# Quantum inspired ML on FPGA

Marco Zanetti

University of Padova and INFN

BOOSTlab

## In the following there won't be

- Answers about what can Quantum Computers do for HEP[1]
- A review of Quantum inspire algorithms for HEP pipelines
- Guidelines on how to use FPGAs to control/simulate quantum hardware

## This is only about

- Tensor Networks are "quantum" tools that can be used for classification/selection and anomaly detection
- Their features make them promising models to be deployed on resurces-limited contexts, notably first Trigger layers
- A case study of the deployment of TN as classifier for HEP tasks is reviewed

[1] See e.g. A. Di meglio et al. «*Quantum Computing for High-Energy Physics: State of the Art and Challenges»,* PRX QUANTUM 5, 037001 (2024)

# Tensor Networs and their usage

- TN: collections of tensors with indeces contracted in specific patterns
  - convenient graphical representation
- Rapresent/solve many-body quantum entangled states, factorizing rank-N tensors into smaller tensors
  - allows linear scaling (vs exp.) on the number of sites
  - allows computation, by reducing number of parameters and algo complexity
- Expressivity of the TN tuned by the bond dimenson
  - the dimension of the index connecting one tensor to the next
- Several topologies are suitable for various tasks
  - Tree TN: the most general loopless architecture



$$\sum_j M_{ij} N_{jkl}$$

$$T^{s_1 s_2 s_3 s_4 s_5 s_6} = \sum_{\{\alpha\}} A_{\alpha_1}^{s_1} A_{\alpha_1 \alpha_2}^{s_2} A_{\alpha_2 \alpha_3}^{s_3} A_{\alpha_3 \alpha_4}^{s_4} A_{\alpha_4 \alpha_5}^{s_5} A_{\alpha_5}^{s_6}$$
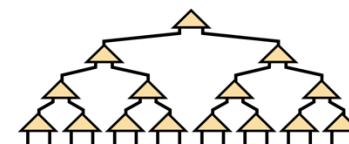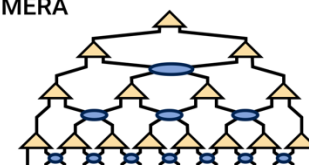
Matrix Product State / Tensor Train

PEPS

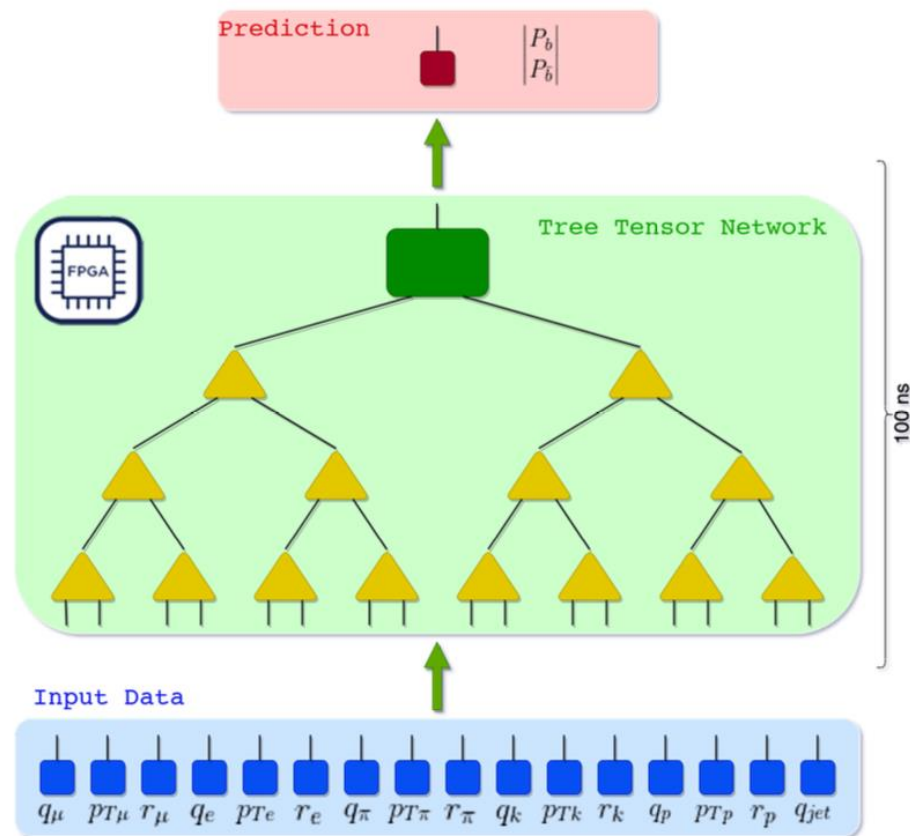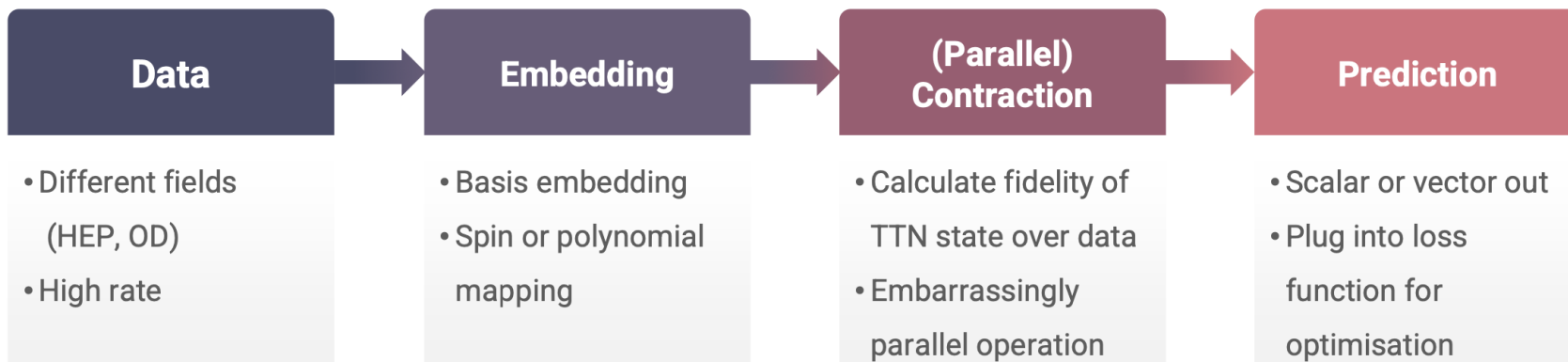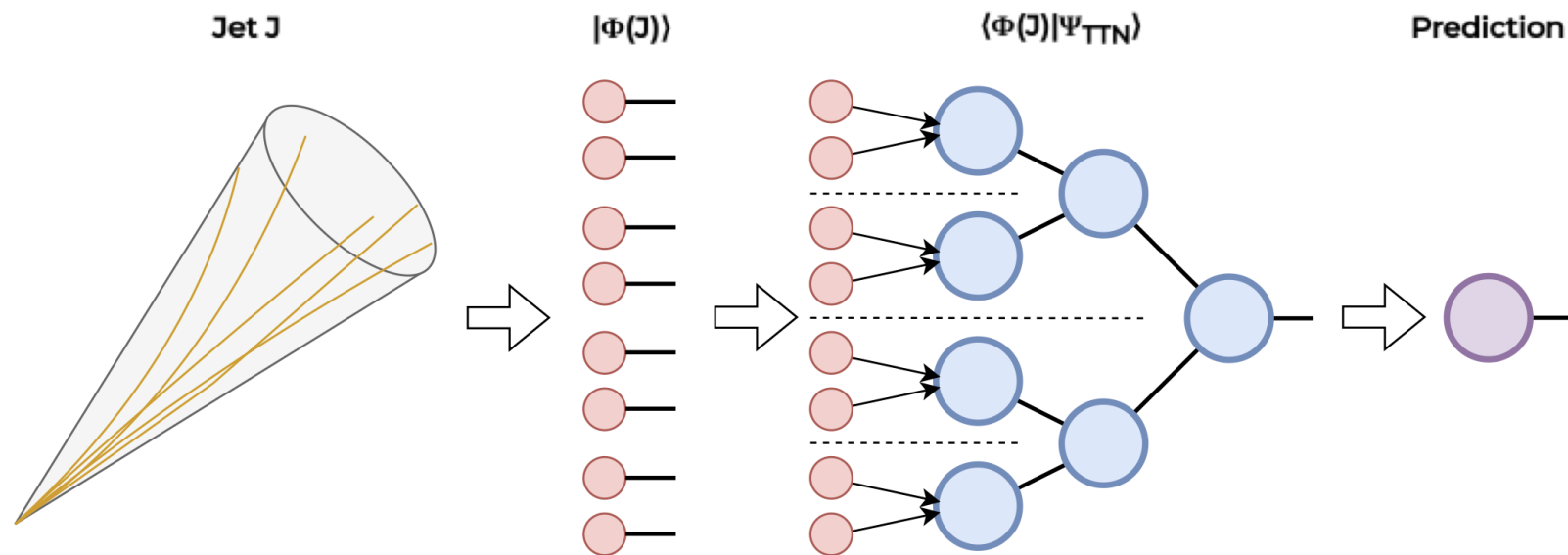Tree Tensor Network / Hierarchical Tucker

MERA

- Being a network, a TN can be trained as any other ML model.

- Methodology:
  - A TN used as "weight tensor" W, acting as a classifier on the input data {x}
  - the sample {x} is encoded into a feature map φ(x)
  - the confidence for a certain label *g* is

$$P_g = W \cdot \varphi(x)$$

- Eventually, the TN architecture encode the learned information representing a quantum entangled state.



Prediction $\begin{vmatrix} P_b \\ P_{\bar b} \end{vmatrix}$

FPGA

Tree Tensor Network

100 ns

Input Data

$q_\mu$  $p_{T\mu}$  $r_\mu$  $q_e$  $p_{Te}$  $r_e$  $q_\pi$  $p_{T\pi}$  $r_\pi$  $q_k$  $p_{Tk}$  $r_k$  $q_p$  $p_{Tp}$  $r_p$  $q_{jet}$

# TN features

## Linearity

Contractions are linear operation → Inference robustness against hallucinations, eases computational representation

## Compression (while learning)

bond dimensions optimized during training: reduction of number of parameters by truncating the size of the hidden links with SVD.

## Quantum correlations

remove redundant information by studying feature correlation and highlighting the ones that are the least correlated.

$$C_{i,j}^l = \frac{\langle\psi_l|\sigma_i^z\sigma_j^z|\psi_l\rangle}{\langle\psi_l|\psi_l\rangle}$$
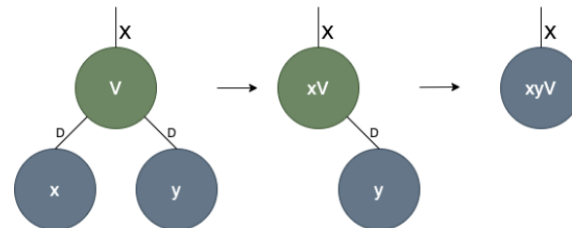
## Von Neumann Entropy

asses the relevance of the learned information encoded in each TTN bipartition → prune useless branches

$$S(\rho_A) = -\mathrm{Tr}[\rho_A \log\rho_A] = -\mathrm{Tr}[\rho_B \log\rho_B] = S(\rho_B)$$

- FPGA programmed with architecture-specific firmware.

- Software-trained weights loaded on static RAM blocks or hardcoded in firware.

- Data to be classified streamed to the FPGA.

- Full contraction with the TTN architecture.

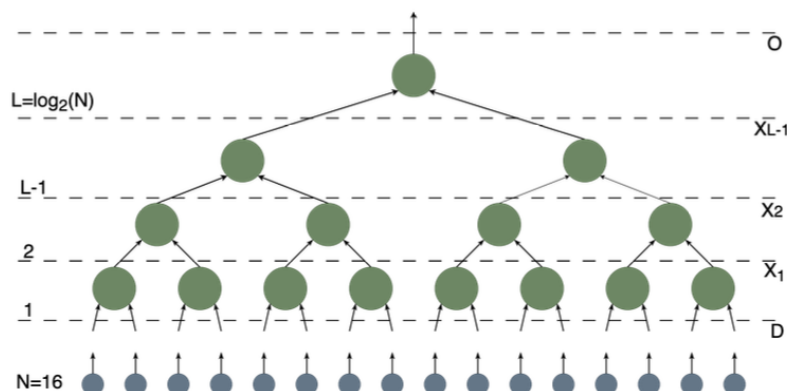- Final probability is retrieved for subsequent steps (e.g. selection)



Tensor contraction is the base operation that needs to be defined on FPGA: choose different degrees of parallelization and iterate it for different layers.



Digital Signal Processors (DSP) exploited for the actual node contraction; just products and sums.
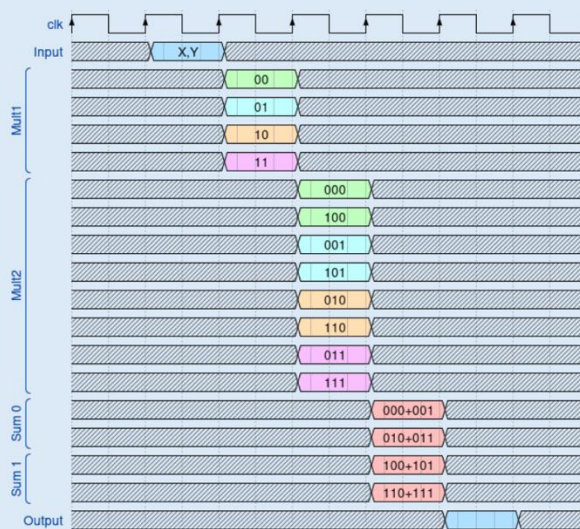
# Case study

- Consider datasets with increasing complexity

- Start of with a couple of ML banchmarks

- Main goal are HEP "standards"

  - b/anti-b classification datasets from LHCb
    - from A. Giannelle et al. «Quantum-inspired machine learning on high-energy physics data» Nature, 2021. https://doi.org/10.1038/s41534-021-00443-w

  - "hlt4ml" jet tagging dataset
    - from Duarte et al. «Fast inference of Deep Neural Networks in FPGA for particle Physics» https://arxiv.org/abs/1804.06913

- Methodology and results described in:

  - L. Borella et al. «*Ultra-low latency quantum-inspired machine learning predictors implemented on FPGA*», arxiv:2409.16075

Maximize number of DSPs used and minimize total algorithmic latency
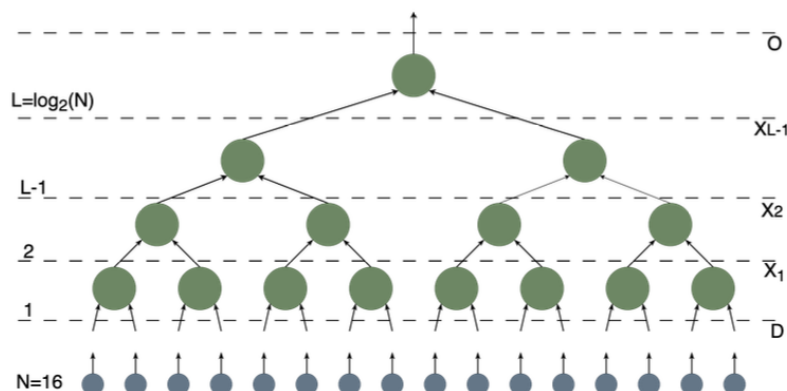


$$latency = \sum_{i=1}^{L} 2 + \log_2(\chi_{i-1}^2)$$

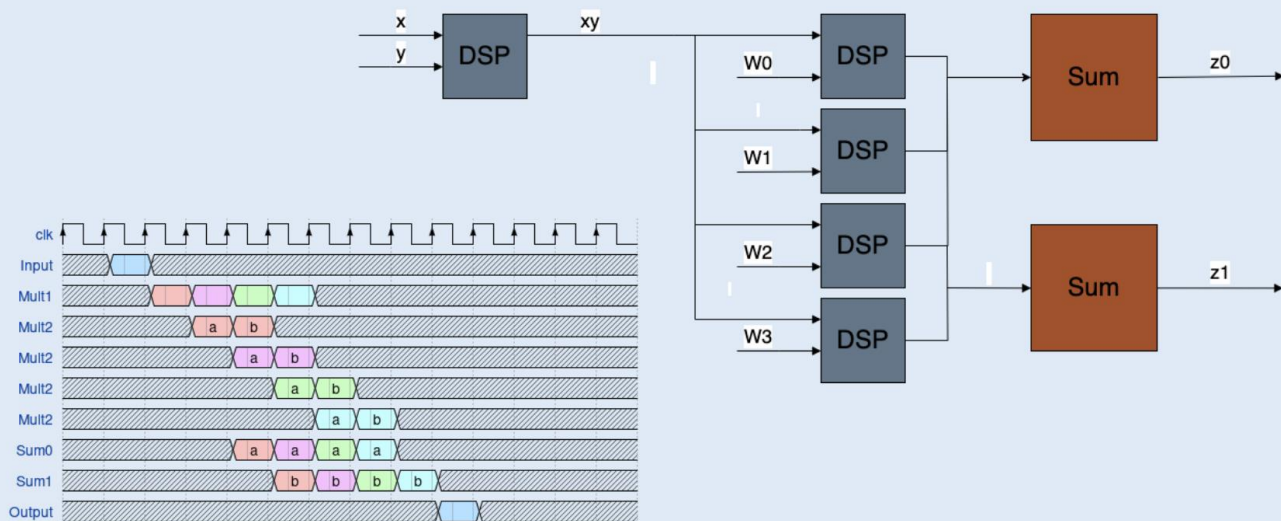$$DSP = \sum_{i=1}^{L} \chi_{i-1}^2 (\chi_i + 1) \frac{N}{2^i}$$

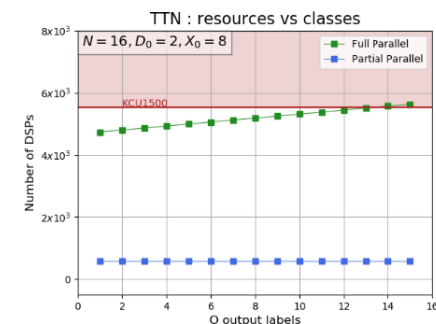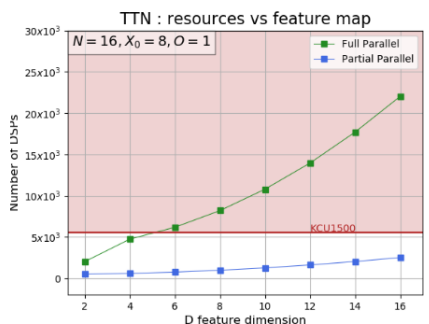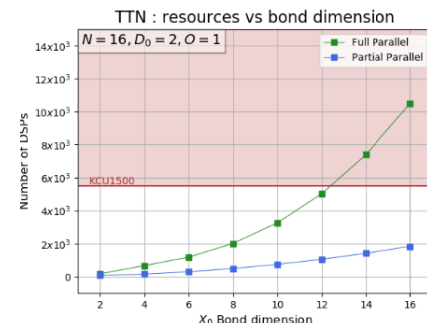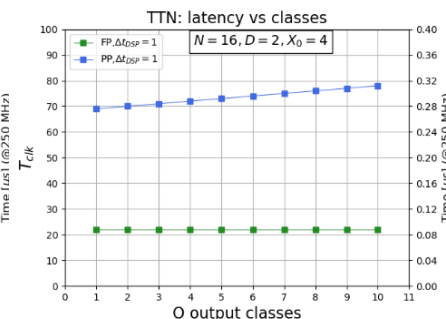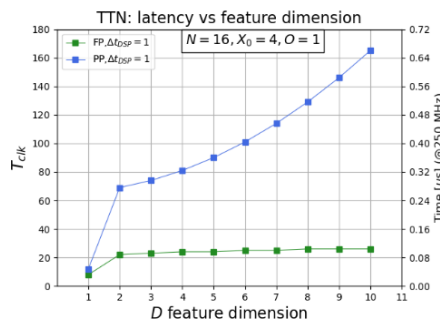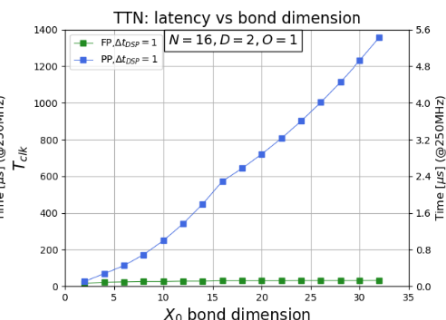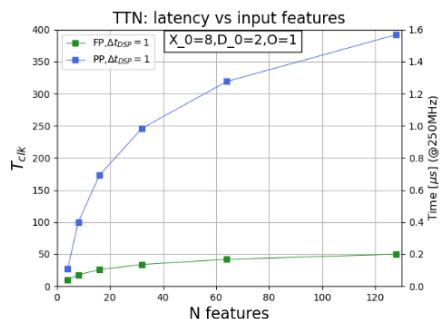Reuse some of the DSPs, with a resulting increase in latency



$$latency = \sum_{i=1}^{L} \chi_{i-1}^2 + \chi_i + 1$$

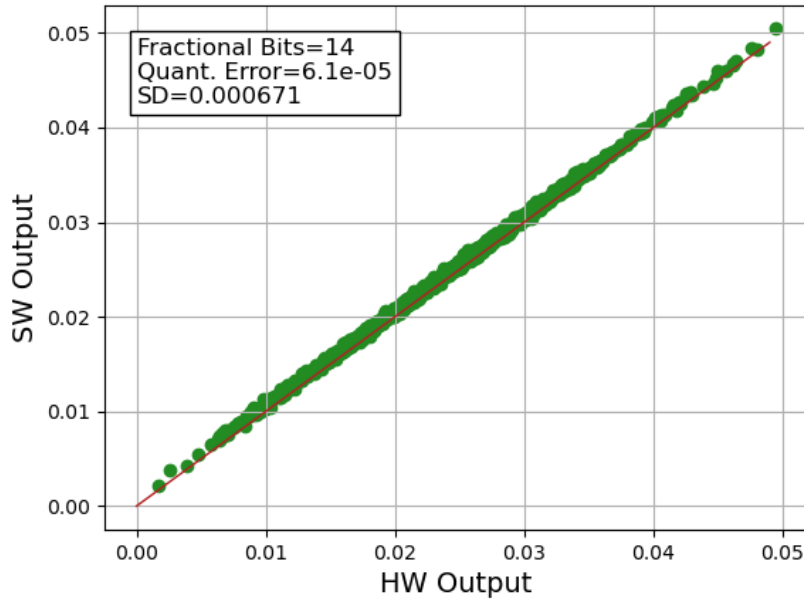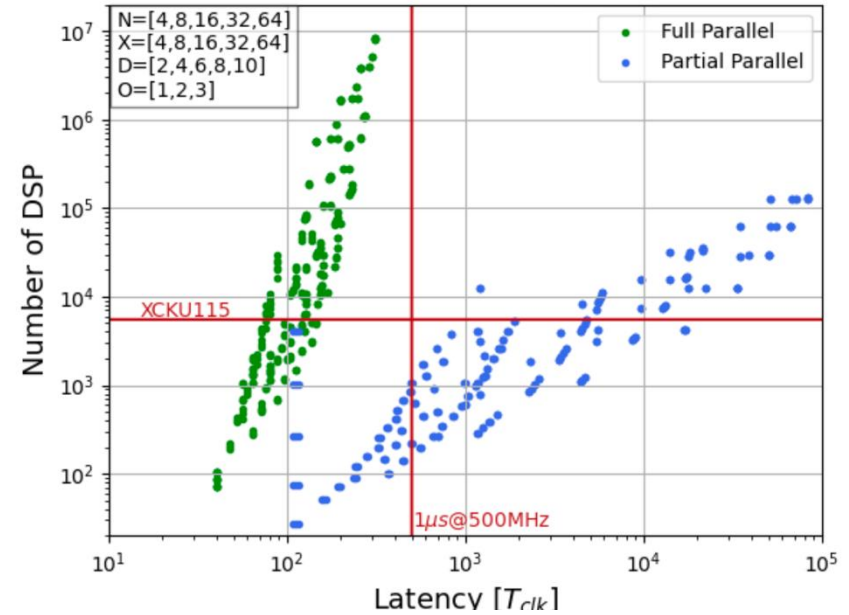$$DSP = \sum_{i=1}^{L} (\chi_{i-1}^2 + 1)\frac{N}{2^i}$$

## Latency

## DSP usage

SW/HW Output Comparison

Fractional Bits=14
Quant. Error=6.1e-05
SD=0.000671



Phase Space TTN

N=[4,8,16,32,64]
X=[4,8,16,32,64]
D=[2,4,6,8,10]
O=[1,2,3]

| Dataset | Iris | Titanic | LHCb [6] | hls4ml [5] |
|---|---|---|---|---|
| Features | 4 | 8 | 16 | 16 |
| Bond dimensions | [2,4] | [2,4,8] | [2,4,8,8] | [2,4,10,10] |
| Classes | 2 | 2 | 2 | 5 |
| Accuracy | 99% | 77% | 62% | 73% |
| Memory | 96 B | 768 B | 3 kB | 6 kB |

# Conclusions and Outlook

- TN as a valid option to tackle ML tasks

- TN's features make them suitable for deployment in hardware to achieve the ultrafast inference needed by HEP trigger systems

- Next steps:
  - Test different TN topologies (MPS, MERA, PEPS etc.) and tasks.
  - Move firmware programming to higher level languages (e.g. from VHDL to HLS4ML).
  - Check Hardware inference on Versal AI Engines.

- Other applications
  - Consider possibility of TN training on FPGA
  - Simulation of quantum circuits