

---

# Project organization

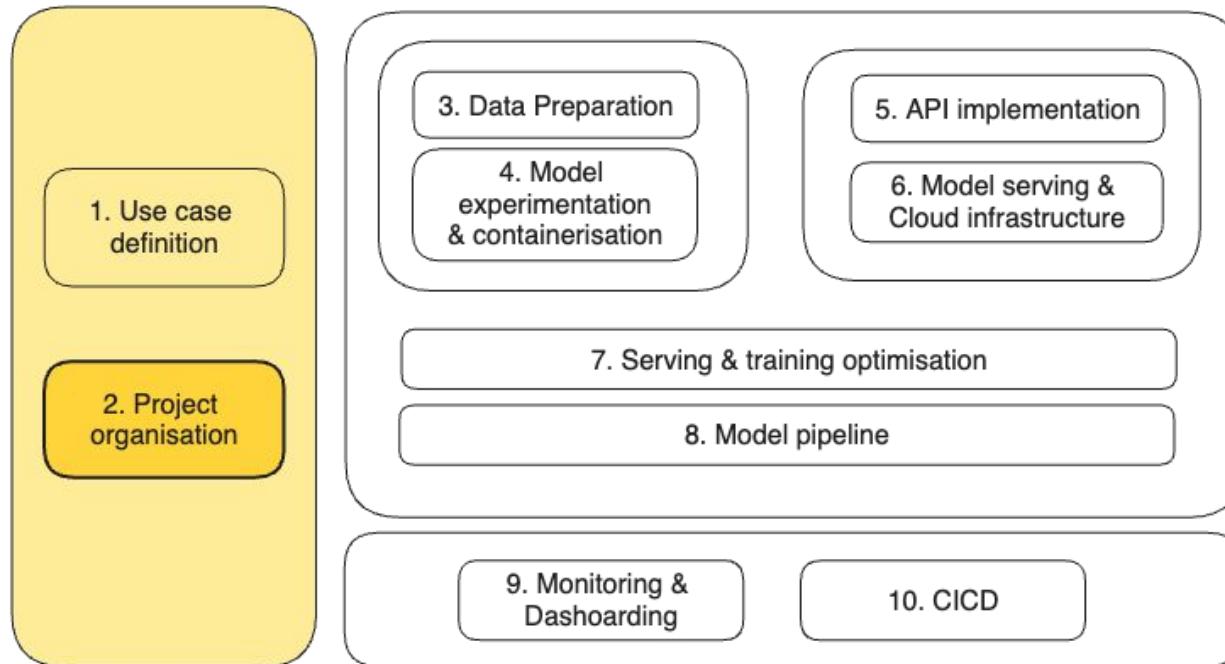
Sprint 1 - Week 2

INFO 9023 - Machine Learning Systems Design

2024 H1

Thomas Vrancken ([t.vrancken@uliege.be](mailto:t.vrancken@uliege.be))  
Matthias Pirlet ([matthias.pirlet@uliege.be](mailto:matthias.pirlet@uliege.be))

# Status on our overall course roadmap



# Hard to organise an IT project...



# Agenda

## What will we talk about today

### Lecture

- Agile way of working
- Data sources
- Code versioning & conventions

### Lab

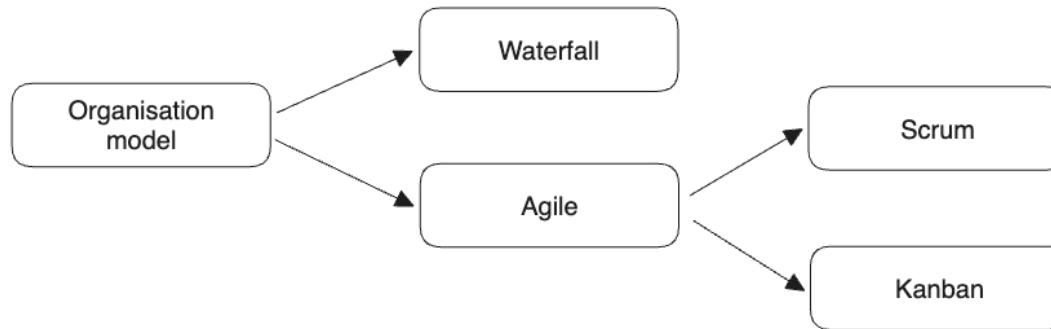
- Git code versioning

---

# Agile way of working



# Different models for organising an IT project

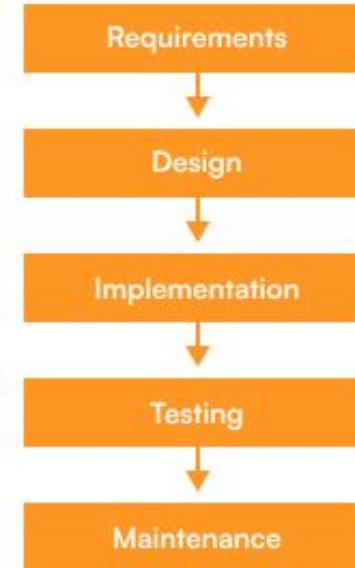


# Waterfall model

## Models for organising IT projects and teams

Organise a project as a **waterfall**, flowing down **sequentially** through 5 phases. It is **rigid, structured and linear** (it is not Agile).

Decision maker team sets a detailed plan for the whole project implementation, with limited flexibility.



Waterfall methodology

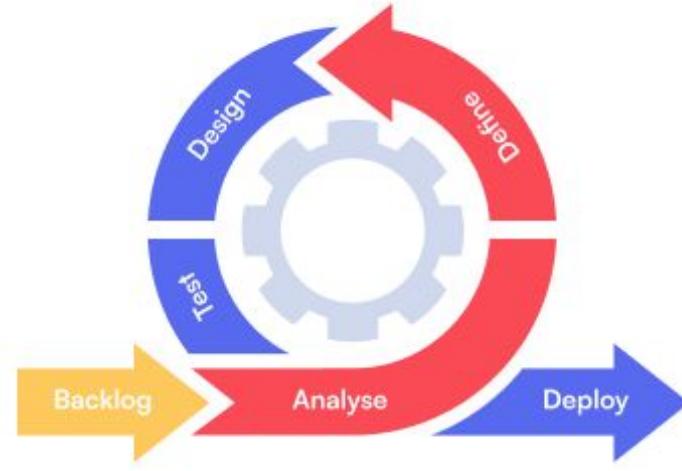
# Agile model

## Models for organising IT projects and teams

Iterative approach to delivering a project, which focuses on **continuous releases** that incorporate **customer feedback**.

Ability to **adjust** during each iteration promotes **velocity** and **adaptability**.

This approach is different from a linear, waterfall project management approach, which follows a set path with limited deviation.



Agile methodology

# Waterfall vs Agile

Criteria	Agile	Waterfall
<b>Approach</b>	Iterative and Incremental	Sequential
<b>Flexibility</b>	Highly flexible and adaptable	Less flexible
<b>Planning</b>	<del>Less</del> Minimal planning is enough	Detailed planning required
<b>Customer Involvement</b>	High level of customer involvement	Low level of customer involvement
<b>Risk management</b>	Continuous risk management	Risk management at the beginning of the project
<b>Documentation</b>	<del>Less</del> Minimal documentation required	Extensive documentation required
<b>Time and cost</b>	Challenging to estimate time and cost	Easy to estimate time and cost

# Agile: Scrum vs Kanban



# Kanban methodology

*In Japanese, kanban literally translates to "visual signal."*

- Originally comes from *lean manufacturing* in Japan (at Toyota)
- Simple and clean application of Agile
- Matches the amount of Work in Progress (WIP) to the team's capacity - *Just In Time (JIT)*
- Minimal set of rules - easy to pick up for new software engineering teams

## Advantages

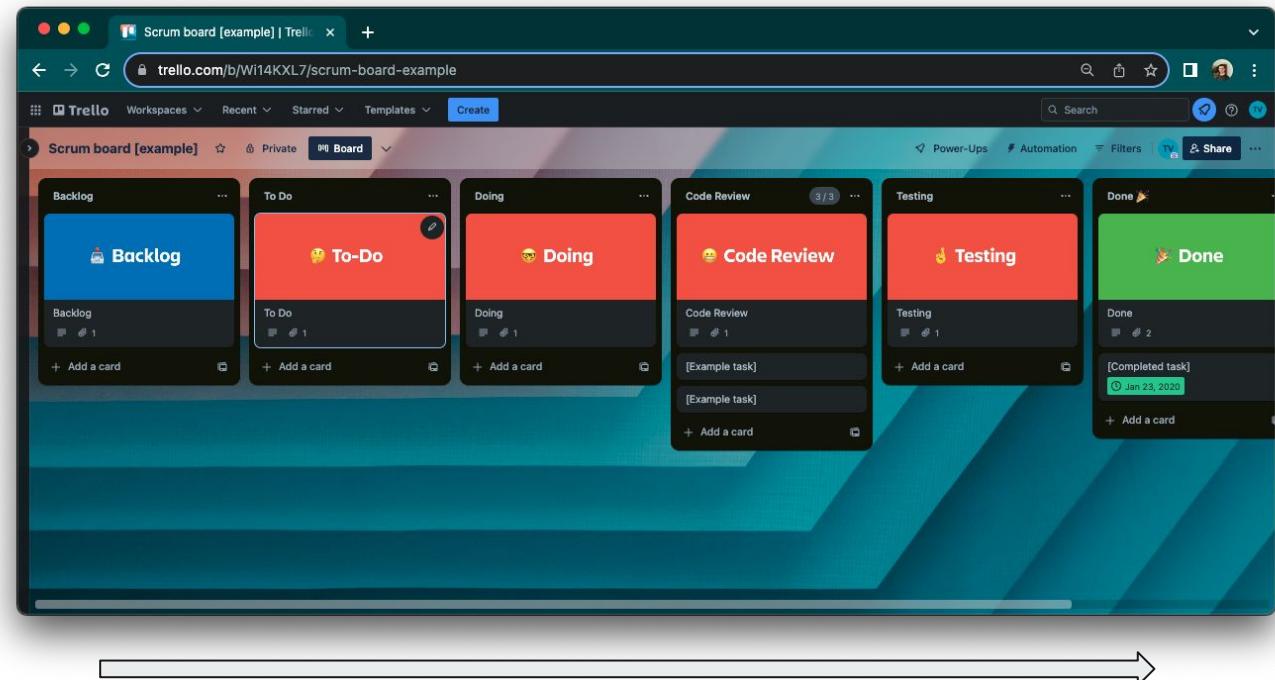
- + Planning flexibility
  - Once done with an item, pick up the next on the backlog
  - Product owner can update the backlog without disrupting the team
- + Short time cycles
  - Time for a unit of work to cycle through the whole process
- + Fewer bottlenecks
- + Visual metrics, transparency

# Kanban board

The **kanban board** is filled with **kanban cards** (aka tickets or tasks).

New tasks are added to the **backlog**.

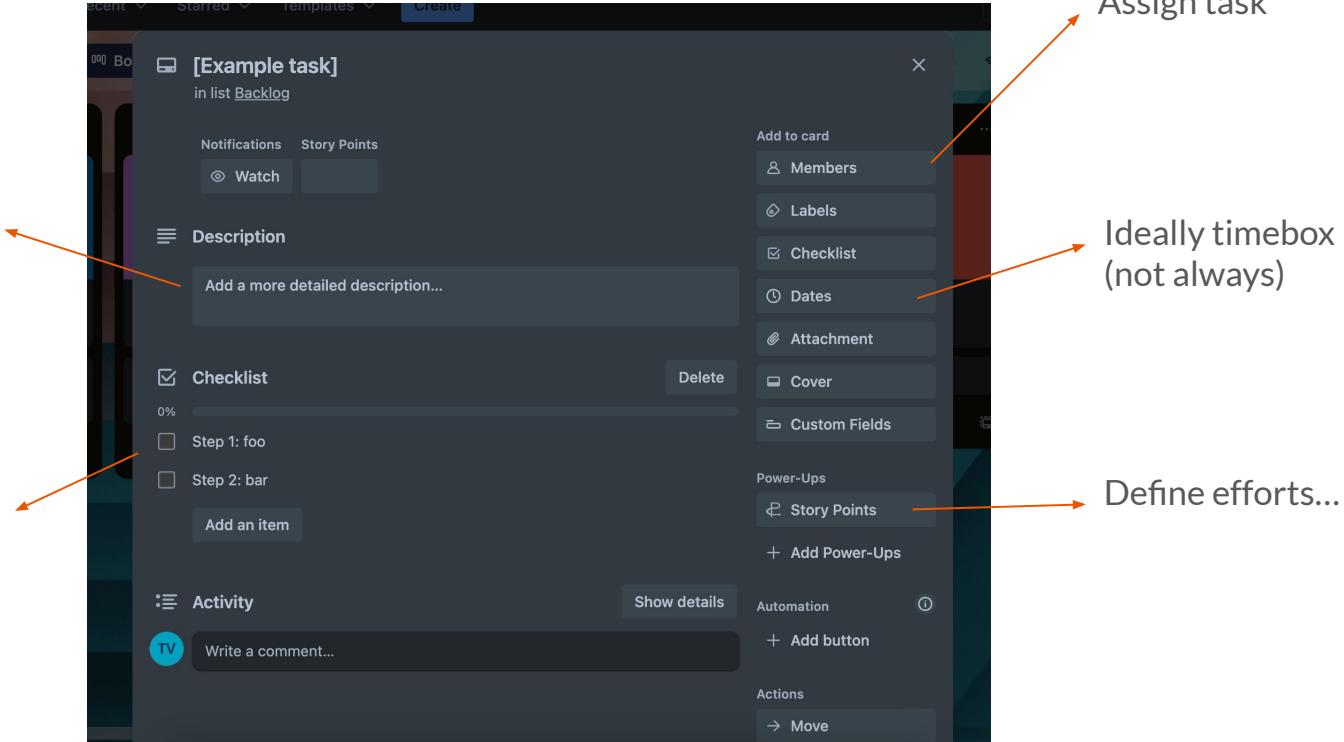
They are then gradually **moved** along the board.



# Kanban card

Define task

Include  
definition of  
done!



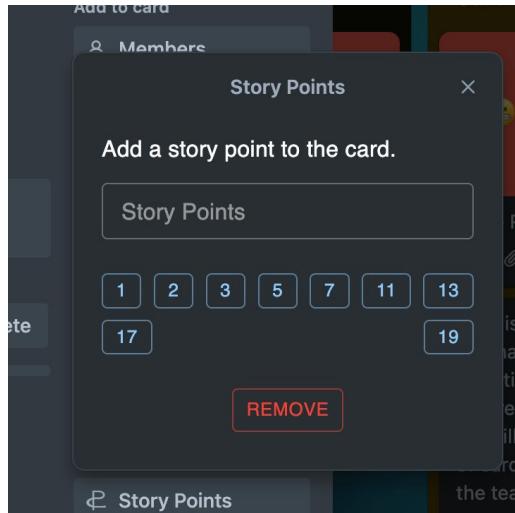
Break up tasks  
into steps

Assign task

Ideally timebox  
(not always)

Define efforts...

# Kanban card



**Story points** are a measure of effort and complexity of a task.

Follows the fibonacci scale.

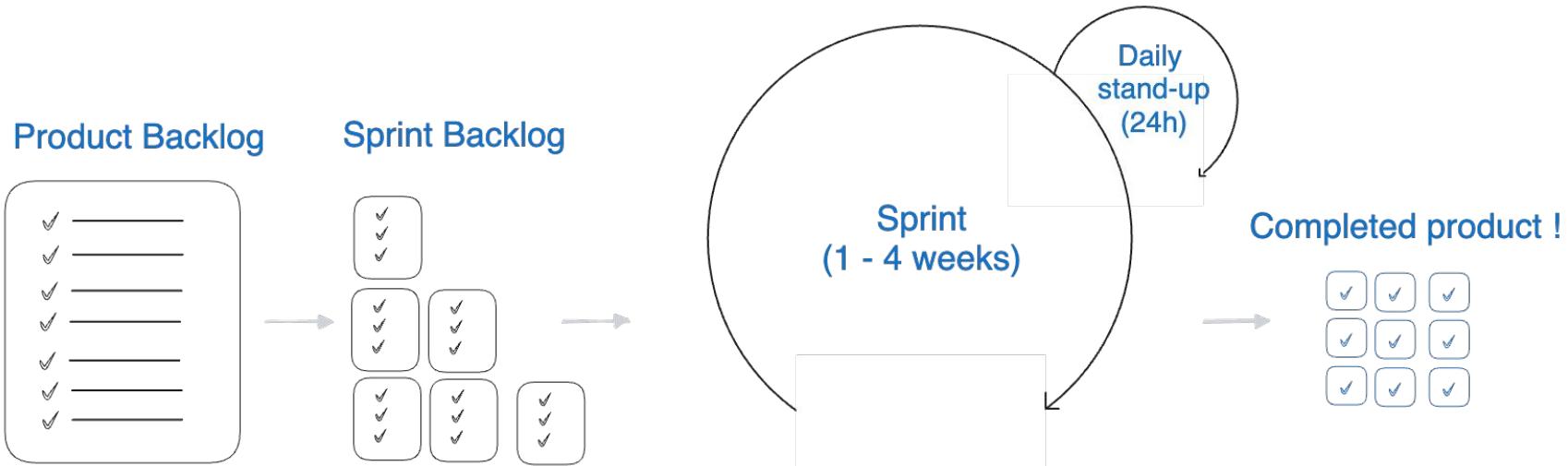
Sometimes voted on (e.g.

<https://www.pointingpoker.com/>)

If you want more information...



# Scrum methodology

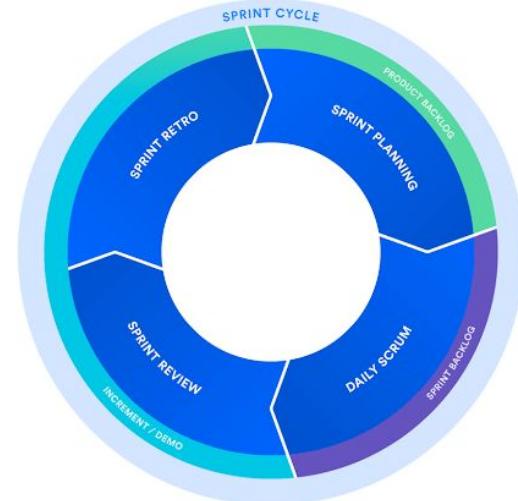


# Scrum methodology

## Sprint organisation

Scrum is organised around **1-4 weeks sprint cycles** (often 2 weeks).

- **Sprint planning:** Before each sprint, plan and add tasks from the product backlog to the sprint backlog.
- **Daily stand-up:** (Aka daily scrum) During the sprint, have dailies to update the **scrum board**, similar to kanban board. Go over each team member's progress.
- **Sprint review:** Casually **present** the work that was done in the last sprint to the rest of the team (definition, celebration and transparency).
- **Sprint retrospective:** **Feedback** on what went well and what can be improved regarding last sprint.



# Why are “daily stand-up” meetings called that way?

# Why are “daily stand-up” meetings called that way?

Goal is to make the meeting efficient.

All team members have to answer a fixed set of questions such as:

- What did you do yesterday?
- What did you do today?
- What, if anything, is blocking your progress?

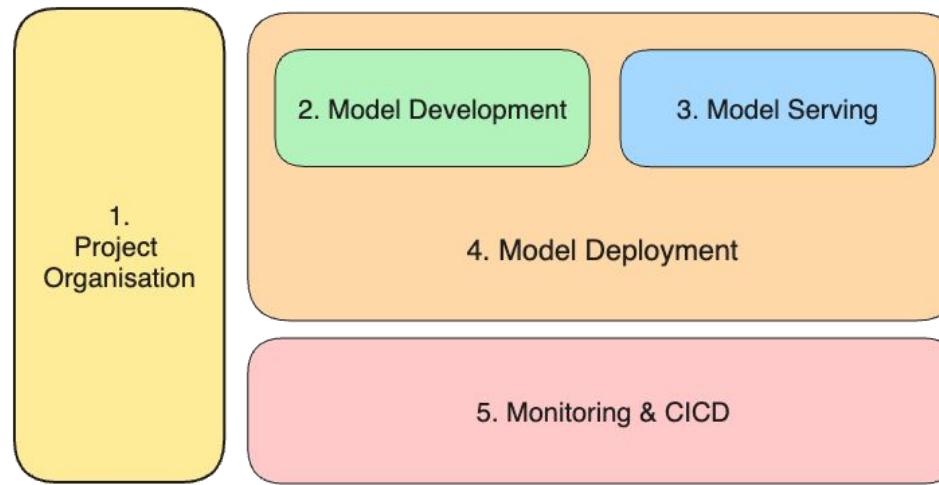


# The scrum team

## Roles & responsibilities

- **Product Owner:** Represents stakeholders' interests. Selects what is relevant to work on considering product objectives. Define user stories, priorities backlog and decide on product's direction.
- **Scrum Master:** Coach for the team. Ensures the best following of the scrum framework. Maintains the scrum board and backlog, facilitates and leads meetings and addresses obstacles.
- **Scrum Development Team:** Cross-functional developer team (data scientists, ML Engineers, data engineers, DevOps, front-end engineers, ...). Delivers a quality product.

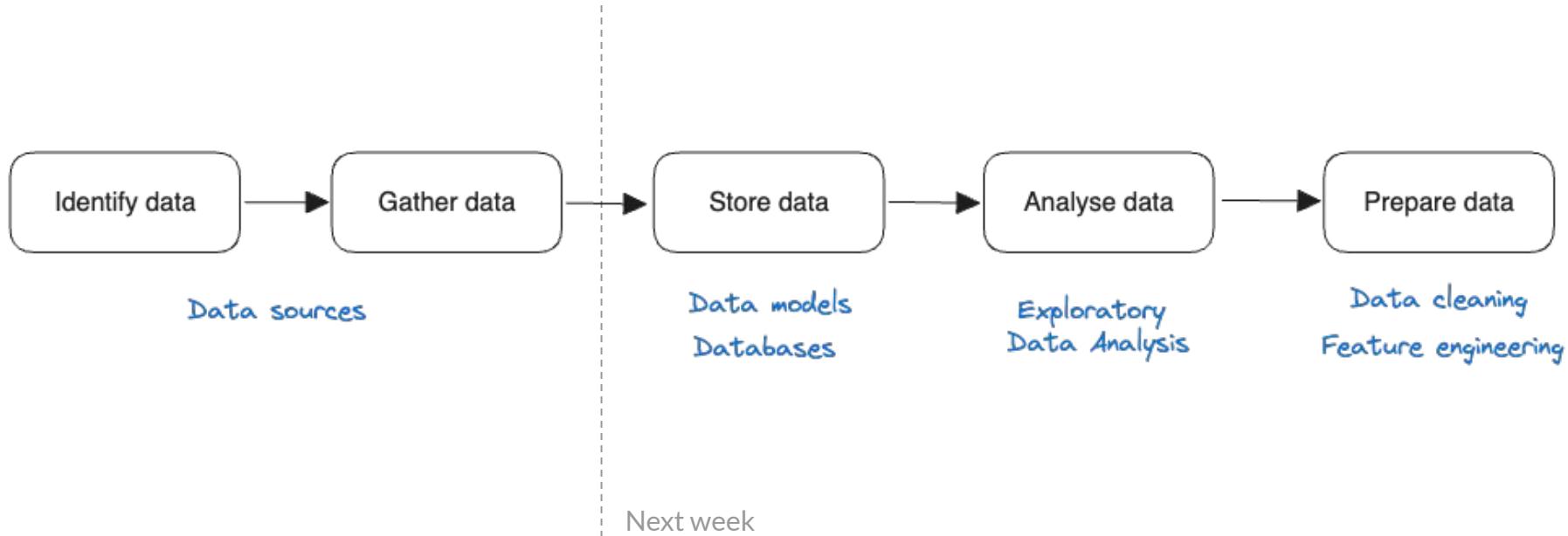
... Sounds kinda familiar...



---

# Data sources

# Overall data process



# Data source

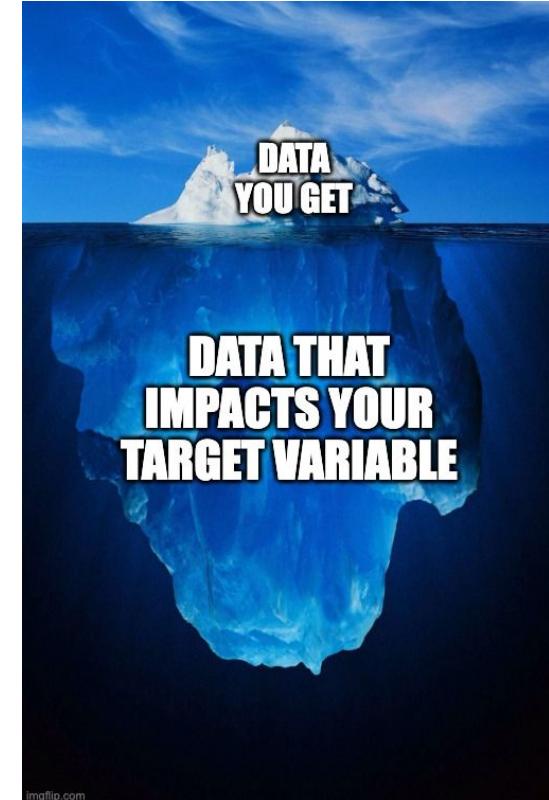
Really important to make sure you understand **all factors** that could be impacting your system.

Spend time discussing with **process experts** and challenge them on what input your model should use.

*Example:*

- Demand forecasting

**What data would you need to predict demand?**



imgflip.com

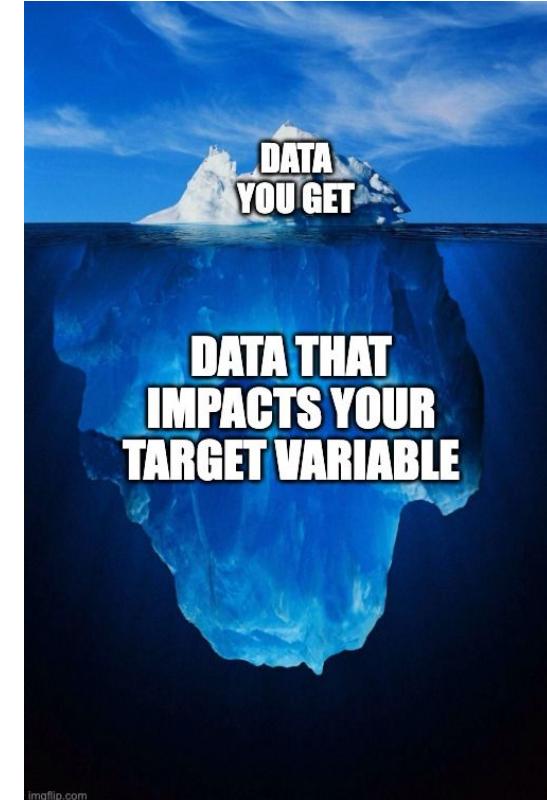
# Data source

Really important to make sure you understand **all factors** that could be impacting your system.

Spend time discussing with **process experts** and challenge them on what input your model should use.

*Example:*

- Demand forecasting
  - It will vary a *lot* depending on industry
    - Fashion → *trends*
    - Ice cream → *weather*
    - Stocks → *macroeconomics*
    - ...
  - Ask sales people how they estimate demand.



# Data source

Split in two dimensions

*Internal data already belong to the organisation building the model.*



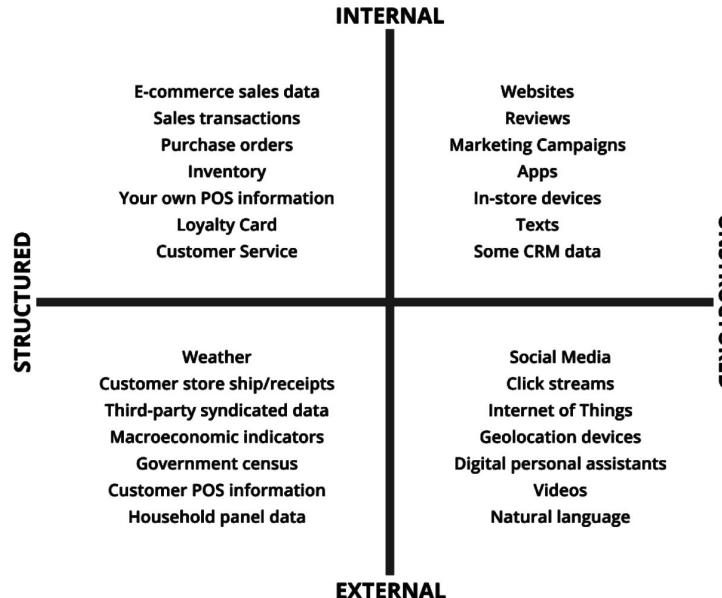
*External data needs to be gathered, scraped or purchased externally.*

# Data source

Split in two dimensions

## Structured data

- Usually tabular organised data
- Schema clearly defined
- Easy to search and analyse
- Schema changes will cause a lot of troubles
- Expensive to setup and store
- Stored in a data warehouse



## Unstructured data

- Data doesn't have to follow a schema
- All non-tabular data (computer vision, NLP, audio, video, ...)
- Fast arrival
- Cheap to store large volumes
- Can handle data from any source
- Harder to extract information
- Stored in a data lake



# Data is the new oil...

NEWS

## This is how much money Facebook earns from your data each year

Facebook may be free to use, but the trade-off isn't just annoying adverts: your data is being harvested and sold for way more money than you might imagine.



By [Jim Martin](#)

Executive Editor, Tech Advisor | JAN 28, 2022 1:30 AM GMT

With 42 million Facebook users, Facebook makes around \$37 billion per year from its UK users alone.

# Labeled data

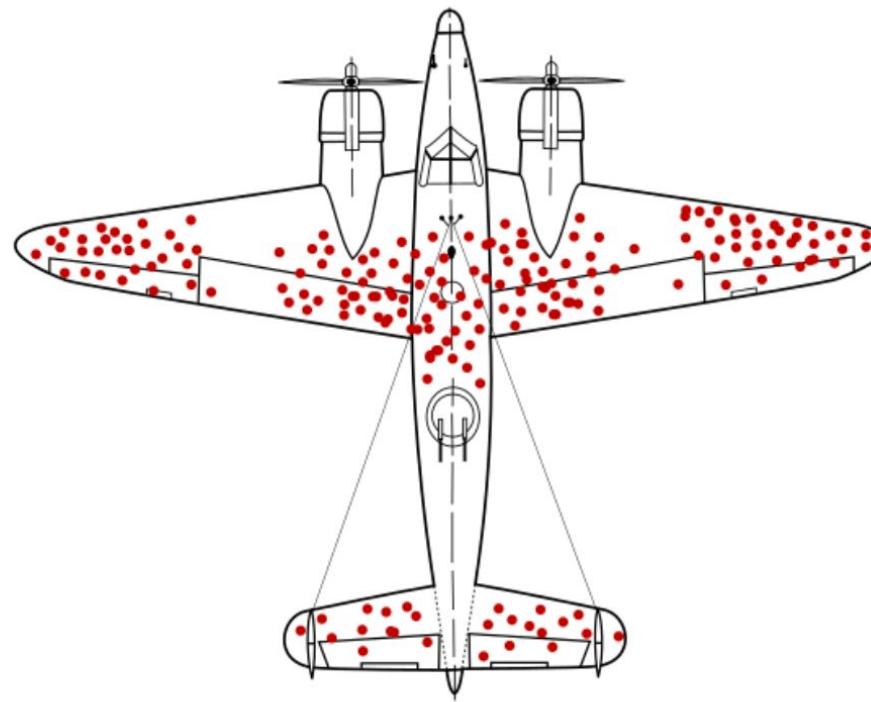
Labeled data is made of observations with a **label** value for what you are trying to predict.

You need

- Enough labeled data
- Good quality labels
- Truly representing what you are trying to model.

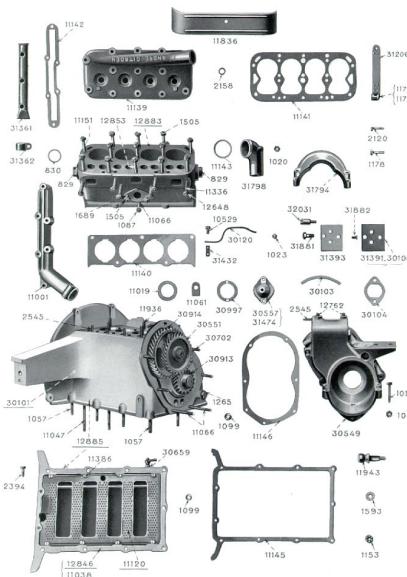
Rarely present in the real world...

# Careful to bias when collecting data



# Setting up a data collection and/or labelling system

## Example 1: Building a physical system



**Context:** Company selling hundreds of different *engine parts*.

**Goal:** Automatically classify physical parts using computer vision.

**How would you approach it?**

# Setting up a data collection and/or labelling system

## Example 1: Building a physical system



**Context:** Company selling hundreds of different *engine parts*.

**Goal:** Automatically classify physical parts using computer vision.

**Approach:**

1. Build a box with clear lighting and multiple cameras
2. Build a UI to easily take a snap picture and enter product code
3. Each new and returning products had to go through it for a summer

**Result:**

>100k product images.

# Setting up a data collection and/or labelling system

Example 2: Build a labeling tool and hire a 3rd party company (or student workers)

Airplane<sup>[1]</sup> Car<sup>[2]</sup>

Results

Nothing selected

Regions (2) Sort ↴

1 Rectangle 155.00 x 155.00  
2 Rectangle 37.00 x 25.00

Relations (0)

No Relations added yet

To have faith is to trust yourself to the water

Positive<sup>[1]</sup>  
 Negative<sup>[2]</sup>  
 Neutral<sup>[3]</sup>

Time Series classification

Growth<sup>[1]</sup>  
 Decay<sup>[2]</sup>

Run<sup>[3]</sup> Walk<sup>[4]</sup>

Results

Nothing selected

Regions (1) Sort ↴

1 TS 13 - 19

Relations (0)

No Relations added yet



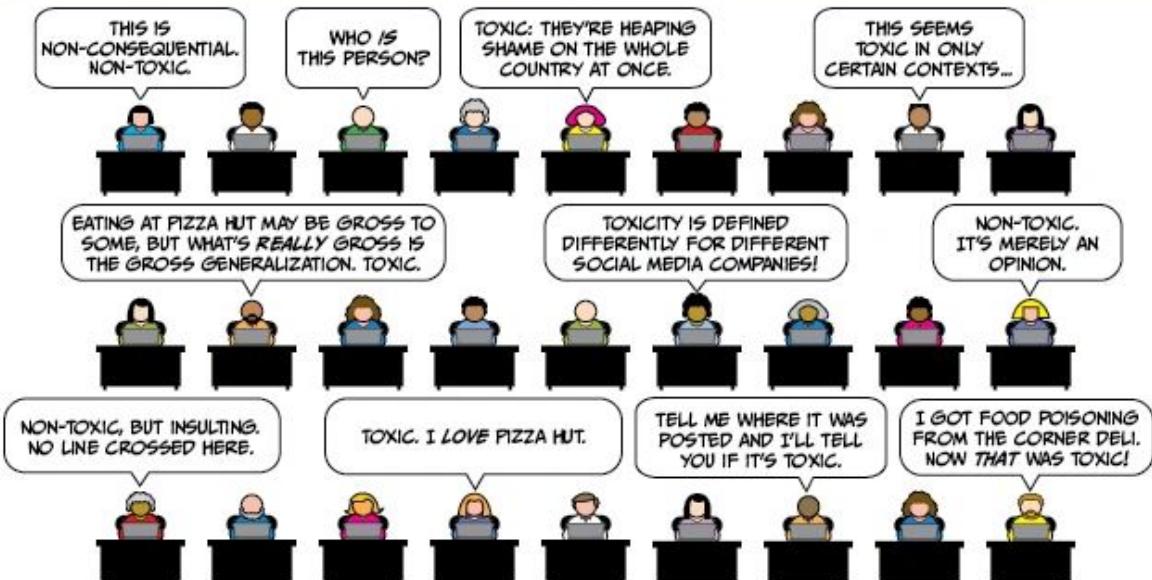
## Label Studio

# Quality challenge in manual data labelling

Beware of bias and data quality when collecting the data! Surprisingly hard to get a lot of people to label in a consistent way - inherent bias...

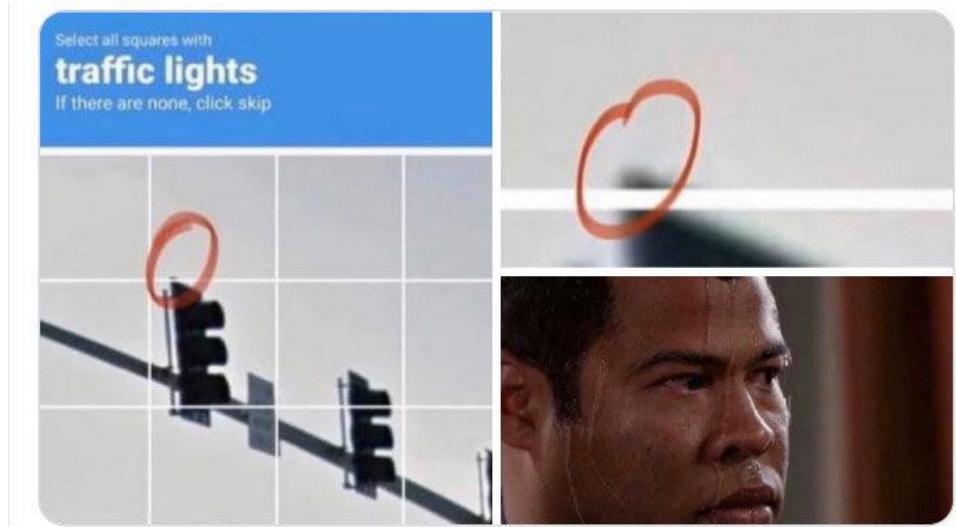
Task: Classify the following online comment as “toxic” or “not toxic.”  
Comment: “1. People still eat at Pizza Hut? Gross. 2. It is shameful how this country[...].”

Toxic  Not toxic



# Quality challenge in manual data labelling

Beware of bias and data quality when collecting the data! Surprisingly hard to get a lot of people to label in a consistent way - inherent bias...



# Quality challenge in manual data labelling

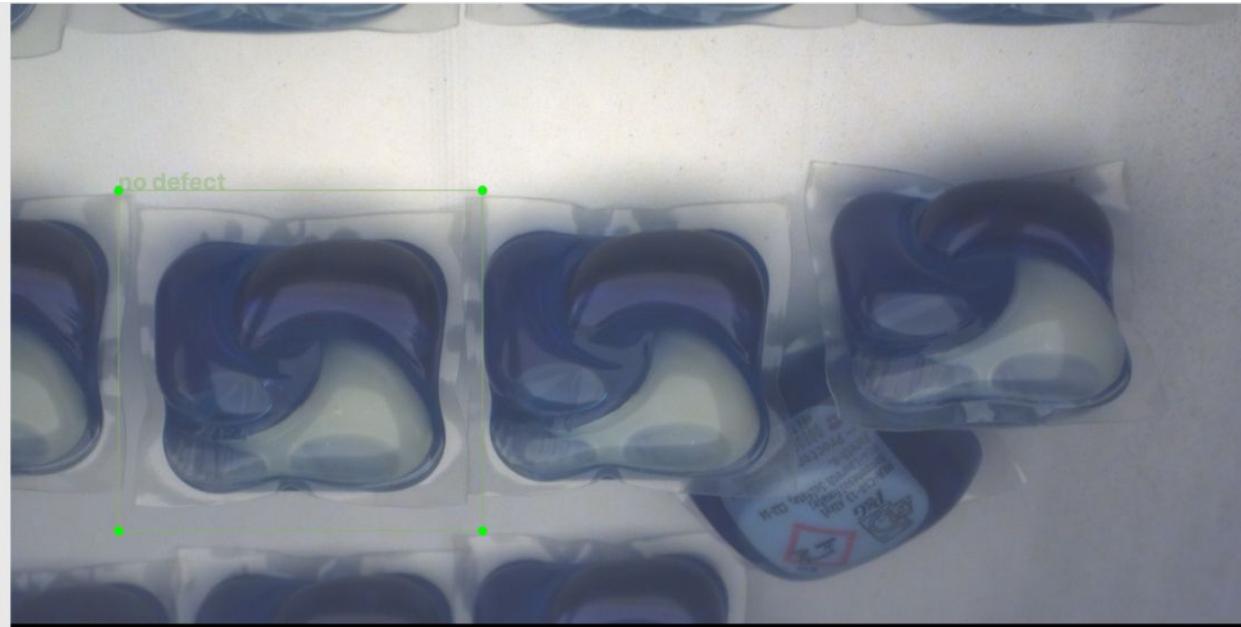
Beware of bias and data quality when collecting the data! Surprisingly hard to get a lot of people to label in a consistent way - inherent bias...

Solutions:

- Set strict **guidelines**
  - Explain in detail what should happen
  - Give a lot of examples
- Review sessions of difficult examples with different labelers
- Overlap samples which are then labeled by multiple labelers
  - Remove or average observations with different labels
- Calculate bias
  - Standard deviation of labels
  - Distribution



# Be really careful with your labeling strategy as it's hard to come back from it...



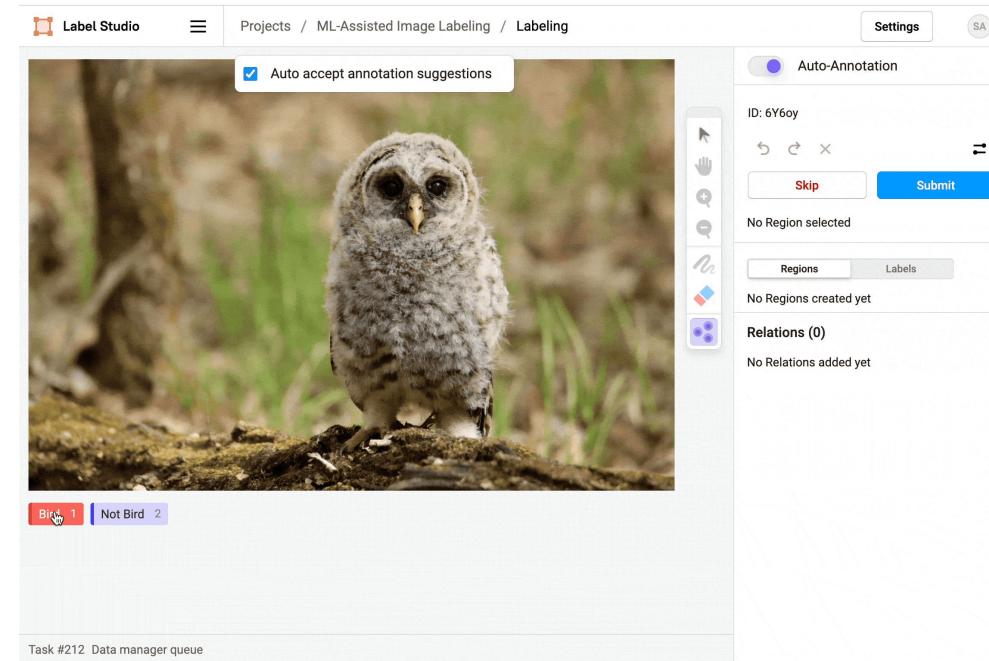
Detecting defect in dishwasher tab pouches.

Did not label overlapping pouches...

# Data pseudo-labelling

Use an unsupervised general model to produce some possible labels and get them verify.

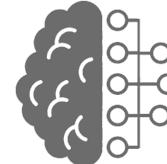
Can lead to overfitting on general model !



# Breakout exercise

Semantic Search works with **text embedding**. An ML model is trained to represent any text in a **vector space** where texts tackling similar topics will be close to each others.

Query:  
“How often should I  
feed my cat?”

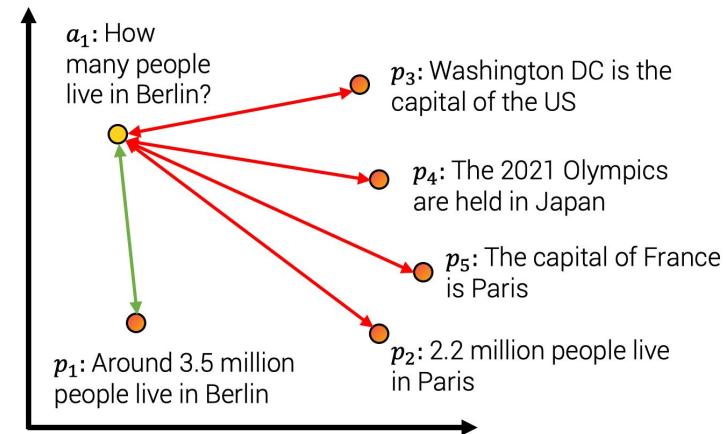


# Breakout exercise

You can find **pre-trained** embedding models (OpenAI, Huggingface, Cohere, ...).

You can **fine-tune** these models on your own data to make them better at understanding **domain specific language**.

What you need: A data set of **Question-Answer pairs**.



# Breakout exercise

## Exercise:

- Team of 3-4 students
- You got 10min

## Context:

- Assume you work for a **publishing company** that has **300k news articles** from the last 30 years
- They want to build a **semantic search engine**
- The documents use domain specific language so you want to **fine-tune a text embedding model**
- For that, you need a labeled data set of **Question and Answer pairs**
  - We want it to specifically use the **language** and **domain** of the newspaper publisher
  - The publisher only has documents

## Goal:

- Try to come up with a couple of strategies to **create a labeled dataset**

# Generative Pseudo Labeling

Loop over documents and ask an LLM to generate relevant questions

## Prompt 1: Generate questions

"""

Your task is to **write** 10 relevant questions on the following documents.  
Write it in {language} using expert terms on [{topic X, Y, ...}]

- Each question must relate to the text
- Each question must be different from other questions
- Each question must be short

Use a JSON format

[{ Question 1: answer}, ...]

Text: {document}

"""

## Prompt 2: Score questions

"""

Your task is to give a relevance **score** on how well these questions are answered by the following document.  
The score must be on a scale from 1 to 3.

- 1: The text does not answer the question at all
- 2: The text partially answers the question
- 3: The text fully answers the question

Use a JSON format

[{ Question 1: score}, ...]

Text: {document}

Questions: [...]

"""

---

# **Code versioning & conventions**

# Why does versioning matter?

Name	Date modified	Type	Size
final-draft	8/1/2017 3:48 PM	Microsoft Word D...	16 KB
final-draft2	8/3/2017 3:49 PM	Microsoft Word D...	12 KB
final-draft-final	8/3/2017 10:22 PM	Microsoft Word D...	13 KB
final-modified	8/8/2017 3:49 PM	Microsoft Word D...	13 KB
real-final-draft	8/9/2017 2:22 PM	Microsoft Word D...	13 KB
real-final-draft-v2	8/10/2017 12:33 PM	Microsoft Word D...	16 KB
final	8/11/2017 11:59 PM	Microsoft Word D...	72 KB



Julien Chaumond  
@julien\_c

Modern ML engineering is 90% knowledge of git

[Traduci il Tweet](#)

4:40 PM · 20 ago 2021 · Twitter for iPhone



# Versioning applied on different parts of your solution.



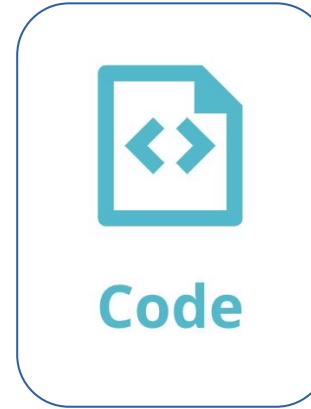
Data

+



Model

+



*Focus for today!*

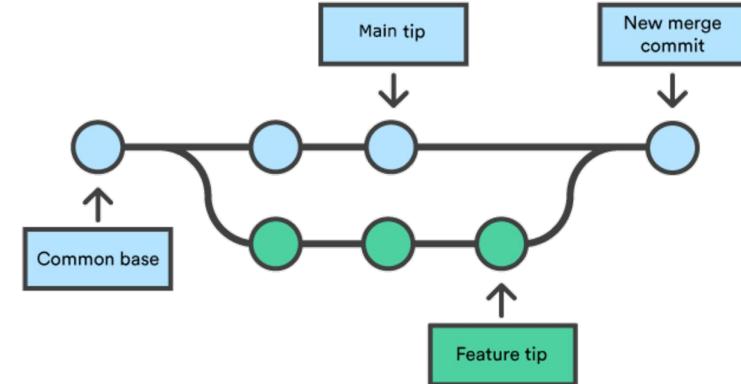
# Code versioning

- Tracking, organising and controlling changes in the code. Enables:
  - Collaboration
  - Track history
  - Revert
  - Versioning
- **Git** is the standard solution for code versioning
- Different Source Code Management (SCM) tools exist such as:
  - Github
  - Gitlab
  - Bitbucket
  - ...



# Git branching and merging

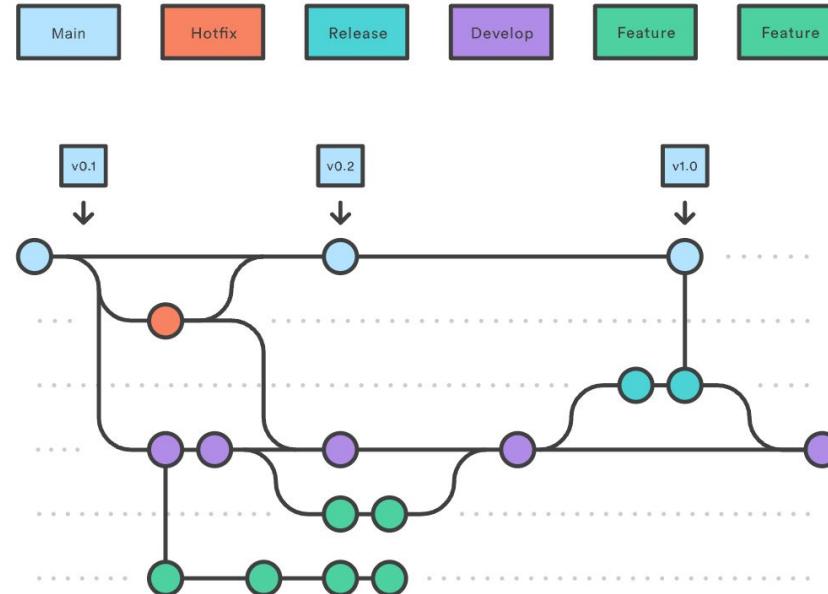
- **Goal:** work at the same time with multiple versions of a repository
- A branch is an independent line of development
  - As default, a repository has a unique branch called **main**
  - Another branch is used to work at changes and extensions before a merge with the main
- For each issue X (feature, bug, etc.):
  - **Create** a new branch X
  - **Implement** X
  - **Merge** X changes with the main



# GitFlow

It uses 4 types of primary branches

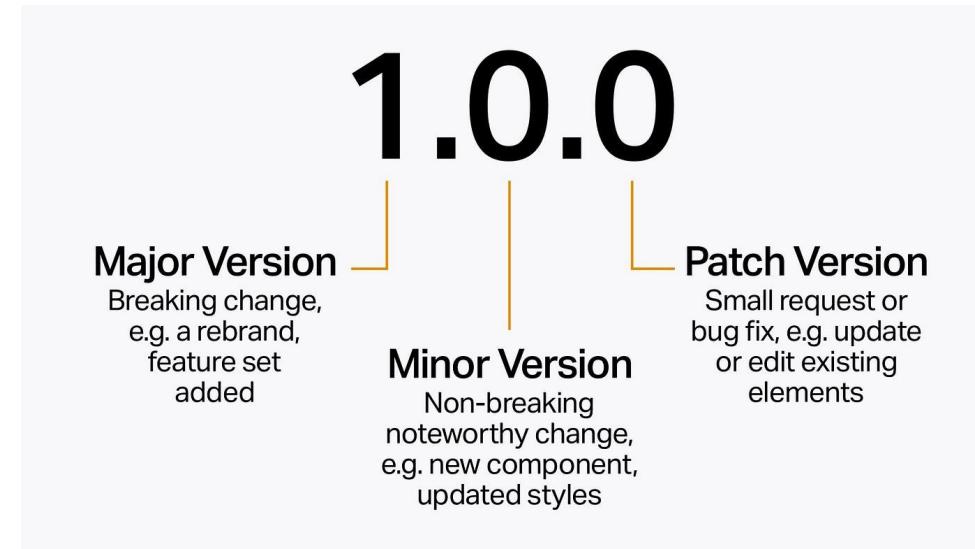
1. **Main** branch stores the official release history (commits with a **version number**)
2. **Develop** branch serves as working branch for new features
3. **Feature** branch gets created each time a specific new features needs to be implemented
4. **Release** cycle - only bug fixes and documentation added to this branch
5. **Hotfix** to quickly patch and maintain production release



<https://nvie.com/posts/a-successful-git-branching-model/>

<https://www.atlassian.com/git/tutorials/comparing-workflows/gitflow-workflow>

# Product versioning convention



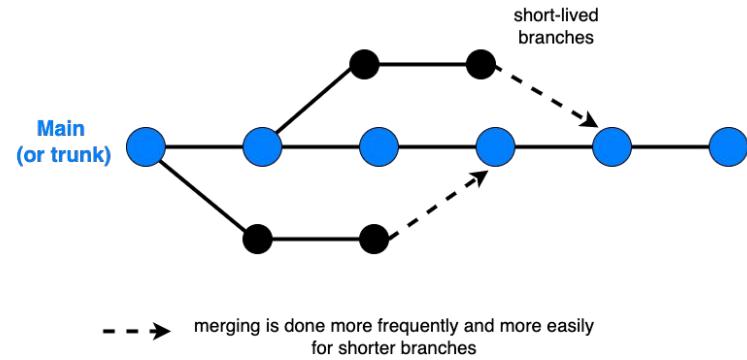
# Git Trunk-based development

- Developers merge **small, frequent** updates to a core “**trunk**” or main branch
- Based on automatic building and testing of the solution
- Requires CICD
- In trunk-based development, code review should be performed immediately (within a day)
- Somewhat aligned with Agile mindset

Some advocates name it as the new thing...

Trunk-based development

StatusNeo



# GitFlow vs Trunk based

GitFlow	Trunk based
Large feature branches	Small, frequent updates to a core “trunk”
Active collaboration	CICD
Multiple base branch	Automatic testing
Large PRs	Small and frequent PRs
Different phases of development (dev, release, hotfix, ...)	No planned code freeze or pause of integration
Still very popular	Documented as the next best thing



*Often, teams operate a bit in between. Traditional Gitflow framework while trying to work with as small branches as possible.*

# GitFlow vs Trunk based

GitFlow	
Large feature branches	Large feature branches
Active collaboration	Active collaboration
Multiple base branch	Multiple base branch
Large PRs	Large PRs
Different phases of development (hotfix, ...)	Pause of integration
Still very popular	First thing

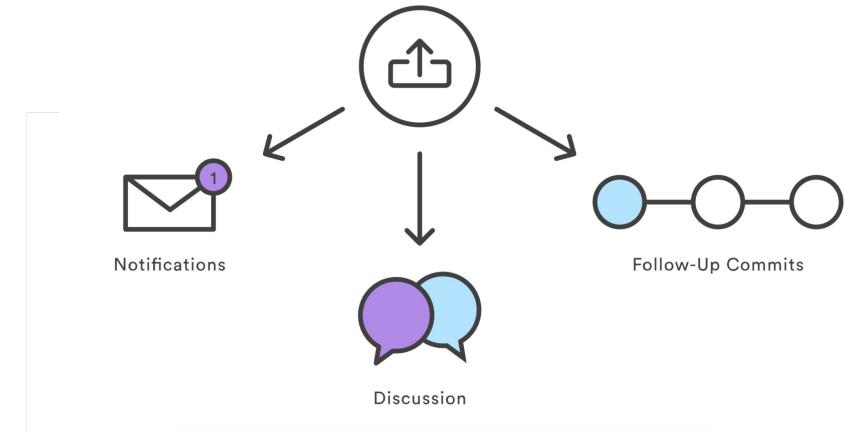
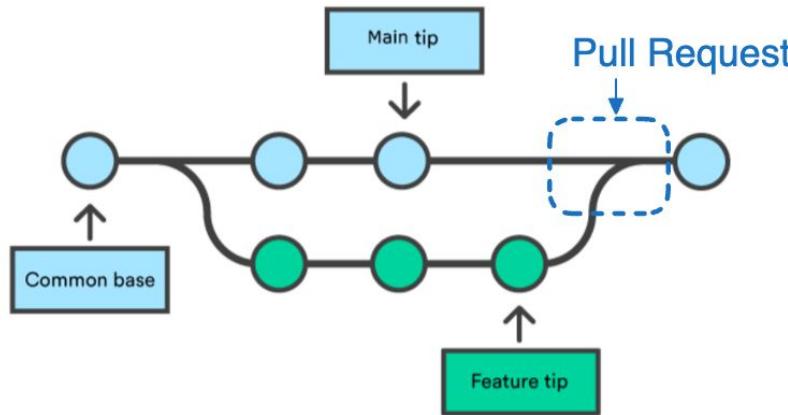
A cartoon illustration of SpongeBob SquarePants from the TV show 'SpongeBob SquarePants'. He is smiling broadly with his arms raised, holding a large, multi-colored rainbow that forms a semi-circle behind him. The rainbow has white sparkles along its edges. Below the rainbow, the text 'KEEP YOUR BRANCHES SHORT AND CONCISE' is written in large, bold, white capital letters. The background is a simple brown gradient.

↔ imgflip.com ↔

Often, teams operate a bit in between. Traditional Gitflow framework while trying to work with as small branches as possible.

# Pull requests

- **Pull Request (aka Merge Request)** ⇒ Process to merge branch into main
- Select members that need to **review the codes** before the branch is merged to main.



# Pull requests: Code review

## Why it is important

- Code reviews are super important:
  - **Quality Assurance:** Is this actually the best implementation?
  - **Knowledge sharing:** Aligned team on how the whole code base works.
  - **Coding standards:** Consistent over code base
  - **Mentorship:** Guidance for less experienced developers
  - **Documentation:** Guarantees that codes are readable
  - **Team morale:** We work as a group, less isolated features
- Tech-lead needs to **dedicate time to code review**

*Strict correlation between developer team maturity and enforcement of code reviews!*

*Be smart - shouldn't hinder team efficiency (especially in early stages of a project)\_*

# Pull requests: Code review

## Best practices

1. Instead of giving orders, make **suggestions**
2. Start with the “**why**”
3. Balance **appreciative** and **constructive** feedback
  - o (Always say something positive)
4. Don’t be spiteful over **mistakes** when reviewing code
5. Don’t be afraid to **comment** during code reviewing



# Merge conflicts

Sometimes multiple developers may try to edit the same content.

If Developer A tries to edit code that Developer B is editing a conflict may occur.

The `git merge` command's primary responsibility is to combine separate branches and resolve any conflicting edits.

```
$ git status
On branch main
You have unmerged paths.
  (fix conflicts and run "git commit")
  (use "git merge --abort" to abort the merge)

Unmerged paths:
  (use "git add <file>..." to mark resolution)

    both modified:  merge.txt
```

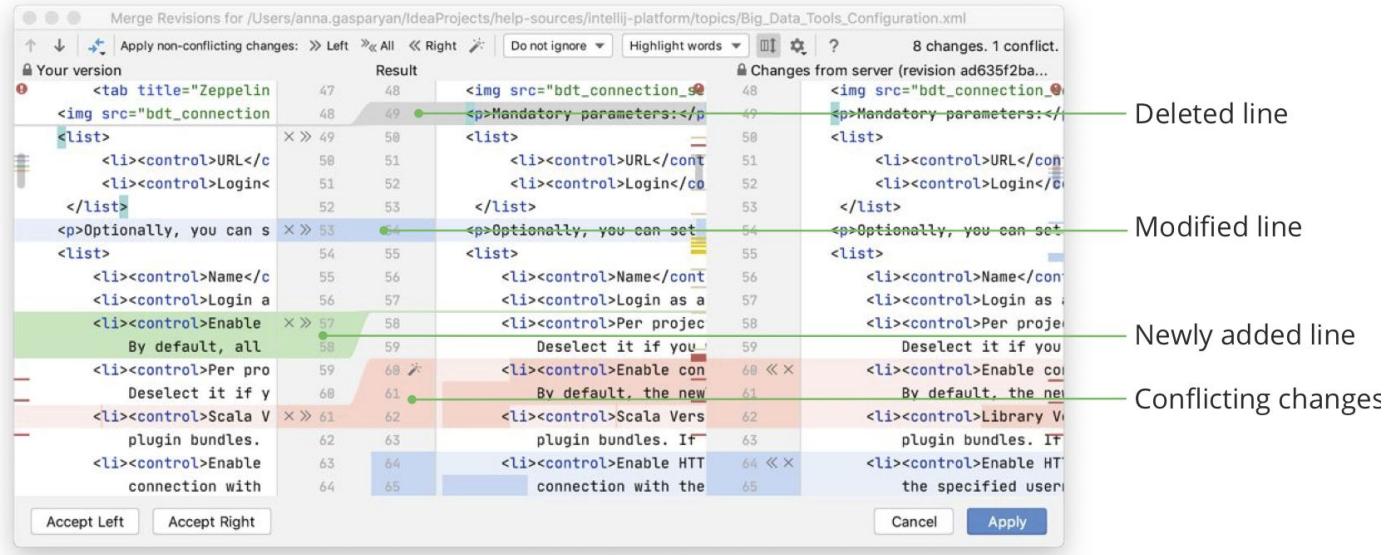
# Merge conflicts

If you tried to merge a branch into another one and experienced a merge conflict, it will then show in your local file.

```
$ cat merge.txt
<<<<< HEAD
this is some content to mess with
content to append
=====
totally different content to merge later
>>>>> new_branch_to_merge_later
```

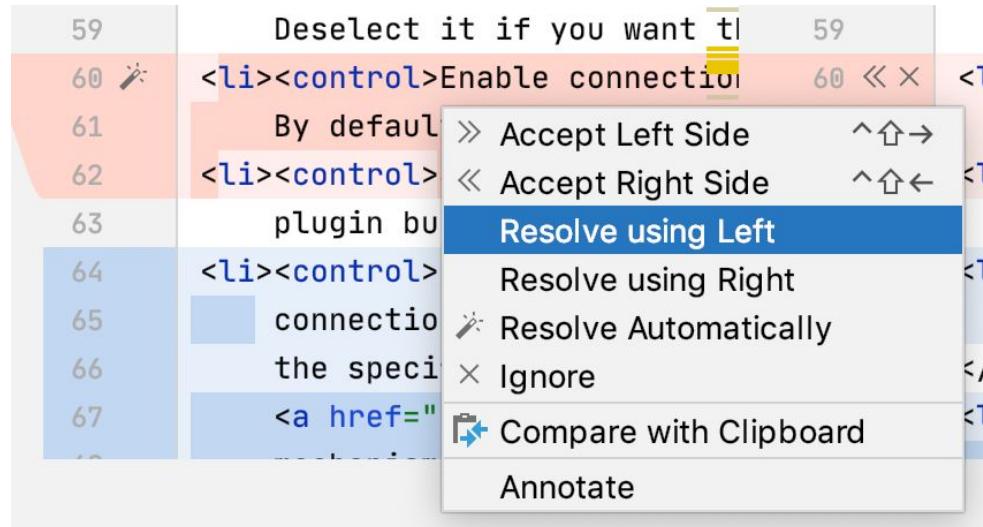
# Merge conflicts

IDEs (PyCharm, VS Code, ...) can help with fixing merge conflicts.



# Merge conflicts

IDEs (PyCharm, VS Code, ...) can help with fixing merge conflicts.



# Importance of documentation

## Docstrings

Used to document your function. “Would I be comfortable that someone else can easily reuse this code?

Describe the general **functionality**, the input **arguments** and what it **returns**. Example styles:

```
"""
This is a reST style.

:param param1: this is a first param
:param param2: this is a second param
:returns: this is a description of what is returned
:raises KeyError: raises an exception
"""
```

[reStructuredText](#) (reST)

```
"""
This is an example of Google style.

Args:
    param1: This is the first param.
    param2: This is a second param.

Returns:
    This is a description of what is returned.

Raises:
    KeyError: Raises an exception.
"""
```

[Google](#)

# Importance

## Docstrings

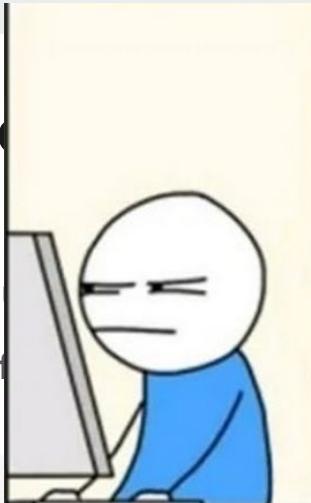
Used to document yo

Describe the general f

```
....  
This is a reST style.
```

```
:param param1: this is a firs  
:param param2: this is a seco  
:returns: this is a descripti  
:raises KeyError: raises an e  
....
```

[reStructuredText](#) (



**When I wrote this code,  
only God & I understood  
what it did.**



use this code?

es:

**Now.....  
only God knows.**

# What is wrong with the docstring in this code snippet?



Pathetic.  
Imgflip.com

```
def calculate_area(length, width):
    '''Calculates the area of a rectangle. This method takes two arguments, the length and
width of the rectangle, and returns the area.
    Arguments:
        length - the length of the rectangle.
        width - the width of the rectangle.
    Returns: The area of the rectangle calculated using the formula length*width.'''
    return length * width
```

# Importance of documentation

## Markdowns

Used as general documentation (e.g. README.md)

Important to clearly document the **structure** of your repository and how to **use it**.

**Rubber ducking:** Tell yourself what exactly you're doing here. Good way to voice a problem. (Usually for debugging but also for documentation).



```
# Foobar
Foobar is a Python library for dealing with word pluralization.

## Installation

Use the package manager [pip](https://pip.pypa.io/en/stable/) to install foobar.

```bash
pip install foobar
```

## Usage

```python
import foobar

# returns 'words'
foobar.pluralize('word')

# returns 'geese'
foobar.pluralize('goose')
# returns 'phenomenon'
foobar.singularize('phenomena')
```

```

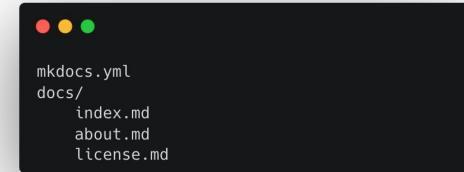
# Importance of documentation

## Documentation page

You can use tools such as **MkDocs** or **docusaurus** to produce a full documentation page from your directory.

It takes in a certain structure of documentation and will generate an HTML page.

Identifies docstrings and other types of documentation directly from your codebase.



The screenshot shows the generated documentation site for the 'MkDocs' project. The top navigation bar includes a search bar, a 'Docs' link, and a 'Home' link. On the right, there are 'Edit on GitHub' and 'Next >' buttons. The main content area has a dark header with the title 'MkDocs' and a subtitle 'Project documentation with Markdown.' Below this, there is a section titled 'Features' with sub-sections 'Great themes available' and 'Easy to customize'. A sidebar on the left contains a navigation menu with links like 'Getting Started', 'USER GUIDE', 'Overview', 'Installation', 'Writing Your Docs', 'Choosing Your Theme', 'Customizing Your Theme', 'Localizing Your Theme', 'Configuration', 'Deploying Your Docs', 'DEVELOPER GUIDE', 'Overview', 'Themes', 'Translations', 'Plugins', 'ABOUT', 'Release Notes', 'Contributing', and 'License'. At the bottom of the sidebar are 'GitHub' and 'Next >' buttons.

**MkDocs**

Project documentation with Markdown.

MkDocs is a fast, simple and downright gorgeous static site generator that's geared towards building project documentation. Documentation source files are written in Markdown, and configured with a single YAML configuration file. Start by reading the introductory tutorial, then check the User Guide for more information.

[Getting Started](#) [User Guide](#)

**Features**

**Great themes available**

There's a stack of good looking themes available for MkDocs. Choose between the built in themes: mkdocs and readthedocs, select one of the third-party themes listed on the MkDocs Themes wiki page, or build your own.

**Easy to customize**

Get your project documentation looking just the way you want it by customizing your theme and/or installing some plugins. Modify Markdown's behavior with Markdown extensions. Many configuration options are available.

# The .gitignore file

## Issue with Git

- Pain to erase memory
  - *Don't ever push any sensitive information on git!*
- Limited storage space (~2 GB)



# The .gitignore file

## Issue with Git

- Pain to erase memory
  - Don't ever push any sensitive information on git!
- Limited storage space (~2 GB)

## Solution

- .gitignore file!
- Automatically ignore some files
- Many examples: <https://github.com/github/gitignore/tree/main>

Remove files per extension



Remove files per location



Any of the characters in square brackets (.pyc, .pyo, .pyd)

```
# These are some examples of commonly ignored file patterns.

# Compiled Python bytecode
*.py[cod]

# Log files
*.log

# Environment
*.env
venv/

# Data
data/

# Secrets
secrets/

# Jupyter Notebook
.ipynb_checkpoints
```

# Scripting vs notebooks

Attributes of scripts:

- **Stateless**
  - Scripts do *not* save variables and what you run in the **global state**. Therefore they are **stateless** (unlike notebooks which are **stateful**)
  - Scripts require passing variables to functions and classes. Healthy to **structure** your implementation.
- **Linear**
  - Scripts run everything **from beginning to end linearly**, which takes away the risk of running cells in the wrong order
- **Reproducibility**
  - Easier to **reproduce** and **run tests** with different parameters with scripts.

"Only a Sith deals in absolutes"

⇒ notebooks aren't bad, they just serve a different purpose ⇒ **Experimentation** (which requires being messy)

# Scripting vs

- **Stateless**
  - Scripts don't have state which are shared between runs.
  - Scripts re-run consistently.
- **Linear**
  - Scripts run sequentially from top to bottom.
- **Reproducibility**
  - Easier to reproduce results.

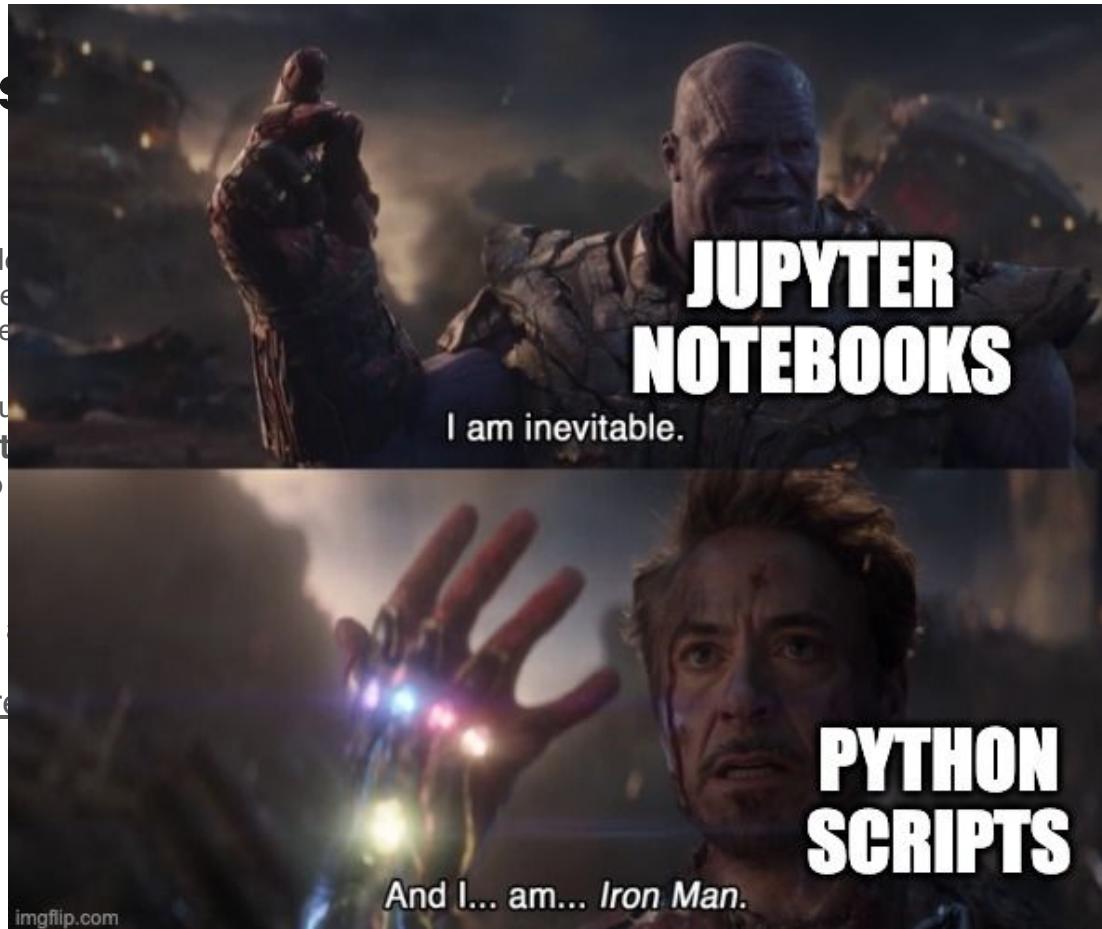
"Only a sith deals in

⇒ notebooks are

(unlike notebooks  
which require manual  
annotation.

in the wrong order

requires being messy)



# Code convention: Python PEP8

- Ensures a consistent code quality across a team.
- Set of rules on styling, such as
  - *Indentation (4 spaces)*
  - *Ordering imports*
  - *Max line length (79 characters)*
  - *Number of blank lines*
    - *Top-level function and class definitions: two blank lines*
    - *Method definitions inside a class: single blank line*
    - *Extra blank lines may be used (sparingly) to separate groups of related functions*
  - ...
- Integrate/automate in your code editor (IDE) to make it easy.  
Often enforced PEP8 during pull request submission.
- 🧑 Developers can be judgy... Make your life easy, adapt clean codes

*Tools to automatically check compliance with  
PEP8 will be covered in the CICD lecture*



# Code conventions: Python PEP8

- Ensures a consistent style
- Set of rules on structure
  - Indentation (4 spaces)
  - Max line length (79 characters)
  - Number of blank lines
    - Top-level
    - Method definitions
    - Extra blank lines between related sections
  - Ordering imports
  - ...
- Integrate/automate
- Often enforced by linters
- Developers can be...



# What is wrong with this code snippet?



Pathetic.  
Imgflip.com

```
import sys, os

def some_function(x,y):
    print('x+y=' ,x+y)
    if (x==0):
        print("x is 0!")
        x=x+1
    for i in range(0,10): print(i)
    try:os.listdir('.')
    except Exception as e:print(e);sys.exit(1)
```

# Cookie cutters: Initiating your repository

Many [cookie cutters](#) exist to easily set up specific kinds of code repositories.

The [Cookiecutter Data Science](#) project sets up a repository structure tailored to data science projects.

 **Disclaimer:** A bit old and subjective. E.g. saving data locally is not best practice.

```
LICENSE           <- The top-level README for developers using this project.  
README.md  
  
notebooks         <- Jupyter notebooks.  
  
references        <- Data schemas, manuals, and all other explanatory materials.  
  
reports           <- Generated analysis as HTML, PDF, LaTeX, etc.  
    └ figures      <- Generated graphics and figures to be used in reporting  
  
requirements.txt  <- The requirements file for reproducing the analysis environment, e.g.  
                    generated with `pip freeze > requirements.txt`  
  
setup.py          <- makes project pip installable (pip install -e .) so src can be imported  
src               <- Source code for use in this project.  
.gitignore        <- Specify rules for code versioning
```

---

# Wrap-up

# Lecture summary

| Topic                    | Concepts  | To know for... |      |
|--------------------------|---|----------------|------|
|                          |   | Project        | Exam |
| Agile / Scrum            | <ul style="list-style-type: none"><li>• Methodology</li><li>• Types of roles &amp; meetings</li><li>• Kanban board</li></ul>                    |                | Yes  |
| Data sources             | <ul style="list-style-type: none"><li>• Internal vs external</li><li>• Structured vs unstructured</li><li>• Data labelling techniques</li></ul> |                | Yes  |
| Code versioning          | <ul style="list-style-type: none"><li>• Git Flow</li><li>• Trunk based</li><li>• Pull requests, feedback giving</li><li>• PEP8</li></ul>        | Yes            |      |
| Lab: Git code versioning |   | Yes            |      |

# Project objective for sprint 1

## Define the use case you will be tackling

| Week | Work package  | Requirement |
|------|---|-------------|
| W01  | Pick a <b>team</b> (3-5 people) <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    |
| W01  | <b>Select</b> a use case<br>Source options <ul style="list-style-type: none"><li>Previous course</li><li><a href="https://www.kaggle.com/datasets">https://www.kaggle.com/datasets</a></li><li>...</li></ul><br>Make sure to pick a use case where <b>data is available</b> . | Required    |
| W01  | <b>Define</b> your use case with the <u>ML Canvas template</u> page   | Required    |
| W02  | Setup <b>communication channel</b> (Discord, Trello board (optional))   | Required    |
| W02  | Setup a <b>code versioning repository</b>   | Required    |
| W02  | Find a cool name for your team ✨  | Required    |



# Project deliverables

## Deliverables

- **3 Milestone (MS) meetings:**
  - MS 1: Present work from sprint 1 & 2
    - Present your general use case (BMC), the data preparation and the result from your model experimentation.
  - MS 2: Present work from sprint 3 & 4
    - Present your architecture for model deployment and automated training.
  - MS 3: Overall presentation
    - Whole project, inc. sprint 5
- Each MS presentation will be accompanied with a *code submission on Github*. Teaching staff will not go over all codes but key information can be provided in README + check that everything is well implemented.

**No handovers before MS meetings.** The work packages per sprints are there to guide you, but you're free to implement it at your pace 😊 Teaching staff is there to provide support every week.

# Project deliverables

 Submit your **team** and **project** by filling in this [project card template](#) (make a copy) and sending it to the teaching staff.

- Purpose is for the teaching staff to validate the project ideas and potentially provide feedback.
- Fill it in before next lesson (26/02/2024).

|  |   |
|--|---|
|  | <p>[Project name]</p> <p>INFO9023 - Project card</p> <p>Purpose of this document is to briefly explain what your project will be about so the teaching staff can validate it and provide feedback. It is not part of grading, no need to spend too much time editing it (just a couple of sentences per section).</p> <p>Make a copy and send it to <a href="mailto:tvranken@uliege.be">tvranken@uliege.be</a> and <a href="mailto:Matthias.Pirlet@uliege.be">Matthias.Pirlet@uliege.be</a> by the <u>26/02/2024</u>.</p> <p><b>Project description</b></p> <p><i>[Short description of your project]</i></p> <p><b>Project data</b></p> <p><i>[Short description of the data you'll use]</i></p> |
|--|---|

---

# Lab: Git code versioning

**MADE IT THROUGH THE LECTURE**



**SEE YOU NEXT WEEK !**

imgflip.com