# Physicists Learning from Machines Learning

## Smart but Interpretable Neural Networks for Physics at the LHC
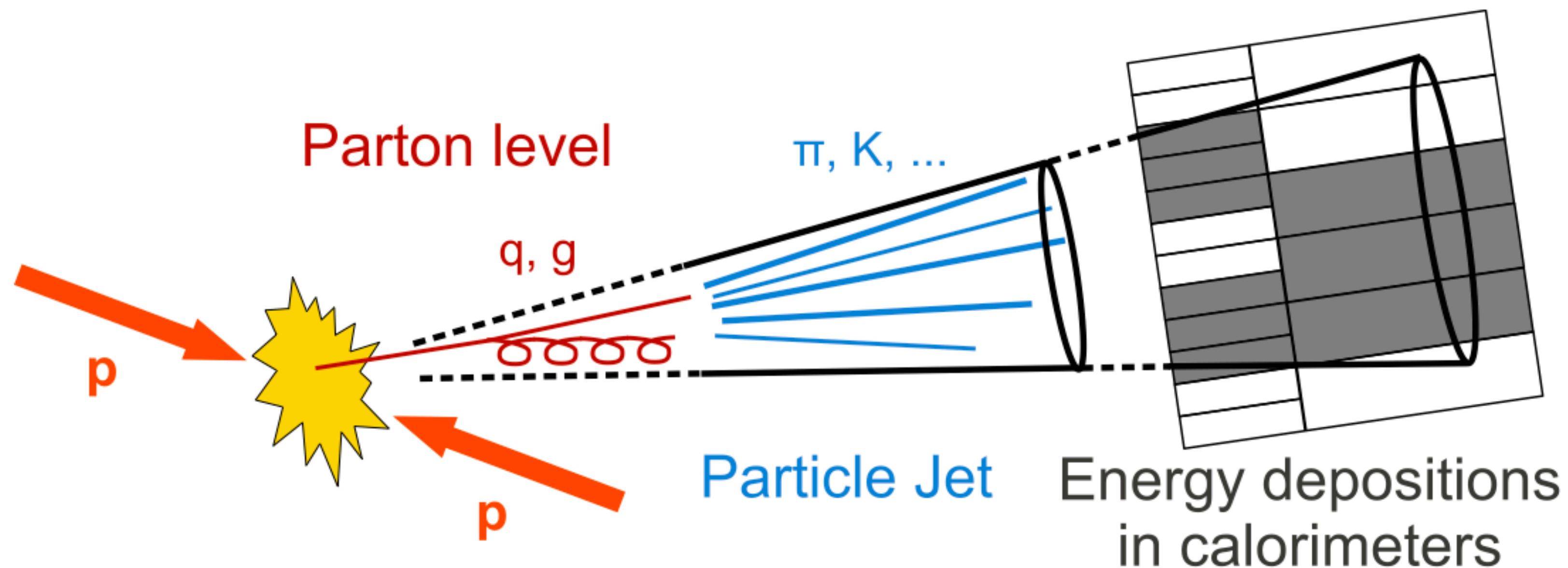
Authors: Taylor Faucett, Daniel Whiteson & Jesse Thaler

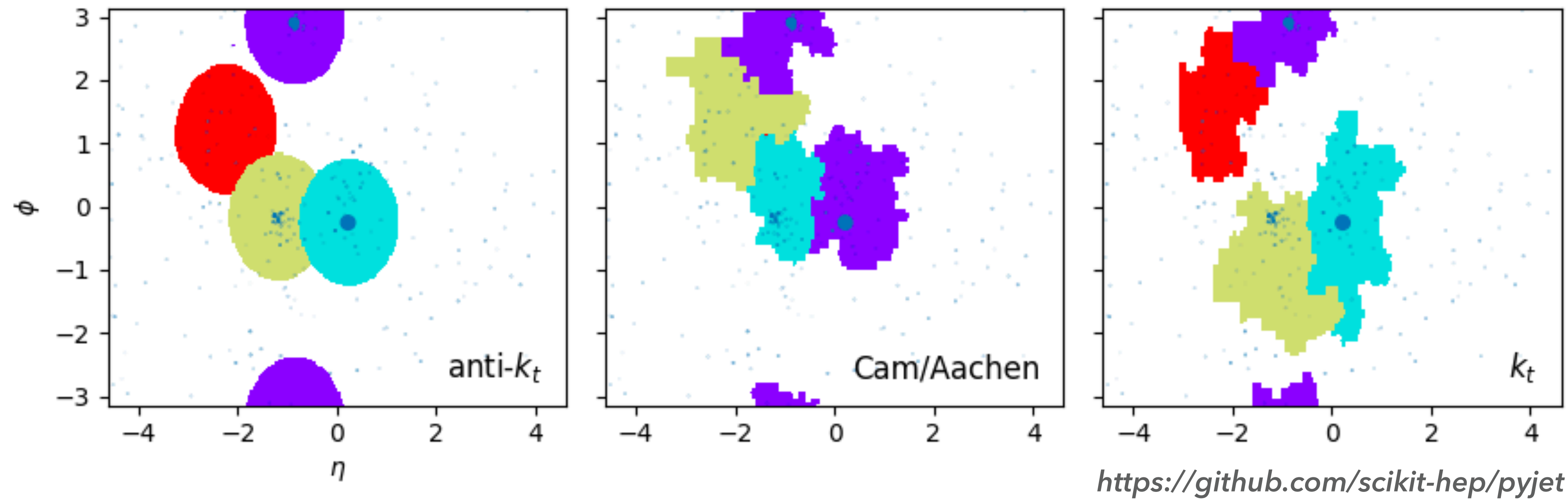**Mapping machine-learned physics into a human-readable space**

**@ Clermont-Ferrand**

▸ Particle collisions create particles with non-zero color charge (i.e. quarks & gluons)

▸ Free quarks/gluons hadronize to produce hadrons (e.g. mesons and baryons)

▸ "Jets" are collimated sprays of many hadrons in a cone.

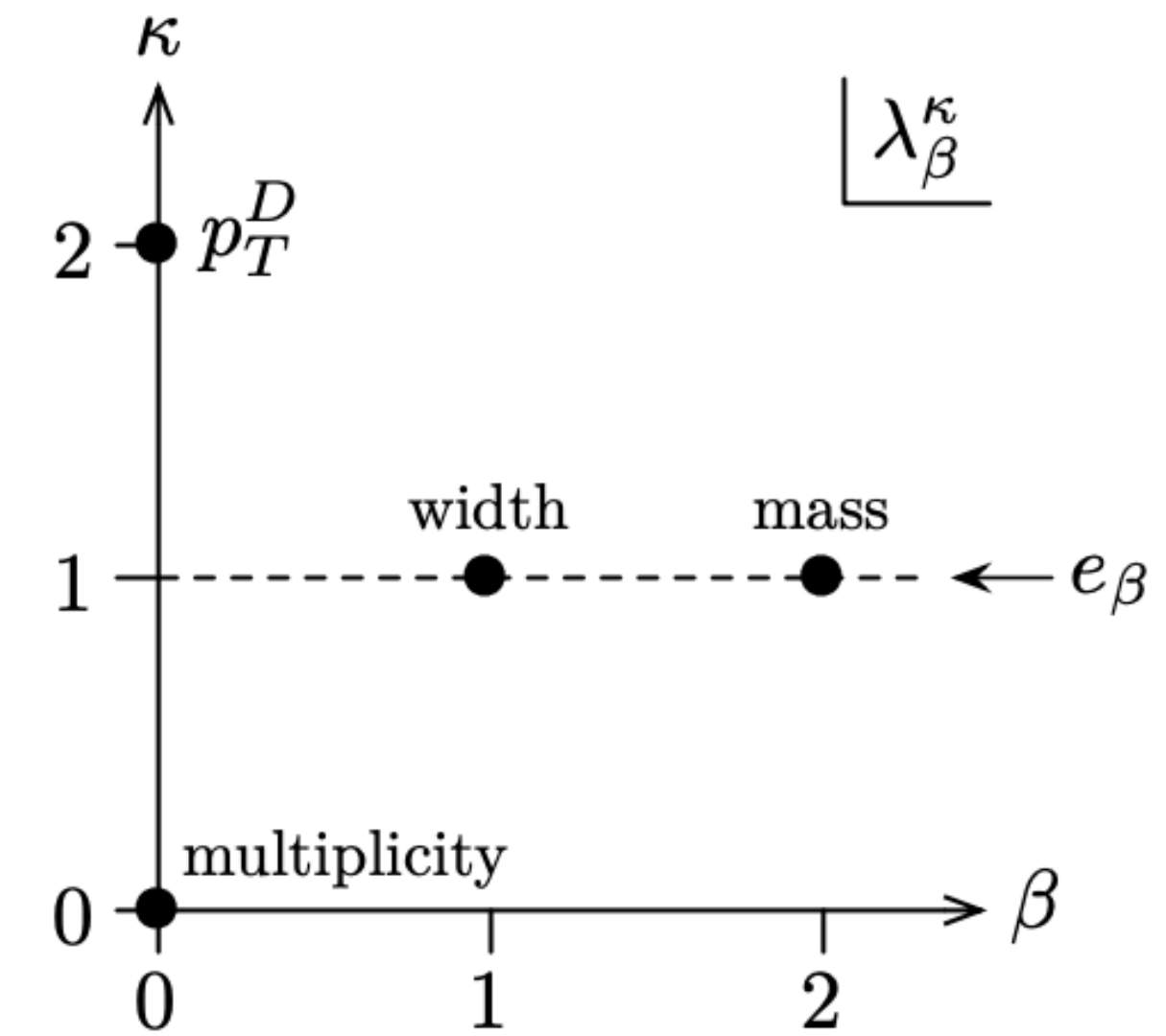▸ Identifying jets and different kinds of jets help distinguish high-energy processes.

Parton level  
π, K, ...  
q, g  
p  
p  
Particle Jet  
Energy depositions in calorimeters

*https://cms.cern/tags/particle-jet*

▸ Jets are isolated in the detector using clustering algorithms

▸ How do we distinguish different types of jets from one another?



*https://github.com/scikit-hep/pyjet*

*What variables would be useful for training ML to classify*

*different jets?*

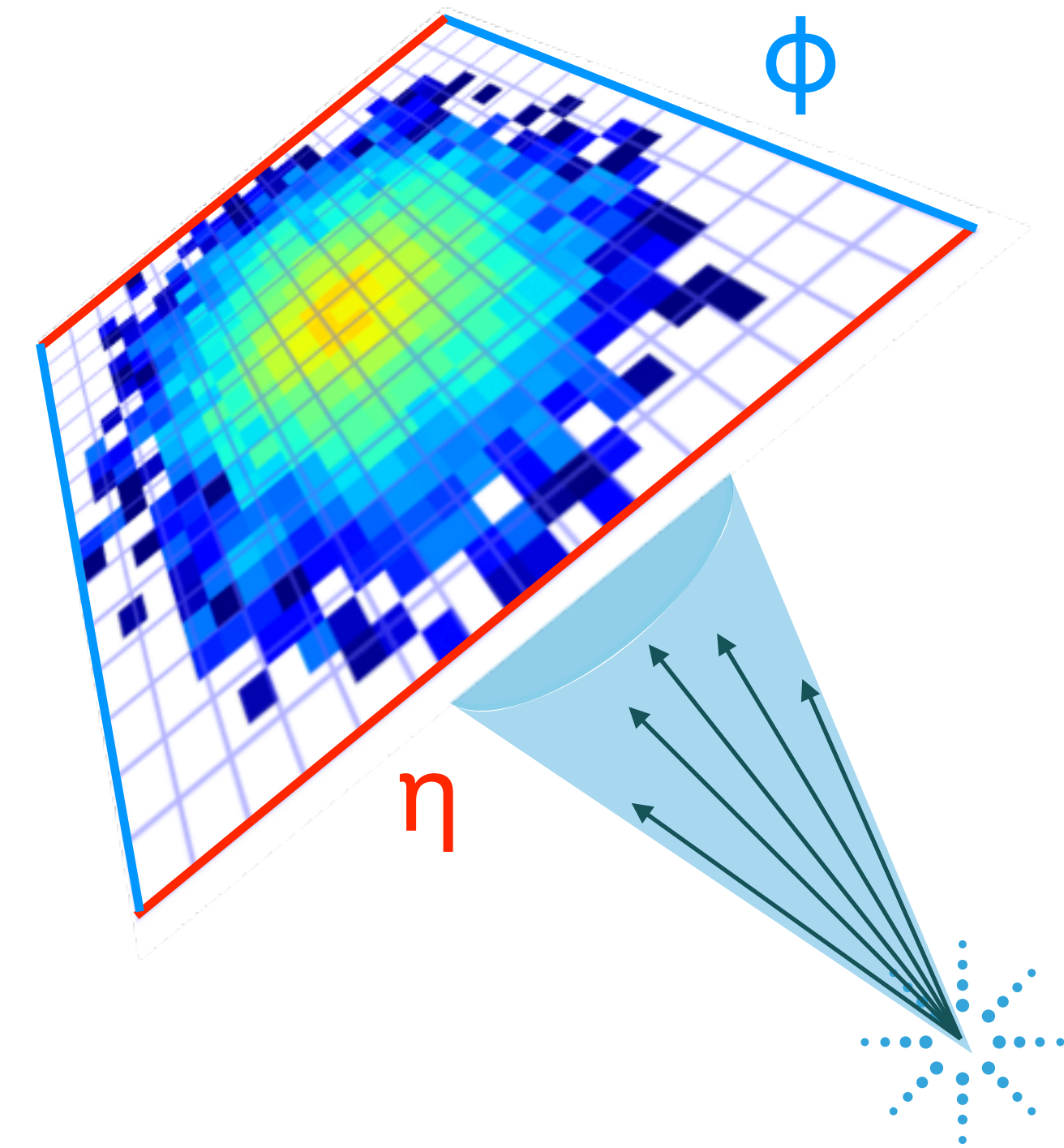$$\lambda_\beta^\kappa = \sum_{i \in jet} z_i^\kappa \theta_i^\beta$$

▸ Quark/Gluon discrimination?

   ▸ Gluons produce more particles.

   ▸ Look at "jet multiplicity" and similar proxies/weightings of momentum/ angular separation

   ▸ From "generalized angularity": multiplicity, LHA, pTD…

▸ What about Jets with multiple sub-jets?

   ▸ W/Z/h decay to 2 jets

   ▸ We invent "N-subjettiness" to quantify separable jet substructure



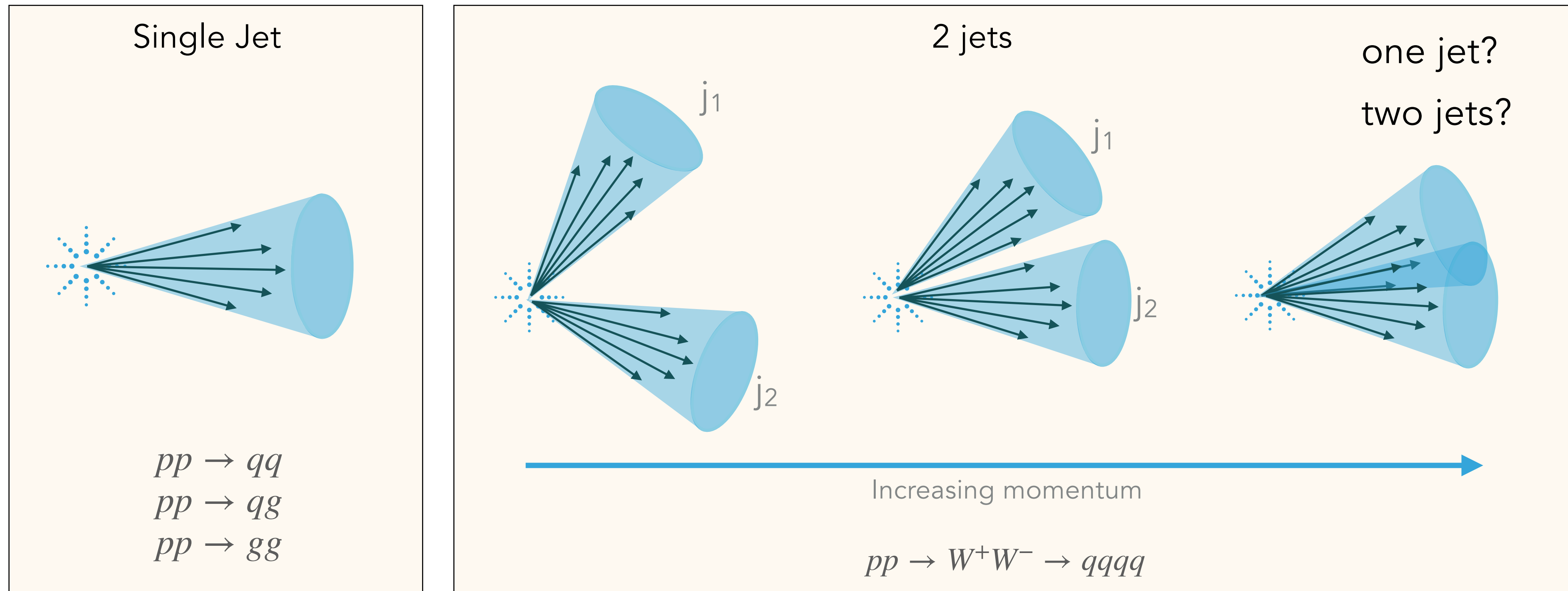$$\tau_N = \frac{1}{d_0} \sum_k p_{T,k} \, min \left\{ \Delta R_{0,k}, \Delta R_{1,k} \ldots \Delta R_{N,k} \right\}$$

▸ Do we need to hand-pick observables in every study?

▸ What if we let a Deep Neural Network learn to solve the problem?

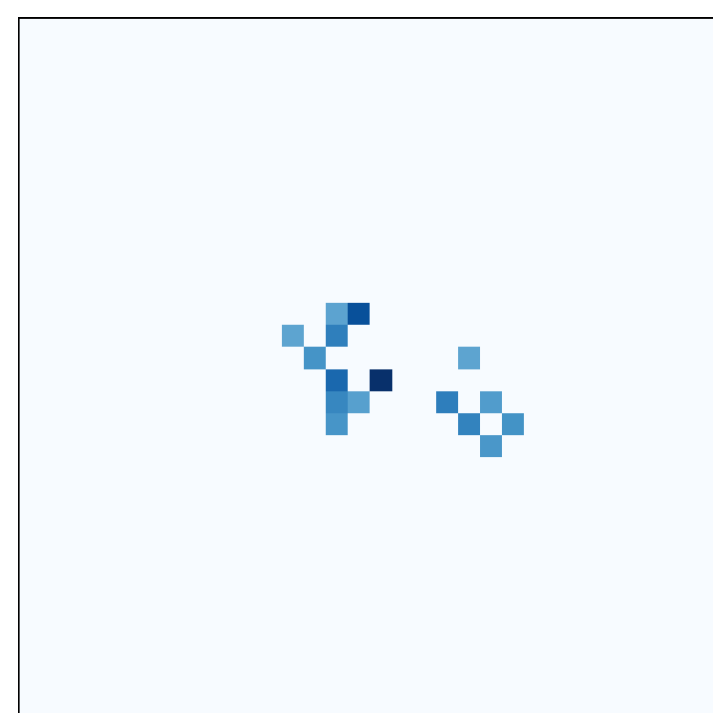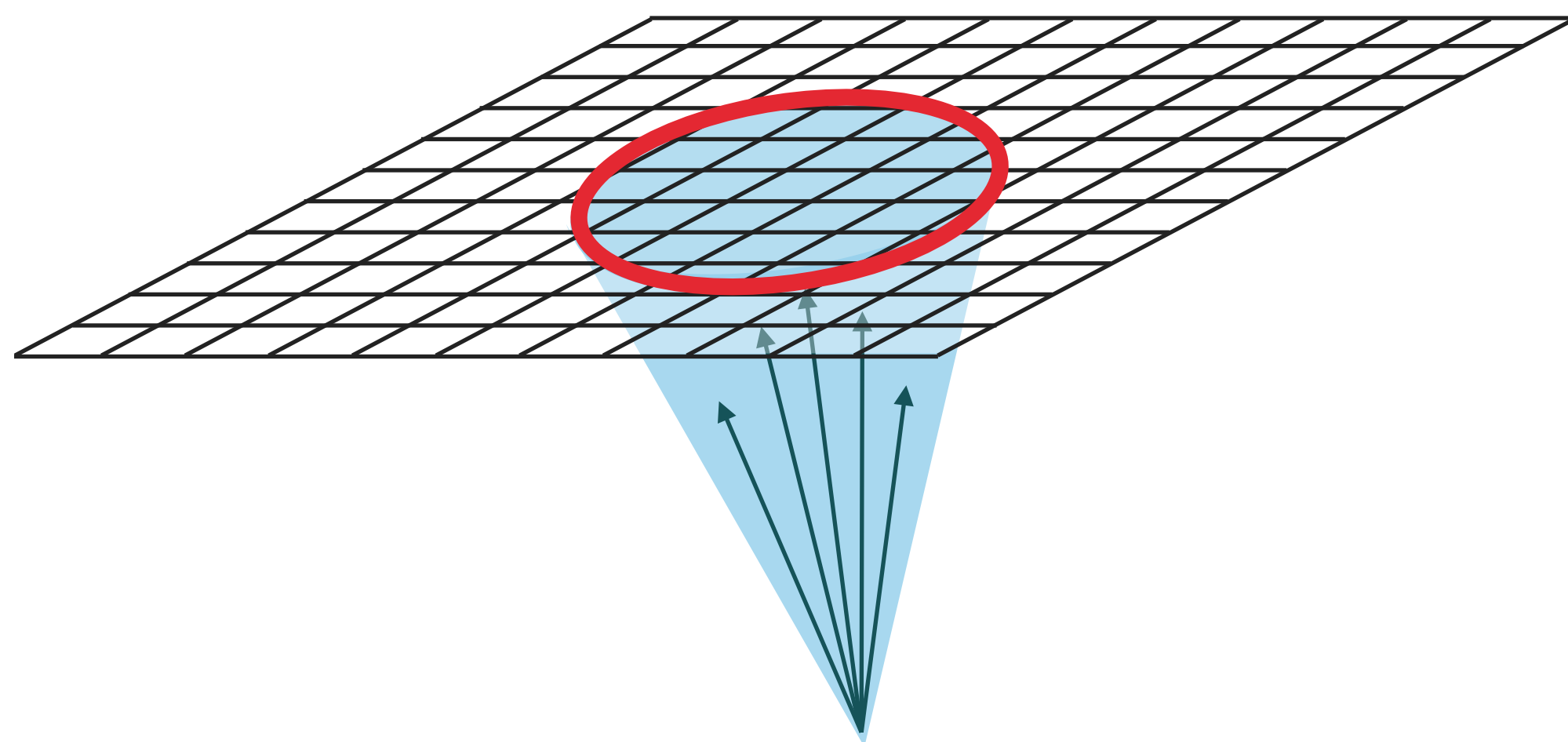▸ Can a CNN to learn to classify directly from the calorimeter data?

- $E_T$ = Transverse Energy
- Position ($\eta$, $\phi$)
- $\eta = -\ln(\tan(\theta/2))$

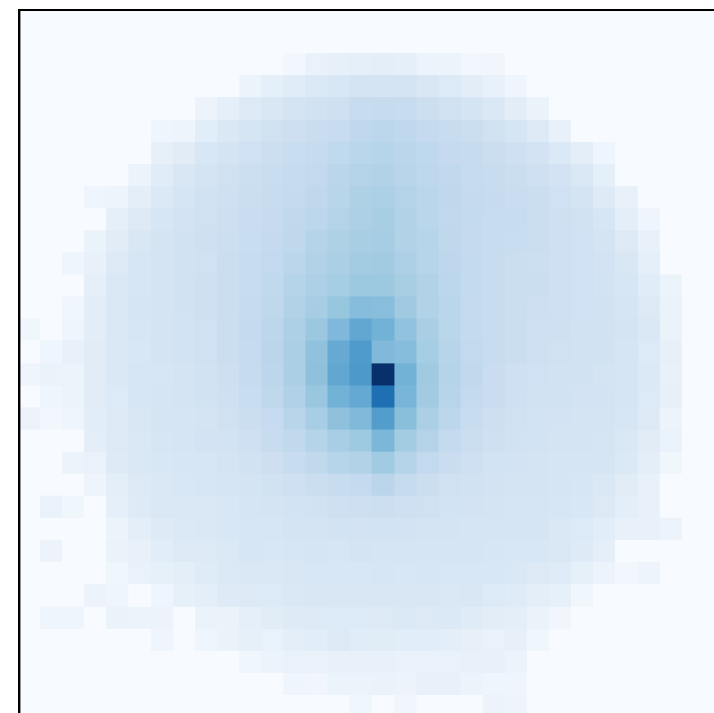Boosted W bosons (W → qq') create highly collimated di-jets.

Can a CNN separate boosted W jets from QCD jets?



Single Jet

$$pp \to qq$$
$$pp \to qg$$
$$pp \to gg$$

2 jets

$j_1$

$j_2$

one jet?
two jets?

Increasing momentum

$$pp \to W^+W^- \to qqqq$$

## QCD Jet (q, g)

## W jet
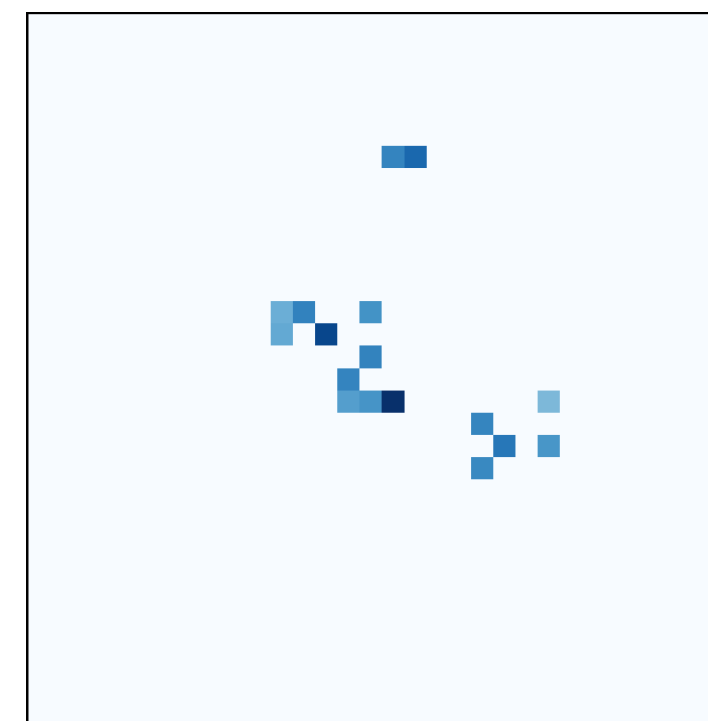


1 Event

Average of all events

1 Event

Average of all events

## 6 HL Variables (HL network)



**VS.**

## Jet Images (LL network)
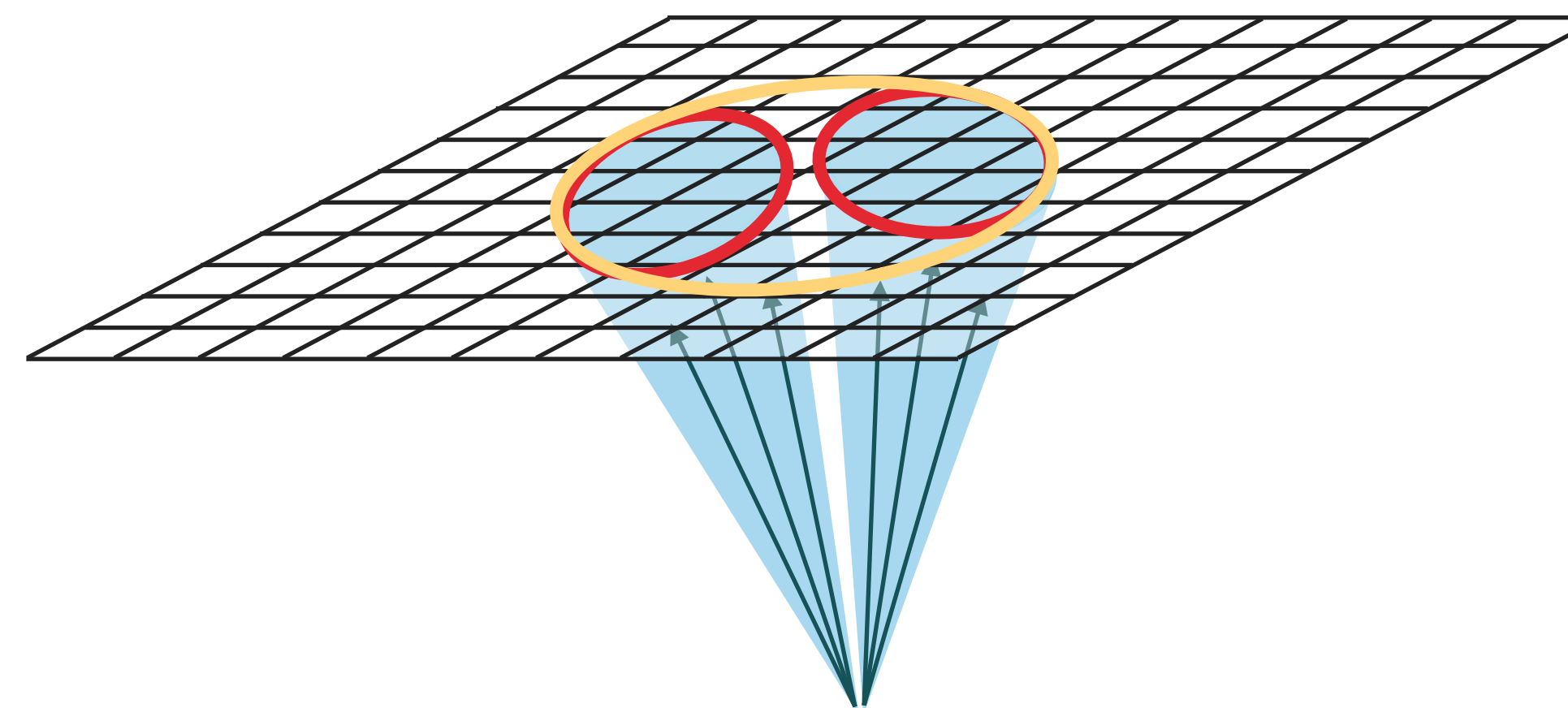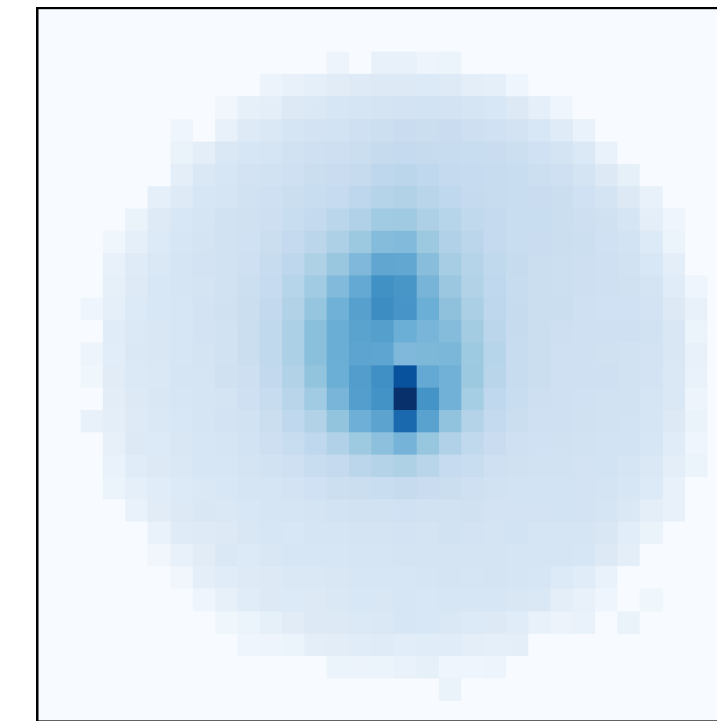


1 Event

Average of all events

1 Event

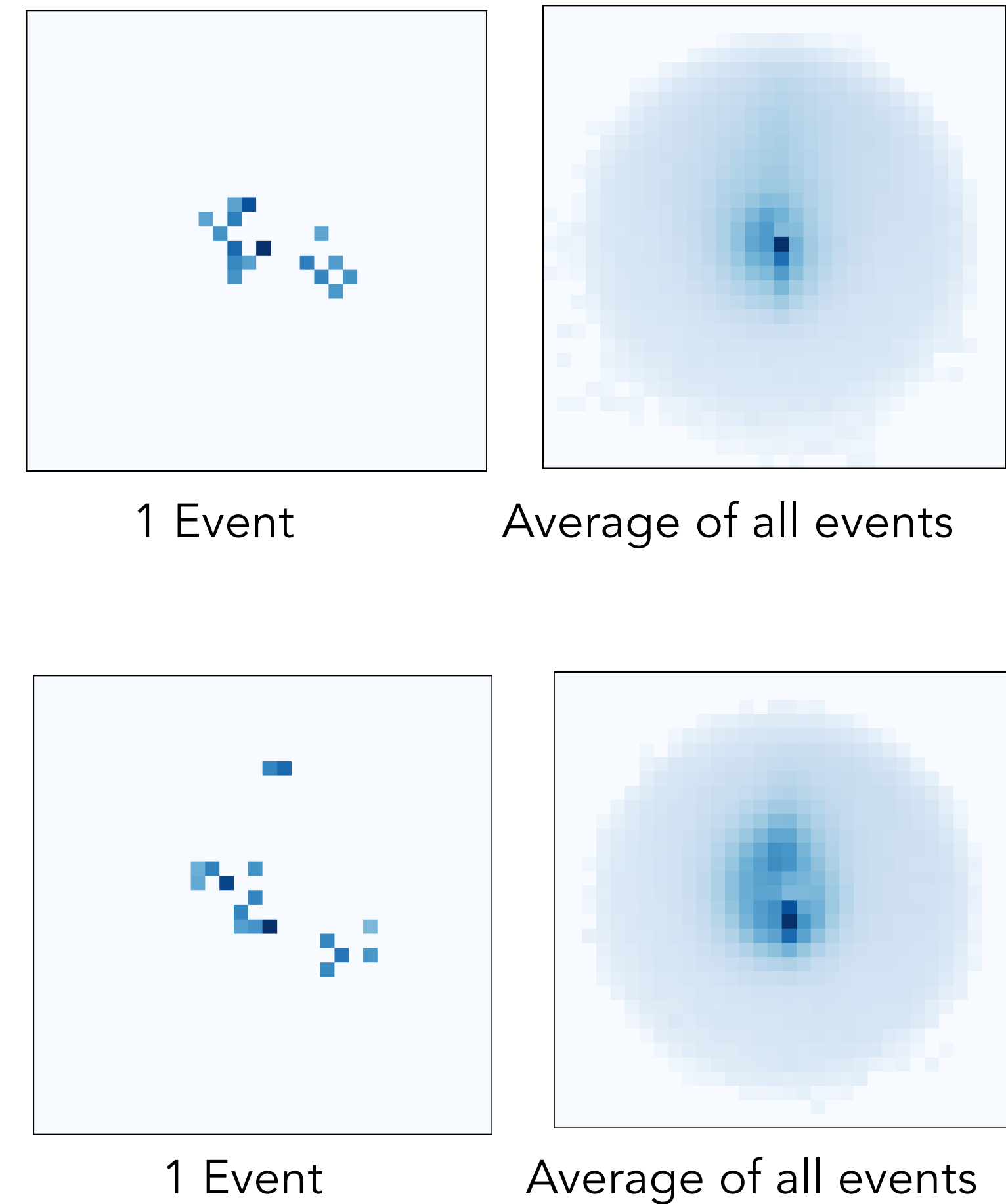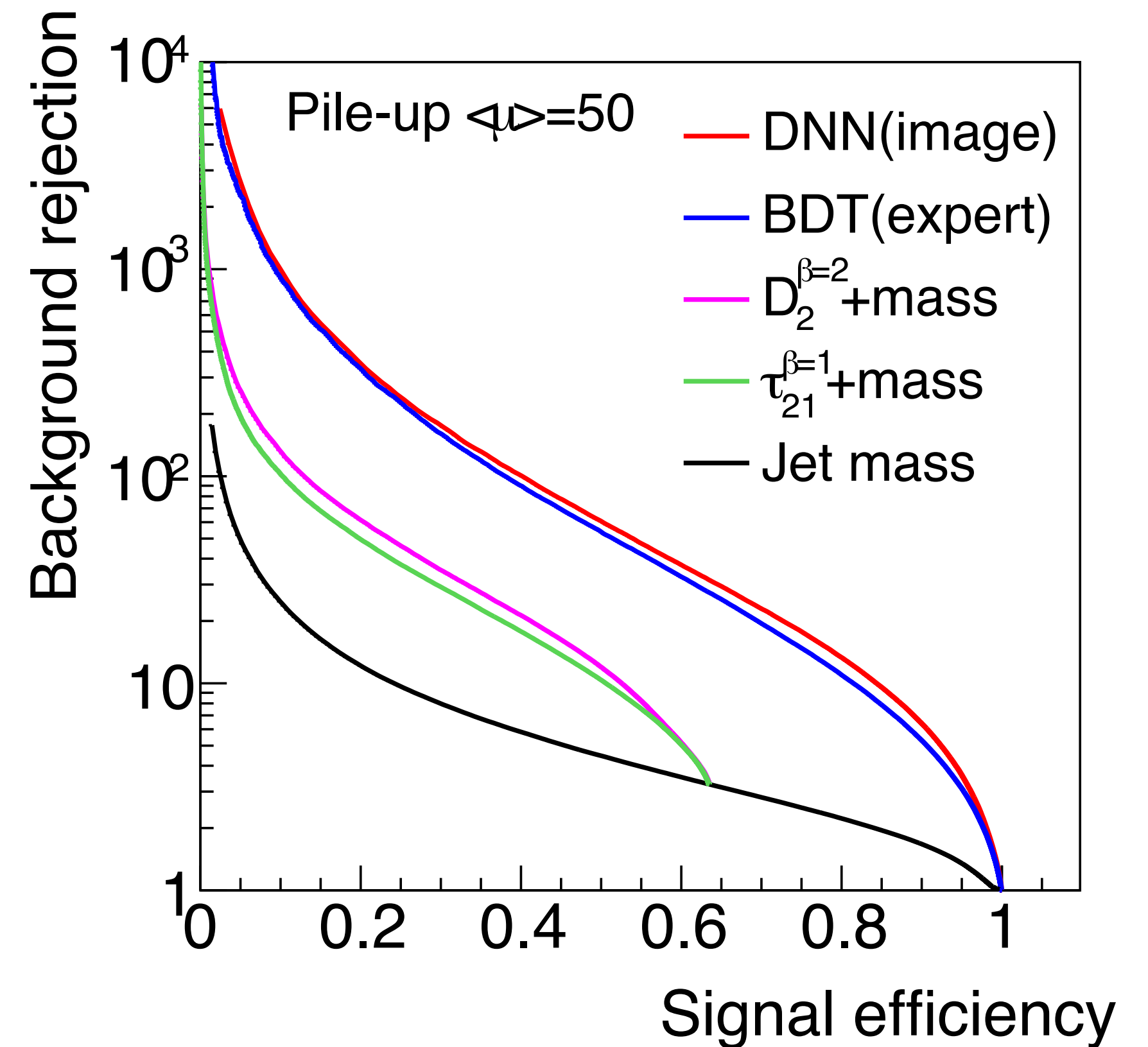Average of all events

- ‣ CNN performs better than Jet Substructure

  - ‣ **LL network (red line)**: AUC = 95.30% ± 0.02%

  - ‣ **HL network (blue line)**: AUC = 95.00% ± 0.02%

- ‣ But wait!

  - ‣ Where is that extra information coming from?

  - ‣ Why don't Jet Substructure observables contain this information?

- ‣ We've used a black box, so how can we investigate?



Graph legend: Pile-up $\langle\mu\rangle$=50, DNN(image), BDT(expert), $D_2^{\beta=2}$+mass, $\tau_{21}^{\beta=1}$+mass, Jet mass. Axes: Background rejection vs Signal efficiency.

*Baldi, P., Bauer, K., Eng, C., Sadowski, P., & Whiteson, D. (2016, March 30). Jet Substructure Classification in High-Energy Physics with Deep Neural Networks. arXiv.org. http://doi.org/10.1103/PhysRevD.93.094034*
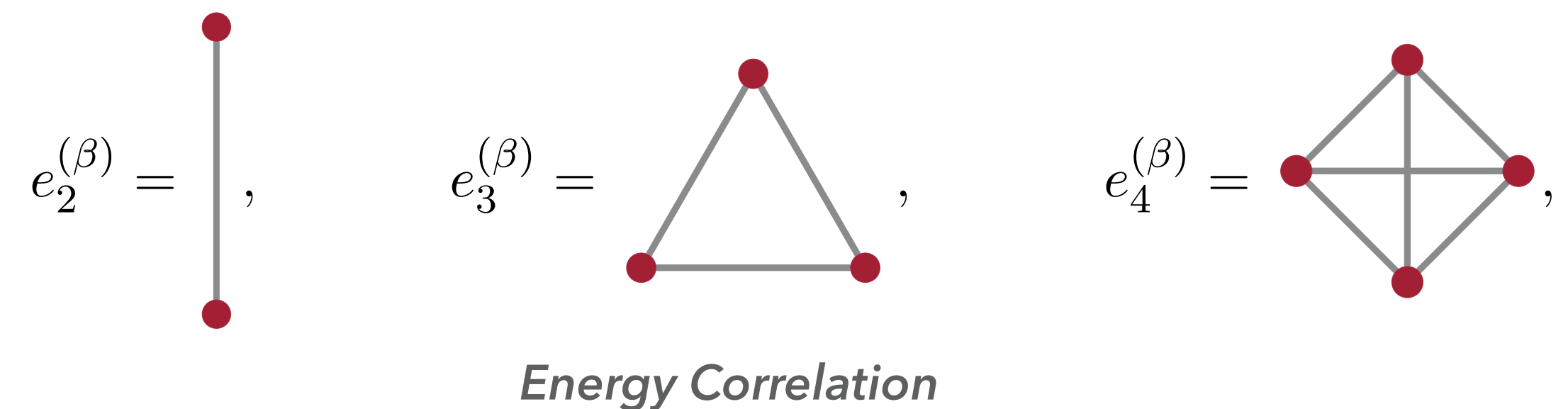
▸ We weren't the only ones thinking about a generalized approach to understanding JSS!
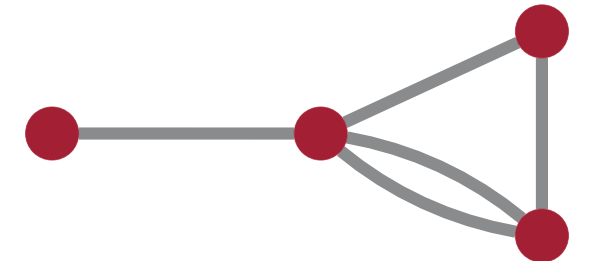
***Jesse Thaler, Patrick Komiske, Eric Metodiev***

*https://arxiv.org/abs/1712.07124*

▸ Energy Flow Polynomials (EFP): A complete linear basis set for jet substructure.

## Existing JSS exists in the EFP space

$$e_2^{(\beta)} = \Big| \quad , \qquad e_3^{(\beta)} = \triangle \quad , \qquad e_4^{(\beta)} = \diamond \quad ,$$

*Energy Correlation*

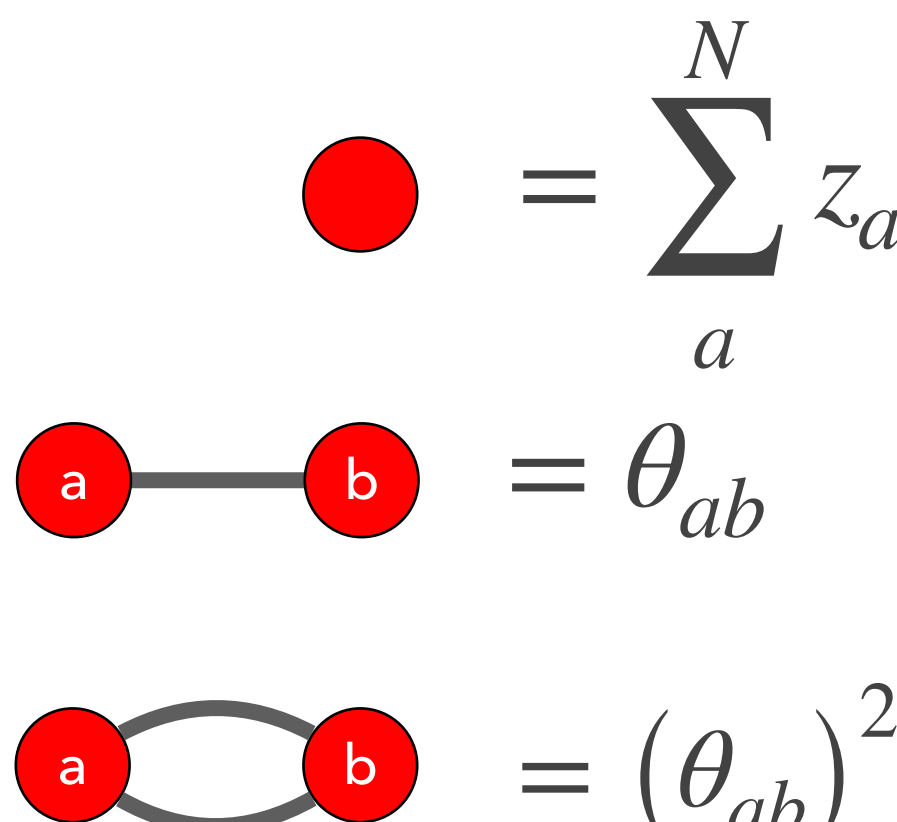## We can also explore more exotic observables.

$$= \sum_{i_1=1}^{M} \sum_{i_2=1}^{M} \sum_{i_3=1}^{M} \sum_{i_4=1}^{M} z_{i_1} z_{i_2} z_{i_3} z_{i_4} \theta_{i_1 i_2} \theta_{i_2 i_3} \theta_{i_2 i_4}^2 \theta_{i_3 i_4}.$$

## Graph Components

$$\bullet = \sum_a^N z_a$$

$$a \!-\! b = \theta_{ab}$$

$$a = b = (\theta_{ab})^2$$

## With variables and parameters (κ,β)

$$z_i^\kappa = \left( \frac{p_{T,i}}{\sum_i p_{T_j}} \right)^\kappa \qquad \theta_{ij}^\beta = \left( \Delta\eta_{ij}^2 + \Delta\phi_{ij}^2 \right)^{\beta/2}$$

1. Black-box Learning

   ‣ **Benefits:** Powerful Performance and no need to hand-pick observables

   ‣ **Drawbacks:** Not interpretable. What are we learning about the problem?

2. Energy Flow Polynomials (EFP)

   ‣ **Benefits:**

      ‣ Physics motivated

      ‣ Modeling can be verified

      ‣ Uncertainties can be defined

      ‣ Compact and efficient

   ‣ **Drawbacks:** It's an infinite space. How do we begin to choose observables and their parameters?

**Combine them! The DNN has learned how to solve the problem. Let the DNN tell us what EFPs to choose!**

▸ Consider boosted W vs QCD jet binary classification

▸ CNN has learned where to draw an ideal decision surface in its feature space.

▸ We want a HL feature space that makes equally good decisions.

▸ How do we compare the decision surface for the CNN to the HL features?

## Comparing pair orderings

- Take a pair of signal (x) and background (x') features,

- Predictions for 2 NN (f(x) and g(x)) of the features will increase/decrease relative to input
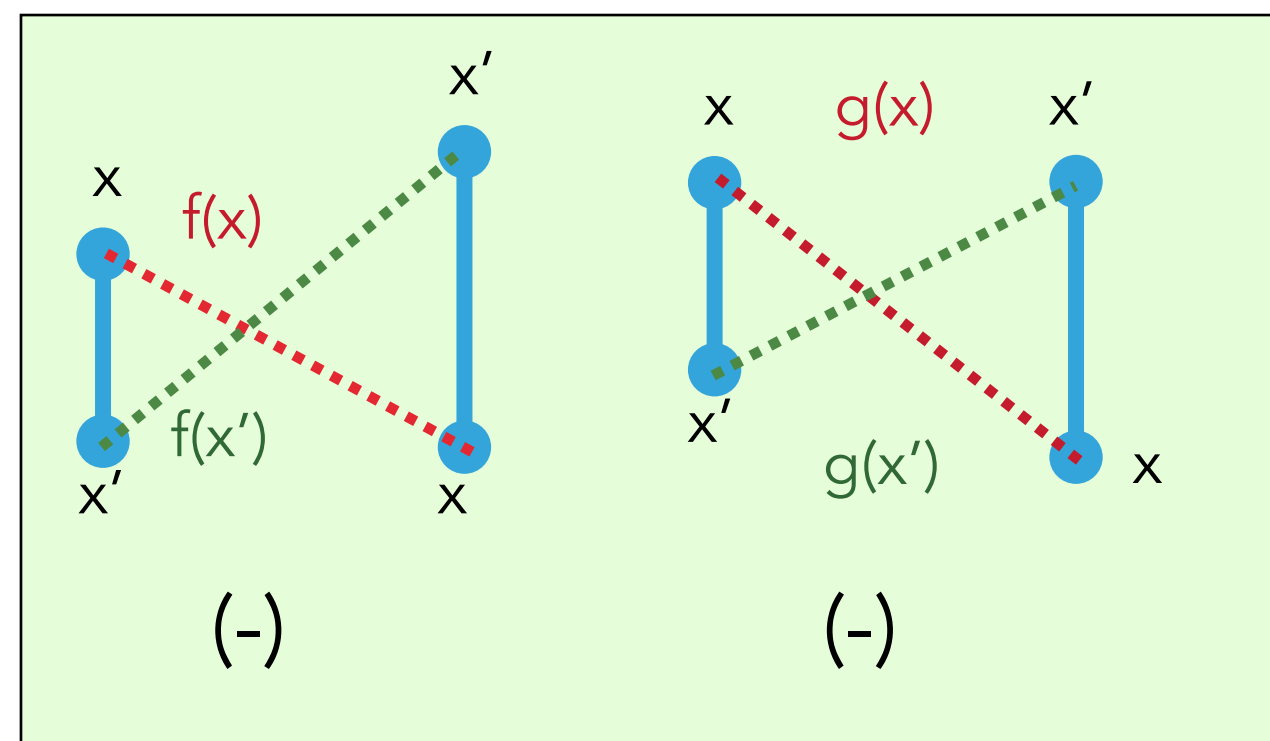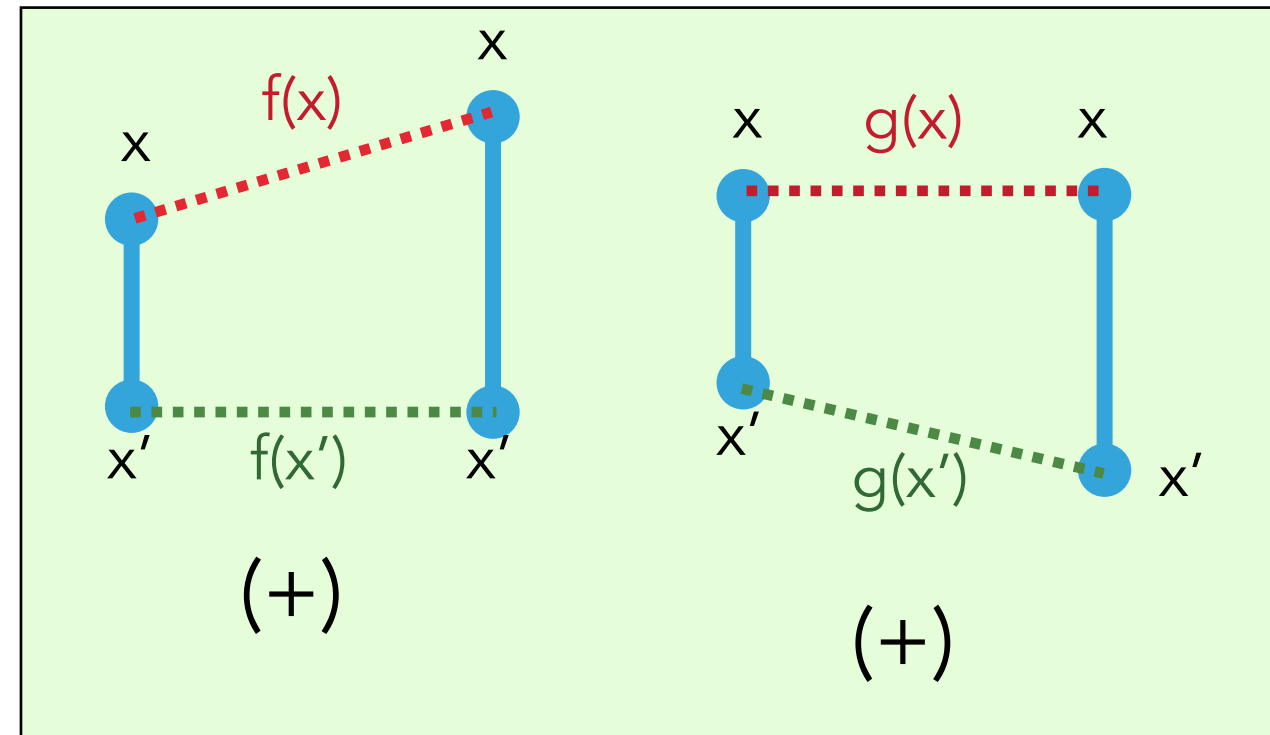
$$\mathrm{DO}(x, x') = \Theta\left[(f(x) - f(x')) \cdot (g(x) - g(x'))\right]$$

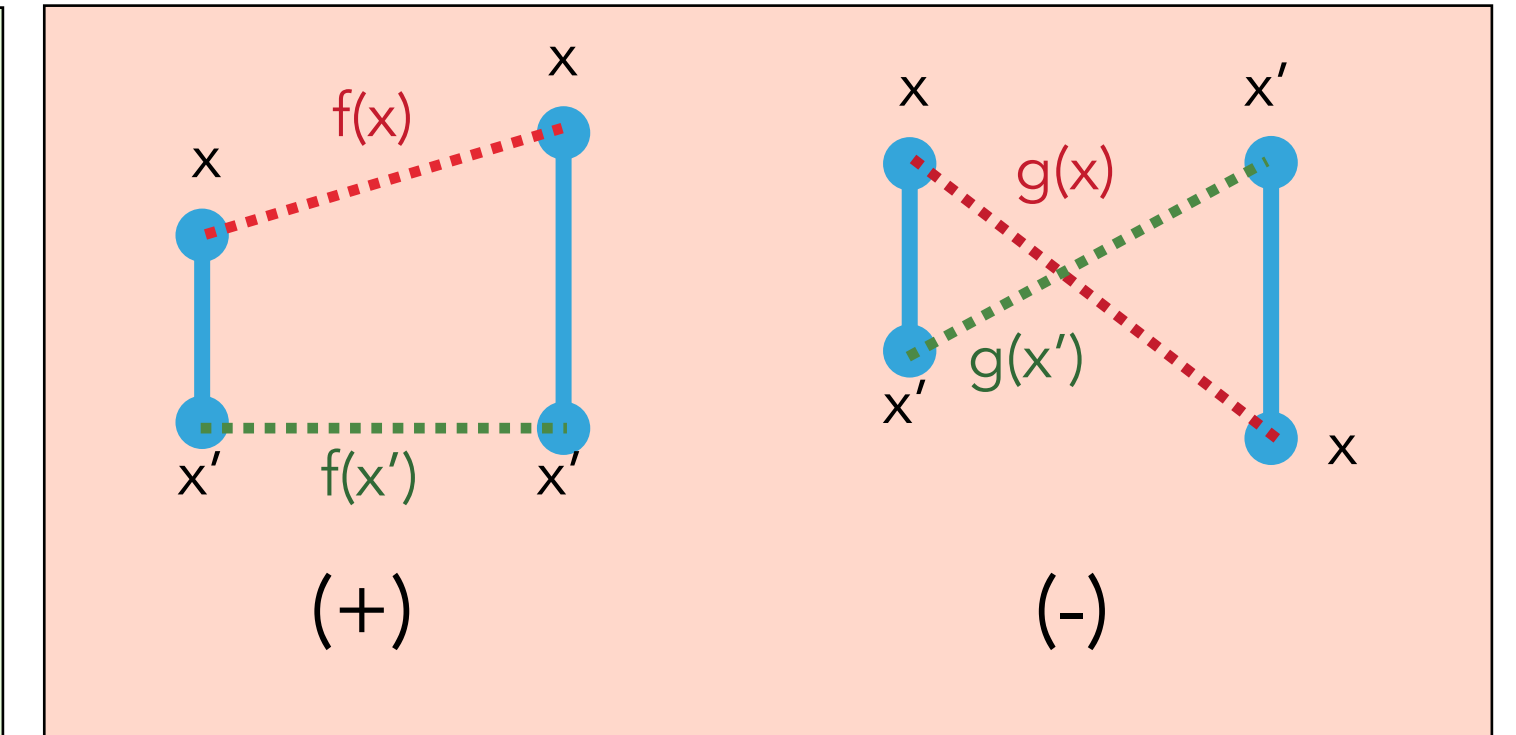Heaviside step function (Θ) sets DO=1 for similar order, DO=0 for dissimilar order.

Average many examples = Average Decision Ordering

$$\mathrm{ADO}' = \sum DO(x, x')$$

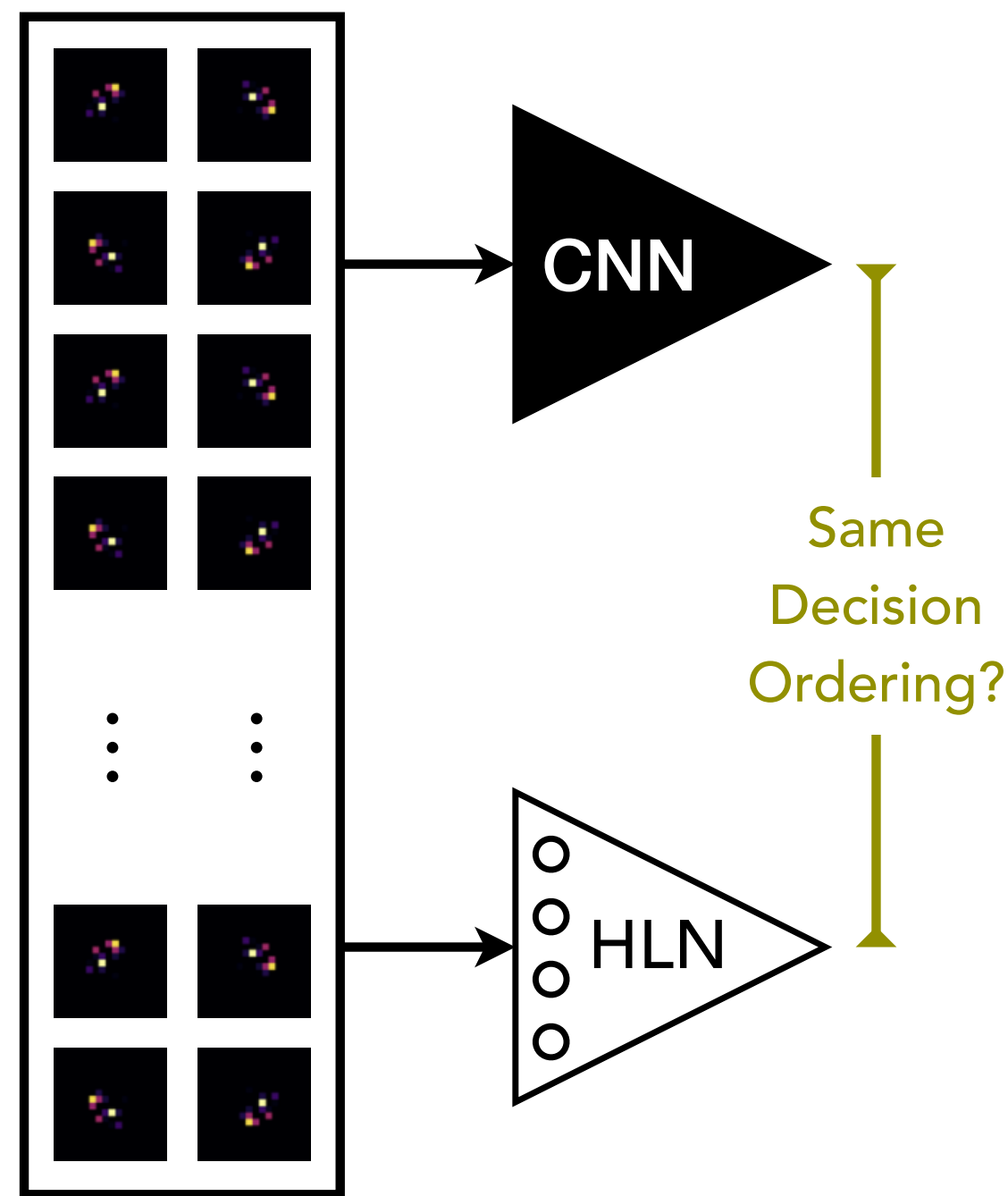## Similar Orderings



## Dissimilar Orderings



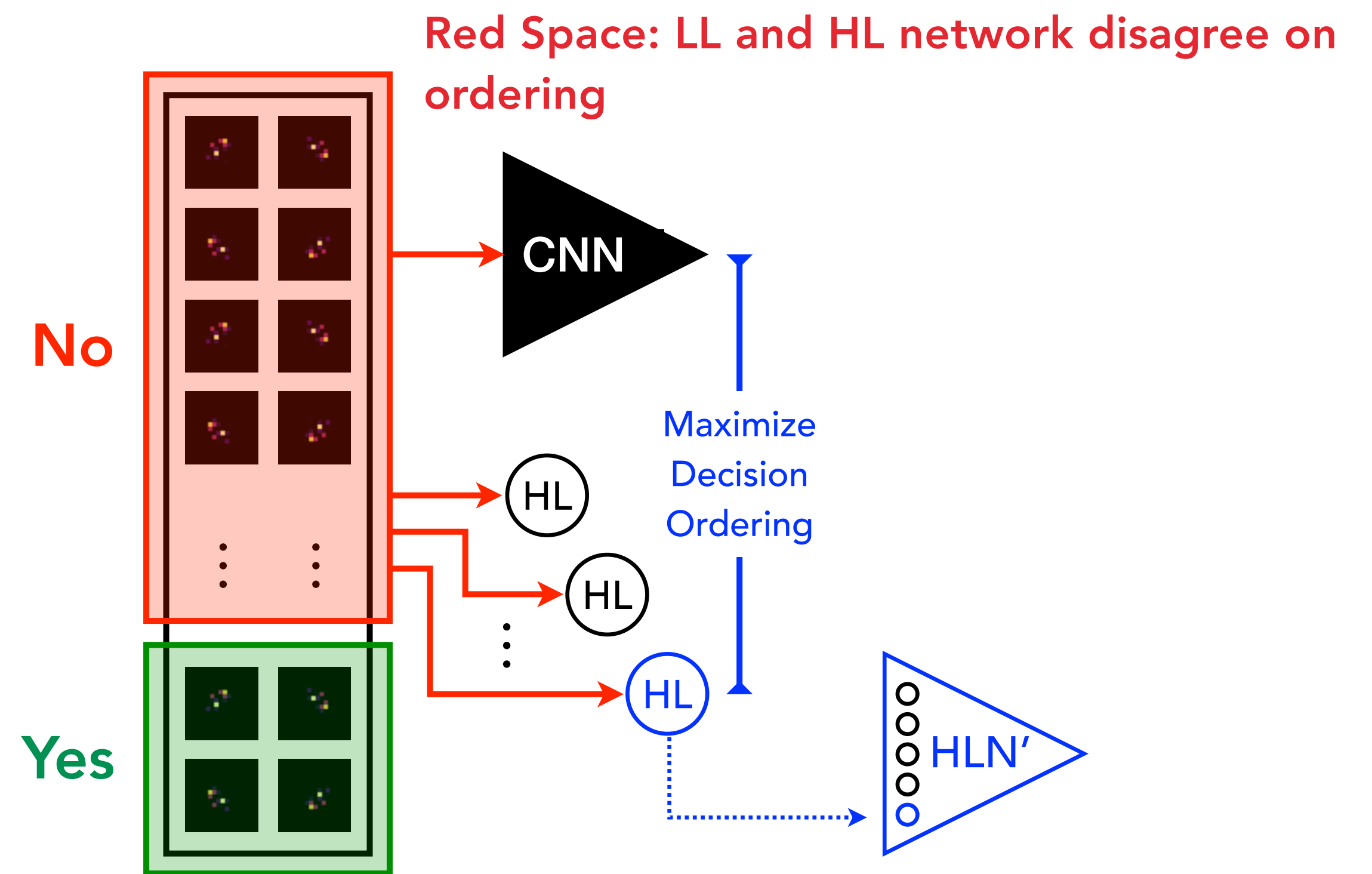**ADO = 1 :** *Identical decisions*

**ADO = 0.5 :** *Random similarity*

*Why is ADO useful?*

We can compare NN decision making. Where does the HL network and LL network disagree?

Signal/Background Pairs

**Red Space: LL and HL network disagree on ordering**

CNN

CNN

Same Decision Ordering?

If NOT–>

Same Decision Ordering?

No

HLN

Maximize Decision Ordering

HL

HL

Yes

HL

HLN'

Use ADO to choose EFP that makes similar choices to the LL network in the "differently ordered" red space
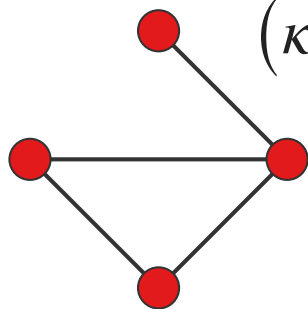
We only need 1 new observable to achieve equal performance with the CNN!

| Observable | AUC | ADO[CNN, Obs.] | |
|---|---|---|---|
| $M_{\text{jet}}$ | $0.898 \pm 0.004$ | $0.807$ | |
| $C_2^{\beta=1}$ | $0.660 \pm 0.006$ | $0.584$ | |
| $C_2^{\beta=2}$ | $0.604 \pm 0.007$ | $0.548$ | |
| $D_2^{\beta=1}$ | $0.790 \pm 0.005$ | $0.743$ | |
| $D_2^{\beta=2}$ | $0.807 \pm 0.005$ | $0.762$ | |
| $\tau_2^{\beta=1}$ | $0.662 \pm 0.006$ | $0.600$ | |
| 6HL | $0.9504 \pm 0.0002$ | $0.971$ | Original HL |
| CNN | $0.9531 \pm 0.0002$ | $1.000$ | |
| 7HL$_{\text{black-box}}$ | $0.9528 \pm 0.0003$ | $0.971$ | |

Original HL + 1 EFP

## Which EFP did we pick?



$$\left(\kappa=2,\ \beta=1/2\right) = \sum_{a,b,c,d=1}^{N} z_a^2 z_b^2 z_c^2 z_d^2 \sqrt{\theta_{ab}\theta_{bc}\theta_{ac}\theta_{ad}}$$

Noteworthy details

‣ EFP is not Infrared-safe (k ≠ 1)

‣ β=1/2 is probing small-angle behaviour

‣ Chromatic #3 graph (probing deviations from 2-prong substructure)

‣ Chromatic Number = Minimum number of prongs to not vanish

Deep networks can identify gaps where low-level data contains unused information

ML Mapping strategies can capture and translate that information into understandable physics

# Questions ?