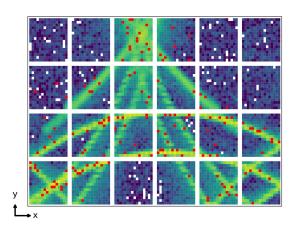
Fast High Fidelity Simulation of Cherenkov Detectors at EIC



James Giroux
Digital Twins for Nuclear and Particle physics
10/07/2025



Outline

- The High Performance DIRC (hpDIRC) at EIC
- Generative Models for Fast Simulation of Cherenkov Detectors at EIC (arXiv:2504.19042, MLST)
 - Hit-level learning for generative AI
 - A collection of SOTA generative models
 - Holistic Simulation pipeline Generating the photon yield
 - o Takeaways
- Towards Foundation Model for Readout Systems combining Discrete and Continuous Data (arXiv:2505.08736)
 - Towards FM in physics
 - Potential Issues with Tokenization
 - Combining Continuous and Discrete Data
 - Conditional Generation through prepended context
 - Class Conditional Generation through conditional computing Mixture of Experts
 - Takeaways



DATA SCIENCE

hpDIRC at EIC

Barrel geometry

- 16-sided polygonal barrel around the beam line (R=1m)
- Divided into optically isolated sectors a bar box and a readout box
- Each bar box contains eleven fused silica radiator bars (~ 4m in length) mirrored ends for photon reflection
- Exiting photons are focused by a 3-layer spherical lens

Pixelated Detector Plane

- o 4x6 PMTs 16x16 pixels per PMT
- Provide spatial and timing information (100ps)

Operation Requirements

 \circ 3σ separation for π/K at 6 GeV/c

Cherenkov Photon
Trajectories
Rediator Bar
Mirror
Particle
Optics
Track

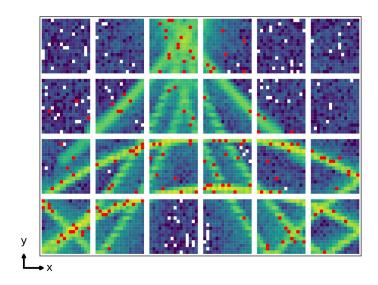
Collider

Image taken from [1]

CHARTERED 1693

[1] Kalicy G 2022 Developing high-performance DIRC detector for the Future Electron Ion Collider Experiment (arXiv:2202.06457) URL https://arxiv.org/abs/2202.06457

Generative Models for Fast Simulation of Cherenkov Detectors at the Electron Ion-Collider



[1] Giroux, James, Michael Martinez, and Cristiano Fanelli. "Generative Models for Fast Simulation of Cherenkov Detectors at the Electron-Ion Collider." *arXiv:2504.19042* (2025). (Accepted into IOP - Machine Learning Science and Technology)



Learning at the Hit Level

- Difficulties in working with Cherenkov detectors for Generative AI
 - Lack of fixed input sizes dynamic photon yield dependent on kinematic parameters
 - Pixelated (discrete) spatial readout system algorithms are designed for continuous spaces

Learning at the Hit Level

- Difficulties in working with Cherenkov detectors for Generative AI
 - Lack of fixed input sizes dynamic photon yield dependent on kinematic parameters
 - Pixelated (discrete) spatial readout system algorithms are designed for continuous spaces
- Abstract away from fixed input sizes
 - o Remain agnostic to the photon yield
- Learning at the hit level, conditional on $< |\mathbf{p}|$, $\boldsymbol{\theta} >$
 - Treating individual Cherenkov photons in a track as ~ independent

Learning at the Hit Level

- Difficulties in working with Cherenkov detectors for Generative AI
 - Lack of fixed input sizes dynamic photon yield dependent on kinematic parameters 0
 - Pixelated (discrete) spatial readout system algorithms are designed for continuous spaces 0
- Abstract away from fixed input sizes
 - Remain agnostic to the photon yield

- $D_{i,j} = \begin{cases} \lfloor M_{PMT.}/6 \rfloor \cdot 16 + \lfloor N_{pixel.}/16 \rfloor \\ (M_{PMT.} \% 6) \cdot 16 + (N_{pixel.} \% 16) \end{cases}$
- Learning at the hit level, conditional on $< |\mathbf{p}|$, $\boldsymbol{\theta} >$
 - Treating individual Cherenkov photons in a track as ~ independent
- Use physical sensor dimensions to remove discrete representation in space
 - DIRC readout has a fixed "row,col" coordinate system (1) 0
 - Transform to x,y coordinate system (mm) (2) 0
 - Smear uniformly over individual PMT pixels 0

$$x = 2 + D_j \cdot p_{width.} + (M_{PMT.} \% 6) \cdot \text{gap}_x + \frac{1}{2} p_{width.}$$

$$y = 2 + D_i \cdot p_{width.} + |M_{PMT.}| / 6| \cdot \text{gap}_x + \frac{1}{2} p_{width.}$$

$$y = 2 + D_i \cdot p_{height.} + \lfloor M_{PMT.} / 6 \rfloor \cdot \text{gap}_y + \frac{1}{2} p_{height.}$$



Learning at the Hit Level Cont'd...

What does this look like during training?

ID	x (mm)	y (mm)	t (ns)	p	$oldsymbol{ heta}$
1				3.0	5.0
1				3.0	5.0
N				4.0	7.0
N				4.0	7.0

Learning at the Hit Level Cont'd...

• What does this look like during training?

ID	x (mm)	y (mm)	t (ns)	p	$oldsymbol{ heta}$
1				3.0	5.0
1				3.0	5.0
N				4.0	7.0
N				4.0	7.0

- What does this look like at inference (generation)?
 - Our models are trained to generate **individual** photons
 - We aggregate multiple forward calls to generate **tracks**
 - \circ These are **not** sequential batch processing of N_{ν}

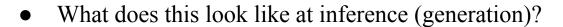
Generated Photon: 1



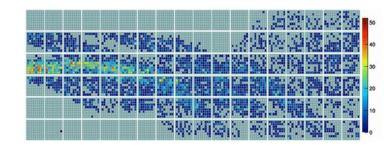
Learning at the Hit Level Cont'd...

What does this look like during training?

ID	x (mm)	y (mm)	t (ns)	$ \boldsymbol{p} $	$\boldsymbol{\theta}$
1				3.0	5.0
1				3.0	5.0
N				4.0	7.0
N				4.0	7.0



- Our models are trained to generate **individual** photons
- We aggregate multiple forward calls to generate **tracks**
- \circ These are **not** sequential batch processing of N_{γ}



Generated Photon: 1

Integrate over tracks to create PDF



A collection	of SOTA	Generative	Models
		The second secon	

A	collection	of SOTA	Generative	Model

Transformation

Objective

$$x_k =$$

$$x_k = f_{\theta}(z, k) = f_{\theta_N} \circ f_{\theta_{N-1}} \circ \dots f_{\theta_1}(z_0, k)$$

Continuous Normalizing Flows

Conditional Flow Matching

Denoising Diffusion Probabilistic Models

Score Based Generative Models

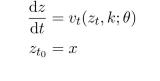
$$\frac{\mathrm{d}z}{\mathrm{d}t} = v_t(z_t, k; \theta)$$

$$min_{\theta} \mathbb{E}_{x \sim p_{data}(x)}[-\log p_{\theta}(x|k)]$$

 $min_{\theta} \mathbb{E}_{x \sim p_{data}(x)}[-\log p_{\theta}(x|k)]$

 $min_{\theta} \mathbb{E}_{t,p(x_{t_1}|k),q_t(x|x_{t_1},k)}||v_t(x,k)-u_t(x|x_{t_1},k)||^2$

 $min_{\theta} \mathbb{E}_{t,x_0,\epsilon}[||\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon,t,k)||^2]$



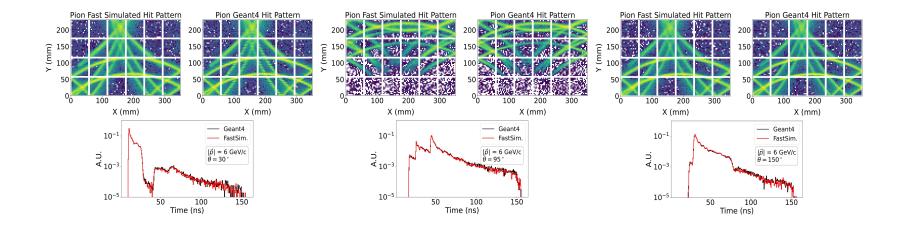
 $p_{\theta}(x_{0:T}|k) := p(x_T) \prod_{t=1}^{T} p_{\theta}(x_{t-1}|x_t, k)$

 $dx = [f(x,t) - g(t)^{2}\nabla_{x} \log p_{t}(x)]dt + g(t)d\bar{w}$

 $z_{t_0} = x$

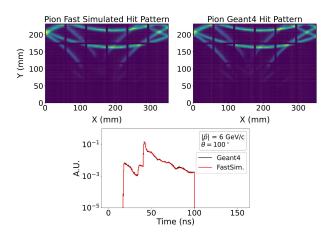
A collection of SOTA Generative Models

Discrete Normalizing Flows





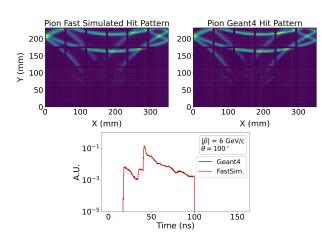
- Translate simulation quality metrics to a more meaningful representation
 - Separation power through KDE based PID method (FastDIRC)





- Translate simulation quality metrics to a more meaningful representation
 - Separation power through KDE based PID method (FastDIRC)

Create large reference populations of π / K (support PDF)



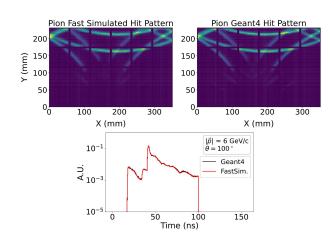


- Translate simulation quality metrics to a more meaningful representation
 - Separation power through KDE based PID method (FastDIRC)

Create large reference populations of π / K (support PDF)

For each photon in a track, calculate likelihood:

$$log \ p(\vec{x}_i|\vec{k})_{\mathcal{K}\pi} \propto (\vec{x}_i - \vec{\mu}_{\mathcal{K}\pi})^T \mathbf{\Sigma}^{-1} (\vec{x}_i - \vec{\mu}_{\mathcal{K}\pi}) \ , \ \mathbf{\Sigma} = \begin{bmatrix} p_{width}^2 & 0 & 0 \\ 0 & p_{height}^2 & 0 \\ 0 & 0 & \sigma_t^2 \end{bmatrix}$$





- Translate simulation quality metrics to a more meaningful representation
 - Separation power through KDE based PID method (FastDIRC)

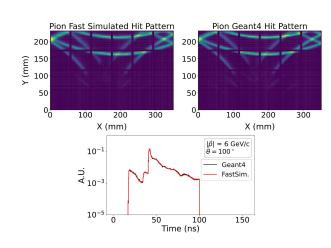
Create large reference populations of π / K (support PDF)

For each photon in a track, calculate likelihood:

$$log \ p(\vec{x}_i|\vec{k})_{\mathcal{K}\pi} \propto (\vec{x}_i - \vec{\mu}_{\mathcal{K}\pi})^T \mathbf{\Sigma}^{-1} (\vec{x}_i - \vec{\mu}_{\mathcal{K}\pi}) \ , \ \mathbf{\Sigma} = \begin{bmatrix} p_{width}^2 & 0 & 0 \\ 0 & p_{height}^2 & 0 \\ 0 & 0 & \sigma_t^2 \end{bmatrix}$$

Perform DLL:

$$\Delta \log \mathcal{L}_{\mathcal{K}\pi} = \sum_{i}^{N} \log p(\vec{x}_{i}|\vec{k}, \mathcal{K}) - \sum_{i}^{N} \log p(\vec{x}_{i}|\vec{k}, \pi)$$

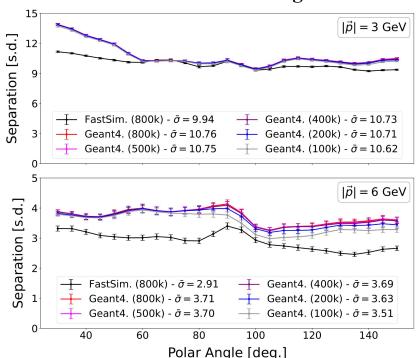


[3] Hardin, John, and Mike Williams. "FastDIRC: a fast Monte Carlo and reconstruction algorithm for DIRC detectors." *Journal of Instrumentation* 11.10 (2016): P10007.



- Evaluate at fixed kinematics
 - o 3,6 GeV/c 5 degree bins over acceptance range
- For each bin, fit a Gaussian distribution to both
 PID's in the DLL space
- Calculate the separation between distributions

Discrete Normalizing Flows

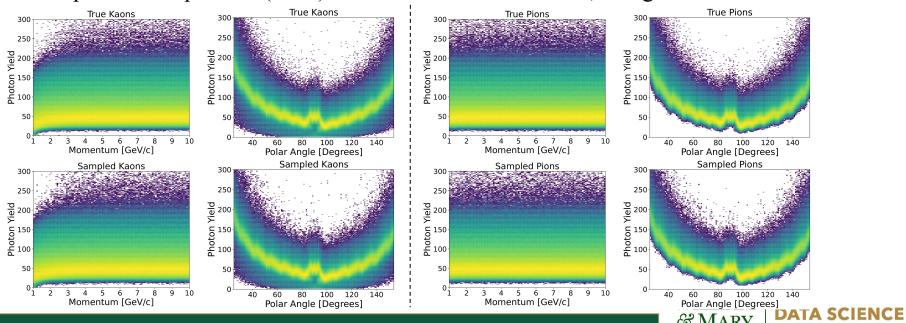




CHARTERED 1693

Holistic Simulation - Photon Yield Generation

- Low dimensional problem $f: \mathbb{R}^2 \to \mathbb{Z}_+$
- Must be fast approximately zero overhead, preferably CPU bound
- A simple Look-Up-Table (LUT) does the trick 100 MeV/c, 1 degree bins

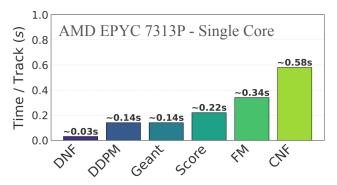


Key Takeaways

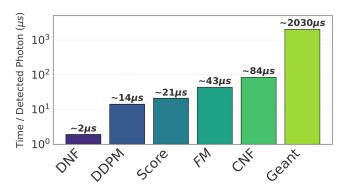
All code is open source and pre-trained models are provided.



- All Cherenkov Rings (underlying PDF's) generated by our models are "correct"
 - Ring and time structures follow correct kinematic dependencies for both PIDs
 - We incur a smoothing effect can cause different PIDs to appear more similar
- Beyond usage in Physics environments
 - We have created an open source suite of SOTA algorithms for the hpDIRC (easily adapted to other detectors)
 - Our fast simulation is self contained, fast and capable of being run on CPU or GPU



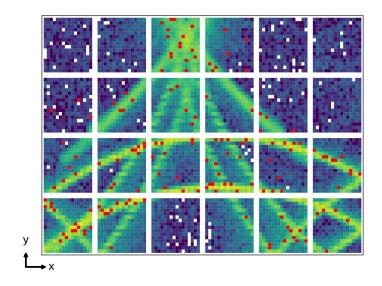
Track Generation (CPU)



PDF Generation (GPU)



Towards Foundation Models for Experimental Readout Systems Combining Discrete and Continuous Data





Foundation Models in Physics

- Foundation Models (FM) are becoming increasingly popular in the Physics community
 - Relatively large pre-trained models
 - Capable of supporting multiple downstream tasks e.g., **fast simulation**, reconstruction, etc.
- Different approaches have emerged (focusing on two recent ones)
 - o Diffusion Transformer (DiT) style (see [3])
 - o GPT style (see [4])
- In both cases, they work with relatively high level features
 - Facilitating generation and classification of Jets through their constituents (4-vector like quantities)

[3] Mikuni, Vinicius, and Benjamin Nachman. "OmniLearn: A method to simultaneously facilitate all jet physics tasks." *arXiv preprint arXiv:2404.16091* (2024).

[4] Birk, Joschka, Anna Hallin, and Gregor Kasieczka. "OmniJet-α: the first cross-task foundation model for particle physics." *Machine Learning: Science and Technology* 5.3 (2024): 035031.



Foundation Models in Physics cont'd...

- More recently [5] has shown very nice, and promising results using GPT style models to generate point clouds in calorimeters
 - Treat cells and energy as **tokenized** representations
 - Generate the shower forward in time, predicting the **next token** given the previous (**context**)
 - This is akin to modern LLMs such as ChatGPT

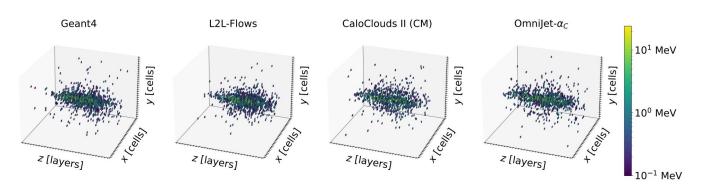
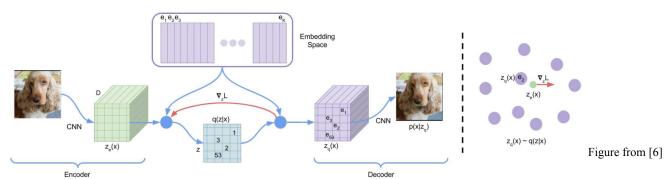


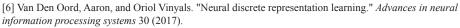
Figure from [5]



Potential Issues with Tokenization

- Tokens in the context of LLM are discrete integers representing a "word"
 - The vocabulary of the LLM is then the discrete set of tokens it is able to learn and generate
- What if we want to use **next token** prediction in continuous domains?
 - For example to generate images, or detector response (location and some value)
 - We can first learn the discretized codebook through some external model e.g., a Vector Quantized Variational Autoencoder (VQ-VAE)







Potential Issues with Tokenization cont'd

- A VQ-VAE like model solves the issue of allowing **next token** prediction style models to operate in continuous domains
- But it does come with drawbacks
 - There is inherent information loss in the encoding procedure
 - The reconstruction is limited by the granularity of the codebook
 - Potential inconsistencies or artifacts crucial in high precision applications

Combining Discrete and Continuous Data

- In attempt to circumvent these potential issues, we devise an alternative strategy
 - o In our Cherenkov data we have a pixel (discrete integer) and a continuous time associated with each hit
- We utilize two vocabularies two prediction heads
 - A discrete set of pixels
 - A discrete set of time bins a linear binning at ½ the timing resolution of the readout system
- As a result, our data structure for a given track is of the form

spatial
$$\rightarrow$$
 { $SOS_p, p_1, \dots, p_n, EOS_p$ }
time \rightarrow { $SOS_t, t_1, \dots, t_n, EOS_t$ }

• We have still discretized the continuous variate - but in a controlled manner



Conditional Generation

- Cherenkov hits in particular are highly dependent on external kinematic parameters
 - o Spatial location (ring structures in PDFs), and time distributions are highly variable
- While these variates are also continuous, we do not need to **tokenize** (discretize them)
- We instead embed them through linear projections and **prepended as context**

spatial
$$\rightarrow \{ |\vec{p}|, \theta, SOS_p, p_1, \dots, p_n, EOS_p \}$$

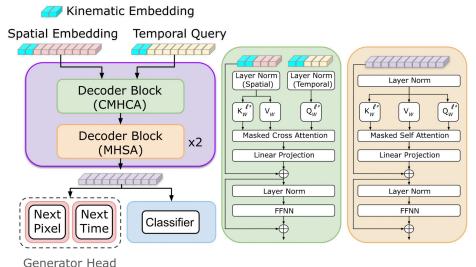
time $\rightarrow \{ |\vec{p}|, \theta, SOS_t, t_1, \dots, t_n, EOS_t \}$

• The prepending strategy allows the kinematics to *guide* sequence generation forward in time



Towards FM for Pixelated Readout Systems

- Our vocabularies operate adjacent to one another providing next **pixel** and **time** through independent prediction heads
- We combine information through Causal Cross Attention
 - Time drives the sequence at a given time, query the pixel space for possible locations

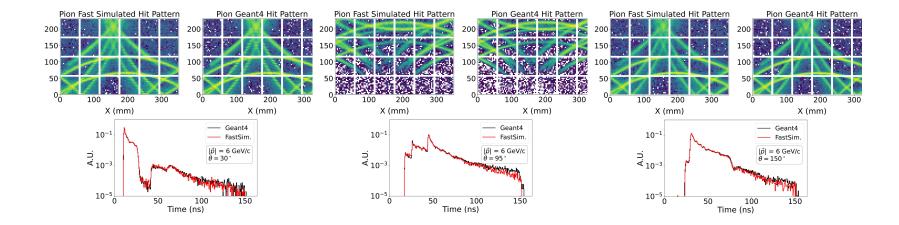


Causal implies we apply a mask to prevent seeing forward in time



Example Generations

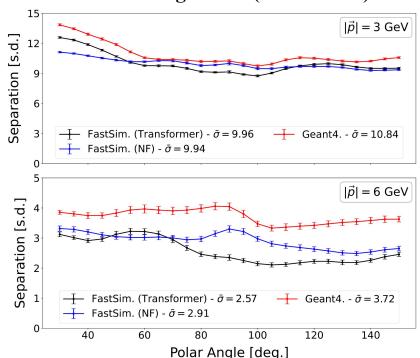
Autoregressive (Next Token)





- Evaluate at fixed kinematics
 - o 3,6 GeV/c 5 degree bins over acceptance range
- For each bin, fit a Gaussian distribution to both
 PID's in the DLL space
- Calculate the separation between distributions

Autoregressive (Next Token)





Class Conditional Generation

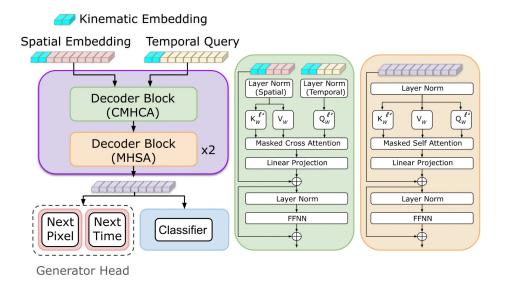
- Generations shown prior are from independent models (π/K)
 - How do we combine multiple classes under a single model?

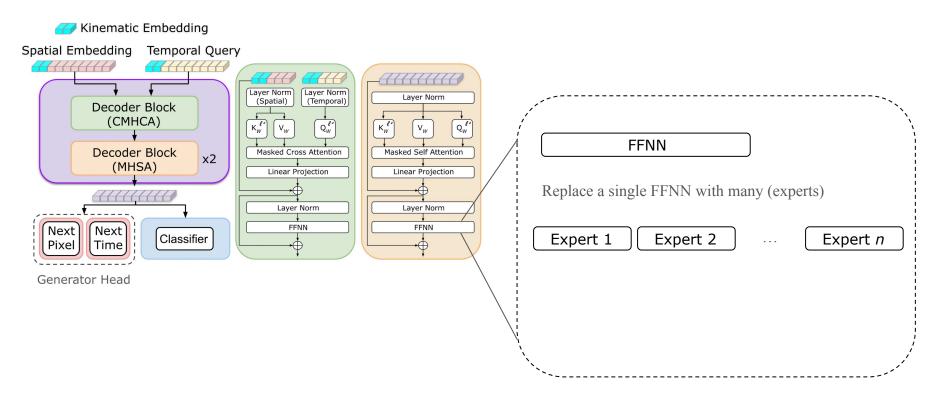
Class Conditional Generation

- Generations shown prior are from independent models (π/K)
 - How do we combine multiple classes under a single model?
- Standard Gen. AI
 - Build conditional probability distribution p(x|k,c)
 - Difficult given merging of π/K PDFs as momentum increases
 - Modes collapse together fidelity decreases

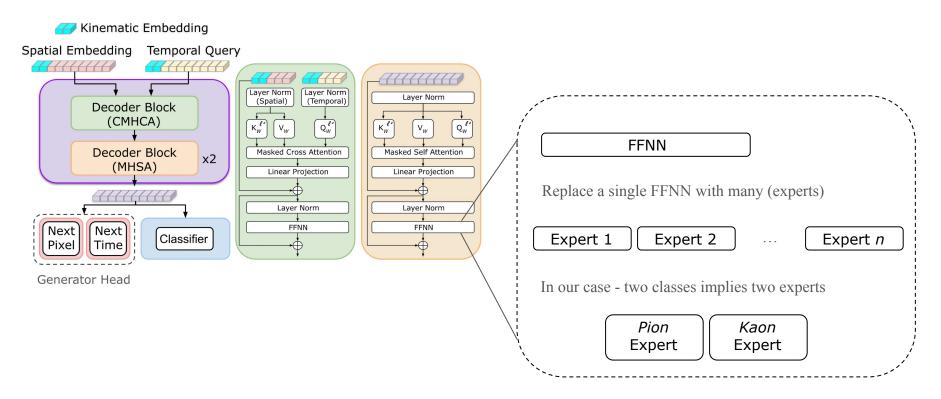
Class Conditional Generation

- Generations shown prior are from independent models (π/K)
 - How do we combine multiple classes under a single model?
- Standard Gen. AI
 - \circ Build conditional probability distribution p(x|k,c)
 - Difficult given merging of π/K PDFs as momentum increases
 - Modes collapse together fidelity decreases
- Autoregressive (next token)
 - We can prepend additional context (class label)
 - Same issues as before modes collapse together









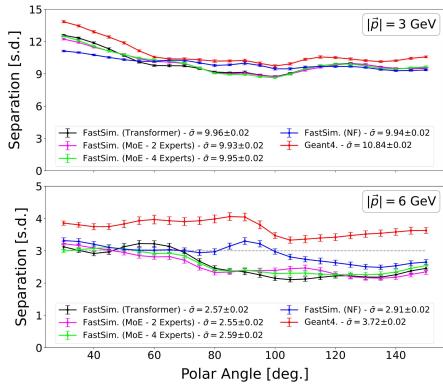


• The general idea

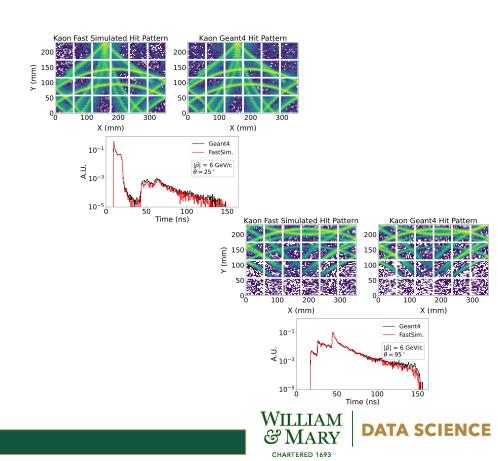
- o Turn of/on parts of the network given a specific type of particle we want to generate fixed routing
- Majority of parameters are shared only experts are aligned with a specific class
- Capture *general* relationship through attention blocks
- Apply *fine grained* corrections through class conditional experts

Extension in multiple ways

- Multiple particles $\pi/K/e$
- Multiple experts per class (see [4] for more details)



Generation Quality does not degrade w.r.t. Independent models



Key Takeaways

- (Proto) Foundation Model operating directly on low level detector signals
 - Generalizable to other detectors with discrete + continuous data streams



• Four core innovations

- All code is open source and pre-trained models are provided.
- Dual Vocabularies for spatial and temporal data fused through Cross Attention
- Scaleable, higher resolution tokenization with joint vocab inflation our resolutions correspond to 36M tokens in a joint vocabulary
- Continuous conditioning using prepending kinematic embeddings
- Class conditional generation through a Mixture of Experts
- Capable of multiple downstream tasks
 - Generation shown here.

