

Unsupervised tagging of semivisible jets with energy-based autoencoders in CMS

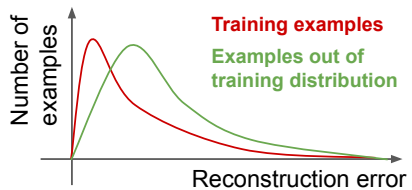
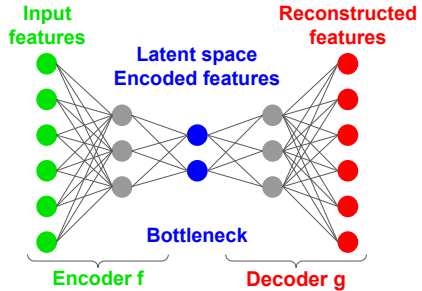
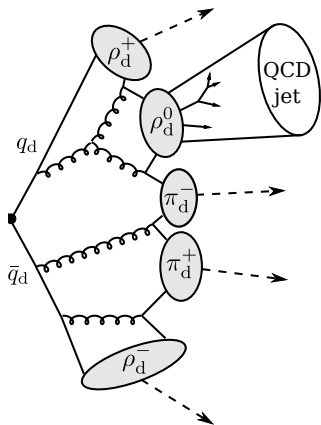
Florian Eble, on behalf of the CMS collaboration

ETH zürich

01/05/2024

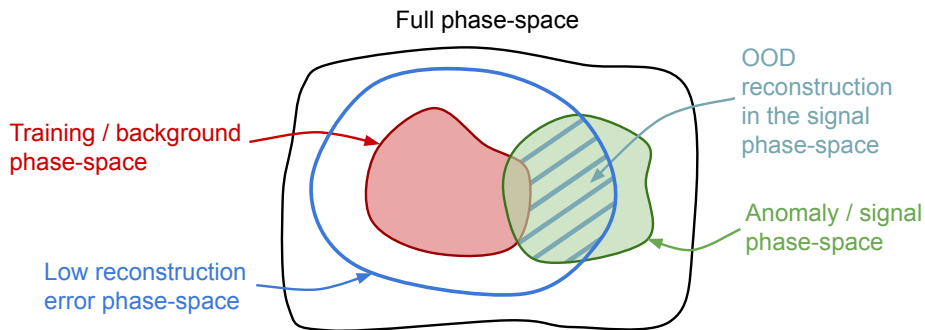
EuCAIFCon 2024

?

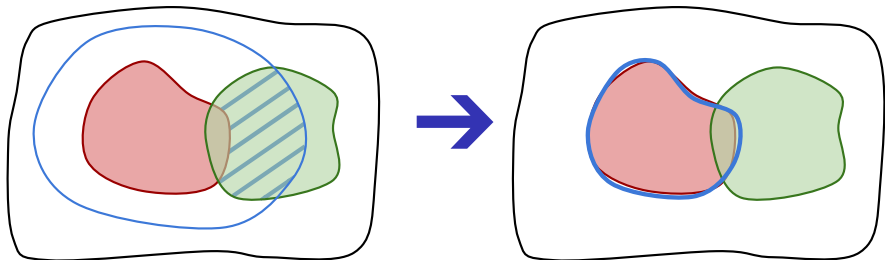


- **AEs are free to minimize reco error outside the background phase space!** including the unknown signal phase space...

→ This is the problem of **OOD reconstruction**:



- NAE features a mechanism to suppress OOD reconstruction
- First introduced in [arXiv:2105.05735](https://arxiv.org/abs/2105.05735) and used in HEP in [arXiv:2206.14225](https://arxiv.org/abs/2206.14225)



- NAE paradigm:
 - Define a probability distribution p_θ so that high probability regions have low reco error
 - Sample from p_θ via a MCMC
 - Minimize the distance between the **background** and p_θ probability distributions
- We propose a different metric to measure this distance, using the Earth Mover's Distance (a.k.a Wasserstein distance) and train NAEs in a fully signal-agnostic fashion

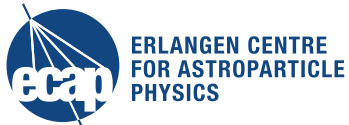
Multi-class classification of gamma-ray sources and the nature of excess of GeV gamma rays near the Galactic center

ERLANGEN CENTRE
FOR ASTROPARTICLE
PHYSICS

Dmitry Malyshev

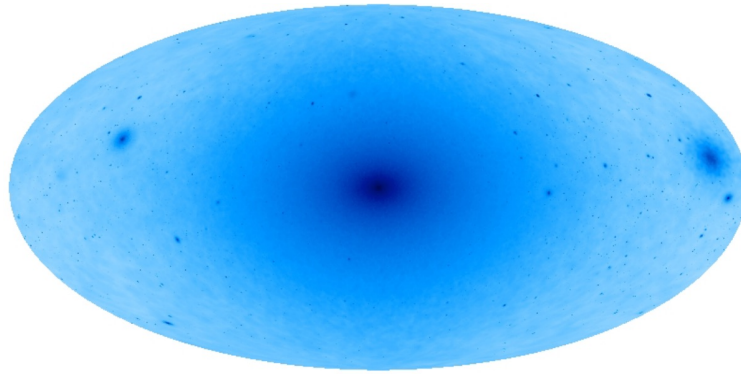
EuCAIFcon, Amsterdam, 30.04 – 03.05 2024

Poster session B, location 116



Dark matter in the Galactic center

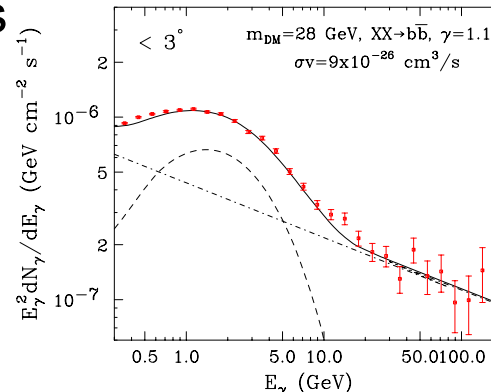
- Galactic center (GC) is the strongest possible source of dark matter (DM) annihilation signal



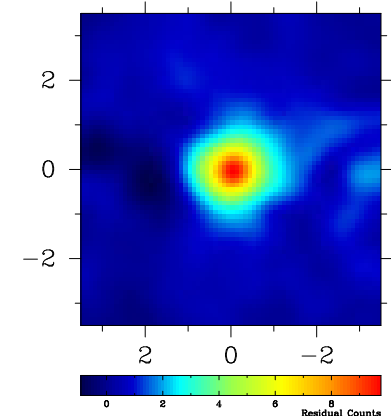
Via Lactea II simulation of a Milky-way-like galaxy. Kuhlen et al. (2009)

- Excess consistent with DM annihilation was detected in Fermi-LAT gamma-ray data two months after the data became public

Goodenough & Hooper (2009), Vitale & Morselli (2009), Hooper & Linden (2011), Abazajian & Kaplinghat (2012), Hooper & Slatyer (2013), Gordon & Macias (2013), Calore et al. (2015), Daylan et al. (2016), Ajello et al. (2016), Ackermann et al. (2017) etc



Goodenough & Hooper (2009)



Abazajian & Kaplinghat (2012)

Astrophysical explanation

- A population of millisecond pulsars (MSPs) near the GC can explain the Galactic center excess (GCE)

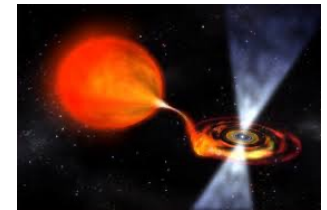
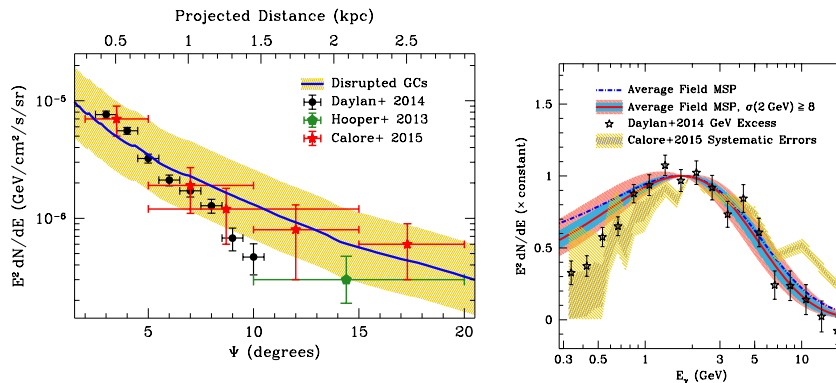


Image credit: NASA



Brandt & Kocsis (2015)

Statistical studies

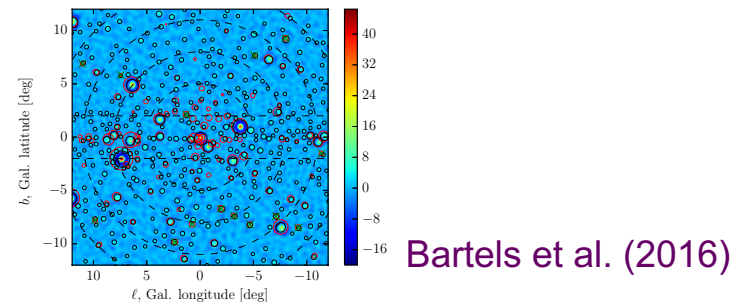
Lee et al. (2015, 2016), Bartels et al. (2016), Leane & Slatyer (2019, 2020), Zhong et al. (2020), List et al. (2020), Calore et al. (2021), Mishra-Sharma & Cranmer (2022), Caron et al. (2023), Manconi et al. (2024) etc.

- Based on statistical properties of the Gamma-ray data

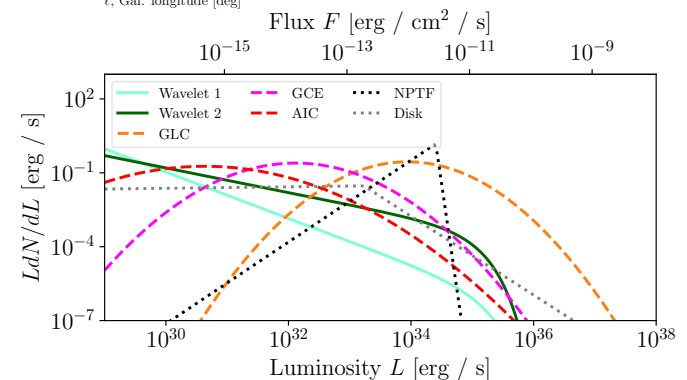
Population studies

Brandt & Kocsis (2015), Hooper & Linden (2016), Bartels et al. (2018), Ploeg et al. (2020), Dinsmore & Slatyer (2022) etc.

- Associated (bright) MSPs are used to constrain the models

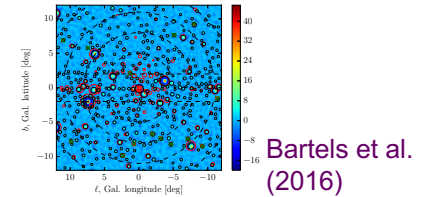
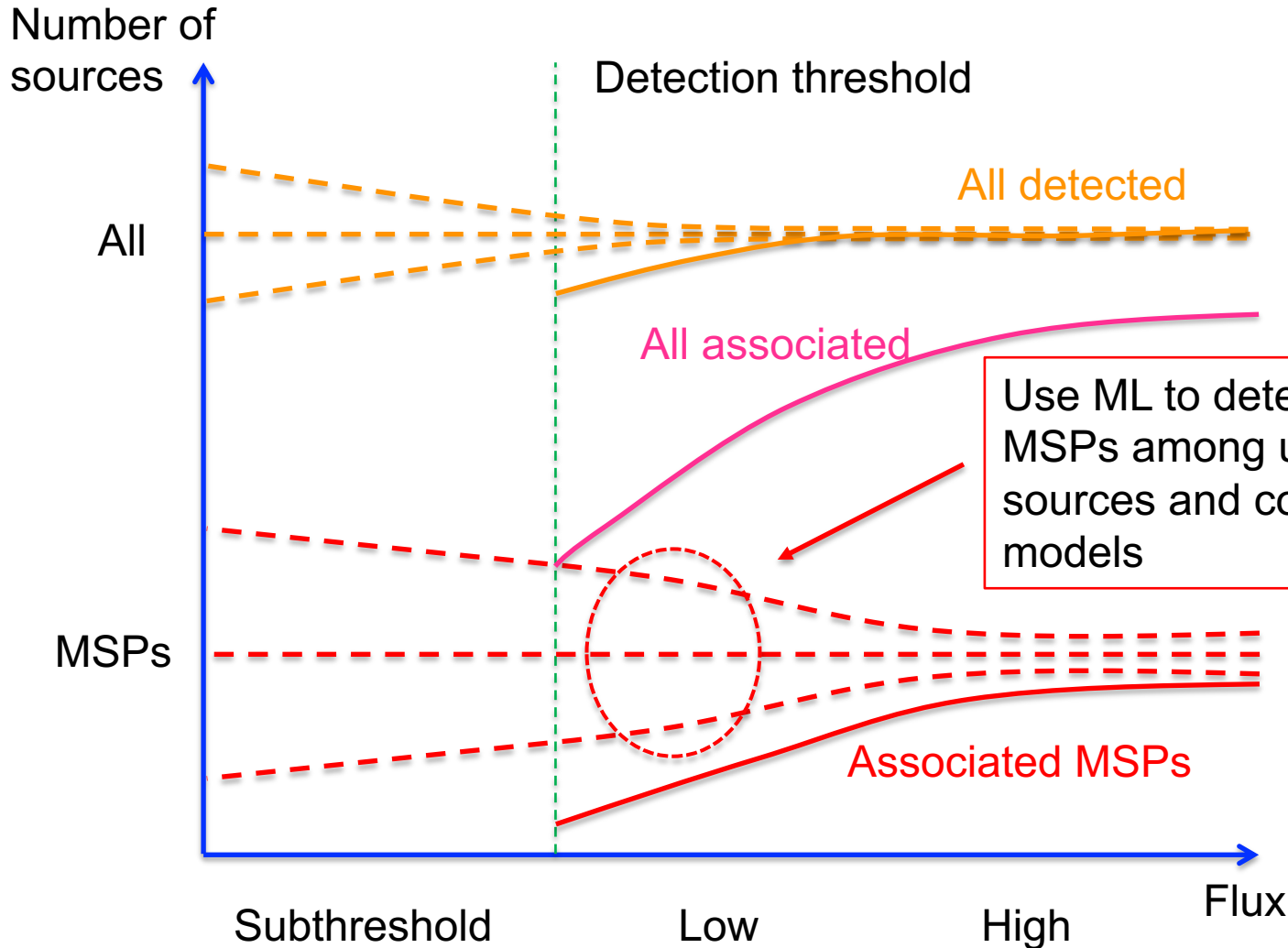


Bartels et al. (2016)



Dinsmore & Slatyer (2022)

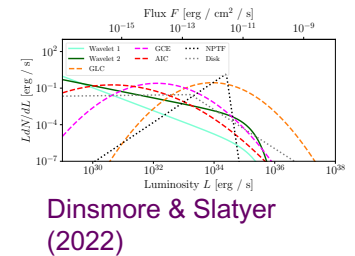
How can machine learning help?



Statistical studies

Use ML to determine MSPs among unassociated sources and constrain MSP models

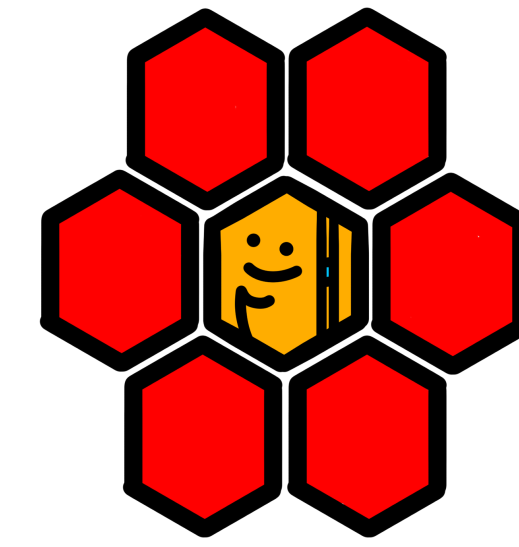
Population studies



Details and results: [poster session B, location 116](#) & [arXiv:2401.04565](#)

b-hive:

a modular training framework for state-of-the-art object-tagging within the Python ecosystem at the CMS experiment



Motivation:

- Everybody wants to do Machine-Learning trainings
- Full end-to-end pipeline is way harder than an example Notebook
 - Big data processing (ROOT files)
 - Conversion into a ML-friendly format (.npy/ .npz)
 - Deploy state-of-the-art models
 - ParticleNet (graph-convolutions)
 - Transformer models

Also:

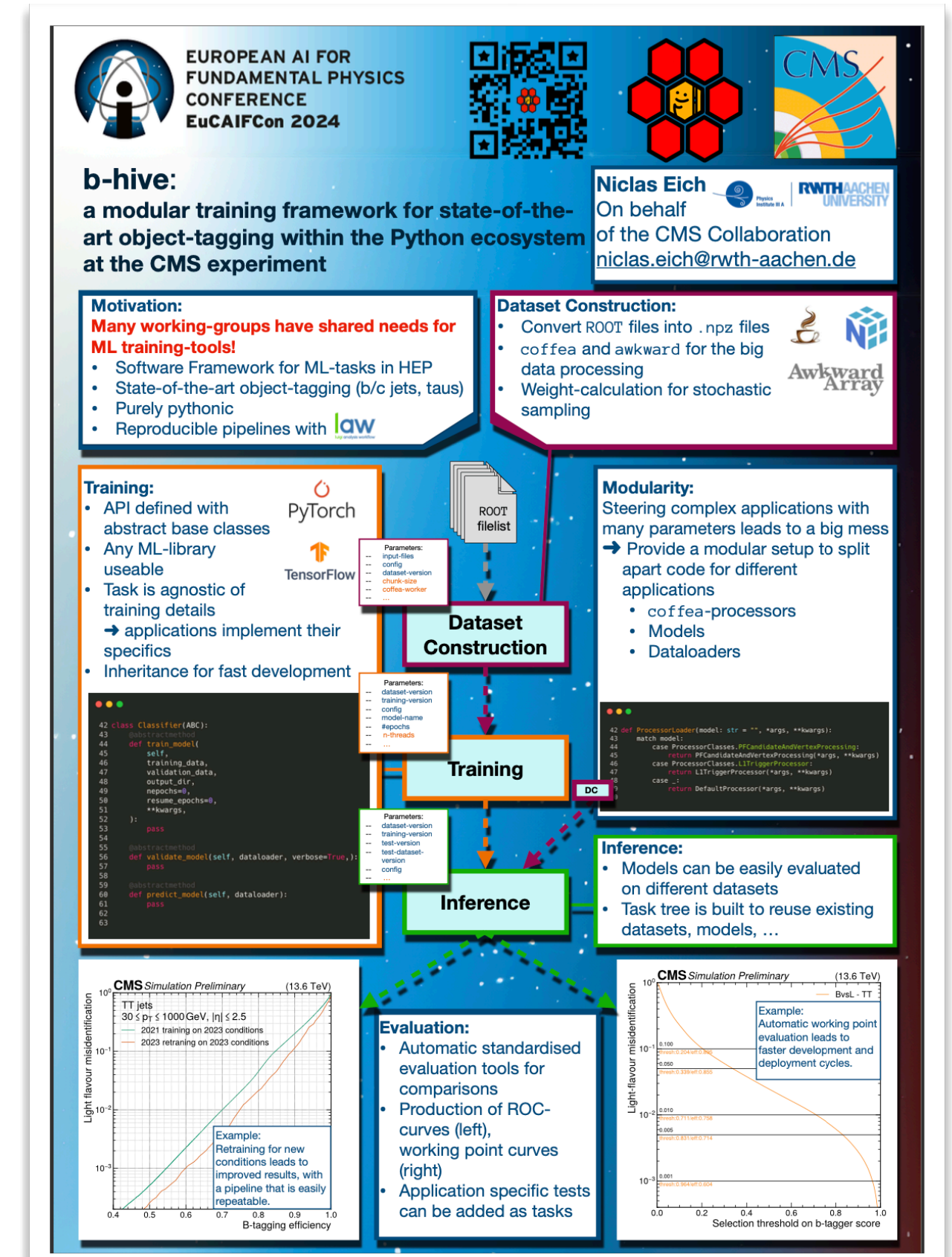
- Have clearly defined workflows (not your 7 bash scripts!)
- Make trainings repeatable
- Standardized evaluation tools

Niclas Eich
On behalf of
the CMS Collaboration

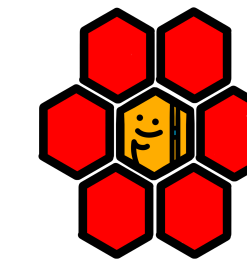


III. Physikalisches
Institut A

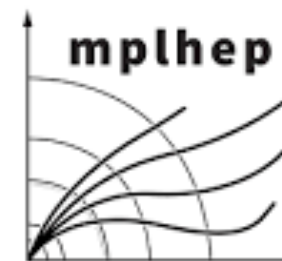
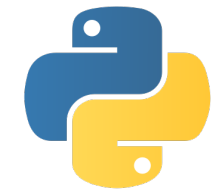
RWTHAACHEN
UNIVERSITY



b-hive attacks these problems

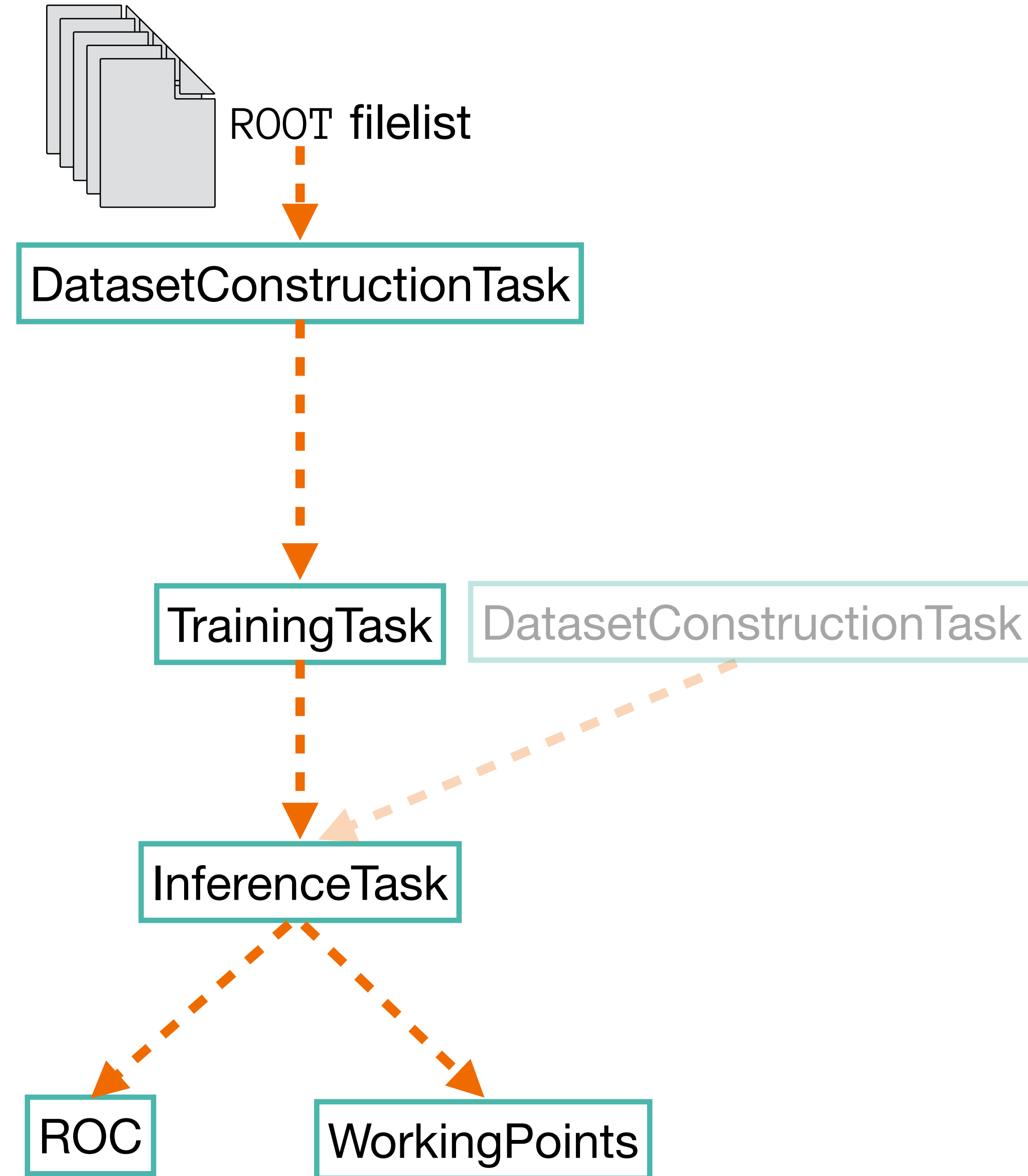


- Pythonic training framework
 - Workflow management with law
 - coffea, awkward, numpy for the heavy data lifting
 - No ML-framework lock in
 - TensorFlow and PyTorch can be used
- Modular Setup



- Easy configuration for different working-groups
 - New applications are embedded in the pipeline
- Knowledge-sharing by code-sharing
- Have a look: [CERN-CMS-DP-2024-020](#)

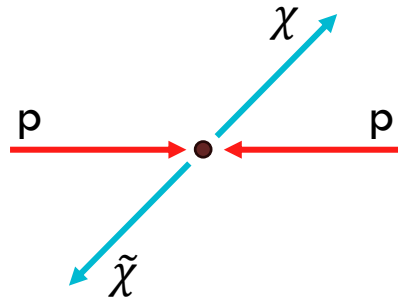
```
42 bins_pt:
43 - 30
44 - 100
45 - 500
46 - 2000
47
48 bins_eta:
49 - -2.5
50 - -2.0
51 - -0.5
52 - 0.5
53 - 2.0
54 - 2.5
55
56 treename:
57 "producer/custom_tree"
58
59 global_features:
60 - "jet_pt"
61 - "jet_eta"
62 - "n_Cpfcand"
63 - "n_Npfcand"
64 - "nsv"
65 - "npv"
66
67 pf_candidate_features:
68 - "pfcand_pt"
69 - "pfcand_eta"
70 - "pfcand_trackPt"
71 - "pfcand_trackEta"
72 - "pfcand_trackDeltaR"
73
74 truths:
75 - "isB"
76 - "isUDS"
77 - "isG"
78
```



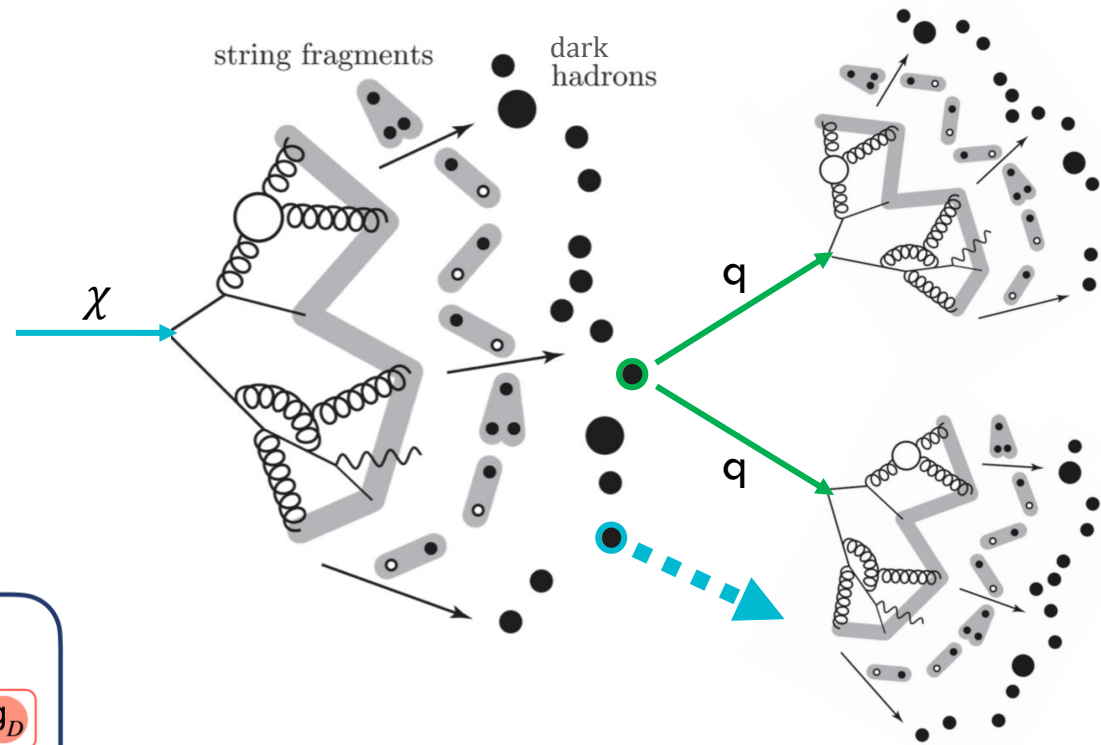
ENERGY-BASED GRAPH AUTOENCODERS FOR SEMIVISIBLE JET TAGGING IN THE LUND REPRESENTATION

Annapaola De Cosa, Roberto Seidita, Florian Eble, Christoph Ribbe¹

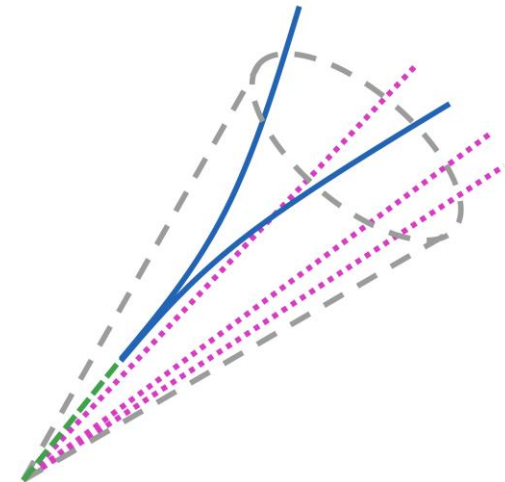
¹Did most of the work but could not join



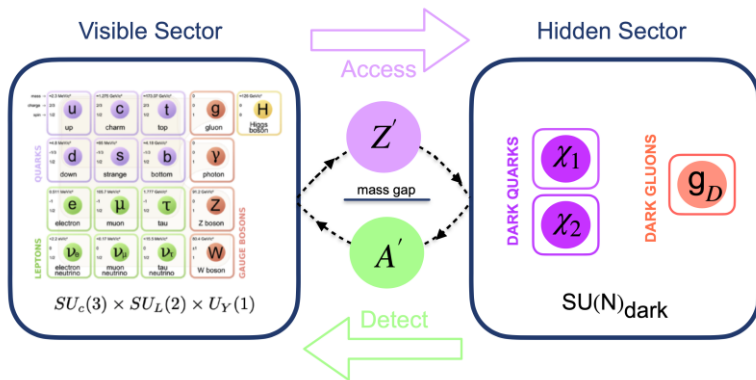
Dark quarks from a proton-proton collision



Shower of **stable (invisible)** and **unstable (visible)** dark particles

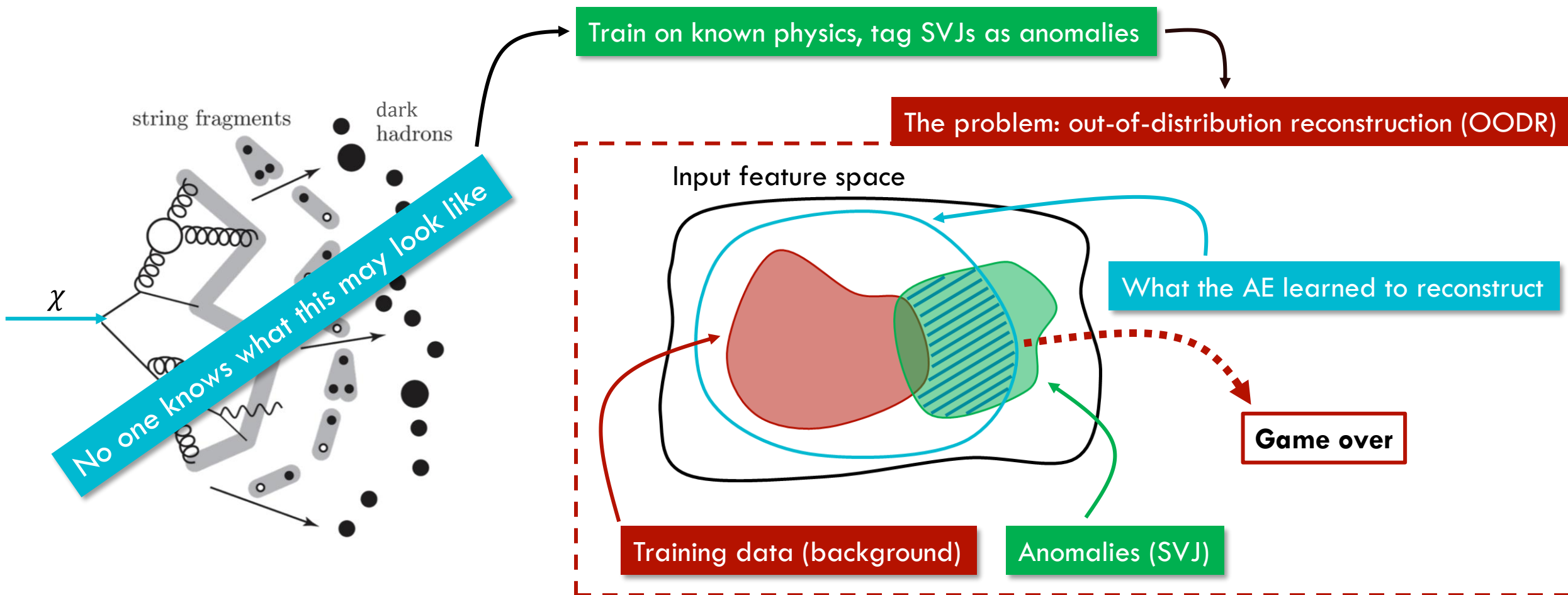


A semivisible jet (SVJ): a spray of particles with missing stuff between them



WHY UNSUPERVISED?

Annapaola De Cosa, Roberto Seidita,
Florian Eble, Christoph Ribbe



JETS AS LUND GRAPHS

Annapaola De Cosa, Roberto Seidita,
Florian Eble, Christoph Ribbe

