Deep Learning for the future of the Large Hadron Collider







SOSC 2018, Perugia

Maurizio Pierini









ML and HEP future challenges





Modern Machine Learning might be the way out

(*)With nowadays software development











Modern Machine Learning might be the way out

(*)With nowadays software development









Too many data, too large data -> need to filter online • Filters based on theoretical bias: we might be loosing

- good events

 - Offline: global, software based, on CPU, @CERN TO



The LHC Big Data Problem

▶ L1 trigger: local, hardware based, on FPGA, @experiment site ▶ HLT: local/global, software based, on CPU, @experiment site Analysis: user-specific applications running on the grid





• We are not seeing new physics: just re-doing what we do today with x10 more data WILL NOT be enough.

• The solution to the HL-LHC problem: modern Machine Learning to be <u>faster</u> and <u>better</u> in what we do today, freeing resources for new ideas

• This ML deployment need to happen in between collisions and data analysis (trigger, reconstruction, ...), where freeing resources will make a difference



The LHC Big Data Problem











Haster Particle Reconstruction Uith Computer Vision erc





Faster/better Jet Reco

- Jets are cone-like showers of quarks and gluons that produce tens of particles, all close to each other
- Jets can be reconstructed from calorimetric deposits, tracks, or full particles
- Depending on the quality of the ingredients, reconstruction can be more or less precise
 - coarse at trigger, when using fast calorimeter reconstructions
 - accurate offline, using all information















• At L1/HLT, raw data	10
objects are used	8

• map of energy deposit on calorimeter

• tracks from local/less accurate tracking

• Low-energy jets are promoted to large energy and not discarded

• So, triggers accept more than what should, or (at fixed budget) one is forced to apply tighter selection

Faster/better Jet Reco





Faster/better Jet Reco

- One could use ML to "guess" the offline jet energy from the online one + additional
- With a simple CNN2D network, substantial improvement observed
- Consequences downstream: better select which events to write /sec
 - We can do the same physics with less resources
 - We can do more physics with same resources







Combinatorics Reduction



• The more tracks we have, the more problematic it becomes (nonlinearity due to combinatoric when connecting dots

• Tracking is the most expensive workflow we have in RECO









• We tried to solve this problem using a ConvNN





13





500 GeVic H,A -> two two tjets + X, 60 fb

• A hit is a window of sensors (16x16 here) with its deposited charge. This can be seen as a sparse digital image.

of hits is a good or bad match







• Given two images, one can train a network to decide if a pair











- inputs:

 - position of the hits in the process



15

Effi	ci	ency	(tpr)	<pre>@ fake rejection</pre>
tpr	6	rej	50% :	0.998996700259
tpr	6	rej	75% :	0.990524391331
tpr	6	rej	908 :	0.922210826719
tpr	6	rej	998 :	0.338669401587







Calorimetry & Computer Vision

- (next generation) digital calorimeters: 3D arrays of sensors with more regular geometry
- Ideal configuration to apply Convolutional Neural Network
 - speed up reconstruction at similar performances
 - and possibly improve performances



See contribution to NIPS workshop









Proof of Principle: Particle ID

- ROC curve for e vs. π^{\pm} classifier 1.0 signal efficiency • We tried particle ID on a sample of simulated events 0.9 0.8• one particle/event (e, γ , π^0 , π) 0.7 Ð Different event representations
 DNN (cells) 0.6DNN (features) BDT 0.5• high-level features related to event 0.3 0.0 0.10.2 0.4 shape (moments of X,Y, and Z π^{\pm} background efficiency projections, etc) ROC curve for γ vs. π^0 classifier signal efficiency 1.0 • raw data (energy recorded in each 0.8*ce11*) 0.6 Pre-filtered pion events to select the
 0.4 nasty ones and make the problem harder DNN (cells) 0.2 DNN (features) BDT 0.0 0.00.2 0.60.8 1.00.4

See contribution to NIPS workshop

17





 π^0 background efficiency



Proof of Principle: Energy Regression

- Correctly reconstruct energy, with physics meaningful performances
 - ECAL performances better than HCAL (as expected)
 - π^0 resolution ~ $\sqrt{2}$ Y resolution (as expected)
- FAST: used only RAW data as <u>inputs -> no pre-processing</u>
 - ▶ Processing time reduced by 10³ wrt traditional approaches
 - Potentially usable both online and offline

See contribution to NIPS workshop





Research





HEP & Language processing networks





Particle Flow & language processing

- CMS uses particle flow for event reconstruction:
 - At some point in the central processing, collision images are turned into a list of particles.
 - From these particles, complex objects (e.g., jets) are formed
- In this framework, Computing vision approaches are not necessarily ideal
- One can instead use language-processing approaches (e.g., recurrent neural networks
 - particles are words in a sentence
 - QCD is the grammar











Recurrent Neural Networks

- A network architecture suitable to process an ordered sequence of inputs
 - words in text processing
 - a time series
 - particles in a list
- Could be used for a single jet or the full event
- Next step: graph networks (active research direction)













A Topology Classifier

<u>A typical example: leptonic triggers</u>

- at the LHC, producing an isolated electron or muon is very rare. Typical smoking gun that something interesting happened (Z,W,top,H production) -> TAKE THEM!
- Triggers like those are very central to ATLAS/CMS physics
- The sample selected is enriched in interesting events, but still contaminated by non-interesting ones
- \odot Can we clean this up w/o biasing the physics? yes, with ML













A Topology Classifier





Abstract Image Classifier

Based on an abstract representation of the reconstructed particles as an image to feed to a convolutional NN.

High-level Feature Classifier

Use high-level features as inputs to a fully connected NN.





(a) Photons



(b) Charged Particles



(c) Neutral Hadrons











Can select 99% of the top events and reduce the fraction of written events by a factor ~ 7

Selection performances







Selection pert



What is the network learning? tt events are more crowded that W events leptons in W and tt events are isolated from other

25

- particles

Imances









Networks erc





Generative Adversarial Training

• Two networks trained in competition

• Generator: creates images starting from random **noise** (and optionally some other information to transform)

• Discriminator: tries to distinguish true from generator-created images

- The loss function to minimise is written as Loss(Gen)-Loss(Disc) • Goes up is discriminator improves

• Goes down if the generator improves

simply training itself to fool the generator



Noise

• The generator learns to make images like a given set it never sees,





Generative Adversarial Training

• Two networks trained in competition

- Generator: creates images starting from random noise (and optionally some other information to transform)
- Discriminator: tries to distinguish true from generator-created images
- The loss function to minimise is written as Loss(Gen)-Loss(Disc) • Goes up is discriminator improves
- - Goes down if the generator improves
 - simply training itself to fool the generator



• The generator learns to make images like a given set it never sees,



28





Generative Adversarial Training

• Two networks trained in competition

- Generator: creates images starting from random noise (and optionally some other information to transform)
- Discriminator: tries to distinguish true from generator-created images
- - Goes up is discriminator improves
 - Goes down if the generator improves
 - simply training itself to fool the generator



Noise

• The loss function to minimise is written as Loss(Gen)-Loss(Disc)

• The generator learns to make images like a given set it never sees,











PROGRESSIVE GROWING OF GANS FOR IMPROVED QUALITY, STABILITY, AND VARIATION

Adversarial training in azione

Submitted to ICLR 2018









• Start from random noise

• Works very well with images

• Applied to electron showers in digital calorimeters as a replacement of GEANT



See contribution to NIPS workshop

<u>Image generation</u>

Shower longitudinal section



31











• Start from random noise

• Works very well with images

• Applied to electron showers in digital calorimeters as a replacement of GEANT



Figure 6: The distributions of image mass m(I), transverse momentum $p_{\rm T}(I)$, and *n*-subjections $\tau_{21}(I)$. See the text for definitions.

Generating full jets













Simulation is half of the problem

33

- Reconstruction involves more than one detector (e.g., tracker + calorimeter) and produces (at least in CMS) a list of particles
- GANs were proved to be useful to emulate the SIM+RECO step in one goa1
 - jets out of full particle reconstruction emulated in a GAN
 - trained on actual SIM+RECO synthetic data by CMS

GENERATION GEANT SIMULATION Tracking+clust ering+... +ParticleFlow RECONSTRU **CTION** Selection **ANALYSIS-SPECIFIC** erc DATASET







Simulation is half of the problem

- Reconstruction involves more than one detector (e.g., tracker + calorimeter) and produces (at least in CMS) a list of particles
- GANs were proved to be useful to emulate the SIM+RECO step in one goal
 - jets out of full particle reconstruction emulated in a GAN
 - It trained on actual SIM+RECO synthetic data by CMS





Fig. 6 Distribution of high level variables used for quark/gluon discrimination (first two rows) and merged jets tagging (last row). Blue histograms are obtained from the input data, while red ones are obtained using the generative model.









An application-specific approach

- CPU is only part of the problem. Storage is the other big one (biggest at the moment) Data Z(vv)+γ/jets W(lv)+γ/jets Typically, multiple copies of the same event are Sinale t stored in different sites worldwide QCD Multiiets Diboson Z(II)+γ/iets ADD $M_n = 2 \text{ TeV}, \delta = 3$ Unparticles d =1.7, $\Lambda_U = 2 \text{ TeV}$ few KB of data: high-level features computed from 10² the 1 MB of raw information (e.g. CMS nanoAOD) 10-1 this dataset? Data - MC MC • PRO: big save both on disk & CPU • CONS: the training setup is analysis specific. 500 600 700 800 300
- A typical LHC event takes 1 MB of disk / tape. • After all processing, a typical analysis uses a • Can we use a generative model to go straight to
- - Several Generative Models will be needed to cover all use cases



35









• We consider a classic GAN setup as baseline (also tried wGAN)

- stabilize the training
- Implemented in keras+TF



• Train a generator an discriminator in an adversarial fashion

• Add regression of meet to the generator cost function, in order to

• Running on server mounting GTX1080 cards + CSCS Piz Daint (project cn01)





- Loss function: cross entropy + c · mse(mee) • fixed c = 0.01
- Dataset size: 2M events
- 100K epochs / 512 events per batch
- Multiple trainings of the same model with randomized starting point, to minimise dependency on initial conditions
- (as normal with GANs) training quite unstable and wGAN didn't really help
 - Cannot use the loss function itself as a guide to the best model
 - instead, work on defining a generator quality assessment based on statistics tests

Training







• Define global- and feature-specific quality tests

- Statistics Score (SS) = $\frac{1}{19^2} \sum_{i=j}^{19} \frac{\Delta \sigma_{ij}}{\sigma_{ij}^{\text{real}}} + \frac{1}{19} \sum_{i=j}^{19} \frac{\Delta \mu_i}{\mu_i^{\text{real}}}$.
 - Roughly the normalized sum of the difference between covariances and averages.
 - $-\Delta\sigma_{ij}$ = difference between the covariance matrix elements in the real and generated data
 - $-\sigma_{ii}^{\text{real}} = \text{covariance matrix element for real data}$
 - $-\Delta \mu_i$ = difference between the average value of the *i*th variable in the real and generated data $-\mu_i^{\text{real}} = \text{average value of the } i\text{th variable in the real data}$
- MLLKS = Kolmogorov-Smirnov test statistic for the MLL
- MetPhiKS = Kolmogorov-Smirnov test statistic for the MET- ϕ
- LepIsoKS = Kolmogorov-Smirnov test statistic for the leading lepton isolation

definitions (KL divergence, earth-mover distance, etc)

• Work in progress: investigating usage of standard pdf-distance erc













Results

k	MLLRank 59 107 79 23 88 306	ScoreRank 43 392 228 67 181 154	StatsScore 27.826635 47.069988 40.682663 30.737935 38.685347 37.152292	MLLKS 0.085580 0.096820 0.090900 0.075540 0.092740 0.122060	MetPhiKS 0.032840 0.024580 0.033440 0.022300 0.041520 0.013740	LepIsoKS 0.236480 0.188500 0.232180 0.305880 0.214840 0.263080	SortKey 560.000000 699.000000 760.000000 800.000000 802.000000 818.000000
	0 200 lep2_e	iatch512_bgbd_mllANDwid	th_NonTC_newdata_mllfix: epoch	81500	-100 0 100 Hep1_py		0 2 lep2_e
	0 100 met	-100) 0 100 lep2_px		-2.5 0.0 2.5 dphi	0 50 nvtxs	0 100 met
	lep2_iso	-5				0 2 lepl_iso	0 lep2_ist
	0 100 jet_pt5	o	100 jet_pt1	0 100 jet_pt2	0 100 jet_pt3	0 100 jet_pt4	0 jet_pt5
	0.0 0.5 1.0		75 100 mll	l lep1_mass	0 1 lep2_mass	0 5 njets	0.0 0.5













- In view of large statistics needs, one can use generative models as statistics augmentation tools
- For instance, could generate expert-feature quantities used in an analysis (muon fourmomenta, jet momenta, etc.)
- Like sampling from histogram with two main advantages
 - no need to bin
 - generalizes to multidimensional problems

with K. Dutta, N. Amin, B. Hashemi and D. Olivito (in preparation) 40

Fast Simulation







• Dimuon events at LHC

momenta, isolation, jet pTs, etc)

• Can use the generator network as a fastsim tool



- Typical analysis would use a few handful of quantities (muon)
- Can learn the N-dim distribution of these quantities with GAN setup











• Dimuon events at LHC

momenta, isolation, jet pTs, etc)

• Can use the generator network as a fastsim tool



- Typical analysis would use a few handful of quantities (muon)
- Can learn the N-dim distribution of these quantities with GAN setup







Analysis-specific unfolding









Fast Decision Taking







The LHC Big Data Problem



	P	×.	
¢.		Ν	
1			
ľ			





- The L1 trigger is a complicated environment
 - decision to be taken in ~10 µsec
 - only access to local portions of the detector
 - processing on Xilinx FPGA, with *limited memory resources*
- Some ML already running @L1
 - CMS has BDT-based regressions coded as look-up tables

• Working to facilitate DL solutions **@L1 with dedicated library**

Brind LJL to







- You have a jet at LHC: spray of hadrons coming from a "shower" initiated by a fundamental particle of some kind (quark, gluon, W/Z/H bosons, top quark)
- You have a set of jet features whose distribution depends on the nature of the initial particle
- You can train a network to start from the values of these quantities and guess the nature of your jet
- To do this you need a sample for which you know the answer

Example: jet tagging





47















- (je
- cor



Make the model cheaper

• **Pruning:** remove really contribute to performances

> possible (regularization)

M<u>ake the model cheap</u>er

• Pruning: remove
parameters that don't
really contribute to
performances

• force parameters
to be as small as
possible
(regularization)

 $L_{\lambda}(\vec{w}) = L(\vec{w}) + \lambda ||\vec{w}_1||$

Remove the small

→ 70% reduction of weights

- Quantization: reduce the number of bits used to represent numbers (i.e., reduce used memory)
 - models are usually trained at 64 or 32 bits
 - this is not necessari 1. needed in real feg_relu ftg_relu

Absolute Relative Weights

In our case, ^{*}
We could to 16 bits w/o loosi precision

peed vs Memory

Fully serial

reuse = 4use 1 multiplier 4 times

mult

reuse = 2use 2 multipliers 2 times each

Reuse factor: how much to parallelize operations in a hidden layer

erc

Performances

will reduce the DSP usage Council 54

LHC

• Neutrino experiments

HL LHC

European Research

Data Quality Monitoring

When taking data, >1 person watches for anomalies in the detector 24/7

- At this stage no global processing of ^b the event
- Instead, local information from detector components available (e.g., detector occupancy in a certain time window)

58

• Given the nature of these data, ConvNN are a natural analysis tool. Two approaches pursued

• Classify good vs bad data. Works if failure mode is known

• Use autoencoders to assess data "typicality". Generalises to unknown failure modes

A. Pol et al., to appear soon

<u>luo approaches</u>

Fully connected

• Given the nature of these data, ConvNN are a natural analysis tool. Two approaches pursued

• Classify good vs bad data. Works if failure mode is known

• Use autoencoders to assess data "typicality". Generalises to unknown failure modes

A. Pol et al., to appear soon

<u>luo approaches</u>

• Autoencoder-based 1-class approach generalises to later stages of quality assessment

- after reconstruction of the events, event reconstruction allows a global assessment (w.g., looking at electrons, muons, etc rather than hits in the detector)
- A global autoencoder can spot all these features
- Monitoring individual contributions to loss function (e.g., MSE) one can track the problem back to a specific physics object/detector component

F. Široký et al., to appear sooner or later

Data Quality Certification

Xilinx Vivado 2017.2

Results are slightly different in other versions of Vivado

Clock frequency: 200 MHz Latency results can vary (~10%) with different clock choices

FPGA: Xilinx Kintex Ultrascale (XCKU115-FLVB2104) Results are slightly different in other FPGAs e.g. Virtex-7 FPGAs are slightly differently optimized

- e.g. 2016.4 optimization is less performant for Xilinx ultrascale FPGAs

Neural network can model non linear functions

• the more complex is the network, the more functions it can approximate

• Neural network are faster to evaluate (inference) than typical reco algorithm.

• This is the speed up we need

• Neural Networks (unlike other kind of ML algorithms) are very good with raw (non-preprocessed) data (the recorded hits in the event)

(pT, η, φ, E)_{OFFLINE} = $f(pT, η, φ, E)_{ONLINE}$

63

• could use them directly on the detector inputs

(**pT**, **n**, **\phi**, **E**)_{OFFLINE} = g(Event hits)

One would have to learn f and g to evaluate them at trigger. Online processing is replaced by offline training

• Approach works in principle

Might not work in absolute

- are not model-agnostic
- be done next

Beyond the toy-model

• With 99% signal efficiency, bias on kinematic variables within the uncertainty of a trigger-efficiency measurement

European

TOPCLASS: do we kill New Physics?

TOPCLASS: do we kill New Physics?

