

Lecture 10

Accessible Multimodal Learning and Transfer Learning with PyKale



[Haiping Lu](#)

YouTube Playlist: <https://www.youtube.com/c/HaipingLu/playlists>

[COM4059/6059: MLAI21@The University of Sheffield](#)

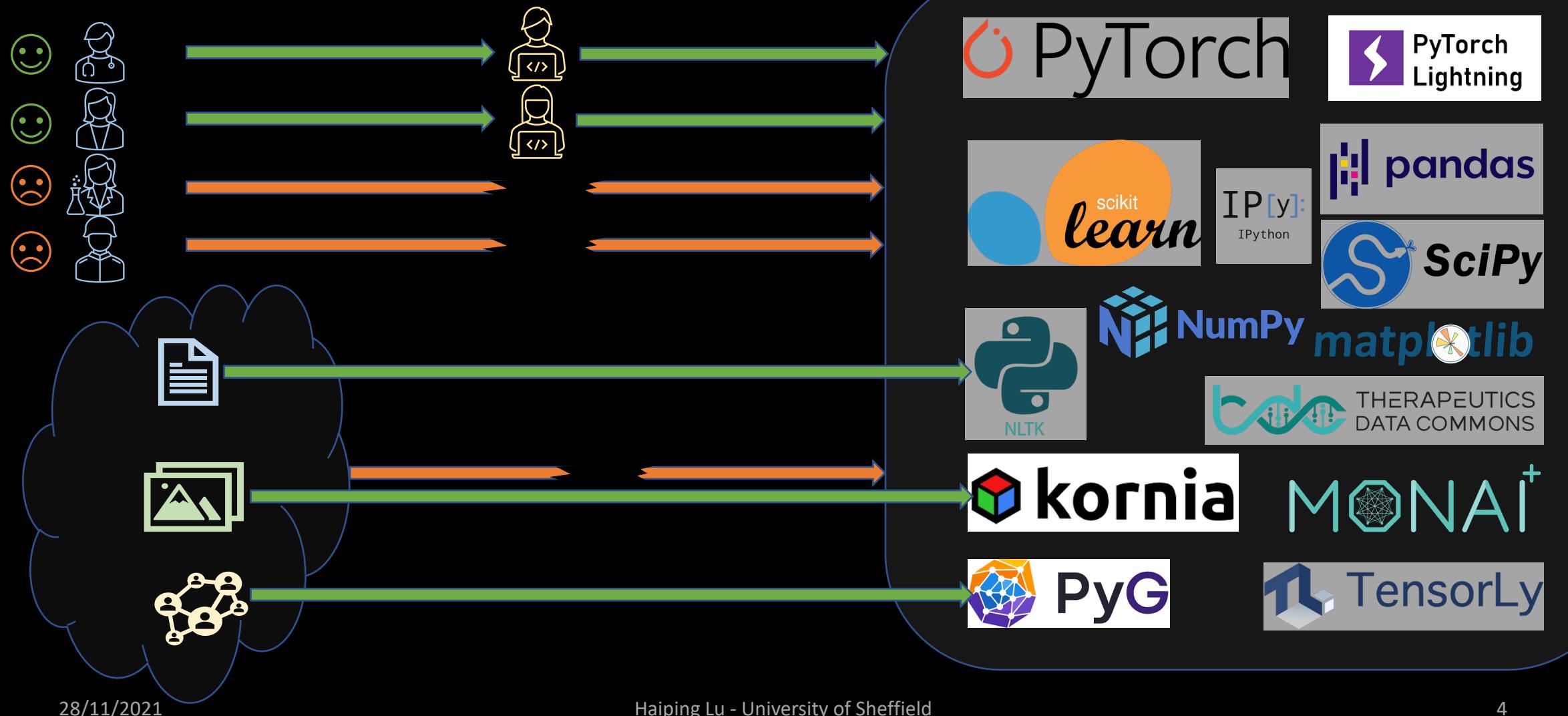
Week 10 Contents / Objectives

- Why PyKale
- How PyKale was Built
- What PyKale Looks Like
- Summary and Review

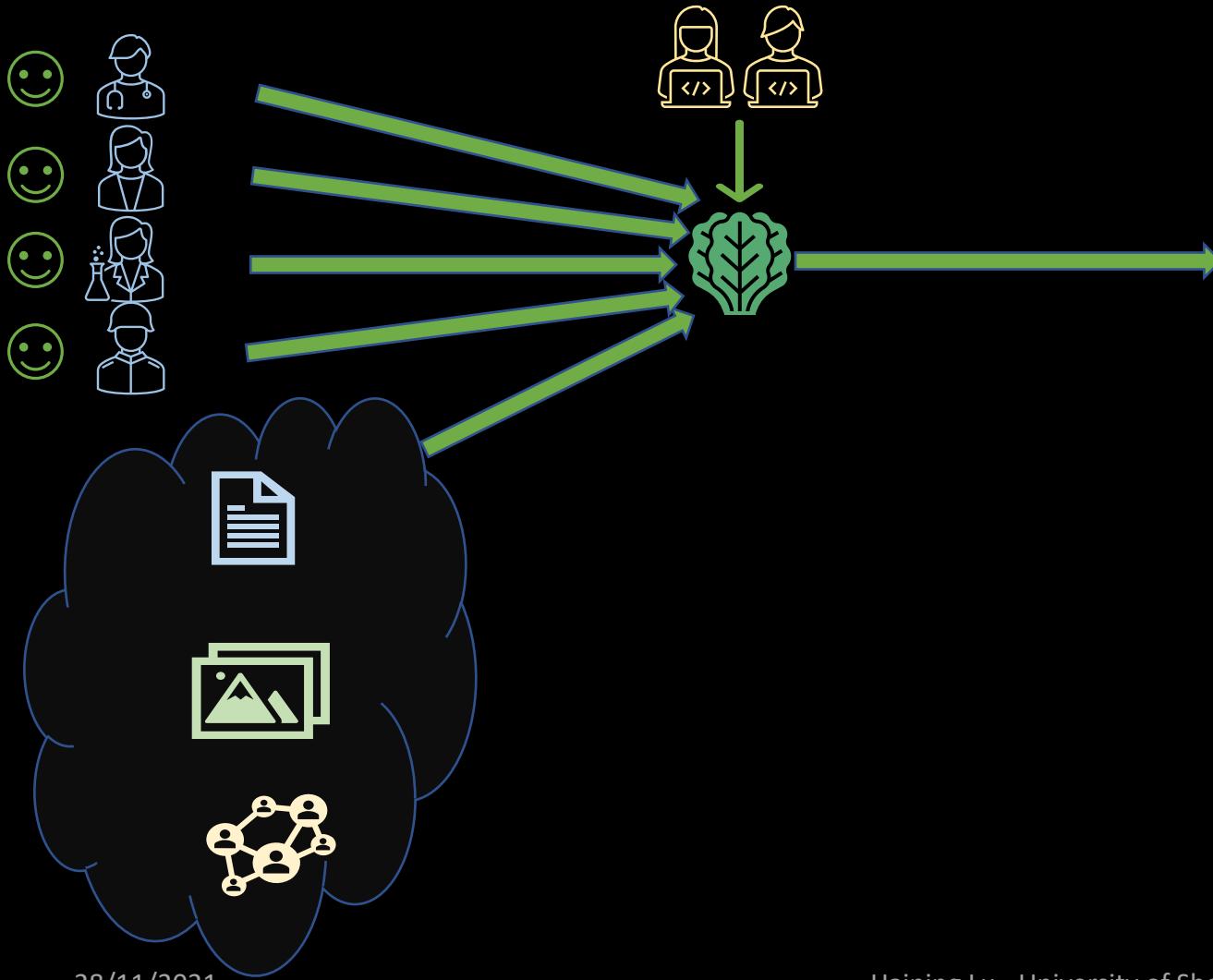
Week 10 Contents / Objectives

- **Why PyKale**
- How PyKale was Built
- What PyKale Looks Like
- Summary and Review

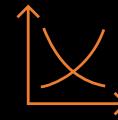
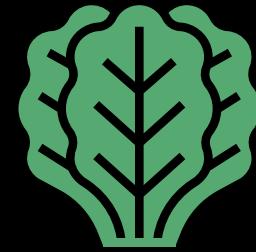
Accessibility issues in interdisciplinary research



Accessibility solution



Why build PyKale?



- Accessibility challenge in interdisciplinary research
- Abundant ML software vs unmet user demands



- Low-level ML software → focus on basics & fundamentals
- Single-domain ML software → tailor for single domains, varying APIs



- Our software → high-level, multi-source, unified API, domain-traversing
- Bridge gaps between data, software and end users for accessible, scalable, and sustainable research in machine learning

Week 10 Contents / Objectives

- Why PyKale
- **How PyKale was Built**
- What PyKale Looks Like
- Summary and Review

Is it FAIR?



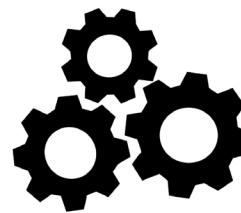
F
indable



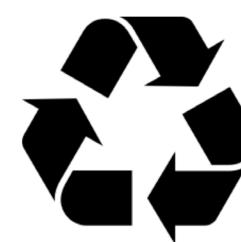
A
ccessible



I
nteroperable



R
Reusable

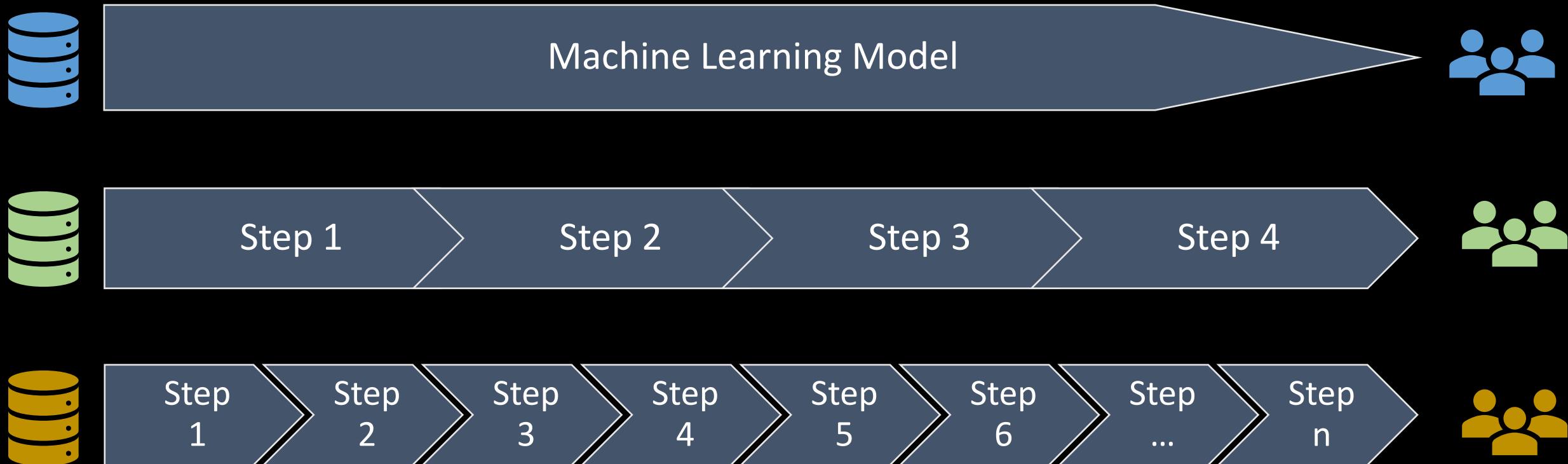


New principles: green machine learning

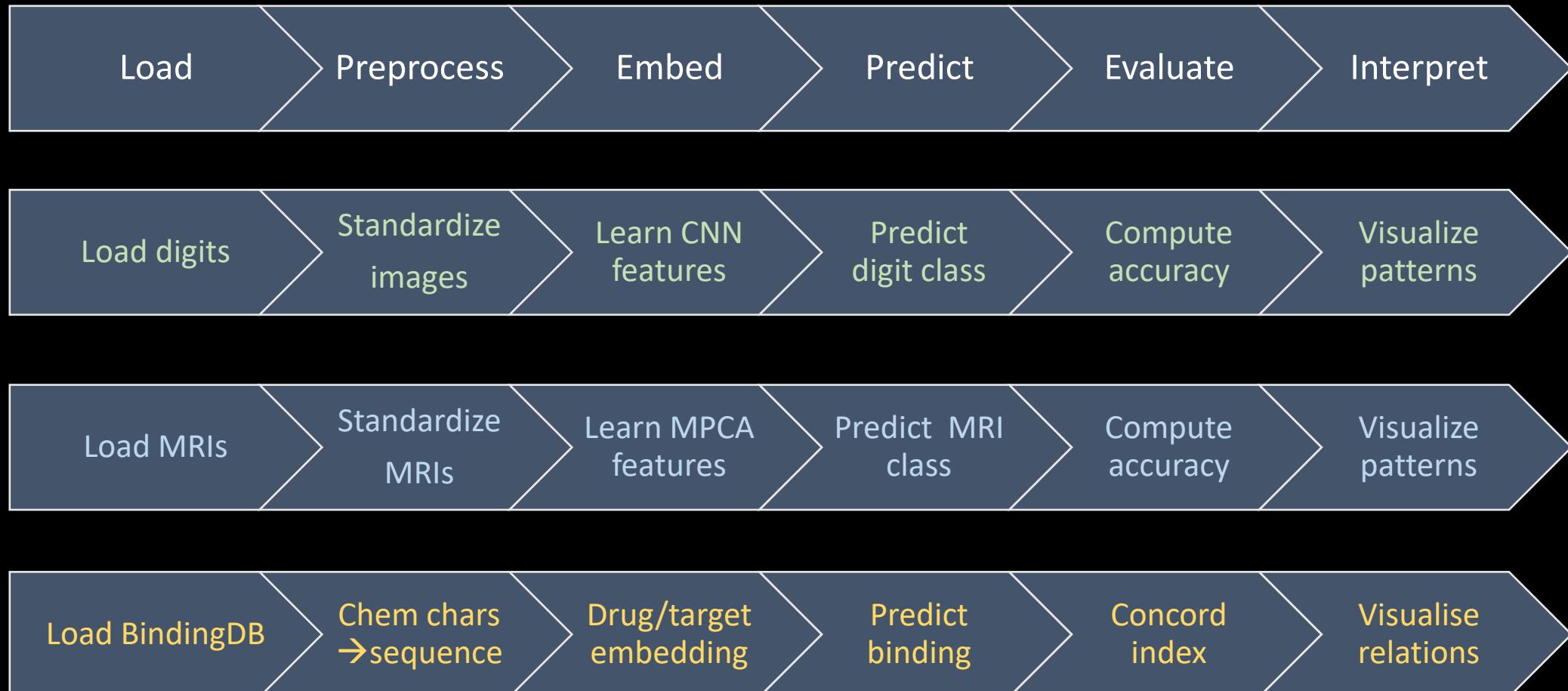
- **Reduce** repetition and redundancy
 - Refactor to standardize workflow and enforce styles
 - Identify and remove duplicated functionalities
- **Reuse** existing resources
 - Reuse the same machine learning pipeline for different data
 - Reuse existing software for available functionalities
- **Recycle** learning models across areas
 - Identify commonalities between applications
 - Recycle models for App A to App B



API without standardization 😞



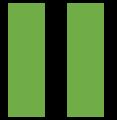
Pipeline-based API



Accessible machine learning



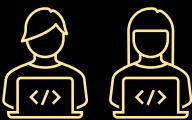
- Bridge the gap: abundant ML software unmet user demands



- Separate code and configuration: pykale/examples

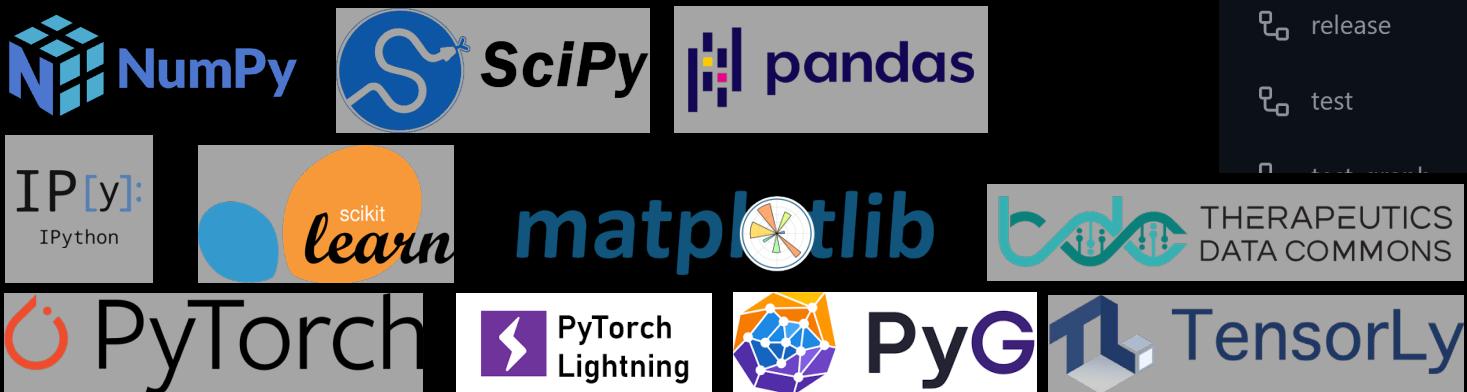


- End users (non-experts in ML)
 - Configure via YAML, no need to write Python/ML → Lower entry barriers
- ML experts
 - Manage config with min coding → Improve reproducibility/efficiency
 - Deliver accessible ML software → **Scalable** collaboration
→ Broad applications



Sustainable software development

- GitHub ecosystem → industrial standard
 - Automation & continuous integration
 - Read the Docs, PyPI release
 - Linting, Pre-commit, PyTest,Codecov
- Python/PyTorch ecosystem



A screenshot of a GitHub repository's Actions page. The top navigation bar shows Code, Issues (5), Pull requests (3), Discussions, and Actions (selected). Below the navigation is a search bar and a button for "New workflow". A sidebar on the left lists workflow names: assign project, changelog, lint, release, and test. The main area displays "All workflows" with 2,617 workflow runs. One run is marked as "test #840: Scheduled" with a green checkmark, and another is marked as "Multi source example" with a red X.

Week 10 Contents / Objectives

- Why PyKale
- How PyKale was Built
- **What PyKale Looks Like**
- Summary and Review

Kale API snapshot: pipeline-based

The screenshot displays a terminal window with four tabs, each showing a list of files and a context menu. The tabs are:

- `pykale / kale /` (highlighted)
- `pykale / kale / loaddata /`
- `pykale / kale / utils /`
- `pykale / kale / pipeline /`

The context menu (visible for the first tab) contains the following items:

- Load Data
- Preprocess Data
- Embed
- Predict
- Evaluate
- Interpret
- Pipeline
- Utilities

The lists of files for each tab are:

- `pykale / kale /`:
 - `..`
 - `README.md`
 - `__init__.py`
 - `cifar_access.py`
 - `dataset_access.py`
 - `digits_access.py`
 - `get_dicom.py`
 - `mnistm.py`
 - `multi_domain.py`
 - `office_access.py`
- `pykale / kale / loaddata /`:
 - `..`
 - `sz144 fix digits_access`
 - `..`
 - `README.md`
 - `__init__.py`
 - `cifar_access.py`
 - `dataset_access.py`
 - `digits_access.py`
 - `get_dicom.py`
 - `mnistm.py`
 - `multi_domain.py`
 - `office_access.py`
- `pykale / kale / utils /`:
 - `..`
 - `XianyuanLiu Delete csv_logger.py`
 - `..`
 - `README.md`
 - `__init__.py`
 - `download.py`
 - `logger.py`
 - `print.py`
 - `seed.py`
- `pykale / kale / pipeline /`:
 - `..`
 - `haipinglu define acronyms`
 - `..`
 - `README.md`
 - `__init__.py`
 - `deep_dti.py`
 - `domain_adapter.py`
 - `mpca_trainer.py`
 - `multi_domain_adapter.py`
 - `video_domain_adapter.py`

At the bottom center of the terminal window, it says "u - University of Sheffield".

Standardized examples & code config separation

The image displays three GitHub commit screenshots illustrating the standardization of examples and code configuration separation.

Commit 1: [haipinglu add help on how to run](#)

- ..
- configs
- README.md
- __init__.py
- config.py
- main.py
- tutorial.ipynb

Commit 2: [haipinglu add help on how to run](#)

- ..
- configs
- figures
- README.md
- __init__.py
- config.py
- main.py
- model.py
- tutorial.ipynb

Commit 3: [haipinglu improve yaml doc](#)

- ..
- configs
- README.md
- __init__.py
- config.py
- main.py
- model.py
- tutorial.ipynb

YAML configuration

```
5 lines (4 sloc) | 49 Bytes

1 DAN:
2   METHOD: "DANN"
3
4 DATASET:
5   SOURCE: "svhn"
```

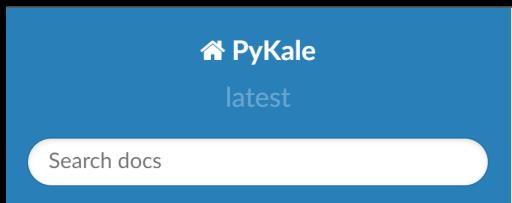
```
14 lines (11 sloc) | 155 Bytes

1 DAN:
2   METHOD: "CDAN"
3
4 DATASET:
5   NUM_REPEAT: 1
6   SOURCE: "svhn"
7   VAL_SPLIT_RATIO: 0.5
8
9 SOLVER:
10  MIN_EPOCHS: 0
11  MAX_EPOCHS: 3
12
13 OUTPUT:
14  PB_FRESH: None
```

```
# -----
# Dataset
# -----
_C.DATASET = CN()
_C.DATASET.ROOT = "../data" # Path to store data and results
_C.DATASET.NAME = "digits" # Name of datasets
_C.DATASET.SOURCE = "mnist" # The source dataset name
_C.DATASET.TARGET = "usps" # The target dataset name
_C.DATASET.NUM_CLASSES = 10
_C.DATASET.NUM_REPEAT = 10
_C.DATASET.DIMENSION = 784
_C.DATASET.WEIGHT_TYPE = "natural"
_C.DATASET.SIZE_TYPE = "source"
_C.DATASET.VAL_SPLIT_RATIO = 0.1
# -----
# Solver
# -----
_C.SOLVER = CN()
_C.SOLVER.SEED = 2020
_C.SOLVER.BASE_LR = 0.001 # Initial learning rate
_C.SOLVER.MOMENTUM = 0.9
_C.SOLVER.WEIGHT_DECAY = 0.0005 # 1e-4
_C.SOLVER.NESTEROV = True

_C.SOLVER.TYPE = "SGD"
_C.SOLVER.MAX_EPOCHS = 120 # "nb_adapt_epochs": 100,
# _C.SOLVER.WARMUP = True
_C.SOLVER.MIN_EPOCHS = 20 # "nb_init_epochs": 20,
_C.SOLVER.TRAIN_BATCH_SIZE = 150 # 150
_C.SOLVER.TEST_BATCH_SIZE = 200 # No difference in ADA
```

Community engagement



GETTING STARTED

Introduction

Installation

Tutorial



Introduction to PyKale, a library for knowledge-aware machine learning from multi

643 views · Apr 17, 2021

1 20 0 SHARE



Haiping Lu

74 subscribers

» PyKale Documentation

PyKale Documentation

Getting Started

- [Introduction](#)
- [Installation](#)

book Tutorials
n using YAML

arXiv.org > cs > arXiv:2106.09756

Computer Science > Machine Learning

[Submitted on 17 Jun 2021]

PyKale: Knowledge-Aware Machine Learning from Multiple Sources in Python

Haiping Lu, Xianyuan Liu, Robert Turner, Peizhen Bai, Raivo E Koot,
Shuo Zhou, Mustafa Chasmai, Lawrence Schobs

→ C https://github.com/pykale/pykale/blob/main/.github/CONTRIBUTING.md + ⚙ | ☆ | ⌂

189 lines (126 sloc) | 18.5 KB

Contributing to PyKale

Light involvements (viewers/users) | Medium involvements (contributors) | Heavy involvements (maintainers)

Ask questions | Report bugs | Suggest improvements | Branch, fork & pull | Coding style | Test | Review & merge | Release & management

A GitHub project board for the PyKale repository. It shows three columns: "Overview & backlog", "To do", and "In progress". The "Overview & backlog" column has 2 items: "Added by haipinglu" and "v0.2.0". The "To do" column has 8 items: "Remove dependency on csv logger and maybe remove csv logger.", "Make pykale.github.io following https://github.com/Project-MONAI/PyKale", and "DA for action recognition TA3N". The "In progress" column has 2 items: "1 of 5 tasks" and "Review required".

Search...

Help | Advanced

New Top: All ▾ Answered Unanswered Label ▾ New discussion

Unique Selling Points
bobturneruk asked 11 days ago in General · Unanswered

2

Should we have tests for tutorial notebooks?
bobturneruk asked on 12 Aug in Q&A · Answered

2

Where should we publish on PyKale?
bobturneruk asked on 12 Aug in Q&A · Unanswered

3

GETTING STARTED

Introduction

Installation

Tutorial

Jupyter Notebook Tutorials

» Jupyter Notebook Tutorials

Jupyter Notebook Tutorials

- Image classification: Digits domain adaptation
- Drug discovery: BindingDB drug-target interaction prediction
- Medical image analysis: Cardiac MRI for MPCA-based diagnosis

← → ⌂ https://colab.research.google.com/github/pykale/pykale/blob/main/examples/digits_dann_lightn/tutorial.ipynb



tutorial.ipynb

File Edit View Insert Runtime Tools Help

+ Code

+ Text

Copy to Drive

PyKale Tutorial: Domain Adaptation on Digits with Lightning

| [Open in Colab](#) (click Runtime → Run all (Ctrl+F9) | [Launch Binder](#) (click Run → Run All Cells) |

If using [Google Colab](#), a free GPU can be enabled to save time via setting Runtime → Change runtime type

← → ⌂ https://hub.gke2.mybinder.org/user/pykale-pykale-77qjut67/doc/tree/examples/digits_dann_lightn/tutorial.ipynb

tutorial.ipynb

File Edit View Run Kernel Tabs Settings Help



Python 3 (ipykernel)

PyKale Tutorial: Domain Adaptation on Digits with Lightning

| [Open in Colab](#) (click Runtime → Run all (Ctrl+F9) | [Launch Binder](#) (click Run → Run All Cells) |

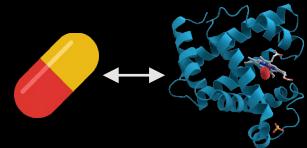
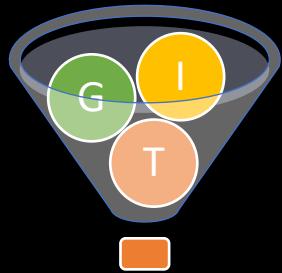
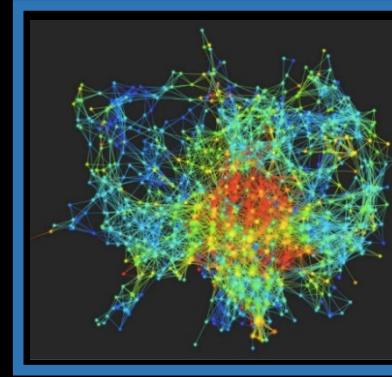
If using [Google Colab](#), a free GPU can be enabled to save time via setting Runtime → Change runtime type → Hardware accelerator: GPU



[livedemo-compressor.jpg \(1450×960\) \(posterno.com\)](#)

Data, model, & app

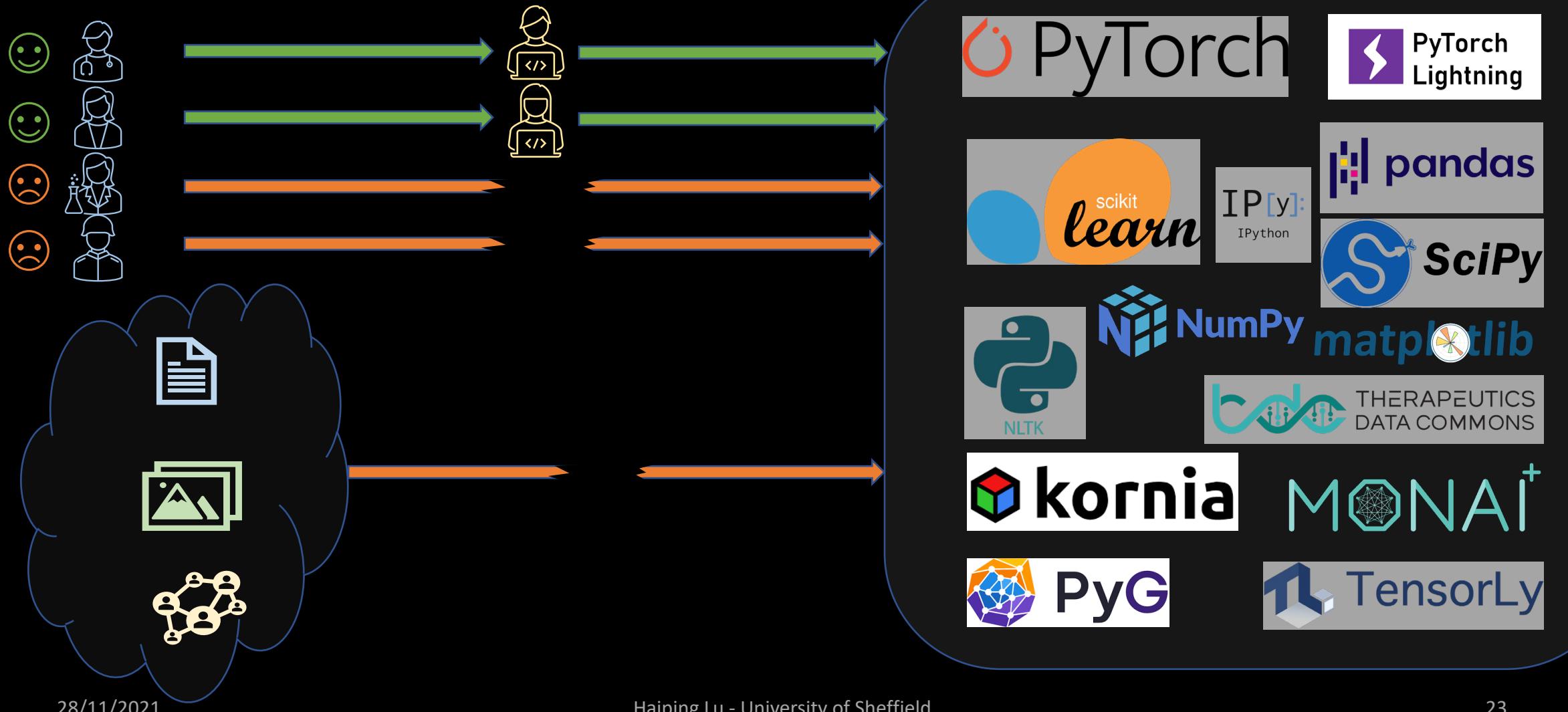
- Data: graph, image, text, video
- Machine learning models
 - Transfer learning: domain adaptation
 - Multimodal learning: integration of heterogeneous data
 - Deep learning: CNN, GNN/GCN, transformers
 - Dimensionality reduction: tensor-based
- Example applications
 - Image/video recognition: CIFAR, digits, office, action videos
 - Bioinformatics/graph analysis: BindingDB, knowledge graphs
 - Medical imaging: cardiac MRI, brain fMRI



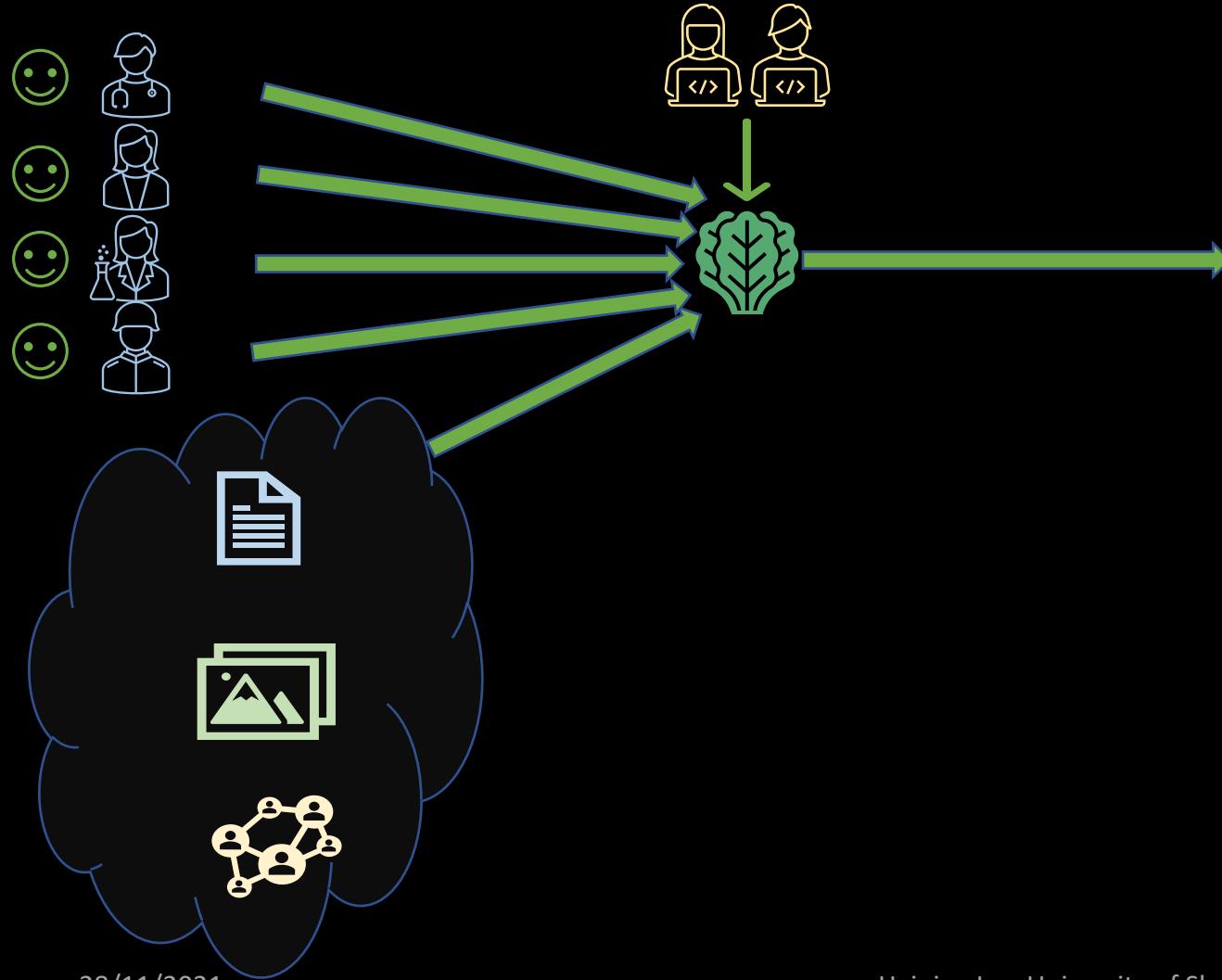
Week 10 Contents / Objectives

- Why PyKale
- How PyKale was Built
- What PyKale Looks Like
- **Summary and Review**

Accessibility issues in interdisciplinary research



PyKale fills the gaps



PyKale: accessible, scalable, sustainable → →

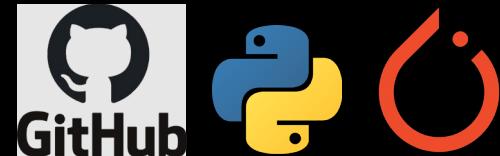
- **Green ML principles:** deliver FAIR ML software



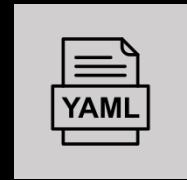
- Pipeline-based API to unify workflow



- Industrial software engineering practices for **sustainability**

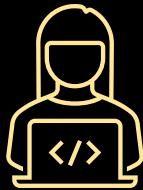


- **Accessible** and **scalable** with code-configuration separation



- Multimodal learning & transfer learning for graphs, images, videos, ...

Targeted audience



ML experts



End users

- Reproducible research
- Scalable collaboration
- Cross-boundary models

- Off-the-shelf ML pipelines
- Friendly configuration editing
- Multi modalities under one roof

Open with strong community engagement

Open & welcome

- Open source on GitHub (MIT License)
 - <https://github.com/pykale/pykale>
- Welcome feedback/contribution!

A screenshot of a web browser window. The address bar shows the URL <https://pytorch.org/ecosystem/>. The page content includes the PyTorch logo and a section titled "PyKale". The text in this section reads: "PyKale is a PyTorch library for multimodal learning and transfer learning with deep learning dimensionality reduction on graphs, images, texts, and videos."



Acknowledgement



- Shef groups + community



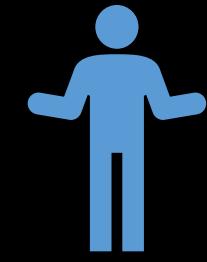
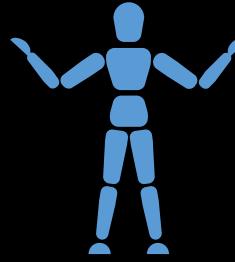
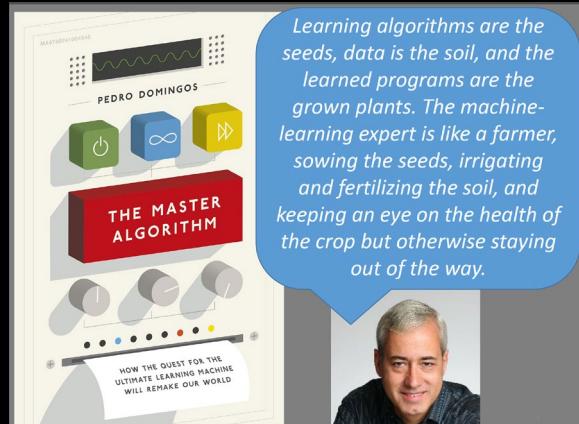
The
University
Of
Sheffield.



- Partially funded by Wellcome Trust via an Innovator Award

Review and Reflect on Machine Learning

- Training: model as well as people
 - You are a **model**
- Reflection and analogy
 - What you want to learn (objective function / goal)
 - How we teach you to learn (training / optimisation)
 - What you have learned (decision model)
- Key challenge: overfitting vs **generalisation**
- Post your questions
- Have fun



[kisspng-fashion-show-computer-icons-man-icon-5acd9cfaa3d9f5.3374529915234245066712.jpg \(900x900\) \(cleanpng.com\)](https://kisspng.com/fashion-show-computer-icons-man-icon-5acd9cfaa3d9f5.3374529915234245066712.jpg)

Data, Model, Metric, Optimisation

- | | |
|--|--|
| 1. Given training data:
$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ | 3. Define goal:
– Objective function
$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$ |
| 2. Choose each of these:
– Decision function/model
$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$ | 4. Train/optimize with SGD: (take small steps opposite the gradient)
$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$ |