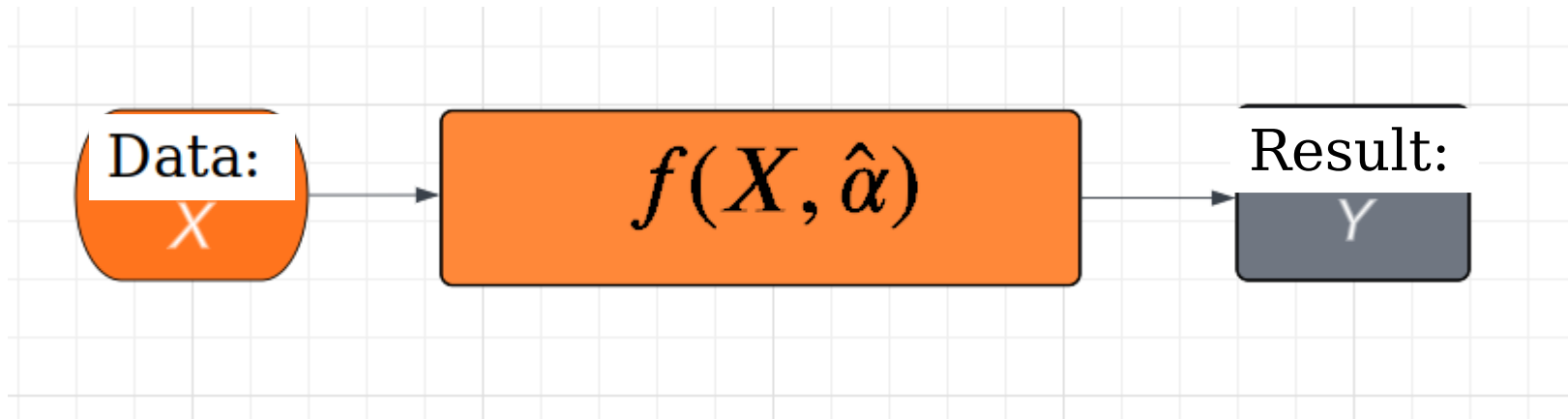


# **Machine learning in HEP**

**Artur Kalinowski**  
University of Warsaw

**X**: n dimensional input data

**Y**: k dimensional output data



Function fitting:

- **best case:  $N \leq 2$ ,  $k=1$**
- **basis functions given *explicit***
- **expansion coefficients *could* be interpretable**

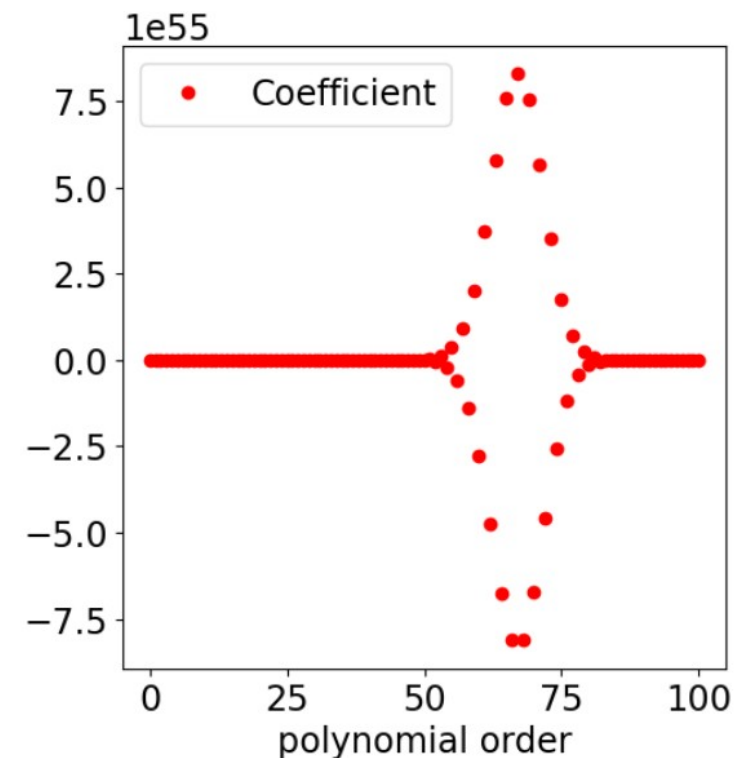
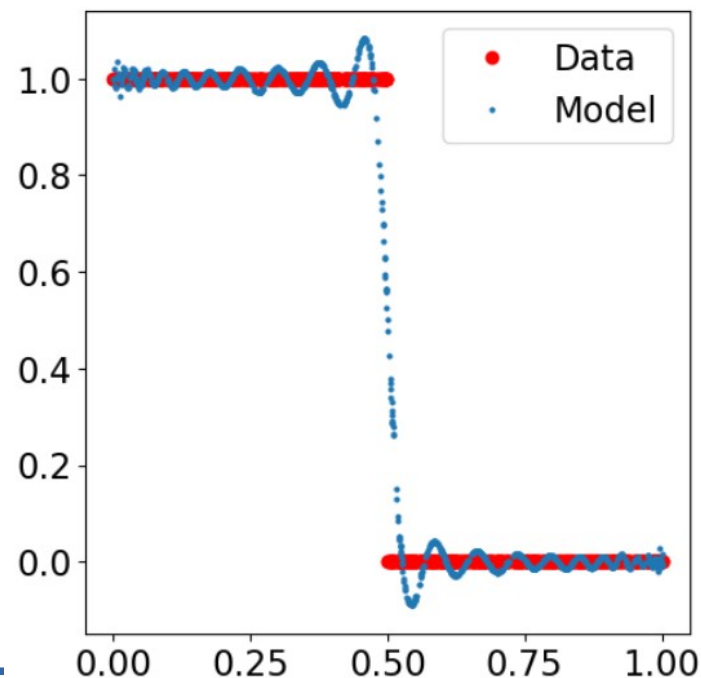
*Taylor theorem (J. Gregory, 1671):*

Every, continuous, differentiable, function  $f(x): \mathbb{R} \rightarrow \mathbb{R}$  can be approximated by a polynomial:

$$f(x, \theta) = \sum_{n=0}^{\infty} \theta_n x^n \simeq \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_n x^n$$

Coefficients  $\theta_i$  are derivatives of  $f(x)$ .

In the case of unknown function (“data”) coefficients can be found by a numerical procedure. Usually...



## Fourier theorem (1807) :

Every continuous, differentiable, and periodic function  $f(x): \mathbb{R} \rightarrow \mathbb{R}$  can be approximated in a basis of sines i cosines:

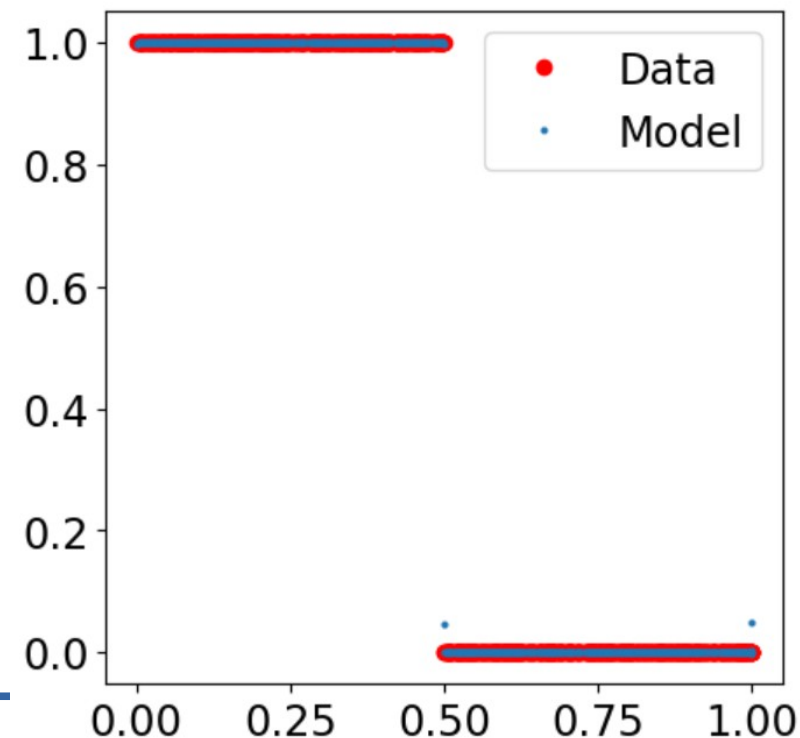
$$f(x, \theta) = \frac{\theta_{0,0}}{2} + \sum_{n=1}^{\infty} \theta_{0,n} \cos(n\omega x) + \theta_{1,n} \sin(n\omega x), \quad \omega = \frac{2\pi}{T}$$

Coefficients  $\theta_i$  can be found analytically.  
In the case of unknown function ("data") coefficients can be found by a numerical procedure.

```

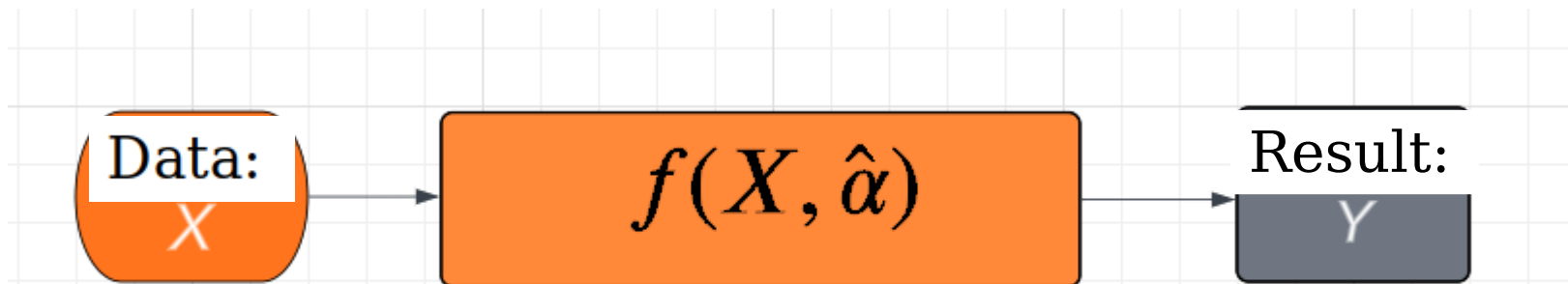
1 nMax = 25000
2 omega = 2*np.pi/1.0
3
4 Y_model = np.full_like(X,0.5)
5 for n in range(1,nMax+1, 2):
6     Y_model += 2.0/np.pi*1/n*np.sin(omega*n*X)

```



**X:** n dimensional input space

**Y:** k dimensional output space



- $n \sim 10^l$ ,  $k \sim 10^m$ ,  $l, m \sim 6$
- basis functions defined *implicite* - through data flow - the network architecture
- expansion coefficients are uninterpretable

A sigmoid function: any non polynomial function fulfilling conditions:

$$A(\theta, x) = A\left(\sum_{i=1}^n \theta_i x_i + b\right) \quad \lim_{x_i \rightarrow -\infty} A(x) \rightarrow 0 \quad \lim_{x_i \rightarrow +\infty} A(x) \rightarrow 1$$

*Universal approximator theorem (Cybenko, 1989):*

Every continuous function  $f(x) \mathbb{R}^n \rightarrow \mathbb{R}$  can be approximated in basis of sigmoidal functions:

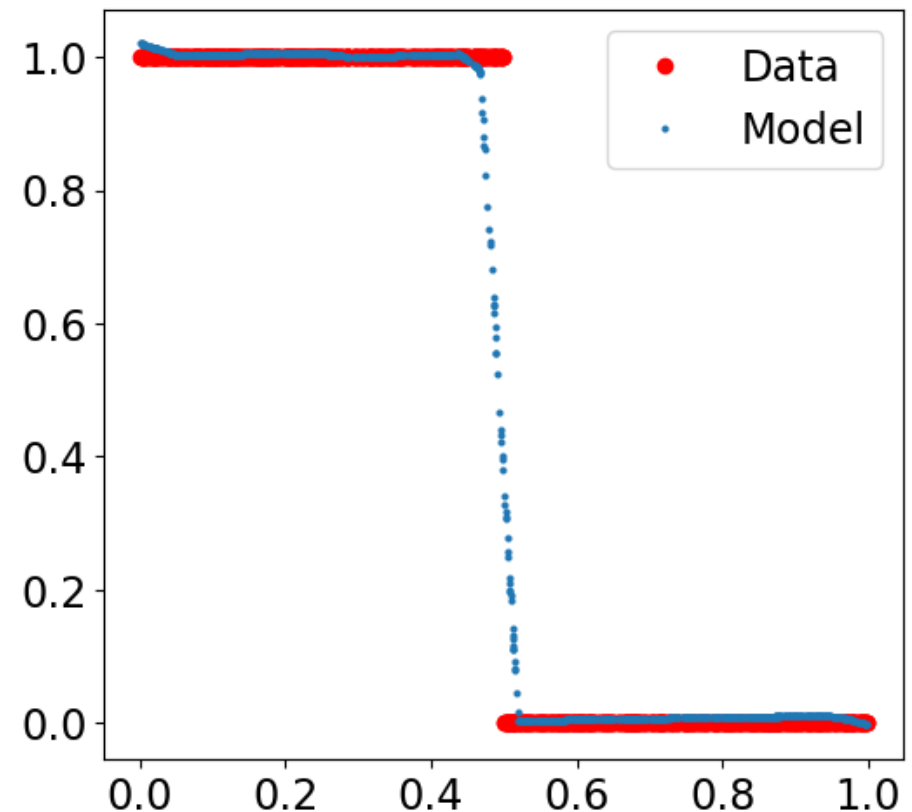
$$f(x, \theta) \simeq \sum_n w_n A(\theta_n, x)$$

Coefficients  $\theta_i, w_i$  do not have in general an analytic form, but can be found using a numerical procedure

```

1 nUnits = 32
2 inputs = tf.keras.Input(shape=(1,))
3 layer1 = tf.keras.layers.Dense(nUnits, activation='relu')(inputs)
4 layer2 = tf.keras.layers.Dense(nUnits, activation='relu')(layer1)
5
6 outputs = tf.keras.layers.Dense(1, activation='linear')(layer2)
7 model = tf.keras.Model(inputs=inputs, outputs=outputs)
8 model.compile(loss = 'mse')
9 ###
10 history = model.fit(X,Y,epochs=150)
11 Y_model = model.predict(X)

```

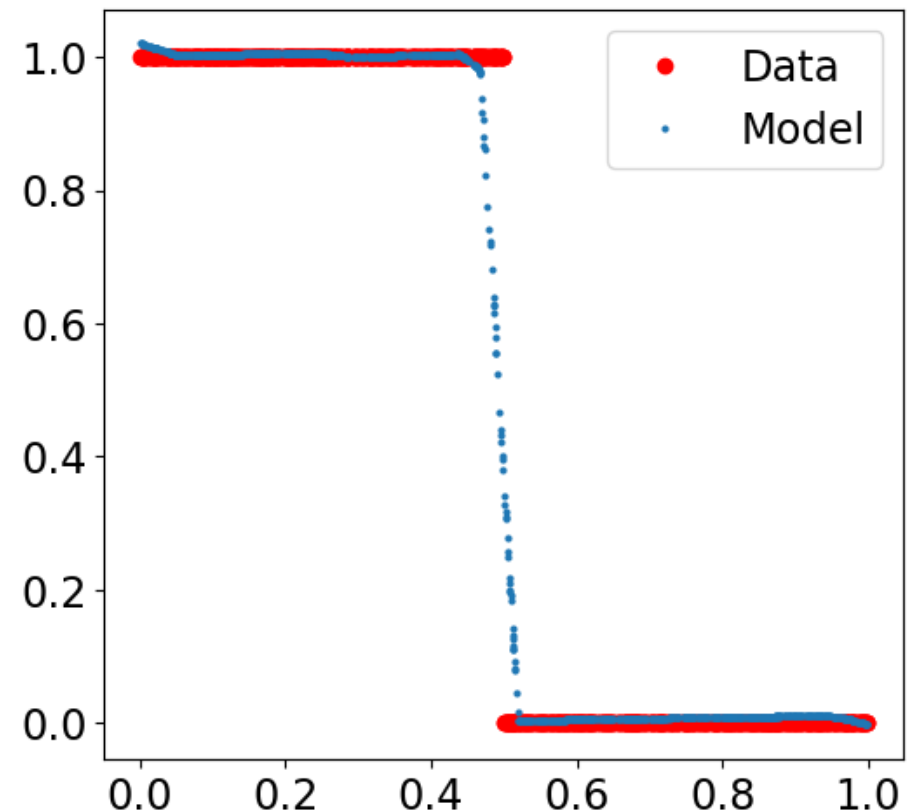


```

1 nUnits = 32
2 inputs = tf.keras.Input(shape=(1,))
3 layer1 = tf.keras.layers.Dense(nUnits, activation='relu')(inputs)
4 layer2 = tf.keras.layers.Dense(nUnits, activation='relu')(layer1)
5
6 outputs = tf.keras.layers.Dense(1, activation='linear')(layer2)
7 model = tf.keras.Model(inputs=inputs, outputs=outputs)
8 model.compile(loss = 'mse')
9 ###
10 history = model.fit(X,Y,epochs=150)
11 Y_model = model.predict(X)

```

**An issue 0:** it is quite hard to get an extremely precise predictions from ML.

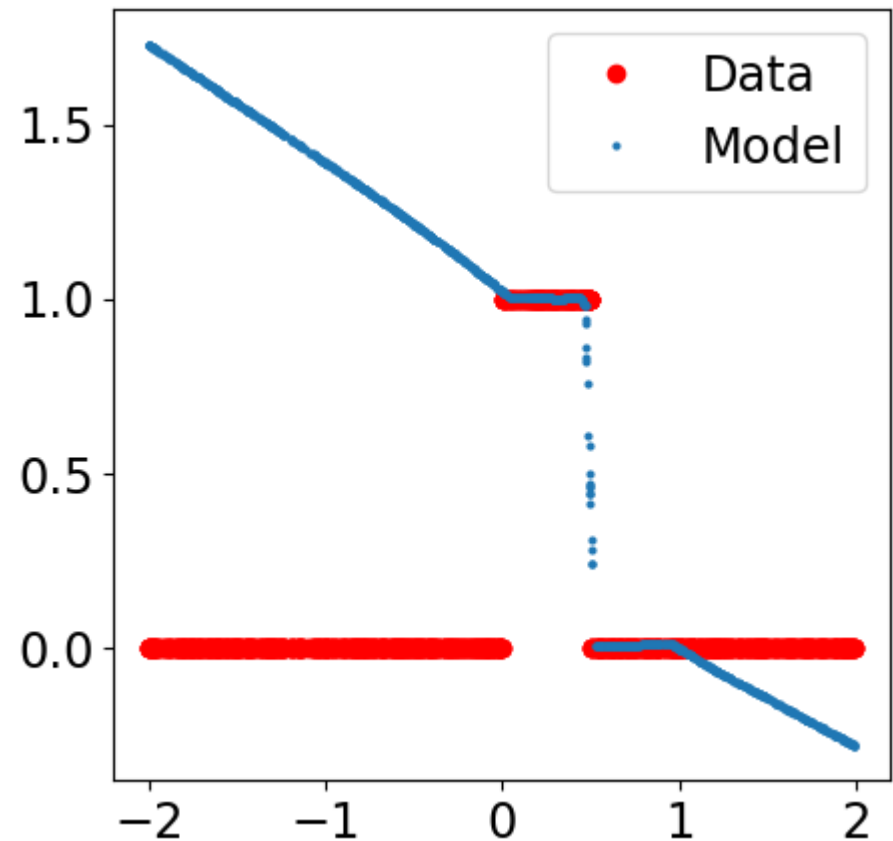




```

1 nUnits = 32
2 inputs = tf.keras.Input(shape=(1,))
3 layer1 = tf.keras.layers.Dense(nUnits, activation='relu')(inputs)
4 layer2 = tf.keras.layers.Dense(nUnits, activation='relu')(layer1)
5
6 outputs = tf.keras.layers.Dense(1, activation='linear')(layer2)
7 model = tf.keras.Model(inputs=inputs, outputs=outputs)
8 model.compile(loss = 'mse')
9 ###
10 history = model.fit(X,Y,epochs=150)
11 Y_model = model.predict(X)

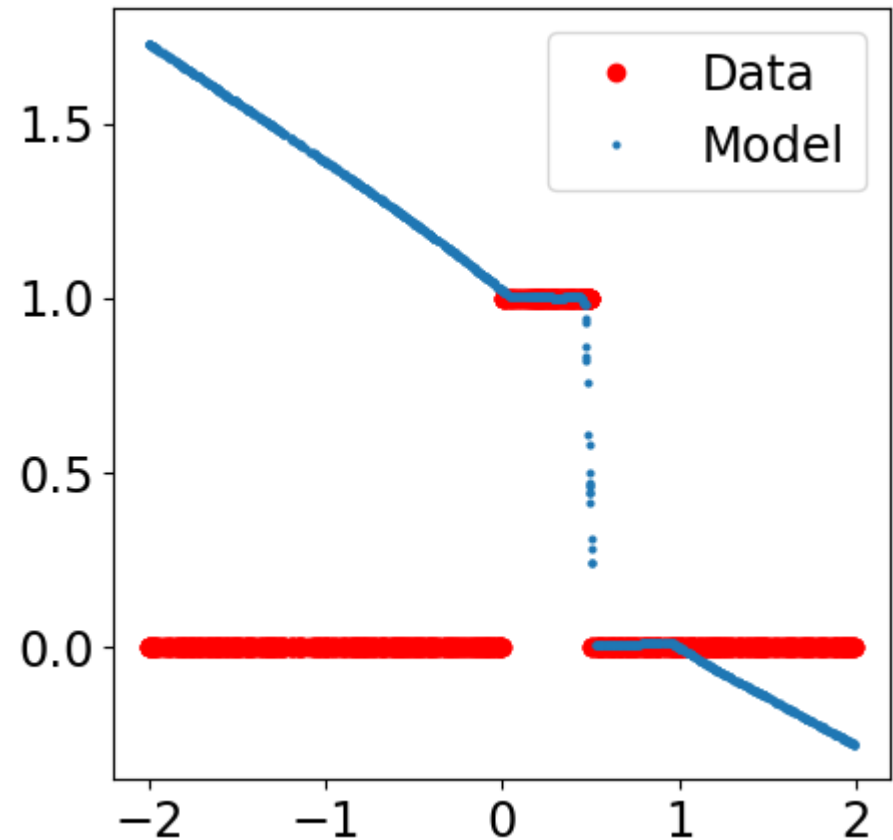
```



```

1 nUnits = 32
2 inputs = tf.keras.Input(shape=(1,))
3 layer1 = tf.keras.layers.Dense(nUnits, activation='relu')(inputs)
4 layer2 = tf.keras.layers.Dense(nUnits, activation='relu')(layer1)
5
6 outputs = tf.keras.layers.Dense(1, activation='linear')(layer2)
7 model = tf.keras.Model(inputs=inputs, outputs=outputs)
8 model.compile(loss = 'mse')
9 ###
10 history = model.fit(X,Y,epochs=150)
11 Y_model = model.predict(X)

```



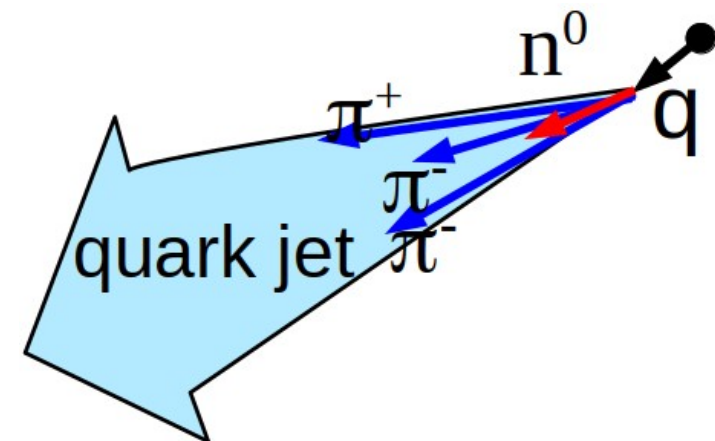
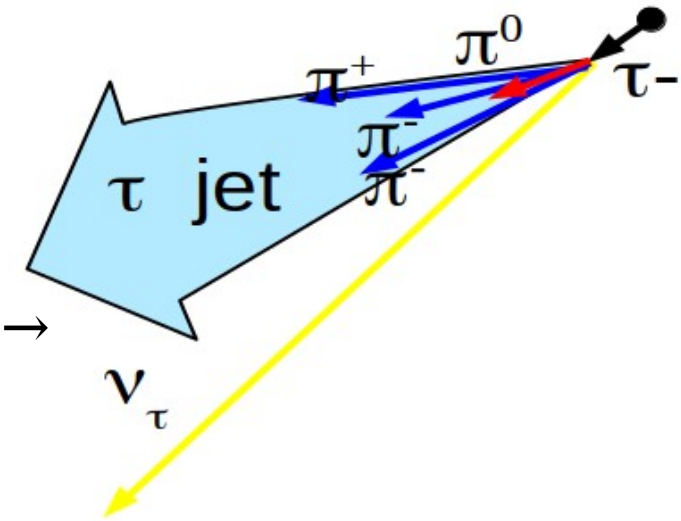
**An issue 1:** behavior for unseen data is hard (impossible?) to predict.

If the model is trained on simulated input the real data might contain unseen parts of the input space

**The task:** jet categorisation:  
hadronic  $\tau$  decay/electron/muon/jet

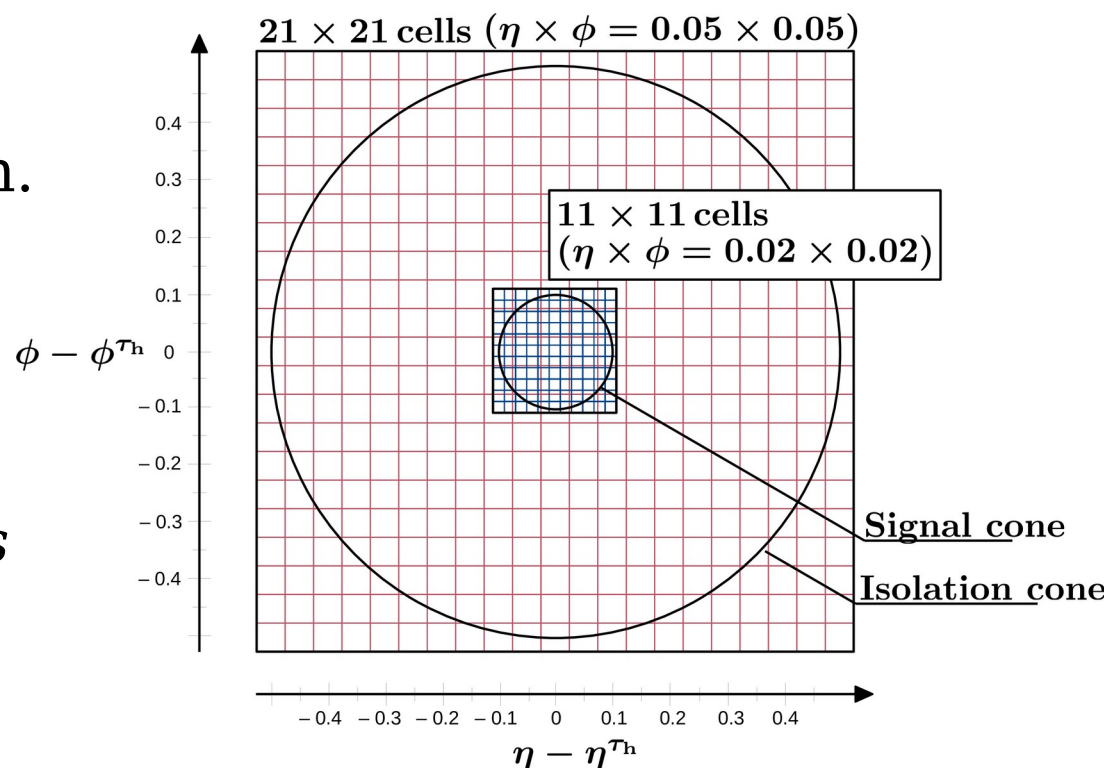
## Context:

- $\tau$  is the heaviest lepton, with short life time  $\rightarrow$  only decay products are observed in the detector
- $\tau$  decays to 1/3 hadrons + neutrinos in 65% cases
- $\tau$  decays look very similar to small jets originating from quarks and gluons



## Input data:

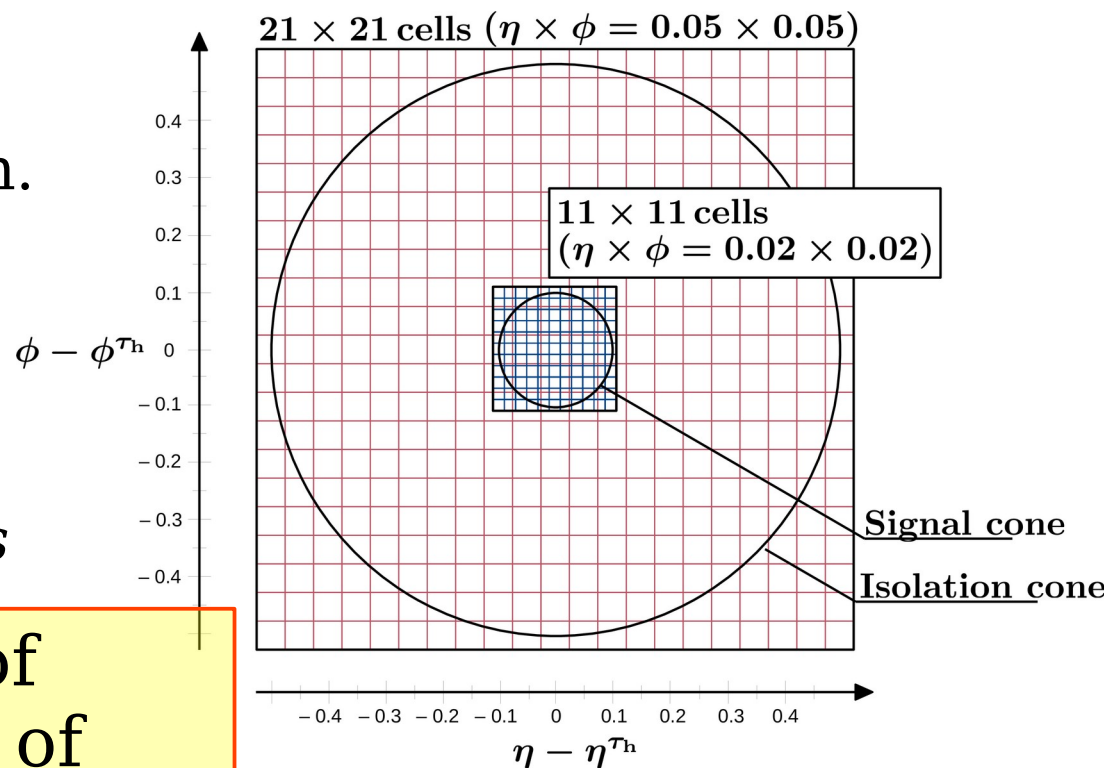
- *low level*:  
highest  $p_T$  particle of given type:  
e/ $\mu$ /charged hadron/neutral hadron  
properties from each of cells  
around the  $\tau$  candidate direction.  
About **20** features per particle
- *high level*:  
**47** human invented features  
derived from particle properties



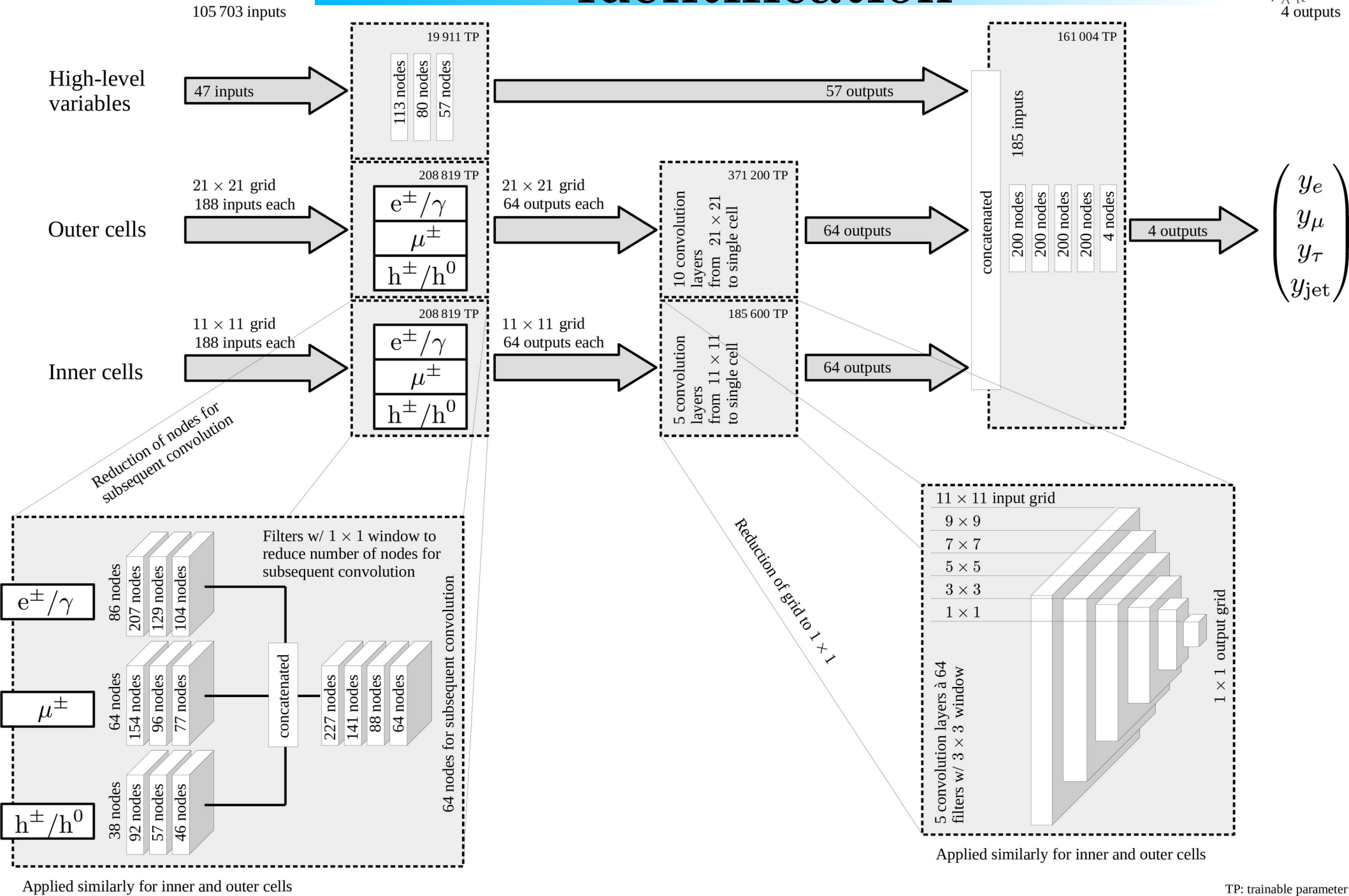
## Input data:

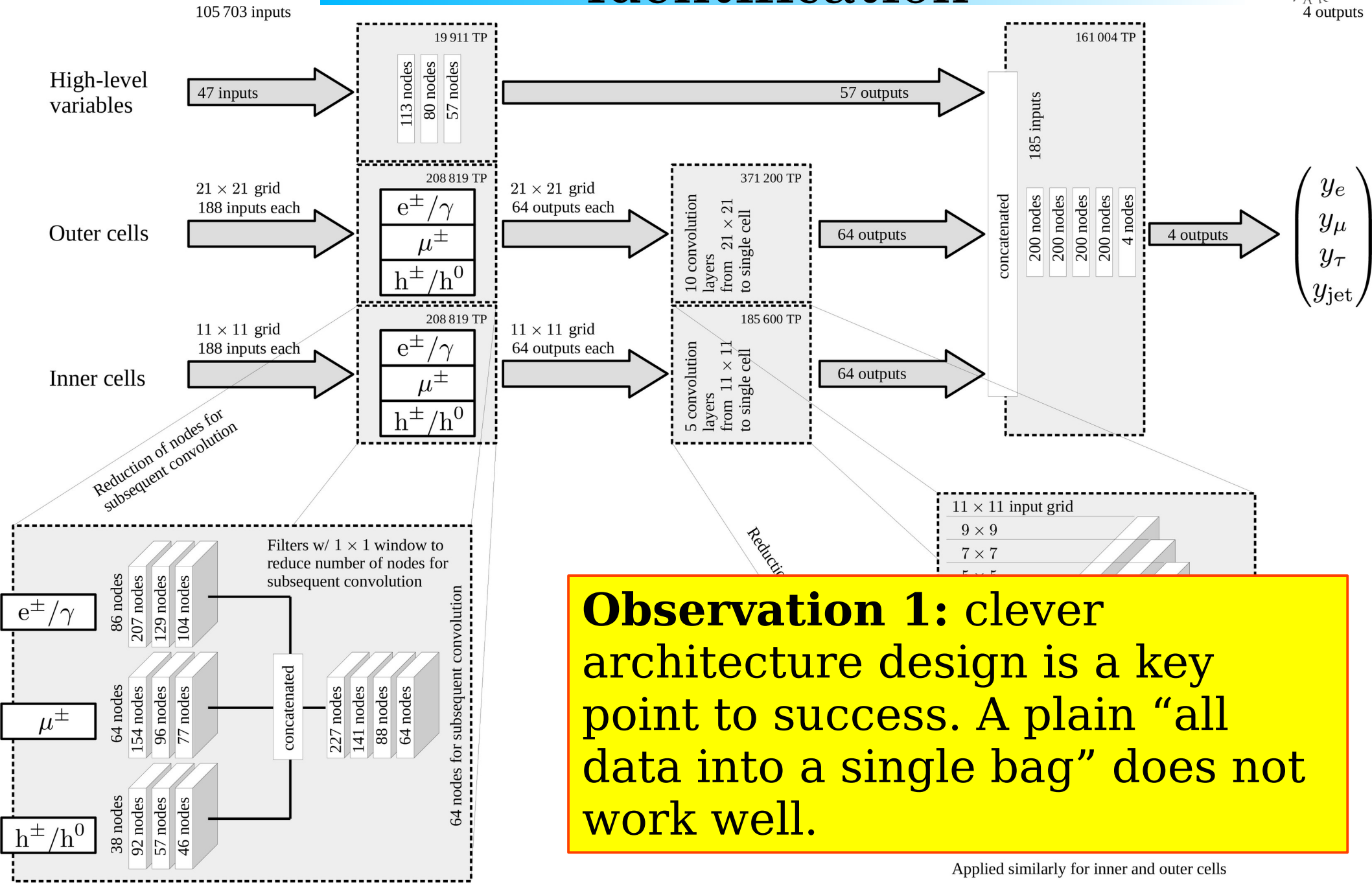
- *low level*:  
highest  $p_T$  particle of given type:  
e/ $\mu$ /charged hadron.neutral hadron  
properties from each of cells  
around the  $\tau$  candidate direction.  
About **20** features per particle
- *high level*:  
**47** human invented features  
derived from particle properties

**Observation 0:** injection of domain knowledge in form of hand crafted features training improves stability and convergence.



# Example: hadronic $\tau$ decay identification





**Observation 1:** clever architecture design is a key point to success. A plain “all data into a single bag” does not work well.

Applied similarly for inner and outer cells

$$L(\mathbf{y}^{\text{true}}, \mathbf{y}; \kappa, \gamma, \omega) = \underbrace{\kappa_\tau H_\tau(\mathbf{y}^{\text{true}}, \mathbf{y}; \omega)}_{\text{(a) Separation of all } \alpha} + \underbrace{(\kappa_e + \kappa_\mu + \kappa_{\text{jet}}) \bar{F}_{\text{cmb}}(1 - y_\tau^{\text{true}}, 1 - y_\tau; \gamma_{\text{cmb}})}_{\text{(b) Focused separation of } e, \mu, \text{ jet from } \tau_h}$$

(a) Separation of all  $\alpha$

(b) Focused separation of  $e, \mu, \text{ jet}$  from  $\tau_h$

$$+ \kappa_F \underbrace{\sum_{i \in \{e, \mu, \text{jet}\}} \kappa_i \hat{\theta}(y_\tau - 0.1) \bar{F}_i(y_i^{\text{true}}, y_i; \gamma_i)}_{\text{(c) Focused separation of } \tau_h \text{ from } e, \mu, \text{ jet for } y_\tau > 0.1}$$

(c) Focused separation of  $\tau_h$  from  $e, \mu, \text{ jet}$  for  $y_\tau > 0.1$

## Focal loss function:

$$F(y^{\text{true}}, y; \gamma) = -y^{\text{true}} (1 - y)^\gamma \log(y) \quad \bar{F}(y^{\text{true}}, y; \gamma) = \mathcal{N} F(y^{\text{true}}, y; \gamma),$$



$$L(\mathbf{y}^{\text{true}}, \mathbf{y}; \kappa, \gamma, \omega) = \underbrace{\kappa_\tau H_\tau(\mathbf{y}^{\text{true}}, \mathbf{y}; \omega)}_{\text{(a) Separation of all } \alpha} + \underbrace{(\kappa_e + \kappa_\mu + \kappa_{\text{jet}}) \bar{F}_{\text{cmb}}(1 - y_\tau^{\text{true}}, 1 - y_\tau; \gamma_{\text{cmb}})}_{\text{(b) Focused separation of } e, \mu, \text{ jet from } \tau_h}$$

(a) Separation of all  $\alpha$                       (b) Focused separation of  $e, \mu, \text{ jet from } \tau_h$

$$+ \underbrace{\kappa_F \sum_{i \in \{e, \mu, \text{jet}\}} \kappa_i \hat{\theta}(y_\tau - 0.1) \bar{F}_i(y_i^{\text{true}}, y_i; \gamma_i)}_{\text{(c) Focused separation of } \tau_h \text{ from } e, \mu, \text{ jet for } y_\tau > 0.1}$$

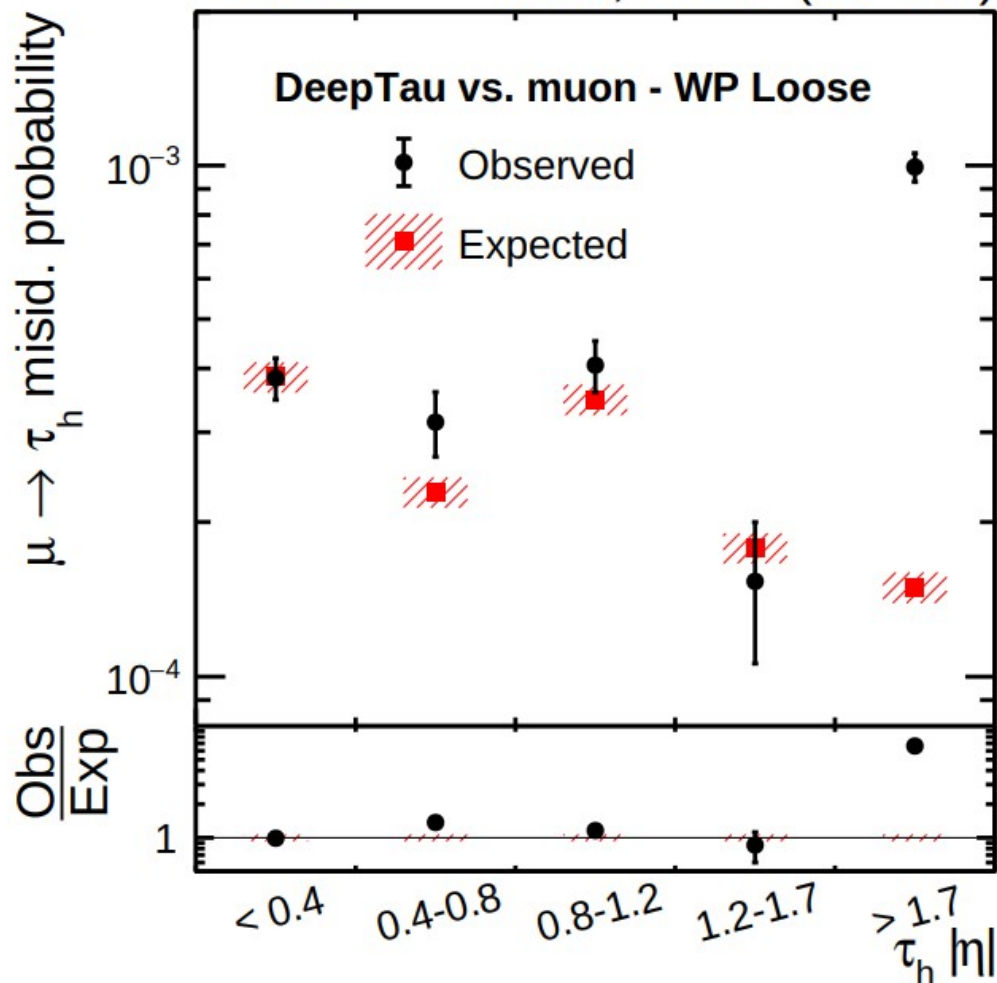
(c) Focused separation of  $\tau_h$  from  $e, \mu, \text{ jet for } y_\tau > 0.1$

**Observation 2:** clever loss function improves model performance in intermediate region of efficiency (aka. recall.)

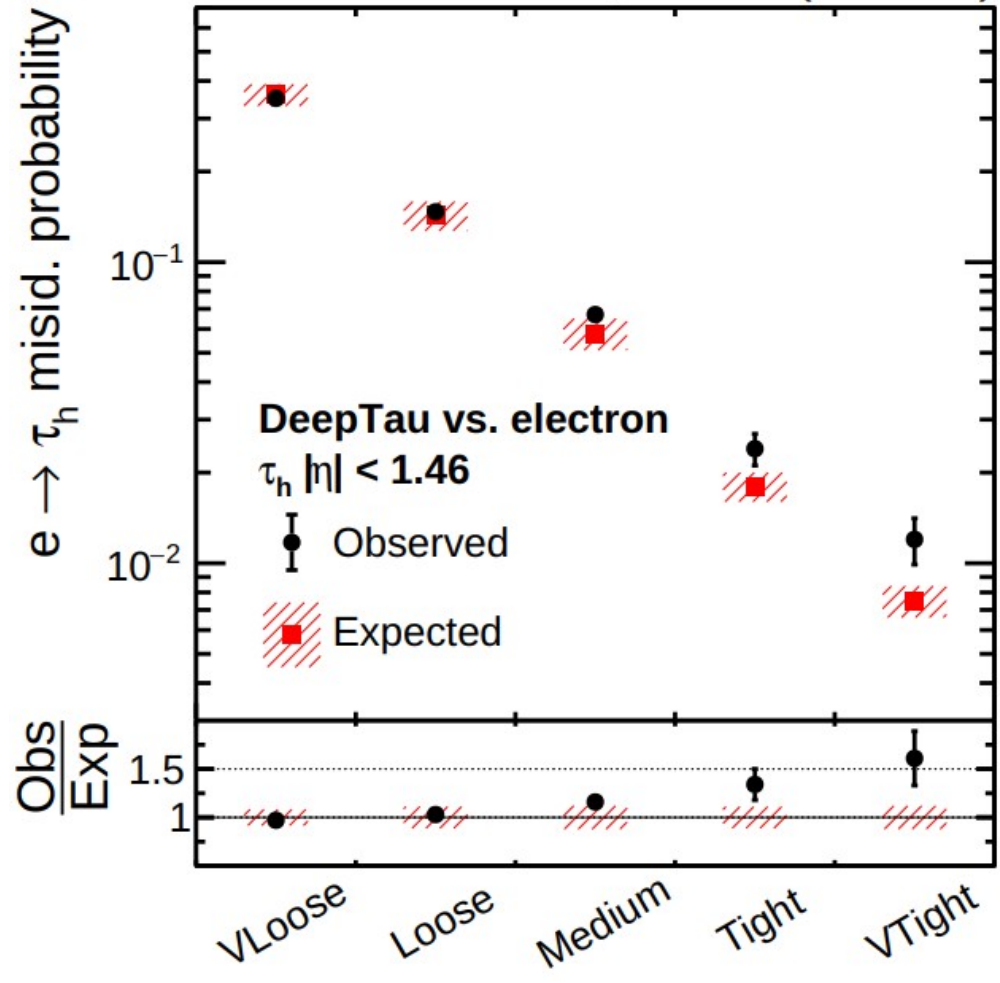
**Focal loss function:**

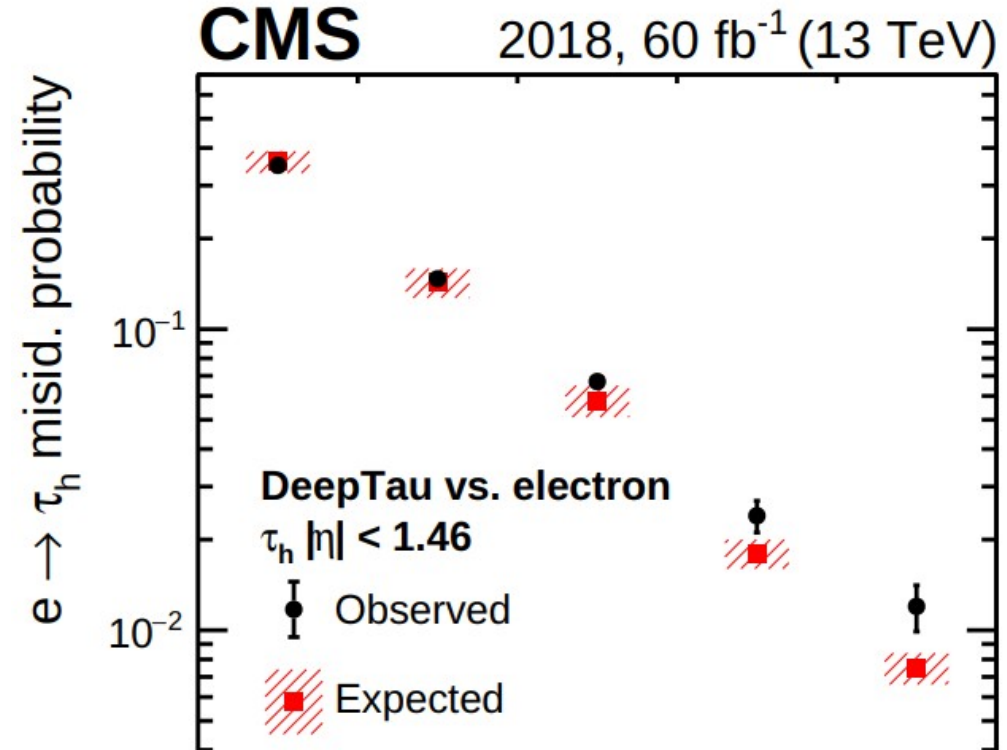
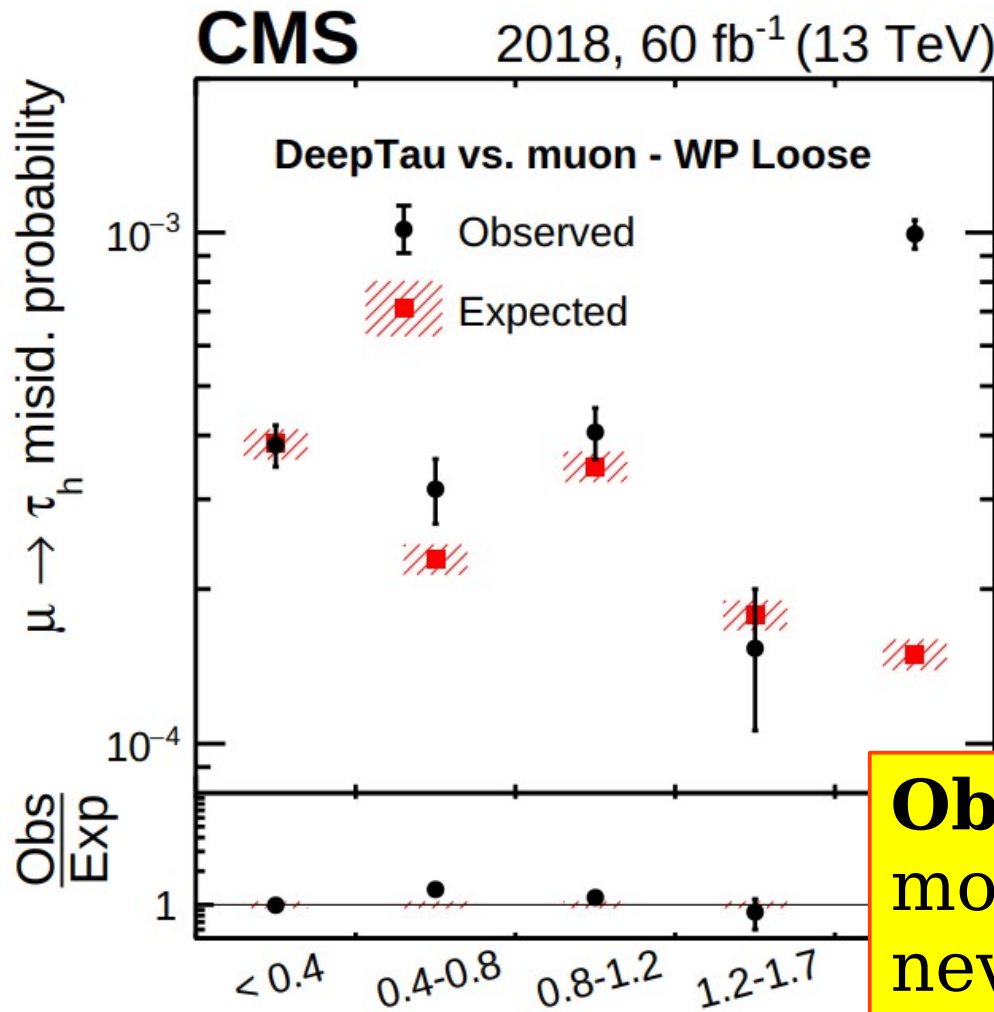
$$F(y^{\text{true}}, y; \gamma) = -y^{\text{true}} (1 - y)^\gamma \log(y) \qquad \bar{F}(y^{\text{true}}, y; \gamma) = \mathcal{N} F(y^{\text{true}}, y; \gamma),$$

**CMS** 2018, 60 fb<sup>-1</sup> (13 TeV)



**CMS** 2018, 60 fb<sup>-1</sup> (13 TeV)





**Observation 3 (trivial):** models trained on simulation never work exactly as expected on a real data. The more rare cases the harder to follow the expectations.

Explaining and Harnessing Adversarial Examples  
arXiv:1412.6572



+

?

=



$x$

“panda”

57.7% confidence

$x + ?$

“gibbon”

99.3 % confidence

Small and smart distortion to input data can lead to dramatic change in model response. This effect is exploited by **bad people** to perform Adversarial Attacks.

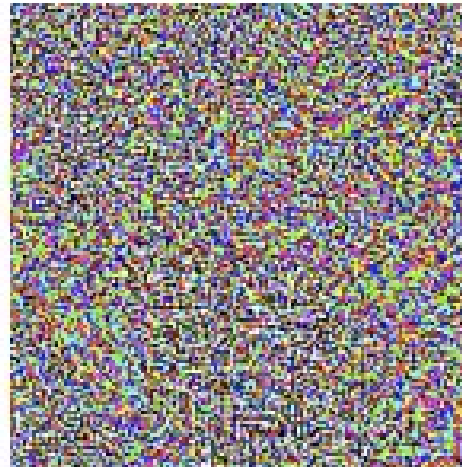


$x$

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

=



$x +$

$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

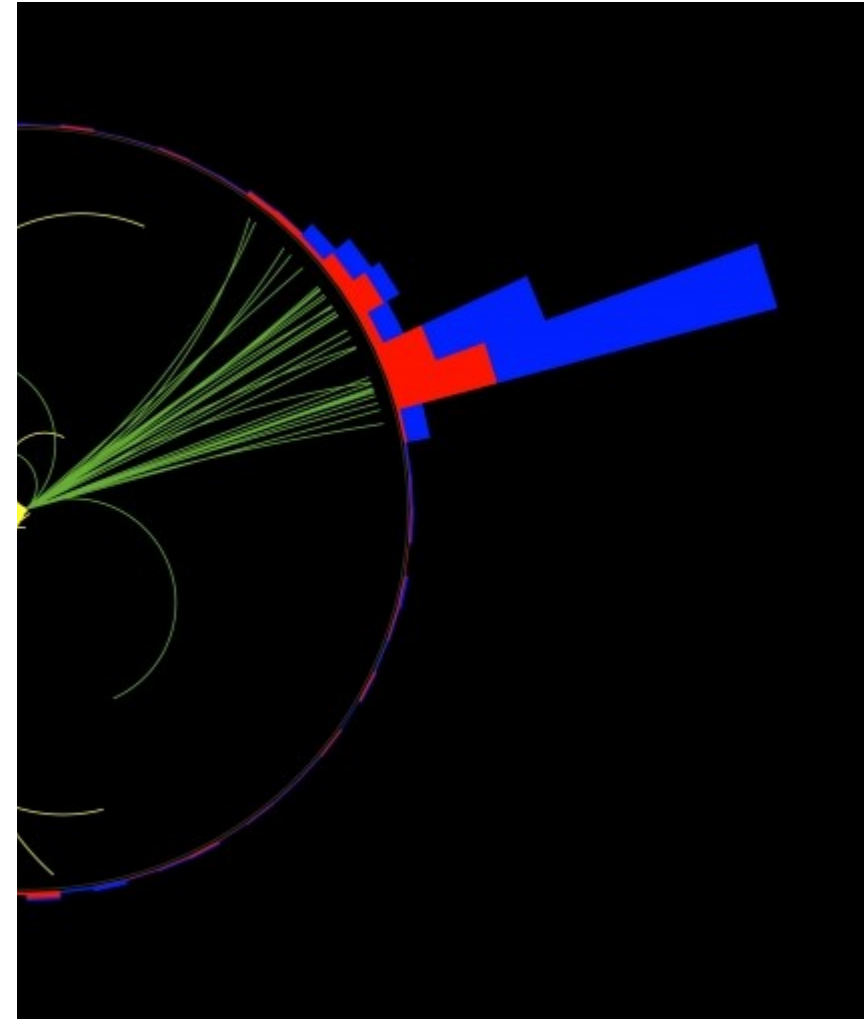
Small and smart distortion to input data can lead to dramatic change in model response. This effect is exploited by **bad people** to make an Adversarial Attacks.

**Good people** can use it to increase robustness of a model to features distortions

**The task:** improve network robustness against Data – Monte Carlo differences.

## Context:

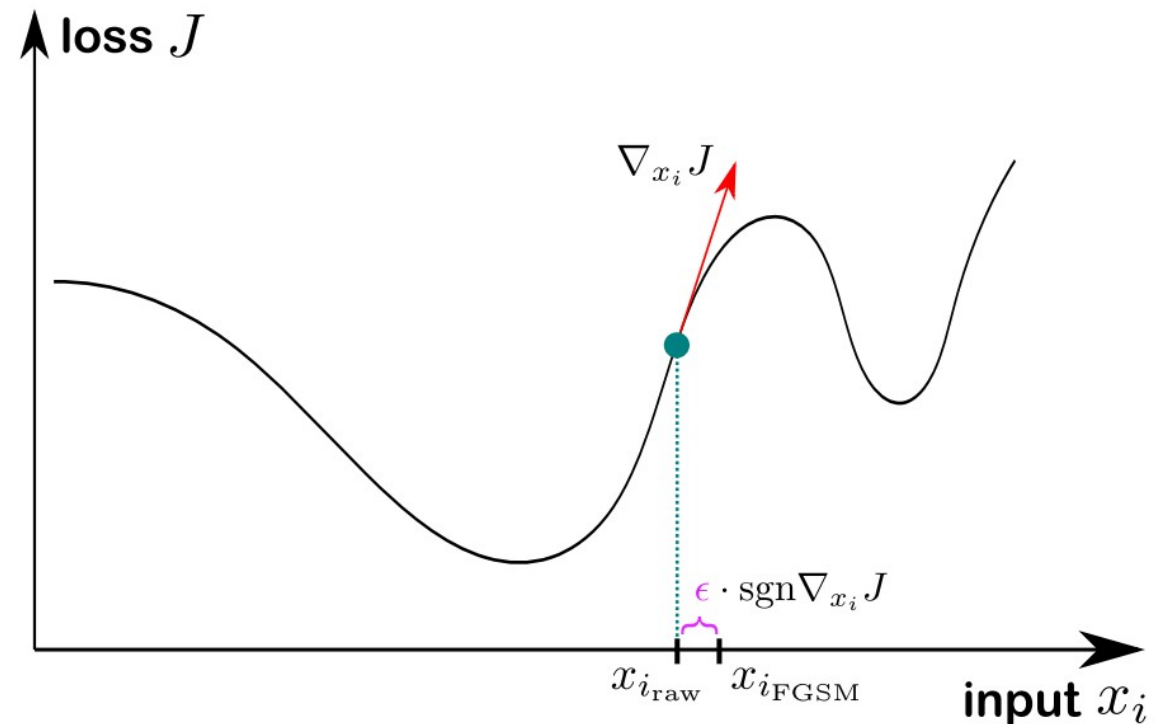
- jet classification:
  - light quark or gluon jets
  - c-quark jets
  - b-quark jets
- features:
  - *low level*: 21 parameters of 6 tracks within a jet
  - *high level*: 14 jet properties

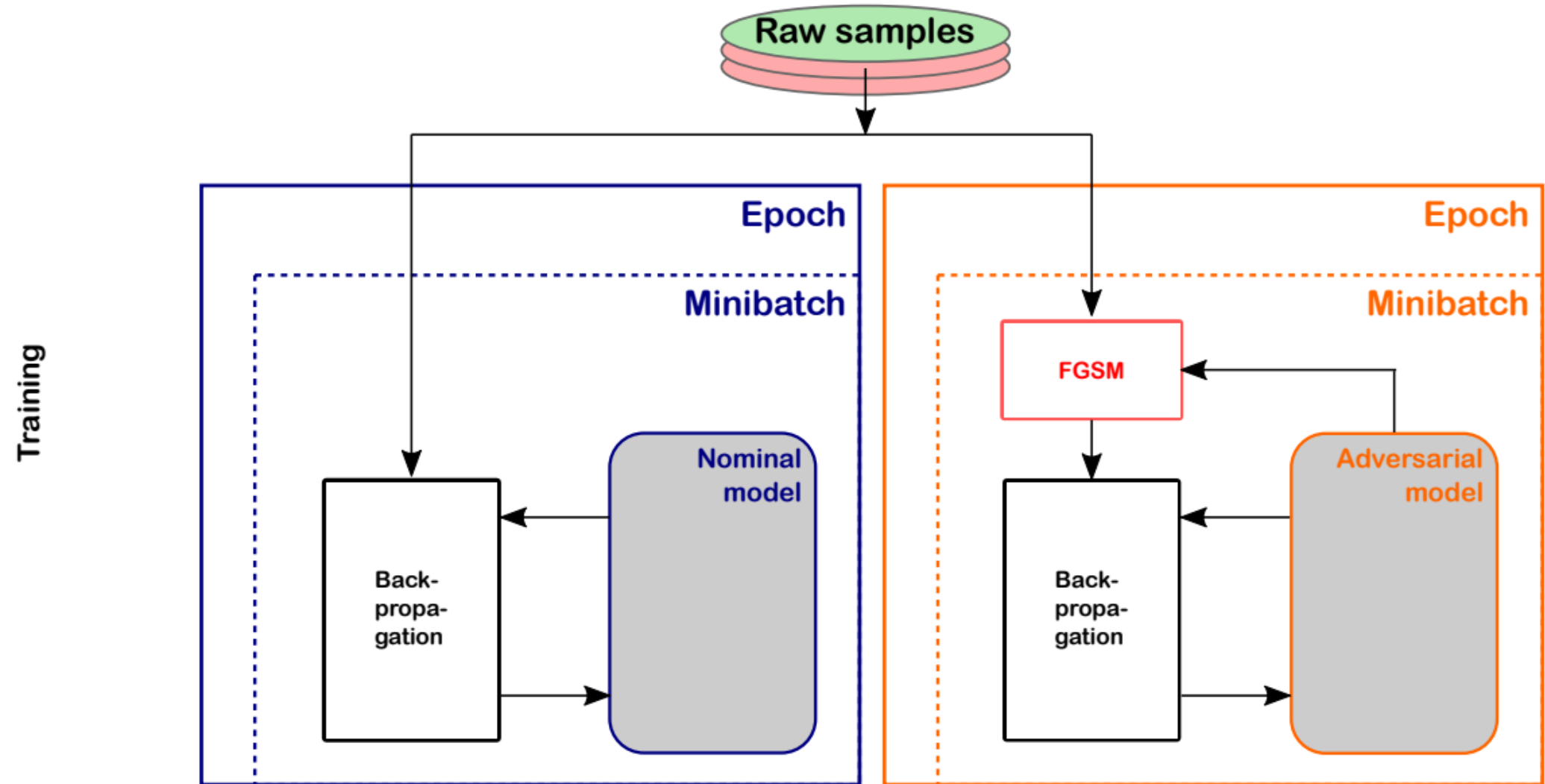


<https://cms.cern/news/machining-jets>

$$x_{\text{FGSM}} = x_{\text{raw}} + \epsilon \cdot \text{sgn}(\nabla_{x_{\text{raw}}} J(x_{\text{raw}}, y))$$

$\epsilon = 0.01$  used  
during the training  
phase

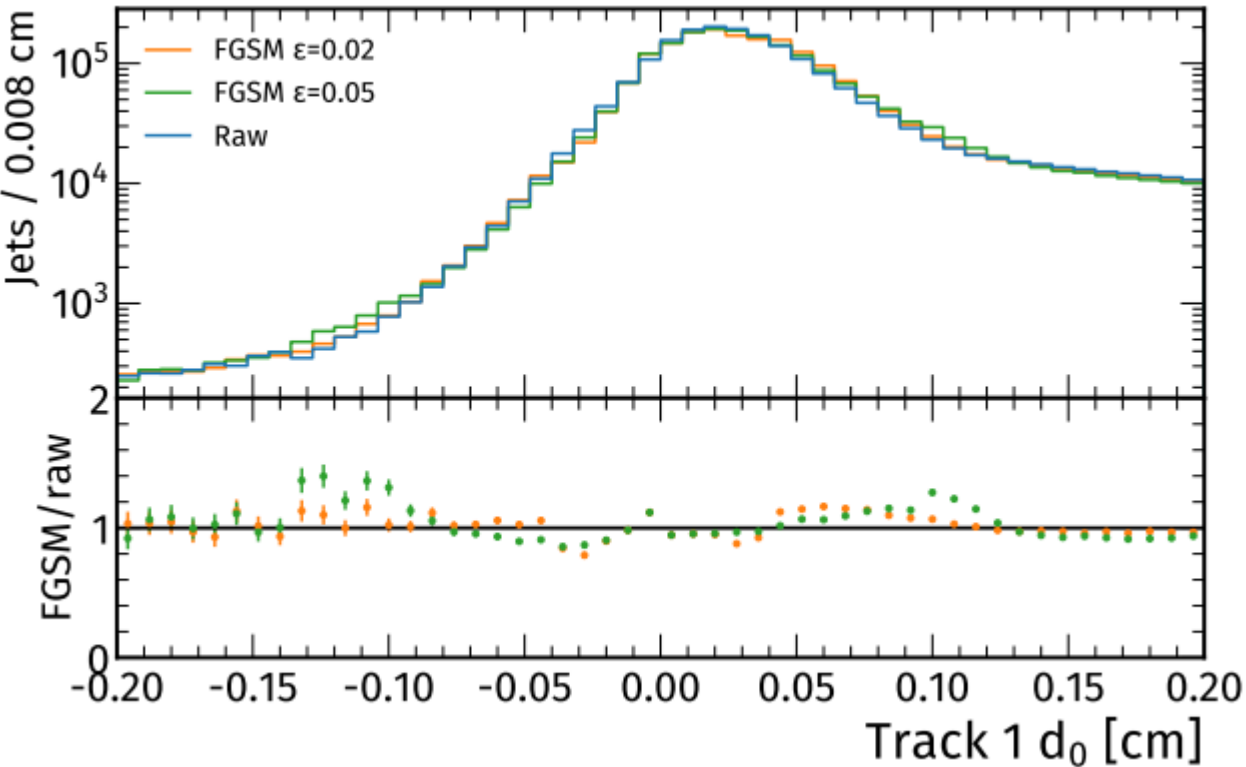
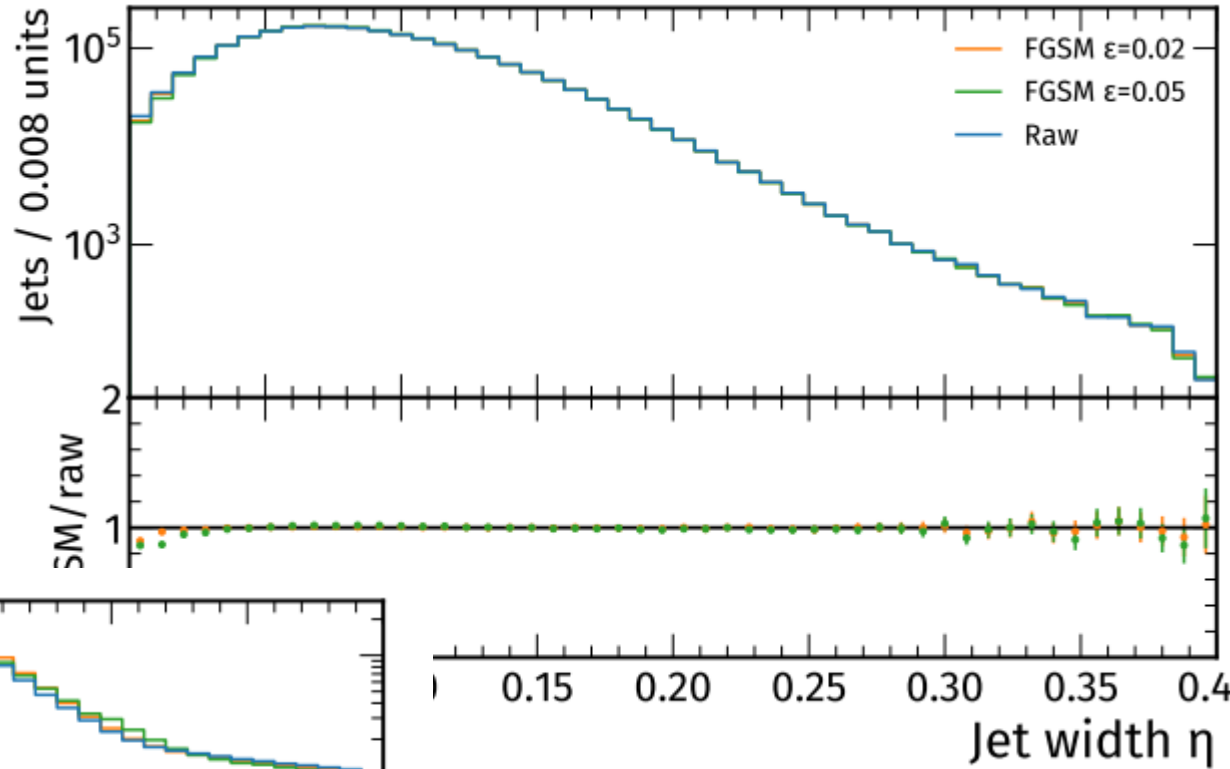




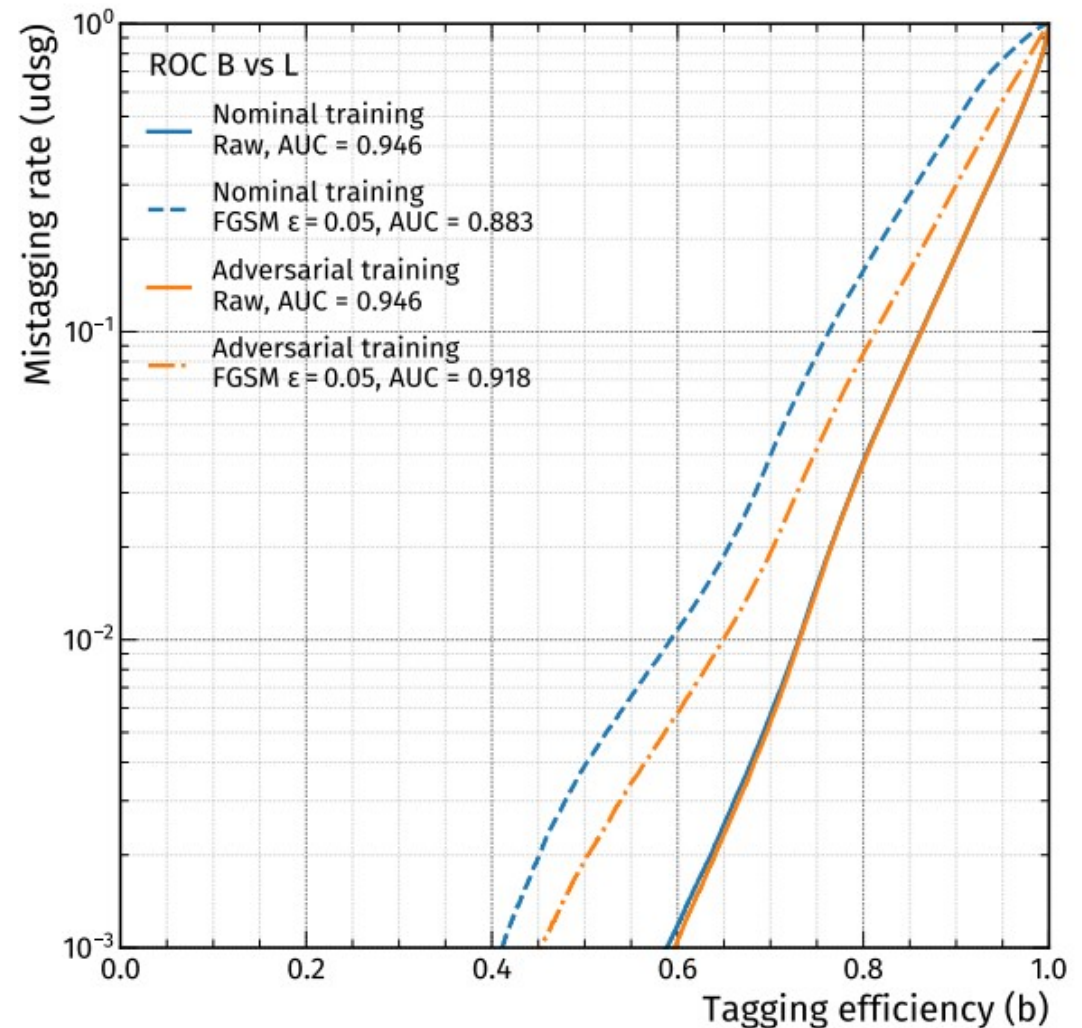


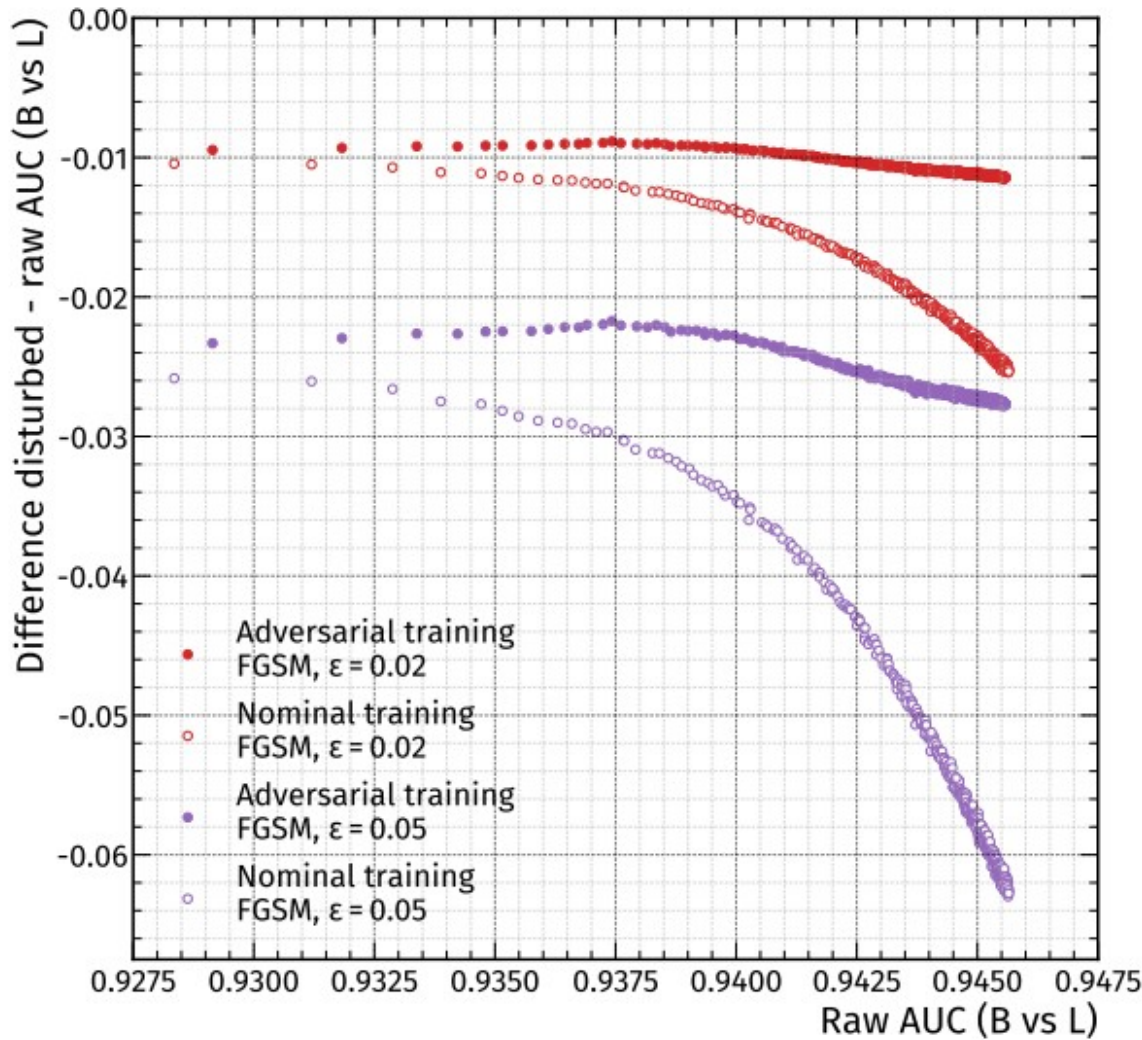
Comput Softw Big Sci 6 (2022) 15

The variable distribution distortion from FGSM is quite small



FGSM acts as a regularisation, which does not affect the performance of the model, but improves its generalization capabilities.





FGSM regularisation, works better than early stopping regularisation.

- Machine Learning is not a magic ward – this is yet another technology

but keep in mind what Arthur C. Clare said: “Any sufficiently advanced technology is indistinguishable from magic.”

- „ordinary” ML users should concentrate on creative problem formulation instead of attempting to invent a new, complicated, architecture
- we should remember about unavoidable difference between features distributions used for training (MC) and inference (Data)



## A Living Review of Machine Learning for Particle Physics

*Modern machine learning techniques, including deep learning, is rapidly being applied, adapted, and developed for high energy physics. The goal of this document is to provide a nearly comprehensive list of citations for those developing and applying these approaches to experimental, phenomenological, or theoretical analyses. As a living document, it will be updated as often as possible to incorporate the latest developments. A list of proper (unchanging) reviews can be found within. Papers are grouped into a small set of topics to be as useful as possible. Suggestions are most welcome.*

### Table of contents

- Reviews
  - Modern reviews
  - Specialized reviews
  - Classical papers
  - Datasets
- Classification
  - Parameterized classifiers



Backup slides