

ML at Work:



Electron Energy regression with CNN

Malte Algren^{*}, Aske Rosted, and **Troels C. Petersen** Niels Bohr Institute, Copenhagen (*now at Univ. of Geneva)

Outline

Outline of talk:

- Motivation
- Context
- Training a CNN for energy reconstruction:
 - The data
 - The selections
 - The input variables
 - The network architecture
 - Feature wIse Linear Modulation (FiLM)
- Results in MC
- Results in data (v1)
- Training in data and "simultaneous training"
 - Results in data (v2)
- Outlook



Motivation

Points of motivation:

- Improve $H \rightarrow ZZ^*$ and $H \rightarrow \gamma\gamma$ analyses
- Optimise searches for:
 - HH $\rightarrow \gamma \gamma bb$
 - $H \rightarrow Z\gamma$
 - $H \rightarrow \gamma^* \gamma$
- Improve resilience to pile-up
- Improve $Z \rightarrow$ ee reconstruction
- Utilise excellent data for testing:
 - CNN and GNN models
 - data+MC simultaneous training
 - $e+\gamma$ simultaneous training
- Improve non-Higgs searches

Goals of lecture:

- Give example of regression with CNN.
- Illustrate concept of attention and FiLM technique.
- Illustrate "target mismatch" and combined training.



Motivation

Points of motivation: (You don't have to care - just know the list is long!)

- Improve $H \rightarrow ZZ^*$ and $H \rightarrow \gamma\gamma$ analyses
- Optimise searches for:
 - HH $\rightarrow \gamma \gamma bb$
 - $-H \rightarrow Z\gamma$
- $H \rightarrow \gamma^* \gamma$
- Improve resilience to pile-up
 Improve Z → ee reconstruction
- Utilise excellent data for testing:
 - CNN and GNN models
 - data+MC simultaneous training
 - $e+\gamma$ simultaneous training
- Improve non-Higgs searches

Goals of lecture:

- Give example of regression with CNN.
- Illustrate concept of attention and FiLM technique.
- Illustrate "target mismatch" and combined training.



ATLAS Detector



ATLAS Detector







The scalars can be seen in table on the right.

The variables are both scalar and cell based.

Туре	Name	Description				
	Eacc	Energy deposit in layer 1-3 of ECAL.				
	η_{index}	η cell index of cluster of layer 2.				
	$f0_{cluster}$	Ratio of energy between layer 0 and E_{acc} in $ \eta < 1.8$ (end of layer 0).				
Energy	R12	Ratio of energy between layer 1 and 2 in the ECAL.				
	p_t^{track}	p_T estimated from tracking for the particle (o e).				
	E_{TG3}	Ratio between the energy in the crack scintillate and E_{acc} within $1.4 < \eta < 1.6$.				
	$E_{tile-gap}$	Sum of the energy deposited in the tile-gap.				
	η	Pseudorapidity of the particle.				
	$\Delta \phi_2^{rescaled}$	Difference between ϕ , as extrapolated by track- ing, use for ECAL momentum estimation and ϕ of the ECAL cluster.				
	$\eta_{ m ModCalo}$	Relative η position w.r.t. the cell edge of layer 2 in the ECAL*.				
Geometric	$\Delta \eta_2$	Difference between η , as extrapolated by tracking, use for ECAL momentum estimation and η of the ECAL cluster (only <i>e</i>).				
	poscs ₂	Relative position of η within cell in layer 2 in ECAL. $2(\eta_{cluster} - \eta_{maxEcell})/0.025 - 1$, $\eta_{cluster}$ is η of the barycenter of the cluster and $\eta_{maxEcell}$ is η of the most energetic cell of the cluster.				
	$\Delta \phi_{TH3}$	Relative position in ϕ in a cell. mod $(2\pi + \phi, \pi/32) - \pi/32$.				
	$\langle \mu \rangle$	Average proton-proton interaction per bunch crossing.				
Misc.	n _{tracks}	# of tracks assigned (only <i>e</i>).				
	$n_{vertexReco}$	Number of reconstructed vertices.				

The variables are both scalar and cell based. The scalars can be seen in table on the right.

We consider the cell energies in the LAr calorimeter as pixels in four images. The cells contain two (used) types of information: • Energy (primary variable)

• Time of cell energy



The variables are both scalar and cell based. The scalars can be seen in table on the right.

We consider the cell energies in the LAr calorimeter as pixels in four images. The cells contain two (used) types of information:

- Energy (primary variable)
- Time of cell energy

In order to have the **same resolution** in each layer, we **upsample** the layers to the lowest common resolution (work by Lucas Erhke).





The variables are both scalar and cell based. The scalars can be seen in table on the right.

We consider the cell energies in the LAr calorimeter as pixels in four images. The cells contain two (used) types of information:

- Energy (primary variable)
- Time of cell energy

Finally, we consider the (up to) 10 nearest tracks in a "TrackNet" input:

Туре	Name	Description				
Energy	p _{t,track} /q _{track}	Transverse momentum of track divided by its charge q				
	d_0/σ_{d0}	d0 is the signed transverse distance between the point of closest approach and the <i>z</i> -axis where σ_{d0} is its uncer- tainty				
	ΔR	$\Delta R = \sqrt{(\phi_0 - \phi)^2 + (\eta_0 - \eta)^2}$				
Geometric	vertex _{track}	Reconstructed vertex of the track				
	<i>z</i> ₀	Longitudinal distance between the point of closest approach and the z-axis.				
	η_{track}	Reconstructed $ \eta $ of tracks.				
	ϕ_{track}	Reconstructed ϕ of tracks.				
	n_{pixel}	Number of hits in the pixel detector				
Misc.	n_{SCT}	Number of hits in the SCT				
	n _{TRT}	Number of hits in the TRT				



The network architecture

There are many ways to combine the input variables, and we have considered the following architectures, where the dashed lines are the considerations.



First, let us consider each part...

The network architecture

The CNN is the main estimator of the overall energy, and is relatively standard. The innovation lies in Feature wIse Linear Modulation (FiLM).

TrackNet is a pile-up-corrective input consisting of a combination of the reconstructed tracks close to the candidate. It can be input to both top layer or FiLM weights.

Finally, the scalars (and possibly tracks) are used in the FiLM.







Feature wIse Linear Modulation



The network architecture

Testing all the different combinations yields the optimal architecture.

We evaluate the performance in the same way as previously done, namely the effective InterQuantile Range (eIQR) of the Relative Error (RE).

еI	$QR = \frac{P_{75}(RE) - P_{25}}{1.240}(RE)$,	R	$E = \frac{E_{cal}}{E_{cal}}$	$\frac{lib}{d}$,	-	Hyperparameter	Parameter
	1.349			Ltru	ıth		Units	(128, 64, 32, 16)
							Normalization	Batch
F							Kernel size & filters	5
		reIÇ	2R75	reIQR95			Connected to	[Top]
3	Basic	-(0.121	-0.025	E.S.		ScalarNet	
1	FiLM: scalar	(0.229	0.257	10		Units	(256)
	FiLM: scalar - top: scalar	(0.220	0.252			Normalization	Batch
	Fil M: scalar - top: scalar track		0.000	0.251	2	1	Connected to	[FiLM]
ſ	TILM. Scalar - top. Scalar track		5.223	0.251	Deel		FiLM gen.	
	FiLM: scalar - top: track		0.226	0.264	Dest	A	Contrecture	(512, 1024)
	FiLM: scalar track	(0.228	0.265	3	5.00	Normalization	Batch
	FiLM: scalar track - top: scalar track	(0 .2 10	0.262	2		CNNnet	
11	FiLM: track - top: scalar	-(0.042	-0.067	3		Down-sampling	MaxPool
	Fil M: track top: track		- 1 40	0.140	57		Globalpooling	MaxPool
	FILM. HACK - top. HACK	(5.140	0.149			Number of blocks	3
	top: scalar	-(0.154	-0.131			Depth of blocks	4
	top: scalar track	(0.213	0.233			Тор	
	top: track	(0.136	0.164			Units	(512, 512, 1)
					1		Output activation	ReLU



The results in 2D - MC

The E_T distribution for truth (x-axis) and reconstruction (y-axis) can be compared for the current ATLAS and the DeepCalo algorithms.

As the figure shows, both algorithms do well, and improve with energy.

As the statistics is largest around 40 GeV, this is where the comparison is most detailed, and here DeepCalo visibly has a significantly reduced lower edge. Thus, the DeepCalo more rarely undershoots the energy.



The results in 1D - MC

Integrating the previous plot into 1D considering the RE distribution, we see a general sharpening. The improvement in relative eIQR (reIQR) is about 22%.



Naively, we would of course love to see a similar number in data!

Result in Zee - MC

On the Zee peak, we evaluate the improvement by fitting with a BW⊗CB fit, considering the CB width (sigmaCB) as the performance parameter. We get:



Result in Zee - MC

On the Zee peak, we evaluate the improvement by fitting with a BW⊗CB fit, considering the CB width (sigmaCB) as the performance parameter. We get:



Results on Zee - data (v1)

The result we get is a much more modest improvement:

$$\langle 1 - rac{\sigma^{DeepCalo}_{CB}}{\sigma^{ATLAS}_{CB}}
angle = 1 - rac{2.058 \pm 0.010}{2.271 \pm 0.019} = 9.4 \pm 0.9\%.$$

Though perhaps a little disappointing, this is not surprising, as we can not expect the MC to mimic data perfectly in the very large space considered. Also, models trained on Zee do not generalise well to all energies (EG, 6.8%).





Training in data

Using Zee events with invariant masses 86-97 GeV, one can get "approximate labels" in data, by assuming the true Z mass: $M^2 = 2p_{T,1}p_{T,2}(\cosh(\eta_1 - \eta_2) - \cos(\phi_1 - \phi_2)), p_T$

Using such labels, we train in data and get...

$$M^{2} = 2p_{T,1}p_{T,2}(\cosh(\eta_{1} - \eta_{2}) - \cos(\phi_{1} - \phi_{2})), \quad p_{T} = E_{T} \updownarrow$$
$$E_{label,data} = \frac{M^{2}}{2E_{T,2}(\cosh(\eta_{1} - \eta_{2}) - \cos(\phi_{1} - \phi_{2}))'}$$
with $E_{T,2} = E$ calib^(BDT) and $M^{2} = 91.19^{2}$

Training in data

Using Zee events with invariant masses 86-97 GeV, one can get "approximate labels" in data, by assuming the true Z mass: $M^2 = 2p_{T,1}p_{T,2}(\cosh(\eta_1 - \eta_2) - \cos(\phi_1 - \phi_2)), p_1$

Using such labels, we train in data and get...

Training in data and MC

Once we have labels in data, there is nothing keeping us from combining the loss functions of MC and data (they even have the same form), and thus training **simultaneously** in data and MC:

$$\mathcal{L}(y, \hat{y}) = \mathcal{L}(y_{(\text{Zee, MC})}, \hat{y}_{(\text{Zee, MC})}) + \mathcal{L}(y_{(\text{Zee, Data})}, \hat{y}_{(\text{Zee, Data})})$$

This allows the model to both use the "strength" of MC, but also learn the differences between MC and real data.

Doing this and trying out the result in MC first yields:

$$\langle reIQR_{75}^{DeepCalo}
angle = 22.1 \pm 0.3\%$$

OK, so at least it doesn't ruin the model for MC. Now let us try data...

Result in data (v2)

The result in data is rather encouraging, and **greater than the sum of the improvements** from training separately in MC (9.4%) and data (5.9%).

Outlook

While this is still "only" an improvement in the electron energy regression, and only for lower energies (Zee range), the simultaneous training allows for extending the energy range, by including the Electron Gun MC.

Furthermore, this training might be extended to include photons, as these behave much the same as electrons, and suffer the same sources of uncertainties and smearing.

For improving the H $\rightarrow \gamma \gamma$ resolution, one might use the following loss function and related training samples:

$$\begin{aligned} \mathcal{L}(y, \hat{y}) &= \mathcal{L}(y_{(\text{Zee, MC})}, \hat{y}_{(\text{Zee, MC})}) + \mathcal{L}(y_{(\text{Zee, Data})}, \hat{y}_{(\text{Zee, Data})}) + \\ \mathcal{L}(y_{(Z\mu\mu\gamma, \text{MC})}, \hat{y}_{(Z\mu\mu\gamma, \text{MC})}) + \mathcal{L}(y_{(Z\mu\mu\gamma, \text{Data})}, \hat{y}_{(Z\mu\mu\gamma, \text{Data})}) + \\ \mathcal{L}(y_{(H\gamma\gamma, \text{MC})}, \hat{y}_{(H\gamma\gamma, \text{MC})}) \end{aligned}$$

Meanwhile, we are trying to write this up somehow (but Malte is now a Ph.D. in Geneva).

Context

We have for the most part worked **only on MC** (step 1 below), comparing our CNN approach to the "**ATLAS BDT**". Here we see significant improvements.

Lacking the remaining "hard work" of corrections and calibration to match data, our performance improvements **in data** have been decent but "mediocre".

While we have lately included data in training also, the following results will almost surely further improve with the subsequent calibration.

The data

We have used "millions" of mainly Zee decays and Electron Gun. The data retained for **testing** is as follows:

1.000.000	450,000
	430.000
350.000	400.000
310.000	No data available
1.100.000	No data available
	350.000 310.000 1.100.000

The selection

We applied a general (loose) selection to the different channels in order to obtain large unbiased samples in both MC and data. In the following, we consider mainly the Zee channel.

For the MC, we furthermore required the truth energy to match the reconstructed energy (by ATLAS) to avoid mis-matches (k = 0.6).

Differential results - MC

Comparing on electron gun MC in a "known" ATLAS figure style, the improvement is isotropic in eta, and decreases slightly with energy.

