# HEP ML Lab

## An end-to-end framework for machine learning application in high energy physics

Jing Li, Hao Sun
Dalian University of Technology, Liaoning, China
ML4Jets (2023)

Preparing for ArXiv
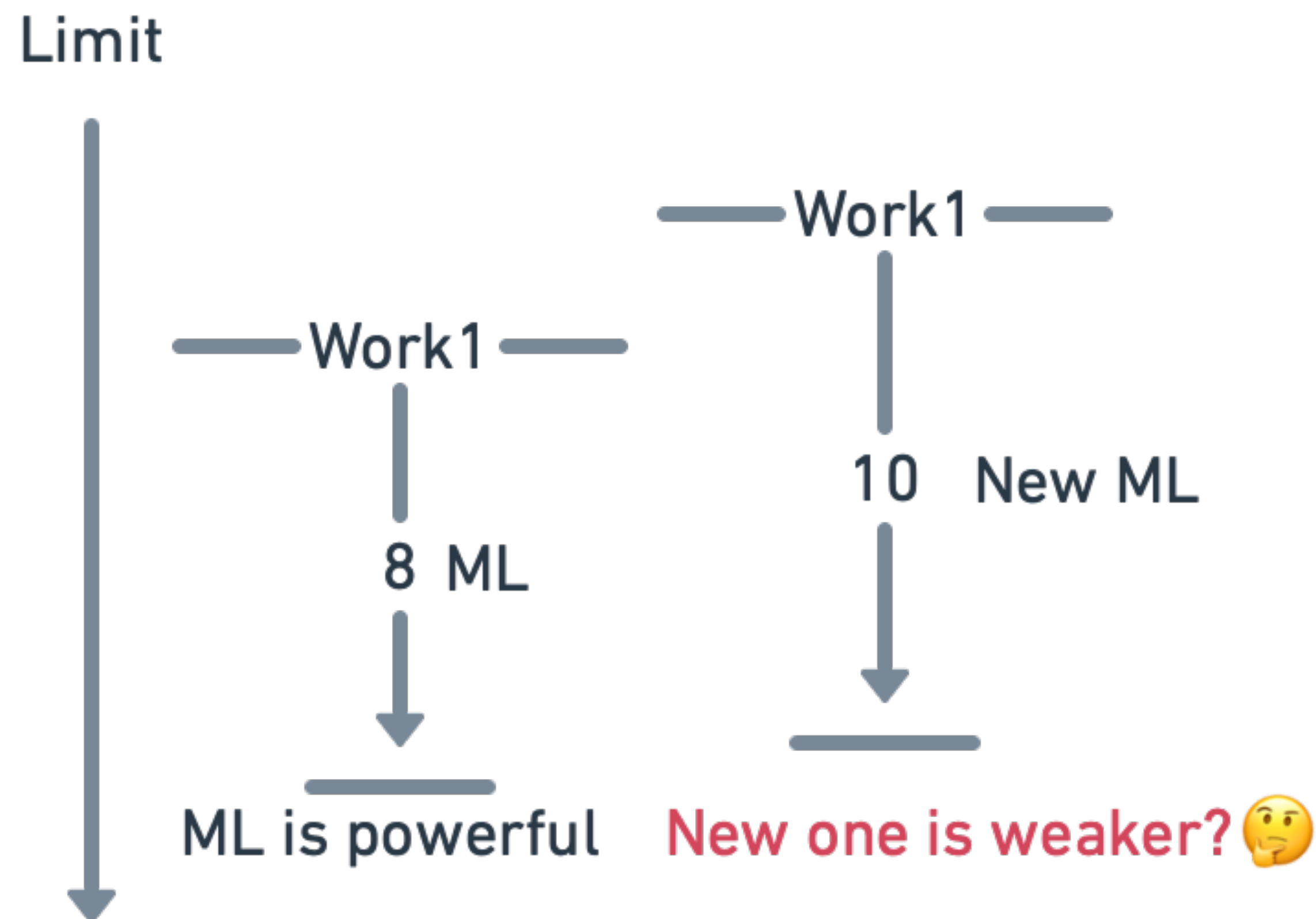https://github.com/Star9daisy/hep-ml-lab

# Table of contents

- Introduction: why we need an end-to-end framework?

- Quick start: generate events, create datasets, apply approaches
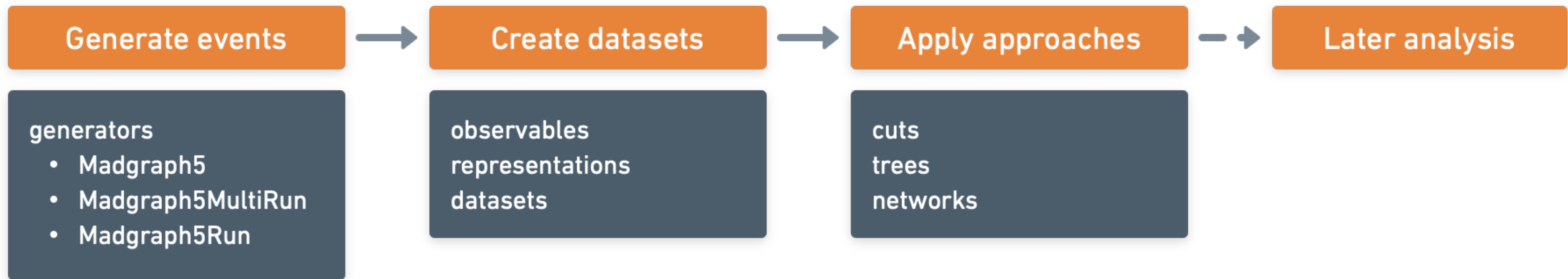
- Future: roadmap

# Introduction
## Reproduction issues

- A lot of works have explore the potential of machine learning approaches.

- The lack of source codes leads to reproduction problems.

- It's crucial to compare different algorithms on the same baseline.
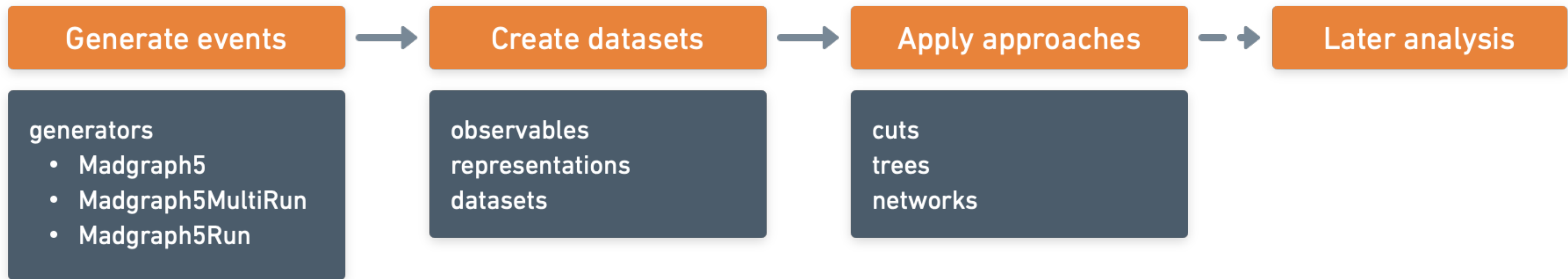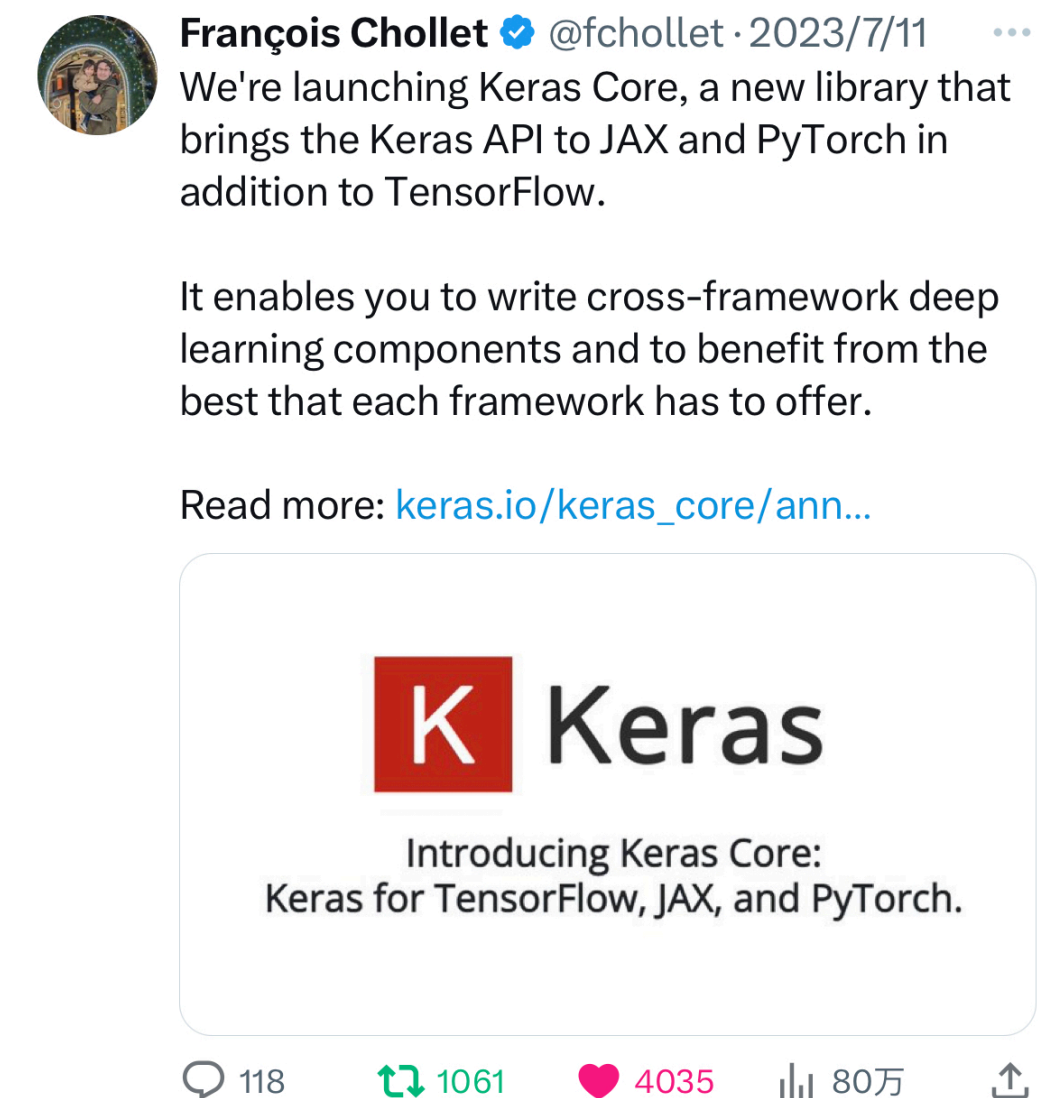
# Introduction
## Control from end to end



- HEP ML Lab (HML) helps apply different approaches into studies.

- An end-to-end framework ensure reproducibility.

# Introduction

## Highlights



- Data control from the very beginning.

- Approaches in **Keras** style.

  - Cross-framework support 👉

# Quick start

## Generate events

- Classical three commands

  - **generate, output, launch**

- Optimized parameters

  - **definitions, processes**

  - **settings**

- **summary** to show all runs' info

```python
1   zjj = Madgraph5(
2       executable="mg5_aMC",
3       model="sm",
4       definitions={},
5       processes=["p p > z z, z > j j, z > vl vl~"],
6       output="data/pp2zz_z2jj_z2vlvl",
7   )
8
9   zjj.launch(
10      shower="pythia8",
11      detector="delphes",
12      settings={
13          "iseed": 42,
14          "nevents": 1000,
15          "htjmin": 400,
16      },
17  )
18
19  zjj.summary()
20
```

# Quick start
## Generate events

```
1    Running Survey
2    Running Pythia8
3    Running Delphes
4    Storing files
5
6    Done
```
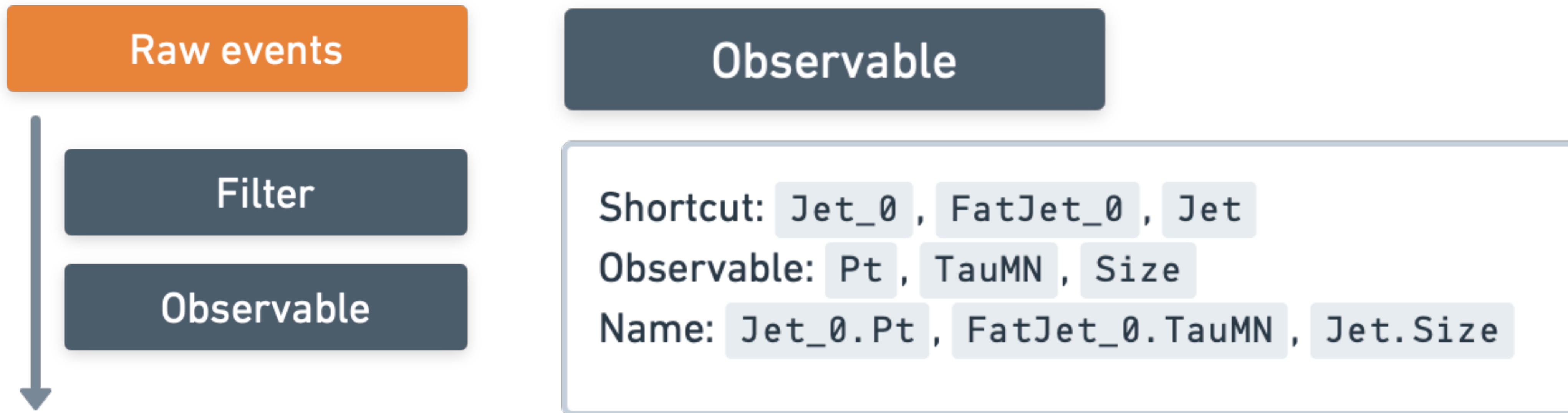
*p p > z z, z > j j, z > vl vl~*

| # | Name | Tag | Cross section (pb) | N events | Seed |
|---|------|-----|--------------------|----------|------|
| 0 | run_01[1] | tag_1 | 2.273e-03 | 1,000 | 42 |

*Output:*
*/root/workspace_ssd/projects/hep-ml-lab/examples/data/pp2zz_z2jj_z2vlvl*

- Simple status checker

- Summary table

7

# Quick start
## Create datasets



Shortcut: `Jet_0` , `FatJet_0` , `Jet`
Observable: `Pt` , `TauMN` , `Size`
Name: `Jet_0.Pt` , `FatJet_0.TauMN` , `Jet.Size`

- Observable parsing system: name = shortcut + observable.
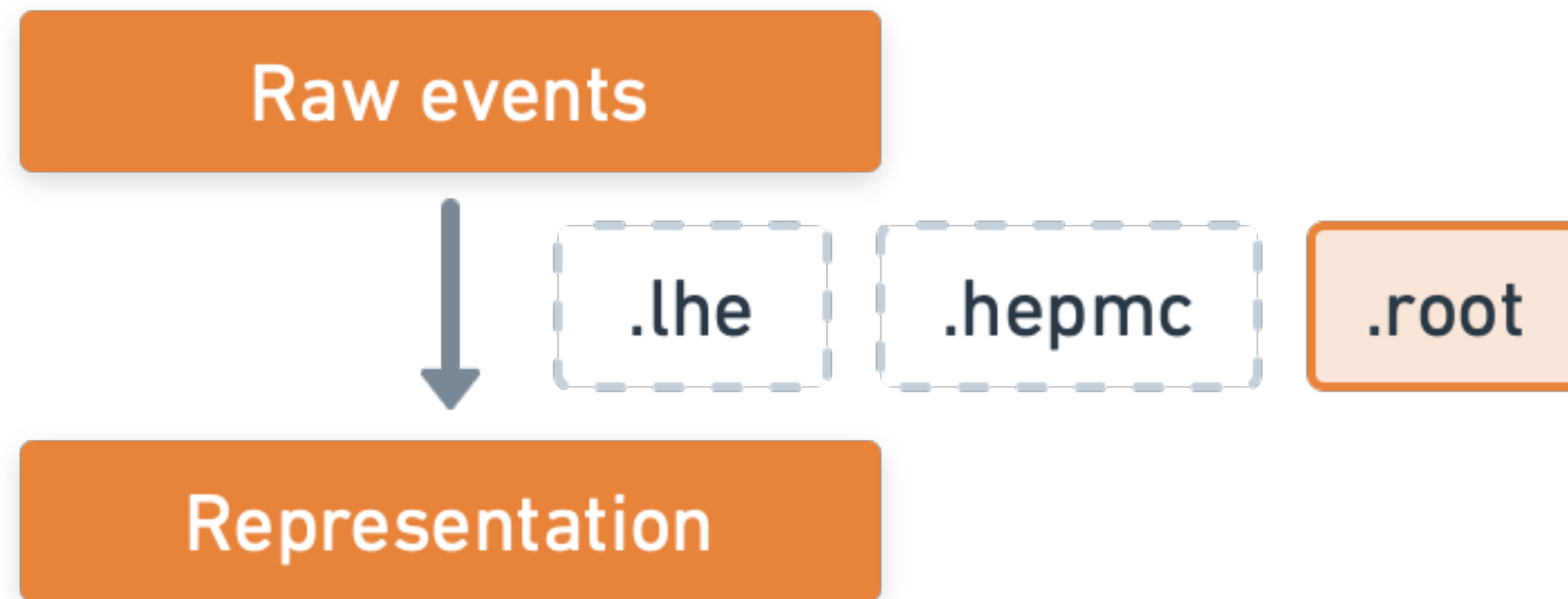
# Quick start
## Create datasets



- **Observable** parsing system: name = shortcut + observable.

- **Filter** accepts a list of logical "and" conditions.
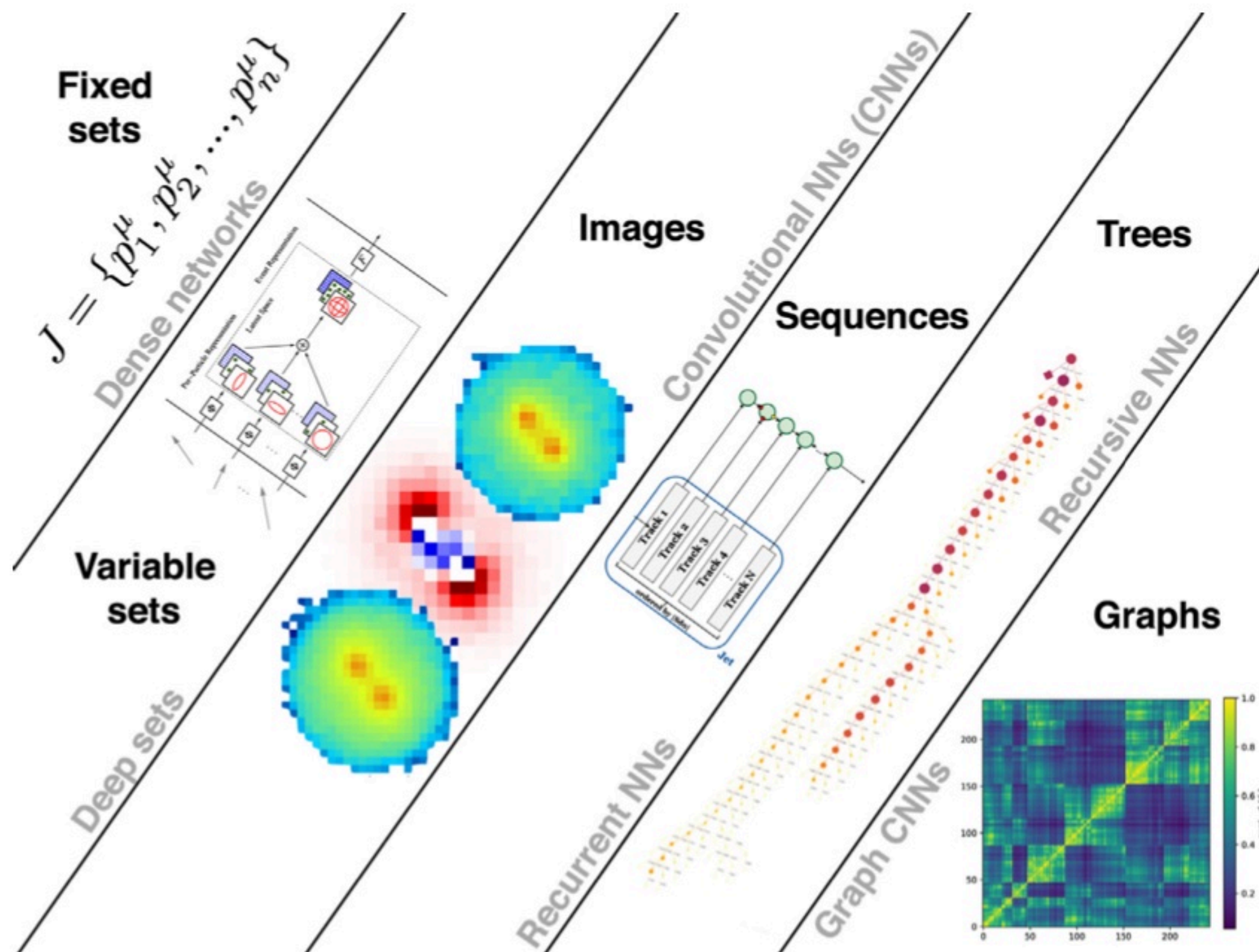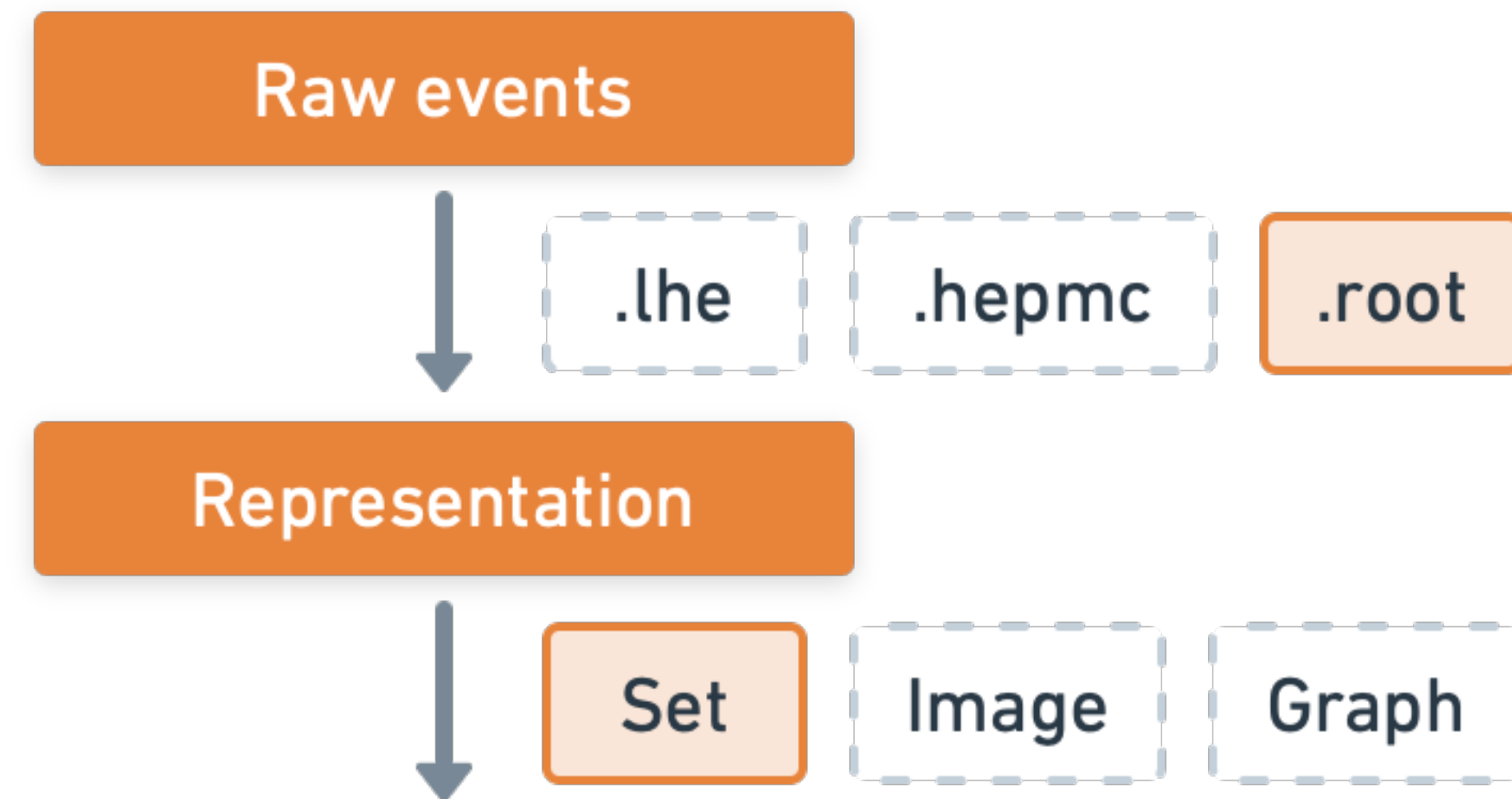
# Quick start
## Create datasets



- **Observable** parsing system: name = shortcut + observable.

- **Filter** accepts a list of logical "and" conditions.

- Currently only supports root format.

# Quick start
## Create datasets



[1709.04464] Jet Substructure at the Large Hadron Collider:
A Review of Recent Advances in Theory and Machine Learning



- Currently only supports **Set**.

- Expanding to **Image** and **Graph**.

# Quick start
## Create datasets

- **from_output** to load the existing runs.

- **get_observable** to parse the name.

  - Could pass strings directly to **Set** here.

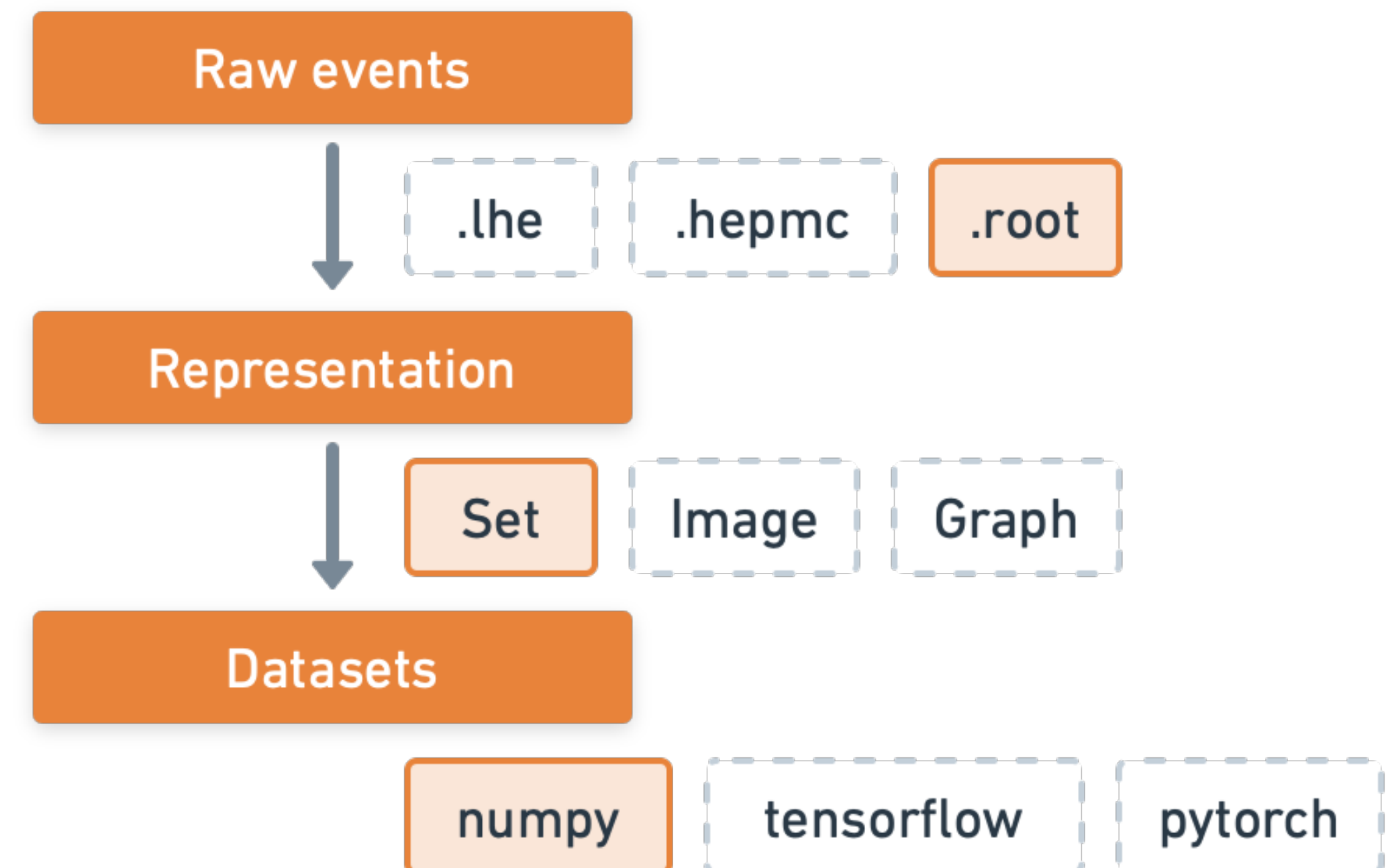- **passed** to check if the event is valid to preselection conditions

```
1   zjj = Madgraph5.from_output("./data/pp2zz_z2jj_z2vlvl")
2
3   preselections = Filter(["FatJet.Size > 0", "Jet.Size > 1"])
4
5   zjj_set = Set(
6       [
7           get_observable("FatJet_0.Mass"),
8           get_observable("FatJet_0.TauMN", m=2, n=1),
9           get_observable("Jet_0-Jet_1.DeltaR"),
10      ]
11  )
12
13  zjj_bar = Progbar(zjj.runs[0].n_events)
14  for i, event in enumerate(zjj.runs[0].events):
15      if preselections.read_event(event).passed():
16          zjj_set.read_event(event)
17
18      zjj_bar.update(i + 1)
19
```

# Quick start
## Create datasets

```python
1   samples = np.array(zjj_set.values + qcd_set.values, "float32")
2   targets = np.array([1] * len(zjj_set.values) + [0] * len(qcd_set.values), "int32")
3   dataset = TabularDataset(
4       samples=samples,
5       targets=targets,
6       feature_names=zjj_set.names,
7       target_names=["Z -> jj", "QCD dijets"],
8       description="Z -> jj vs QCD dijets",
9   )
10
11  dataset.save("./data/zjj_vs_qcd")
12
```
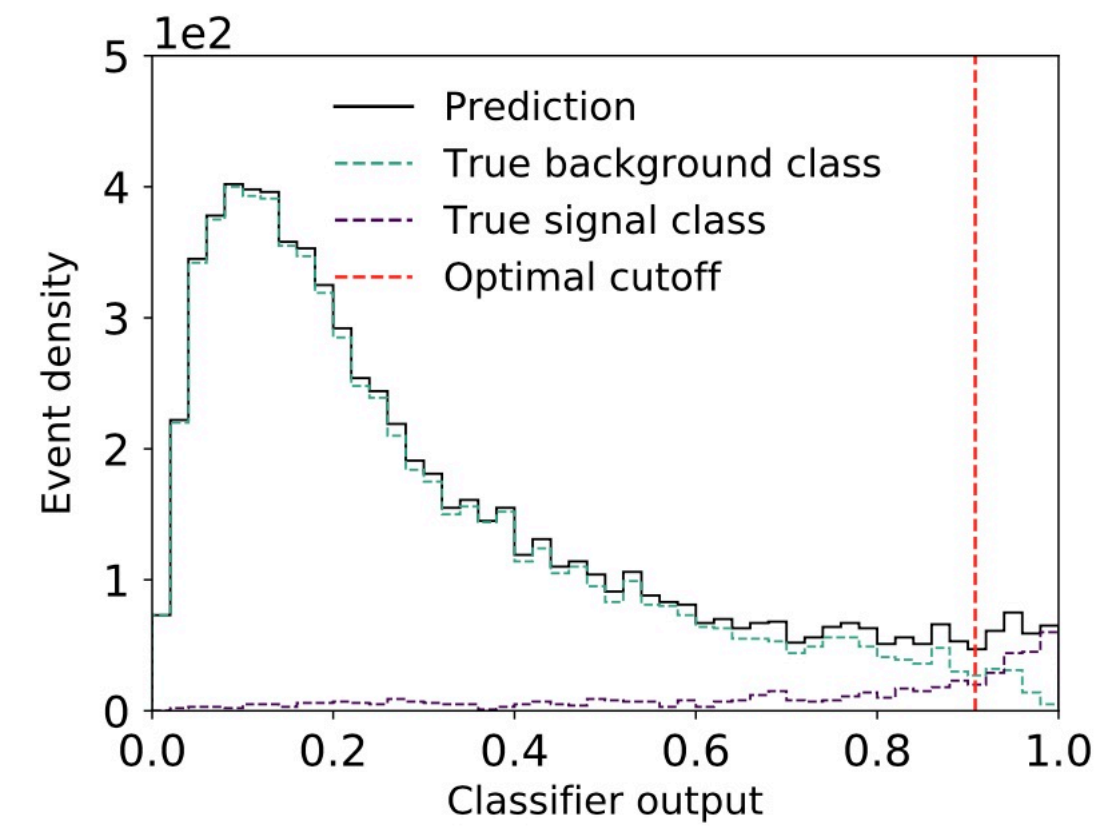


- **TabularDataset** for **Set**.

- Currently supports saving to **.npz** files
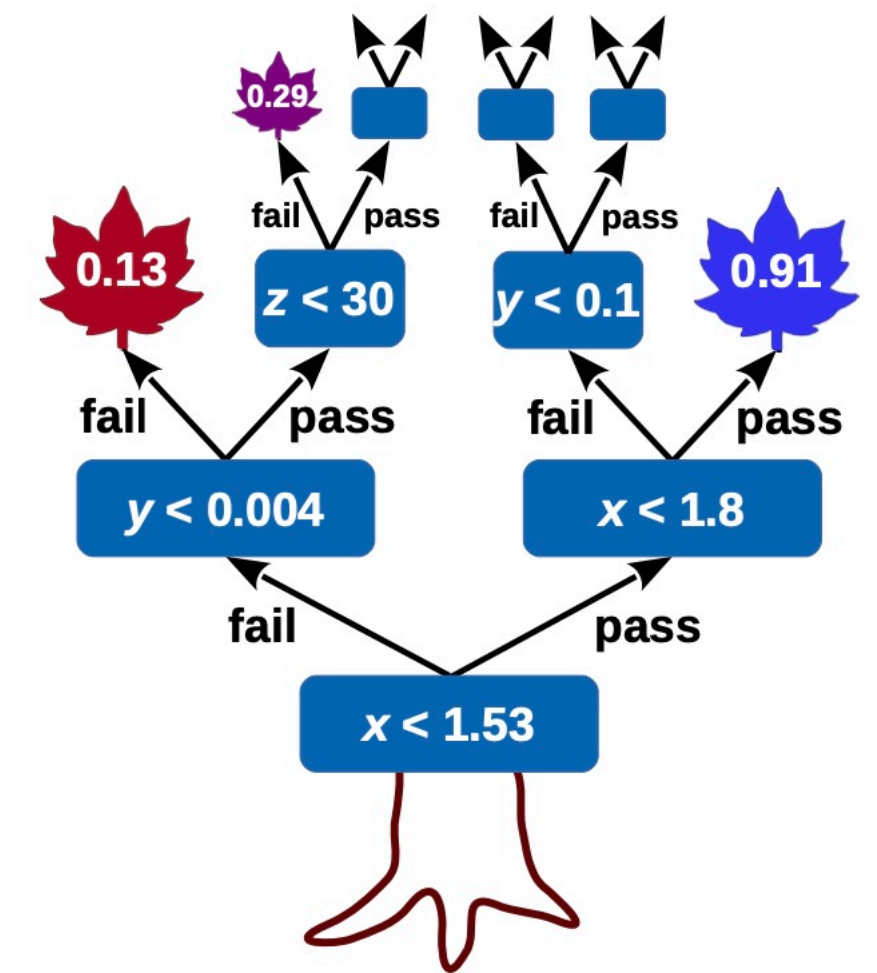
# Quick start
## Apply approaches

- Three most-used approaches
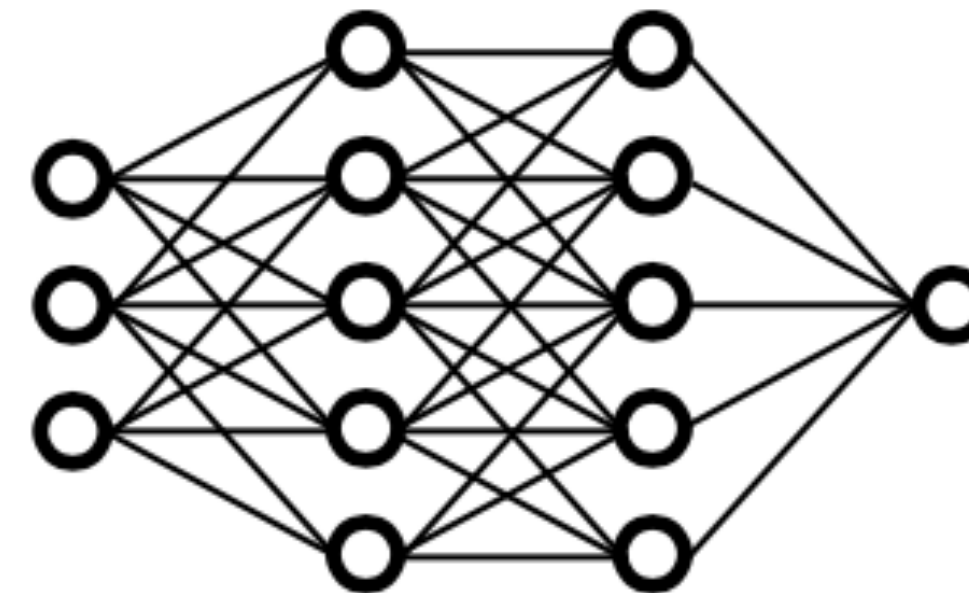
  - cuts

  - trees

  - neural networks



(a) XGBoost with optimized cutoff at 0.9081.
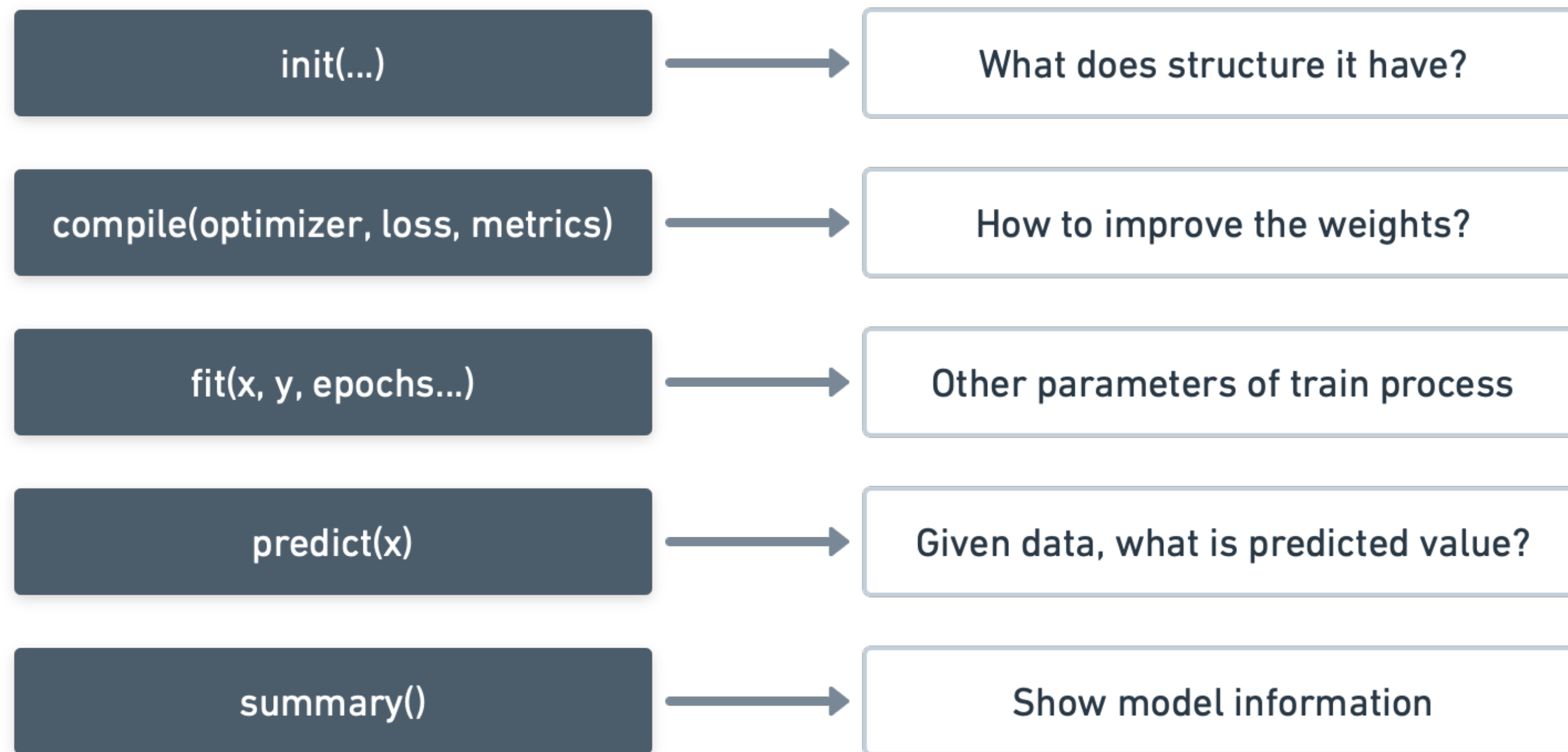
[2108.03125] Beyond Cuts in Small Signal Scenarios



[2206.09645] Boosted decision trees



[1709.04464] Jet Substructure at the Large Hadron Collider

# Quick start

## Apply approaches

| | |
|---|---|
| init(...) | What does structure it have? |
| compile(optimizer, loss, metrics) | How to improve the weights? |
| fit(x, y, epochs...) | Other parameters of train process |
| predict(x) | Given data, what is predicted value? |
| summary() | Show model information |

scikit learn

K Keras

TensorFlow

PyTorch

**Approach** protocol accepts any models with these member functions

# Quick start
## Apply methods

```
1  dataset = load_dataset("./data/zjj_vs_qcd.npz")
2
3  x_train, x_test, y_train, y_test = train_test_split(
4      dataset.samples, dataset.targets, test_size=0.3, random_state=42
5  )
6  x_train, x_val, y_train, y_val = train_test_split(
7      x_train, y_train, test_size=0.2, random_state=42
8  )
9
```

- Load dataset from previous saved location.

- Split train/val/test sets with fixed random seed.

# Quick start
## Apply methods

- **CutAndCount** as a cut-based analysis (CBA).

- **batch_size** should be the whole train set.

- **epochs** have no effect.

```python
approach1 = CBA()
approach1.compile(
    optimizer="adam",
    loss="sparse_categorical_crossentropy",
    metrics=["accuracy"],
)
history = approach1.fit(
    x_train,
    y_train,
    batch_size=len(x_train),
    validation_data=(x_val, y_val),
)
```

# Quick start

## Apply methods

- **GradientBoostedDecisionTree** as a boosted decision tree (BDT).

- **optimizer**, **loss** have no effect in **compile**.

- **batch_size**, **epochs** have no effect in **fit**.

```
1  approach2 = BDT()
2  approach2.compile(
3      metrics=["accuracy"],
4  )
5  history = approach2.fit(
6      x_train,
7      y_train,
8      validation_data=(x_val, y_val),
9  )
```

# Quick start
## Apply methods

- A built-in **ToyMultilayerPerceptron** as MLP for demonstration.

- Normal **Keras** models.

```python
1   approach3 = MLP()
2   approach3.compile(
3       loss="sparse_categorical_crossentropy",
4       metrics=["accuracy"],
5   )
6   approach3.fit(
7       x_train,
8       y_train,
9       batch_size=128,
10      epochs=20,
11      validation_data=(x_val, y_val),
12  )
```

# Quick start
## Apply methods

```
1    from keras.metrics import Accuracy, AUC
2    from sklearn.metrics import roc_curve
3    from hml.metrics import MaxSignificance, RejectionAtEfficiency
```

- **MaxSignificance** calculates the maximum significance under uniform distributed thresholds.

$$\text{significance} = \sqrt{S/(S+B)}$$

- **RejectionAtEfficiency** $(1/\varepsilon_b$ at $\varepsilon_s = 50\,\%)$ calculates the background rejection at a given signal efficiency.
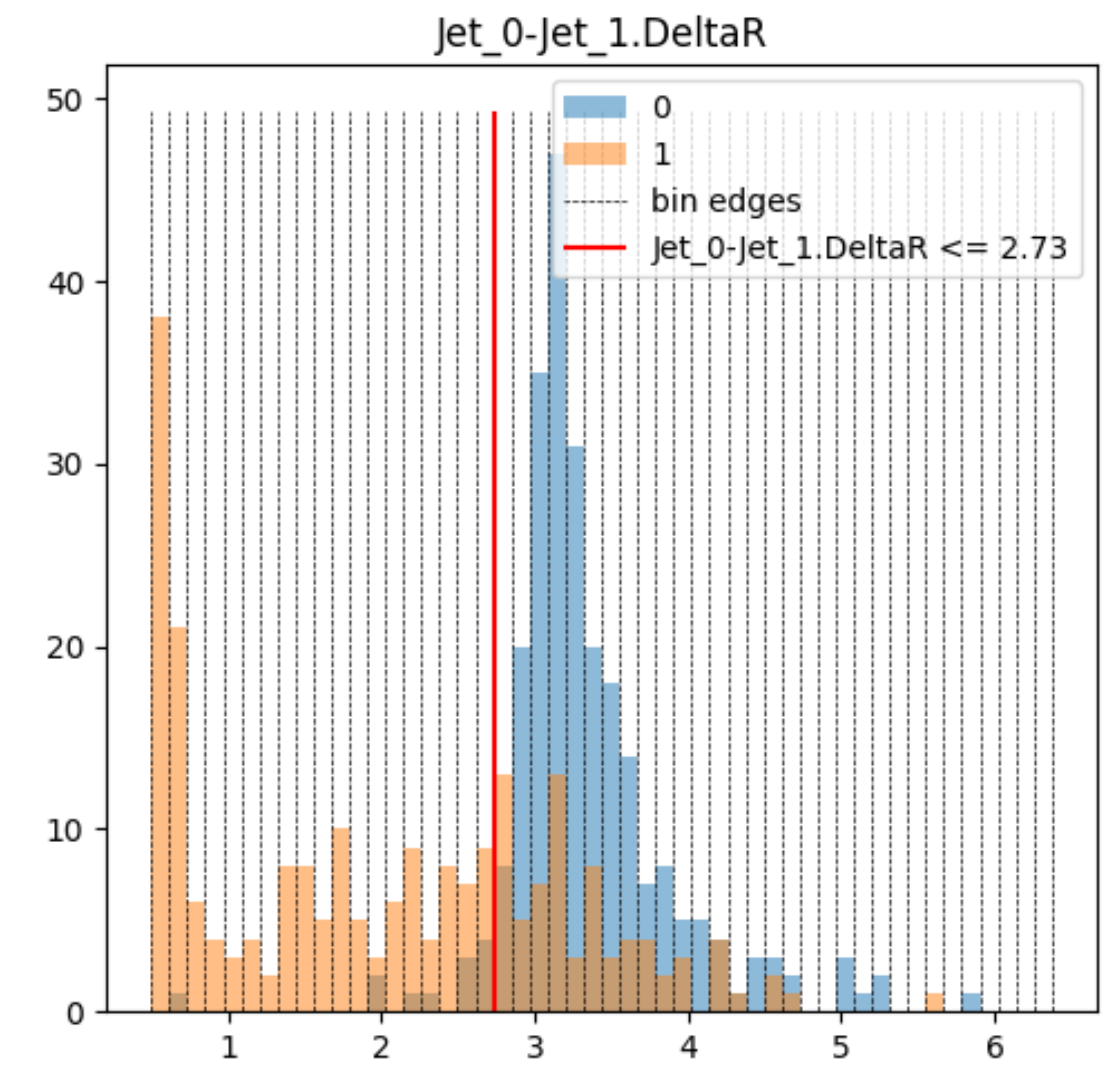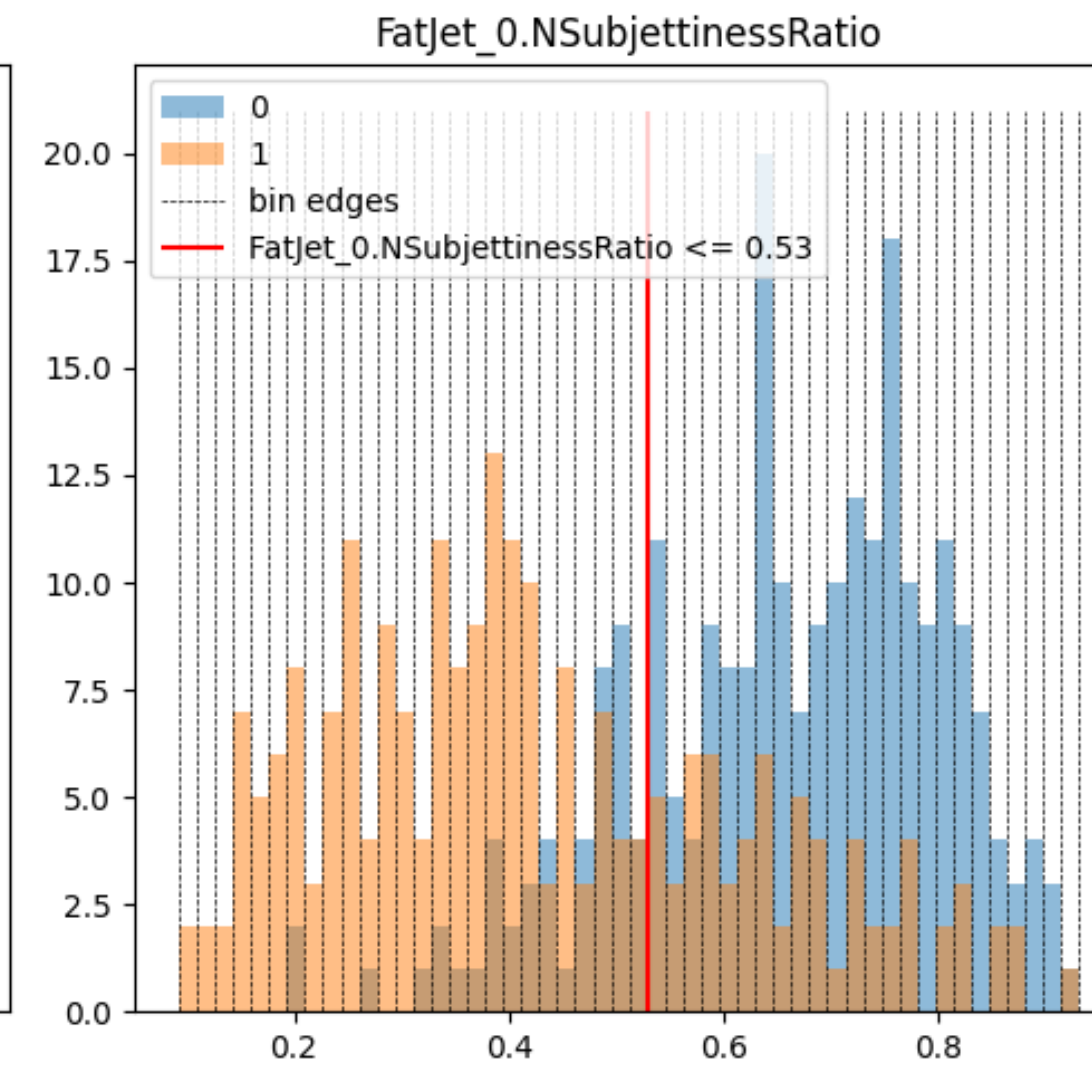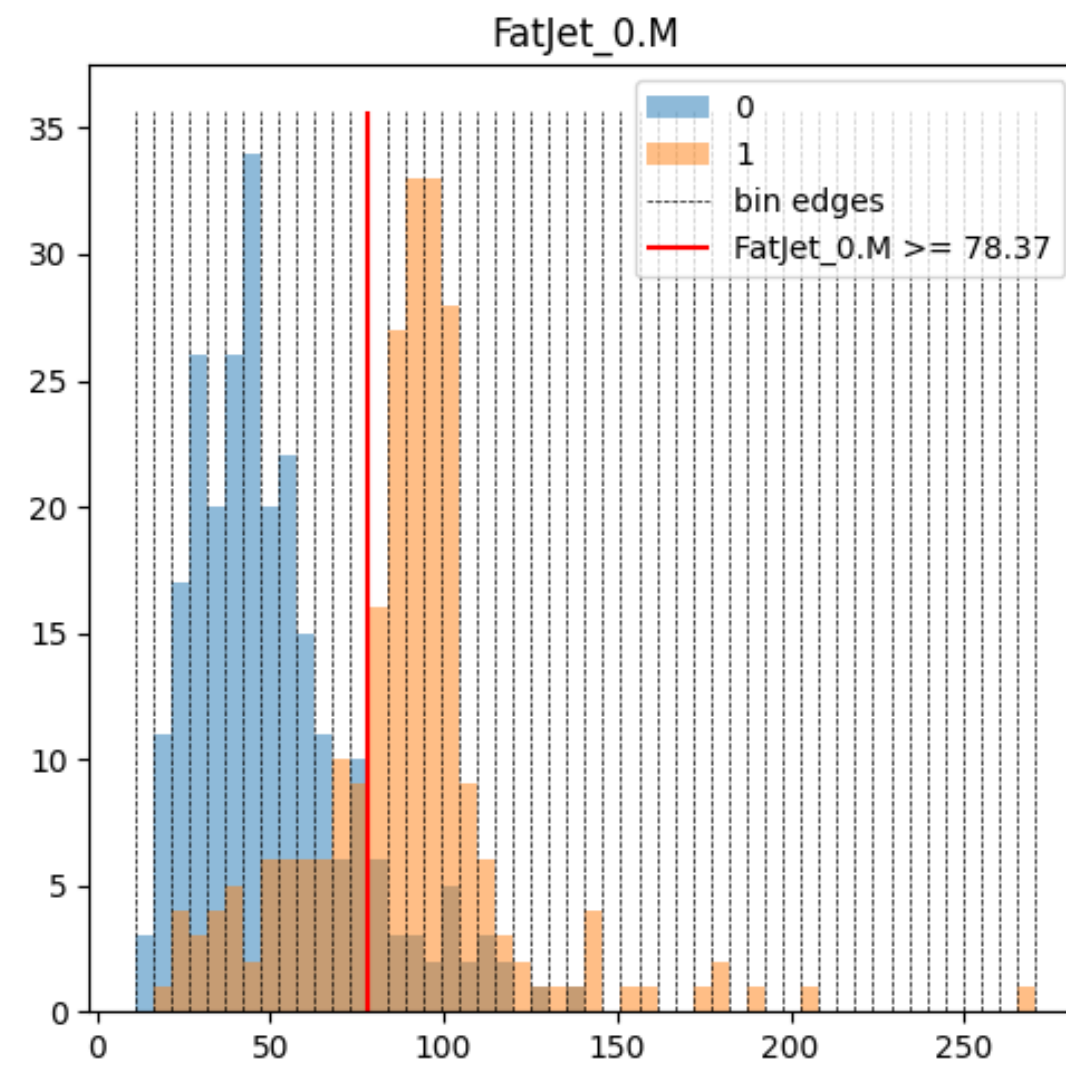
20

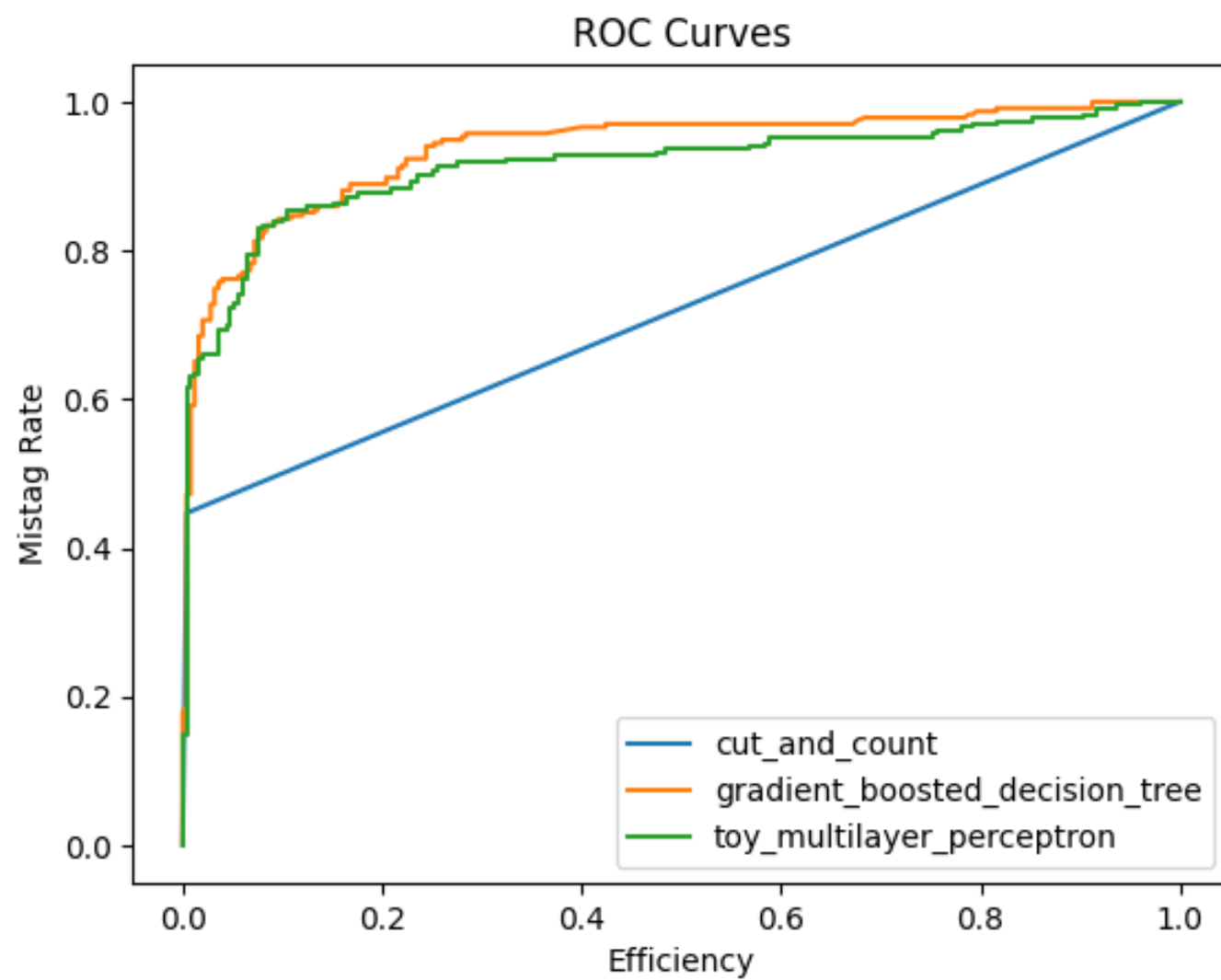# Quick start

## Apply methods

```
1   y_prob = approach.predict(x_test, verbose=0)
2   y_pred = y_prob.argmax(axis=1)
```

- Get class probability and predictions to evaluate all approaches.

# Quick start

## Apply methods

| Name | ACC | AUC | MaxSignificance | RejectionAtEfficiency |
|------|-----|-----|-----------------|-----------------------|
| cut_and_count | 0.729897 | 0.721404 | 10.1985 | 1 |
| gradient_boosted_decision_tree | 0.876289 | 0.93703 | 13.3053 | 125 |
| toy_multilayer_perceptron | 0.610309 | 0.915157 | 6.78387 | 249.999 |

# Future
## Roadmap

- 0.3.x

  - Support loading parts of existing datasets from Zenodo, Hugging Face, GitHub.

  - Support image and graph representation and ToyCNN, ToyGNN to test.

- 0.4.*

  - Support for frameworks from **Scikit-HEP**

**zenodo**

**GitHub**

🤗 **Hugging Face**

**Scikit HEP**

**uproot**

**Awkward Array**

**VECTOR**

**F⭐STJET**

# Future
## Roadmap

- WELCOME your priceless contributions ❤️

- WELCOME any comments to my email: star9daisy@outlook.com

- Check the documents for more details: https://star9daisy.github.io/hep-ml-lab/

- Find source code here: https://github.com/Star9daisy/hep-ml-lab

- Have a try today! pip install hep-ml-lab

# THANK YOU!