

# Multi-Scale Transformer Encoder for Di-Tau Invariant Mass Reconstruction

Valentina Camagni

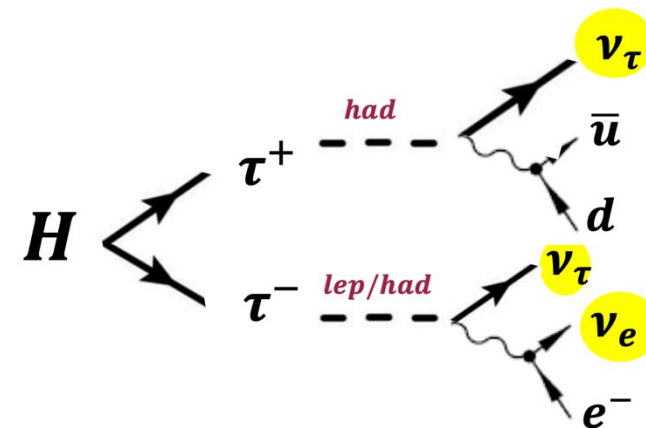
*on behalf of the CMS collaboration*

**EPS-HEP 2025**

July 9, 2025

# Introduction (1)

Reconstructing the di-tau invariant mass ( $m_{\tau\tau}$ ) is crucial for precise SM measurements and BSM searches.  
Neutrinos from tau decays escape detection, worsening mass resolution and complicating resonance identification.



The CMS experiment currently employs the **Secondary Vertex Fit (SVFit)** algorithm [\[1\]](#)

✗ High computational time

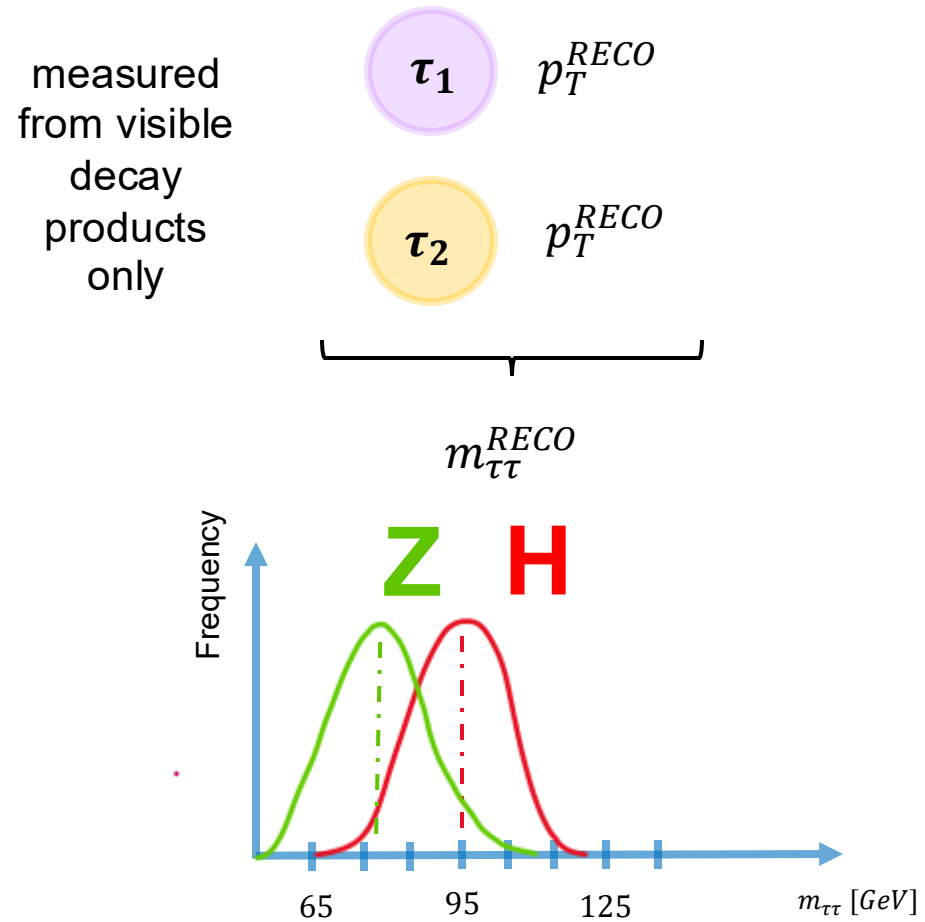


New strategies based  
on deep learning  
**Tau Pair Mass  
Transformer (TPMT)**

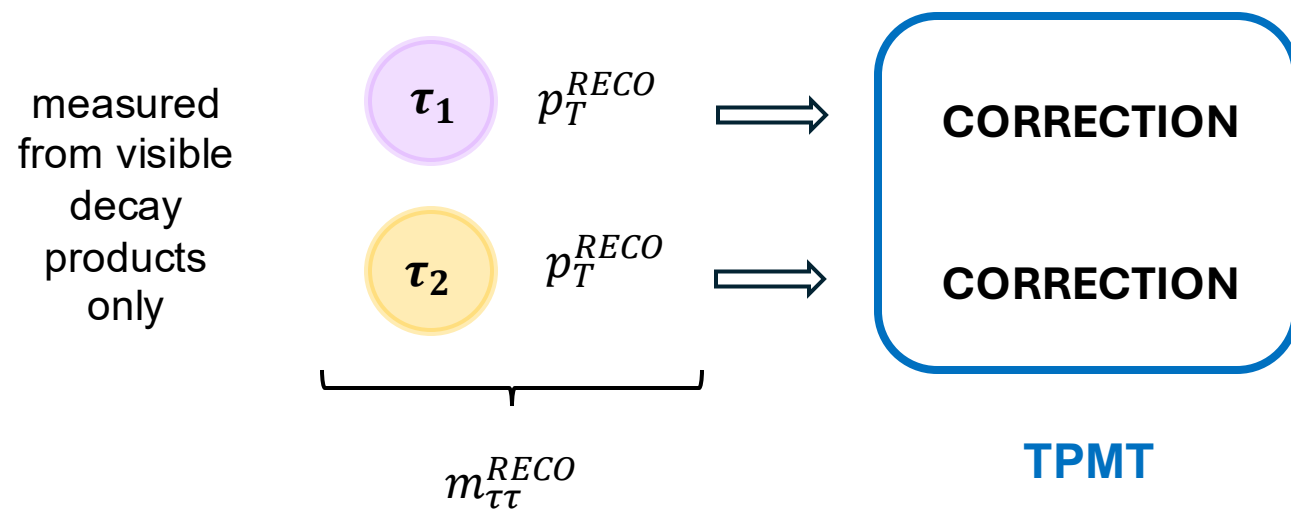
**Aim**

Reconstruct the four-momentum of each  $\tau$  lepton prior to its decay, in order to recover the kinematics of the parent particle and accurately estimate the invariant mass

# Introduction (2)

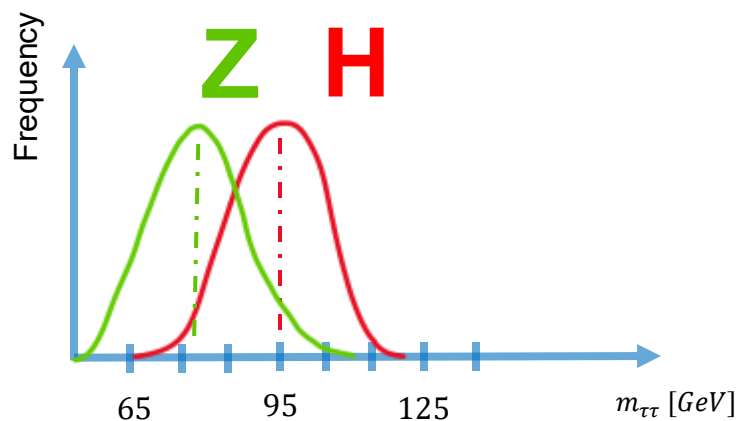


# Introduction (2)

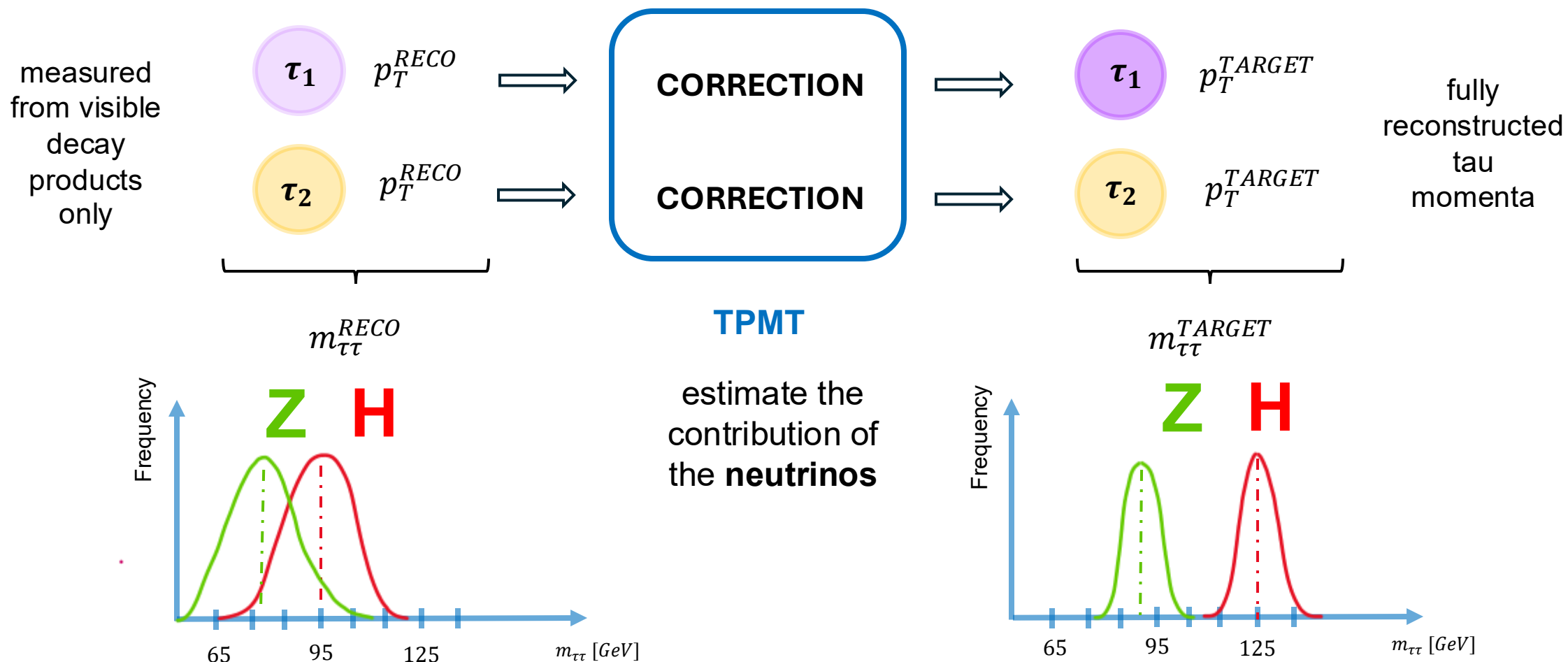


TPMT

estimate the contribution of the **neutrinos**



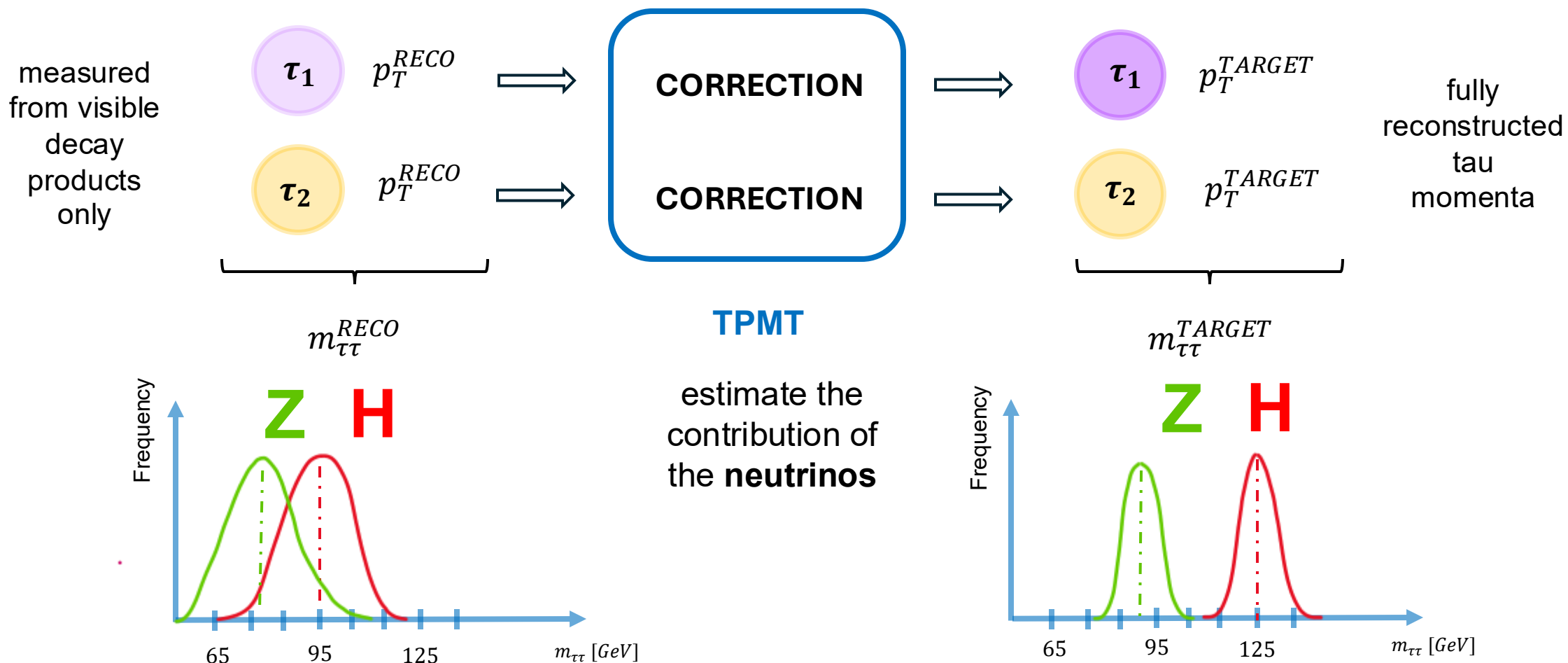
# Introduction (2)



# Introduction (2)

## Assumption

- Uses  $\eta$  and  $\phi$  from reconstructed taus  
 $\rightarrow$  *Collinear approximation* valid for taus  
 with  $p_T^{RECO} > 20 \text{ GeV}$



# First Strategy: Flat-mass samples

**Goal:** improve the reconstruction of the di-tau invariant mass by correcting the visible  $p_T$  of tau candidates for the momentum carried away by neutrinos, without biasing the model toward any specific mass value

## Training configuration:

- Events from  $X \rightarrow \tau\tau$  decays, generated via gluon-gluon fusion and vector boson fusion
- The invariant mass of the parent particle was drawn from a flat distribution in the range 30–300 GeV
- Includes both hadronic ( $\tau_h\tau_h$ ) and semileptonic tau decays ( $\tau_h\tau_\mu$ )

## Evaluation configuration:

After training, the model was tested on realistic resonant processes to assess performance:

- **Higgs boson production** ( $H \rightarrow \tau\tau$  with  $m_{\tau\tau} = 125 \text{ GeV}$ )
- **Drell–Yan** events ( $Z \rightarrow \tau\tau$  with  $m_{\tau\tau} \approx 91 \text{ GeV}$ )

## Motivation

Using flat-mass training allows the model to learn from kinematics alone, avoiding dependence on Z/H mass peaks and minimizing sculpting effects

# TPMT inputs

Feature importance analysis with **Random Forest** on the full set of taus, jets, MET variables to maximize H / Z classification

[2], [3]

## Taus

$\tau_h \tau_h$

**H/Z/SUSY/X:** 2 genuine  $\tau_h$   
(opposite charge,  $p_T^{RECO} > 20 \text{ GeV}$ )  
 **$t\bar{t}$ :** 2 fake  $\tau_h$  that match ( $\Delta R < 0.4$ )  
to top daughters ( $b$  or  $W$  decay  
products)

$\tau_h \tau_\mu$

**H/Z/SUSY/X:** 1 genuine  $\tau_h$  + 1  $e/\mu$  from tau  
decay (opposite charge,  $p_T^{RECO} > 20 \text{ GeV}$ )  
 **$t\bar{t}$ :** 1 fake  $\tau_h$  + 1  $e/\mu$  not from tau decay, both  
that match ( $\Delta R < 0.4$ ) to top daughters ( $b$  or  $W$   
decay products)

## TauProd

Up to 10 decay products from the selected tau pair are included, sorted by  $p_T$ . In the semi-leptonic case, only the  $\tau_h$  decay products are used. Zero-padding ensures fixed input size.

## Jets

Up to 3 leading jets (with  $\Delta R > 0.4$  from the selected taus) are provided as input. If fewer than 3 are found, zero-padding is applied.

## MET

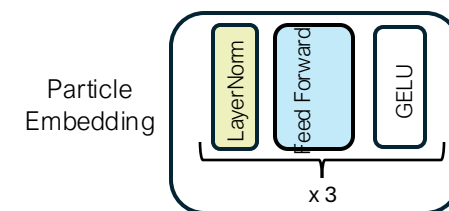
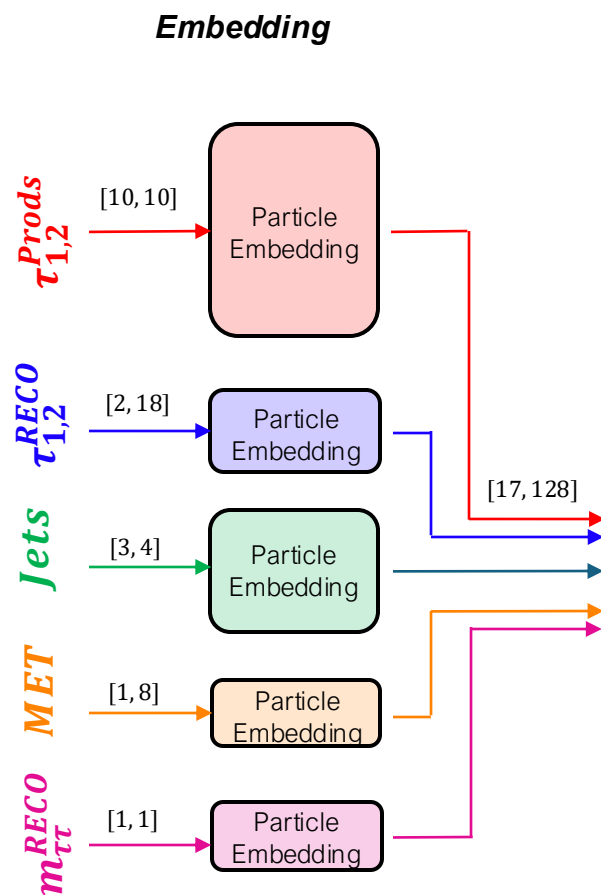
Particle Flow MET is included as a global feature to account for undetected neutrinos.

$m_{\tau\tau}^{RECO}$

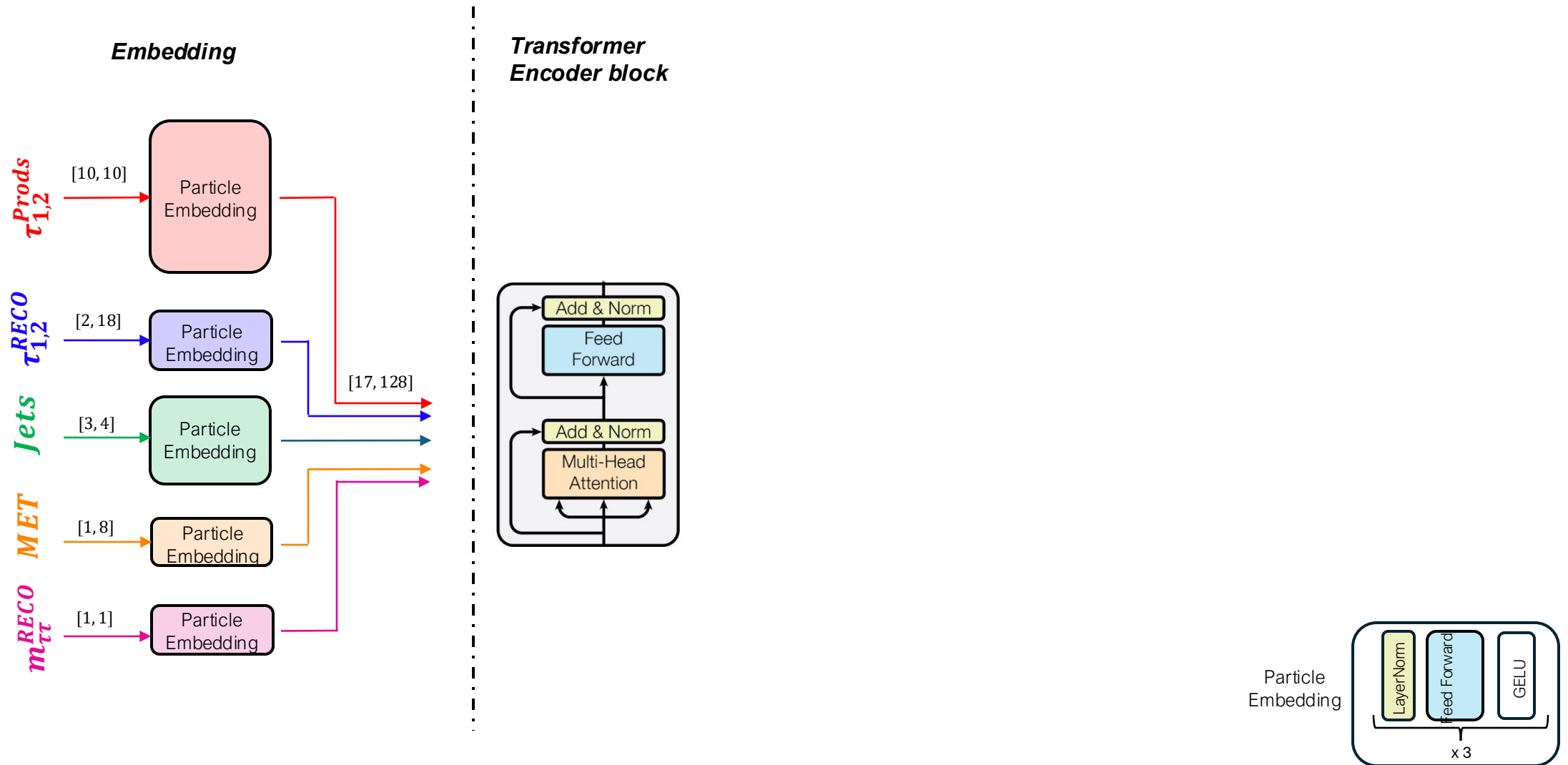
A scalar input with the visible di-tau mass.



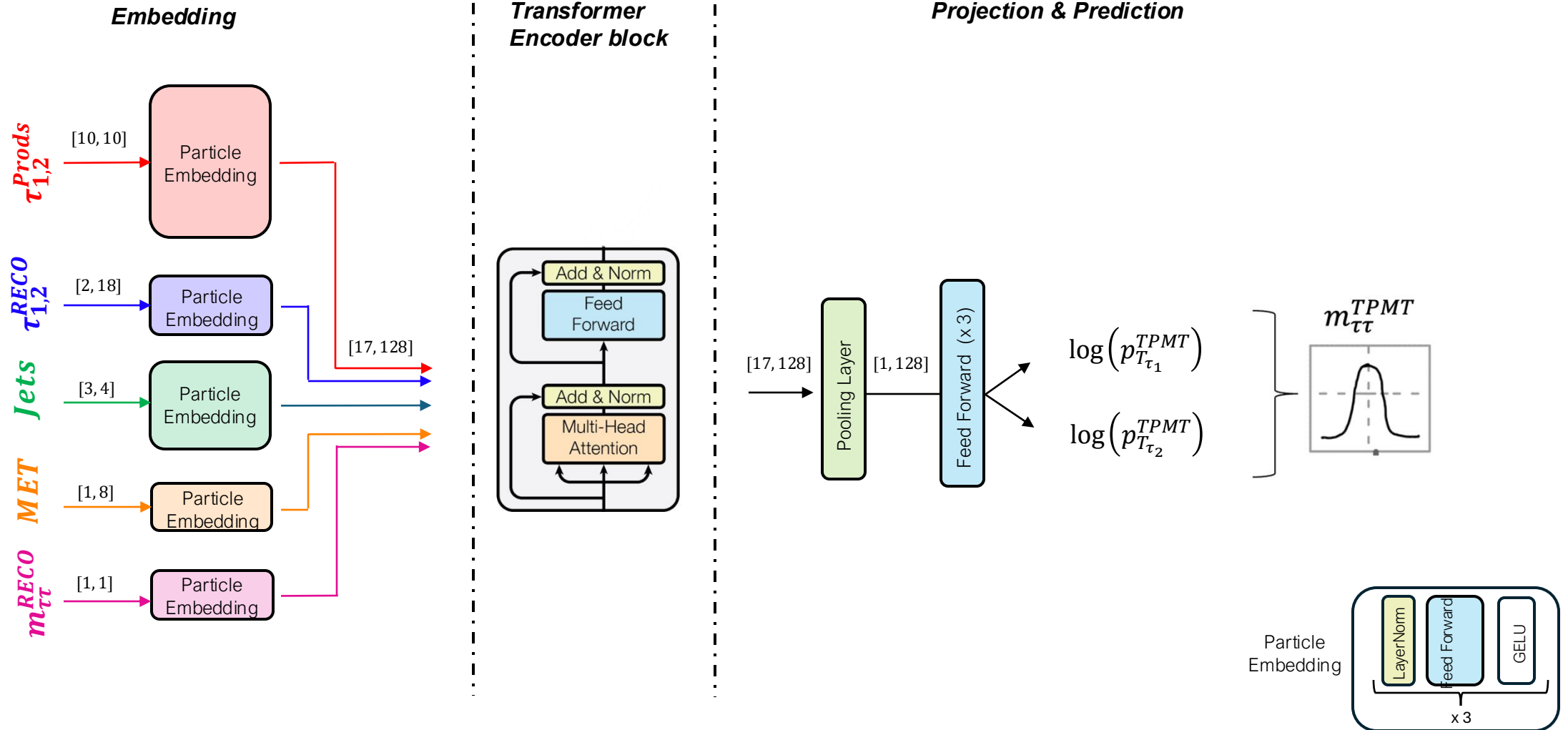
# TPMT Architecture [4]



# TPMT Architecture [4]



# TPMT Architecture [4]



# Training Hyperparameters

- **Batch size:** 1024
- **Initial learning rate:**  $1 \times 10^{-4}$   
*Reduced to  $1 \times 10^{-6}$  using ReduceLROnPlateau*
- **N° Encoder blocks:** 1
- **N° Attention heads:** 8
- **d\_model** (Embedding dimension): 512
- **EarlyStopping patience:** 15 epochs
- **ReduceLROnPlateau patience:** 10 epochs
- **Optimizer:** Adam
- **Trainable parameters:**  $\sim 800,000$
- **Hardware:** NVIDIA Tesla T4 (16GB)
- **Inference time:**  $10^{-3}s$  per event  
(for SVFit,  $\mathcal{O}(s)$  per event)
- **GPU memory usage:**  $\sim 50\%$

*Flat Mass Trainings*

<i>Pairtype</i>	<i>Total Training Events</i>
$\tau_h \tau_h$	$\sim 800$ k

VBF: 400k, GGF: 400k

# Loss Design

To guide the network toward physically meaningful predictions, the loss function is the weighted sum of two terms, both based on the MAE:

$$\mathcal{L}_{\text{total}} = \lambda_{\tau} \cdot \mathcal{L}_{\tau} + \lambda_{m_{\tau\tau}} \cdot \mathcal{L}_{m_{\tau\tau}}$$

- **Tau  $p_T$  loss ( $\mathcal{L}_{\tau}$ ):**

$$MAE(p_T^{TPMT}, p_T^{TARGET}) \text{ for the two taus}$$

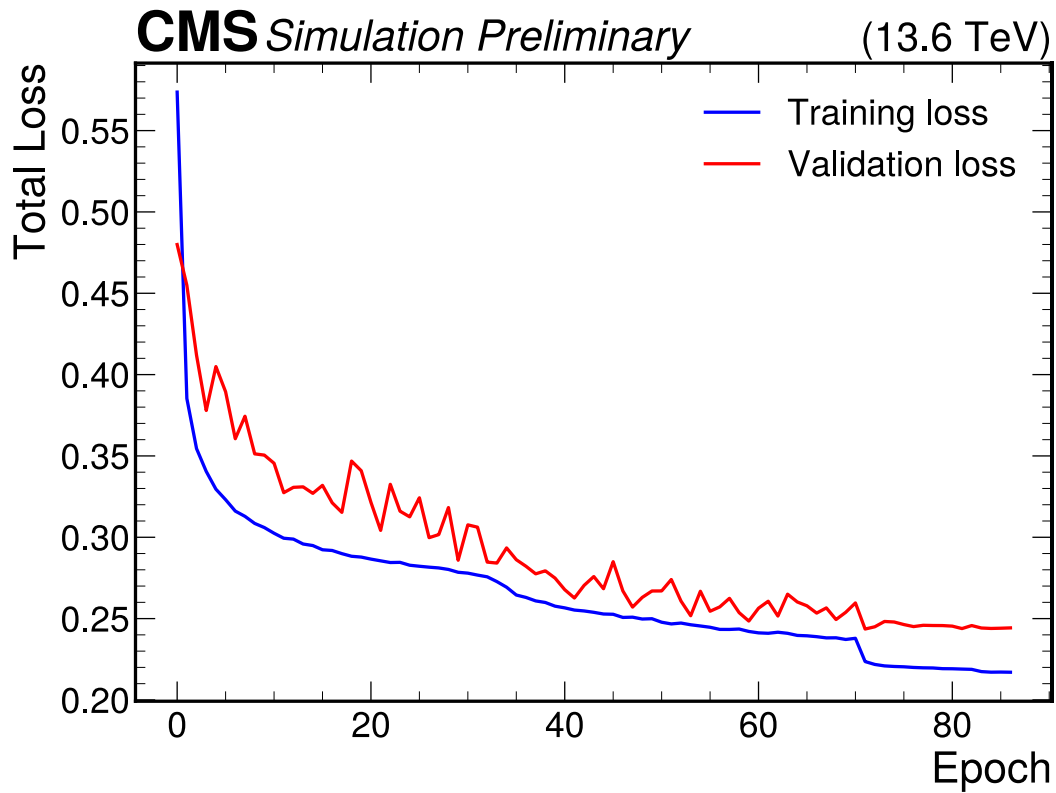
- **Invariant Mass Loss ( $\mathcal{L}_{m_{\tau\tau}}$ ):**

$$MAE(m_{\tau\tau}^{TPMT}, m_{\tau\tau}^{TARGET})$$

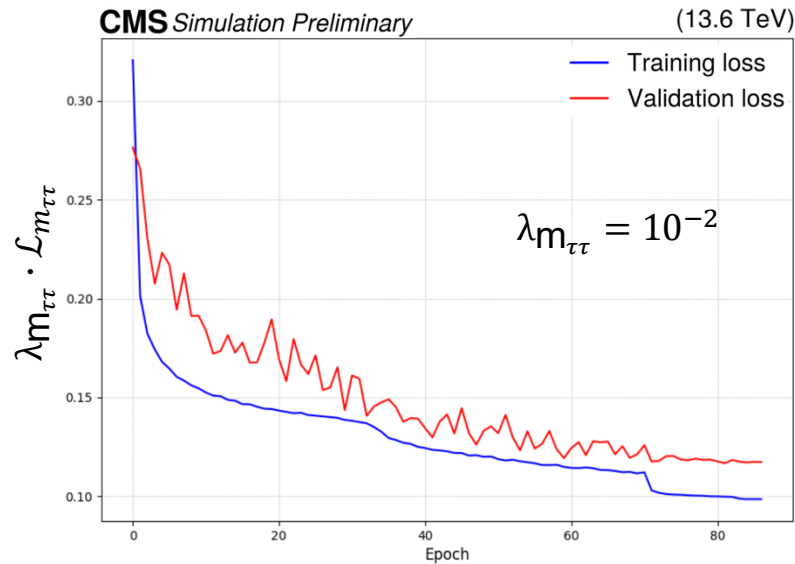
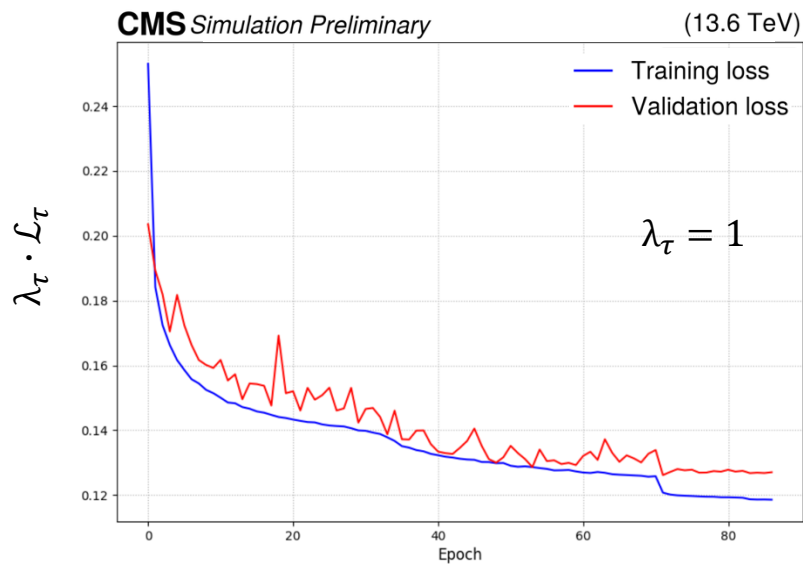
While  $\mathcal{L}_{\tau}$  ensures per-object accuracy, small errors in  $p_T$  can cause large deviations in  $m_{\tau\tau}$  due to its non-linear dependence on kinematics:  $\mathcal{L}_{m_{\tau\tau}}$  term helps correct for this and encourages physically consistent predictions

$$\lambda_{\tau} = 1$$
$$\lambda_{m_{\tau\tau}} = 10^{-2}$$

chosen to  
balance the  
two  
contributions  
to the same  
order of  
magnitude  
during training



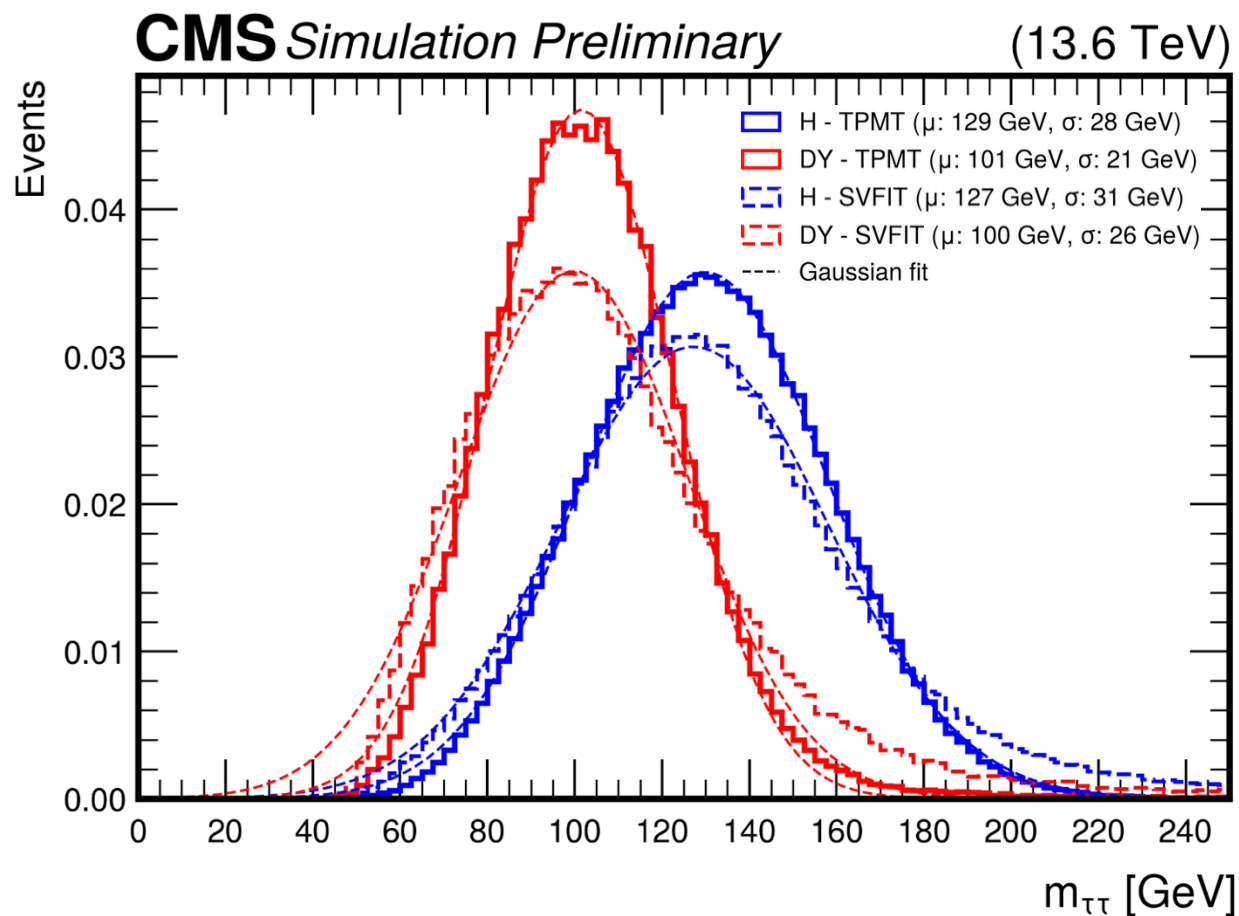
- Rapid decrease in early epochs
- Learning rate drops visible (triggered by ReduceLROnPlateau)
- Losses plateau around epoch  $\sim 75$  ( $lr = 10^{-6}$ )
- Training vs validation: good agreement  $\rightarrow$  good generalization



- Convergence still under study
  - Loss components follow similar trends
- ↓
- smooth and stable total loss

# Results

## $m_{\tau\tau}^{TPMT}$ distribution from Flat Mass Training



**TPMT vs SVFit** comparison on test samples:

➤  $H \rightarrow \tau\tau$  and  $Z \rightarrow \tau\tau$

**Solid lines:** TPMT prediction from predicted  $\log(p_T^{TARGET})$  of the two taus

**Dashed lines:** SVFit  $m_{\tau\tau}$  reconstruction + Gaussian fits

TPMT matches SVFit resolution on  $H \rightarrow \tau\tau$  and  $Z \rightarrow \tau\tau$  (fully hadronic), with  $> 1000x$  faster inference ( $ms$  vs  $s/event$ )

# Second Strategy: Resonant-mass samples

Specialized model, tailored for  $H \rightarrow \tau\tau$ , not general-purpose  $\tau$  reconstruction

Training and evaluation configuration:

- SM-only:
  - $DY \rightarrow \tau\tau$  ( $m_{\tau\tau} > 50 \text{ GeV}$ ),
  - $H \rightarrow \tau\tau$  ( $m_{\tau\tau} = 125 \text{ GeV}$ ),
  - $t\bar{t} \rightarrow (W^+b)(W^-\bar{b})$
- SM+BSM:
  - SM-only processes
  - + SUSY 2HDM signals
  - ( $ggH/bbH$ ,  $m_{\tau\tau} = 350 \text{ GeV}$ )

## Motivation

Resonant samples to better reflect the conditions expected at application time, where a resonance at a fixed mass is assumed



# Second Strategy: Resonant-mass samples

Specialized model, tailored for  $H \rightarrow \tau\tau$ , not general-purpose  $\tau$  reconstruction

**Training and evaluation configuration:**

- SM-only:
    - $DY \rightarrow \tau\tau$  ( $m_{\tau\tau} > 50 \text{ GeV}$ ),
    - $H \rightarrow \tau\tau$  ( $m_{\tau\tau} = 125 \text{ GeV}$ ),
    - $t\bar{t} \rightarrow (W^+b)(W^-\bar{b})$
  - SM+BSM:
    - SM-only processes
    - + SUSY 2HDM signals
    - ( $ggH/bbH$ ,  $m_{\tau\tau} = 350 \text{ GeV}$ )
- ! The choice of the 350 GeV benchmark is not driven by any specific theoretical motivation

## Motivation

Resonant samples to better reflect the conditions expected at application time, where a resonance at a fixed mass is assumed

**Channel-specific trainings** for better adaptation to physics scenario and decay mode

# Second Strategy: Resonant-mass samples

Specialized model, tailored for  $H \rightarrow \tau\tau$ , not general-purpose  $\tau$  reconstruction

## Training and evaluation configuration:

- SM-only:

- $DY \rightarrow \tau\tau$  ( $m_{\tau\tau} > 50 \text{ GeV}$ ),
- $H \rightarrow \tau\tau$  ( $m_{\tau\tau} = 125 \text{ GeV}$ ),

➤  $p_T^{TARGET} = p_T^{MC}$

- $t\bar{t} \rightarrow (W^+b)(W^-\bar{b})$

➤ used as **fake- $\tau$**  background

➤  $p_T^{TARGET} = p_T^{RECO}$

*enables the model to preserve the original  $p_T^{RECO}$  input for fake candidates, ensuring consistent treatment across genuine and misidentified taus*

- SM+BSM:

SM-only processes

+ SUSY 2HDM signals

(**ggH**/**bbH**,  $m_{\tau\tau} = 350 \text{ GeV}$ )

! The choice of the 350 GeV benchmark is not driven by any specific theoretical motivation

## Motivation

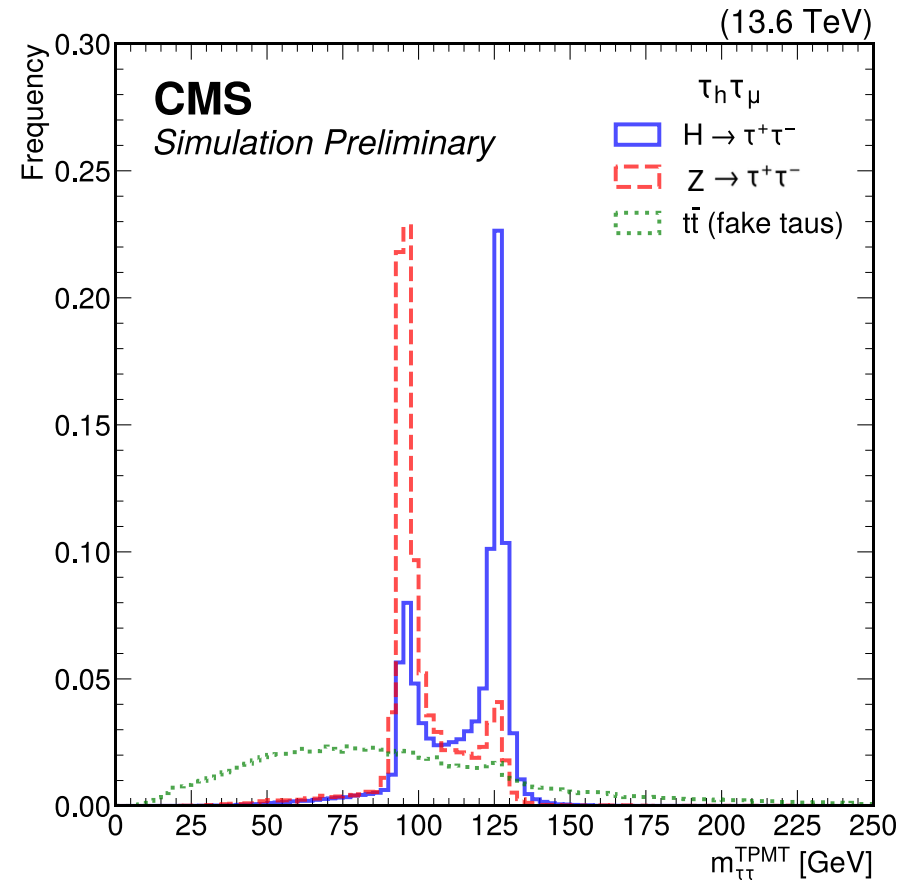
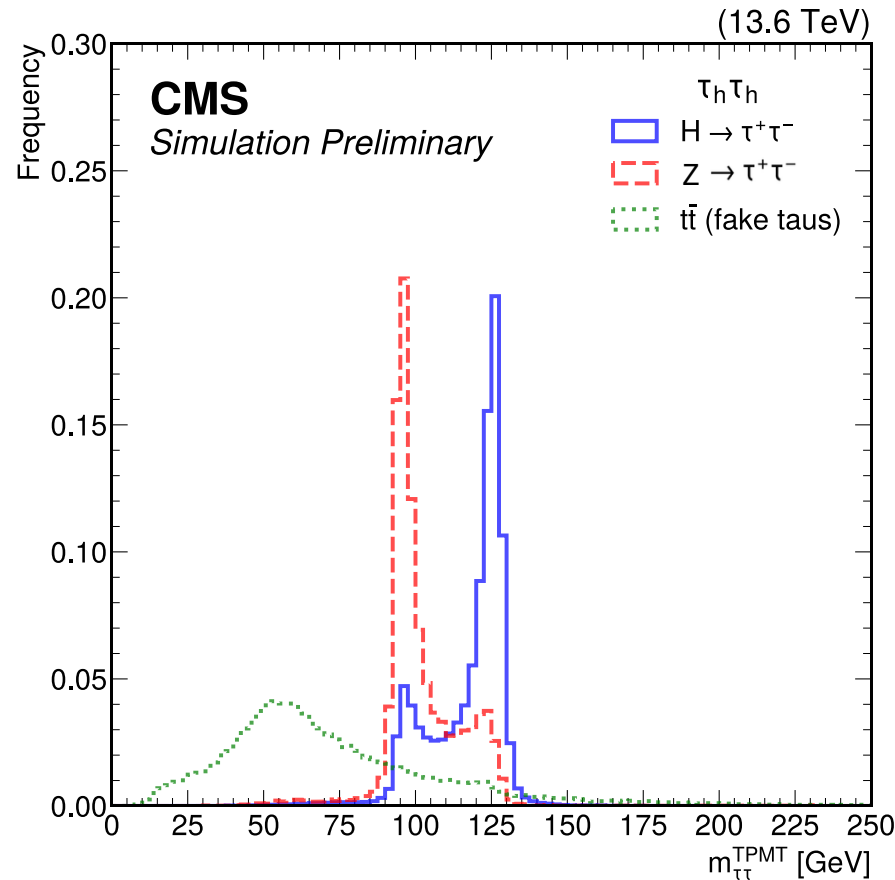
Resonant samples to better reflect the conditions expected at application time, where a resonance at a fixed mass is assumed

**Channel-specific trainings** for better adaptation to physics scenario and decay mode

# Results

$m_{\tau\tau}^{TPMT}$  distribution from **SM-only scenario**  
**Resonance Mass Training**

Differences across final states  $\rightarrow$  under study



*TPMT shows lower  $Z/t\bar{t}$  contamination in a signal-specific mass window designed to retain 90% signal efficiency*

## Two-mode behavior:

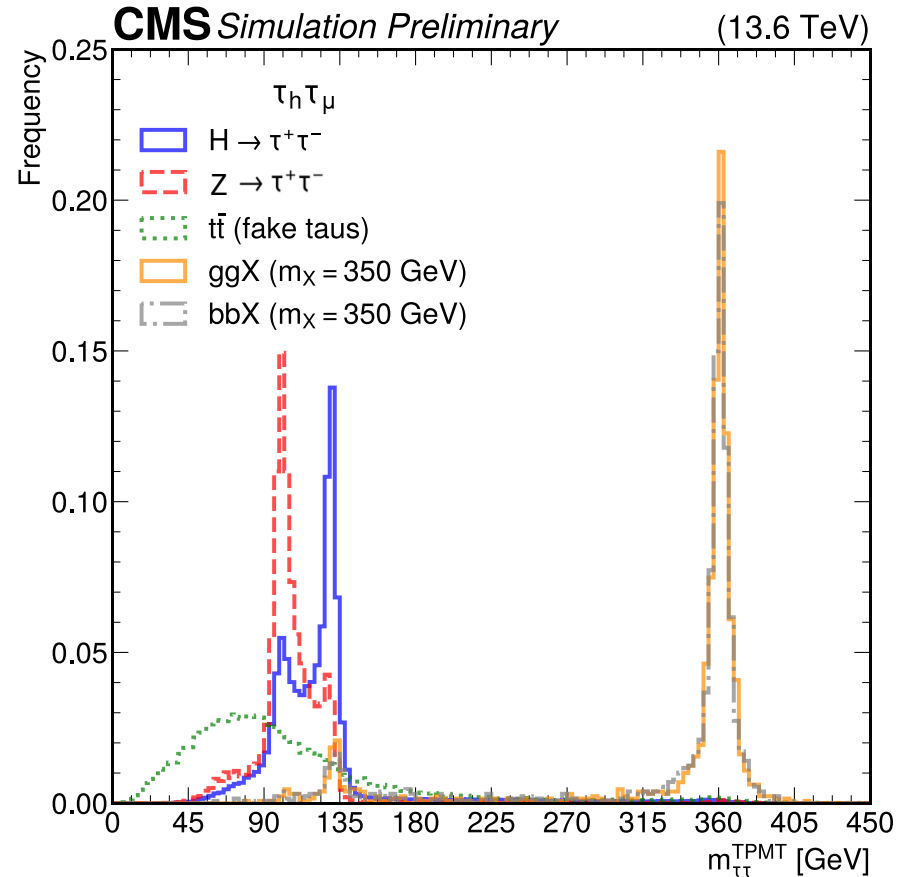
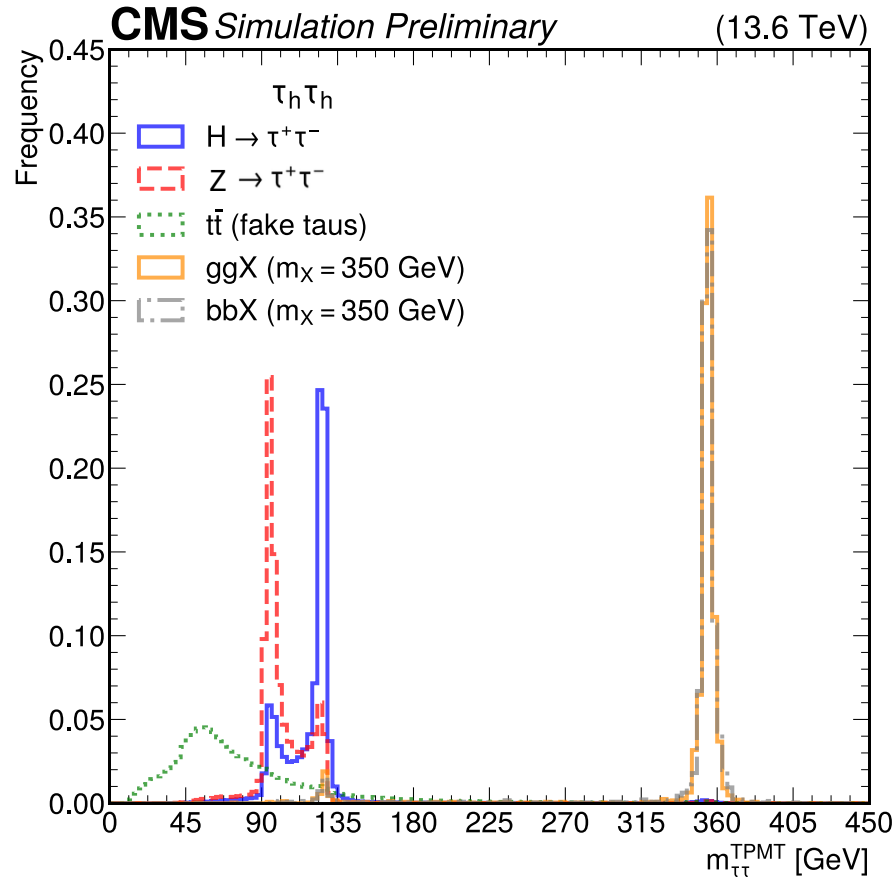
H and Z each show primary + secondary peaks

$\rightarrow$  Model favors either  $\sim 125$  or  $\sim 91$  GeV, reflecting learned resonance structure

# Results

$m_{\tau\tau}^{TPMT}$  distribution from **SM+BSM scenario**  
**Resonance Mass Training**

Differences across final states  $\rightarrow$  under study



**Z vs misreco H  
 (~ 91 GeV):**

Features overlap  
 $\rightarrow$  hard to separate  
 (Random Forest  
 AUC = 0.60)

## Two-mode behavior:

H and Z each show primary + secondary peaks

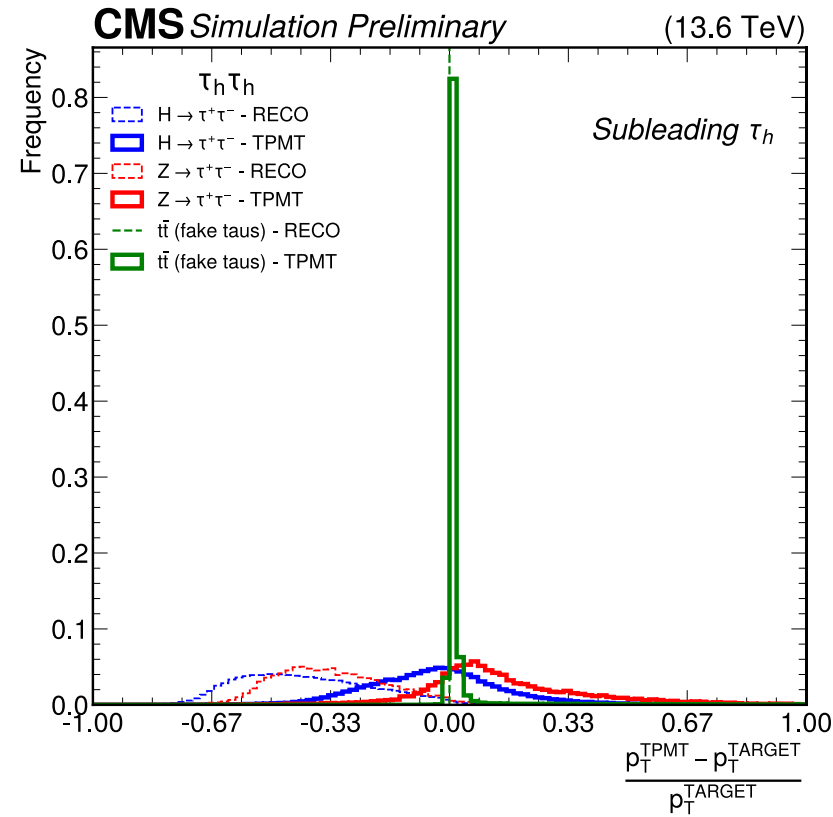
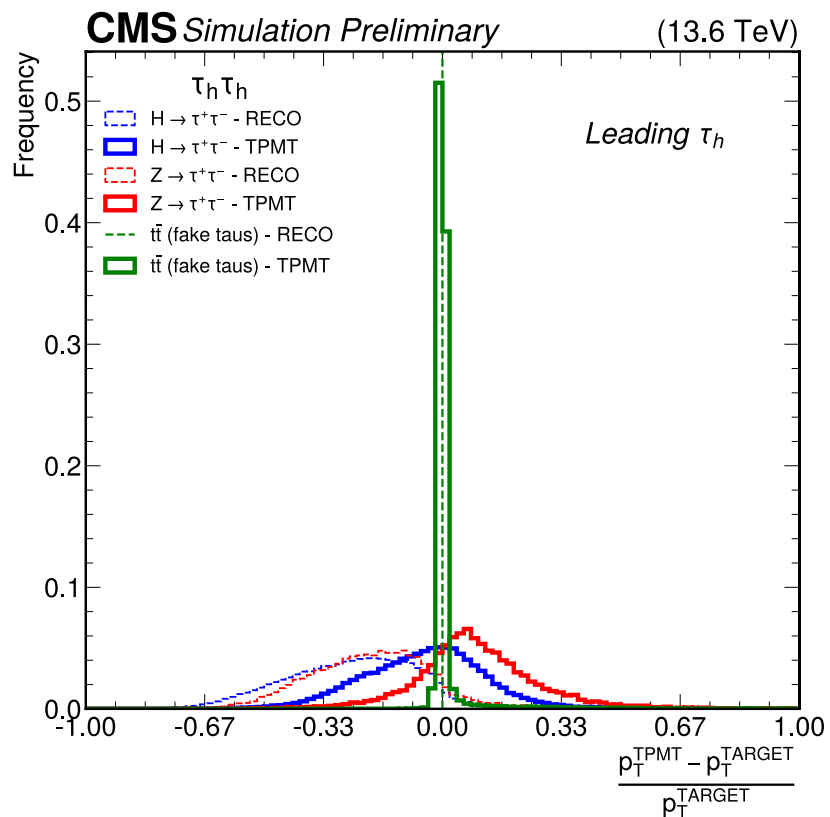
$\rightarrow$  Model favors either  $\sim 125$  or  $\sim 91$  GeV, reflecting learned resonance structure

# Results

## $p_T$ distribution from Resonance Mass Training

SM-only  
scenario

Leading vs subleading  $\tau$ :  
differences still under study



TPMT predictions  
↓  
better match to target  
vs. visible  $p_T^{\text{RECO}}$

- in  $t\bar{t}$  (fakes):
- $p_T^{\text{TARGET}} = p_T^{\text{RECO}}$
  - model learns to keep them unchanged

Broader distributions for  
H and Z  
↓  
residual bias remains

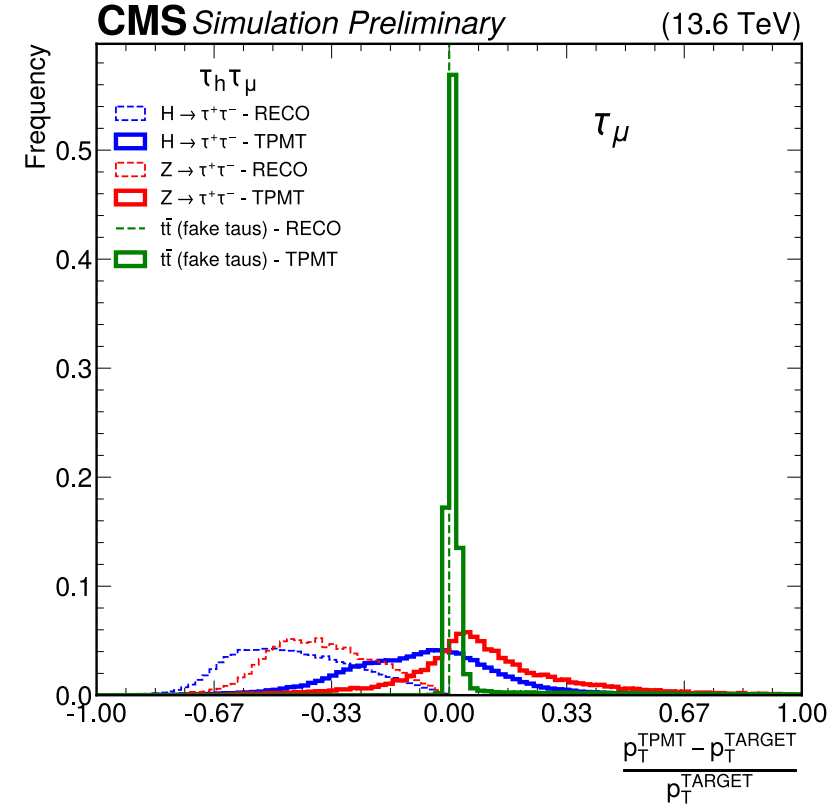
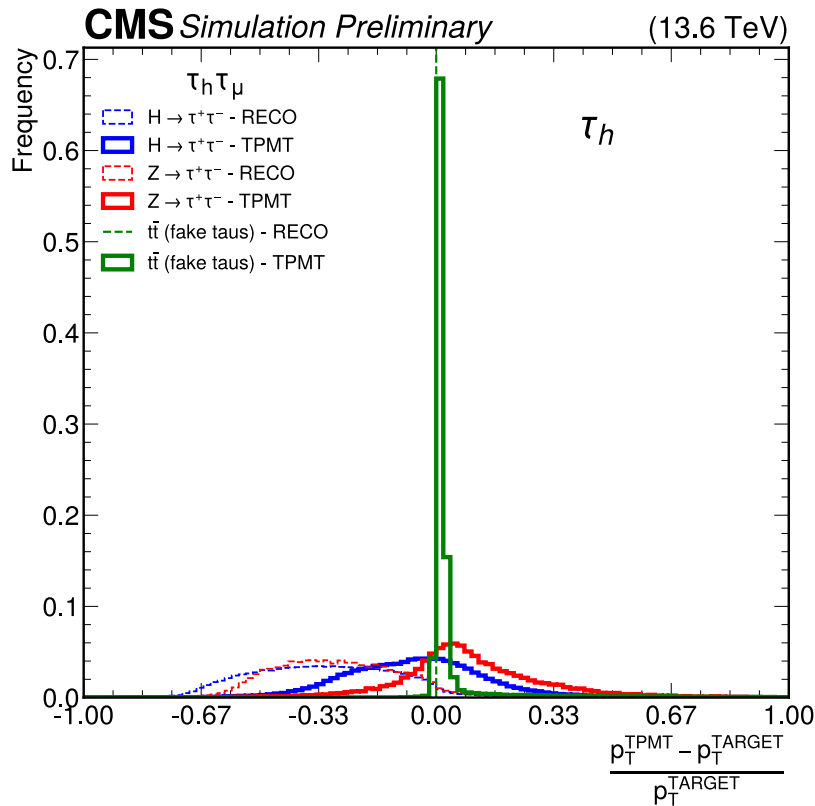
# Results

## $p_T$ distribution from Resonance Mass Training

SM-only  
scenario

$\tau_h \tau_\mu$

$\tau_h$  and  $\tau_\mu$ :  
differences still  
under study



TPMT predictions  
↓  
better match to target  
vs. visible  $p_T^{\text{RECO}}$

- in  $t\bar{t}$  (fakes):
- $p_T^{\text{TARGET}} = p_T^{\text{RECO}}$
  - model learns to keep them unchanged

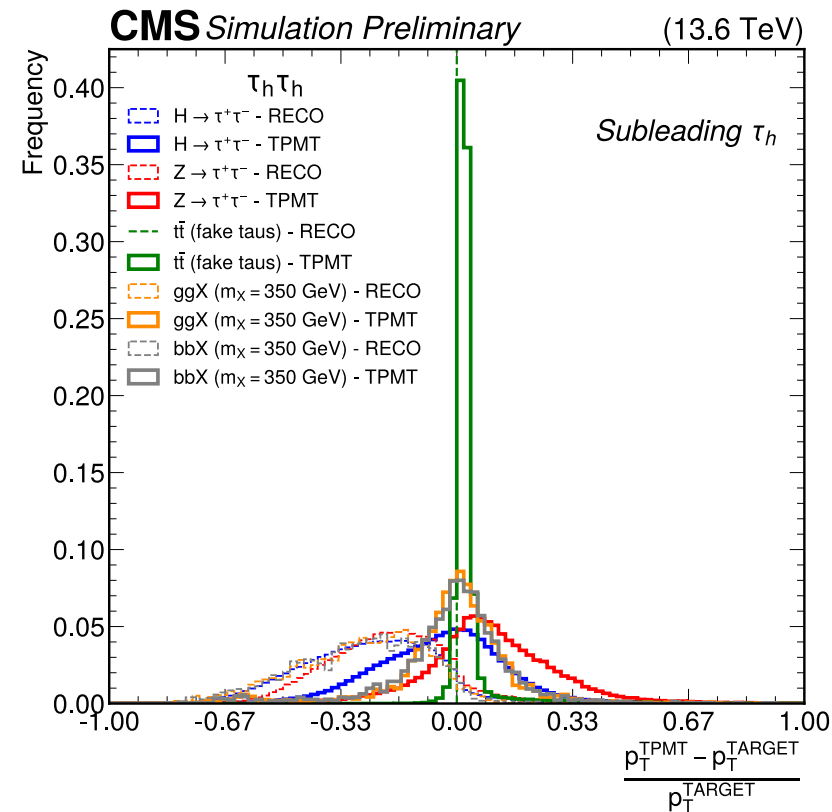
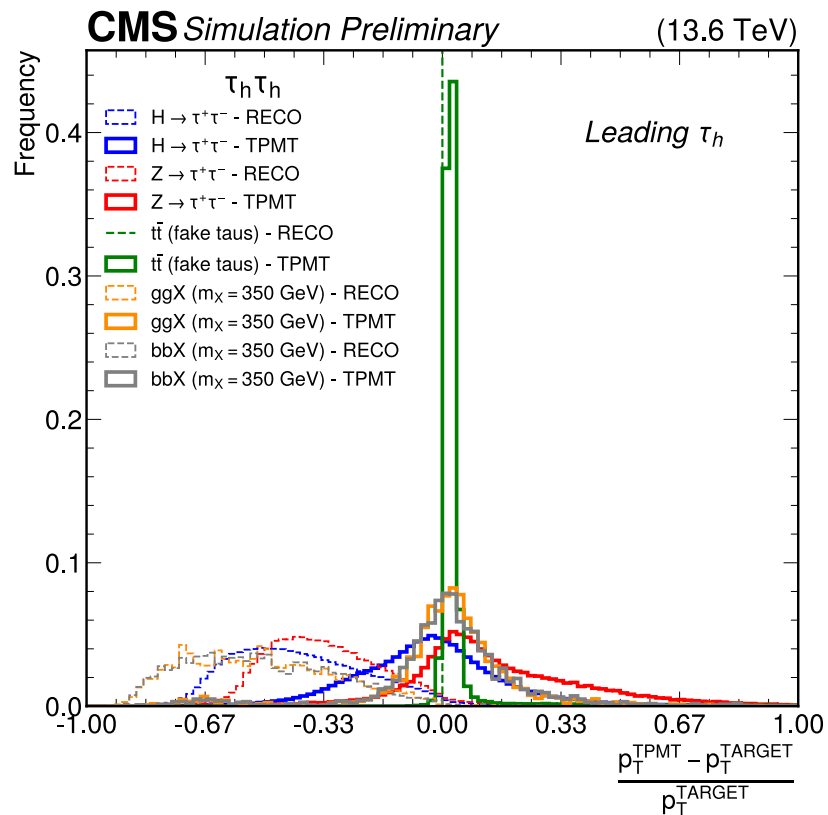
Broader distributions for  
H and Z  
↓  
residual bias remains

# Results

$p_T$  distribution from Resonance  
Mass Training

SM+BSM  
scenario

Leading vs subleading  $\tau$ :  
differences still under study



TPMT predictions  
↓  
better match to target  
vs. visible  $p_T^{\text{RECO}}$

- in  $t\bar{t}$  (fakes):
- $p_T^{\text{TARGET}} = p_T^{\text{RECO}}$
  - model learns to keep them unchanged

Broader distributions for  
H and Z  
↓  
residual bias remains

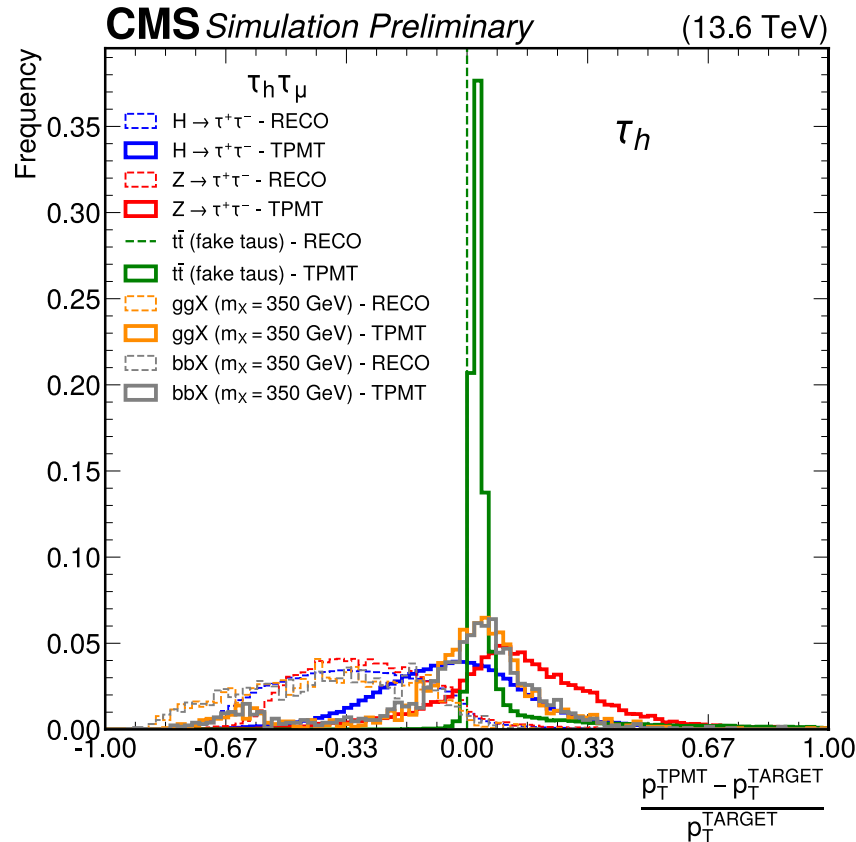
$\tau_h$  and  $\tau_\mu$ :  
differences still  
under study

SM+BSM  
scenario

17

# Results

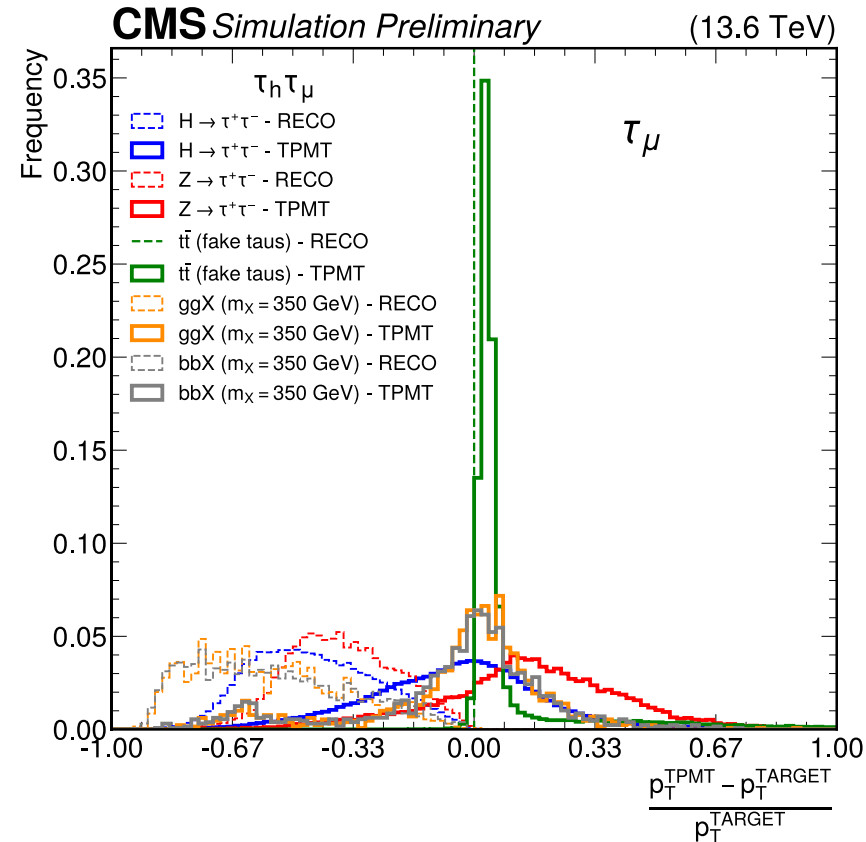
$p_T$  distribution from Resonance  
Mass Training



TPMT predictions



better match to target  
vs. visible  $p_T^{\text{RECO}}$



in  $t\bar{t}$  (fakes):

- $p_T^{\text{TARGET}} = p_T^{\text{RECO}}$
- model learns to keep them unchanged

Broader distributions for  
H and Z



residual bias remains



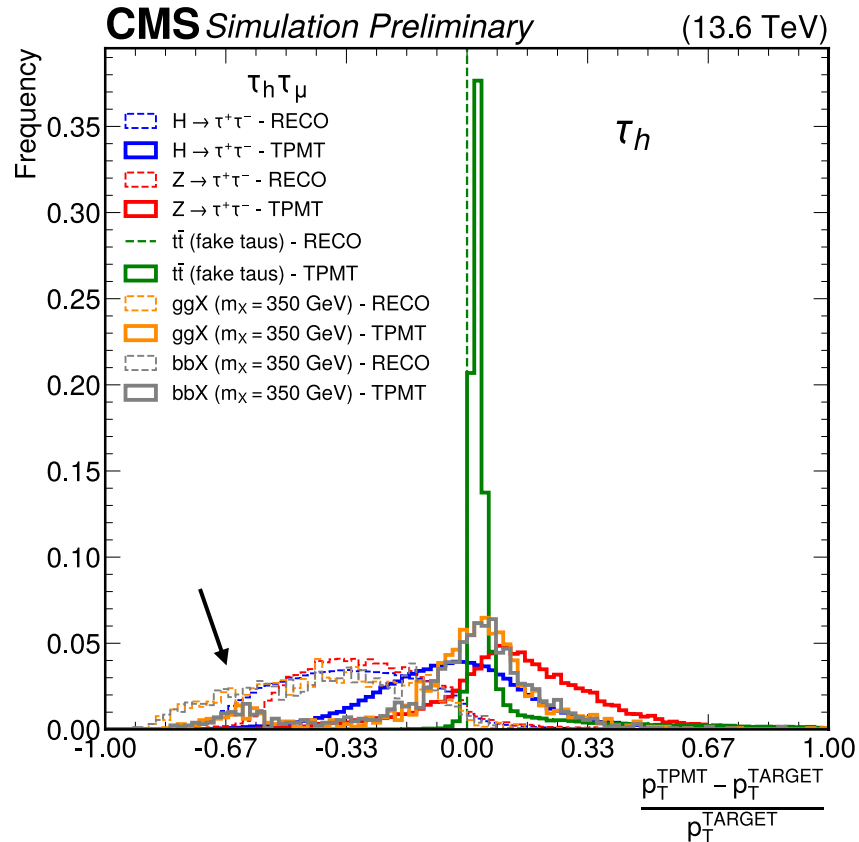
$\tau_h$  and  $\tau_\mu$ :  
differences still  
under study

17

# Results

$p_T$  distribution from Resonance  
Mass Training

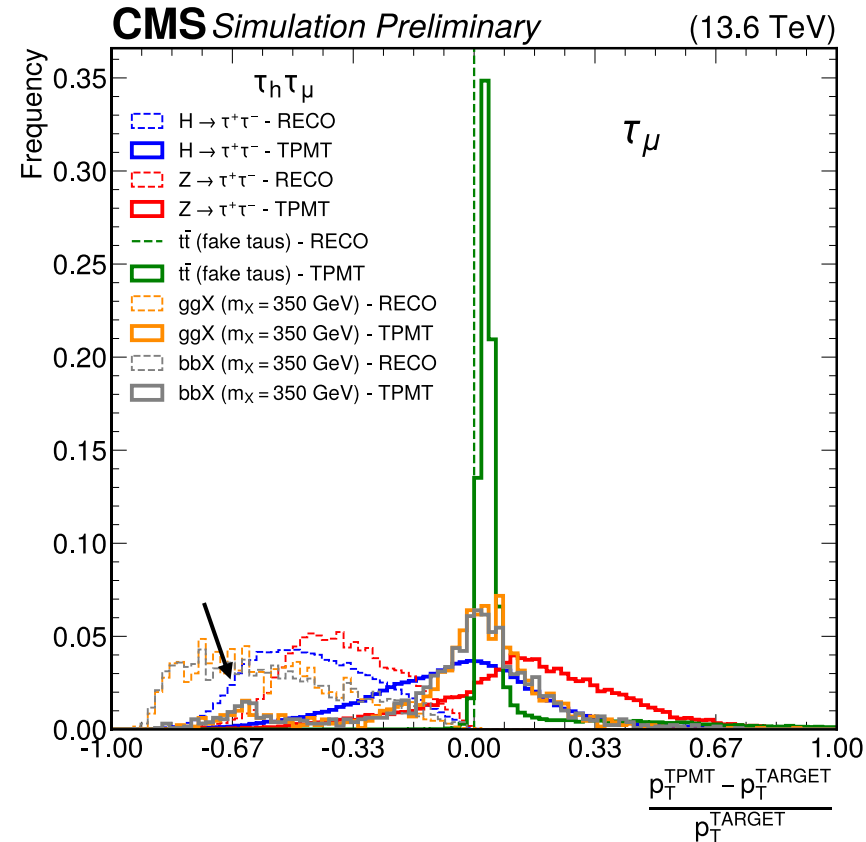
SM+BSM  
scenario



TPMT predictions



better match to target  
vs. visible  $p_T^{\text{RECO}}$



in  $t\bar{t}$  (fakes):

- $p_T^{\text{TARGET}} = p_T^{\text{RECO}}$
- model learns to keep them unchanged

Broader distributions for  
H and Z



residual bias remains

This structure  
reflects a systematic  
underestimation of  
 $p_T$ , causing SUSY  
events to be  
misreconstructed  
near 125 GeV

# Conclusions & Outlook

- Explored a Transformer-based approach for di- $\tau$  identification and kinematic reconstruction in SM and BSM scenarios
- The model shows promising performance, but further refinement is needed.

## Key Observations:

- Binary-like behavior in mass reconstruction, especially H vs Z separation
- Limited smooth interpolation across overlapping mass regions.

## Future developments:

- Introduce a parametric variant (mass hypothesis as input) to improve flexibility
- Focus on better generalization across mass spectra and transition regions

**Thank you for the attention**

## References

- [1] Bianchini, Lorenzo, et al. "Reconstruction of the Higgs mass in  $H \rightarrow \tau\tau$  events by dynamical likelihood techniques." *Journal of Physics: Conference Series*. Vol. 513. No. 2. IOP Publishing, 2014.
- [2] Qu, Huilin, and Loukas Gouskos. "Jet tagging via particle clouds." *Physical Review D* 101.5 (2020): 056019.
- [3] CMS Collaboration. "Reconstruction and identification of  $\tau$  lepton decays to hadrons and  $\nu_\tau$  at CMS. " *Journal of Instrumentation* 11.01 (2016): P01019.
- [4] Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017)
- [5] CMS Collaboration. "Particle-flow reconstruction and global event description with the CMS detector." *Journal of Instrumentation* 12.10 (2017): P10003.
- [6] CMS Collaboration, "Identification of hadronic tau lepton decays using a deep neural network," *Journal of Instrumentation*, 17 (2022): P07023