

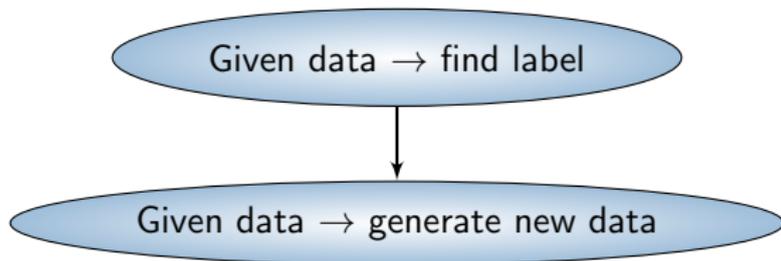
Generative networks for particle physics

Lecture at QU Data Science Basics (Hamburg)

June 2021

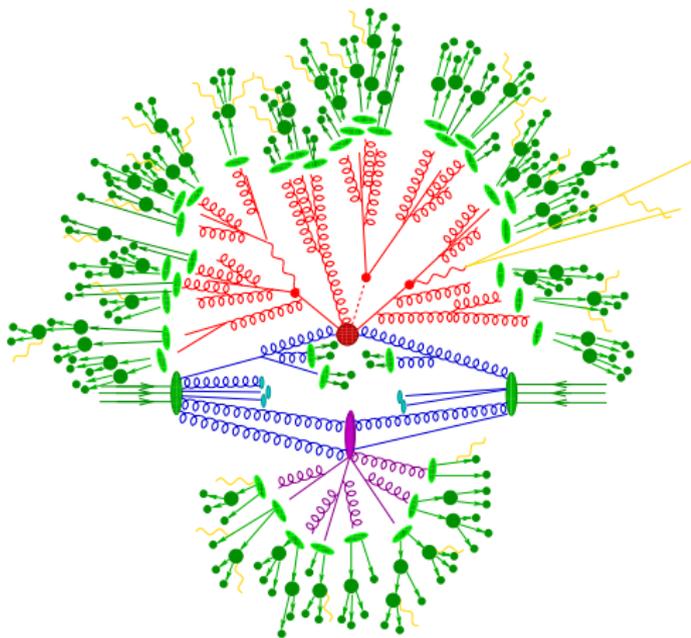
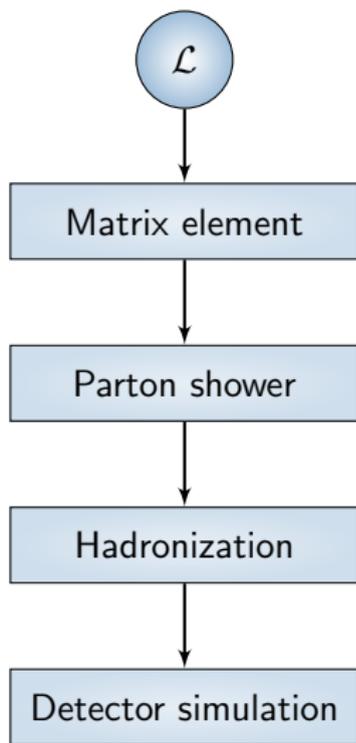
Anja Butter

Paradigm shift



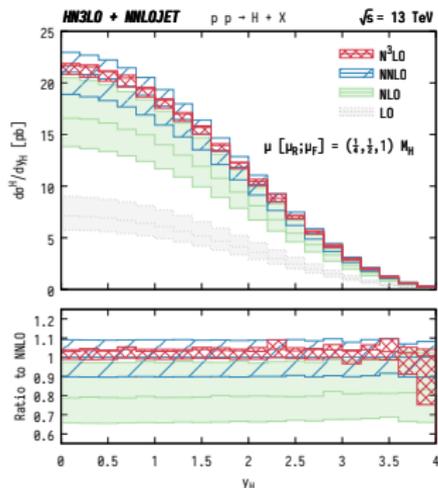
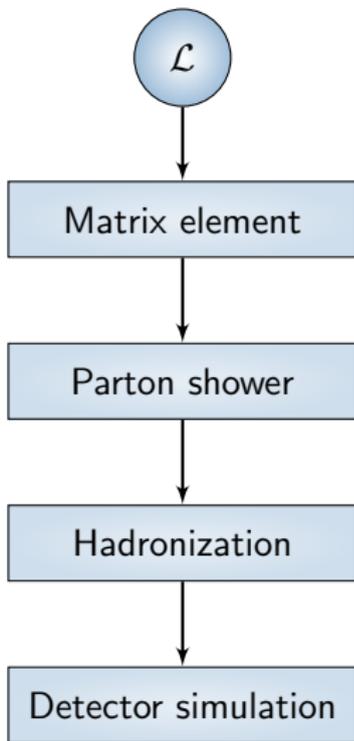
What do we generate in high energy physics?

First principle based event generation



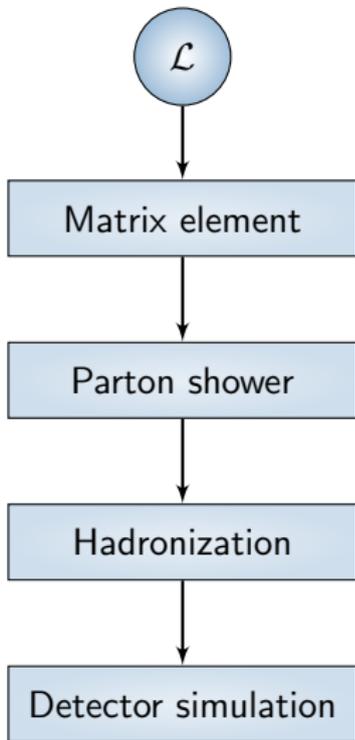
a sherpa artist

First principle based event generation

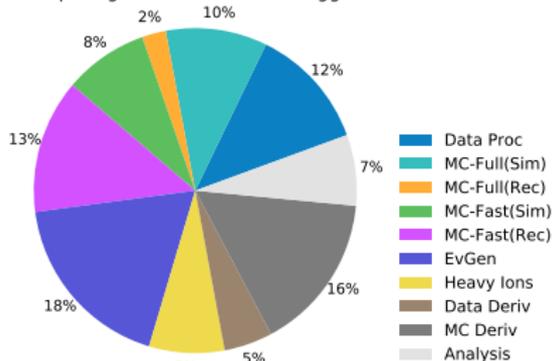


[1807.11501] Cieri, Chen, Gehrmann, Glover, Huss

First principle based event generation

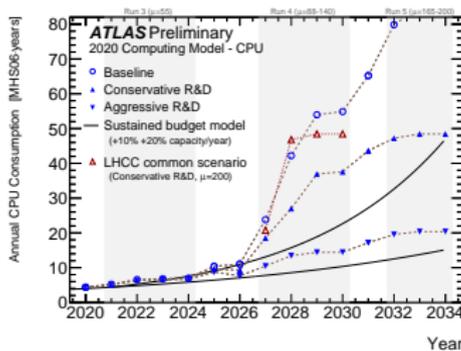
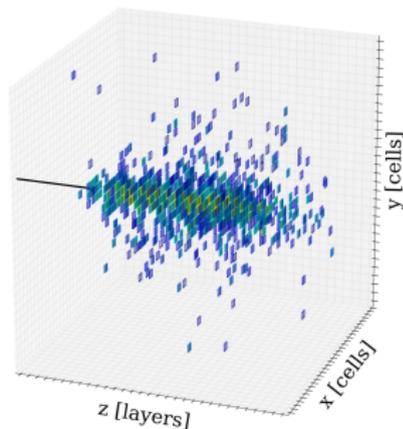


ATLAS Preliminary
2020 Computing Model -CPU: 2030: Aggressive R&D



Fast detector simulations

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



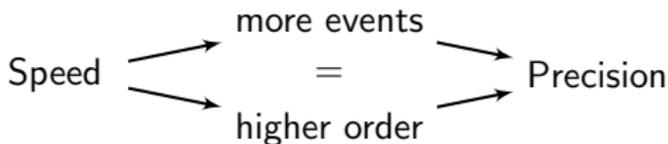
- Challenge:
High-dimensional output
 $\leftarrow 30 \times 30 \times 30$

Why do we need machine learning for data simulation?

precision
event generation

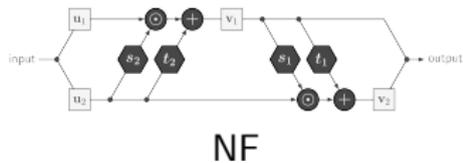
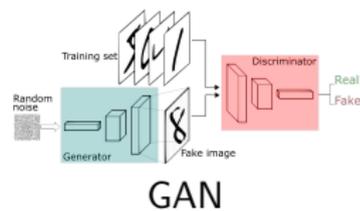
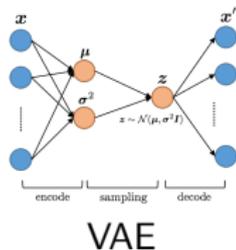
detailed
detector simulation

limited computing resources

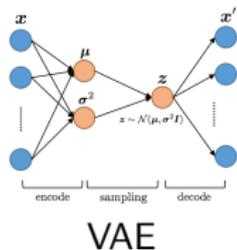


Use NN to speed up simulations!

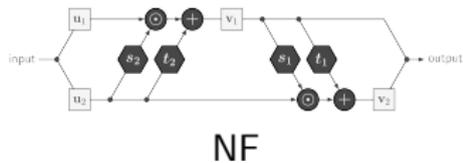
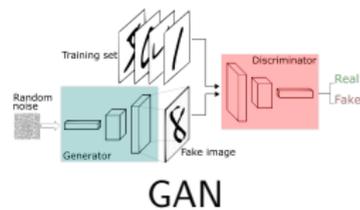
Neural network based generative networks



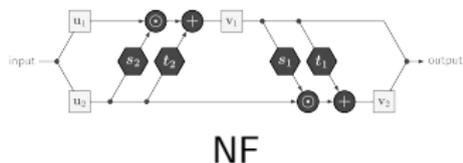
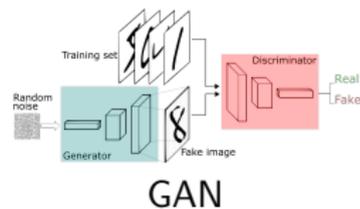
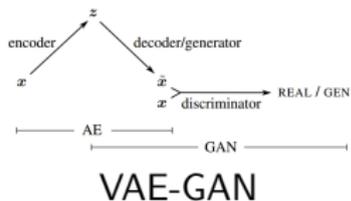
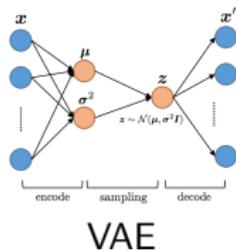
Neural network based generative networks



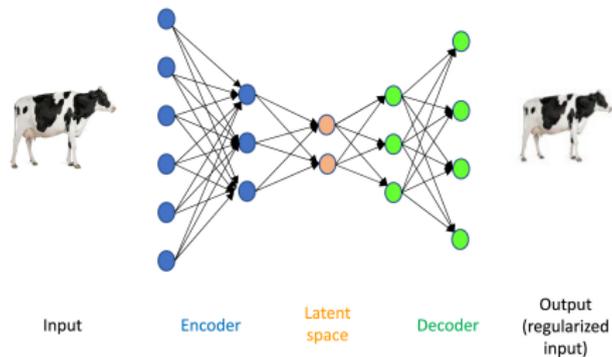
all kinds of hybrids



Neural network based generative networks

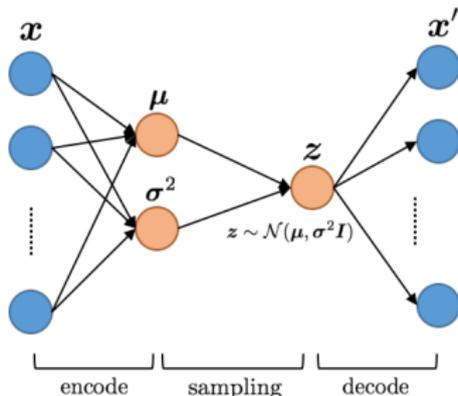


From autoencoders to variational autoencoders



$$\text{AE: } \mathbf{x} \xrightarrow{\text{encoder}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{AE}} = (\mathbf{x} - \mathbf{x}')^2$$

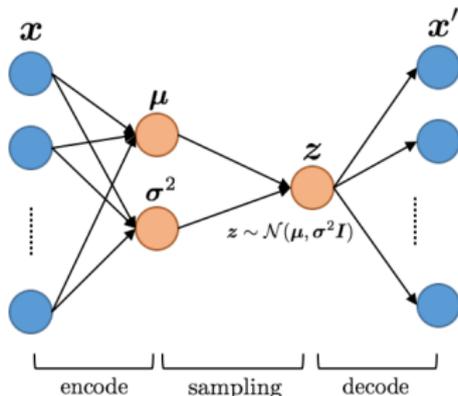
From autoencoders to variational autoencoders



$$\text{AE: } \mathbf{x} \xrightarrow{\text{encoder}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{AE}} = (\mathbf{x} - \mathbf{x}')^2$$

$$\text{VAE: } \mathbf{x} \xrightarrow{\text{encoder}} \begin{pmatrix} \mu \\ \sigma \end{pmatrix} \xrightarrow[\mathbf{z} \sim \mathcal{N}(\mu, \sigma)]{\text{sample}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \mathcal{L}_{\text{lat}}$$

From autoencoders to variational autoencoders



$$\text{AE: } \mathbf{x} \xrightarrow{\text{encoder}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{AE}} = (\mathbf{x} - \mathbf{x}')^2$$

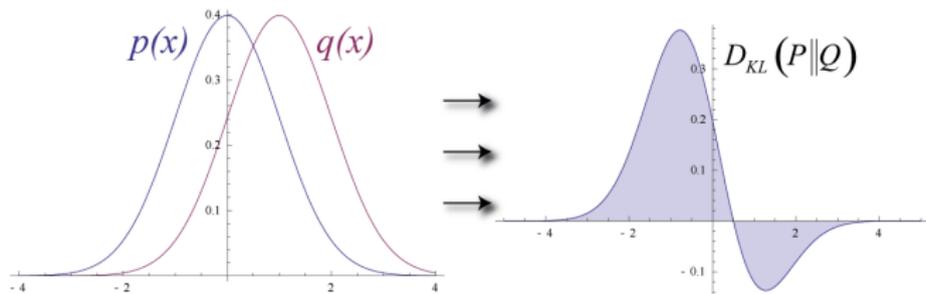
$$\text{VAE: } \mathbf{x} \xrightarrow{\text{encoder}} \begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\sigma} \end{pmatrix} \xrightarrow[\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})]{\text{sample}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \mathcal{L}_{\text{lat}}$$

Loss enforces Gaussian latent space

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \beta \cdot \text{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}) | \mathcal{N}(0, 1)) \quad \leftarrow \text{similarity measure}$$

Interlude: KL Divergence

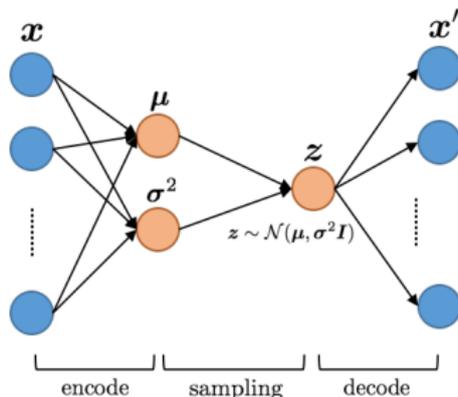
- Distance measure for probability distributions P and Q
- $D_{\text{KL}}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$



By Mundhenk at English Wikipedia, CC BY-SA 3.0

$$D_{\text{KL}}(\mathcal{N}(\mu, \sigma) \parallel \mathcal{N}(0, 1)) = \frac{1}{2}(1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$

From autoencoders to variational autoencoders



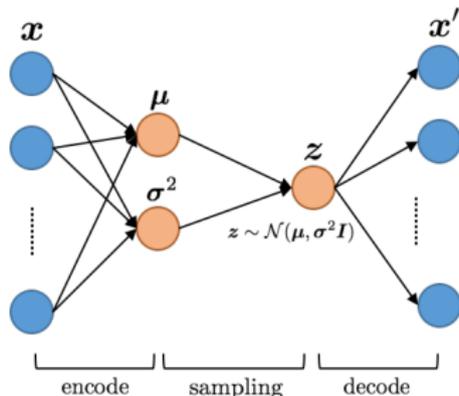
$$\text{AE: } \mathbf{x} \xrightarrow{\text{encoder}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{AE}} = (\mathbf{x} - \mathbf{x}')^2$$

$$\text{VAE: } \mathbf{x} \xrightarrow{\text{encoder}} \begin{pmatrix} \mu \\ \sigma \end{pmatrix} \xrightarrow[\mathbf{z} \sim \mathcal{N}(\mu, \sigma)]{\text{sample}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \mathcal{L}_{\text{lat}}$$

Loss enforces Gaussian latent space

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \beta \cdot \text{KL}(\mathcal{N}(\mu, \sigma) | \mathcal{N}(0, 1)) \quad \leftarrow \text{similarity measure}$$

From autoencoders to variational autoencoders



$$\text{AE: } \mathbf{x} \xrightarrow{\text{encoder}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{AE}} = (\mathbf{x} - \mathbf{x}')^2$$

$$\text{VAE: } \mathbf{x} \xrightarrow{\text{encoder}} \begin{pmatrix} \boldsymbol{\mu} \\ \boldsymbol{\sigma} \end{pmatrix} \xrightarrow[\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma})]{\text{sample}} \mathbf{z} \xrightarrow{\text{decoder}} \mathbf{x}' \quad \mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \mathcal{L}_{\text{lat}}$$

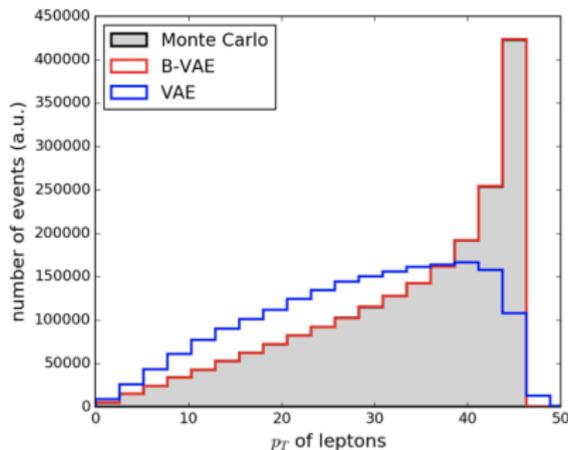
Loss enforces Gaussian latent space

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{AE}} + \beta \cdot \text{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}) | \mathcal{N}(0, 1)) \quad \leftarrow \text{similarity measure}$$

$$= \mathcal{L}_{\text{AE}} + \frac{\beta}{2} \sum_j 1 + \log(\sigma_j^2) - \mu_j^2 - \sigma_j^2$$

$e^+e^- \rightarrow Z \rightarrow l^+l^-$ with VAE

[1901.00875] S. Otten et al.



naive VAE fails to reproduce distributions

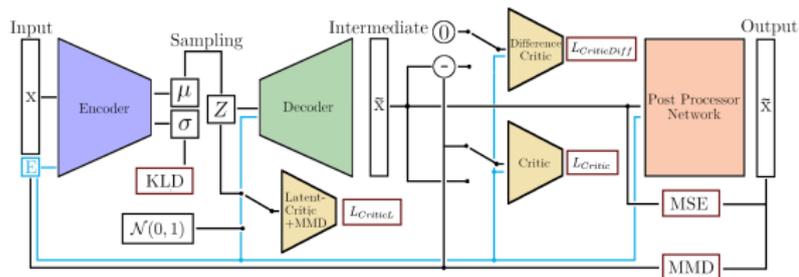
Why? \rightarrow latent space not perfectly Gaussian

Fix: insert information **buffer** to sample from real latent distribution

\rightarrow **B-VAE** shows excellent performance

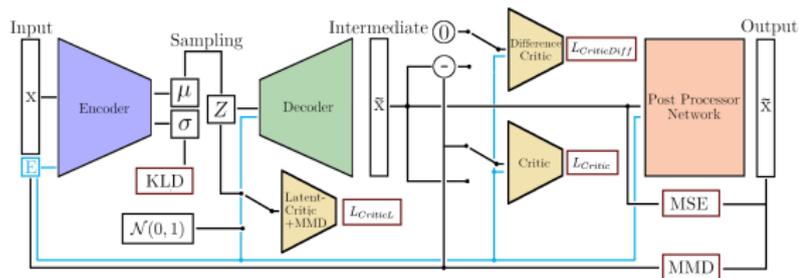
Detector simulation with BIB-AE PP

Bounded-Information-Bottleneck autoencoder with post processing

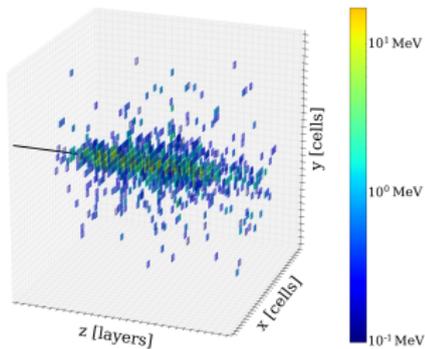


Detector simulation with BIB-AE PP

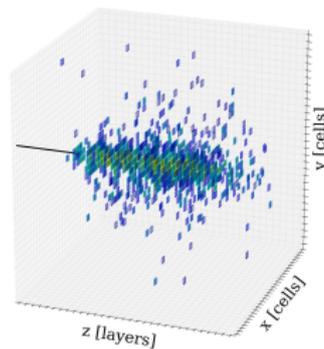
Bounded-Information-Bottleneck autoencoder with post processing



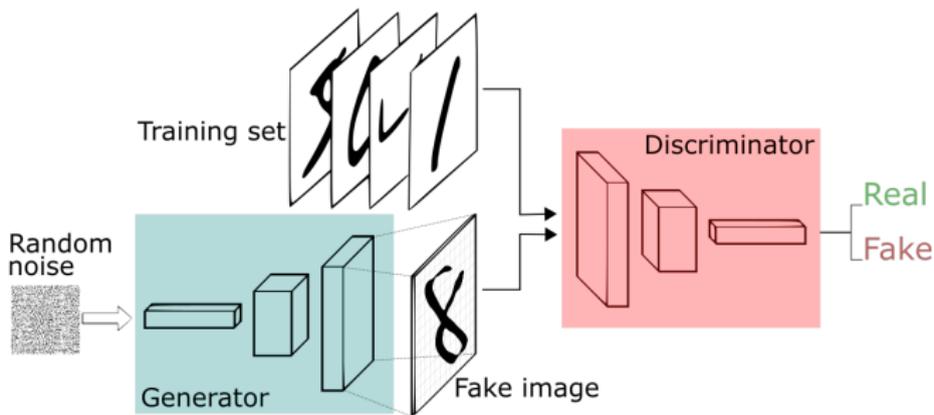
GEANT4



Simulation



Generative Adversarial Networks



Discriminator

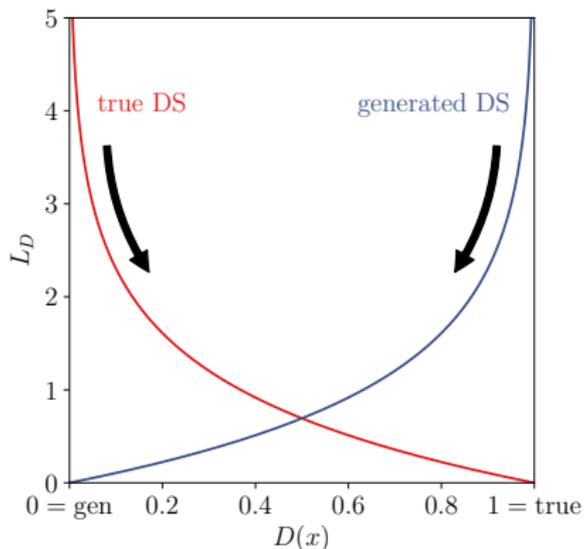
$$L_D = \langle -\log D(x) \rangle_{x \sim P_{Truth}} + \langle -\log(1 - D(x)) \rangle_{x \sim P_{Gen}}$$

Generator

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

Training the Discriminator

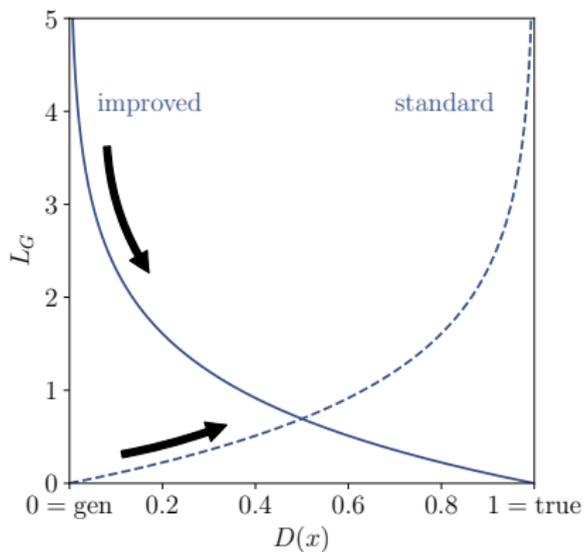
Discriminator loss



$$\text{Minimize } L_D = \langle -\log D(x) \rangle_{x \sim P_T} + \langle -\log(1 - D(x)) \rangle_{x \sim P_G}$$

Training the Generator

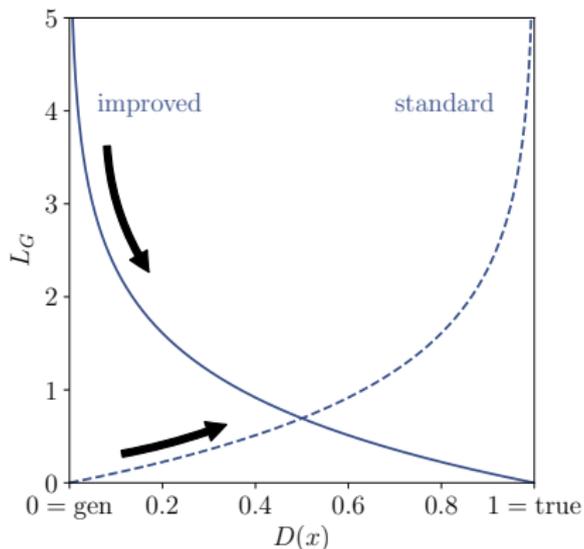
Generator loss



$$\text{Maximize } L_G = \langle -\log(1 - D(x)) \rangle_{x \sim P_G}$$

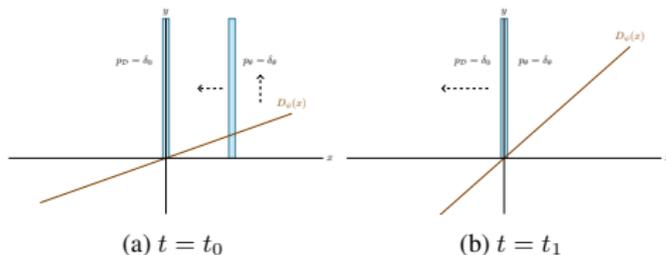
Training the Generator

Generator loss



$$\text{Minimize } L_G = \langle -\log D(x) \rangle_{x \sim P_G}$$

Regularization



[1801.04406]

Adding gradient penalty

$$\phi(x) = \log \frac{D(x)}{1 - D(x)} \quad \Rightarrow \quad \frac{\partial \phi}{\partial x} = \frac{1}{D(x)} \frac{1}{1 - D(x)} \frac{\partial D}{\partial x}$$

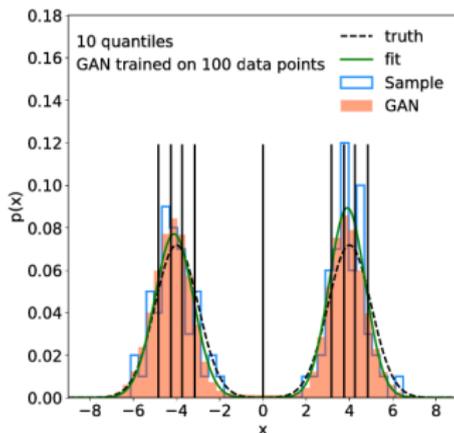
$$L_D \rightarrow L_D + \lambda_D \langle (1 - D(x))^2 |\nabla \phi|^2 \rangle_{x \sim P_T} + \lambda_D \langle D(x)^2 |\nabla \phi|^2 \rangle_{x \sim P_G}$$

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

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

Evaluation on quantiles:

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

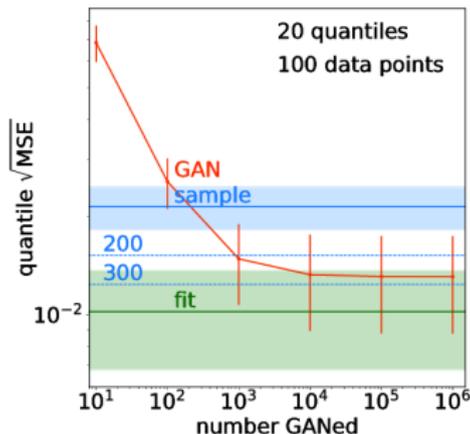


What is the statistical value of GANned events? [2008.06545]

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

Evaluation on quantiles:

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



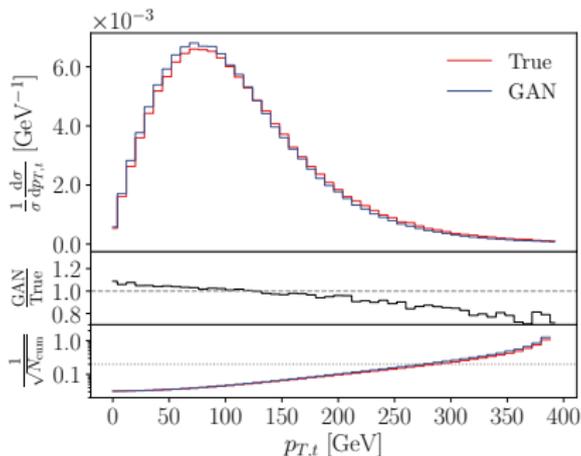
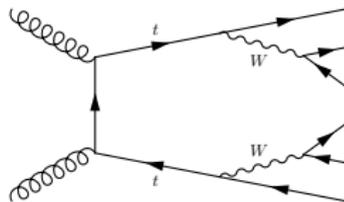
→ Amplification factor 2.5

Sparser data → bigger amplification

How to GAN LHC events

[1907.03764]

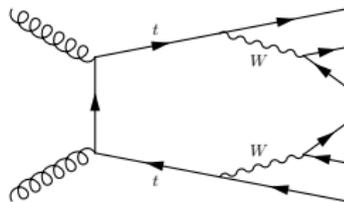
- $t\bar{t} \rightarrow 6$ quarks
 - 18 dim output
 - external masses fixed
 - no momentum conservation
- + Flat observables ✓
- Systematic undershoot in tails



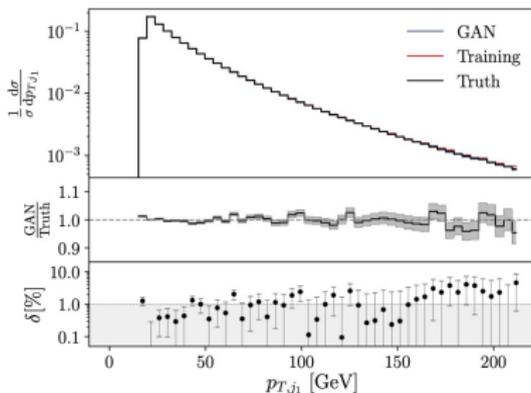
How to GAN LHC events

[1907.03764]

- $t\bar{t} \rightarrow 6$ quarks
- 18 dim output
 - external masses fixed
 - no momentum conservation



- + Flat observables ✓
- Systematic undershoot in tails
→ improve network (symmetries, preprocessing, ...) ✓



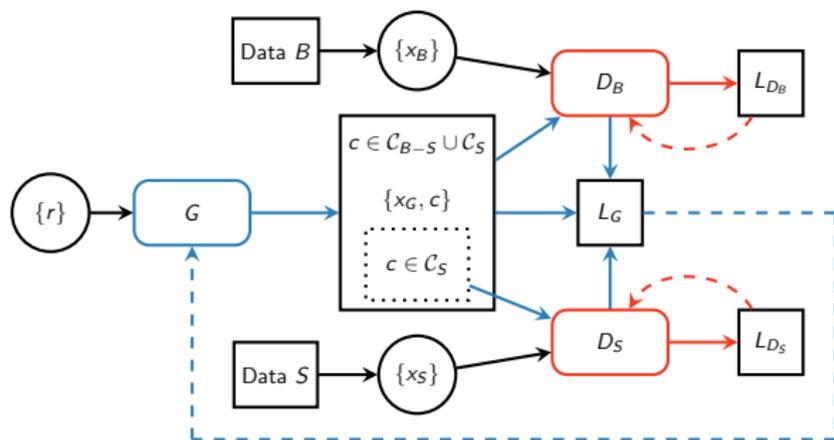
Generating the high-dim. difference of distributions

[1912.08824]

- Necessary to include negative events
- Beat bin-induced uncertainty

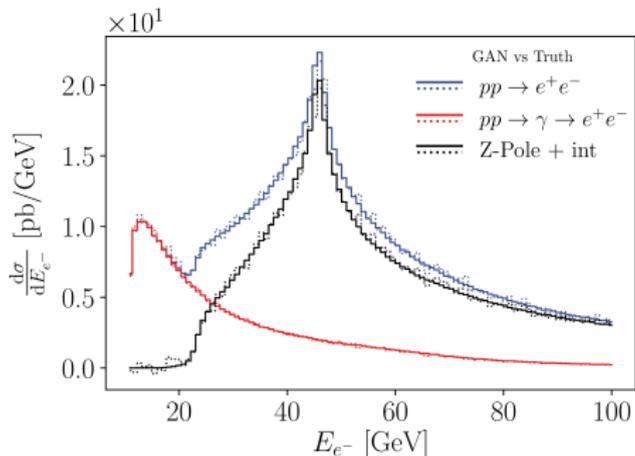
$$\Delta_{B-S} > \max(\Delta_B, \Delta_S)$$

- Applications:
 - Background subtraction, soft-collinear subtraction, ...



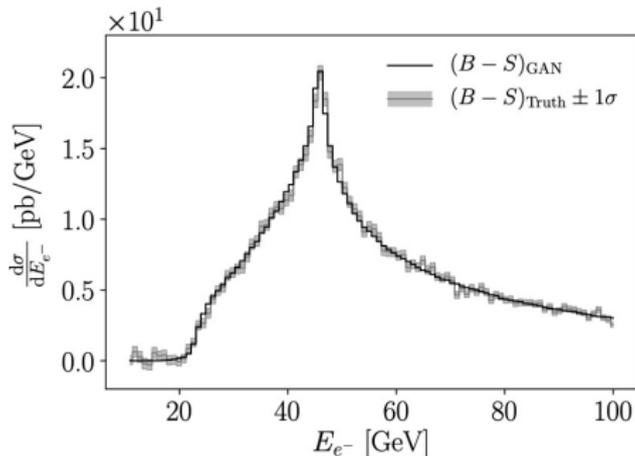
Generative background subtraction

- Training data:
 - $pp \rightarrow e^+e^-$
 - $pp \rightarrow \gamma \rightarrow e^+e^-$
- Generated events: Z-Pole + interference

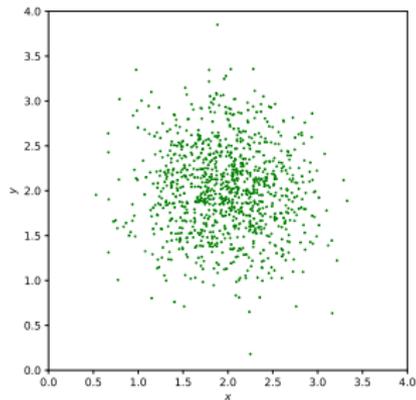


Generative background subtraction

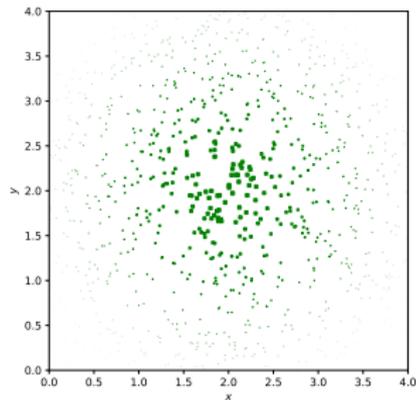
- Training data:
 - $pp \rightarrow e^+e^-$
 - $pp \rightarrow \gamma \rightarrow e^+e^-$
- Generated events: Z-Pole + interference



Information in distributions



=

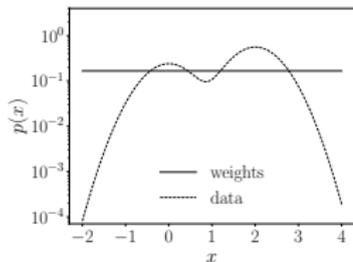


Information in space distribution
(what we want)

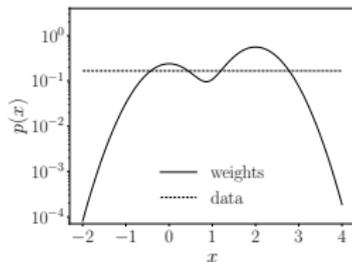
Information in weight
(what we have)

Training on weighted events

Information contained in distribution or event weights

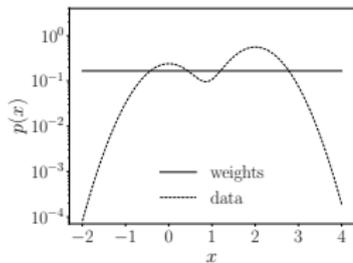


Train on
weighted events

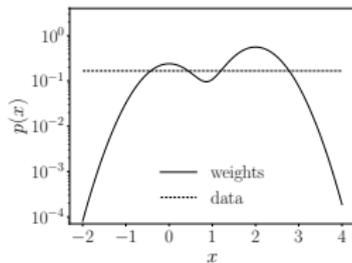


Training on weighted events

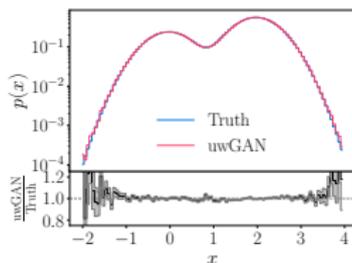
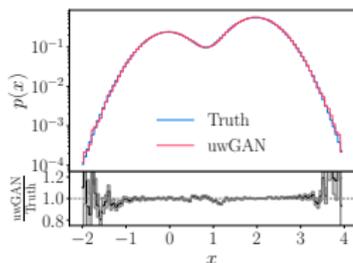
Information contained in distribution or event weights



Train on
weighted events



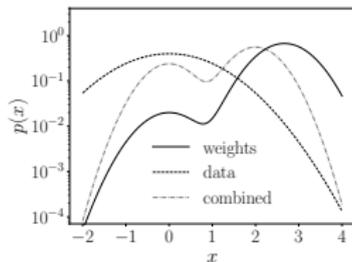
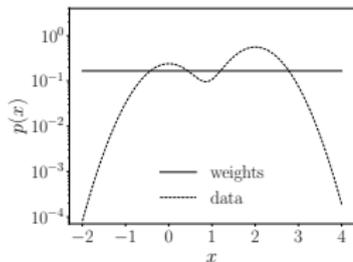
Generate
unweighted events



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

Training on weighted events

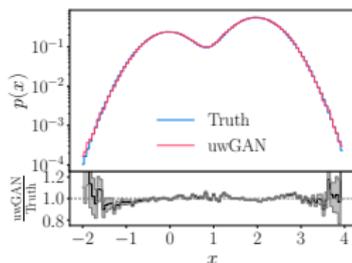
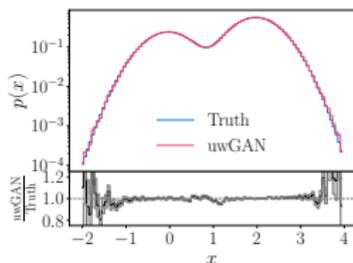
Information contained in distribution or event weights



Train on
weighted events



Generate
unweighted events

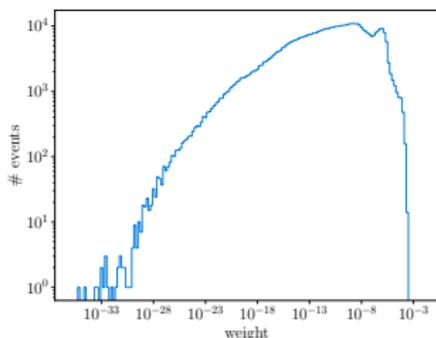


$$L_D = \langle -w \log D(x) \rangle_{x \sim P_{Truth}} + \langle -\log(1 - D(x)) \rangle_{x \sim P_{Gen}}$$

The unweighting bottleneck

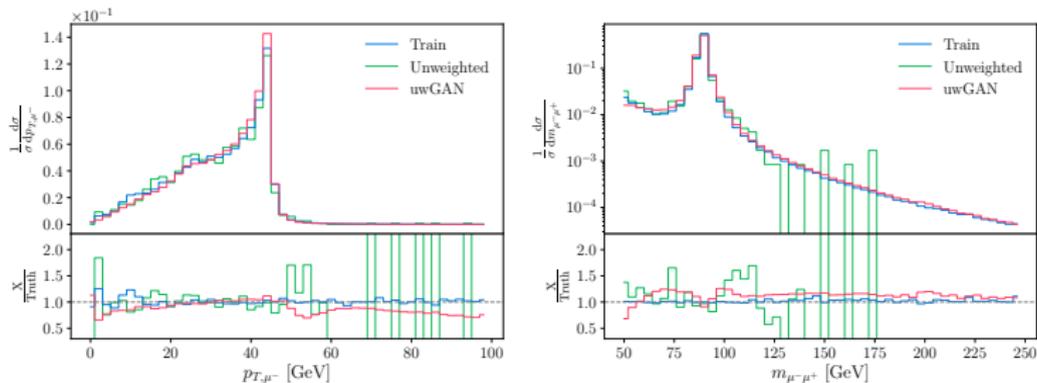
- High-multiplicity / higher-order \rightarrow large variation of weights
 \rightarrow unweighting efficiencies $< 1\%$
- \rightarrow Simulate conditions with naive Monte Carlo generator
ME by Sherpa, parton densities from LHAPDF, Rambo-on-diet

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



\rightarrow unweighting efficiency 0.2%

uwGAN results

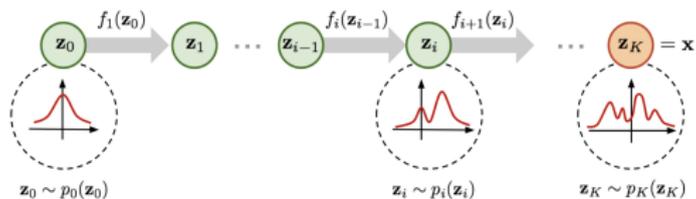


Populates high energy tails

Large amplification wrt. unweighted data!

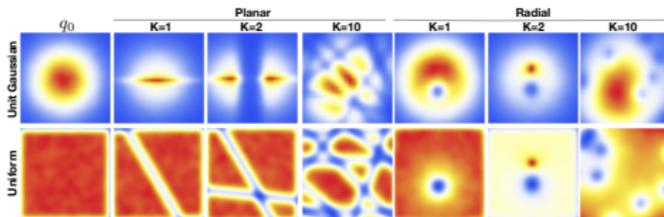
Normalizing flow

Transform input distribution into target distribution via invertible layers

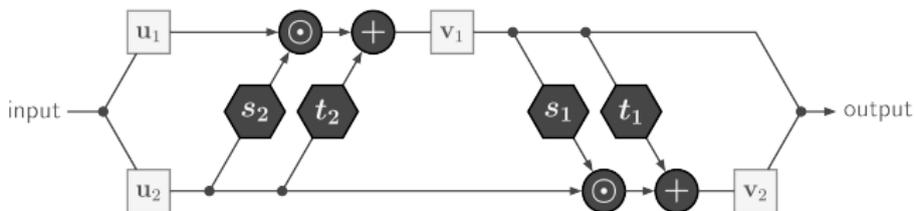


<https://lilianweng.github.io/lil-log/2018/10/13/flow-based-deep-generative-models.html>

- planar flow: $\mathbf{g}(\mathbf{x}) = \mathbf{x} + \mathbf{u}h(\mathbf{w}^T \mathbf{x} + b)$ [2015 Rezende, Mohamed]



Advanced coupling blocs



Features:

+ invertible

+ efficient evaluation in both directions

+ tractable Jacobian ($J(L_2(L_1)) = J(L_2) \cdot J(L_1)$)

+ trainable on samples OR probability

Standard event generation in a nutshell

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

→ keep events if $w_{event}/w_{max} > r \in [0, 1]$

→ optimal for $w \approx 1$

Standard event generation in a nutshell

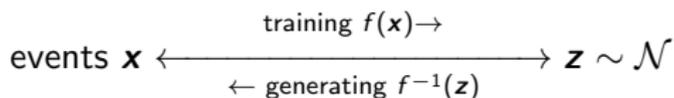
$$w_{event} = f(x_1, Q^2)f(x_2, Q^2) \times \mathcal{M}(x_1, x_2, p_1, \dots, p_n) \times J(p_i(r))^{-1}$$

Diagram illustrating the components of the event weight equation:

- $f(x_1, Q^2)f(x_2, Q^2)$ is labeled as PDF.
- $\mathcal{M}(x_1, x_2, p_1, \dots, p_n)$ is labeled as Matrix element.
- $J(p_i(r))^{-1}$ is labeled as Phase space mapping.

Find phase space mapping p_i such that $w \approx 1$

Training on samples

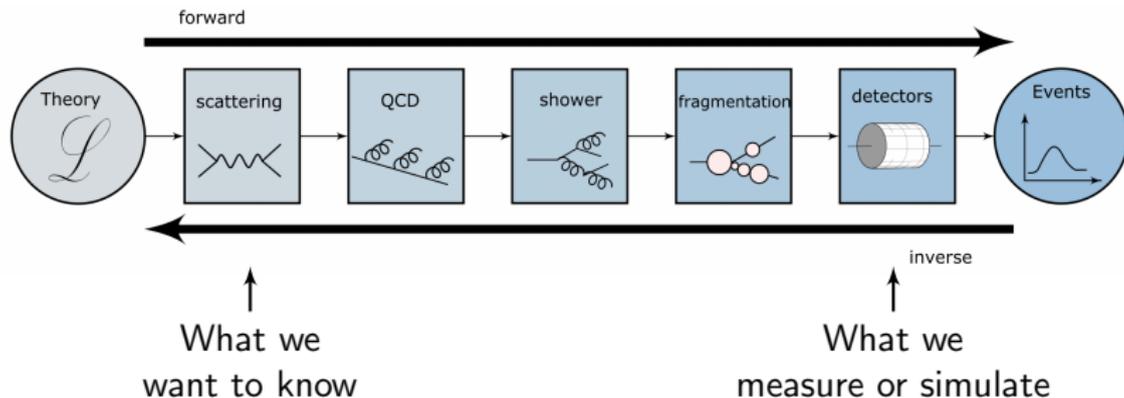


How to formulate the loss?

→ Maximize posterior over network parameters

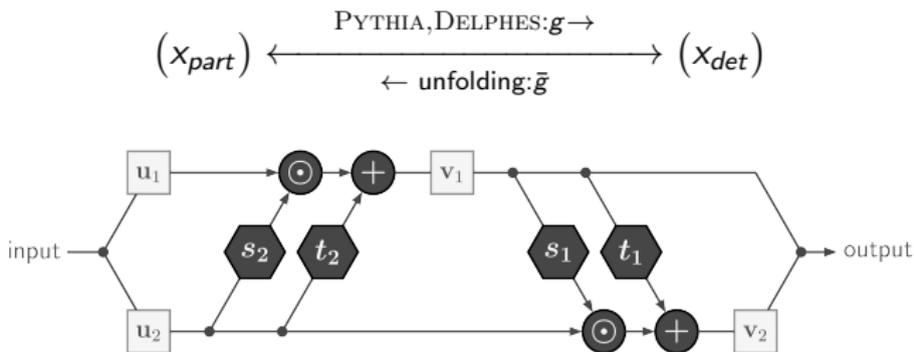
$$\begin{aligned} L &= -\langle \log p(\theta|x) \rangle_{x \sim P_x} \\ &= -\langle \log p(x|\theta) \rangle_{x \sim P_x} - \log p(\theta) + \text{const.} \quad (\text{Bayes' theorem}) \\ &= -\left\langle \log p(f(x)) + \log \left| \frac{\partial f(x)}{\partial x} \right| \right\rangle_{x \sim P_x} - \log p(\theta) + \text{const.} \\ &= -\langle -f(x)^2 + \log |J| \rangle_{x \sim P_x} - \log p(\theta) + \text{const.} . \end{aligned}$$

Can we invert the simulation chain?



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

Invertible networks



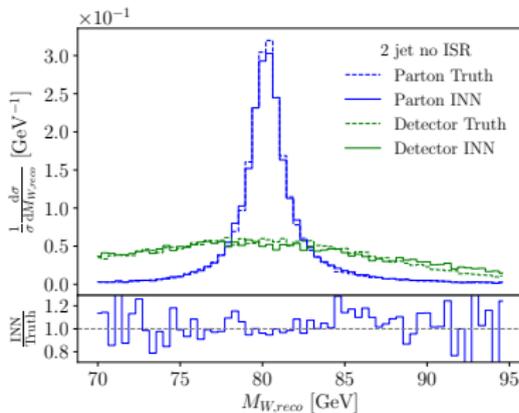
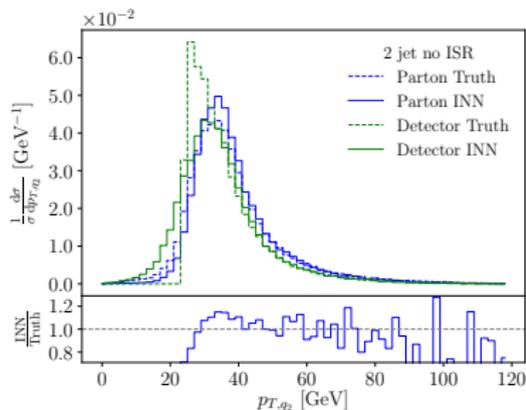
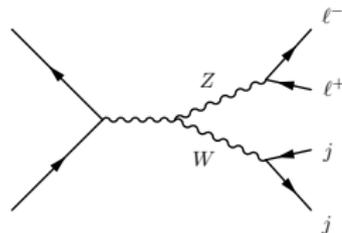
[1808.04730] L. Ardizzone, J. Kruse, S. Wirkert, D. Rahner,

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

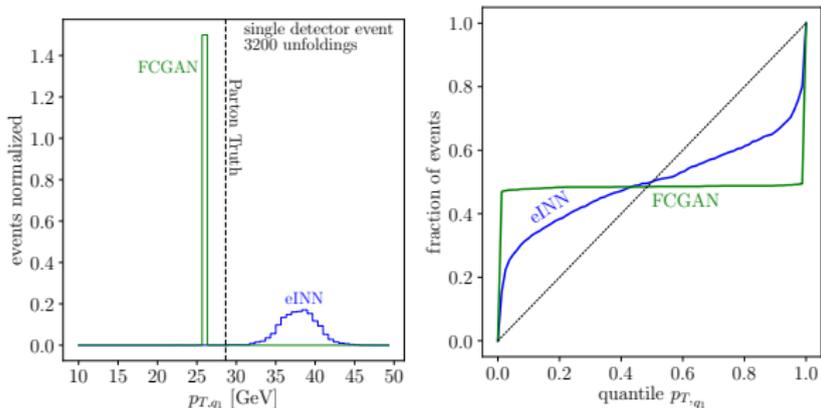


multi-dimensional ✓ bin independent ✓ statistically well defined ?

Including stochastic effects

$$\begin{array}{ccc} \left(\begin{array}{c} x_p \\ r_p \end{array} \right) & \begin{array}{c} \xrightarrow{\text{PYTHIA, DELPHES: } g \rightarrow} \\ \xleftarrow{\text{unfolding: } \bar{g}} \end{array} & \left(\begin{array}{c} x_d \\ r_d \end{array} \right) \end{array}$$

Sample r_d for fixed detector event
How often is Truth included in distribution quantile?



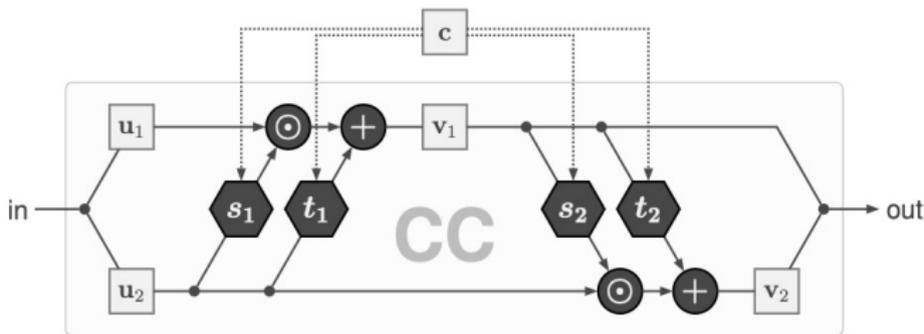
- Problem: arbitrary balance of many loss functions

Taking a different angle

Given an event x_d , what is the probability distribution at parton level?

→ sample over r , condition on x_d

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



Taking a different angle

Given an event x_d , what is the probability distribution at parton level?

→ sample over r , condition on x_d

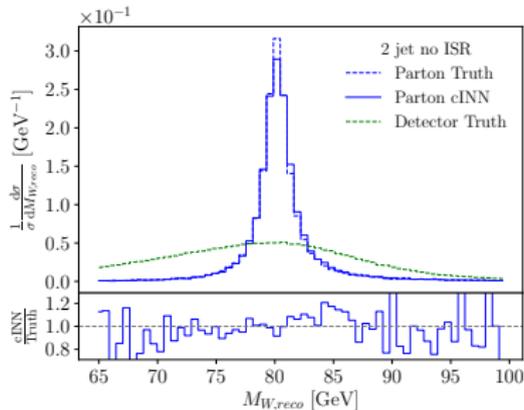
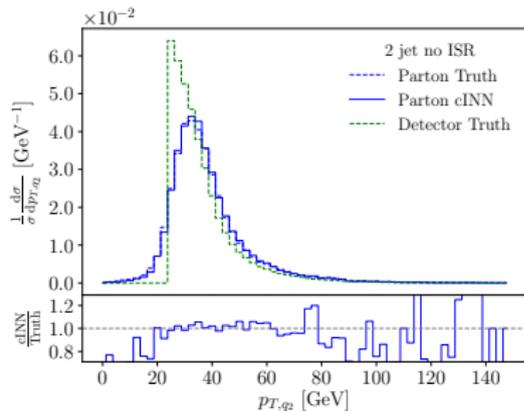
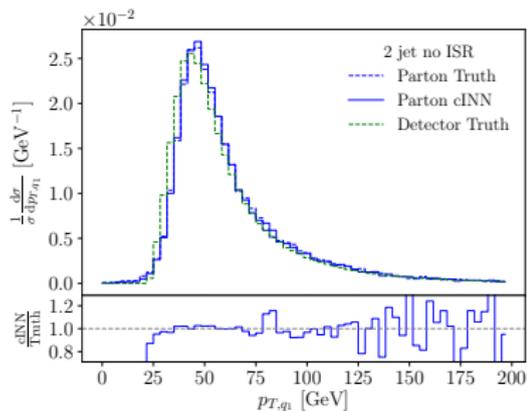
$$x_p \leftarrow \begin{array}{c} \xrightarrow{g(x_p, f(x_d))} \\ \xleftarrow{\text{unfolding: } \bar{g}(r, f(x_d))} \end{array} r$$

→ 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

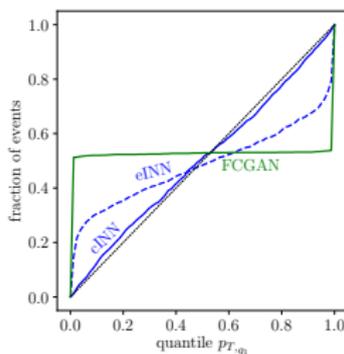
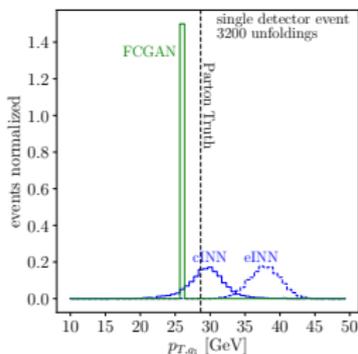
Cross check distributions



Condition INN on detector data [2006.06685]

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

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



multi-dimensional ✓ bin independent ✓ statistically well defined ✓

Summary

- Three types of generative models VAE, GAN, NF
- VAE: latent space encoding, KL loss can limit performance
- GAN: based on simple classifier, efficient training if stabilized
- NF: invertible, useable to train directly on probability

Now it's your turn 😊