# The OmniJet-α Foundation Model

**Extending to Calorimeter Showers** 

# 8th Round table on Deep Learning

November 14th 2025

Joschka Birk, Frank Gaede, Anna Hallin, Gregor Kasieczka, Martina Mozzanica, **Henning Rose** 



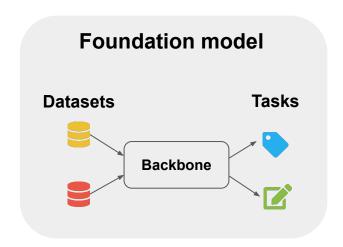
**CLUSTER OF EXCELLENCE** 

QUANTUM UNIVERSE

### What is OmniJet-α?

First cross-task Foundation model for jet-physics<sup>1</sup>

- But we aim to build a foundation model for all particle physics
- Applied on Calorimeter showers to show Adaptability of the architecture<sup>2</sup>



<sup>1</sup> Birk et al. "OmniJet-α: The first cross-task foundation model for particle physics" 2403.05618 (2024)

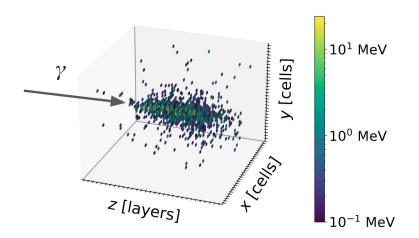
<sup>2</sup> Birk et al. "OmniJet- $\alpha_c$ : Learning point cloud calorimeter simulations using generative transformers" 2501.05534 (2025)

## Why Calorimeter Showers?

- γ-showers in the ECAL in the proposed
  International Large Detector (ILD) [4]
  = long sequences, sparse spatial structure
- Perfect stress test for architecture scaling

#### **Shower Properties**

- $\gamma$ -energy 10-100 GeV
- x, y, z in 3D grid (30<sup>3</sup> voxels)
- hit energy 0-13 MeV
- 100-1700 hits per shower



3

~30 times longer than jets

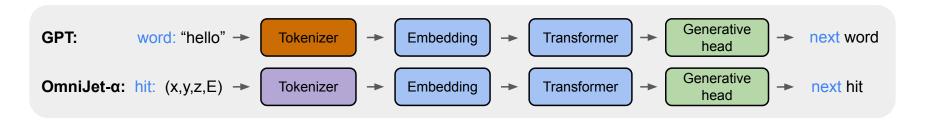
[5]

<sup>4</sup> ILD Collaboration "International Large Detector: Interim Design Report" 2003.01116 (2020)

<sup>5</sup> Buhmann et al. "Getting High: High Fidelity Simulation of High Granularity Calorimeters with High Speed" 2005.05334 (2020)

### OmniJet-α architecture

- Autoregressive decoder-only Transformer (GPT-1 style)<sup>3</sup>
- Predicts the next hit-token, one at a time, forming showers sequentially
- GPTs use Byte-Level Tokenizers (codebooksize ~40k)
- VQ-VAE acts as a Tokenizer: it turns continuous calorimeter data into discrete hit-tokens.

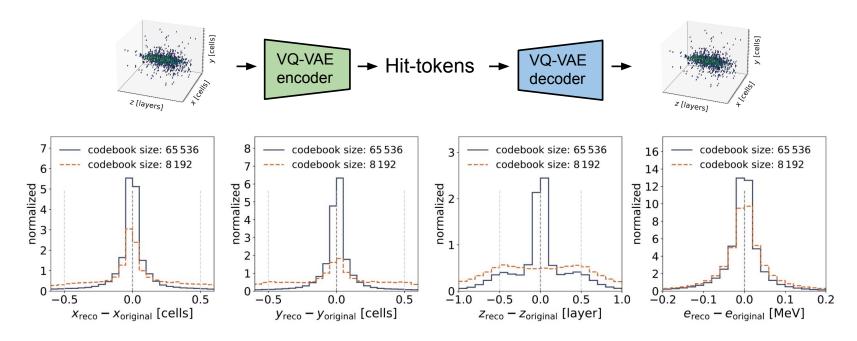


**Problem:** Input data is continuous → VQ-VAE introduces loss of information

<sup>3</sup> Radford et al, "Improving language understanding by generative pre-training" (2018)

## Loss of information

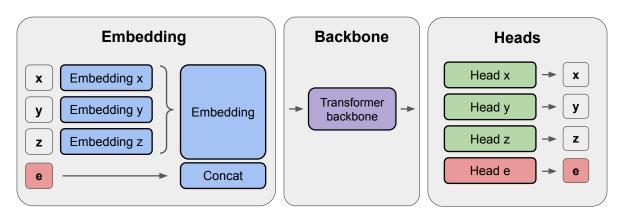
Going from shower to token-space and back, introduces an error



## **Upgrades to OmniJet-α**

#### **Direct Spatial Tokenization:**

- Instead of compressing hits to tokens, encode coordinates directly in the embedding
- Full spatial precision preserved
- Simpler pipeline → closer to a universal architecture

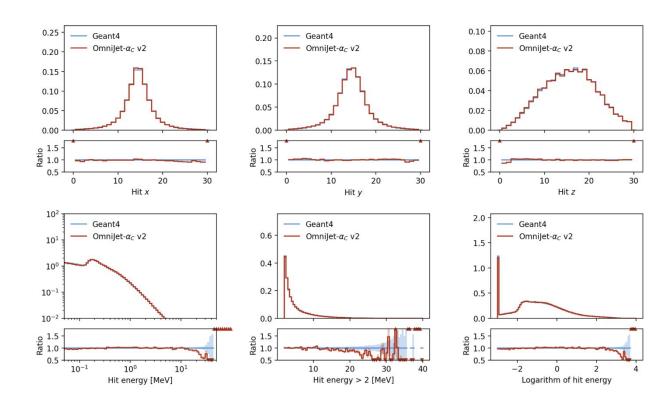


Use Mixture of Gaussians to sample continuous energy values

## Results

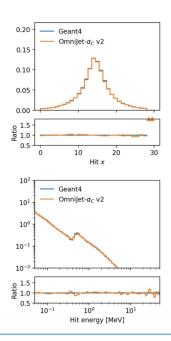
#### **Hit-Level Distributions:**

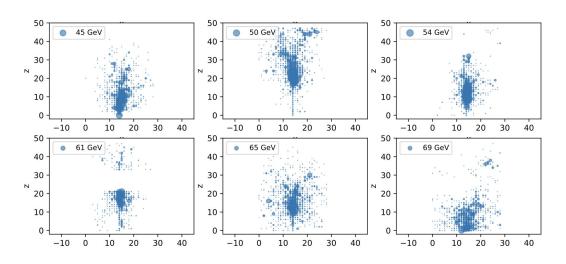
- Spatial distributions are well reproduced
- Hit-energy spectrum matches Geant4 over3 orders of magnitude



## **Results - HCAL**

OmniJet-α can also be trained to model pion showers in an HCAL





Plots produced by Martina Mozzanica

## Conclusion

- Direct spatial tokenization removes Tokenization bottlenecks
  Full spatial precision, simpler pipeline
- High-fidelity generation across ECAL & HCAL
  Spatial and energy distributions match Geant4
- Architecture remains flexible for future extensions
  Same framework can incorporate additional detectors without major changes

9

## Results

#### **Global-Level Distributions:**

 Global distributions are well reproduced

