

Generative Models in High Energy Physics

Sofia Vallecorsa

February 3rd, 2020



Conclusions

Generative Models are powerful techniques

Object of an intense R&D , require careful validation

Many applications in High Energy Physics

Simulation is one of the main applications

High level of accuracy & speed

Multiple features/capabilities

Deep Learning is a driving force toward changing our computing model

DL workloads are HPC friendly

Benefit from dedicated hardware and accelerators







Deep Learning in HEP

DL recognize patterns in large complicated data sets Better performances if applied directly to raw data

Re-cast physics problems as "DL problems"

Detector output as **images** and apply **computer vision** techniques Physics events as **sentences** and apply NLP techniques

Intense R&D activity

i cern **openlab**

Performance requirements Model interpretability

Results Validation against classical methods Detailed Systematics studies

"New" computing models

Accelerators and dedicated hardware HPC integration Cloud environment & Big Data platforms



Generative Models

Basic concepts

Some references:

CERN Openlab

- NIPS 2016 Tutorial on Generative Adversarial Networks, I. Goodfellow
- <u>2018 Generative Models tutorial,</u> <u>D. Rezende</u>
- <u>2019 Stanford course on Deep</u> Generative Models



"GM Taxonomy" Image from NIPS 2016 Tutorial, I. Goodfellow

5

GM as Density Estimators

Generative models learn probability distributions from data

Given a collection of samples x_i and the underlying p_{data} distribution Choose parameters { θ } so that $p_{\{\theta\}}(x) \approx p_{data}$



6



Prior knowledge?

Generative models use *some* prior knowledge



Extract meaningful **representations directly from data** but need:

Loss functions Learning principle Optimisation algorithms Domain knowledge

CERN openlab

Probability divergence D(p,q)

Name	Formula
f-divergences	$D(p;q) = \mathbb{E}_q[f(rac{q}{p})]$
Relative entropy (KL)	$D(p;q) = \mathbb{E}_q[\ln rac{p}{q}]$
Jensen-Shannon (JS)	$D(p;q) = rac{1}{2}\mathrm{KL}(p;m) + rac{1}{2}\mathrm{KL}(q;m)$
Stein divergence	$D(p;q) = \sup_{f} \mathbb{E}_q [abla \ln pf + abla f]^2$
Energy distance	$D(p;q) = \mathbb{E}[2\ x - y\ - \ x - x'\ - \ y - y'\]$
Wasserstein distance	$D_lpha(p;q) = [\inf_ ho \mathbb{E}_ ho[\ x-x'\ ^lpha]]^{rac{1}{lpha}}$
Max-min dis. (MMD)	$D(p;q) = \sup_{f} (\mathbb{E}[f]_p - \mathbb{E}[f]_q)$

D. Rezende: https://danilorezende.com/2018/07/12/short-notes-on-divergence-measures/

O. Cerri, ACAT2019, arXiv:1811.10276

Multiple Applications

Event sampling (**Simulation**) Find Underlying Factors (**Discovery**) Detect Rare events (**Anomaly Detection**) Predict future events (**Planning**) Find Analogies (**Transfer Learning & Style Transfer**)







(Deep) Generative Models

Internal representations learned by shallow systems are simple (Bengio & LeCun 2007, Bengio 2009)

- Incapable of learning complex hidden structures
- Large amounts of labeled data
- \rightarrow Deep Generative Models
 - \rightarrow Higher levels of abstraction
 - → Improved generalisation and transfer
- \rightarrow Categorized according to
 - → Feature representation
 - → Fully observed, or latent variables based
 - \rightarrow Learning principles
 - → MLE, variational, or likelihood-free



CERN CERN

Fully Observed Models

Directly observe data without introducing new local (latent) variables

$$p(x_{1,...,N}) = \prod_{i=1}^{N} p(x_i | x_{1,...,(i-1)})$$

Ex. Pixel Recurrent Neural Networks:



$$p(x_t \mid x_{1:t-1}) = p(x_t^{red} \mid x_{1:t-1})p(x_t^{green} \mid x_{1:t-1}, x_t^{red})p(x_t^{blue} \mid x_{1:t-1}, x_t^{red}, x_t^{green})$$

 x_2

 X_1

 x_3

 X_{1}





Latent Variables models

Unsupervised representation learning

Introduce an **unobserved random variable** for every observed data point to explain features. Prescribed models: Use observer likelihoods and assume observation noise.

Implicit models: Likelihood-free models.

$$x \in \mathbb{R}^{d_x} \quad z \in \mathbb{R}^{d_z} \quad \theta \in \mathbb{R}^{d_\theta}$$
$$\mathcal{D} = \{x_i\} \quad i \in \{1, ..., N\}$$

$$\log p_{\theta}(x) = \log \int p_{\theta}(x|z) p(z) dz = \log \mathbb{E}_{p(z)}[p_{\theta}(x|z)]$$

$$\log p_{\theta}(\mathcal{D}) = \sum_{i=1}^{N} \log \mathbb{E}_{p(z)}[p_{\theta}(x_i|z)]$$



$$p(z)$$

$$z$$

$$p(x|z)$$

Variational AutoEncoders

Describe training dataset in latent space

Explicit constraints on the encoded representations

Learn the **latent variable** distribution

Two components in the loss function

reconstruction loss

CERN

🚅 openlab

KL divergence between the learned latent distribution and the prior regularization





Likelihood-free learning

Density estimation by comparison

CERN Openlab

Compare the **estimated distribution** q(x) to the **true distribution** $p^*(x)$ using samples.

Build an **auxiliary model** to indicate how data simulated from the generative model differs from observed data.

Adjust model parameters to better match the data distribution



Adversarial training

Assume a deterministic generator: $\mathbf{x} = G_{\theta}(\mathbf{z})$

A prior over latent space:

$$\mathbf{z} \sim p_{\lambda}(\mathbf{z})$$

Define a discriminator:

CERN Openlab $D_{\psi}(\mathbf{x}) \in [0, 1]$

min max $\mathbf{E}_{x \sim \mathcal{D}_{real}}[D_{\psi}(x)] - \mathbf{E}_{h}[D_{\psi}(G_{\theta}(h))]$

A learnable loss function from the min-max game

$$\min_{\theta} \max_{\psi} \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\ln D_{\psi}(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{z} \sim p_{\lambda}(\mathbf{z})} \left[\ln \left(1 - D_{\psi}(G(\mathbf{z})) \right) \right]$$



Wasserstein loss

14

arXiv:1406.2661v1

Generative adversarial networks

Simultaneously train two networks that compete and cooperate with each other

Generator G generates data from random noise Discriminator D learns how to distinguish real data from generated data





https://arxiv.org/pdf/1701.00160v1.pdf

Multiple GAN flavors

Different base layers technologies: Original GAN was based on MLP in 2014 Deep Convolutional GAN in 2015 Graph GAN Recurrent -GAN

Different "information paths"

CERN

ī: 🚅 openlab

Conditional GAN Extended to learn a parameterized generator p_{model}(x|θ) Useful to obtain a single generator object for all θ configurations Interpolate between distribution Auxiliary Classifier GAN D can assign a class to the image Progressive growing GAN VAE-GAN Stack GAN BiGAN



Conditional GAN (Mirza & Osindero, 2014)

arXiv: 1411.1784



AC-GAN (Present Work)

<u>arXiv:1610.0958</u>

16







https://thispersondoesnotexist.com/







GAN for Art?

Image-to-Image Translation with Conditional Adversarial Networks



http://www.memo.tv/portfolio/learning-to-see/

https://github.com/phillipi/pix2pix¹⁸

Evaluating performance is tricky

Different applications

Density estimation

Sampling/generation

Latent representation learning

Custom task (i.e. image translation, compressions..)

Tractable likelihood models

Split dataset into train, validation, test sets

- Evaluate gradients on train set
- Tune hyperparameters on validation set
- Evaluate generalization by reporting likelihoods on test set
- Non-tractable likelihood or likelihood-free models Use lower bounds or approximations (I,e ELBO, KDE)



https://deepgenerativemodels.github.io/syllabus.html



GAN for simulation

Measure convergence

Evaluate difference between model PDF and real PDF

At convergence evaluate sample quality

- Mixing and coverage (diversity)
- Saliency
- Mode collapse or mode dropping
- Overfitting

Need quantities that are **invariant** to small translation, rotation, intensity changes

Simple pixel space Euclidean distances don't work

Define a way to map input into a feature space

- **Inception score**
- Maximum Mean Discrepancy
- **Fréchet Inception Distance**

openlab

Structural Similarity Index

+ Physics Quantities Validation

Generated support space

Size of the dataset

GANs produce distributions with limited support

Support size grows ~linearly with discriminator size Training dataset size does not help much for a given discriminator

Example: BiGAN support size is around 1M (training set ~200k)

Birthday paradox test*

If a sample of size **s** has near-duplicate images with prob > 1/2, then distribution has only **s**² distinct images.



Birthday Paradox: Suppose a distribution is supported on N images. Then P[sample of size \sqrt{N} has a duplicate image] > $\frac{1}{2}$.

Search for nearest neighbor



One extreme case: Mode collapse

Goal is to generate fake examples imitating real sample Simple solution is to just generate easy modes (classes).





Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks (2016).

Duality Gap

A natural metric in minimax games

For a given set of Dicriminator and Generator: $(u^*, v^*) \in K_1 \times K_2$

$$DualGap := \max_{v \in K_2} M(u^*, v) - \min_{u \in K_1} M(u, v^*)$$

Fix generator, find most adversarial discriminator

Fix discriminator, find most adversarial generator



Conclusions

Generative Models are powerful techniques

Object of an intense R&D , require careful validation

Many applications in High Energy Physics

Simulation is one of the main applications

High level of accuracy & speed

Multiple features/capabilities

Deep Learning is a driving force toward changing our computing model

DL workloads are HPC friendly

Benefit from dedicated hardware and accelerators



3DGAN







https://openlab.cern/project/fast-simulation



25

Compact Linear Collider

High-luminosity linear e+e- collider Three energy stages up to 3 TeV Associated detector studies





Electromagnetic calorimeter detector design 1.5 m inner radius 5 mm×5 mm segmentation 25 tungsten absorber layers + silicon sensors

http://cds.cern.ch/record/2254048# 26

High Granularity Calorimeter

Open data set developed for ML applications ⁽¹⁾

Pixelized datector volume to fully contain shower

1M single particle samples (e,γ,π) with flat energy spectrum (10-500) GeV

 α =+/- 30° random incident angle

Detector response as 3D images (51x51x25 pixels)









3D Generative Adversarial Networks

Condition training on input variables, Custom losses

Auxiliary regression tasks assigned to the discriminator



Convergence

and discriminator performance

Stable test loss

Discriminator Real/Fake probability peaks at ~50%

Correct incident angle











Single cell energy





G4

0.03

31



32



Internal Correlations



Energy shower shapes





Position along Z axis

α= 120°



Transfer learning: extending the energy range



- Transfer learning from 100-200 GeV pre-trained network
- Double training dataset statistics, 3 epochs training

Improved correlation description!



Further Validation (I)

Structural Similarity Index (SSIM) is used to assess image similarity

Measure diversity in GAN generated images

Run Triforce on GAN/GEANT4 events

Trained on Geant4 data

500 1.0 400 0.8 ш° Predicted 300 0.6 0.4 100 0.2 0.0 CERN GAN GEANT openiau

electron identification

SSIM (scale)	SSIM G4 vs G4	SSIM GAN vs GAN
1	0.94	0.95
1e-2	0.21	0.25
1e-4	0.045	0.061
1e-6	0.045	0.051



Triforce, Matt Zhang, https://github.com/BucketOfFish/Triforce_CaloML

Further validation (II)

Train TriForce classifier on GAN data

Similar to **Classification Accuracy Score** idea⁽¹⁾ Train Triforce with a mix of GAN (electrons) and Geant4 (pions) data

Validate performance on Geant4:

CERN

Obtain very good results (99% AUC)







Some electrons are not classified correctly

Computing performance

20000x faster than Geant4



Convergence in 15min

CERN High Energy Physics: 3D GANS Training Performance Intel 2S Xeon(R) Cluster, OPA Fabric Xeon(R) 8268 (2019) vs Xeon(R) 8160 (2018)





Other simulation examples at the LHC





ATLAS LAr calorimeter

Complex geometrical structure

Wasserstein GAN: One generator against two critics







DijetGAN

QCD dijet events at the LHC

- Investigate regions of the phase spaces with **low**
- cross-section → interesting kinematic region for BSM searches
- GAN can produce "**reasonable**" output beyond the training data phase space!





CMS HGCAL prototype

High Granularity and Hexagonal cells for CMS upgrade

Wasserstein conditional GAN (convolutions)

Train a generator against a critic and 2 constrainer networks reconstructing energy and impact point coordinates 1/N dN/dX [a.u. Geant4 70 GeV

Good agreement to Geant4

Some problems at low energy



0.15

0.1



Geonté, 70 GeV (

eant4 90 GeV WGAN 32 GeV e

WGAN 70 GeV e

I/N dN/dX [a

Seant4 90 GeV

WGAN 70 GeV

WGAN 90 GeV

0.06

0.04

Other applications in HEP (simulation)

- Generative models for ALICE TPC simulation (ACAT2019)
- GAN in LHCb: calorimeter and RICH simulation (CHEP2019)
- Graph-GAN for CMS HGCAL (ACAT2019)
- Variational AutoEncoders to simulate ATLAS LAr calorimeter (CHEP19)
- Wasserstein GANs to generate high-level physics variables based on Monte Carlo ttH (superfast-simulation) (IML WG 04/18)
- Particle-GAN for Full Event Simulation at the LHC (ACAT2019)
- **Refining Detector Simulation using Adversarial Networks (CHEP2019)**

Model-Assisted GANs for the optimisation of simulation parameters (AISIS2019)

,;;;;;; CERN



Conclusions

Generative Models are powerful techniques

Object of an intense R&D , require careful validation

Many applications in High Energy Physics

Simulation is one of the main applications

High level of accuracy & speed

Multiple features/capabilities

Deep Learning is a driving force toward changing our computing model

DL workloads are HPC friendly

Benefit from dedicated hardware and accelerators

CERN III CERN



Addressing Computing challenges

Millions of operations

Mostly matrix-multiplications

HEP models are designed and **optimised for specific tasks**

Generally custom models

~**Fewer weights** and operations than out-of-the-box tools

Higher accuracy

Depending on the task, we might need:

Fast inference

Online training capability

Fast training for large optimisations



Accelerators

Deep Learning workloads are naturally accelerator friendly

Large number of frameworks and ecosystems to simplify deployments Can work with half precision arithmetic (16FP, ...)

GPUs are de-facto standard to run DL

R&D on reducing bottlenecks (memory size, I/O. ...)

FPGAs can provide low latency inference

Network compression/quantization/parallelisation

Different programming approaches Hardware Description Language vs High Level Synthesis

Frameworks exist that "compile" ML code for different hardware Customisation

Available in cloud environments for **on-demand access**

Initial tests to time inference on cloud vs local

CERN

🖉 openlab





Heterogeneous computing



UF

FPGA-acceleration of **3DGAN** inference

Intel PAC Arria 10 card; OpenVINO; DLA design suite Implement missing components:

- 3D convolutions (done)
- 3D upsampling (inprogress)



DLA customization



Distributing the training process

Data distribution

openlab

- Compute gradients on several batches independently
- Update the model synchronously or asynchronously
- Applicable to large dataset

Data is sent and received from a single node

Can have high communication costs Smart update strategies Computation or bandwidth bottleneck

Gradients Model Gradients **Training Process** Data Store Averaged Gradients Model Gradients **Training Process** Averaged Model Gradients Gradients 1. Read Data 3. Average Gradients 4. Update Model 2. Compute Model Updates (Gradients) Send GPU 1 GPU 2 GPU 0 GPU 3 Reducer GPU 4 ← Receive

Averaged

47

Training Process

HPC resources

HPC "creates" new DL models

Most powerful systems are hybrid Ease access to the resources

Integration in HEP infrastructure

Software stack deployment ?

Containers

Most DL framework provide containerised versions

Workflow management?

Native **DL platforms** are natural choice (i.e. Kubeflow) Adapt HTC job schedulers ?

Data access/management?

S3 works very well Supported widely in **commercial clouds** Not all HPC centers allow access to it



SURF SARA

Distributed training on Cascade Lake

Training 3DGAN can take days

Implements data parallel approach using Horovod.

128 Nodes:

Xeon 8160 (Sky Lake): ~2 Mins/Epoch

Xeon 8268 (Cascade Lake): < 1 Min/Epoch

256 Nodes:

Time to Convergence: ~15 Mins

CERN High Energy Physics: 3D GANS Training Performance Intel 2S Xeon(R) Cluster, OPA Fabric Xeon(R) 8268 (2019) vs Xeon(R) 8160 (2018)



Intel Endeavour cluster: Xeon® 8268 Cascade Lake, 2 Sockets /node, 24 cores per socket, Intel® Omni-Path Architecture Software: Tensorflow 1.14 (Intel optimized), MKL-DNN 0.18, Horovod 0.16.4, Keras 2.2.4

Containerized deployment on HPC



Transition AI algorithms from laptops to supercomputers

Deploy 3DGAN training on SuperMuc-NG using CharlieCloud

Intel Xeon Platinum 8174 (Skylake)

48 cores

Intel® Omni-Path Architecture





Cloud resources

A solution to consider for offline applications

Not always feasible/effective to **buy specialized hardware**

Most MLaaS solutions are not customizable enough for scientific use cases

Great opportunities for R&D with industry New initiatives to increase access to commercial clouds deploy hybrid models (OCRE in the

context of the EOSC)

3DGAN training on public cloud: Docker + Kubernetes/Kubeflow



Bonus Material

Refining Simulation using GANs : Ultra High Energy Cosmic Rays Pushing the boundaries: Quantum GAN





52

Refining Simulation using GANs

Pierre Auger Observatory

- Detection of UHECR E>10^{17.5} eV
- Hybrid Technique
 - 27 Nitrogen Fluorescence (ultra-violet) telescopes
 - 1660 Surface detectors (water tanks)
- 3000 km² array size
- Simulation/Data mistmatch Refine simulation using WGAN Train DNN on refined simulation









arXiv:1802.03325

www.auger.org

Refining Simulation using GANs

DNNs very sensitive to simulation / data mismatches Simulation underestimates the muon component Generate 2 datasets "data-like" and "simulation-like Run DNN regression on them

 \rightarrow DNN does not fit the "data-like" events because it is trained on simulation



Refining Simulation using GANs





Why Quantum Machine Learning?

Quantum approach to ML could solve more complicated problems... faster

ML based tool can recognize complicated (hidden) patterns in data

Quantum processors can produce statistical patterns that are computationally difficult to produce with classical approaches

→ Could quantum processors recognize more complicated patterns in data?

Defining what quantum speed-up means is a complicated task Need to compare to the "best available" classical algorithm





Quantum ML

... and ML for Quantum Computing

QML introduces quantum algorithms as part of a larger implementation Fully quantum or hybrid classical/quantum approaches Input data could be quantistic \rightarrow ML for QC

How do we construct Quantum Neural Networks (QNN)?

Direct association between neurons and qubits

- Encode information into amplitudes of a quantum state
- How do we represent learning rules?
 - Need association rule between NN activation patterns and pure quantum states

How do we address data loading?

- Quantum state preparation
- Direct access through qRAM ?

Possible to train on large datasets by only loading a small number of samples!



Principle of Quantum Computing

Circuit model

CERN openlab $|0\rangle = \begin{pmatrix} 1\\ 0 \end{pmatrix}, \ |1\rangle = \begin{pmatrix} 0\\ 1 \end{pmatrix}$

Based on quantum mechanics \rightarrow Encoding in Quantum Bits (= Qubit) Superposition of states ; 2ⁿ possible states $|0\rangle$, $|1\rangle$, ..., $|2^n-1\rangle$ for n qubit quantum system

$$|\psi\rangle = \sum_{k=0}^{2^{n}-1} \alpha_{k} |k\rangle, \sum_{k=0}^{2^{n}-1} |\alpha_{k}|^{2} = 1$$

Able to manipulate computations in parallel

Computation described in unitary quantum logic gates \rightarrow **Reversible Exponentially reduced** in space and time

> **Example** : Pauli-Y-Rotation

$$R_{y}(\theta) = \begin{pmatrix} \cos\frac{\theta}{2} & -\sin\frac{\theta}{2} \\ \frac{\theta}{\sin\frac{\theta}{2}} & \cos\frac{\theta}{2} \end{pmatrix}$$

 $R_{y}(\theta)|0\rangle = \cos(\theta/2)|0\rangle + \sin(\theta/2)|1\rangle$





Quantum GAN



Based on IBM QGAN implementation

Hybrid model : Quantum Generator + Classical Discriminator

Efficient in **loading and learning a probability** over discrete values

Increase resolution by adding qubit







60

A simplified 3DGAN

25x25x25 too big for initial test: reduce to 1D distribution Focus on longitudinal energy profile



Binned into $2^n = N$ pixels \rightarrow Map to 8 values expressed by 3-qubit generator Probability of getting state $|k\rangle =$ (Relative) Energy at pixel k



61

Reproducing longitudinal energy shape

- Generator: 3 qubits, 3 layers
- **Classical Discriminator:** 512 nodes + Leaky ReLU \rightarrow 216 nodes + Leaky ReLU \rightarrow single-node + sigmoid
- AMSGRAD optimizer for both generator and discriminator





Preparation of the initial state





Understanding performance

Kolmogrov-Smirnov Statistics : Measures equality between P(x) & Q(x)

 \rightarrow With 95% confidence level & 1,000 samplings, null hypothesis is accepted if $D_{KS} \leq 0.0547$

Initialization	depth	D_{KL}	$D_{KS} \; (\times 10^{-2})$	$\operatorname{Accept}/\operatorname{Reject}$
normal	1	0.031 ± 0.042	0.033 ± 0.024	Accept
	2	$(1.9 \pm 1.7) \cdot 10^{-3}$	0.024 ± 0.002	Accept
	3	$(5.0 \pm 3.0) \cdot 10^{-4}$	0.038 ± 0.008	Accept
uniform	1	0.044 ± 0.003	0.048 ± 0.007	Accept
	2	0.12 ± 0.04	0.072 ± 0.018	Reject
	3	0.070 ± 0.039	0.054 ± 0.017	Accept
random	1	0.051 ± 0.096	0.042 ± 0.015	Accept
	2	0.096 ± 0.162	0.042 ± 0.026	Accept
	3	0.072 ± 0.078	0.059 ± 0.017	Reject



First attempts at 2D (6 qubits)

Using more powerful discriminator (4 hidden layers) Uniform initialization + depth 1 After 3,000 epochs, able to reproduce similar shape



Conclusions... again

Generative Models are powerful techniques

Fast simulation applications are reaching a high level of accuracy

Performance validation remains a key issue

- Need to start designing common procedures
 - **Results validations**

Integration

Computing resources availability is a driving factor of the problem size we can solve

65

Benefit from dedicated hardware and accelerators (GPUs, FPGAs, TPUs)

DL development is accelerated by a diversified community

Continuous R&D and collaboration with external partner is essential

Thanks!

Questions?





https://openlab.cern

CERN OPENLAB

Driving Innovation since 2001



