# Toward a General-Purpose Foundation Model for Neutrino Physics



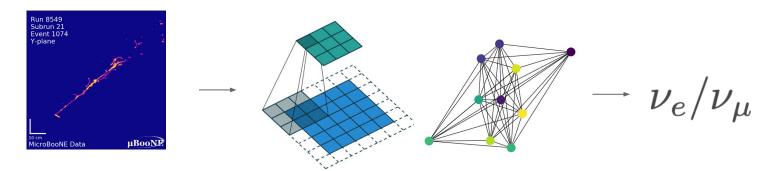
Samuel Young // NPML 2025
youngsam@stanford.edu | youngsm.com





# Deep learning in neutrino physics

**Current approach:** deep neural network(s) trained used simulated datasets with labels via fully supervised learning.

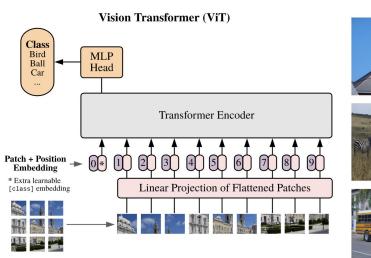


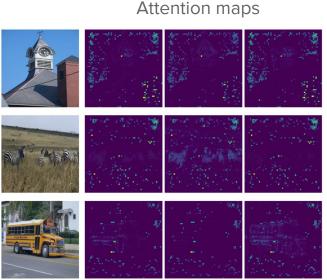
#### Extremely successful, but:

- Vulnerable to data-simulation discrepancies ⇒ heavy calibration efforts
- "Smart vs. big" models: can removing hand-crafted physical priors result in better performance with the introduction of more compute? [1]
- Task-specific, trained separately from scratch.

# Task specificity

Cat/dog classifier will not learn anything about the difference between trees and flowers.



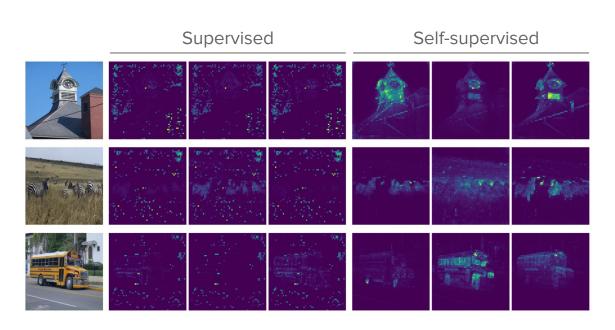


Attention maps from image classification in a vision transformer DINO (2104.14294)

$$Attention(Q, K, V) = \underbrace{\text{softmax}(\frac{QK^T}{\sqrt{d_k}})}_{\text{scores}} V$$

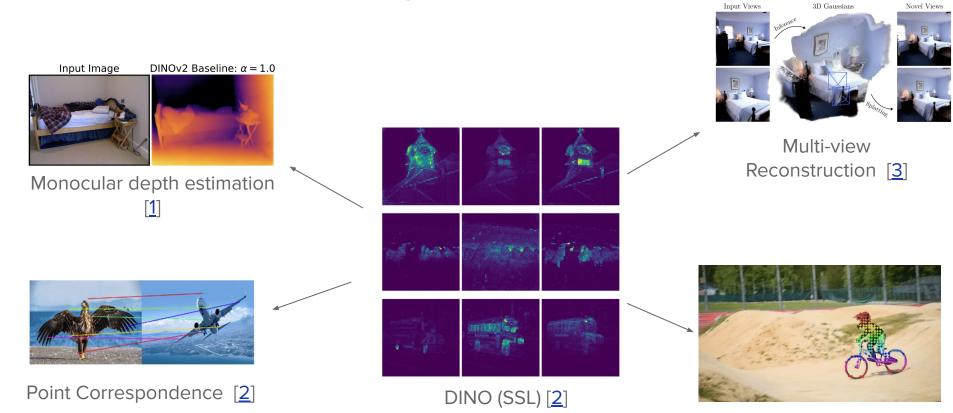
# Foundation models are generalists

#### **FM** = learn more than the task requires so you can reuse it later

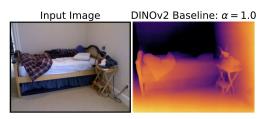


Attention maps from image classification and self-supervised tasks in a vision transformer DINO (2104.14294)

# Foundation models are generalists



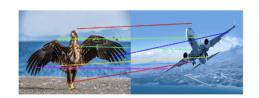
# Foundation models are generalists



Monocular depth estimation [1]



Let's apply to LArTPC data



Point Correspondence [2]



DINO (SSL) [2]



Multi-view Reconstruction [3]



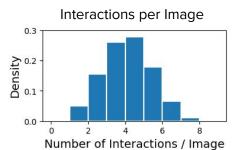
Video Tracking [4]

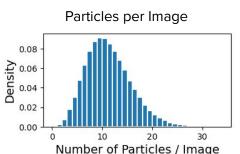
#### Dataset: PILArNet-Medium

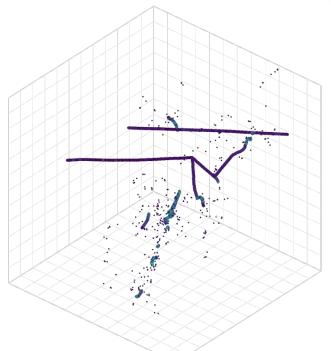




- Simulated dataset of 1.2M 3D events
- $(2.3 \text{ m})^3$  cube  $(768 \text{ px})^3$ . ~5B non-zero voxels.
- +1M events on top of previous open dataset, <u>PILArNet (2020)</u>.
- Simply 3D energy depositions, equivalent to "digital hits" from a LArTPC (e.g., DUNE Near Detector)
- 1024 30,000 voxels/event





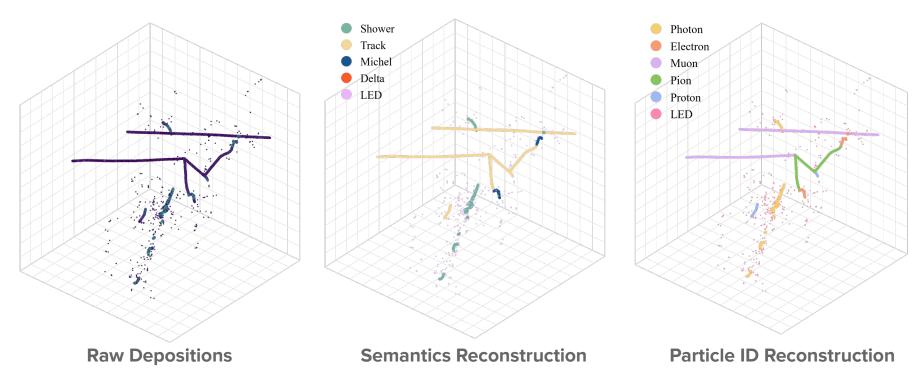


**Raw Depositions** 

# Semantic Segmentation Labels



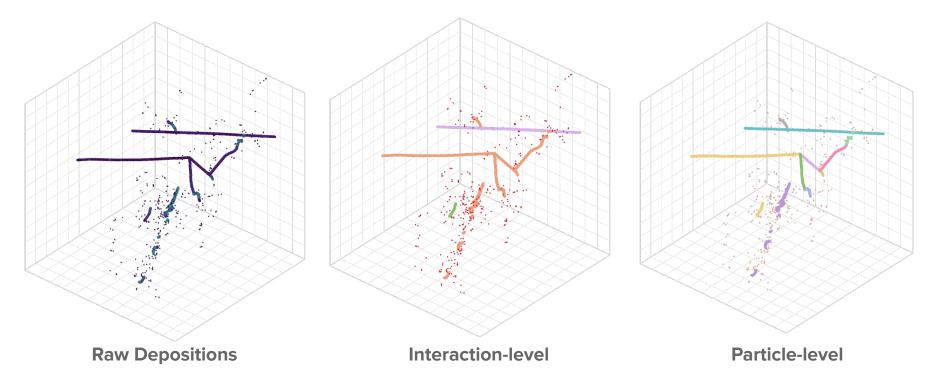




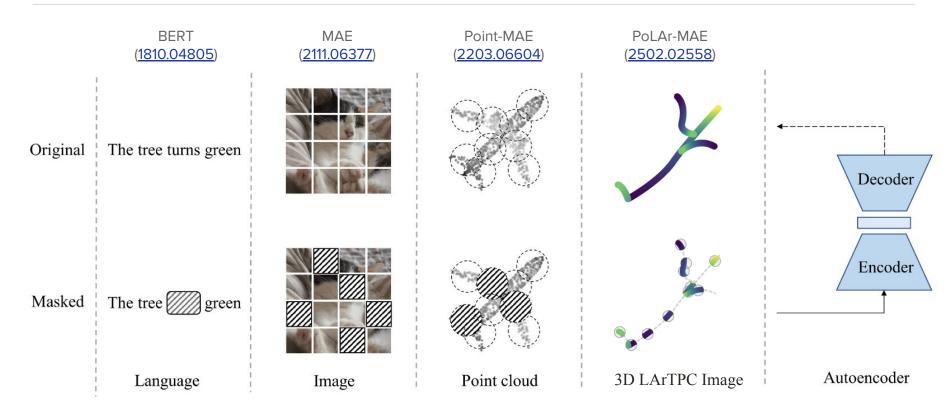
### Instance Labels



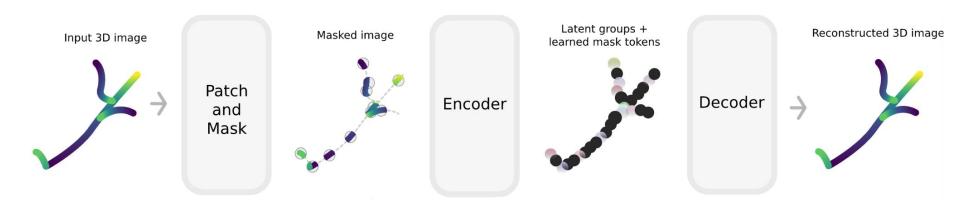




#### Masked Autoencoders

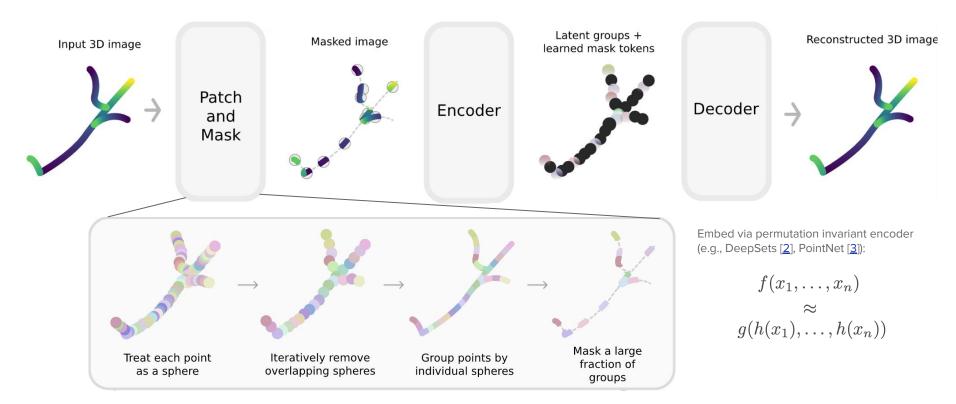


### Polar-MAE: Point-based LAr Masked Autoencoder [1]



Encoder-decoder is asymmetric, i.e. encoder params ≫ decoder params.

### Polar-MAE: Point-based LAr Masked Autoencoder [1]



### Polar-MAE: Point-based LAr Masked Autoencoder [1]

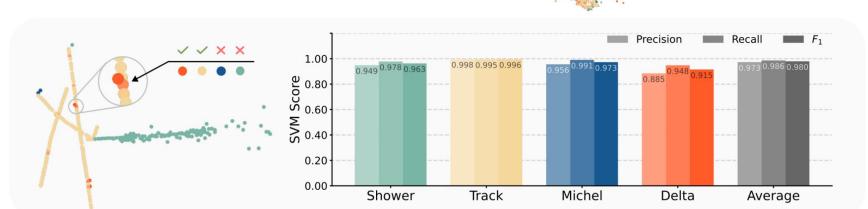


# Patch Representations

A look at patch representations.

Remember: one patch contains a group of pixels, so can contain >1 particle type.

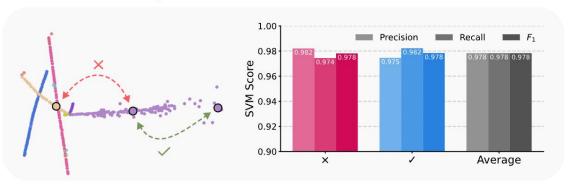
#### Patch makeup semantic segmentation





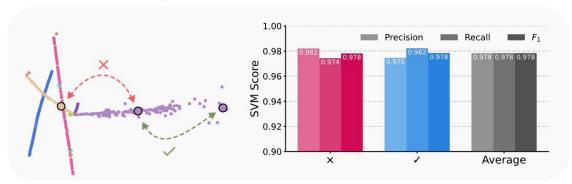
### Instance and Vertex Patch Classification

#### **Instance Sharing**

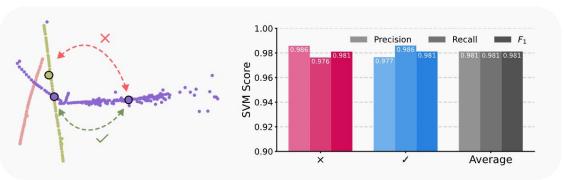


#### Instance and Vertex Patch Classification

#### **Instance Sharing**

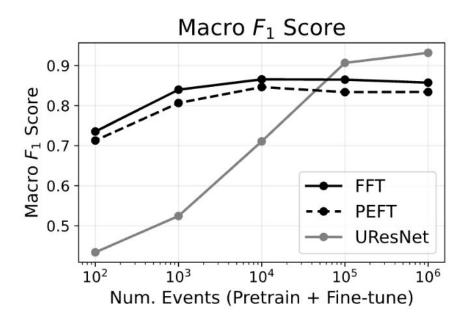


#### **Vertex Sharing**



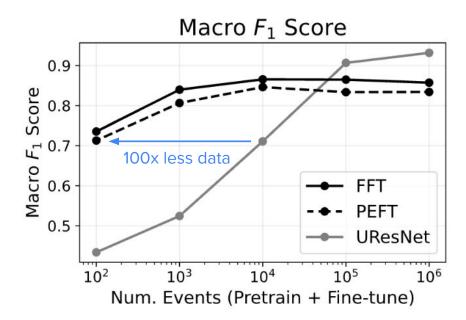
#### What we care about: per-pixel classification

 Beats state-of-the-art in data-constrained environment, but not in the limit of many events.



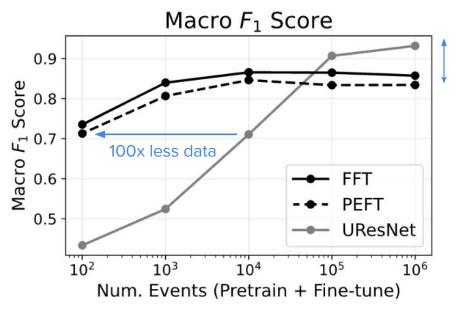
#### What we care about: per-pixel classification

 Beats state-of-the-art in data-constrained environment, but not in the limit of many events.



#### What we care about: per-pixel classification

 Beats state-of-the-art in data-constrained environment, but not in the limit of many events.



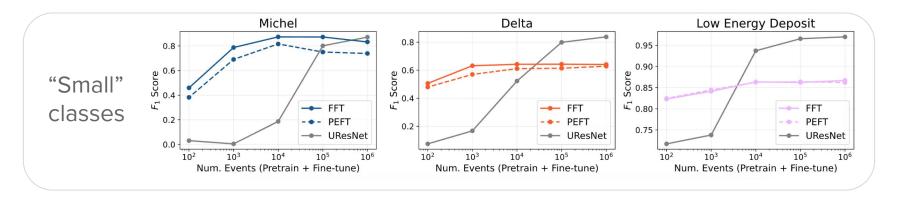
does not beat UResNet at high event counts.

→ fundamental limit in PoLAr-MAE architecture.

- Small features poorly modeled, i.e. "paint brush" classification.
- This is due to single-scale patches being used, which smears tiny structures.



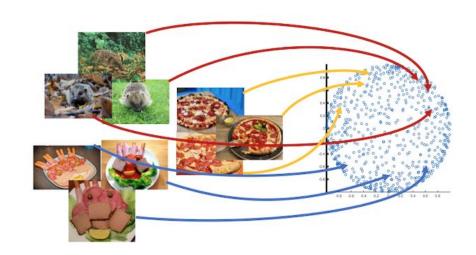


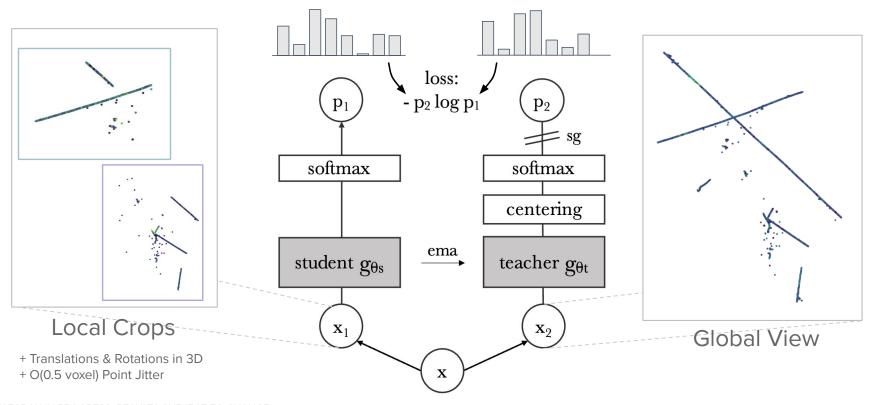


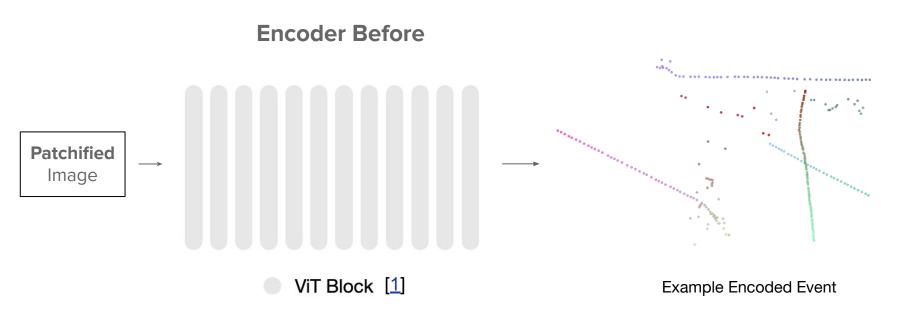
instead of reconstructing masked portions of image directly,

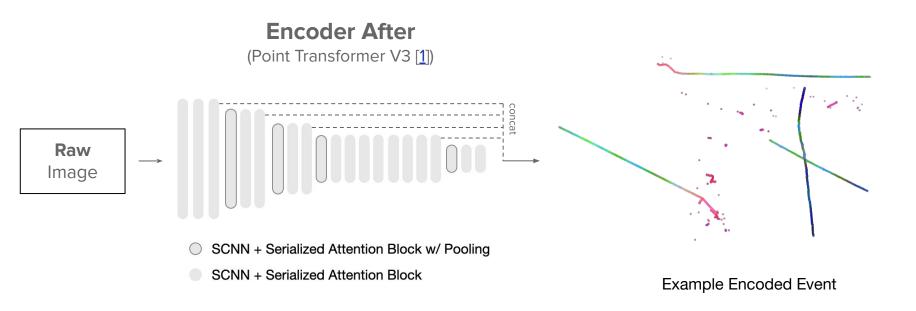
let's predict where they would would end up on a unit sphere
(i.e., classify them),

by enforcing consistency between global and local views of the same image.



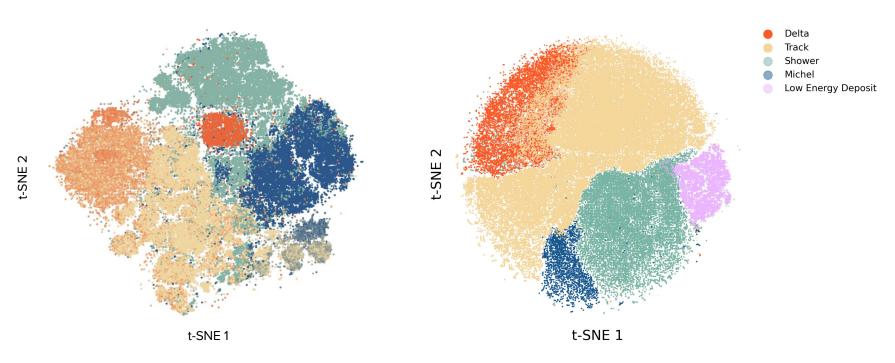




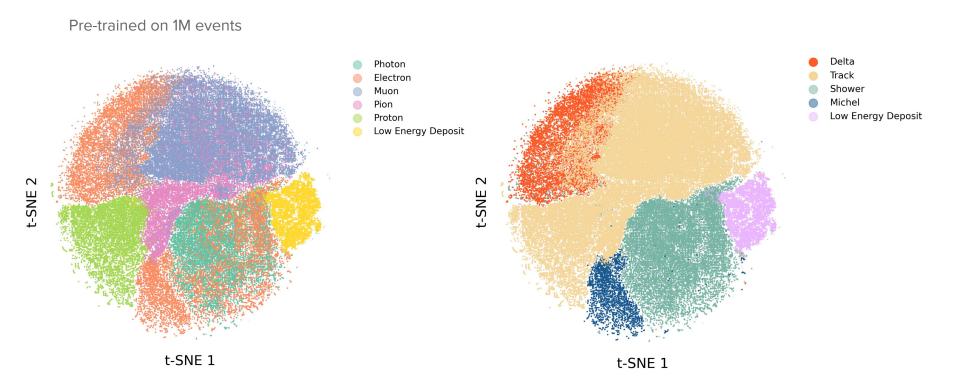


### Structure within the Feature Manifold

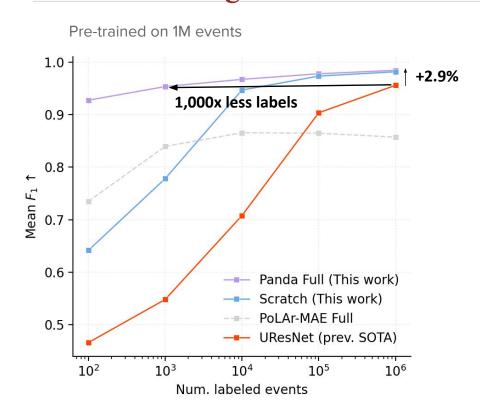
#### Pre-trained on 1M events

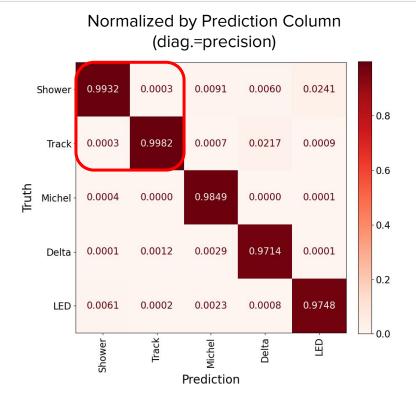


### Structure within the Feature Manifold

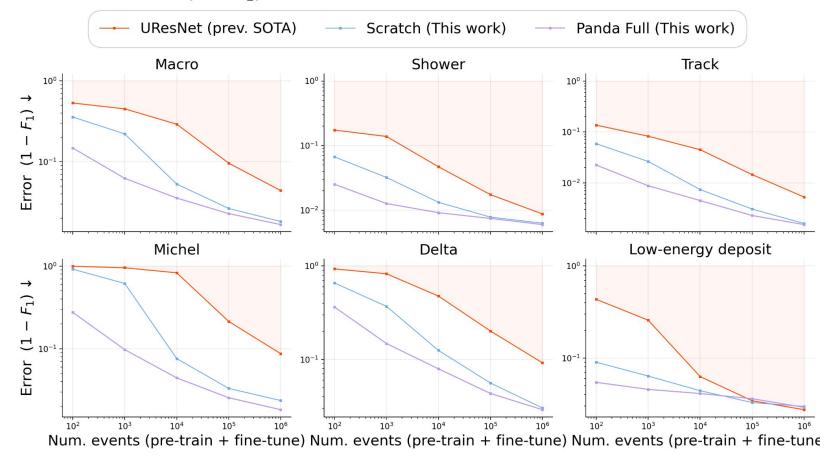


# Semantic Segmentation: Motif



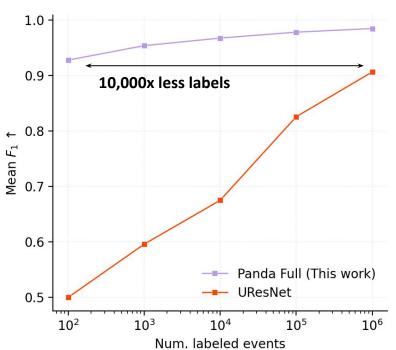


#### Error $(1 - F_1)$ vs. fine-tune set size — Pre-train + Fine-tune

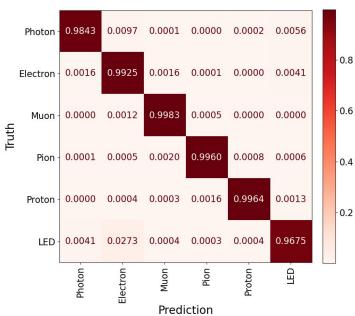


# Semantic Segmentation: Particle

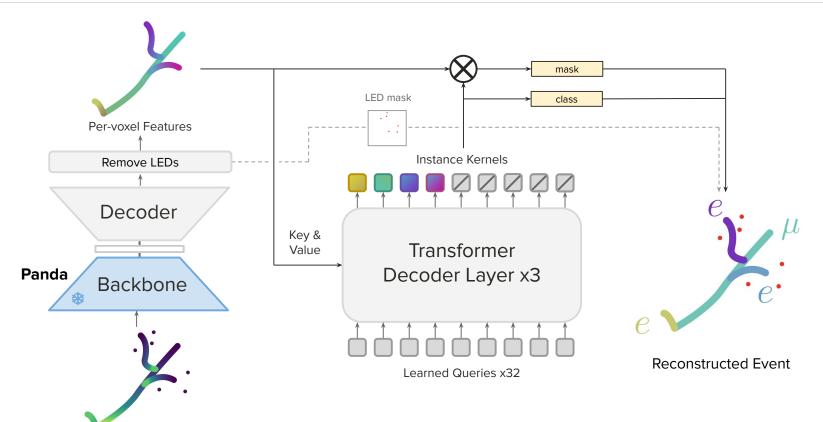




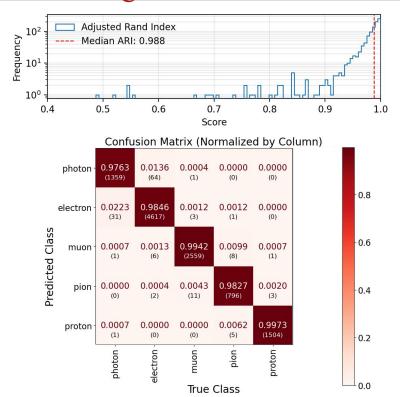
# Normalized by Truth Column (diag.=recall/efficiency), 1M fine-tune

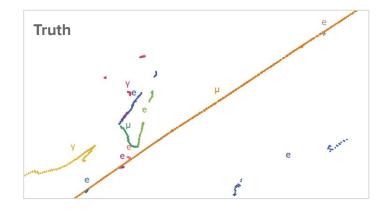


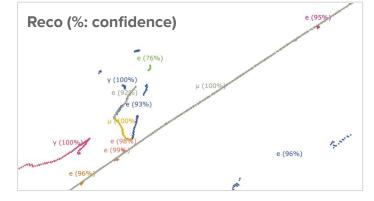
# Instance Segmentation: separating particles from one another



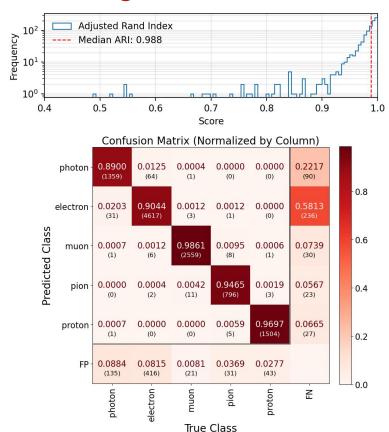
# Instance Segmentation: Particle

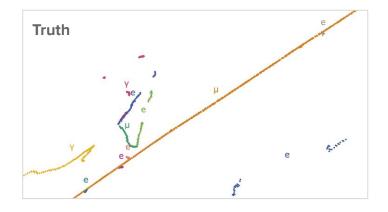


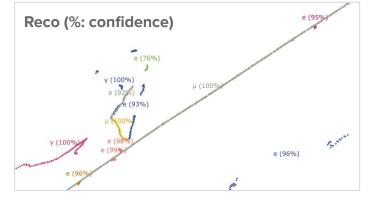




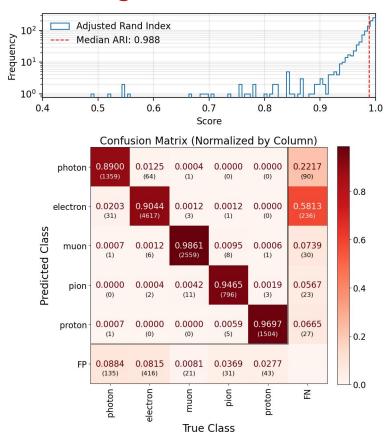
# Instance Segmentation: Particle

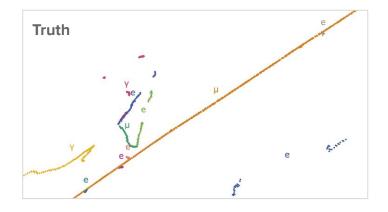


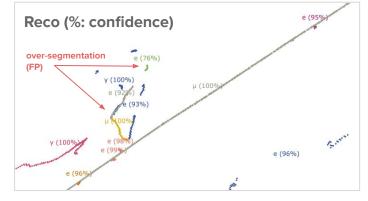




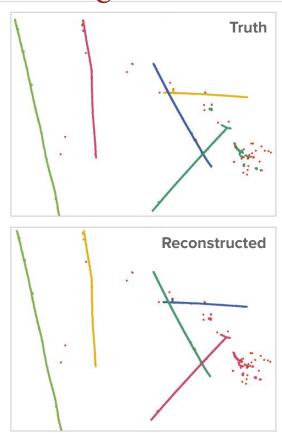
# Instance Segmentation: Particle

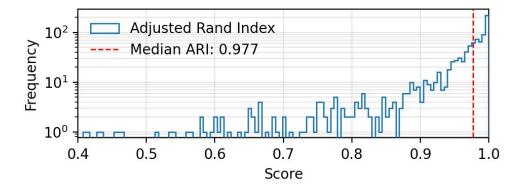






# Instance Segmentation: Interaction





# Takeaways

- Self-supervised learning works, and we're just getting started.
- A generic feature extractor unlocks new possibilities that were simply not possible before:
  - Few-shot learning w/o well-calibrated sim: track/shower, Michel tagging, particle ID, ...
  - **Reasoning over images/captioning** with language (human-in-the-loop)
  - Content-retrieval at scale: "find events like this" in this dataset.
  - **Cross-experiment datasets →** invariant embeddings across detector conditions, easy adaptation.
- **Future paths**: multi-modal sensors (optical, 2D projection, ...), incorporating (neuro-)symbolic reasoning, cross-experiment datasets, detector-related effects.
- Reproducing the results of other models is really hard and time-consuming (specific data formats, dependencies, closed-source). Common benchmark(s) would go a long way in accelerating the work we do.

#### The future is exciting!

### Extras

## Sharpening + Centering

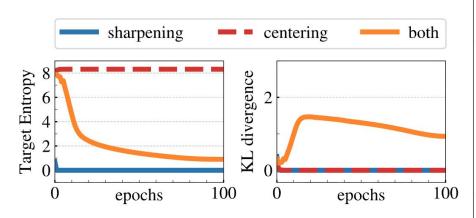


Figure 7: **Collapse study.** (**left**): evolution of the teacher's target entropy along training epochs; (**right**): evolution of KL divergence between teacher and student outputs.

There are two forms of collapse: regardless of the input, the model output is uniform along all the dimensions or dominated by one dimension. The centering avoids the collapse induced by a dominant dimension, but encourages an uniform output. Sharpening induces the opposite effect. We show this complementarity by decomposing the cross-entropy H into an entropy h and the Kullback-Leibler divergence ("KL")  $D_{KL}$ :

$$H(P_t, P_s) = h(P_t) + D_{KL}(P_t|P_s).$$
 (5)

A KL equal to zero indicates a constant output, and hence a collapse.

#### **Sharpening:**

$$P_s(x)^{(i)} = \frac{\exp(g_{\theta_s}(x)^{(i)}/\tau_s)}{\sum_{k=1}^K \exp(g_{\theta_s}(x)^{(k)}/\tau_s)},\tag{1}$$

with  $\tau_s > 0$  a temperature parameter

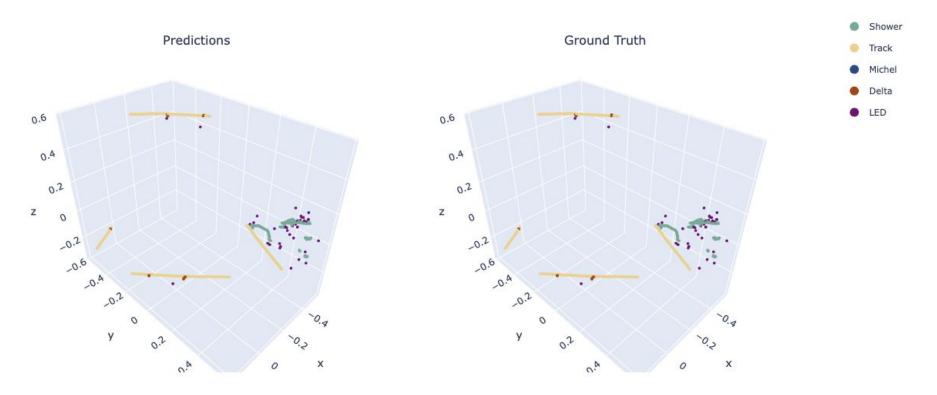
### Augmentations

What about diffusion/attenuation?

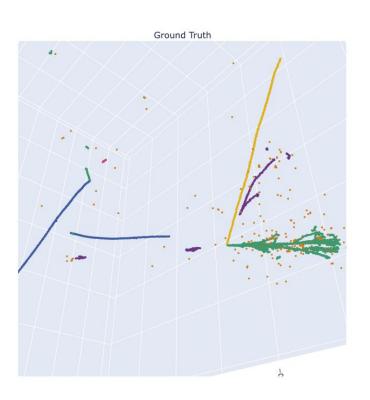
#### Linear Evaluation on pre-training – 1M events

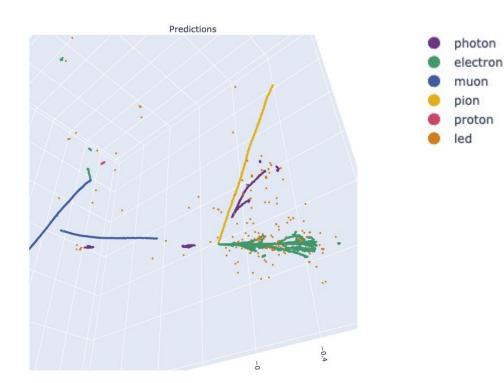


## Example Sem Seg - Motif

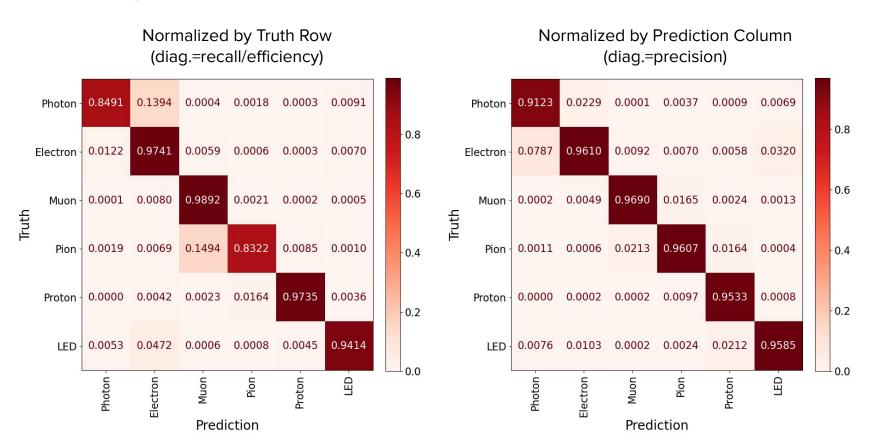


# Example Sem Seg - PID

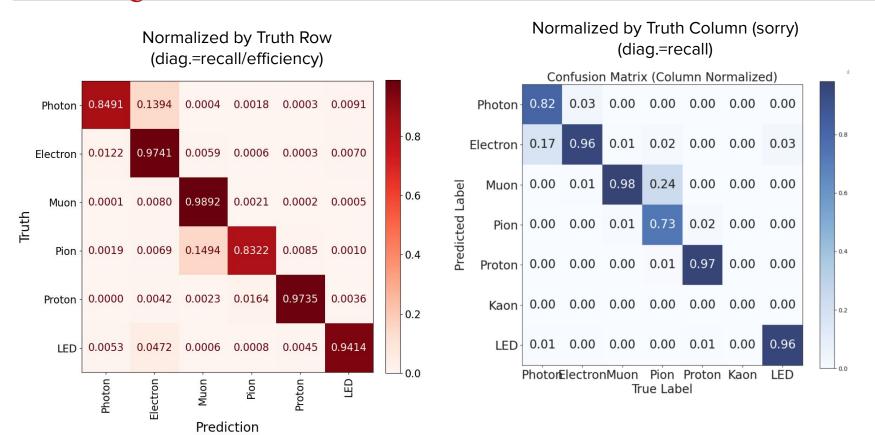




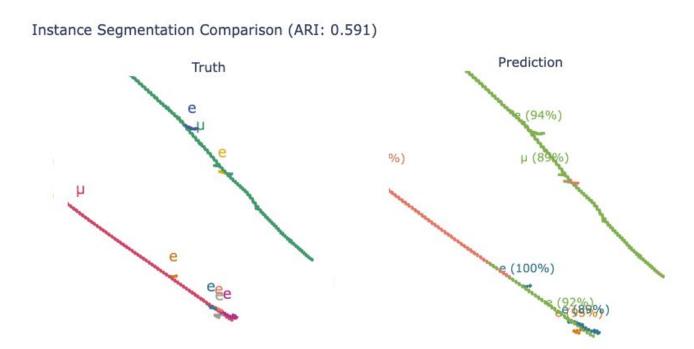
## Sem Seg PID – 100 events



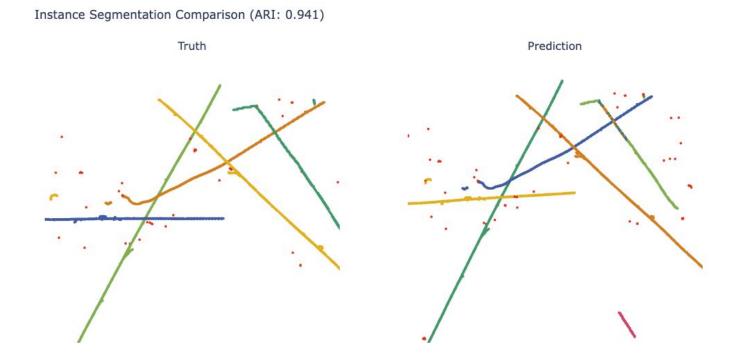
### Sem Seg PID – UResNet Confusion



### Poor instance reconstruction – particle

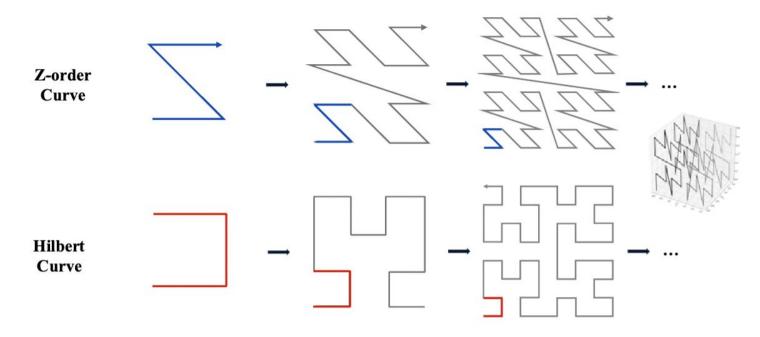


### Poor instance reconstruction - interaction



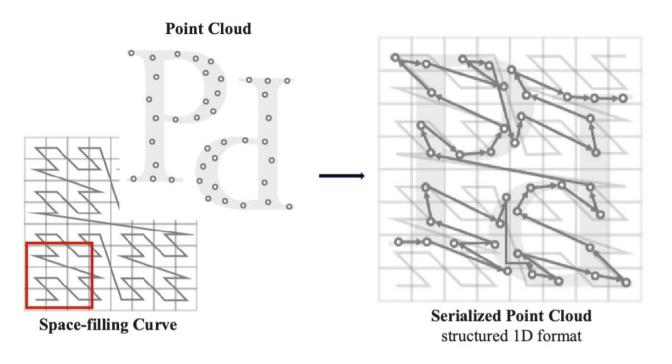
### Serialized Attention

See Point Transformer V3 paper (arXiv:2312.10035) for more detail



### Serialized Attention

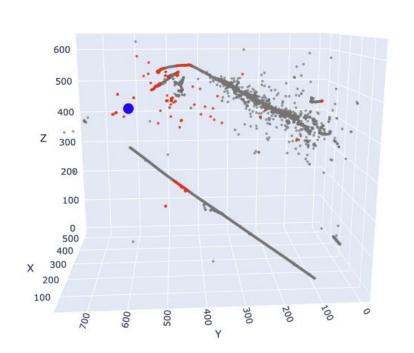
See Point Transformer V3 paper (arXiv:2312.10035) for more detail

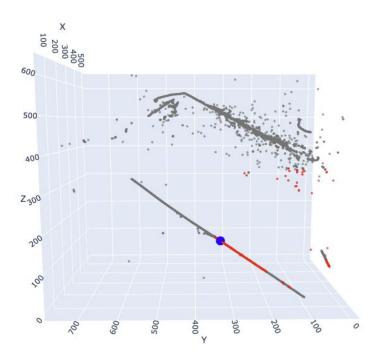


### Serialized Attention – 256 voxel patches

Receptive field at a single point at stage 0

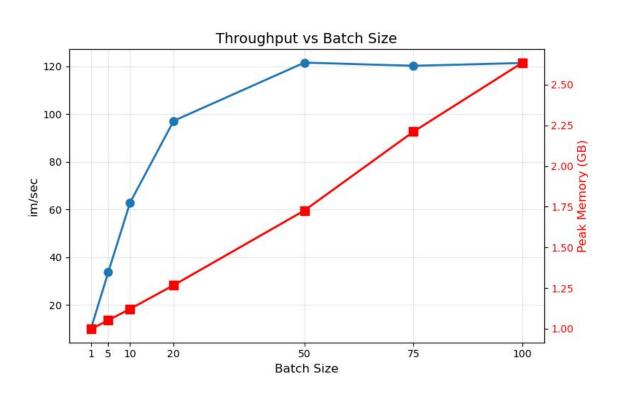
not perfect, but good enough with enough depth.





## Scalability

Semantic Segmentation (measured on single A100)



#### Peak throughput:

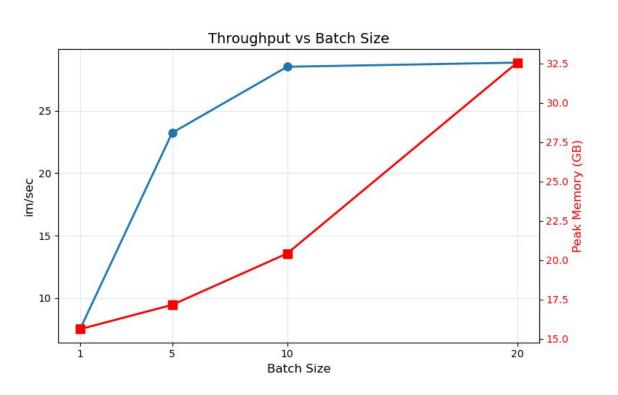
8.3 ms/image 120 img/sec

### Mem@peak:

~1.75 GB

### Scalability

Instance/Panoptic Segmentation (model forward + NMS post-processing)



### Peak throughput:

34.5 ms / image 29 img/sec

### Mem@peak:

~20 GB

NMS post-processing is serial, but parallelized via multiprocessing

## Scaling Model Params

