

## **BitHEP – The Limits of Low-Precision ML in HEP**

#### Daohan Wang

#### Institute of High Energy Physics (HEPHY), Austrian Academy of Sciences (OeAW)

July 8, 2025

Collaborated with Claudius Krause and Ramon Winterhalder

arXiv 2504.03387





## **Motivation**

- In NLP and LLM, a recent proposal of only using 2 or 3 discrete states in the weights matrix called BITNET has gained significant
  attention. In this way, the matrix multiplication reduces to addition only. This significantly reduces storage and computational costs of
  transformer-based large language models while maintaining same level performance.
- Transformer-based LLMs: large size, high energy consumption BitNet: 1-bit Transformer





SCIENCES



# **Bitlinear Laver**



- Weight Quantization:  $\widetilde{W}_{1b} = \text{sign}(W \langle W \rangle) \Rightarrow \{1, -1\}, \quad \widetilde{W}_{1.58b} = \max\left(-1, \min\left(1, \text{round}\left(\frac{W}{b}\right)\right)\right), \beta = \langle |W| \rangle \Rightarrow \{1, 0, -1\}.$ ۲
- Pre-Activation Quantization:  $\tilde{x} = \max\left(-Q_b, \min\left(Q_b, \operatorname{round}\left(\frac{xQ_b}{\gamma}\right)\right)\right)$ ,  $\gamma = \max(|x|)$ . ۲
- Quantized outputs during forward propagation:  $y = \widetilde{W}\widetilde{x}$ . ۲
- Rescaled outputs during back propagation:  $y = \widetilde{W}\widetilde{x} \times \frac{\beta \gamma}{\Omega_{L}}$ ۰
- We employ the 1.58-bit weights and choose an 8-bit input quantization, i.e. b = 8 and  $Q_b = 128$ . ۲





## **Model Implementation**

Binarizing other architectures also holds significant potential. We explore this potential by applying 1.58b-BITNET to benchmark the performance of various HEP applications.



Ensure that the gradient calculation in back propagation is stable and accurate

Performance

Timing

Energy Cost



Daohan Wang (HEPHY Vienna)

BitHEP

July 8, 2025





# **HEP Applications**

#### Benchmark models with 1.58b-BITNET implemented:

 $Linear \ Layer \Rightarrow BitLinear \ Layer$ 





AUSTRIAN

CUENCES



# **Classification Application: P-DAT**

#### Particle-Dual Attention Transformer (2307.04723)

M. He & D. Wang

#### **Quark/Gluon Discrimination**

AUSTRIAN ACADEMY OF SCIENCES



## **Dual Attention Mechanism**



BitHEP





# **P-DAT Model Architecture**

- Input features: log *E*, log *p*<sub>T</sub>,  $\frac{p_T}{p_{TT}}$ ,  $\frac{E}{E_I}$ ,  $\Delta \eta \Delta \phi$ ,  $\Delta R$ , PID of leading 100 particles.
- The particle attention module ( $P \times P$  attention map) and the channel attention module ( $C \times C$  attention map) are stacked while maintaining a consistent feature dimension of N = 64 and they can complement each other.
- Particle Dual Attention Transformer: 2 Feature Extractor (1 EdgeConv + 3 Conv2D + 1 AvgPool) + 2 Particle Attention modules + 2 Channel Attention modules + 1D CNN + MLP.







# **P-DAT-BIT Model Architecture**

- Replacing all the linear layers with BitLinear layers in the four attention modules of the P-DAT model (60% of the total parameters).
- All hyperparameters are identical to the non-binarized version.







### Performance

	Accuracy	AUC	Rej <sub>50%</sub>	Rej <sub>30%</sub>	Parameters	FLOPS
ParticleNet PCT LorentzNet ParT	0.840 0.841 0.844 0.849	0.9116 0.9140 0.9156 0.9203	$39.8 \pm 0.2$ $43.2 \pm 0.7$ $42.4 \pm 0.4$ $47.9 \pm 0.5$	$\begin{array}{c} 98.6 \pm 1.3 \\ 118.0 \pm 2.2 \\ 110.2 \pm 1.3 \\ 129.5 \pm 0.9 \end{array}$	370k 193.3k 224k 2.13M	540M 266M - 260M
P-DAT	0.839	0.9092	$39.2\pm0.6$	$95.1\pm1.3$	498k	144M
P-DAT-Bit	0.834	0.9040	$35.0\pm0.3$	$83.3 \pm 1.2$	498k	144M

Table: Comparison among the performance reported for P-DAT, P-DAT-Bit, and some existing classification algorithms on the quark-gluon discrimination dataset. The uncertainties in rejection rates are calculated by taking the standard deviation of 5 training runs with different random weight initialization.





### **Calibration Curves**







# **Regression Application: SMEFTNet**

#### IRC-safe and Rotation-Equivariant Graph Neural Network (2401.10323)

S. Chatterjee, S. S. Cruz, R. Schöfbeck & D. Schwarz

**Decay Plane Angle Regression** 





# **Motivation**



- Equivariant SMEFTNet targets the linear SM–SMEFT interference by using dedicated angular analyses of W/Z decay planes to resolve helicity structures modified by CP-even (O<sub>W</sub>) and CP-odd (O<sub>W</sub>) gauge self-interaction operators.
- $pp \rightarrow W(\rightarrow q\bar{q})Z(\rightarrow l\bar{l})$  MG5+Pythia+Delphes Events with  $H_T > 300$  GeV are retained. anti- $k_T$  algorithm with R=0.8
- The decay plane angle changes as the W jet rotates. To study the hadronic final states of W boson, SMEFTNet is constructed to be equivariant to azimuthal rotations of the boosted jet's constituents around the jet axis, maintaining SO(2) symmetry.
- The particle features of each event inherently encode information about the decay plane for each event, hidden within the radiation patterns mapped to the variable-length constituent vector.

SMEFTNet's goal is to serve as a surrogate model to provide an optimal observable for detecting SMEFT effects from LHC collision data. Our focus, however, is on testing BITNET's performance in regression tasks, so we limit our study to the regression of the decay plane angle.







• Input features for l=0: Four-vectors  $p_i$  of the particles.  $\mathbf{h}_{\varphi,i}^{(0)} = \varphi_i$ ,  $\mathbf{h}_i^{(0)} = \Delta R_i$ 

- Message passing function:  ${}^{i}\mathbf{m}_{j}^{(l)} = \omega_{j}^{(\mathcal{N}(i))} f_{m}^{(l)}(\hat{p}_{i}, \hat{p}_{j})$  with  $\omega_{j}^{\mathcal{N}} = \frac{p_{T,j}}{\Sigma_{k \in \mathcal{N}} p_{T,k}}$ . Particle  $j \in \mathcal{N}(i)$  of a particle i with  $\Delta R_{ij} \leq \Delta R$ .
- We demand SO(2) equivariance:  $S_{\Delta \varphi}(\mathbf{h}_{\varphi}, \mathbf{h}) = (\mathbf{h}_{\varphi} + \Delta \varphi, \mathbf{h})$ :

$$\begin{aligned} \mathbf{h}_{i}^{(l+1)} &= \sum_{j \in N(i)} \omega_{j}^{(N(i))} \mathbf{f}_{\mathbf{h}}^{(l)} \left( \mathbf{h}_{i}^{(l)}, \mathbf{h}_{j}^{(l)}, h_{\phi,i}^{(l)} - h_{\phi,j}^{(l)} \right) & \text{Invariance} \\ e^{ih_{\phi,i}^{(l+1)}} &= e^{ih_{\phi,i}^{(l)} + i\sum_{j \in N(i)} \omega_{j}^{(N(i))} f_{\Phi}^{(l)} \left( \mathbf{h}_{i}^{(l)}, \mathbf{h}_{j}^{(l)}, h_{\phi,i}^{(l)} - h_{\phi,j}^{(l)} \right) } & \text{Equivariance} \end{aligned}$$

• After *L* iterations, the global pooling is applied to sum over all the constituents with the energy-weighting. The results along with the global features x<sub>global</sub> are fed into a final MLP.

AUSTRIAN





## **Loss Function**

- Regression Target: Decay plane angle of the W boson's parton-level decay products φ<sub>i, decay</sub>
- Inputs:  $\mathbf{x}_j = \{p_{T,i}, \phi_i, \Delta R_i\}_{i=1}^{N_j}$  with AK8 jet  $p_T > 500$  GeV. 80%/20% of WZ data as train/test dataset.
- $L = \sum_{\mathbf{x}_j \in \mathcal{D}_{sim}} \sin^2 \left( \hat{f}(\mathbf{x}_j) \varphi_{j,decay} \right).$

Since the simulated data lacks the information to distinguish constituents originating from up-type and down-type quarks, although switching the positions of the two partons alters the decay plane angle by  $\pi$ , the underlying simulated data remains unchanged. Consequently, the sine function is specifically employed to speed up learning the periodicity.







### **Scatter Plots**





## **Probability Density Histograms of Residuals**



Model Comparison	Wasserstein Distance	Separation Power
SMEFTNet	0.0021	0.0001
SMEFTNet-Bit	0.4546	1.4860
SMEFTNet-Bit70	0.2528	0.6138
SMEFTNet-Bit30	0.1040	0.1503





# **Generation Application: CALOINN and CALODREAM**

#### Normalizing Flows for High-Dimensional Detector Simulations (2312.09290)

F. Ernst, L. Favaro, C. Krause, T. Plehn & D. Shih

#### Detector Response Emulation via Attentive flow Matching (2405.09629)

L. Favaro, A. Ore, S. Schweitzer & T. Plehn

#### **Fast Calorimeter Shower Simulations**





## **Motivation**

Deep generative networks based on normalizing flow provide fast and accurate surrogates for simulations in high-dimensional phase spaces by learning the underlying probability distribution of calorimeter showers from a reference dataset and then generating new samples based on this learned distribution.

The CaloChallenge was a data challenge in the HEP community, with four different datasets, increasing in their dimensionality. The goal of the challenge was to train generative networks on the datasets and to generate artificial samples as fast and precise as possible.

#### We consider two well-performing submissions of the CaloChallenge

- CALOINN based on normalizing flow.
- CALODREAM based on conditional flow matching.

#### **Evaluation Metrics**

- Low-level classifier: Phase space of the voxels in each layer.
- High-level classifier:  $\mathcal{I}_{ia}, E_i = \sum_a \mathcal{I}_{ia}, \frac{E_{dep}}{E_{inc}} = \frac{\sum_{a,i} \mathcal{I}_{ia}}{E_{inc}}, \langle l \rangle = \frac{\sum_a l_a \mathcal{I}_{ia}}{\sum_a \mathcal{I}_{ia}}, \sqrt{\frac{\sum_a l_a^2 \mathcal{I}_{ia}}{\sum_a \mathcal{I}_{ia}} \left(\frac{\sum_a l_a \mathcal{I}_{ia}}{\sum_a \mathcal{I}_{ia}}\right)^2}, \lambda_i.$





# CALOINN

- INN (coupling layer based normalizing flow) is applied for dataset 1 & 2 to sample  $p_{model}(x)$  from  $p_{latent}(r)$ .
- INN uses rational quadratic splines for dataset 1 and cubic splines for dataset 2 & 3.
- Loss Function:

$$\mathcal{L}_{\text{INN}} = -\langle \log p_{model}(x) \rangle_{\mathcal{P}_d} = -\left\langle \log p_{latent}(\overline{G}_{\theta}(x)) + \log \left| \frac{\partial \overline{G}_{\theta}(x)}{\partial x} \right| \right\rangle_{\mathcal{P}_d}.$$
(1)

- The INN is trained on the full data, conditioned on the logarithm of the incident energies.
- CaloINN consists of multiple coupling layers, each containing an independent small neural network that models an invertible transformation for high-dimensional distributions.





# **Quantization Strategies for CALOINN**

Strategy	Quantized Layers	Permutation Scheme	Notes	
Default	None	Random	Baseline from original paper. Confirms consistency with prior work.	
Exchange Permutation	None	Fixed (Exchange)	Shares permutation with BlockCentral. Ensures each dimension is transformed exactly once.	
NNCentral	Central layers in each NN	Random	Only quantizes central layers in each NN. Minimal impact on model performance.	
BlockCentral	All layers in central bijectors	Fixed (Exchange)	Float $\rightarrow$ quantized $\rightarrow$ float structure. Balances expressiveness and compression.	
All	All layers in all bijectors	Random	Most aggressive compression strategy. Highest risk of performance degradation.	



## Performance of default and BITNET CALOINN

Dataset	Setup		Quantization Low-Level AUC		High-Level AUC	
ds1–γ	reg.	Default Exchange Perm.	-	0.633(3) 0.640(4)	0.656(3) 0.651(3)	
	quant.	NNCentral BlockCentral All	8.4% 66.6% 99.9%	0.640(3) 0.680(3) 0.759(2)	0.650(2) 0.669(3) 0.828(2)	
ds1– $\pi^+$	reg.	Default Exchange Perm.	-	0.793(3) 0.784(2)	0.742(3) 0.736(3)	
	quant.	NNCentral BlockCentral All	5.9% 66.6% 99.9%	0.801(2) 0.852(1) 0.882(2)	0.751(3) 0.807(2) 0.907(2)	
ds2	reg.	Default Exchange Perm.	-	0.738(4) 0.728(6)	0.859(2) 0.857(3)	
	quant.	NNCentral BlockCentral All	0.3% 71.4% 99.9%	0.780(3) 0.950(2) 0.993(1)	$\begin{array}{c} 0.876(4) \\ 0.979(1) \\ 0.998(0) \end{array}$	





# CALODREAM

CaloDream combines Conditional Flow Matching (CFM) with transformers, consisting of:

- Energy Network: an autoregressive transformer predicting 45 layer energies.
- Shape Network: a vision transformer learning normalized shower shapes.
- Uses patching to group nearby voxels, enabling scalability to large detectors like DS3.
- Self-attention captures spatial correlations and sparsity in calorimeter showers.
- CFM enables efficient training of continuous normalizing flows.
- Unlike CaloINN, which uses many small NNs, CaloDream uses two large networks.
- This design makes it **less sensitive to quantization**, as fewer components are involved.
- The two networks can be quantized **independently and flexibly**.





# **Quantization Strategies for CALODREAM**

- Embedding: 4-head transformer, 64-dim embeddings.
- CFM network: 8-layer MLP with 256 hidden units.
- Quantization Options:
  - Regular: No quantization (baseline).
  - **Quantized**: Quantize 6 out of 8 hidden layers in the MLP (CFM), resulting in 66.09% of the energy network being quantized, corresponding to 5.54% of total CALODREAM parameters.
- Operates on 135 patches of 48 voxels (total: 6480 voxels).
- Uses time, position, and conditional embeddings  $\rightarrow$  6 ViT blocks  $\rightarrow$  final MLP.
- Quantization Options:
  - **Regular**: No quantization (baseline).
  - **No Embedding**: Quantize QKV, projections, and MLPs in ViT blocks (63.8% quantized).
  - Full: Additionally quantize embedding layers (66.22% quantized).

AUSTRIAN ACADEMY OF SCIENCES



# Performance of regular and BITNET CALODREAM

Energy net	Shape net	Quantization	low-level AUC	high-level AUC
regular	regular	-	0.531(3)	0.523(3)
quantized		5.5%	0.532(3)	0.525(3)
regular	no embedding	58.5%	0.611(2)	0.543(3)
quantized		63.9%	0.610(5)	0.545(2)
regular	full	60.7%	0.735(4)	0.942(2)
quantized		66.2%	0.738(4)	0.944(3)

Table: Performance of regular and BITNET CALODREAM using the classifier metric of the CaloChallenge. Uncertainties show the standard deviation over 10 random initializations and trainings of the classifier on the same CALODREAM sample.





# **Summary and Outlook**

- We demonstrated that **quantization-aware training (QAT)** enables competitive performance across classification, regression, and generative tasks in HEP, especially for classification tasks such as quark–gluon tagging.
- Larger networks can be quantized more easily, maintaining performance even with 60% of weights quantized.
- Selective quantization (targeting specific layers) significantly improves performance compared to full model quantization.
- Self-attention layers in transformer-based architectures are particularly robust to low-bit quantization.
- QAT aligns with future **hardware and energy constraints**, offering a promising direction for scalable and efficient ML at the HL-LHC and beyond.
- Future directions include **optimizing quantization configurations**, **integrating QAT into large-scale models**, and adapting BITNET for fully quantized, low-precision inference—with **measurable gains in energy, memory, and latency**. Extending QAT to **real-time tasks**, such as the LHC trigger system, may enable fast, accurate inference on resource-constrained hardware like FPGAs.