



Unsupervised Learning for Fun and Profit Precision

Anja Butter & Gregor Kasieczka
(gregor.kasieczka@uni-hamburg.de)
(butter@thphys.uni-heidelberg.de)

KITP - PRECISION2 I



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

CLUSTER OF EXCELLENCE
QUANTUM UNIVERSE



Partnership of
Universität Hamburg and DESY

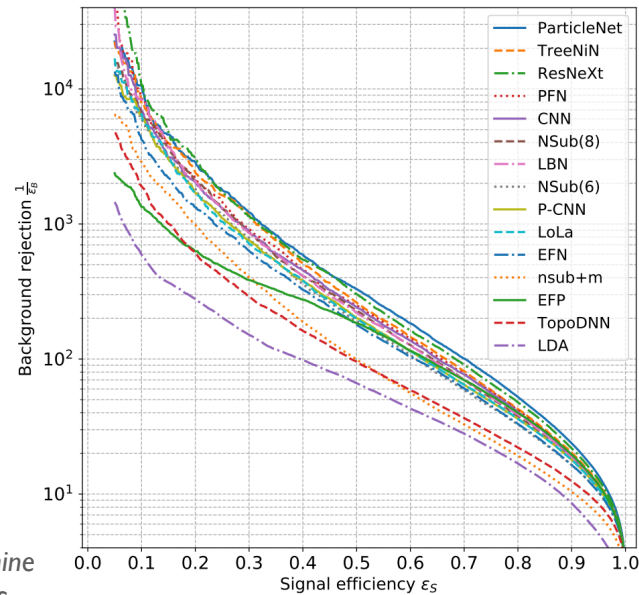


Bundesministerium
für Bildung
und Forschung

Motivation

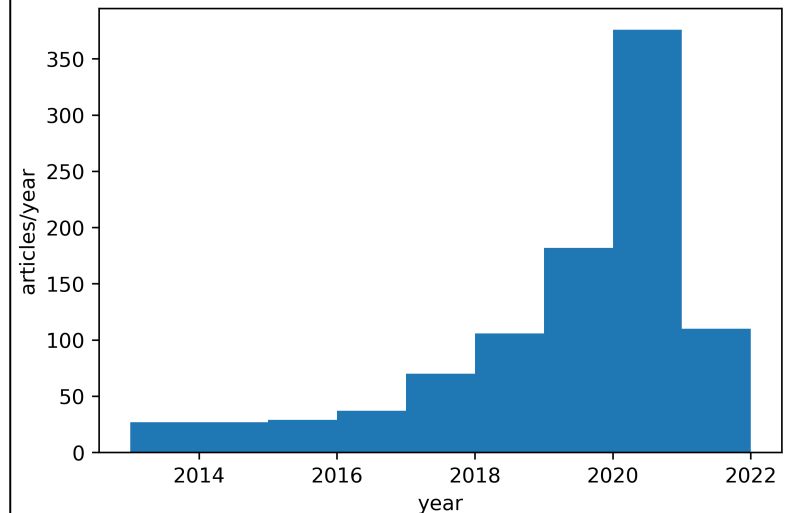
- Large excitement for machine learning in particle physics:
 - Particle tagging / signal selection
 - Low level reconstruction / calibration
 - Simulation

... and many more



Inspire Search:

("machine learning" or "deep learning" or neural) and (hep-ex or hep-ph or hep-th)

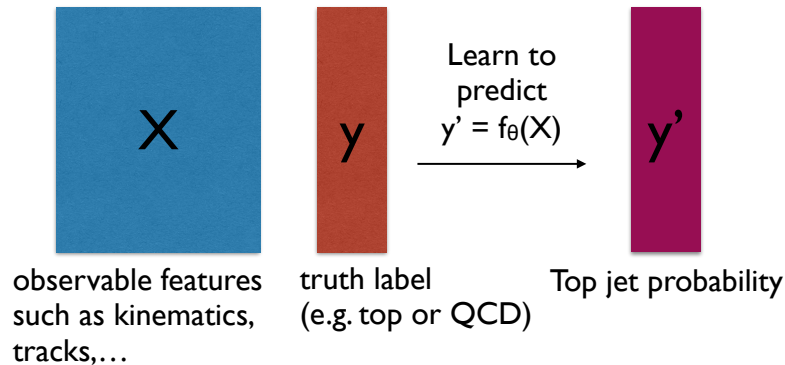


Two* types of problem:

Supervised Learning

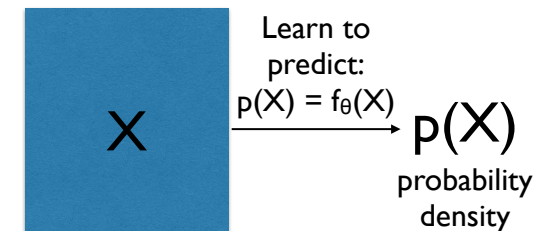
Attempt to infer some target (*truth label*):
classification (*jet flavour tagging*) or
regression (*energy calibration*)

Use training data with known labels
(often from Monte Carlo simulation)



Unsupervised

No target, learn the probability
distribution (directly from data)



**Maximize likelihood $p(X)$
(minimize $-\log p(x)$)**

***There also exists a number of other less-than-supervised approaches (weakly supervised learning, semi-supervised learning, ...) Not so important for now.**

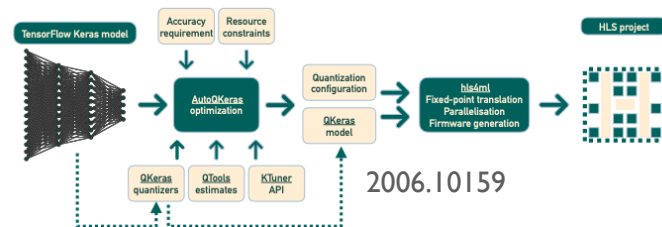
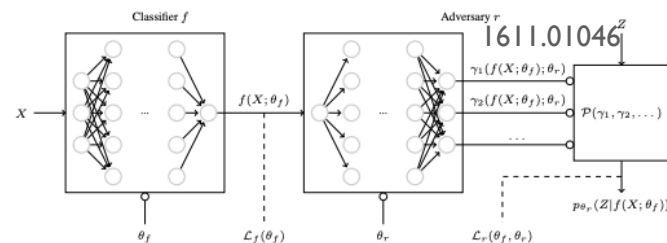
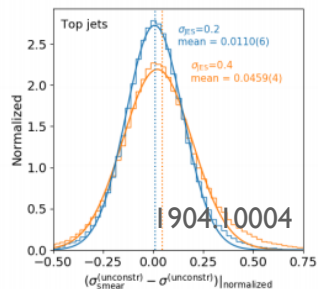
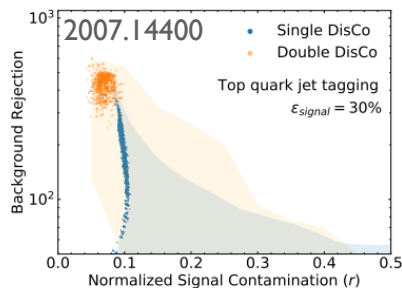
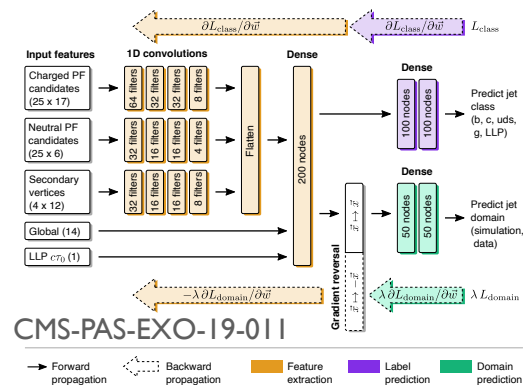
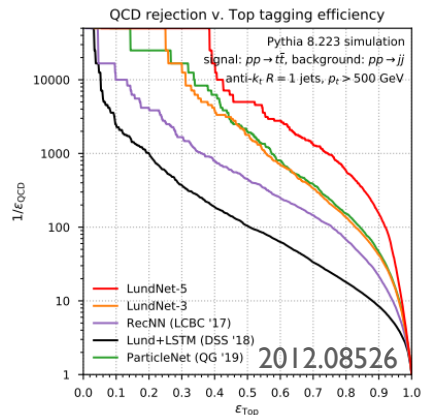
Supervised

Most early works fall under this category.

Crucially important for large number of tasks.

Need:

- Higher accuracy (easy to measure, many results)
- Better stability (domain adaptation issue)
- More control over uncertainties
- Resource efficient implementations
- Experimental integration
- ...



Unsupervised

Exciting space for developing new ideas
(also including all other forms of less-than-supervised learning).

Topic of our talks today.

Part I (Gregor):

How can we use a learned $p(X)$ to find new physics?
Anomaly detection //
Model independent
searches



Part II (Anja):

*How can we efficiently
sample from $p(X)$?*
Generative Models // Fast
simulation

Anomaly Searches

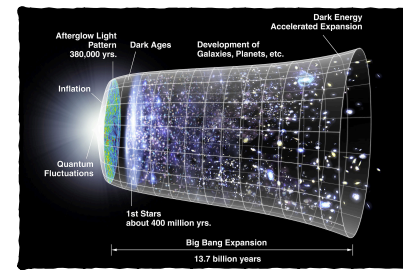
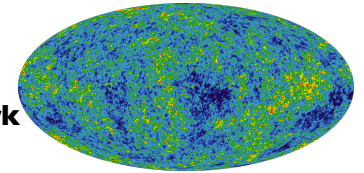
Physics Motivation

- Theoretical and experimental reasons to expect new physics beyond the Standard Model
- However, only negative results in searches
- Make sure that we do not miss potential discoveries at the LHC:
Supplement traditional searches with model-independent anomaly searches*

***Tricky term, will discuss meaning of model-independence later**

Why are neutrinos massive?

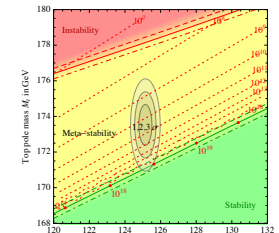
What is the nature of dark matter & dark energy?



What are the origins of the LHCb flavour anomaly?

Why is there more matter than anti-matter?

Why is there more matter than anti-matter?



Is the electroweak vacuum stable?

How can the Higgs boson be light when the mass receives large quantum corrections?

What are the details of cosmic inflation?

Dissecting the problem

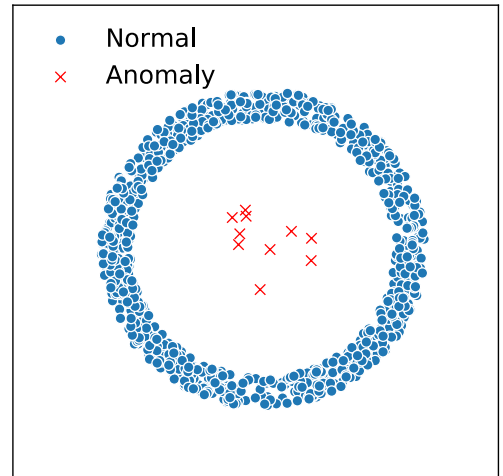
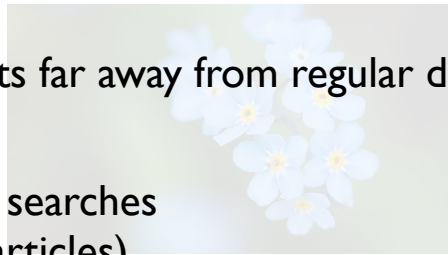
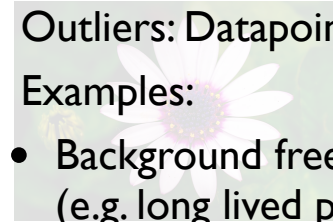
- Zero: *What are anomalies*
- First: *Build an anomaly scoring function $\mathbf{a}(X)$*
- Second: *Design analysis strategy*
- Third: *Interpret result*

What is an anomaly?

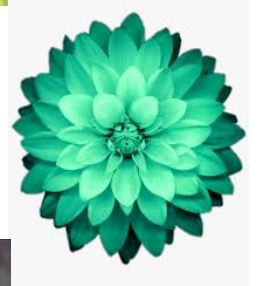


Point anomaly

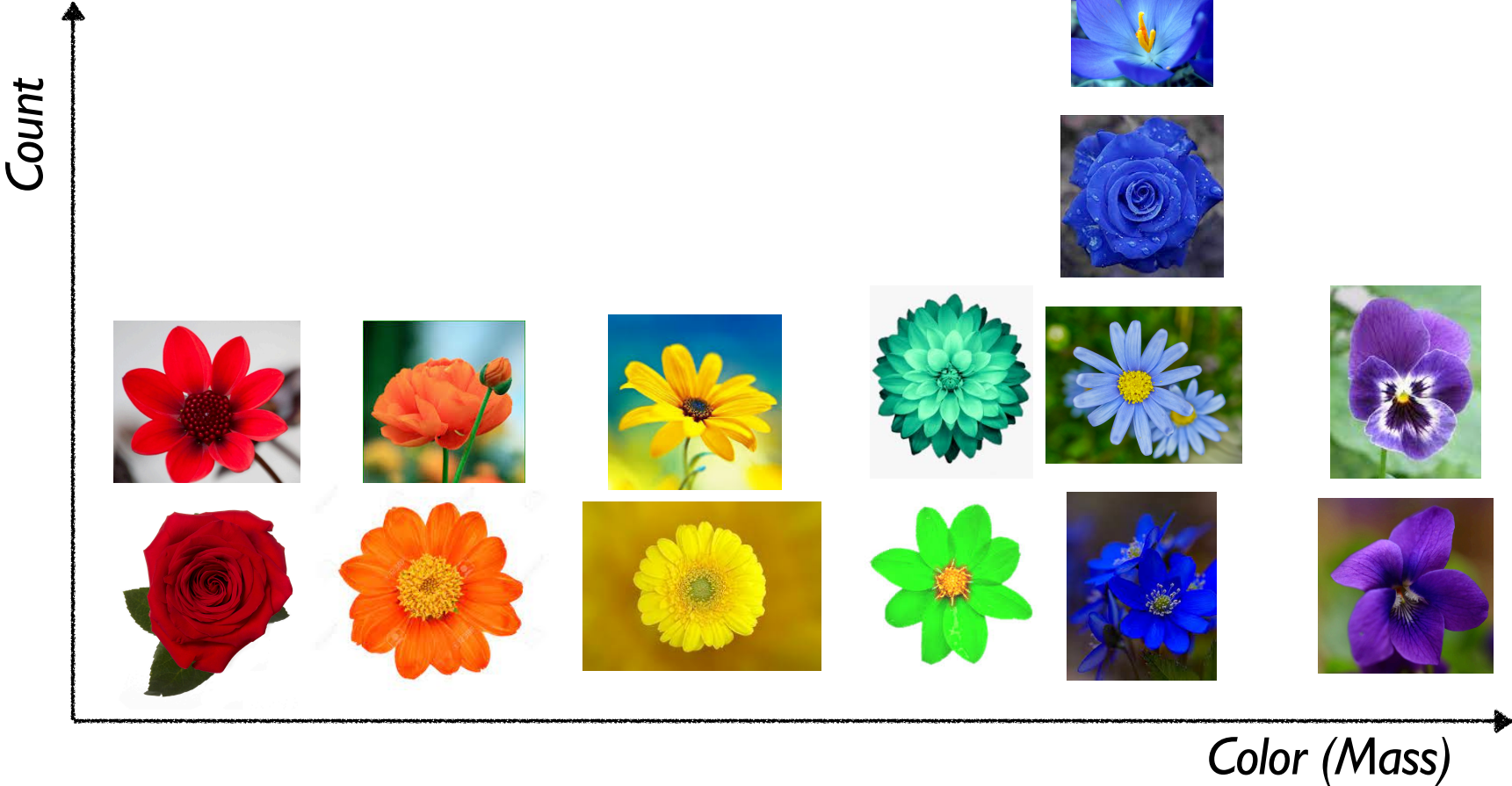
- Outliers: Datapoints far away from regular distribution
- Examples:
 - Background free searches (e.g. long lived particles)
 - Detector malfunctions



And now?



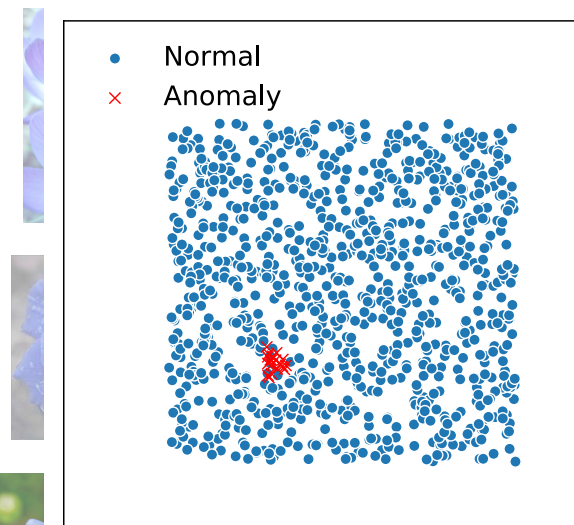
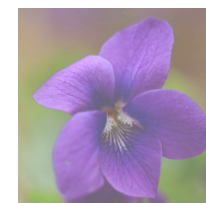
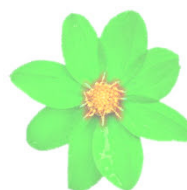
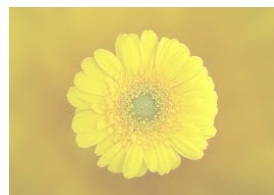
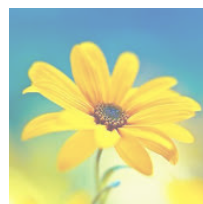
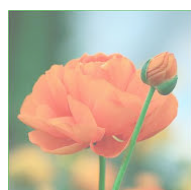
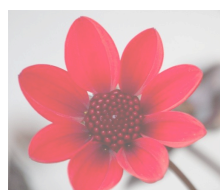
Group anomaly



Group anomaly

- Individual examples not anomalous, but interesting collective behaviour
- Examples:
 - New physics searches, e.g. resonances

Count



Color (Mass)

How to build anomaly score?

- Anomaly score a should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $a(x) = \text{(Semi-) Supervised}$

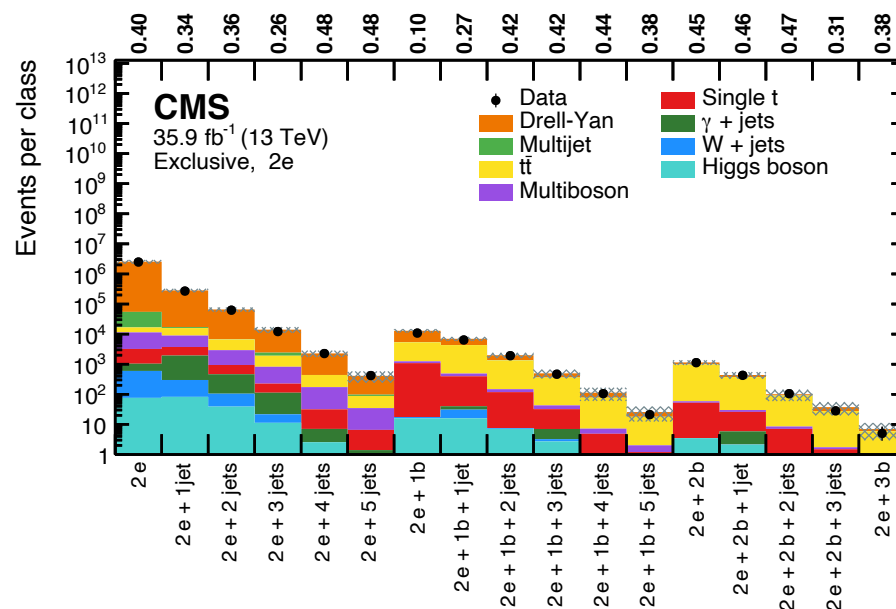
How to build anomaly score?

- Anomaly score a should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $a(x) = \text{(Semi-) Supervised}$
- Train binary classifier network (using simulation) to discriminate:
 - Standard Model background vs
 - Cocktail of new physics models
- Pros:
 - Close to known methods, simple training
 - Clear trade-off: width vs sensitivity
- Cons:
 - Ambiguity on mixture choice
 - Needs to account for residual difference between data/simulation

How to build anomaly score?

- Anomaly score \mathbf{a} should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $\mathbf{a}(\mathbf{x}) =$ (Semi-) Supervised
 - $\mathbf{a}(\mathbf{x}) = 1 / \mathbf{p}(\mathbf{x}|\mathbf{Background})$ (from simulation)

- Systematically look for differences between background simulation and data
- MUSIC / General search



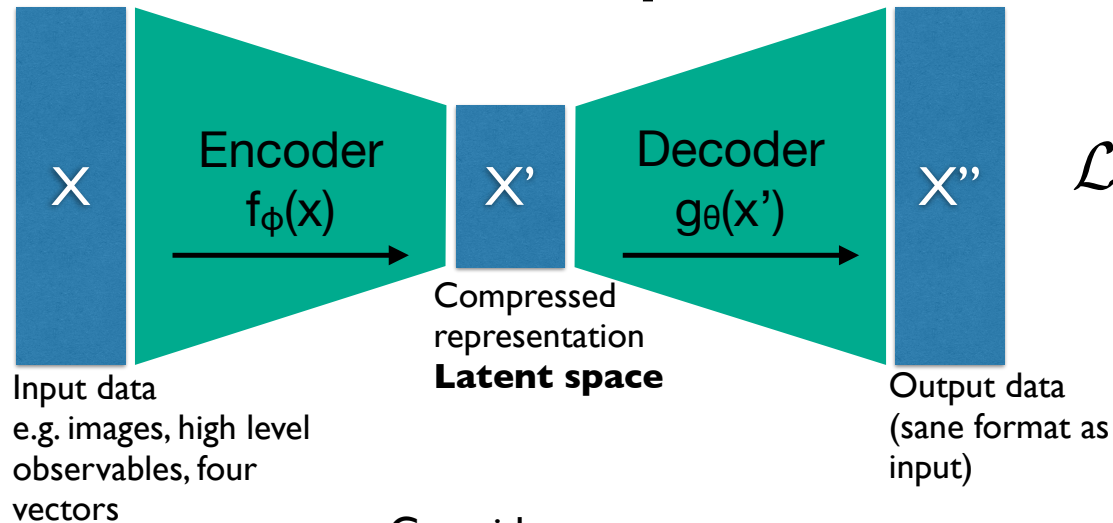
How to build anomaly score?

- Anomaly score \mathbf{a} should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $\mathbf{a}(\mathbf{x}) =$ (Semi-) Supervised
 - **$\mathbf{a}(\mathbf{x}) = \mathbf{1} / \mathbf{p}(\mathbf{x}|\mathbf{Background})$**
(from simulation)
- Systematically look for differences between background simulation and data
- MUSIC / General search
 - Use histograms of many variables in many dimensions to estimate
 - Potentially also improve via ML
- Pros:
 - Very signal model independent
 - Already delivering results
- Cons:
 - Strongly depends on background simulation
 - Large penalty from many histogram bins

How to build anomaly score?

- Anomaly score \mathbf{a} should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $\mathbf{a}(\mathbf{x}) =$ (Semi-) Supervised
 - $\mathbf{a}(\mathbf{x}) = \mathbf{I} / \mathbf{p}(\mathbf{x}|\mathbf{Background})$
(from simulation)
(from data)
- Search differences between different phase space regions in data
- Show example using **autoencoders**

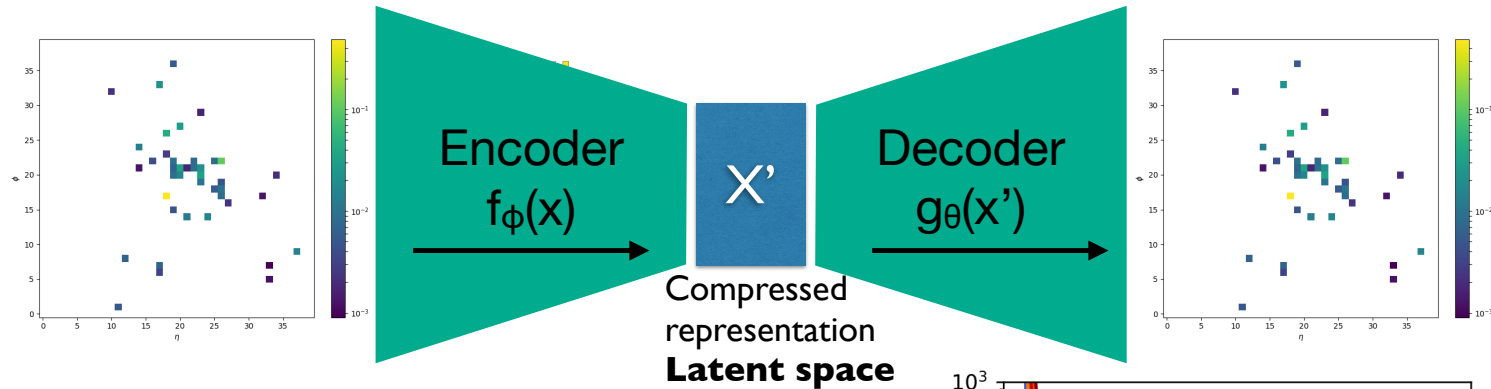
Example: Autoencoder



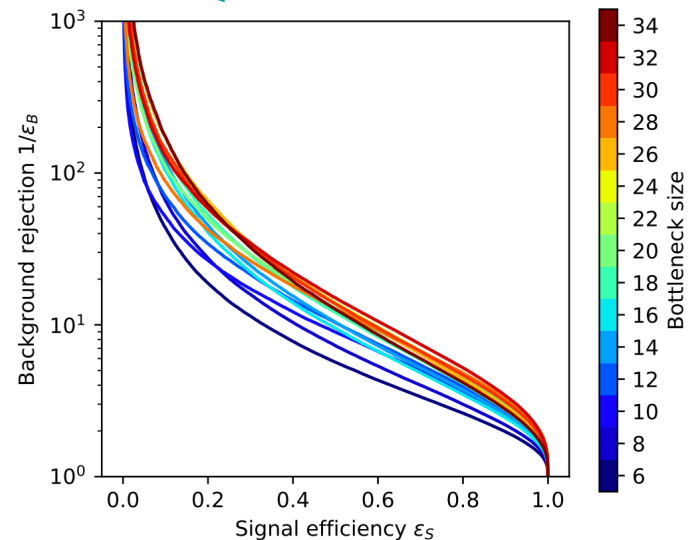
$$\mathcal{L}(x) = \|x - g_{\theta}(f_{\phi}(x))\|_2$$
$$a(x) = \mathcal{L}(x)$$

- Core idea:
 - Train lossy compression algorithm on anomaly-free data (minimise L)
 - Apply to data containing potential anomalies
 - Expect quality to decrease for atypical examples: anomaly score

Apply to jet images



- Represent data as images
 - Boosted top vs QCD jets (~ 600 GeV)
1 jet = 1 image (40x40 pixels, color=energy)
- Train QCD only sample
- Evaluate on mixed top/QCD jet sample
 - **Tops detected as anomaly**



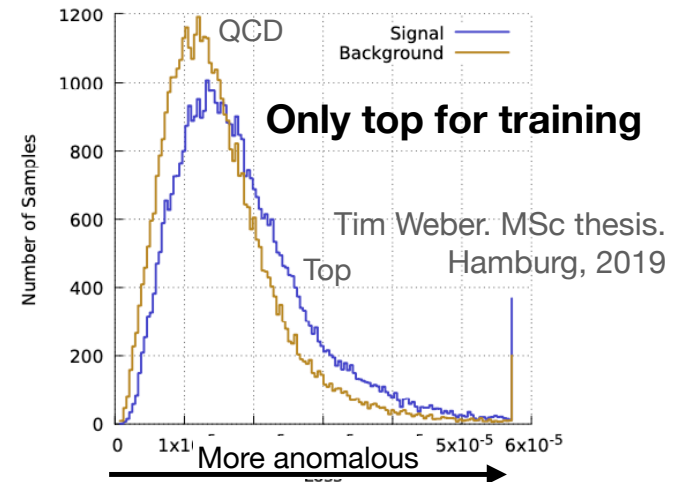
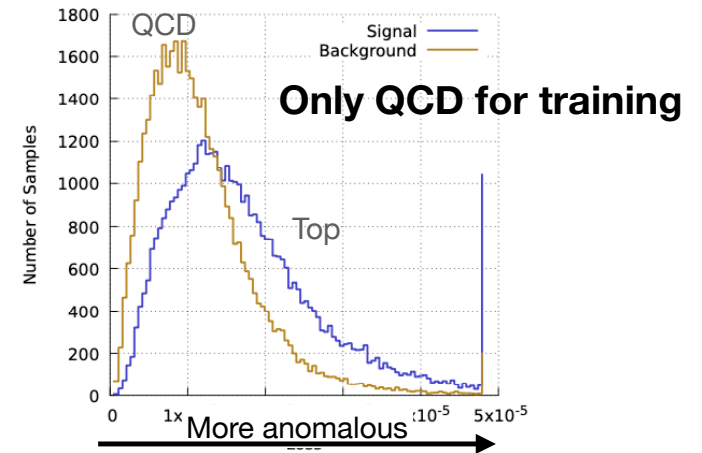
Heimel, GK, Plehn, Thompson, *QCD or What?*, 1808.08979

Farina, Nakai, Shih, *Searching for New Physics with Deep Autoencoders*, 1808.08992

Limitations

Complexity

- If anomalies are much simpler (therefore easier to reconstruct):
a(x) will still be lower, despite never encountered in training
- Observed with naive AE in QCD vs top
 - Train on tops only; top still considered anomaly wrt/ QCD

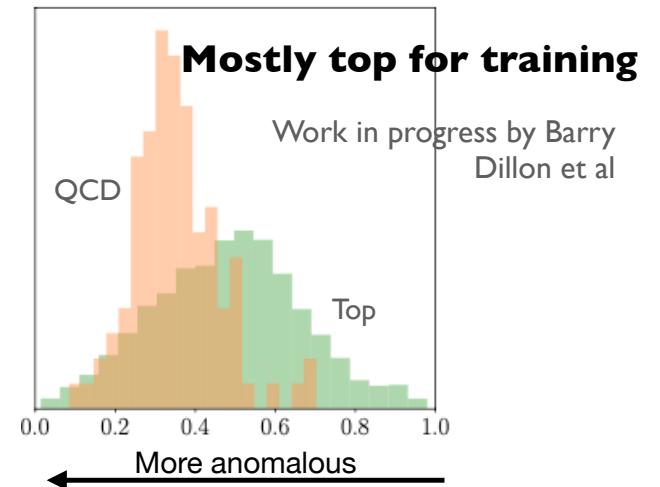
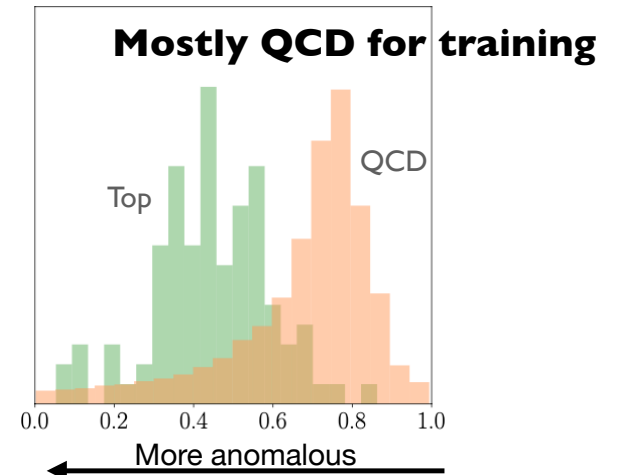


Limitations

Complexity

- If anomalies are much simpler (therefore easier to reconstruct):
a(x) will still be lower, despite never encountered in training
- Observed with naive AE in QCD vs top
 - Train on tops only; top still considered anomaly wrt/ QCD

Hope that this can be overcome with alternative AE trainings: Stay tuned for update by Heidelberg group using mixture model latent space!



Limitations

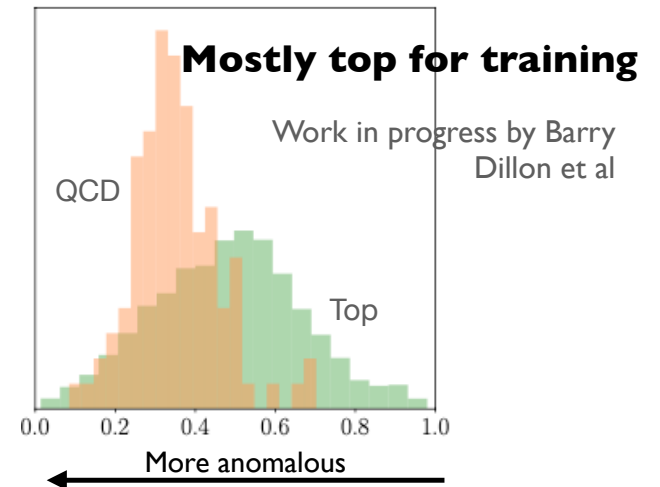
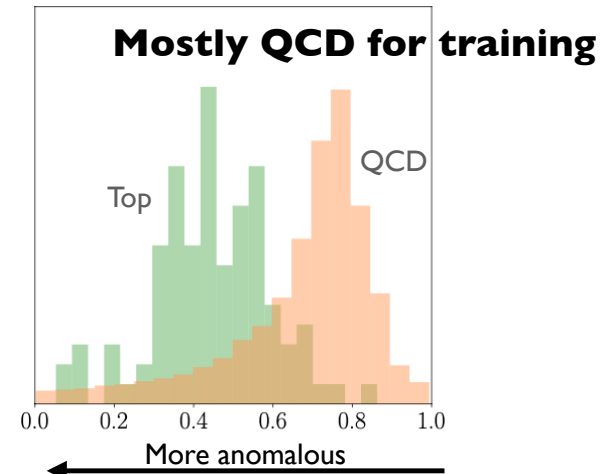
Complexity

- If anomalies are much simpler (therefore easier to reconstruct):
a(x) will still be lower, despite never encountered in training
- Observed with naive AE in QCD vs top
 - Train on tops only; top still considered anomaly wrt/ QCD

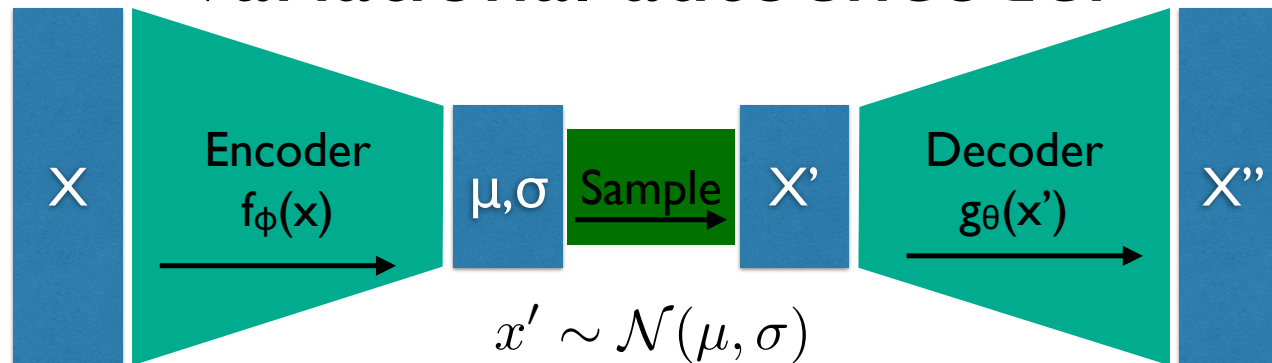
Hope that this can be overcome with alternative AE trainings: Stay tuned for update by Heidelberg group using mixture model latent space!

Topology

- Additional potential difficulty if data space has a non-trivial global topology.
See 2102.08380 for more



Brief aside on generative models: Variational autoencoder



- The decoder maps a latent space distribution X' to realistic examples
- Control over latent space:
 - Decode X' to generate new examples
- Achieve by:
 - Make X' Gaussian, encoder learns parameters μ, σ
 - Add term to loss so that (μ, σ) approach standard normal $(0, 1)$

How to build anomaly score?

- Anomaly score \mathbf{a} should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $\mathbf{a}(\mathbf{x}) =$ (Semi-) Supervised
 - $\mathbf{a}(\mathbf{x}) = \mathbf{1} / \mathbf{p}(\mathbf{x}|\mathbf{Background})$
(from simulation)
(from data)
- Search differences between different phase space regions in data
- Show example using **autoencoders**
- Pros:
 - Relatively signal model independent
 - Intuitive to construct and train
- Cons:
 - Little control over sensitivity
 - Some model assumptions needed for construction

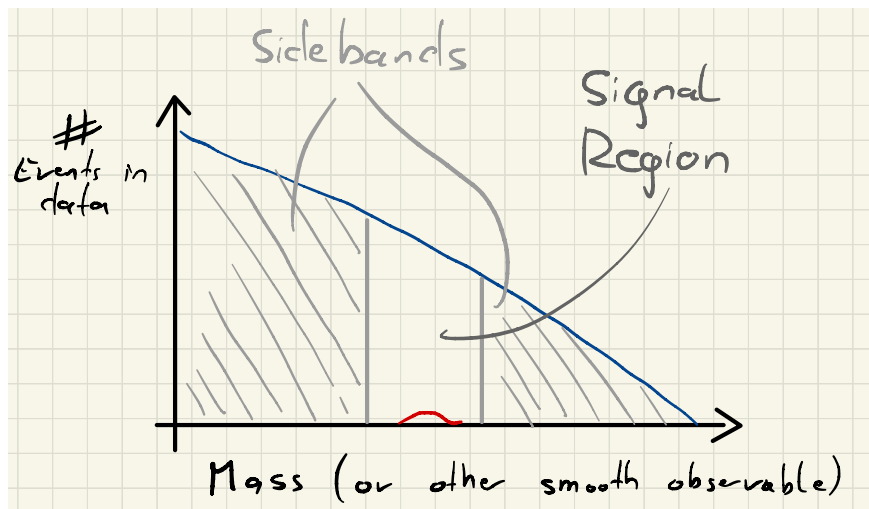
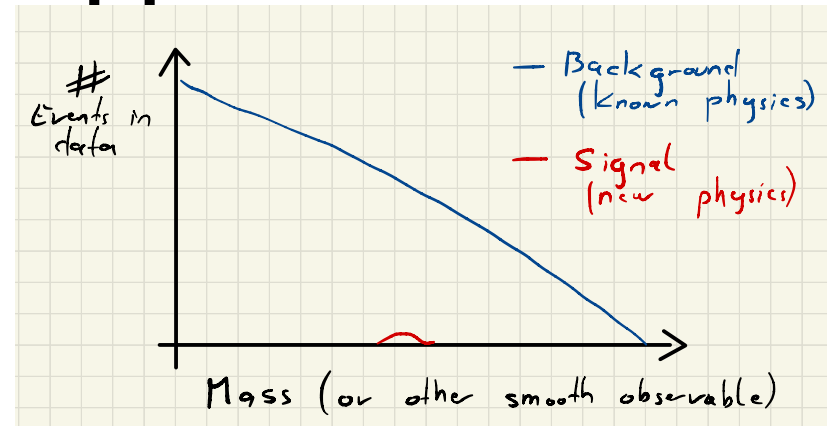
How to build anomaly score?

- Anomaly score **a** should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $a(x) =$ (Semi-) Supervised
 - $a(x) = 1 / p(x|\text{Background})$
(from simulation)
(from data)
 - **$a(x) = p(x|\text{Signal}) / p(x|\text{Background})$**
- Systematically look for differences between different phase space regions in data
- Show two examples:
 - **Mixed sample training**

Sideband approach

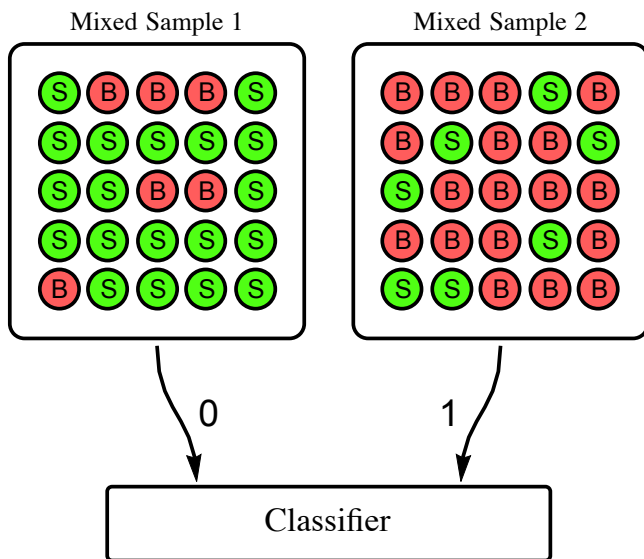
Key assumptions:

- There exists one feature so that:
 - Background distribution is smooth
 - Signal distribution is localised (and very small wrt/ background)



- Use sidebands to train anomaly score.
- Test signal region for new physics.
- Scan over different signal regions (trial factor)
- (Other ways to define anomaly-free regions in data possible as well. Not thoroughly explored yet)

Example: Mixed Sample training (aka CWola hunting)



$$L_{M_1/M_2} = \frac{p_{M_1}}{p_{M_2}} = \frac{f_1 p_S + (1 - f_1) p_B}{f_2 p_S + (1 - f_2) p_B} = \frac{f_1 L_{S/B} + (1 - f_1)}{f_2 L_{S/B} + (1 - f_2)}$$

- Distinguishing mixed samples is equivalent to signal/background classification assuming
- Signal/background in both mixed samples are from same source
- Sufficiently different mixed samples
- Translated to anomaly detection:
 - Train to distinguish signal region and sideband
 - Only use inputs independent of variable used to define these regions

Metodiev, Nachman, Thaler, *Classification without labels: Learning from mixed samples in high energy physics*, 1708.02949

Collins, Howe, Nachman, *Anomaly Detection for Resonant New Physics with Machine Learning*, 1805.02664

How to build anomaly score?

- Anomaly score a should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $a(x) =$ (Semi-) Supervised
 - $a(x) = 1 / p(x|\text{Background})$
(from simulation)
(from data)
 - **$a(x) = p(x|\text{Signal}) / p(x|\text{Background})$**
- Systematically look for differences between different phase space regions in data
- Show two examples:
 - **Mixed sample training**
 - density estimation
- Pros:
 - Relatively signal model independent
 - Cheap to train
- Cons:
 - Sensitive to correlations
 - Some model assumptions needed for construction

How to build anomaly score?

- Anomaly score a should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $a(x) =$ (Semi-) Supervised
 - $a(x) = 1 / p(x|\text{Background})$
(*from simulation*)
(*from data*)
 - **$a(x) = p(x|\text{Signal}) / p(x|\text{Background})$**
- Systematically look for differences between different phase space regions in data
- Show two examples:
 - Mixed sample training
 - **density estimation**

Example: Density Estimation

Per Neyman-Pearson: Likelihood-ratio is optimal test statistic

Unfortunately, $p(x|anomaly)$ is not available

$$L_{S/B} = \frac{p(x|anomaly)}{p(x|normal)}$$

Example: Density Estimation

Per Neyman-Pearson: Likelihood-ratio is optimal test statistic

Unfortunately, $p(x|anomaly)$ is not available

$$L_{S/B} = \frac{p(x|anomaly)}{p(x|normal)}$$

Build data/background ratio:

$$L_{D/B} = \frac{p(x)}{p(x|normal)}$$

Example: Density Estimation

Per Neyman-Pearson: Likelihood-ratio is optimal test statistic

Unfortunately, $p(x|\text{anomaly})$ is not available

$$L_{S/B} = \frac{p(x|\text{anomaly})}{p(x|\text{normal})}$$

Build data/background ratio:

$$L_{D/B} = \frac{p(x)}{p(x|\text{normal})}$$

Approximate background density using control measurement (e.g. sideband)

$$L_{D/B} \approx \frac{p(x)}{\tilde{p}(x|\text{normal})}$$

Example: Density Estimation

Per Neyman-Pearson: Likelihood-ratio is optimal test statistic

Unfortunately, $p(x|\text{anomaly})$ is not available

$$L_{S/B} = \frac{p(x|\text{anomaly})}{p(x|\text{normal})}$$

Build data/background ratio:

$$L_{D/B} = \frac{p(x)}{p(x|\text{normal})}$$

Approximate background density using control measurement (e.g. sideband)

$$L_{D/B} \approx \frac{p(x)}{\tilde{p}(x|\text{normal})}$$

Expand $p(x) = f_{\text{normal}} p(x|\text{normal}) + f_{\text{anomaly}} p(x|\text{anomaly})$

Example: Density Estimation

Per Neyman-Pearson: Likelihood-ratio is optimal test statistic
Unfortunately, $p(x|\text{anomaly})$ is not available

$$L_{S/B} = \frac{p(x|\text{anomaly})}{p(x|\text{normal})}$$

Build data/background ratio:

$$L_{D/B} = \frac{p(x)}{p(x|\text{normal})}$$

Approximate background density using control measurement (e.g. sideband)

$$L_{D/B} \approx \frac{p(x)}{\tilde{p}(x|\text{normal})}$$

Expand $p(x) = f_{\text{normal}} p(x|\text{normal}) + f_{\text{anomaly}} p(x|\text{anomaly})$

And insert: $L_{D/B} \approx f_{\text{normal}} + f_{\text{anomaly}} \frac{p(x|\text{anomaly})}{\tilde{p}(x|\text{normal})}$

Example: Density Estimation

Per Neyman-Pearson: Likelihood-ratio is optimal test statistic
 Unfortunately, $p(x|\text{anomaly})$ is not available

$$L_{S/B} = \frac{p(x|\text{anomaly})}{p(x|\text{normal})}$$

Build data/background ratio:
 • Data-Background likelihood is monotonous to Signal-Background likelihood if we can approximate background.

$$L_{D/B} = \frac{p(x)}{p(x|\text{normal})}$$

Approximate background density using control measurement (e.g. sideband)
 • **We can use this to construct an anomaly score**

$$L_{D/B} \approx \frac{p(x)}{\tilde{p}(x|\text{normal})}$$

Expand $p(x) = f_{\text{normal}} p(x|\text{normal}) + f_{\text{anomaly}} p(x|\text{anomaly})$

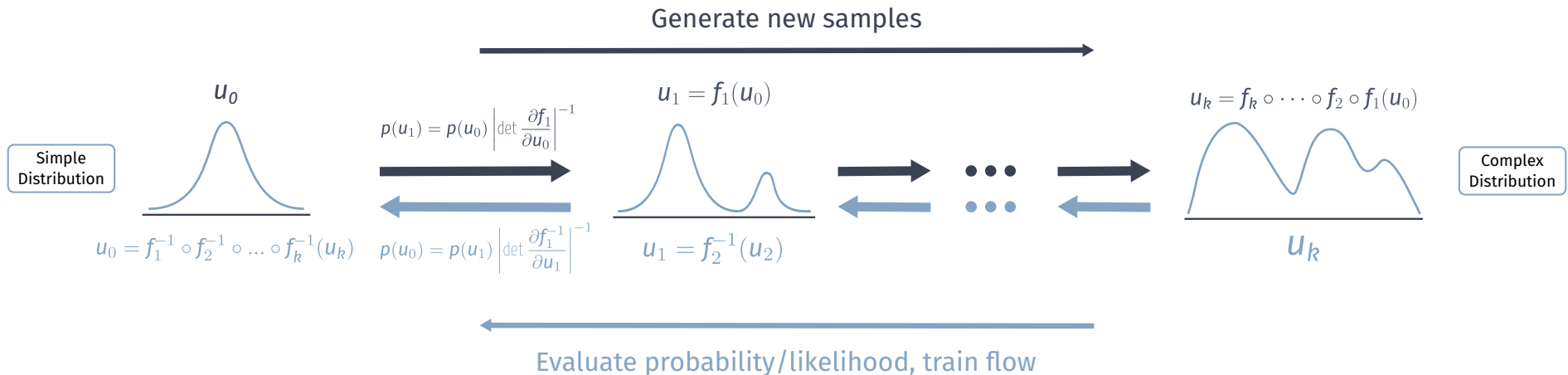
And insert: $L_{D/B} \approx f_{\text{normal}} + f_{\text{anomaly}} \frac{p(x|\text{anomaly})}{\tilde{p}(x|\text{normal})}$

Normalising Flows

- **Goal:** assign probability density to each datapoint
- Learn bijective transformation between data and a latent space with tractable probability
- Build from simple invertible transformations, tractable Jacobian

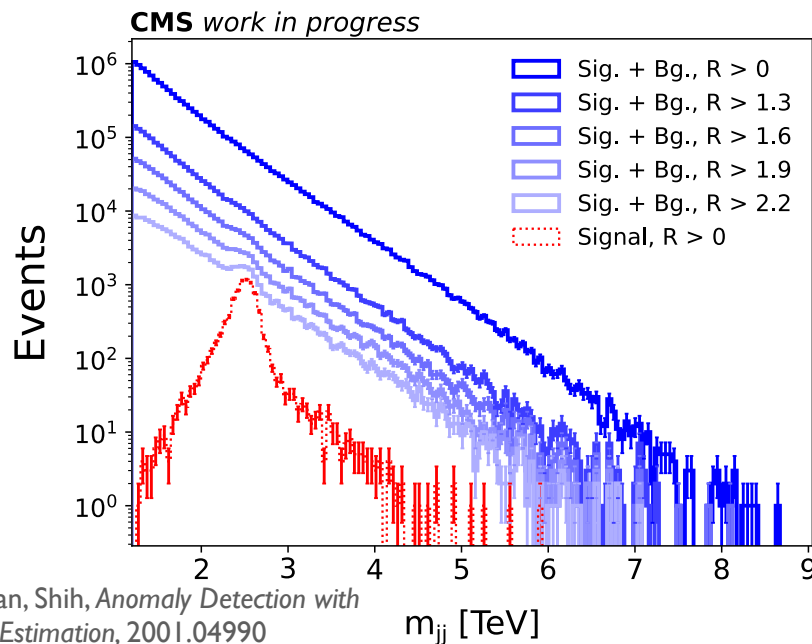
$$p(\mathbf{x}) = p(\mathbf{f}^{-1}(\mathbf{x})) \prod_i \left| \det \left(\frac{\partial \mathbf{f}_i^{-1}}{\partial \mathbf{x}} \right) \right| =$$

$$p(\mathbf{u}) \prod_i \left| \det \left(\frac{\partial \mathbf{f}_i}{\partial \mathbf{u}} \right) \right|^{-1}$$

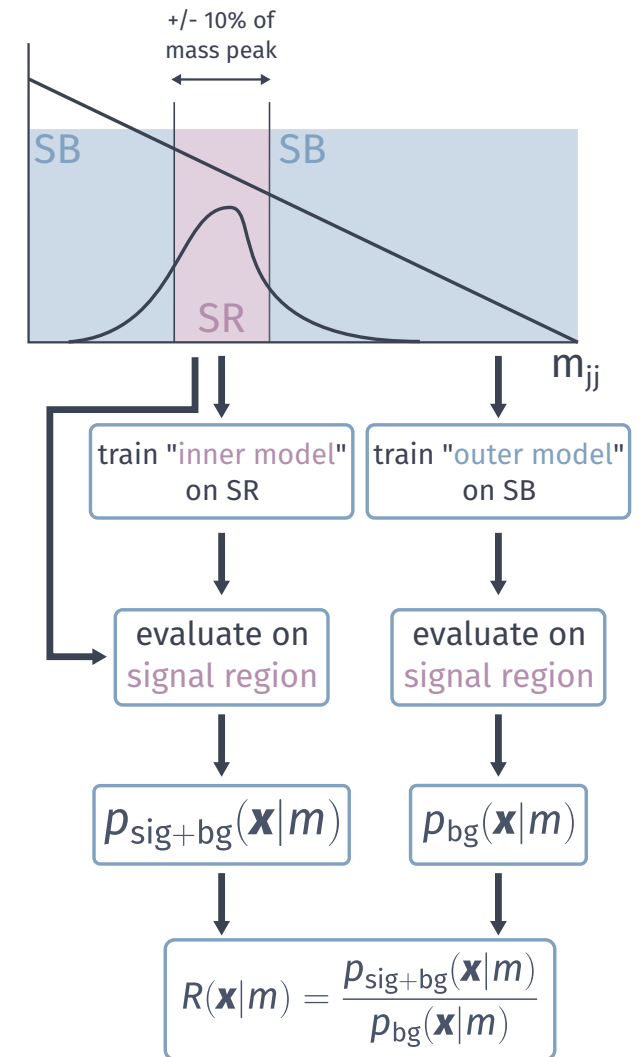


ANODE

- Use Masked Autoregressive Flow (1705.07057) to learn p (easy to invert NN with simple Jacobian)
- Compare extrapolated and in-region probability densities
- $a(\mathbf{x}) = R(\mathbf{x}|m) = \text{ratio of densities} = L_{D/B}$



Nachman, Shih, *Anomaly Detection with Density Estimation*, 2001.04990
Thanks to T. Loesche



How to build anomaly score?

- Anomaly score a should be high for anomalous (signal-like) and low for background-like events
- Some options:
 - $a(x) =$ (Semi-) Supervised
 - $a(x) = 1 / p(x|\text{Background})$
(*from simulation*)
(*from data*)
 - **$a(x) = p(x|\text{Signal}) / p(x|\text{Background})$**
- Systematically look for differences between different phase space regions in data
- Show two examples:
 - Mixed sample training
 - **density estimation**
- Pros:
 - Relatively signal model independent
 - Powerful
- Cons:
 - Expensive to train
 - Some model assumptions needed for construction

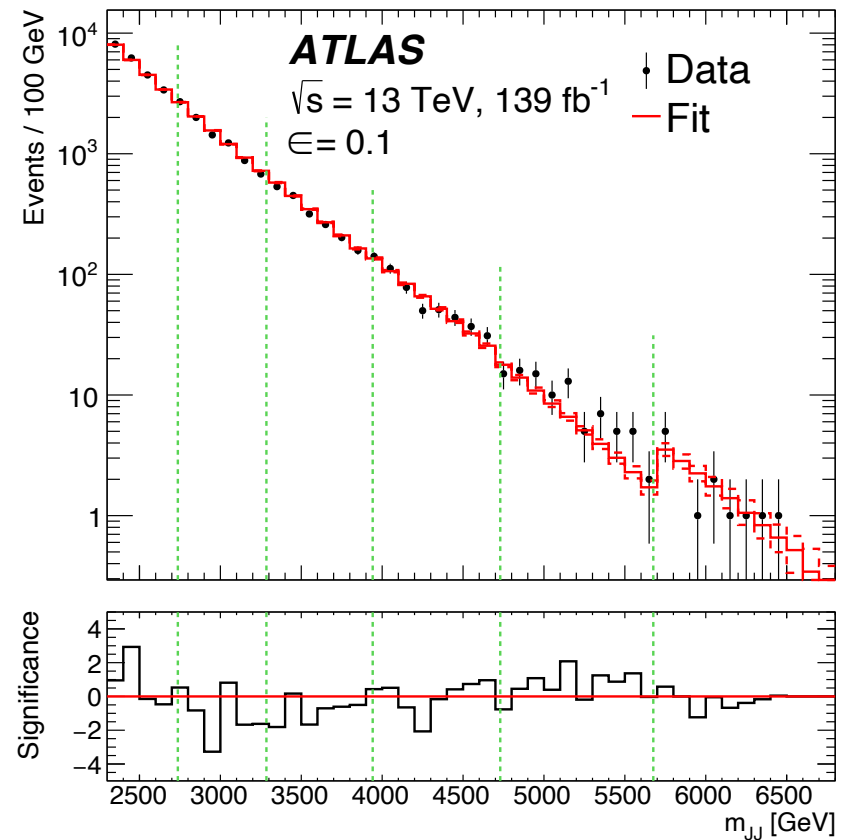
Moving on

- Many strategies exist to construct anomaly scores
 - $a(x) = "p(x|\text{Signal})"$ → **Semi-Supervised Cocktails**
 - $a(x) = 1 / p(x|\text{Background})$ → **General Search, Autoencoders**
 - $a(x) = p(x|\text{Signal}) / p(x|\text{Background})$ → **CWoLA, Density Estimation**
- **How can we use them in a search?**

MANY more ideas exist, see e.g. 2101.08320

Application: ATLAS Di-Jet Search

- ATLAS carried out a search following CWoLa approach
 - $A \rightarrow BC$ resonance search
no assumption on masses of A,B,C
 - Resonance search in di-jet invariant mass using $R=1.0$ jets for B,C
 - Split spectrum into discrete signal regions
 - Use CWoLa method, cut on 10% and 1% most anomalous events
 - Fit spectrum from sidebands
 - Interpret results in W' model



ATLAS Collaboration, *Dijet resonance search with weak supervision using $\sqrt{s}=13 \text{ TeV}$ pp collisions in the ATLAS detector*, 2005.02983

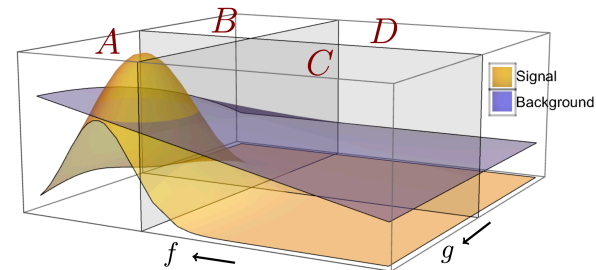
General

Simulation Driven

- Similar to classical analyses
- high signal model independence
- downside of data/simulation difference

Data Driven

- Straightforward idea to combine with bump hunt (see ATLAS example)
- Other data-driven techniques (ABCD?) should be possible as well, currently less explored
- Big advantage of data-only search:
No systematic uncertainties
(except background estimation from data)



How to interpret results?

Positive Result (some anomaly found)

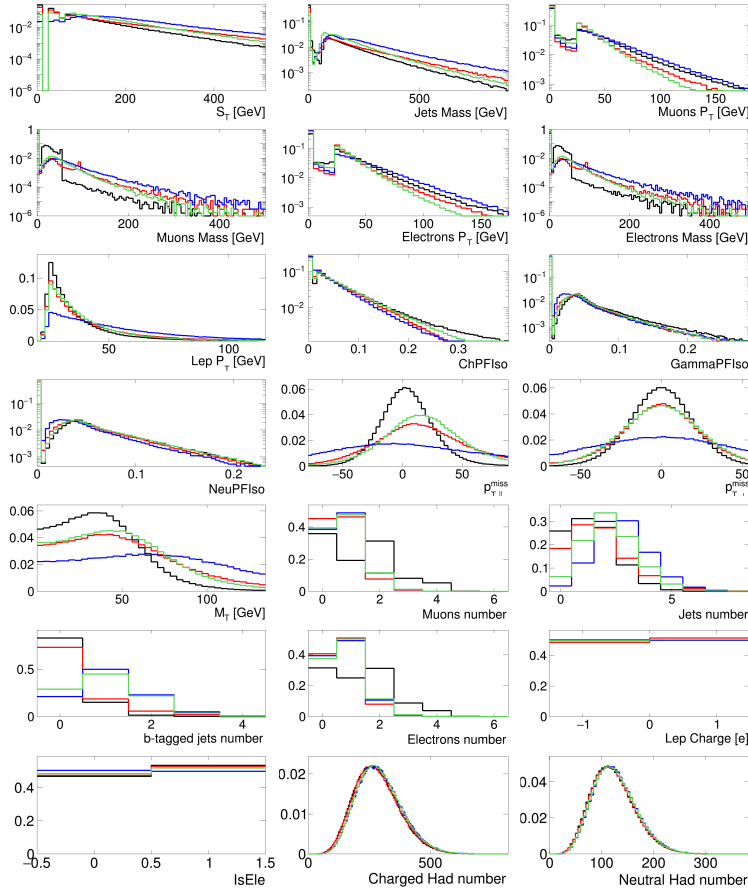
- Characterise what we found.
- Of course, can compare with different models and see which one fits.
(Also test for detector effects, of course)
- Ideas to systematise this needed
- Publish events?

Negative Result (no anomaly found)

- Need to interpret resulting exclusion
- Of course, can run different models and test (systematic uncertainties enter here!)
- Again, strategy for interpretation needed.
Publish anomaly score for recasting?

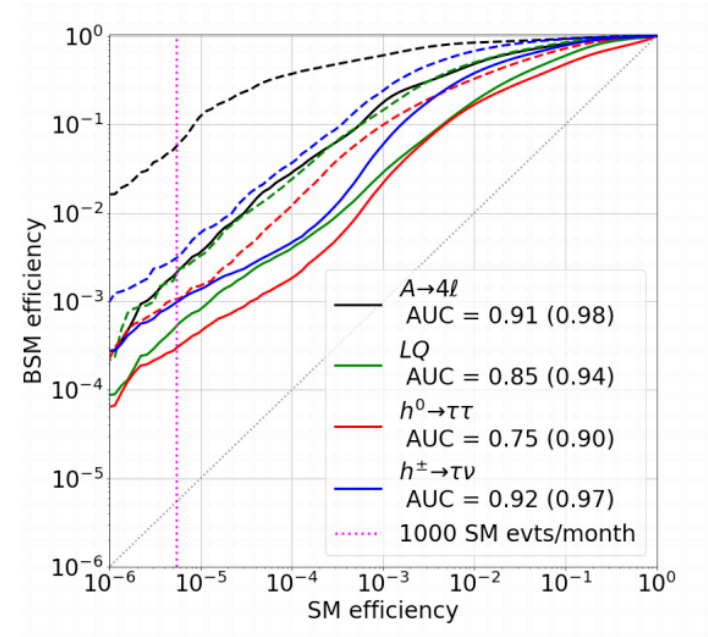
Aside: Trigger

- Consider anomaly detection for CMS/ATLAS triggers
- Strategy: (Variational) autoencoder trained on SM cocktail



21 high level observables as input

Cerri, Nguyen, Pierini, Spiropulu, Vlimant,
Variational Autoencoders for New Physics Mining at the Large Hadron Collider,
 1811.10276



Advertisement

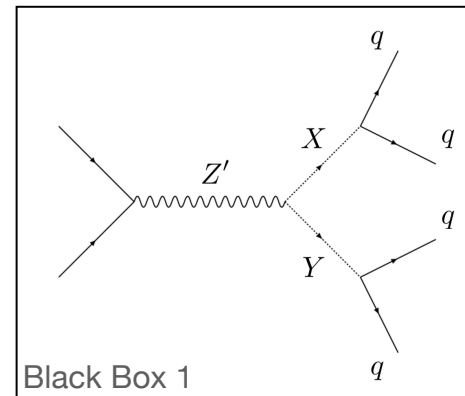
The LHC Olympics 2020

A Community Challenge for Anomaly
Detection in High Energy Physics



Gregor Kasieczka (ed),¹ Benjamin Nachman (ed),^{2,3} David Shih (ed),⁴ Oz Amram,⁵ Anders Andreassen,⁶ Kees Benkendorfer,^{2,7} Blaz Bortolato,⁸ Gustaaf Brooijmans,⁹ Florencia Canelli,¹⁰ Jack H. Collins,¹¹ Biwei Dai,¹² Felipe F. De Freitas,¹³ Barry M. Dillon,^{8,14} Ioan-Mihail Dinu,⁵ Zhongtian Dong,¹⁵ Julien Donini,¹⁶ Javier Duarte,¹⁷ D. A. Faroughy,¹⁰ Julia Gonski,⁹ Philip Harris,¹⁸ Alan Kahn,⁹ Jernej F. Kamenik,^{8,19} Charanjit K. Khosa,^{20,30} Patrick Komiske,²¹ Luc Le Pottier,^{2,22} Pablo Martín-Ramiro,^{2,23} Andrej Matevc,^{8,19} Eric Metodiev,²¹ Vinicius Mikuni,¹⁰ Inês Ochoa,²⁴ Sang Eon Park,¹⁸ Maurizio Pierini,²⁵ Dylan Rankin,¹⁸ Veronica Sanz,^{20,26} Nilai Sarda,²⁷ Uroš Seljak,^{2,3,12} Aleks Smolkovic,⁸ George Stein,^{2,12} Cristina Mantilla Suarez,⁵ Manuel Szewc,²⁸ Jesse Thaler,²¹ Steven Tsan,¹⁷ Silviu-Marian Udrescu,¹⁸ Louis Vaslin,¹⁶ Jean-Roch Vlimant,²⁹ Daniel Williams,⁹ Mikaeel Yunus¹⁸

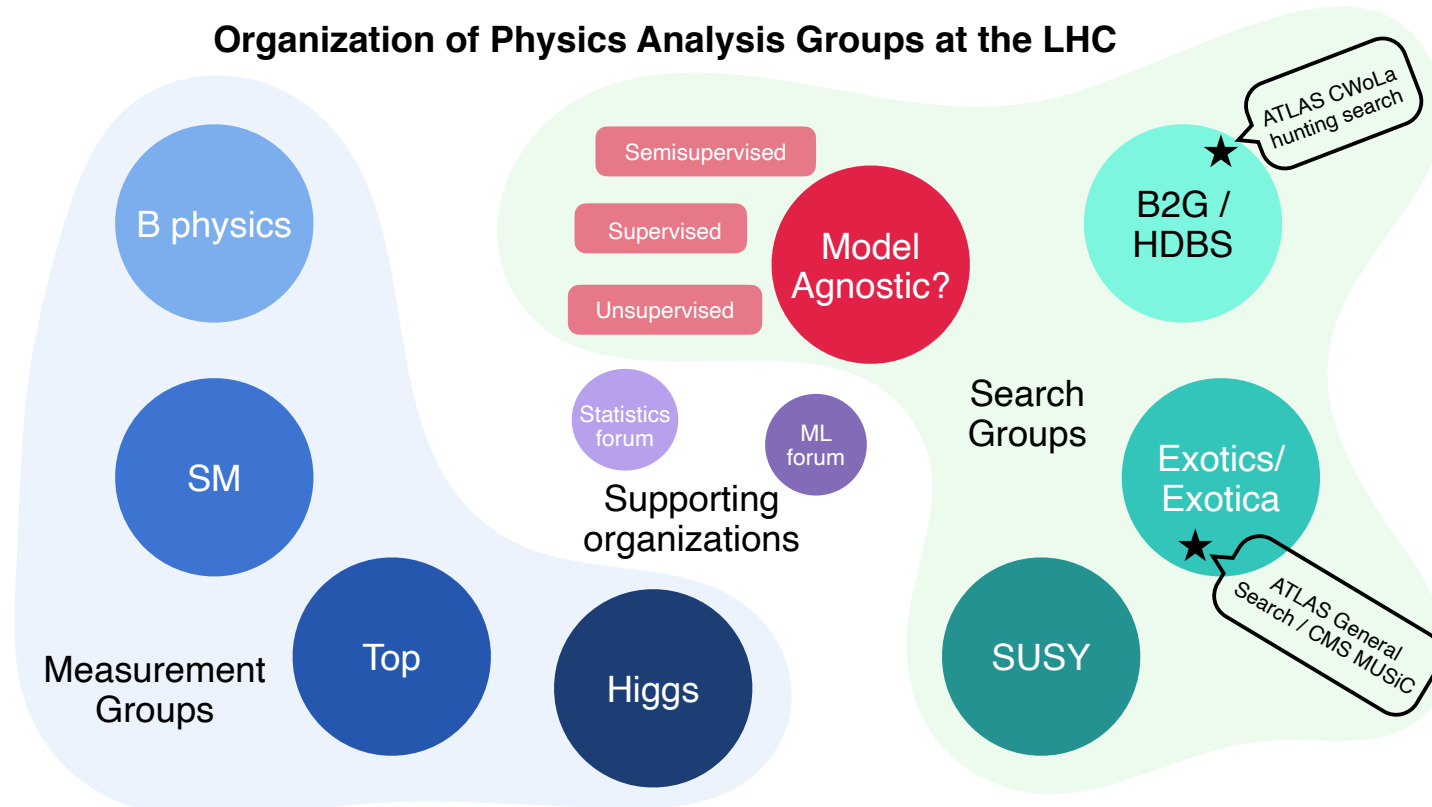
- For more on anomaly detection see: <https://indico.desy.de/e/anomaly2020>
- Public datasets available: <https://lhco2020.github.io/homepage/>
- Community paper with ~20 methods



Kasieczka, Nachman, Shih (eds), et al, *The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics*, 2101.08320

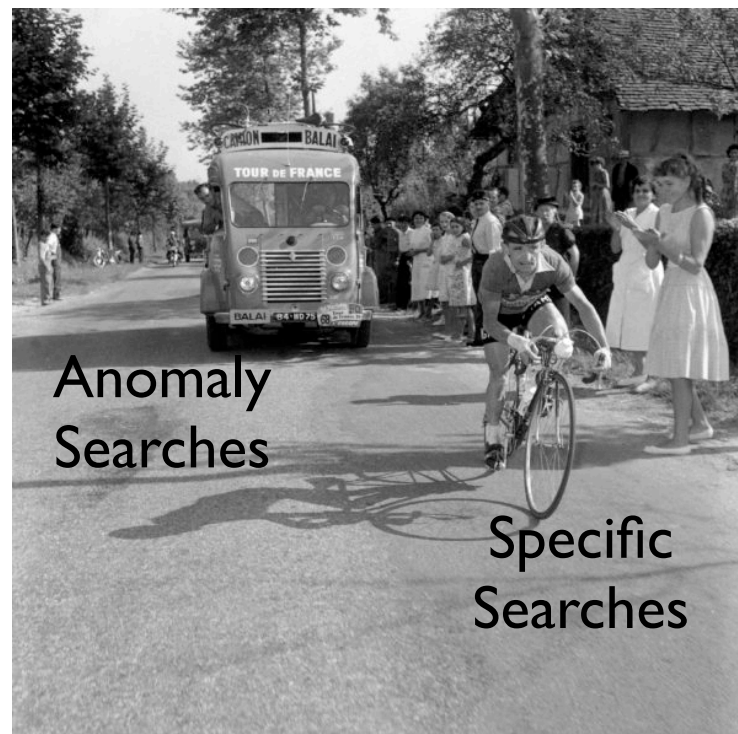
Advertisement

Organization of Physics Analysis Groups at the LHC



Closing

- Focused on new physics searches.
Anomaly detection also considered for data quality monitoring, detector control, computing monitoring
- Improve power of anomaly detectors
- Extends to higher number of features
 - Beyond images / high-level observables
- How to properly encode normal physics / anomalous physics?
- Systematically understand sensitivity of different approaches
- Develop interpretation strategies
- Widely apply to experimental data



Thank you!

Backup

MADE/MAF

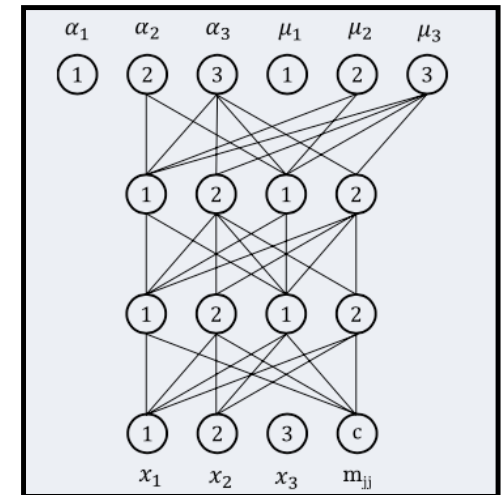
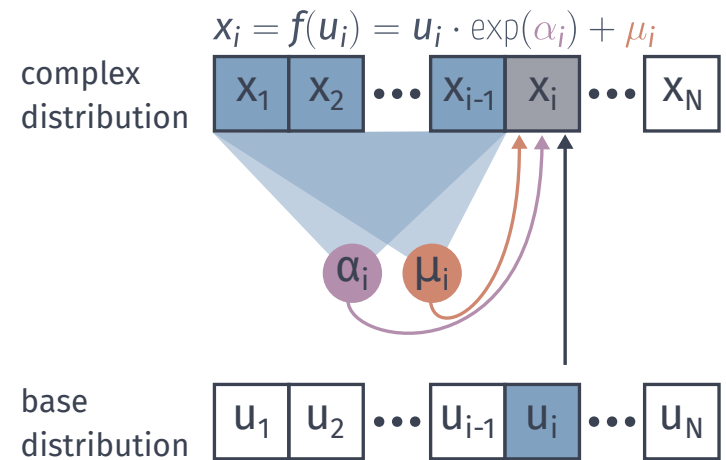
- Masked Autoregressive flow (1502.03509/1705.07057)
- Start with fully connected network, but drop connections so output a_j / μ_j are only connected to input x_1, \dots, x_{j-1}
- Autoregressive: no dependence of early features on late features
 - -> Jacobian is upper triangular matrix and easily invertible
- Combine multiple such blocks

$$p(x) = \prod_i p(x_i | x_{1:i-1})$$

$$p(x_i | x_{1:i-1}) = \mathcal{N}(x_i | \mu_i, (\exp \alpha_i)^2)$$

$$\mu_i = f_{\mu_i}(x_{1:i-1})$$

$$\alpha_i = f_{\alpha_i}(x_{1:i-1})$$



Unsupervised Learning for Fun and Precision

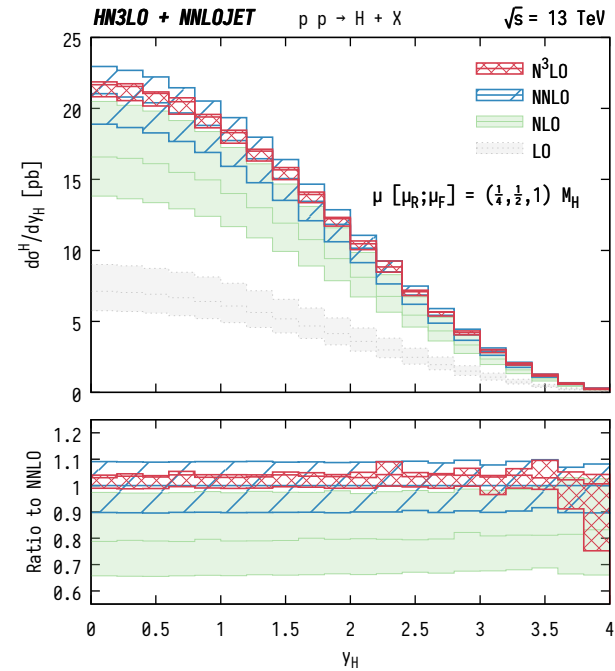
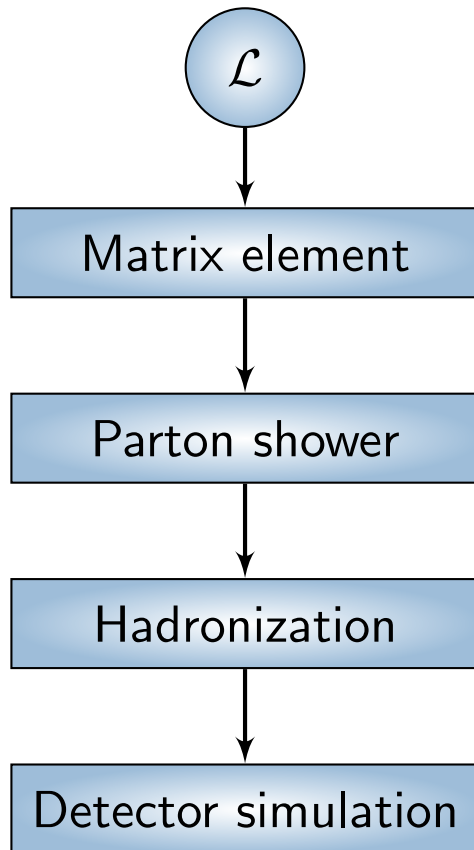
KITP - Precision21

Anja Butter & Gregor Kasieczka

ITP, Universität Heidelberg

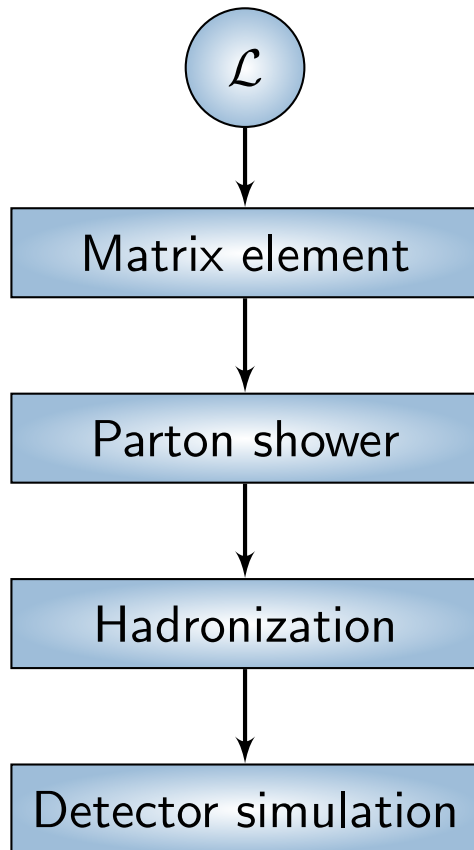


Precision simulations with limited resources

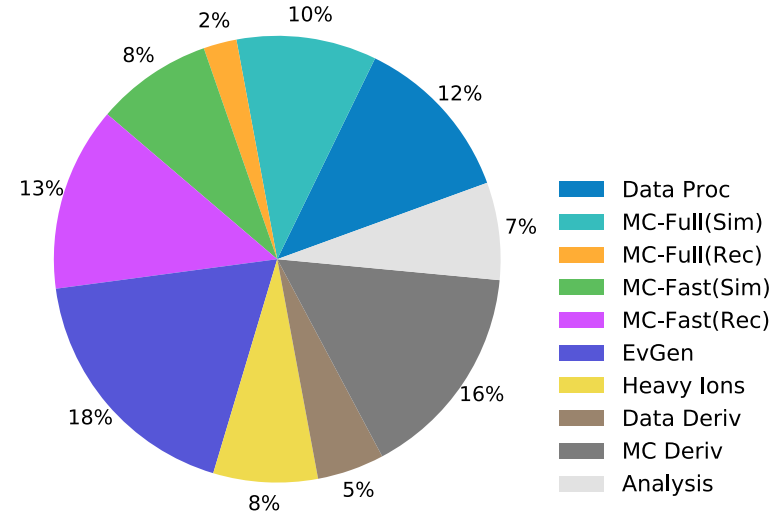


[1807.11501] Cieri, Chen, Gehrmann, Glover, Huss

Precision simulations with limited resources



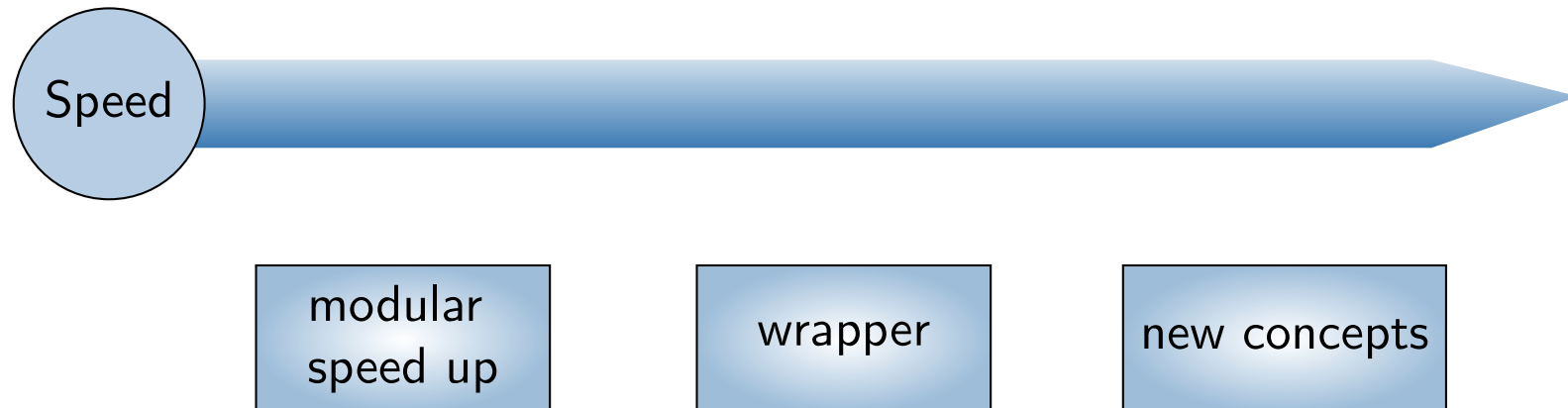
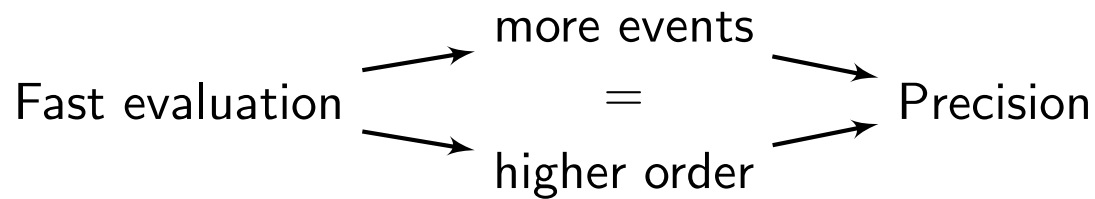
ATLAS Preliminary
2020 Computing Model -CPU: 2030: Aggressive R&D



Speed = Precision

How can we boost MC simulations

- ML 2.0 Generative models
 - Can we simulate new data?



Boosting standard event generation...

1. Generate phase space points

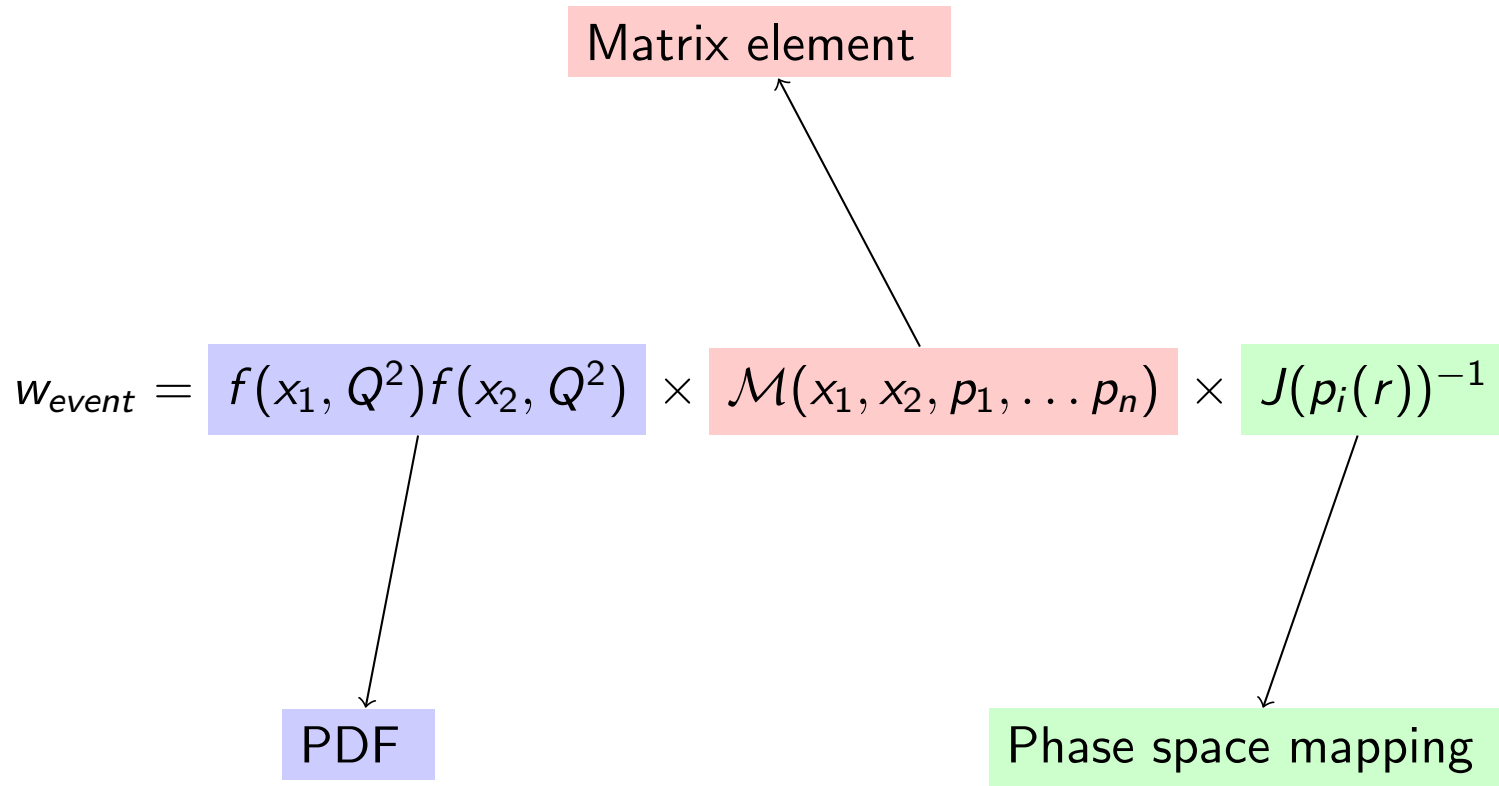
2. Calculate event weight

$$W_{event} = f(x_1, Q^2)f(x_2, Q^2) \times \mathcal{M}(x_1, x_2, p_1, \dots, p_n) \times J(p_i(r))^{-1}$$

3. Unweighting via importance sampling

→ optimal for $w \approx 1$

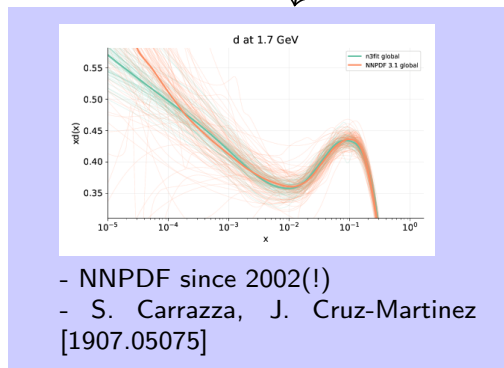
Boosting standard event generation...



Boosting standard event generation...

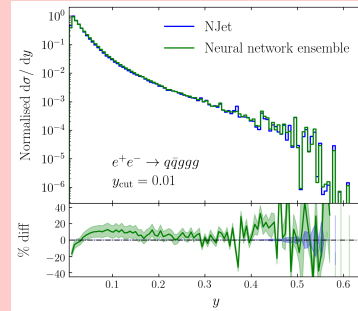
Matrix element

$$W_{event} = f(x_1, Q^2)f(x_2, Q^2) \times \mathcal{M}(x_1, x_2, p_1, \dots, p_n) \times J(p_i(r))^{-1}$$



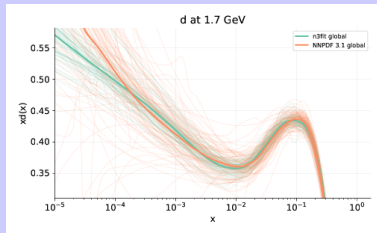
Phase space mapping

Boosting standard event generation...



- Amplitude estimation
- S. Badger, J. Bullock [2002.07516]
- J. Bendavid [1707.00028]

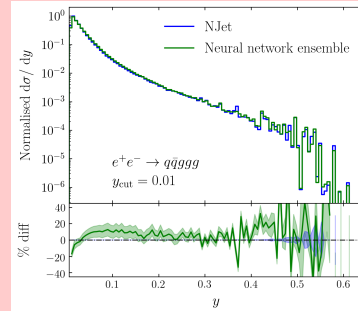
$$W_{event} = f(x_1, Q^2)f(x_2, Q^2) \times \mathcal{M}(x_1, x_2, p_1, \dots, p_n) \times J(p_i(r))^{-1}$$



- NNPDF since 2002(!)
- S. Carrazza, J. Cruz-Martinez [1907.05075]

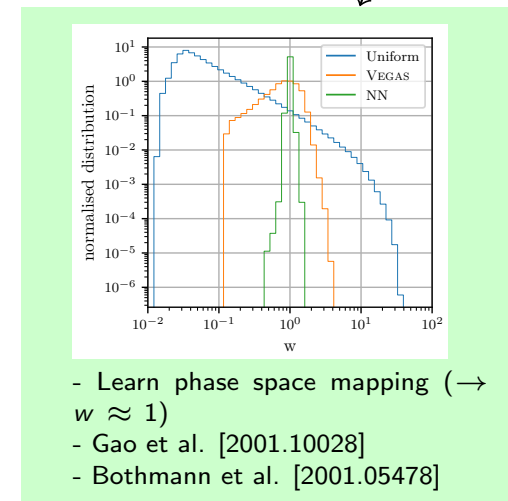
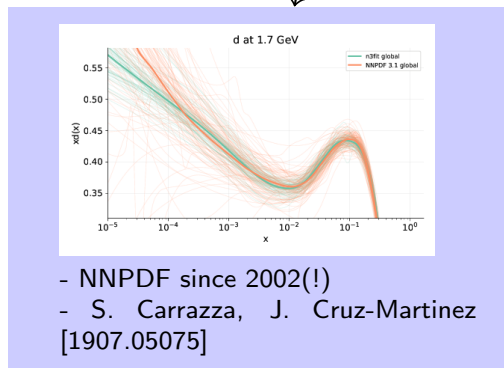
Phase space mapping

Boosting standard event generation...



- Amplitude estimation
- S. Badger, J. Bullock [2002.07516]
- J. Bendavid [1707.00028]

$$W_{event} = f(x_1, Q^2)f(x_2, Q^2) \times \mathcal{M}(x_1, x_2, p_1, \dots, p_n) \times J(p_i(r))^{-1}$$



... or training directly on event samples

Event generation

- Generating 4-momenta
- $Z > ll, pp > jj, pp > t\bar{t} + \text{decay}$

[1901.00875] Otten et al. **VAE & GAN**

[1901.05282] Hashemi et al. **GAN**

[1903.02433] Di Sipio et al. **GAN**

[1903.02556] Lin et al. **GAN**

[1907.03764, 1912.08824] Butter et al. **GAN**

[1912.02748] Martinez et al. **GAN**

[2001.11103] Alanazi et al. **GAN**

[2011.13445] Stienen et al. **NF**

[2012.07873] Backes et al. **GAN**

[2101.08944] Howard et al. **VAE**

Detector simulation

- Jet images
- Fast calorimeter simulation

[1701.05927] de Oliveira et al. **GAN**

[1705.02355, 1712.10321] Paganini et al. **GAN**

[1802.03325, 1807.01954] Erdmann et al. **GAN**

[1805.00850] Musella et al. **GAN**

[ATL-SOFT-PUB-2018-001, ATLAS-SIM-2019-004, ATL-SOFT-PROC-2019-007] ATLAS **VAE & GAN**

[1909.01359] Carazza and Dreyer **GAN**

[1912.06794] Belayneh et al. **GAN**

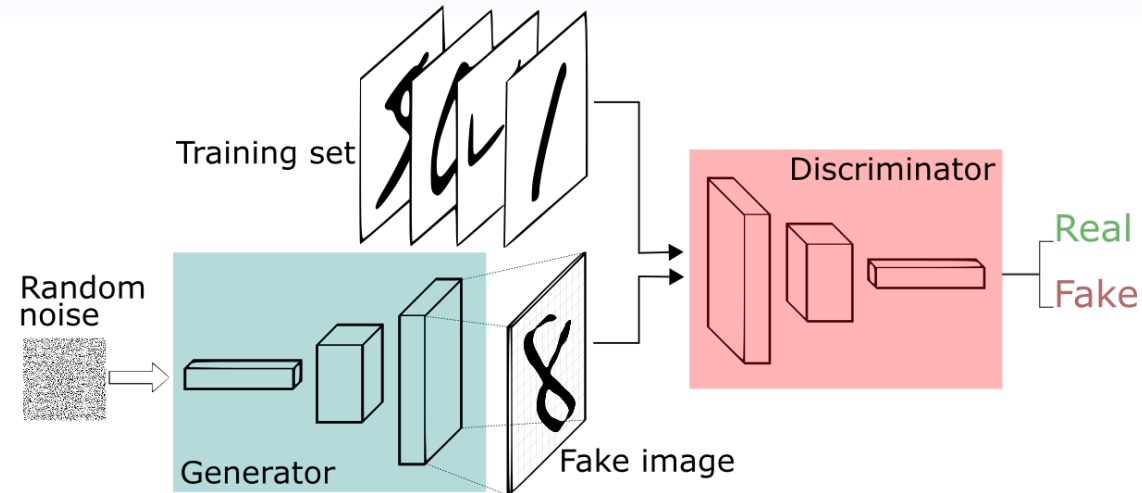
[2005.05334, 2102.12491] Buhmann et al. **VAE**

[2009.03796] Diefenbacher et al. **GAN**

[2009.14017] Lu et al.

NO claim to completeness!

Generative Adversarial Networks



Discriminator $[D(x_T) \rightarrow 1, D(x_G) \rightarrow 0]$

$$L_D = \langle -\log D(x) \rangle_{x \sim P_{Truth}} + \langle -\log(1 - D(x)) \rangle_{x \sim P_{Gen}} \rightarrow -2 \log 0.5$$

Generator $[D(x_G) \rightarrow 1]$

$$L_G = \langle -\log D(x) \rangle_{x \sim P_{Gen}}$$

\Rightarrow **Nash Equilibrium**

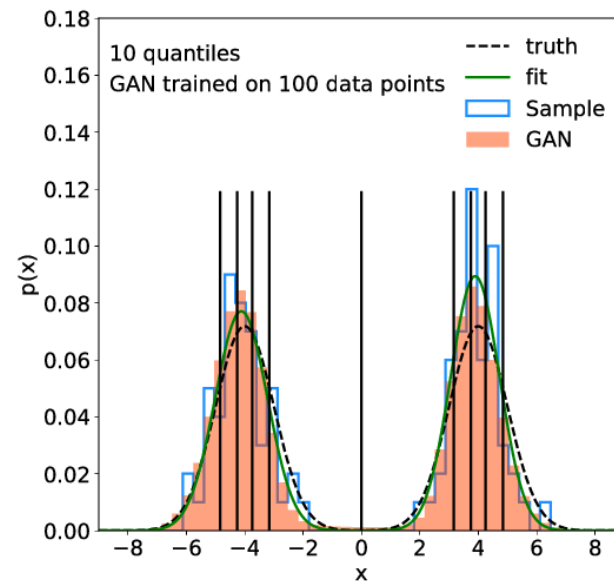
\Rightarrow **New statistically independent samples**

What is the statistical value of GANned events? [2008.06545]

- Camel function
- Sample vs. GAN vs. 5 param.-fit

Evaluation on quantiles:

$$\text{MSE}^* = \sum_{j=1}^{N_{\text{quant}}} \left(p_j - \frac{1}{N_{\text{quant}}} \right)^2$$

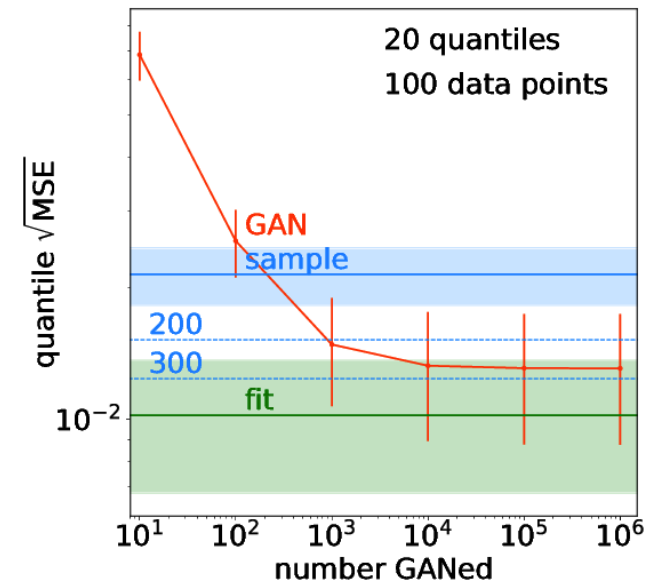


What is the statistical value of GANned events? [2008.06545]

- Camel function
- Sample vs. GAN vs. 5 param.-fit

Evaluation on quantiles:

$$\text{MSE}^* = \sum_{j=1}^{N_{\text{quant}}} \left(p_j - \frac{1}{N_{\text{quant}}} \right)^2$$

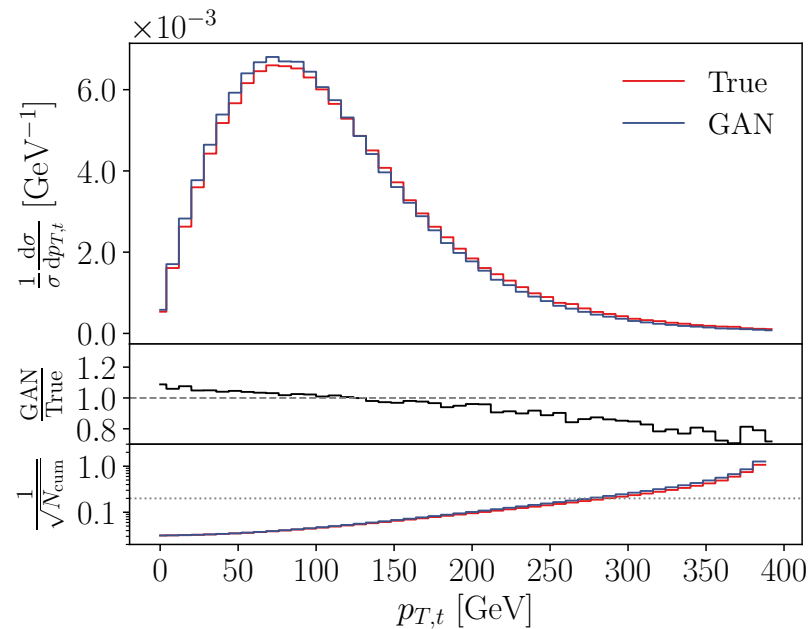
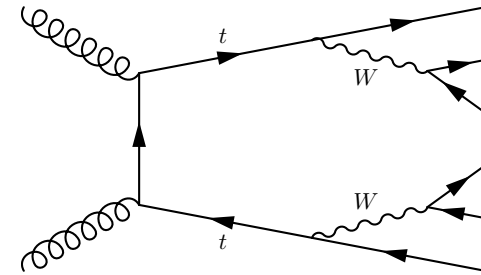


→ Amplification factor 2.5

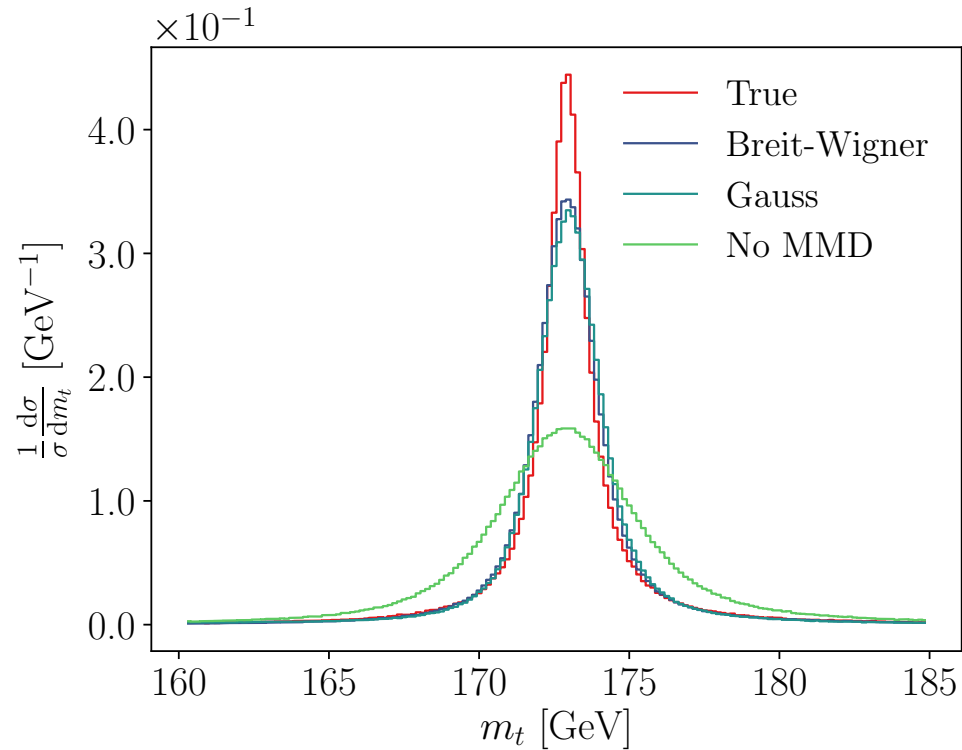
Sparser data → bigger amplification

How to GAN LHC events [1907.03764]

- $t\bar{t} \rightarrow 6$ quarks
- 18 dim output
 - external masses fixed
 - no momentum conservation
- + Flat observables ✓
- Systematic undershoot in tails [10-20% deviation]



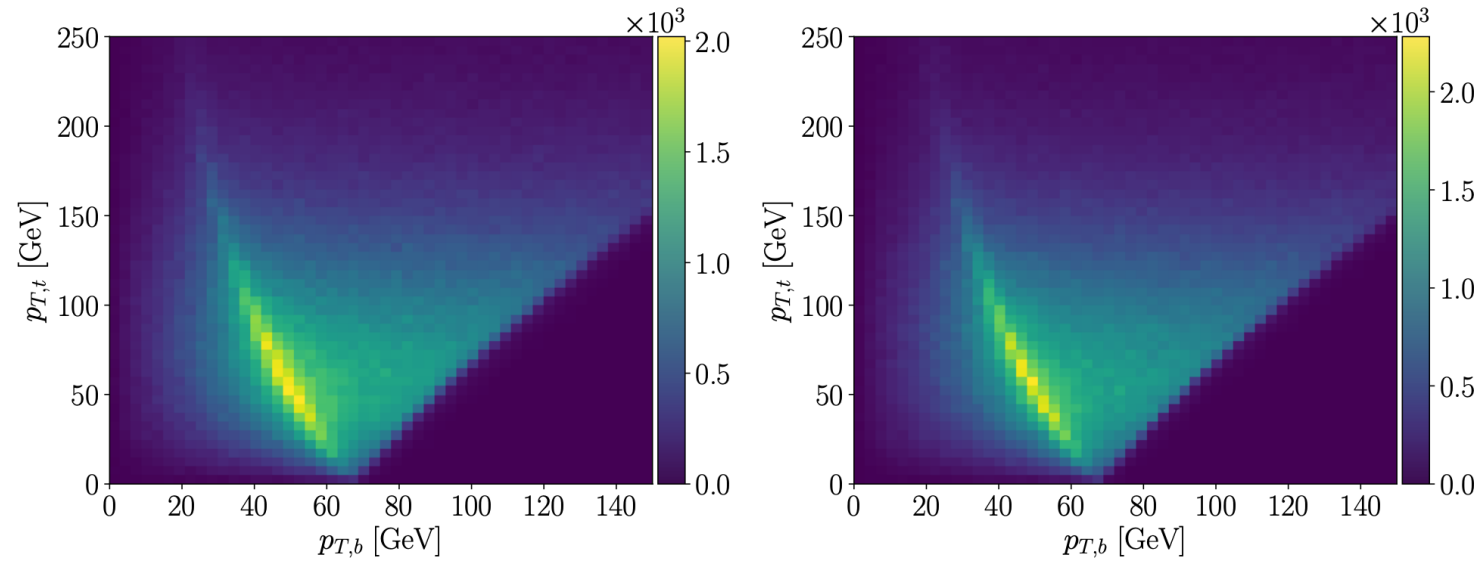
Special features



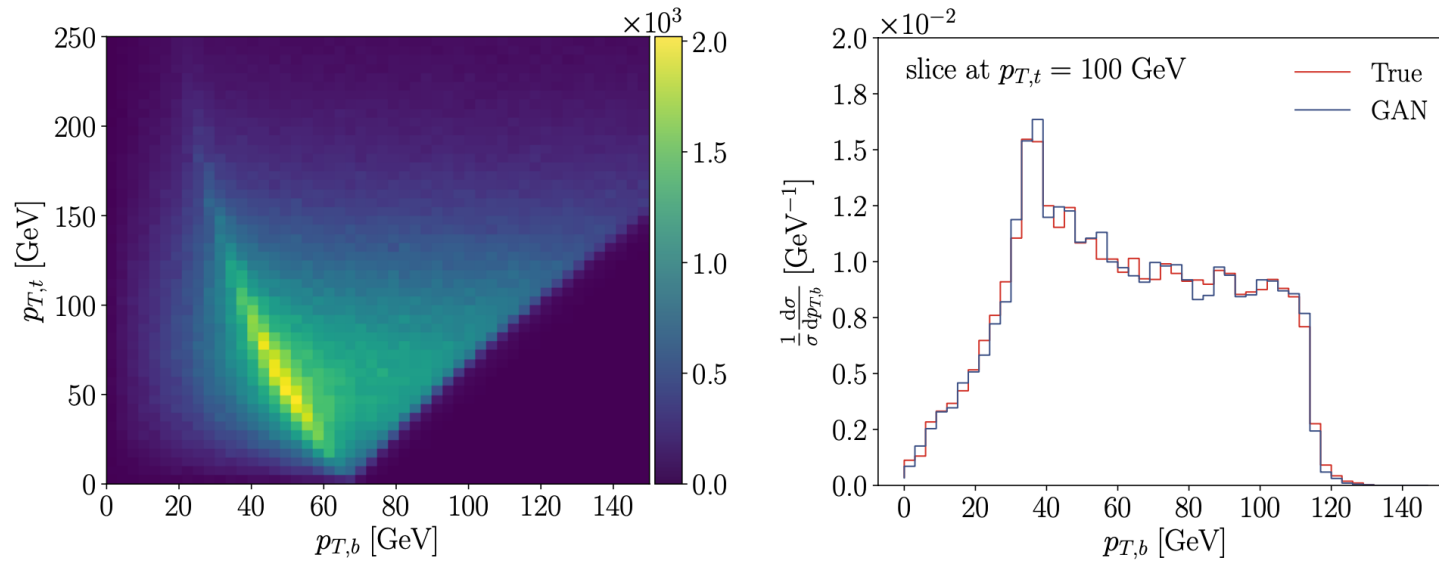
Solution: MMD kernel

$$\text{MMD}^2(P_T, P_G) = \langle k(x, x') \rangle_{x, x' \sim P_T} + \langle k(y, y') \rangle_{y, y' \sim P_G} - 2 \langle k(x, y) \rangle_{x \sim P_T, y \sim P_G}$$

Correlations



Correlations

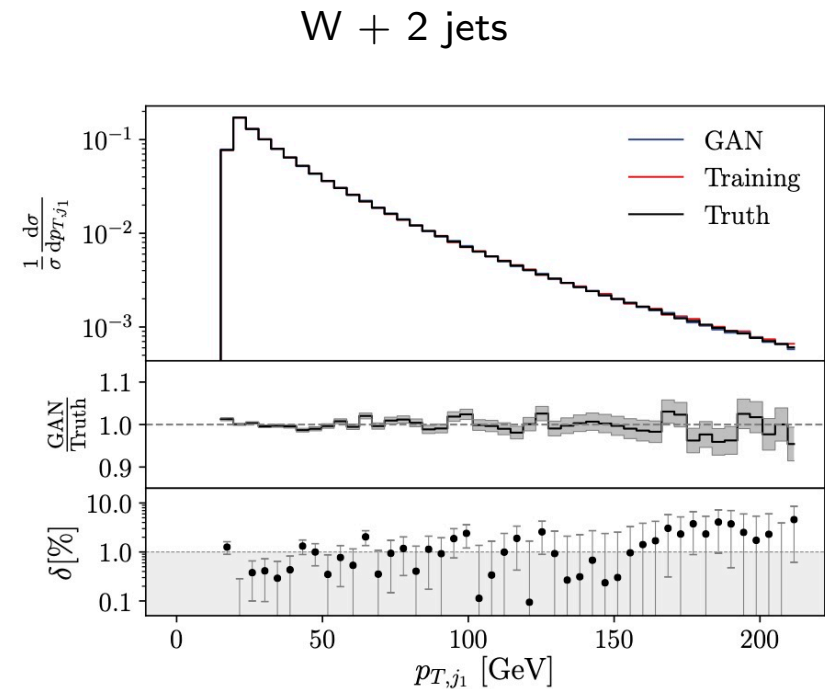


Reaching precision (preliminary)

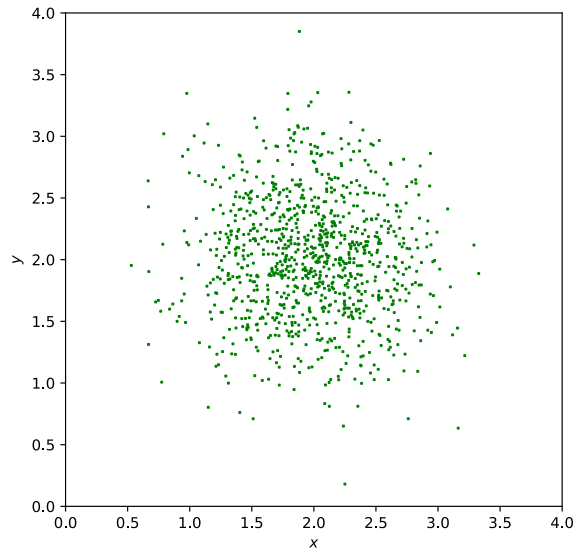
1. Representation p_T, η, ϕ
2. Momentum conservation
3. Resolve $\log p_T$
4. Regularization: spectral norm
5. Batch information

→ 1% precision ✓

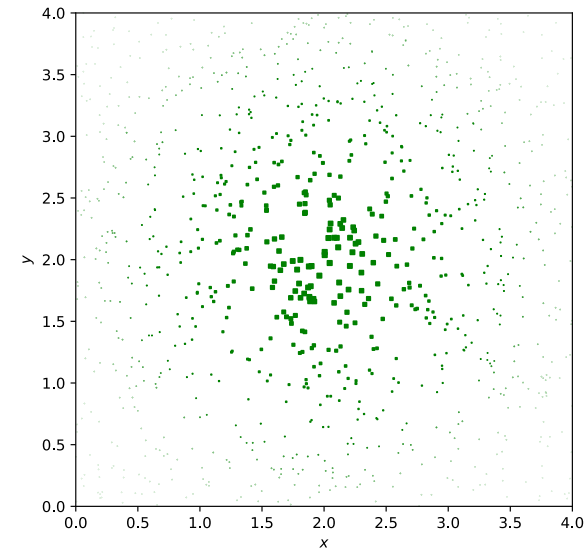
Next step automatization



Information in distributions



==



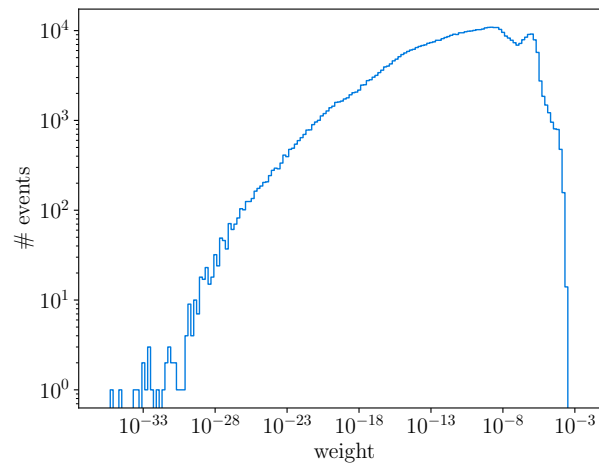
Information in space distribution
(what we want)

Information in weight
(what we have)

The unweighting bottleneck

- High-multiplicity / higher-order \rightarrow unweighting efficiencies $< 1\%$
- \rightarrow Simulate conditions with naive Monte Carlo generator
ME by Sherpa, parton densities from LHAPDF, Rambo-on-diet

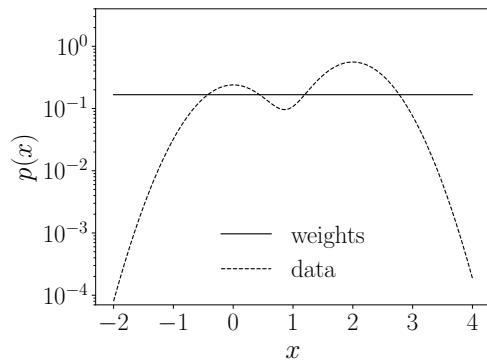
$pp \rightarrow \mu^+ \mu^-$ with $m_{\mu\mu} > 50$ GeV



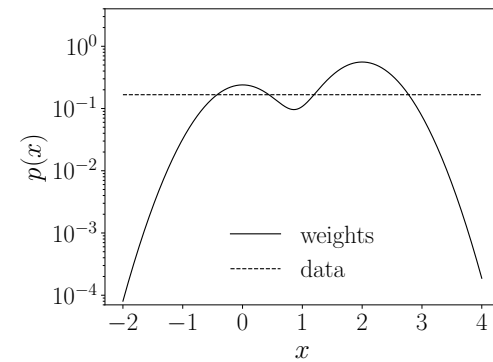
\rightarrow unweighting efficiency 0.2%

Training on weighted events

Information contained in distribution or event weights

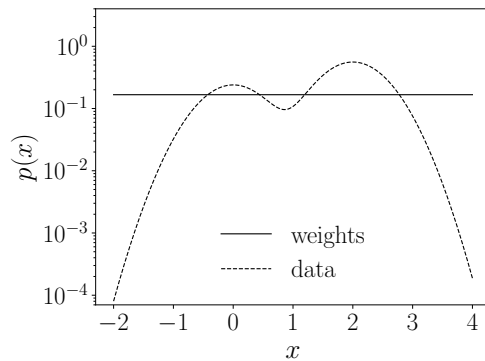


Train on
weighted events

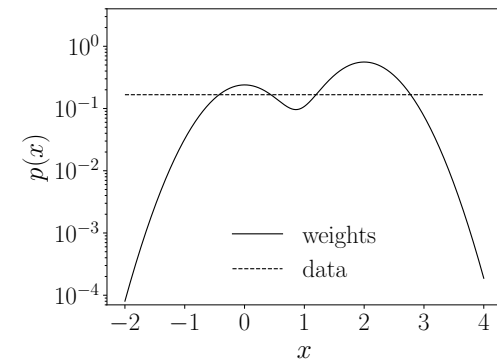


Training on weighted events

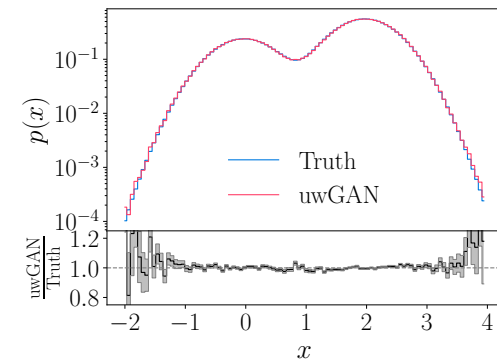
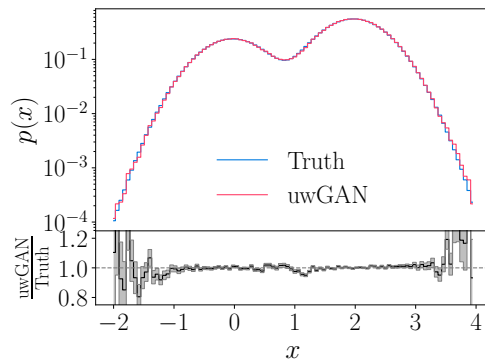
Information contained in distribution or event weights



Train on
weighted events



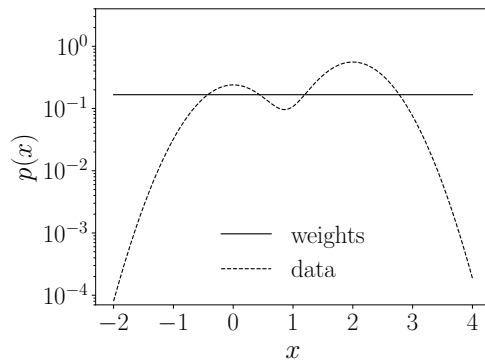
Generate
unweighted events



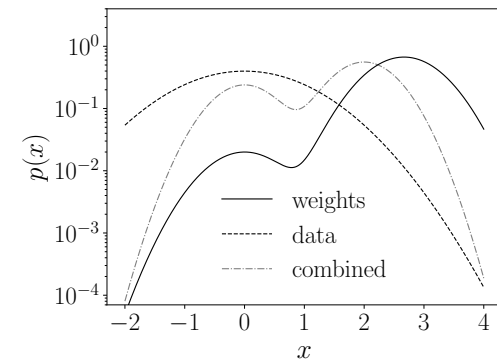
$$L_D = \langle -w \log D(x) \rangle_{x \sim P_{Truth}} + \langle -\log(1 - D(x)) \rangle_{x \sim P_{Gen}}$$

Training on weighted events

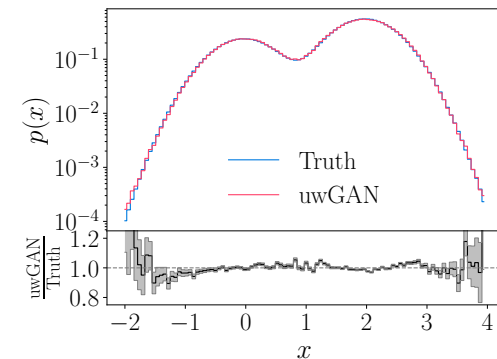
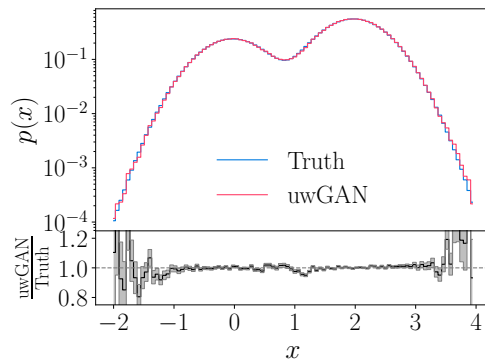
Information contained in distribution or event weights



Train on
weighted events



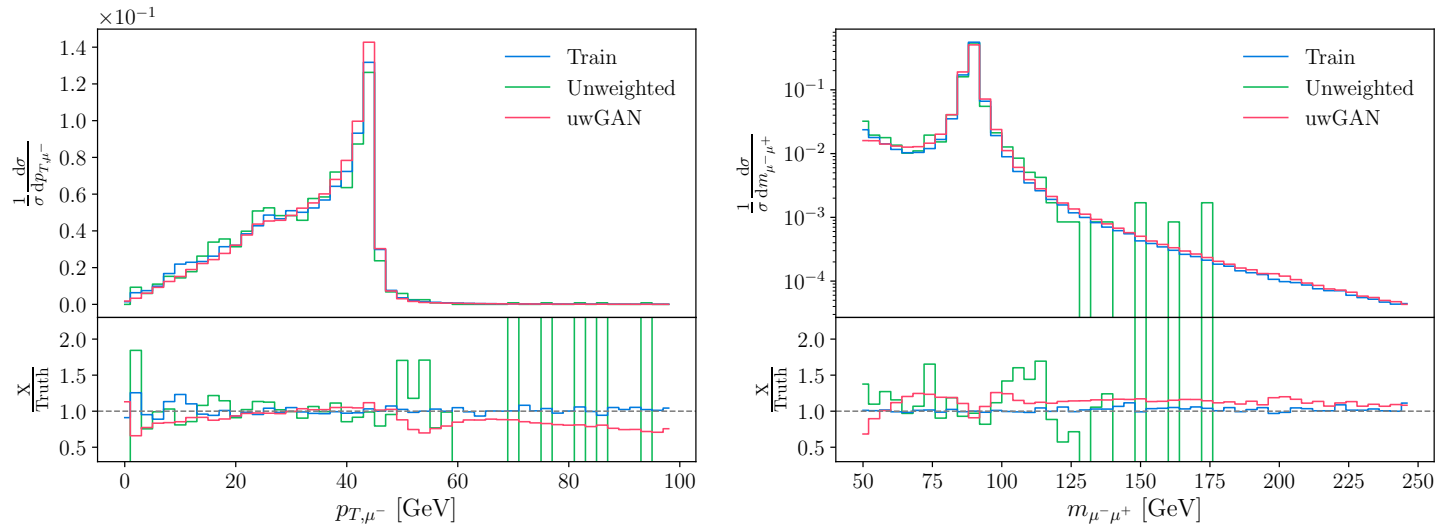
Generate
unweighted events



$$L_D = \langle -w \log D(x) \rangle_{x \sim P_{Truth}} + \langle -\log(1 - D(x)) \rangle_{x \sim P_{Gen}}$$

normalizing flow: B. Stienen, R. Verheyen [2011.13445]

uwGAN results

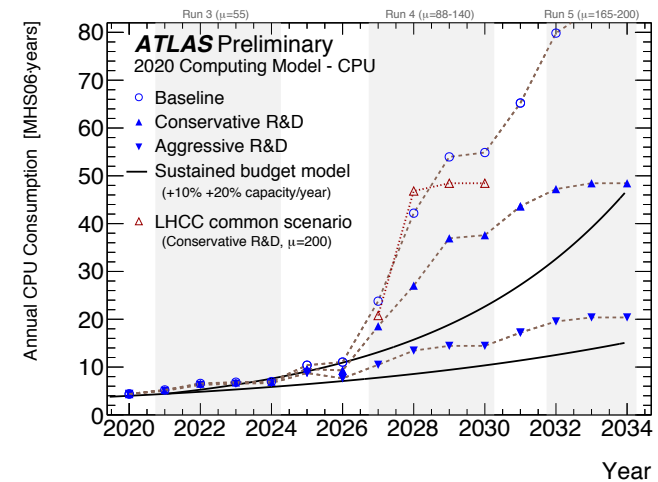
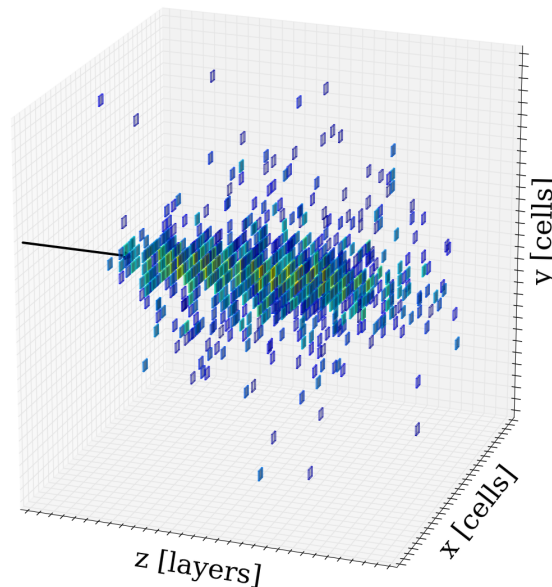


Populates high energy tails

Large amplification wrt. unweighted data!

Fast detector simulations

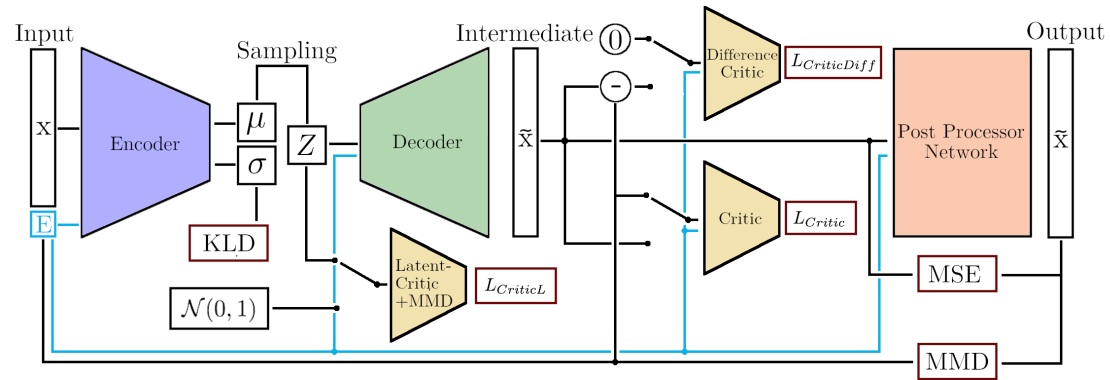
- Important R&D potential
NN evaluation $\times 100-1000$
faster than GEANT4



- Same underlying techniques [GAN, VAE, (NF)]
- Challenge: High-dimensional output $\leftarrow 30 \times 30 \times 30$

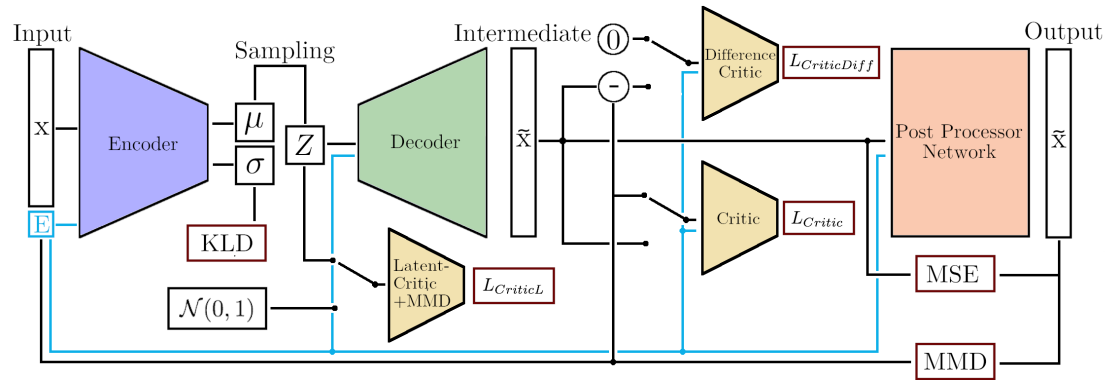
BIB-AE PP

Bounded-Information-Bottleneck autoencoder with post processing

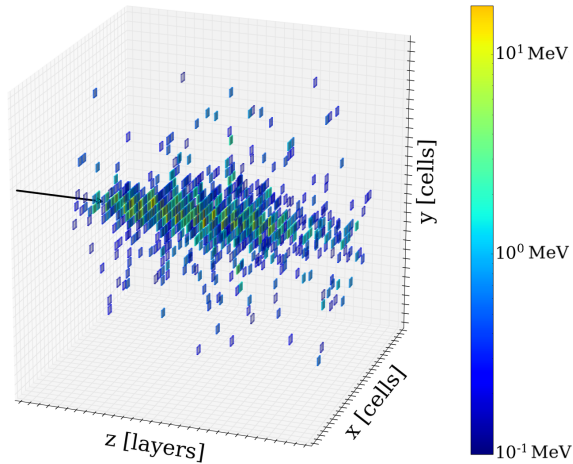


BIB-AE PP

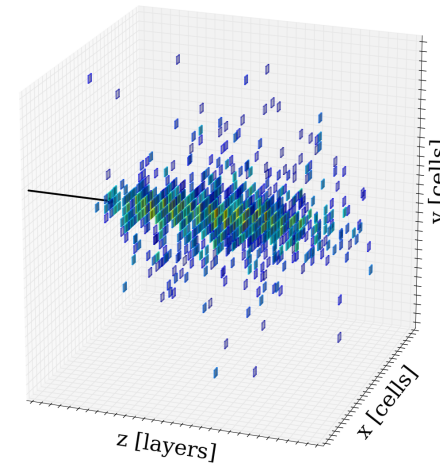
Bounded-Information-Bottleneck autoencoder with post processing



GEANT4

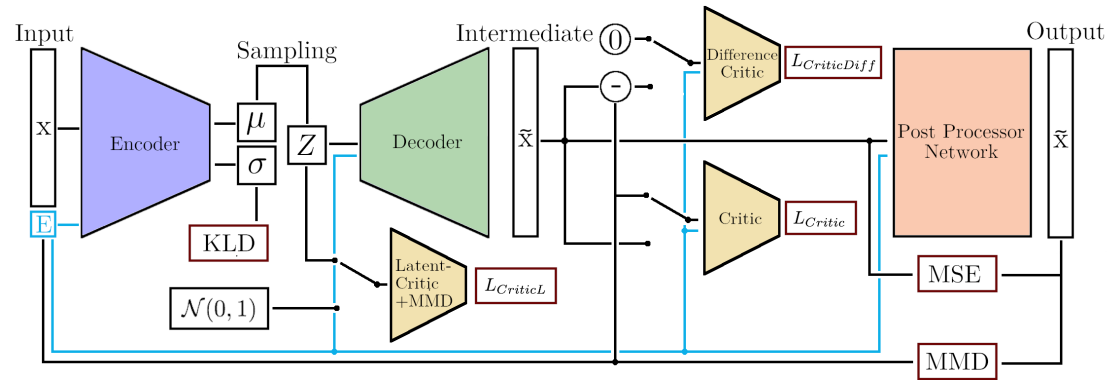


Simulation

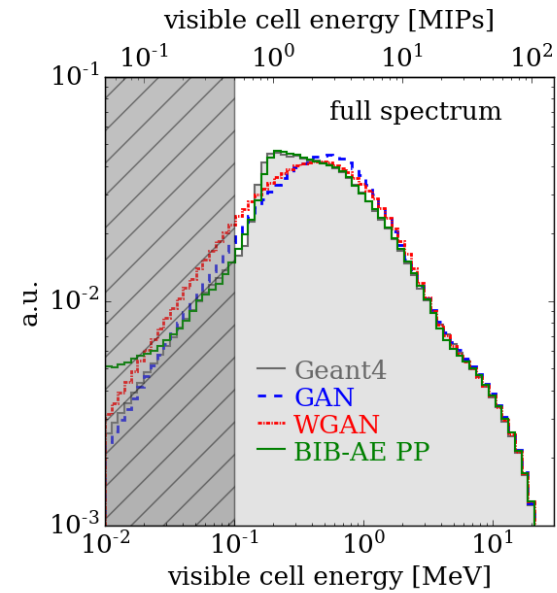


BIB-AE PP

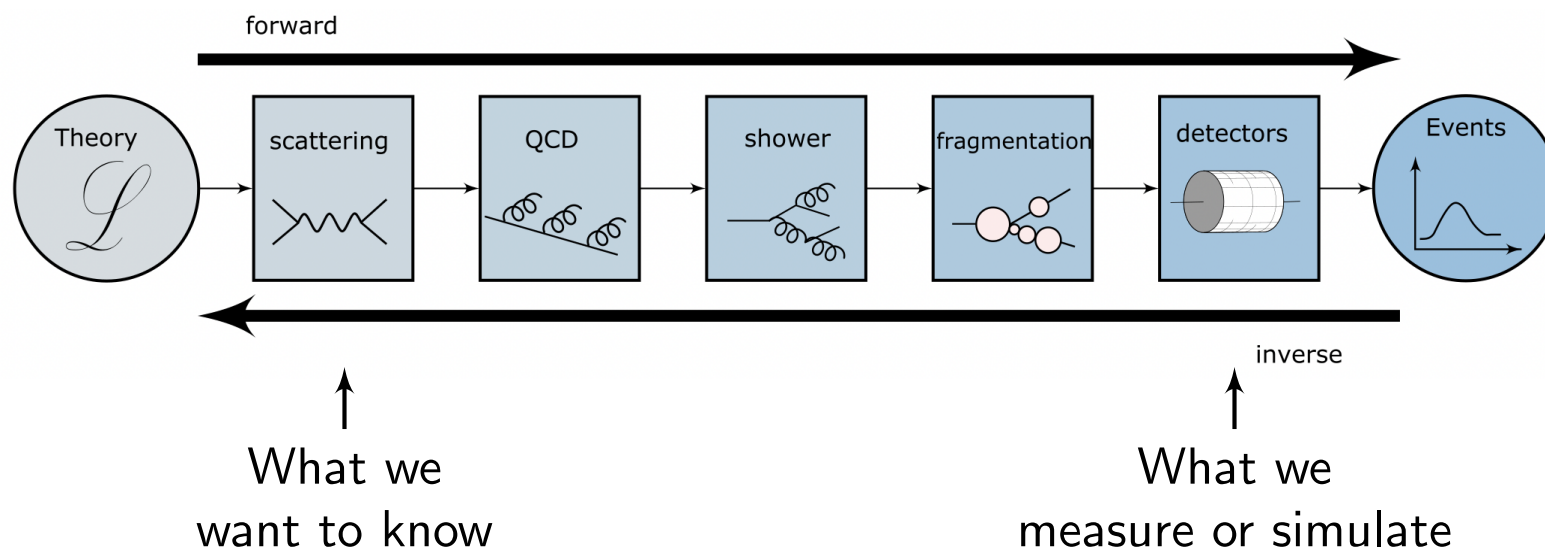
Bounded-Information-Bottleneck autoencoder with post processing



Post Processor Network adjusts energy to recover spectrum
→ MIP bump

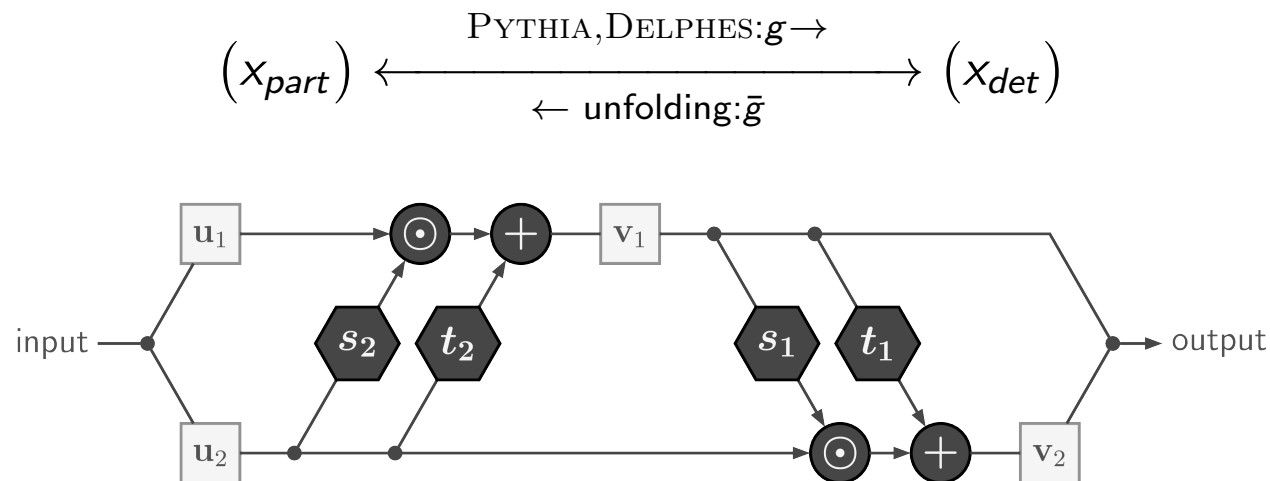


Can we invert the simulation chain?



- wish list:
- multi-dimensional
 - bin independent
 - statistically well defined

Invertible networks



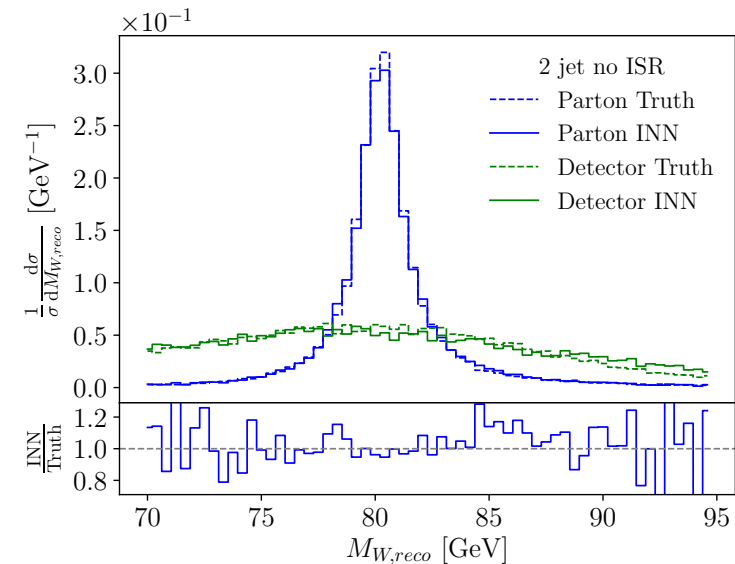
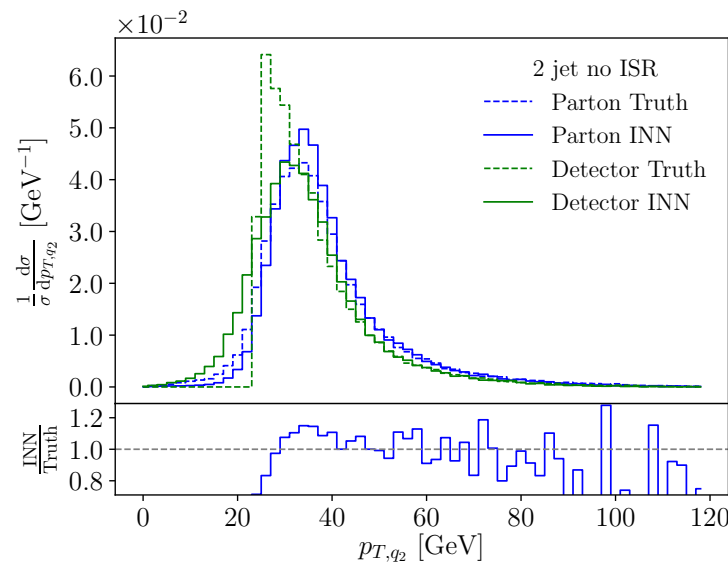
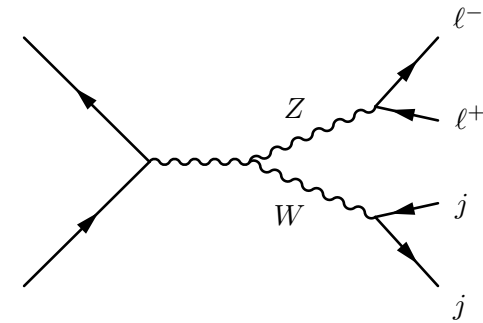
[1808.04730] L. Ardizzone, J. Kruse, S. Wirkert, D. Rahner,

E. W. Pellegrini, R. S. Klessen, L. Maier-Hein, C. Rother, U. Köthe

- + Bijective mapping
- + Tractable Jacobian
- + Fast evaluation in both directions
- + Arbitrary networks s and t

Inverting detector effects

- $pp \rightarrow ZW \rightarrow (\ell\ell)(jj)$
- Train: parton \rightarrow detector
- Evaluate: parton \leftarrow detector



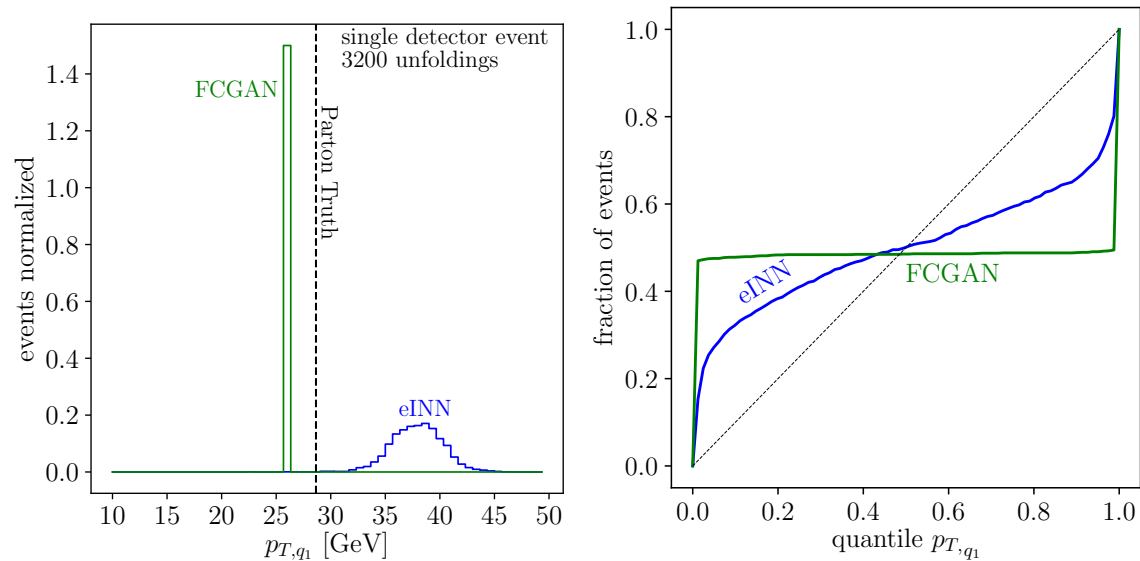
multi-dimensional ✓ bin independent ✓ statistically well defined ?

Including stochastic effects

$$\begin{pmatrix} x_p \\ r_p \end{pmatrix} \xleftarrow{\text{PYTHIA, DELPHES: } g \rightarrow} \begin{pmatrix} x_d \\ r_d \end{pmatrix} \xleftarrow{\text{unfolding: } \bar{g}}$$

Sample r_d for fixed detector event

How often is Truth included in distribution quantile?

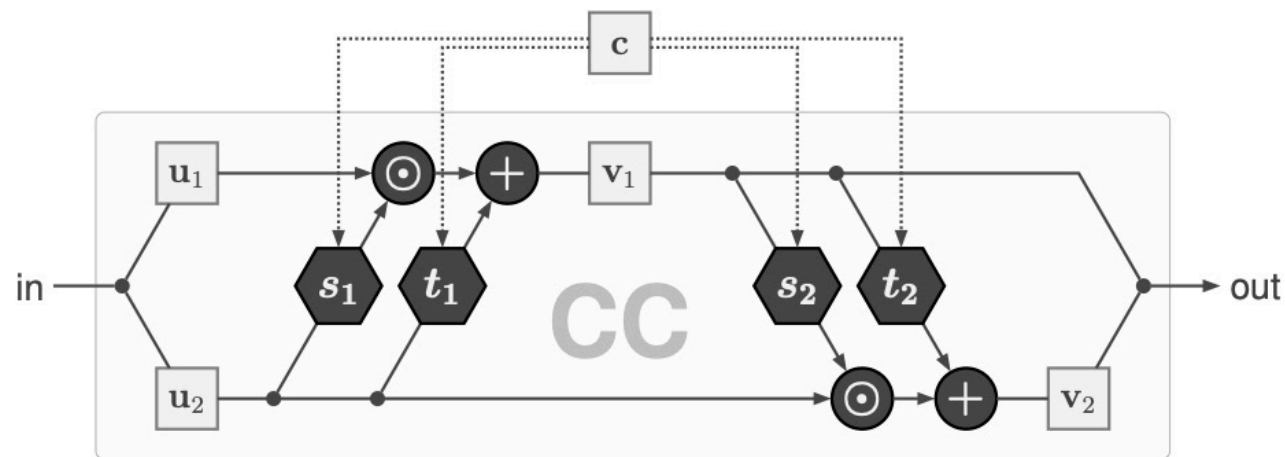


- Problem: arbitrary balance of many loss functions

Taking a different angle

Given an event x_d , what is the probability distribution at parton level?
 \rightarrow sample over r , condition on x_d

$$x_p \leftarrow \begin{array}{c} \xrightarrow{g(x_p, f(x_d))} \\ \xleftarrow{\text{unfolding: } \bar{g}(r, f(x_d))} \end{array} r$$



Taking a different angle

Given an event x_d , what is the probability distribution at parton level?
→ sample over r , condition on x_d

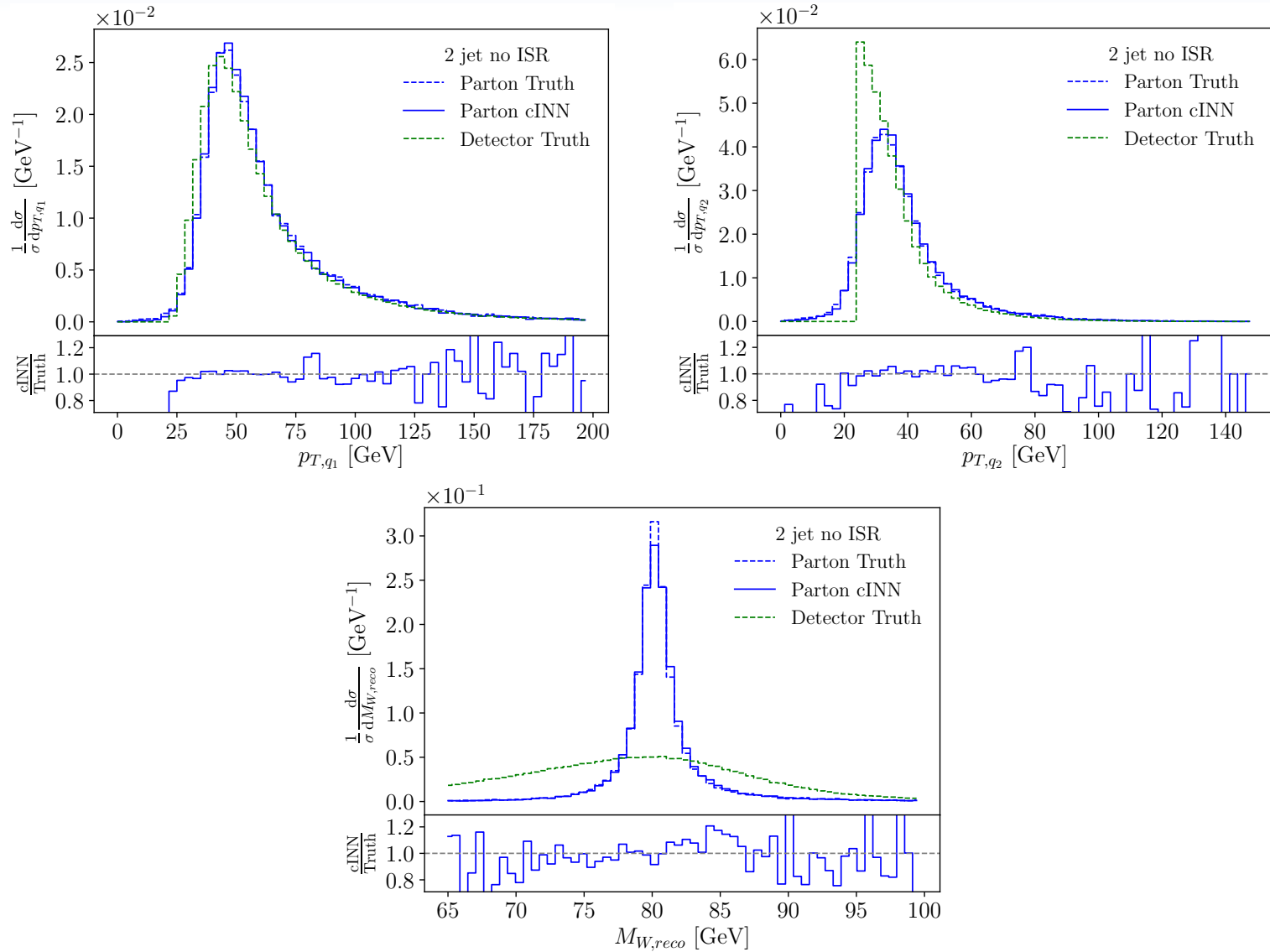
$$x_p \leftarrow \begin{array}{c} \xrightarrow{g(x_p, f(x_d))} \\ \xleftarrow{\text{unfolding: } \bar{g}(r, f(x_d))} \end{array} r$$

→ Training: Maximize posterior over model parameters

$$\begin{aligned} L &= - \langle \log p(\theta | x_p, x_d) \rangle_{x_p \sim P_p, x_d \sim P_d} \\ &= - \langle \log p(x_p | \theta, x_d) \rangle_{x_p \sim P_p, x_d \sim P_d} - \log p(\theta) + \text{const.} \quad \leftarrow \text{Bayes} \\ &= - \left\langle \log p(\bar{g}(x_p, x_d)) + \log \left| \frac{\partial \bar{g}(x_p, x_d)}{\partial x_p} \right| \right\rangle - \log p(\theta) \quad \leftarrow \text{change of var} \\ &= \langle 0.5 \| \bar{g}(x_p, f(x_d)) \|_2^2 - \log |J| \rangle_{x_p \sim P_p, x_d \sim P_d} - \log p(\theta) \end{aligned}$$

→ Jacobian of bijective mapping

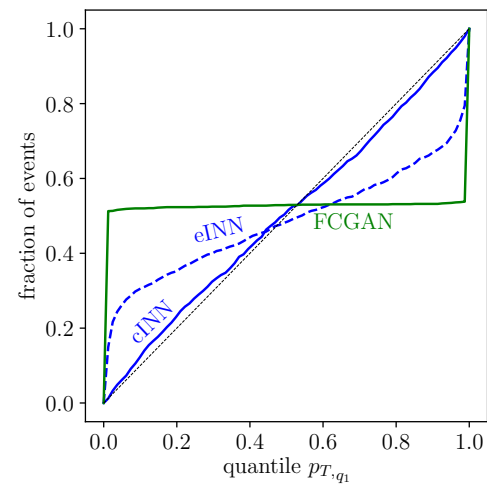
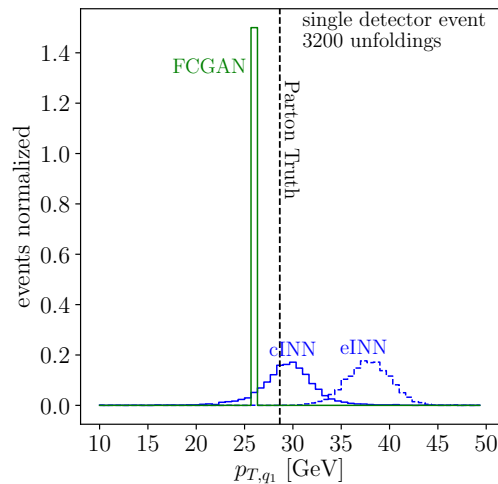
Cross check distributions



Condition INN on detector data [2006.06685]

$$\begin{array}{c}
 g(x_p, f(x_d)) \rightarrow \\
 \leftarrow \text{unfolding: } \bar{g}(r, f(x_d)) \\
 x_p \leftarrow \text{-----} \rightarrow r
 \end{array}$$

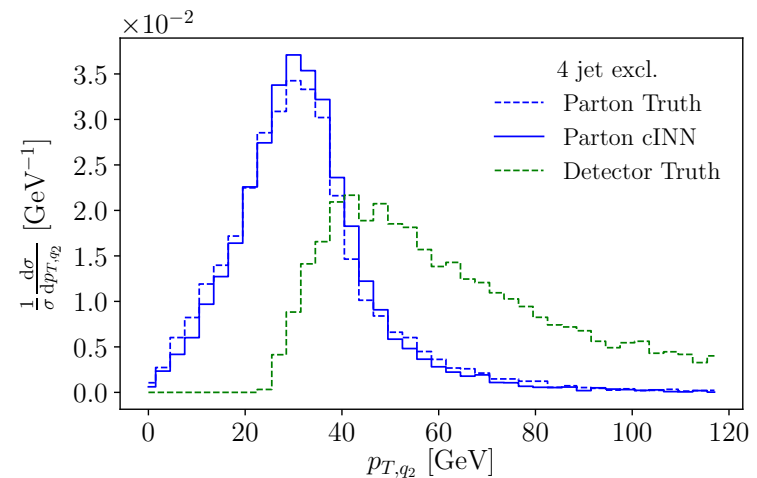
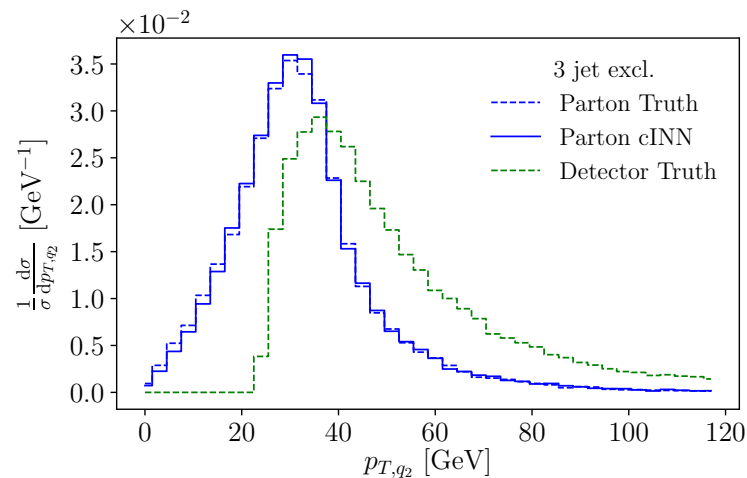
$$\text{Minimizing } L = \langle 0.5 \|\bar{g}(x_p, f(x_d))\|_2^2 - \log |J| \rangle_{x_p \sim P_p, x_d \sim P_d} - \log p(\theta)$$



multi-dimensional ✓ bin independent ✓ statistically well defined ✓

Inverting the full event I

- $pp > WZ > q\bar{q}l^+l^- + \text{ISR}$
- ISR leads to large fraction of 2/3/4 jet events
- Train and test on exclusive channels

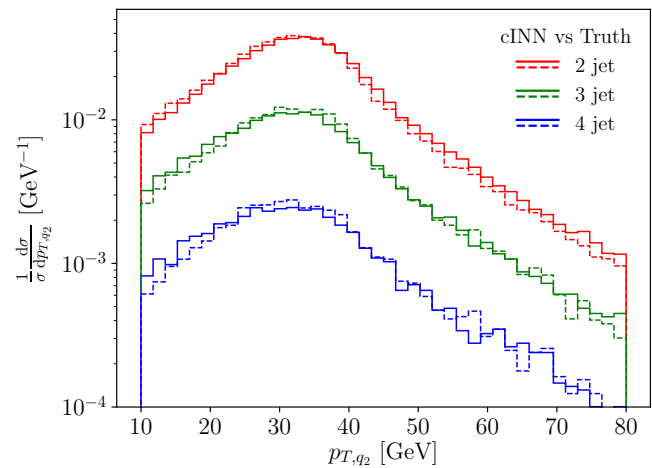
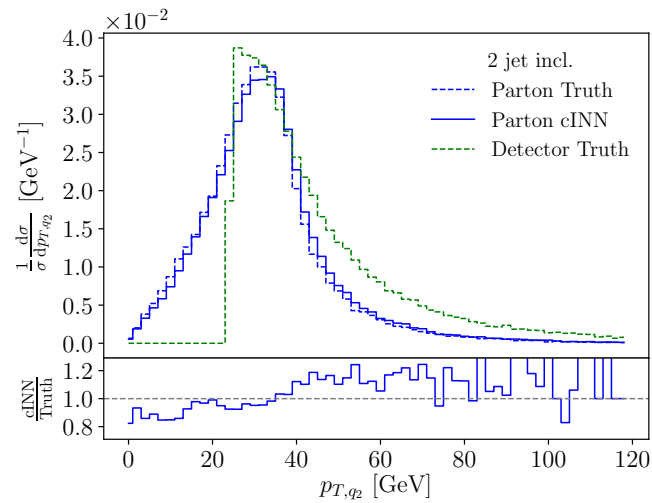


Inverting the full event II

$pp > WZ > q\bar{q}l^+l^- + \text{ISR}$

Train on inclusive dataset

Evaluate
exclusive 2/3/4 jet channels



We can use ML ...

... to enable precision simulations in forward direction

... to turn weighted into unweighted events

... to invert the simulation chain statistically

... for fun and precision :)