# **Transformer-inspired**

# ML models for particle track reconstruction

Yue Zhao

High performance machine learning group SURF, national HPC centre, the Netherlands yue.zhao@surf.nl

# The collaboration





UNIVERSITEIT VAN AMSTERDAM

- Sascha Caron
- Nadezhda Dobreva
- José Martín-Guerrero
- Uraz Odyurt
- Slav Pshenov
- Roberto Ruiz de Austri Bazan
- Evgeniy Shalugin
- Zef Wolffs
- Yue Zhao

- High-Energy Physics, Radboud University, The Netherlands National Institute for Subatomic Physics (Nikhef)
- Radboud University, The Netherlands
- Intelligent Data Analysis Laboratory (IDAL) University of Valencia, Spain ValgrAI, Spain
- University of Twente, The Netherlands Nikhef, The Netherlands

Nikhef, The Netherlands

Instituto de Física Corpuscular, University of Valencia

Radboud University

- Institute of Physics, University of Amsterdam Nikhef, the Netherlands
- SURF (national HPC centre), the Netherlands

# The collaboration

Present at the conference!

- Sascha Caron
- Nadezhda Dobreva
- José Martín-Guerrero
- Uraz Odyurt
- Slav Pshenov
- Roberto Ruiz de Austri Bazan
- Evgeniy Shalugin
- Zef Wolffs
- Yue Zhao

High-Energy Physics, Radboud University, The Netherlands National Institute for Subatomic Physics (Nikhef)

Radboud University, The Netherlands

Intelligent Data Analysis Laboratory (IDAL) University of Valencia, Spain ValgrAI, Spain

University of Twente, The Netherlands Nikhef, The Netherlands

Nikhef, The Netherlands

Instituto de Física Corpuscular, University of Valencia

**Radboud University** 

Institute of Physics, University of Amsterdam Nikhef, the Netherlands

SURF (national HPC centre), the Netherlands



Universiteit van Amsterdam

Nik



Image source: CERN

3

#### The data: 3D point cloud of hits



Image source: Schoonheid, L. (2023). Fit-assisted DNN model for HEP Particle Tracking. [BSc thesis]. UvA, VU Amsterdam.

#### The task:

#### Reconstruct tracks from hits



Image source: arXiv:1904.06778 [hep-ex]

#### The task:

#### Reconstruct tracks from hits





Image source: arXiv:1904.06778 [hep-ex]

#### The task:

#### Reconstruct tracks from hits





-1000

750

500

250

0

-250

-500

-750

-1000

1000 750

500

250

0

Image source: arXiv:1904.06778 [hep-ex]

#### The data challenge

- 40 million collision events per second at LHC
- Hundreds of terabytes of data per second (online)
- Tens of petabytes saved per year (offline)
- High Luminosity LHC conditions with  $<\mu> \sim 200$  creates a substantial CPU time challenge for current reconstruction techniques due to combinatorics



Image source: Fast tracking for the HL-LHC ATLAS detector  $<\mu$ >

# 1<sup>st</sup> phase of our project transformed luminosity

Caron, S. *et al.* Trackformers: in search of transformer-based particle tracking for the high-luminosity LHC era. *Eur. Phys. J. C* **85**, 460 (2025). https://doi.org/10.1140/epjc/s10052-025-14156-3



#### Results of 1<sup>st</sup> phase

10-50 tracks (up to 500 hits) per event



Kaggle score: 94% Inference speed: ~2 miliseconds Kaggle score: double majority

- >50% hits from one particle
- <50% of its hits outside the track

200-500 tracks (up to 5k hits) per event



#### Kaggle score: 60-70% Inference speed: ~50 miliseconds

#### Results of 1<sup>st</sup> phase: inference speed of EncReg



#### Current efforts in our project (2<sup>nd</sup> phase)

- To address memory issue
  - FlexAttention and masking

attention matrix ~  $(\# hits)^2$ 





#### Current efforts in our project (2<sup>nd</sup> phase)

- To address memory issue
  - FlexAttention and masking
  - Domain decomposition or convolutions over  $\varphi?$
- To improve accuracy
  - Projections of hits to surfaces  $\rightarrow$  Clustering
  - Clustering in latent space of EncCla/EncReg model

#### Projections of hits to surfaces $\rightarrow$ Clustering

- Tracks of interest with  $p_t > 0.9 \ GeV$ are local in space
- Project pixel detector hits onto 3 surfaces (straight line to the origin)
- Choose surfaces to minimize hits spread





#### Projections of hits to surfaces $\rightarrow$ Clustering



#### Projections of hits to surfaces $\rightarrow$ Clustering

BARREL					
Recon. Ratio	Track DM eff.	# clusters	# tracks	Attn. matrix reduction	
0.987	0.14	1119	779	0.009138	
0.955	0.004	499	779	0.000285	
0.982 0.908	0.024 0.01	684 1176	779 779	0.002237 0.000998	
	<b>Recon.</b> <b>Ratio</b> 0.987 0.955 0.982 0.908	Recon. Track DM eff.   Ratio 0.14   0.987 0.14   0.955 0.004   0.982 0.024   0.908 0.01	Recon. Track DM #   Ratio eff. clusters   0.987 0.14 1119   0.955 0.004 499   0.982 0.024 684   0.908 0.01 1176	Recon. Track DM #   Ratio eff. clusters tracks   0.987 0.14 1119 779   0.955 0.004 499 779   0.982 0.024 684 779   0.988 0.01 1176 779	

	ENDCAP						
	Recon. Ratio	Track DM eff.	# clusters	# tracks	Attn. matrix reduction		
DBSCAN	0.985	0.832	756	267	0.00052		
Iterative window	0.993	0.468	871	267	0.002063		
Density Peaks Grid Window	0.993 0.989	0.831 0.577	1684 1330	267 267	0.000653 0.001099		

#### Current efforts in our project (2<sup>nd</sup> phase)

- To address memory issue
  - FlexAttention and masking
  - Domain decomposition or convolutions over  $\varphi?$
- To improve accuracy
  - Projections of hits to surfaces  $\rightarrow$  Clustering
  - Clustering in latent space of EncCla/EncReg model
- Simulations for more training data: 40k events collaborate?

#### Current efforts in our project (2<sup>nd</sup> phase)

- To address memory issue
  - FlexAttention and masking
  - Domain decomposition or convolutions over  $\varphi?$
- To improve accuracy
  - Projections of hits to surfaces  $\rightarrow$  Clustering
  - Clustering in latent space of EncCla/EncReg model
- Simulations for more training data: 40k events collaborate?
- Tokenization

#### From Hits to Tracks: Tokens in Physics



#### From Hits to Tracks: Tokens in Physics

- Goal is to convert a 4-vector to integer
- VQ-VAE replaces the continuous latent sampling of a standard VAE with a discrete codebook, trained end-to-end



Image modified from : arXiv:1711.00937

#### Key messages

- Transformer-based architectures are a promising direction for tracking
- FlexAttention and masking help address quadratic scaling of attention matrix
- Projecting the data to 'barrels' and 'endcaps' improves accuracy
- New TrackML-like simulations provide more training data
- If we see tracks as tokens, can we use VQ-VAE to tackle tracking?

# The collaboration

Present at the conference!

- Sascha Caron
- Nadezhda Dobreva
- José Martín-Guerrero
- Uraz Odyurt
- Slav Pshenov
- Roberto Ruiz de Austri Bazan
- Evgeniy Shalugin
- Zef Wolffs
- Yue Zhao

High-Energy Physics, Radboud University, The Netherlands National Institute for Subatomic Physics (Nikhef)

Radboud University, The Netherlands

Intelligent Data Analysis Laboratory (IDAL) University of Valencia, Spain ValgrAI, Spain

University of Twente, The Netherlands Nikhef, The Netherlands

Nikhef, The Netherlands

Instituto de Física Corpuscular, University of Valencia

**Radboud University** 

Institute of Physics, University of Amsterdam Nikhef, the Netherlands

SURF (national HPC centre), the Netherlands



UNIVERSITEIT VAN AMSTERDAM

Nik

#### Kaggle score: double majority Additional info Results of 1<sup>st</sup> phase: 10-50 tracks



(a) FitAccuracy versus the transverse momentum  $p_T$ .

(b) FitAccuracy versus pseudorapidity  $\eta$ .

>50% hits from one particle

<50% of its hits outside the track

FitAccuracy = Kaggle score without weights for hits

# Model performance $\rightarrow$ Not enough training data in trackML challenge ! (model can still improve)



(a) FitAccuracy versus the transverse momentum  $p_T$ .