02476 Machine Learning Operations
Nicki Skafte Detlefsen

# Reproducibility and software

# What is in this presentation?

♻ The reproducibility crisis

  What is it?

  How bad is it?

🖥 How can software help?

  What to use when?

```
 8    // Dear programmer:
 9    // When I wrote this code, only god and
10    // I knew how it worked.
11    // Now, only god knows it!
12    //
13    // Therefore, if you are trying to optimize
14    // this routine and it fails (most surely),
15    // please increase this counter as a
16    // warning for the next person:
17    //
18    // total_hours_wasted_here = 254
19    //
20
```
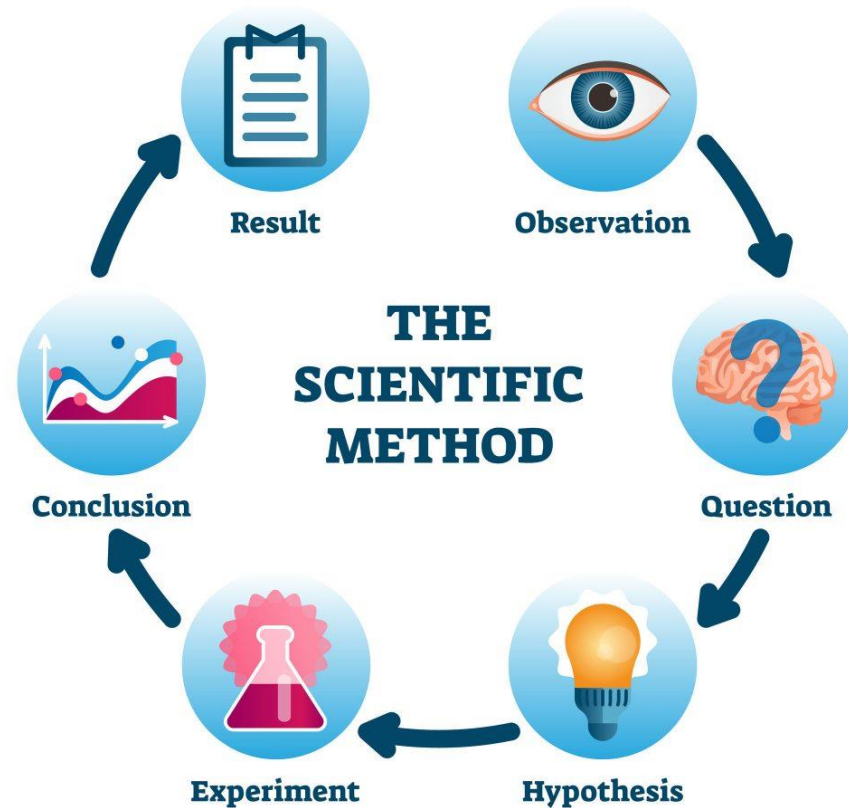
# What is reproducibility?

💡 Reproducibility is the ability of **an entire experiment** or study to be duplicated, either by the same researcher or **by someone else working independently**.

💡 Reproducible data - **repeatability** which is the degree of agreement of tests or measurements on replicate specimens by the same observer in the same laboratory.

💡 Computationally reproducible research - the idea that the ultimate product of **academic research** is the paper along with the **full computational environment** used to produce the results in the paper such as the code, data, etc. that can be used to reproduce the results and create new work based on the research.

# Why do we need it?

I would argue, because of this

# What does reproducibility have to do with MLOps?

⚑ Knowledge Preservation

      If other cannot reproduce your work, knowledge can be lost


⚑ Transparency and Accountability

      To make sure that others can verify your claims before going into production


⚑ Regulatory Compliance

      To secure the correct documentation to make sure everything is compliant
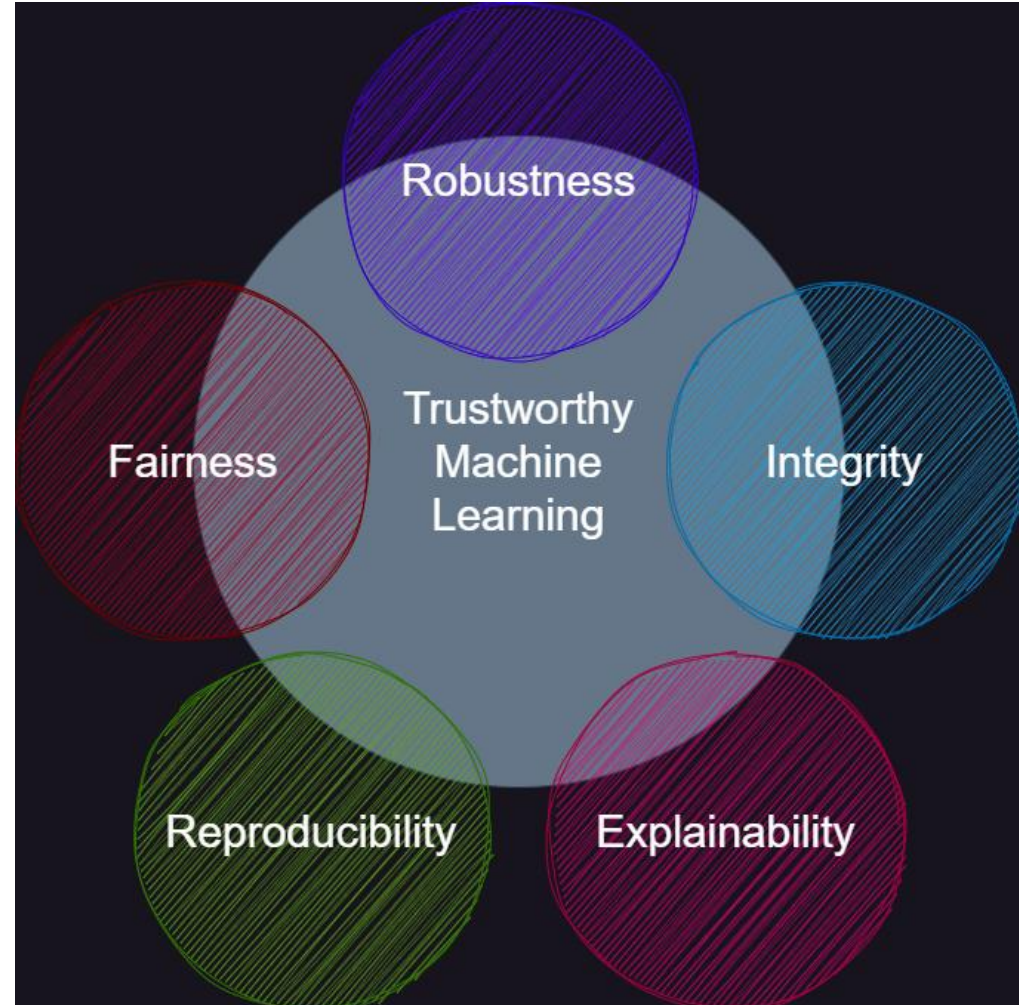

⚑ Continuous Improvement

      Making sure that improvements to the pipeline are real and not artifacts of random effects

# Trustworthy ML

Reproducibility is a key component in *Trustworthy ML*

⚠ Case:

Imaging an AI agent used for diagnostics. Without reproducibility two persons with the exact same symptoms could get different diagnosis

# We are in a crisis!

🏳 There is growing alarm about results that cannot be reproduced. Explanations include

- Increased levels of scrutiny
- Complexity of experiments and statistics
- Pressures on researchers



IS THERE A REPRODUCIBILITY CRISIS?

7% Don't know

52% Yes, a significant crisis

3% No, there is no crisis

1,576 researchers surveyed

38% Yes, a slight crisis

©nature

# An example: Trouble in the lab

The biotech company Amgen had a team of about 100 scientists trying to reproduce the findings of 53 "landmark" articles in cancer research published by reputable labs in top journals.
Only 6 of the 53 studies were reproduced (about 10%).

**REPRODUCIBILITY OF RESEARCH FINDINGS**
Preclinical research generates many secondary publications, even when results cannot be reproduced.

| Journal impact factor | Number of articles | Mean number of citations of non-reproduced articles* | Mean number of citations of reproduced articles |
|---|---|---|---|
| >20 | 21 | 248 (range 3–800) | 231 (range 82–519) |
| 5–19 | 32 | 169 (range 6–1,909) | 13 (range 3–24) |

Results from ten-year retrospective analysis of experiments performed prospectively. The term 'non-reproduced' was assigned on the basis of findings not being sufficiently robust to drive a drug-development programme.
*Source of citations: Google Scholar, May 2011.

# Another example: We are only humans

💭 The (subjective) choices we take during research may impact the conclusion



## ONE DATA SET, MANY ANALYSTS

Twenty-nine research teams reached a wide variety of conclusions using different methods on the same data set to answer the same question (about football players' skin colour and red cards).

Dark-skinned players four times more likely than light-skinned players to be given a red card.

- Statistically significant effect
- Non-significant effect

Twice as likely

Equally likely

Point estimates and 95% confidence intervals. *Truncated upper bounds.

# What are we trying to do within research?

https://paperswithcode.com/

Checklist for conferences

# A closer look at machine learning

⚡ Re-Implementation of 255 paper. Hypothesis testing on what "paper features" have an effect on reproducibility.

## A Step Toward Quantifying Independently Reproducible Machine Learning Research

**Edward Raff**
Booz Allen Hamilton
raff_edward@bah.com
University of Maryland, Baltimore County
raff.edward@umbc.edu

### Abstract

What makes a paper independently reproducible? Debates on reproducibility center around intuition or assumptions but lack empirical results. Our field focuses on releasing code, which is important, but is not sufficient for determining reproducibility. We take the first step toward a quantifiable answer by manually attempting to implement 255 papers published from 1984 until 2017, recording features of each paper, and performing statistical analysis of the results. For each paper, we did not look at the authors code, if released, in order to prevent bias toward discrepancies between code and paper.

Table 1: Significance test of which paper properties impact reproducibility. Results significant at $\alpha \leq 0.05$ marked with "*".

| Feature | p-value |
|---|---|
| Year Published | 0.964 |
| Year First Attempted | 0.674 |
| Venue Type | 0.631 |
| Rigor vs Empirical* | $1.55 \times 10^{-9}$ |
| Has Appendix | 0.330 |
| Looks Intimidating | 0.829 |
| Readability* | $9.68 \times 10^{-25}$ |
| Algorithm Difficulty* | $2.94 \times 10^{-5}$ |
| Pseudo Code* | $2.31 \times 10^{-4}$ |
| Primary Topic* | $7.039 \times 10^{-4}$ |
| Exemplar Problem | 0.720 |
| Compute Specified | 0.257 |
| Hyperparameters Specified* | $8.45 \times 10^{-6}$ |
| Compute Needed* | $8.75 \times 10^{-5}$ |
| Authors Reply* | $6.01 \times 10^{-8}$ |
| Code Available | 0.213 |
| Pages | 0.364 |
| Publication Venue | 0.342 |
| Number of References | 0.740 |
| Number Equations* | 0.004 |
| Number Proofs | 0.130 |
| Number Tables* | 0.010 |
| Number Graphs/Plots | 0.139 |
| Number Other Figures | 0.217 |
| Conceptualization Figures | 0.365 |
| Number of Authors | 0.497 |

# Levels of reproducibility
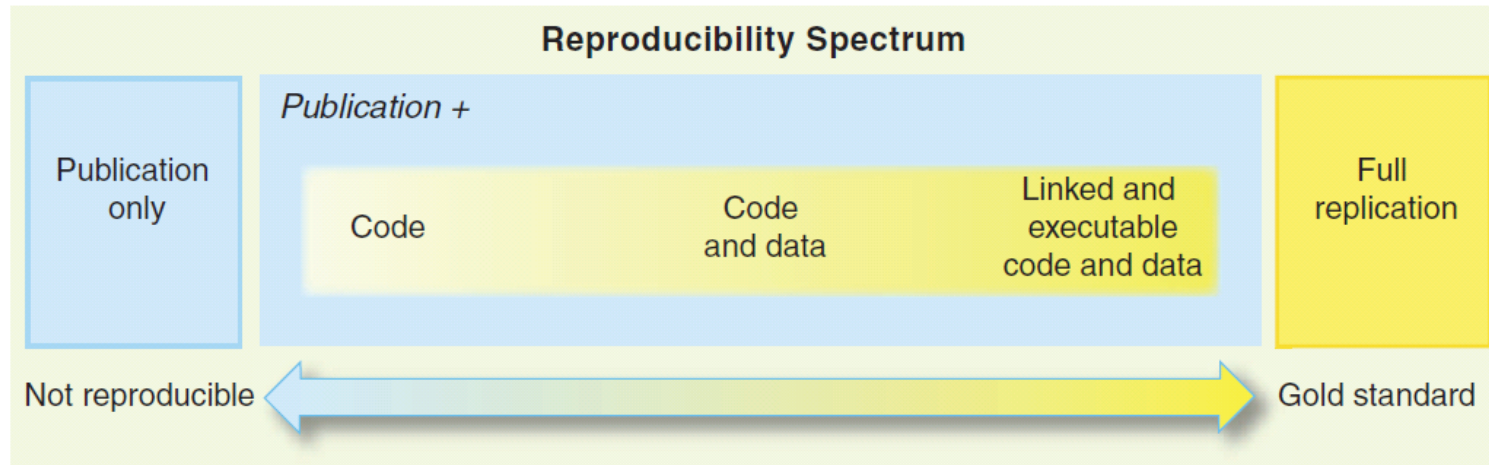
Reproducibility is not binary, it's a spectrum



**Fig. 1.** The spectrum of reproducibility.

Example from ML

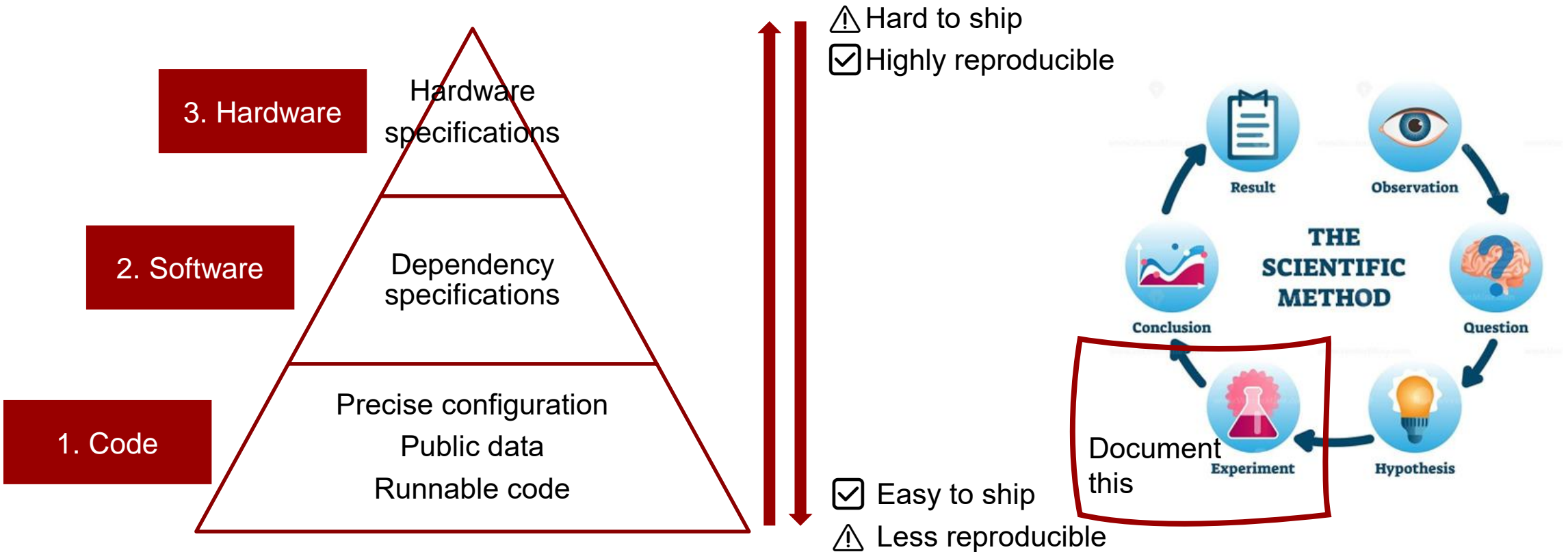$$\begin{Bmatrix} w_1 \\ \vdots \\ w_n \end{Bmatrix} == \begin{Bmatrix} w'_1 \\ \vdots \\ w'_n \end{Bmatrix}$$

m1.ckpt     m2.ckpt

$\boxed{\text{vs}}$

| Dataset | Model Architecture | Random Init | Transfer | Parameters | IMAGENET Top5 |
|---------|-------------------|-------------|----------|------------|---------------|
| RETINA | Resnet-50 | 96.4% ± 0.05 | 96.7% ± 0.04 | 23570408 | 92.% ± 0.06 |
| RETINA | Inception-v3 | 96.6% ± 0.13 | 96.7% ± 0.05 | 22881424 | 93.9% |
| RETINA | CBR-LargeT | 96.2% ± 0.04 | 96.2% ± 0.04 | 8532480 | 77.5% ± 0.03 |
| RETINA | CBR-LargeW | 95.8% ± 0.04 | 95.8% ± 0.05 | 8432128 | 75.1% ± 0.3 |
| RETINA | CBR-Small | 95.7% ± 0.04 | 95.8% ± 0.01 | 2108672 | 67.6% ± 0.3 |
| RETINA | CBR-Tiny | 95.8% ± 0.03 | 95.8% ± 0.01 | 1076480 | 73.5% ± 0.05 |

# What can we do about it ?

Make sure that you document everything about your experiments
🩸 Nicki's hierarchy of reproducibility of ML 🩸



**3. Hardware** — Hardware specifications

**2. Software** — Dependency specifications

**1. Code** — Precise configuration / Public data / Runnable code

⚠ Hard to ship
☑ Highly reproducible

☑ Easy to ship
⚠ Less reproducible

THE SCIENTIFIC METHOD

Result — Observation — Question — Hypothesis — Experiment — Conclusion

Document this

# Reproducibility level 1

⚡ There is a lot of subjective choices that we do when running experiments in machine learning, most notable the hyperparameters.

| Parameters in scripts | Argument parser | Config files |
|---|---|---|
| ```python<br>class hparams:<br>    lr = 0.1<br>    batch_size = 16<br>    num_layers = 5``` | ```python<br>python my_script.py \<br>    --lr 0.1 \<br>    --batch_size 16 \<br>    --num_layers 5``` | ```yaml<br>experiment1.yaml<br><br>lr: 0.001<br>batch_size: 16<br>num_layers: 5<br><br>python my_script.py \<br>    config=experiment1.yaml``` |
| ☑Easy to code<br>⚠Not easy to configure on the run<br>⚠Experimental info may be lost if not careful | ☑Easy to configure<br>⚠Falls on user to save the config | ☑Highly configurable<br>☑Configuration is systematically saved (and version controlled) |

# Reproducibility level 1

Hydra is a framework for elegantly configuring complex (ML) applications

https://github.com/facebookresearch/hydra

Example:

```
├── conf
│   ├── config.yaml
│   └── dataset
│       ├── cifar10.yaml
│       └── imagenet.yaml
└── my_app.py
```

```python
import hydra
from omegaconf import DictConfig

@hydra.main(config_path="config.yaml")
def my_app(cfg: DictConfig) -> None:
    print(cfg.pretty())

if __name__ == "__main__":
    my_app()
```

Other options
- 💡 https://github.com/IDSIA/sacred
- 💡 https://mlflow.org/

# Reproducibility level 2

For python: Just use a package management system

Examples:

💡 [Conda](#) (what I like)

💡 [Pipenv](#)

💡 [venv](#)

💡 [pyenv](#)

# Reproducibility level 3

💡 The easiest way for someone to reproduce your work, would be to just hand over your computer.

💡 Instead of doing this, lets hand over a virtual copy of our machine

Virtual machine works by taking hardware from the host and creates virtual CPU, RAM, storage for each virtual machine. The virtual machines are completely independent from the host

# Reproducibility level 3

💡 The core advantage of a VM is that it in principal can run on any host without changes, because it is independent.

💡 Docker can be seen as an lightweight version of full VMs.

💡 *VM is isolation of machines, while Containers is isolation of processes*



| App 1 | App 2 | App 3 |
|-------|-------|-------|
| Bins/Lib | Bins/Lib | Bins/Lib |
| Guest OS | Guest OS | Guest OS |

Hypervisor

Infrastructure

**Virtual Machines**

| App 1 | App 2 | App 3 |
|-------|-------|-------|
| Bins/Lib | Bins/Lib | Bins/Lib |

Container Engine

Operating System

Infrastructure

Containers

# Reproducibility level 3

A way to create containerized applications = low overhead virtual machines (VMs)



| Dockerfile | Docker Image | Docker Container |
|---|---|---|
| Instructions for creating the container | File containing the everything | A running instance of our application |

# A 6-step process for reproducible software



Use version control

# A 6-step process for reproducible software

Use version control

Use templates



```
├── LICENSE
├── Makefile           <- Makefile with commands like `make data` or `make train`
├── README.md          <- The top-level README for developers using this project.
├── data
│   ├── external       <- Data from third party sources.
│   ├── interim        <- Intermediate data that has been transformed.
│   ├── processed      <- The final, canonical data sets for modeling.
│   └── raw            <- The original, immutable data dump.
│
├── docs               <- A default Sphinx project; see sphinx-doc.org for details
│
├── models             <- Trained and serialized models, model predictions, or model summaries
│
├── notebooks          <- Jupyter notebooks. Naming convention is a number (for ordering),
│                         the creator's initials, and a short `-` delimited description, e.g.
│                         `1.0-jqp-initial-data-exploration`.
│
├── references         <- Data dictionaries, 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.
│   ├── __init__.py    <- Makes src a Python module
│   │
│   ├── data           <- Scripts to download or generate data
│   │   └── make_dataset.py
│   │
│   ├── features       <- Scripts to turn raw data into features for modeling
│   │   └── build_features.py
│   │
│   ├── models         <- Scripts to train models and then use trained models to make
│   │   │                 predictions
│   │   ├── predict_model.py
│   │   └── train_model.py
│   │
│   └── visualization  <- Scripts to create exploratory and results oriented visualizations
│       └── visualize.py
│
└── tox.ini            <- tox file with settings for running tox; see tox.readthedocs.io
```
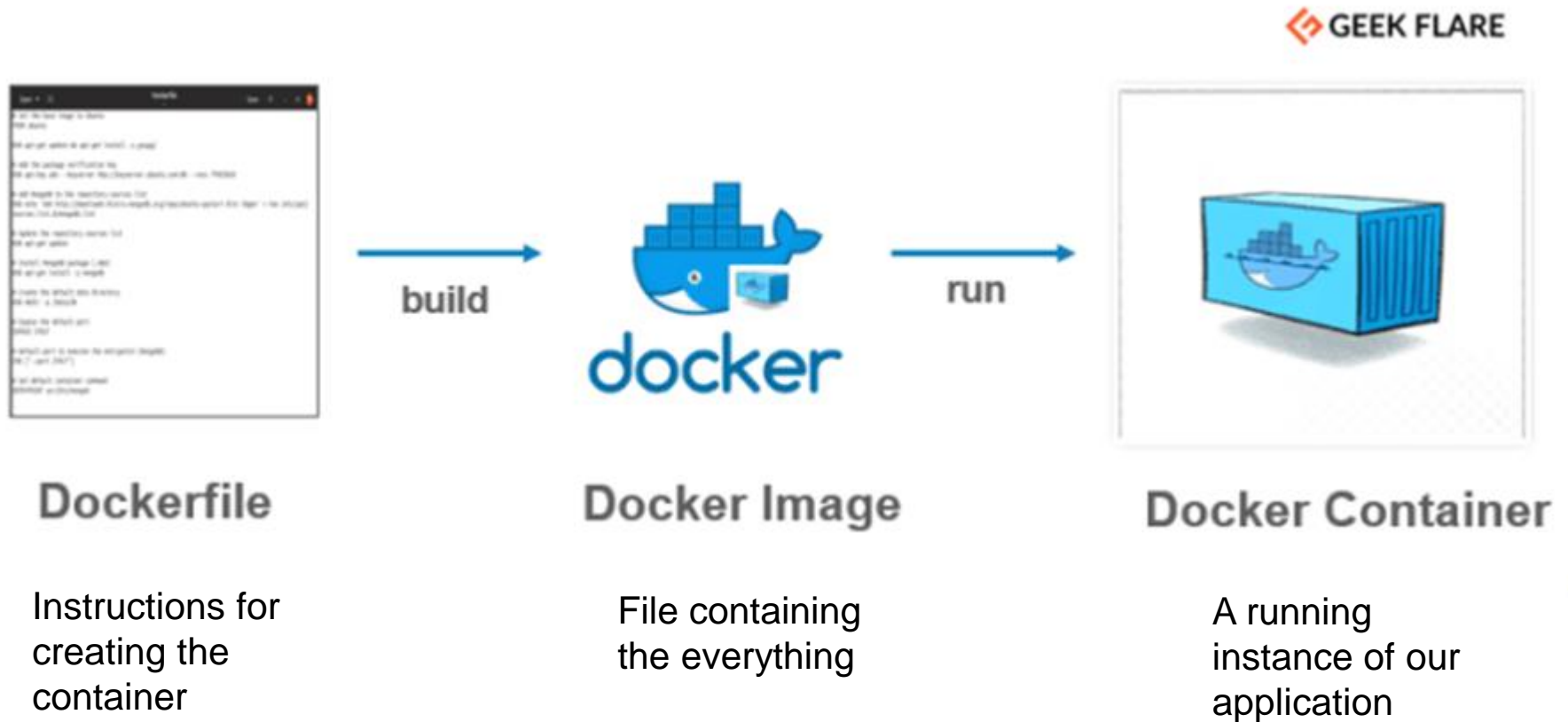
# A 6-step process for reproducible software

Use version control

Use templates

Write down your dependencies
(and use virtual environments)

```
≡ requirements.txt ✕

≡ requirements.txt
    1    Click==7.0
    2    Flask==1.1.1
    3    gunicorn==19.9.0
    4    itsdangerous==1.1.0
    5    Jinja2==2.10.1
    6    MarkupSafe==1.1.1
    7    Werkzeug==0.15.6
    8
    9
```

# A 6-step process for reproducible software

Use version control

Use templates

Write down your dependencies
(and use virtual environments)

Document your code!



**Train & Test**

To train and test, you can call:

```
cd src;
CUDA_VISIBLE_DEVICES=0 python trainer_[INSERT MODEL].py --config PATH_TO_CONFIG
```

For example to train online LAE on mnist and evaluate on kmnist

```
cd src;
CUDA_VISIBLE_DEVICES=0 python trainer_lae_elbo.py --config ../configs/ood_experiments/mnist/linear/lae_elbo.y
```

and try train a VAE

```
cd src;
CUDA_VISIBLE_DEVICES=0 python trainer_vae.py --config ../configs/ood_experiments/mnist/linear/vae.yaml
```

You can monitor training on tensorboard

```
tensorboard --logdir lightning_log --port 6006
```

To test on missing data imputation experiments, you can call. This require that you have a trained model.

```
cd src/data_imputation;
CUDA_VISIBLE_DEVICES=0 python lae.py
```

or

```
cd src/data_imputation;
CUDA_VISIBLE_DEVICES=0 python vae.py
```

# A 6-step process for reproducible software

Use version control

Use templates

Write down your dependencies
(and use virtual environments)

Document your code!

Test your code

# A 6-step process for reproducible software

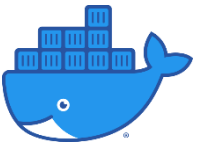Use version control

Use templates

Write down your dependencies
(and use virtual environments)

Document your code!

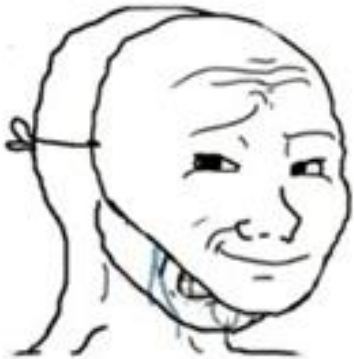Test your code

Containerize your code

```
1     # Import a base image so we don't have to start from scratch
2     #FROM python:3.10-slim
3     FROM huggingface/transformers-pytorch-cpu
4
5     # Use EXPOSE so we can give docker run the appropriate commandline argument (PORT) as:
6     # docker run predict:latest -e PORT=8000
7     EXPOSE $PORT
8     ENV LC_ALL=C.UTF-8
9     ENV LANG=C.UTF-8
10
11    # Run a bunch of linux commands
12    RUN apt update && \
13        apt install --no-install-recommends -y build essential gcc & \
14        apt clean & rm -rf /var/lib/apt/lists/*
15
16    # Copy the essential files from our folder to docker container.
17    COPY src/ src/
18    COPY requirements_predict.txt requirements_predict.txt
19    COPY setup.py setup.py
20
21    RUN pip install -r requirements_predict.txt --no-cache-dir
22
23    RUN dvc init --no-scm
24    RUN dvc remote add -d gcloud_storage gs://mlops-dataset-small
25    RUN dvc pull
26
27    # Set working directory as / and install dependencies
28    WORKDIR /
29
30    RUN mkdir app
31
32    # Set entry point, i.e. which file we run with which argument when running the docker container.
33    # The -u flag makes it print to console rather than the docker log file.
34    #ENTRYPOINT ["python", "-u", "src/models/predict_model.py"]
35    CMD exec uvicorn src.models.predict_model:app --host 0.0.0.0 --workers 1 --port $PORT
36    #ENTRYPOINT ["uvicorn", "src.models.predict_model:app", "--host", "0.0.0.0", "--workers", "1", "--port", $PORT]
```
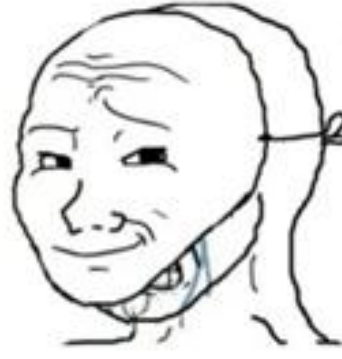
# Meme of the day