

# Tau-Lepton Identification and Decay Mode Classification using Graph Neural Networks

Robert McNulty

Supervisors:

*Dr. Nikolaos Rompotis, Prof. Monica D'Onofrio*

With contributions from: *Dr. Joseph Carmignani*



# Introduction

- 2<sup>nd</sup> Year LIV.INNO Student
- Dual Funded PhD – 50/50 split between LIV.INNO and ARO (Industry Placement)
- Physics Project (ATLAS): Tau-Lepton Identification and Decay Mode Classification using Graph Neural Networks
- Industry Project (ARO/UoL): Use of AI models for predictions in different eye-related disease progressions

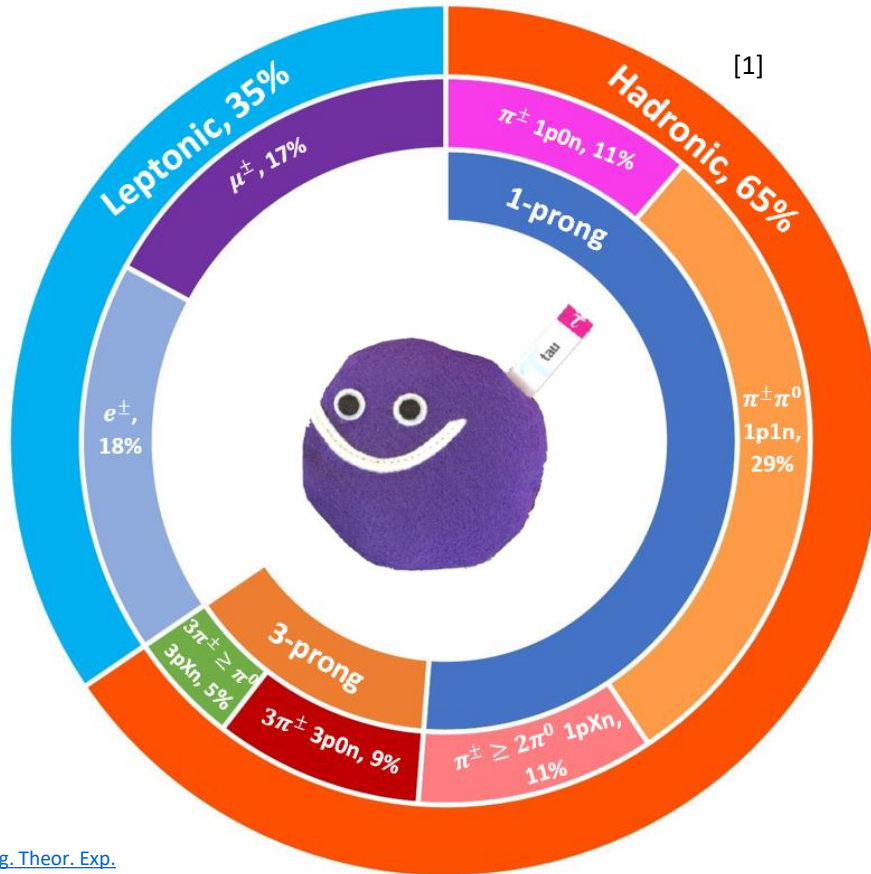


# TauJetGraphs

*GNN Developed by Dr. Joseph Carmignani*

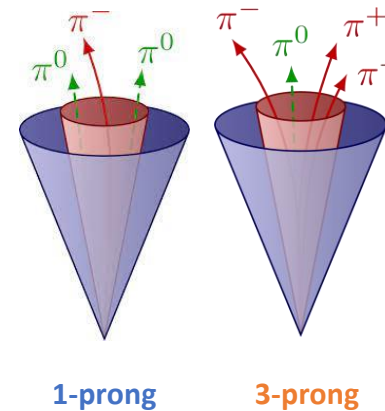
# Tau-Leptons and QCD Dijets

- Leptonic Tau decays,  $\tau_{lep}$ , (35%) produce  $e^\pm, \mu^\pm$  and corresponding  $\nu$ 's
- Hadronic Tau decays,  $\tau_{had}$ , (65%) produce 1 or 3  $\pi^\pm$  (1- & 3-prong decays) and maybe a few  $\pi^0$



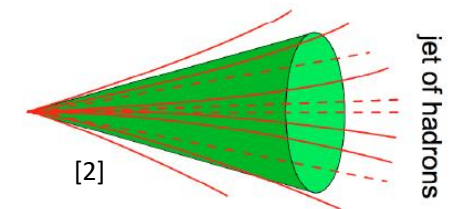
## $\tau_{had}$ Decay Cones:

- Highly collimated – narrow cone
- Small cross-section
- Low multiplicity



## Dijet Production & Cone

- Main background source of fake  $\tau_{had}$  are jets from QCD
- Shower shape can mimic/drown-out the shower shape of  $\pi$ 's from  $\tau_{had}$
- Fragment into multiple hadrons (high multiplicity)
- High production cross-section for dijets
- Wider cone area



[1] R.L. Workman et al. (Particle Data Group), Prog. Theor. Exp. Phys. **2022**, 083C01 (2022) and 2023 update

[2] Joern Mahlstedt and the ATLAS collaboration 2014 J. Phys.: Conf. Ser. **513** 012021, DOI [10.1088/1742-6596/513/1/012021](https://doi.org/10.1088/1742-6596/513/1/012021)

# Motivations and Goals

## Motivations:

- BR for  $\tau_{\text{had}}$  almost twice as much as  $\tau_{\text{lep}}$
- Identification (ID) important for several areas of research, such as:
  - $H \rightarrow \tau\tau$  production [CERN-EP-2021-217](#)
  - Di-Higgs searches with  $b\bar{b}\tau^+\tau^-$

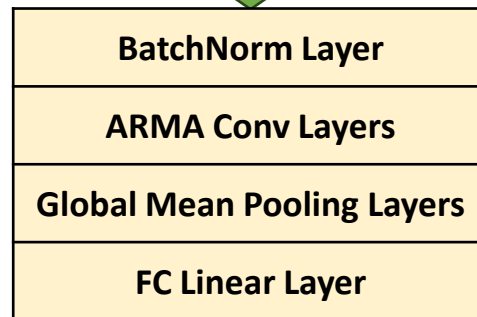
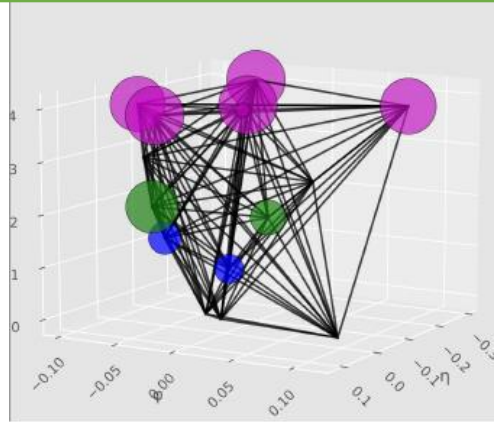
## Goals:

- To study further the unification of the Decay Mode Classification (DMC) and ID Neural Network into a single Graph Neural Network (GNN) algorithm
- It should handle  $\tau_{\text{had}}$ -candidates with 1 & 3 tracks
- The final classifier should also be able to classify 5 decay modes (1p0n, 1p1n, 1pXn, 3p0n, 3pXn) and background QCD jets (dijets)

# TauJetGraphs – Graph & Model Structure

- Each node has 74 attributes – both local (physics object) and global (jet) variables
- For each layer, nodes that are within a predefined distance,  $\Delta R = 0.4$ , are connected by an edge
- Edges are then added between nodes across layers, constructing a 3D graph

**TauJetGraphs NN Model Pipeline\***



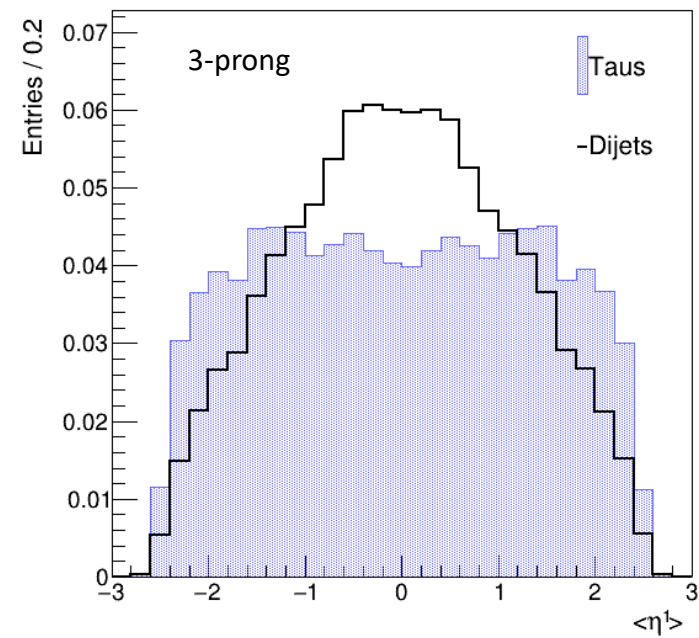
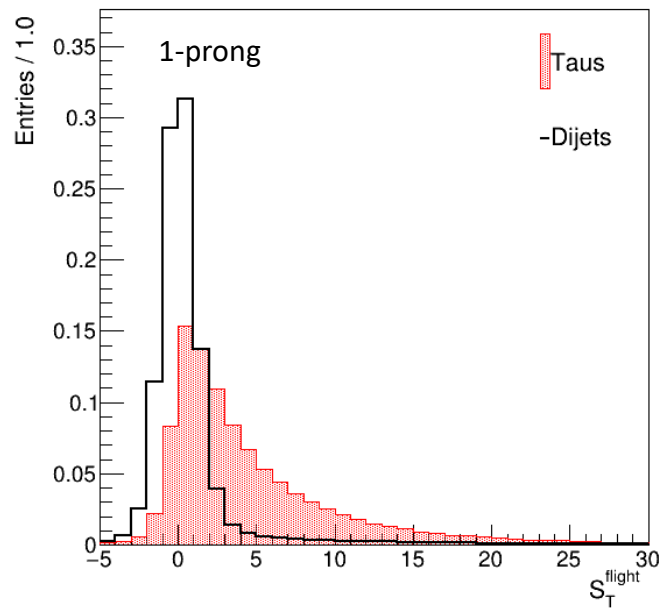
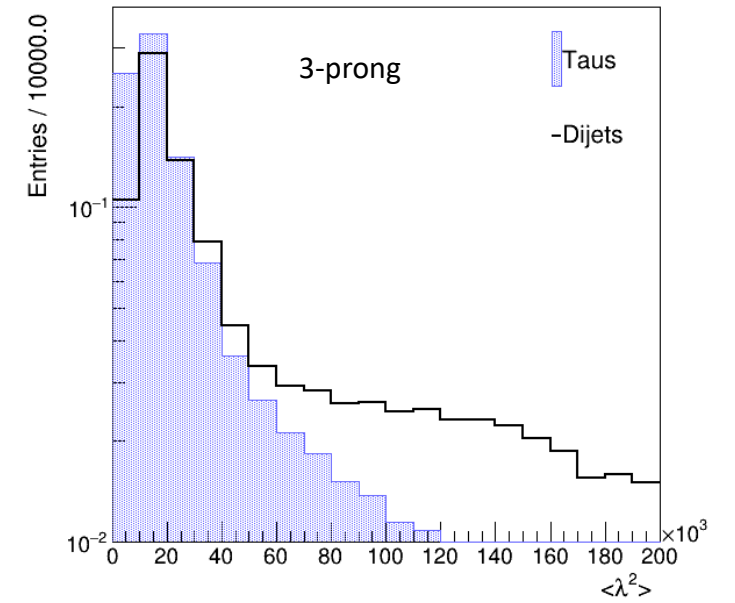
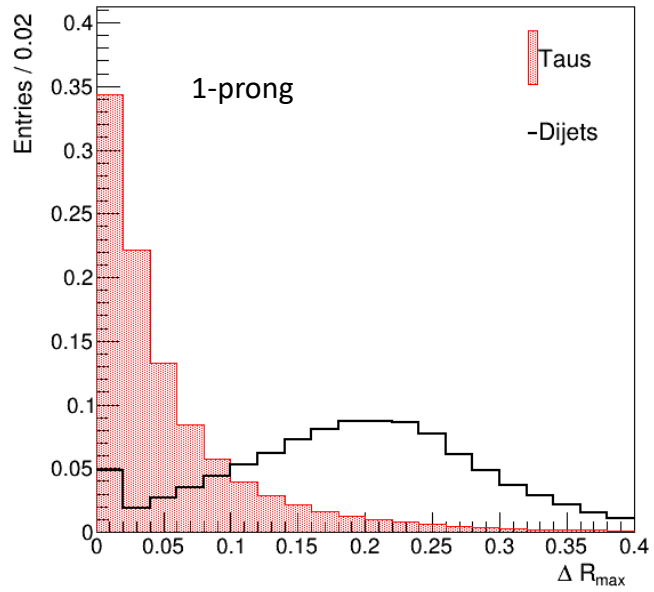
**Output Score**

Layer	Description
0	$h^\pm$ candidates (from $\tau_{\text{had}}$ tracks)
1	$\pi^0$ candidates (from $\tau_{\text{had}}$ decay)
2	$\gamma$ -energy deposits in EM Calorimeters (originating from $\pi^0$ )
3	$e^-e^+$ tracks (from $\gamma \rightarrow e^+e^-$ )
4	Energy deposits, $E_T$ , in the Calorimeter Layers

- Each node represents a reconstructed physics object within a given layer of the graph, i.e. each node in Layer 0 is a reconstructed  $h^\pm$  candidate, and so on

\*  $p_T$  can also be represented as the size of a node

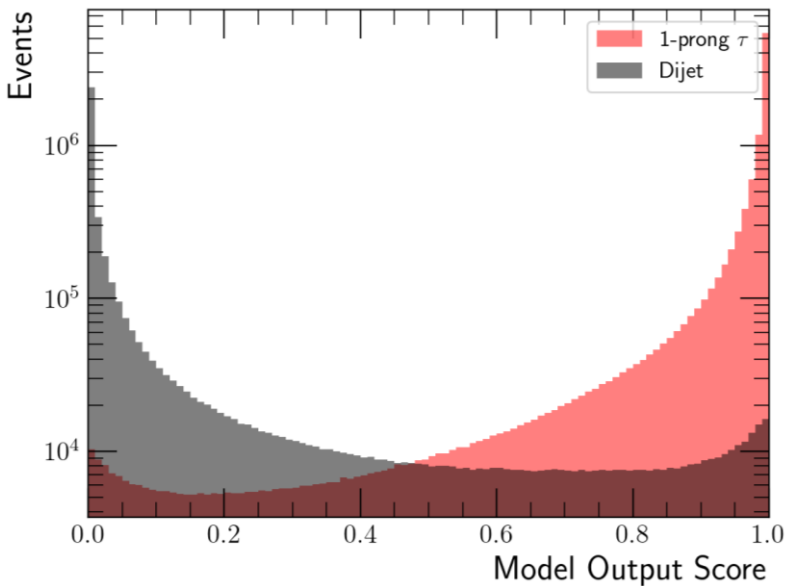
# Some Example Distributions



# Results (Identification)

- 3-prong performance shows similar performance to 1-prong
- Model shows good general separation between true and fake  $\tau_{\text{had}}$  candidates

1-prong GNN Output Scores



- Rejection = inverse of background selection efficiency
- Receiver Operating Characteristic (ROC) Curve displays the rejection of misidentified background samples as a function of the signal efficiency

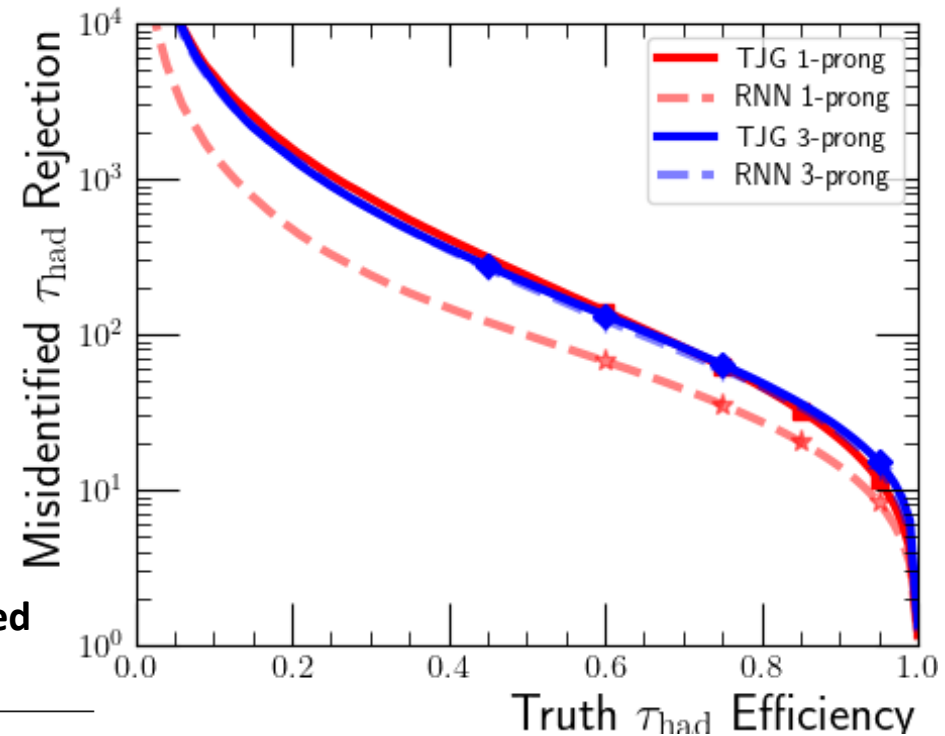
**At 60% Efficiency**

- 1-prong rejection improves in GNN by order of 10
- 3-prong rejection has some but no significant improvement

**TauID GNN Background Rejection at Specified Efficiencies**

	1-prong				3-prong			
Efficiency	60 %	75 %	85 %	95 %	45 %	60 %	75 %	95 %
Rejection	140.96	62.24	32.40	11.62	275.47	132.81	63.95	15.06

**TauID GNN vs RNN ROC Curve**





# Results (Decay Mode Classification)

- Results shown are for 1-prong at 75% Efficiency
- Efficiency matrix: each column sums to 100%
- Similar performance for 1-prong candidates at same efficiency
- **Efficiency (Recall)** =  $\frac{tp}{tp+fn}$ 
  - $tp$  = true positive (signal sample identified as signal)
  - $fn$  = false negative (signal sample identified as background)

GNN (1-prong)

GNN tau decay mode	1pXn	0.54	6.91	45.46
	1p1n	17.15	84.86	52.61
	1p0n	82.31	8.23	1.93
		1p0n	1p1n	1pXn
		Truth tau decay mode		

Current Method - DSNN

DeepSet NN tau decay mode	3pXn	0.02	0.13	0.21	3.46	69.87
	3p0n	0.39	0.11	0.05	96.47	29.54
	1pXn	0.50	5.93	55.25	0.00	0.11
	1p1n	10.14	86.01	42.85	0.05	0.47
	1p0n	88.95	7.82	1.63	0.01	0.01
		1p0n	1p1n	1pXn	3p0n	3pXn
		Truth tau decay mode				

[ATL-PHYS-PUB-2022-044](#)

# Industry Placement and Upcoming Project

# Work with ARO and New Project



- Role with ARO is Research Developer.
- This has required working with ARO's clients on their projects.
- As such, I have received training e.g. in the use of the following:
  - PRTG Dashboards
  - REDCap Databases
  - Enovacom Integration Engine
- N.B. the project involving Enovacom has since moved onto using Apache Airflow software.

## Upcoming Project

with Dr. Philip Burgess and Prof. Yalin Zheng

- Use of an AI model for prediction of progression of Age-related Macular Degeneration<sup>1</sup>. [3, 4]
- Validation of AI model for progression of diabetic retinopathy<sup>2</sup> to treatment.

<sup>1</sup>Age-related macular degeneration (AMD) is a common condition that affects the middle part of your vision. More information can be found here on the [NHS Website](#).

<sup>2</sup>Diabetic retinopathy is a complication of diabetes, caused by high blood sugar levels damaging the back of the eye (retina). This can cause blindness if left untreated. More information can be found here on the [NHS Website](#).

[3] Bridge J, Harding S, Zheng Y. Development and validation of a novel prognostic model for predicting AMD progression using longitudinal fundus images. *BMJ Open Ophthalmology* 2020;5:e000569. doi: [10.1136/bmjophth-2020-000569](https://doi.org/10.1136/bmjophth-2020-000569)

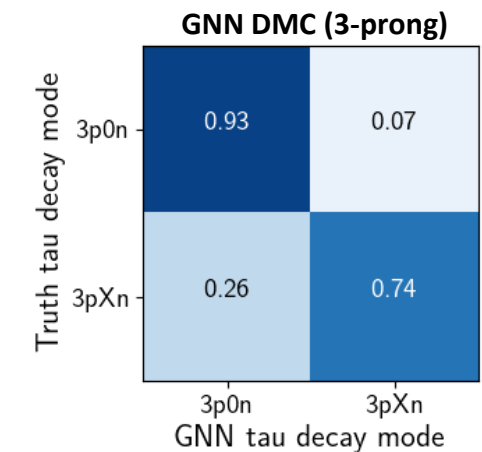
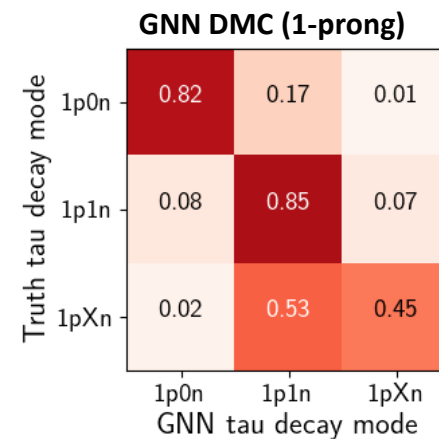
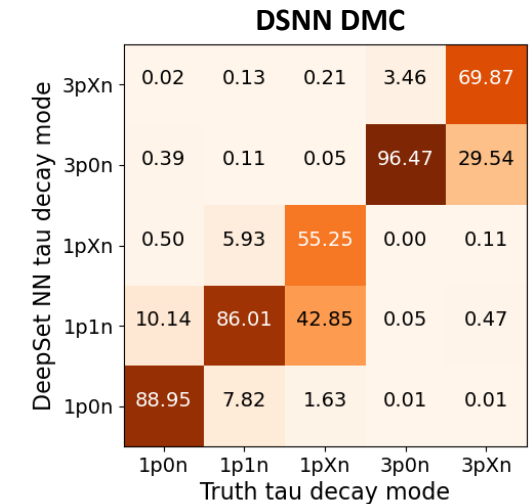
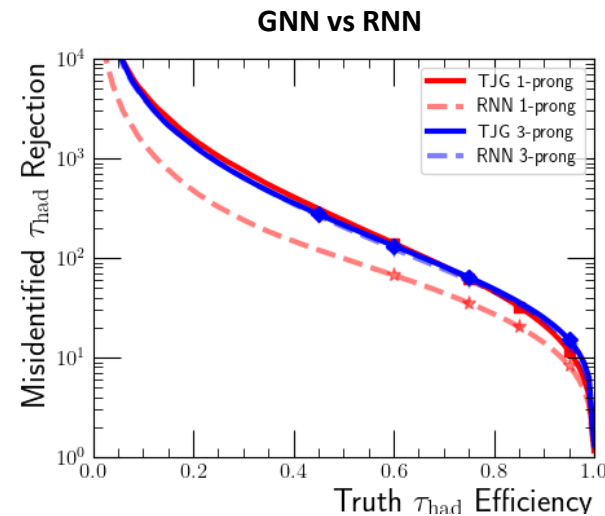
[4] Bridge, J., Harding, S., Zheng, Y. (2021). End-to-End Deep Learning Vector Autoregressive Prognostic Models to Predict Disease Progression with Uneven Time Intervals. In: Papież, B.W., Yaqub, M., Jiao, J., Namburete, A.I.L., Noble, J.A. (eds) *Medical Image Understanding and Analysis*. MIUA 2021. Lecture Notes in Computer Science(), vol 12722. Springer, Cham. [https://doi.org/10.1007/978-3-030-80432-9\\_38](https://doi.org/10.1007/978-3-030-80432-9_38)

# Summary

# Summary and Future Steps

- TauJetGraphs progressing nicely:
  - Showing improvements over current methods
  - Further analysis currently taking place (see Mehul's 1<sup>st</sup> Year talk, as well as Monica's and Jordy's ATLAS talks)
- Currently writing up TauJetGraphs work for thesis
- Work at ARO proving to be good work/industrial experience
- Project on the use of AI models for predictions in different eye-related disease progressions to start soon

	1-prong				3-prong			
Efficiency	60%	75%	85%	95%	45%	60%	75%	95%
Rejection	140.96	62.24	32.40	11.62	275.47	132.81	63.95	15.06



Thank you

# Backup

# Backup: Input Variables



Inputs used for TauID RNN, [ATL-PHYS-PUB-2019-033](#),  
also used in TauJetGraphs GNN

Track inputs			Cluster inputs			High-level inputs		
Observable	1-prong	3-prong	Observable	1-prong	3-prong	Observable	1-prong	3-prong
$p_T^{\text{seed jet}}$	•	•	$p_T^{\text{jet seed}}$	•	•	$p_T^{\text{uncalibrated}}$	•	•
$p_T^{\text{track}}$	•	•	$E_T^{\text{cluster}}$	•	•	$f_{\text{cent}}$	•	•
$\Delta\eta^{\text{track}}$	•	•	$\Delta\eta^{\text{cluster}}$	•	•	$f_{\text{leadtrack}}^{-1}$	•	•
$\Delta\phi^{\text{track}}$	•	•	$\Delta\phi^{\text{cluster}}$	•	•	$\Delta R_{\text{max}}$	•	•
$ d_0^{\text{track}} $	•	•	$\lambda_{\text{cluster}}$	•	•	$ S_{\text{leadtrack}} $	•	
$ z_0^{\text{track}} \sin \theta $	•	•	$\langle \lambda_{\text{cluster}}^2 \rangle$	•	•	$S_T^{\text{flight}}$		•
$N_{\text{IBL hits}}$	•	•	$\langle r_{\text{cluster}}^2 \rangle$	•	•	$f_{\text{track}}^{\text{iso}}$	•	•
$N_{\text{Pixel hits}}$	•	•				$f_{\text{track}}^{\text{EM}}$	•	•
$N_{\text{SCT hits}}$	•	•				$p_T^{\text{EM+track}} / p_T$	•	•
						$m^{\text{EM+track}}$	•	•
						$m^{\text{track}}$		•

Variable	Description
$p_T(\tau_{\text{had}})$	$p_T$ of the $\tau_{\text{had}}$ (using calorimeter based $\tau_{\text{had-vis}}$ energy scale)
$p_T(\text{object})$	$p_T$ of the object
$\Delta\phi(\text{object}, \tau_{\text{had}})$	Distance between the object and $\tau_{\text{had}}$ in $\phi$
$\Delta\eta(\text{object}, \tau_{\text{had}})$	Distance between the object and $\tau_{\text{had}}$ in $\eta$
$\Delta\phi(\text{object}, \text{trackECal})$	Distance between the object and the extrapolation of highest- $p_T$ $\tau_{\text{had}}$ track to EM calorimeter in $\phi$
$\Delta\eta(\text{object}, \text{trackECal})$	Distance between the object and the extrapolation of highest- $p_T$ $\tau_{\text{had}}$ track to EM calorimeter in $\eta$
$\langle\eta^1\rangle$	First moment in $\eta$ in cluster shower axis
$\log(\langle r^2\rangle)$	Second moment in the radial distance of cluster cells from the shower axis
$\Delta\theta$	Distance in $\theta$ between the EM shower axis and the vector pointing from the primary vertex to the centre of the shower
$\log(\lambda_{\text{centre}})$	Distance of the cluster shower centre from the calorimeter front face measured along the shower axis
$\langle\lambda^2\rangle$	Mean distance of a cell from the shower centre along the shower axis
$\log(\langle\rho^2\rangle)$	Second moment in the cluster energy density, where $\rho = E^{\text{cluster}}/V^{\text{cluster}}$
$f_{\text{core}}$	Sum of energy fractions in the most energetic cells per sampling
$f_{\text{core}}^{\text{EM1}}$	Same as $f_{\text{core}}$ but only consider EM1
$N_{\text{pos,EM1}}$	Number of cells with positive energy in EM1
$N_{\text{pos,EM2}}$	Number of cells with positive energy in EM2
$E_{\text{EM1}}$	Energy in the EM1 layer
$E_{\text{EM2}}$	Energy in the EM2 layer
$\langle\eta_{\text{EM1}}^1\rangle$ w.r.t. cluster	First moment in $\eta$ in EM1 with respect to the cluster
$\langle\eta_{\text{EM2}}^1\rangle$ w.r.t. cluster	First moment in $\eta$ in EM2 with respect to the cluster
$\log(\langle\eta_{\text{EM1}}^2\rangle)$ w.r.t. cluster	Second moment in $\eta$ in EM1 with respect to the cluster
$\log(\langle\eta_{\text{EM2}}^2\rangle)$ w.r.t. cluster	Second moment in $\eta$ in EM2 with respect to the cluster

Physics object kinematic variables

Variables used in Decay Mode Classification DSNN, [ATL-PHYS-PUB-2022-044](https://arxiv.org/abs/2204.044), also utilised in TauJetGraphs GNN

Neutral pion cluster variables

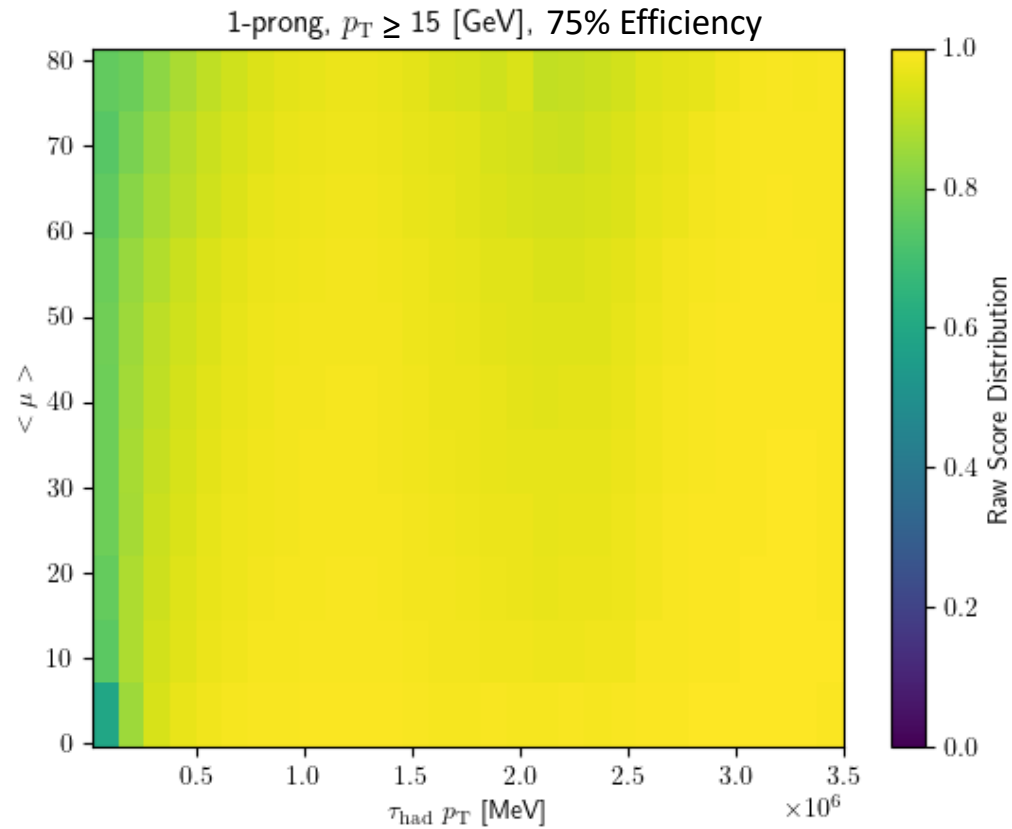
# Backup: Score Bias

# Accounting for Bias (1)

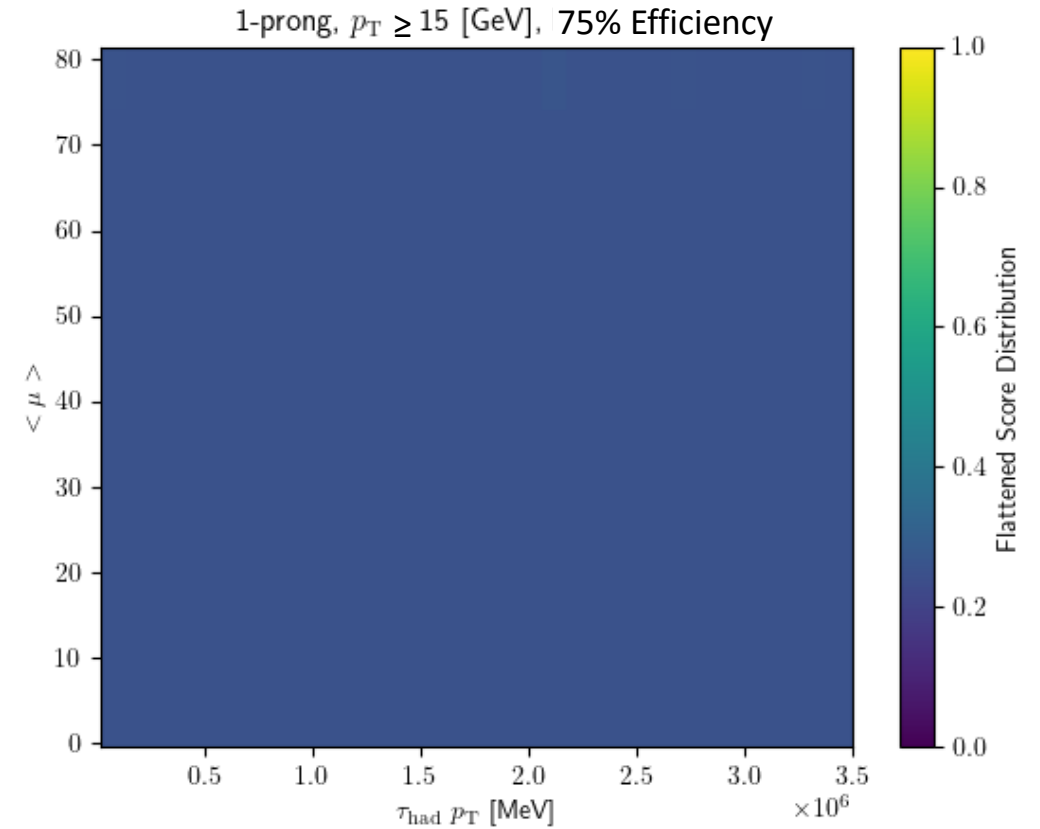
- Performance is typically measured through efficiency, purity, and background rejection
- Observations have shown that classification methods typically favour high- $p_T$  events for a higher score than for those with a lower  $p_T$
- This bias is accounted for by a transformation on the output scores for a given efficiency:
  - A 2D histogram of transverse momentum,  $p_T$ , and the average number of interactions per bunch-crossing,  $\mu$ , is created from the signal dataset (slide [17](#))
  - The top percentage (given by the desired efficiency) of samples in each bin is taken, and the score threshold of this bin is determined by the sample with the lowest score in this selection.
  - The scores of each sample is then transformed w.r.t. the desired efficiency.
  - Background samples are then treated in a similar fashion, however the bin threshold scores from the signal histogram are used as the threshold for the background histogram
    - Samples which pass this selection are labelled as signal in the predictions for the given efficiency (referred to as “False Positives”)
    - From here, the background rejection for the given efficiency is determined

# Accounting for Bias (2)

Signal sample score distribution before transformation

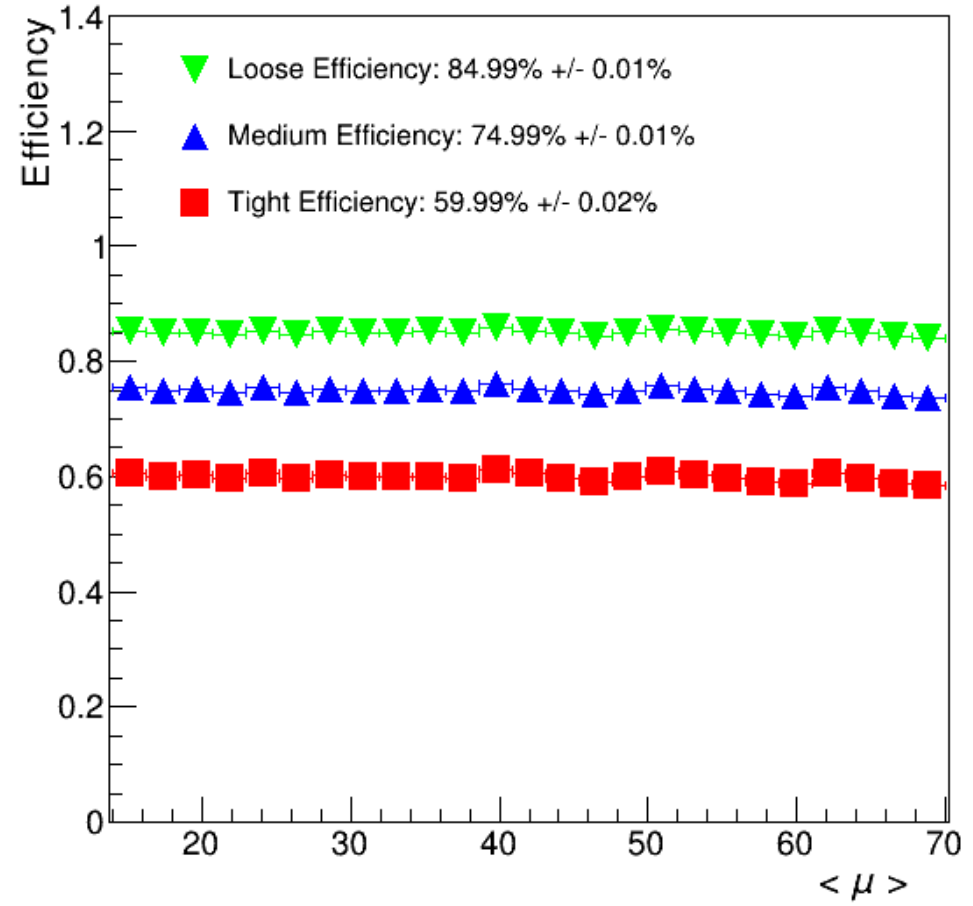
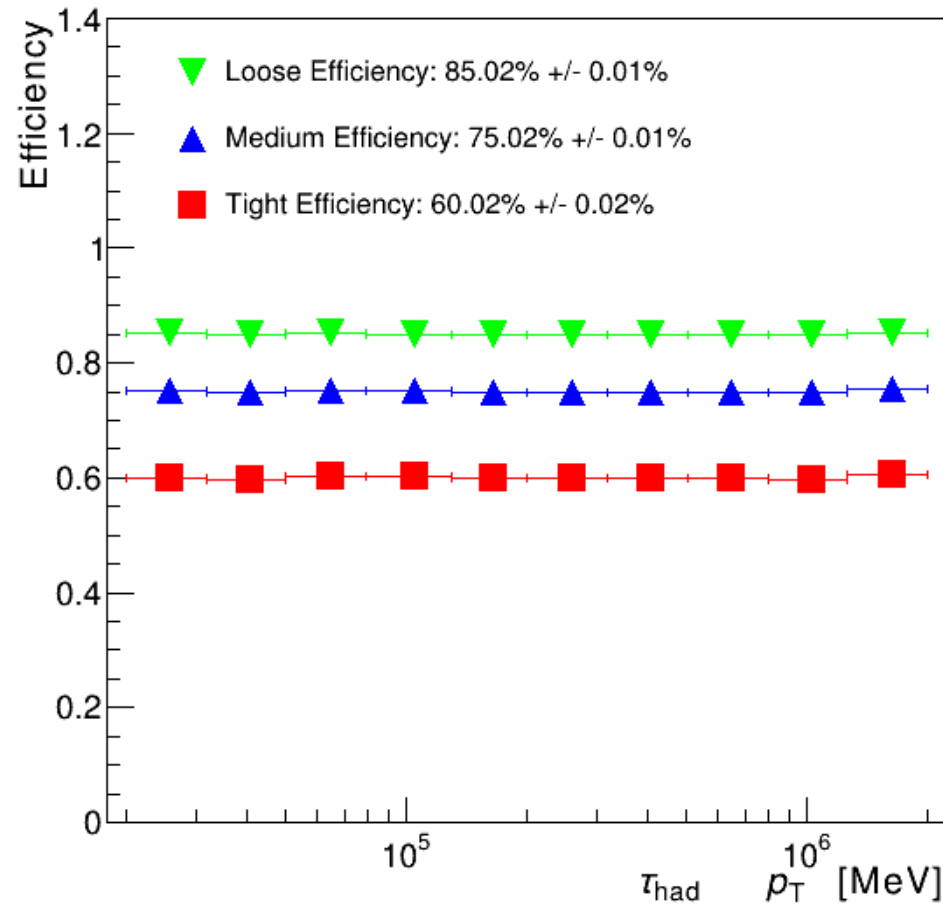


Signal sample score distribution after transformation

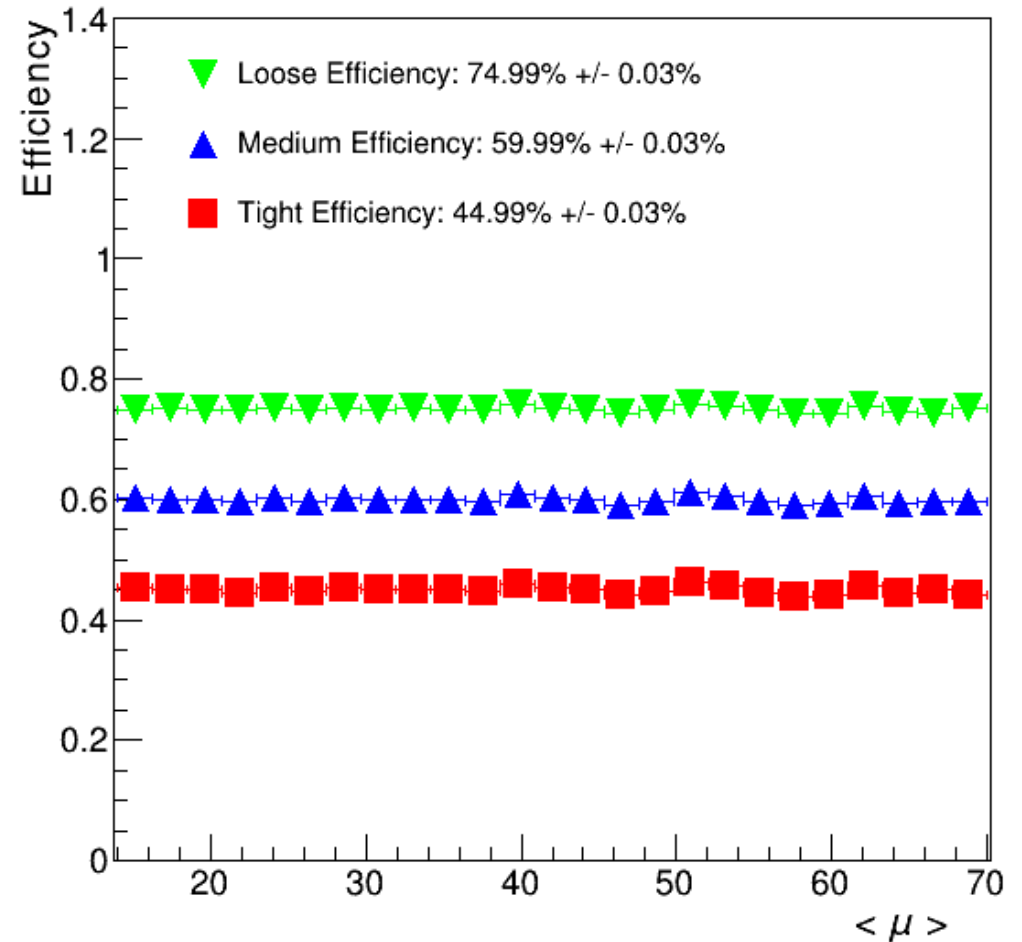
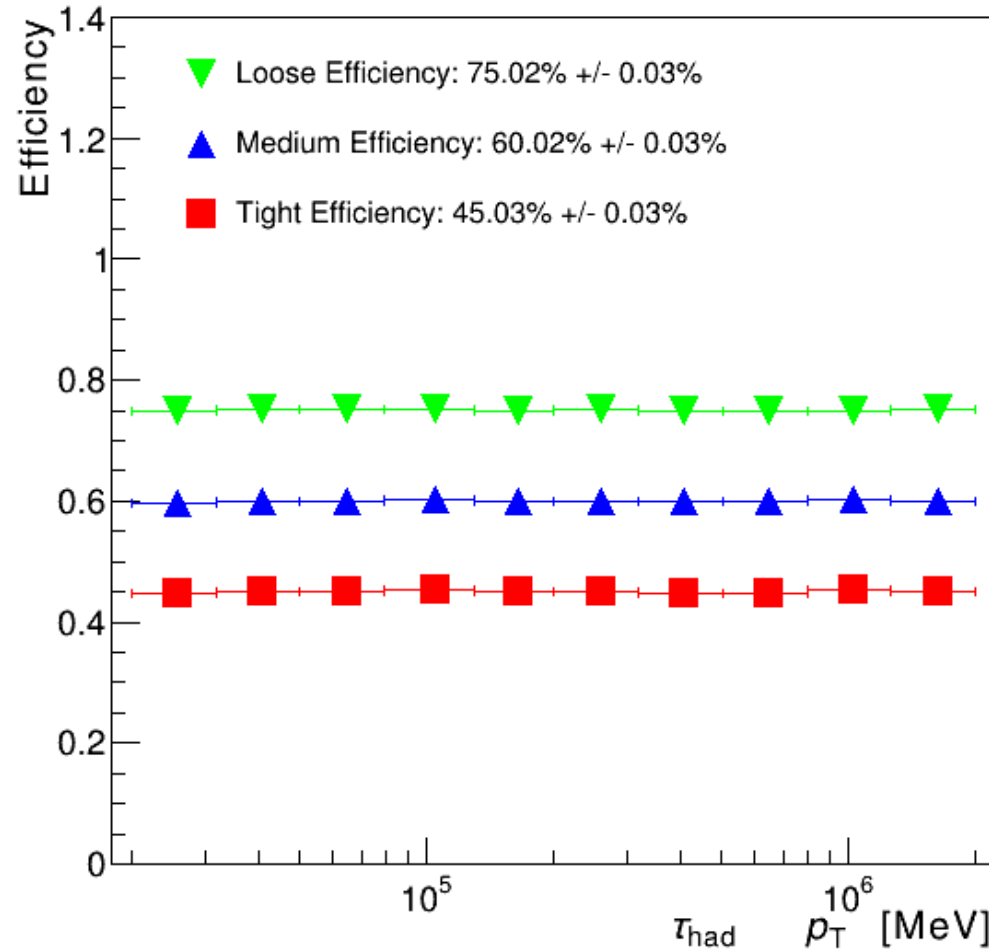


# Backup: TauJetGraphs Signal Efficiency and Background Rejection

# GNN Efficiency Plots (1-prong)

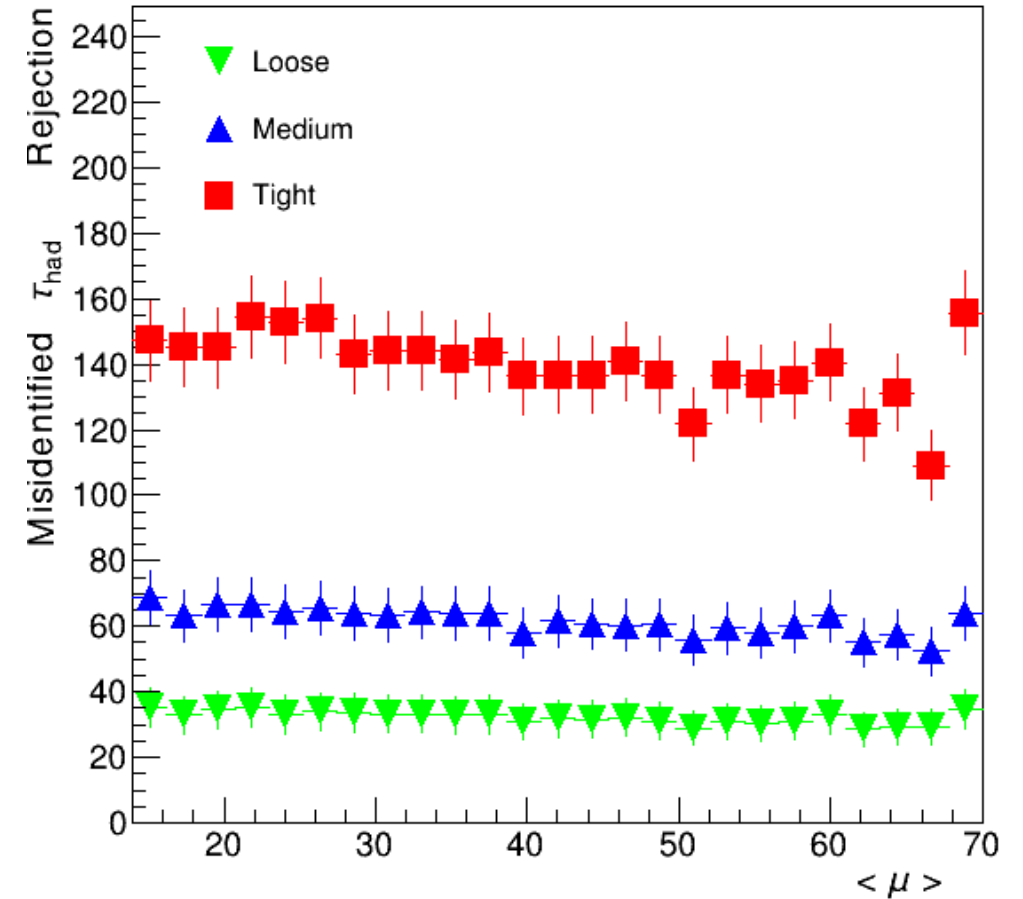
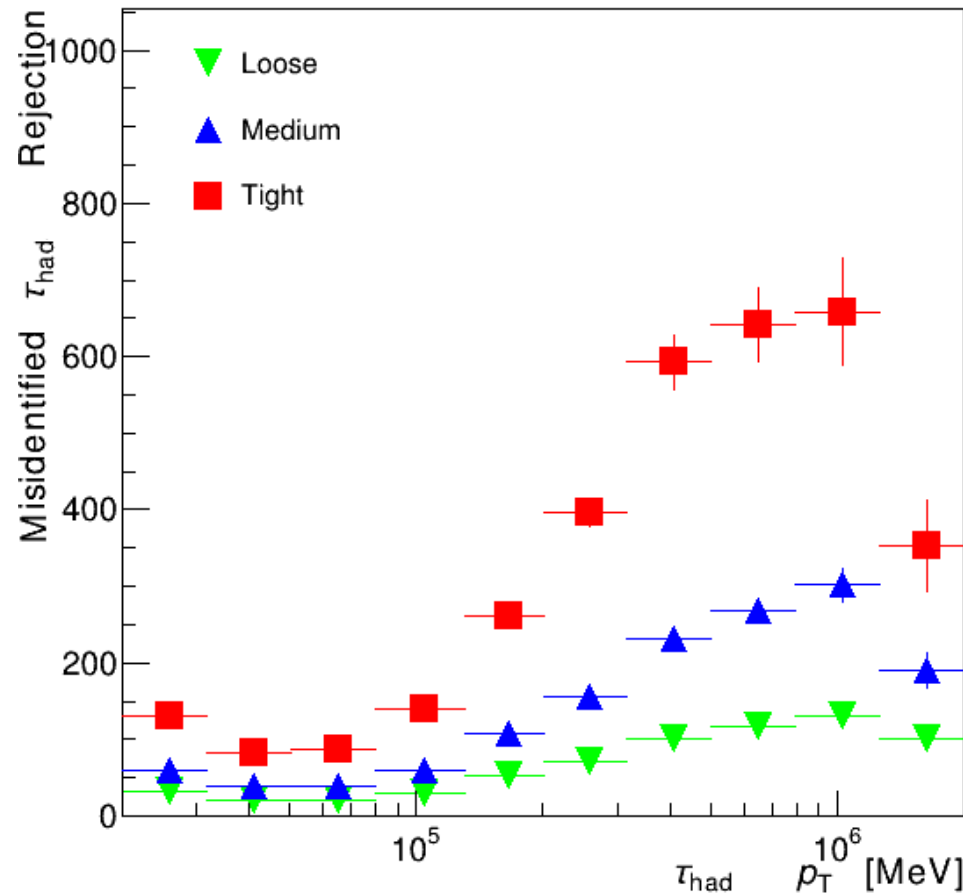


# GNN Efficiency Plots (3-prong)

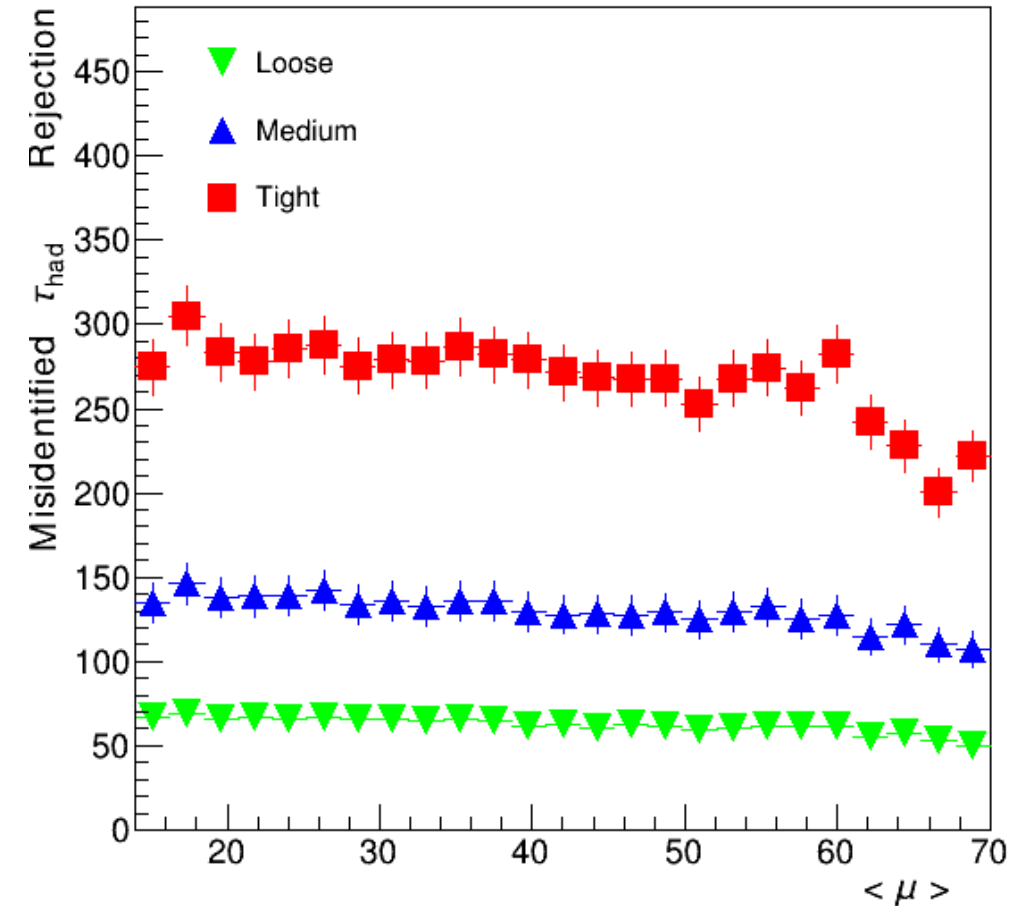
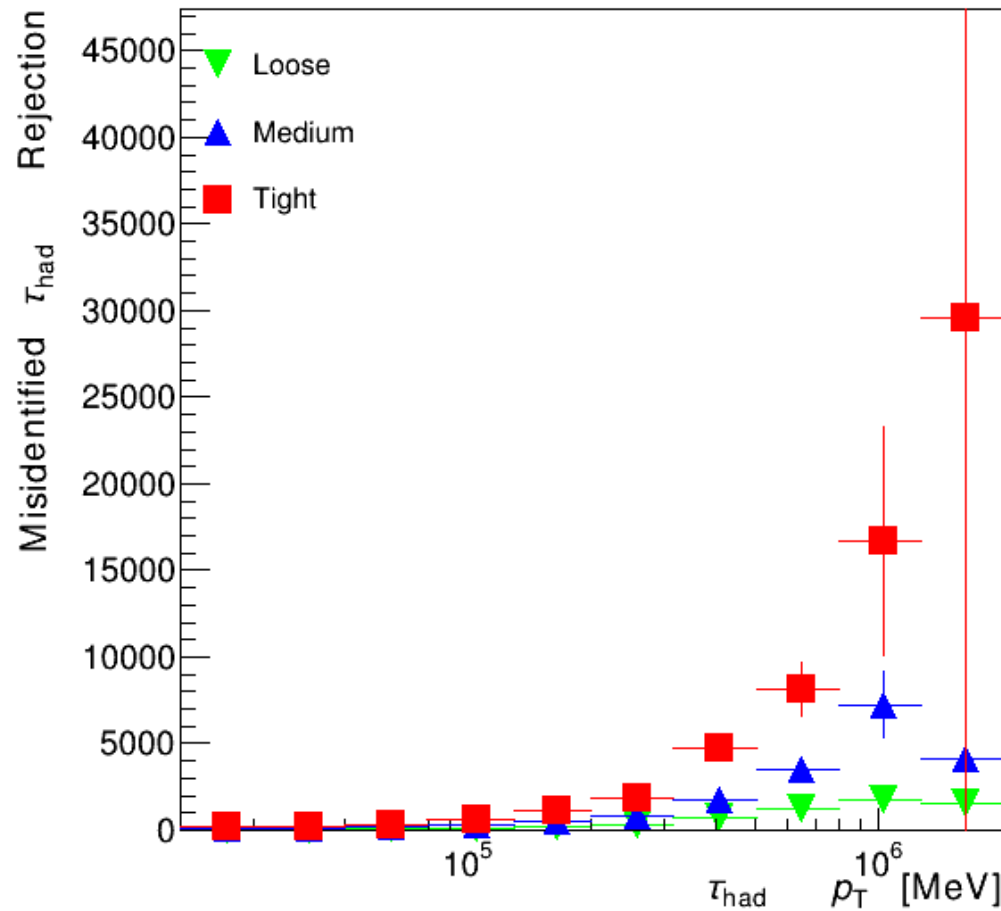




# GNN Rejection (1-prong)



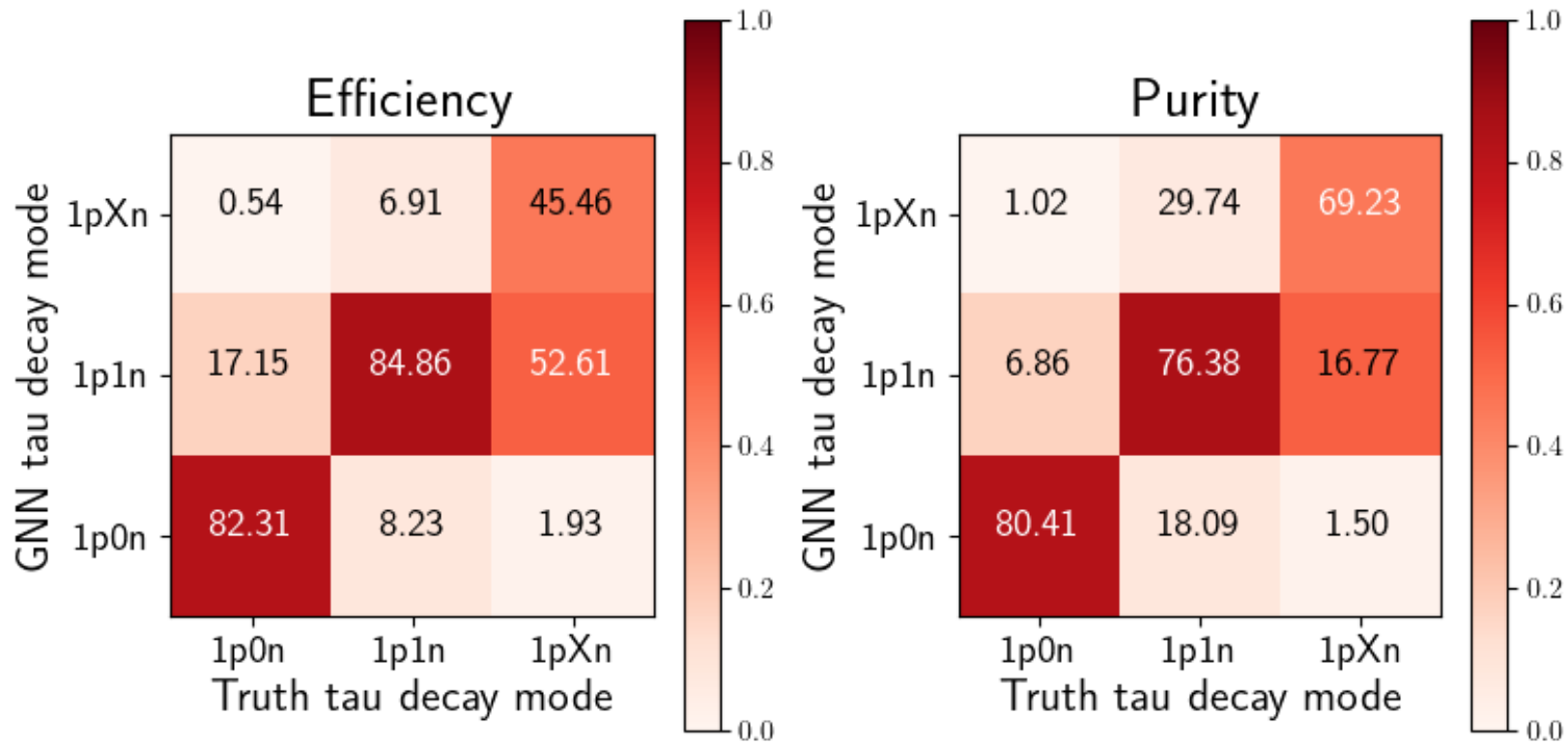
# GNN Rejection (3-prong)



# Backup: GNN Confusion Matrices

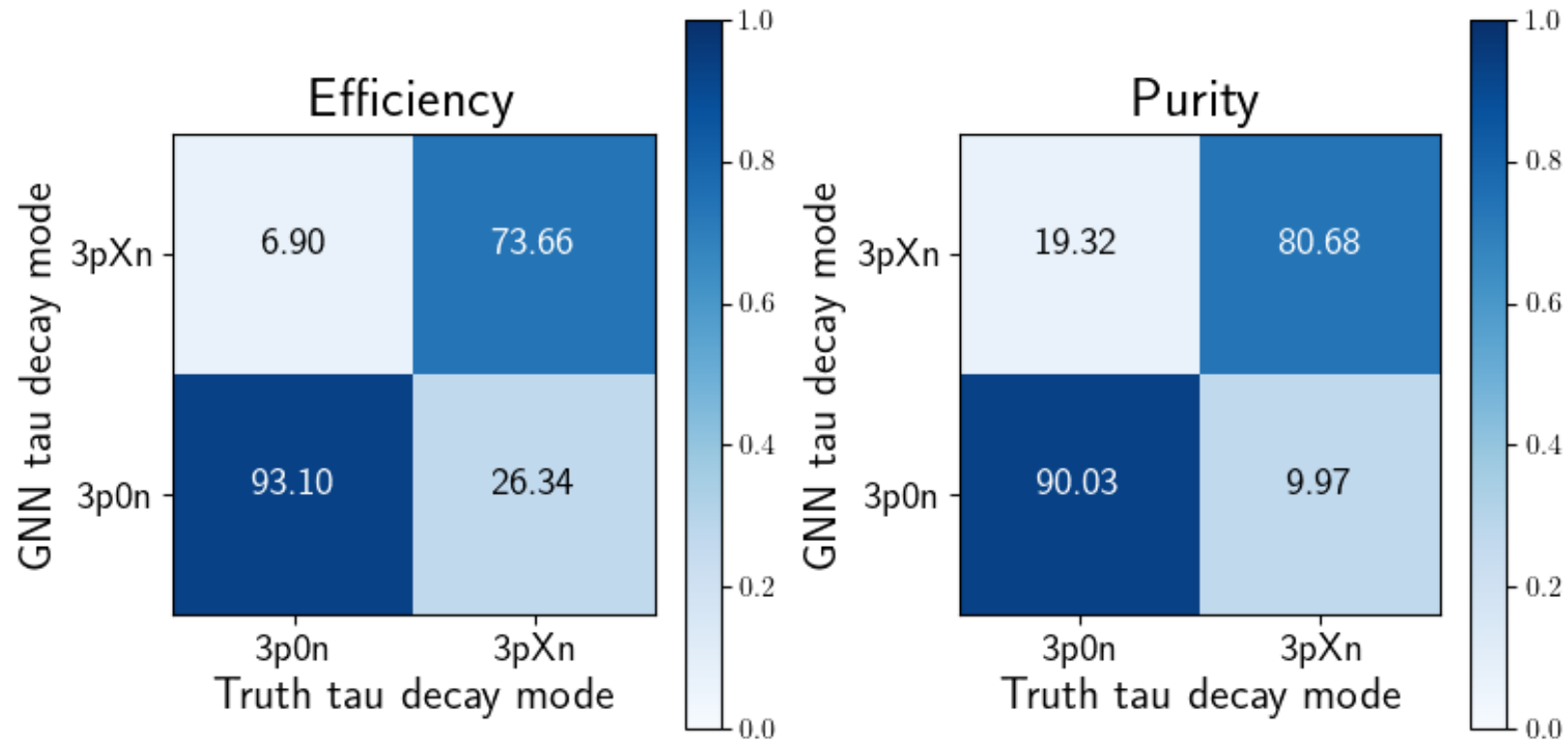
# GNN Confusion Matrices (1-prong)

For 75% Signal Efficiency



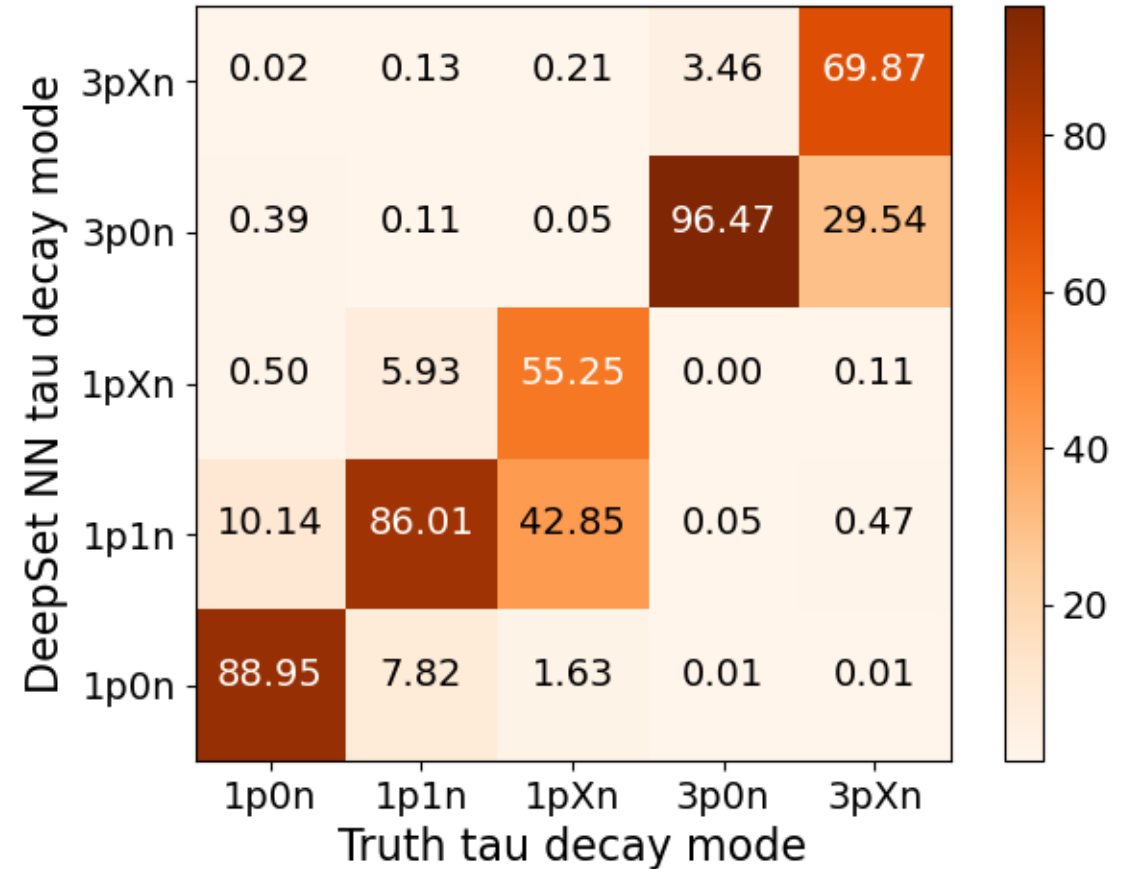
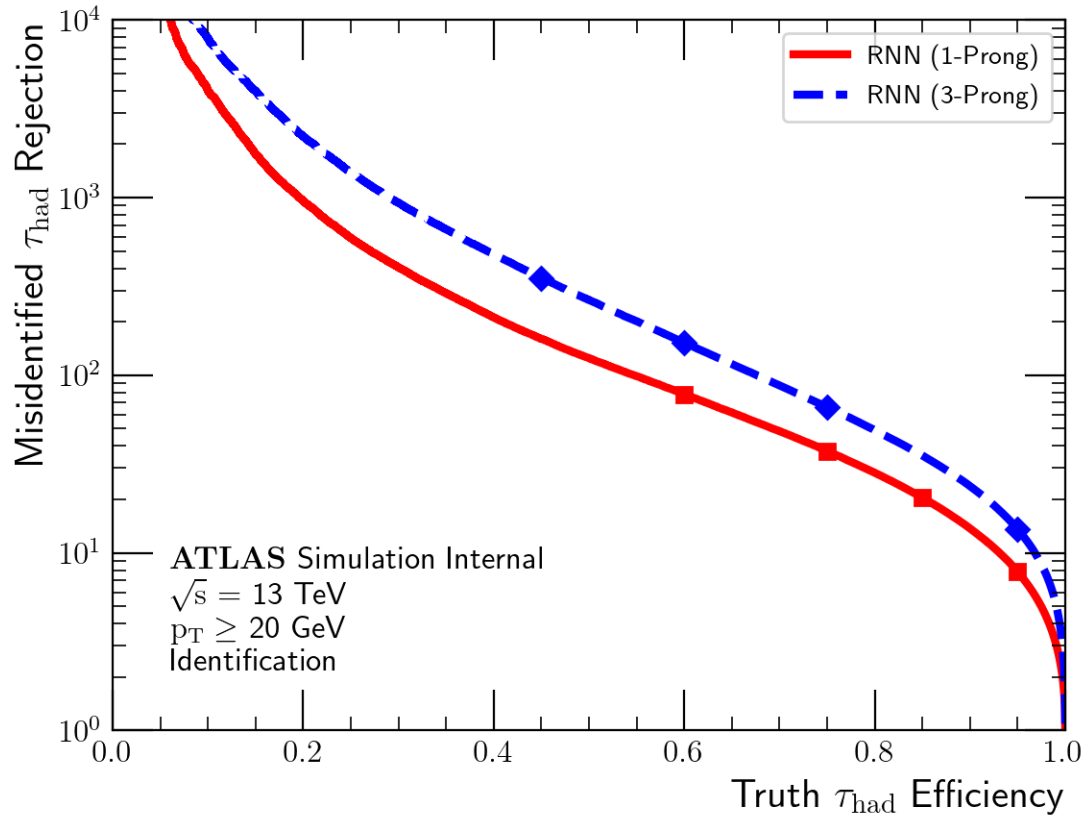
# GNN Confusion Matrices (3-prong)

For 60% Signal Efficiency



# Backup: Performance of Current Methods

# Tau ID RNN and Decay Mode Classifier DSNN



# Backup: Definitions and Glossary



# Metric Definitions

- **Accuracy** – The fraction of correctly classified samples (if normalised = True)
- **Purity (Precision)** – Purity is the measure of how well a classifier avoids incorrectly labelling a sample as positive. It's calculated as true positives divided by true positives plus false positives:
  - $\frac{tp}{tp+fp}$  where  $tp$  is true positive and  $fp$  is false positive
- **Efficiency (Recall)** – Efficiency measures how well a classifier finds all the true positives. It's calculated as true positives divided by true positives plus false negatives:
  - $\frac{tp}{tp+fn}$  where  $tp$  is true positive and  $fn$  is false negative
- **Background Rejection** – The inverse of the Background Selection Efficiency, depending on the Signal Selection Efficiency

# Glossary

- **ID** – Identification
- **DMC** – Decay Mode Classification
- $\tau_{\text{had}}$  - Hadronically decaying  $\tau$ -lepton
- **RNN** – Recurrent Neural Network
- **DSNN** - DeepSet Neural Network
- **GNN** – Graph Neural Network
- **ROC Curve** - Receiver Operator Characteristic Curve