

Generative Models in High Energy Physics

Fast Simulation examples

Sofia Vallecorsa CERN openlab AI and Quantum Research CERN IT Department



13th Terascale Detector Workshop 06/04/2021

Content

- Introduction
- Architecture Trends
- Results Validation
- Computing Resources
- Summary and Conclusion

... a short teaser !



Deep Generative Models

- Shallow models learn simple internal representations
- Deep models allow higher levels of abstractions and improve generalisation
- Multiple applications
 - Simulation
 - Anomaly Detection
 - Data manipulation
- Different use cases have different requirements:
 - Fast inference
 - Real time training capability
 - Fast training for large optimizations





Simulation applications

S. Shirobokov *et al., Workshop on Real World Experiment Design and Active Learning at ICML* 2020, arxiv:2002.04632

Fast detector simulation ..

Examples in ATLAS, CMS, LHCb, ALICE.

Mostly calorimeters, but also RICH, TPCs..

but also:

Detector design

Optimisation

Domain Adaptation

Monte Carlo agreement to real data





Magnet system for muon beam optimisation





Multiple Architectures



Variational AutoEncoders

- Variational (KL Divergence), arxiv:1312.6114
- Wasserstein (MMD), arxiv:1711.01558
- Sinkhorn arxiv:1810.01118



A. Ghosh, Journal of Physics: Conference Series. Vol. 1525. No. 1. IOP, 2020.







	MNIST		HEP	
Model	FID	Sinkhorn	MAE	Sinkhorn
cond VAE	6.61	30.13	23.13 (±65.53)	14.92
cond WAE (MMD)	34.73	30.44	43.54 (±55.23)	34.46
cond e2eSAE (ours)	4.11	24.92	13.50 (± 29.82)	7.91
cond DCGAN	0.93	22.23	68.27 (±180.45)	6.95
original data	0.33	0	6.59	2.89

Generative Adversarial Networks

- Vanilla, Conditional, Auxiliary Classifier GAN
- Wasserstein, Cramer (arxiv:1704.00028)
- Complex topologies (ex. BiGAN, BIB-AE)





Wasserstein GAN for Cherenkov detector simulation:

D. Derkach et al., Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 952 (2020): 161804.



Detector response as images



G. Khattak et al., ICMLA 2019

F. Rehm *et al.*, AAAI 2020, arXiv:2103.13698

Convolutional layers for 51x51x25 pixels image: sparse, large dynamic range



Detector response as graphs

Calorimeter energy deposits as graphs

Data as a graph of connected hits

Connect hits using **geometric** constraints

Embedding requires large graphs (~10⁵ nodes)

GrapSAGE: W. Hamilton, R. Ying, J. Leskovec, NIPS (2017), arXiv:1706.02216:

Inductive approach to generate embedding for unseen nodes





1. Sample neighborhood

2. Aggregate feature information 3. from neighbors

Predict graph context and labe using aggregated information





Detector Size..

- FastCaloGAN: full calorimeter fast simulation using 300 GANs
- 100 η slices (0.05 wide)
 - 15 energy points Electrons, Photons, Pions
- Generated hits in Athena





M. Faucci Giannelli, 4th IML workshop, 2020

INFN

The GAN

06/04/2021

• Train the GANs on voxelised hits



Michele Faucci Giannelli

Need to take training time into account!

Results validation



Energy patterns along detector axes

> G4 GAN

Х

Х

G4 GAN

deposited energy [GeV

Physics validation

- Compare GAN images against Monte Carlo
- **Depending on the application need few percent accurate representation** of all relevant physics variables

G. Khattak et al., ICMLA 2019





G. Khattak et al., ICMLA 2019

Physics validation (II)

 Triforce* DNN has been developed to distinguish different kind of particles and measure their energy



*D. Belayneh et al., "Calorimetry with deep learning: Particle simulation and reconstruction for collider physics," 2019, https://inspirehep.net/literature/1770936



Systematic effects and interpretability





Systematics: image similarity

GAN can exhibit mode-collapse or mode-drop

How much **diversity** in the generated sample?

• Use the Structural Similarity Index

SSIM(\mathbf{x}, \mathbf{y}) = $\frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}$

where x, y are two samples to be compared

- Calculated on sliding windows, then averaged.
- Ours is a 3D problem: SSIM computed in *xy* plane, 3rd dimension is channel
- Adjust C1-C2 to the pixel dynamic range





Systematics: Support size

Empirical evidence of the GAN low support size (Arora and Sanjeev, 2017)

• Learnt distribution not representative enough

Birthday paradox test (Brink, 2012):

How many people need to be in one room so that P(at least two people have same birthday) > 0.5 ?

• 365 days in a year \rightarrow 23 people is enough

Generalized problem:

How many samples is it necessary to generate to have P(at least one pair of duplicates among the samples) > 0.5 ?

• (The answer)² = estimate of the support size







Birthday paradox for GANs

Original birthday paradox problem

- Days in a year finite set of possible values with discrete uniform distribution
- Unique duplicates definition people born on the same day

GAN distribution

- Images pixels of continuous values
- Multivariate continuous distribution → occurrence of exact duplicates has zero probability
- Duplicates as "similar enough" images

Similarity metrics depend on the use case and data type

Exact duplicates



Not exact duplicates But similar enough?







Support size estimates

- GAN samples significantly more similar \rightarrow **smaller** support size
- Test depends strongly on duplicates definition



Not adapted to our problem?







Systematics: rare events

In some cases it is important to reproduce correctly the topoloov and occurrence of rare even

Percentage of events with multiple peaks

250

350



"Standard"

10-1



Computing resources



Faster then Monte Carlo?

Post training quantization (INT 8) using Intel DLBoost and iLoT tool

FP32: 3DGAN is **38000x faster** than Monte Carlo (on Intel Xeon processors) INT8: quantized 3DGAN is **68000x faster** than Monte Carlo (FP32)

CERN 3D-GANS Inference FP32 & INT8 (DL Boost) Operation Times per Batch on 1S Intel(R) Xeon(R) Scalable Processor 8280







Reduced data representation reduce inference time but reduce phsyics performance

Need ad-hoc optimisation strategy





Training time

CERN

penlab

- Training the 3D convolutional GAN model (3M parameters) takes about 7 days on a GPU
 - Distributed training is essential
 - Need to keep physics under control
- Tested different data parallel approach on different hardware on HPC and Cloud









Access to Cloud resources through CloudBank EU project

Summary

- High Energy Physics experiments are heavily investigating Generative Models for fast simulation
 - High level of customization
 - Increasing complexity
- Domain-specific knowledge is key
 - Model design and architecture **otpimisation**
 - Results validation
 - Models interpretability and performance systematic studies
- Efficient use of computing resources broadens the scope





Teaser: Quantum Generative Adversarial Networks

Explore quantum advantage in terms of:

- Compressed data representation in quantum states
- Faster training with smaller number of parameters
- Support space of the learned distribution

Simplify 3DGAN simulation problem

1D & 2D energy profiles from detector Train a **hybrid classical-quantum** GAN eggine a give a second second



S. Y. Chang et al., arXiv:2103.15470 arXiv:2101.11132



QUANTUM TECHNOLOGY

INITIATIVE



Thanks!

Sofia. Vallecorsa@cern.ch

https://openlab.cern/

https://home.cern/



A. Butter et al., arxiv:2008.06545

Systematics: increasing statistics

0.18

0.16

0.14

0.12

0.10 2 0.08

0.06

0.04

0.02

0.00

- If a GAN is trained on **N** data points, how many **new** points can be drawn?
- GAN can describe distribution better than training data
- Needs 10,000 GAN points to match 150 true points
- In terms of information:
 - **sample**: only data points

penlab

- fit: data + true function
- GAN: data + smooth, continuous function



Most physics data sets described by continuous function \rightarrow GAN can interpolate