**University Milano Bicocca**

Physics and Astronomy Doctoral School (39[th] cycle)

# Development of high-throughput machine learning techniques on FPGAs

Valentina Camagni

**Tutor**: Pietro Govoni[12]
**Supervisors**: Simone Gennai [2], Federico De Guio [12]

## End of the year seminar

September 20, 2024

# *FPGAs for CMS Level-1 Trigger*

CMS Phase II Level-1 Trigger system intends to perform precise physics selection using a global event reconstruction already at <u>hardware</u> level

Improved triggering with full detector view:
Trigger decision includes calorimeters, muons & tracker (~5us latency)

→ <u>L1Rate</u> = **750 kHz**
→ <u>Latency</u> = **12.5 us latency**
→ <u>Bandwidth</u>: ~ **50 Tb/s**
            (1.8 Tb/s in Phase I)

## FPGAs

✓ low-latency processing

✓ ability to handle highly parallel tasks

✓ reconfigurable nature allows for customization to meet specific requirements

✓ superior performance for real-time data processing, with lower power consumption

*Deploying ML on FPGAs*

*New trigger algorithms*

## Challenges

- meet the stringent latency requirements ($\mu s$)

- FPGA resources are limited: ML models need to be compressed and optimized through *quantization* and *pruning*

- Model optimization: tools like **hls4ml**, which facilitate high-level synthesis.

# *Planned activities*

**1** Starting from the Master Thesis work implement a DNN for the di $-\tau$ mass regression to replace SVFit algorithm in all Run III analyses

*Particle Transformer for $\tau$ lepton pair invariant mass reconstruction for the $HH \to b\bar{b}\tau^+\tau^-$ CMS analysis*

**Tau Pair Mass Transformer**
**TPMT**

✓ Tau costituents
✗ b-jets information

**2** Model distillation optimized for Phase-II implementation on FPGAs. Incorporating invariant mass information could lower the tau trigger threshold, currently set at 40 GeV, thereby recovering the corresponding phase space

*As CERN Doctoral student*

**3** Level-1 Trigger Scouting on soft taus.
Improvement of the trigger acceptance of tau leptons, specifically extending the coverage towards lower pT

UNIVERSITÀ DEGLI STUDI DI MILANO
BICOCCA

# PhD courses, Workshops and Schools

- ✓ Introduction to FPGAs *(November 2023)*

- ✓ ML@L1 Trigger Workshop at CERN  *(December 2023)*

- ✓ 6th Inter-experiment Machine Learning Workshop
    + poster presentation  *(February 2024)*

- ✓ Mandatory interdisciplinary courses:
    1. Communicating research in the era of social media
    2. Productivity tool for (young) researchers
    3. Surfing the academic job marketing

- ✓ Tutor activity for Laboratory II  *(March-June 2024)*

- ✓ AI-INFN 1° User Form (talk) *(June 2024)*

- X Internal courses:
    Deep Learning for Physicists (**to attend**)
    Physics at Colliders (**to attend**)
    Particle Physics II   (**ongoing**)

- ~ AI-PHY school *(October 2024)*
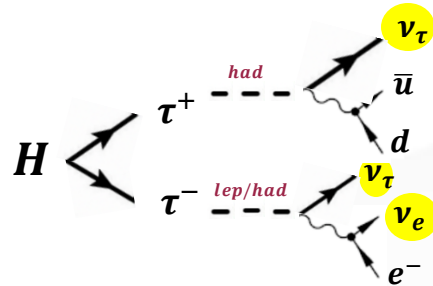
*Best presentation award - 109th SIF Conference*

Article publication on *Nuovo Cimento* Journal

*Open Access*
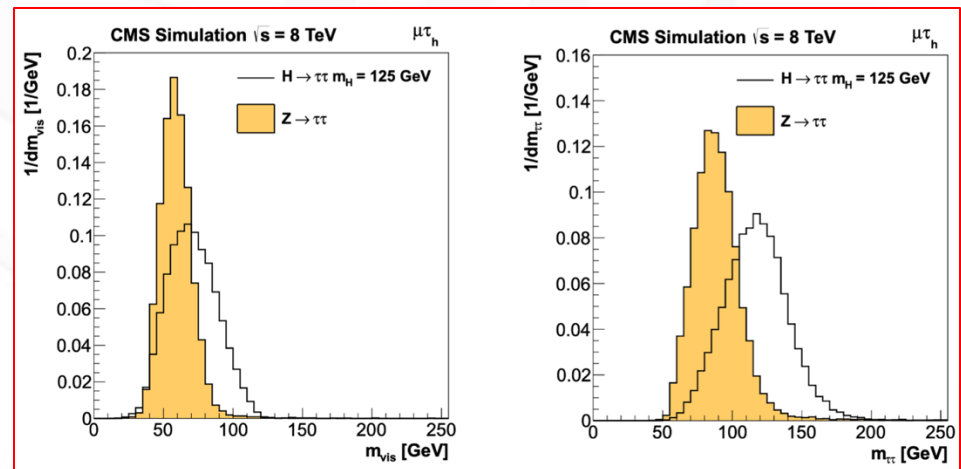
# *Di-$\tau$ invariant mass reconstruction*



The presence of neutrinos from tau decay prevent the full reconstruction of the di-tau system invariant mass, allowing only the reconstruction of the visible tau-decay products ($m_{\tau\tau}^{VIS}$) whose low resolution doesn't help in the signal discrimination task

**SVFit algorithm**

- Improves the $m_{\tau\tau}$ resolution only marginally
- High computational time

⇩

**Tau Pair Mass Transformer**
*TPMT*



**Objective:** Reconstruct the four-momentum of each $\tau$ particle before decay to accurately estimate the invariant mass and retrieve the kinematics of the parent particle

**1° GOAL**
Understand the model functionality on $H \to \tau^+\tau^-$ and $Z \to \tau^+\tau^-$
and considering only taus that decay hadronically so far

# Input features

## ① TauProd

**Taus' decay products**

Shape: $(10, 12) = (num\_tauprods , num\_features)$

*padding* if an event has less than 10 tau products

$logp_t$
$\eta$
$\phi$
$m$
$log\left(\dfrac{p_T}{p_{T(\tau)}}\right)$
$charge$
$tauIdx$

$is\_electron$
$is\_muon$
$is\_pion$
$is\_kaon$
$is\_photon$

*Categorical variables from particle ID*

## ② Tau

Shape: $(6, 3) =$
$(num\_particles , num\_features)$

$\tau_1$     $logp_T$
$\tau_2$     $\eta$
$MET$     $\phi$
$jet_1$
$jet_2$
$jet_3$

## GenPart

Not the full 4-momentum since eta and phi does not change

Shape: $(2, 1) =$
$(num\_particles , num\_features)$

$\tau_1$     $logp_T$
$\tau_2$

+ di-$\tau$ invariant mass at generator level

$m_{\tau\tau}^{GEN}$

UNIVERSITÀ DEGLI STUDI DI MILANO BICOCCA

# *Pre-processing steps*

Data sets   **GluGluHToTauTau_M125**          **DYJetsToLL_M-50-madgraphMLM**

**TAU SELECTION**

At least 2 taus
- Gen matched
- Hadronic decay
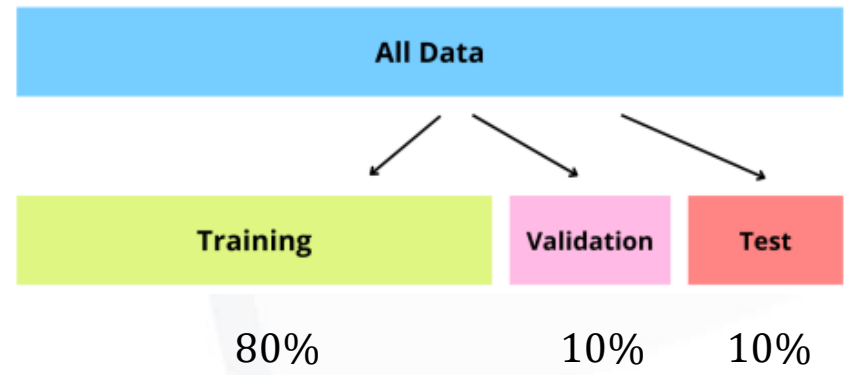- $p_T \geq 20$ GeV

**JETS SELECTION**

First 3 leading jets with
$\Delta R(jet, tau) > 0.4$

(minimum $p_T$: 10 GeV )

**VARIABLE ENCODING & FEATURE ENGENEERING**

- Definition of new variables
- Order TauProd with respect to their $p_T$ and padding with max_$len = 10$
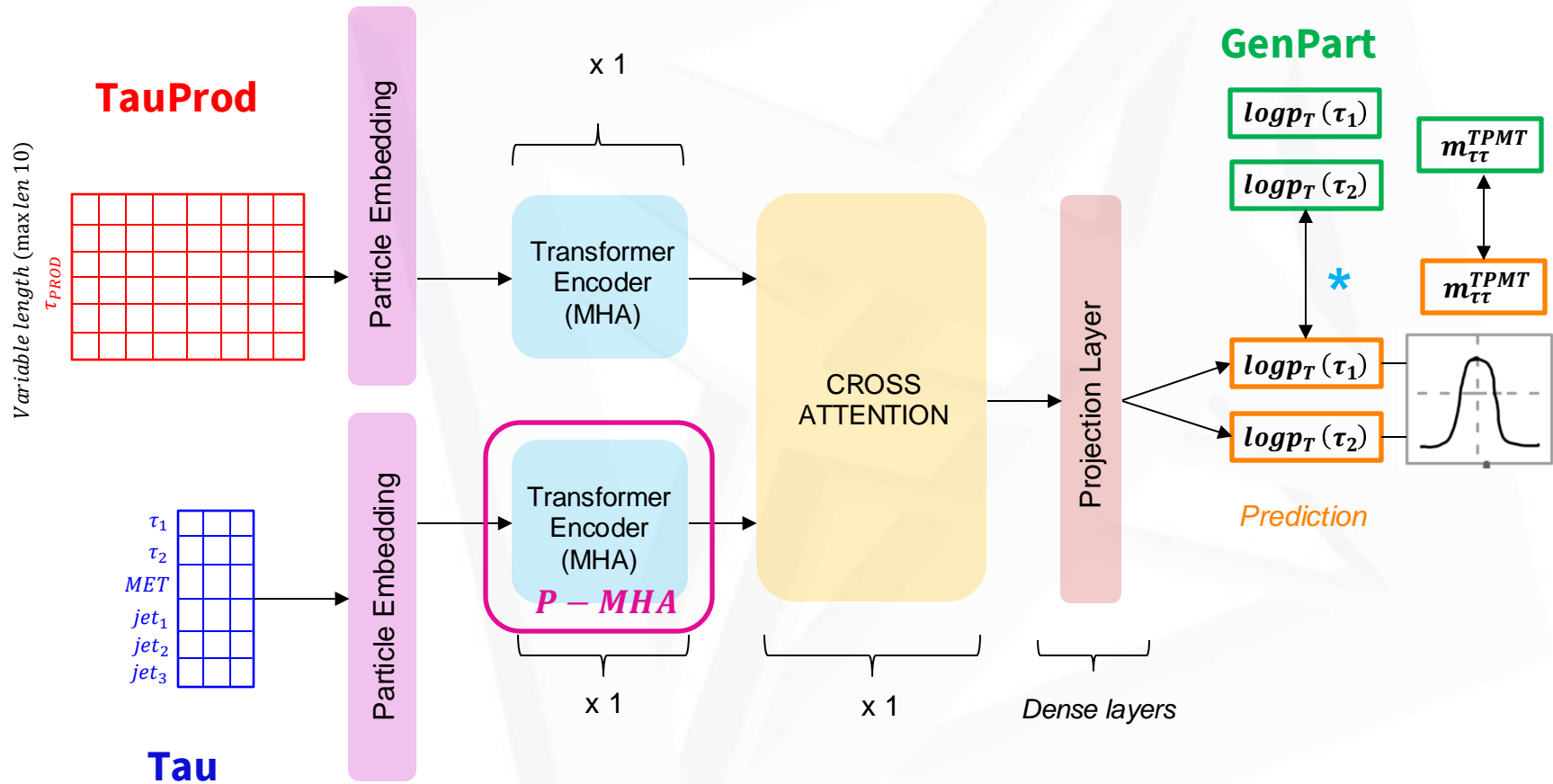
**SPLIT IN TRAIN, TEST AND VALIDATION**



All Data

Training     Validation     Test

80%          10%        10%

# Model Architecture

PyTorch

**TauProd**

*Variable length (max len 10)*
$\tau_{PROD}$

Particle Embedding

x 1

Transformer Encoder (MHA)

**GenPart**

$logp_T(\tau_1)$

$logp_T(\tau_2)$

$m_{\tau\tau}^{TPMT}$

$m_{\tau\tau}^{TPMT}$

CROSS ATTENTION

Projection Layer

$logp_T(\tau_1)$

$logp_T(\tau_2)$

*Prediction*

$\tau_1$
$\tau_2$
$MET$
$jet_1$
$jet_2$
$jet_3$

Particle Embedding

Transformer Encoder (MHA)
$P - MHA$

x 1

x 1

*Dense layers*

**Tau**

*

**Training time:** $\sim 1.5\ min\ per\ epoch$
**Inference time:** $\sim 2 \times 10^{-3}\ s\ per\ event$
**Number of parameters:** $\sim 0.5\ M$

UNIVERSITÀ DEGLI STUDI DI MILANO BICOCCA

# $m_{\tau\tau}$ results



**GluGluHToTauTau_M125**

**DYJetsToLL_M-50-madgraphMLM**

**AUC as evaluation metric**

Train ratio
$H : Z = 2 : 1$

Best training:
*AUC score* of 0.84

# *Preliminary considerations*

✓ AUC suggests that TPMT has a better separation capability

✗ The wrong peak is slightly higher for H than for DY (due to the different response)

✓ Training time: $1.5 \ min \ per \ epoch \ (\sim 80 \ epochs)$
Inference time: $2 \cdot 10^{-3} \ s$

✗ Inference on any other resonance would have worked worse
(if not added in the train set composition)

⇩

Training on flat mass samples
**GluGlutoXto2Tau_M-30to300**
**VBFtoXto2Tau_M-30to300**
and inference on H and Z samples

! **No more jet information**

! **Leptonic decaying taus in addition to hadronic ones $(\tau_h + l\,(e, \mu))$**

UNIVERSITÀ DEGLI STUDI DI MILANO BICOCCA

# Overall training on ggF sample – tau_tau, ele_tau, mu_tau



| Fit Type | | Mean | Std |
|---|---|---|---|
| $m_{\tau\tau}^{SVFit} - \mathbf{H}$ | | 127.26 | 30.79 |
| $m_{\tau\tau}^{TPMT} - \mathbf{H}$ | tau_tau | 129.81 | 28.08 |
| $m_{\tau\tau}^{SVFit} - \mathbf{DY}$ | | 99.83 | 26.28 |
| $m_{\tau\tau}^{TPMT} - \mathbf{DY}$ | | 101.47 | 21.18 |
| $m_{\tau\tau}^{SVFit} - \mathbf{H}$ | | 164.98 | 47.23 |
| $m_{\tau\tau}^{TPMT} - \mathbf{H}$ | ele_tau | 139.72 | 30.22 |
| $m_{\tau\tau}^{SVFit} - \mathbf{DY}$ | | 152.69 | 53.85 |
| $m_{\tau\tau}^{TPMT} - \mathbf{DY}$ | | 116.53 | 23.47 |
| $m_{\tau\tau}^{SVFit} - \mathbf{H}$ | | 166.28 | 47.63 |
| $m_{\tau\tau}^{TPMT} - \mathbf{H}$ | mu_tau | 139.91 | 29.51 |
| $m_{\tau\tau}^{SVFit} - \mathbf{DY}$ | | 155.12 | 53.94 |
| $m_{\tau\tau}^{TPMT} - \mathbf{DY}$ | | 116.81 | 23.50 |

# $p_T^{\tau 1}, p_T^{\tau 2}, m_{\tau\tau}$ resolution results for tau_tau pairType      *H* & *Z*



**p_T resolution - $\tau_h$**

**m_$\tau\tau$ resolution**

tau_tau

**p_T resolution - $\tau_h$**

| RECO | | |
|---|---|---|
| **Fit Type** | **Mean** | **Std** |
| $p_T\ \tau_1$ - **H** | -0.26 | 0.2 |
| $p_T\ \tau_2$ - **H** | -0.42 | 0.35 |
| $m_{\tau\tau}$ - **H** | -0.35 | 0.15 |
| $p_T\ \tau_1$ - **DY** | -0.22 | 0.18 |
| $p_T\ \tau_2$ - **DY** | -0.33 | 0.41 |
| $m_{\tau\tau}$ - **DY** | -0.28 | 0.13 |

| TPMT | | |
|---|---|---|
| **Fit Type** | **Mean** | **Std** |
| $p_T\ \tau_1$ - **H** | 0.04 | 0.23 |
| $p_T\ \tau_2$ - **H** | 0.04 | 0.35 |
| $m_{\tau\tau}$ - **H** | 0.04 | 0.3 |
| $p_T\ \tau_1$ - **DY** | 0.13 | 0.23 |
| $p_T\ \tau_2$ - **DY** | 0.12 | 0.28 |
| $m_{\tau\tau}$ - **DY** | 0.12 | 0.23 |

# $p_T^{\tau 1}, p_T^{\tau 2}, m_{\tau\tau}$ resolution results for ele_tau pairType    *H* & *Z*



ele_tau

| RECO | | |
|---|---|---|
| **Fit Type** | **Mean** | **Std** |
| $p_T$ $\tau_1$ - **H** | -0.26 | 0.2 |
| $p_T$ $\tau_2$ - **H** | -0.42 | 0.35 |
| $m_{\tau\tau}$ - **H** | -0.35 | 0.15 |
| $p_T$ $\tau_1$ - **DY** | -0.22 | 0.18 |
| $p_T$ $\tau_2$ - **DY** | -0.33 | 0.41 |
| $m_{\tau\tau}$ - **DY** | -0.28 | 0.13 |

| TPMT | | |
|---|---|---|
| **Fit Type** | **Mean** | **Std** |
| $p_T$ $\tau_1$ - **H** | 0.04 | 0.23 |
| $p_T$ $\tau_2$ - **H** | 0.04 | 0.35 |
| $m_{\tau\tau}$ - **H** | 0.04 | 0.3 |
| $p_T$ $\tau_1$ - **DY** | 0.13 | 0.23 |
| $p_T$ $\tau_2$ - **DY** | 0.12 | 0.28 |
| $m_{\tau\tau}$ - **DY** | 0.12 | 0.23 |

# $p_T^{\tau1}, p_T^{\tau2}, m_{\tau\tau}$ resolution results for tau_tau pairType    *ggF* & *VBF*



**RECO**

| Fit Type | | Mean | Std |
|---|---|---|---|
| $p_T\ \tau_1$ - | **ggF** | -0.26 | 0.2 |
| $p_T\ \tau_2$ - | **ggF** | -0.44 | 0.35 |
| $m_{\tau\tau}$ - | **ggF** | -0.35 | 0.18 |
| $p_T\ \tau_1$ - | **VBF** | -0.29 | 0.21 |
| $p_T\ \tau_2$ - | **VBF** | -0.37 | 0.39 |
| $m_{\tau\tau}$ - | **VBF** | -0.37 | 0.18 |

**TPMT**

| Fit Type | | Mean | Std |
|---|---|---|---|
| $p_T\ \tau_1$ - | **ggF** | -0.02 | 0.21 |
| $p_T\ \tau_2$ - | **ggF** | -0.02 | 0.26 |
| $m_{\tau\tau}$ - | **ggF** | -0.02 | 0.19 |
| $p_T\ \tau_1$ - | **VBF** | -0.04 | 0.2 |
| $p_T\ \tau_2$ - | **VBF** | -0.03 | 0.28 |
| $m_{\tau\tau}$ - | **VBF** | -0.04 | 0.18 |

# $p_T$ ratio versus $p_T^{RECO}$ for tau_tau pairType



➢ **Different response between resonances and flat mass samples**
➢ **More differences between H and Z compared to ggF and VBF responses**

**Due to convolution of tau resolution and $p_T^{GEN}$ distribution**

**Studying new
training strategies**

## *Conclusions*

➤ **Training on H and DY**
- TPMT behavies as a classifier
- Good mass resolution but strong dependent on the training samples

➤ **Training on ggF sample**
- Resolution and fits much worst, still better than SVFit but suboptimal

For optimal training, it is essential to include samples that reflect the true underlying distributions of the events whose mass we aim to estimate, rather than using flat distributions that can lead to suboptimal performance

## *Future plans*

- **Add a loss term regarding MET**   $\mathcal{L}_{MET} = |MET_{observed} - \left(p_T^{neutrinos}\right)|$

- **Train TPMT with the TauProd matrix divided by the two taus**

# New Model Architecture



**TauProd**

$\tau_1$

*Variable length (max len 5)*

$\tau_{PROD}$

$\tau_2$

*Variable length (max len 5)*

$\tau_{PROD}$

$\tau_1$
$\tau_2$
MET

**Tau**

Particle Embedding

Transformer Encoder (MHA)

Transformer Encoder (MHA)

Transformer Encoder (Powered-MHA)
$P - MHA$

x 1

x 1

+

CROSS ATTENTION

x 1

Projection Layer

*Dense layers*

**GenPart**

$logp_T(\tau_1)$

$logp_T(\tau_2)$

$m_{\tau\tau}^{TPMT}$

*

$m_{\tau\tau}^{TPMT}$

$logp_T(\tau_1)$

$logp_T(\tau_2)$

*Prediction*

**Loss function**

① Mean between $MAE_{logp_T}$ for the two taus

② MAE between $m_{\tau\tau}^{TRANS}$ and $m_{\tau\tau}^{MC}$   (7% of the total loss)

$$MAE = \frac{1}{N}\sum_{i=1}^{N}|y_i - \hat{y}_i|$$

③ MET term $= |MET_{observed} - (p_T^{neutrinos})|$

**Number of parameters:** $\sim 0.9\ M$

UNIVERSITÀ DEGLI STUDI DI MILANO
BICOCCA

*Thank you for your attention!*

*BACKUP*

# *Powered MHA*

Particles → Embedding → $\mathbf{x}^0$ → Transformer Encoder (P-MHA)

**NODES**
Array
$(6, 4)$

Interactions → Embedding → $\mathbf{U}$

**EDGES**
Array
$(6, 6, 4)$

**4 pairwise features**

$$(\Delta, k_T, z, m^2)$$

from Particle Transformer for jet tagging

**P-MHA**

MatMul

SoftMax

$V$

$\mathbf{U}$ → $\oplus$

Mask

Scale

MatMul

$Q$    $K$

Linear   Linear   Linear

$\mathbf{x}$

$$\Delta = \sqrt{(y_a - y_b)^2 + (\phi_a - \phi_b)^2},$$
$$k_{\mathrm{T}} = \min(p_{\mathrm{T},a}, p_{\mathrm{T},b})\Delta,$$
$$z = \min(p_{\mathrm{T},a}, p_{\mathrm{T},b})/(p_{\mathrm{T},a} + p_{\mathrm{T},b}),$$
$$m^2 = (E_a + E_b)^2 - \|\mathbf{p}_a + \mathbf{p}_b\|^2,$$

# $logp_T^{\tau_1}, logp_T^{\tau_2}$ results

There is a $logp_T$ transition region:
for $logp_T < 4$, H and Z taus' $logp_T$ need a different scale factor



logp_T histograms - $\tau_1$ - H sample

Leads to the wrong peak mass

logp_T histograms - $\tau_2$ - H sample

logp_T histograms - $\tau_1$ - Z sample

logp_T histograms - $\tau_2$ - Z sample

# $p_T^{\tau 1}, p_T^{\tau 2}, m_{\tau\tau}$ resolution results for mu_tau pairType



|  |  |  |  |  |
| --- | --- | --- | --- | --- |
| **RECO** | | | | |
| **Fit Type** | **Mean** | **Std** | | |
| $p_T\ \tau_1$ - **H** | -0.26 | 0.2 | | |
| $p_T\ \tau_2$ - **H** | -0.42 | 0.35 | | |
| $m_{\tau\tau}$ - **H** | -0.35 | 0.15 | | |
| $p_T\ \tau_1$ - **DY** | -0.22 | 0.18 | | |
| $p_T\ \tau_2$ - **DY** | -0.33 | 0.41 | | |
| $m_{\tau\tau}$ - **DY** | -0.28 | 0.13 | | |

| **TPMT** | | |
| --- | --- | --- |
| **Fit Type** | **Mean** | **Std** |
| $p_T\ \tau_1$ - **H** | 0.04 | 0.23 |
| $p_T\ \tau_2$ - **H** | 0.04 | 0.35 |
| $m_{\tau\tau}$ - **H** | 0.04 | 0.3 |
| $p_T\ \tau_1$ - **DY** | 0.13 | 0.23 |
| $p_T\ \tau_2$ - **DY** | 0.12 | 0.28 |
| $m_{\tau\tau}$ - **DY** | 0.12 | 0.23 |

# $m_{\tau\tau}^H, m_{\tau\tau}^Z$ quartiles



Box Plot for SVFIT - TPMT H mass distribution comparison

tau_tau

| Distribution | Q1 | Q2 | Q3 |
|---|---|---|---|
| SVFIT - H | 109.19 | 130.00 | 153.64 |
| TPMT - H | 111.37 | 130.18 | 149.17 |
| SVFIT - DY | 85.31 | 103.24 | 124.69 |
| TPMT - DY | 88.09 | 102.31 | 117.06 |

ele_tau

| Distribution | Q1 | Q2 | Q3 |
|---|---|---|---|
| SVFIT - H | 147.15 | 191.63 | 270.05 |
| TPMT - H | 120.89 | 140.39 | 160.59 |
| SVFIT - DY | 132.44 | 184.51 | 279.27 |
| TPMT - DY | 101.99 | 117.80 | 134.57 |

mu_tau

| Distribution | Q1 | Q2 | Q3 |
|---|---|---|---|
| SVFIT - H | 148.48 | 193.65 | 272.63 |
| TPMT - H | 120.80 | 140.67 | 160.89 |
| SVFIT - DY | 134.61 | 186.57 | 279.32 |
| TPMT - DY | 102.45 | 117.99 | 134.84 |

# $p_T$ ratio versus $p_T^{RECO}$ for all pairTypes

# $p_T$ ratio versus $p_T^{RECO}$ for all pairTypes

Ratio $p_T^{RECO}/p_T^{GEN}$ vs $p_T^{GEN}$

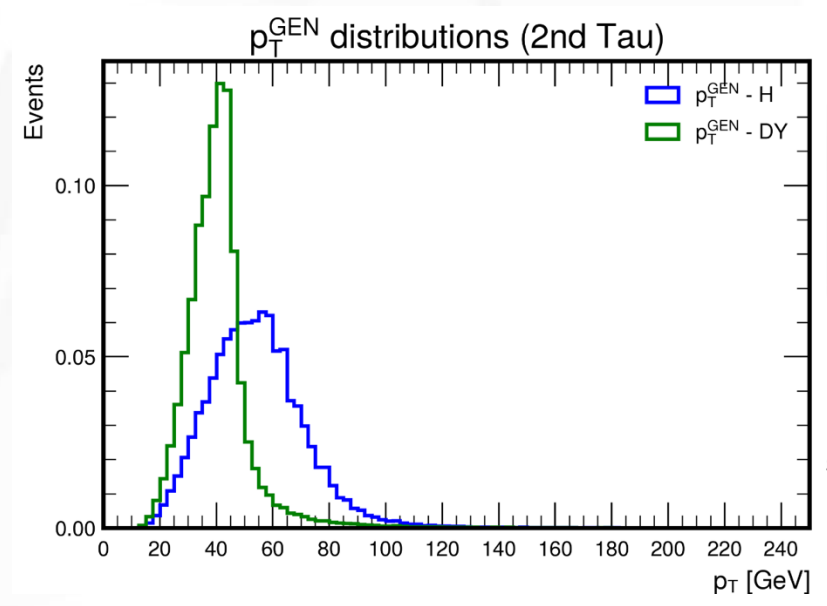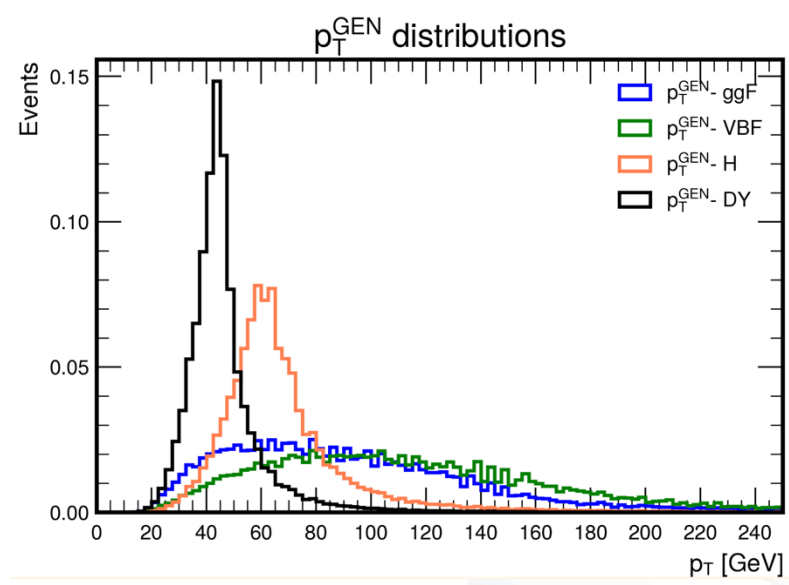# $p_T^{GEN}$ *distributions*

**Before resampling on the first tau**

**After resampling on the first tau**



$p_T^{GEN}$ distributions

- $p_T^{GEN}$ - ggF
- $p_T^{GEN}$ - VBF
- $p_T^{GEN}$ - H
- $p_T^{GEN}$ - DY



$p_T^{GEN}$ distributions

- $p_T^{GEN}$ - H
- $p_T^{GEN}$ - DY



$p_T^{GEN}$ distributions (2nd Tau)

- $p_T^{GEN}$ - ggF
- $p_T^{GEN}$ - VBF
- $p_T^{GEN}$ - H
- $p_T^{GEN}$ - DY



$p_T^{GEN}$ distributions (2nd Tau)

- $p_T^{GEN}$ - H
- $p_T^{GEN}$ - DY

## Scaled Dot - Product

$$Multihead(Q, K, V) = Concat(head_1, head_2, ..., head_h)W^O$$

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



**Scaled Dot-Product Attention**

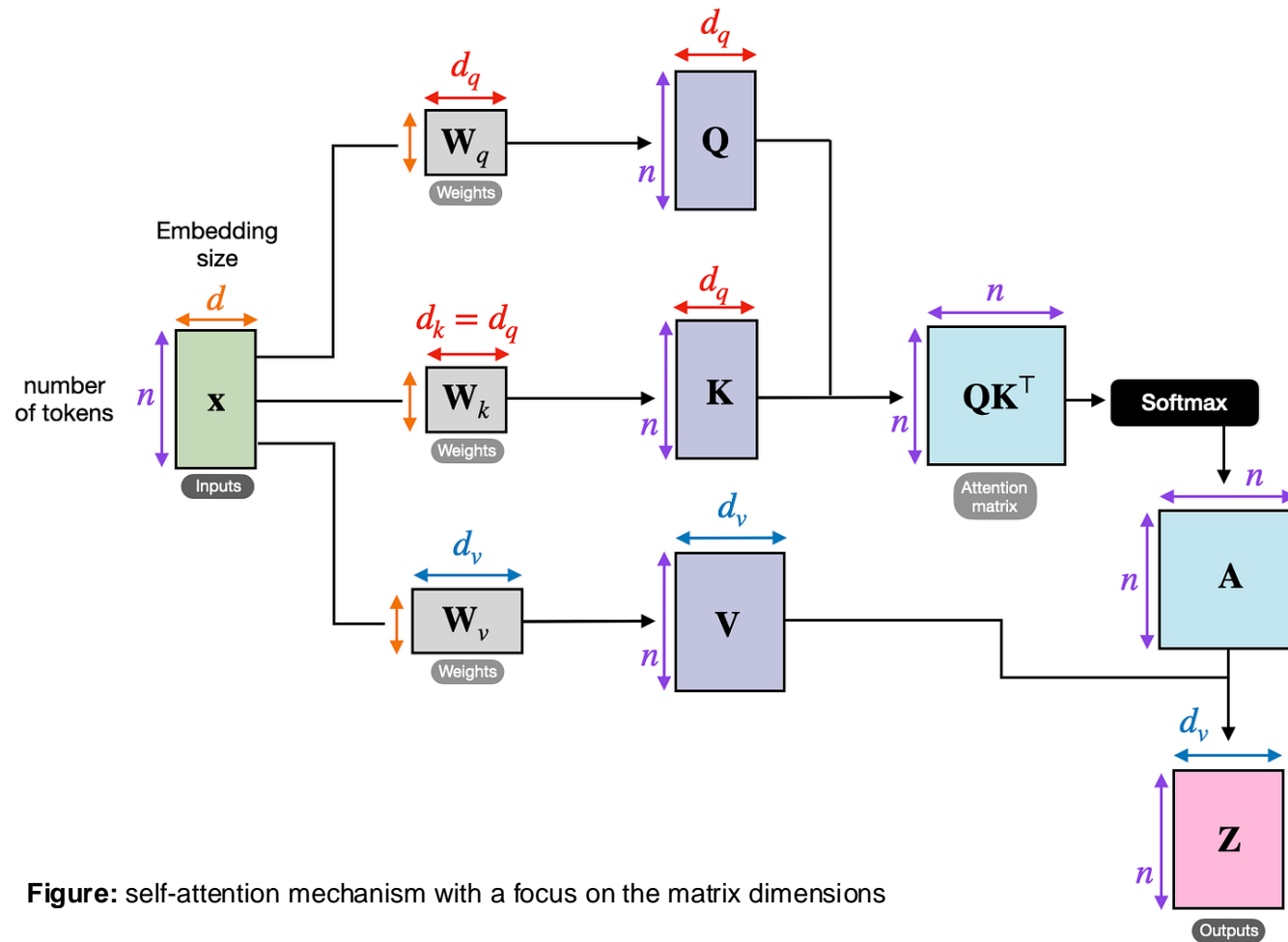**Multi-Head Attention**

# Self-Attention



**Figure:** self-attention mechanism with a focus on the matrix dimensions
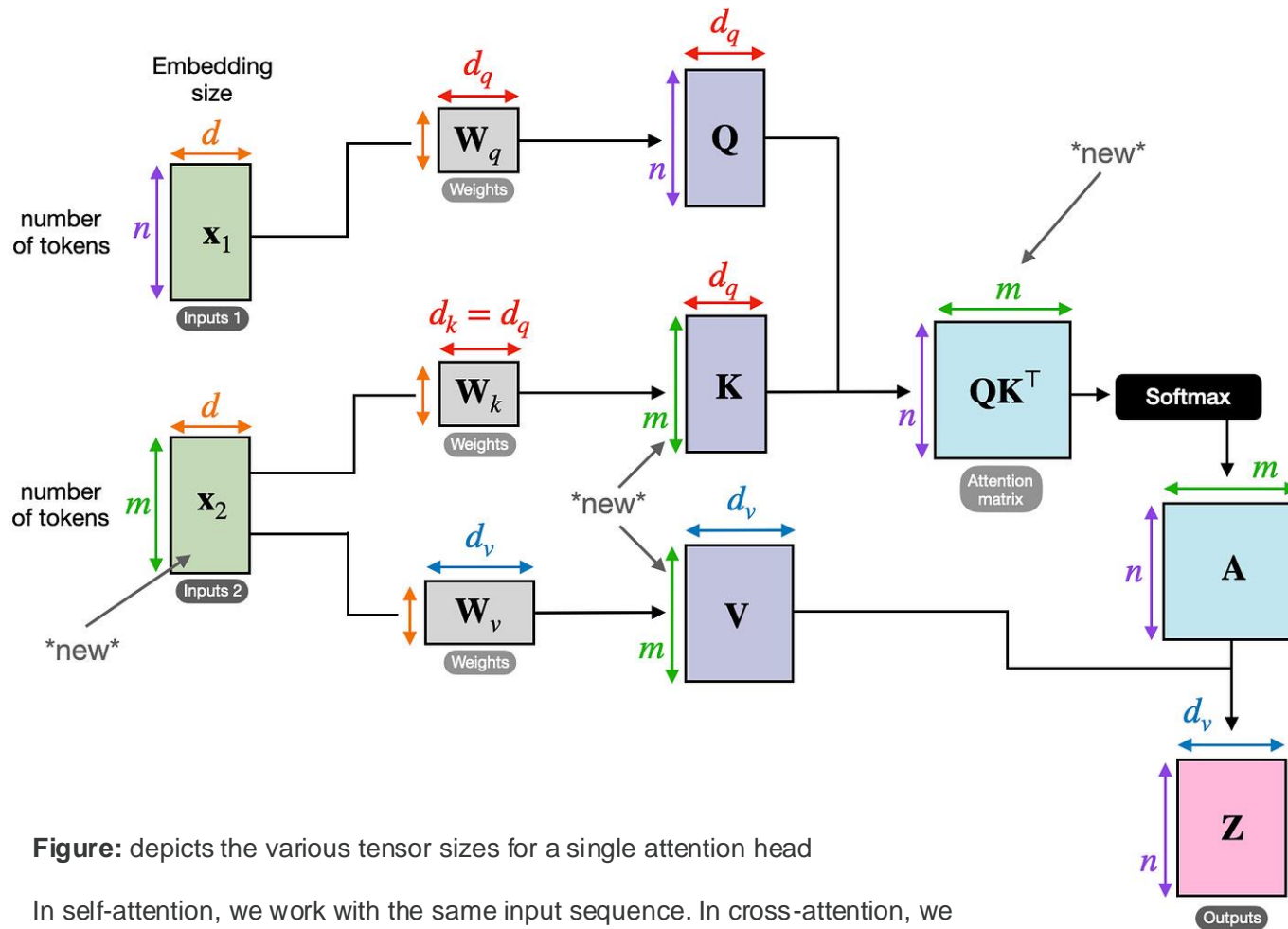
# Cross-Attention



**Figure:** depicts the various tensor sizes for a single attention head

In self-attention, we work with the same input sequence. In cross-attention, we mix or combine two *different* input sequences. In the case of the original transformer architecture, that's the sequence returned by the encoder module and the input sequence being processed by the decoder part on the right. The two input sequences and can have different numbers of elements. However, their embedding dimensions must match.

# Multi-scale cross-attention transformer encoder for event classification

A. Hammad[a], S. Moretti[b,c] and M. Nojiri[a,d,e]

[a]Theory Center, IPNS, KEK, 1-1 Oho, Tsukuba, Ibaraki 305-0801, Japan.
[b]School of Physics and Astronomy, University of Southampton, Highfield, Southampton, UK.
[c]Department of Physics & Astronomy, Uppsala University, Box 516, SE-751 20 Uppsala, Sweden.
[d]The Graduate University of Advanced Studies (Sokendai), 1-1 Oho, Tsukuba, Ibaraki 305-0801, Japan
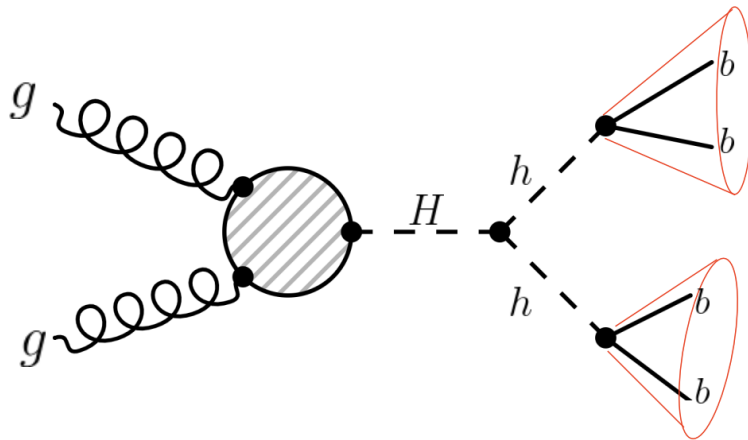[e]Kavli IPMU (WPI), University of Tokyo, 5-1-5 Kashiwanoha, Kashiwa, Chiba 277-8583, Japan

Figure 2: Feynman diagram for the signal process.