MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC

Ambre Visive Nikhef & University of Amsterdam

Worked realized in collaboration with Roberto R de Austri, Polina Moskvitina, Sascha Caron and Clara Nellist



van Amsterdam

TABLE OF CONTENT



- Large-Language Model & Tokenization
- Dataset & Input Data
- Reflections on the Global Strategy
- Tokenization Strategies
- Results & Outlooks
- Conclusion

LARGE-LANGUAGE MODEL & TOKENIZATION

LLMs, Tokens, and Other Things We Should Probably Define

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025



LARGE A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC LANGUAGE MODEL

Transformer model that uses tokens as input

Inputs (= sequence of tokens)



LARGE AL VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC LANGUAGE MODEL

> Transformer model that uses tokens as input

Very good at learning the data distribution

Inputs (= sequence of tokens)



Modified from [



Thanks to this process, the model can process complicated data which have meaning to us.

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



 $[jet_1, pt_{jet1}, \varphi_{jet1}\eta_{jet1}, jet_2, pt_{jet2}, \varphi_{jet2}\eta_{jet2}, e_1, pt_{e-1}, \varphi_{e-1}\eta_{e-1}, ...]$

 $[jet_1, pt_{jet1}, \varphi_{jet1}\eta_{jet1}, bjet_1, pt_{bjet1}, \varphi_{bjet1}\eta_{bjet1}, e_1, pt_{e-1}, \varphi_{e-1}\eta_{e-1}, ...]$







ttbarZ,ttbarWW,ttbarW,ttbarHiggs

4 top

Token Index

[jet₁,pt_{jet1}, $\varphi_{jet1}\eta_{jet1}$, $jet_2, pt_{jet2}, \varphi_{jet2}\eta_{jet2}, e^-$ Tokenization [1,103,385...] $_{1}, \mathsf{pt}_{\mathsf{e}-1}, \varphi_{\mathsf{e}-1}\eta_{\mathsf{e}-1}, ...]$ [jet₁,pt_{jet1}, φ_{jet1} , η_{jet1} , [23,199,385,...] bjet₁, pt_{bjet1}, $\varphi_{bjet1}\eta_{bjet1}$, e⁻₁, pt_{e-1}, $\varphi_{e-1}\eta_{e-1}$, ...] Repartition of the tokens between signal and background





TRAINING PROCESS

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC





TRAINING PROCESS



A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



TRAINING PROCESS





A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



TRAINING PROCESS





Repeat for every batch in an epoch, until early stopping conditions are reached.





TRAINING PROCESS



TRAINING PROCESS





TRAINING PROCESS



TRAINING PROCESS



DATASET & INPUT DATA

Which cats should we use to train our model?

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC

7

INPUT OF THE MODEL

- a sequence of tokens

can represent a particle-type (and charge), its pt, its φ and its η , the MET, MET φ of the event...



Output

Probabilities

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC





(QUICK) TYPOLOGY

MET = (magnitude of) missing transverse energy

 MET_{φ} = azimuthal angle of MET



4-VECT = (energy, p_{\dagger} , φ , η)

 p_t = transverse mometum

id = particle-type + its charge

4-VECT = (id*, p_t , φ , η)

REFLECTIONS ON THE GLOBAL STRATEGY

How to build an anomaly detector that will take us to Stockholm

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025





4VECT-model

MET-model

 $MET\varphi$ -model

During Training







At Inference

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



OUR IDEA > REFLEXIONS & CHOICES

4VECT-model

- Should MET/MET φ be included in the events?
- How to tokenize the 4-vectors?

MET-model

- Should MET φ be • included in the events?
- How to tokenize the MET?

$MET\phi$ -model

How to tokenize • the MET φ ?

$\mathsf{MET}\varphi\ \mathsf{MODEL}$



Conf	usion Mat	rix for Ba	ackgroun	d Event F	Reconstru	ction	
881	0.60	0.18	0.03	0.04	0.14	- 0.6	
S 882 -	- 0.17	0.59	0.17	0.04	0.02	- 0.4	
rue Token 883	- 0.03	0.16	0.61	0.17	0.02	- 0.3	
TI 884	- 0.03	0.05	0.18	0.59	0.15	- 0.2	
885	- 0.17	0.04	0.03	0.20	0.56	- 0.1	sus
	881	882 Pre	883 dicted Tok	884 ens	885		Toke

-								
	Co	onfusion I	Matrix for	⁻ Signal E	vent Rec	onstructi	on	
	881	0.58	0.19	0.03	0.04	0.15		- 0.5
	5 882 -	0.18	0.57	0.18	0.04	0.03		- 0.4
	rue Tokens 883	0.03	0.20	0.57	0.18	0.02		- 0.3
	Tr 884	0.04	0.05	0.19	0.57	0.15		- 0.2
	885	0.20	0.05	0.04	0.19	0.53		- 0.1
	ľ	881	⁸⁸² Pre	883 dicted Tok	884 ens	885		

num epochs	dropout rate	ROC-AUC
6	0.1	0.5191
7	0.05	0.5207
7	0.05	0.5200
7	0.05	0.5202
7	0.05	0.5201
8	0.1	0.5191
9	0.1	0.5195
11	0.1	0.5187
11	0.1	0.5203
11	0.05	0.5198
12	0.05	0.5200
14	0.05	0.5197

ROC curves calculated at inference for different models

MET MODEL



num epochs	dropout rate	ROC-AUC
6	0.1	0.5474
6	0.1	0.5475
6	0.05	0.5478
6	0.05	0.5468
14	0.1	0.5438
14	0.05	0.5456
28	0.05	0.5451
31	0.05	0.5445
35	0.05	0.5467

ROC curves calculated at inference for different models

Confu	usion Mat	trix for Ba	ackground	d Event P	leconstru	ction
876 '	0.77	0.05	0.05	0.09	0.03	- 0.7
877	0.63	0.08	0.07	0.17	0.04	- 0.6 - 0.5
ue Tokens 878	0.48	0.08	0.08	0.24	0.11	- 0.4
Tr 879	0.36	0.06	0.08	0.27	0.24	- 0.3 - 0.2
- 88	0.21	0.02	0.05	0.22	0.50	- 0.1
	876	⁸⁷⁷ Pre	878 dicted Toke	879 ens	880	

Сс	onfusion	Matrix for	⁻ Signal E	vent Rec	onstructio	on
876	0.75	0.05	0.07	0.10	0.03	- 0.7
S 877	0.65	0.05	0.09	0.16	0.05	- 0.5
rue Token 878	0.55	0.05	0.10	0.21	0.10	- 0.4
П 879	0.44	0.04	0.08	0.26	0.17	- 0.3 - 0.2
880	0.27	0.02	0.07	0.26	0.38	- 0.1
	876	⁸⁷⁷ Pre	878 dicted Tok	879 ens	880	

OURIDEA REFLECTIONS & CHOICES

Should abandon the MET and MET φ specific models

A Not performing well on their own → diminish the performance of the "global" model



TOKENIZATION STRATEGIES

Slicing collisions like a good comté: a token at a time

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025





OVERVIEW OF OUR TOKENIZATION STRATEGIES

Should we add MET, MET φ tokens?



• GNN with VQ layer

• 8 "simple" (by-hand) tokenization





OVERVIEW OF OUR TOKENIZATION STRATEGIES

Should we add MET, MET φ tokens?



• GNN with VQ layer



• 8 "simple" (by-hand) tokenization

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



OVERVIEW OF OUR TOKENIZATION STRATEGIES

Should we add MET, MET φ tokens?

• VQVAE

GNN with VQ layer
into how many bins?
by binning the data
6 8 "simple" (by-hand) tokenization







MET+MET φ included in the sequence as a token









used tokens in test dataset: 452 codebook_size = 512

used tokens in train dataset: 453





GNN

*Graph Attention Convolution (layer from Graph Attention Network Architecture [7])



GNN







- Codebook size: 512
- Num tokens used in test: 452
- Num tokens used in train: 453



*Graph Attention Convolution (layer from Graph Attention Network Architecture [7])





- 5 bins of 20% of p_T
- 5 bins of 20% of $|\eta|$
- 5 geometrical bins of φ
- 1 bin per id
- 5 bins of 20% of MET
- 5 geometrical bins of MET φ



MET and MET φ tokens are included at the end of the sequence

	je [.]	t					b	je	ŧ				e	э-	ł					(e								μ	ι+	-							μ	ļ-						γ	
				(+	12	F	2	2	H.	st D-							<u>_</u>	2	F		DT DT)	0 _T		N PT									÷	Σ	L.		JM DT		L L	 st p _T					
						Low	ii/07V		HIGH	Hiche) 									2	Medi		High	Ц:СРО	пуле											Low p	-	Mediu		High	Highe					
																																												1		
 at Lu		ull,	, a [1]	1111		0.	'¶\$\$	wl	1					1					Įų, l				11.		1		1		.1							Į,I		J		1.						
											4			1								ľ				1.1				ľ		11	' '								1	1.1				

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



- 5 bins of 20% of p_T
- 5 bins of 20% of $|\eta|$
- 5 geometrical bins of φ
- 🔹 1 bin per id
- 5 bins of 20% of MET
- 5 geometrical bins of $MET\varphi$



MET and MET φ tokens are included at the end of the sequence





A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC



- 5 bins of 20% of p_T
- 5 bins of 20% of $|\eta|$
- 5 geometrical bins of φ
- 1 bin per id
- 5 bins of 20% of MET
- 5 geometrical bins of MET φ



MET and MET φ tokens are included at the end of the sequence

RESULTS & OUTLOOKS

Finding the needle in the datastack

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025

4-TOP EVENTS





APPLIED TO RARE EVENTS





APPLIED TO **BSM EVENTS**



0.8

1.0



How good is the token reconstruction

OUTLOOKS



- Not a fine-tuned model yet
- Still investigating best tokenization method
- However, method sounds promising

Stay tuned, paper will be published soon



CONCLUSION

From anomaly detection to foundation models: tokenization will reshape how we search for the unknown.

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025

THANK YOU FOR LISTENING

Ambre Visive Nikhef & University of Amsterdam ambre.visive@cern.ch





MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025

Х

REFERENCES

[1] Vaswani, A., Shazeer, N. et al. (2023). Attention Is All You Need. arXiv preprint arXiv:1706.03762.

[2] ATLAS detector schematics. Atlas experiment website. https://atlas.cern/Discover/Detector

[3] Aarrestad, T., van Beekveld, M., Bona, M. et al. (2022). The Dark Machines Anomaly Score Challenge: Benchmark data and model independent event classification for the Large Hadron collider. SciPost Physics, 12(1). https://doi.org/10.21468/scipostphys.12.1.043.

[4] Builtjes, L., Caron, S., Moskvitina P. et al. (2025). Attention to the strengths of physical interactions: Transformer and graph-based event classification for particle physics experiments. arXiv preprint arXiv:2211.05143.

[5] Van den Oord, A., Vinyals, O., & Kavukcuoglu, K. (2018). Neural Discrete Representation Learning. arXiv preprint arXiv:1711.00937.

[6] Birk, J., Hallin, A., & Kasieczka, G. (2024). OmniJet-a: the first cross-task foundation model for particle physics. Machine Learning Science and Technology, 5(3), 035031. https://doi.org/10.1088/2632-2153/ad66ad

[7] Veličković, P., Cucurull, G., Casanova, et al. (2018). Graph attention networks. . arXiv preprint arXiv:1710.10903

[8] Caron, S., Garcia Navarro, J., Moreno Llacer M. et al. (2025). Universal anomaly detection at the LHC: transforming optimal classifiers and the DDD method. 85(4). https://doi.org/10.1140/epjc/s10052-025-14087-z.

BACK-UPS

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC, EPS-HEP 2025

MET MODEL

ROC curves calculated at inference for different model

num epochs	dropout rate	ROC-AUC
6	0.1	0.5474
6	0.1	0.5475
6	0.05	0.5478
6	0.05	0.5468
14	0.1	0.5438
14	0.05	0.5456
28	0.05	0.5451
31	0.05	0.5445
35	0.05	0.5467







$MET\phi MODEL$

ROC curves calculated at inference for different model

num epochs	dropout rate	ROC-AUC
6	0.1	0.5191
7	0.05	0.5207
7	0.05	0.5200
7	0.05	0.5202
7	0.05	0.5201
8	0.1	0.5191
9	0.1	0.5195
11	0.1	0.5187
11	0.1	0.5203
11	0.05	0.5198
12	0.05	0.5200
14	0.05	0.5197









VQVAE - inspired by [d

MET+MET φ included in the sequence as a token





A. VISIN FOR AN codebook_size = 512 .

3

S/B

*Grapl

18

used tokens in test dataset: 452 e = 512 { used tokens in train dataset: 453







Output (2 classes)

GNN

Tokens for our model

Training of the GNN



*Graph Attention Convolution (layer from Graph Attention Network Architecture [7])





descriptionin the transverse planslightly forward particleforward particleforward particlewery forward particlemost forward particlearound 20% in each bin

• 5 bins of 20% of MET

• 5 bins of 20% of p_T

A. VISIVE, MASKED-TOKEN PREDICTION

FOR ANOMALY DETECTION AT THE LHC

+ 5 geometrical bins of ${\rm MET}\varphi$

MET and μ MET φ tokens are included at the end of the sequence

 $token_{4VECT} =$ $(cat_{id} - 1) \times 125$ $+(cat_{pt}-1)\times 25$ $+(cat_n-1)\times 5+cat_{\phi}$

20.1

Focus on

 $MET/MET\varphi$ bins

5 bins of 20% of p_T

- 5 bins of 20% of $|\eta|$
- 5 geometrical bins of φ

- 1 bin per id
- 5 bins of 20% of MET

A. VISIVE, MASKED-TOKEN PREDICTION

FOR ANOMALY DETECTION AT THE LHC

• 5 geometrical bins of $MET\varphi$







- 5 bins of 20% of p_T
- 5 bins of 20% of $|\eta|$
- 5 geometrical bins of φ
- 1 bin per id

A. VISIVE, MASKED-TOKEN PREDICTION FOR ANOMALY DETECTION AT THE LHC

- 5 bins of 20% of MET
- 5 geometrical bins of $MET\varphi$

20.3



- 5 bins of 20% of p_T
- 5 bins of 20% of $|\eta|$
- 5 geometrical bins of φ
- 1 bin per id

A. VISIVE, MASKED-TOKEN PREDICTION

FOR ANOMALY DETECTION AT THE LHC

- 5 bins of 20% of MET
- 5 geometrical bins of MET φ

Focus on BSM events (SUSY gluino-gluino production)

20.4





