

When (not) to use uncertainties for ML in particle physics

Bayesian Deep Learning for Cosmology and Time Domain Astrophysics

Anja Butter, ITP Heidelberg



Setting the stage for LHC physics

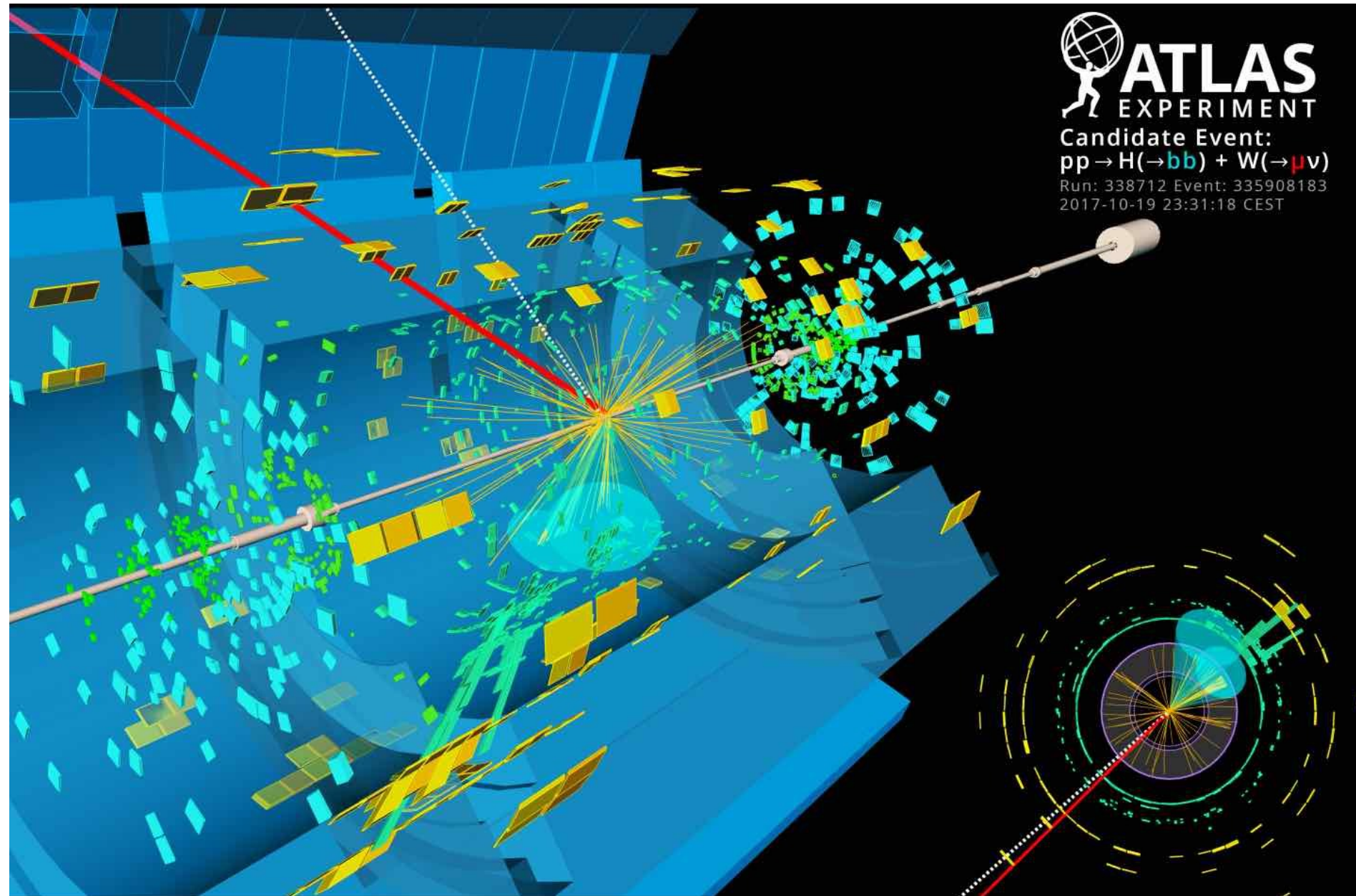


Before the particle accelerator

Setting

- Large Hadron Collider at CERN
- Independent proton collisions at 13 TeV
- Recorded by ATLAS, CMS, LHCb, ALICE
- **Huge** dataset $\sim 1\text{Pb/s}$ before trigger selection

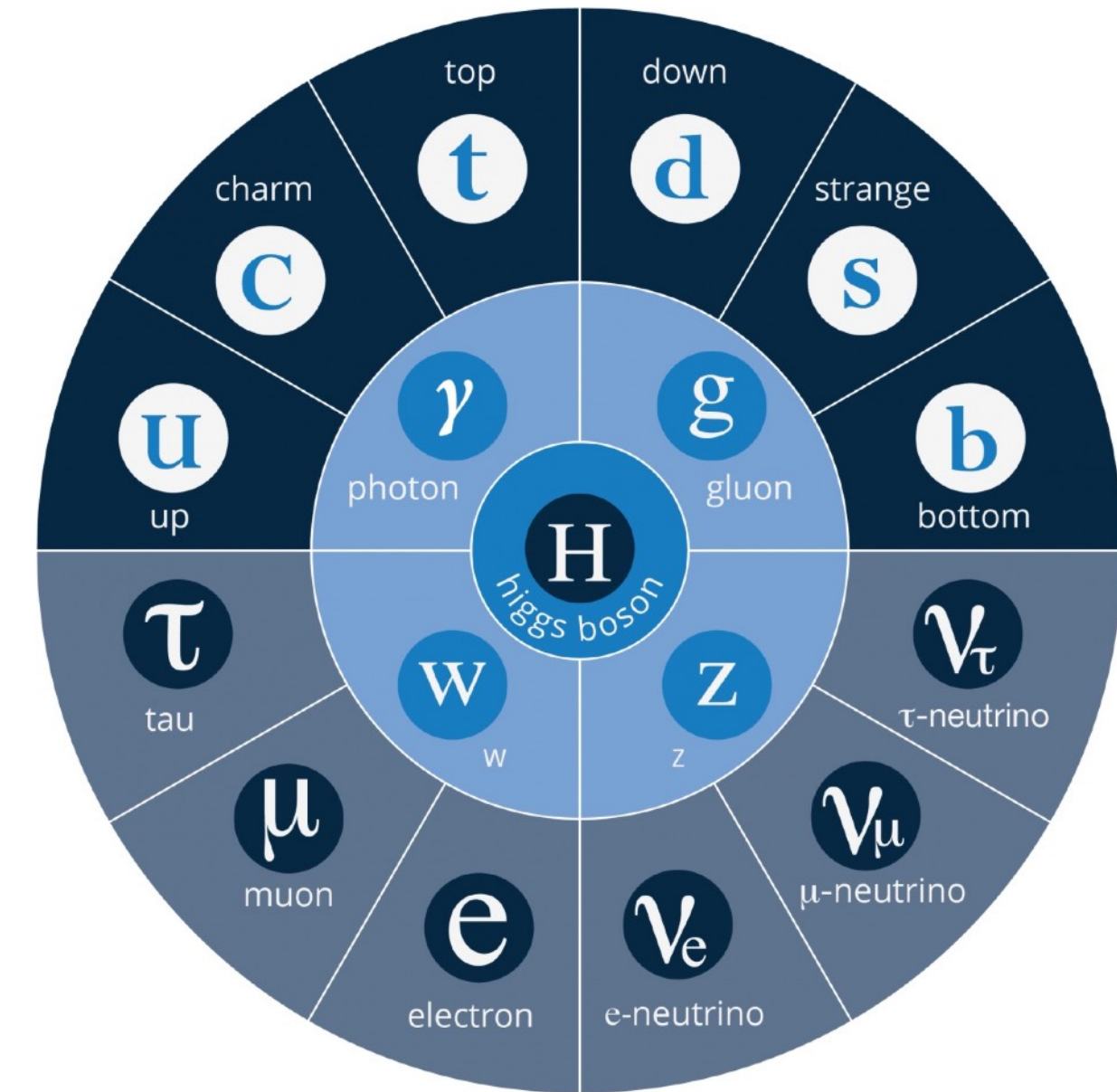
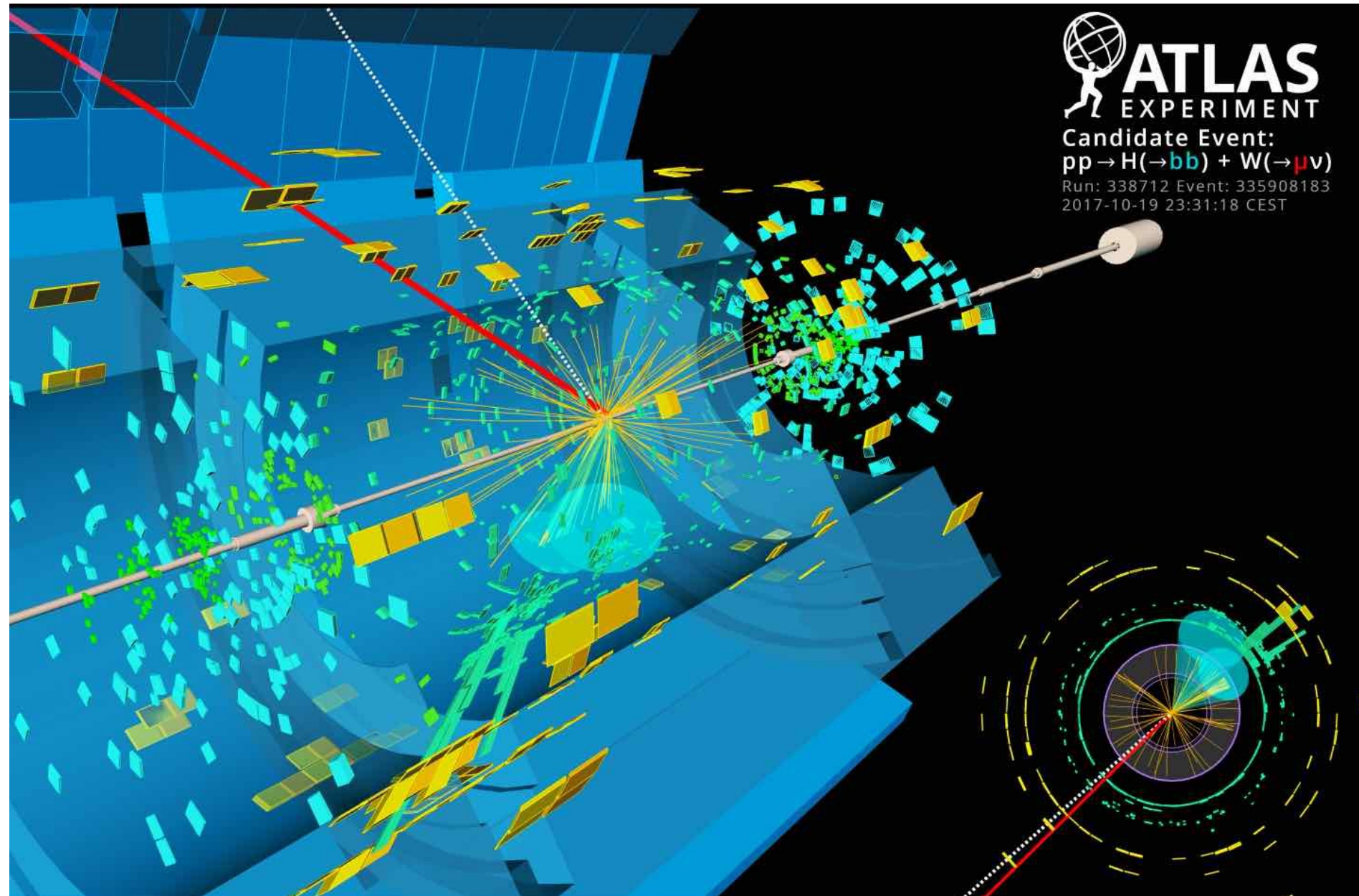
Setting the stage for LHC physics



Setting

- Large Hadron Collider at CERN
- Independent proton collisions at 13 TeV
- Recorded by ATLAS, CMS, LHCb, ALICE
- **Huge** dataset $\sim 1\text{Pb/s}$ before trigger selection

Setting the stage for LHC physics



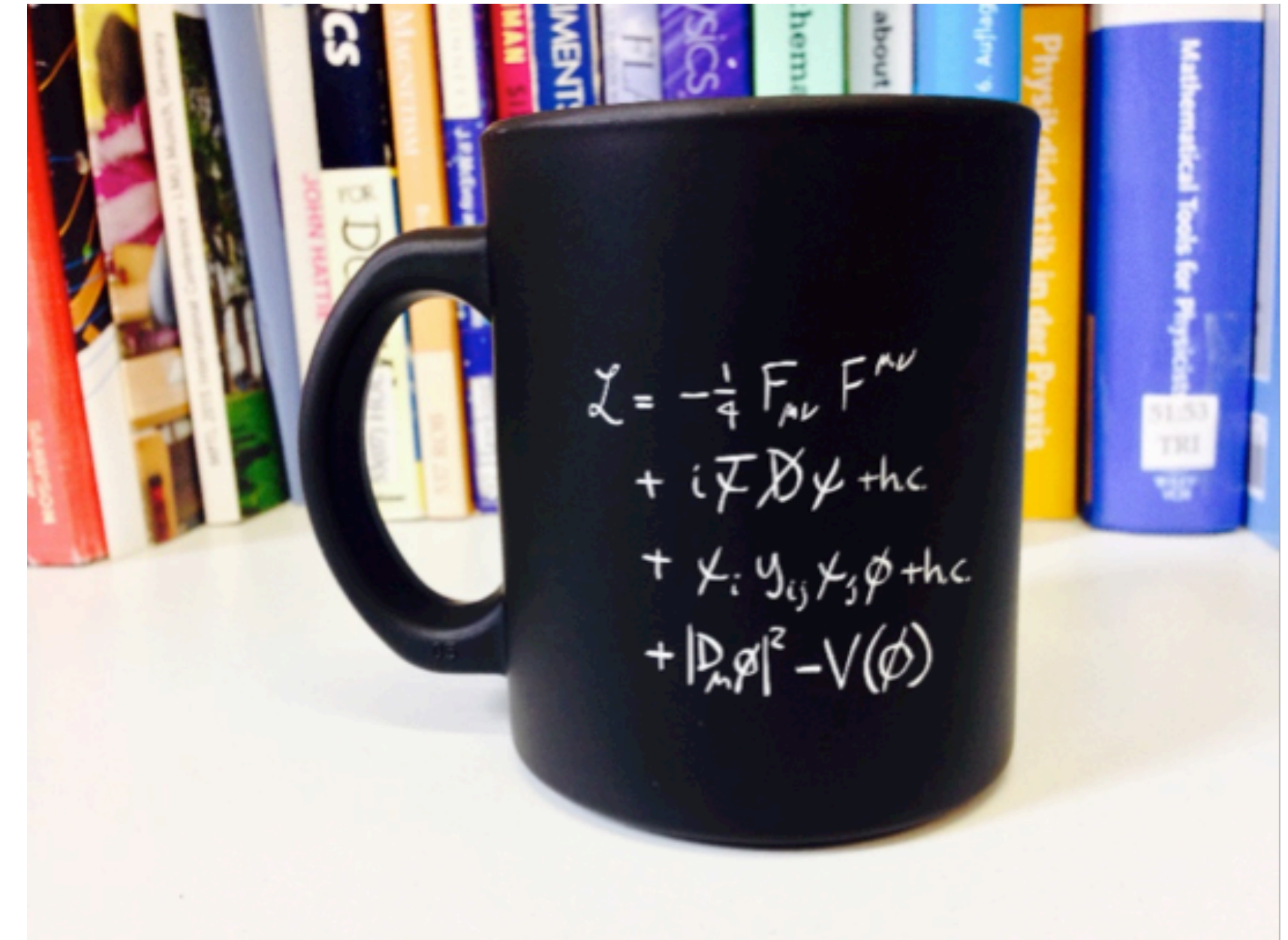
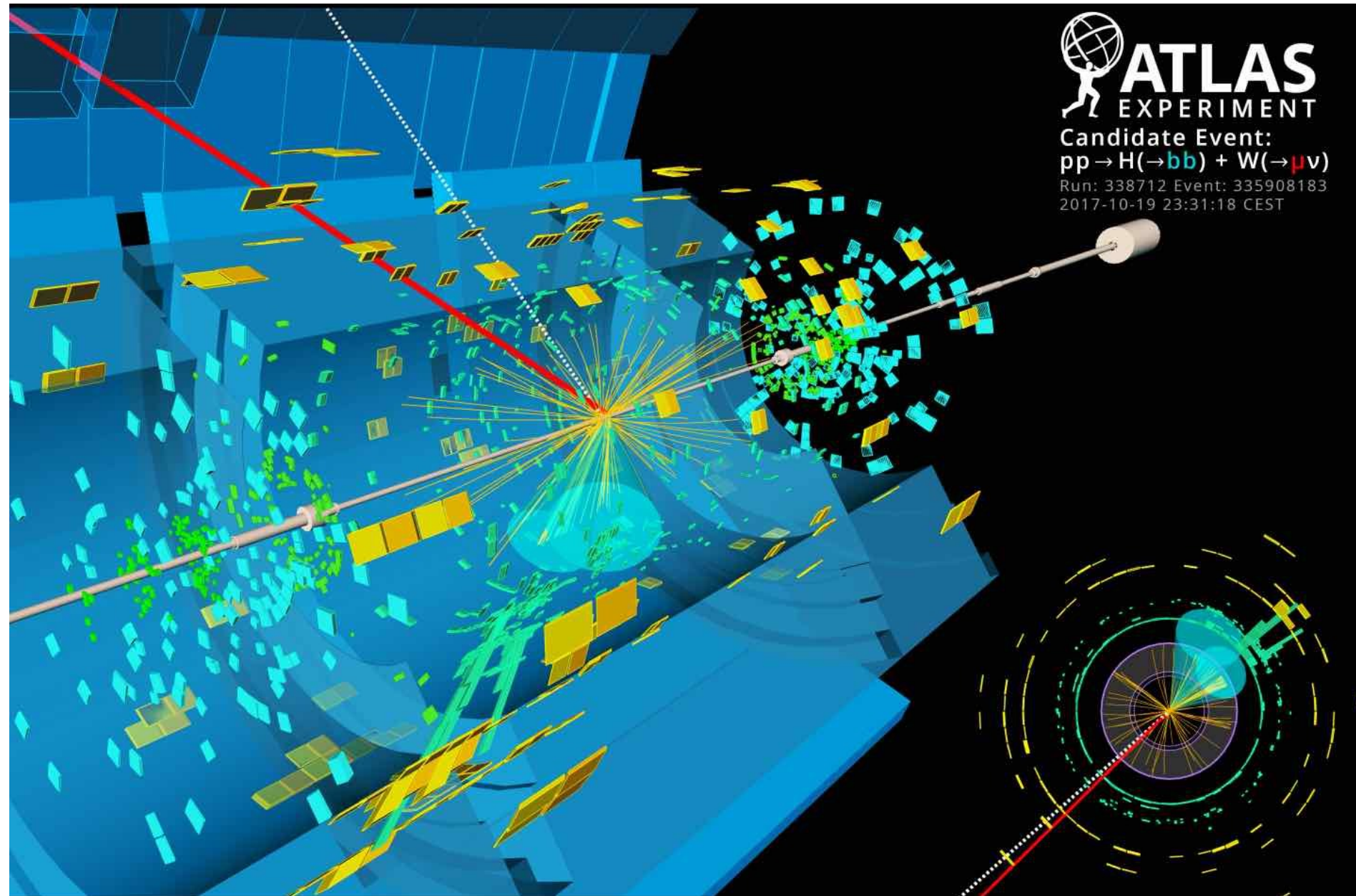
Setting

- Large Hadron Collider at CERN
- Independent proton collisions at 13 TeV
- Recorded by ATLAS, CMS, LHCb, ALICE
- **Huge** dataset $\sim 1\text{Pb/s}$ before trigger selection

Goal

- Understand full dataset from **1st principles**
- Precision measurements of the SM
- Find signs of new physics (eg dark matter)

Setting the stage for LHC physics



Setting

- Large Hadron Collider at CERN
- Independent proton collisions at 13 TeV
- Recorded by ATLAS, CMS, LHCb, ALICE
- **Huge** dataset $\sim 1\text{Pb/s}$ before trigger selection

Goal

- Understand full dataset from **1st principles**
- Precision measurements of the SM
- Find signs of new physics (eg dark matter)

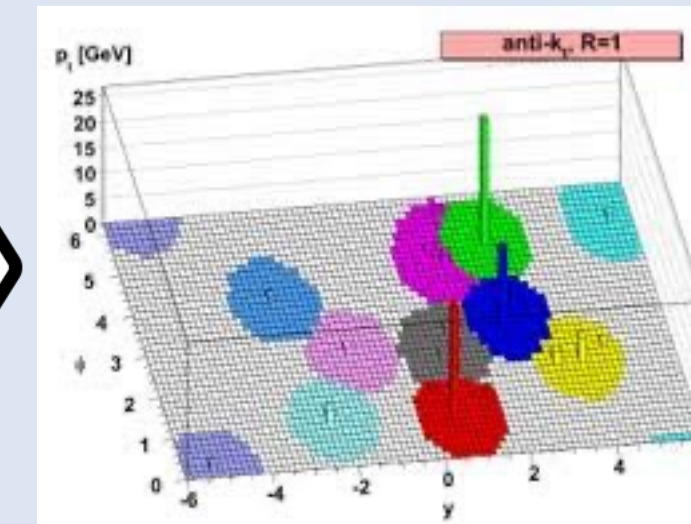
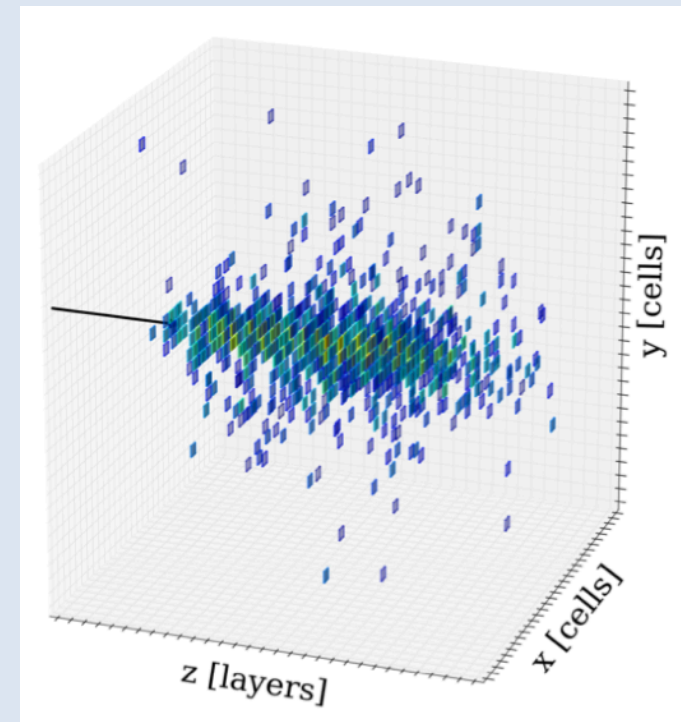
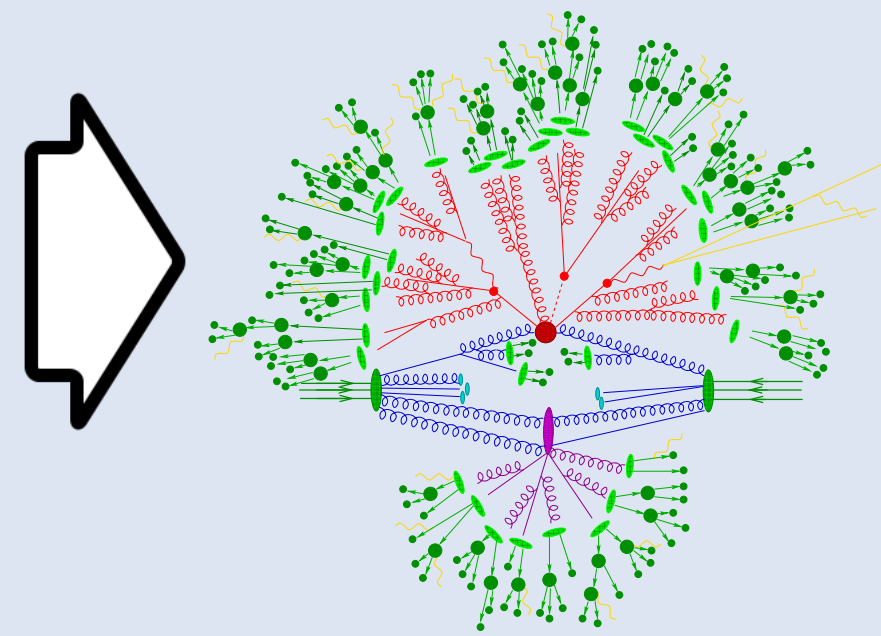
LHC analyses in a nutshell

Theory/Hypothesis

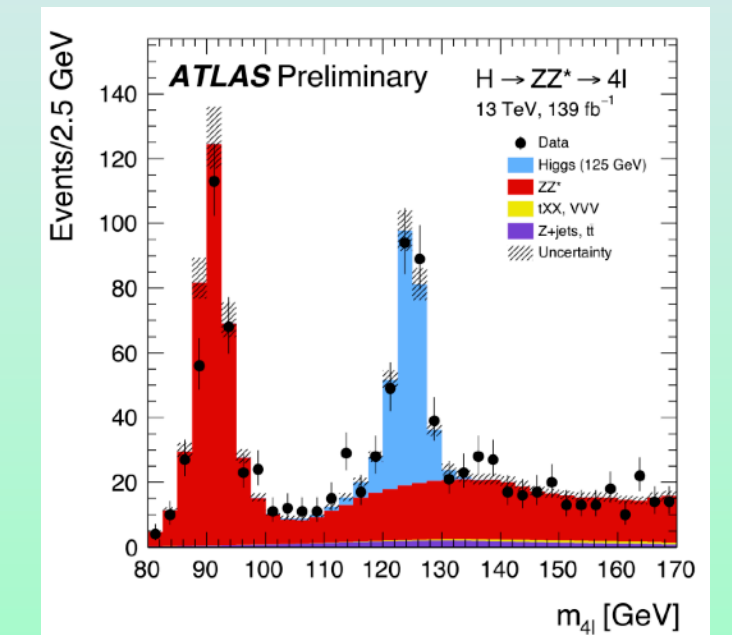
Event generation

Detector simulation

Reconstruction

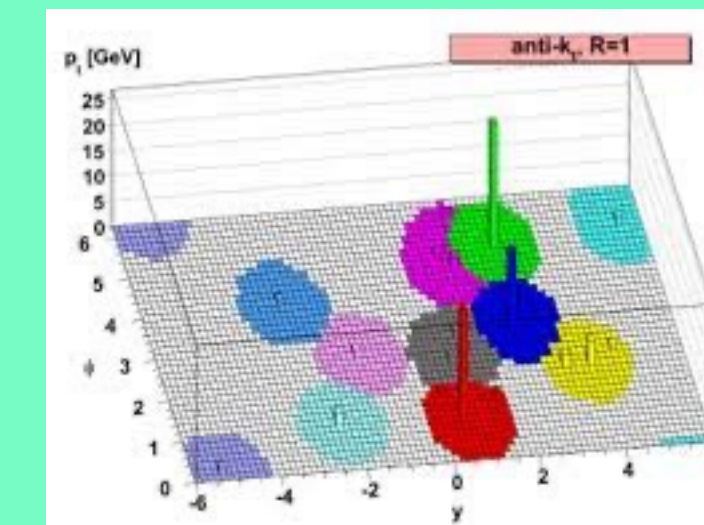
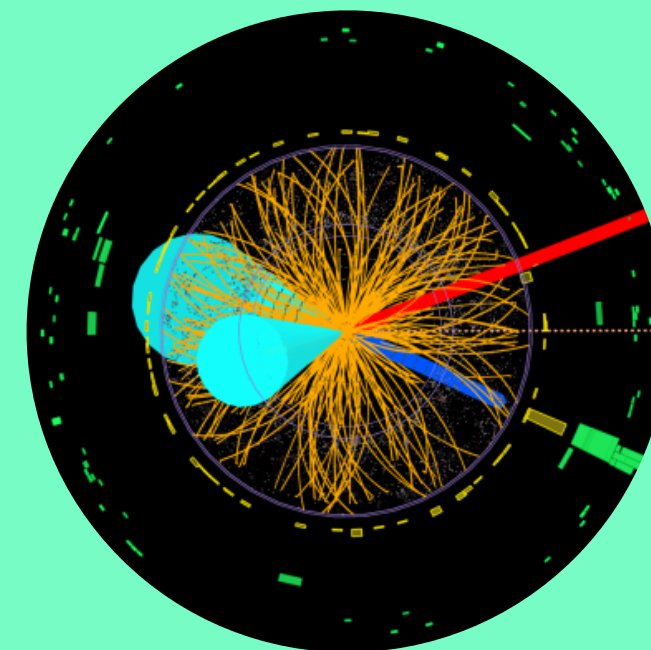


Statistical analysis



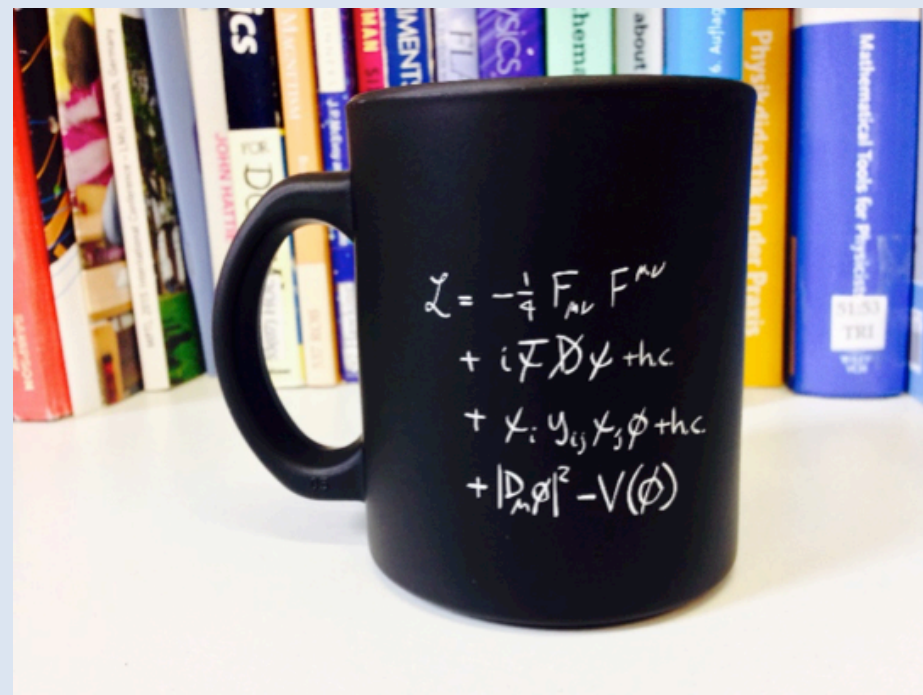
Data

Reconstruction

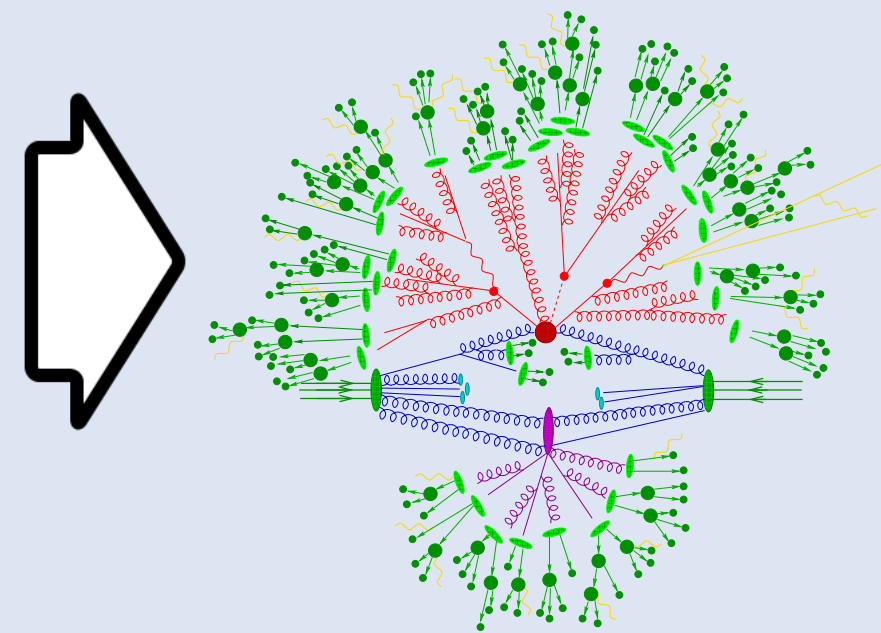


LHC analyses in a nutshell

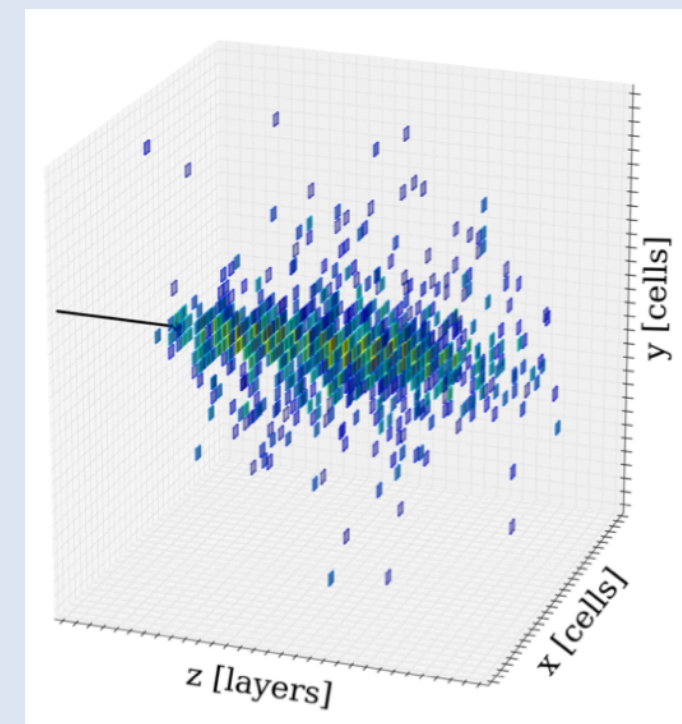
Theory/Hypothesis



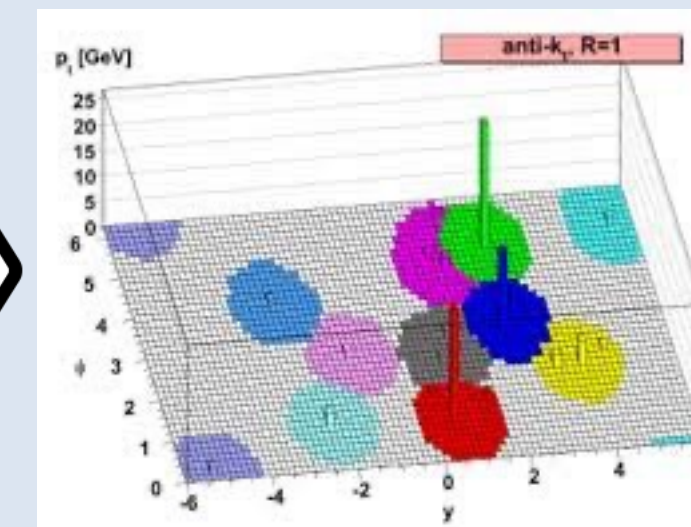
Event generation



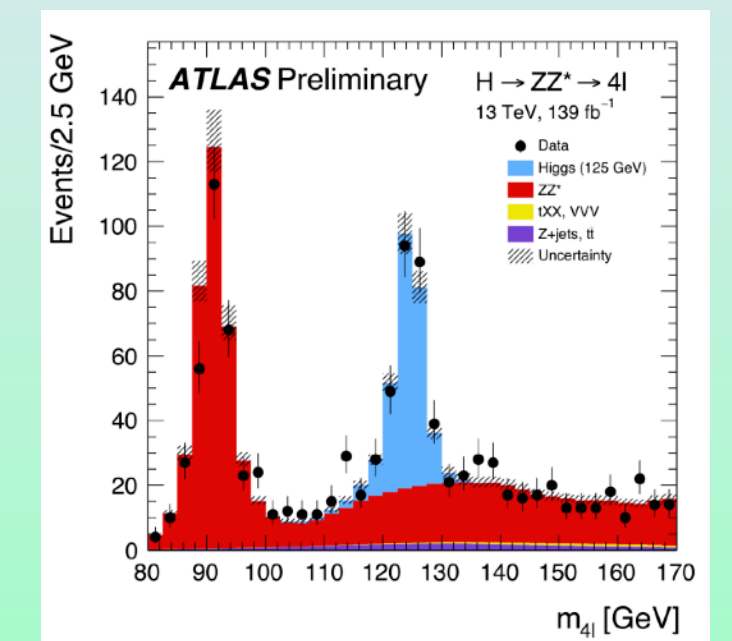
Detector simulation



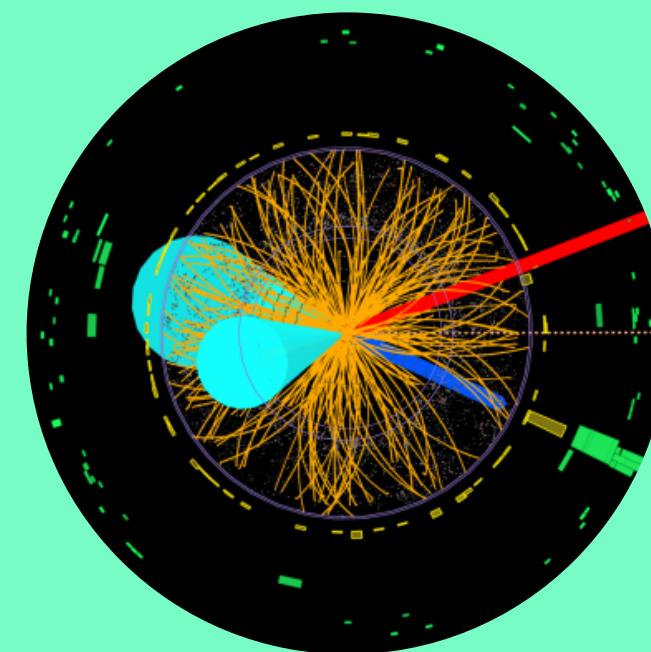
Reconstruction



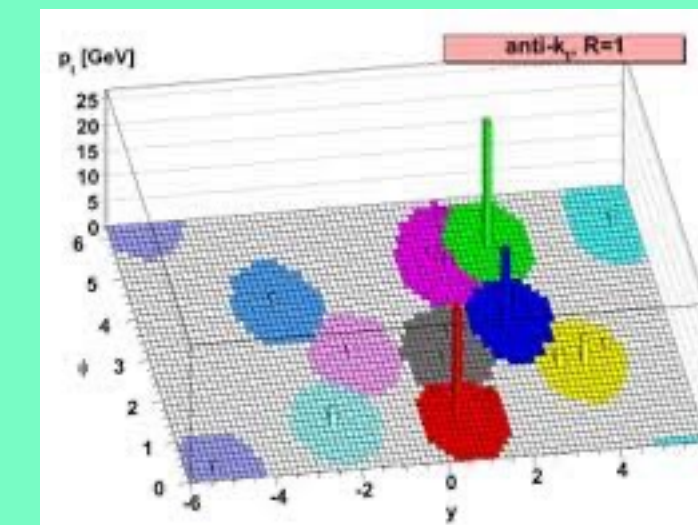
Statistical analysis



Data



Reconstruction



Limit setting for standard analyses

We want to set a limit on new physics parameter θ

$$p(\text{data} | \theta, \nu) = \prod_{i=1}^{n_{bins}} P(n_i | s_i(\theta, \nu) \cdot \epsilon_i(\nu) + b_i(\nu)) P(\nu | \text{aux. data})$$

Likelihood $\rightarrow L(\theta) = p(\text{data} | \theta)$

Best estimate $\rightarrow \hat{\theta} = \arg \max_{\theta} L(\theta)$

p-value $\rightarrow p_{\theta} = \int_{n_{obs}}^{\infty} dn p(n | \theta)$

68% confidence limit $\rightarrow [\theta_{min}, \theta_{max}]$ with $p_{\theta_{min}} < 0.16$

LHC discovery = 5 sigma $\rightarrow p < 3 \times 10^{-7}$

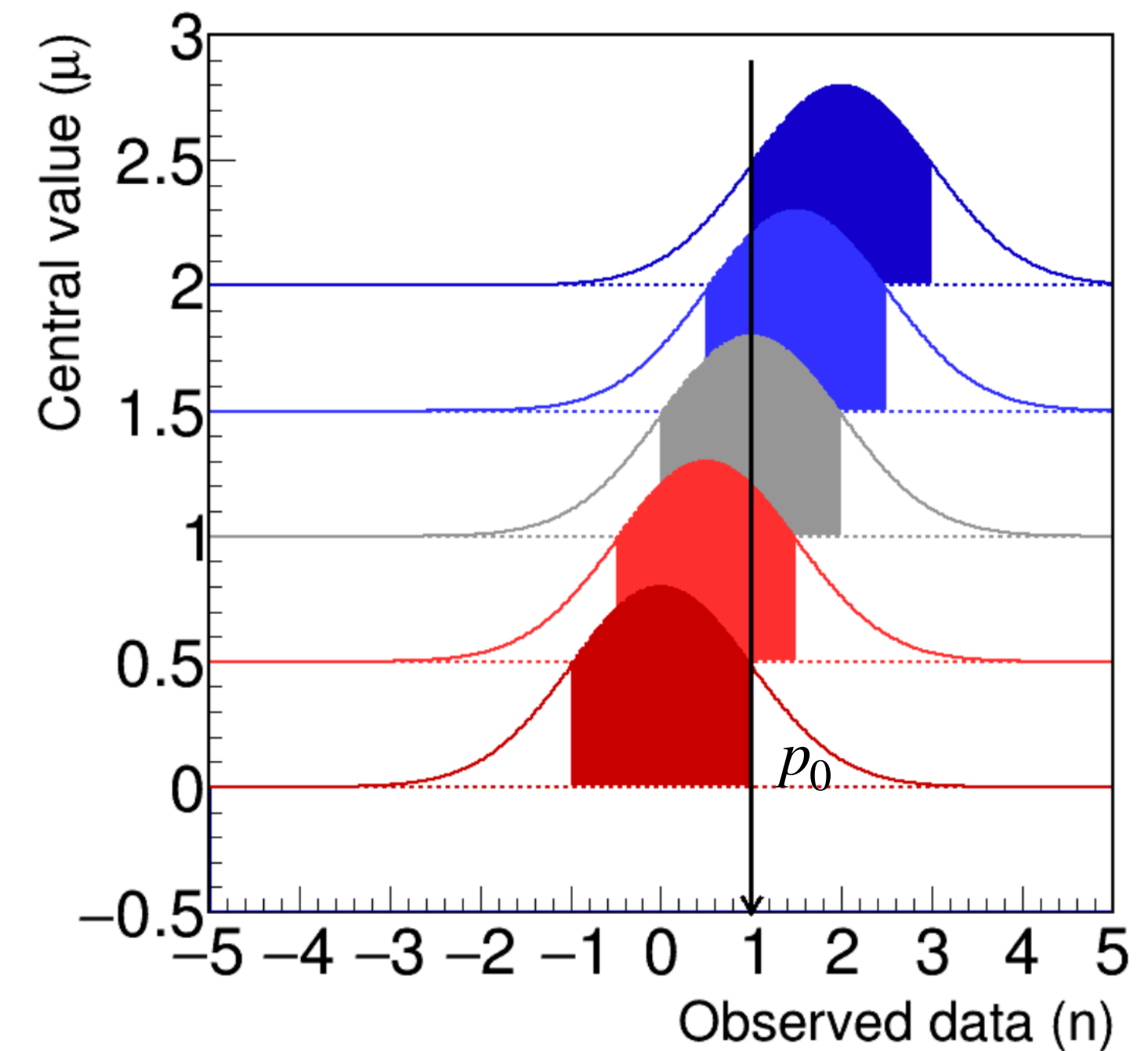


Illustration by Nicolas Berger

Not a statement on the probability of θ !!!

Types of uncertainties to build the likelihood

	Statistical uncertainties	Systematic uncertainties	Theory uncertainties
What for?	Counting $n_{event} = p_{event} \cdot N_{collisions}$	Anything that is calibrated Efficiencies, Luminosity,...	Scale dependence prediction Two point correlations (Eg. Sherpa vs MadGraph)
Shape	Poisson $\frac{\lambda^x}{x!} \exp(-\lambda)$	Gaussian $\frac{1}{\sigma\sqrt{2\pi}} \exp - (x - \mu)^2/2\sigma^2$	Flat/Gaussian $\theta(x - \mu_{min})\theta(\mu_{max} - x)$
Correlation	Independent between bins	Correlated Covariance matrix from toys	Correlated Eg. Nuisance parameters

How does ML enter in this?

A short history of ML in HEP

First HEP NN papers

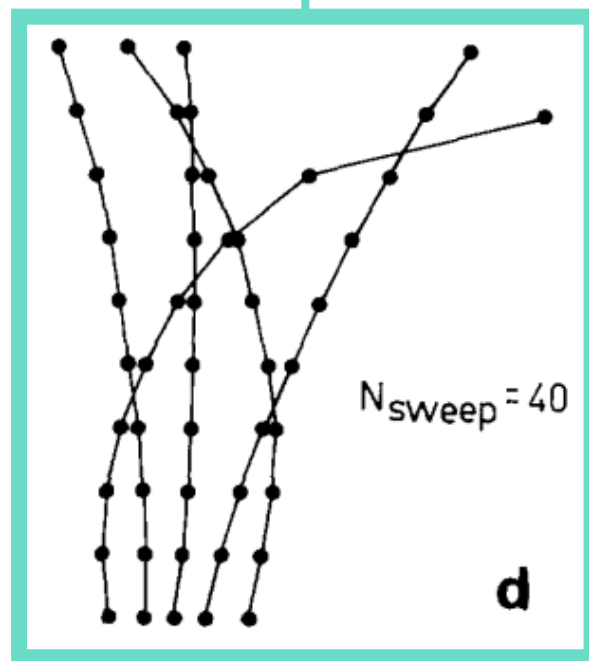
Track finding

Denby (LAL, Orsay) '87

Peterson (Lund) '88

→ Jet identification

1987/88

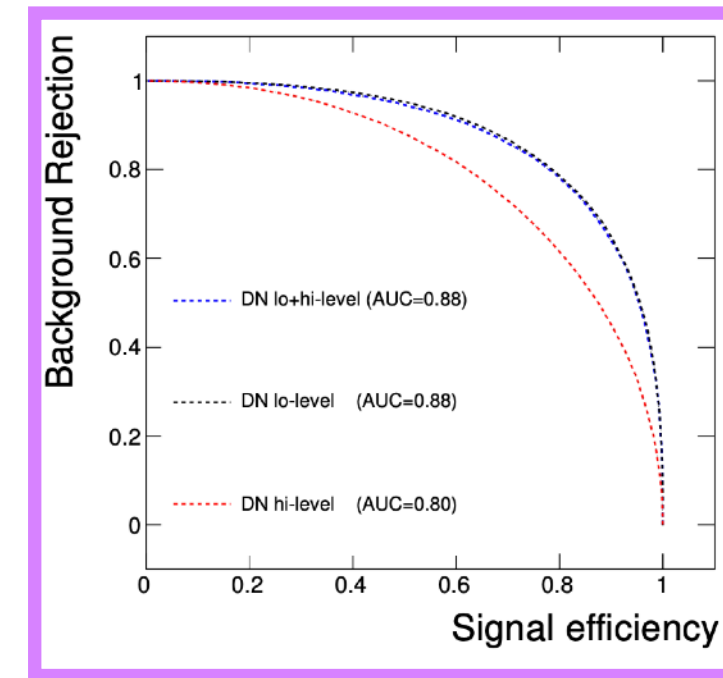
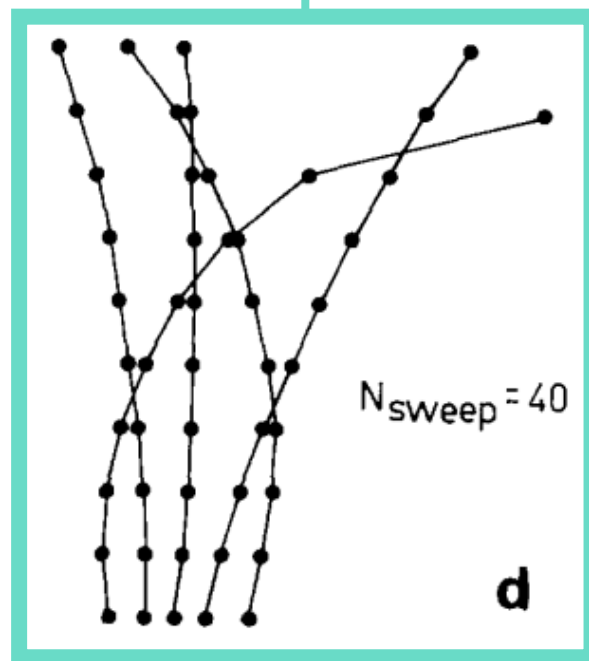


A short history of ML in HEP

First HEP NN papers

Track finding
Denby (LAL, Orsay) '87
Peterson (Lund) '88
→ Jet identification

1987/88



2014

Relaunch

Deep Learning in HEP
Signal vs Background

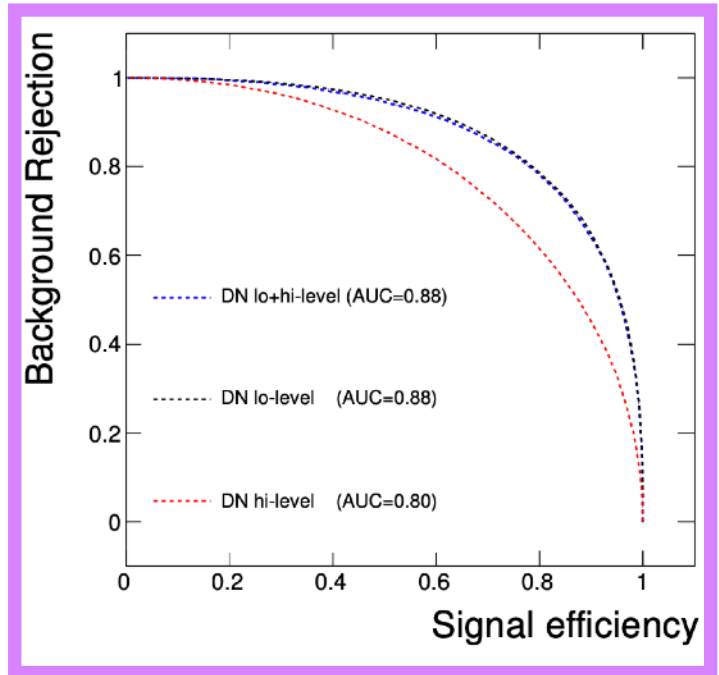
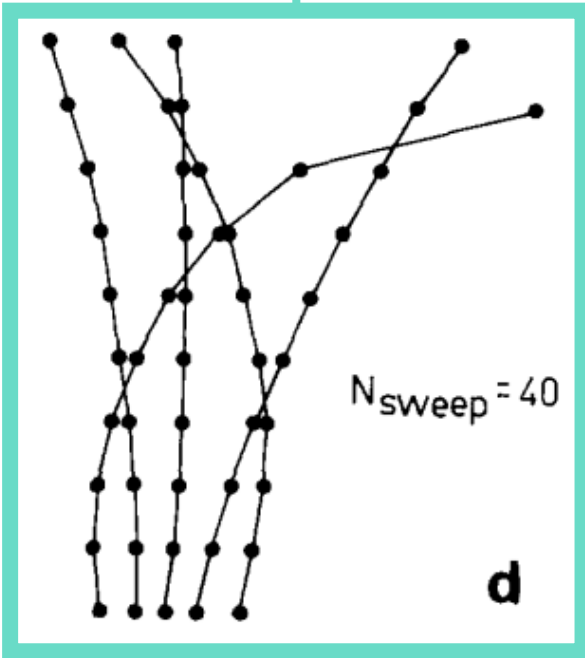
P. Baldi, P. Sadowski,
D. Whiteson

A short history of ML in HEP

First HEP NN papers

Track finding
Denby (LAL, Orsay) '87
Peterson (Lund) '88
→ Jet identification

1987/88



2014



Relaunch

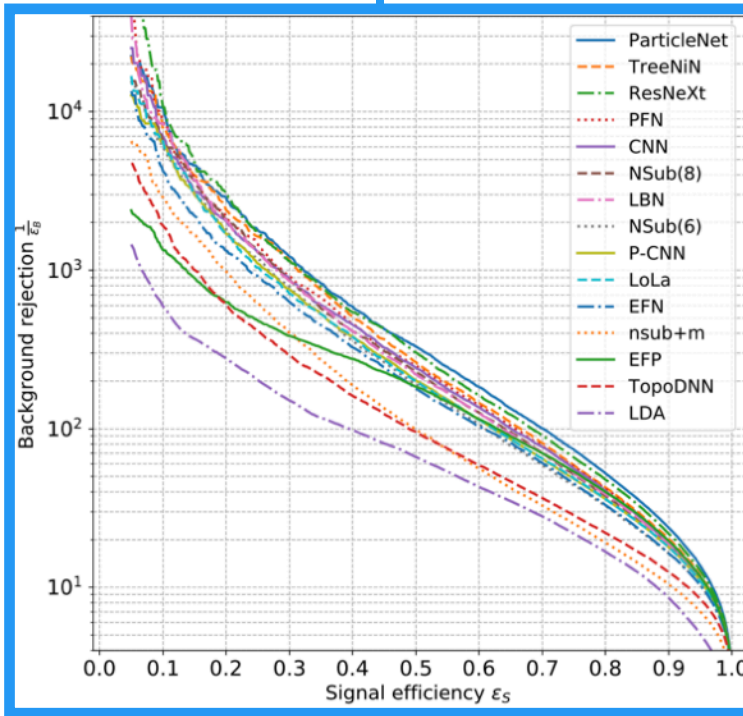
Deep Learning in HEP
Signal vs Background

P. Baldi, P. Sadowski,
D. Whiteson

First community Paper

Machine Learning
Landscape of
TopTagging
G. Kasieczka, et al.

2019

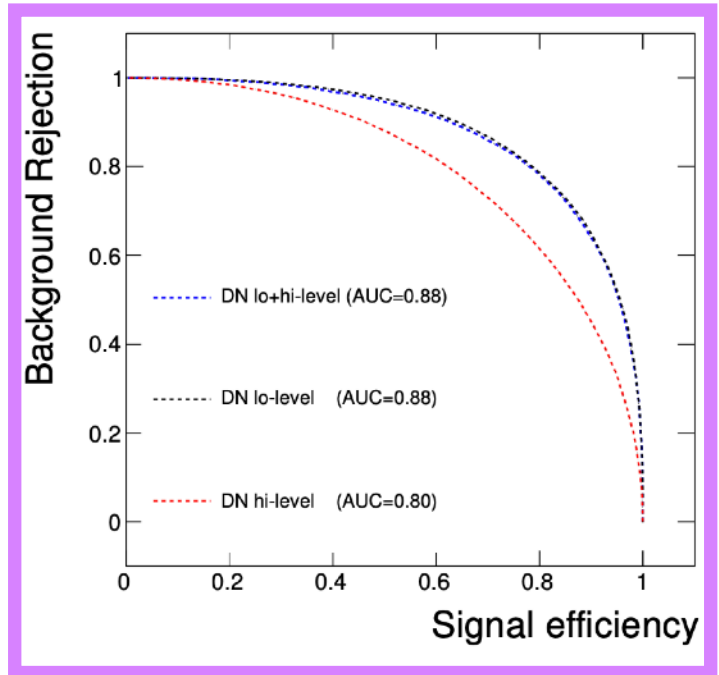
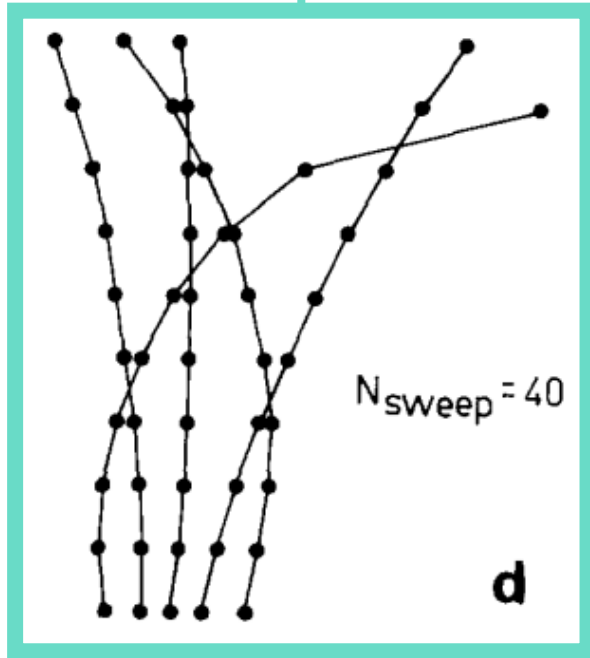


A short history of ML in HEP

First HEP NN papers

Track finding
 Denby (LAL, Orsay) '87
 Peterson (Lund) '88
 → Jet identification

1987/88



2014

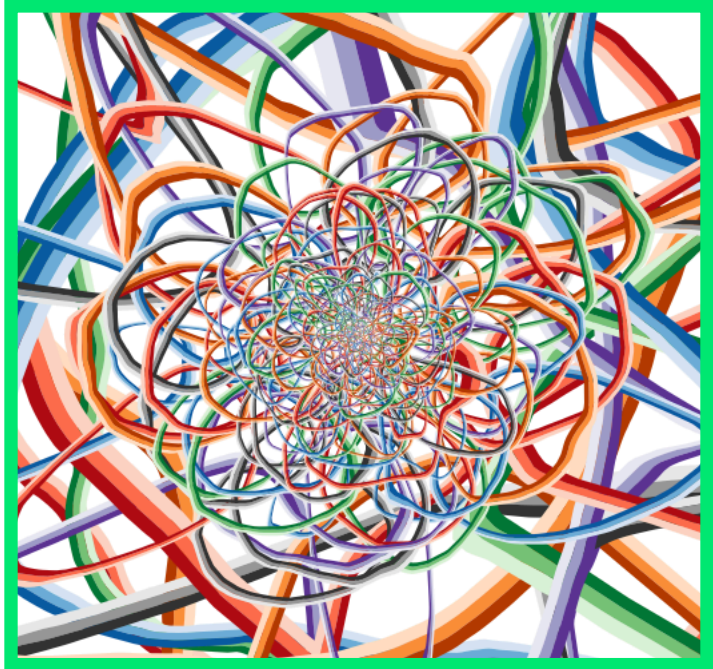
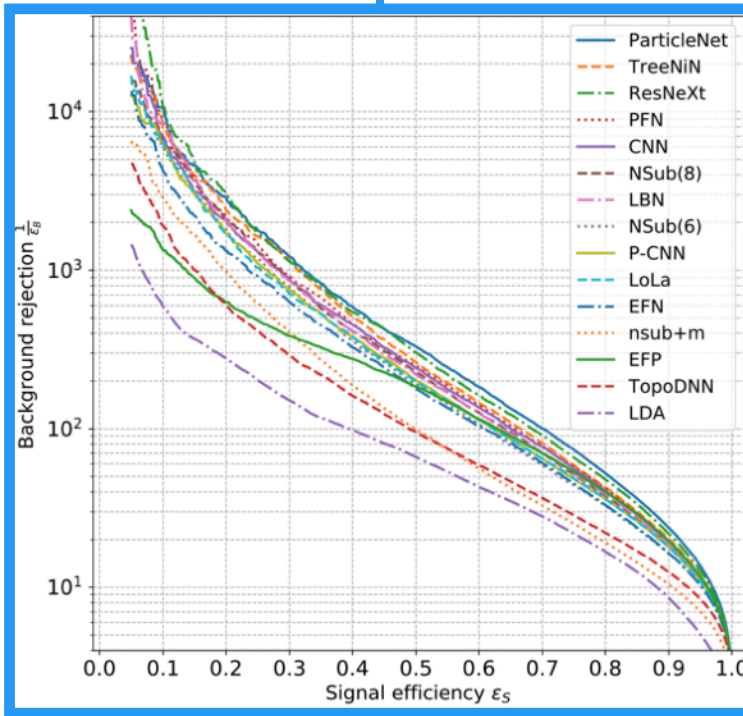


Relaunch
 Deep Learning in HEP
 Signal vs Background
 P. Baldi, P. Sadowski,
 D. Whiteson

First community Paper

Machine Learning
 Landscape of
 TopTagging
 G. Kasieczka, et al.

2019



2019



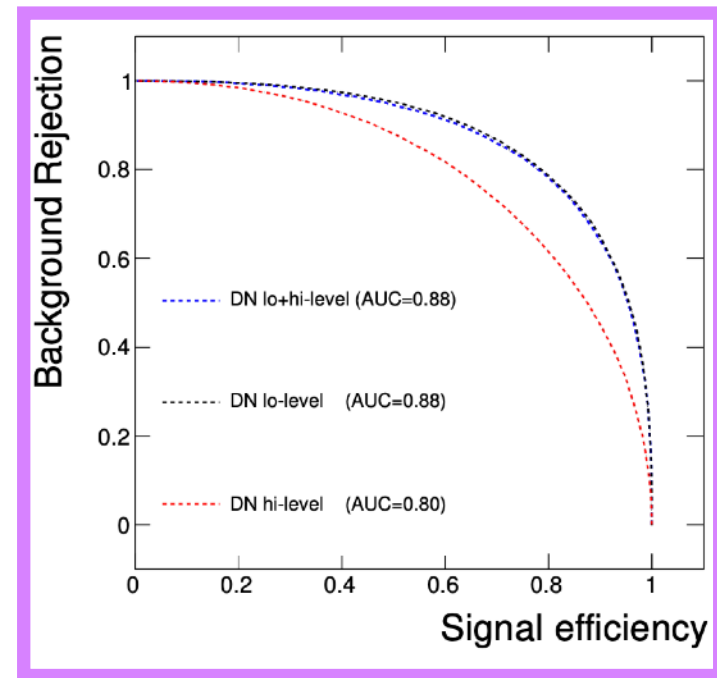
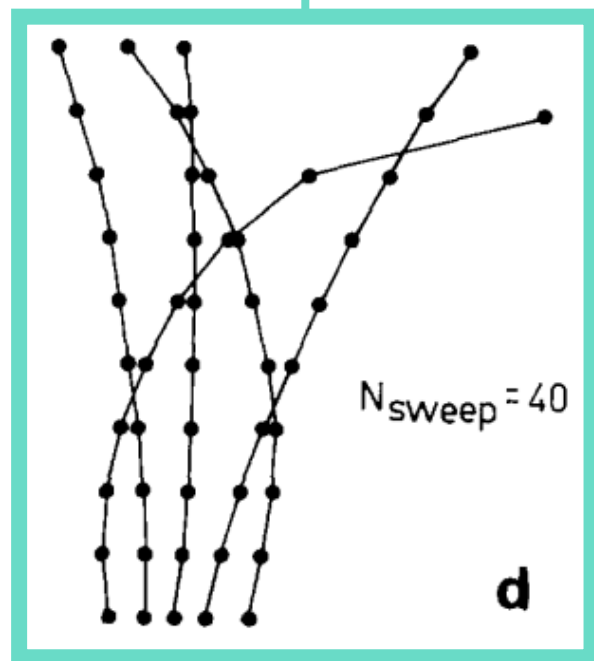
First attempts at "understanding" neural networks
 Deep Thinking
 J. Thaler

A short history of ML in HEP

First HEP NN papers

Track finding
Denby (LAL, Orsay) '87
Peterson (Lund) '88
→ Jet identification

1987/88



2014

Relaunch

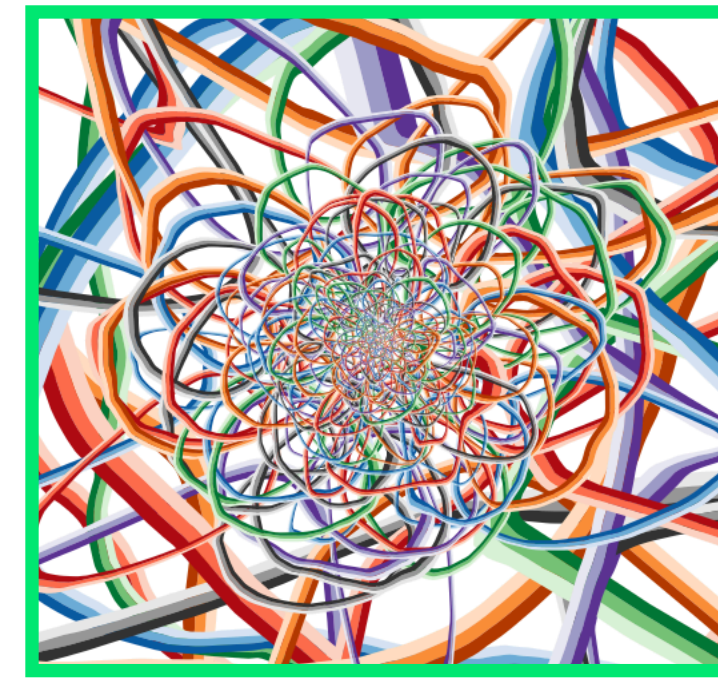
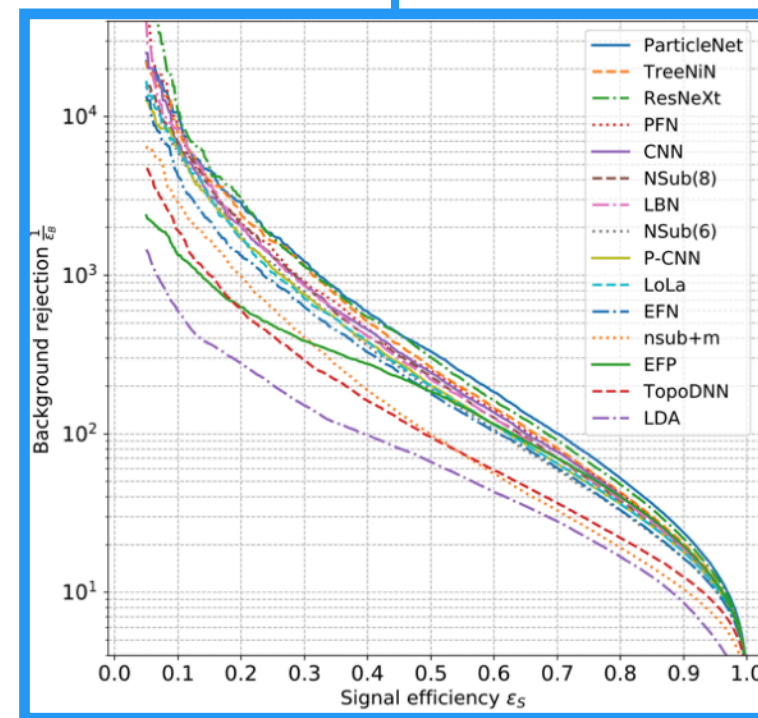
Deep Learning in HEP
Signal vs Background

P. Baldi, P. Sadowski,
D. Whiteson

First community Paper

Machine Learning
Landscape of
TopTagging
G. Kasieczka, et al.

2019



2019

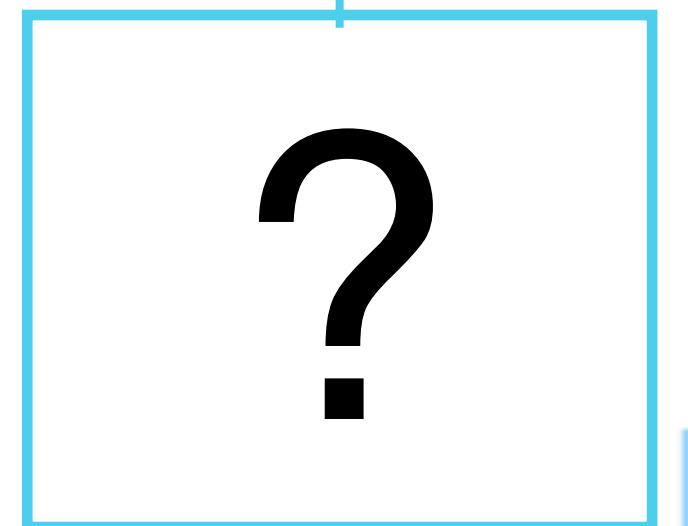
First attempts at
“understanding”
neural networks

Deep Thinking

J. Thaler

Today

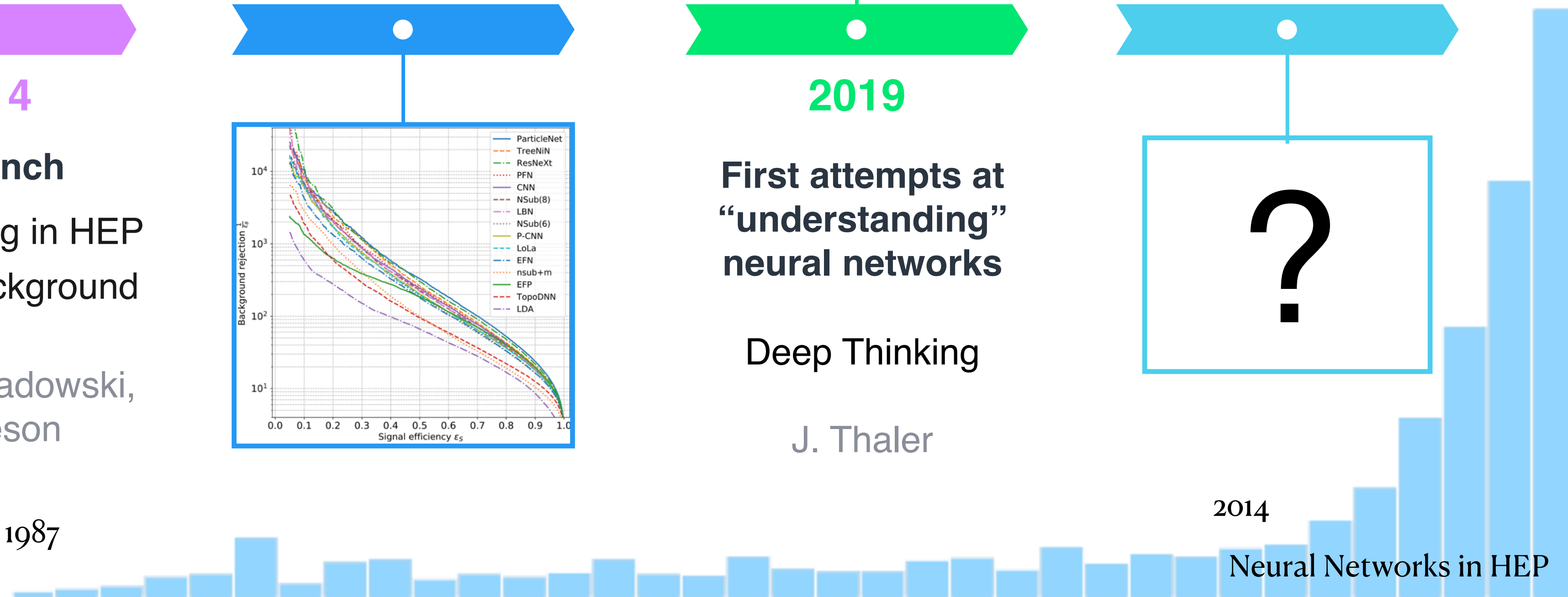
2021



2014

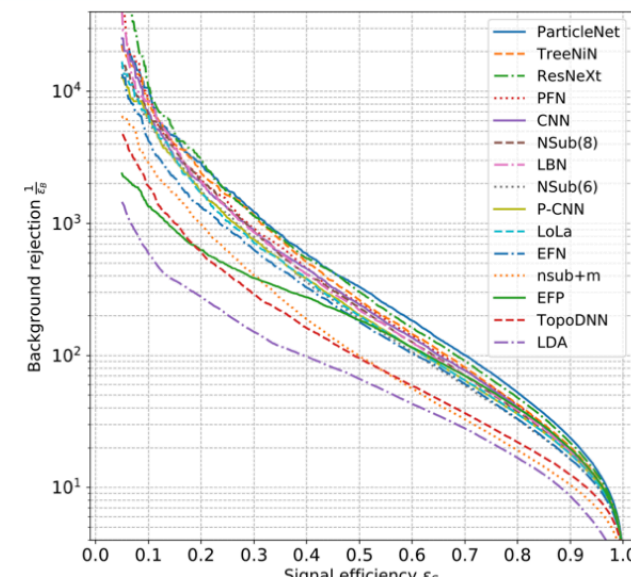
Neural Networks in HEP

1987



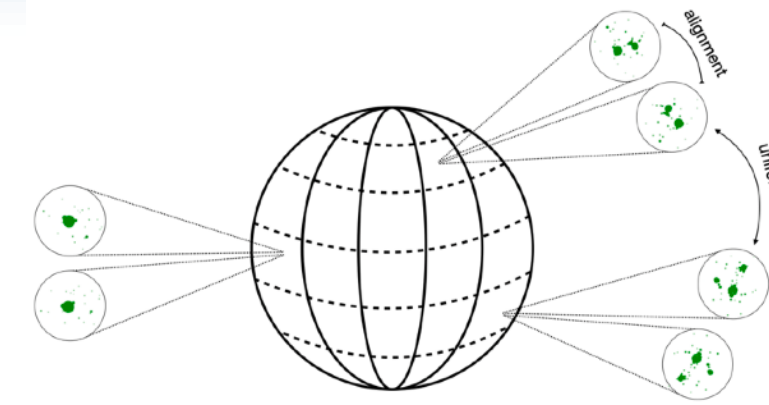
ML for big data in particle physics

Top tagging



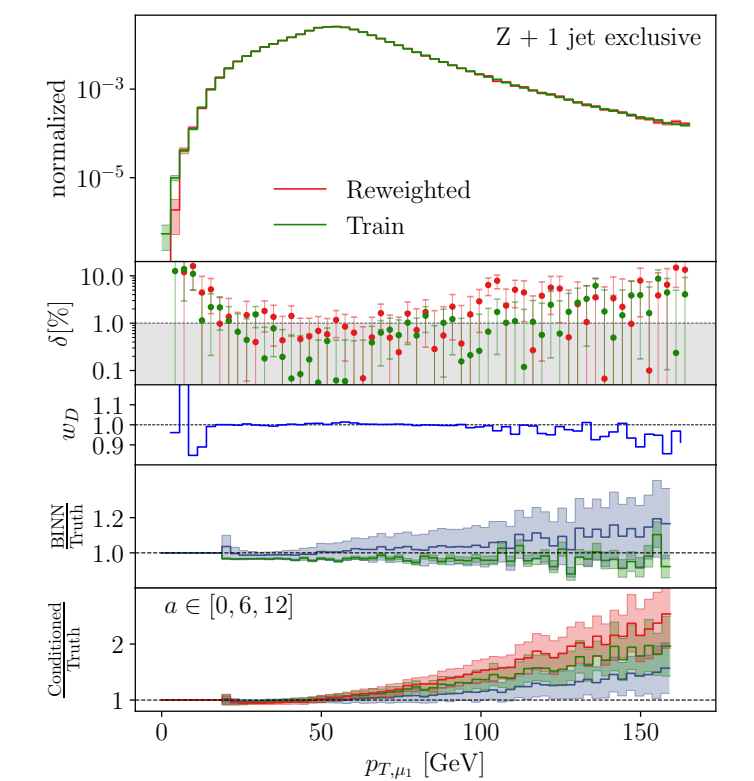
G. Kasieczka et al. [1902.09914]

Anomaly detection



B. Dillon et al. [2108.04253]

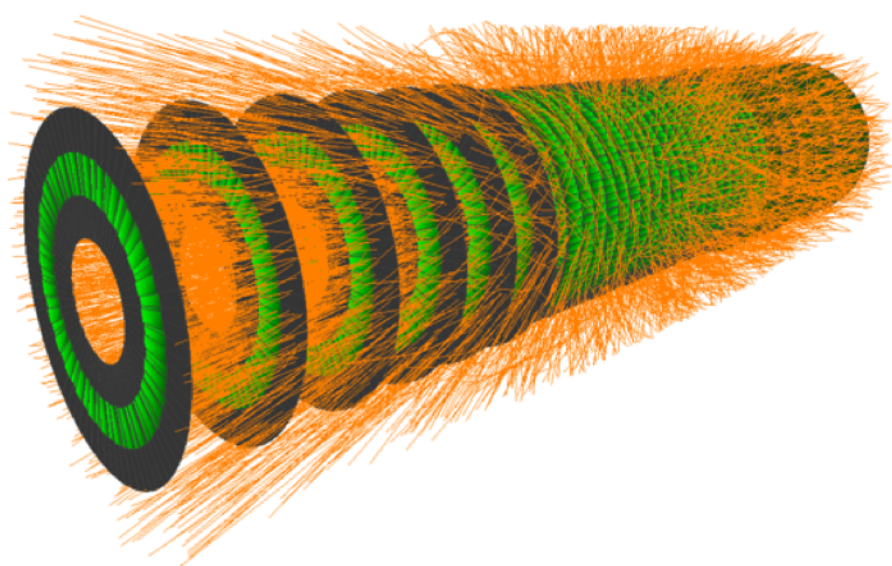
Event generation



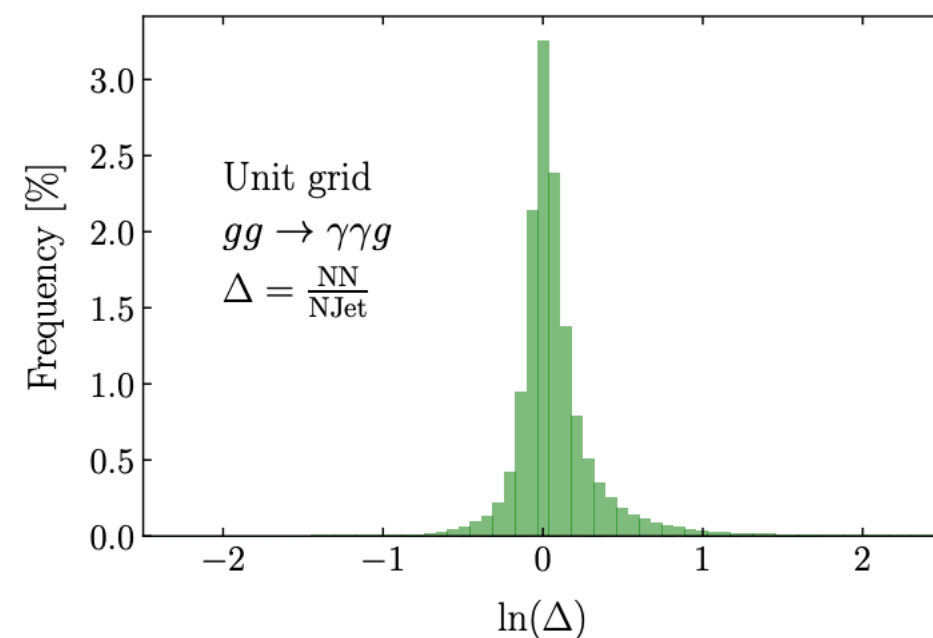
A. Butter et al. [2110.13632]

Track reconstruction

Kaggle challenge



Amplitude estimation



J. Aylett-Bullock, et al. [2106.09474]

Classification

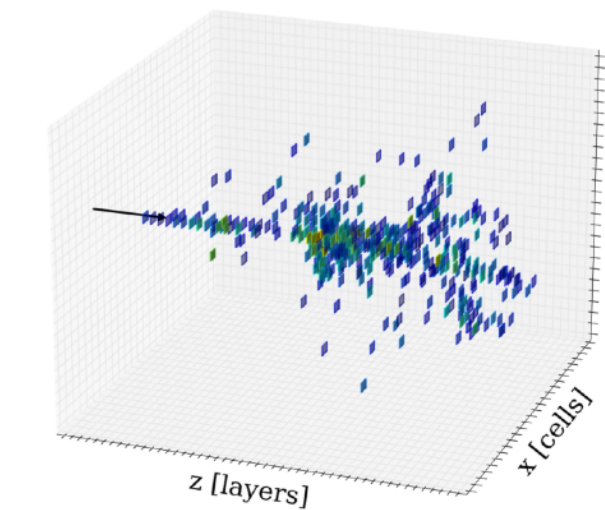
Generative models

Graph networks

Bayesian networks

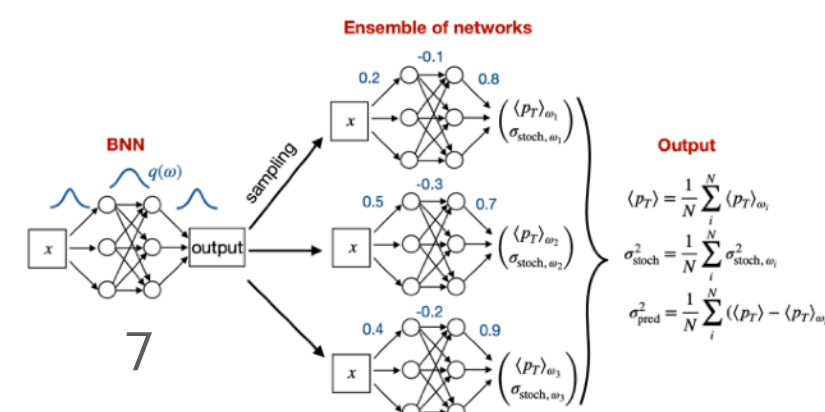
Regression

Detector simulation



E. Buhmann et al. [2112.09709]

Jet calibration & uncertainties



G. Kasieczka et al. [2003.11099]

Complete citations $\mathcal{O}(800)$
<https://iml-wg.github.io/HEPML-LivingReview/>

ML examples and their uncertainties

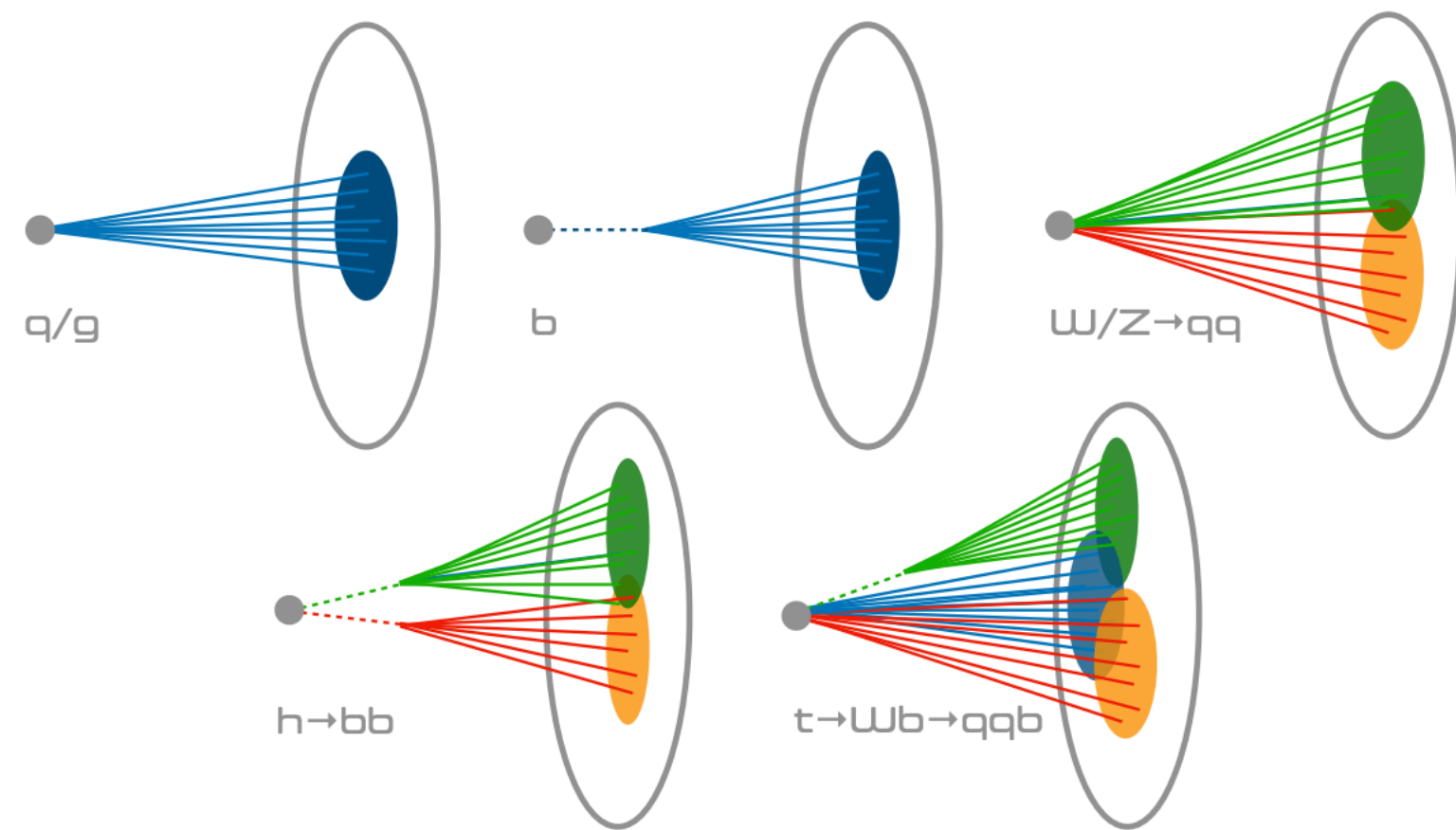
Classification

Simulations

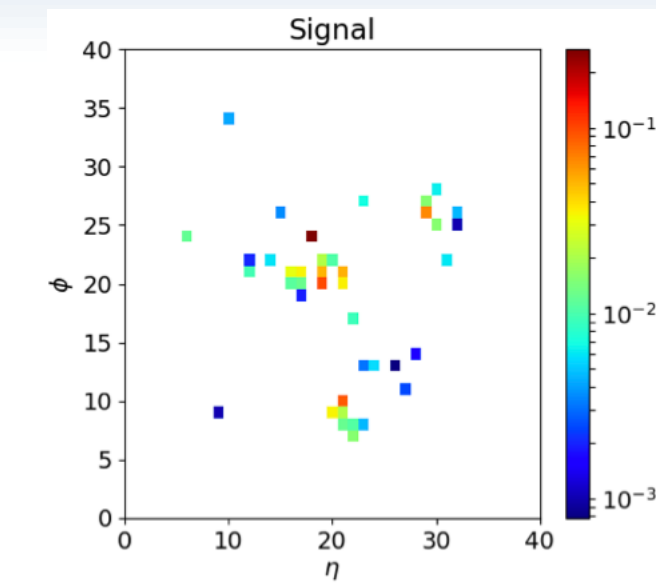
Unfolding

Jet classification

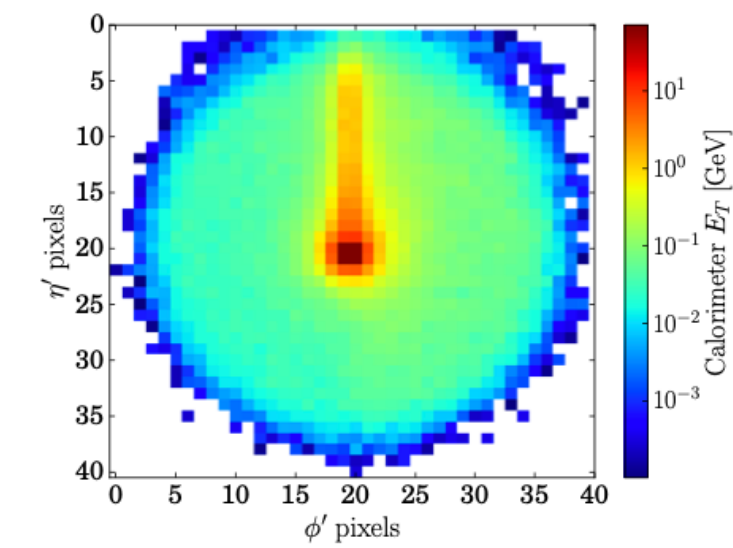
- How to distinguish **top** from **QCD** jets?
- Immensely important for top & Higgs physics studies



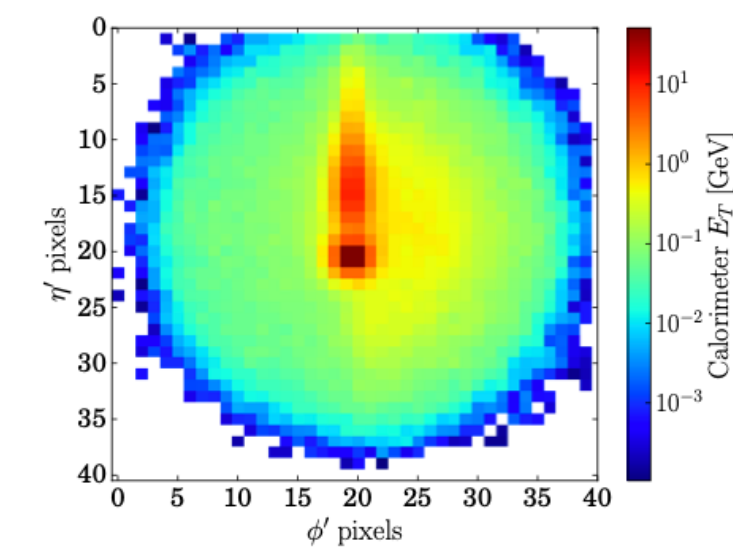
E. Moreno et al. [1909.12285]



Single jet



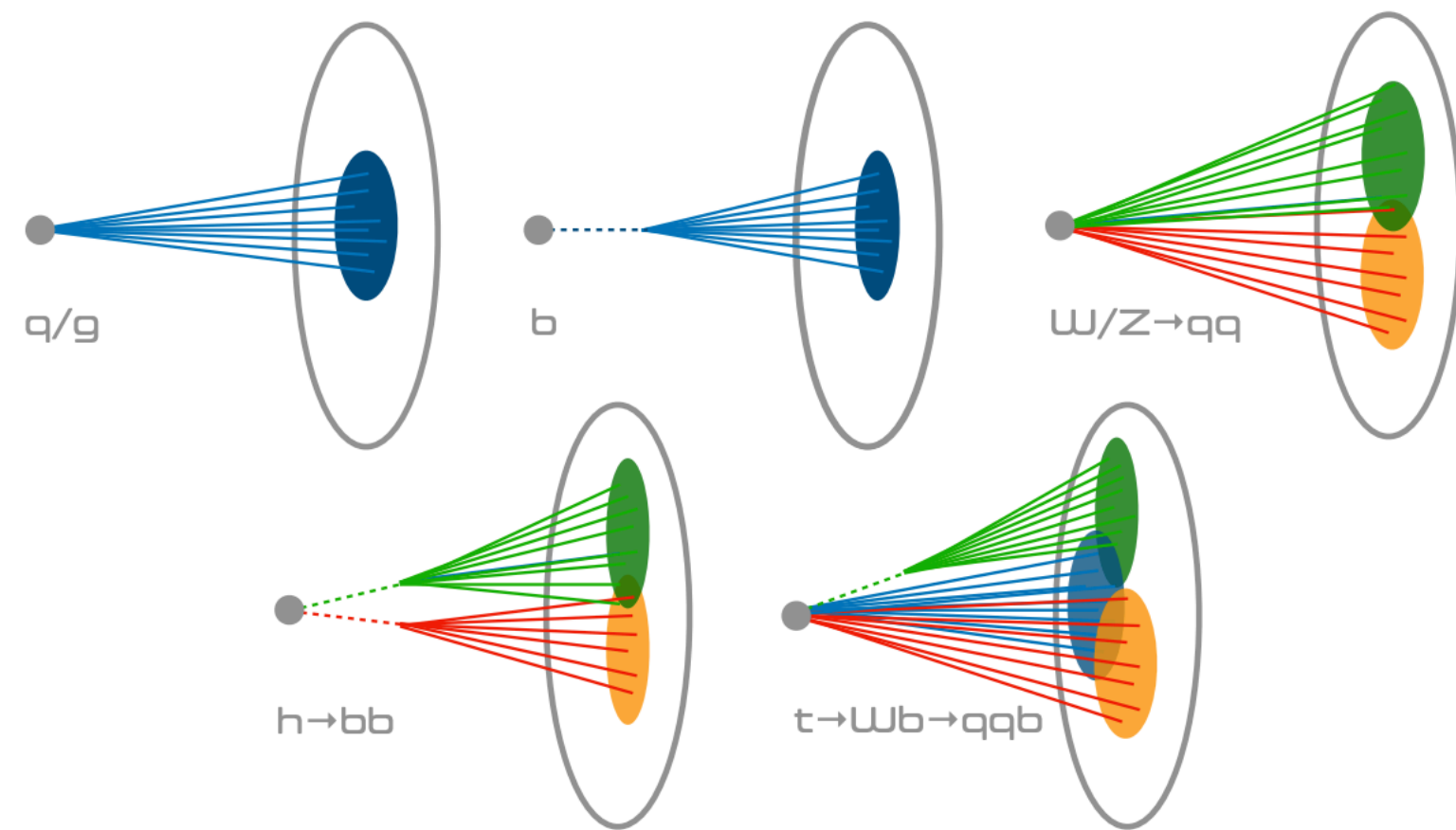
QCD jet averaged



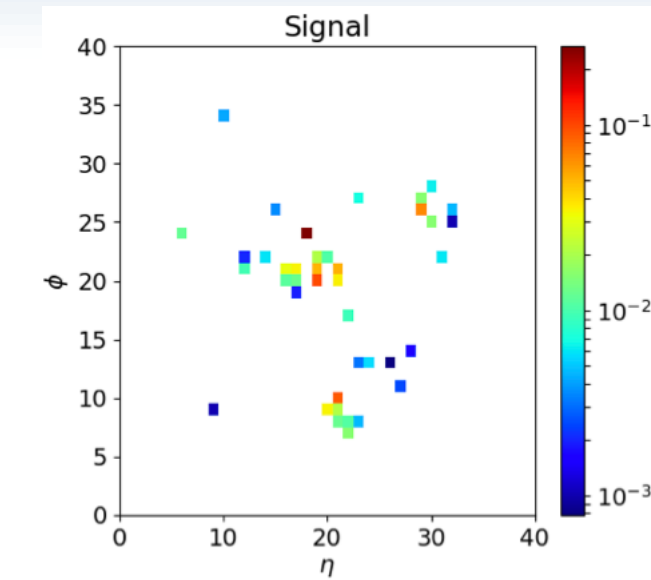
top jet Averaged

Jet classification

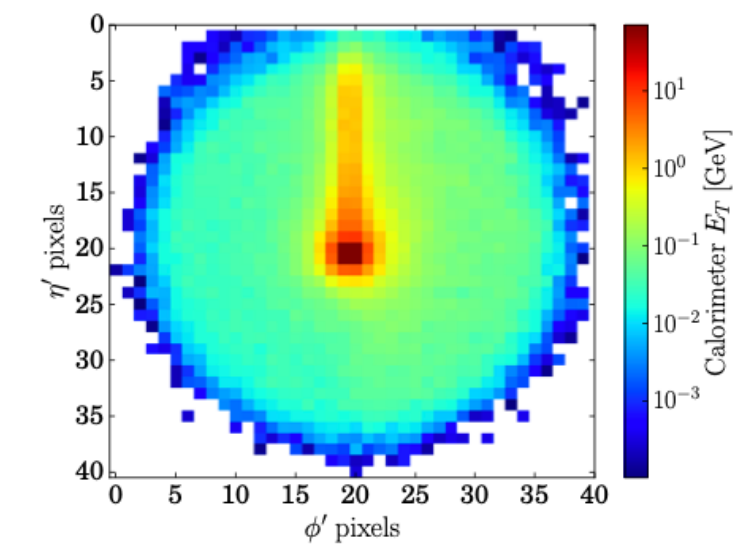
- How to distinguish **top** from **QCD** jets?
- Immensely important for top & Higgs physics studies
- Standard supervised classification task



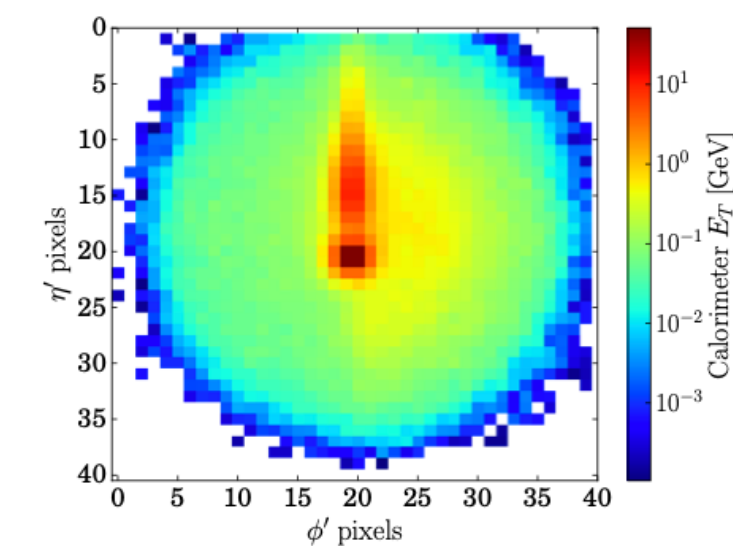
E. Moreno et al. [1909.12285]



Single jet



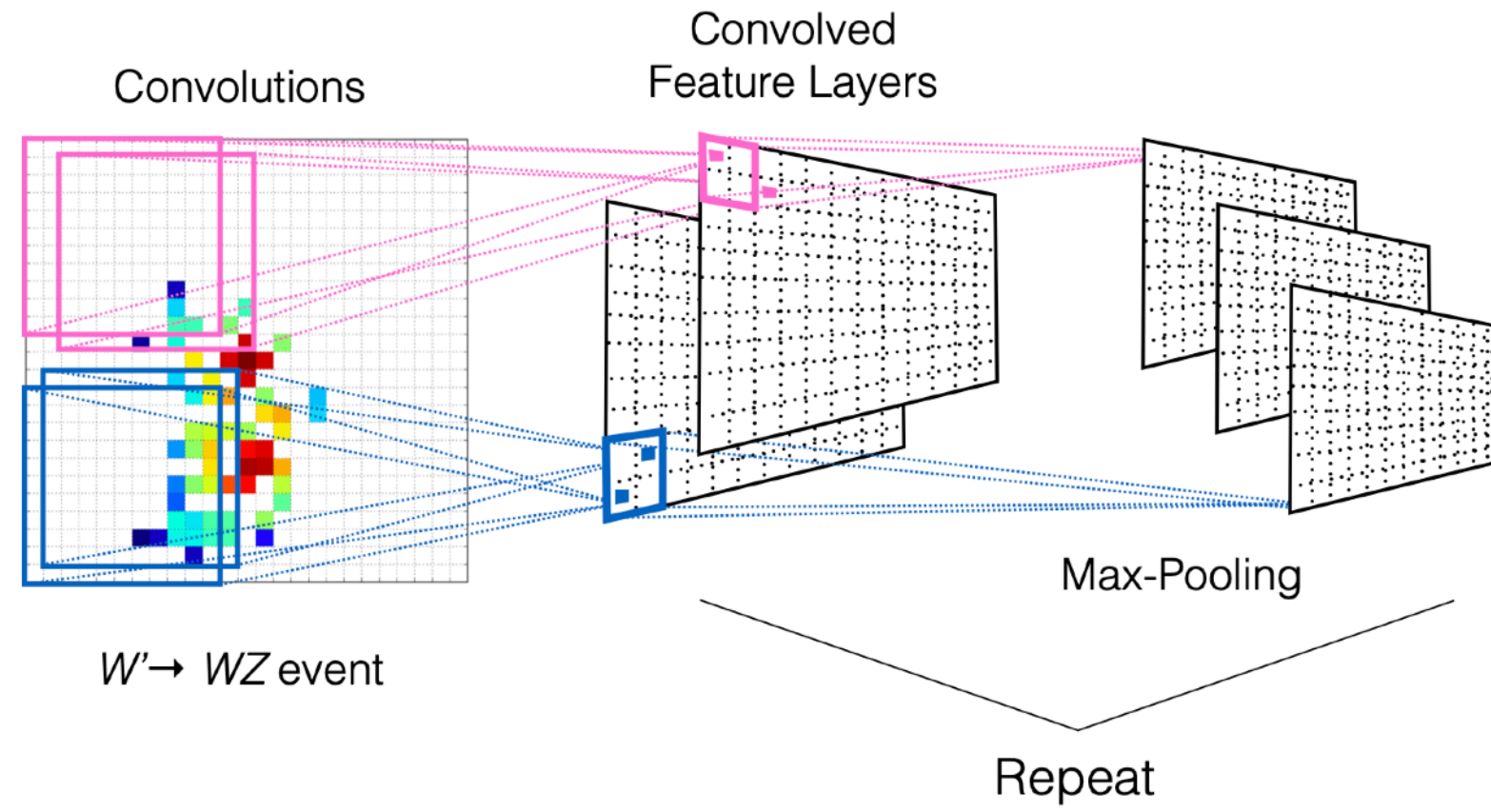
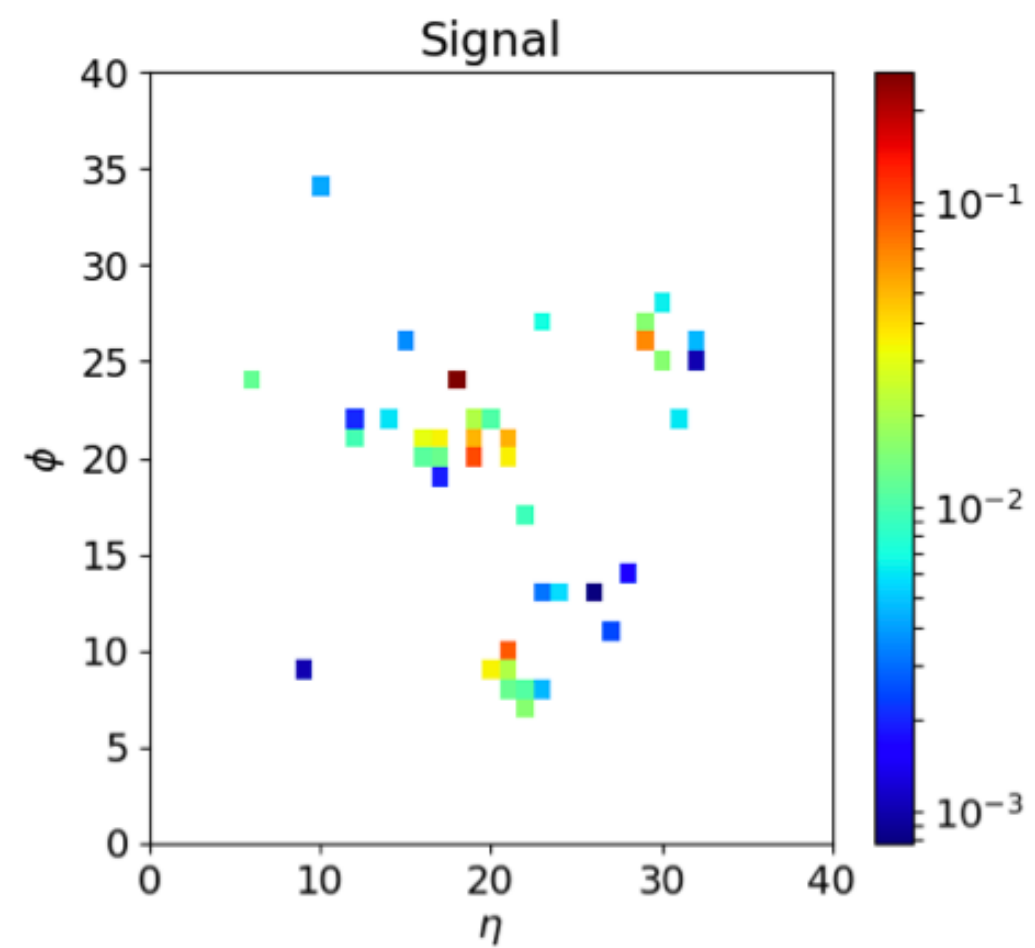
QCD jet averaged



top jet Averaged

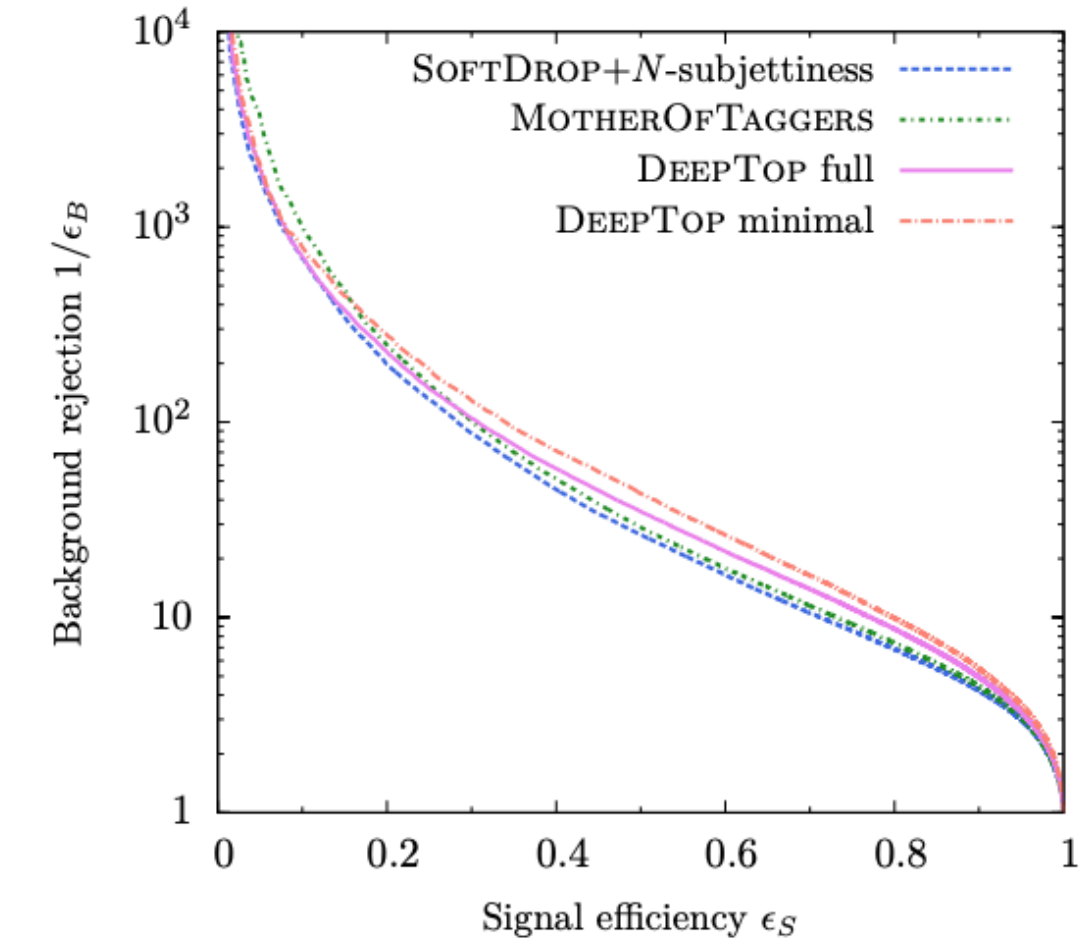
Jet representation I

Single jet



L. Oliveira et al. [1511.05190]

CNNs

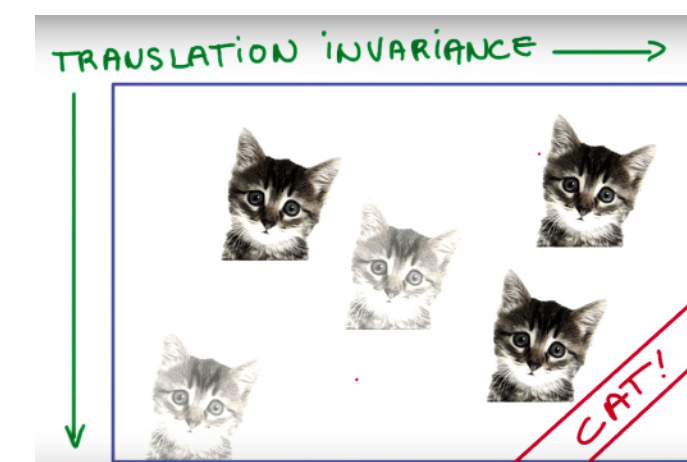


G. Kasieczka et al. [1701.08784]

Use **equivariance** principle for translation t :

$$\text{CNN}(t(x)) = t(\text{CNN}(x))$$

Pooling layer
→ **Invariance**

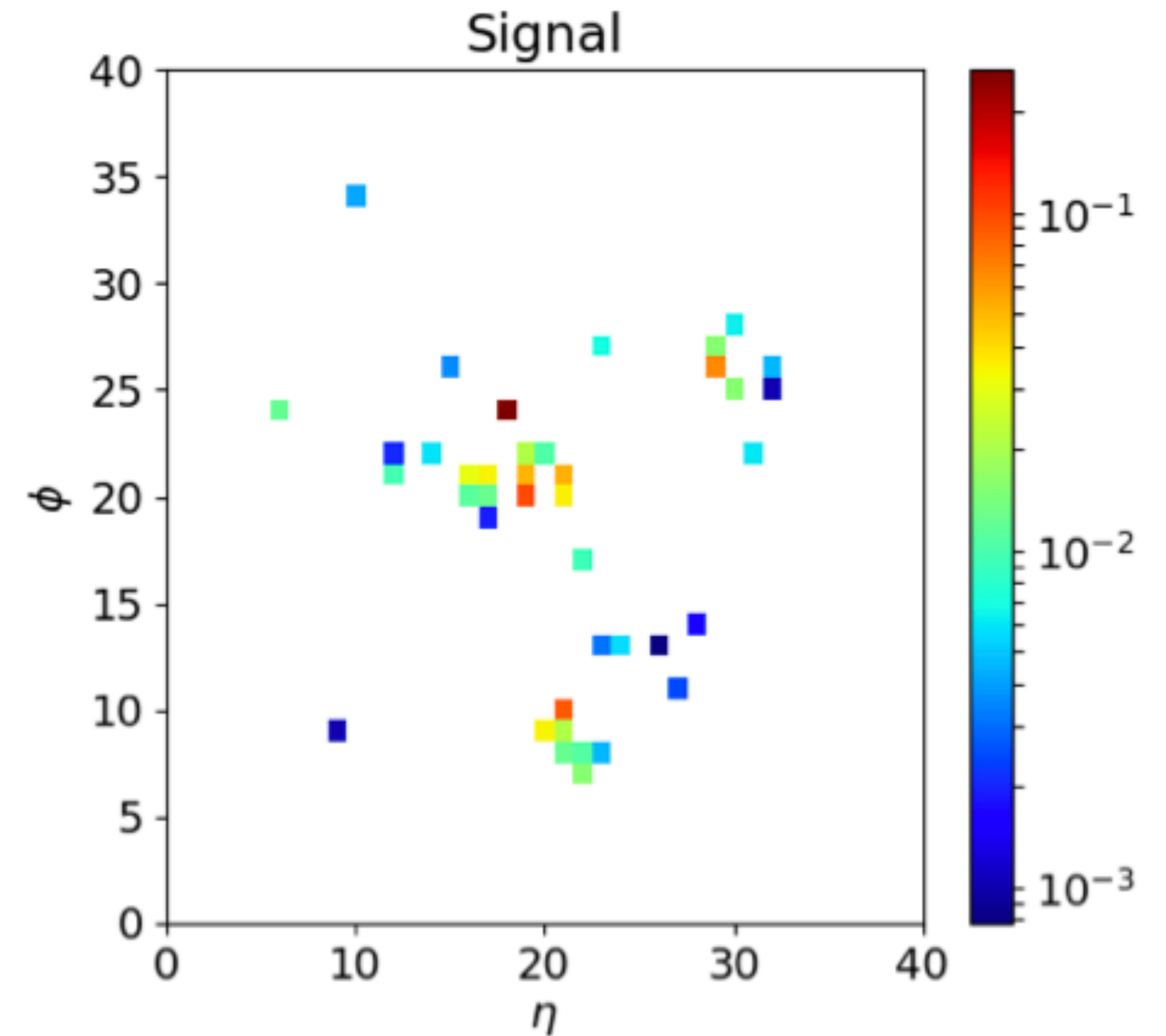


Our image ain't a very good image...

No continuity, no edges, no cats....



VS



Jet representation II

How to represent a graph

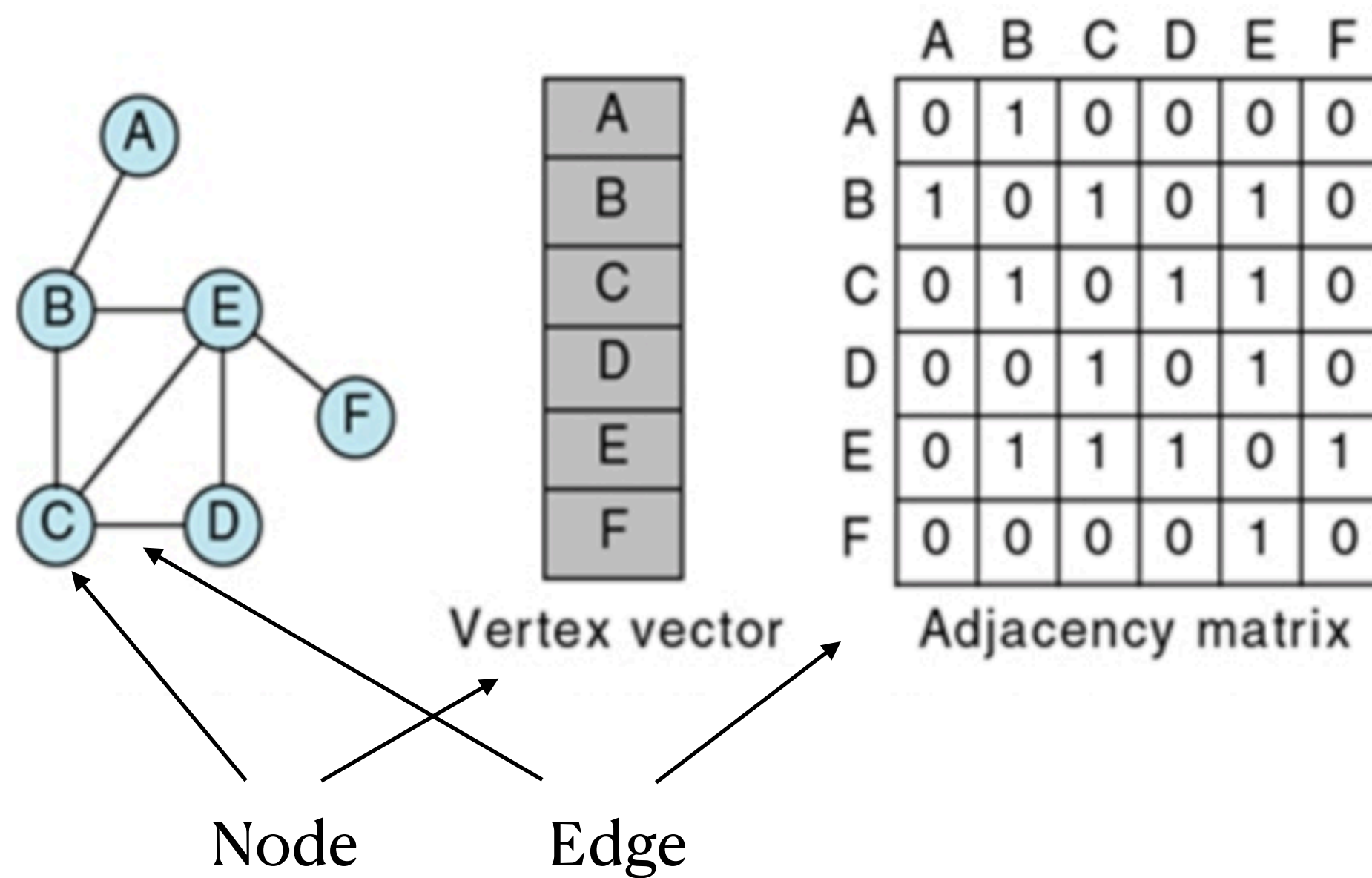
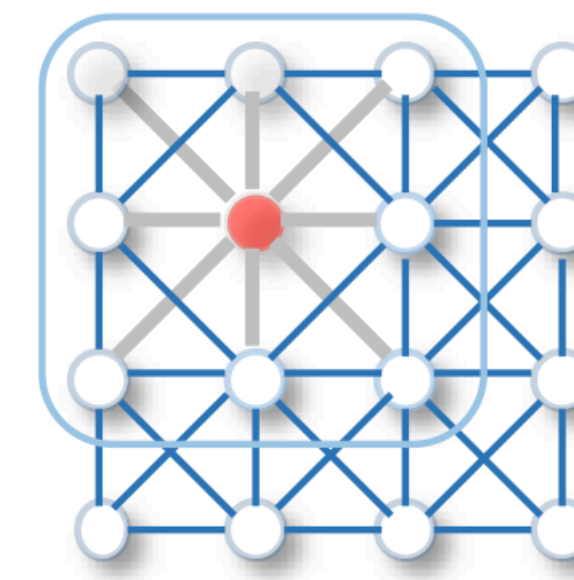
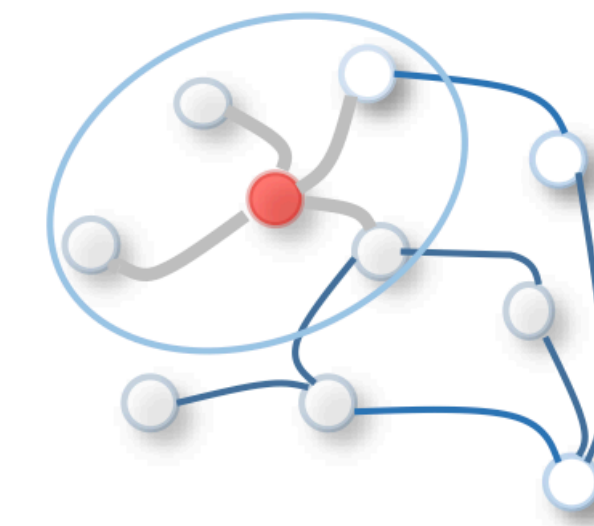


Image vs Graph



pixels
neighbouring pixel

CNN

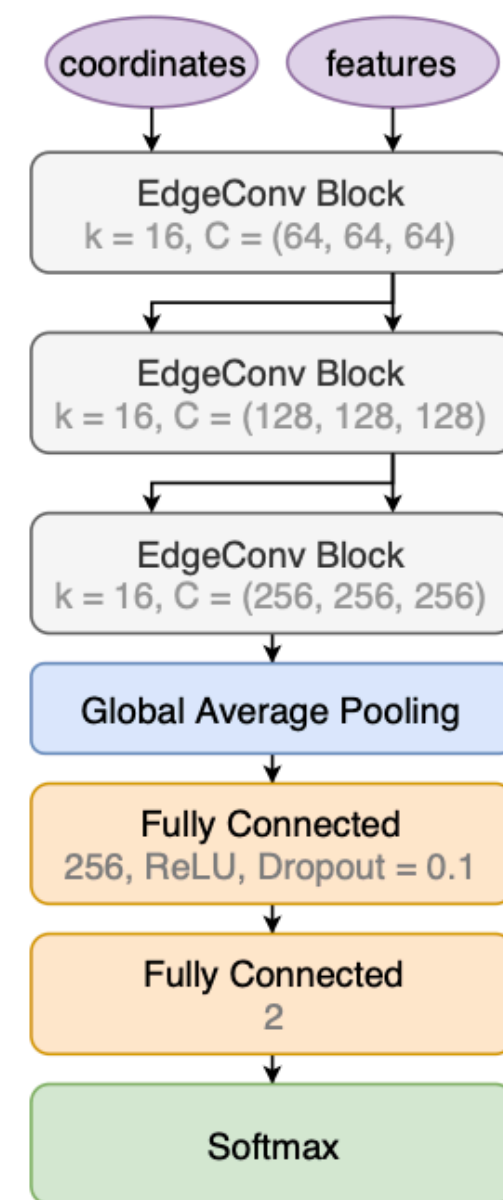


→ node
→ neighbouring node (graph edges)

→ edge convolution

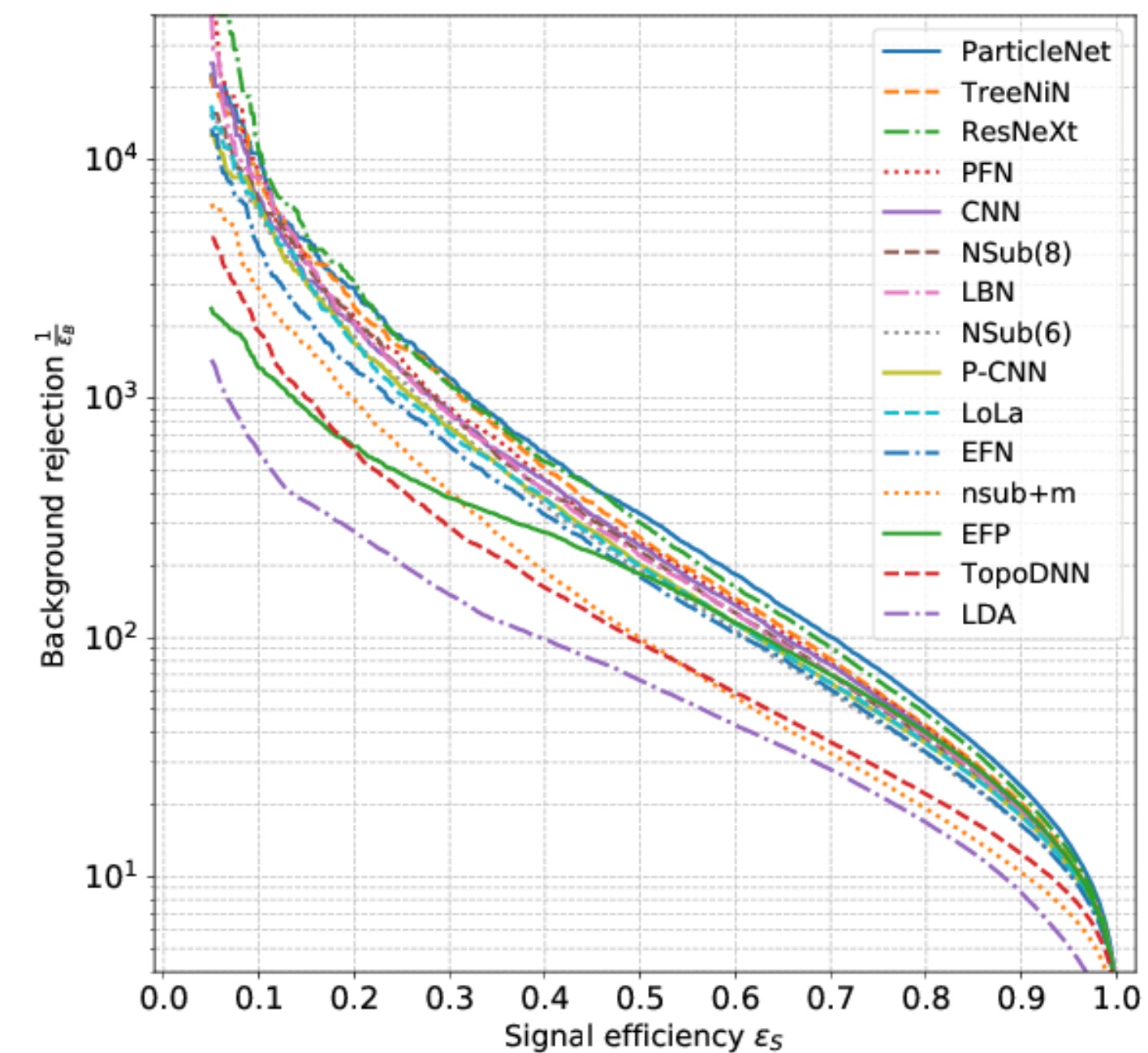
$$\vec{x}'_i = \frac{1}{k} \sum_{j=1}^k h_{\Theta}(\vec{x}_i, \vec{x}_{i_j} - \vec{x}_i)$$

The ML landscape of top taggers



(a) ParticleNet

H. Qu, L. Gouskos [1902.08570]



Difference in performance for various different approaches!

Where are the uncertainties?

What if the network is not perfect?

A perfect network has learned the likelihood ratio $\frac{p(x | top)}{p(x | QCD)}$

What if the network is not perfect?

Classification loss function

$$\begin{aligned}\mathcal{L} &= \sum_{x_i} -\log C(x_i) y_i - \log(1 - C(x_i)) (1 - y_i) \\ &= - \int dx p_{top}(x) \log C(x) + p_{QCD} \log(1 - C(x))\end{aligned}$$

$$\text{Variance yields } \rightarrow \frac{p_{top}(x)}{p_{QCD}(x)} = \frac{C(x)}{1 - C(x)}$$

A perfect network has learned the likelihood ratio $\frac{p(x | top)}{p(x | QCD)}$

What if the network is not perfect?

A suboptimal network will label more tops "wrong"

Classification loss function

$$\begin{aligned}\mathcal{L} &= \sum_{x_i} -\log C(x_i) y_i - \log(1 - C(x_i)) (1 - y_i) \\ &= - \int dx p_{top}(x) \log C(x) + p_{QCD} \log(1 - C(x))\end{aligned}$$

$$\text{Variance yields } \rightarrow \frac{p_{top}(x)}{p_{QCD}(x)} = \frac{C(x)}{1 - C(x)}$$

A perfect network has learned the likelihood ratio $\frac{p(x | top)}{p(x | QCD)}$

What if the network is not perfect?

Classification loss function

$$\begin{aligned}\mathcal{L} &= \sum_{x_i} -\log C(x_i) y_i - \log(1 - C(x_i)) (1 - y_i) \\ &= - \int dx p_{top}(x) \log C(x) + p_{QCD} \log(1 - C(x))\end{aligned}$$

$$\text{Variance yields } \rightarrow \frac{p_{top}(x)}{p_{QCD}(x)} = \frac{C(x)}{1 - C(x)}$$

A suboptimal network will label more tops "wrong"

Applies to prediction & data!

Equivalent to **poor efficiency** in

$$p(\text{data} | \theta, \nu) = \prod_{i=1}^{n_{bins}} P(n_i | s_i(\theta, \nu) \cdot \epsilon_i(\nu) + b_i(\nu)) P(\nu | \text{aux. data})$$

A perfect network has learned the likelihood ratio $\frac{p(x | top)}{p(x | QCD)}$

What if the network is not perfect?

What is the network supposed to learn?

Classification loss function

$$\begin{aligned}\mathcal{L} &= \sum_{x_i} -\log C(x_i) y_i - \log(1 - C(x_i)) (1 - y_i) \\ &= - \int dx p_{top}(x) \log C(x) + p_{QCD} \log(1 - C(x))\end{aligned}$$

Variance yields $\rightarrow \frac{p_{top}(x)}{p_{QCD}(x)} = \frac{C(x)}{1 - C(x)}$

A suboptimal network will label more tops "wrong"

Applies to prediction & data!

Equivalent to **poor efficiency** in

$$p(data | \theta, \nu) = \prod_{i=1}^{n_{bins}} P(n_i | s_i(\theta, \nu) \cdot \epsilon_i(\nu) + b_i(\nu)) P(\nu | \text{aux. data})$$

A perfect network has learned the likelihood ratio $\frac{p(x | top)}{p(x | QCD)}$

What if the network is not perfect?

What is the network supposed to learn?

Classification loss function

$$\begin{aligned}\mathcal{L} &= \sum_{x_i} -\log C(x_i) y_i - \log(1 - C(x_i)) (1 - y_i) \\ &= - \int dx p_{top}(x) \log C(x) + p_{QCD} \log(1 - C(x))\end{aligned}$$

$$\text{Variance yields } \rightarrow \frac{p_{top}(x)}{p_{QCD}(x)} = \frac{C(x)}{1 - C(x)}$$

A suboptimal network will label more tops "wrong"

Applies to prediction & data!

Equivalent to poor efficiency in

$$p(\text{data} | \theta, \nu) = \prod_{i=1}^{n_{bins}} P(n_i | s_i(\theta, \nu) \cdot \epsilon_i(\nu) + b_i(\nu)) P(\nu | \text{aux. data})$$

A perfect network has learned the likelihood ratio $\frac{p(x | top)}{p(x | QCD)}$

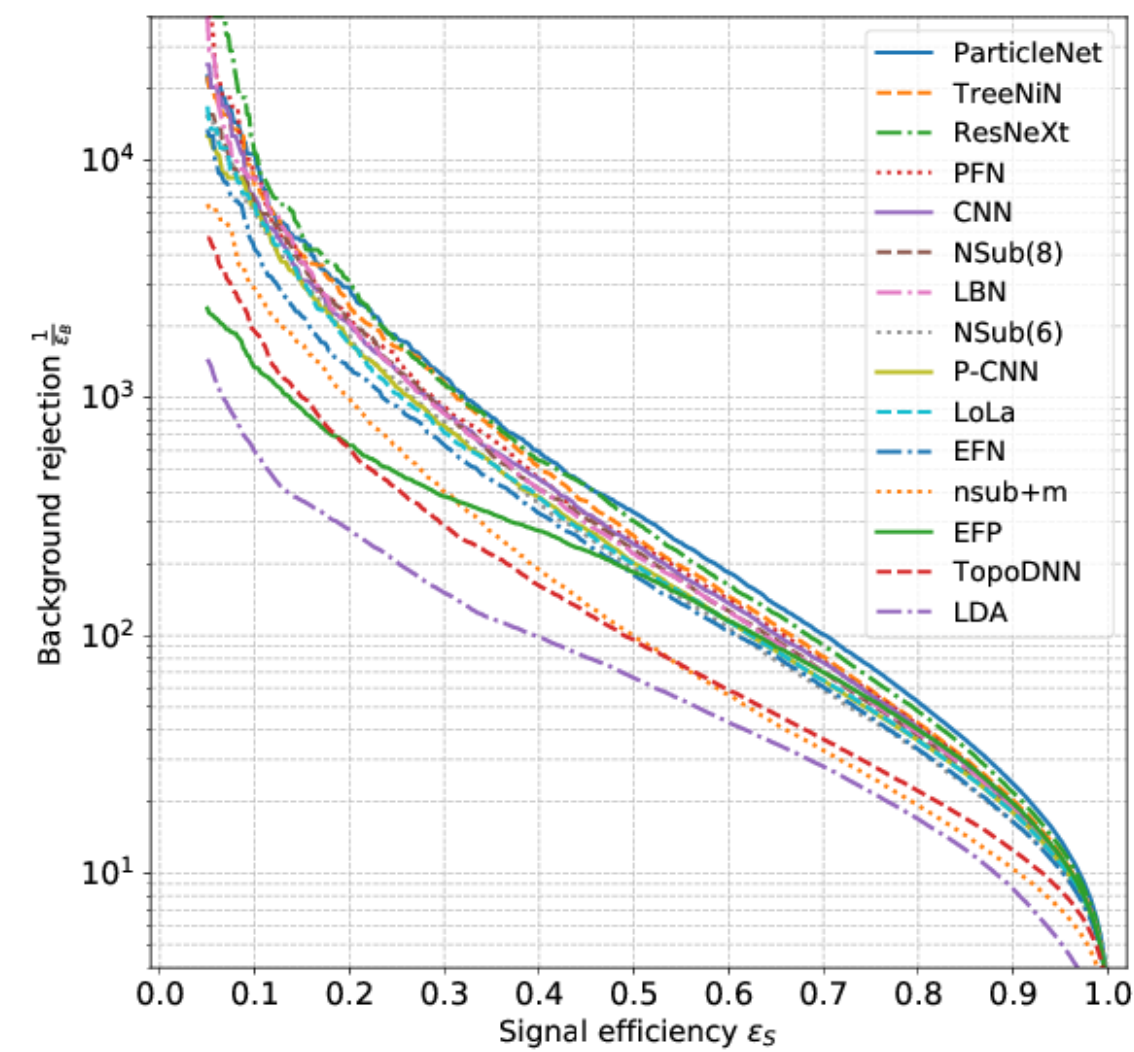
The result is not optimal - but still correct!

ML examples and their uncertainties

Classification

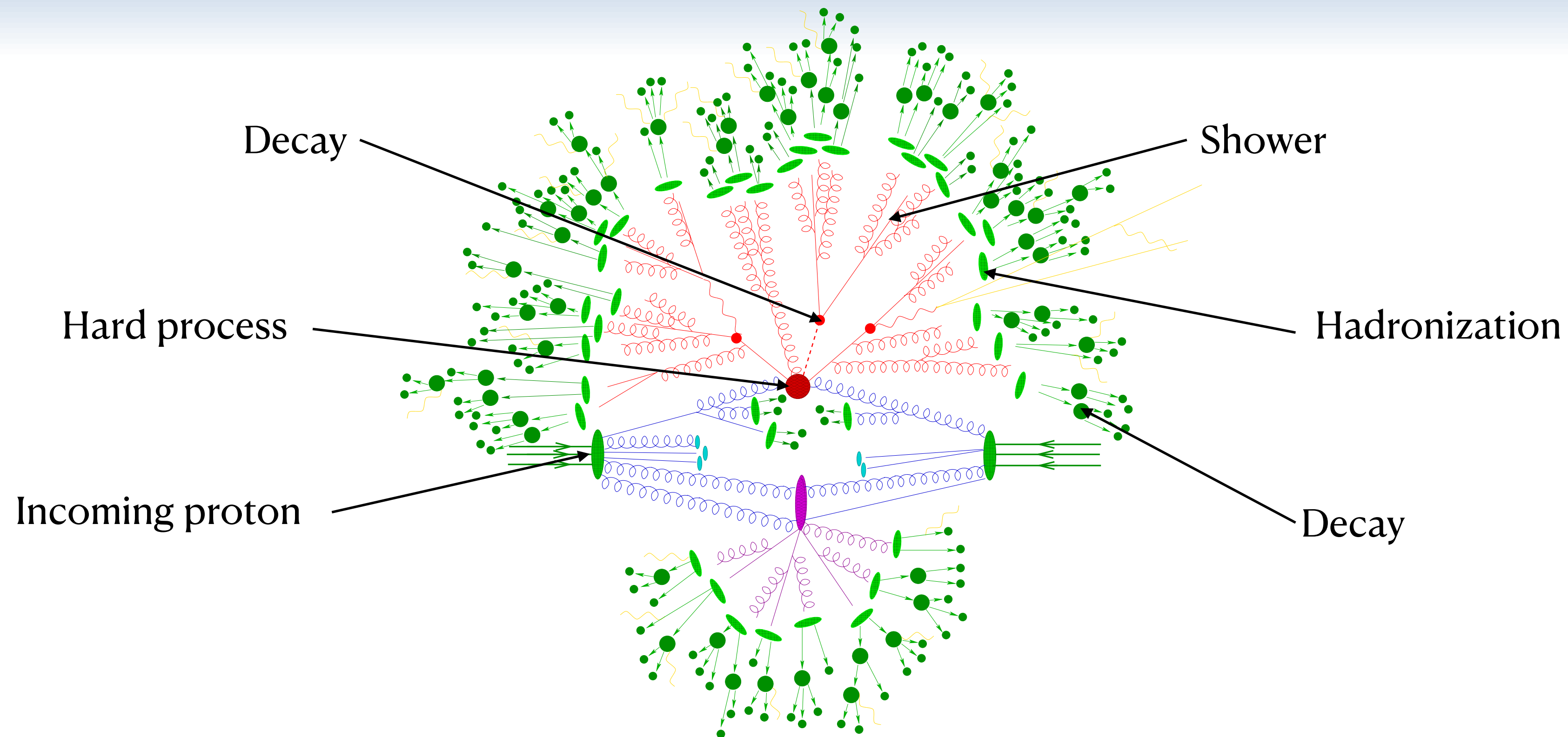
Simulations

Unfolding

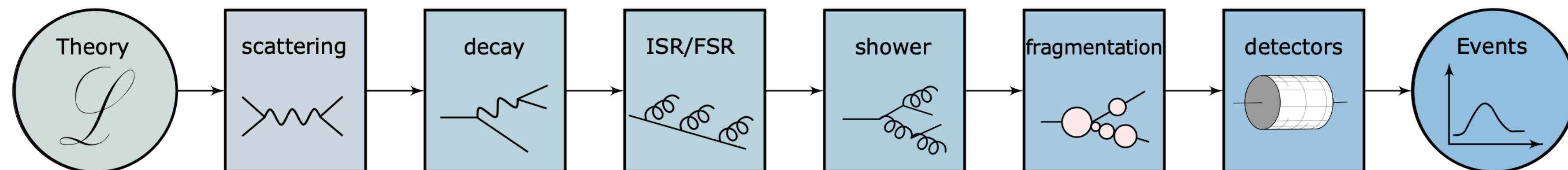


No uncertainty needed

Event generation at the LHC



Simulation tool chain



Monte carlo event generation

1. Generate phase space points

→ set of four-momenta p_i

2. Calculate event weight

$$w_{\text{event}} = \underbrace{f(x_1, Q^2)f(x_2, Q^2)}_{\text{PDF}} \times \underbrace{\mathcal{M}(x_1, x_2, p_1, \dots, p_n)}_{\text{Matrix element †}} \times \underbrace{J(p_i(r))}_{\text{Phase space mapping}}$$

3. Unweighting †

keep events with $\frac{w_i}{w_{\text{max}}} > r \in [0,1]$

Monte carlo event generation

1. Generate phase space points

→ set of four-momenta p_i

2. Calculate event weight

$$w_{\text{event}} = \underbrace{f(x_1, Q^2)f(x_2, Q^2)}_{\text{PDF}} \times \underbrace{\mathcal{M}(x_1, x_2, p_1, \dots, p_n)}_{\text{Matrix element †}} \times \underbrace{J(p_i(r))}_{\text{Phase space mapping}}$$

3. Unweighting †

keep events with $\frac{w_i}{w_{\text{max}}} > r \in [0,1]$

† Bottlenecks

1. Slow **matrix element** calculation
 - ◆ Complexity grows exponentially with
 - # final state particles
 - Precision (LO, NLO, NNLO, ...)
2. Low **unweighting** efficiency
 - ◆ Discard most events if $w_i \ll w_{\text{max}}$
 - ◆ Optimize phase space mapping
 - ➔ $J(p_i(r)) = (f \times \mathcal{M})^{-1}$

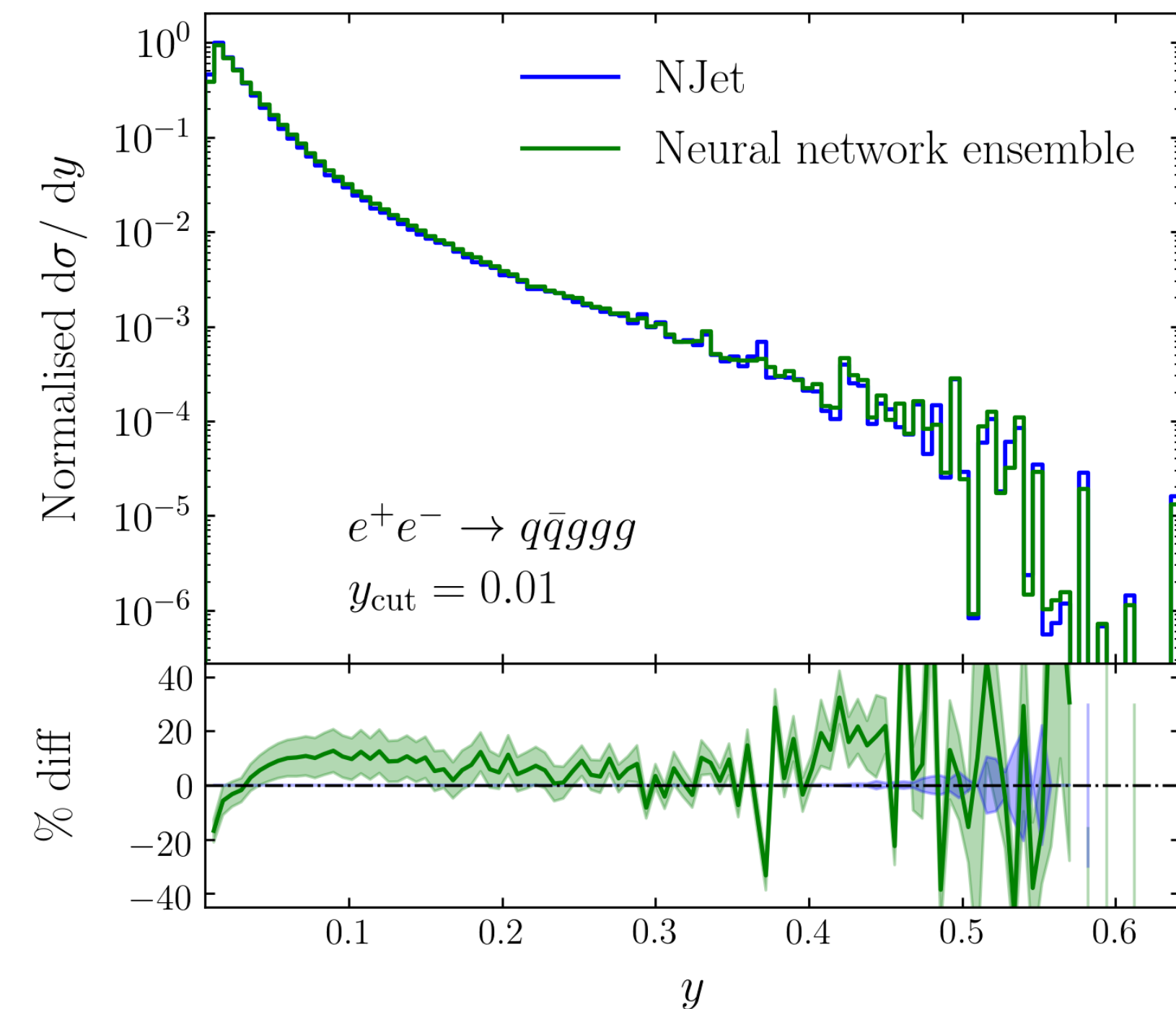
Approximating Amplitudes

- Approximate squared matrix element with NN
- **Regression** problem
- Minimize distance between prediction and truth
 $\Rightarrow \mathcal{L} = (NN(p_i) - \mathcal{M}(p_i))^2$
- + Generalization of interpolation
- + Better **scaling** than grids for large dimensions
- Open questions
 - Limited precision ?
 - Overtraining vs interpolation ?

Problem

Wrong estimation leads to wrong prediction!

→ assign uncertainties



Badger, Bullock [2002.07516]

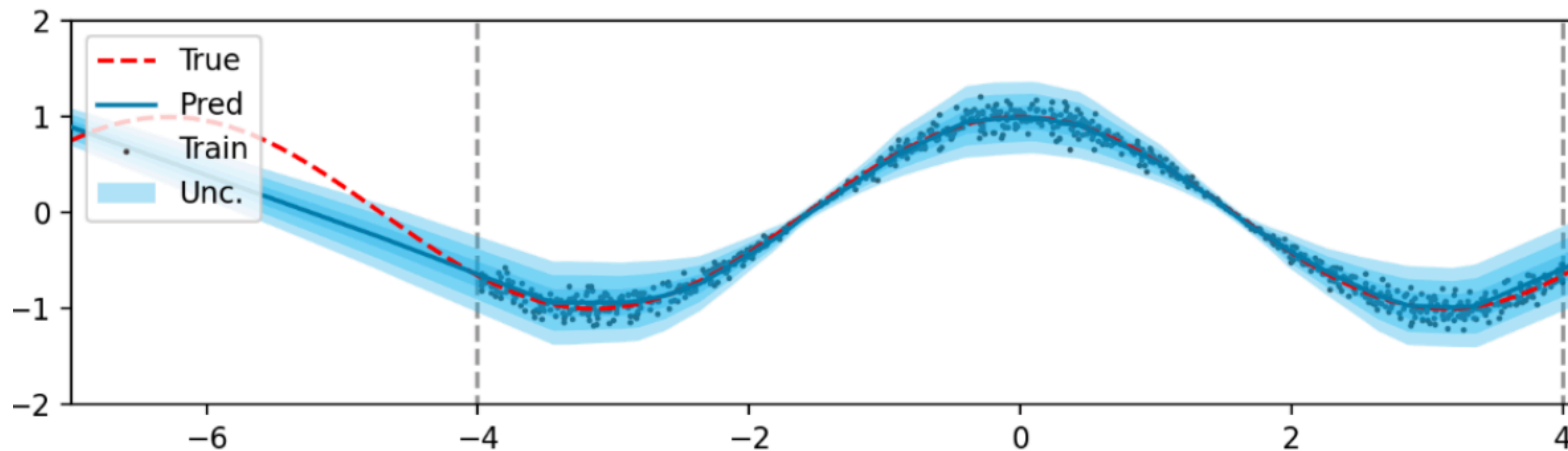
Estimating uncertainties on amplitude predictions

1. **Extend** standard network **output** to include uncertainty

→ $(\mu(x), \sigma(x))$

- Gaussian approximation

- $\mathcal{L}_{\text{Gauss}} = -\log(\sqrt{2\pi}\sigma(x)) - \frac{1}{2} \frac{(\mu(x) - y)^2}{\sigma(x)^2}$



- Captures only $\mathbf{p}(y | \mathbf{x}, \mathbf{w})$ for fixed network weights
- w varies for different trainings!

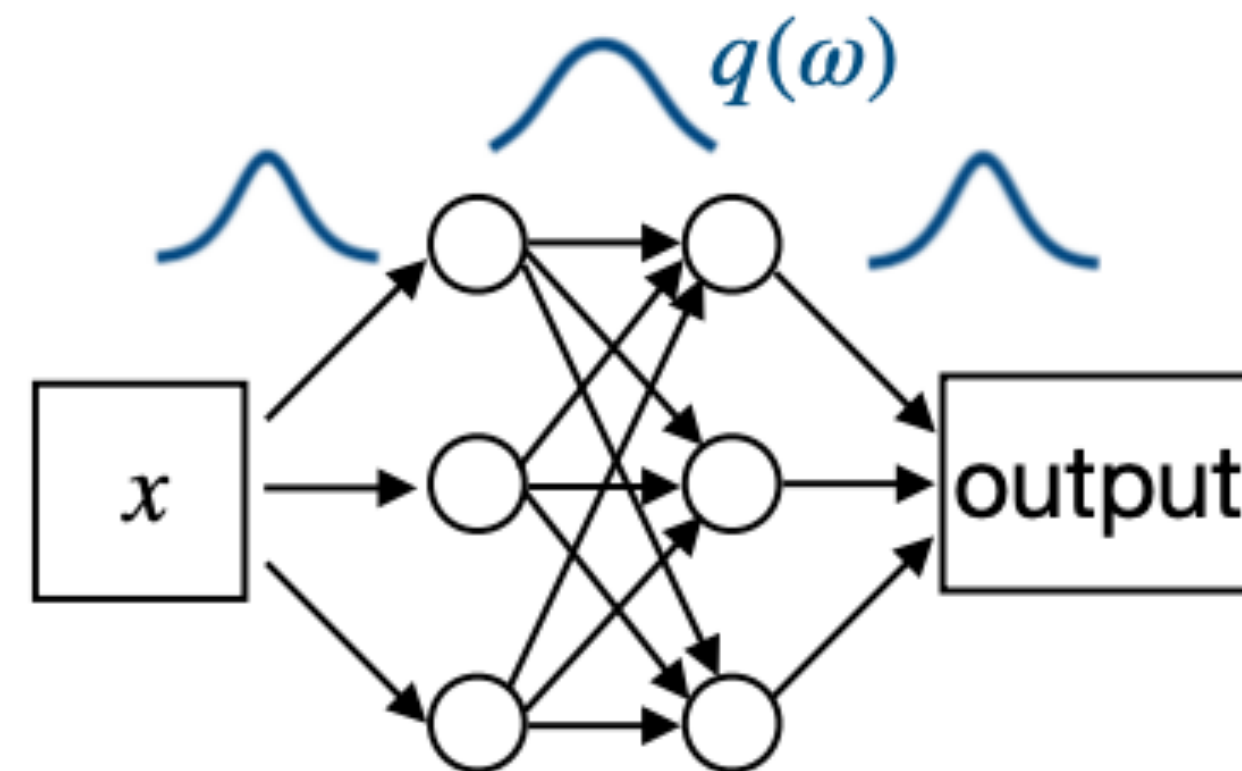
Estimating uncertainties on amplitude predictions

2. Estimating $\mathbf{p}(\mathbf{y} \mid \mathbf{x}, \mathbf{D})$ with training dataset D

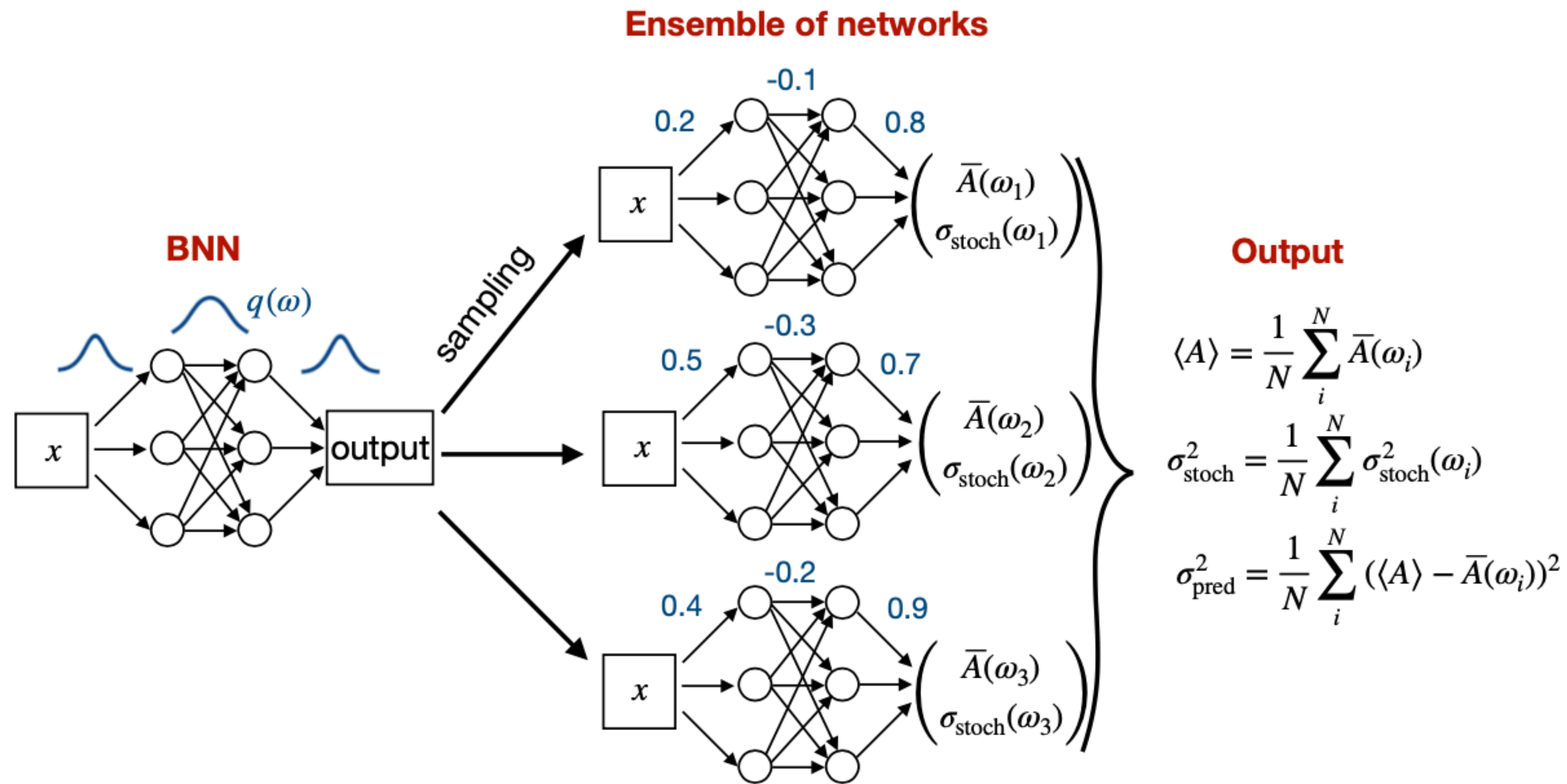
$$\cdot p(y \mid x, D) = \int dw p(y \mid x, w) p(w \mid D)$$

$\mathcal{L}_{\text{Gauss}}$

BNN



Bayesian Neural Network



$$\mathcal{L}_{\text{BNN}} = \int d\omega q(\omega) \sum_{\text{points } j} \left[\frac{\left| \bar{A}_j(\omega) - A_j^{(\text{truth})} \right|^2}{2\sigma_{\text{stoch},j}(\omega)^2} + \log \sigma_{\text{stoch},j}(\omega) \right] + \text{KL}[q(\omega), p(\omega)]$$

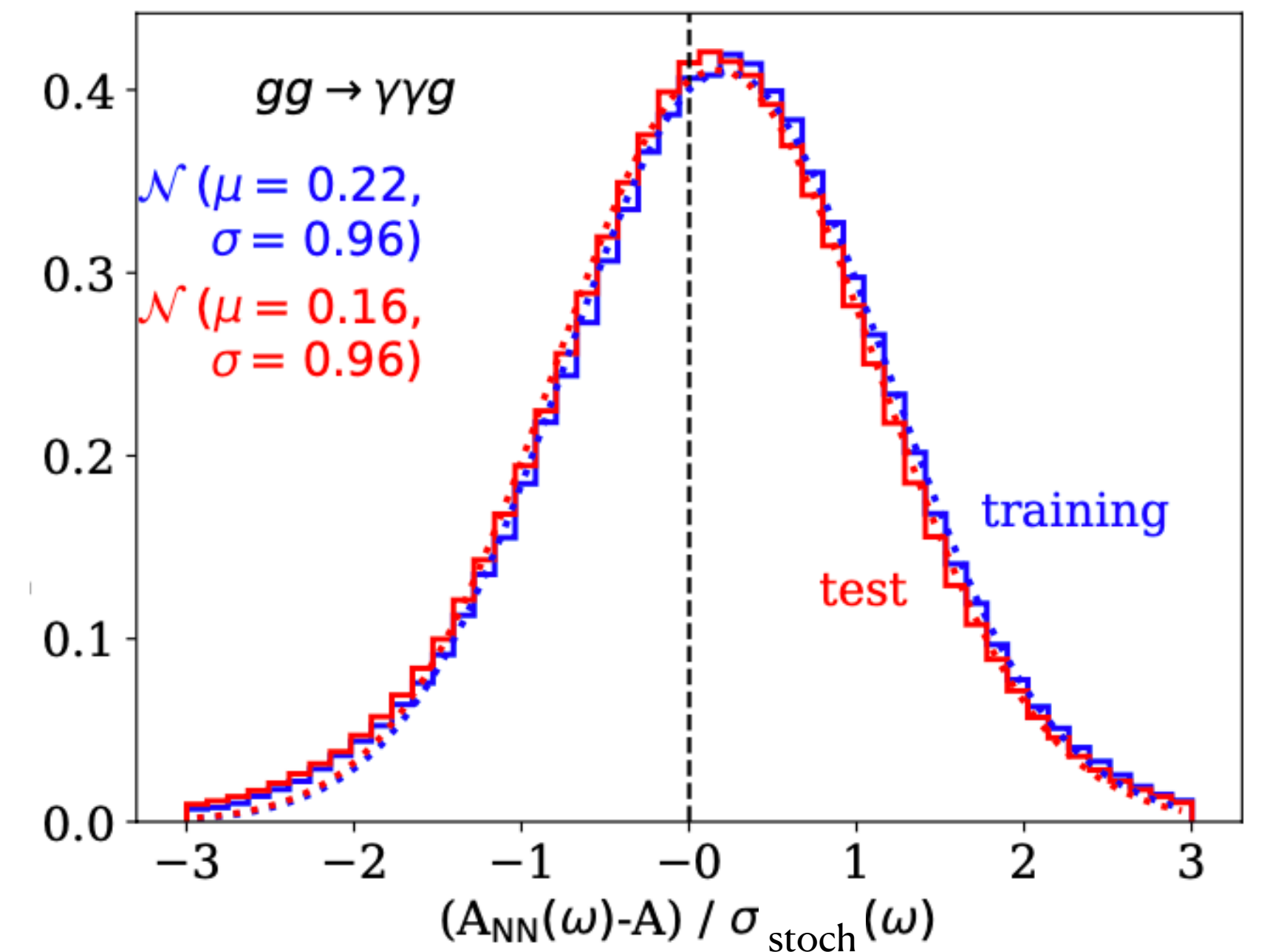
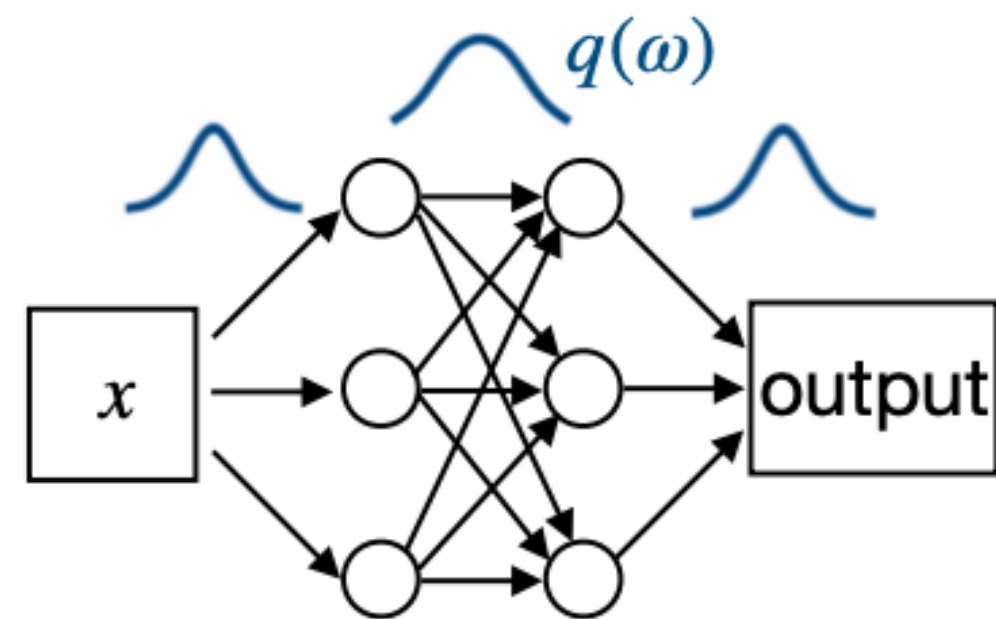
Results BNN

2. Estimating $\mathbf{p}(y | \mathbf{x}, \mathbf{D})$ with training dataset D

$$p(y | x, D) = \int dw p(y | x, w) p(w | D)$$

$\mathcal{L}_{\text{Gauss}}$

BNN

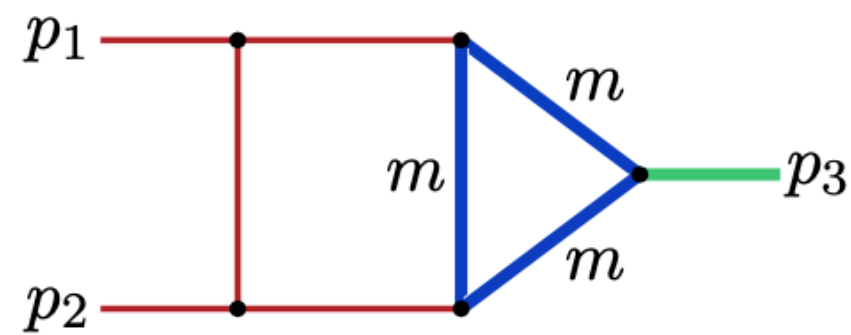


[preliminary]

$$\mathcal{L}_{\text{BNN}} = \int d\omega q(\omega) \sum_{\text{points } j} \left[\frac{|\bar{A}_j(\omega) - A_j^{(\text{truth})}|^2}{2\sigma_{\text{stoch},j}(\omega)^2} + \log \sigma_{\text{stoch},j}(\omega) \right] + \text{KL}[q(\omega), p(\omega)]$$

Multi-loop calculations with NNs

Precision predictions based on loop diagrams



Analytic expression for loop amplitude

$$G = \int_{-\infty}^{\infty} \left(\prod_{l=1}^L \frac{d^D k_l}{i\pi^{\frac{D}{2}}} \right) \prod_{j=1}^N \frac{1}{(q_j^2 - m_j^2 + i\delta)^{\nu_j}}$$

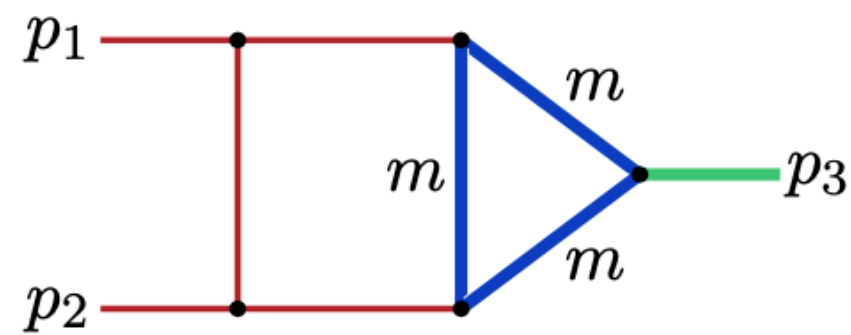
$$= \int_0^1 \prod_{j=1}^{N-1} dx_j x_j^{\nu_j-1} \frac{U^{\nu-(L+1)D/2}}{F^{\nu-LD/2}} = \int_0^1 \prod_{j=1}^{N-1} dx_j I(\vec{x})$$

Rewrite with
Feynman parameters

Still contains singularities

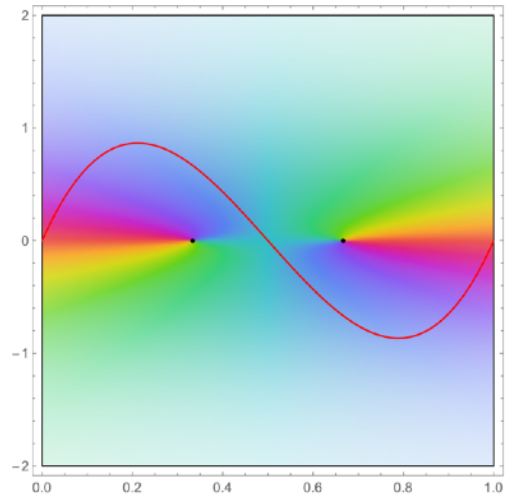
Multi-loop calculations with NNs

Precision predictions based on loop diagrams



Solved by contour deformation due to Cauchy's theorem

$$\int_0^1 \prod_{j=1}^N dx_j I(\vec{x}) = \int_0^1 \prod_{j=1}^N dx_j \det\left(\frac{\partial \vec{z}(\vec{x})}{\partial \vec{x}}\right) I(\vec{z}(\vec{x}))$$



Analytic expression for loop amplitude

$$G = \int_{-\infty}^{\infty} \left(\prod_{l=1}^L \frac{d^D k_l}{i\pi^{D/2}} \right) \prod_{j=1}^N \frac{1}{(q_j^2 - m_j^2 + i\delta)^{\nu_j}}$$

$$= \int_0^1 \prod_{j=1}^{N-1} dx_j x_j^{\nu_j-1} \frac{U^{\nu-(L+1)D/2}}{F^{\nu-LD/2}} = \int_0^1 \prod_{j=1}^{N-1} dx_j I(\vec{x})$$

Rewrite with Feynman parameters

Still contains singularities



Optimal parametrization = minimal variance

Turn it into an ML Problem

Integration with normalizing flows

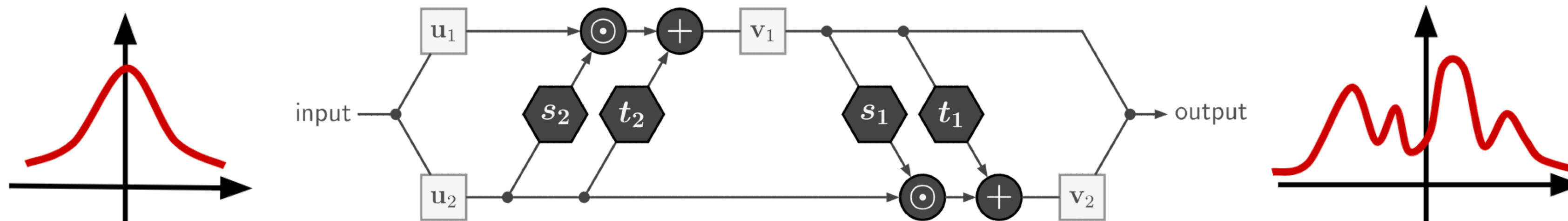
Numeric evaluation of integral $G = \int_0^1 dx_j \det\left(\frac{\partial \vec{z}(\vec{x})}{\partial \vec{x}}\right) I(\vec{z}(\vec{x}))$

Parametrization $\rightarrow z = \text{INN}(x)$

Minimize **variance** $\rightarrow \text{loss } \mathcal{L} = \sigma_n^2 = \frac{1}{n-1} \sum_{i=1}^n \left| \det\left(\frac{\partial \vec{z}(\vec{x}_{(i)})}{\partial \vec{x}_{(i)}}\right) I(\vec{z}(\vec{x}_{(i)})) - \langle I \rangle \right|^2$

Normalizing flow networks

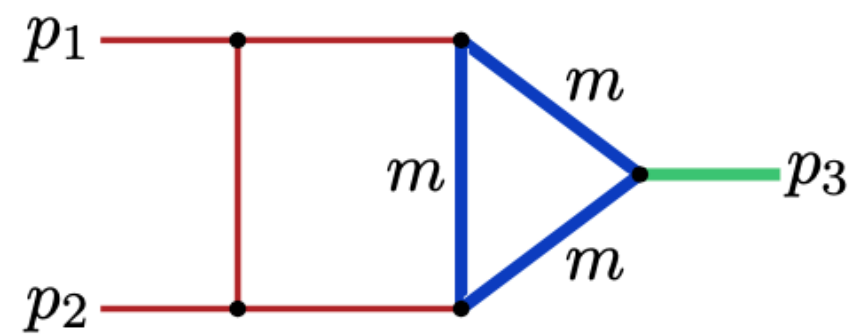
- + Bijective mapping
- + Tractable Jacobian
- + Combine many blocks



Multi-loop calculations with INN

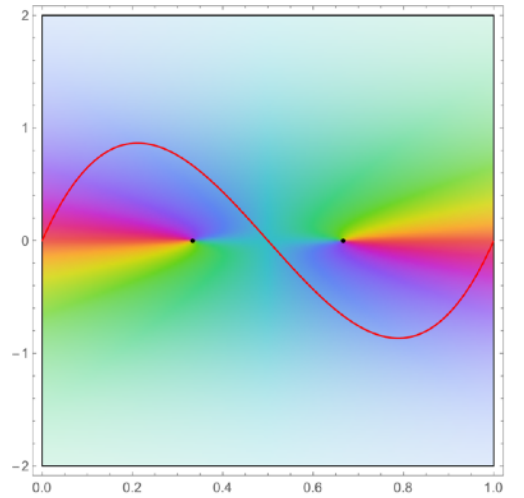
Profiting from the Jacobian

Precision predictions based on loop diagrams



Solved by contour deformation due to Cauchy's theorem

$$\int_0^1 \prod_{j=1}^N dx_j I(\vec{x}) = \int_0^1 \prod_{j=1}^N dx_j \det\left(\frac{\partial \vec{z}(\vec{x})}{\partial \vec{x}}\right) I(\vec{z}(\vec{x}))$$



Analytic expression for loop amplitude

$$G = \int_{-\infty}^{\infty} \left(\prod_{l=1}^L \frac{d^D k_l}{i\pi^{D/2}} \right) \prod_{j=1}^N \frac{1}{(q_j^2 - m_j^2 + i\delta)^{\nu_j}}$$

$$= \int_0^1 \prod_{j=1}^{N-1} dx_j x_j^{\nu_j-1} \frac{U^{\nu-(L+1)D/2}}{F^{\nu-LD/2}} = \int_0^1 \prod_{j=1}^{N-1} dx_j I(\vec{x})$$

Rewrite with Feynman parameters

Still contains singularities



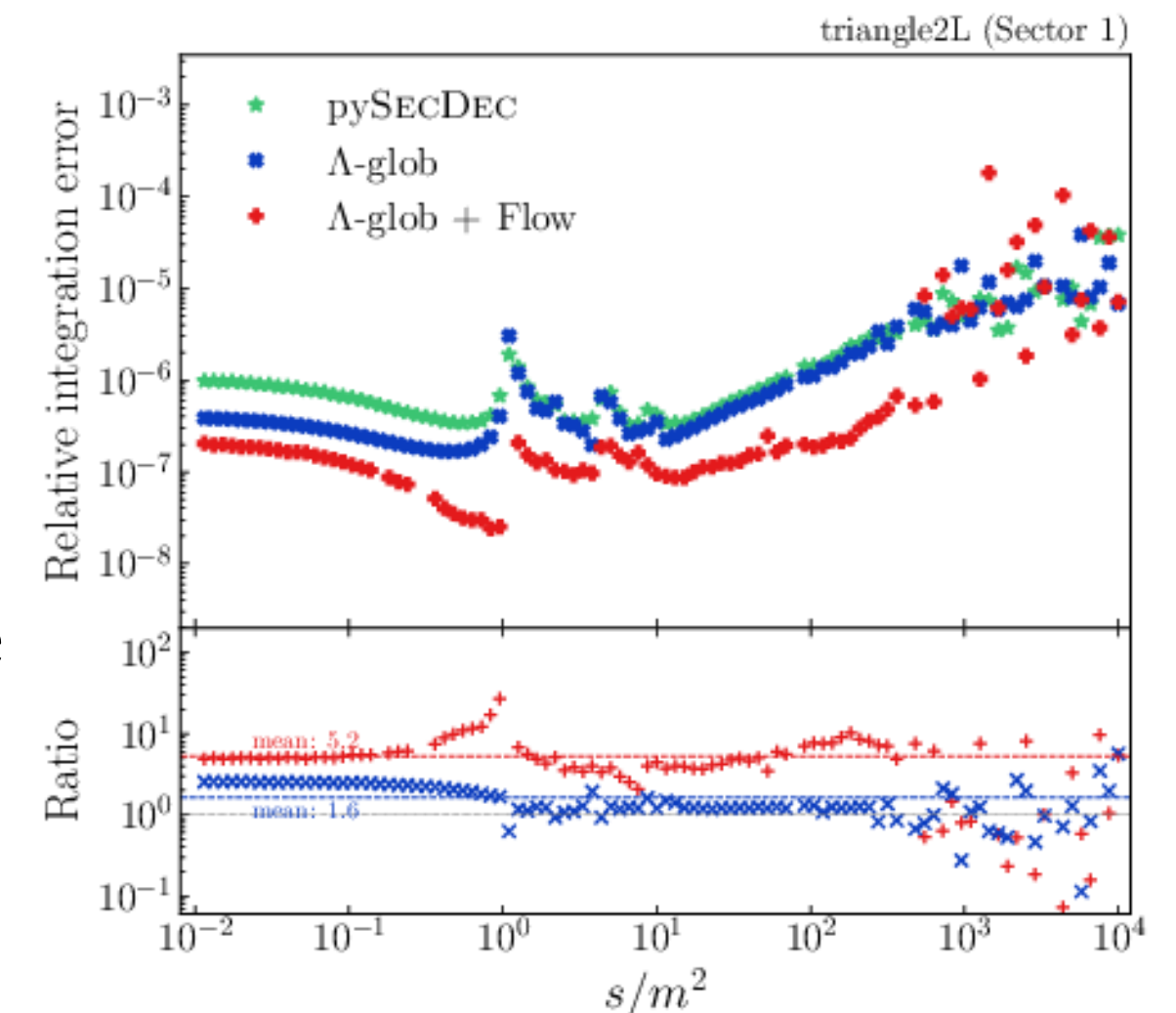
Optimal parametrization = minimal variance

Turn it into an ML Problem

Parametrization $\rightarrow z = \text{INN}(x)$

Variance $\rightarrow \mathcal{L}$

Better network \rightarrow smaller variance



Monte carlo event generation

1. Generate phase space points

→ set of four-momenta p_i

2. Calculate event weight

$$w_{\text{event}} = \underbrace{f(x_1, Q^2)f(x_2, Q^2)}_{\text{PDF}} \times \underbrace{\mathcal{M}(x_1, x_2, p_1, \dots, p_n)}_{\text{Matrix element †}} \times \underbrace{J(p_i(r))}_{\text{Phase space mapping}}$$

3. Unweighting †

keep events with $\frac{w_i}{w_{\text{max}}} > r \in [0,1]$

† Bottlenecks

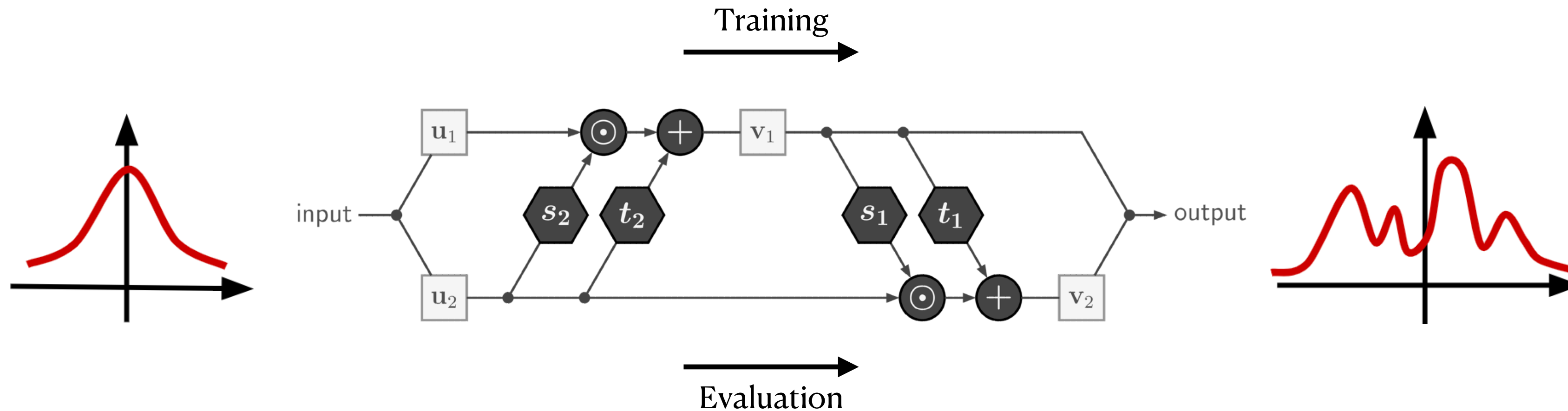
1. Slow **matrix element** calculation
 - ◆ Complexity grows exponentially with
 - # final state particles
 - Precision (LO, NLO, NNLO, ...)
2. Low **unweighting** efficiency
 - ◆ Discard most events if $w_i \ll w_{\text{max}}$
 - ◆ Optimize phase space mapping
 - ➔ $J(p_i(r)) = (f \times \mathcal{M})^{-1}$

Phase space sampling with generative networks (GAN, VAE, NF)

Normalizing flows

Invertible networks for complex transformations

- + Bijective mapping
- + Tractable Jacobian $\rightarrow p_x(x) = p_z(z) \cdot J_{NN}$
- + Fast evaluation in both direction



Training on density $t(x)$
 \rightarrow Minimize difference

$$\begin{aligned}\mathcal{L} &= \log p_x(x) / t(x) \\ &= \log p_z(z(x)) J_{NN} / t(x)\end{aligned}$$

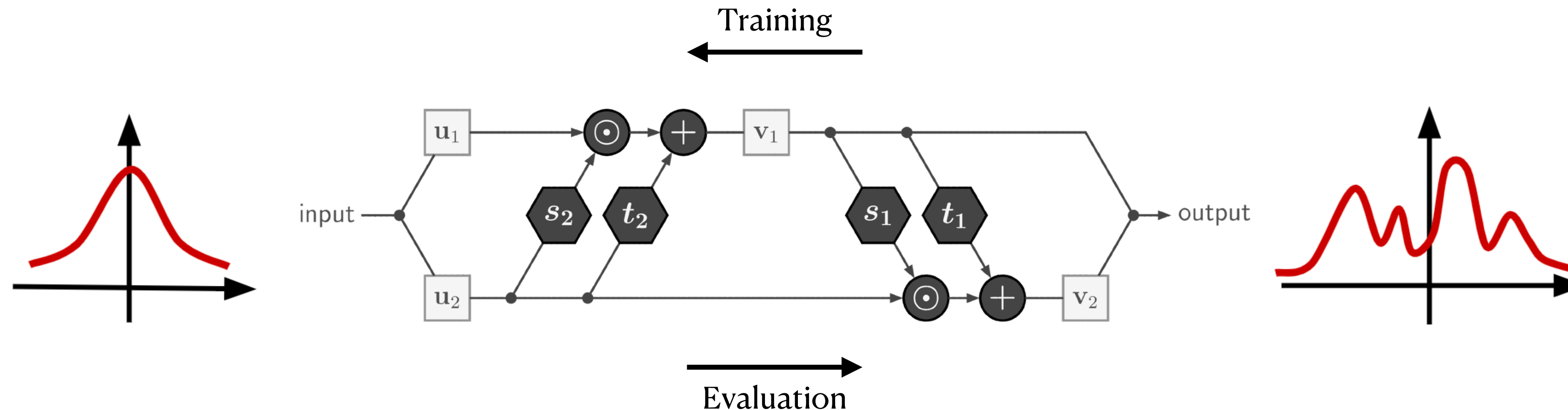
Training on samples x
 \rightarrow Maximize the log-likelihood

$$\begin{aligned}\mathcal{L} &= \log p(\theta | x) \\ &= \log p(z | \theta) + \log J_{NN} + p(\theta)\end{aligned}$$

Normalizing flows

Invertible networks for complex transformations

- + Bijective mapping
- + Tractable Jacobian $\rightarrow p_x(x) = p_z(z) \cdot J_{NN}$
- + Fast evaluation in both direction



Training on density $t(x)$
 \rightarrow Minimize difference

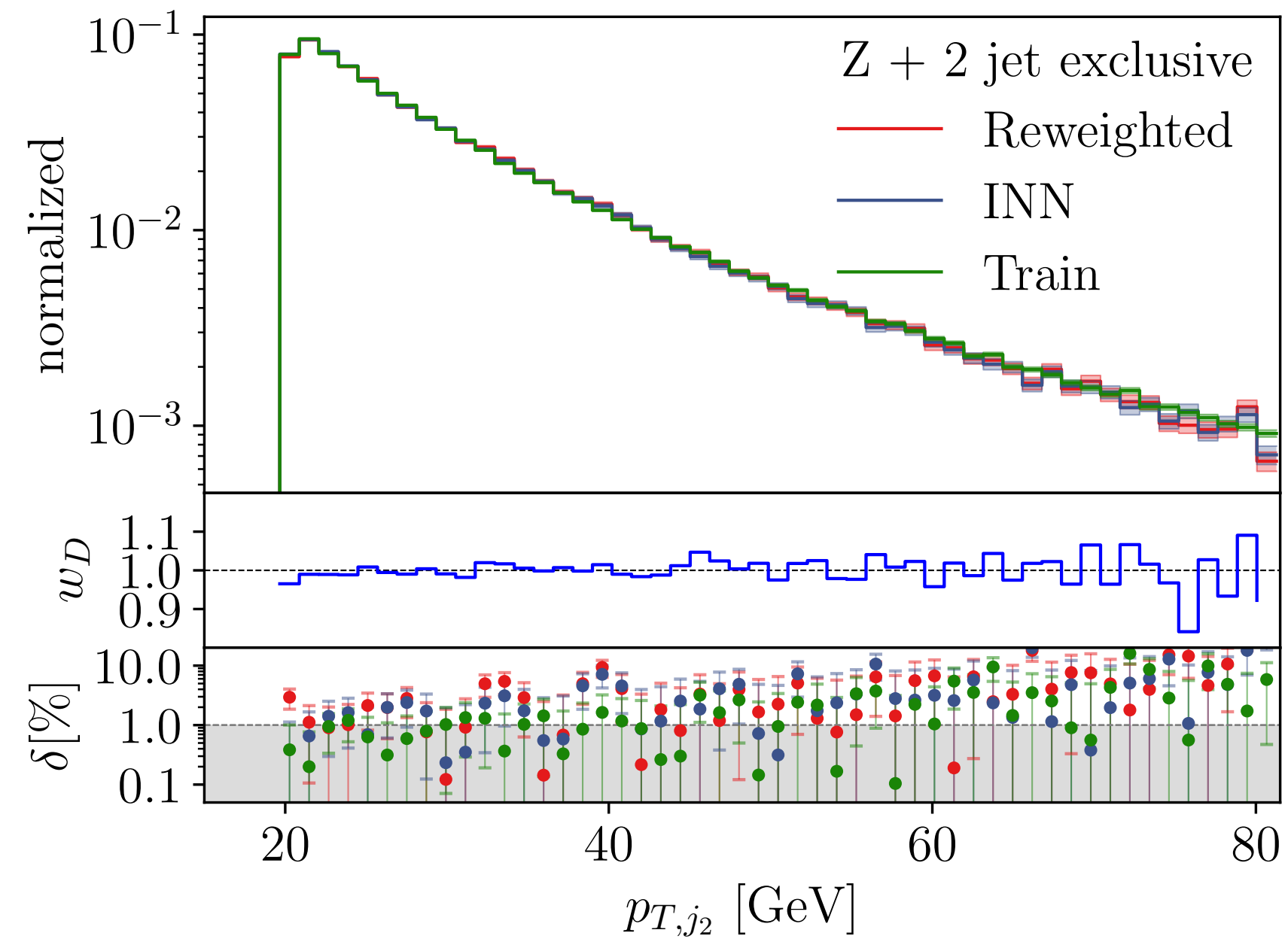
$$\begin{aligned}\mathcal{L} &= \log p_x(x)/t(x) \\ &= \log p_z(z(x)) J_{NN} / t(x)\end{aligned}$$

Training on samples x
 \rightarrow Maximize the log-likelihood

$$\begin{aligned}\mathcal{L} &= \log p(\theta | x) \\ &= \log p(z | \theta) + \log J_{NN} + p(\theta)\end{aligned}$$

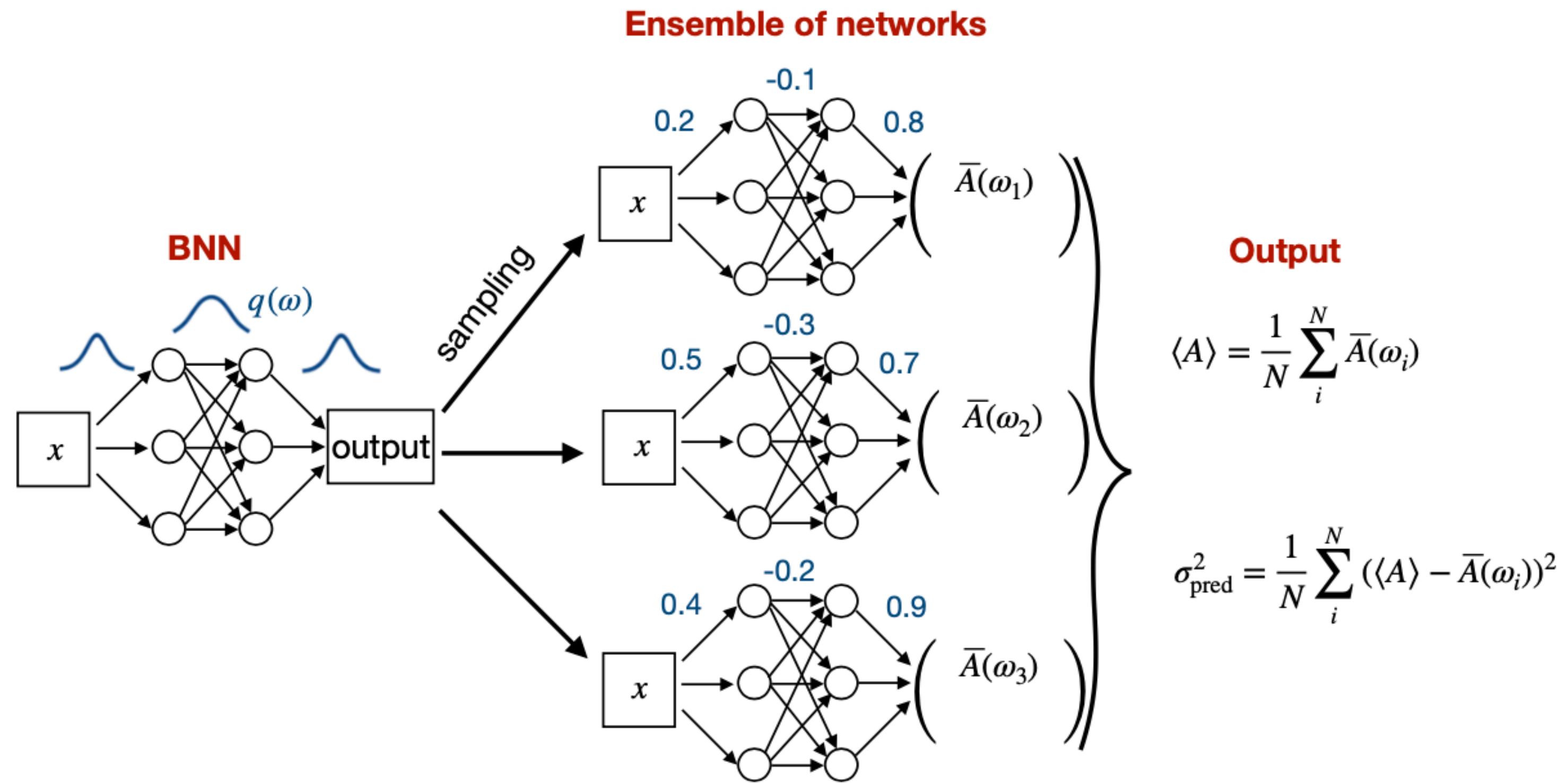
Putting flows to work

Event generation



- Train normalizing flow on 4-momenta
- Include symmetries in feature representation
- Excellent performance for direct output

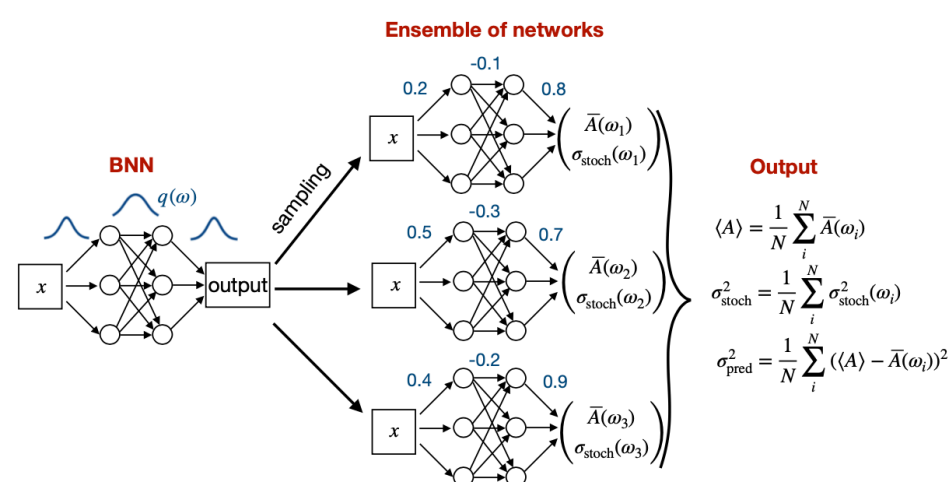
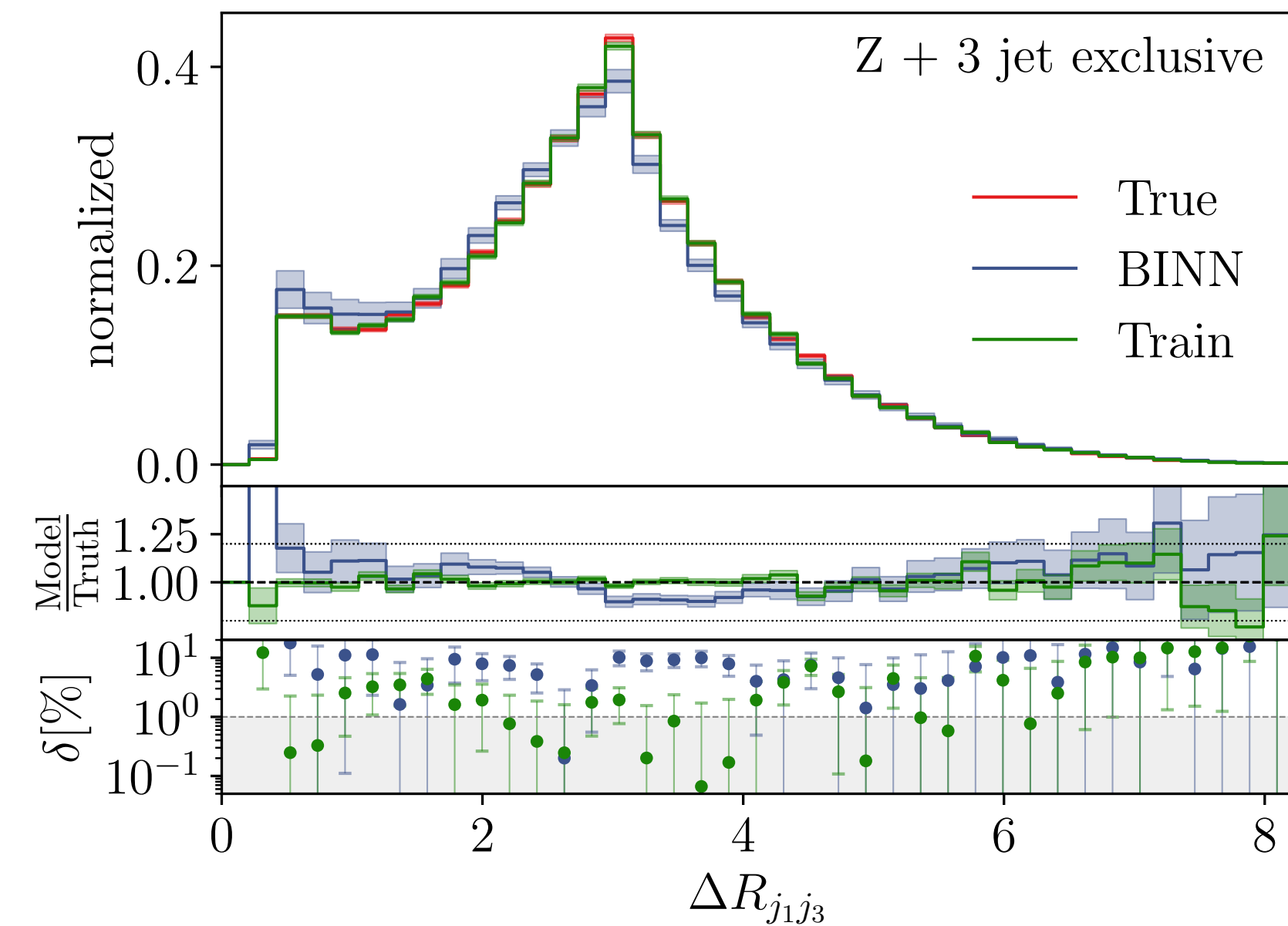
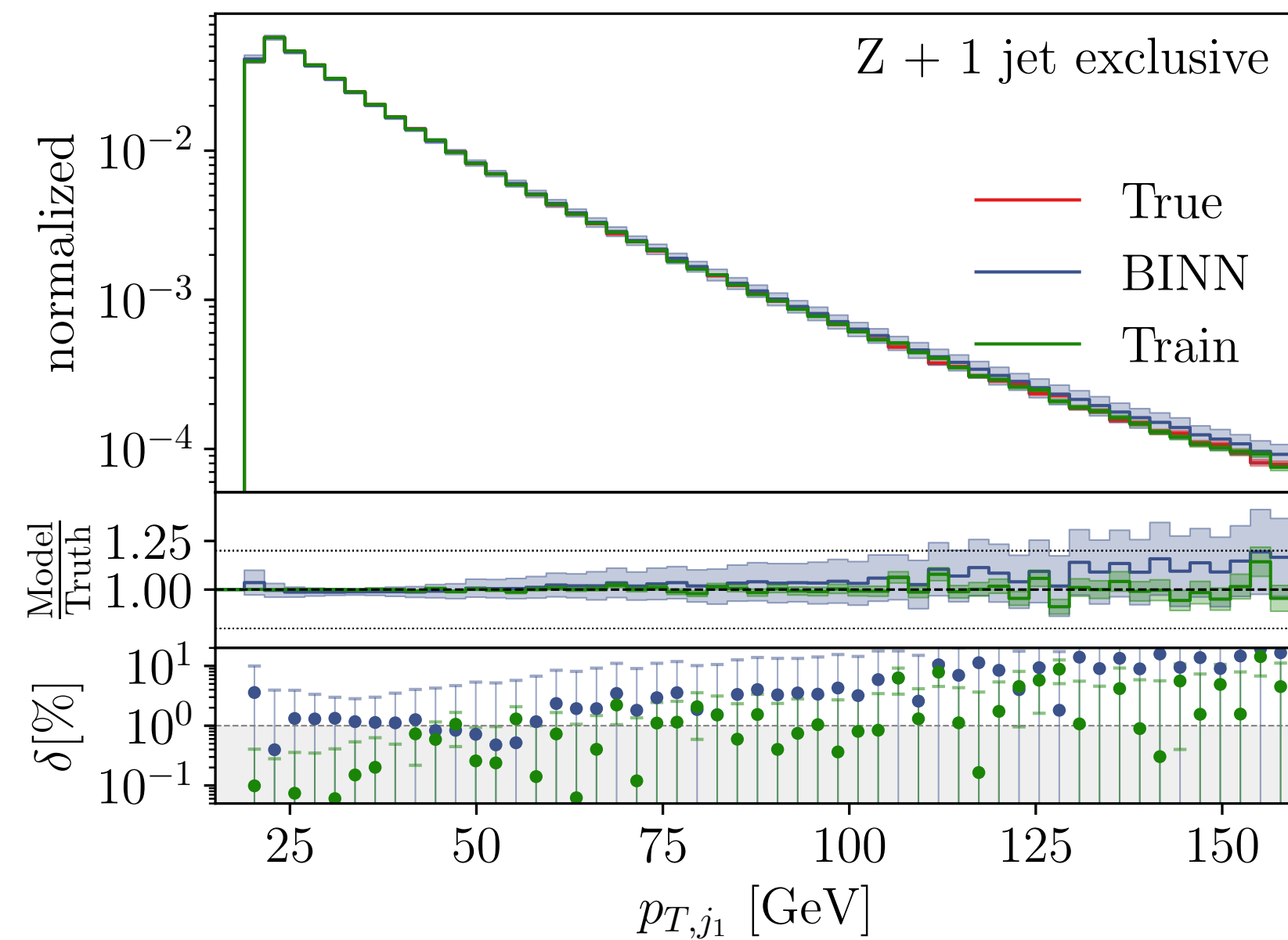
Bayesian Neural Network



$$\mathcal{L} = \mathcal{L}_{INN} + KL_{\text{prior}}$$

$$= \sum_{n=1}^N \langle \log p_X(x_n | \theta) \rangle_{\theta \sim q_{\Phi}(\theta)} - KL(q_{\Phi}(\theta), p(\theta))$$

Bayesian generative networks

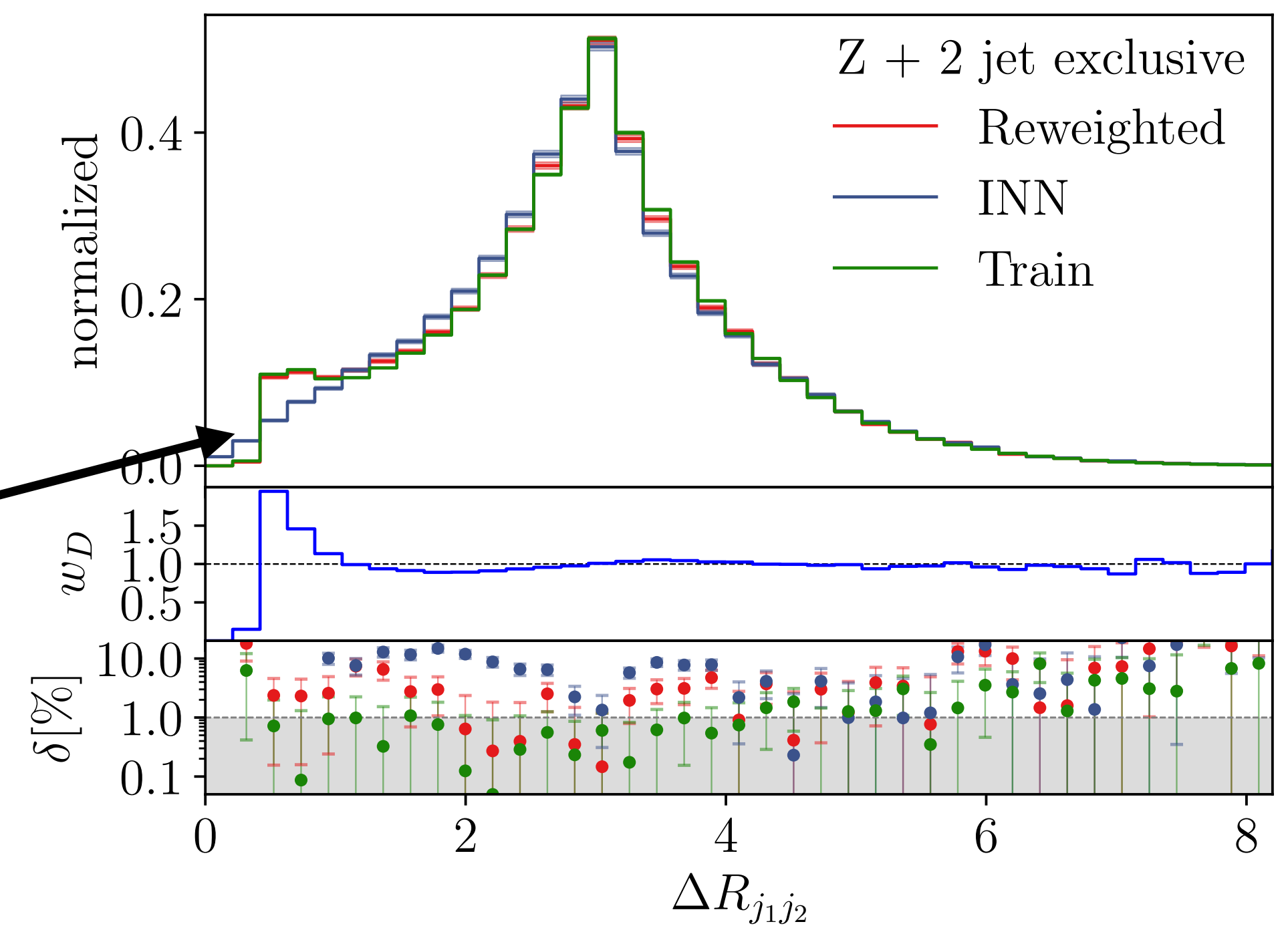
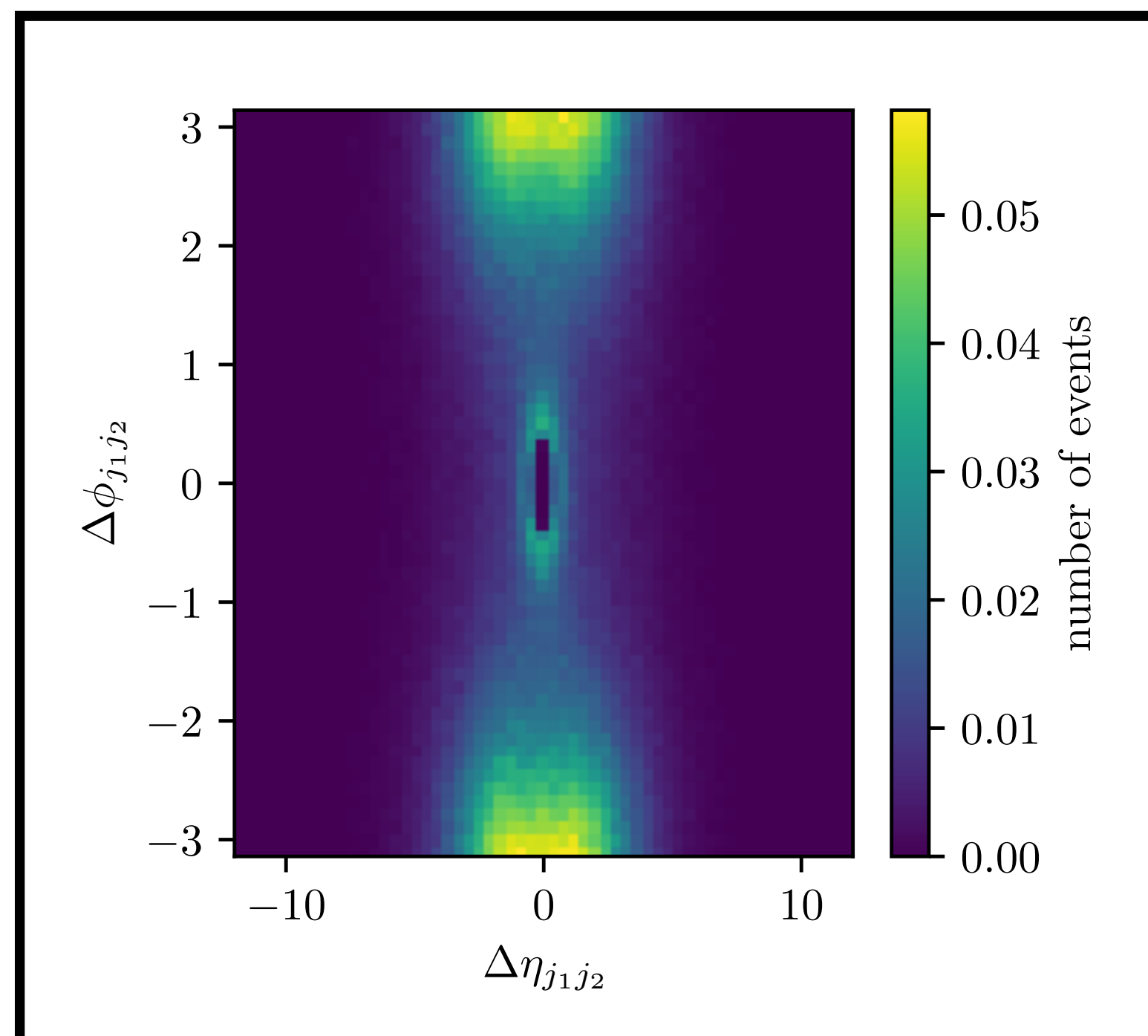


⇒ BINN captures uncertainty related to convergence and statistical uncertainties

⇒ BINN does not capture lack of expressiveness

Challenges for normalizing flows

- Narrow features
- Topological holes (eg ΔR cuts)
 - no bijective mapping possible
 - can only be approximated



Reweighting for Precision

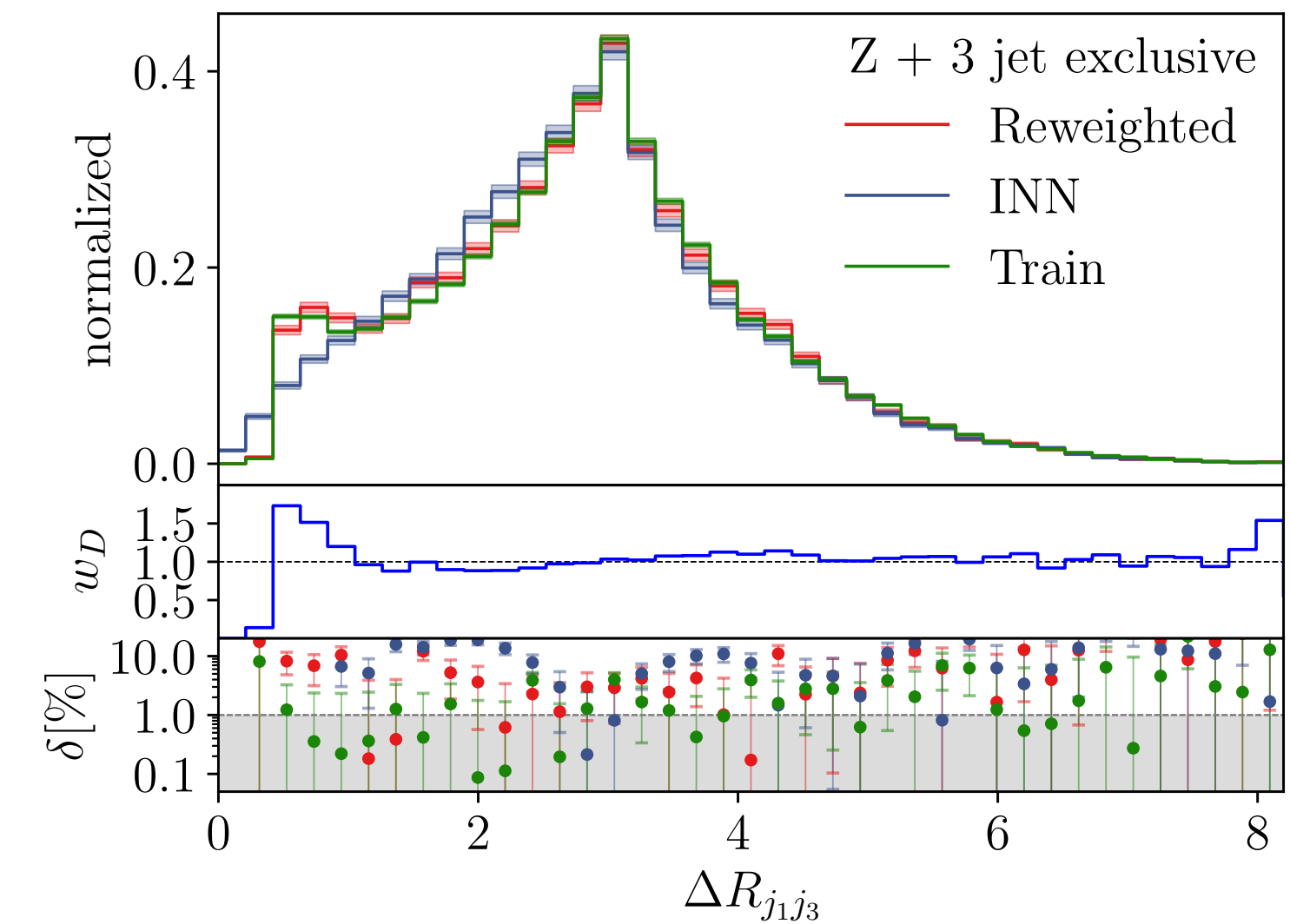
- Classifier loss

$$\begin{aligned}\mathcal{L} &= - \sum_{x \sim p_{data}} \log(D(x)) - \sum_{x \sim p_{INN}} \log(1 - D(x)) \\ &= - \int dx p_{data}(x) \log(D(x)) + p_{INN}(x) \log(1 - D(x))\end{aligned}$$

- Upon convergence obtain **reweighting factor**

$$\Rightarrow \frac{p_{data}(x)}{p_{INN}(x)} = \frac{D(x)}{1 - D(x)} = w_D$$

- Improve precision through reweighting
- Quantifies deviation



Reweighting for Precision

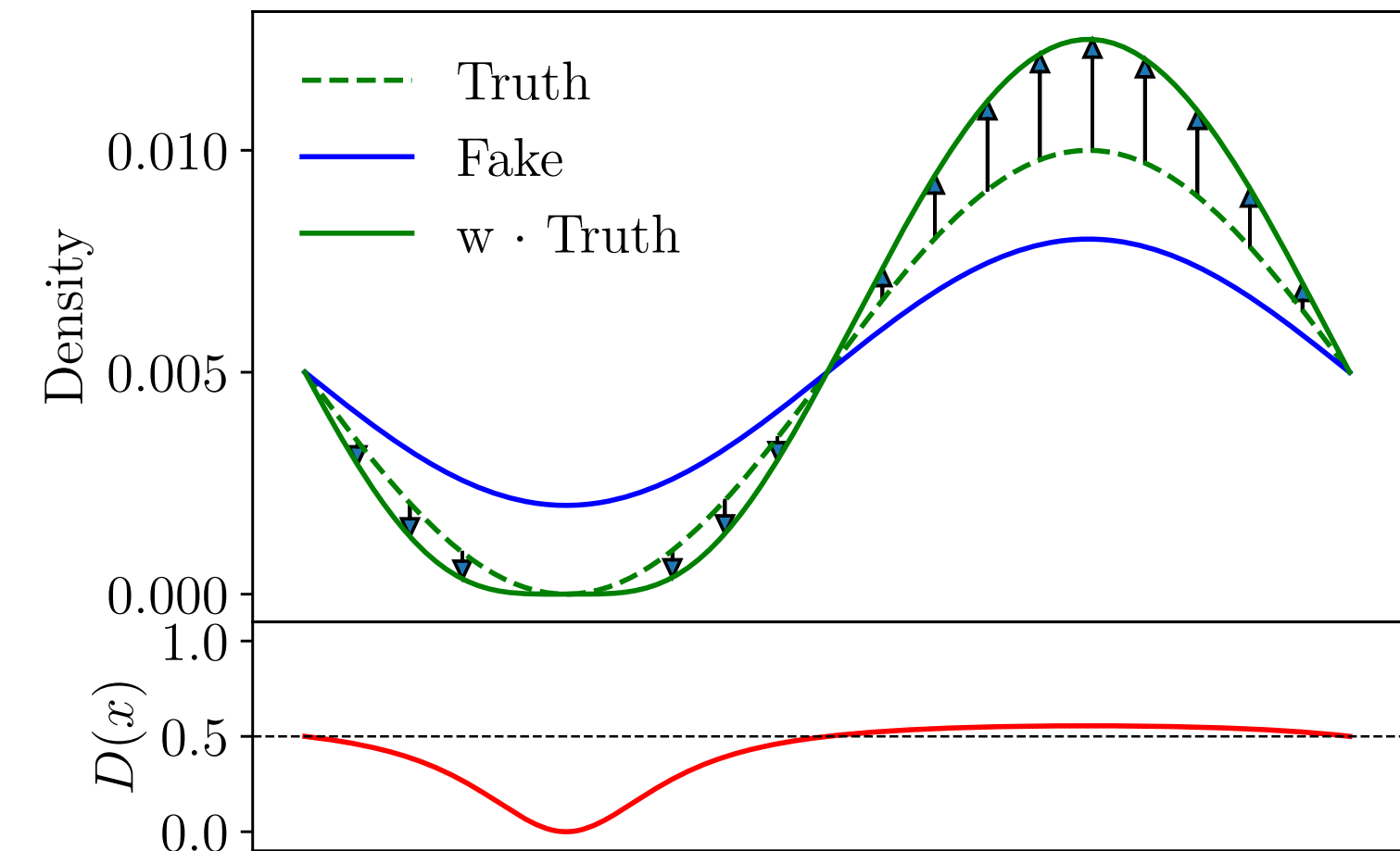
- Classifier loss

$$\begin{aligned}\mathcal{L} &= - \sum_{x \sim p_{data}} \log(D(x)) - \sum_{x \sim p_{INN}} \log(1 - D(x)) \\ &= - \int dx p_{data}(x) \log(D(x)) + p_{INN}(x) \log(1 - D(x))\end{aligned}$$

- Upon convergence obtain **reweighting factor**

$$\Rightarrow \frac{p_{data}(x)}{p_{INN}(x)} = \frac{D(x)}{1 - D(x)} = w_D$$

- Improve precision through reweighting
- Quantifies deviation



Reweighting for Precision

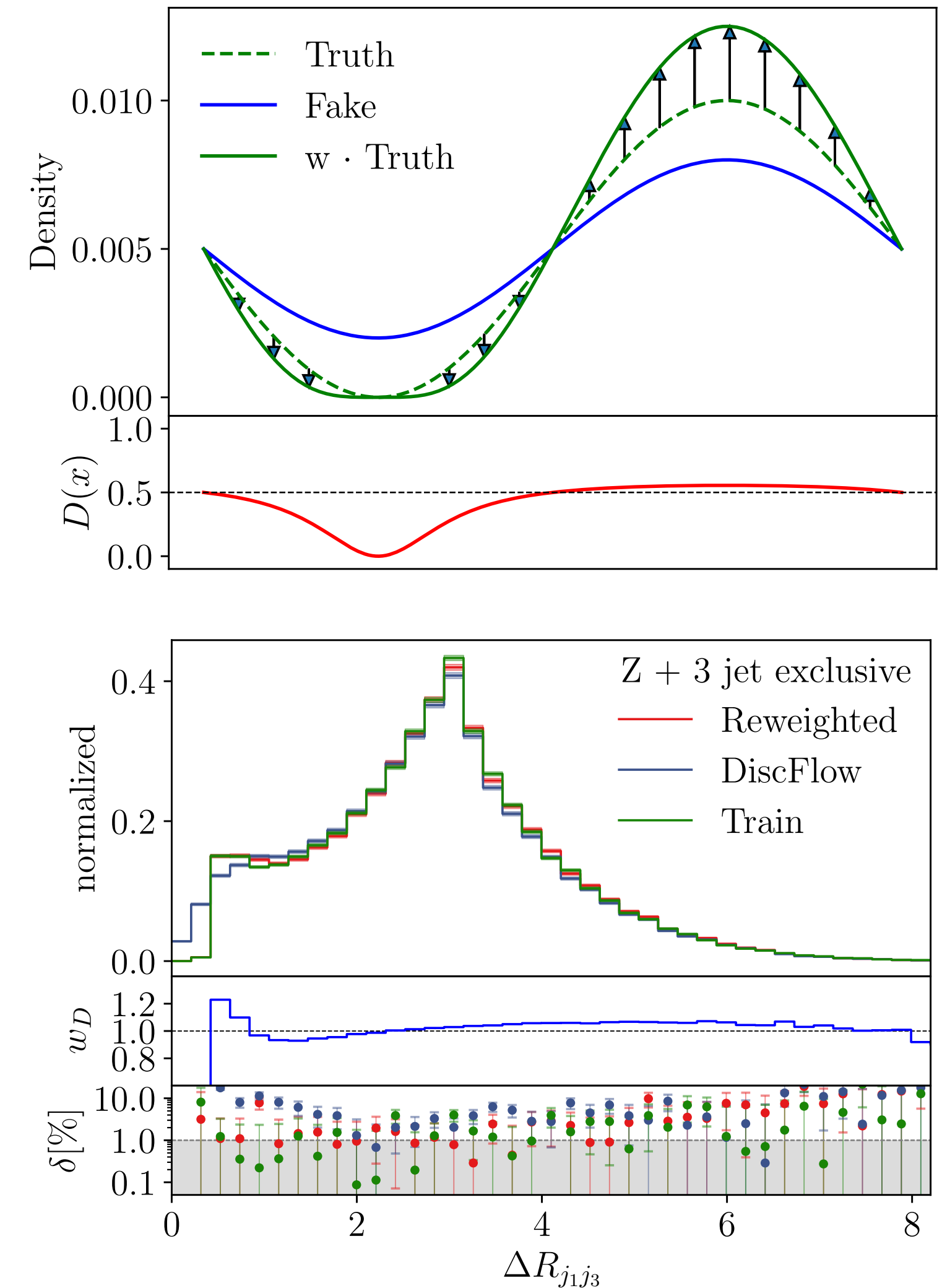
- Classifier loss

$$\begin{aligned} \mathcal{L} &= - \sum_{x \sim p_{data}} \log(D(x)) - \sum_{x \sim p_{INN}} \log(1 - D(x)) \\ &= - \int dx p_{data}(x) \log(D(x)) + p_{INN}(x) \log(1 - D(x)) \end{aligned}$$

- Upon convergence obtain **reweighting factor**

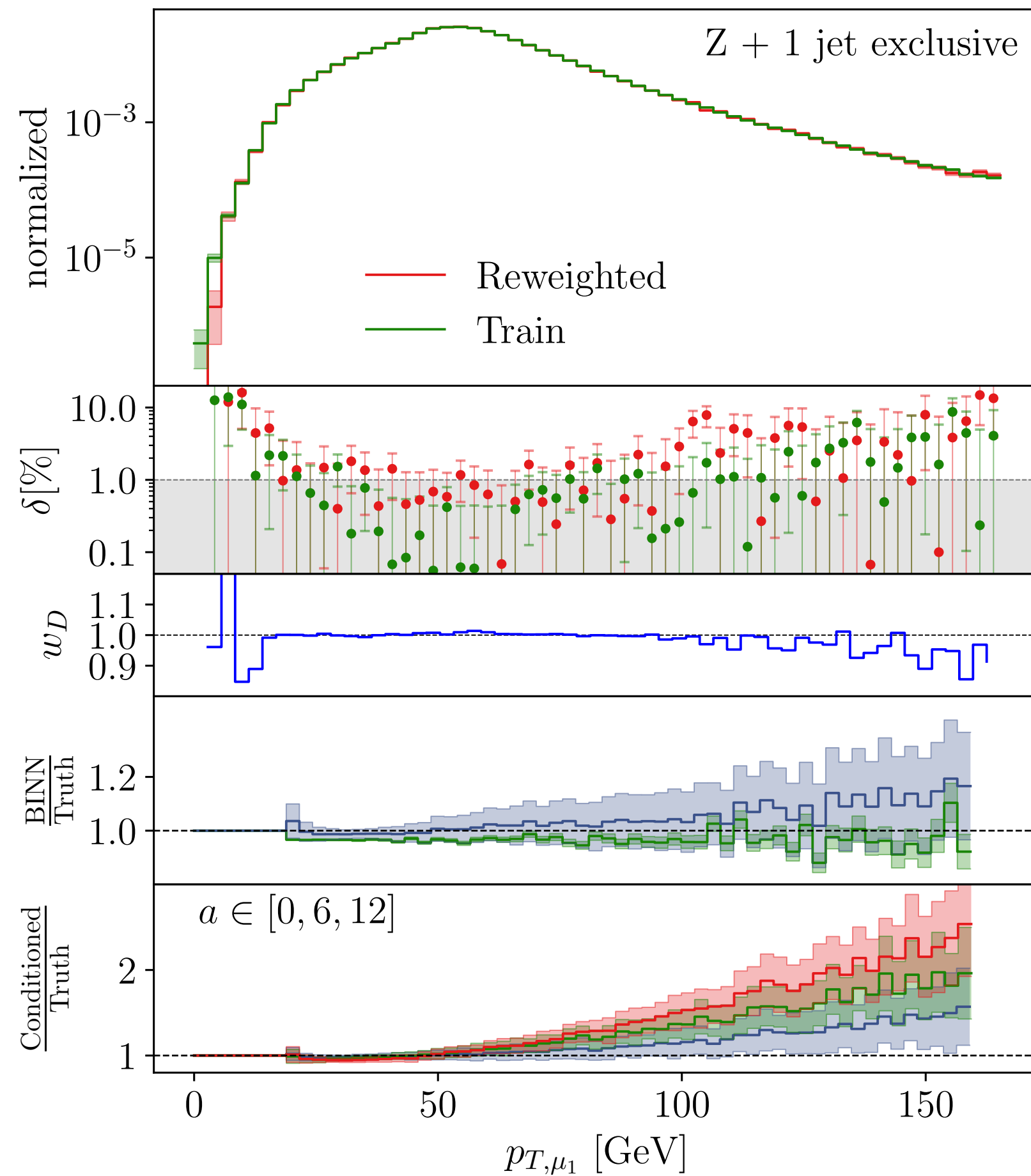
$$\Rightarrow \frac{p_{data}(x)}{p_{INN}(x)} = \frac{D(x)}{1 - D(x)} = w_D$$

- Improve precision through reweighting
- Quantifies deviation



Putting flows to work

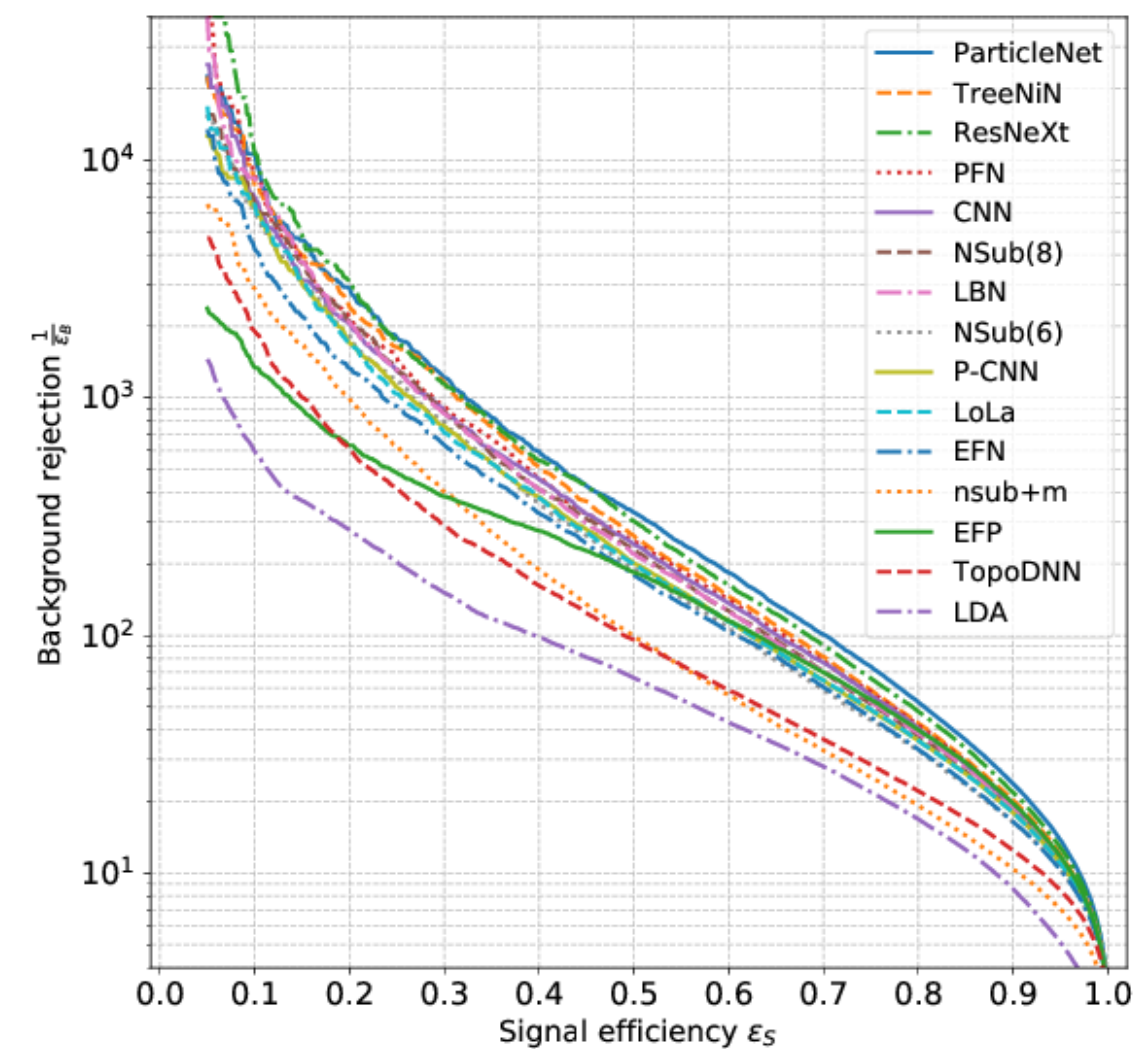
Event generation



- Basis: INN
 - Phase space symmetries in architecture
- Control via classifier D
 - $\frac{p_{\text{truth}}(x)}{p_{\text{INN}}(x)} = \frac{D(x)}{1 - D(x)}$
- Precision via reweighting
 - Correct deviations of p_{INN}
- ➔ Uncertainty estimation via Bayesian NN
- ➔ Uncertainty propagation via conditioning

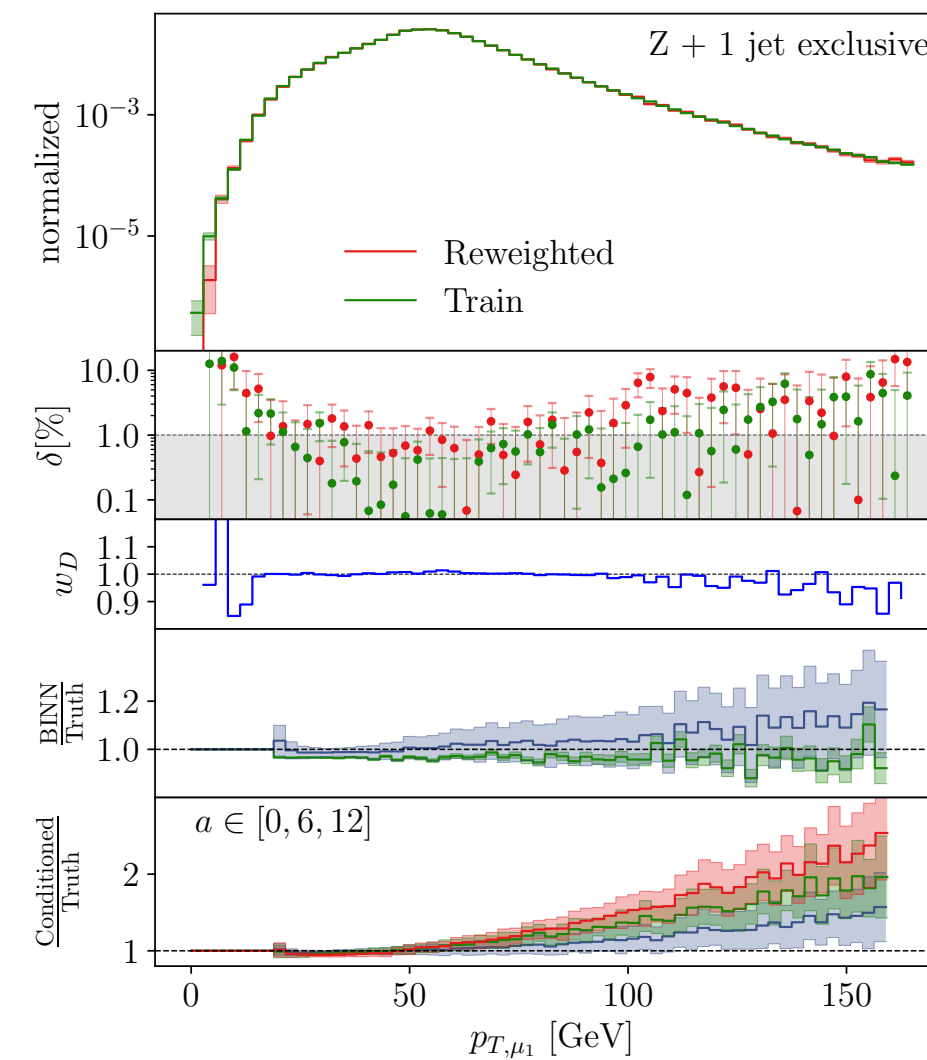
ML examples and their uncertainties

Classification



No uncertainty needed

Simulations



Amplitude estimation -> yes (BNN)

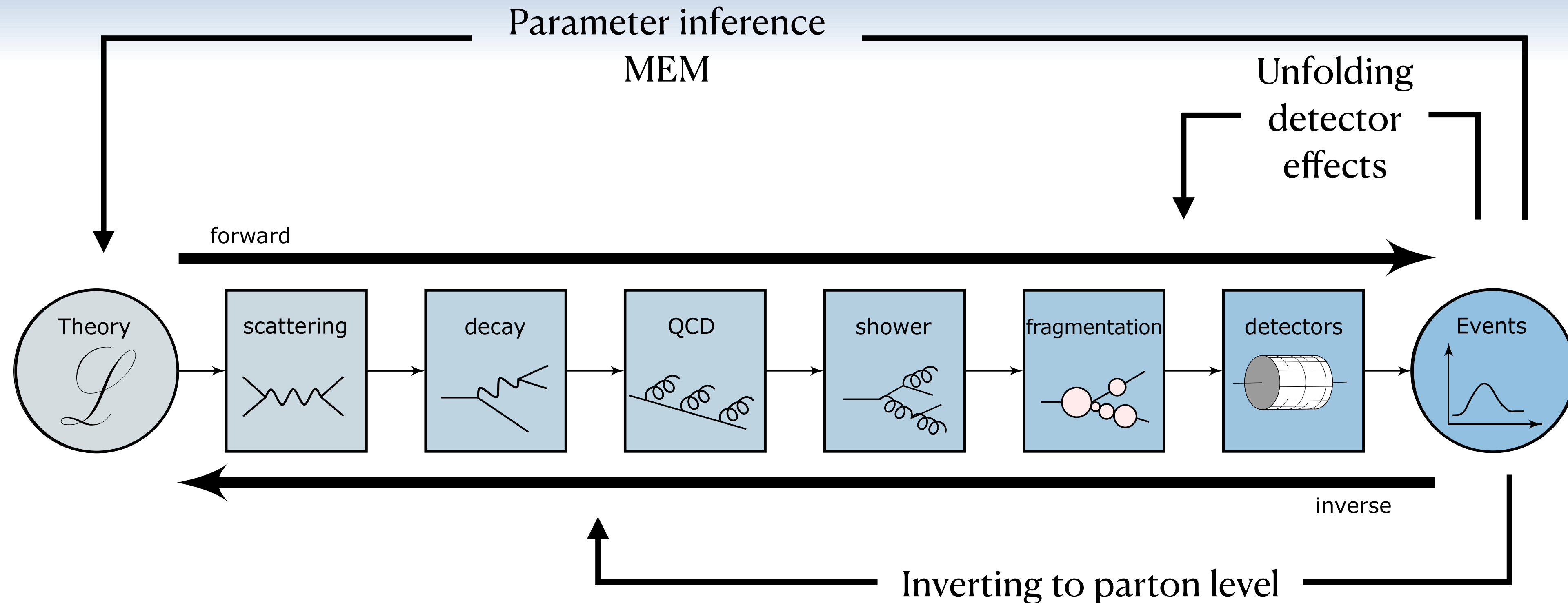
Loop integration -> no

Phase space sampling -> no

Data compression -> yes (BNN & classifier)

Unfolding

Inverting the simulation chain

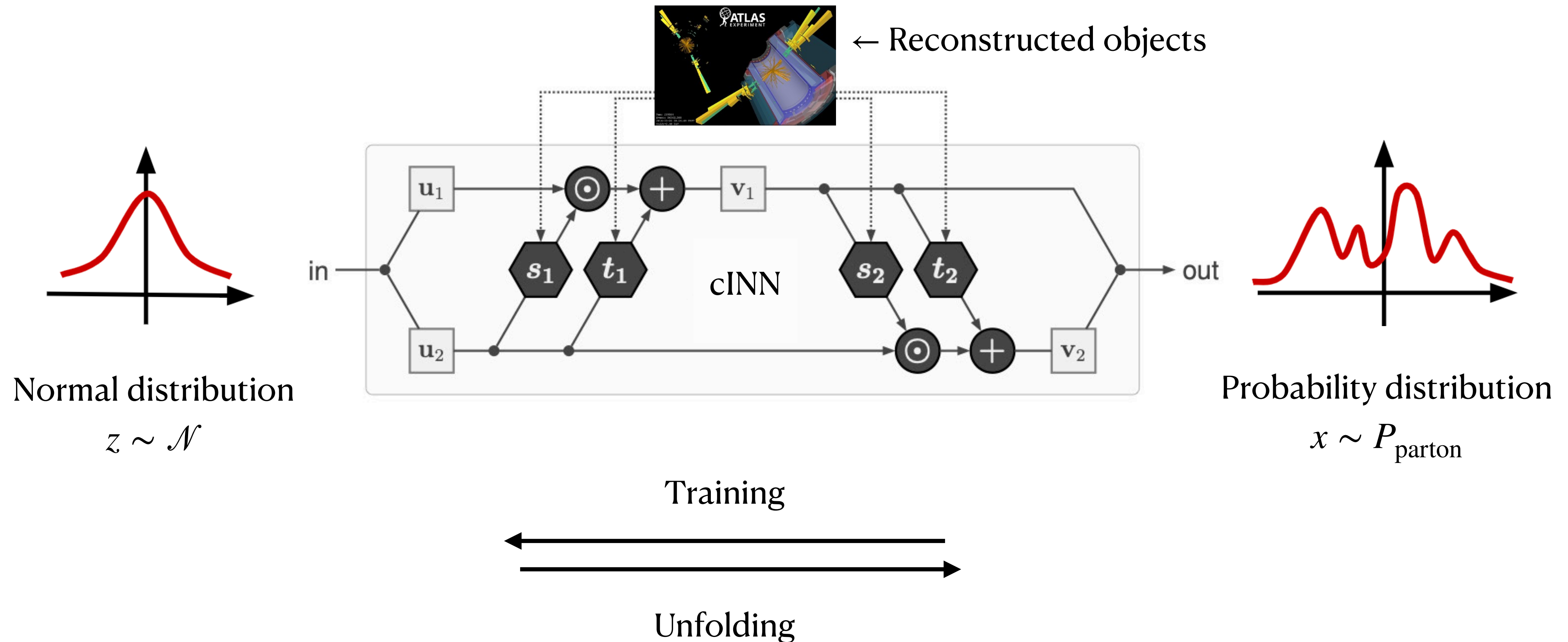


Requirements

- Highdimensional
- Bin independent
- Statistically well defined

cINN unfolding

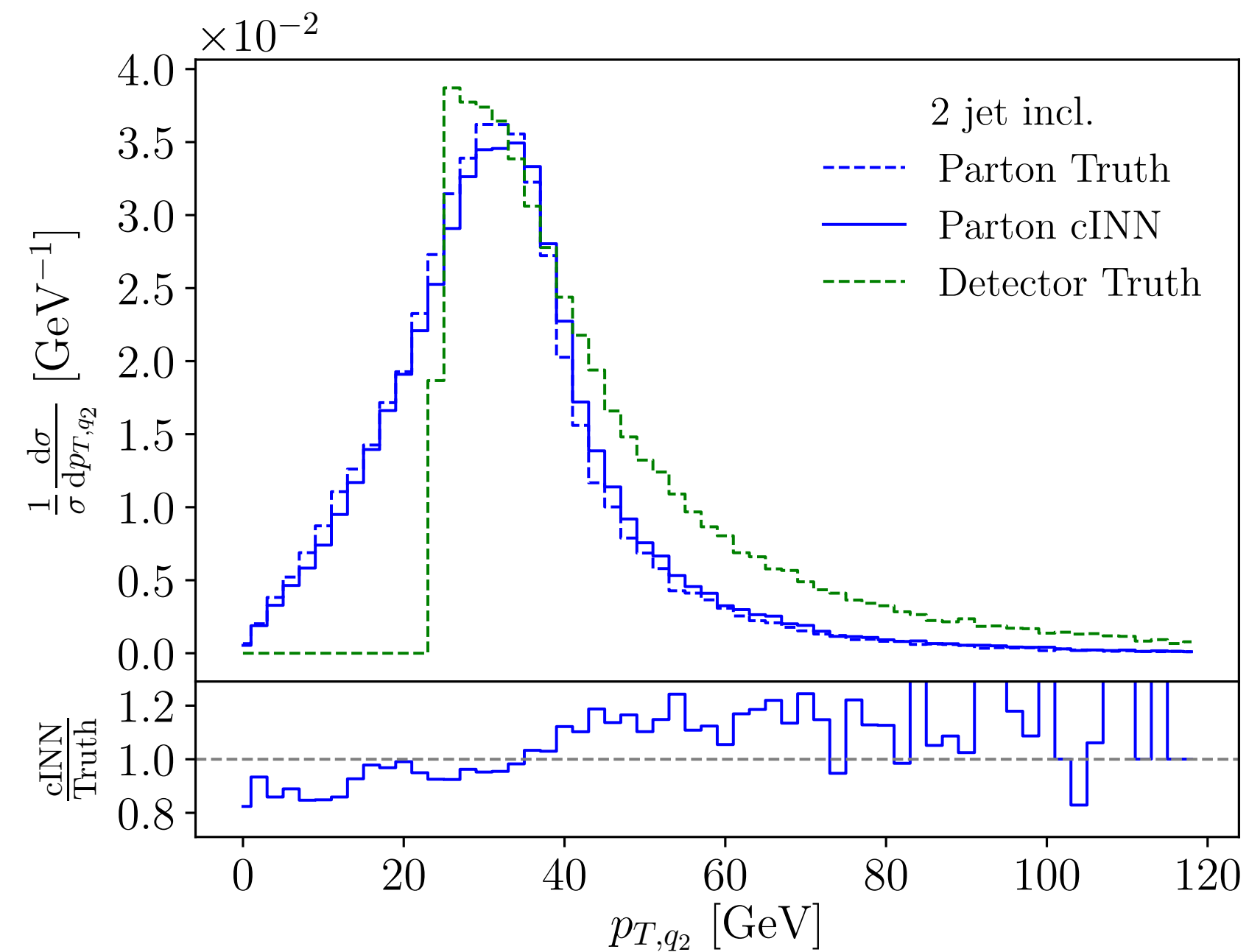
Given a reconstructed event:
What is the probability distribution at particle level?



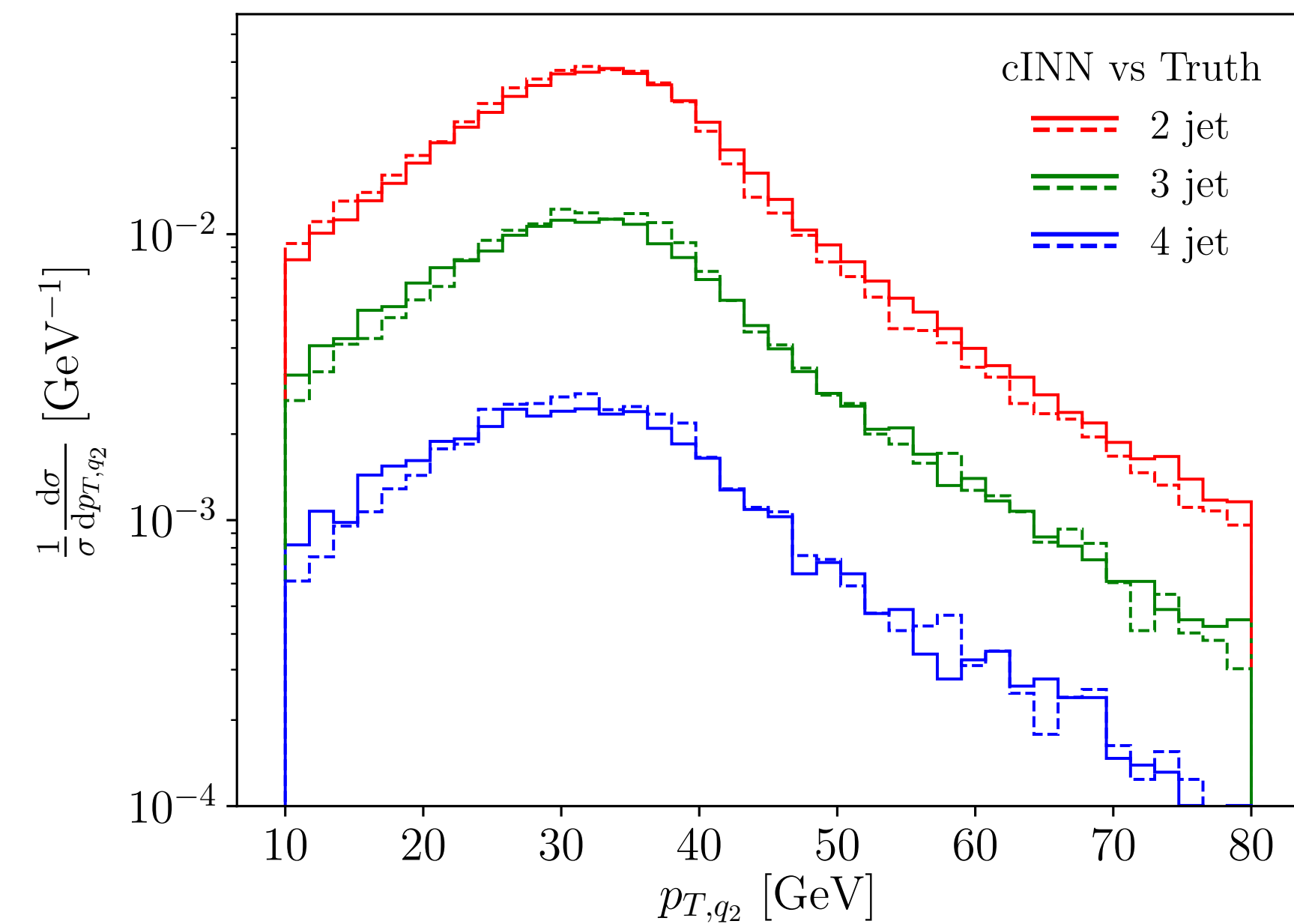
Inverting inclusive distributions

$$pp > WZ > q\bar{q}l^+l^- + \text{ISR} \rightarrow 2/3/4 \text{ jet events}$$

Training on inclusive dataset



Evaluate exclusive 2/3/4 jet events



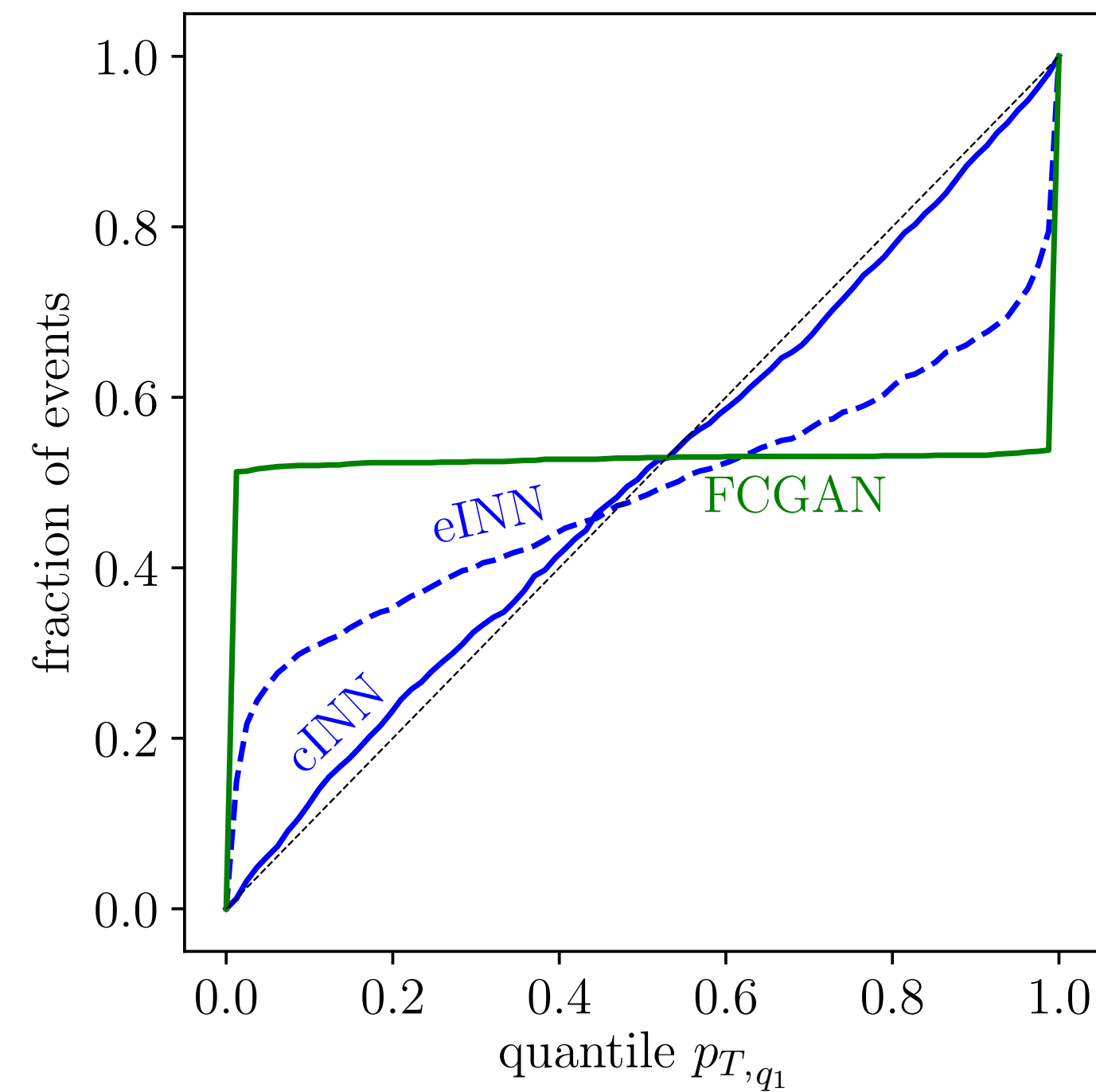
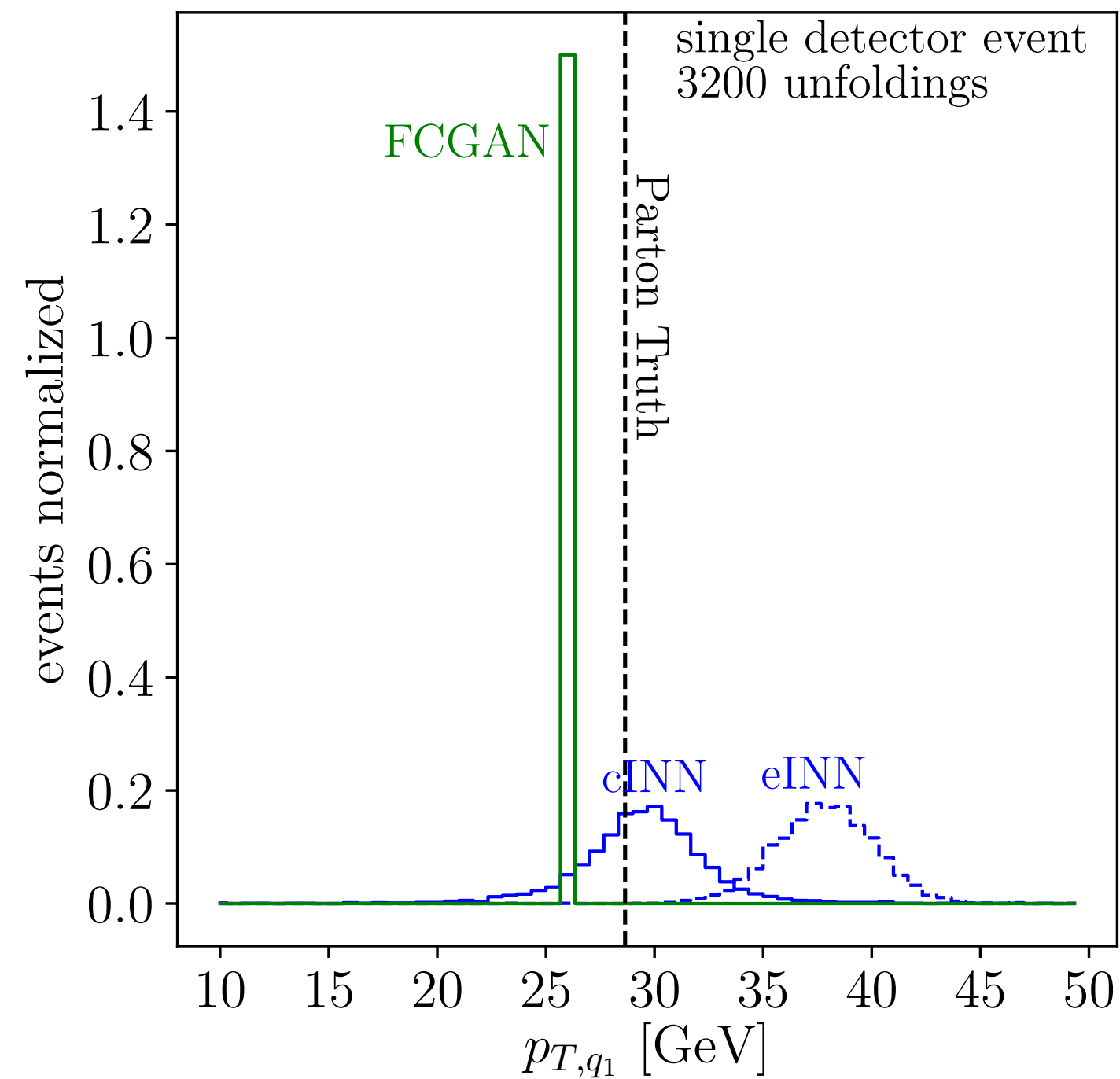
- High-dimensional
- Bin-independent
- Statistically well defined?

M. Bellagente et al. [2006.06685]

Event-wise unfolding

No deterministic mapping!

Check calibration of probability density for individual event unfolding

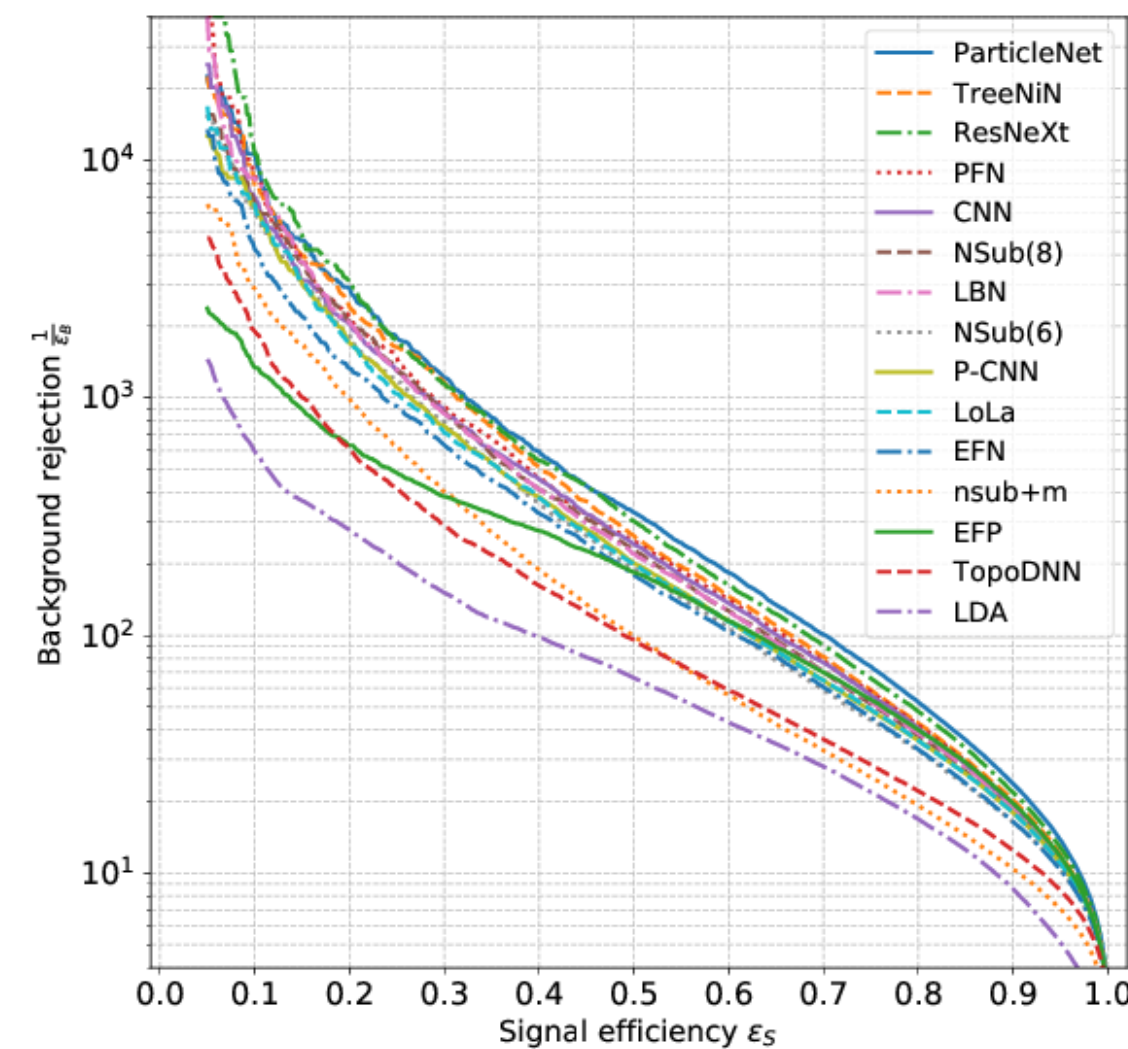


- High-dimensional
- Bin-independent
- Statistically well defined

M. Bellagente et al. [2006.06685]

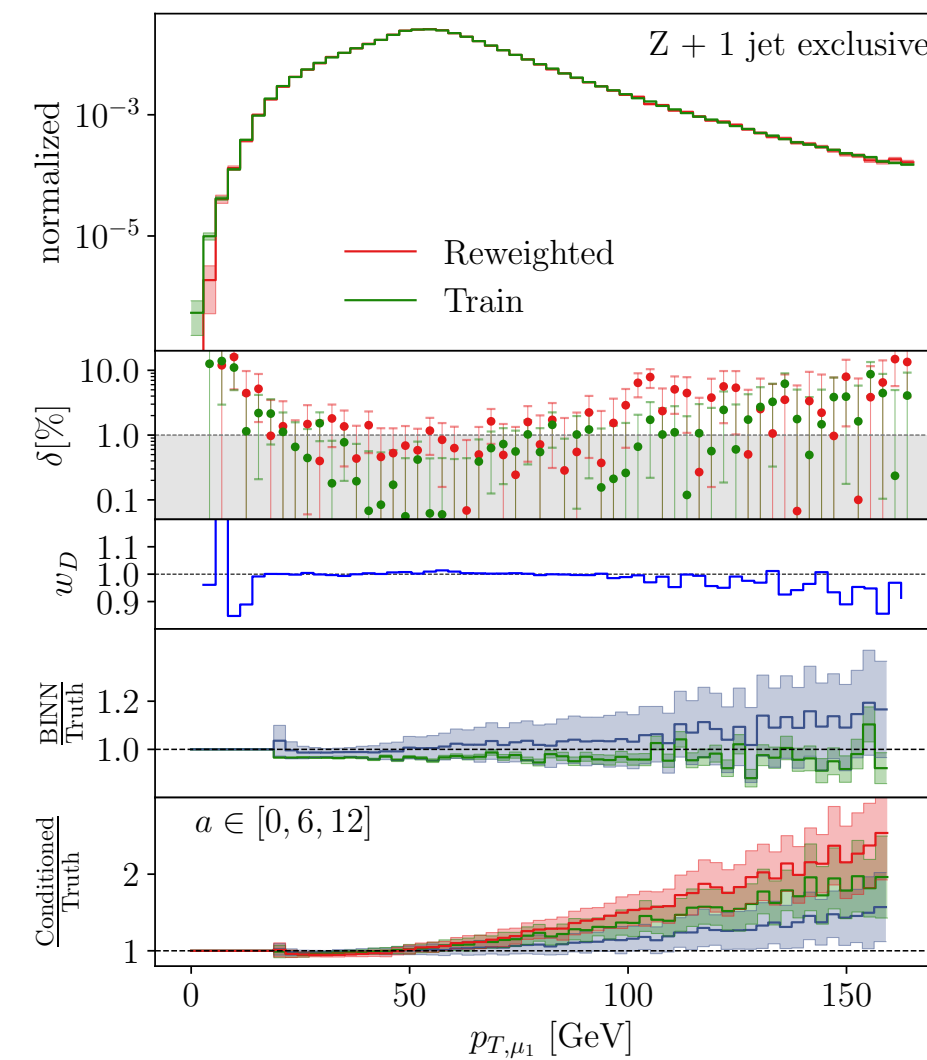
ML examples and their uncertainties

Classification



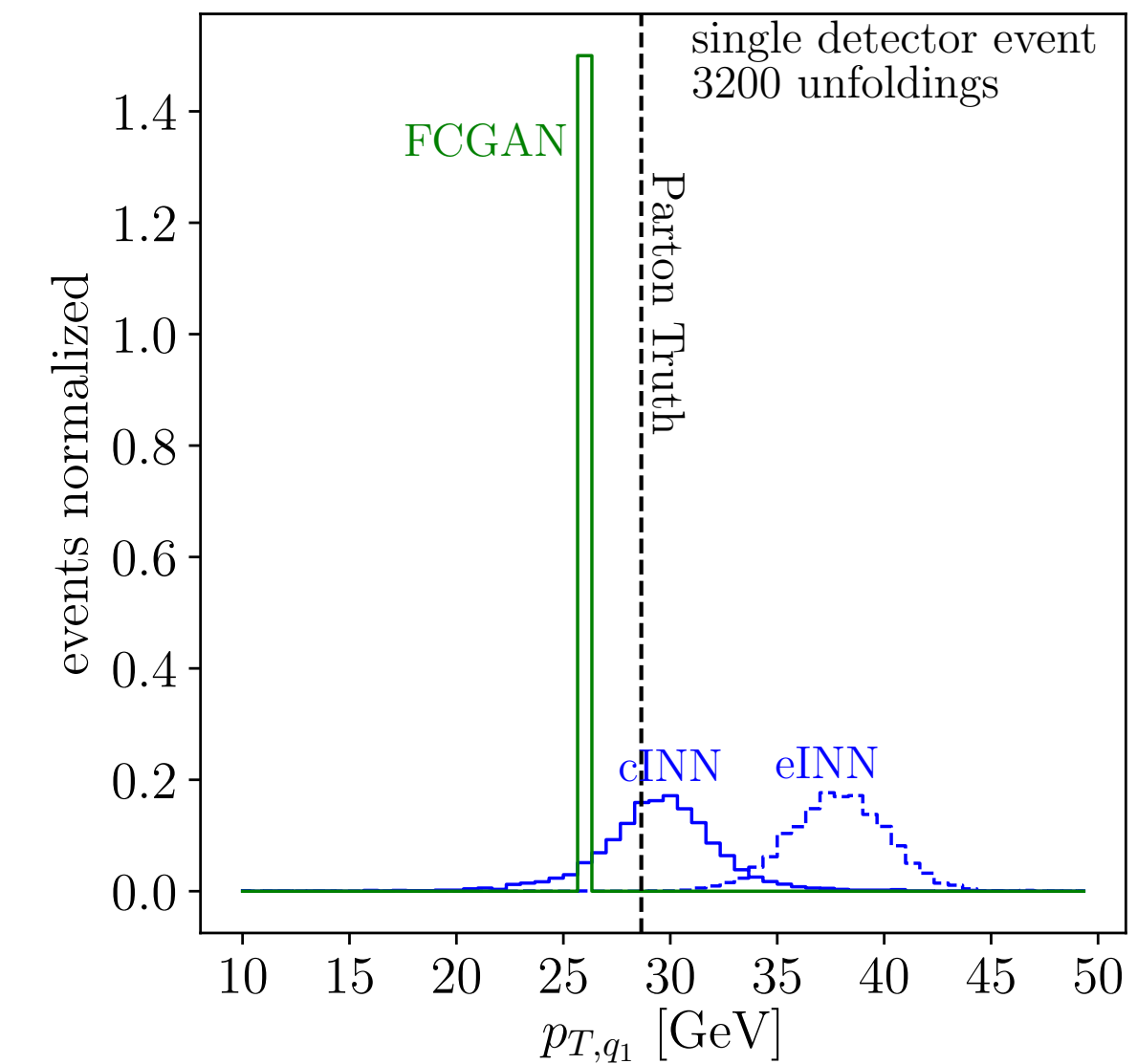
No uncertainty needed

Simulations



Amplitude estimation -> yes (BNN)
 Loop integration -> no
 Phase space sampling -> no
 Data compression -> yes (BNN & classifier)

Unfolding

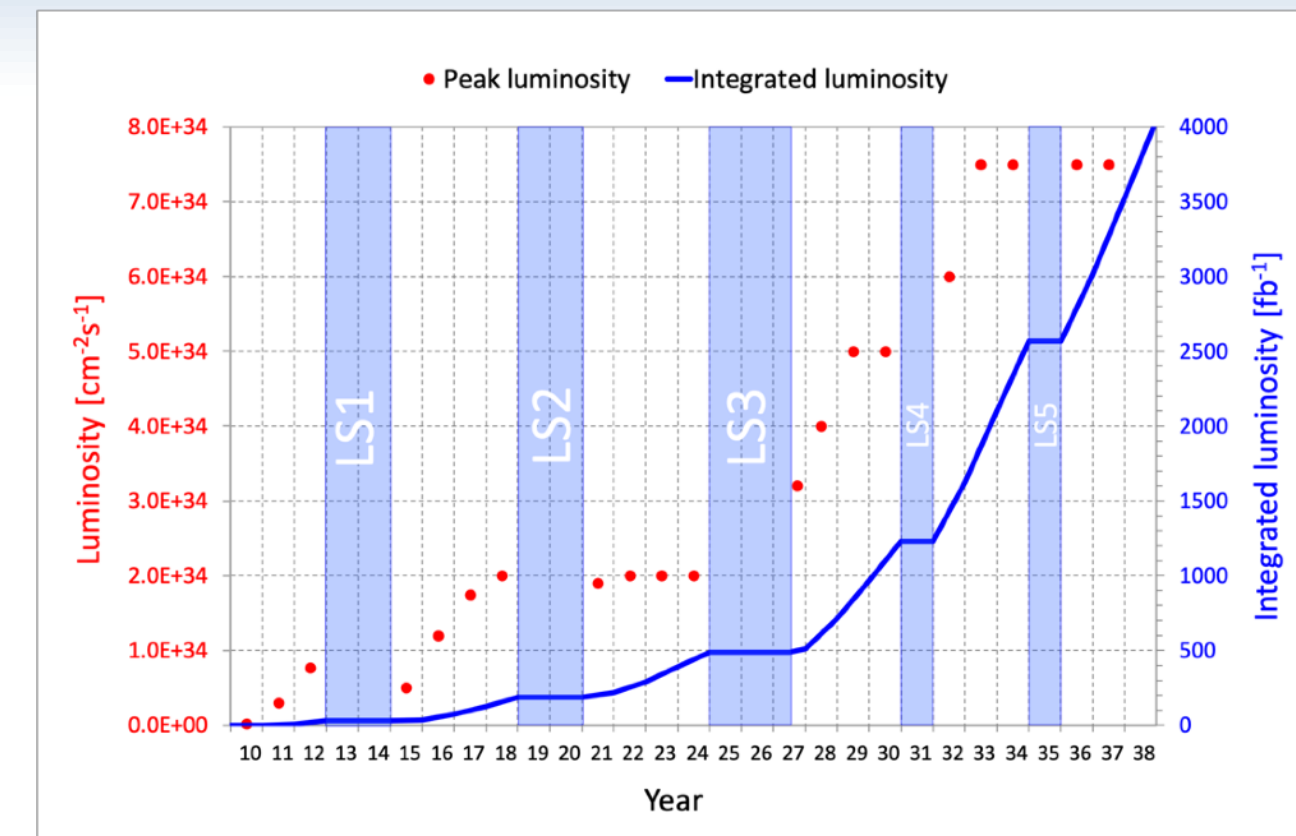


Probability distributions from generative networks
 Uncertainties on pdfs?

Open questions towards HL-LHC

A biased selection

- Facing **25 times** the amount of data
- What do we need to understand the data? (*read*: find new physics)



- **Precision predictions**
 - Higher order amplitudes
 - Event generation
 - Shower
 - Detector simulation

- **Optimized analysis for high-dimensional data**
 - Likelihood free inference
 - Optimal Observables, **Unfolding**
 - Anomaly detection
 - **Uncertainty treatment**

ML can help tackle all of these problems. Uncertainties included.