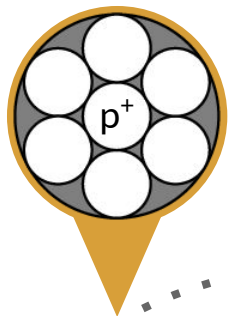
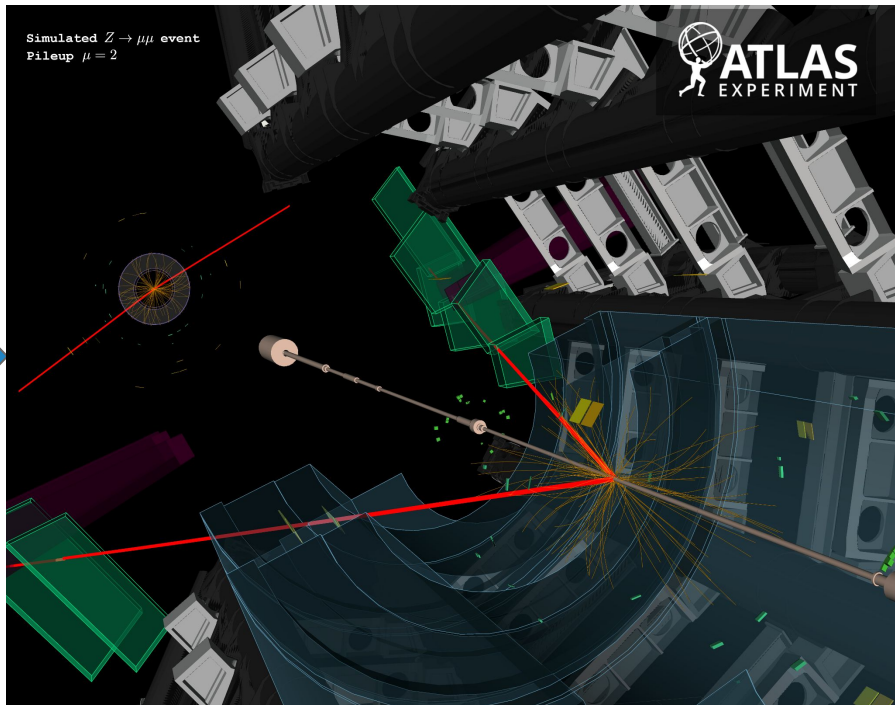


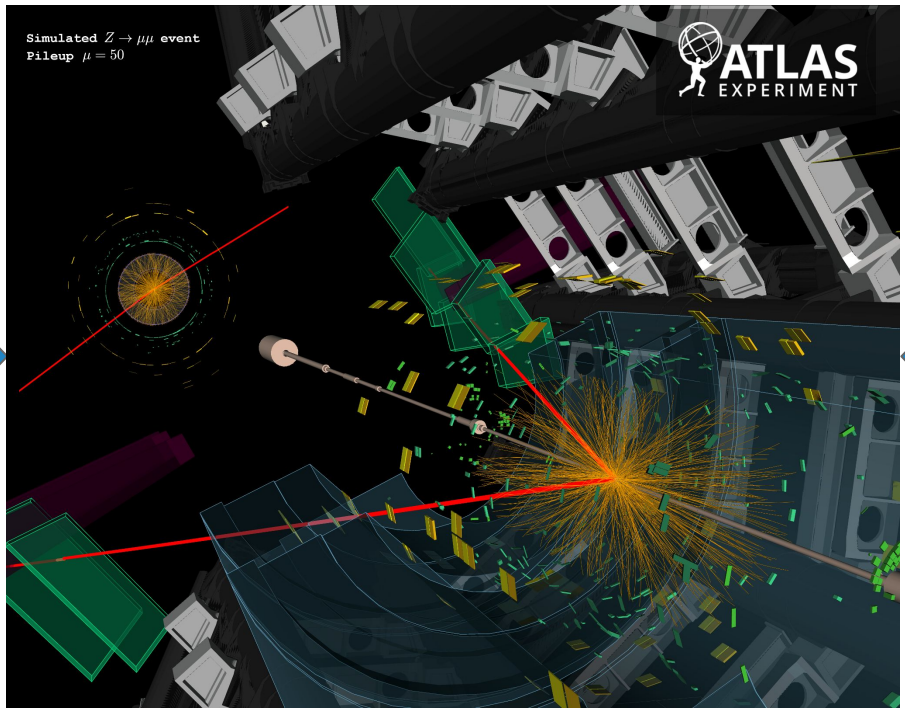
WOTAN: Weakly-supervised Optimal Transport Attention-based Noise Mitigation

Nathan Suri, Vinicius Mikuni, Benjamin
Nachman
ML4Jets 2024

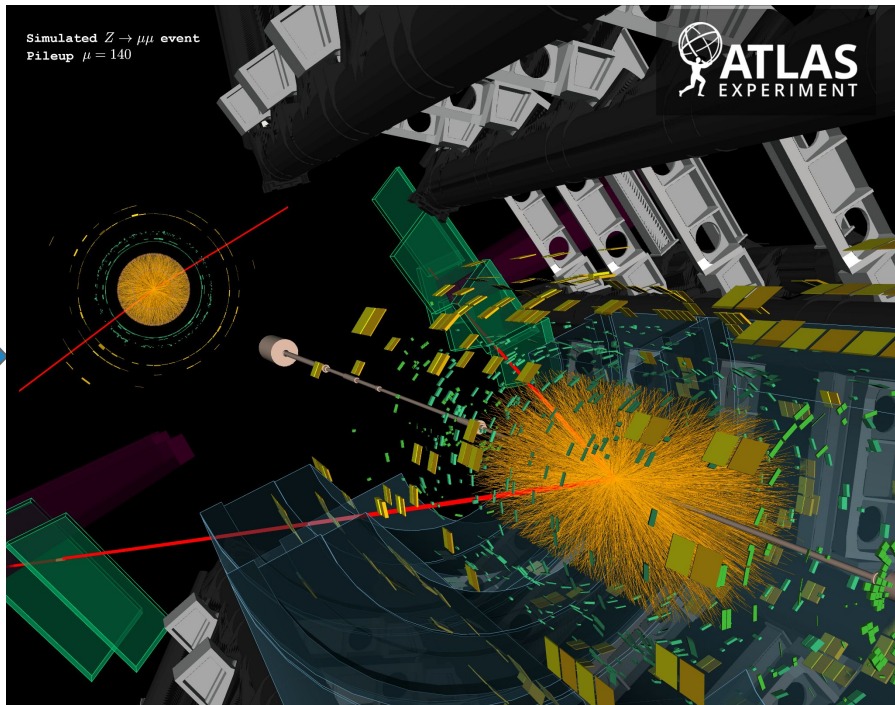




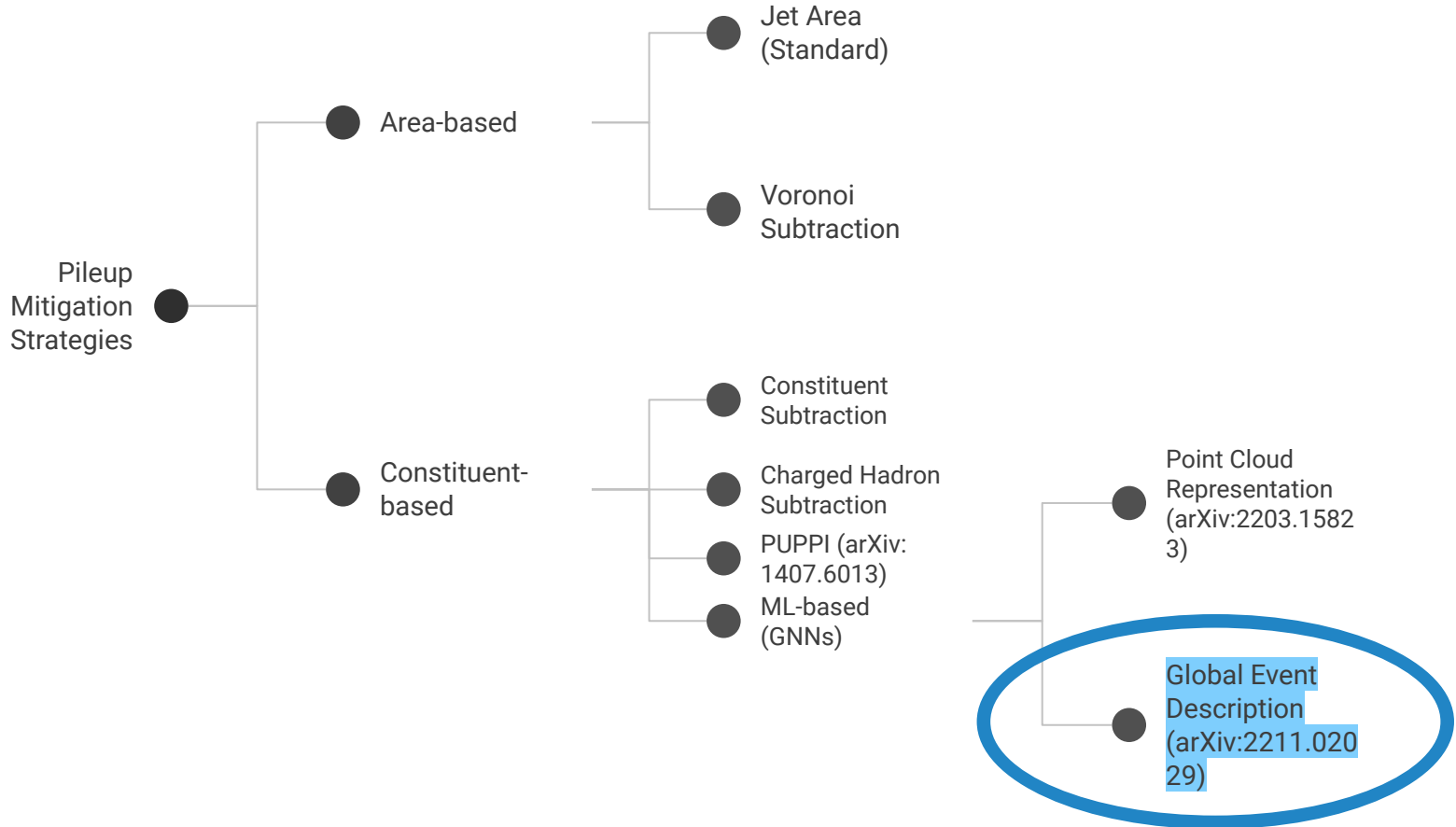
Charged + neutral pileup
In-time + out-of-time pileup



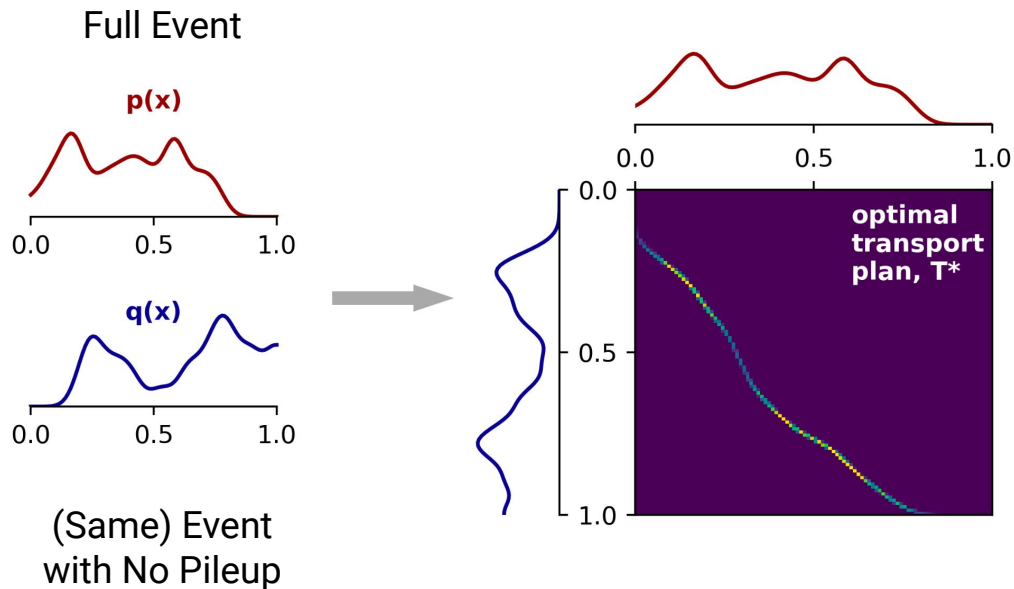
Charged + neutral pileup
In-time + out-of-time pileup



Charged + neutral pileup
In-time + out-of-time pileup



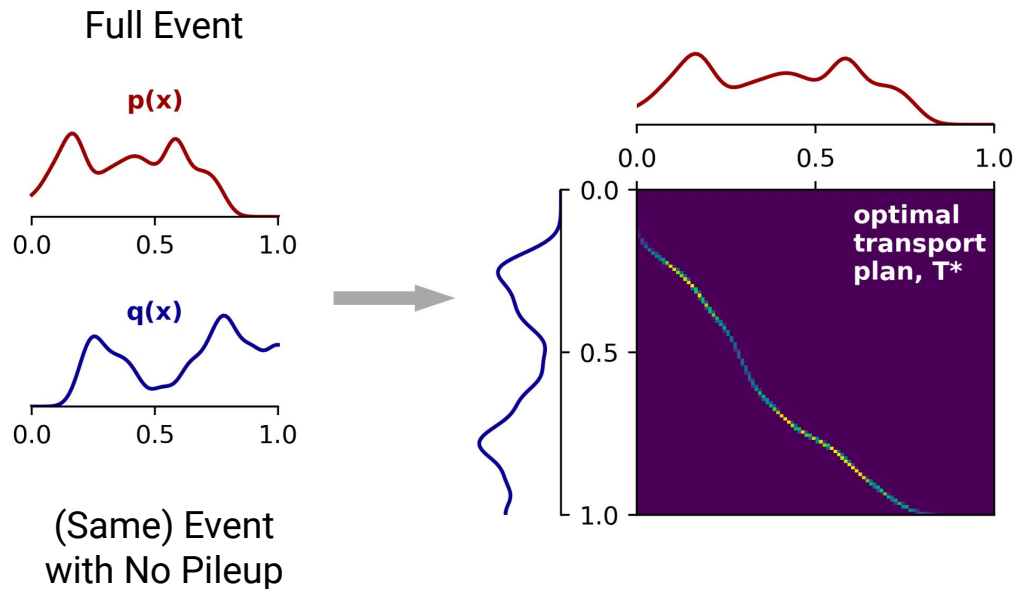
TOTAL: Training Optimal Transport with Attention Learning



$$\mathcal{L} = \text{SWD}(x'_p, x_{np}) + \lambda \text{MSE}(p_T^{\text{miss}}(x'_p), p_T^{\text{miss}}(x_{np}))$$

TOTAL: Training Optimal Transport with Attention Learning

- The probability density is intractable, but we can approximate the density
- Realizations of the density are accessible
- Optimal transport over the space of inputs allows for approximation



$$\mathcal{L} = \text{SWD}(x'_p, x_{np}) + \lambda \text{MSE}(p_T^{\text{miss}}(x'_p), p_T^{\text{miss}}(x_{np}))$$


TOTAL: Training Optimal Transport with Attention Learning

$$\mathcal{L} = \text{SWD}(x'_p, x_{np}) + \lambda \text{MSE}(p_T^{\text{miss}}(x'_p), p_T^{\text{miss}}(x_{np}))$$

TOTAL: Training Optimal Transport with Attention Learning

$$\mathcal{L} = \text{SWD}(x'_p, x_{np})$$

- Wasserstein distance (WD): Finds the transport function that keeps hard scattering particles and removes those from simultaneous vertices
- Sliced WD to compensate for poor scaling of computational costs of calculating WD at high dimensions

$$x'_p = \omega x_p$$


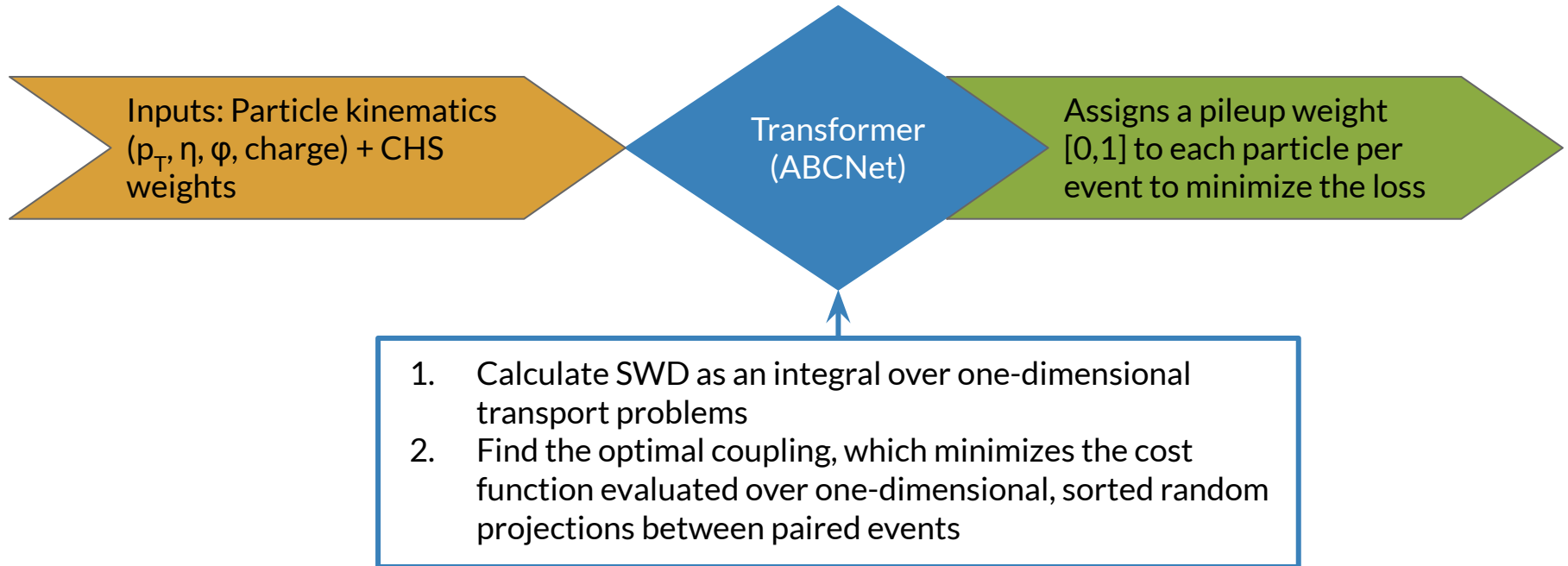
TOTAL: Training Optimal Transport with Attention Learning

$\mathcal{L} =$

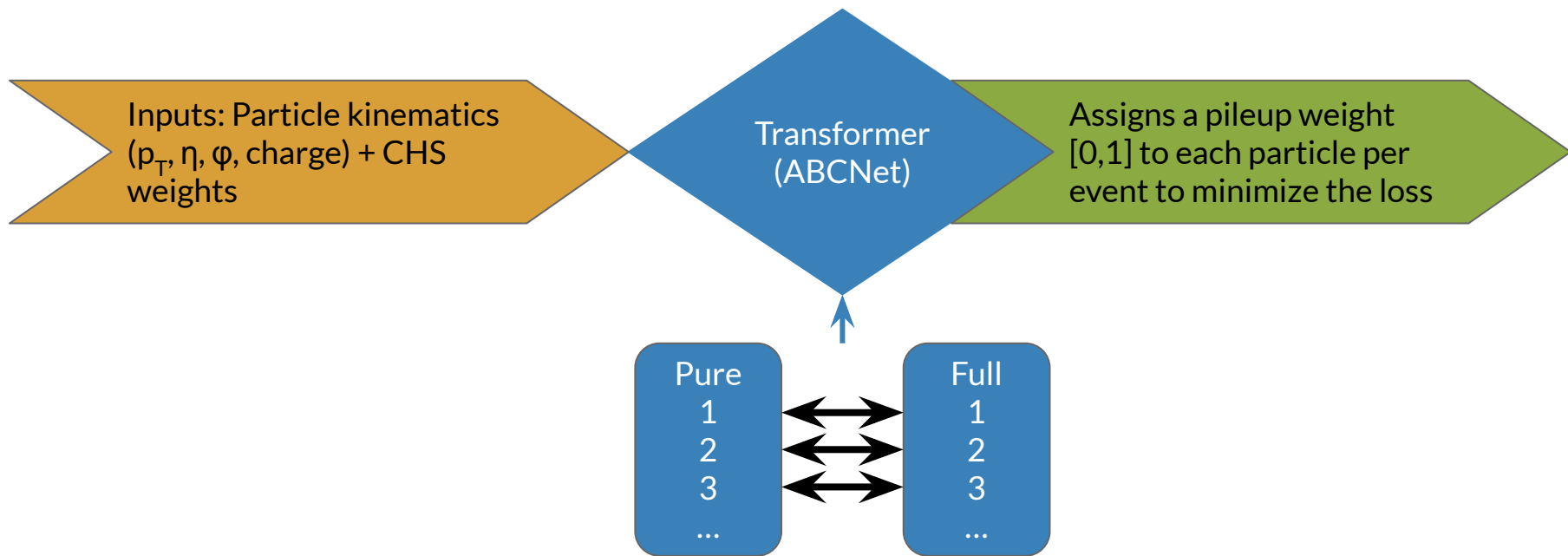
- Scaled Mean Square Error of missing p_T
- Forces energy conservation between the pure and full samples

$$+ \lambda \text{MSE}(p_T^{\text{miss}}(x'_p), p_T^{\text{miss}}(x_{np}))$$

TOTAL: Training Optimal Transport with Attention Learning

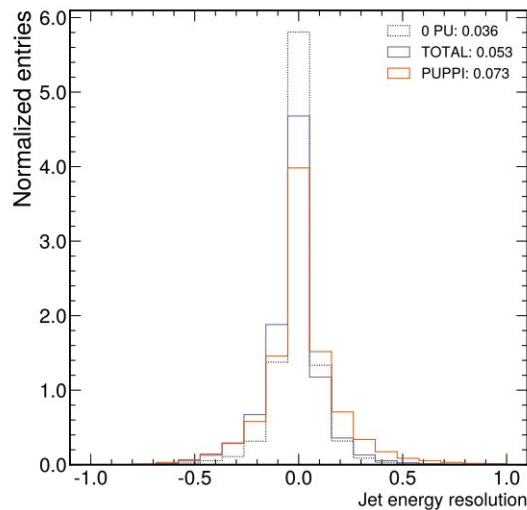


TOTAL: Training Optimal Transport with Attention Learning

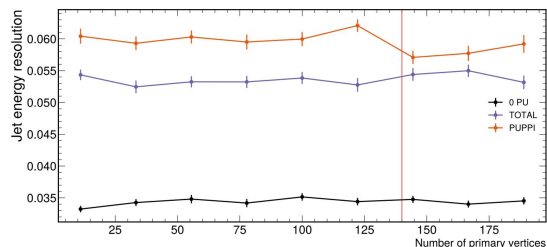


TOTAL: Training Optimal Transport with Attention Learning

- + Outperforms traditional and ML-based alternatives
- + Relies on global event descriptions
- + Robustly learns pileup characteristics as a transport function



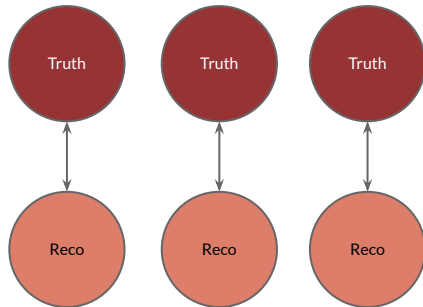
- Requires direct matching of events
- Overall limited due to supervision



TOTAL: Training Optimal Transport with Attention Learning

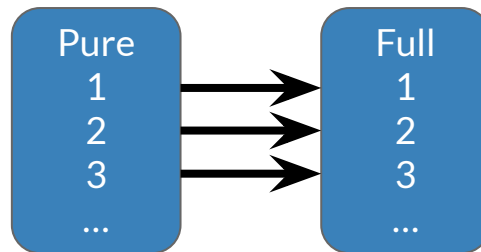
ML Competitors

- Matching between truth and reco at particle level (MC correction)



TOTAL

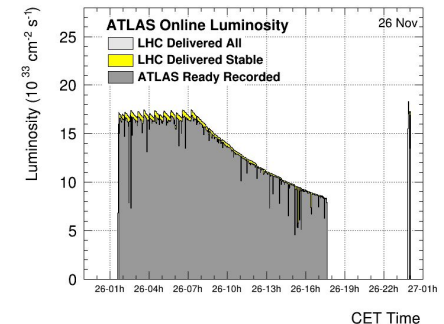
- Matching between pileup events and same event without pileup vertices (data-driven*)



**Such supervision is not physically realizable*

WOTAN

- Matching between ensembles of events with different relative pileup densities (fully data-driven)



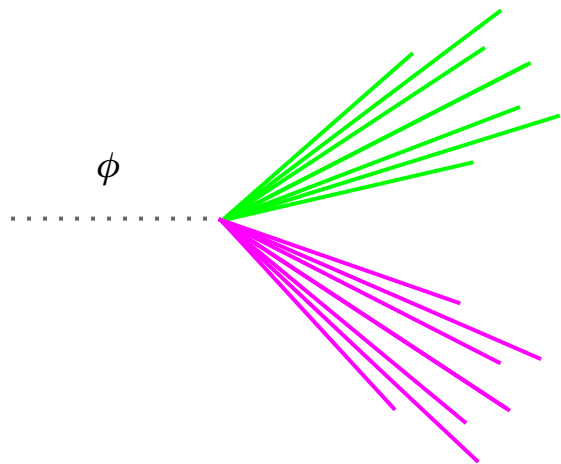


How can we improve TOTAL's flexibility?



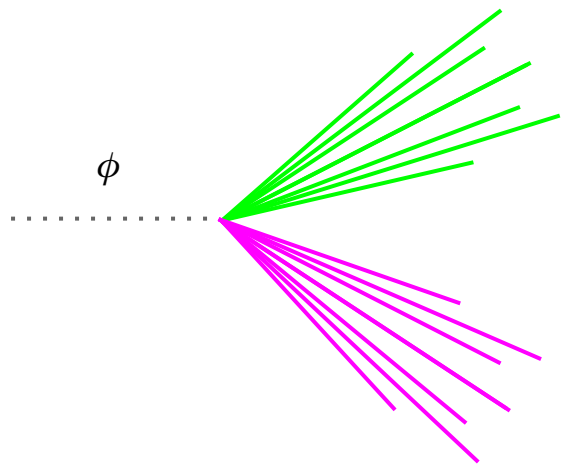
Reduce supervision!

Physics Example: High p_T Jets



- ▷ PUMML Dataset:
<https://zenodo.org/records/2652034>
- ▷ Process: q-qbar
light-quark-initiated jets from the
from the decay of a Higgs-like scalar
- ▷ Pileup was generated by
overlaying soft QCD on top of
signal

Physics Example: High p_T Jets

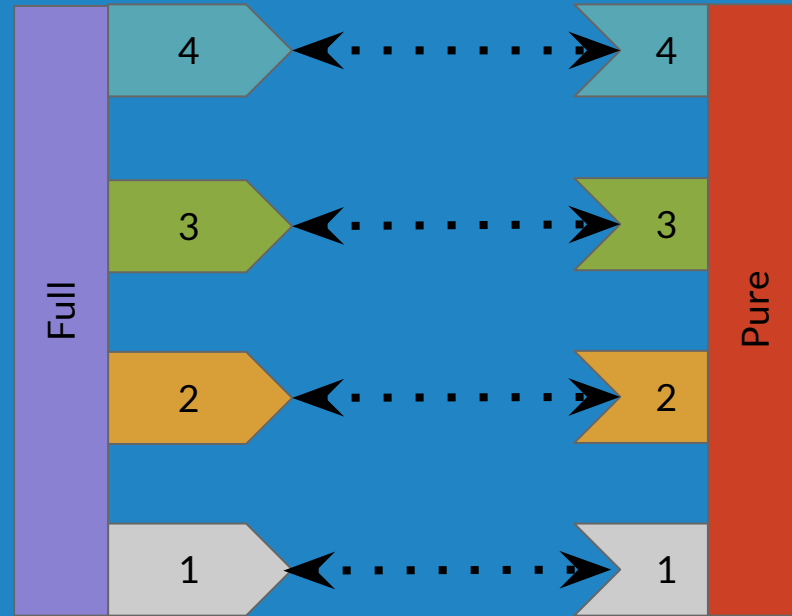


- ▷ PUMML Dataset:
<https://zenodo.org/records/2652034>
- ▷ Datasets
 - $\mu = 140, \Delta m$: Set pileup vertex count, varied scalar mass
 - $\Delta\mu, m = 500$ GeV: Varied pileup vertex (PV) count, set scalar mass
 - PV: 130-141

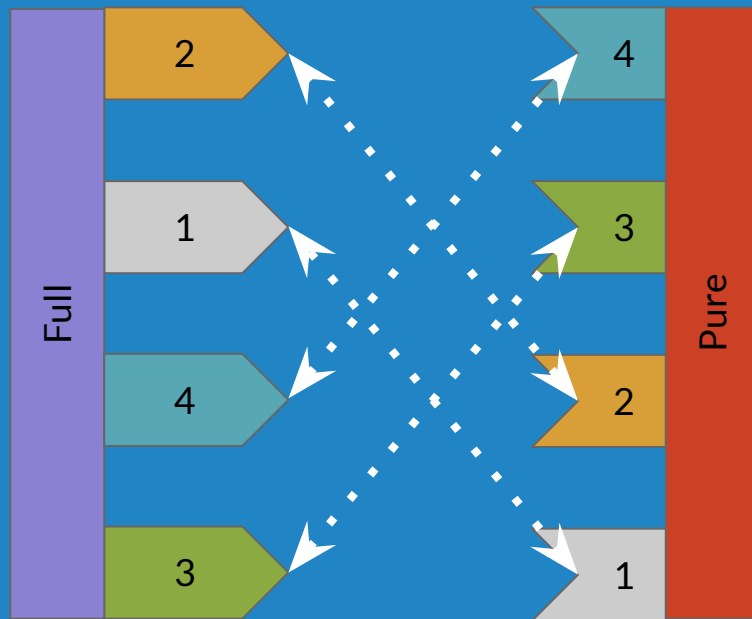


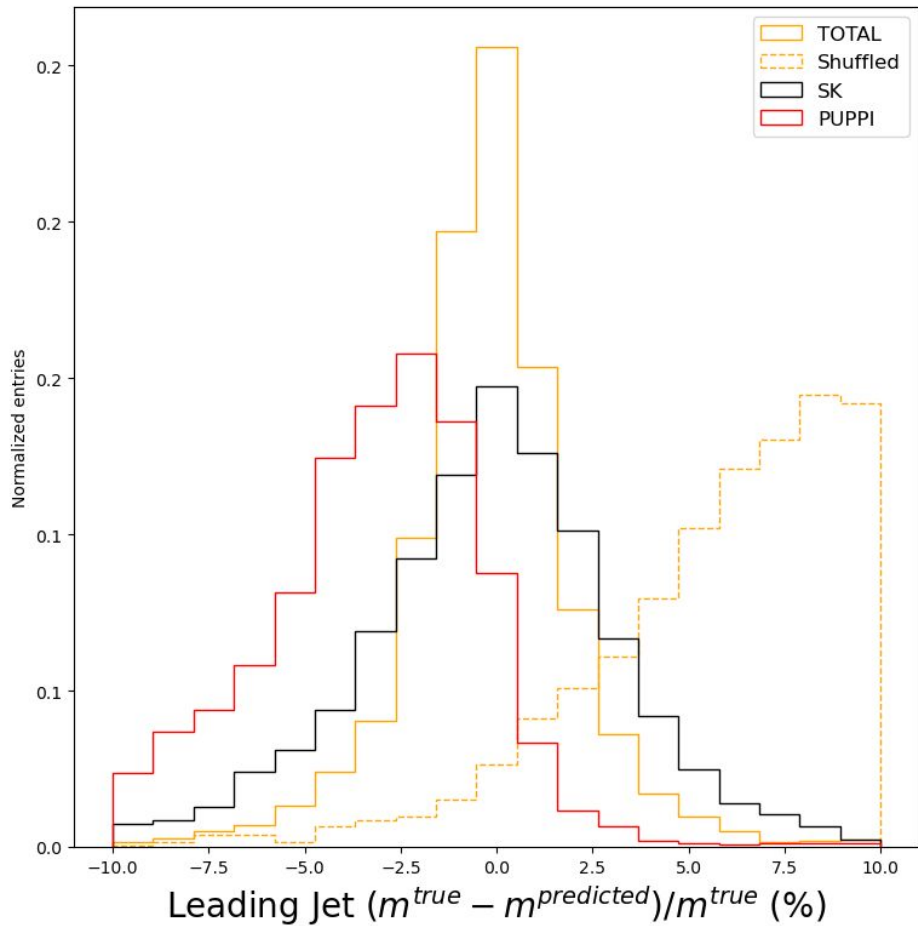
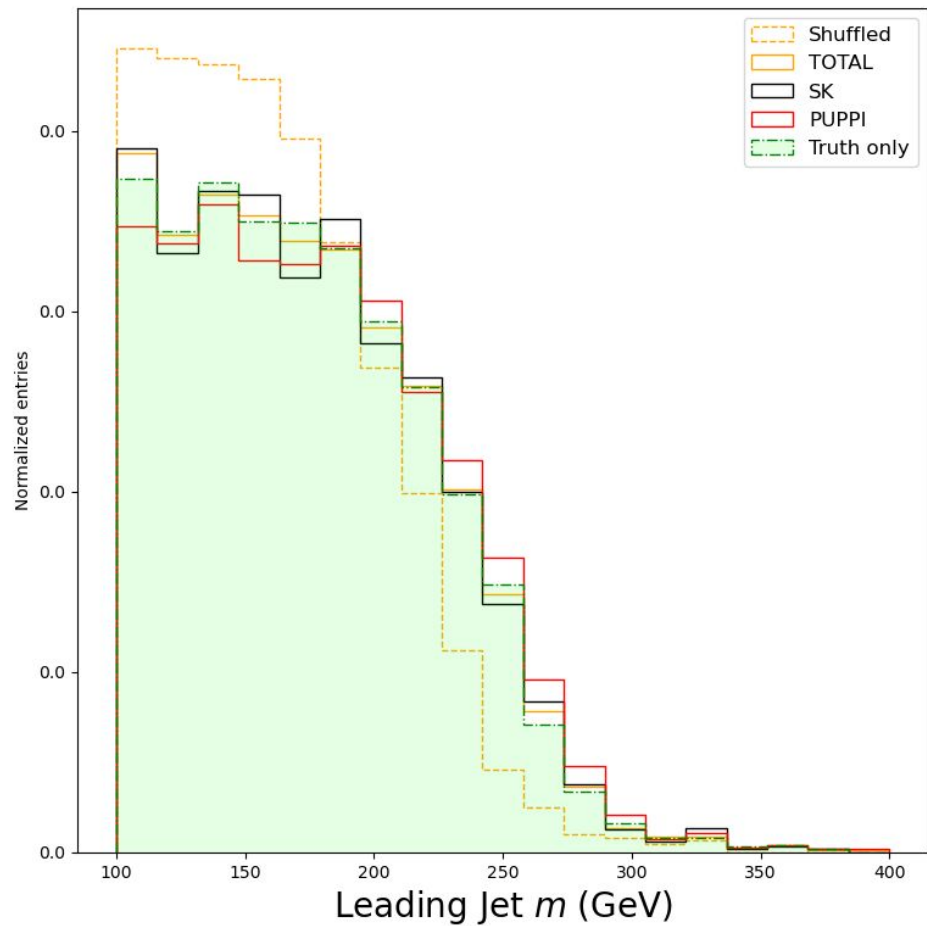
*What happens if we do not require
direct matching?*

TOTAL

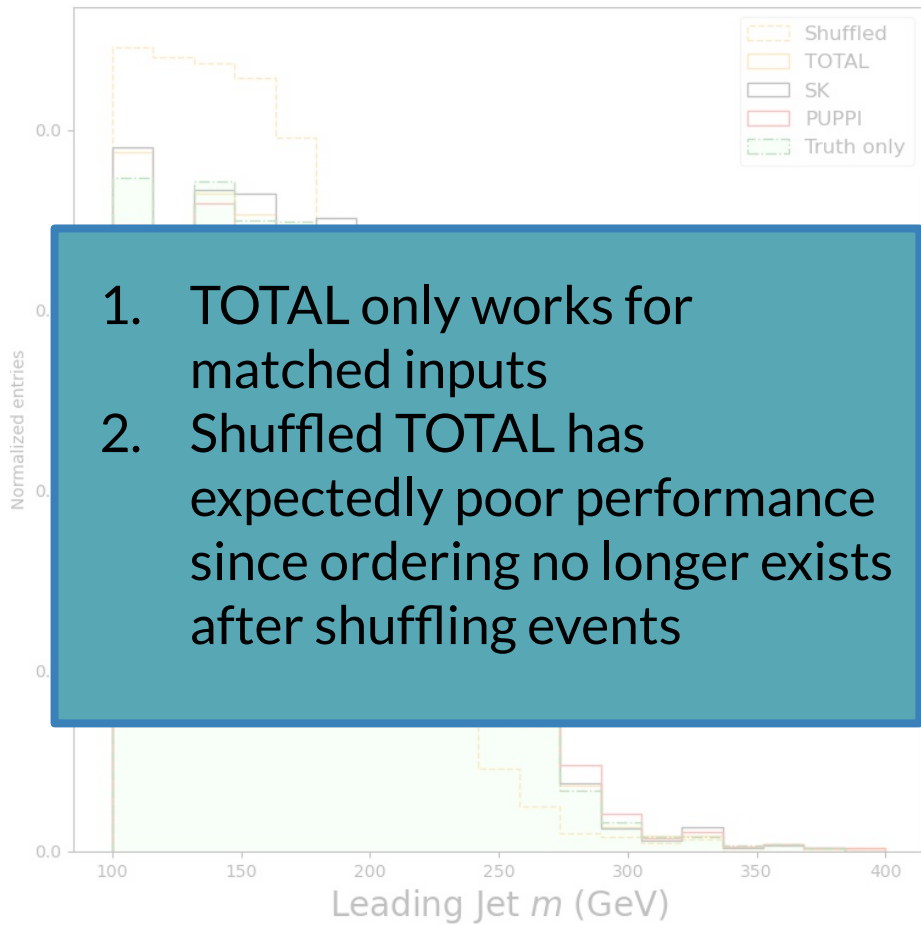


Shuffled

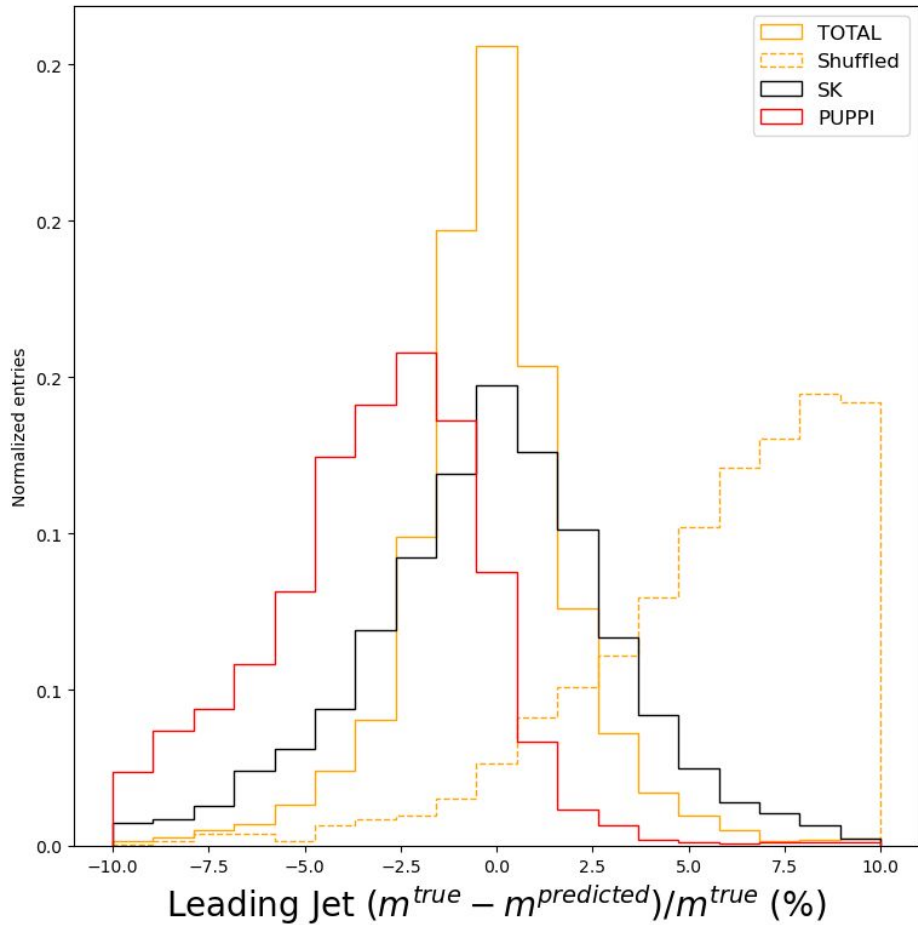




*Each result averaged over 5 trainings



1. TOTAL only works for matched inputs
2. Shuffled TOTAL has expectedly poor performance since ordering no longer exists after shuffling events

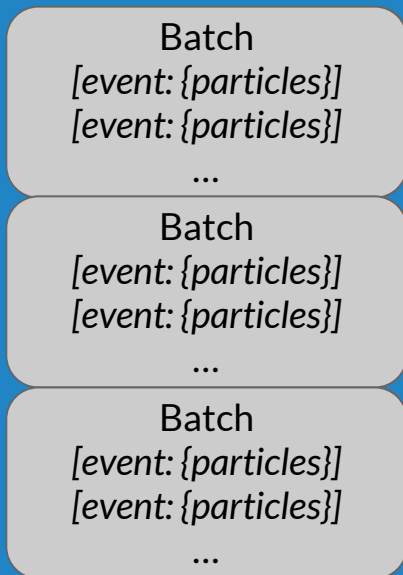




*How can we mitigate the
information loss of not matching
events?*

TOTAL

$[n_{\text{batch}} \times n_{\text{particles}} \times n_{\text{features}}]$

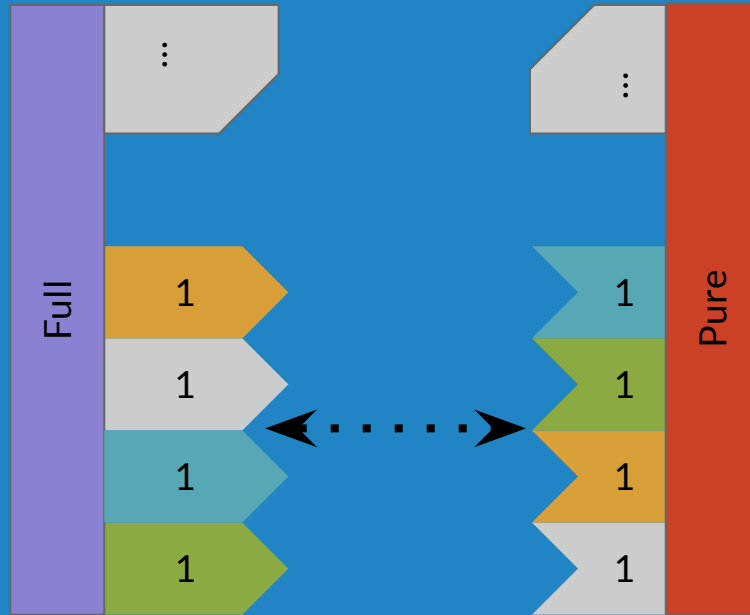


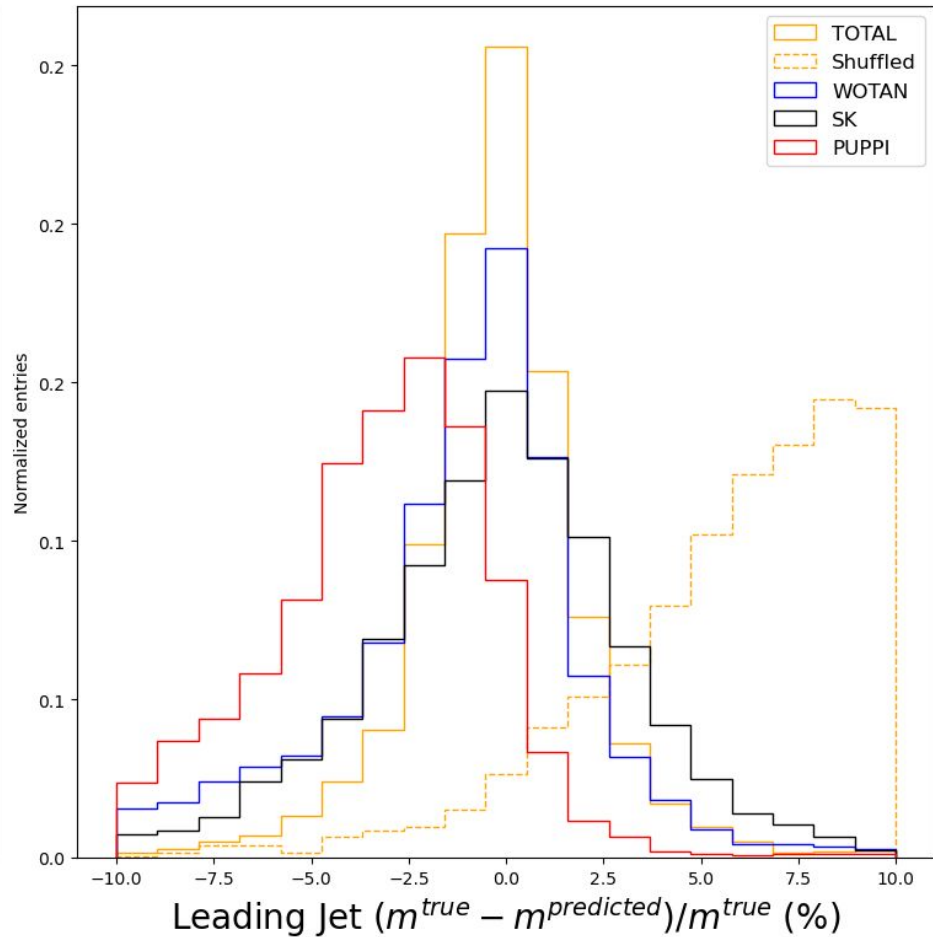
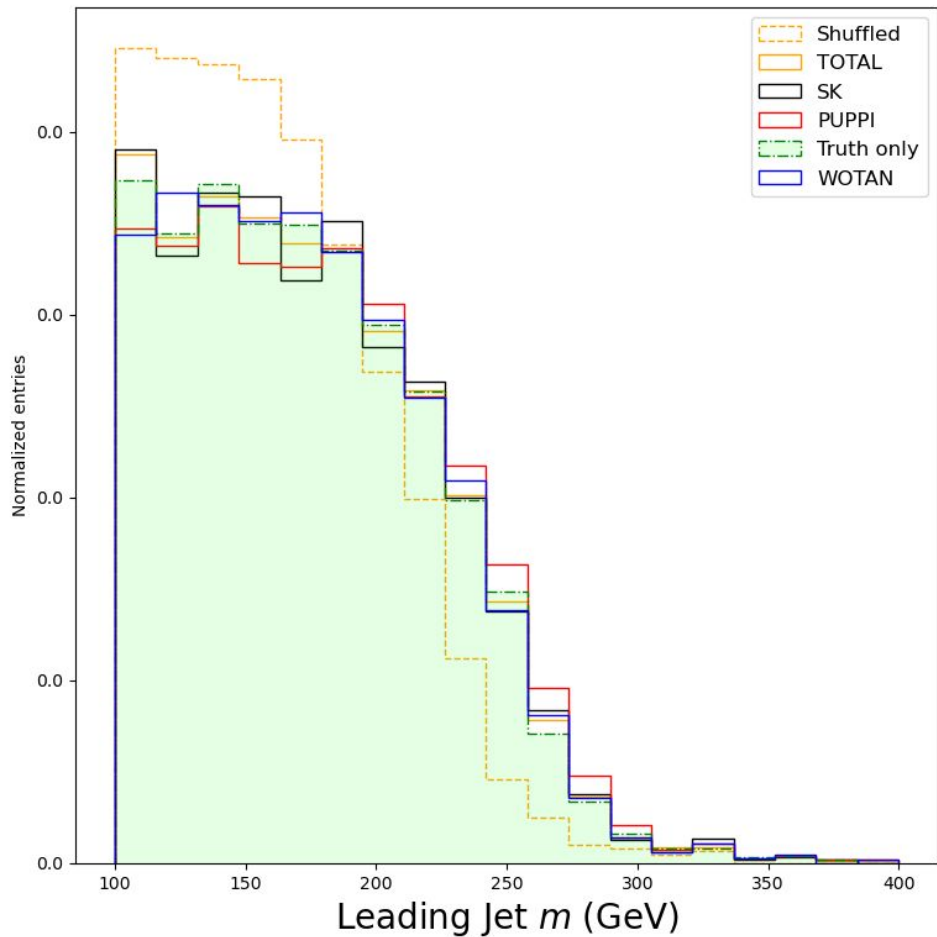
WOTAN

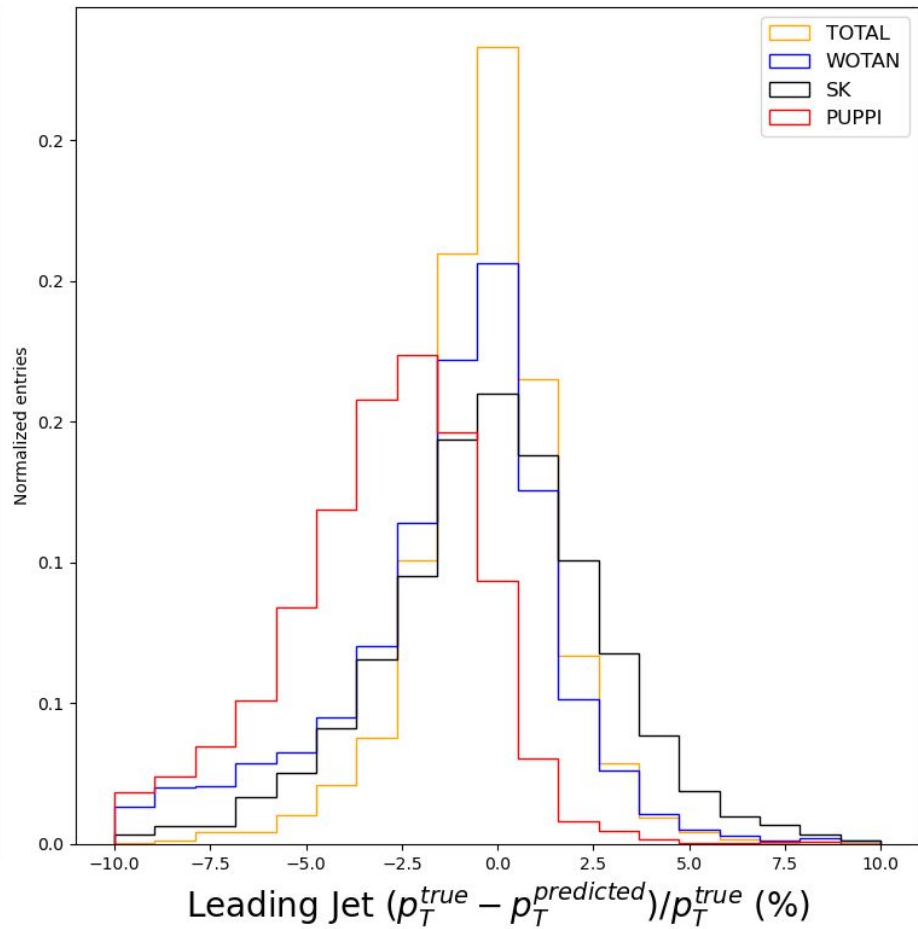
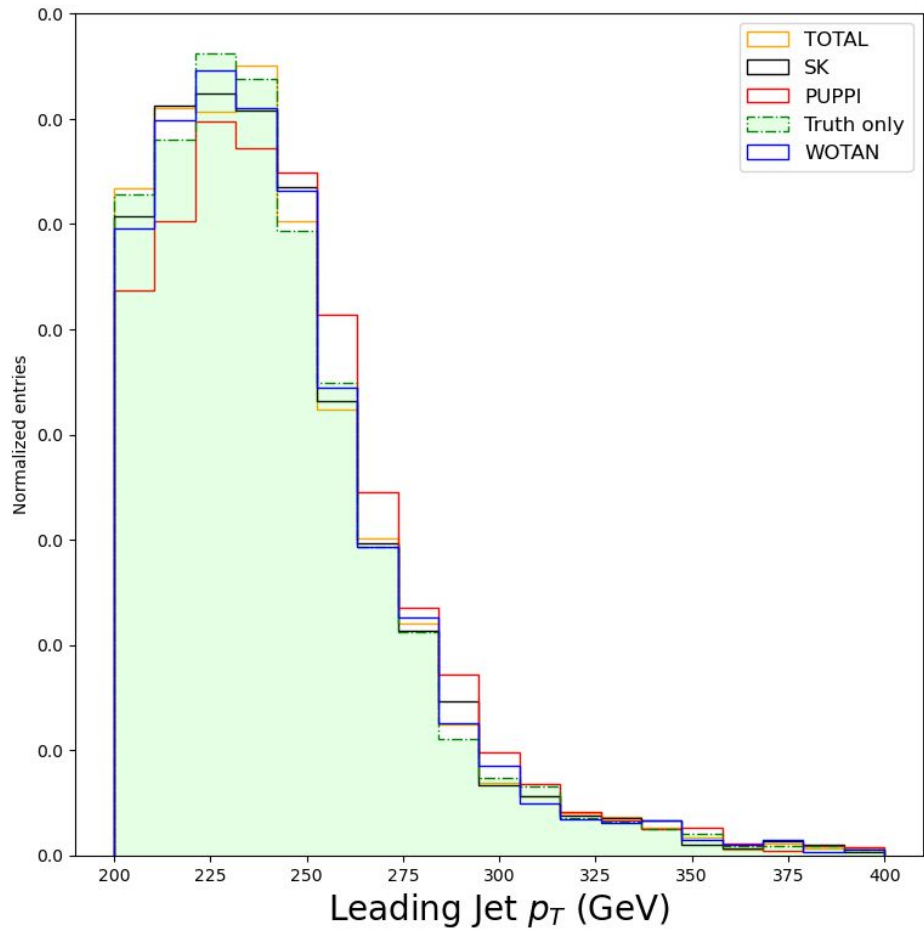
$[(n_{\text{batch}} \times n_{\text{particles}}) \times n_{\text{features}}]$

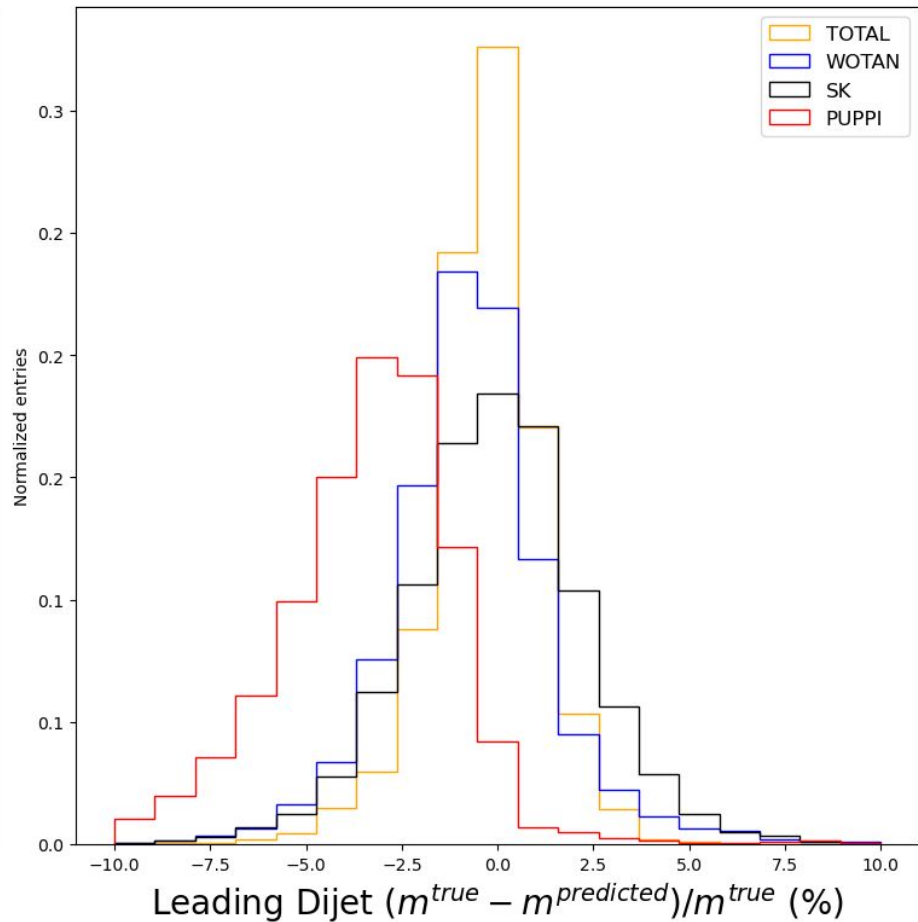
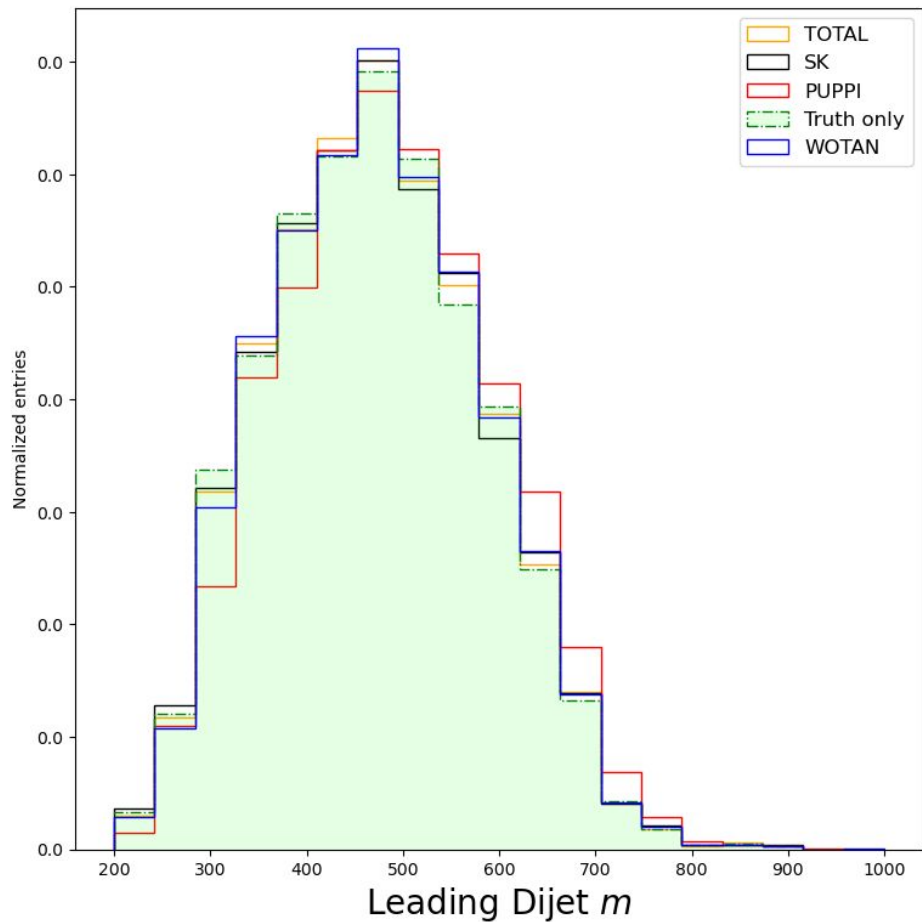


WOTAN

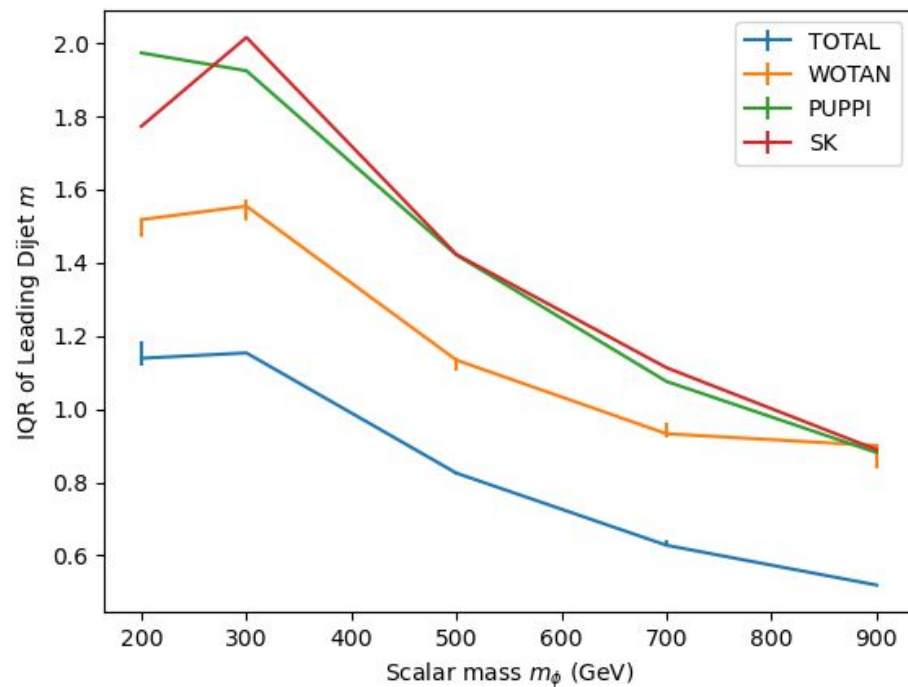
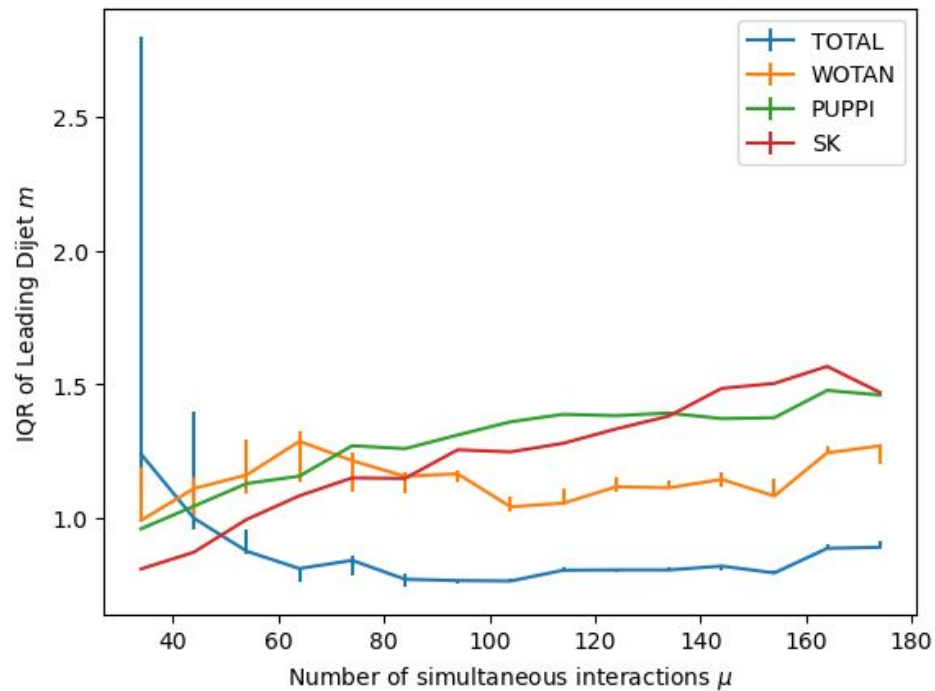








IQR	TOTAL	WOTAN	SoftKiller	PUPPI
Leading Jet p_T Resolution	0.967	1.704	1.749	1.714
Leading Jet m Resolution	1.117	1.922	2.112	2.143
Leading Dijet m Resolution	0.793	1.117	1.422	1.421



Key Takeaways

01

WOTAN is a completely data-driven pileup mitigation technique

02

WOTAN outperforms conventional pileup mitigation strategies without requiring unphysical supervision

03

WOTAN is generalizable to any denoising problem that matches the outline discussed in this talk



Stay tuned!
Final results to be released soon
arXiv 2411.XXXXX

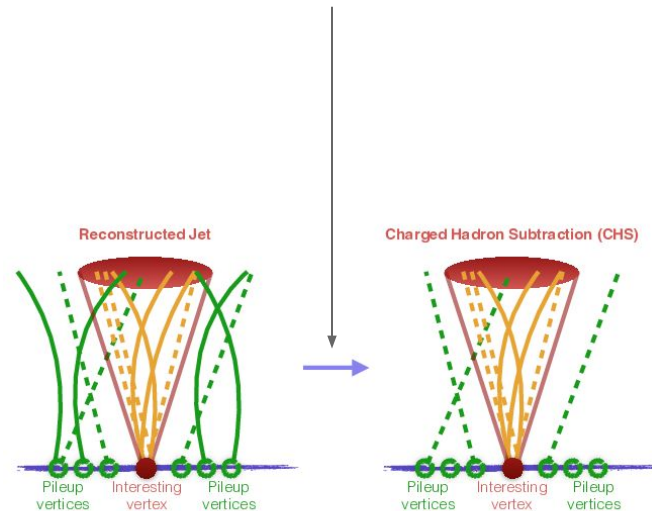


Backup Slides

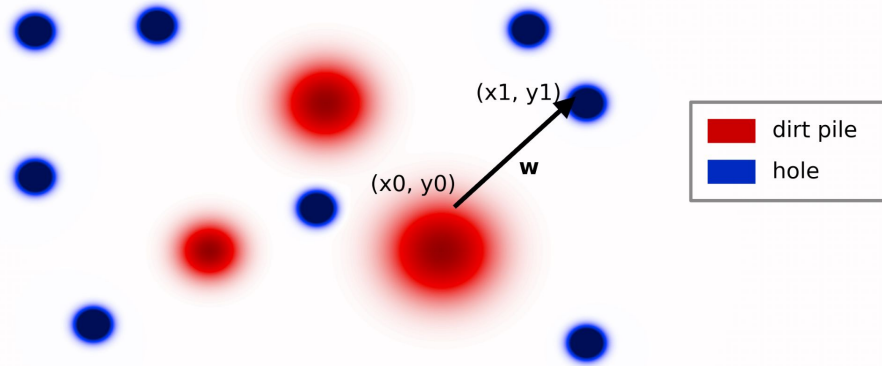
Charged Hadron Subtraction

$$\text{CVF} = \frac{\sum_{\text{tracks,HS}} p_T^{\text{track}}}{\sum_{\text{tracks,HS}} p_T^{\text{track}} + \sum_{\text{tracks,PU}} p_T^{\text{track}}}$$

- ▷ Benefits
 - Very effective at removing charged pileup due to track information
- ▷ Drawbacks
 - Inapplicable to neutral pileup
- ▷ (arXiv:2012.06271)



Earth Mover's Distance = W_1



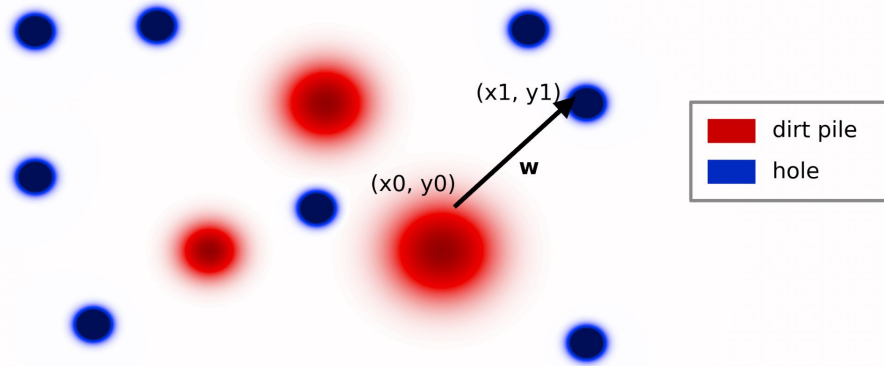
- ▷ Assumption: Total volume of the holes = total volume of the dirt piles
- ▷ Piles as the probability density function of P and holes as the probability density function of Q
- ▷ Per unit transportation cost:

$$C(x_0, y_0, x_1, y_1) = (x_0 - x_1)^2 + (y_0 - y_1)^2$$

- ▷ Transportation Plan:

$$T(x_0, y_0, x_1, y_1) = w$$

Earth Mover's Distance = W_1



$$\iint T(x_0, y_0, x, y) dx dy = p(x_0, y_0)$$

$$\iint T(x, y, x_1, y_1) dx dy = q(x_1, y_1)$$

$$\text{Total Cost} = \iiint \int C(x_0, y_0, x_1, y_1) \cdot T(x_0, y_0, x_1, y_1) dx_0 dy_0 dx_1 dy_1$$

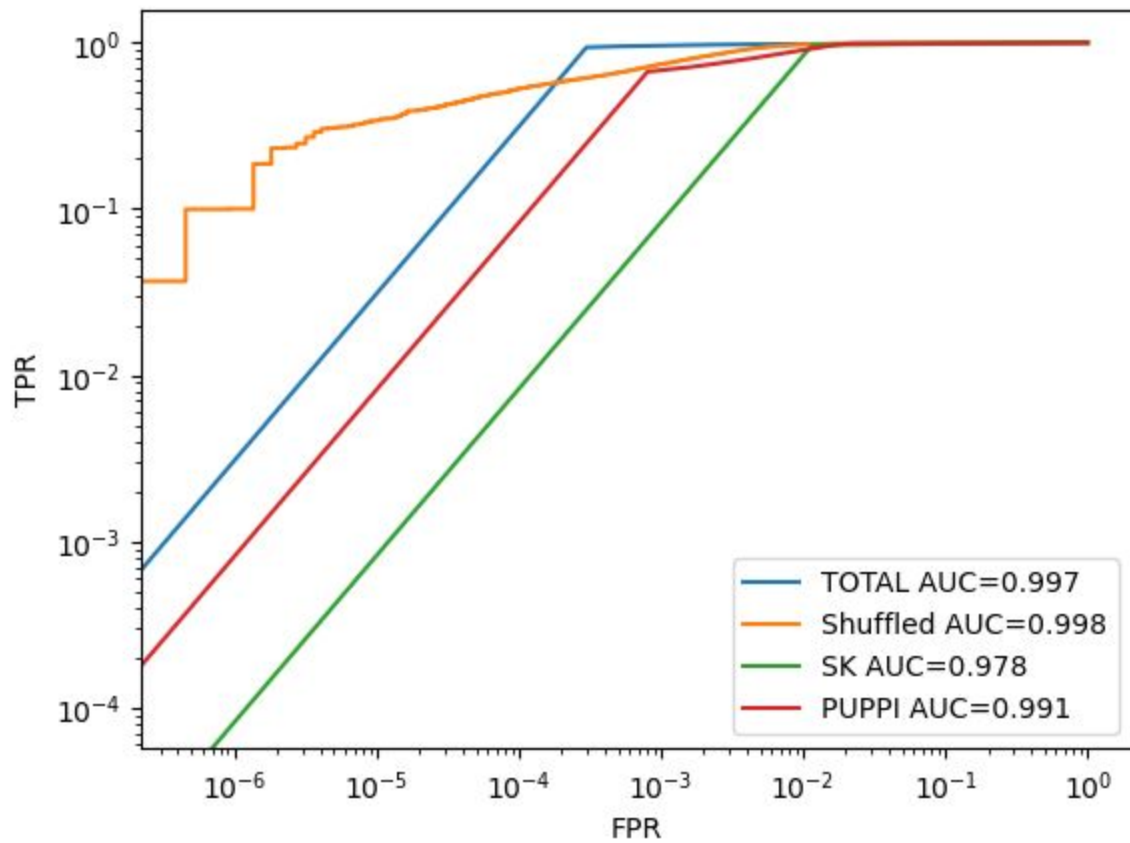
TOTAL: Training Optimal Transport with Attention Learning

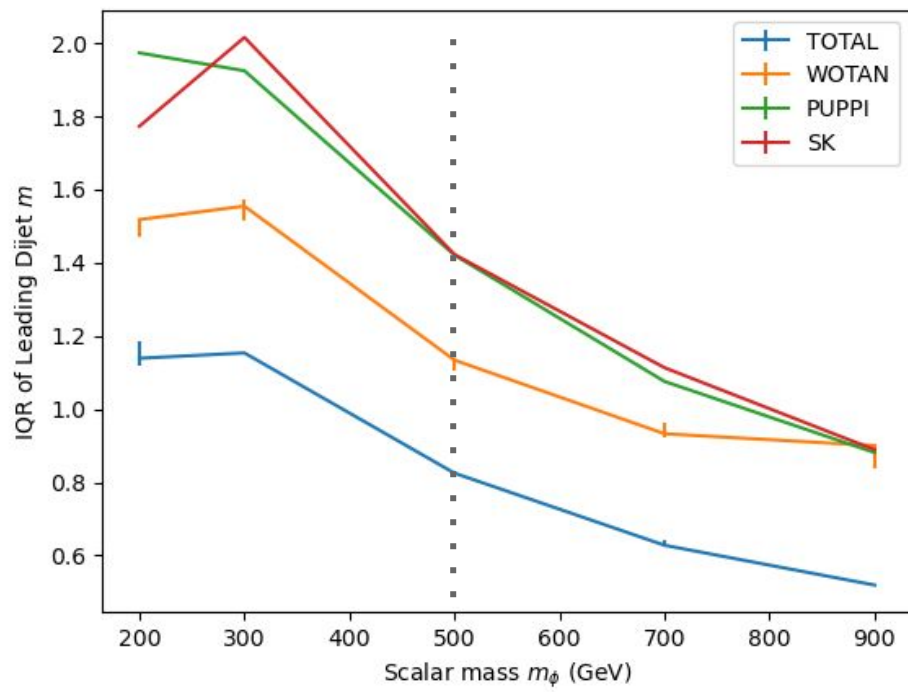
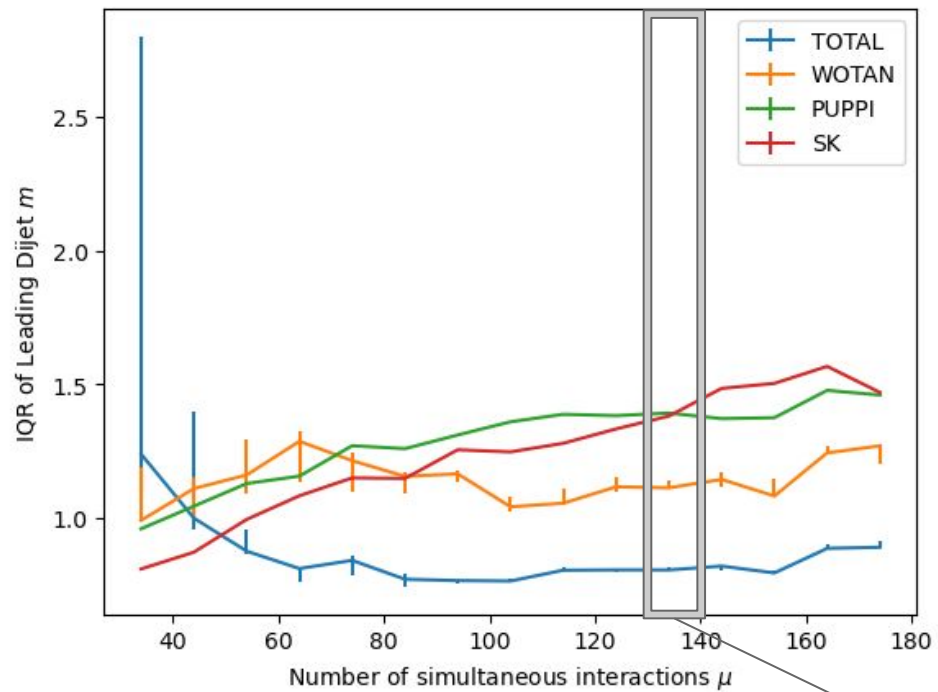
$$\mathcal{L} = \text{SWD}(x'_p, x_{np}) + \lambda \text{MSE}(p_T^{\text{miss}}(x'_p), p_T^{\text{miss}}(x_{np}))$$



Modification for jet-based dataset (PUMML)

$$\mathcal{L} = \text{SWD}(x'_p, x_{np}) + \lambda \text{MSE}(E_T(x'_p), E_T(x_{np}))$$





**Initial training ranges*