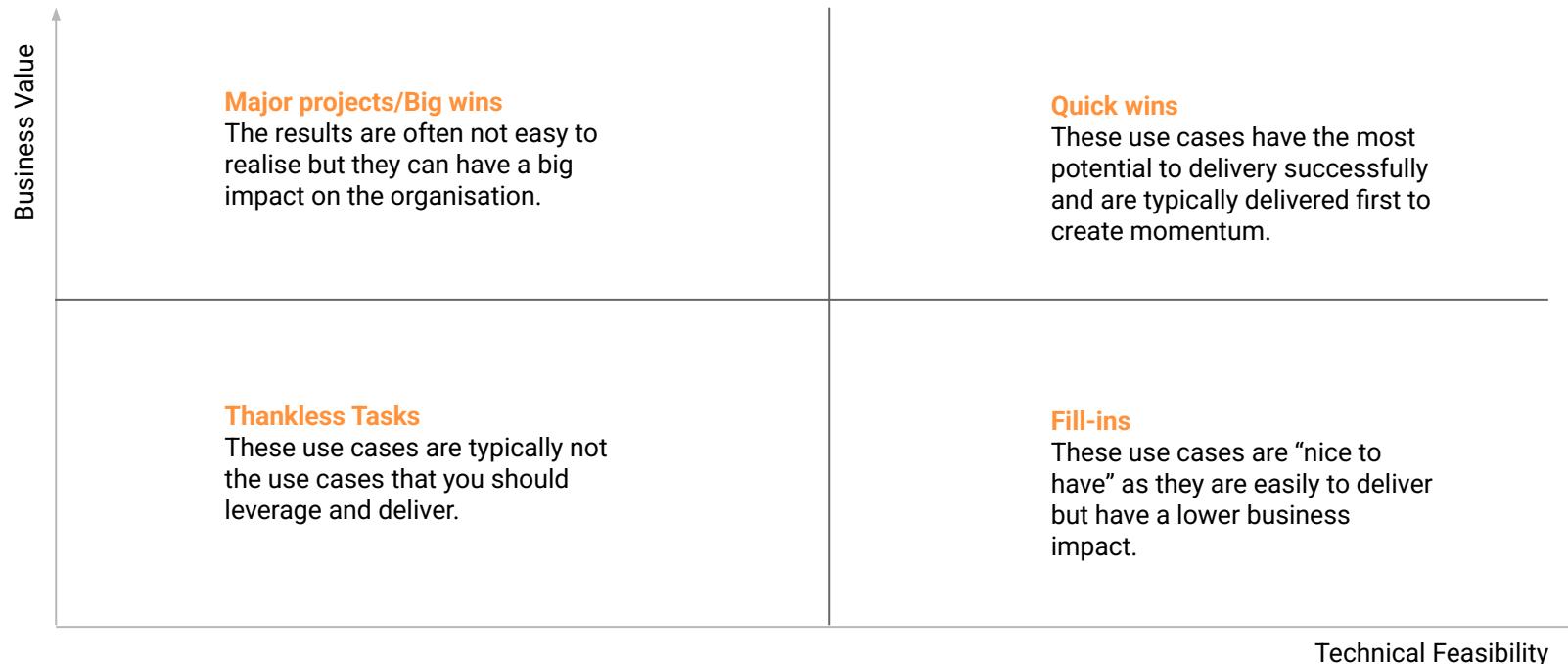
# Use case template.

- How to quickly iterate over a few use cases?
- How to efficiently capture the point of view of different people?
- How to set the vision on a specific use case?

|  |   |
|--|---|
| <b>Use Case:</b> [Cool Name]   |   |
| <b>What?</b><br>[Describe the use case in 2 sentences]   | <b>Value</b><br>[Score out 5 - flash vote]<br>       |
|  | <b>Feasibility</b><br>[Score out 5 - flash vote]<br> |
| <b>Why?</b><br>[Purpose of the solution - e.g. reducing costs, helping users, climate, ...]              |   |
|  | <b>Who?</b><br>[Stakeholders benefiting from the solution (e.g. customers, users, role X, ...)]   |
|  | <b>How?</b><br>[Approach, simplified]   |
| <b>Challenges?</b> <ul style="list-style-type: none"> <li>• ...</li> <li>• ...</li> <li>• ...</li> </ul> |   |
| <b>Evaluation?</b><br>[Metrics and success criteria]   |   |

# Prioritisation matrix.

How to evaluate the different use cases?



# Define and scope your project.

Which questions to answer before getting started with the selected project?  
(Often done offline, after the workshop)



Define value  
drivers



Set success  
criteria



Identify  
challenges



Define building  
blocks

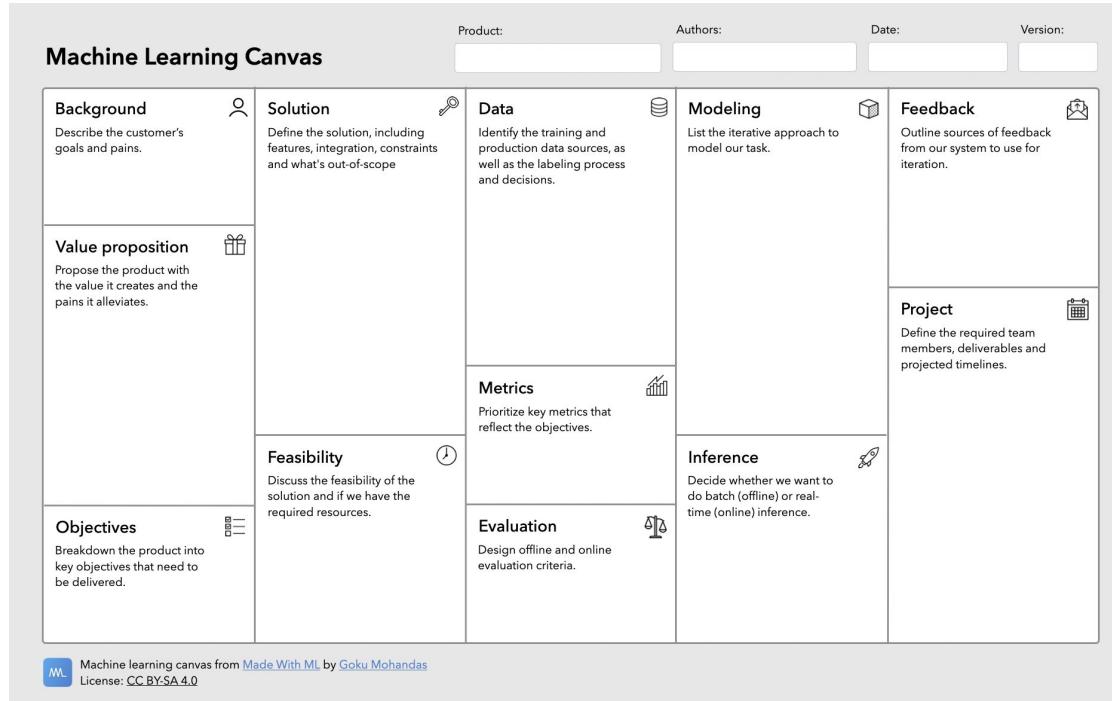


Estimate time  
& budget

Think about  
intermediate  
milestones that  
show value

# Define and scope your project.

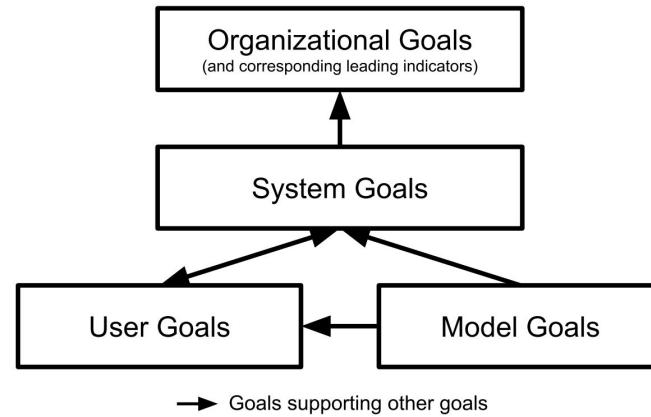
## Product design template



# Aligning your solution with goals on different levels.

- **Organizational goals:** Innate/overall goals of the organization.
- **System goals:** Goals of the software system/feature to be built.
- **User goals:** How well the system is serving its users, from the user's perspective.
- **Model goals:** Quality of the model used in a system, such as performance.

These goals should be aligned with each other



# User adoption

“You can have the best model with the best data, success always depends on how users will adopt it.”

Ways to ensure user adoptions:

- **Power users:** Work with users since day 1. Throughout the use case ideation and during development. You receive critical feedback and can get champions who fully understand the solution to spread its usage once developed.
- **Change management strategy:** From executives and process experts.
- **Integration:** Make sure it works with users favorite tools (a new board in existing platform has much higher chances of being utilised than a new program/website).
- **Documentation:** Clear explanation of *how the model works, performs and should be used*. Training program, videos, tutorials, FAQs, support line, ...
- **Monitor usage:** ... and improve the solution from it.

# When not to use Machine Learning?

It's not always the right solution...

- Clear specifications are available
- Simple heuristics are good enough
- Cost of building and maintaining the ML system outweighs its benefits
- Correctness is of utmost importance
- ML is used only for the hype (e.g., to attract funding)

Examples of these?

# (Really) accurate predictions might not even be that important

## The over-optimizing paradox

- "Good enough" may be good enough
- Prediction critical for system success or just an gimmick?
- Better predictions may come at excessive costs
  - Data is often the bottleneck
  - Cost of producing more data (labeling, infra, collection, ...)
- Better user interface ("experience") may mitigate many problems
  - Explain decisions to users with Explainable AI (XAI)
- Use only high-confidence predictions?

# Critical thinking when doing the project definition

Ask the right questions - make sure you have a solid use case before you start building anything.

- **Baseline:** What is the performance of an alternative to ML? How do simple heuristics or human guess-predictions perform?
- **Probabilistic:** ML is by definition not deterministic. Are probabilities/ranges fine for this use case? E.g. for demand forecasting the model can make errors, for self-driving cars not...
- **Precision / recall:** Are both important? If not, can I make it a success by sacrificing one? E.g. for fraud detection we can raise a warning on false positive, but cannot have false negative...
- **Interpretability:** Do we need to explain why the model makes specific decisions? If yes, can we?
- **Do not reinvent the wheel:** Are there existing open source or 3rd party solutions? Did anybody in my organisation work on something like this?

# Course organisation

# Objective for this course.

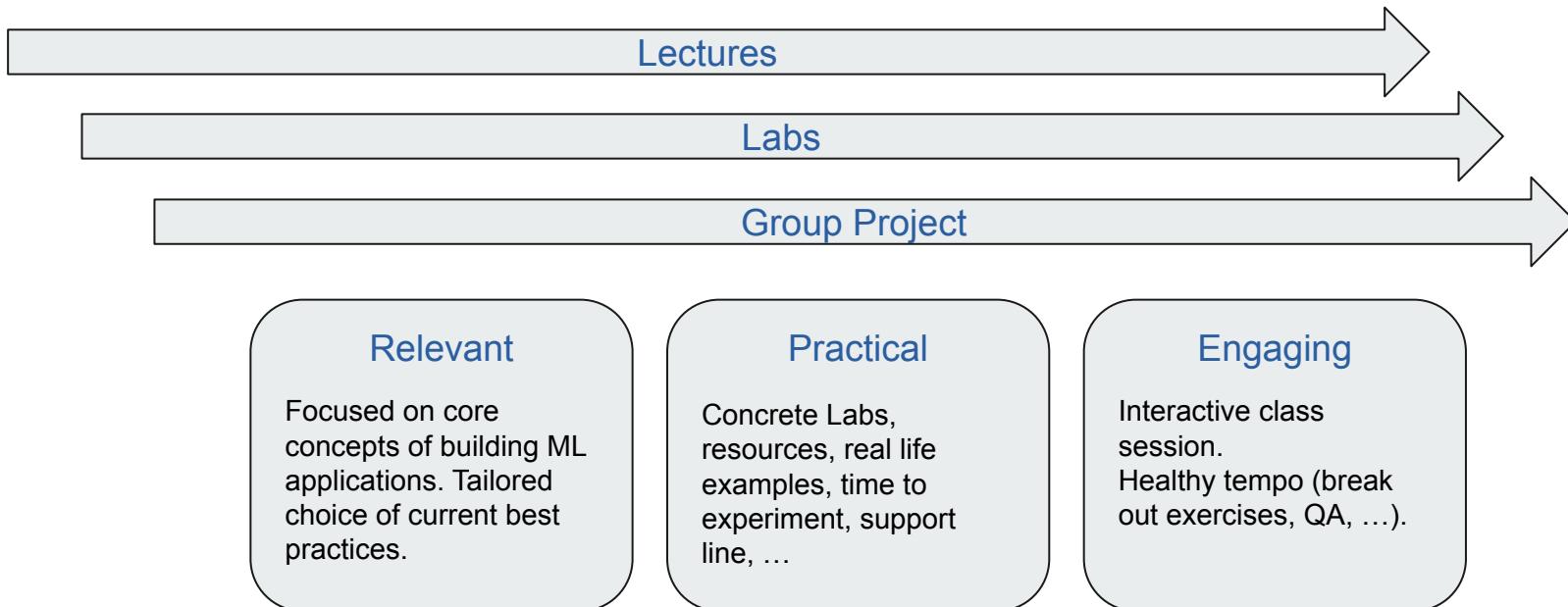
We want to enable you with skills to **design** and **build ML application** 

We selected core **topics** of MLSD to be tackled in this course. Tools are selected based on usability, performance, popularity and accessibility.

Goal is to provide

- **Theoretical** concepts
- **Technical** tools & skills
- Practical real world **practices**

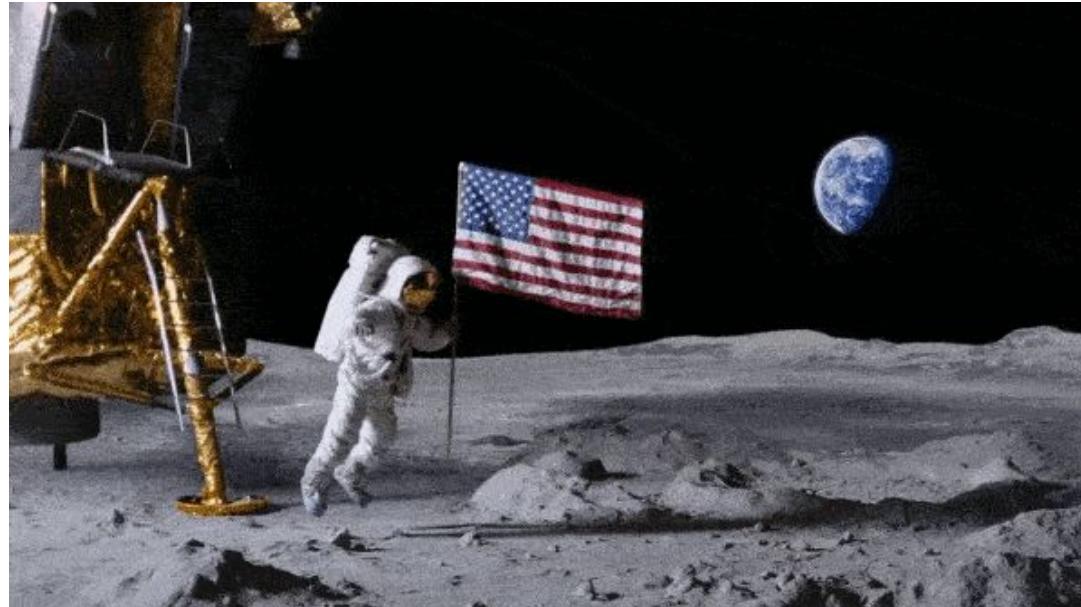
# Structure of the course



# We're (again) making history.

This is the second edition of this class

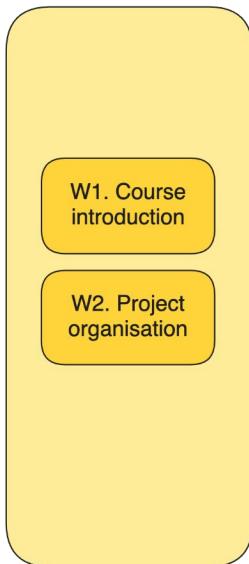
- Quick feedback cycles
- Open communication
- Enthusiasm for trying new things 
- Active support from teaching staff



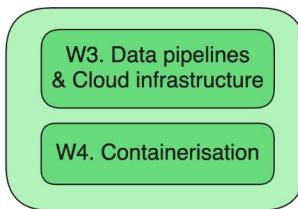
# Course outline

## Overview of sprints & classes

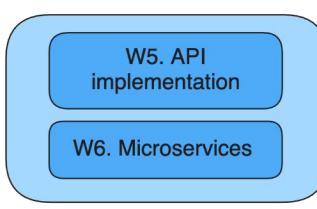
Sprint 1:  
Project organisation



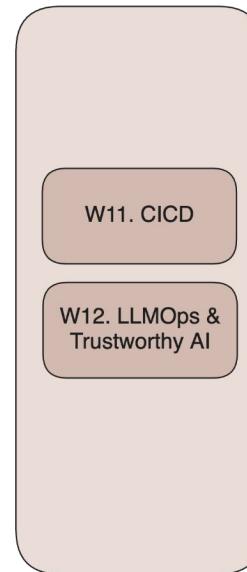
Sprint 2:  
Cloud & containerisation



Sprint 3:  
API implementation



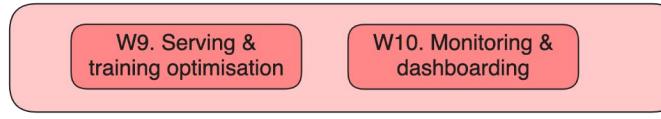
Sprint 6:  
CICD



Sprint 4: Model deployment



Sprint 5: Optimisation & monitoring



# Overall organisation & communication

## Class organisation

- We meet every Monday from **9:00** to **12:30**
- Typically you'll have about 2h of lecture + labs. Remaining of the time can be spent working on your project.

## Useful links:

- All info on the Github page: <https://github.com/ThomasVrancken/info9023-mlops>
  - Project info
  - Sample exam
  - Lecture & labs (before the class)
- Discord: <https://discord.gg/kY6B3cchkr>
- Open office hours on **Monday afternoons** (office Number I 77 B in Montefiore)

# Project

## Organisation

Build one ML system throughout the course. The application is picked by yourself.

- **Teams:** 2 - 4 students
  - Form group by next week!
  - Let the teaching staff know if you don't have a group and you'll be assigned one
- **Structure**
  - The building blocks to be implemented in the project follow the course's **6 sprints**.
- **Handovers**
  - There will be **3 milestone meetings** where you can present your results
  - **Code submission** - make sure to document clearly anything you want the teaching staff to read
- **Support**
  - Often lectures/labs will be shorter than the time slot for this course. You can spend the extra time working with your team. Teaching staff will be in the room to provide support.
  - Open office hours on Monday afternoon in office Number I 77 B in Montefiore
  - Feel free to reach out by email if you have any question/struggle
- **You're in the driving seat!**
  - Many building blocks are optional. You are free to choose the overall design and tools used for your project. Experiment and ask questions if you have any.

# Project

## Guiding principles

- Learn, learn and learn!
  - Find an interesting project to work on - ideally with a real world usage
  - Come up with your own design and toolstack
  - Focus on relevant parts of your specific system
  - Motivate your choices
- 
- ... And pick a cool name for it



# Example projects from last year

- Hessian: <https://github.com/alexandre-eymael/HESSIAN>
- ClipMorph: <https://github.com/iSach/clipmorph>
- ...



# Project objective for sprint 1

| Week | Work package   | Requirement |
|------|--|-------------|
| W01  | <p>Pick a <b>team</b></p> <ul style="list-style-type: none"><li>• Try to mix skills and experience</li><li>• If you didn't find one let one of the teachers know and we'll allocate you to one</li></ul>   | Required    |
| W02  | <p>Select a <b>use case</b></p> <ul style="list-style-type: none"><li>• Previous course</li><li>• <a href="#">Kaggle Datasets</a></li><li>• ...</li></ul> <p>Make sure to pick a use case where <b>data is available</b>. Ideally pick something with interesting data and a real world application.</p> | Required    |
| W02  | <p><b>Define</b> your use case. Fill in a ML Canvas <a href="#">template page</a> (You can skip the <i>Inference</i> part as we will tackle that in a later sprint.)</p>   | Required    |
| W02  | <p>Find a <b>cool name</b> for your project ✨</p>  | Required    |
| W02  | <p>Submit your project by sending a filled in <a href="#">project card</a> to the teaching staff with basic information about your project. We might give you some feedback and ask for parts to be changed.</p>   | Required    |
| W02  | <p>Setup <b>communication channel</b> (Discord, trello)</p>  | Required    |
| W02  | <p>Setup a <b>code versioning repository</b></p> <ul style="list-style-type: none"><li>• We recommend Github as we will cover Github Actions during this course</li></ul>  | Required    |

# Resources

## Similar courses

- University of Bari
  - Paper: "[Teaching MLOps in Higher Education through Project-Based Learning.](#)" arXiv preprint arXiv:2302.01048 (2023)
  - Lanubile, Filippo, Silverio Martínez-Fernández, and Luigi Quaranta
- Stanford University
  - CS 329S: Machine Learning Systems Design ([link](#))
  - Chip Huyen
- Carnegie-Mellon University
  - Machine Learning in Production / AI Engineering ([link](#))
  - Christian Kästner

## Interesting resources

- [Machine Learning Engineering for Production \(MLOps Specialization\)](#) (Coursera, Andrew Ng)
  - [GitHub](#), [Youtube](#)
- Made with ML ([link](#))
- Marvelous MLOps ([link](#))

## Books

- Designing Machine Learning Systems: An Iterative Process for Production-Ready Applications (Chip Huyen)
- Building Machine Learning Powered Applications: Going from Idea to Product (Emmanuel Ameisen)
- Introducing MLOps (Mark Treveil, Nicolas Omont, Clément Stenac et al.)
- Machine Learning Design Patterns (Valliappa Lakshmanan, Sara Robinson, Michael Munn)

**That's it for today!**

