



# Tracking our Treasures

Punit Sharma

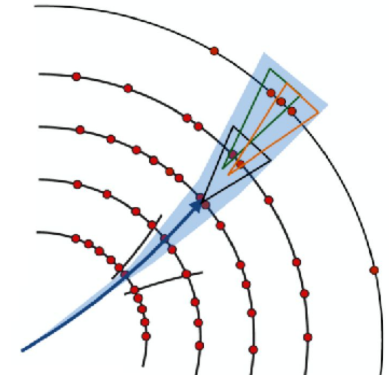
April 28, 2026



# Overview

Some overview of the tracking

- Dataset Generation
- Pipeline Overview



Low Level (Level 4?) data not available in the ATLAS Open Data!

# ACTS + ODD

## ACTS- A Common Tracking Software

ACTS provides algorithms for clusterization, track finding, track fitting, vertexing, ...

Benefits of being experiment-independent: knowledge transfer, common solutions for various domain specific problems

Already used in production by various experiments

## ODD- Open Data Detector ...

provides a template (HL-)LHC style particle detector for algorithm research and development

An evolution of the TrackML detector with more realistic material budget and full simulation support

Full silicon prototype inner detector based on DD4hep

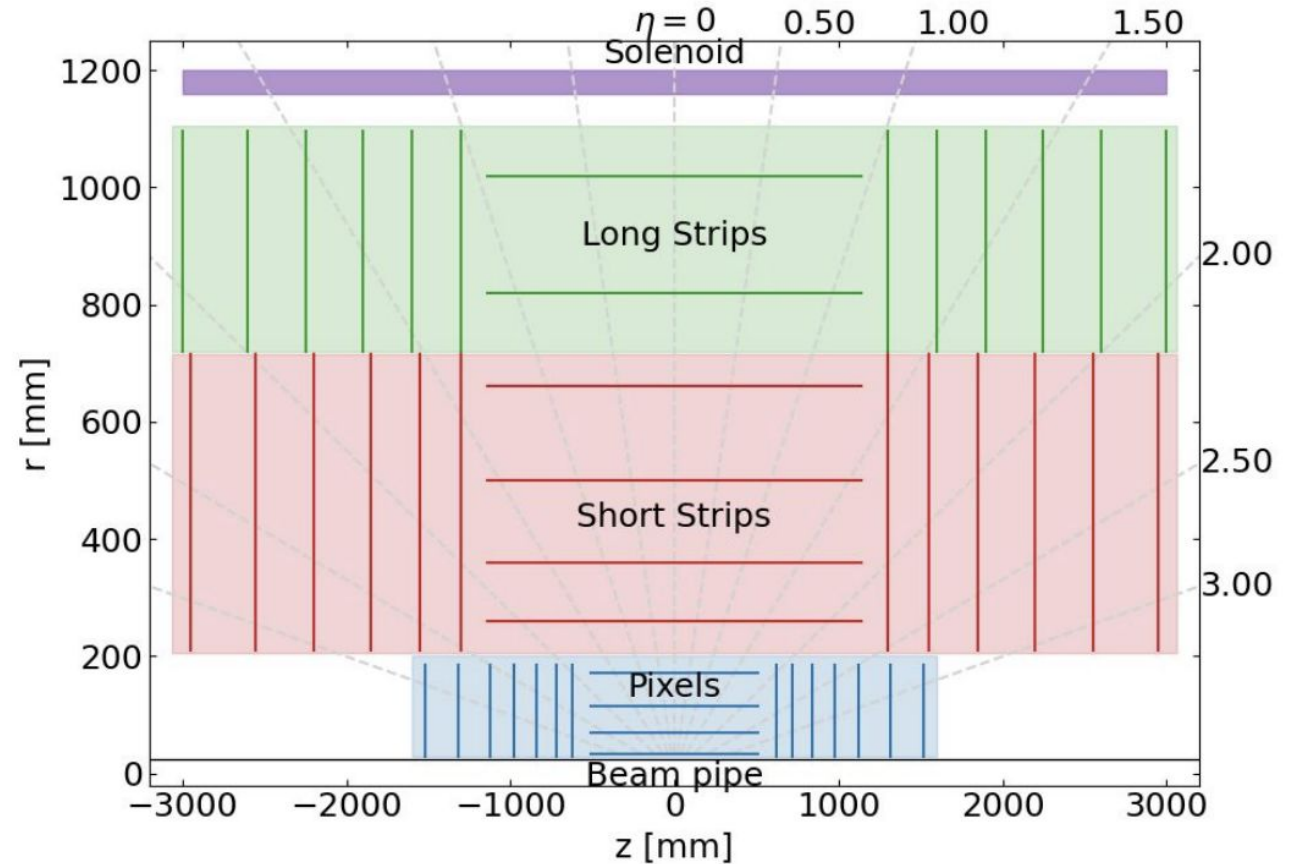
The detector description does not specify digitization

Used in Acts for validation and performance studies

For our studies only the silicon part was important

# ODD layout

- Pixel
  - 2D + time
  - Resolution: 15  $\mu\text{m}$  spatial, 25 mm time
  - 4 barrel layers
  - 7 endcap disks
- Short Strip
  - 2D
  - Resolution: 43  $\mu\text{m}$  / 1.2 mm spatial
  - 4 barrel layers
  - 6 endcap disks
- Long Strip
  - 1D, two sided, with stereo angle
  - Resolution: 72  $\mu\text{m}$  spatial
  - 2 barrel layers
  - 6 endcap disks



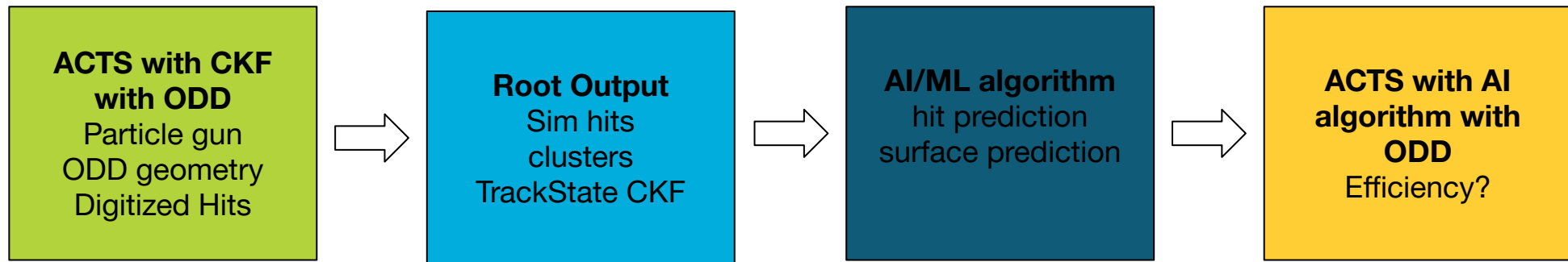
# Collider ML dataset

Pipeline

Simulation

	ID	Process Label	Process Type	Comments
$\langle \mu \rangle = 0,200$	1	ttbar	$pp \rightarrow t\bar{t}$	Top quark production; important for SM tests and complex final states.
$\langle \mu \rangle = 200$	2	zee	$pp \rightarrow Z \rightarrow e^+e^-$	Drell-Yan process; detector calibration and electroweak probe.
$\langle \mu \rangle = 200$	3	zmumu	$pp \rightarrow Z \rightarrow \mu^+\mu^-$	Drell-Yan process; detector calibration and tracking probe.
$\langle \mu \rangle = 0,200$	4	dihiggs	$gg \rightarrow HH$	Gluon fusion di-Higgs; benchmark for high-Luminosity studies.
	5	diphoton	$pp \rightarrow \gamma\gamma$	Diphoton production; loop-induced QCD/EW process.
	6	multijet	$pp \rightarrow \text{jets}$	QCD production; dominant background for hadronic signatures.
$\langle \mu \rangle = 0,200$	7	ggf	$pp \rightarrow H$	Gluon-gluon fusion; main Higgs production channel.
	8	susy	$pp \rightarrow \tilde{g}\tilde{g}$	GMSB gluino production; $\tilde{\chi}_1^0$ NLSP provides displaced vertex benchmark.
$\langle \mu \rangle = 0$	9	zprime	$pp \rightarrow Z'_{SSM}$	Heavy resonance; standard candle for bump hunting.
$\langle \mu \rangle = 0$	10	hiddenvally	$pp \rightarrow Z' \rightarrow \text{dark}$	Dark QCD sector; semi-visible jets for anomaly detection.

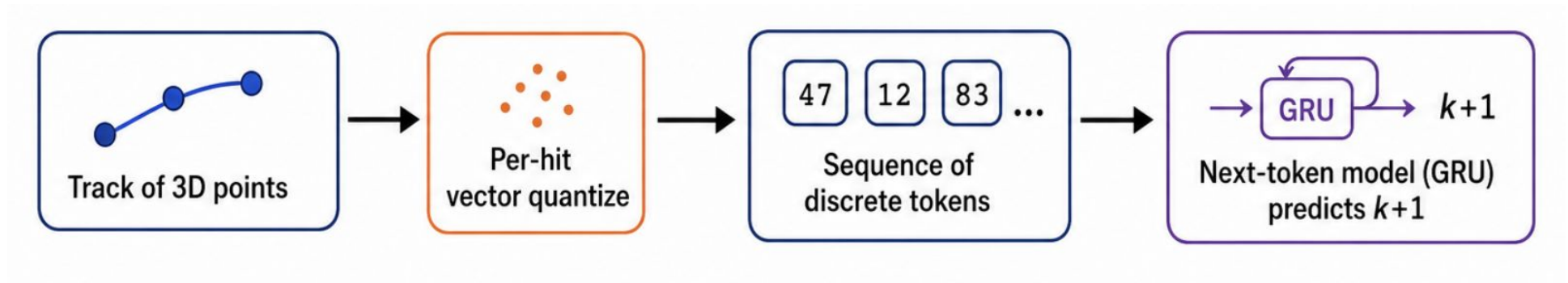
# Dataset generation



[ColliderML](#) [Arxiv](#)

# What is our problem statement?

The first order problem that I started looking into was if I can tokenize the tracks and use the token prediction as the way of track extrapolation.



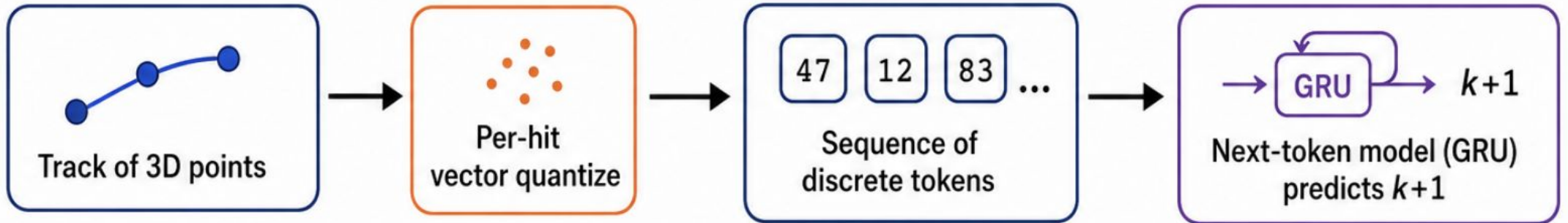
Some things to think about

1. Detector is big,  $\sim 1000 \times 1000 \times 3000$
2. How much resolution do we want with tokenization?
3. IS predicting next token enough?

A question to ask,

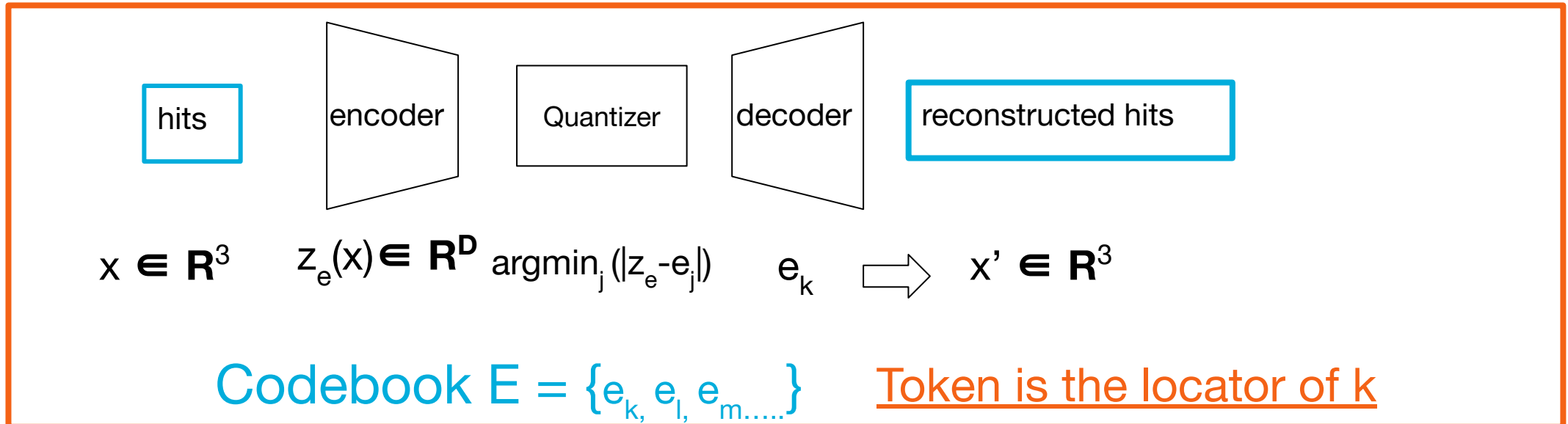
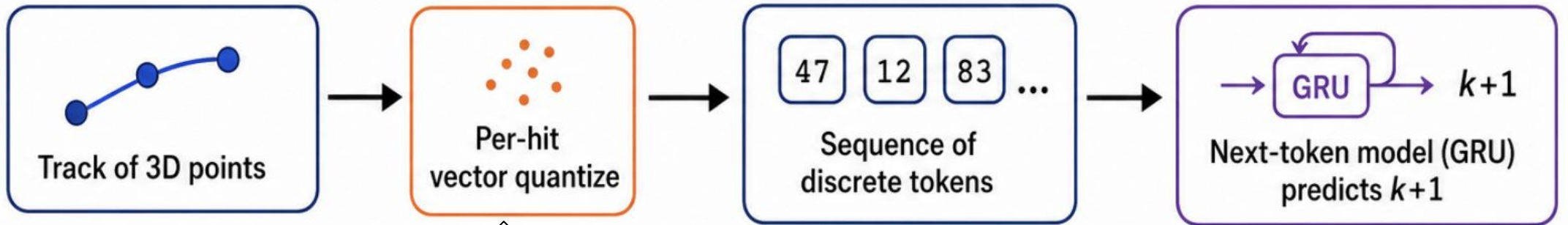
Let's say I want 10mm resolution on Z prediction, what would be the codebook size?

# Preprocessing

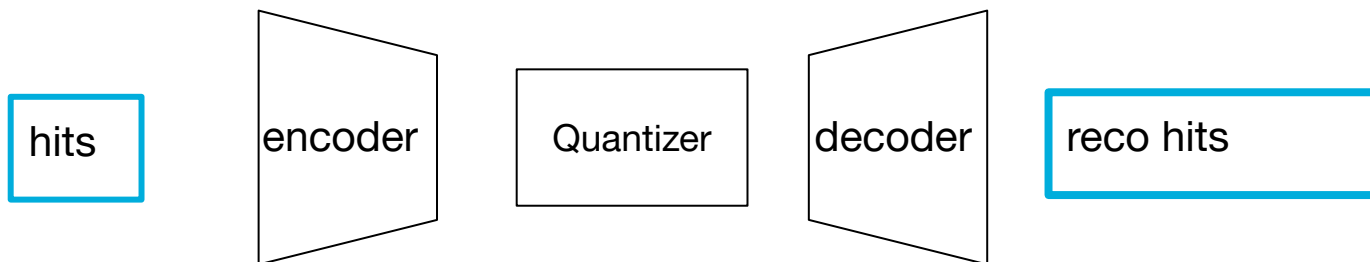


! Preprocessing of the dataset is needed.  
Plot some tracks to see if they make sense.  
Normalizations, splitting in test/train etc.

# VQ-VAE and Tokenization



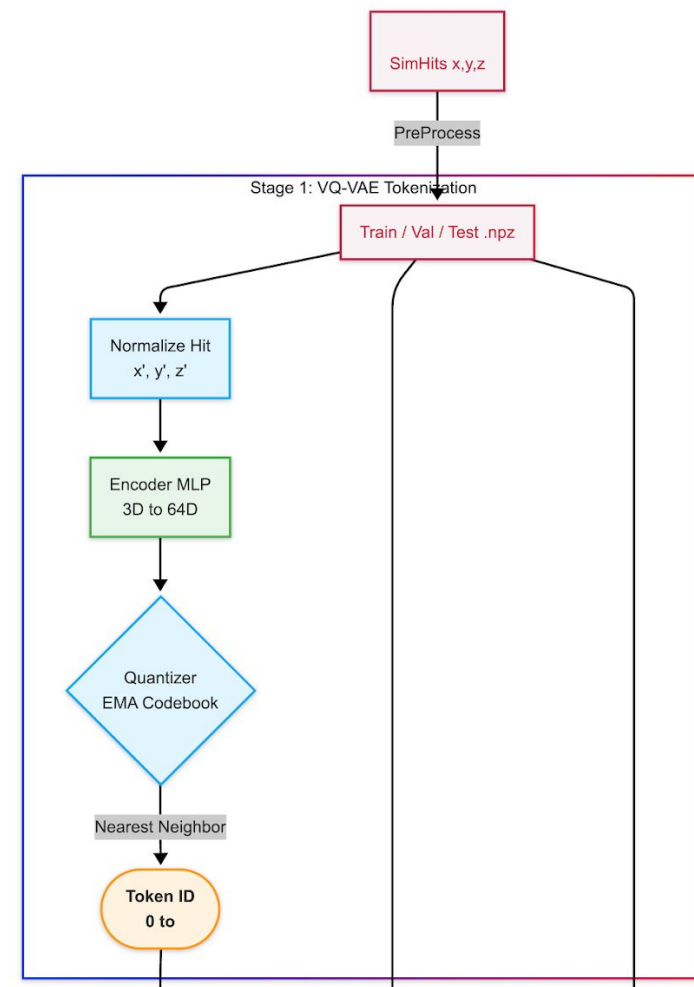
# VQ-VAE and Tokenization



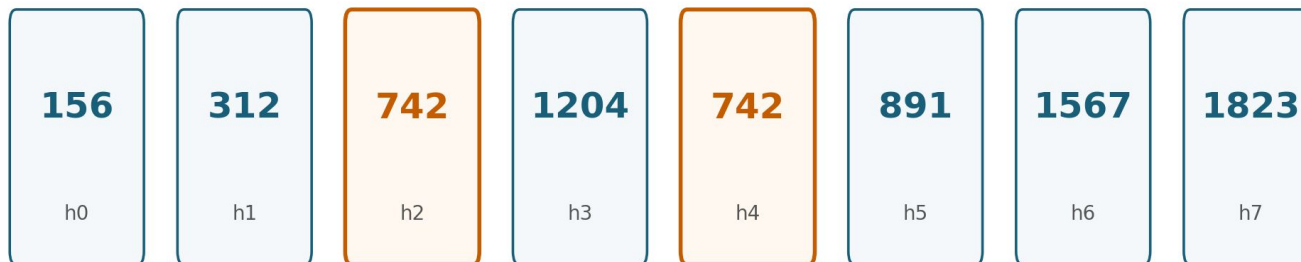
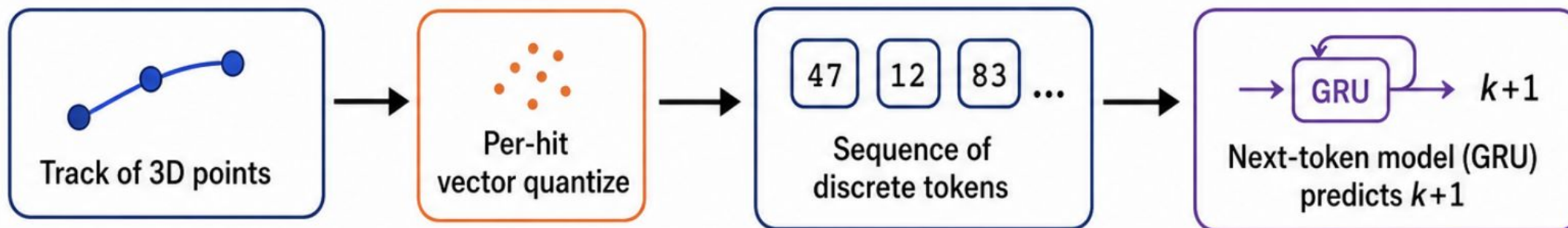
$$x \in \mathbb{R}^3 \quad z_e(x) \in \mathbb{R}^D \quad \operatorname{argmin}_j (|z_e - e_j|) \quad e_k \Rightarrow x' \in \mathbb{R}^3$$

Codebook is  $\{e_k, e_l, e_m, \dots\}$

In my case  
D= 64,  
codebook size of 2048



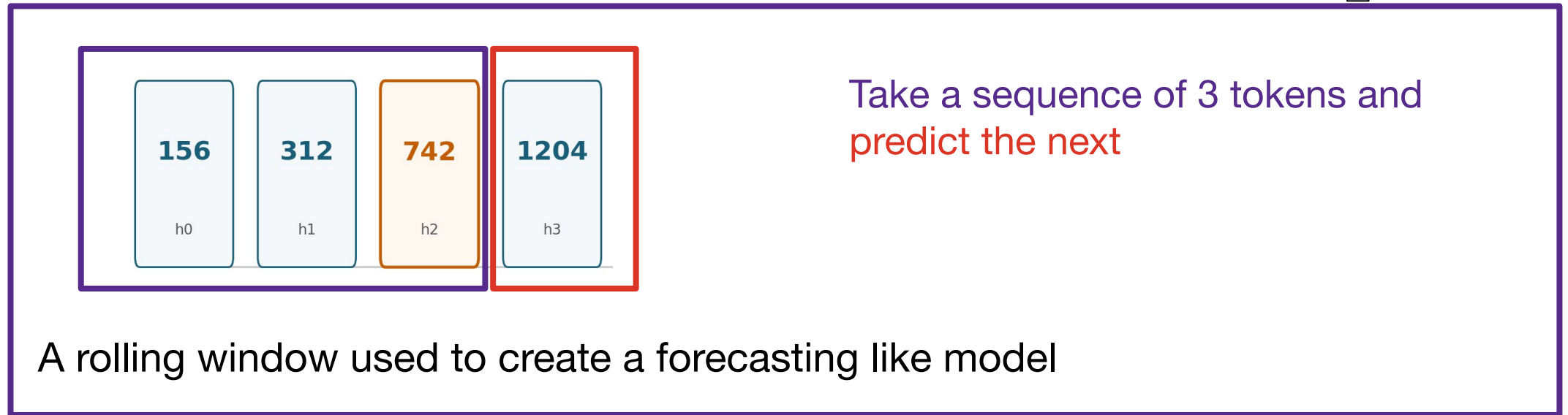
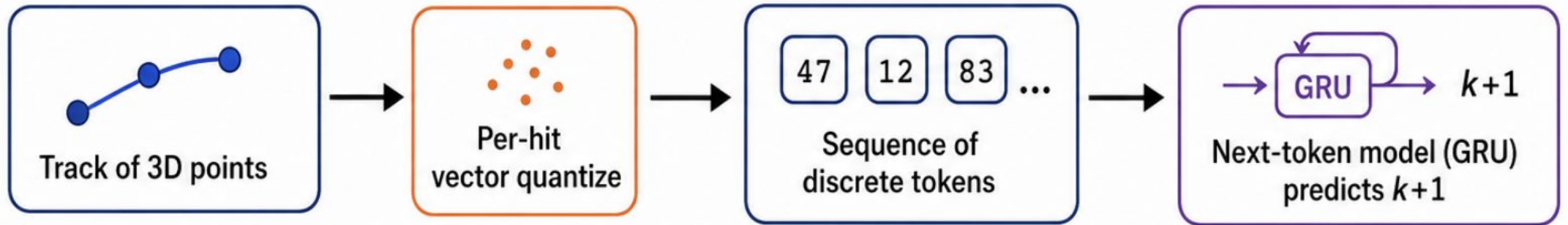
# Tracks to Tokens



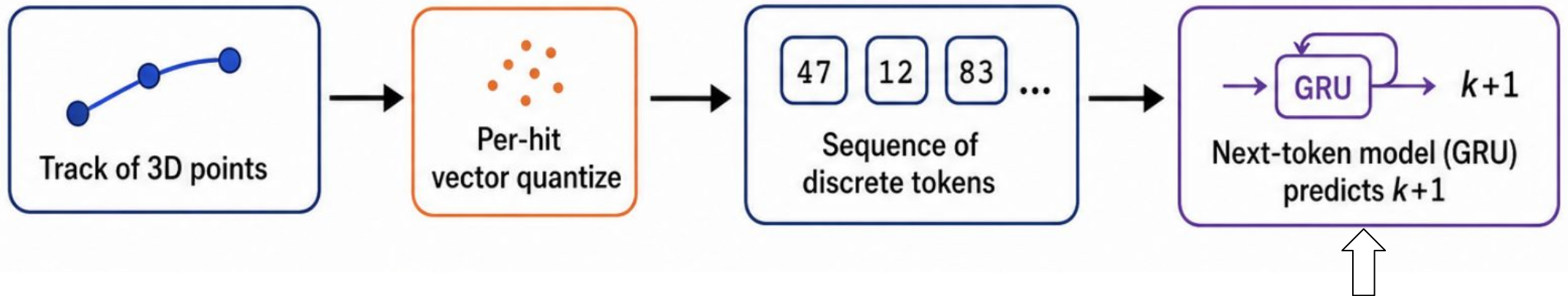
*hits 2 & 4 → same token 742 (different xyz map into one Voronoi cell)*

156 → 312 → 742 → 1204 → 742 → 891 → 1567 → 1823

# Prediction of Next Token (Ideal)



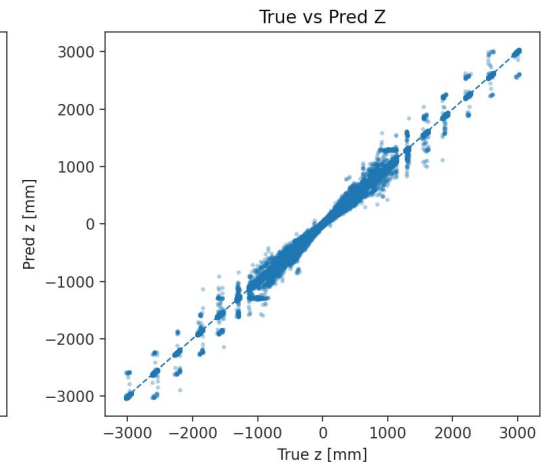
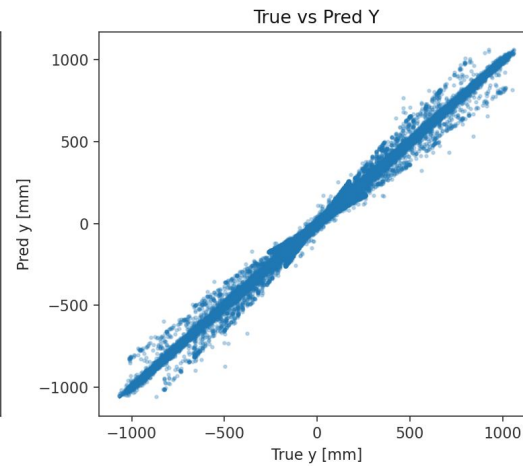
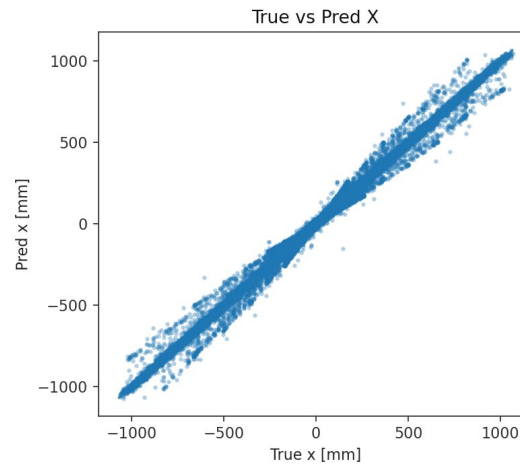
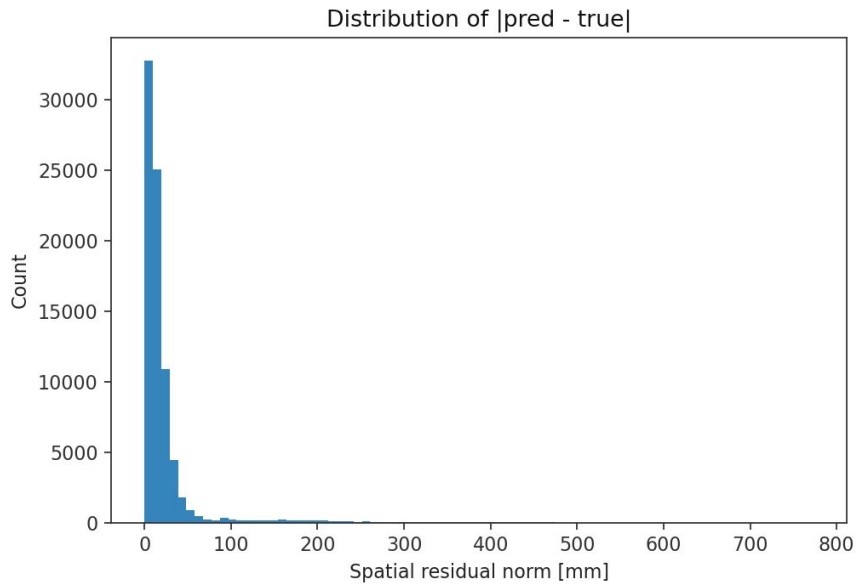
# Is Prediction of Next Token enough?



I found the answer to this question as no! Most likely the model + dataset I had isn't sufficient for this.

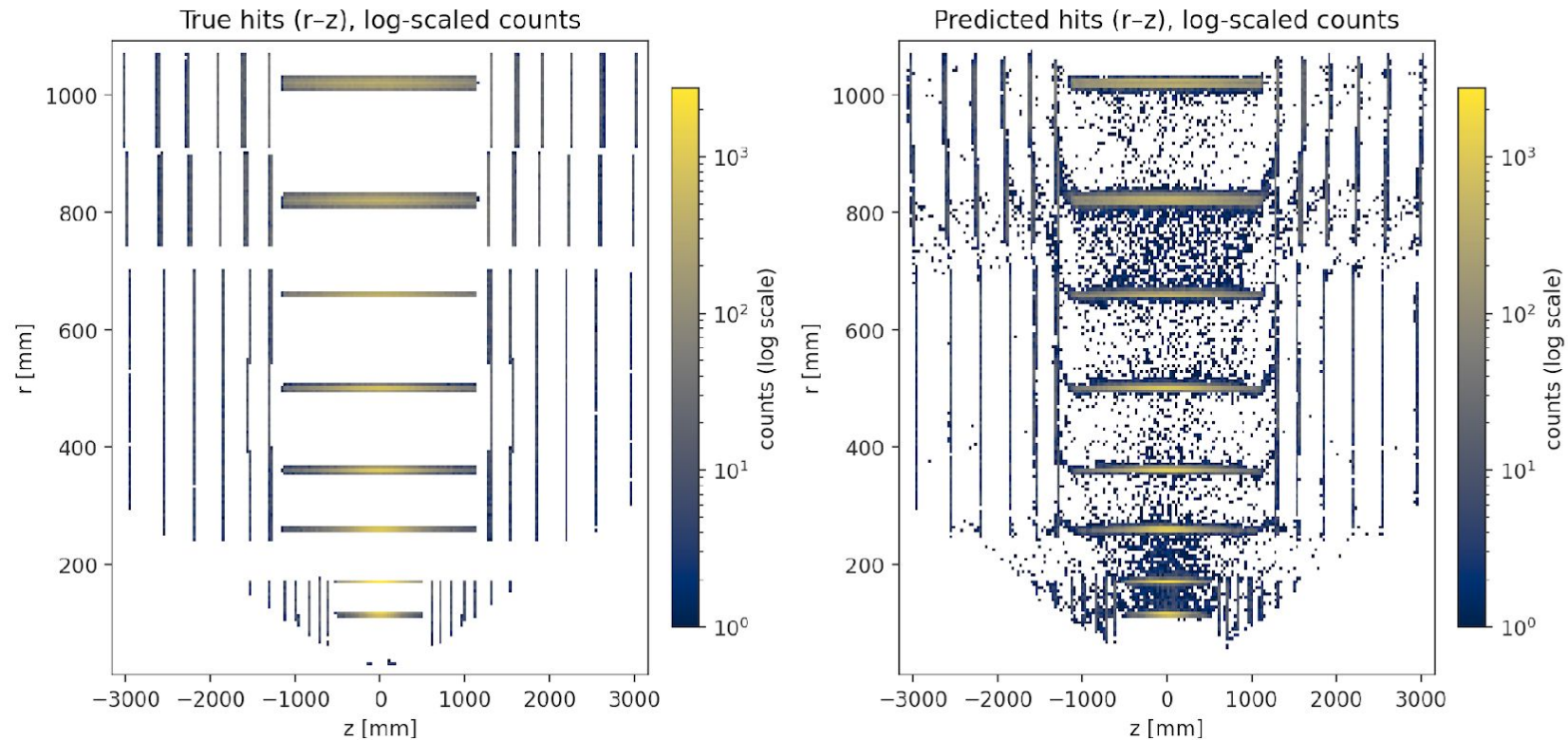
I had to use a concatenation of tokens with hits and layer information to predict the next token, hit and layer.

# Efficiency



- 100mm uncertainty on predicted hit.

# Early Results



- 100mm uncertainty on predicted hit.

# Summary and Next steps

- It's an interesting way of thinking the same problem.
- How would we add the ODD to the dataformat we discussed yesterday?
- Many use-cases, first one being using it with CKF?
- Are just having tokenization of the hits enough? Most likely not, we need track states on surface to be tokenized as well?
- Many other SOTA models can be used to fancy tracking tasks on the tokenized dataset.