

Generative Models for Natural Science

3rd IML Machine Learning Workshop

CERN 15-18 April 2019

Fedor Ratnikov



SCHOOL OF DATA ANALYSIS

NRU Higher School of Economics,
Yandex School of Data Analysis

Approaches

Generative models [\[edit \]](#)

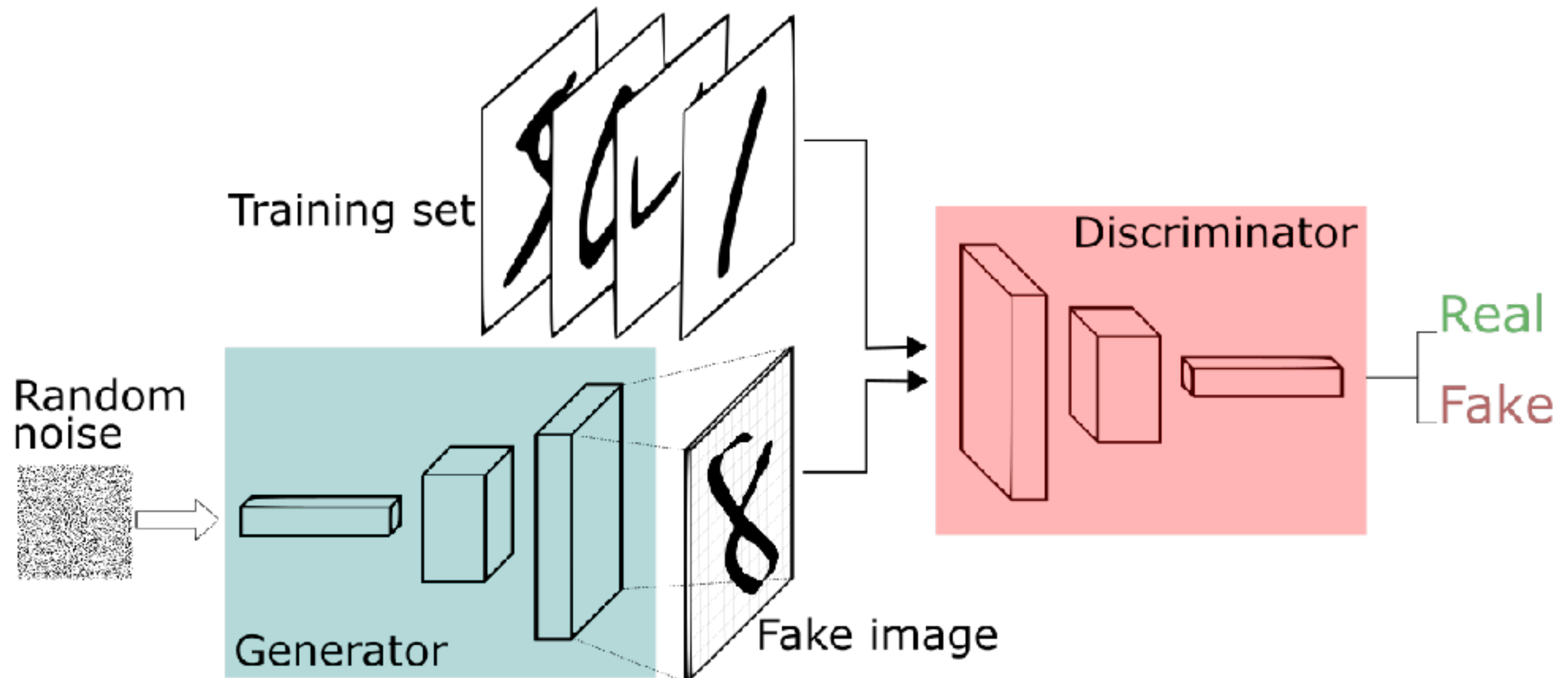
Types of generative models are:



- Gaussian mixture model (and other types of mixture model)
 - Hidden Markov model
 - Probabilistic context-free grammar
 - Bayesian network (e.g. Naive bayes, Autoregressive model)
 - Averaged one-dependence estimators
 - Latent Dirichlet allocation
 - Boltzmann machine (e.g. Restricted Boltzmann machine, Deep belief network)
 - Variational autoencoder
 - Generative adversarial network
 - Flow-based generative model
-
- GAN and VAE are mostly used nowadays for generating complicated objects

Generative Adversarial Network (GAN)

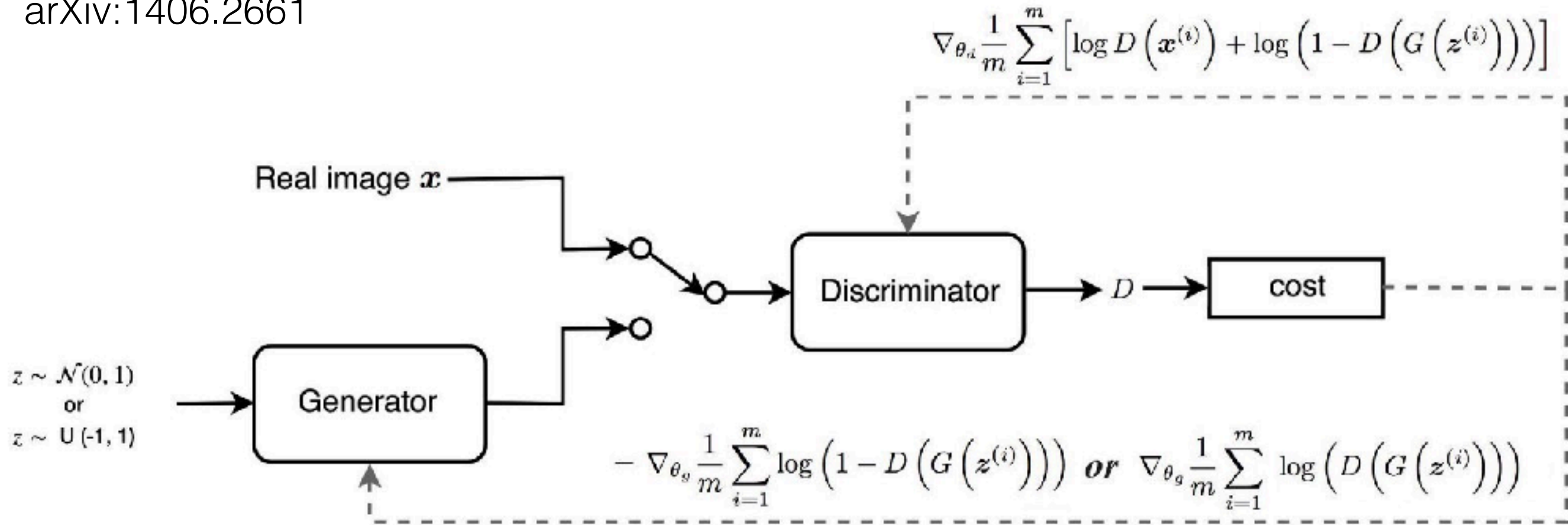
<https://medium.freecodecamp.org/an-intuitive-introduction-to-generative-adversarial-networks-gans-7a2264a81394>



- Implicit $p(x|y)$, sampling only

Classic GAN

arXiv:1406.2661



- Discriminator approaches Jensen–Shannon divergence
 - vanishing gradients for poor generator
 - mode collapse

Illustration: Jonathan Hui

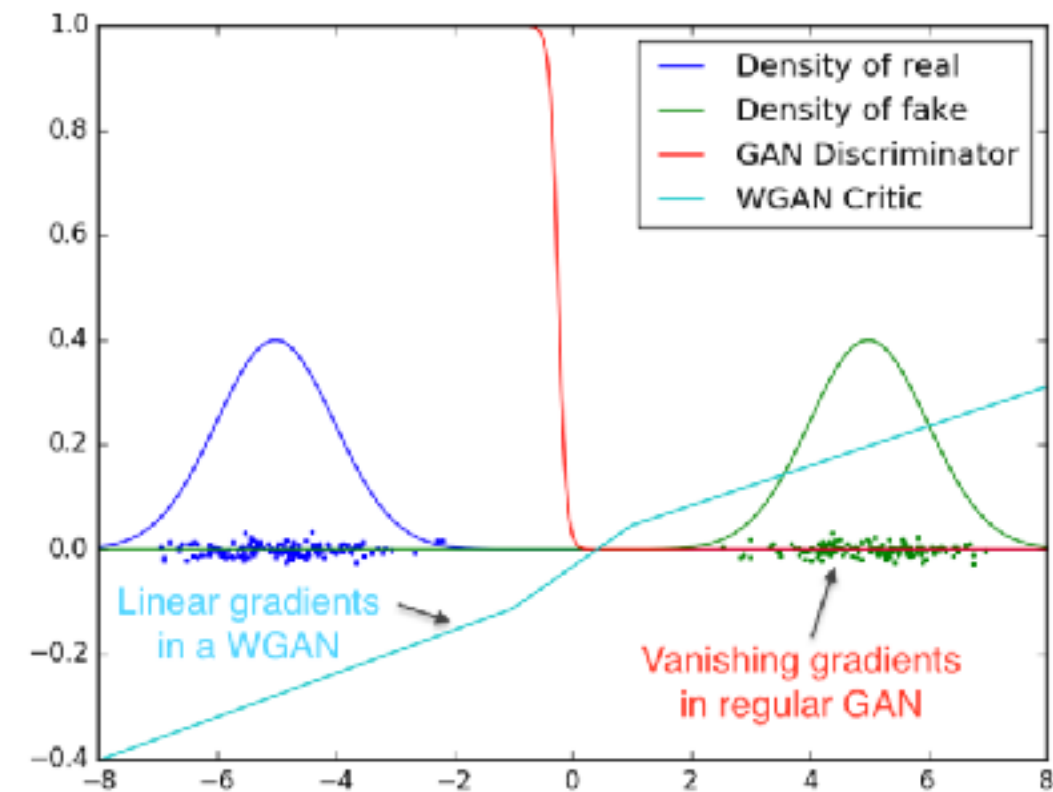
Wasserstein GAN

- Uses “earth mover’s distance”:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

- Kantorovich-Rubinstein duality:

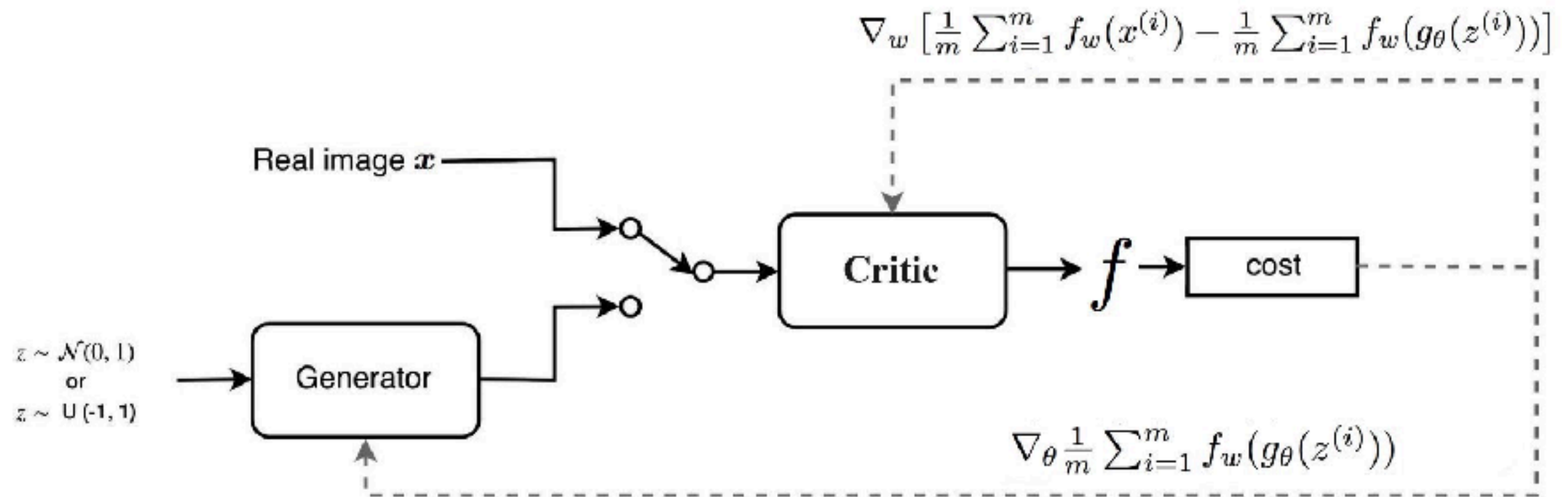
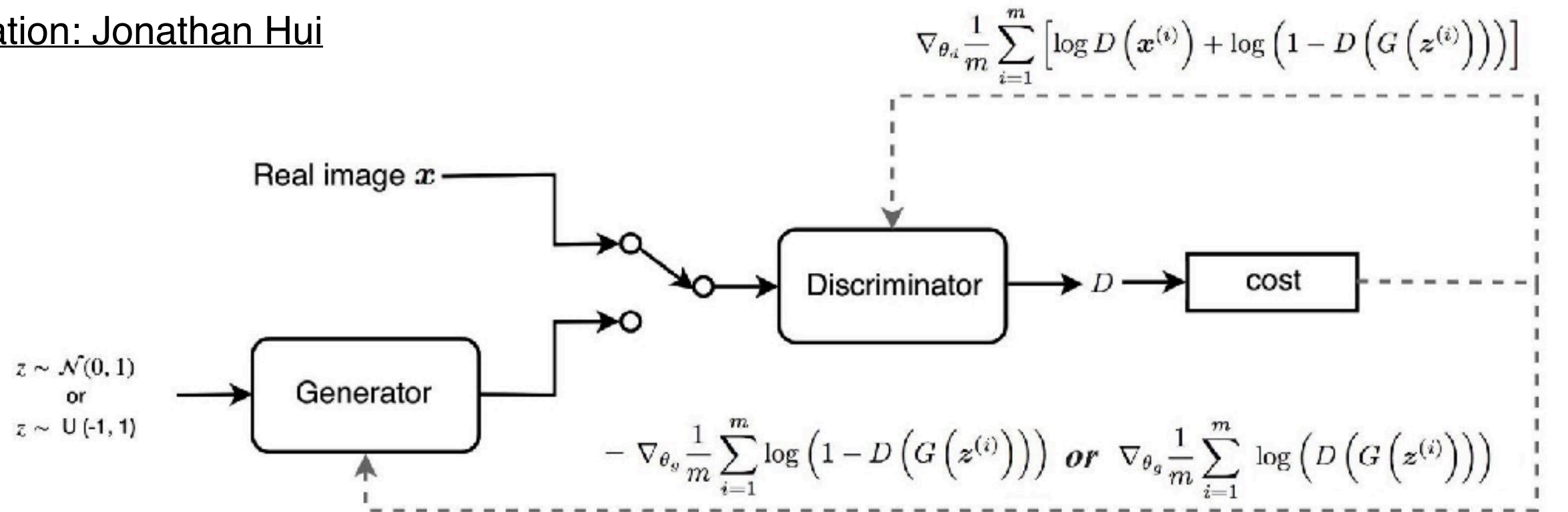
$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\substack{\|f\|_L \leq 1 \\ |f(x_1) - f(x_2)| \leq |x_1 - x_2|}} \mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)]$$



arXiv:1701.07875

- “Discriminator” function $f(x)$ may be approached using deep network
 - output is not probability, but any scalar number
 - need to satisfy 1-Lipschitz condition

Illustration: Jonathan Hui



CramerGAN

- WGAN produces biased gradients, that makes converging slower, sometimes never reaching optimum
- CramerGAN uses energy distance as a critic (discriminator):

$$\mathcal{E}(X, Y) := 2 \mathbb{E} \|X - Y\|_2 - \mathbb{E} \|X - X'\|_2 - \mathbb{E} \|Y - Y'\|_2 \quad \text{arXiv:1705.10743}$$

- where X, X', Y, Y' are statistically independent samples from two distributions
- corresponds to the Cramer distance in 1D case:

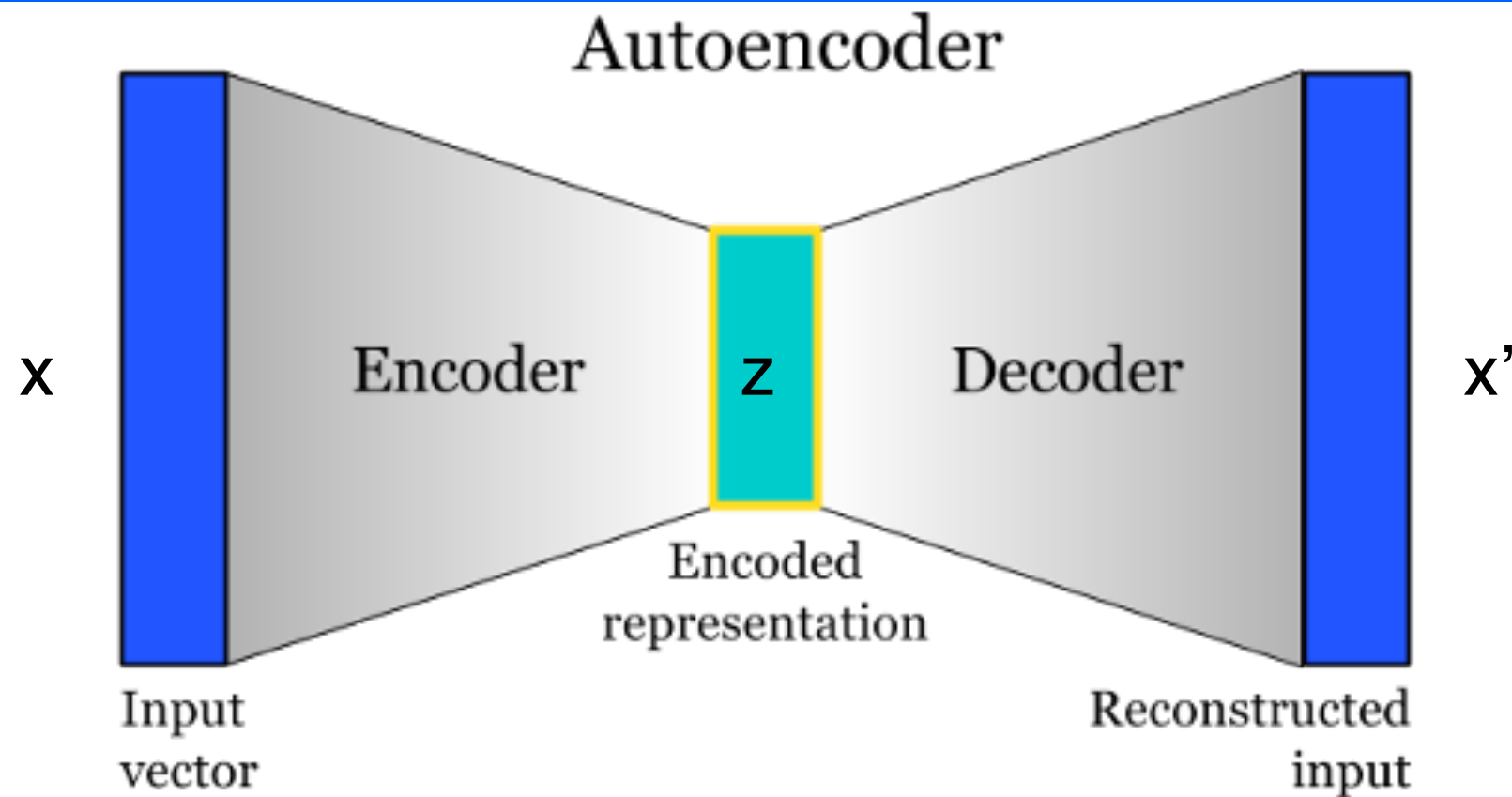
$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$

- generator loss is therefore more complicated:

$$L_g = 2\|h(x_r) - h(x_g)\|_2 - \|h(x_r) - h(x'_r)\|_2 - \|h(x_g) - h(x'_g)\|_2$$

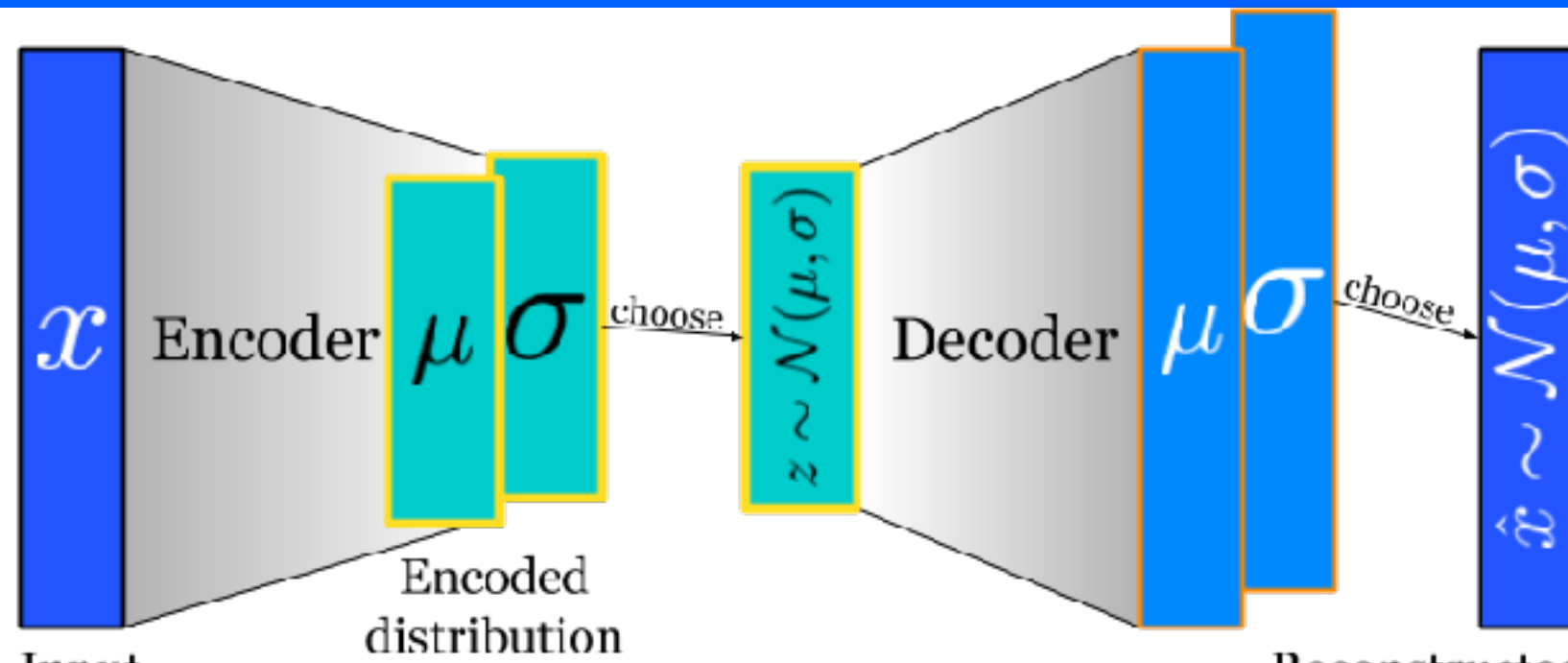
- critic is trained to maximize the energy distance
- CramerGAN demonstrates better convergency indeed

Variational Autoencoder



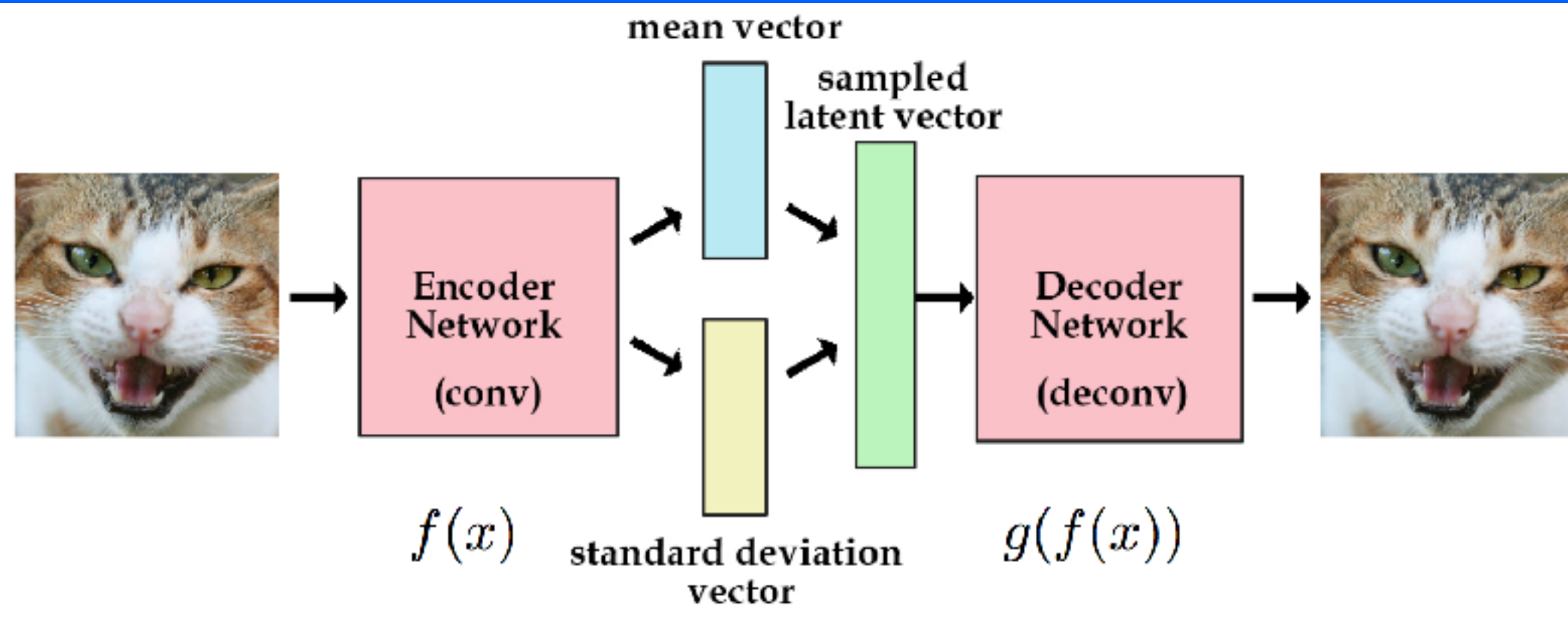
- Autoencoder can be trained to sample realistic objects
- $x \rightarrow \text{encoder} \rightarrow z \rightarrow \text{decoder} \rightarrow x'$
 - require $x' \sim x$
- Decoder part of the AE can generate realistic objects
 - ... providing correct prior distribution in the latent space $p(z)$

Variational Autoencoder



- Decoder part of the AE can generate realistic objects
 - ... providing correct prior distribution in the latent space $p(z)$
- Put extra requirement into the loss
 - latent distribution $p(z/X)$ must approach some standard one, e.g. $\mathcal{N}(\mathbf{0}, \mathbf{I})$
 - make $z(x)$ variational
 - (make $x'(z)$ variational)

Variational Autoencoder



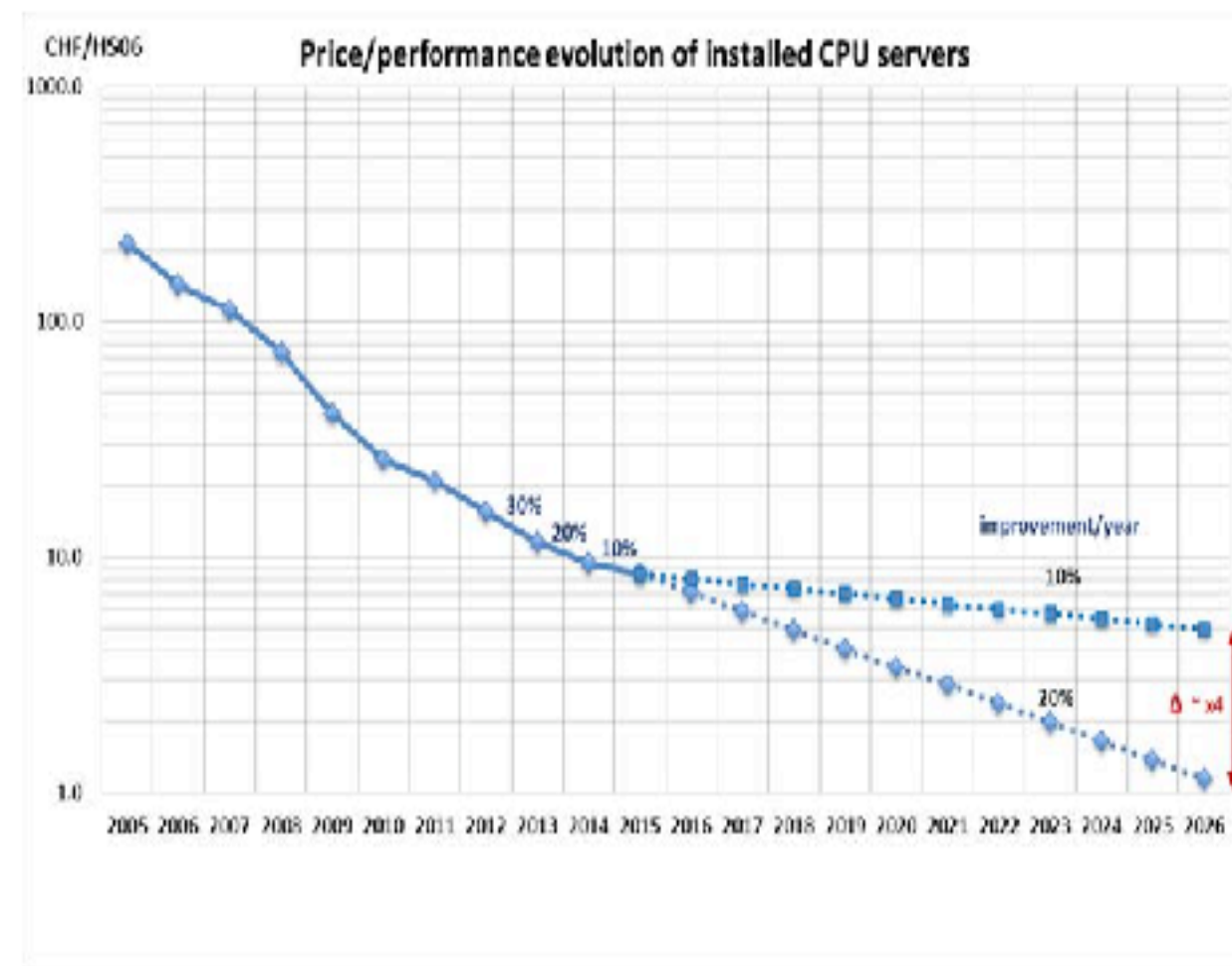
- VAE allows calculating $p(x|y)$
 - NB: GAN only allows sampling from $p(x|y)$
- ... but smaller number of dimensions in the latent space
 - blurry objects

Library Approach

- We have train sample for the generative model anyway
 - consistency with this train sample is a figure of merit for the generative model
- Objects of the train sample may be used for generation directly
 - similar to KNN classification algorithm
 - $k=1$: search for the object with appropriate conditions in the (presumably huge) data library
 - $k>1$: need to interpolate between objects
 - short distance objects interpolation, more robust than global generation
- NB: library approach **by construction** uses full information which is contained in the training sample

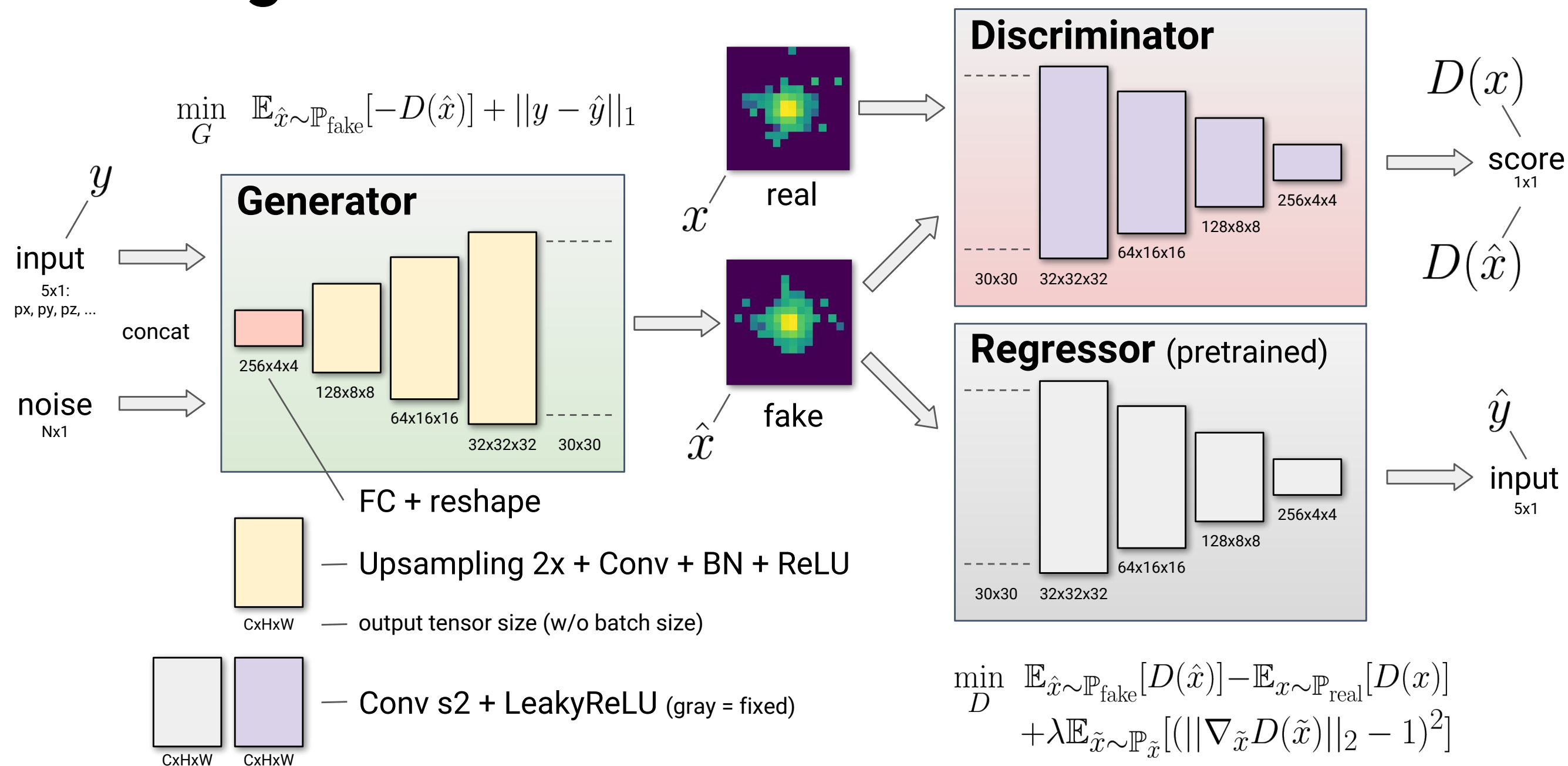
Generative Models at LHC

- About 80% of computing resources are used for MC simulation in HEP experiments
 - Calorimeter simulation is one of bottlenecks
 - RICH is the next in the row for LHCb detector
 - > 85% of simulation is taken by these
- Can not expect exponential rise of CPU performance
- Need a work around for Run3 and HL-LHC
- Generative models trained on the detailed GEANT simulation may be a solution



Example: ECAL Conditional Fast Simulation

Training scheme

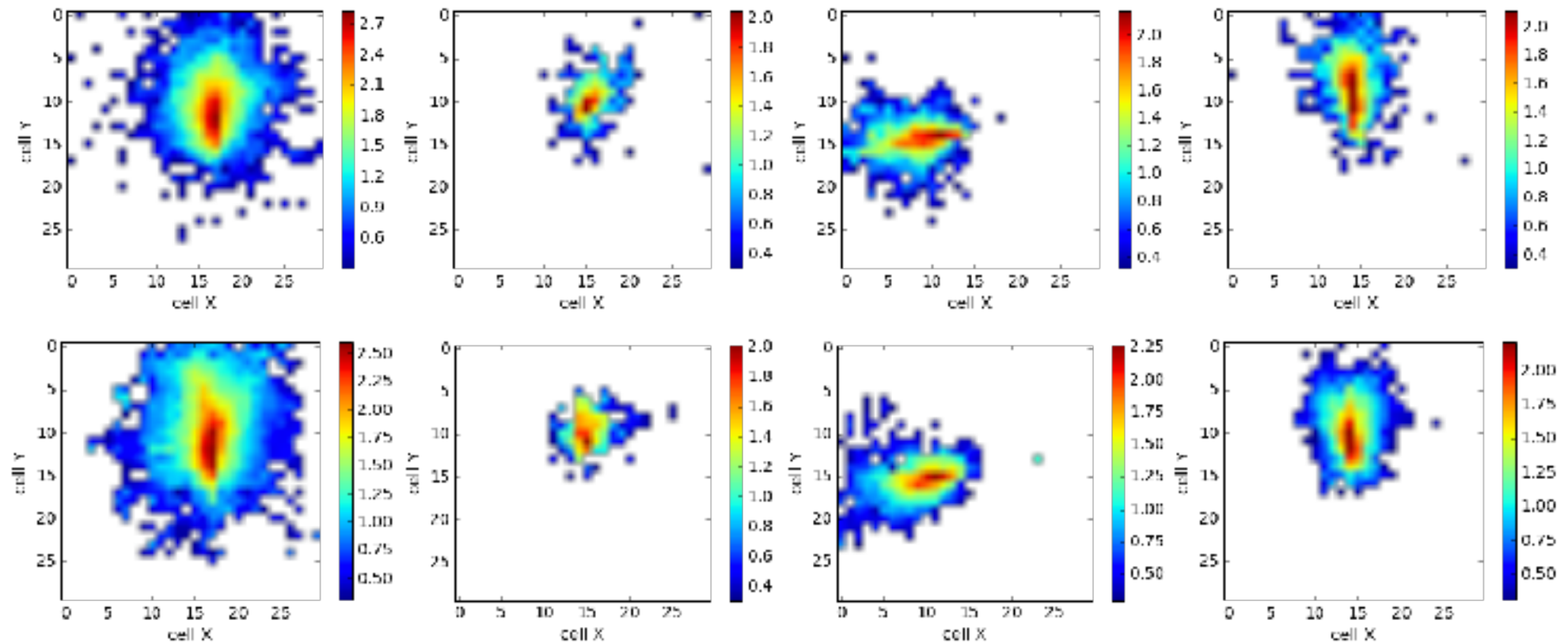


LHCb ECAL Simulation

GEANT Simulated

$\log_{10}(\text{cell energy})$

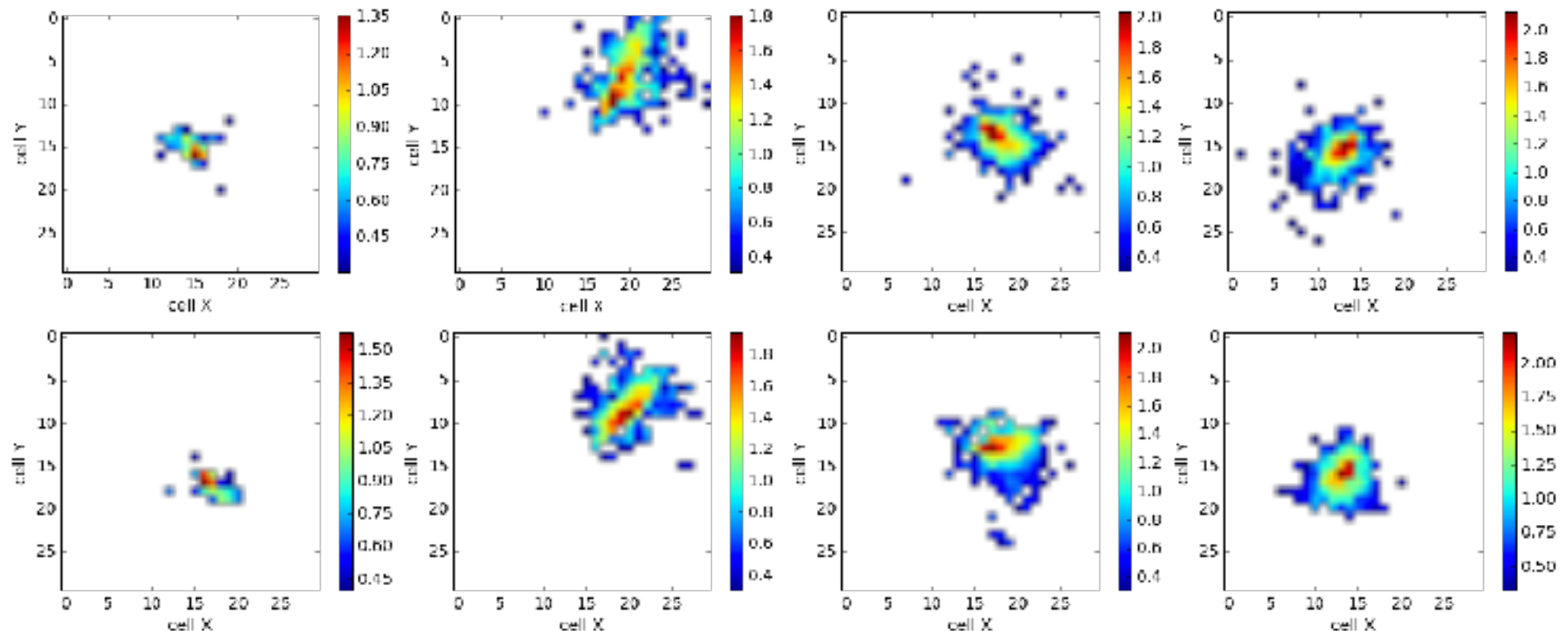
GAN Generated



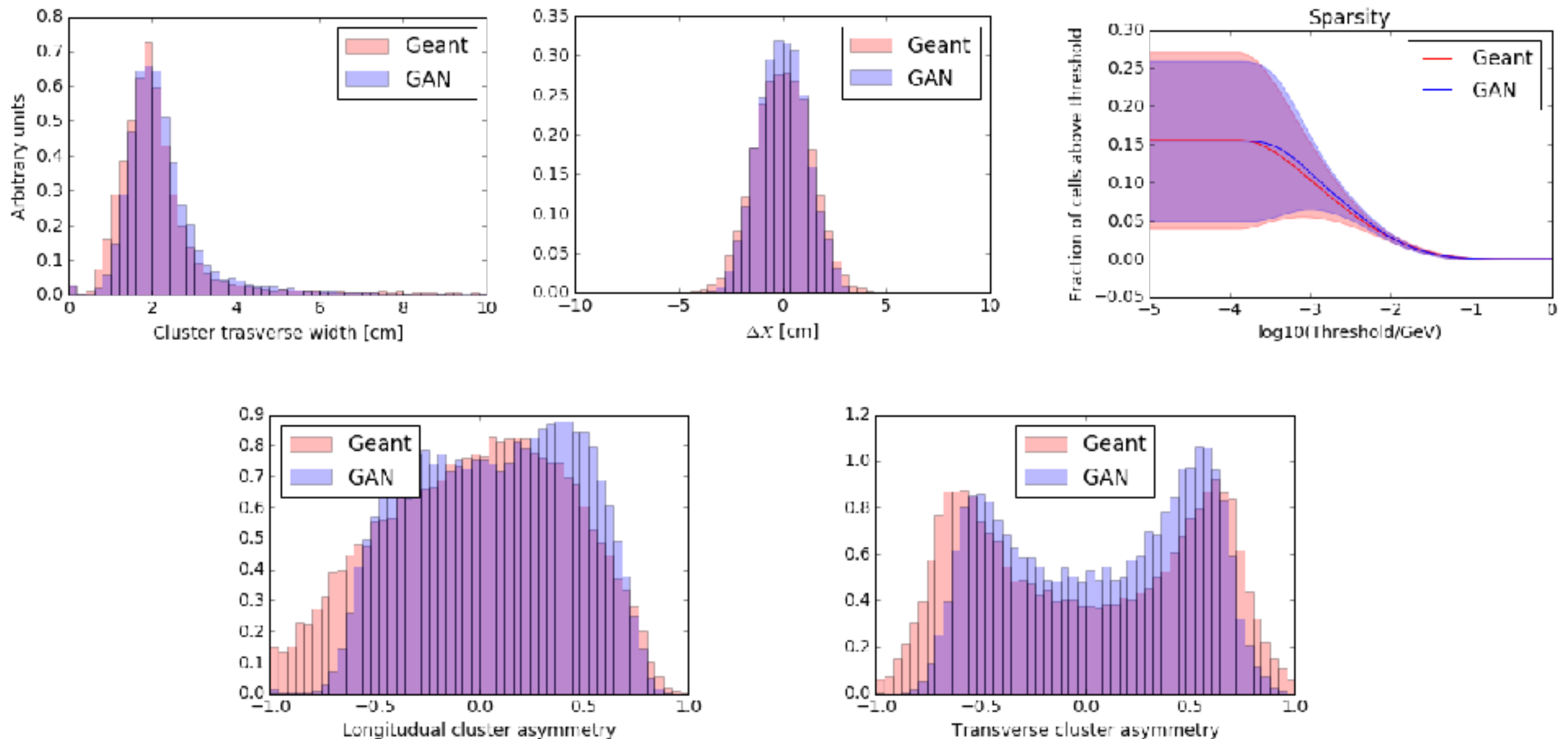
GEANT Simulated

$\log_{10}(\text{cell energy})$

GAN Generated



Primary and Marginal Distributions



- Is hard to fit marginal distributions
 - unless the model is aware that those are important for us

Scientific Requirements

- For image generation we are usually happy if the result **looks** like it is desired
- In science we need the result to reasonably well match the given set of requirements. Requirements are driven by **scientific considerations** closely connected to the ultimate scientific goal

Enforcing Important Statistics

- No generative model is ideal
 - some deviations from the original distribution remain
- Models tend to learn primary statistics of generated objects
- In physics applications, we often need for our model to learn particular statistics which are marginal for the generated object
 - e.g. cluster shape fluctuations for fast calorimeter simulation
- Can enforce these statistics by explicit adding them to the loss
 - can't we?

Enforcing Important Statistics

- Can enforce statistics by explicit adding them to the loss
 - can't we?
- By adding necessary statistics to the loss we do enforce match for these statistics
 - most likely by the price of overtraining these particular statistics
 - ... and we lose handle to validate quality of generator on this statistics
- Still can remove those statistics from loss, and see how far they would deviate
 - figure of merit for generating this statistics

Generating Tails

arXiv:1903.02433v2 [hep-ex] 22 Mar 2019

DijetGAN: A Generative-Adversarial Network Approach for the Simulation of QCD Dijet Events at the LHC

Riccardo Di Sipio,¹ Michele Fucci Giannelli,² Sana Ketabchi Haghighat,¹ Serena Palazzo.²

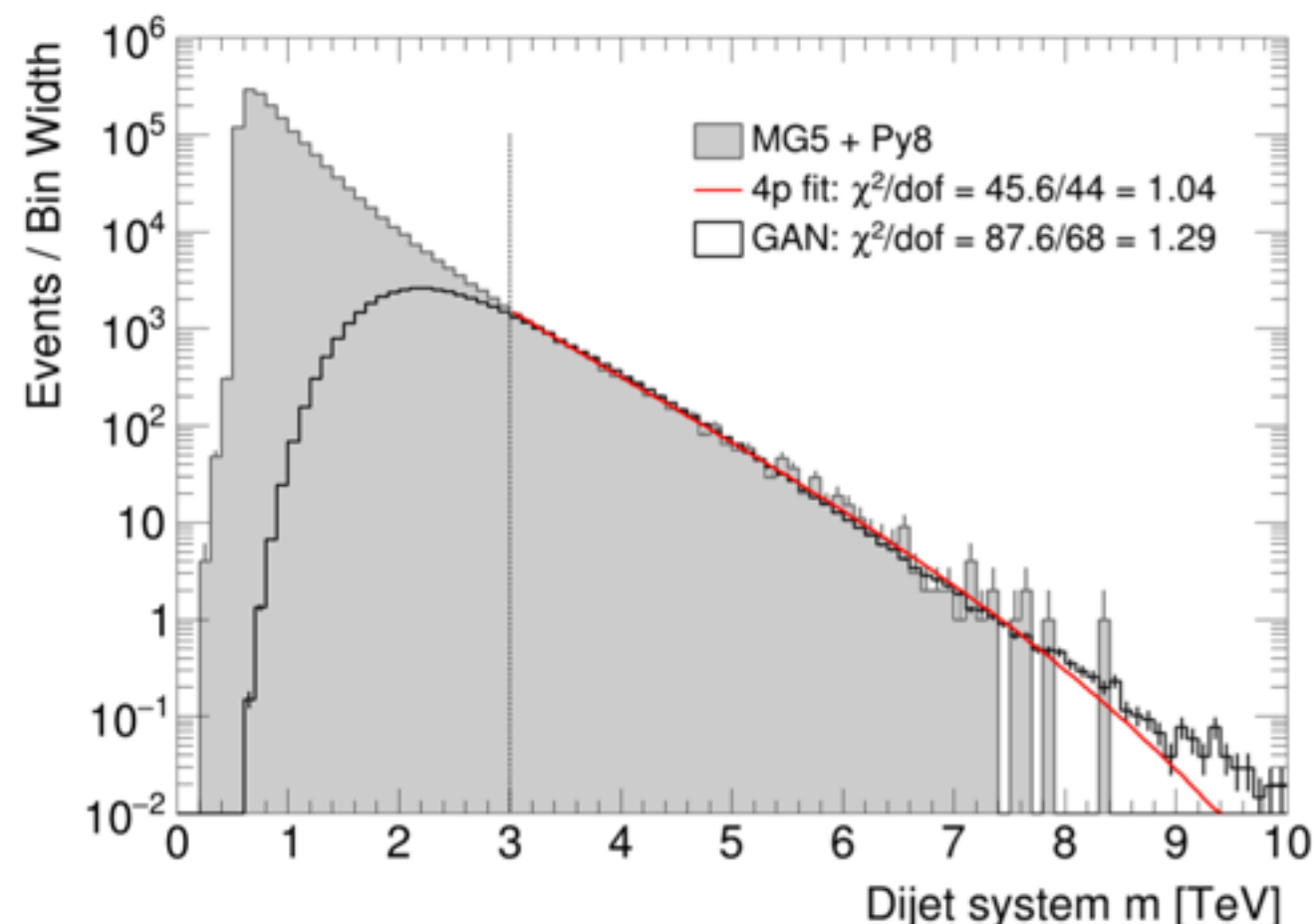
¹University of Toronto, Canada

²University of Edinburgh, UK

E-mail: riccardo.disipio@utoronto.ca,
michele.fucci.giannelli@ed.ac.uk,
sana.ketabchi.haghighat@mail.utoronto.ca, serena.palazzo@ed.ac.uk

ABSTRACT: A Generative-Adversarial Network (GAN) based on convolutional neural networks is used to simulate the production of pairs of jets at the LHC. The GAN is trained on events generated using MADGRAPH5, PYTHIA8, and DELPHES3 fast detector simulation. We demonstrate that a number of kinematic distributions both at Monte Carlo truth level and after the detector simulation can be reproduced by the generator network with a very good level of agreement.

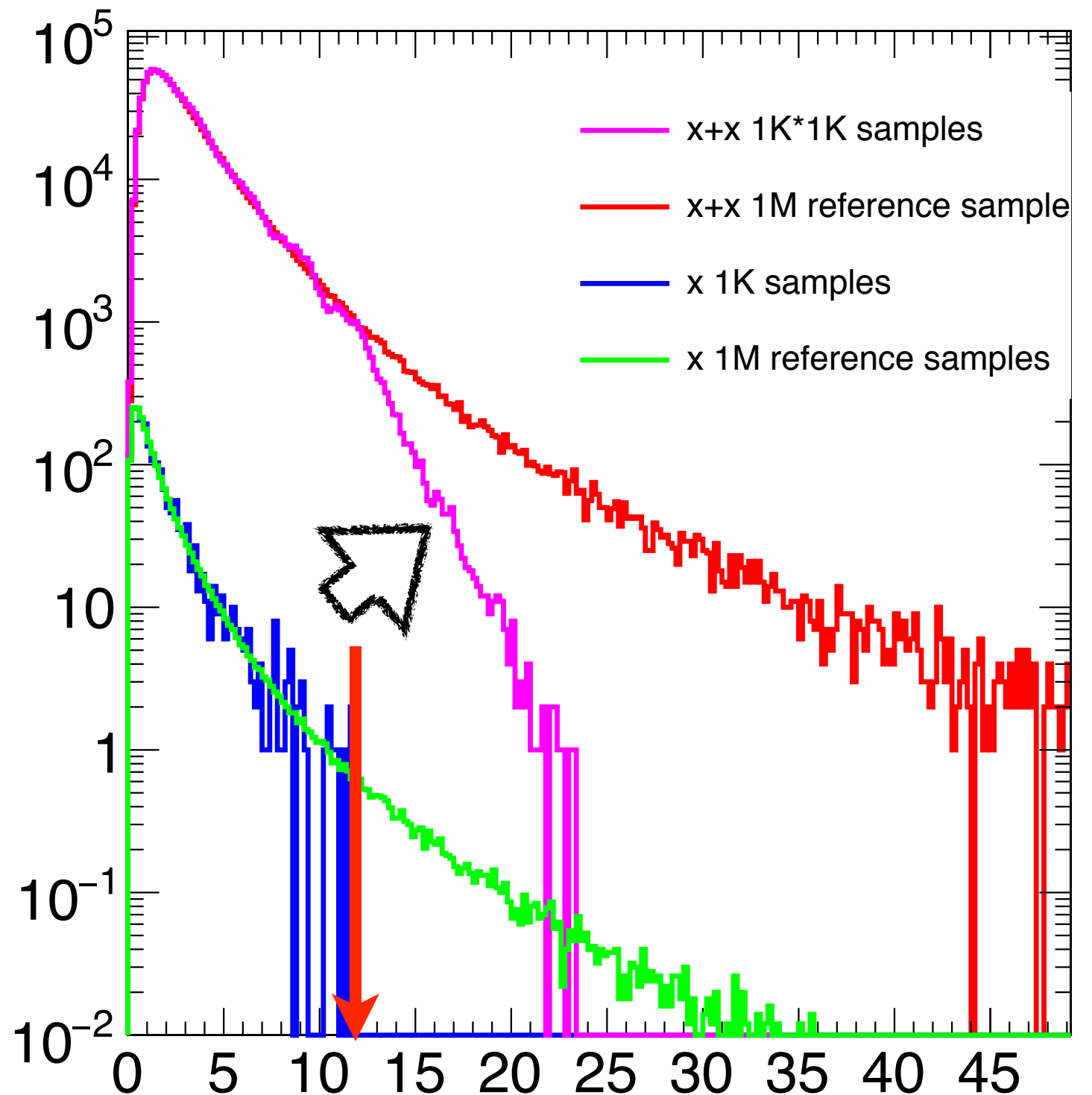
The code can be checked out or forked from the publicly accessible online repository <https://gitlab.cern.ch/disipio/DiJetGAN>.



- If the model is trained on the limited sample, how reliable are predictions beyond the training domain?

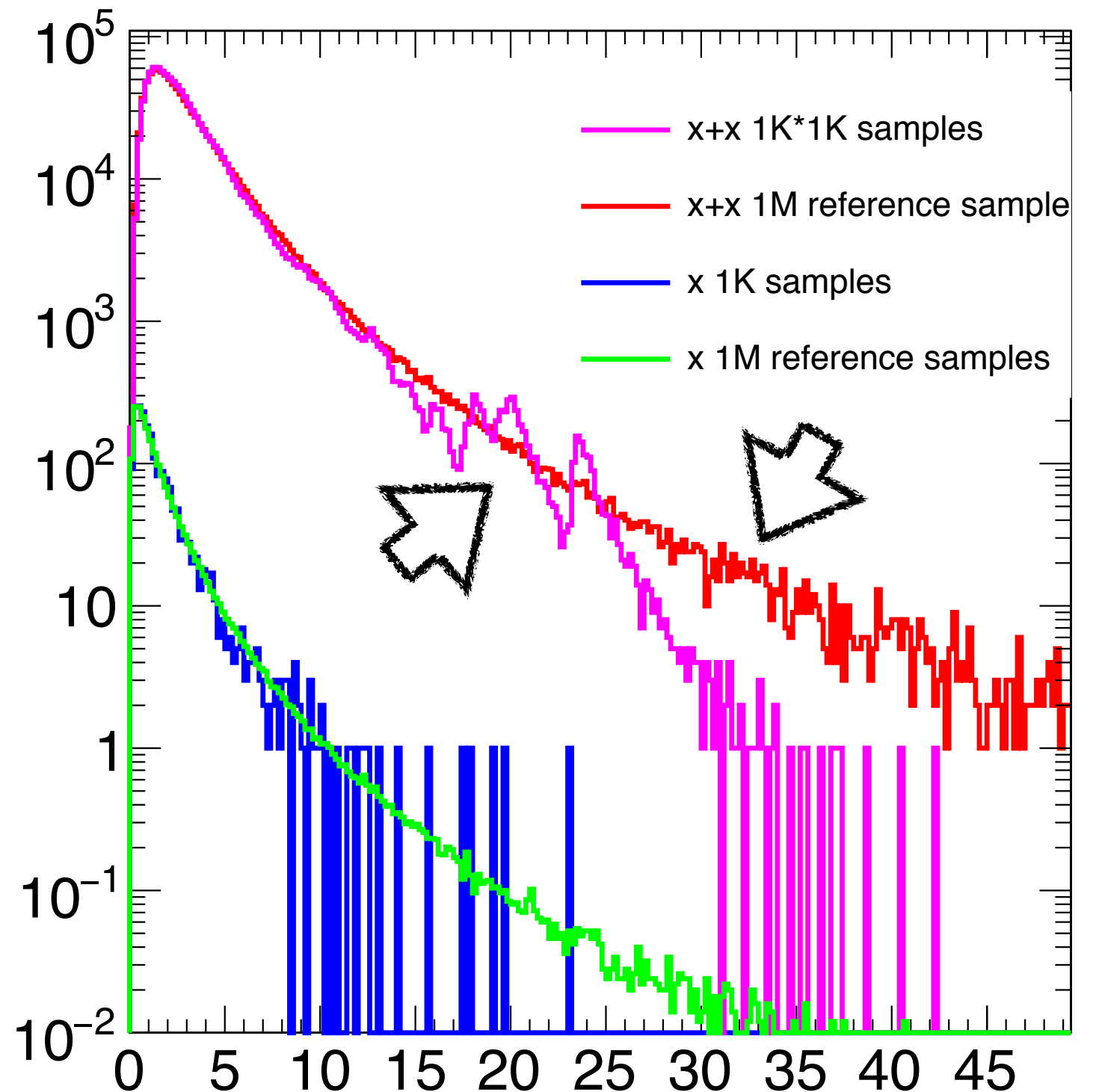
Asymptotic Properties

- Toy model
 - two variables distributed LogNormal
 - training sample 1K events ($x < 12$)
 - target sample 1M events $x_1 + x_2$
 - use 1K samples with permutations
- Systematics due to marginal cut off



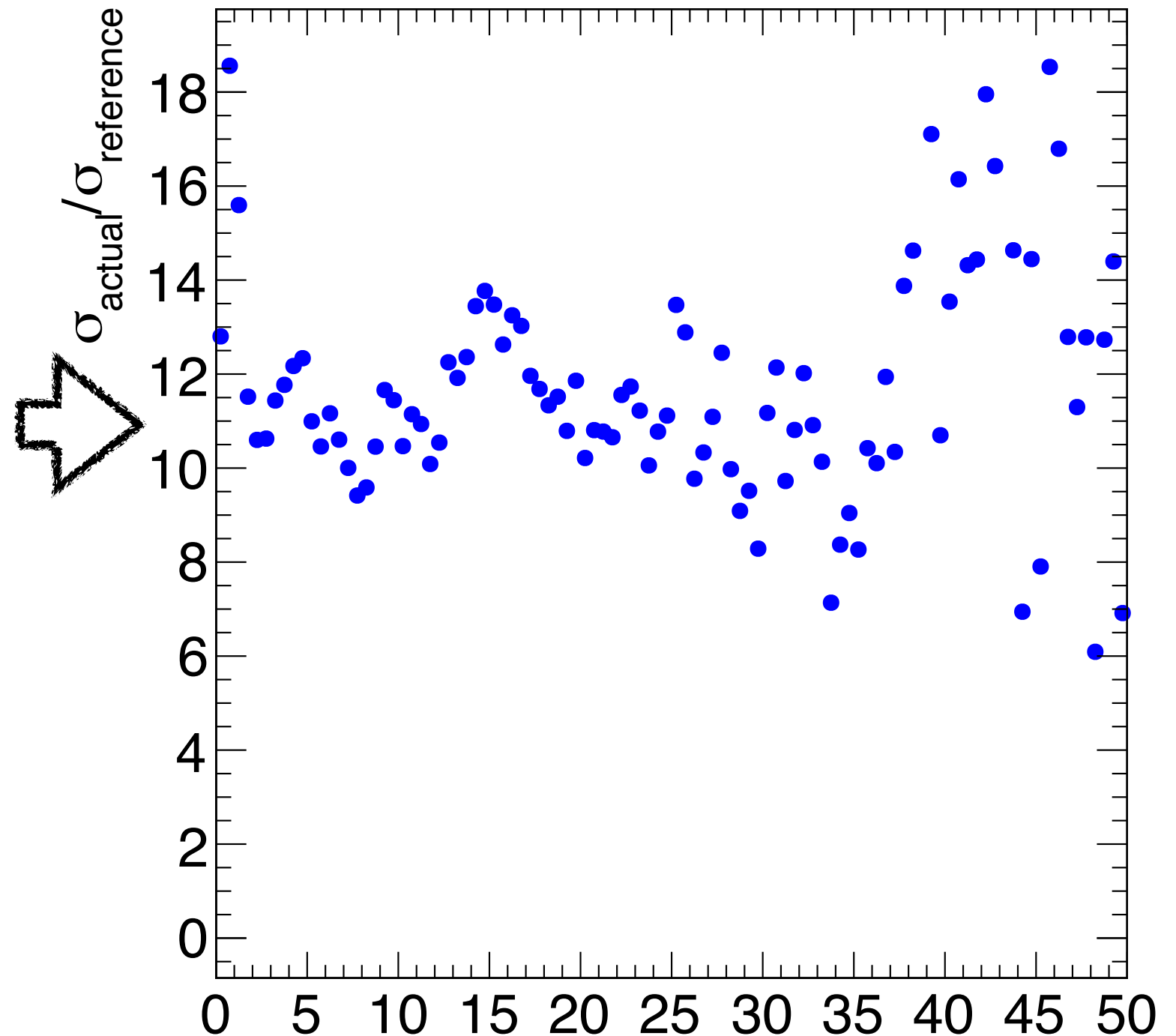
Asymptotic Properties

- Toy model
 - two variables distributed LogNormal
 - training sample 1K events
 - target sample 1M events x_1+x_2
 - use 1K samples with permutations
- Systematics due to fluctuation in tails



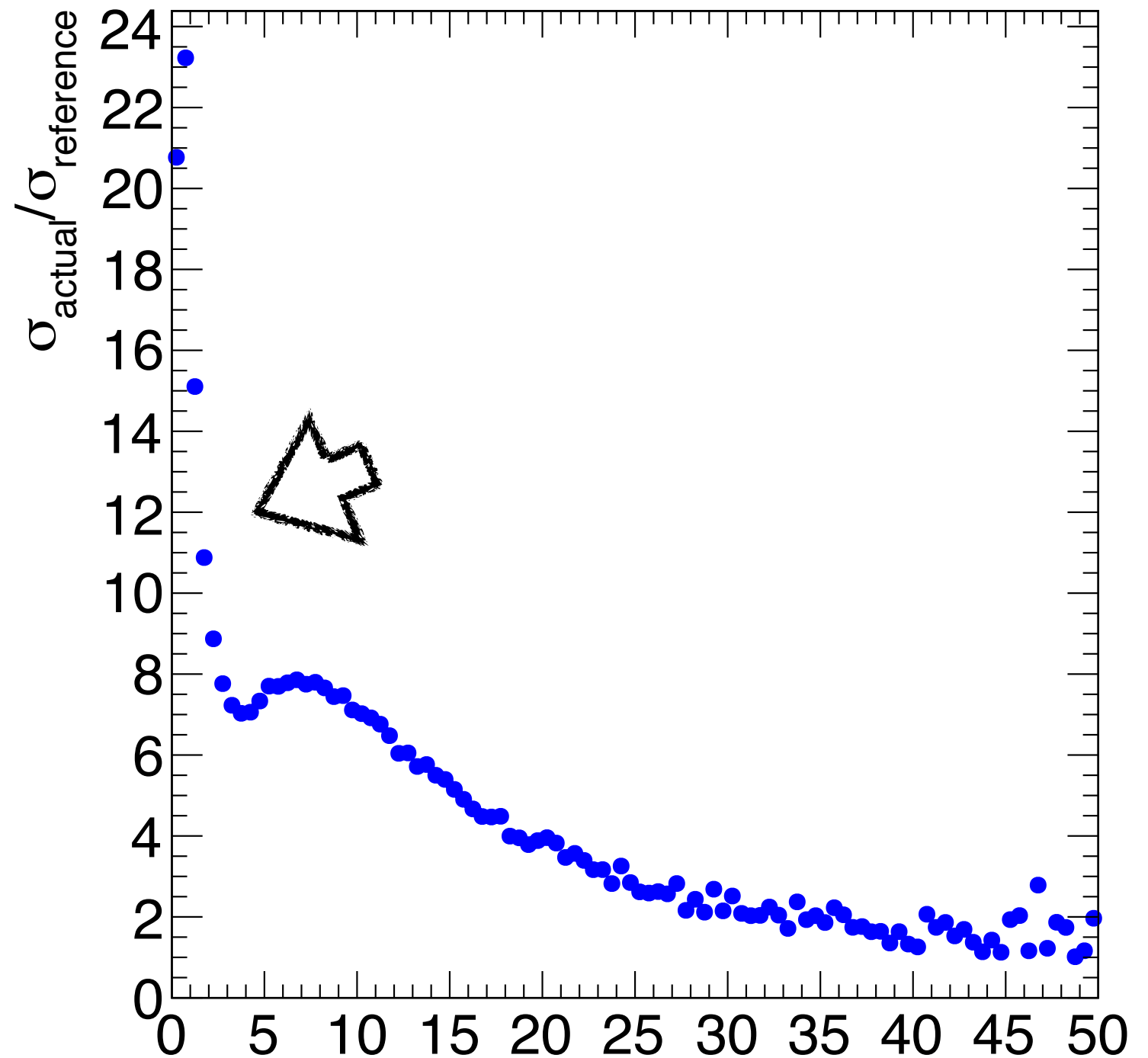
Statistics Properties

- Toy model
 - two variables distributed LogNormal
 - training sample 1K events
 - use 1K samples with permutations
 - variation for x_1+x_2
- Systematics due to the sample intrinsic correlation



Statistics Properties

- Toy model
 - two variables distributed LogNormal
 - training sample 1K events
 - fit LogNormal to 1K samples
 - variation for x_1+x_2
- Systematics from the model systematics

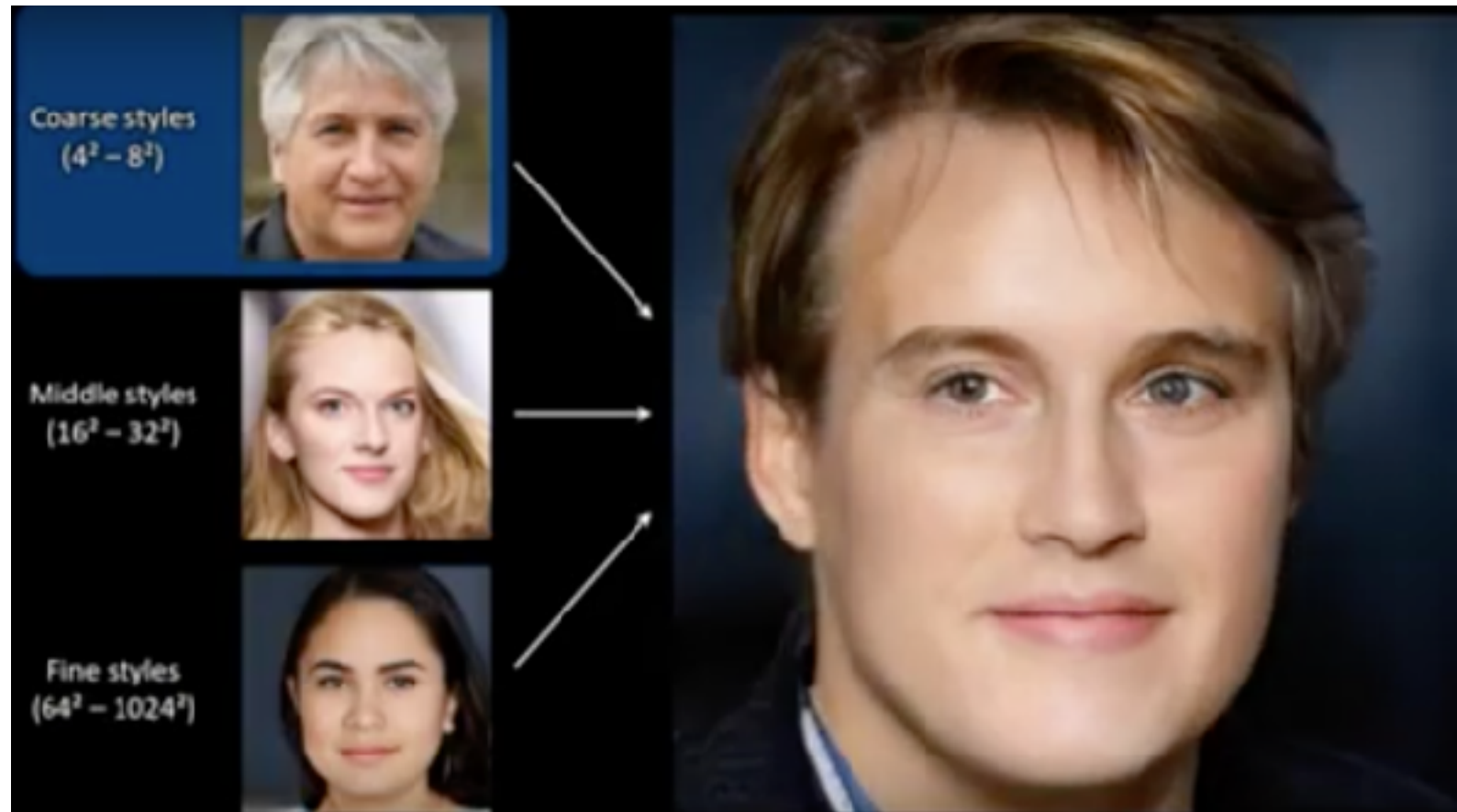


Decomposition

- Quality of the generative models is limited by the size of the train data sample
- generative models may not give profit for producing statistically correct big data sets
 - no information beyond the train sample is available
 - model systematics corresponds to the train sample statistics

Decomposition

- No information beyond the train sample is available



- Not quite if we can decompose generative model into separate components
 - random combinations of different components may drastically increase variability

Decomposition

- Quality of the generative models is limited by the size of the train data sample
 - generative models may not give profit for producing statistically correct big data sets
 - no information beyond the train sample is available
- Not quite if we can decompose generative model into separate components
 - random combinations of different components may drastically increase variability
- E.g. fast simulation of the calorimeter response
 - generator is trained on 10^6 incident particles
 - ~ 50 particles in the calorimeter per event
 - total variability $\sim (10^6)^{50} = 10^{300}$! (NB intrinsic correlation)

Quality Metric

- No generative model is ideal
 - some deviations from the original distribution remain
- Minor deviations are not that important e.g. for image generation
- Minor deviations may be a big deal for generative models in physics
 - e.g. we could want $E^2 - p^2 = m^2$ for generated particles to be precise
- Ultimate generative model quality metric is a comparing the final physics result obtained using generative model with the one obtained using the test data
 - accuracy is limited by the size of the test data

Conclusions

- Surrogate generative models demonstrate extraordinary progress in current years
- There are many applications for use in natural science research
- Generative models need attention to ensure scientifically solid results
 - satisfying boundary conditions, control of scientifically important but marginal statistics
 - appropriate evaluating the quality of the model
 - propagating model intrinsic systematics to the systematic uncertainties of the final scientific result

Generative Model. ML Perspective

- Generative models look very different from regression/classification models
 - actually they are not that different
- Consider set of objects each of which is described by a vector of parameters
 - we arbitrary split this vector into “features” \mathbf{x} and “labels” \mathbf{y}
- For classification/regression problem we search for deterministic function f which approximates dependency \mathbf{y} from \mathbf{x} : $\mathbf{y} = f(\mathbf{x})$
 - in probabilistic approach we search for probability $p(\mathbf{y}|\mathbf{x})$
- For generation problem we want to sample objects for a given label
 - we search for probability $p(\mathbf{x}|\mathbf{y})$
 - \mathbf{y} for generative model is called “condition”
 - condition may be absent - unconditional generative model

Generative Model. ML Perspective

- In both discriminative model and generative model we want to get probability for subset of object parameters conditioned by another subset of object parameters
- Discriminative models:
 - evaluate distributions for few, usually redundant, parameters conditioned by many features
 - can discriminate basing on this parameters
- Generative models:
 - evaluate many features conditioned by few parameters (conditions)
 - can sample these features
- NB: logistic regression + binomial distribution = generative model
 - for the binary objects