



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

# Tree Tensor Network supervised classifier for High Energy Physics

Rocco Ardino

Mat: 1231629

[rocco.ardino@studenti.unipd.it](mailto:rocco.ardino@studenti.unipd.it)

Alessandro Valente

Mat: 1234429

[alessandro.valente.4@studenti.unipd.it](mailto:alessandro.valente.4@studenti.unipd.it)

Quantum Information and Computing  
a.y. 2020/21

March 19, 2021

## 1 Introduction

## 2 Theory

- Kinematics of high energy collisions
- HIGGS dataset
- Tree Tensor Networks
- Advanced techniques

## 3 Code Implementation

- Input data preprocessing
- TTN layer: framework and initialisation
- Model building and training

## 4 Results

- Comparison between “pure” and advanced TTN models
- TTN advanced model characterisation
- TTN advanced model final results

## 5 Conclusions

## 6 Back-up

## High Energy Physics (HEP):

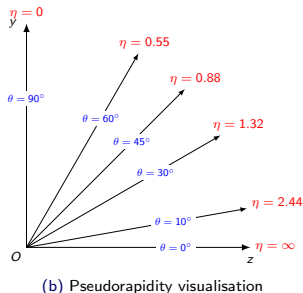
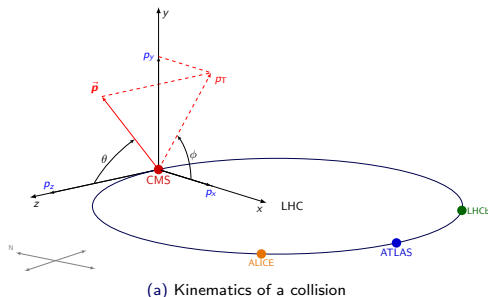
- purpose of understanding the nature of the particles that constitute matter and radiation
- Standard Model as theoretical climax
- still lots of unanswered questions
- ⇒ still New Physics to discover

## Application of Machine Learning techniques in HEP:

- current techniques used in HEP fail to capture all the available information
- as proved by Baldi et al ([1]), Machine Learning can overcome this issue
- many other advantages, such as solution to the curse of high dimensionality of data
- e.g., Artificial Neural Networks exploited to build powerful classifiers

## New approach with Tree Tensor Networks (TTN):

- quantum-inspired version of Biological Neural Networks
- structure based on Tensor Network methods
- possibility to solve non-convex optimisation tasks, such as loss functions minimisation



**Figure:** LHC structure and kinematics of a product particle, emitted with a polar angle  $\theta$  and azimuth angle  $\phi$ , in **1a**. In **1b**, visualisation of the linking between the polar angle and the pseudorapidity  $\eta$ .

## Product particles of $pp$ collisions

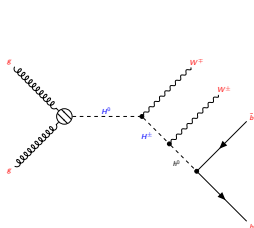
- Leptons ( $\ell^\pm, \bar{\nu}_\ell$ )
- Bosons ( $h^0, Z^0, W^\pm$ )
- Quarks ( $q$ ) and Gluons ( $g$ )
- Jets from quark/gluon hadronisation ( $j$ )

## Main features for the analysis

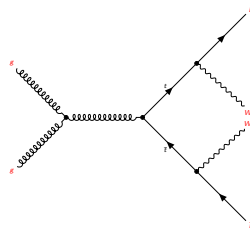
- Transverse momentum  $p_T$
- Pseudorapidity  $\eta$
- Azimuth angle  $\phi$
- $b$ -tag (for jets)

## Monte Carlo generated sample of gluon fusion process:

- a mass of the Higgs  $m_{h^0} = 125$  GeV is assumed
- new exotic Higgs bosons  $H^0$  and  $H^\pm$  are introduced ( $m_{H^0} = 425$  GeV,  $m_{H^\pm} = 325$  GeV)
- 21 low level features (such as  $p_T$ ,  $\eta$ ,  $\phi$ ,  $b$ -tags, ...)
- 7 high level features (invariant masses distributions)



(a) Signal process



(b) Background process

**Figure:** Feynman diagrams for the Monte Carlo simulated events. In 2a, the diagram of the signal channel is portrayed with the exotic Higgs bosons  $H^0$  and  $H^\pm$ . In 2b, the background diagram is represented.

$$\text{Signal: } gg \rightarrow H^0 \rightarrow W^\mp H^\pm \rightarrow W^\mp W^\pm h^0 \rightarrow W^\mp W^\pm b\bar{b} \quad (1)$$

$$\text{Background: } gg \rightarrow g \rightarrow t\bar{t} \rightarrow W^\mp W^\pm b\bar{b} \quad (2)$$

## Tensor Networks (TNs) as core of TTNs:

- factorisations of high rank tensors into networks of smaller rank tensors
- intuitive graphical language
  - tensors are notated by solid shapes, tensor indices are notated by lines emanating from these shapes
  - connecting two index lines implies a contraction, or summation over the connected indices

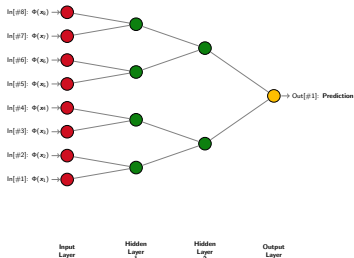
$$\begin{aligned}
 \begin{array}{c} \bullet \\ | \\ i \end{array} &= v_i \longrightarrow \text{Vector} \\
 i \text{ --- } \bullet \text{ --- } j &= M_{ij} \longrightarrow \text{Matrix} \\
 \begin{array}{c} i \text{ --- } \bullet \text{ --- } k \\ | \\ j \end{array} &= T_{ijk} \longrightarrow \text{Rank-3 Tensor}
 \end{aligned} \tag{3}$$

- graphical notation for some of the most common operations between tensors

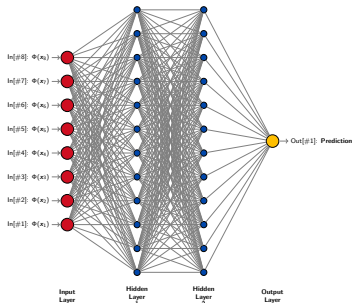
$$\begin{aligned}
 i \text{ --- } \bullet \text{ --- } j \text{ --- } \bullet &= \sum_j M_{ij} v_j \longrightarrow \text{Vector-Matrix} \\
 i \text{ --- } \bullet \text{ --- } j \text{ --- } \bullet \text{ --- } k &= \sum_j A_{ij} B_{jk} \longrightarrow \text{Matrix-Matrix}
 \end{aligned} \tag{4}$$

## TNs for Machine Learning in a nutshell:

- input features  $x$ , mapped into  $\Phi(x)$  through an apposite feature map
- TN with a tree-like structure  $T(w; \chi)$ , with:
  - $w$  the entries of the tensors, namely the weights to be tuned by the learning algorithm
  - $\chi$  the bond dimension, which controls the complexity of the structure
- decision function:  $f(x; w) = T(w; \chi) \cdot \Phi(x)$



(a) Example of TTN structure



(b) Example of DNN structure

**Figure:** Comparison between examples of structures for a TTN, in 3a, and a DNN, in 3b. In particular, in the TTN the dimension of the bonds between the tensors in the hidden layers is the bond dimension  $\chi$ .

## Loss function:

- need to minimise the ladder on a fraction of dataset reserved for training the TTN
- quantification of TTN misclassification on new fraction of dataset reserved for validation
- **binary cross-entropy** possible loss for signal-versus-background discrimination with:
  - $y_{\text{true}}^{(i)}$ : true label of  $i^{\text{th}}$  sample of validation set
  - $y_{\text{pred}}^{(i)}$ : TTN predicted label for the  $i^{\text{th}}$  sample of validation set

$$L(y_{\text{true}}, y_{\text{pred}}) = \sum_i^{n_{\text{val}}} y_{\text{true}}^{(i)} \log(y_{\text{pred}}^{(i)}) + (1 - y_{\text{true}}^{(i)}) \log(1 - y_{\text{pred}}^{(i)}) \quad (5)$$

## Metrics for performances quantification:

- **accuracy**: fraction of correctly classified samples (using a decision threshold  $\xi \in (0, 1)$ )
- **AUC**: area under the Receiver Operating Characteristic (ROC) Curve
  - True Positive Rate (TPR) in function of the False Positive Rate (FPR)
  - obtained by sweeping the decision threshold  $\xi$  in  $(0, 1)$
  - directly linked to concept of discovery significance commonly used in HEP

## Optimiser for weights update:

- **ADAM**: adapts the correction to the weights at each step parameter per parameter
- update after the TTN has processed a **mini batch** of input training set of dimension  $m$
- after processing all mini batches in the training set, a **training epoch** has ended



## Activation function:

- function applied at the output of a layer
- source of non-linearity to enlarge the space of functions that the TTN can approximate
- for this work, **ELU** (for inner layers) and **sigmoid** (for output layer) tested:

$$\text{ELU}(x) = \begin{cases} x & x \geq 0 \\ a(e^x - 1) & x < 0 \end{cases} \quad \sigma(x) = \frac{1}{1 + e^{-x}} \quad (6)$$

## Kernel regularisation:

- add a penalty for large weights to the loss function to avoid overfit
- **$\ell_2$  regularisation**:

$$J(y_{\text{true}}, y_{\text{pred}}) = L(y_{\text{true}}, y_{\text{pred}}) + \lambda \sum_{i=1}^{n_{\text{weights}}} |w_i|^2 \quad (7)$$

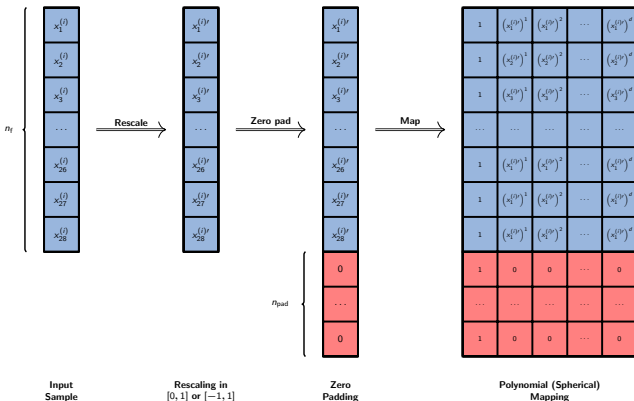
## Batch normalisation:

- training networks with lots of layer is very challenging
- calculate mean and standard deviation of the output of each layer for every mini batch
- $\Rightarrow$  perform a **standardisation**
- $\Rightarrow$  **speed up of learning algorithm convergence**

# Input data preprocessing

## Input features preprocessing workflow:

- **rescaled** in  $[0, 1]$  or  $[-1, 1]$  (depending on the feature)
- **padded** to match the number of input “legs” of the TTN
- **mapped** through a feature map to enhance the performances



**Figure:** Workflow of preprocessing procedure, starting from the rescaling of data, going through the zero padding in order to get proper dimensions for the input of the TTN, and lastly the polynomial or spherical mapping.

# Input data preprocessing: rescaling and padding

## Rescaling

- rescale “positive” and “negative” features in  $[0, 1]$  and  $[-1, 1]$ , respectively

$$x_j^{(i)} \longrightarrow x_j^{(i)'} = \frac{x_j^{(i)}}{m_j} \quad \text{with} \quad m_j = \max_{i \in \mathcal{D}_{\text{train}}} |x_j^{(i)}| \quad (8)$$

## Padding

- total number of features  $n_f$  must be divisible by features contracted in each site  $n_{\text{con}}$
- dataset is “padded” with fictitious features to reach  $\tilde{n}_f$  features

$$\tilde{n}_f = \min_{n \in \mathbb{N}} \{n \geq n_f : n = (n_{\text{con}})^m, m \in \mathbb{N}\} \quad (9)$$

```
# Rescaling
def Standardize(x, nt):
    for j in range(x.shape[1]):
        vec = x[:, j]
        vec_norm = vec[:nt]
        x[:, j] = vec / np.max(np.abs(vec_norm))
    return x

# Padding
def PadToOrder(x, con_order):
    # compute number of padded features
    n_pad = int( con_order** ( math.ceil( math.log(x.shape[1], con_order) ) ) - x.shape[1] )
    # pad dataset
    x = np.append(x, np.zeros((x.shape[0], n_pad)), axis=1)
    return x
```

# Input data preprocessing: feature map

## Feature map:

- applying a map to the input features enhances the classifier performances
- **polynomial map**, common in standard ML classification tasks

$$\Phi_d^{\text{pol}}(x) = [1, x, \dots, x^d] \quad (10)$$

- **spherical map**, quantum inspired, at order 2 maps features to spins:

$$\Phi_d^{\text{sph}}(x) = [\phi_d^{(1)}(x), \dots, \phi_d^{(d)}(x)] \quad (11)$$

$$\phi_d^{(s)}(x) = \sqrt{\binom{d-1}{s-1}} \left( \cos\left(\frac{\pi}{2}x\right) \right)^{d-s} \left( \sin\left(\frac{\pi}{2}x\right) \right)^{s-1}$$

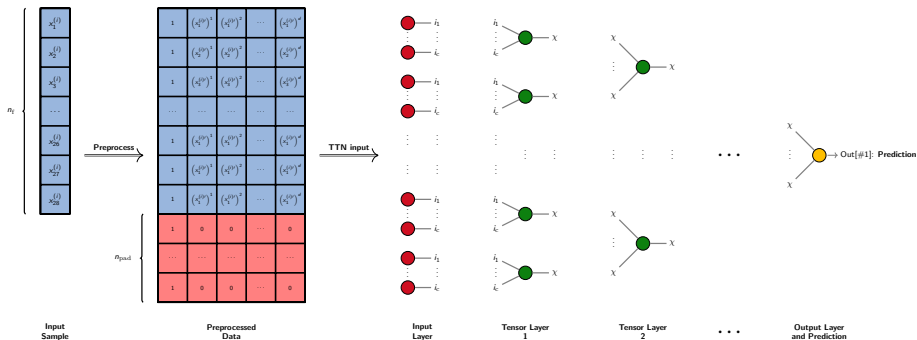
```
# Spherical map
def SphericalMap(x, order=2, dtype=np.float32):
    x_map = np.zeros((x.shape[0], x.shape[1], order), dtype=dtype)
    for i in range(order):
        comb_coef = np.sqrt(scipy.special.comb(order-1, i))
        x_map[:, :, i] = comb_coef * np.power(np.cos(x), order-1-i) * np.power(np.sin(x), i)
    return x_map

# Polynomial map
def PolynomialMap(x, order=2, dtype=np.float32):
    x_map = np.zeros((x.shape[0], x.shape[1], order+1), dtype=dtype)
    for i in range(order+1):
        x_map[:, :, i] = np.power(x, i)
    return x_map
```

# TTN framework: layer

## Framework and workflow:

- TTN classifier is a series of layers built using **TensorFlow** and **TensorNetwork** libraries
- possibility of running on **both CPU and GPU** through Keras API
- layer object main input parameters:
  - **n\_contraction**: number of features to contract in each node site, namely  $n_{\text{con}}$
  - **bond\_dim**: the dimension  $\chi$  of the bonds between the tensor nodes inside the TTN
  - **activation**: string carrying the name of the activation function to use, if specified
  - **use\_batch\_norm**: introduce the batch normalisation of the layers if true is specified



**Figure:** Schematic representation of workflow of the data preprocessing and of the TTN structure, with  $i_c = i_{\text{con}}$ .

## Node creation and contraction:

- the computations inside a TTN layer are divided in three logical steps:
  - **nodes creations**: input and weight tensors are initialised as TensorNetwork nodes
  - **edge connection**: edges of each tensor node are connected to the corresponding input node
  - **edge contraction**: node contractions along connected edges are executed
- In each step the same loop structure is used but separating allows better GPU parallelisation

```
# tensor nodes initialisation
for i in range(len(nodes)):
    for j in range(n_contr):
        # create feature nodes
        x_nodes.append(tn.Node(x[n_contr*i+j], name='xnode', backend="tensorflow"))
# create ttn node
tn_nodes.append(tn.Node(nodes[i], name=f'node_{i}', backend="tensorflow"))
```

```
# tensor nodes edges connection
for i in range(len(nodes)):
    for j in range(n_contr):
        # make connections
        x_nodes[n_contr*i+j][0] ^ tn_nodes[i][j]
```

```
# tensor nodes edges contraction
for i in range(len(nodes)):
    result.append(
        tn.contractors.greedy([x_nodes[n_contr*i+j] for j in range(n_contr)]+[tn_nodes[i]])
    )
```

## Model structure:

- using the TTN layers, a classifier can be created using the **Sequential Keras API**
- when instantiating the layers, it is possible to specify many parameters, such as:
  - **input\_shape**: shape of input samples
  - **n\_contraction**: number of features to contract at each weight node
  - **bond\_dim**: the dimension  $\chi$  of the bonds between the tensor nodes inside the TTN
  - **activation**: string carrying the name of the activation function to use, if specified
  - **use\_bias**: introduce bias weights if true is specified
  - **use\_batch\_norm**: introduce the batch normalisation of the layers if true is specified
  - **kernel\_regularizer**: TensorFlow regulariser object to introduce the regularisation

```
# create keras sequential model
tn_model = Sequential()

# first layer, input shape must be specified
tn_model.add( TTN_SingleNode(bond_dim=10, activation='elu',          n_contraction=2,
                             input_shape=(x_train.shape[1:])) )

# intermediate layers, input shape computed from previous layers output
tn_model.add( TTN_SingleNode(bond_dim=10, activation='elu',          n_contraction=2 ) )
tn_model.add( TTN_SingleNode(bond_dim=10, activation='elu',          n_contraction=2 ) )
tn_model.add( TTN_SingleNode(bond_dim=10, activation='elu',          n_contraction=2 ) )

# last layer, bond dim 1 and sigmoid function to interpret output as probability
tn_model.add( TTN_SingleNode(bond_dim=1, activation='sigmoid', n_contraction=2 ) )
```

# TTN framework: model training

## Model compilation:

- after layers instantiation, the model is compiled inside TensorFlow framework
- optimiser, loss function and metrics are specified at compilation time
  - **ADAM** optimiser
  - **binary cross entropy** loss
  - **accuracy** and **AUC** metrics

## Model training:

- TTN classifier is trained using TensorFlow framework
- training and validation sets are provided to the training function, alongside with training parameters:
  - number of **epochs** of training
  - **batch size**, namely after how many samples processed the weights should be updated by the optimiser

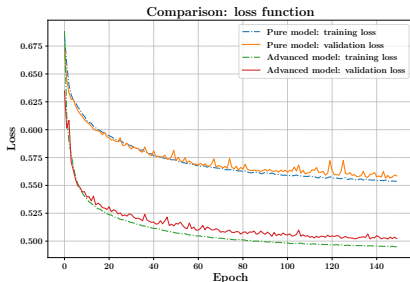
```
# model compilation
tn_model.compile(
    optimizer = 'adam',           # optimizer for training
    loss       = 'binary_crossentropy', # loss function to minimize
    metrics    = ['accuracy', 'AUC']  # metrics to monitor
)

# model training
with tf.device('/device:gpu:0'):  # execution on GPU
    history = tn_model.fit(
        x_train, y_train,         # training set
        validation_data = (x_val, y_val), # validation set
        epochs          = 150,      # training epochs
        batch_size       = 5000     # batch size
    )
```



## “Pure” TTN model

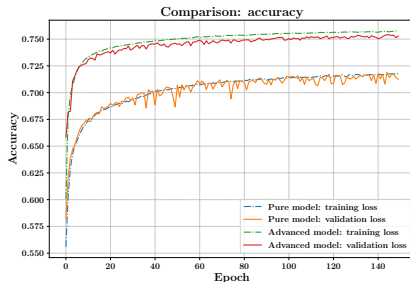
- Standard TTN structure
  - no advanced techniques
  - no batch normalisation, no regularisation
  - no activation for inner layers
  - only sigmoid activation for final prediction



(a) Loss metric over training epochs

## Advanced TTN model

- ML optimisations added:
  - batch normalisation
  - $\ell_2$  regularisation
  - ELU activation function for inner layers
  - sigmoid activation for final prediction



(b) Accuracy metric over training epochs

**Figure:** Comparison between the “pure” TTN model and the advanced one. The trend of the cost function and of the accuracy score during the training are showed in **6a** and **6b**, respectively.

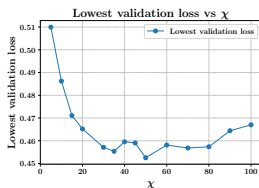
# Model Characterisation: parameters and metrics

## Bond dimension $\chi$ :

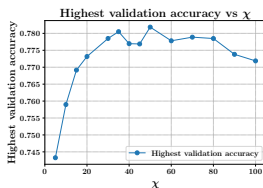
- determines the **number of parameters** and consequently the **model complexity**

$\chi$	Parameters	$\chi$	Parameters	$\chi$	Parameters	$\chi$	Parameters
5	$\approx 2.73 \cdot 10^3$	30	$\approx 3.84 \cdot 10^5$	55	$\approx 2.34 \cdot 10^6$	80	$\approx 7.18 \cdot 10^6$
10	$\approx 1.60 \cdot 10^4$	35	$\approx 6.08 \cdot 10^5$	60	$\approx 3.03 \cdot 10^6$	85	$\approx 8.62 \cdot 10^6$
15	$\approx 5.03 \cdot 10^4$	40	$\approx 9.05 \cdot 10^5$	65	$\approx 3.86 \cdot 10^6$	90	$\approx 1.02 \cdot 10^7$
20	$\approx 1.16 \cdot 10^5$	45	$\approx 1.28 \cdot 10^6$	70	$\approx 4.82 \cdot 10^6$	95	$\approx 1.20 \cdot 10^7$
25	$\approx 2.24 \cdot 10^5$	<b>50</b>	<b><math>\approx 1.76 \cdot 10^6</math></b>	75	$\approx 5.92 \cdot 10^6$	100	$\approx 1.40 \cdot 10^7$

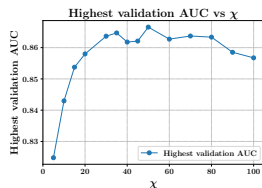
**Table:** Number of parameters of the TTN depending on  $\chi$ , namely the bond dimension of the tensor nodes.



(a) Model lowest loss value



(b) Model highest accuracy score

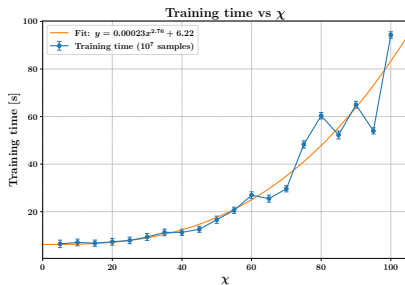


(c) Model highest AUC score

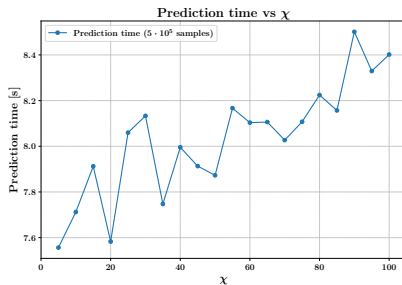
**Figure:** Performances of TTN model depending on the bond dimension. The results are obtained from the validation set, on which the best values over training epochs for the metrics are computed. In particular, in **7a** the lowest loss, in **7b** the highest accuracy and in **7c** the highest AUC, depending on the bond dimension  $\chi$ .

## Time scaling with $\chi$ :

- number of parameters grows with  $\chi$  as  $O(\chi^3)$
- training and prediction time will be higher for more complex models
  - training time follows a **power law with exponent  $k = 2.76$**
  - expectation for the latter is  $k_{th} = 3$  (for rank-3 tensor-vector multiplication)
  - prediction time grows linearly (with a limited number of samples)



(a) Model training time per epoch for  $10^7$  samples



(b) Model prediction time for  $5 \cdot 10^5$  samples

**Figure:** Timing analysis and scaling of TTN classifier depending on the bond dimension  $\chi$ . In **8a** the training time per epoch for  $10^7$  samples is showed, while in **8b** the prediction time for  $5 \cdot 10^5$  samples is visualised.

## Feature map and order:

- performance dependence on the applied map during data preprocessing and on map order
- **polynomial map**, common in standard ML classification tasks

$$\Phi_d^{\text{pol}}(x) = [1, x, \dots, x^d] \quad (12)$$

- **spherical map**, quantum inspired, at order 2 maps features to spins:

$$\Phi_d^{\text{sph}}(x) = [\phi_d^{(1)}(x), \dots, \phi_d^{(d)}(x)] \quad (13)$$

- **map orders**  $d \in \{2, 3, 4, 5\}$  tested

	Spherical Map				Polynomial Map			
Map Order $d$	2	3	4	5	2	3	4	5
Min loss	0.532	0.523	0.512	<b>0.511</b>	0.526	0.519	0.512	0.514
Max accuracy	0.777	0.781	0.781	<b>0.782</b>	0.781	0.781	0.781	0.781
Max AUC	0.862	0.866	0.866	<b>0.867</b>	0.866	0.866	0.866	0.865
Training time [s]	120	127	134	140	125	133	141	150
Prediction time [s]	8.65	8.80	8.91	9.13	8.61	8.80	8.93	9.30

**Table:** Results of the feature map and map order characterisation analysis. The training time is obtained for  $10^7$  input samples, while the prediction time for  $5 \cdot 10^5$  input samples.

## Time scaling with batch size $m$ :

- $m$  is the number of samples processed before weight update by optimiser
- strongly influences training time
- a higher batch size needs more training epochs to reach same performances
- bigger batch sizes more efficient with GPU acceleration

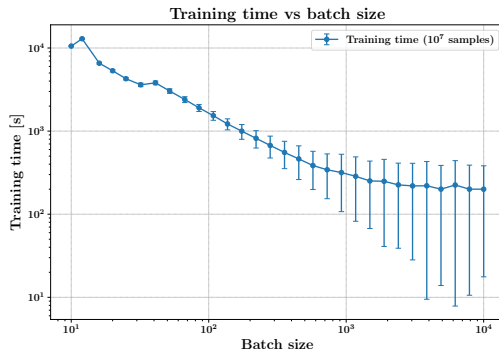


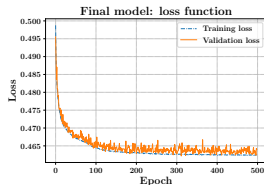
Figure: TTN classifier training time per epoch for  $10^7$  input samples depending on the batch size.

## Best model creation:

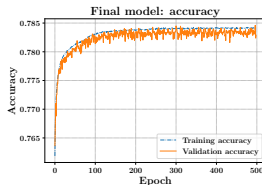
- **optimal configuration** (from characterisation) **trained and tested using the full dataset**

Feature Map		Architecture	
Type	Spherical	$\chi$	50
Order	5	Activation	Elu
Training		Dataset size	
Batch size	$10^4$	Train	$10^7$
Epochs	500	Validation	$5 \cdot 10^5$
Optimizer	Adam	Test	$5 \cdot 10^5$

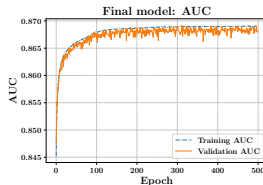
**Table:** Hyperparameters for final TTN model, from which the best results obtained in this work are extracted.



(a) Loss over training epochs



(b) Accuracy over training epochs

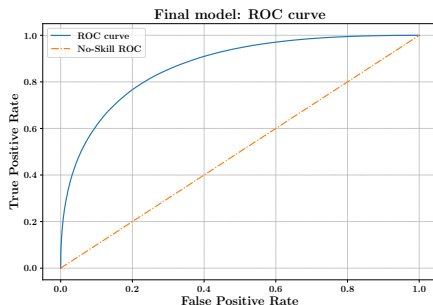


(c) AUC over training epochs

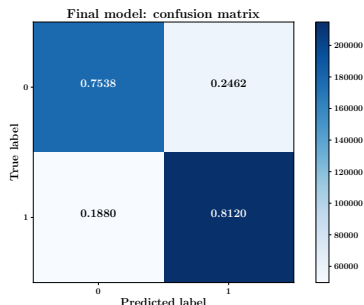
**Figure:** Final TTN model metrics computed on training and validation sets over training epochs, with the loss in **10a**, the accuracy in **10b** and the AUC in **10c**.

## Performance evaluation:

- final performance evaluation done on a test set of  $5 \cdot 10^5$  samples
- from ROC curve:  $AUC_{\text{final}} = 0.8694$
- from confusion matrix:  $\text{Accuracy}_{\text{final}} = 78.46\%$



(a) Final model ROC curve



(b) Final model confusion matrix

**Figure:** Final TTN classifier predictions on test set. In **11a**, the ROC curve of the model prediction is showed, in **11b** the confusion matrix of the predicted classes is visualised. In particular, in the ladder the colormap represents the number of samples belonging to a certain prediction class.

**HIGGS dataset from [1] taken as benchmark:**

- common machine learning dataset for HEP classification
- signal-over-background discrimination task

**We have seen how to properly preprocess the dataset for the TTN:**

- features are **rescaled**, **padded** and **mapped** through a feature map

**We have seen how to implement a TTN classifier:**

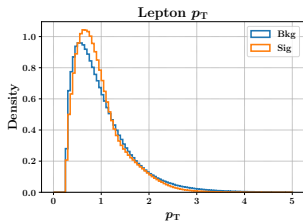
- using **TensorFlow** and **TensorNetwork**, a TTN layer has been implemented
- Keras API is used to create a full model starting from the TTN layer
- many advanced techniques also implemented, e.g. batch normalisation, regularisation, ...

**Model trained on the full training set:**

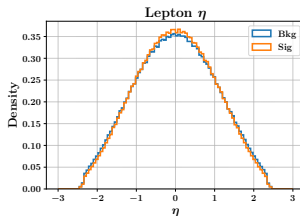
- **performance and time scaling** studied depending on:
  - bond dimension  $\chi$
  - feature map
  - batch size  $m$
- best achieved model:
  - best test accuracy of 78.46%
  - best test AUC of 0.8694



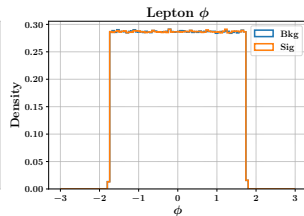
**Thank you for the  
attention!**



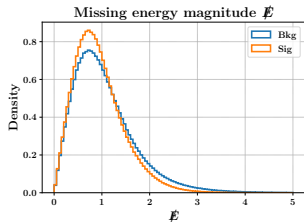
(a) Lepton  $p_T$



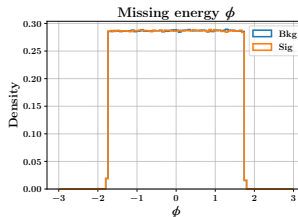
(b) Lepton  $\eta$



(c) Lepton  $\phi$

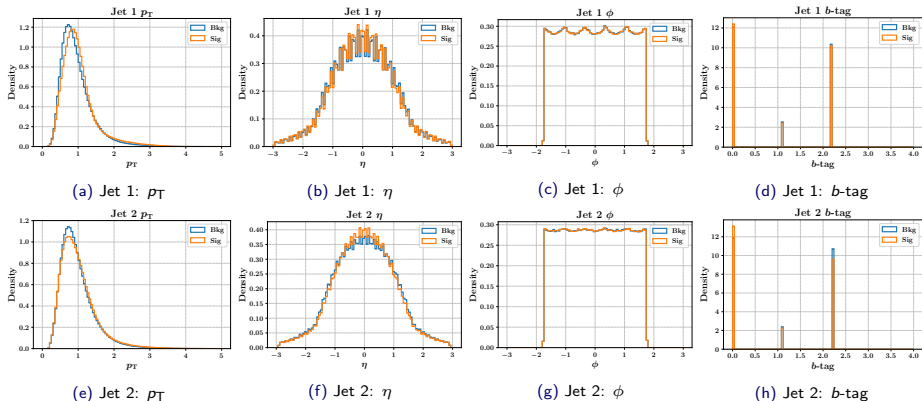


(d) Missing energy  $E$

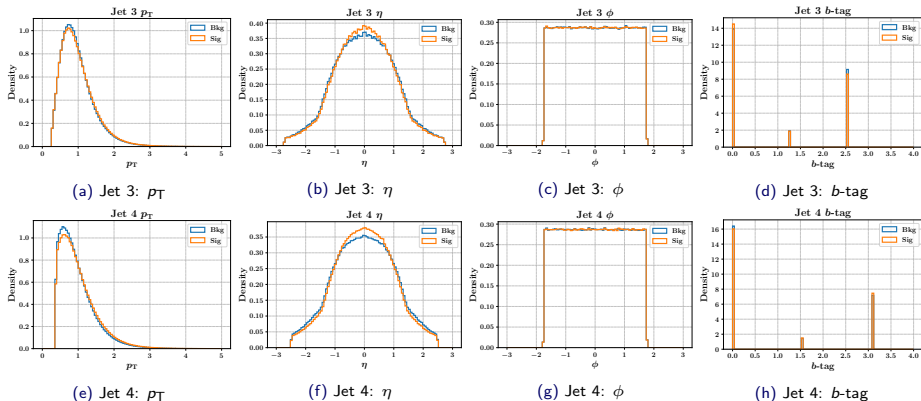


(e) Missing energy  $\phi$

**Figure:** Distributions of low-level features of leptonic part for background (blue line) and signal (orange line) events. The physical units of measure are omitted due to the fact that the dataset is not available in non-standardised form.

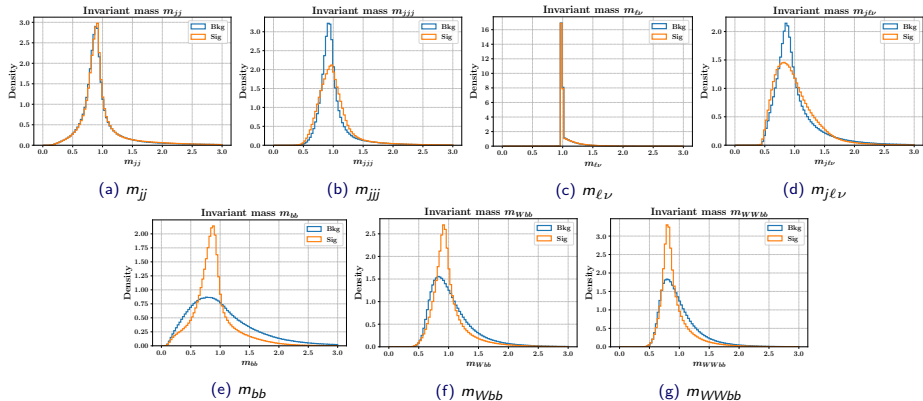


**Figure:** Distributions of low-level features of hadronic part for background (blue) and signal (orange) events. The physical units of measure are omitted due to the fact that the dataset is not available in non-standardised form.



**Figure:** Distributions of low-level features of hadronic part for background (blue) and signal (orange) events. The physical units of measure are omitted due to the fact that the dataset is not available in non-standardised form.

# Back-up: high level features



**Figure:** Distributions of high-level features for background (blue) and signal (orange) events. The physical units of measure are omitted due to the fact that the dataset is not available in non-standardised form.

## Algorithm ADAM algorithm

**Require:** Step size  $\varepsilon$  (suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1, \rho_2 \in [0, 1)$  (suggested defaults:  $\rho_1 = 0.9$ ,  $\rho_2 = 0.999$ )

**Require:** Small constant  $\delta$ , usually  $10^{-8}$ , used to stabilise division by small numbers

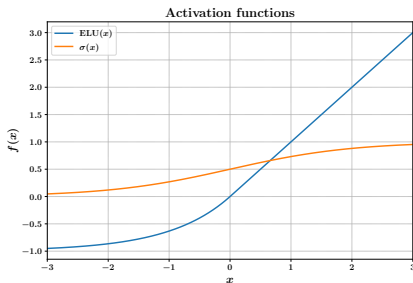
**Require:** Initial weights  $\mathbf{w}$

**Require:** Minibatch dimension  $m$

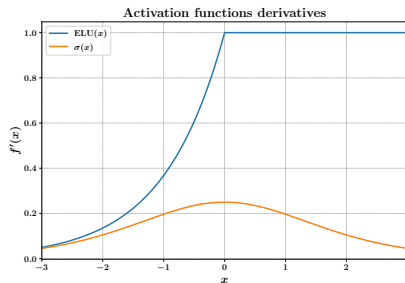
```
1: procedure ADAM( $\{x^{(i)}\}_{i=1,\dots,n}$ )                                ▷ Input: training dataset of  $n$  elements
2:   Initialize 1st and 2nd moment variables,  $s = 0$ ,  $r = 0$ 
3:   Initialize time step  $t = 0$ 
4:   while stopping criterion not met do
5:     Sample minibatch  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ 
6:      $g \leftarrow +\frac{1}{m} \nabla_{\mathbf{w}} \sum_i L(f(x^{(i)}; \mathbf{w}), y^{(i)})$                 ▷ Compute gradient estimate
7:      $t \leftarrow t + 1$ 
8:      $s \leftarrow \rho_1 s + (1 - \rho_1) g$                                     ▷ Update biased 1st moment estimate
9:      $r \leftarrow \rho_2 r + (1 - \rho_2) g \odot g$                             ▷ Update biased 2nd moment estimate
10:     $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$                                         ▷ Correct bias in 1st moment
11:     $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$                                         ▷ Correct bias in 2nd moment
12:     $\Delta \mathbf{w} = -\varepsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$                                 ▷ Compute parameter update
13:     $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$ 
14:  end while
15:  return  $\mathbf{w}$ 
16: end procedure
```

$$\text{ELU}(x) = \begin{cases} x & x \geq 0 \\ a(e^x - 1) & x < 0 \end{cases}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



(a) Activation functions



(b) Derivatives of activation functions

**Figure:** Comparison between ELU and sigmoid activation functions in **16a**, and between their derivatives, in **16b**, which are used also for loss function gradient calculation.

```
# prepare input data for the vectorization of the contract function
input_shape = list(inputs.shape)
inputs      = tf.reshape(inputs, (-1, input_shape[1], input_shape[2]))
# vectorize the contraction over all the input samples
result = tf.vectorized_map( # vectorize
    lambda vec: contract(   # create a lambda function to be vectorized
        vec                 , # input sample
        self.nodes          , # weight tensors
        self.n_contraction  , # number of feaqtue to contract
        self.use_bias       ,
        self.bias_var       , # bias tensor
    ),
    inputs                  # input dataset over which vectorize the lambda function
)
```

```
# apply activation, if specified
if self.activation is not None:
    result = self.activation(result)
```