

Precision Simulations Using Machine Learning

Tilman Plehn

Universität Heidelberg

Edinburgh, April 2023



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Data denoising [Cars]

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Data denoising [Cars]

- How to look for BSM physics?

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Data denoising [Cars]

- How to look for BSM physics?

Autoencoders [SAP]

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Data denoising [Cars]

- How to look for BSM physics?

Autoencoders [SAP]

- How to analyse LHC data?



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Data denoising [Cars]

- How to look for BSM physics?

Autoencoders [SAP]

- How to analyse LHC data?

Simulation-based inference [Edinburgh Ultra-Mini 2007]



LHC physics vs data scientist

LHC questions

- How to trigger from 3 PB/s to 300 MB/s?

Data compression [Netflix]

- How to analyze events with 4-vectors?

Graph neural networks [Cars]

- How to incorporate symmetries?

Contrastive learning [Google]

- How to combine tracker and calorimeter?

Super-resolution [Gaming]

- How to remove pile-up?

Data denoising [Cars]

- How to look for BSM physics?

Autoencoders [SAP]

- How to analyse LHC data?

Simulation-based inference [Edinburgh Ultra-Mini 2007]

- [How to treat uncertainties??](#)



Shortest ML-intro ever

Fit-like approximation [ask NNPDF]

- approximate known $f(x)$ using $f_\theta(x)$
- no parametrization, just very many values θ
- new representation/latent space θ

Construction and control

- define loss function
- minimize loss to find best θ
- compare $x \rightarrow f_\theta(x)$ for training/test data

LHC applications

- regression $x \rightarrow f_\theta(x)$
- classification $x \rightarrow f_\theta(x) \in [0, 1]$
- generation $r \sim \mathcal{N} \rightarrow f_\theta(r)$
- conditional generation $r \sim \mathcal{N} \rightarrow f_\theta(r|x)$
- ...

→ Transforming numerical science



Networks with error bar

Training-related uncertainties

- different trainings
- different initializations
- different data sets
- histogram network output: $f_{\theta}(x) \pm \Delta f(x)$

→ **Bayesian network: $\Delta f_{\theta}(x)$ from $\Delta \theta$** [Yarin Gal (2016)]

Energy measurement with NN

- expectation value from probability distribution

$$\langle E \rangle = \int dE \ E \ p(E) \rightarrow \int dE \ E \ p_{\theta}(E)$$

- energy $p(E|\theta)$ encoded in network parameters
- parameters $p(\theta|T)$ trained on data T

$$p_{\theta}(E) = \int d\theta \ p(E|\theta) \ p(\theta|T)$$

→ **Prediction by sampling weights**

$$\langle E \rangle = \int dE \ d\theta \ E \ p(E|\theta) \ p(\theta|T) = \int dE \ d\theta \ E \ p(E|\theta) \ q(\theta)$$



Constructing the loss function

Training means encoding $p(\theta|T)$

- so-called variational approximation [think $q(\theta)$ as Gaussian with mean and width]

$$p(E) = \int d\theta \, p(E|\theta) \, p(\theta|T) \stackrel{!}{=} \int d\theta \, p(E|\theta) \, q(\theta)$$

- similarity through minimized KL-divergence

$$D_{\text{KL}}[q(\theta), p(\theta|T)] = \int d\theta \, q(\theta) \log \frac{q(\theta)}{p(\theta|T)}$$



Constructing the loss function

Training means encoding $p(\theta|T)$

- so-called variational approximation [think $q(\theta)$ as Gaussian with mean and width]

$$p(E) = \int d\theta \, p(E|\theta) \, p(\theta|T) \stackrel{!}{=} \int d\theta \, p(E|\theta) \, q(\theta)$$

- similarity through minimized KL-divergence

$$D_{\text{KL}}[q(\theta), p(\theta|T)] = \int d\theta \, q(\theta) \log \frac{q(\theta)}{p(\theta|T)}$$

- Bayes' theorem to replace $p(\theta|T)$

$$\begin{aligned} D_{\text{KL}}[q(\theta), p(\theta|T)] &= \int d\theta \, q(\theta) \log \frac{q(\theta)p(T)}{p(T|\theta)p(\theta)} \\ &= D_{\text{KL}}[q(\theta), p(\theta)] - \int d\theta \, q(\theta) \log p(T|\theta) + \log p(T) \int d\theta \, q(\theta) \end{aligned}$$

- normalize distributions, ignore irrelevant terms, so minimize

$$D_{\text{KL}}[q(\theta), p(\theta|T)] \approx D_{\text{KL}}[q(\theta), p(\theta)] - \int d\theta \, q(\theta) \log p(T|\theta)$$



Constructing the loss function

Training means encoding $p(\theta|T)$

- so-called variational approximation [think $q(\theta)$ as Gaussian with mean and width]

$$p(E) = \int d\theta \, p(E|\theta) \, p(\theta|T) \stackrel{!}{=} \int d\theta \, p(E|\theta) \, q(\theta)$$

- similarity through minimized KL-divergence

$$D_{\text{KL}}[q(\theta), p(\theta|T)] = \int d\theta \, q(\theta) \log \frac{q(\theta)}{p(\theta|T)}$$

- Bayes' theorem to replace $p(\theta|T)$

$$\begin{aligned} D_{\text{KL}}[q(\theta), p(\theta|T)] &= \int d\theta \, q(\theta) \log \frac{q(\theta)p(T)}{p(T|\theta)p(\theta)} \\ &= D_{\text{KL}}[q(\theta), p(\theta)] - \int d\theta \, q(\theta) \log p(T|\theta) + \log p(T) \int d\theta \, q(\theta) \end{aligned}$$

- normalize distributions, ignore irrelevant terms, so minimize

$$D_{\text{KL}}[q(\theta), p(\theta|T)] \approx D_{\text{KL}}[q(\theta), p(\theta)] - \int d\theta \, q(\theta) \log p(T|\theta)$$

→ Loss combining likelihood and regularization

$$L = - \int d\theta \, q(\theta) \log p(T|\theta) + D_{\text{KL}}[q(\theta), p(\theta)]$$

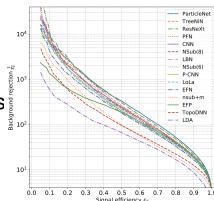


ML-applications for analysis

Top tagging [supervised classification]

- 'hello world' of LHC-ML
- end of QCD-taggers
- different NN-architectures

→ Non-NN left in the dust...



SciPost Physics

Submission

The Machine Learning Landscape of Top Taggers

G. Kasieczko^{1(a)}, T. Plehn^(a), A. Butter², K. Craner³, D. DeLauter⁴, B. M. Ertel⁵, M. Fairhead⁶, D. A. Farrelly⁷, W. Fickel⁸, C. Gay⁹, L. Goulet¹⁰, J. F. Kerner¹¹, P. T. Komodo¹², S. Lott¹³, A. Lister¹⁴, S. Maciunas¹⁵, E. M. Metodiev¹⁶, L. Moore¹¹, B. Nussens^{1,11}, K. Nussens^{1,11}, J. Puck¹⁶, H. Qiu⁹, Y. Ratz¹⁶, M. Rieger¹⁶, D. Shtyl¹⁶, J. M. Thompson¹⁶, and S. Varma¹⁶

¹ Institut für Experimentelle Physik, Universität Hamburg, Germany

² Institut für Theoretische Physik, Universität Hamburg, Germany

³ Center for Cosmology and Particle Physics and Center for Data Science, NYU, USA

⁴ NHETC, Dept. of Physics and Astronomy, Rutgers, The State University of NJ, USA

⁵ Joint Institute for Nuclear Research, Dubna, Russia

⁶ Theoretical Particle Physics and Cosmology, King's College London, United Kingdom

⁷ Department of Physics and Astronomy, The University of British Columbia, Canada

⁸ Department of Physics, University of California, Santa Barbara, USA

⁹ Faculty of Mathematics and Physics, University of Ljubljana, Ljubljana, Slovenia

¹⁰ Center for Theoretical Physics, MIT, Cambridge, USA

¹¹ CPJ, Universitat Catòlica de Leuven, Leuven-la-Neuve, Belgium

¹² Physics Division, Lawrence Berkeley National Laboratory, Berkeley, USA

¹³ Simons Inst. for the Theory of Computing, University of California, Berkeley, USA

¹⁴ National Institute for Subatomic Physics (NINHEP), Amsterdam, Netherlands

¹⁵ LPTHE, CNRS & Sorbonne Université, Paris, France

¹⁶ III. Physikalisches Institut A, RWTH Aachen University, Germany

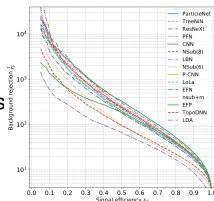


ML-applications for analysis

Top tagging [supervised classification]

- 'hello world' of LHC-ML
- end of QCD-taggers
- different NN-architectures

→ Non-NN left in the dust...



SciPost Physics

Submitted

The Machine Learning Landscape of Top Taggers

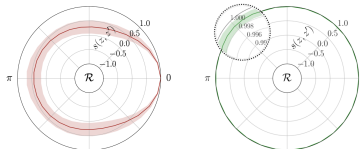
G. Kasieczko (ed.), T. Plehn (ed.), A. Butter*, K. Craner*, D. DeLoraine*, B. M. Dillon*, M. Fairhead*, D. A. Farrelly*, W. Fickel*, C. Gay*, L. Gossiaux*, J. F. Kerner*, P. T. Komarek*, S. Loefer*, A. Luter*, S. Maclean*, E. M. Metodiev*, L. Moore*, B. Raduana*, S. J. R. Nandorini*, J. Puckler*, H. Qiu*, Y. Ruan*, M. Rieger*, D. Shih*, J. M. Thompson*, and S. Varron*

- 1 Institut für Experimentelle Physik, Universität Heidelberg, Germany
- 2 Institut für Theoretische Physik, Universität Heidelberg, Germany
- 3 Center for Cosmology and Particle Physics and Center for Data Science, NYU, USA
- 4 NHETC, Dept. of Physics and Astronomy, Rutgers, The State University of NJ, USA
- 5 Joint Institute for Nuclear Research, JINR, Dubna, Russia
- 6 Theoretical Particle Physics and Cosmology, King's College London, United Kingdom
- 7 Department of Physics and Astronomy, The University of British Columbia, Canada
- 8 Department of Physics, University of California, Santa Barbara, USA
- 9 Faculty of Mathematics and Physics, University of Ljubljana, Ljubljana, Slovenia
- 10 Center for Theoretical Physics, MIT, Cambridge, USA
- 11 CPJ, Universitè Catholique de Louvain, Louvain-la-Neuve, Belgium
- 12 Physics Division, Lawrence Berkeley National Laboratory, Berkeley, USA
- 13 Simons Institute for the Theory of Computing, University of California, Berkeley, USA
- 14 National Institute for Subatomic Physics (NINHEP), Amsterdam, Netherlands
- 15 LPTHE, CNRS & Sorbonne Université, Paris, France
- 16 III. Physikalisches Institut A, RWTH Aachen University, Germany

Symmetric networks [contrastive learning, transformer network]

- rotations, translations, permutations, soft splittings, collinear splittings
- learn symmetries/augmentations

→ Symmetric latent representation



SciPost Physics

Submitted

Symmetries, Safety, and Self-Supervision

Berry M. Dillon*, Gregor Kasieczko*, Hans Othelidag*, Tilman Plehn*, Peter Sorrensen*, and Lorenz Vogt*

- 1 Institut für Theoretische Physik, Universität Heidelberg, Germany
- 2 Institut für Experimentelle Physik, Universität Heidelberg, Germany
- 3 Heidelberg Collaboratory for Image Processing, Universität Heidelberg, Germany

August 11, 2021

Abstract

Collider searches face the challenge of defining a representation of high-dimensional data such that physical symmetries are manifest, the distributional features are retained, and the choice of representation is non-physics agnostic. We introduce JetCLR to solve the mapping from low-level data to optimized observables through self-supervised contrastive learning. As an example, we construct a data representation for top and QCD jets using a permutation-invariant transformer-encoder network and visualize its symmetry properties. We compare the JetCLR representation with alternative representations using linear classifier tests and find it to work quite well.



Events and amplitudes

Speeding up Sherpa and MadNIS [sampling]

- precision simulations limiting factor for Runs 3&4
- unweighting critical

→ Phase space sampling

	$gg \rightarrow H_{\text{eff}}$	$u\bar{u} \rightarrow H_{\text{eff}}$	$s\bar{s} \rightarrow H_{\text{eff}}$	$b\bar{b} \rightarrow H_{\text{eff}}$
σ_{tot}	$1.1\text{e-}2$	$7.3\text{e-}3$	$6.6\text{e-}3$	$6.6\text{e-}4$
$\sigma_{H_{\text{eff}}}$	$8.7\text{e-}3$	$5.8\text{e-}3$	$4.7\text{e-}3$	$3.6\text{e-}4$
$(\sigma_{\text{tot}}/\sigma_{H_{\text{eff}}})$	38812	2417	199	64
$\mu_{H_{\text{eff}}}^{\text{stat}}$	52.03	32.52	69.75	326.19
$\mu_{H_{\text{eff}}}^{\text{th}}$	$2.4\text{e-}2$	$3.5\text{e-}2$	$2.1\text{e-}2$	$5.5\text{e-}3$
$\mu_{H_{\text{eff}}}^{\text{stat}}$	0.0669	0.3904	0.3904	0.0681
$\mu_{H_{\text{eff}}}^{\text{th}}$	2.21	4.89	1.47	0.19
$\mu_{H_{\text{eff}}}^{\text{stat}}$	20.40	19.14	27.75	35.34
$\mu_{H_{\text{eff}}}^{\text{th}}$	$4.3\text{e-}2$	$6.4\text{e-}2$	$5.1\text{e-}2$	$7.1\text{e-}2$
$\mu_{H_{\text{eff}}}^{\text{stat}}$	0.0683	0.0906	0.0903	0.0321
$\mu_{H_{\text{eff}}}^{\text{th}}$	3.96	8.26	5.91	2.22

Table 6: Performance measures for partonic channels contributing to $H \rightarrow 3$ jets production at the LHC.

SciPost Physics

Submissions

MCNET-21-13

Accelerating Monte Carlo event generation – rejection sampling using neural network event-weight estimates

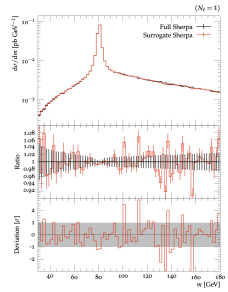
K. Danziger¹, T. Jocher², S. Schaefer², F. Siegel¹

¹ Institut für Kern- und Teilchenphysik, TU Dresden, Dresden, Germany
² Institut für Theoretische Physik, Georg-August-Universität Göttingen, Göttingen, Germany

September 27, 2021

Abstract

The generation of unit-weight events for complex scattering processes presents a severe challenge to modern Monte Carlo event generators. Even when using sophisticated phase-space sampling techniques adapted to the underlying transition matrix elements, the efficiency for generating unit-weight events from weighted samples can become a limiting factor in practical applications. Here we present a novel two-stage unweighting procedure that makes use of a neural-network surrogate for the full event weight. The algorithm can significantly accelerate the unweighting process, while it still guarantees unbiased sampling from the correct target distribution. We apply, validate and benchmark the new approach in high-multiplicity LHC production processes, including $2W+4$ jets and $2t+3$ jets, where we find speed-up factors up to ten.



Events and amplitudes

Speeding up Sherpa and MadNIS [sampling]

- precision simulations limiting factor for Runs 3&4
- unweighting critical

→ Phase space sampling

	$gg \rightarrow H_{\text{eff}}$	$gg \rightarrow t\bar{t}gg$	$gg \rightarrow t\bar{t}gg$	$gg \rightarrow H_{\text{eff}}$
r_{full}	$1.1e-2$	$7.3e-3$	$6.8e-3$	$6.6e-4$
$r_{\text{full,LO}}$	$8.7e-3$	$5.8e-3$	$4.7e-3$	$3.6e-4$
$(r_{\text{full}}/r_{\text{full,LO}})$	30033	3017	149	66
$r_{\text{full}}^{\text{MC}}$	52.03	32.52	69.75	206.19
$r_{\text{full}}^{\text{MC,LO}}$	2.4e-2	$3.8e-2$	$3.1e-2$	$5.6e-3$
$r_{\text{full}}^{\text{MC,LO}}$	0.0689	0.0884	0.0904	0.0961
$r_{\text{full}}^{\text{MC,LO}}$	2.21	1.89	1.47	0.19
$r_{\text{full}}^{\text{MC,LO}}$	20.01	18.14	27.78	35.34
$r_{\text{full}}^{\text{MC,LO}}$	4.3e-2	$6.4e-2$	$5.1e-2$	$7.1e-2$
$r_{\text{full}}^{\text{MC,LO}}$	0.0563	0.0900	0.0943	0.0921
$r_{\text{full}}^{\text{MC,LO}}$	3.50	8.20	3.91	2.22

Table 6: Performance measures for partonic channels contributing to $gg \rightarrow 3$ jets production at the LHC.

RePost Physics

Subsimulation

MCNET-21-13

Accelerating Monte Carlo event generation – rejection sampling using neural network event-weight estimates

K. Dauterle¹, T. Jausen¹, S. Schmeiser², F. Singer¹

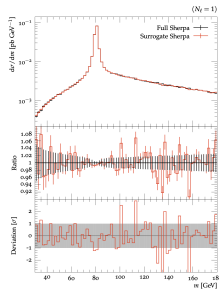
¹ Institut für Kern- und Teilchenphysik, TU Dresden, Dresden, Germany

² Institut für Theoretische Physik, Georg-August-Universität Göttingen, Göttingen, Germany

September 27, 2023

Abstract

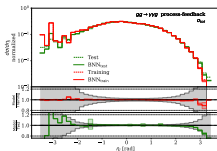
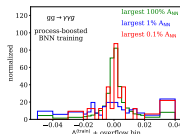
The generation of unit-weight events for complex scattering processes presents a severe challenge to modern Monte Carlo event generators. Even when using sophisticated phase-space sampling techniques adapted to the underlying transition matrix elements, the efficiency for generating unit-weight events from weighted samples can become a limiting factor in practical applications. Here we present a novel two-stage unweighting procedure that makes use of a neural-network surrogate for the full event weight. The algorithm can significantly accelerate the unweighting process, while it still guarantees unbiased sampling from the correct target distribution. We apply, validate and benchmark the new approach in high-multiplicity LHC production processes, including $2W \rightarrow 4$ jets and $t\bar{t} + 3$ jets, where we find speed-up factors up to ten.



Speeding up amplitudes [precision regression]

- loop-amplitudes expensive
- interpolation standard

→ Precision NN-amplitudes



PREPARED FOR SUBMISSION TO JHEP

JHEP09(2019)138

Optimising simulations for diphoton production at hadron colliders using amplitude neural networks

Joseph Aylott-Gullick^{a,b}, Simon Badger^a, Ryan Meebley^a

^aInstitute for Particle Physics Phenomenology, Department of Physics, Durham University, Durham, DH1 1TA, United Kingdom

^bInstitute for Data Science, Durham University, Durham, DH1 1TA, United Kingdom

^cDepartment of Physics and Astronomy, University of Toronto, and TRIUMF, Science at Toronto, Via P. O'Brien, 1, 6030 Keele, Toronto, Italy

E-mail: j.p.gullick@durham.ac.uk, simonbadger@utoronto.ca, ryan.meebley@durham.ac.uk

ABSTRACT: Machine learning technology has the potential to dramatically optimise event generation and simulation. We continue to investigate the use of neural networks to approximate matrix elements for high-multiplicity scattering processes. We focus on the case of loop-induced diphoton production through gluon fusion, and develop a modular simulation method that can be applied to hadron collider observables. Neural networks are trained using the one-loop amplitudes implemented in the Rivet++ library, and interfaced to the Sherpa Monte Carlo event generator, where we perform a detailed study for $2 \rightarrow 3$ and $2 \rightarrow 4$ scattering problems. We also consider how the trained networks perform when varying the kinematic cuts affecting the phase space and the reliability of the neural network simulations.



Invertible event generation

Precision NN-generators [Bayesian discriminator-flows]

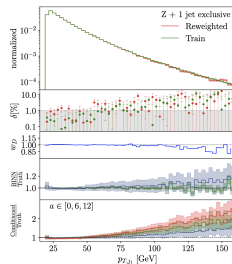
- control through discriminator [GAN-like]
- uncertainties through Bayesian networks

→ Discussed later



Abstract

Generative networks are opening new avenues in fast event generation for the LHC. We show how generative flow networks can reach percent-level precisions for Monte Carlo distributions, how they can be trained jointly with a discriminator, and how this discriminator improves the generation. Our joint training relies on a novel coupling of the two networks which does not require a Nash equilibrium. We then estimate the generation uncertainty through a Bayesian network setup and through conditional data augmentation, while the discriminator ensures that there are no systematic inconsistencies compared to the training data.



Invertible event generation

Precision NN-generators [Bayesian discriminator-flows]

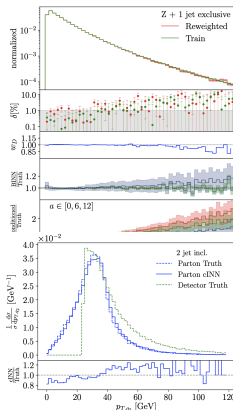
- control through discriminator [GAN-like]
- uncertainties through Bayesian networks

→ Discussed later



Abstract

Generative networks are opening new avenues in fast event generation for the LHC. We show how generative flow networks can reach percent-level precision for kinematic distributions, how they can be trained jointly with a discriminator, and how this discriminator improves the generation. Our joint training relies on a novel coupling of the two networks which does not require a Nash equilibrium. We then estimate the generation uncertainties through a Bayesian network using and through conditional data augmentation, while the discriminator ensures that there are no systematic inconsistencies compared to the training data.



Unfolding and inversion [conditional normalizing flows]

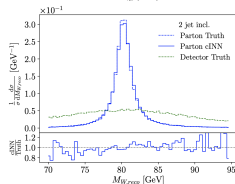
- shower/hadronization unfolded by jet algorithm
- detector/decays unfolded e.g. in tops
- calibrated inverse sampling

→ Discussed later



Abstract

For simulations where the forward and the inverse directions have a physics meaning, invertible neural networks are especially useful. A conditional INN can learn a detector simulation in terms of high-level observables, specifically for ZW production at the LHC. It allows for a per-event statistical interpretation. Next, we allow for a variable number of QCD jets. We model detector effects and QCD radiation to a pre-defined hard process, again with a per-event probabilistic interpretation over parton-level phase space.



Modern generative networks

Generative networks

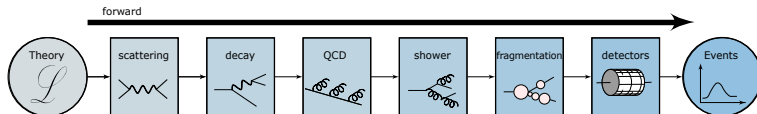
- generate **new** images, text blocks, etc
- encode density in target space
sample from Gaussian into target space
- reproduce training data, statistically independently



Modern generative networks

Generative networks

- generate **new** images, text blocks, etc
 - encode density in target space
sample from Gaussian into target space
 - reproduce training data, statistically independently
 - Variational Autoencoder
→ low-dimensional physics, high-dimensional objects
 - Generative Adversarial Network
→ generator trained by classifier
 - Normalizing Flow/Diffusion Model
→ stable bijective mapping
 - Generative Pre-trained Transformer
→ learning all structures
- **Pick best model for purpose**



Modern generative networks

Generative networks

- generate **new** images, text blocks, etc
 - encode density in target space
sample from Gaussian into target space
 - reproduce training data, statistically independently
 - Variational Autoencoder
→ low-dimensional physics, high-dimensional objects
 - Generative Adversarial Network
→ generator trained by classifier
 - Normalizing Flow/Diffusion Model
→ stable bijective mapping
 - Generative Pre-trained Transformer
→ learning all structures
- **Pick best model for purpose**

Fundamental question: GANplification

- first generated instances reproducing structures
- too many generated instances reproducing noise?

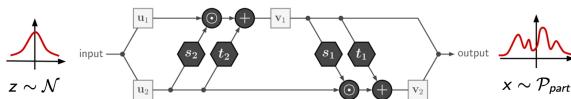


Modern generative networks

Normalizing flows — INN

- Gaussian latent space
- bijective mapping
- known Jacobian
- likelihood loss
- variety of coupling layers

→ Perfect for speed and precision



Modern generative networks

Normalizing flows — INN

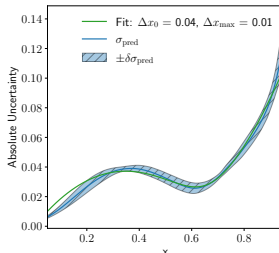
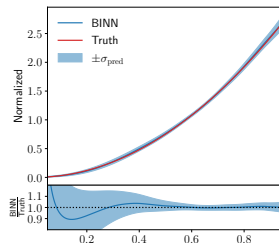
- Gaussian latent space
- bijective mapping
- known Jacobian
- likelihood loss
- variety of coupling layers

→ Perfect for speed and precision

INNs with uncertainties

- Bayesian NN for density estimation
- events with error bars
- density & uncertainty maps cross-talking

→ Bayesian INNs just fits with error bars



Precision generator

ML-event generators

- useful ML-playground
- training from event samples
no momentum conservation
no detector effects [sharper structures]

1- top-quark pairs $t\bar{t} \rightarrow 6 \text{ jets}$ [resonance peaks]

2- $Z_{\mu\mu} + \{1, 2, 3\} \text{ jets}$ [Z-peak, variable jet number, jet-jet topology]

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg



Precision generator

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg

ML-event generators

- useful ML-playground
- training from event samples
- no momentum conservation
- no detector effects [sharper structures]

1- top-quark pairs $t\bar{t} \rightarrow 6$ jets [resonance peaks]

2- $Z_{\mu\mu} + \{1, 2, 3\}$ jets [Z-peak, variable jet number, jet-jet topology]

INN-generator [Butter, Heime, Hummerich, Krebs, TP, Rousselot, Vent]

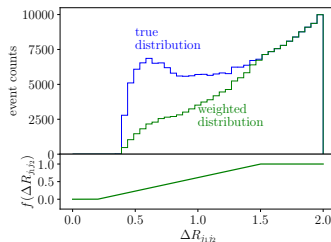
- challenging ΔR_{jj} features
- opposite of importance sampling

$$w^{(1\text{-jet})} = 1$$

$$w^{(2\text{-jet})} = f(\Delta R_{j_1, j_2})$$

$$w^{(3\text{-jet})} = f(\Delta R_{j_1, j_2})f(\Delta R_{j_2, j_3})f(\Delta R_{j_1, j_3})$$

$$f(\Delta R) = \frac{\Delta R - R_-}{R_+ - R_-} \quad (\Delta R \in [R_-, R_+])$$



Precision generator

ML-event generators

- useful ML-playground
- training from event samples
- no momentum conservation
- no detector effects [sharper structures]

- 1- top-quark pairs $t\bar{t} \rightarrow 6 \text{ jets}$ [resonance peaks]
- 2- $Z_{\mu\mu} + \{1, 2, 3\} \text{ jets}$ [Z-peak, variable jet number, jet-jet topology]

INN-generator [Butter, Heimes, Hummerich, Krebs, TP, Rousselot, Vent]

- challenging ΔR_{jj} features
- opposite of importance sampling

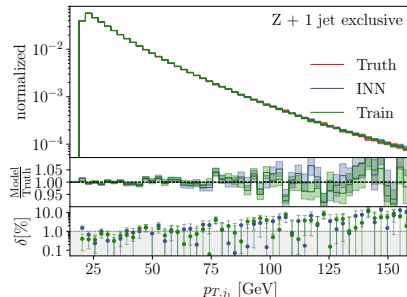
$$w^{(1\text{-jet})} = 1$$

$$w^{(2\text{-jet})} = f(\Delta R_{j_1, j_2})$$

$$w^{(3\text{-jet})} = f(\Delta R_{j_1, j_2}) f(\Delta R_{j_2, j_3}) f(\Delta R_{j_1, j_3})$$

$$f(\Delta R) = \frac{\Delta R - R_-}{R_+ - R_-} \quad (\Delta R \in [R_-, R_+])$$

→ Per-cent precision in reach

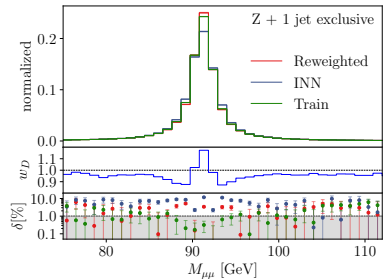
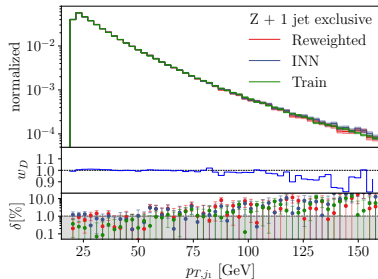


Controlled precision generator

Discriminator: training vs generated

- probability output $D = 0(\text{generator}), 1(\text{truth})$
- decent generator $D \approx 0.5$
- additional event weight $w_D = D/(1 - D)$

→ Dual use — control & reweight



Controlled precision generator

Discriminator: training vs generated

- probability output $D = 0(\text{generator}), 1(\text{truth})$
- decent generator $D \approx 0.5$
- additional event weight $w_D = D/(1 - D)$

→ Dual use — control & reweight

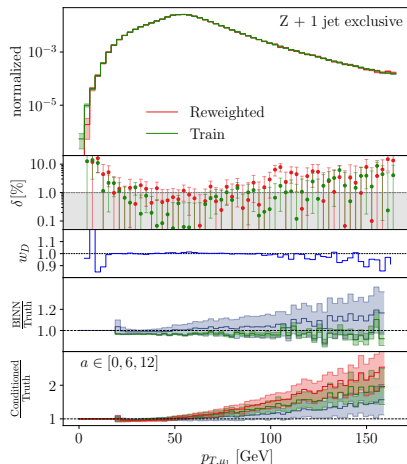
Uncertainties

- training uncertainties from BINN
- low statistics challenging
- systematics from data augmentation
- adjust data in tails $[a = 0 \dots 30]$

$$w = 1 + a \left(\frac{p_{T,j_1} - 15 \text{ GeV}}{100 \text{ GeV}} \right)^2$$

- train conditionally on smeared a
- error bar from sampling a

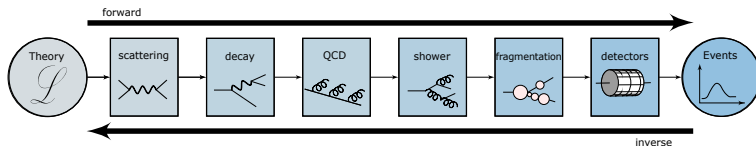
→ INNs for LHC standards



Inverse simulation

Invertible ML-simulation

- forward: $r \rightarrow$ events trained on model
- inverse: $r \rightarrow$ anything trained on model, conditioned on event



inverse



Inverse simulation

Invertible ML-simulation

- forward: $r \rightarrow$ events trained on model
- inverse: $r \rightarrow$ anything trained on model, conditioned on event
- individual steps known problems

detector unfolding

unfolding to QCD parton means jet algorithm

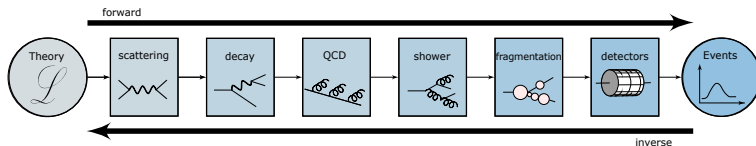
unfolding jet radiation known combinatorics problem

unfolding to hard process standard in top groups [needed for global analyses]

matrix element method an old dream

- improved through coherent ML-method

→ Free choice of data-theory inference point



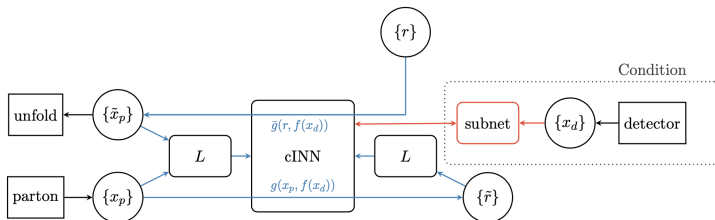
Inverting to hard process

Conditional INN

- partonic events x_p from $\{r\}$, given reco-event x_r
- loss based on likelihood

$$\begin{aligned}
 L &= - \langle \log p(\theta | x_p, x_r) \rangle_{x_p, x_r} \\
 &= - \langle \log p(x_p | x_r, \theta) + \log p(\theta | x_r) - \log p(x_p | x_r) \rangle_{x_p, x_r} \\
 &= - \langle \log p(x_p | x_r, \theta) \rangle_{x_p, x_r} - \log p(\theta) + \text{const.} \\
 &= - \left\langle \log p(g(x_p | x_r)) + \log \left| \frac{\partial g(x_p | x_r)}{\partial x_p} \right| \right\rangle_{x_p, x_r} - \log p(\theta) + \text{const.}
 \end{aligned}$$

→ Stable and statistically calibrated



Inverting to hard process

Conditional INN

- partonic events x_p from $\{r\}$, given reco-event x_r
- loss based on likelihood

$$L = - \langle \log p(\theta | x_p, x_r) \rangle_{x_p, x_r}$$

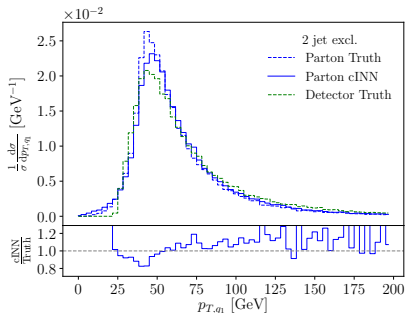
$$= - \left\langle \log p(g(x_p | x_r)) + \log \left| \frac{\partial g(x_p | x_r)}{\partial x_p} \right| \right\rangle_{x_p, x_r} - \log p(\theta) + \text{const.}$$

→ Stable and statistically calibrated

Undo detector and QCD jet radiation in $pp \rightarrow ZW + \text{jets}$

- hard process given
- detector and reconstruction universal
- jet radiation (approximately) universal
- model-independence: Butter-Malaescu

→ Stable and statistically calibrated



Inverting to hard process

Conditional INN

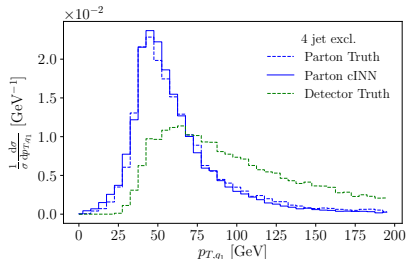
- partonic events x_p from $\{r\}$, given reco-event x_r
- loss based on likelihood

$$\begin{aligned}
 L &= - \langle \log p(\theta | x_p, x_r) \rangle_{x_p, x_r} \\
 &= - \left\langle \log p(g(x_p | x_r)) + \log \left| \frac{\partial g(x_p | x_r)}{\partial x_p} \right| \right\rangle_{x_p, x_r} - \log p(\theta) + \text{const.}
 \end{aligned}$$

→ Stable and statistically calibrated

Undo detector and QCD jet radiation in $pp \rightarrow ZW + \text{jets}$

- hard process given
 - detector and reconstruction universal
 - jet radiation (approximately) universal
 - model-independence: Butter-Malaescu
- Stable and statistically calibrated



Inverting to hard process

Conditional INN

- partonic events x_p from $\{r\}$, given reco-event x_r
- loss based on likelihood

$$L = - \langle \log p(\theta | x_p, x_r) \rangle_{x_p, x_r}$$

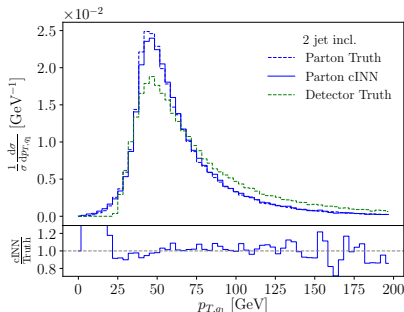
$$= - \left\langle \log p(g(x_p | x_r)) + \log \left| \frac{\partial g(x_p | x_r)}{\partial x_p} \right| \right\rangle_{x_p, x_r} - \log p(\theta) + \text{const.}$$

→ Stable and statistically calibrated

Undo detector and QCD jet radiation in $pp \rightarrow ZW + \text{jets}$

- hard process given
- detector and reconstruction universal
- jet radiation (approximately) universal
- model-independence: Butter-Malaescu

→ Stable and statistically calibrated



Optimal observables

Measure model parameter θ optimally

- single-event likelihood

$$p(x|\theta) = \frac{1}{\sigma_{\text{tot}}(\theta)} \frac{d^m \sigma(x|\theta)}{dx^m}$$

- expanded in θ around θ_0 , define score

$$\log \frac{p(x|\theta)}{p(x|\theta_0)} \approx (\theta - \theta_0) \left. \nabla_{\theta} \log p(x|\theta) \right|_{\theta_0} \equiv (\theta - \theta_0) \, t(x|\theta_0) \equiv (\theta - \theta_0) \, \phi^{\text{opt}}(x)$$

- leading order parton level

$$p(x|\theta) \approx |\mathcal{M}|_0^2 + \theta |\mathcal{M}|_{\text{int}}^2 \quad \Rightarrow \quad t(x|\theta_0) \sim \frac{|\mathcal{M}|_{\text{int}}^2}{|\mathcal{M}|_0^2}$$



Optimal observables

Measure model parameter θ optimally

- single-event likelihood

$$p(x|\theta) = \frac{1}{\sigma_{\text{tot}}(\theta)} \frac{d^m \sigma(x|\theta)}{dx^m}$$

- expanded in θ around θ_0 , define score

$$\log \frac{p(x|\theta)}{p(x|\theta_0)} \approx (\theta - \theta_0) \left. \nabla_{\theta} \log p(x|\theta) \right|_{\theta_0} \equiv (\theta - \theta_0) t(x|\theta_0) \equiv (\theta - \theta_0) \phi^{\text{opt}}(x)$$

- leading order parton level

$$p(x|\theta) \approx |\mathcal{M}|_0^2 + \theta |\mathcal{M}|_{\text{int}}^2 \quad \Rightarrow \quad t(x|\theta_0) \sim \frac{|\mathcal{M}|_{\text{int}}^2}{|\mathcal{M}|_0^2}$$

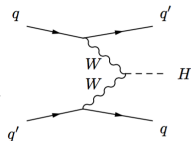
CP-violating Higgs production

- unique CP-observable

$$t \propto \epsilon_{\mu\nu\rho\sigma} k_1^{\mu} k_2^{\nu} q_1^{\rho} q_2^{\sigma} \text{sign}[(k_1 - k_2) \cdot (q_1 - q_2)] \xrightarrow{\text{lab frame}} \sin \Delta\phi_{jj}$$

- CP-effect in $\Delta\phi_{jj}$
D6-effect in $p_{T,j}$

⇒ Established LHC task



Analytic formula for score

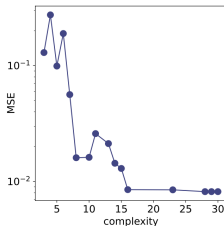
- function to approximate $t(x|\theta)$
- phase space parameters $x_p = p_T/m_H, \Delta\eta, \Delta\phi$ [node]
- operators $\sin x, x^2, x^3, x + y, x - y, x * y, x/y$ [node]
- represent formula as tree [complexity = number of nodes]

⇒ Figures of merit

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n [g_i(x) - t(x, z|\theta)]^2 \rightarrow \text{MSE} + \text{parsimony} \cdot \text{complexity}$$

Score around Standard Model

compl	dof	function	MSE
3	1	$a \Delta\phi$	$1.30 \cdot 10^{-1}$
4	1	$\sin(a\Delta\phi)$	$2.75 \cdot 10^{-1}$
5	1	$a\Delta\phi x_{p,1}$	$9.93 \cdot 10^{-2}$
6	1	$-x_{p,1} \sin(\Delta\phi + a)$	$1.90 \cdot 10^{-1}$
7	1	$(-x_{p,1} - a) \sin(\sin(\Delta\phi))$	$5.63 \cdot 10^{-2}$
8	1	$(a - x_{p,1}) x_{p,2} \sin(\Delta\phi)$	$1.61 \cdot 10^{-2}$
14	2	$x_{p,1}(a\Delta\phi - \sin(\sin(\Delta\phi)))(x_{p,2} + b)$	$1.44 \cdot 10^{-2}$
15	3	$-(x_{p,2}(a\Delta\eta^2 + x_{p,1}) + b) \sin(\Delta\phi + c)$	$1.30 \cdot 10^{-2}$
16	4	$-x_{p,1}(a - b\Delta\eta)(x_{p,2} + c) \sin(\Delta\phi + d)$	$8.50 \cdot 10^{-3}$
28	7	$(x_{p,2} + a)(bx_{p,1}(c - \Delta\phi) - x_{p,1}(d\Delta\eta + ex_{p,2} + f) \sin(\Delta\phi + g))$	$8.18 \cdot 10^{-3}$



PySR

Brief intro

Bayesian nets

ML examples

Generation

Inversion

Symbolic reg

Analytic formula for score

- function to approximate $t(x|\theta)$
- phase space parameters $x_p = p_T/m_H, \Delta\eta, \Delta\phi$ [node]
- operators $\sin x, x^2, x^3, x + y, x - y, x * y, x/y$ [node]
- represent formula as tree [complexity = number of nodes]

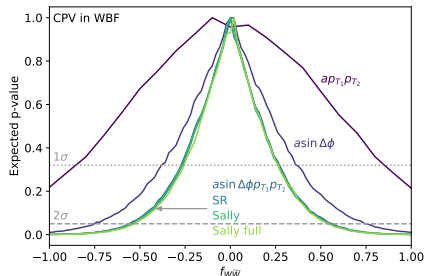
⇒ Figures of merit

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n [g_i(x) - t(x, z|\theta)]^2 \rightarrow \text{MSE} + \text{parsimony} \cdot \text{complexity}$$

Score around Standard Model

- expected limits:
very wrong formula
wrong formula
right formula
MadMiner
- same within statistical limitation

⇒ New optimal observables next



ML for LHC Theory

ML-applications

- just another numerical tool for a numerical field
- driven by money from data science and medical research
- goals are...
 - ...improve established tasks
 - ...develop new tools for established tasks
 - ...transform through new ideas
- xAI through...
 - ...precision control
 - ...uncertainties
 - ...symmetries
 - ...formulas

→ Fun with good old LHC problems

Modern Machine Learning for LHC Physicists

Tilman Plehn^a, Anja Butter^{a,b}, Barry Dillon^a, and Claudius Krause^{a,c}

^a Institut für Theoretische Physik, Universität Heidelberg, Germany

^b LPNHE, Sorbonne Université, Université Paris Cité, CNRS/IN2P3, Paris, France

^c NHETC, Dept. of Physics and Astronomy, Rutgers University, Piscataway, USA

November 2, 2022

Abstract

Modern machine learning is transforming particle physics, faster than we can follow, and bullying its way into our numerical tool box. For young researchers it is crucial to stay on top of this development, which means applying cutting-edge methods and tools to the full range of LHC physics problems. These lecture notes are meant to lead students with basic knowledge of particle physics and significant enthusiasm for machine learning to relevant applications as fast as possible. They start with an LHC-specific motivation and a non-standard introduction to neural networks and then cover classification, unsupervised classification, generative networks, and inverse problems. Two themes defining much of the discussion are well-defined loss functions reflecting the problem at hand and uncertainty-aware networks. As part of the applications, the notes include some aspects of theoretical LHC physics. All examples are chosen from particle physics publications of the last few years. Given that these notes will be outdated already at the time of submission, the week of ML4jets 2022, they will be updated frequently.



Inverting to QCD

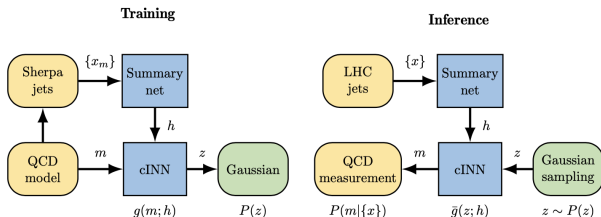
cINN for inference [Bieringer, Butter, Heimes, Höche, Köthe, TP, Radev]

- condition jets with QCD parameters
- train model parameters \rightarrow Gaussian latent space
- test Gaussian sampling \rightarrow parameter measurement
- beyond C_A vs C_F [Kluth et al]

$$P_{qq} = C_F \left[D_{qq} \frac{2z(1-y)}{1-z(1-y)} + F_{qq}(1-z) + C_{qq}yz(1-z) \right]$$

$$P_{gg} = 2C_A \left[D_{gg} \left(\frac{z(1-y)}{1-z(1-y)} + \frac{(1-z)(1-y)}{1-(1-z)(1-y)} \right) + F_{gg}z(1-z) + C_{gg}yz(1-z) \right]$$

$$P_{gq} = T_R \left[F_{gq} (z^2 + (1-z)^2) + C_{gq}yz(1-z) \right]$$



Inverting to QCD

cINN for inference [Bieringer, Butter, Heimes, Höche, Köthe, TP, Radev]

- condition jets with QCD parameters
- train model parameters \rightarrow Gaussian latent space
- test Gaussian sampling \rightarrow parameter measurement

- beyond C_A vs C_F [Kluth et al]

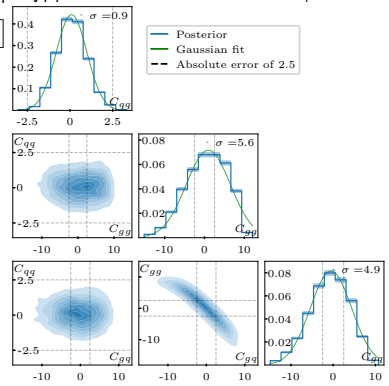
$$P_{qq} = C_F \left[D_{qq} \frac{2z(1-y)}{1-z(1-y)} + F_{qq}(1-z) + C_{qq}yz(1-z) \right]$$

$$P_{gg} = 2C_A \left[D_{gg} \left(\frac{z(1-y)}{1-z(1-y)} + \frac{(1-z)(1-y)}{1-(1-z)(1-y)} \right) + F_{gg}z(1-z) + C_{gg}yz(1-z) \right]$$

$$P_{gq} = T_R \left[F_{gq} (z^2 + (1-z)^2) + C_{gq}yz(1-z) \right]$$

- idealized shower [Sherpa]

- More ML-opportunities...

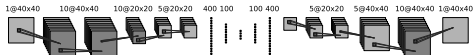


Learning background only

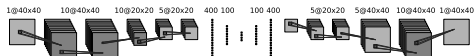
Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
- reconstruct QCD jets \rightarrow top jets hard to describe
- reconstruct top jets \rightarrow QCD jets just simple top-like jet

\rightarrow Symmetric performance $S \leftrightarrow B?$



Learning background only

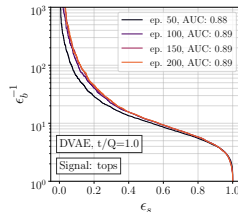
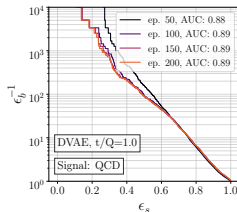
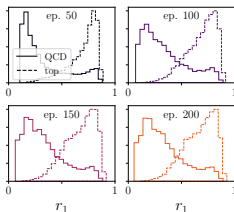


Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets \rightarrow top jets hard to describe
 - reconstruct top jets \rightarrow QCD jets just simple top-like jet
- \rightarrow Symmetric performance $S \leftrightarrow B?$

Moving to latent space

- anomaly score from latent space?
- VAE \rightarrow does not work
- GMVAE \rightarrow does not work
- Dirichlet VAE \rightarrow works okay
- density estimation \rightarrow does not work



Learning background only



Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets \rightarrow top jets hard to describe
 - reconstruct top jets \rightarrow QCD jets just simple top-like jet
- \rightarrow Symmetric performance $S \leftrightarrow B?$

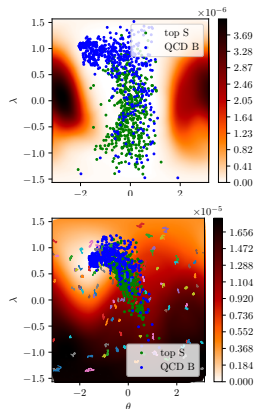
Normalized autoencoder [penalize missing features]

- normalized probability loss
- Boltzmann mapping [$E_\theta = \text{MSE}$]

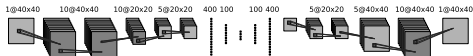
$$p_\theta(x) = \frac{e^{-E_\theta(x)}}{Z_\theta}$$

$$L = -\langle \log p_\theta(x) \rangle = \langle E_\theta(x) + \log Z_\theta \rangle$$

- inducing background metric
 - small MSE for data, large MSE for model
 - Z_θ from (Langevin) Markov Chain
- \rightarrow Symmetric autoencoder, at last



Learning background only



Unsupervised classification

- train on background only
extract unknown signal from reconstruction error
 - reconstruct QCD jets → top jets hard to describe
 - reconstruct top jets → QCD jets just simple top-like jet
- Symmetric performance $S \leftrightarrow B?$

Normalized autoencoder [penalize missing features]

- normalized probability loss
- Boltzmann mapping [$E_\theta = \text{MSE}$]

$$p_\theta(x) = \frac{e^{-E_\theta(x)}}{Z_\theta}$$

$$L = -\langle \log p_\theta(x) \rangle = \langle E_\theta(x) + \log Z_\theta \rangle$$

- inducing background metric
 - small MSE for data, large MSE for model
 - Z_θ from (Langevin) Markov Chain
- Symmetric autoencoder, at last

