# Optimizing Tree Tensor Networks for classification on hardware accelerators
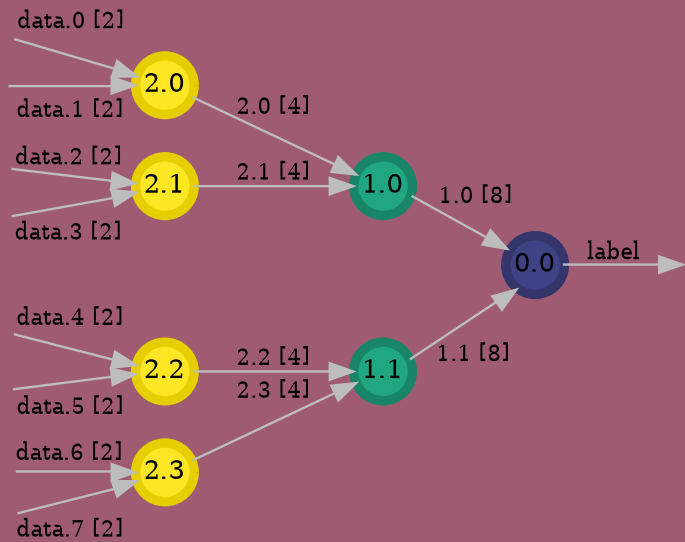
**PhD Student**

*Alberto Coppi*

*Workshop on Tensor Networks and*
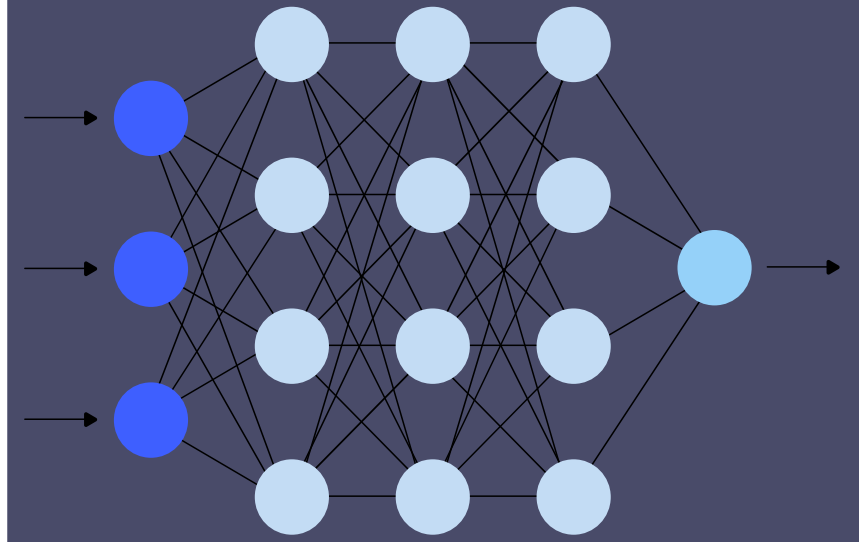*(Quantum) Machine Learning for High-Energy Physics*
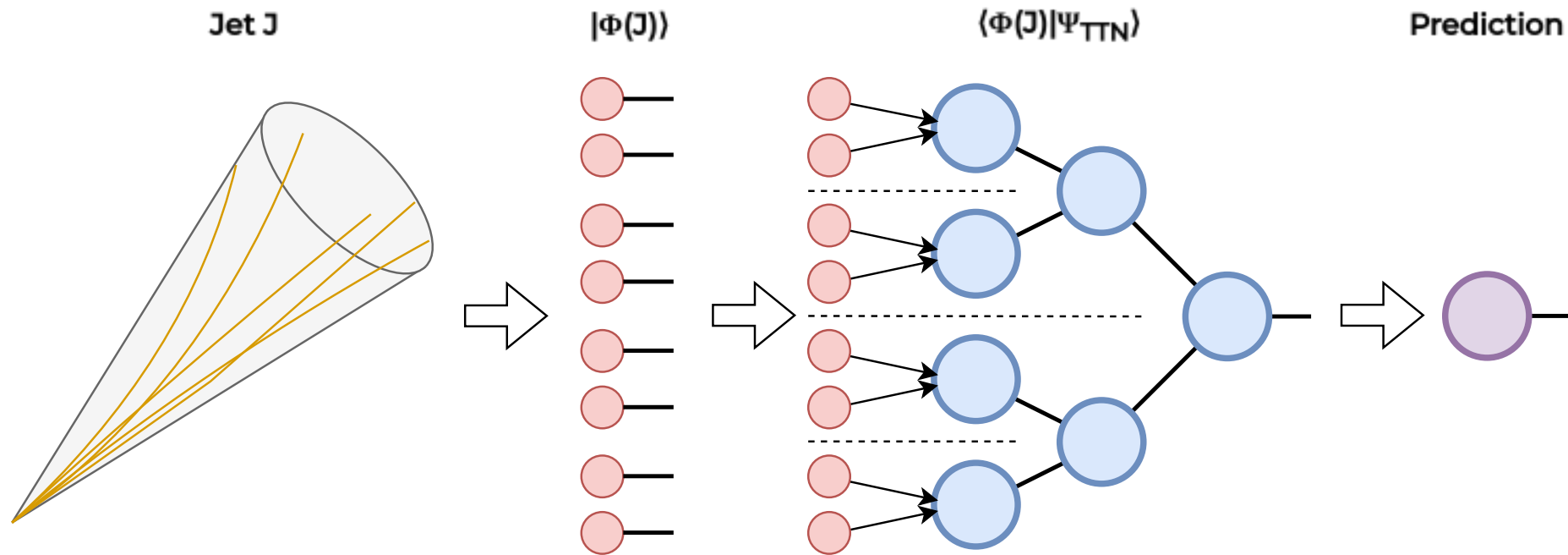
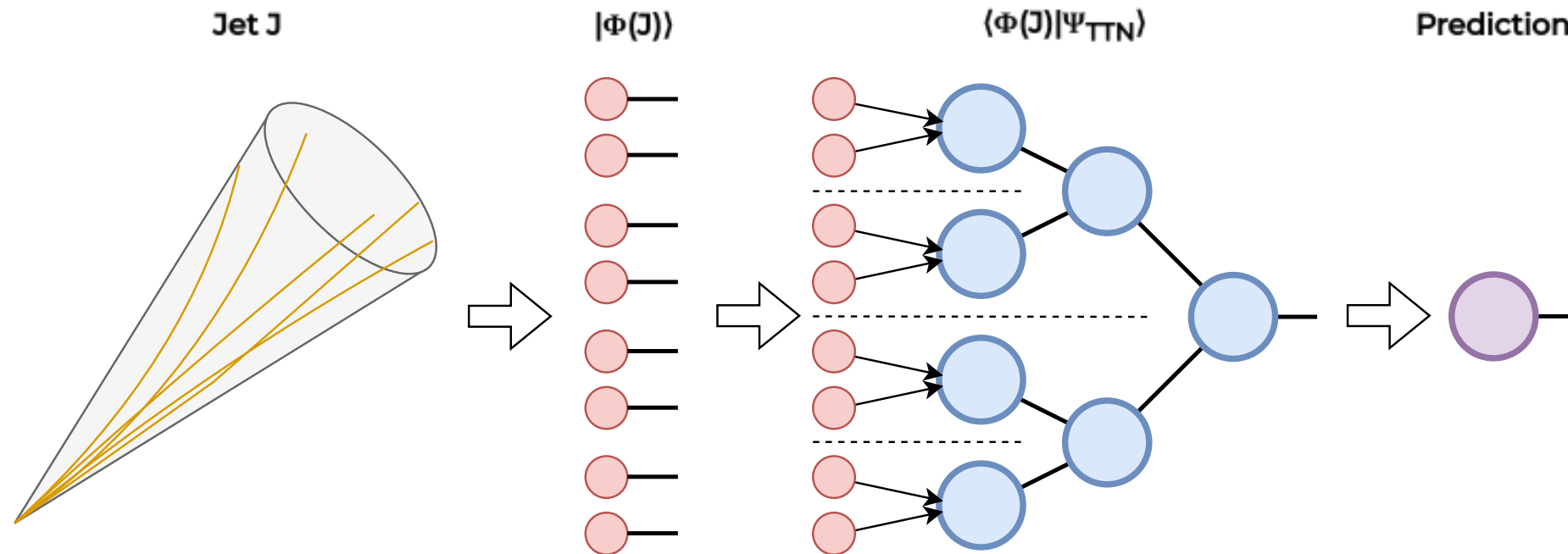*Date: 05/11/2024*

# Tree Tensor Networks

# Tree Tensor Networks classifiers

# Tree Tensor Networks classifiers



**Jet J**

$|\Phi(J)\rangle$

$\langle\Phi(J)|\Psi_{TTN}\rangle$

**Prediction**

**Data**

- Different fields
  (HEP, OD)
- High rate

3

# Tree Tensor Networks classifiers

# Tree Tensor Networks classifiers

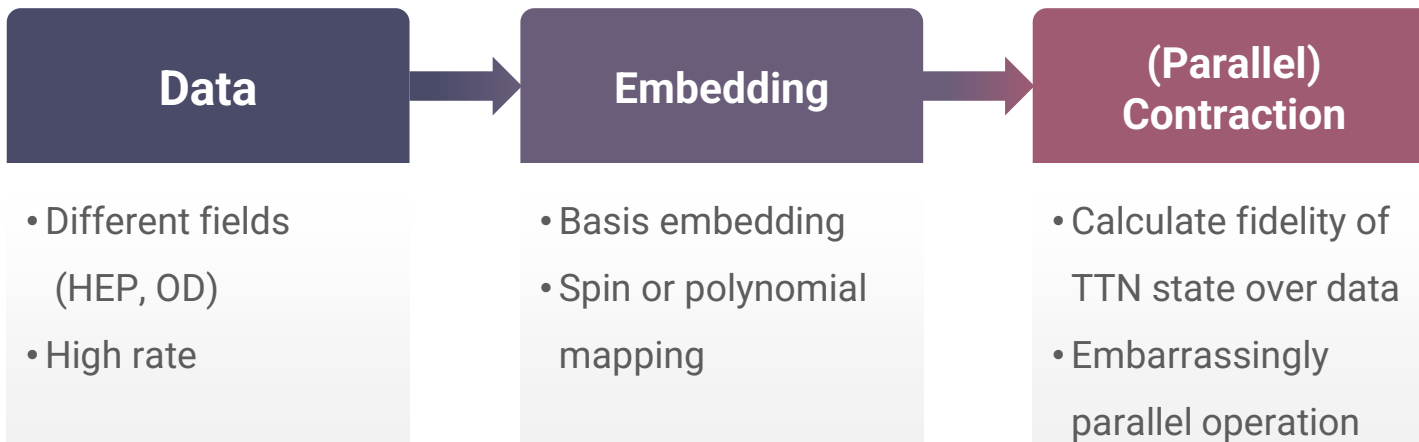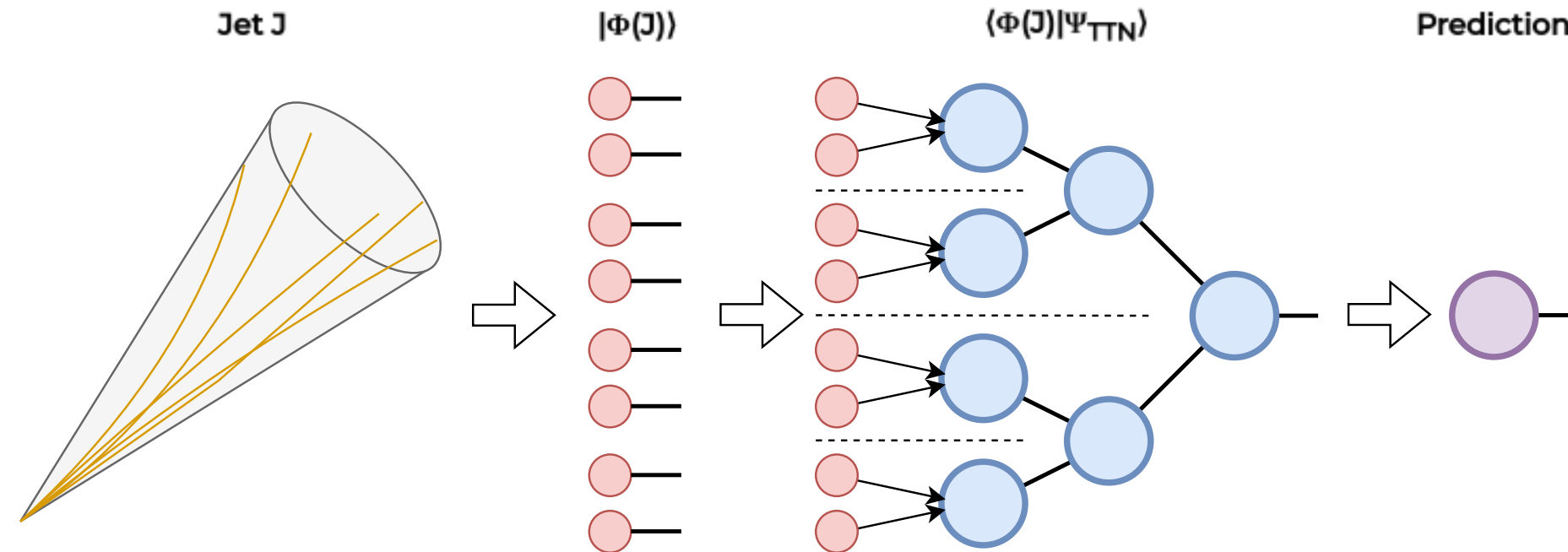# Tree Tensor Networks classifiers



| **Data** | | **Embedding** | | **(Parallel) Contraction** | | **Prediction** |
|---|---|---|---|---|---|---|
| • Different fields (HEP, OD)<br>• High rate | | • Basis embedding<br>• Spin or polynomial mapping | | • Calculate fidelity of TTN state over data<br>• Embarrassingly parallel operation | | • Scalar or vector out<br>• Plug into loss function for optimisation |

# Tree Tensor Networks classifiers:
## Why?

### Software

**Training**
- Initialisation
- Optimisation:
  - Global SGD
  - Sweeping

**Explainability**
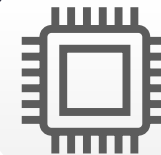- Measurement of physical quantities

### Hardware

**FPGA**
- Deterministic latency
- Limited resources

**Use case**
- **Jet tagging** for HEP experiments, e.g. CMS
- Currently done **offline** by complex ML models like ParticleNet on **already filtered data**
- Objective: deploy a **tagger online** in the L1 trigger of CMS experiment to improve selection efficiency

# Tree Tensor Networks classifiers:
## Why?

### Software

**Training**
- Initialisation
- Optimisation:
  - Global SGD
  - Sweeping

**Explainability**
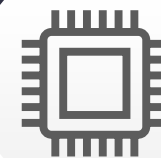- Measurement of physical quantities

### FPGA

**Explainability**
- The model can be "explained" through physical measurements
- **No black-boxes** for filtering out ~98% of data

**Compressibility**
- The number of parameters can be reduced, post-learning
- **Fit** the model **to limited resources** hardware

**Speed**
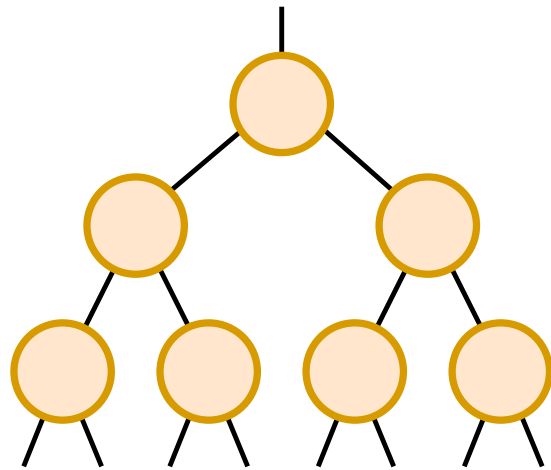- Based on simple operations and **parallelizable**
- **Compliance with latency limits**

# Tree Tensor Networks classifiers:
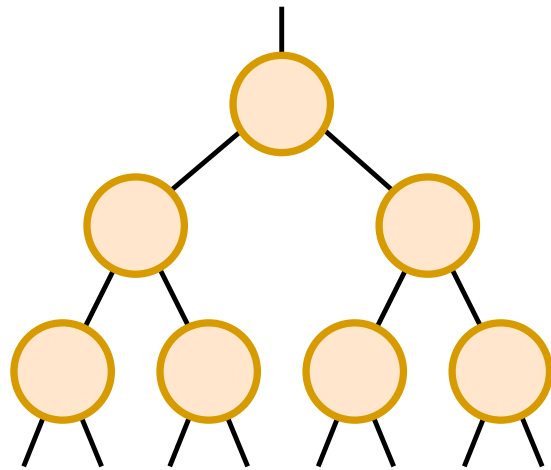## Training

**Initialisation**

**Optimisation**

# Tree Tensor Networks classifiers:
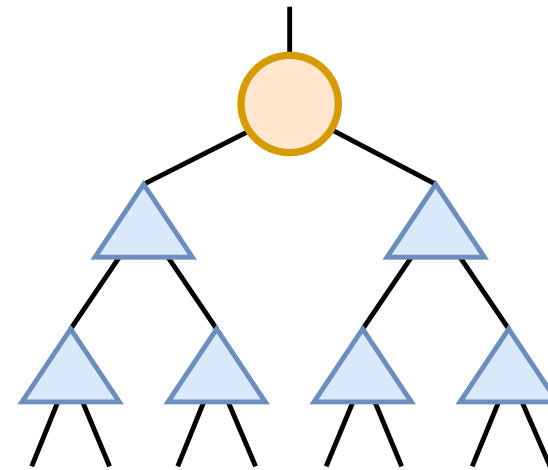## Initialisation



Randomly initialised
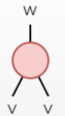
# Tree Tensor Networks classifiers:
## Initialisation



Randomly initialised

Each layer progressively project data into a lower dimensional space[1]

[1] E Miles Stoudenmire, Learning relevant features of data with multi-scale tensor networks, *Quantum Sci. Technol.* **3** 034003 (2018). https://doi.org/10.1088/2058-9565/aaba1a

# Tree Tensor Networks classifiers:
## Optimisation

### Global SGD

**Pros:**

- Treat TTN parameters as classical NN parameters

- Exploit **gradient tracking automation** of PyTorch

**Cons:**

- Suffers from **barren plateaus**

**VS**

### Sweeping

**Pros:**

- Decompose a large problem in many **smaller problems**

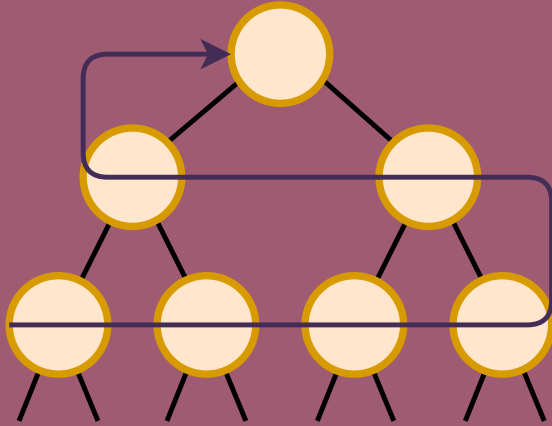- More **stable and robust**, enabling training of larger models

**Cons:**

- Must be **manually implemented**, thus can be **slower**

# Tree Tensor Networks classifiers:
## Optimisation



**Sweeping**

**VS**

**Pros:**

- Decompose a large problem in many **smaller problems**

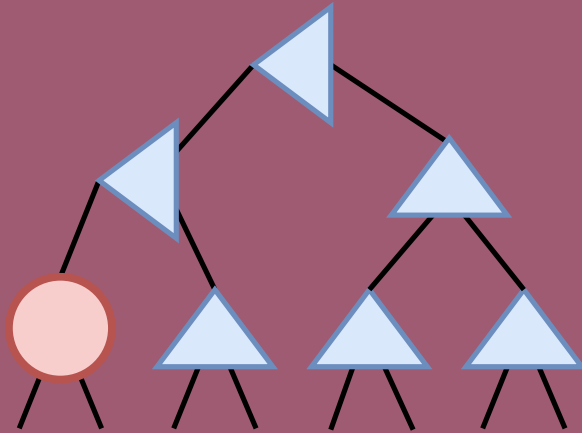- More **stable and robust**, enabling training of larger models

**Cons:**

- Must be **manually implemented**, thus can be **slower**

9

# Tree Tensor Networks classifiers:
## Optimisation



**Sweeping**

**VS**

**Pros:**

- Decompose a large problem in many **smaller problems**

- More **stable and robust**, enabling training of larger models

**Cons:**

- Must be **manually implemented**, thus can be **slower**

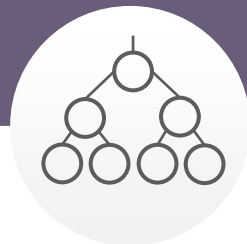# Tree Tensor Networks classifiers:
## **Optimisation**

### Global SGD

**Pros:**

- Treat TTN parameters as classical NN parameters

- Exploit **gradient tracking automation** of PyTorch
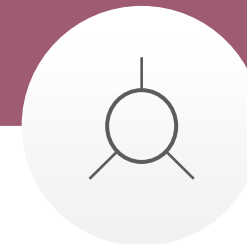
**Cons:**

- Suffers from **barren plateaus**

**VS**

### Sweeping

**Pros:**

- Decompose a large problem in many **smaller problems**

- More **stable and robust**, enabling training of larger models

**Cons:**

- Must be **manually implemented**, thus can be **slower**

# Software implementation

```python
1    class TIndex:
2        def __init__(self,
3                      name: str,
4                      inds: Sequence[str] |
5                          np.ndarray):
6            self.__name = name
7            self.__tindices = np.array(inds,
8                                  dtype=np.str_)
9            self.__ndims = len(inds)
```

```python
1    class TTN:
2        def __init__(
3                self,
4                n_features,
5                n_phys=2,
6                n_labels=2,
7                label_tag="label",
8                bond_dim=4,
9                dtype=torch.double,
10               device="cpu",
11               quantizer = None
12       ):
```

```python
1    class TTNModel(torch.nn.Module, TTN):
2        def __init__(
3                self,
4                n_features,
5                n_phys=2,
6                n_labels=2,
7                label_tag="label",
8                bond_dim=8,
9                dtype=torch.double,
10               device="cpu",
11               quantizer = None
12       ):
```

## Characteristics
- **Open-source**
- Developed **from scratch**
- Based on **PyTorch** 🔥
- **Enables to train and explain a TTN ML model**

## Structure

### TIndex
- Base class for **tensor indexing**
- Can be used as **key** in dictionaries

### TTN
- Fundamental class to construct a **full TTN**
- Equipped with methods to enable **TTN functionality**: **contraction,** initialization, derivative, **expectation value**, entanglement **entropy**, drawing
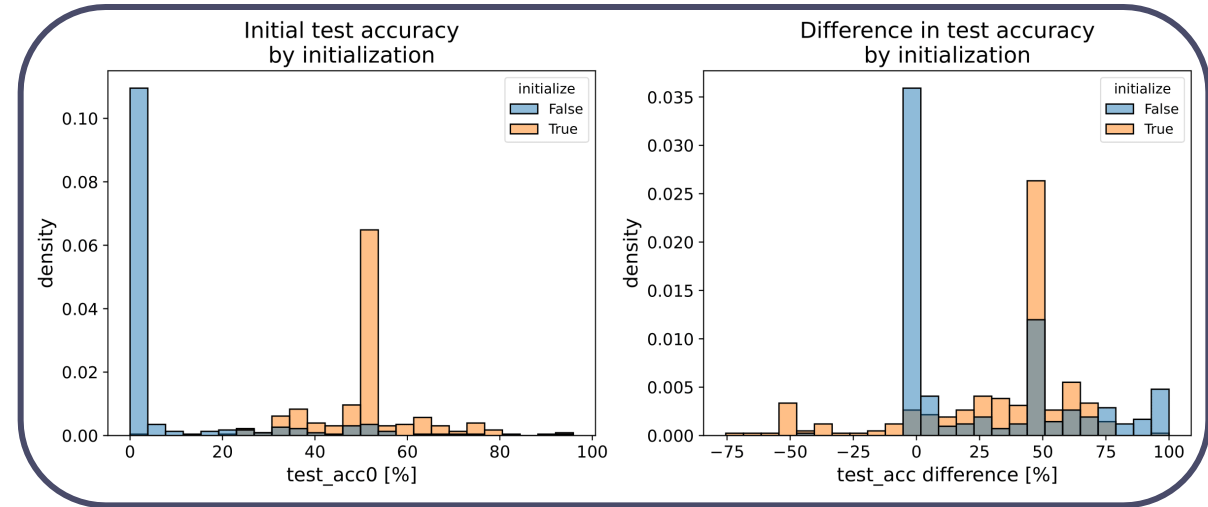
### TTNModel
- Derived class of both `TTN` and `torch.nn.Model`
- This provides a ML based approach to optimise the TTN, with (an almost free) **SGD**
- Implements the **Sweeping** optimisation algorithm
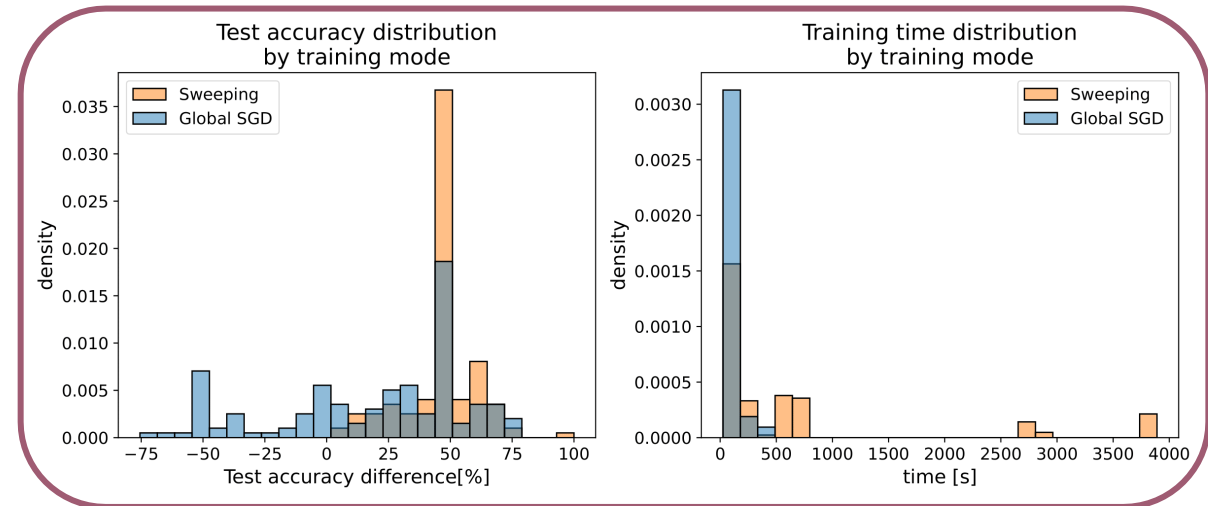
# Tests on hyperparameters

**Initialisation on/off:**

- The initialisation procedure moves the model towards optimum
- From there, the **training** procedure **is facilitated** to find the optimum
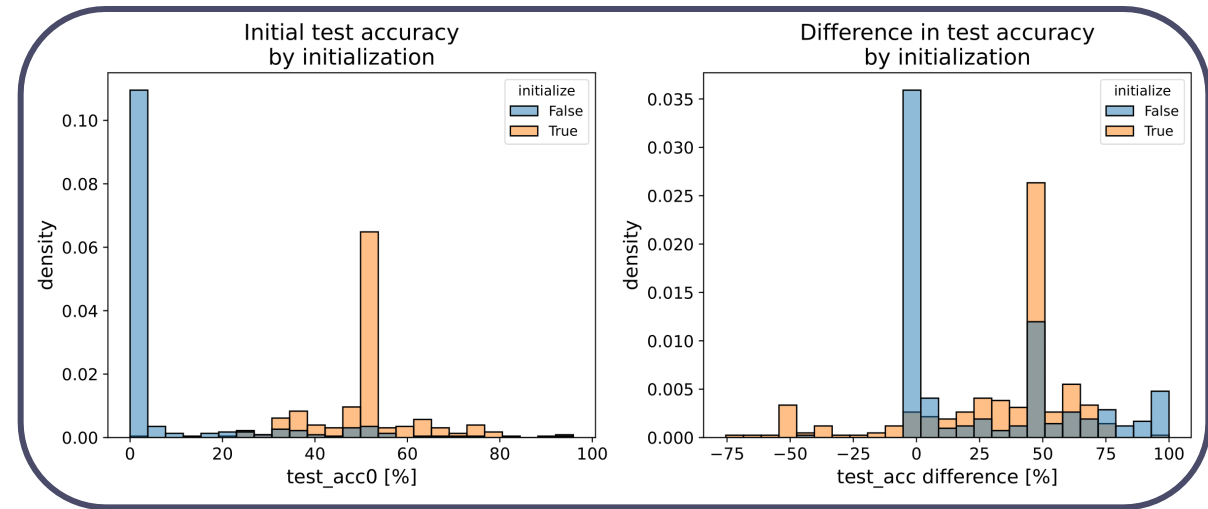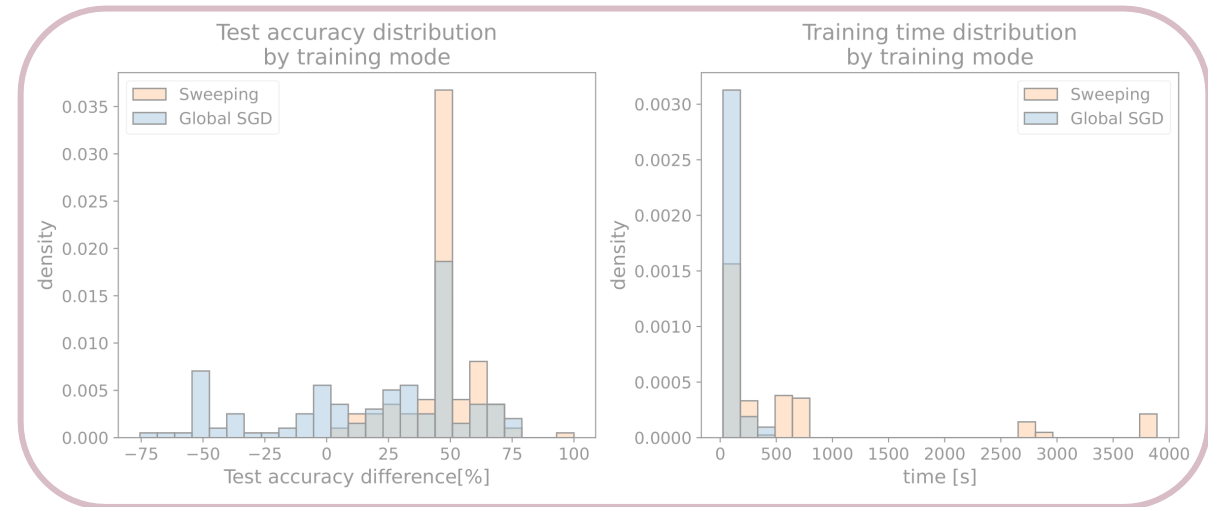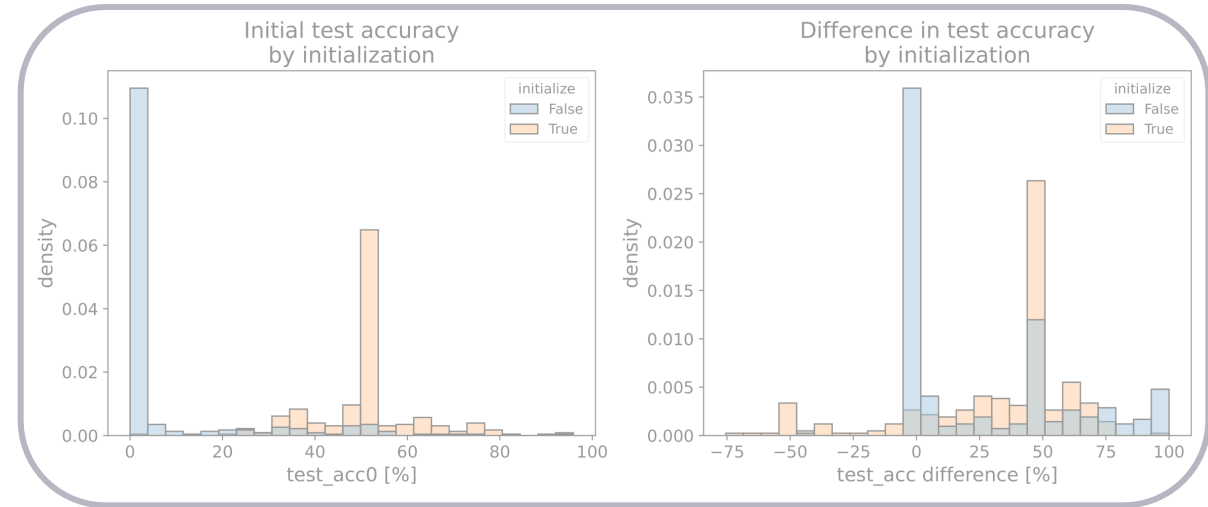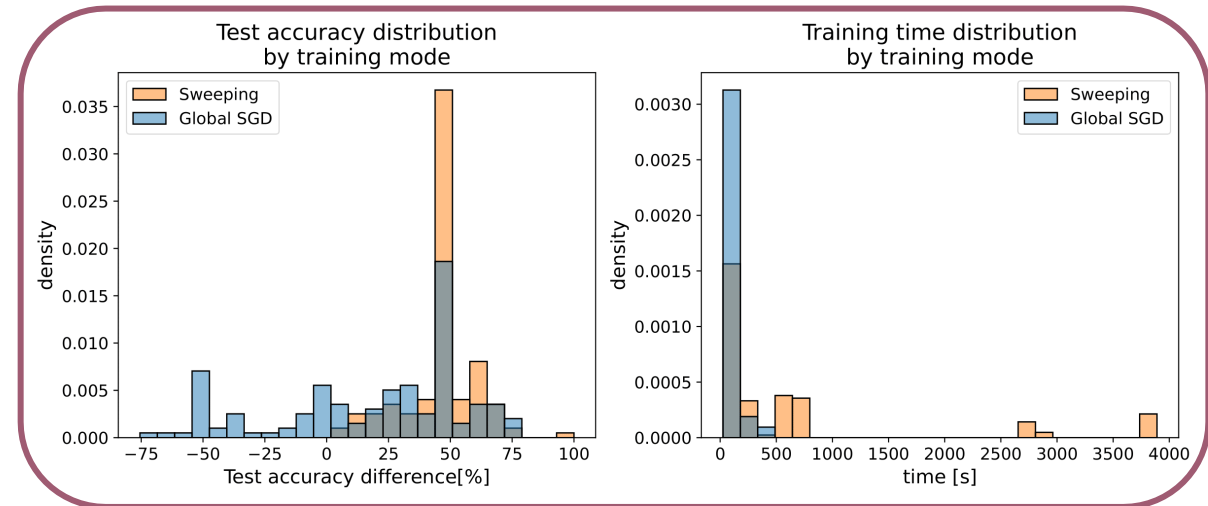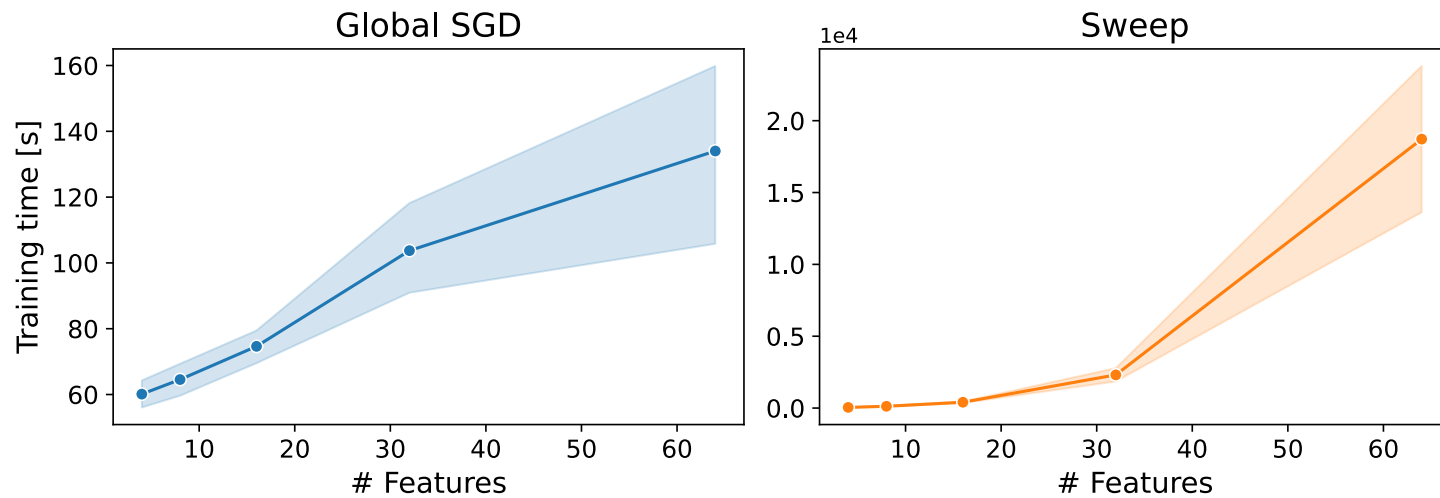
**Global SGD vs. Sweeping:**

- **Sweeping is more stable** and robust, reaching the optimum more frequently, **but slower**
- Global SGD doesn't work for large models

# Tests on hyperparameters

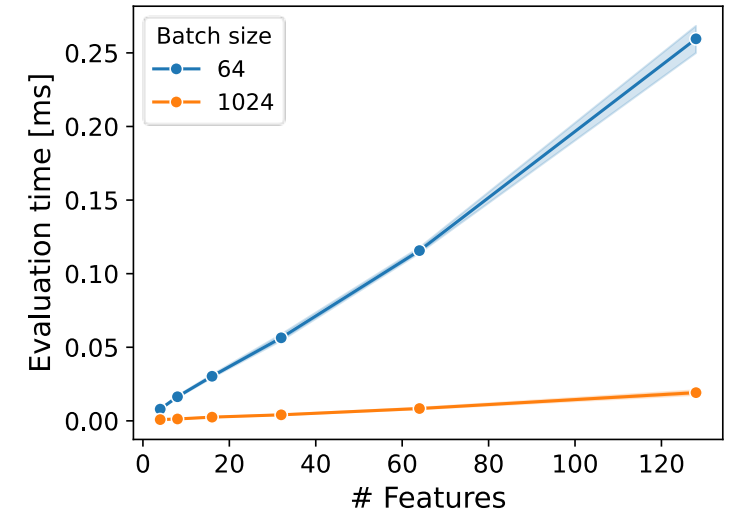**Initialisation on/off:**

- The initialisation procedure moves the model towards optimum
- From there, the **training** procedure **is facilitated** to find the optimum

**Global SGD vs. Sweeping:**

- **Sweeping is more stable** and robust, reaching the optimum more frequently, **but slower**
- Global SGD doesn't work for large models

# Tests on hyperparameters

**Initialisation on/off:**

- The initialisation procedure moves the model towards optimum
- From there, the **training** procedure **is facilitated** to find the optimum

**Global SGD vs. Sweeping:**

- **Sweeping is more stable** and robust, reaching the optimum more frequently, **but slower**
- Global SGD doesn't work for large models

# Tests on hyperparameters:
## Performance

### Training time vs. # Features

#### Global SGD



#### Sweep



### Evaluation time vs. # Features



- Approximately linear dependence (thanks PyTorch)

- Expected at least quadratic dependence (the number of nodes is $\mathcal{O}(N^2)$)
- Actually worse

- Roughly linear dependence (NVIDIA GTX 1050)
- Greater batch size -> better exploitation of parallel computing capabilities (thanks PyTorch)

# Tests on synthetic datasets

## Iris

Commonly used in ML for benchmarking purposes

**Characteristics**

- 4 features

- 150 samples

- 3 classes:
  Iris-setosa, Iris-versicolor, Iris-virginica

## Titanic

Commonly used in ML for benchmarking purposes

**Characteristics**

- 8 features (originally 13)

- 1043 samples (originally 1309)

- 2 classes:
  (not) survived

## LHCb

Dataset coming from LHCb open data, already used in literature to test TTNs in ML [2]

**Characteristics**

- 16 features

- $\sim 1.1 \times 10^6$ samples

- 2 classes:
  $b$ vs $\overline{b}$

[2] Felser, T., Trenti, M., Sestini, L. *et al.* Quantum-inspired machine learning on high-energy physics data. *npj Quantum Inf* **7**, 111 (2021). https://doi.org/10.1038/s41534-021-00443-w

# Tests on synthetic datasets

| | Iris | Titanic | LHCb |
|---|---|---|---|
| **Accuracy (Training/Test)** | 99,1% / 96,7% | 80,9% / 77,0% | 61,7% / 61,8% |
| **AUC** | 1,0 | 0,83 | 0,66 |
| **Accuracy (smaller)** | 96,7% / 96,7% | 79,3% / 74,1% | 60,0% / 60,3% |
| **AUC (smaller)** | 0,99 | 0,84 | 0,63 |



Predictions distribution titanic, BD=8



Predictions distribution bbdata, BD=10

# Tests on synthetic datasets:
## Explainability



Two-sites $\sigma_z$ correlation



Entanglement entropy

# Tests on synthetic datasets:
## Explainability



$\sigma_z$ correlations between features

Two-sites $\sigma_z$ correlation

# Tests on synthetic datasets:
## Explainability



Entanglement entropy

# Tests on synthetic datasets:
## Explainability



Entanglement entropy

# Tests on synthetic datasets:
## Explainability



$\sigma_z$ correlations between features

Two-sites $\sigma_z$ correlation



Titanic dataset features entropy

Entanglement entropy

# Tests on synthetic datasets:
## Explainability



$\sigma_z$ correlations between features

Two-sites $\sigma_z$ correlation

Titanic dataset features entropy

Entanglement entropy

# Tests on synthetic datasets:
## Explainability



$\sigma_z$ correlations between features

Two-sites $\sigma_z$ correlation

Titanic dataset features entropy

Entanglement entropy

# Tests on synthetic datasets:
## Explainability



$\sigma_z$ correlations between features

Two-sites $\sigma_z$ correlation



Titanic dataset features entropy

Entanglement entropy

# Tests on synthetic datasets:
## Towards hardware

- Different hardware implement **different logic**
- **Simulate model behaviour** in different numerical representations
- On FPGA, **APFP** representation

- Studied model performance **dependence on the number of bits** used, by means of QPyTorch
- Further compression



Test accuracy study for fixed precision, striped



Test accuracy study for fixed precision, titanic

# Future perspectives

**Hardware optimisation**

Optimise these algorithms exploiting TN properties for faster execution on specialized hardware.

**Online jet tagging**

Deploy in the online selection algorithm a b-tagger.

**TTN based emulation and ML**

Implement the TTN ansatz both in the Quantum MATCHA and CHAI TEA applications to improve performance and accuracy of quantum computer emulations.

**Extension to other ansatzes**

TTN may be not powerful enough for some specific tasks. Other ansatzes like MERA are worth exploring.

# Hardware optimisation



FPGA



GPU



Gauge freedom

**What**

- Exploit the characteristics of different hardware for faster execution of computations

- Parallelism, different numerical representations (APFP, TF32)

**How**

- Compressibility of TNs

- Gauge freedom

**Why**

- Enable fast, low-resource simulations, accelerating research into quantum algorithms and making it widely feasible.

# TTN based emulation



## MPS based emulation

**Pros**

- **Well developed** technology, integrated in popular emulation libraries (Quantum MATCHA TEA, qiskit)
- **Straightforward** application of one and two qubits gates

**Cons**

- Supports only **linear topology** circuits
- Difficult representation of **long distance** interactions

# TTN based emulation



**TTN based emulation**

**Pros**

- Can handle **longer distance** interactions

- **Non-linear** circuits

P. Seitz, I. Medina, E. Cruz, Q. Huang, and C. B. Mendl, *Simulating quantum circuits using tree tensor networks,* Quantum 7, 964 (2023).

# TTN based emulation



## TTN based emulation

**Pros**

- Can handle **longer distance** interactions

- **Non-linear** circuits

P. Seitz, I. Medina, E. Cruz, Q. Huang, and C. B. Mendl, *Simulating quantum circuits using tree tensor networks,* Quantum 7, 964 (2023).

# TTN based emulation



## TTN based emulation

**Pros**

- Can handle **longer distance** interactions

- **Non-linear** circuits

**Cons**

- **More complex** operations

- Still difficulty for circuits with **high connectivity**

P. Seitz, I. Medina, E. Cruz, Q. Huang, and C. B. Mendl, *Simulating quantum circuits using tree tensor networks,* Quantum 7, 964 (2023).

# TTN based emulation



## TTN based emulation

**Pros**

- Can handle **longer distance** interactions
- **Non-linear** circuits

**Cons**

- **More complex** operations
- Still difficulty for circuits with **high connectivity**

**Tests**

- Computation **time** and **resource** usage
- Final state **accuracy** for common quantum algorithms

Figures from: P. Seitz, I. Medina, E. Cruz, Q. Huang, and C. B. Mendl, *Simulating quantum circuits using tree tensor networks,* Quantum 7, 964 (2023).

# TTN based ML



**Quantum-Enhanced ML**

**What**

- Quantum computer as part of a **hybrid quantum-classical** model

- Classical part optimizes **classical parameters** controlling the quantum part

- ➤ Expand CHAI TEA to use MATCHA TEA for hybrid quantum-classical ML. Expose as PyTorch layers

**Why**

- Offer a complete framework

- Competitive performances w.r.t. Tensorflow Quantum, possibly improving on TorchQauntum

Top figure from: D. Peral-García, J. Cruz-Benito, F. J. García-Peñalvo, *Systematic literature review: Quantum machine learning and its applications*, Computer Science Review, Volume 51, 2024, 100619, ISSN 1574-0137.
Bottom figure from: Tensorflow Quantum.

Thank you for your attention

BACKUP

Software studies

# Sweeping analysis



Difference in test accuracy
after and before training

# Barren plateaus



Median gradient magnitude
for different number of features

# Tensors

# Tensors



**Multi-dimensional arrays**

- Generalisation of the idea of vectors and matrices.
- Connected to software representation.

Scalar    Vector    Matrix

Order-$N$ tensor

**Multi-linear maps**

- Extending the idea of linear maps (matrices).
- Useful to interpret some TNs algorithms.

W

V    V

**Elements of tensor product space**

- Define properties such as order and shape.
- Connected to quantum many-body systems:

$$|\Psi\rangle = \sum_{\{i\}^N} \Psi_{\{i\}^N} |i_1\rangle \otimes \cdots \otimes |i_N\rangle$$

T

# Tensor operations



**Manipulations**

**Index fusing & splitting**

**Contraction**

**Decompositions**

**SVD / eigenvalue**

**QR**

# Isometrisation

# Isometrisation

# Differentiation

# Derivative

$$\frac{\partial}{\partial T^{[k]}} \left[ \text{tree with } T^{[k]} \right] = \text{tree}$$

# Quantum TEA

# Quantum TEA library

# Quantum TEA library



**Quantum RED TEA**

Provide the interfaces to BLAS/LAPACK and CUDA for the higher-level applications.

# Quantum TEA library



**Quantum RED TEA**

Provide the interfaces to BLAS/LAPACK and CUDA for the higher-level applications.

**Quantum TEA LEAVES**

Python solutions for common TN geometries, ground state search algorithms, time evolution via TDVP, Python-FORTRAN interfaces.

# Quantum TEA library



**Quantum RED TEA**

Provide the interfaces to BLAS/LAPACK and CUDA for the higher-level applications.

**Quantum TEA LEAVES**

Python solutions for common TN geometries, ground state search algorithms, time evolution via TDVP, Python-FORTRAN interfaces.

**Quantum MATCHA TEA**

Quantum computer emulator powered by **matrix product states**

# Quantum TEA library



**Quantum RED TEA**

Provide the interfaces to BLAS/LAPACK and CUDA for the higher-level applications.

**Quantum TEA LEAVES**

Python solutions for common TN geometries, ground state search algorithms, time evolution via TDVP, Python-FORTRAN interfaces.

**Quantum MATCHA TEA**

Quantum computer emulator powered by **matrix product states**

**Quantum GREEN TEA**

Solves the static and time-dependent Schrödinger equation and Lindblad equation

# Quantum TEA library



**Quantum RED TEA**

Provide the interfaces to BLAS/LAPACK and CUDA for the higher-level applications.

**Quantum TEA LEAVES**

Python solutions for common TN geometries, ground state search algorithms, time evolution via TDVP, Python-FORTRAN interfaces.

**Quantum MATCHA TEA**

Quantum computer emulator powered by **matrix product states**

**Quantum GREEN TEA**

Solves the static and time-dependent Schrödinger equation and Lindblad equation

**Quantum CHAI TEA**

Contains the machine learning applications using TNs.

# Quantum TEA library



**Quantum MATCHA TEA**

Quantum computer emulator powered by **matrix product states**

**Quantum CHAI TEA**

Contains the machine learning applications using TNs.

# Quantum TEA library



**Quantum MATCHA TEA**

Quantum computer emulator powered by **matrix product states**

**Quantum CHAI TEA**

Contains the machine learning applications using TNs.

# TN ansatzes

# Extend to other ansatzes



PEPS

MERA

## What

- Explore other, more complex ansatez like PEPS and MERA

- Integrate them in the Quantum TEA framework making the switch between different ansatzes effortless

## Why

- Different geometries are guaranteed to capture higher amount of entropy and longer-range interactions

- Thought for non-linear systems → can stand higher connectivity

# Results on synthetic datasets

# Titanic

# Titanic

spin

poly



Predictions distribution
titanic, BD=8

# Titanic



Train Confusion Matrix
titanic, BD=8

spin

poly

Train Confusion Matrix
titanic, BD=4

# Striped images



Training Loss
stripe, BD=8



ROC Curve
stripe, BD=8

# Striped images



Predictions distribution
stripe, BD=8



Test Confusion Matrix
stripe, BD=8

# Striped images



$\sigma_z$ correlations between features

# Iris

# Iris

# Iris



$\sigma_z$ correlations between features

# Iris

LHCb

# LHCb



Training Loss
bbdata, BD=10

1152.581

ROC Curve
bbdata, BD=10

AUC = 0.66

# LHCb

# LHCb
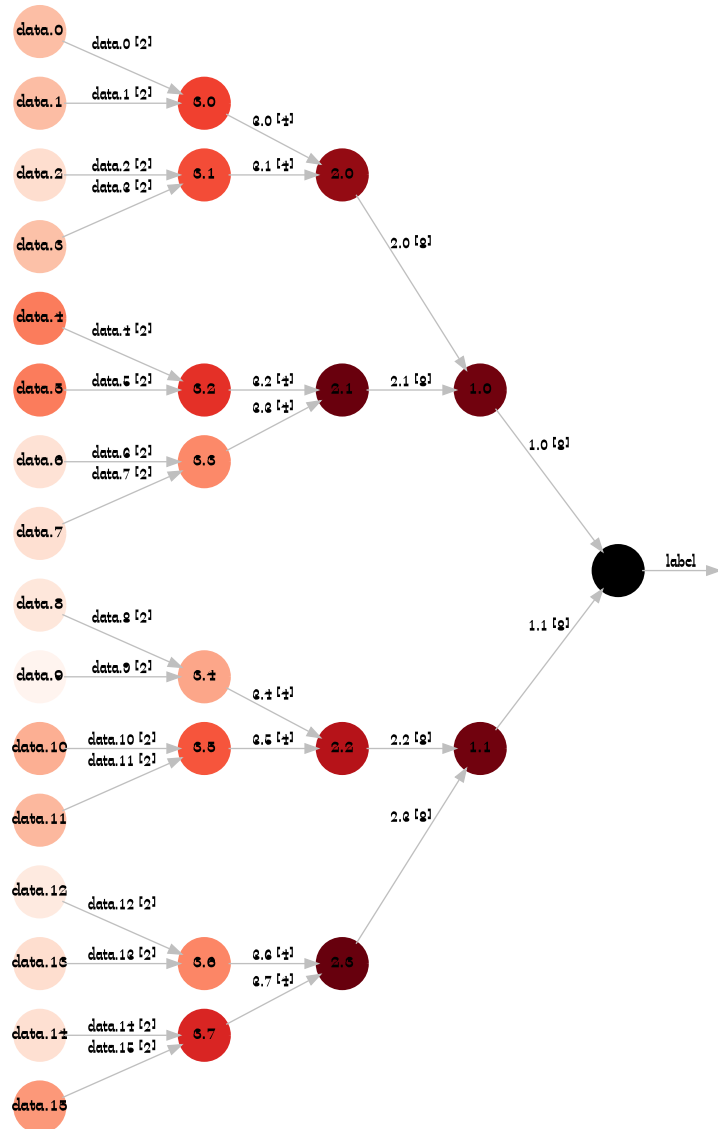


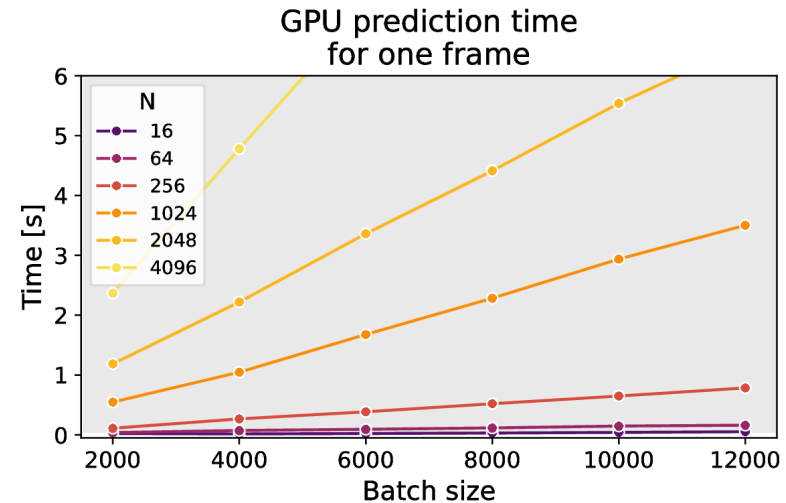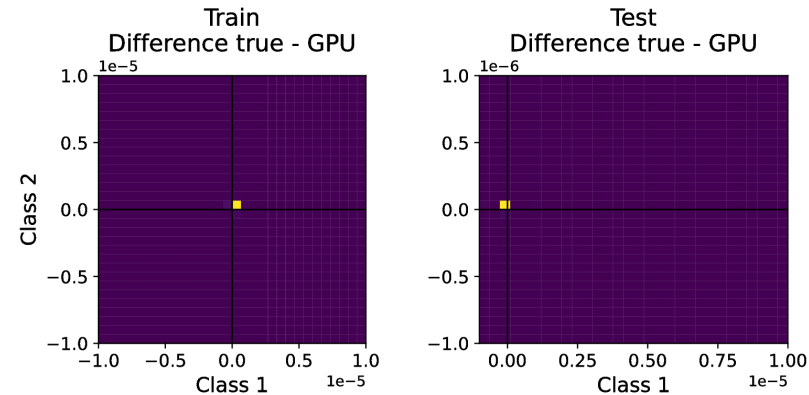$\sigma_z$ correlations between features

# LHCb

GPU

# Classifiers on hardware accelerators:
## GPU

- GPU predictor implemented as part of **internship** for Tensor AI Solutions

- Based on **cuTENSOR**, a CUDA library for tensor contraction on NVIDIA GPUs

- **Tested** on the **trained models** mentioned before

- Tested on **FSOCO dataset** for traffic cones detection

### Results:

➤ Perfect **match between** software and hardware **outputs**

➤ **Partial compliance** with video frame-rates

# Classifiers on hardware accelerators:

## GPU

- Trained a model on **FSOCO dataset**, containing high-quality camera images with cones delimited by bounding boxes

- Traffic cones detection performed through **sliding windows** technique

- **Features entropy** explain what the model learned

- Model shows **promising results** in object identification, but further refinements are needed





| Original image | Class blue | Class yellow | Class large_orange | Class orange |