

Lorentz-GATr

Lorentz-Equivariant
Geometric Algebra Transformers
for High-Energy Physics

Jonas Spinner*, Victor Breso*,
Pim de Haan, Tilman Plehn,
Jesse Thaler, Johann Brehmer

ML Journal Club
RWTH Aachen

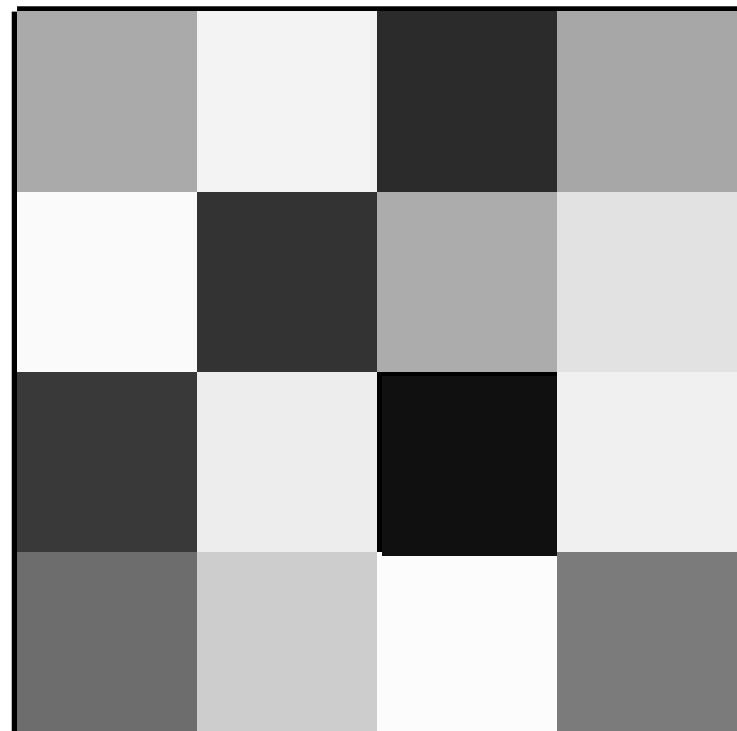


UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Lorentz symmetry is key in
high-energy physics...

$$\begin{aligned}\mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i \bar{\psi} D^\mu \psi + h.c \\ & + \bar{\chi}_i \gamma_{ij} \chi_j \phi + h.c \\ & + |\nabla_\mu \phi|^2 - V(\phi)\end{aligned}$$

... so let's build it into
our neural networks



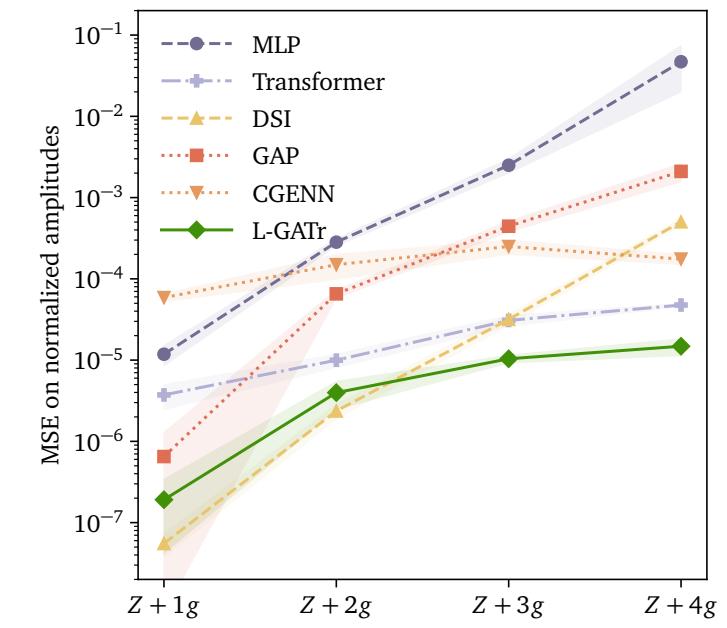
Transformer
architecture



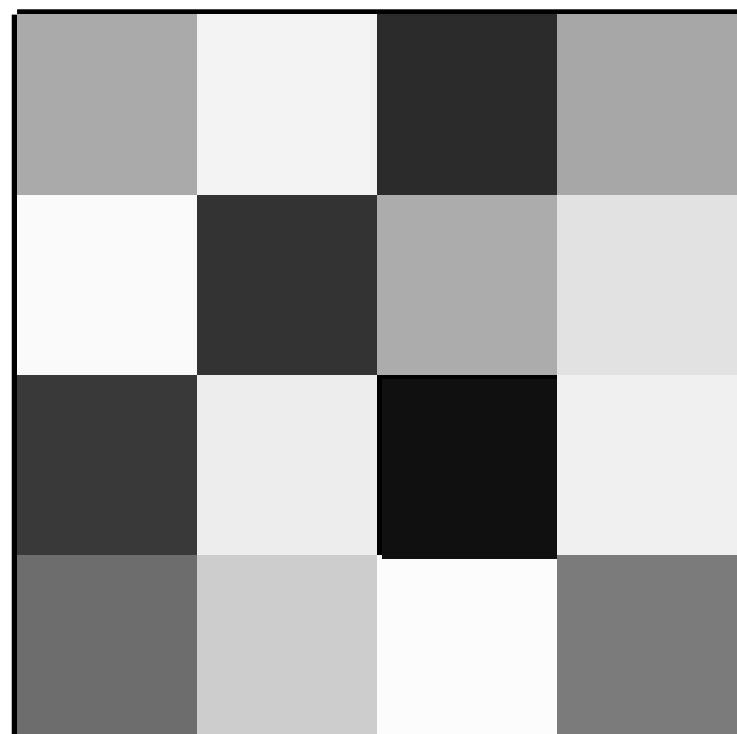
Geometric algebra
representations

**Lorentz-Equivariant
Geometric **Algebra**
Transformer**

GATr was originally
developed for E(3)
arXiv:2305.18415



Strong performance
on diverse problems

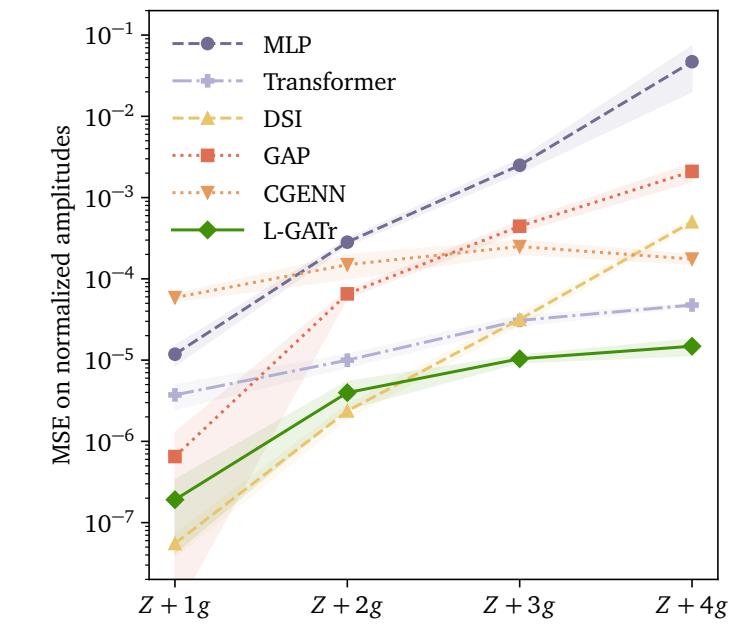


Transformer
architecture



Geometric algebra
representations

Lorentz-Equivariant
Geometric Algebra
Transformer

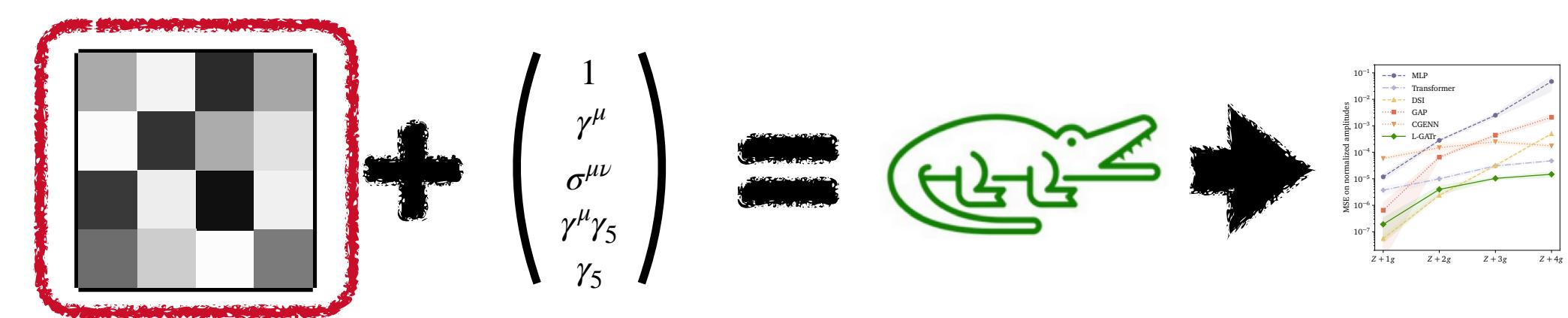


Strong performance
on diverse problems

GATr was originally
developed for E(3)
arXiv:2305.18415

Transformers

Point clouds



UNSTRUCTURED

Structure



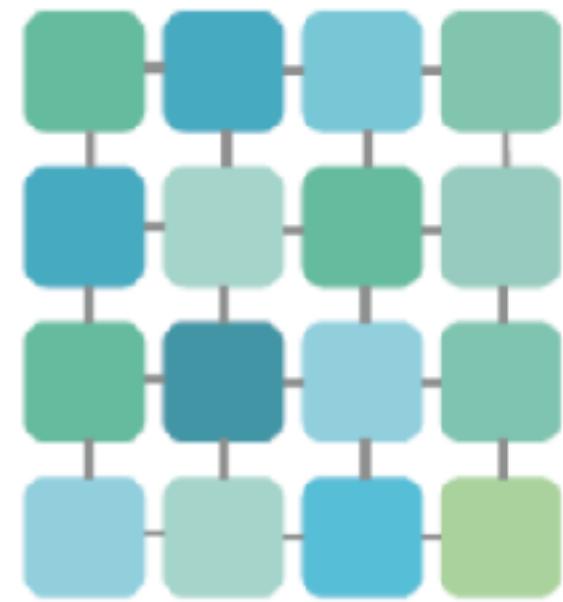
Network
Architecture

MLP

Inductive Bias

-

GRID



Convolutional
NN (CNN)

Locality

POINT CLOUD



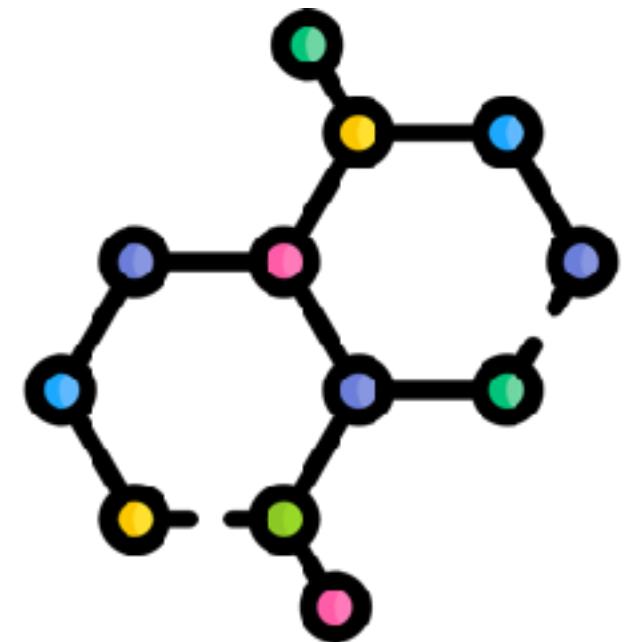
Graph NN (GNN)/
Transformer

Permutation symmetry
Variable-size inputs
Pairwise relations

Transformers

Point clouds

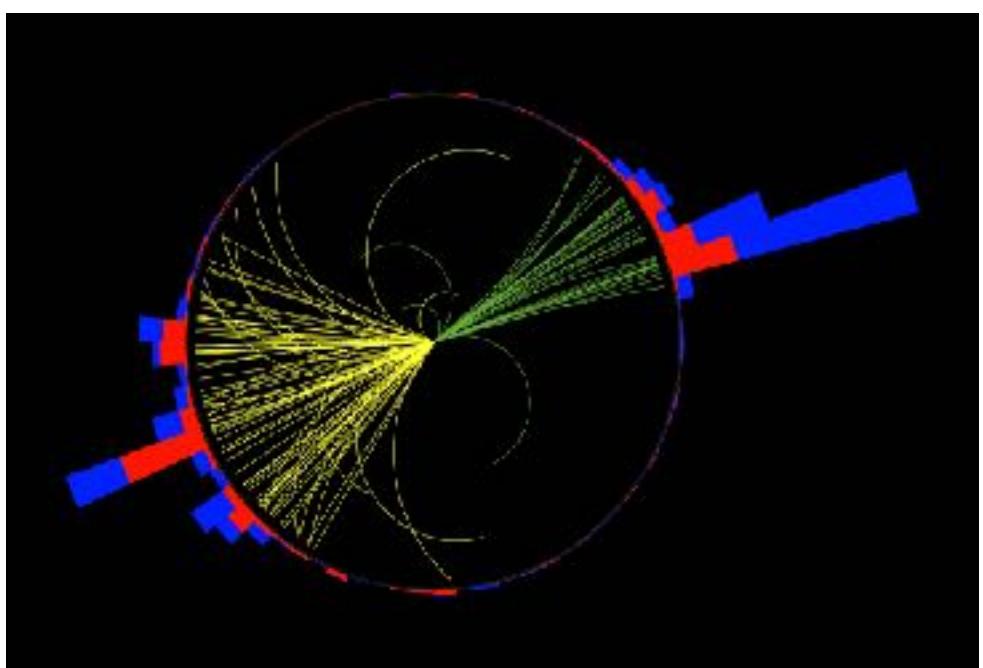
Molecules



Words

Attention is all you need .

Jets



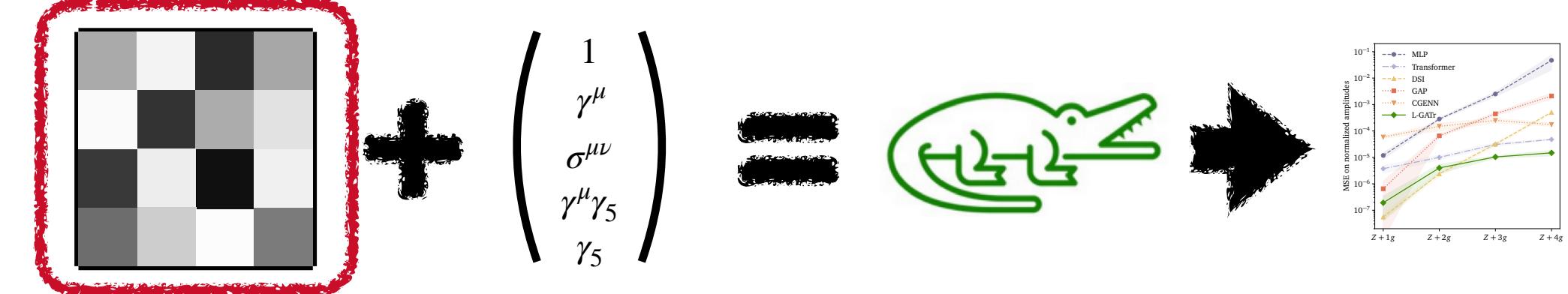
Events

$$q\bar{q} \rightarrow e^+e^-\gamma$$

Calorimeter cells



Patches in
an image



POINT CLOUD



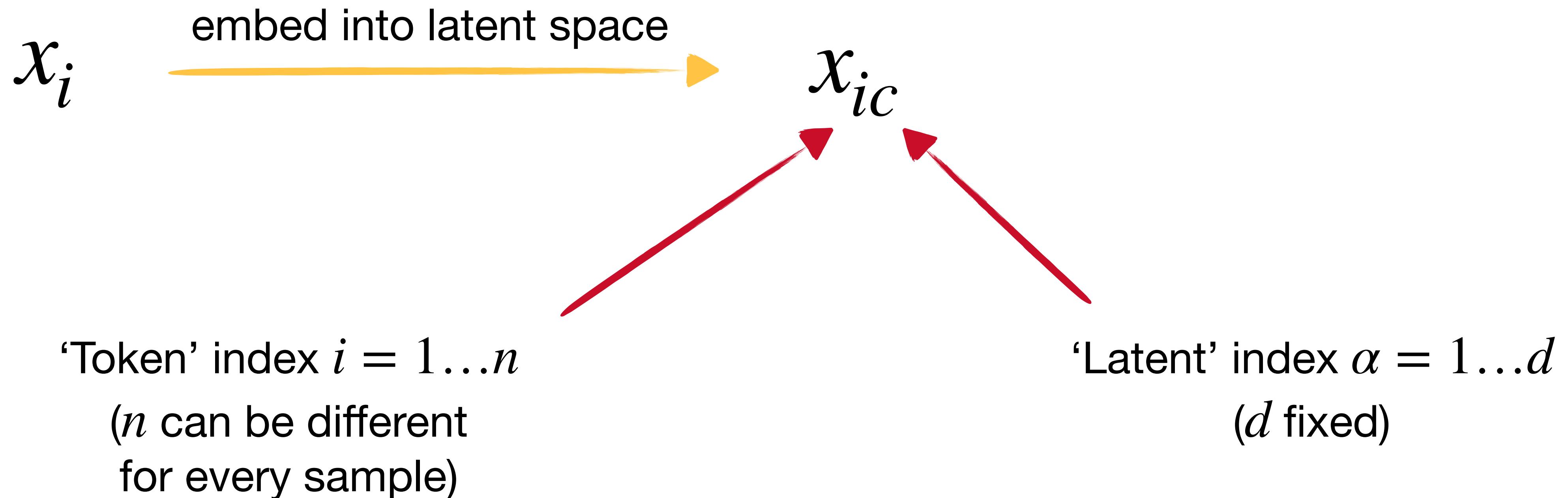
Graph NN (GNN)/
Transformer

Permutation symmetry
Variable-size inputs
Pairwise relations

Transformers

Point cloud representations

A diagram illustrating a mathematical operation. On the left, a 4x4 grid of gray and black squares is enclosed in a red-bordered box. To its right is a plus sign. Next is a vector in parentheses: $\begin{pmatrix} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu \gamma_5 \\ \gamma_5 \end{pmatrix}$. This is followed by an equals sign. Then there is a green hand icon pointing to a large black arrow pointing right. Finally, a small graph titled 'MSB on normalized amplitudes' is shown, plotting various models against Z values.



Data has shape
(batchsize, #tokens, #latent dims)

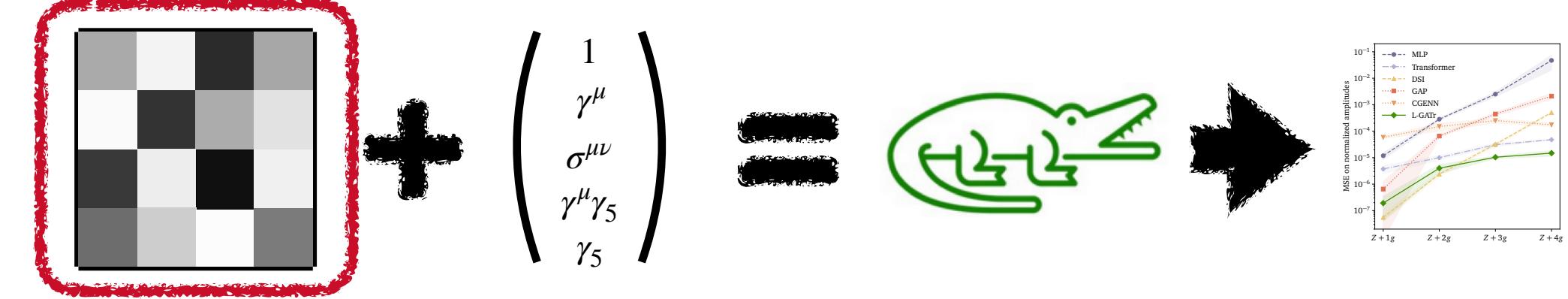
Transformers

Attention

Want a simple update operation on point clouds $x'_{ic} = f_{ic}(x)$

Exchange information
along ‘latent’ index c
→ Linear layer

$$x'_{ic} = W_{cc'}^V x_{ic'}$$



Transformers

Attention

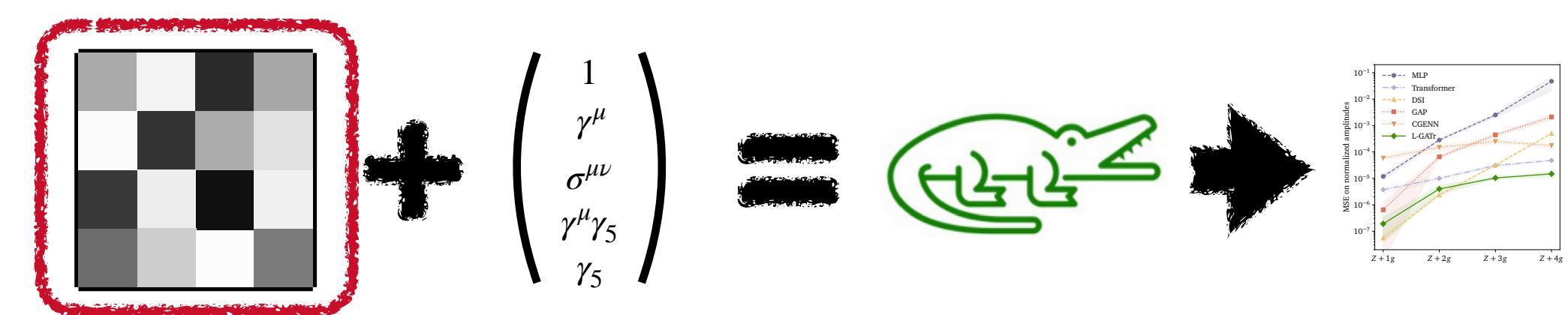
Want a simple update operation on point clouds $x'_{ic} = f_{ic}(x)$

Exchange information
along '**token**' index i
(without breaking
permutation symmetry)
→ Attention matrix

$$A_{ij}(x) = \text{Softmax}_j\left(\frac{x_{ic} W_{cc'}^Q W_{c'c''}^K x_{jc''}}{\sqrt{d}}\right)$$

$$x'_{ic} = A_{ij}(x) W_{cc'}^V x_{jc'}$$

More sophisticated update operations:
Message-Passing Graph Networks

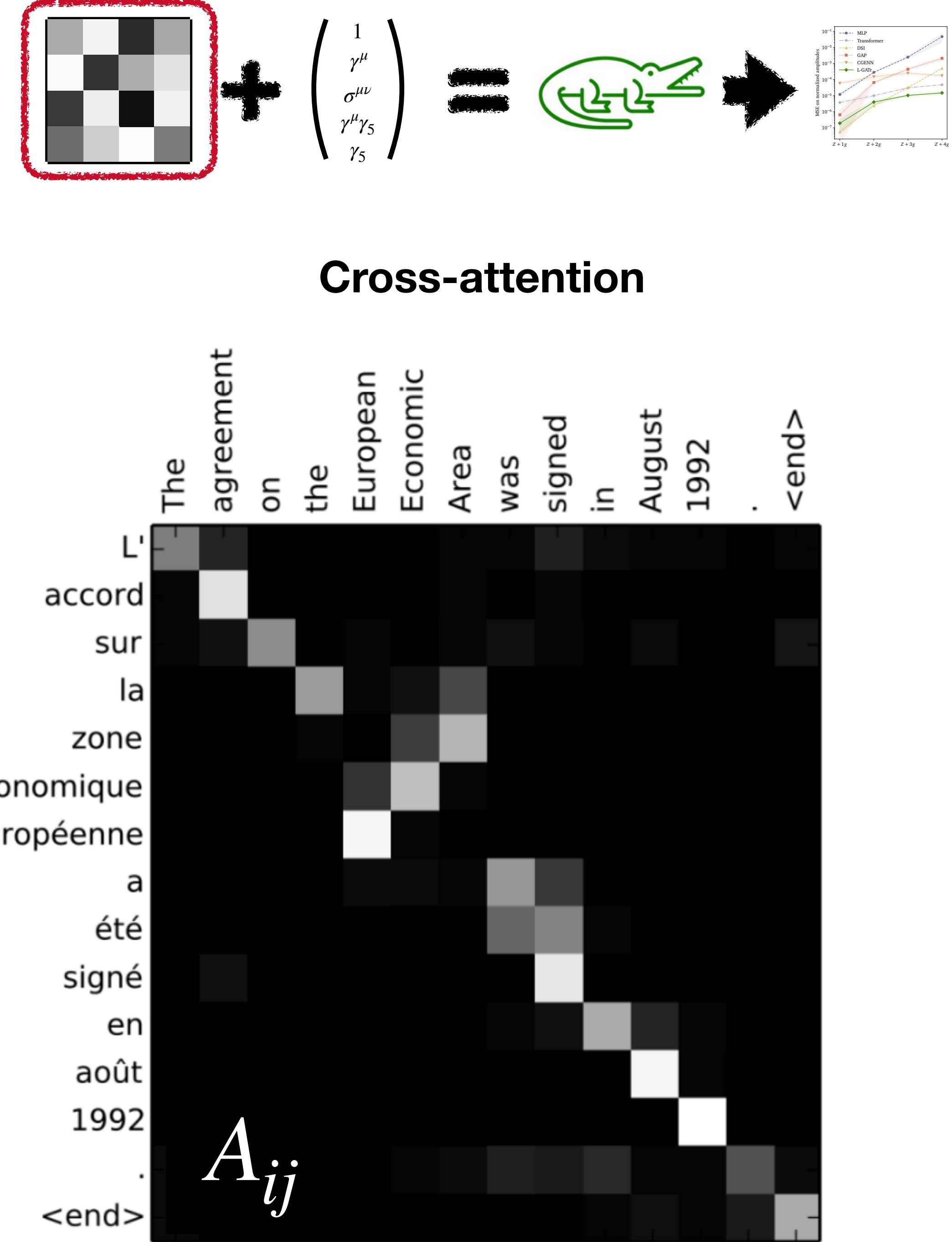


Transformers

Attention matrix A_{ij} - visualised

Self-attention

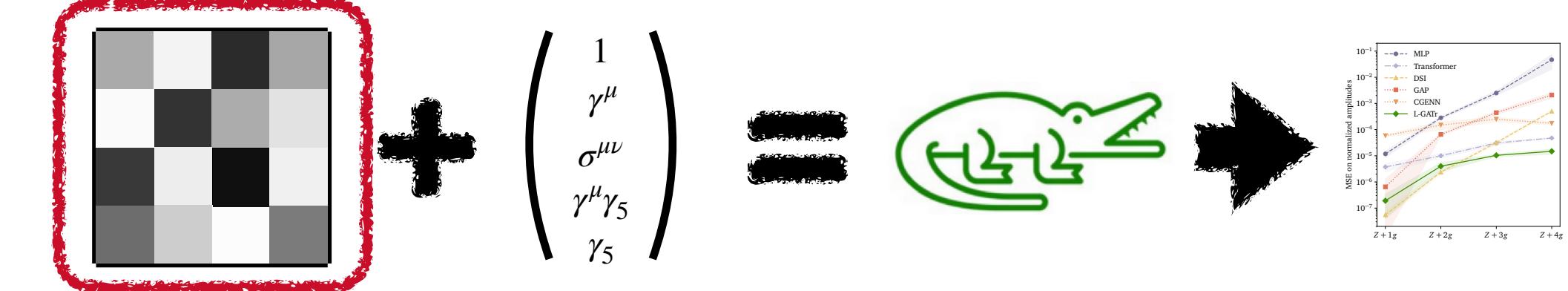
	Hello	I	love	you
Hello	0.8	0.1	0.05	0.05
I	0.1	0.6	0.2	0.1
love	0.05	0.2	0.65	0.1
you	0.2	0.1	0.1	0.6



Transformers

Attention is all you need

- Transformer = Network architecture based solely on attention mechanisms
- Proposed 2017 in ‘Attention is all you need’ (already 130k citations)
- Established as the most scalable point cloud architecture



Attention Is All You Need

1706.03762

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

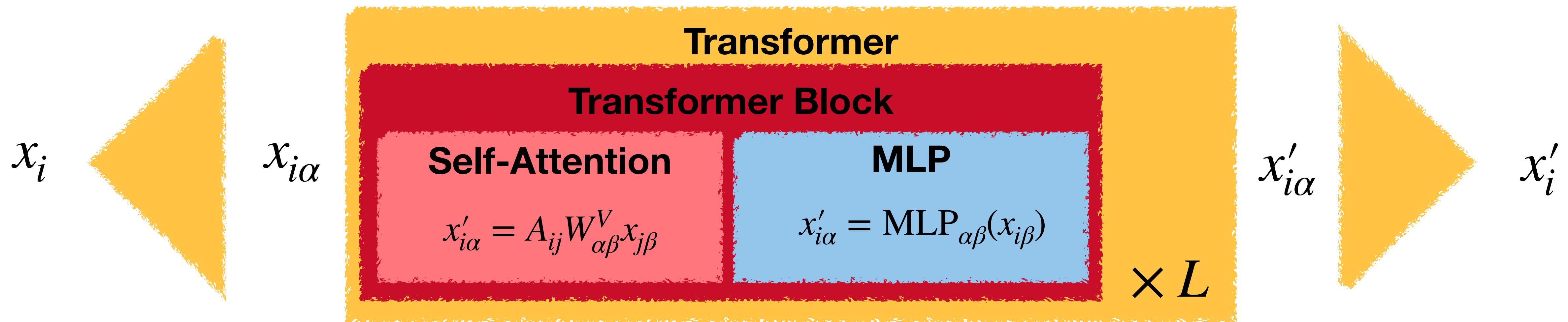
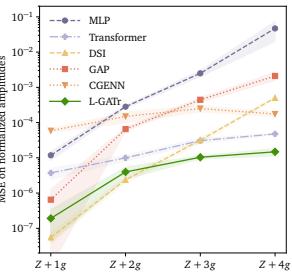
Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to

Transformers

Architecture overview

$$\begin{pmatrix} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu\gamma_5 \\ \gamma_5 \end{pmatrix} = \text{Matrix} + \begin{pmatrix} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu\gamma_5 \\ \gamma_5 \end{pmatrix}$$



Embed into
latent space

‘Exchange
information
between tokens’

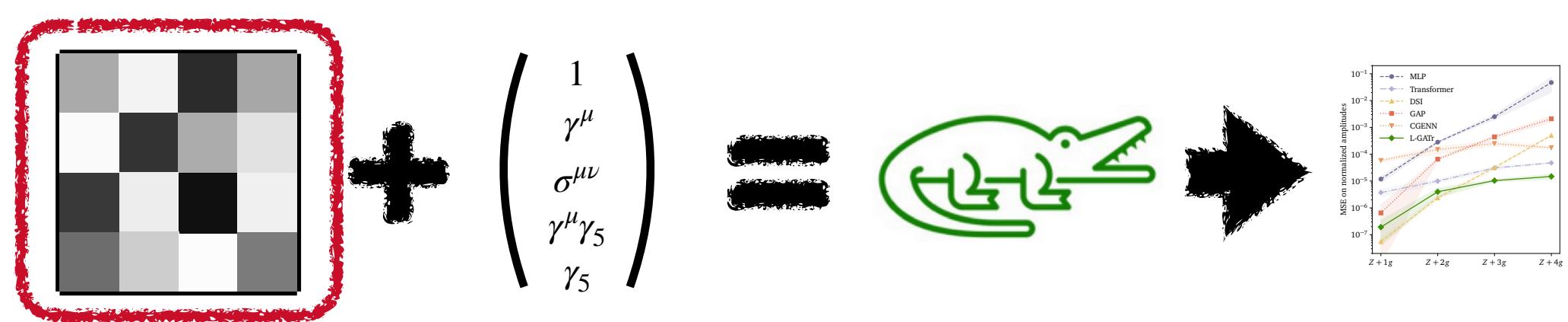
Nonlinearity +
‘Clean up
latent space’

Extract from
latent space

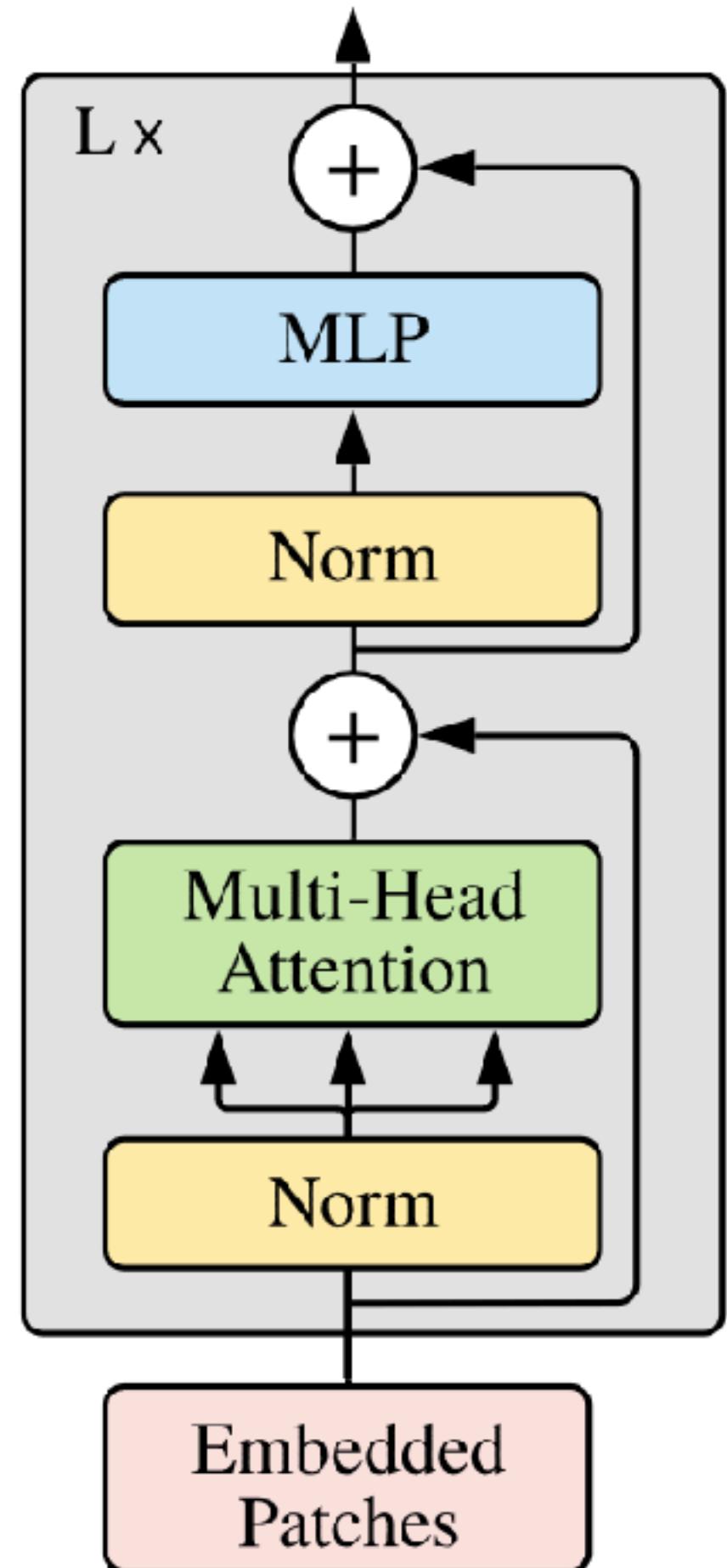
Transformers

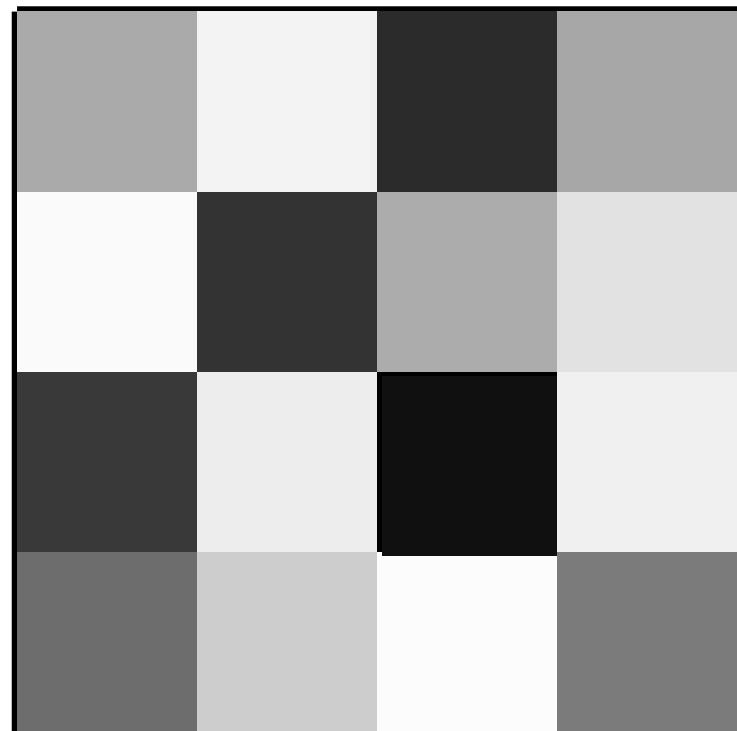
Implementation

- **Residual connections** \Rightarrow Helps for deep networks
- **LayerNorm** \Rightarrow Numerical stability
- **Dropout** before residual connections \Rightarrow Regularisation
- Off-the-shelf implementation:
`torch.nn.TransformerEncoderLayer`


$$\text{Input Grid} + \begin{pmatrix} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu \gamma_5 \\ \gamma_5 \end{pmatrix} = \text{Output}$$

Transformer Encoder



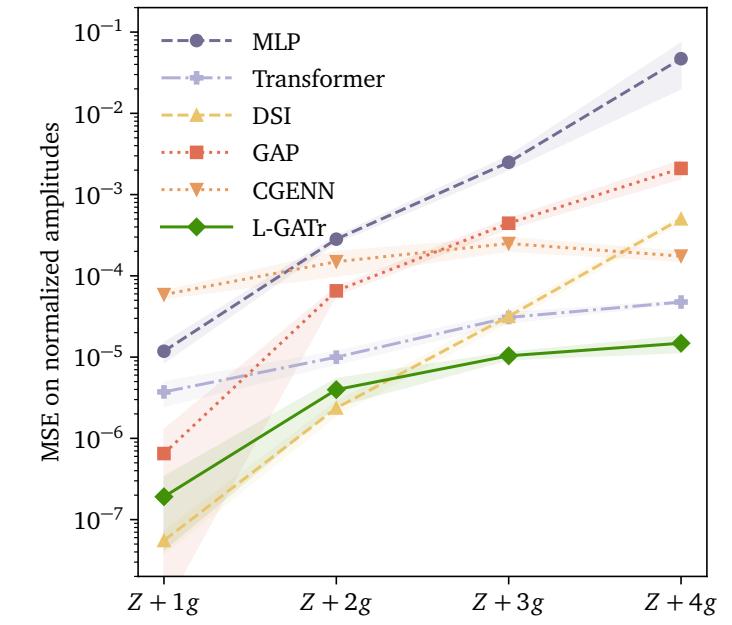


Transformer
architecture



Geometric algebra
representations

**Lorentz-Equivariant
Geometric **Algebra**
Transformer**



Strong performance
on diverse problems

Geometric algebra

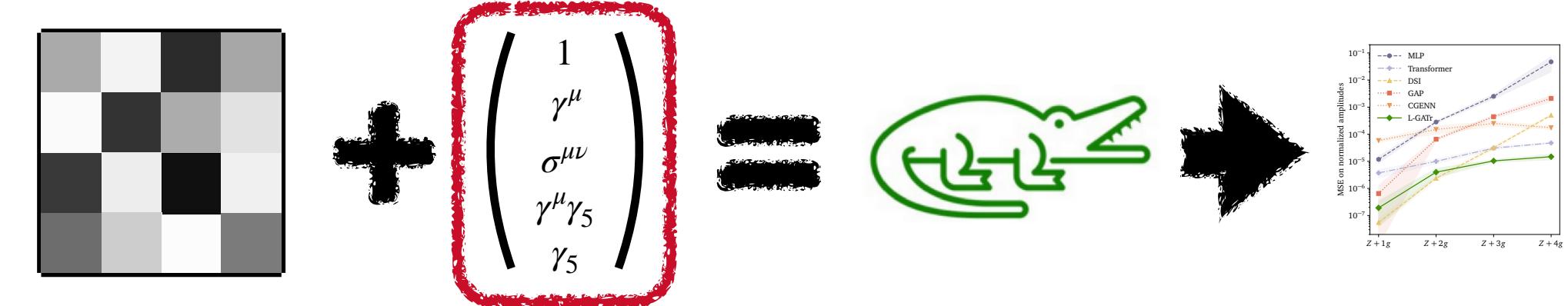
Geometric algebra = Clifford algebra

Geometric algebra = Vector space + geometric product $xy = \frac{\{x, y\}}{2} + \frac{[x, y]}{2}$

- Symmetric part $\{x, y\}$: scalar/inner product
- Antisymmetric part $[x, y]$: outer product (yields higher-order objects)

Spacetime algebra: Geometric algebra over vector space \mathbb{R}^4

- Basis elements γ^μ are orthonormal: $\{\gamma^\mu, \gamma^\nu\} = 2g^{\mu\nu}$
- The Dirac algebra is the same up to $\mathbb{R} \rightarrow \mathbb{C}$

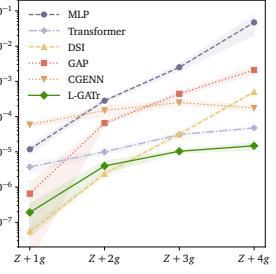
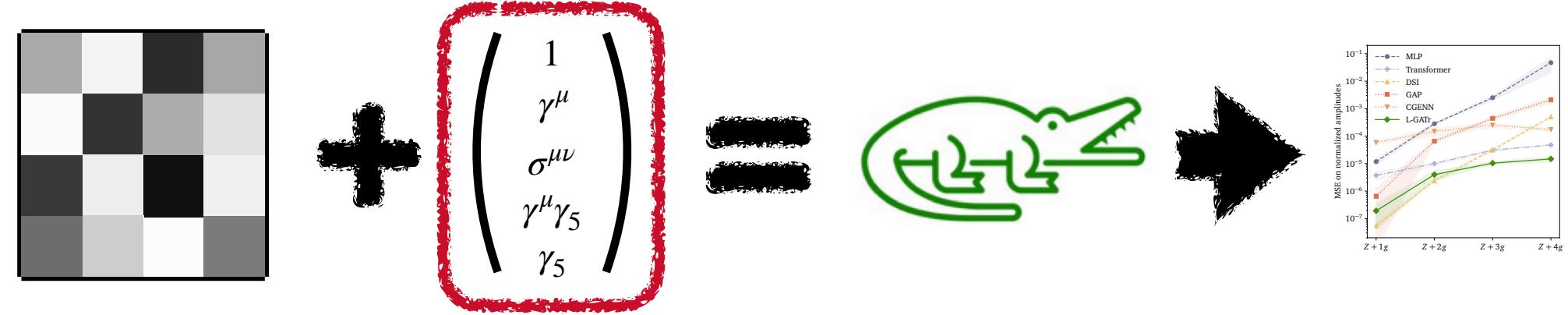


Geometric algebra

Building multivectors

- Scalar and vectors $1, \gamma^\mu$ (1+4 objects)
- Product of two vectors: $\gamma^\mu \gamma^\nu = \frac{\{\gamma^\mu, \gamma^\nu\}}{2} + \frac{[\gamma^\mu, \gamma^\nu]}{2} = g^{\mu\nu} + \sigma^{\mu\nu}$ (6 new objects)
- Axial vector: $\epsilon_{\mu\nu\rho\sigma} \gamma^\nu \gamma^\rho \gamma^\sigma$ (4 new objects)
- Pseudoscalar: $\gamma^5 = \gamma^0 \gamma^1 \gamma^2 \gamma^3 = \frac{1}{4!} \epsilon_{\mu\nu\rho\sigma} \gamma^\mu \gamma^\nu \gamma^\rho \gamma^\sigma$ (1 new object)

Multivector: $x = x^S 1 + x_\mu^V \gamma^\mu + x_{\mu\nu}^B \sigma^{\mu\nu} + x_\mu^A \gamma^\mu \gamma^5 + x^P \gamma^5$ with $(x^S, x_\mu^V, x_{\mu\nu}^B, x_\mu^A, x^P) \in \mathbb{R}^{16}$



Geometric algebra

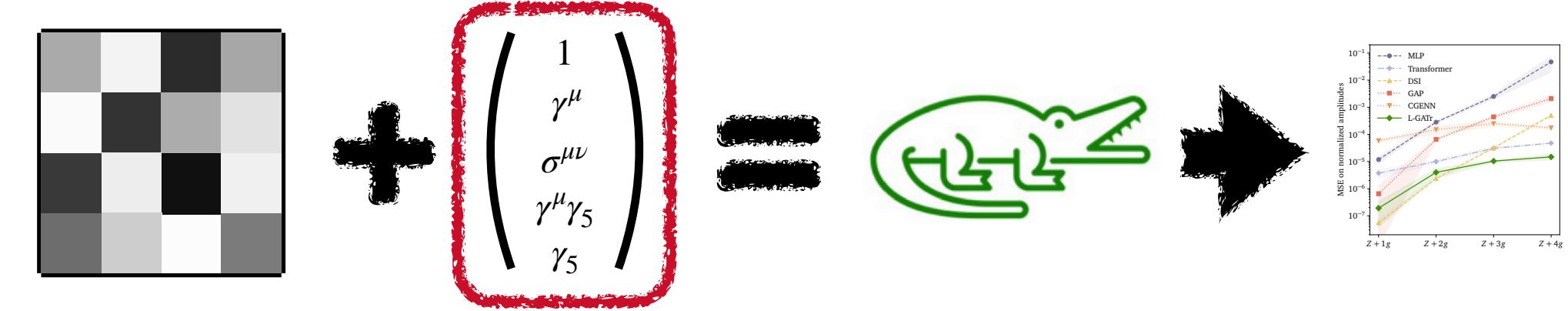
Physics with multivectors

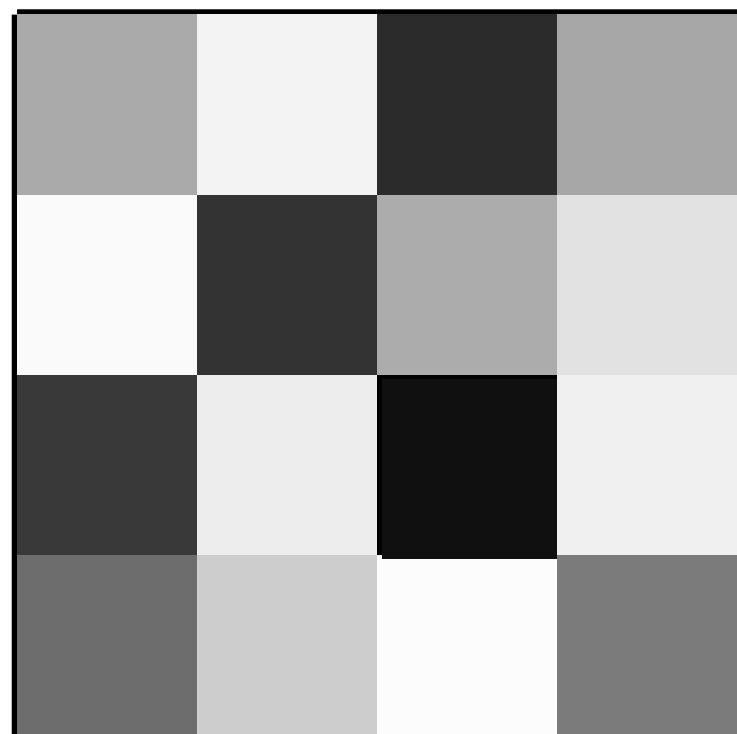
- Particle: $x = p_\mu \gamma^\mu$
- Boost in z direction: $v = e^{\omega \sigma^{03}} = 1 \cosh \omega + \sigma^{03} \sinh \omega$
- Boost applied to particle moving in z direction $x = E\gamma^0 + p_z \gamma^3$:

$$vxv^{-1} = (1 \cosh \omega + \sigma^{03} \sinh \omega)(E\gamma^0 + p_z \gamma^3)(1 \cosh \omega - \sigma^{03} \sinh \omega)$$

$$= (E \cosh \omega - p_z \sinh \omega)\gamma^0 + (p_z \cosh \omega - E \sinh \omega)\gamma^3$$

$$= E^{\text{boosted}}\gamma^0 + p_z^{\text{boosted}}\gamma^3$$



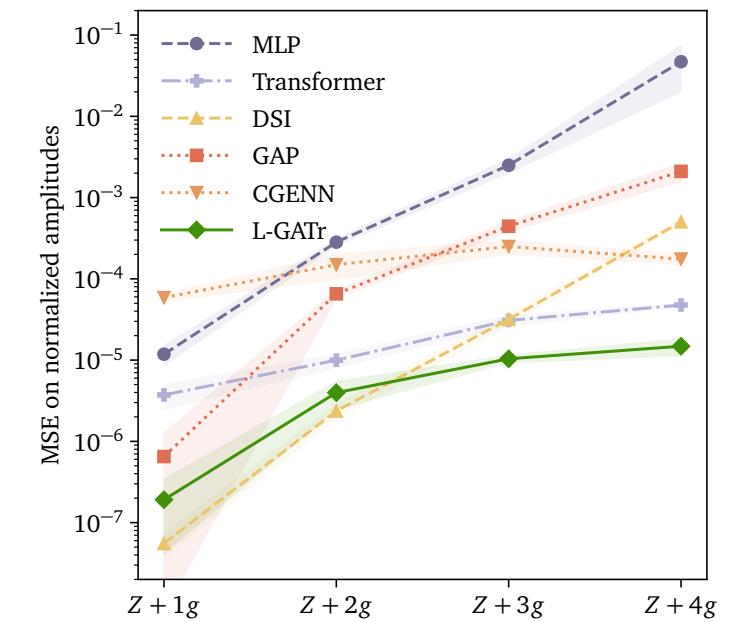


Transformer
architecture



Geometric algebra
representations

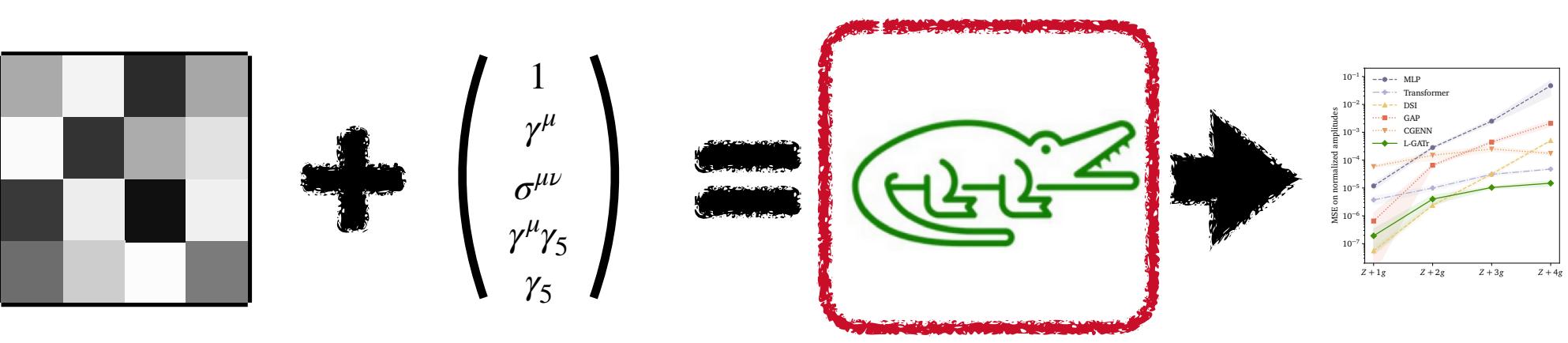
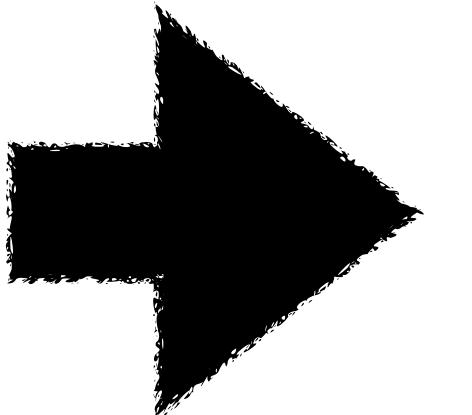
Lorentz-Equivariant
Geometric **A**lgebra
Transformer



Strong performance
on diverse problems

L-GATr

Geometric algebra representations

 x^S 

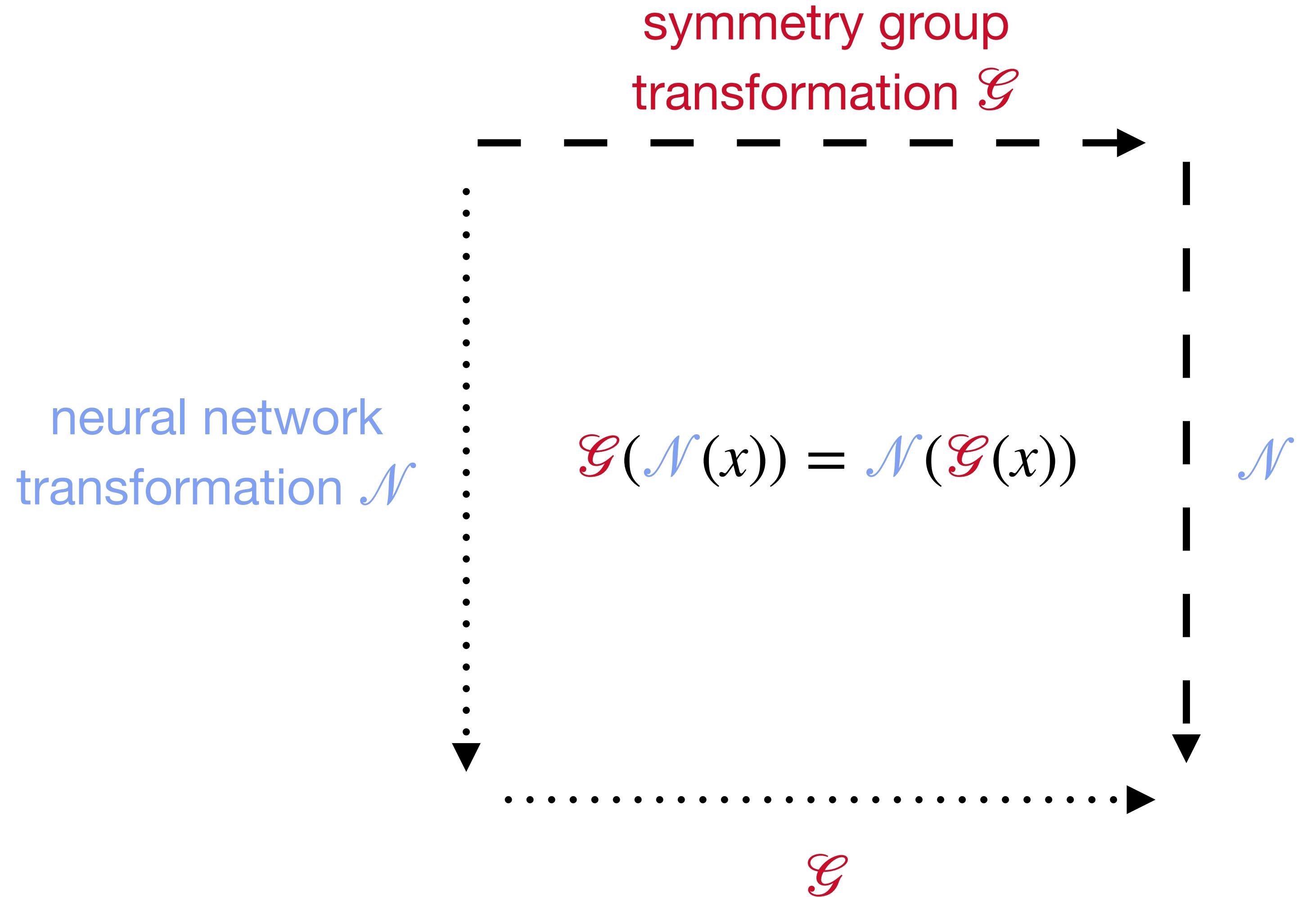
scalars

 x^S
 x_μ^V
 $x_{\mu\nu}^B$
 x_μ^A
 x^P $\in \mathbb{R}^{16}$

multivectors
(+ extra scalars)

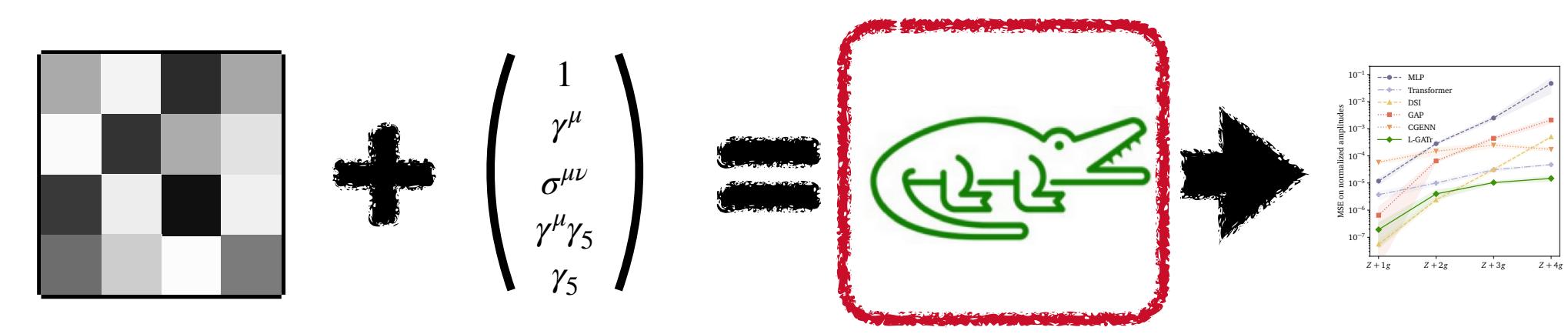
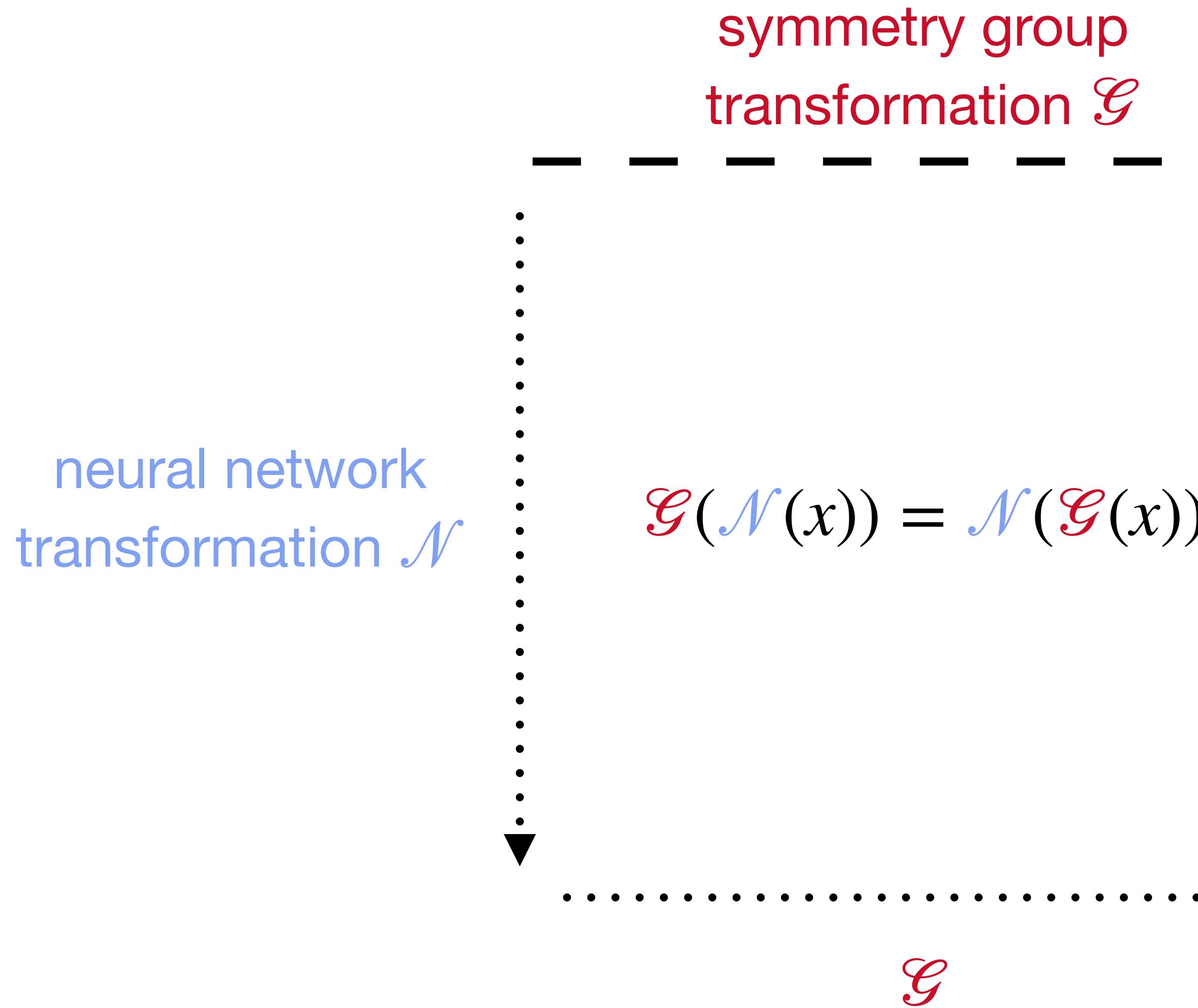
L-GATr

Equivariance



L-GATr

Equivariance



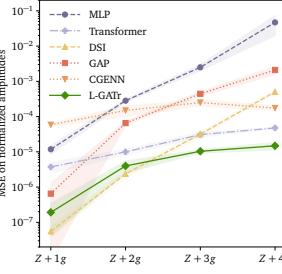
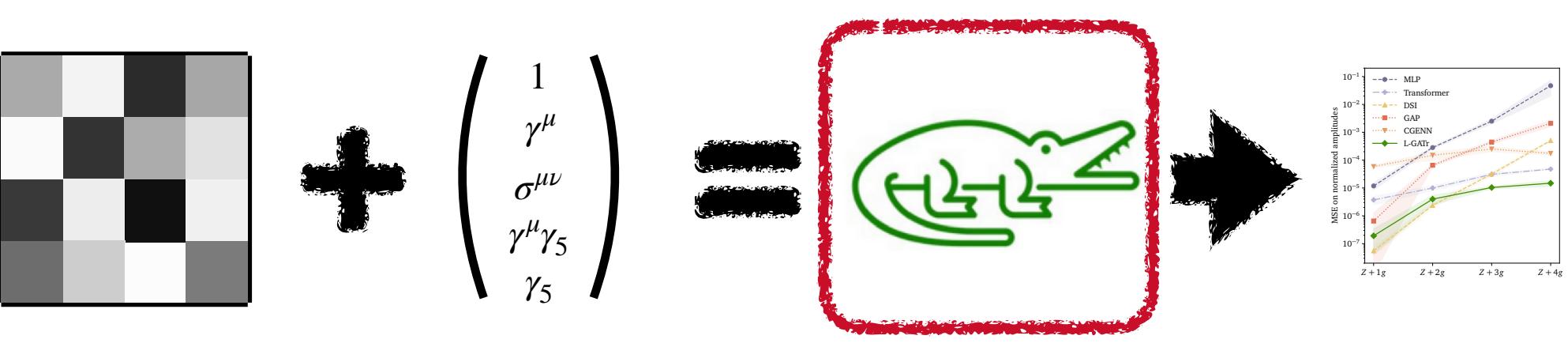
Equivariance on geometric algebra representations:

Multivector grades form subrepresentations of the Lorentz group → All components within the same grade have to transform equally

$$\begin{pmatrix} x^S \\ x_\mu^V \\ x_{\mu\nu}^B \\ x_\mu^A \\ x_\mu^P \end{pmatrix} \xrightarrow{\mathcal{G}} \begin{pmatrix} w^S x^S \\ w^V x_\mu^V \\ w^B x_{\mu\nu}^B \\ w^A x_\mu^A \\ w^P x_\mu^P \end{pmatrix}$$

L-GATr

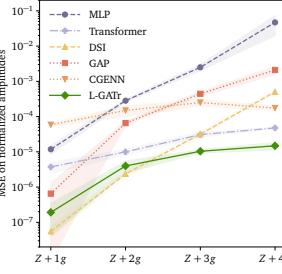
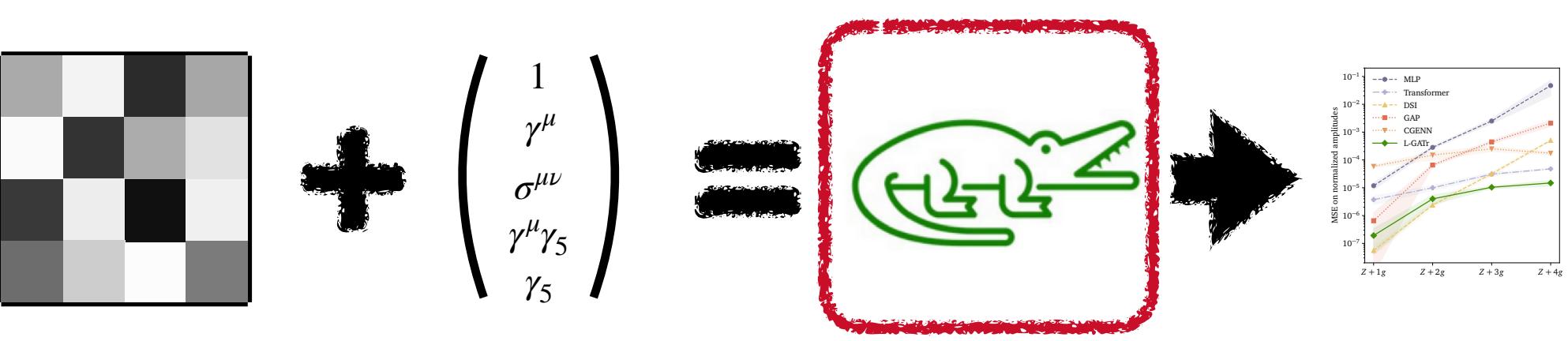
GATr-ing all transformer layers



Layer type	Transformer	L-GATr
Linear(x)	$vx + w$	
Attention(q, k, v) _{ic}	$\sum_{j,c'} \text{Softmax}_j \left(\frac{q_{ic'} k_{jc'}}{\sqrt{n_c}} \right) v_{jc}$	
LayerNorm(x)	$x \left[\frac{1}{n_c} \sum_{c=1}^{n_c} x_c^2 + \epsilon \right]^{-1/2}$	
Activation(x)	GELU(x)	

L-GATr

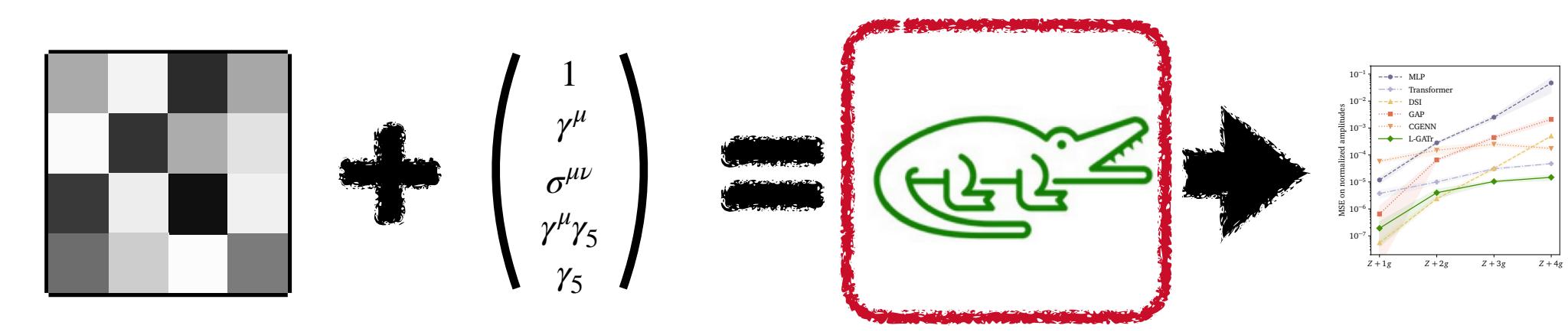
GATr-ing all transformer layers



Layer type	Transformer	L-GATr
Linear(x)	$vx + w$	$\sum_{k=0}^4 v_k \langle x \rangle_k + \sum_{k=0}^4 w_k \gamma^5 \langle x \rangle_k$
Attention(q, k, v) _{ic}	$\sum_{j,c'} \text{Softmax}_j \left(\frac{q_{ic'} k_{jc'}}{\sqrt{n_c}} \right) v_{jc}$	$\sum_{j,c'} \text{Softmax}_j \left(\frac{\langle q_{ic'}, k_{jc'} \rangle}{\sqrt{16n_c}} \right) v_{jc}$
LayerNorm(x)	$x \left[\frac{1}{n_c} \sum_{c=1}^{n_c} x_c^2 + \epsilon \right]^{-1/2}$	$x \left[\frac{1}{n_c} \sum_{c=1}^{n_c} \sum_{k=0}^4 \left \langle \langle x_c \rangle_k, \langle x_c \rangle_k \rangle \right + \epsilon \right]^{-1/2}$
Activation(x)	GELU(x)	GELU($\langle x \rangle_0$) x
GP(x, y)	—	xy

L-GATr

Full architecture

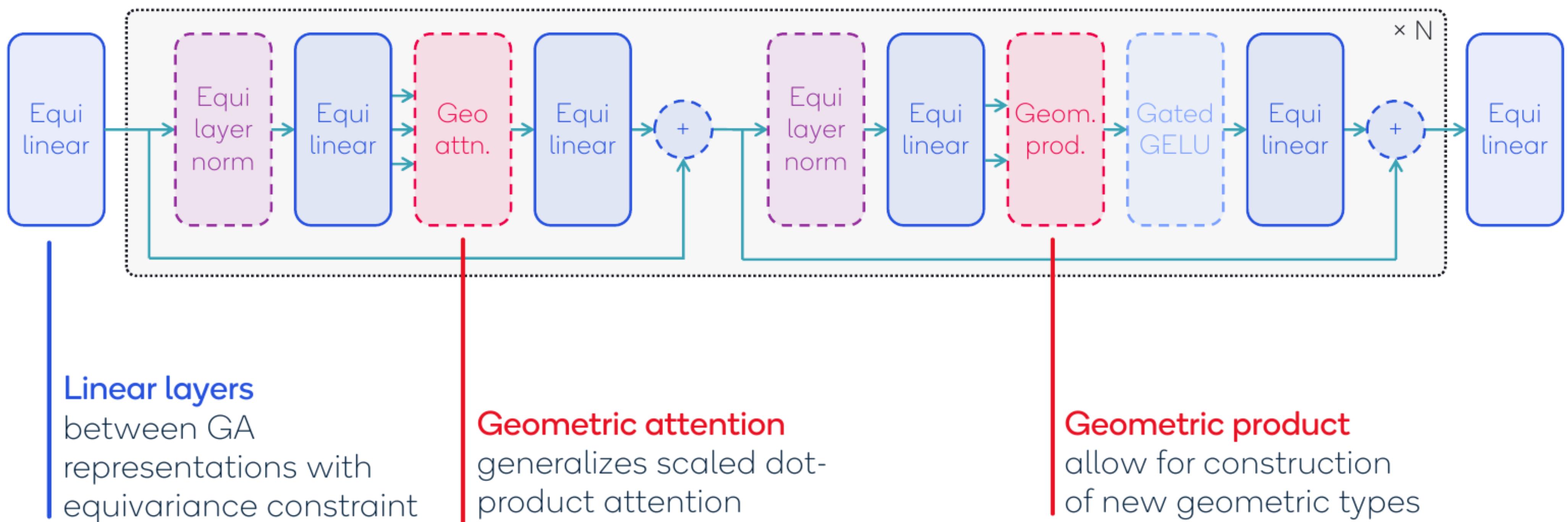


Input and output data

can have one or multiple token dimensions

Attention blocks

can be stacked to large depth, gradients are propagated efficiently



Linear layers

between GA representations with equivariance constraint

Geometric attention

generalizes scaled dot-product attention

Geometric product

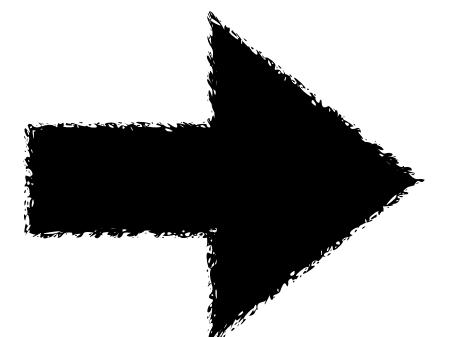
allow for construction of new geometric types

L-GATr

Symmetry breaking with spurious

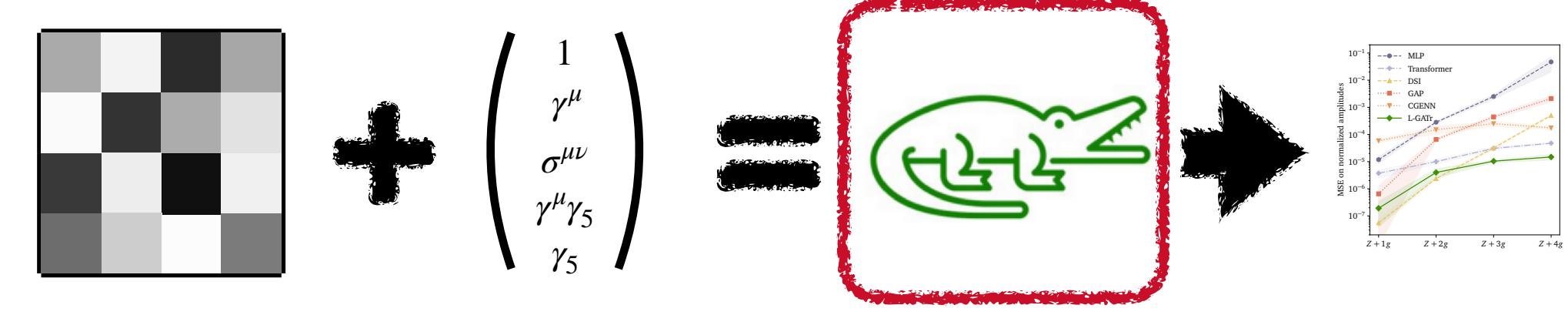
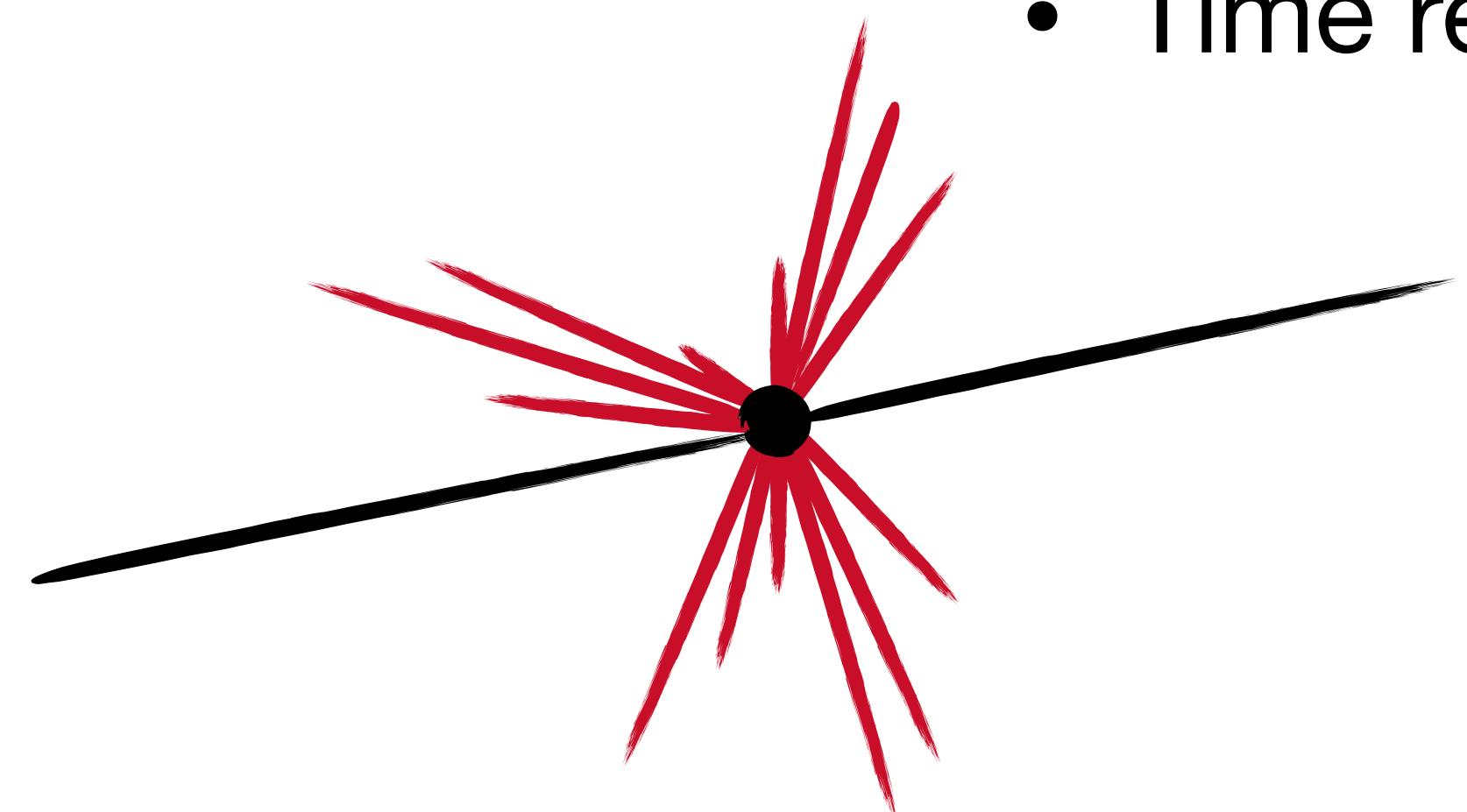
Lorentz symmetry is rarely exact

- Beam direction in collider
- Detector effects
- ...?



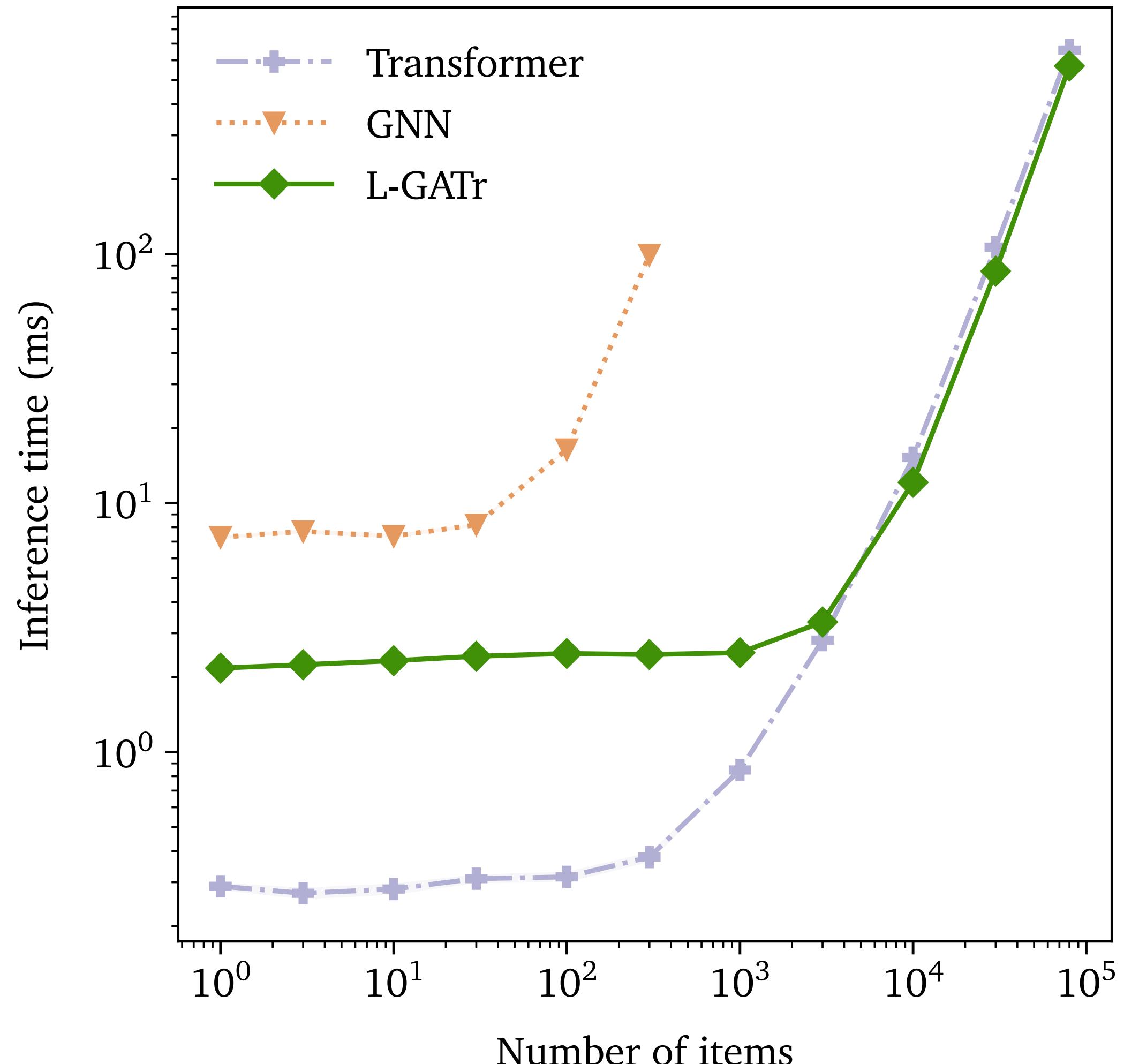
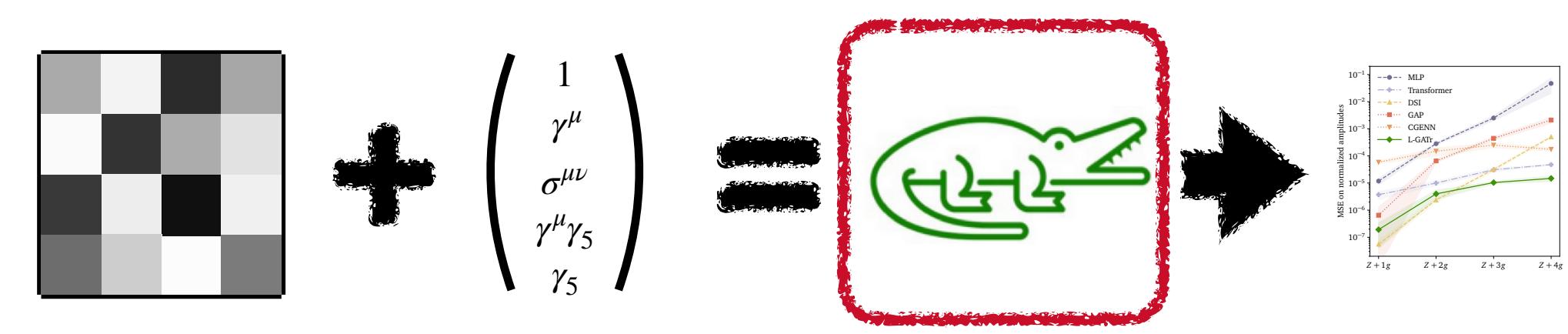
Add a **spurion** to the particle list
(either as token or channel)

- Beam reference: $p^\mu = (0,0,0, \pm 1)$
- Time reference: $p^\mu = (1,0,0,0)$

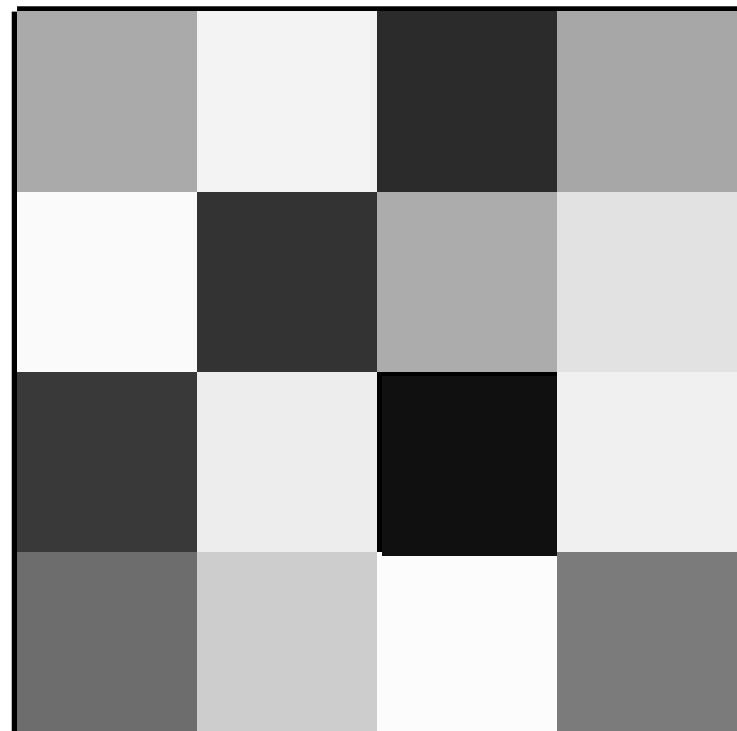


L-GATr

Processing thousands of particles



Transformers scale
better than graph networks

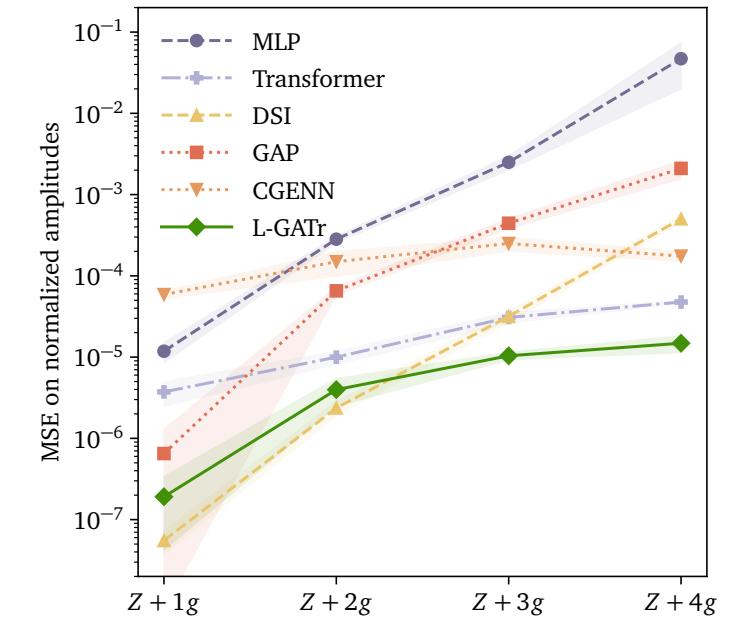


Transformer
architecture



Geometric algebra
representations

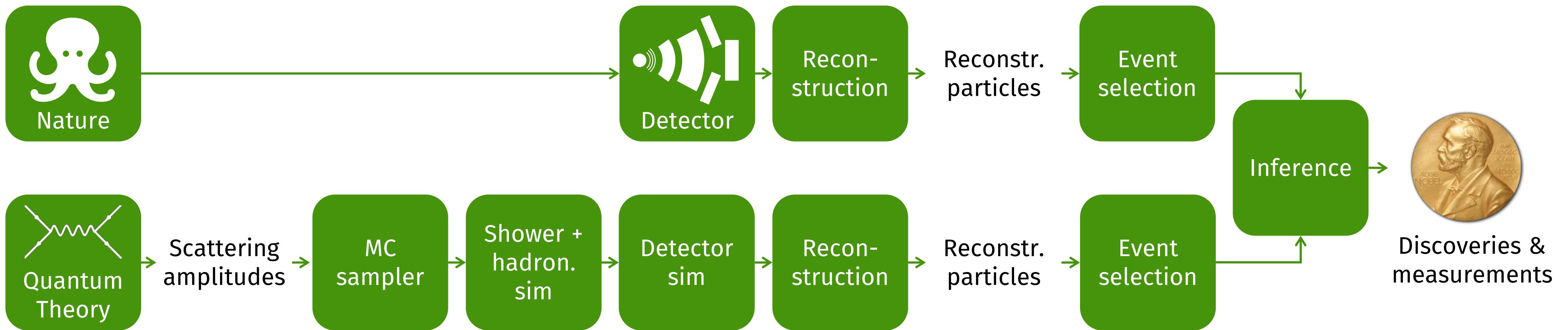
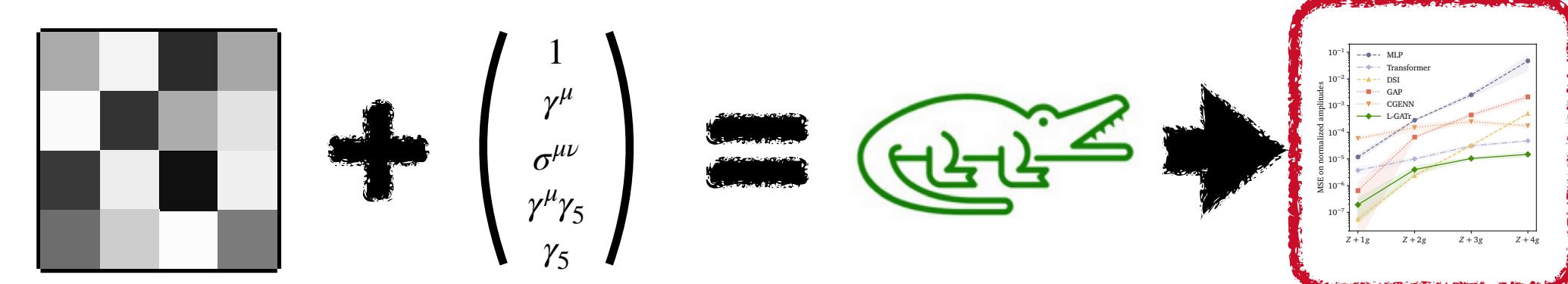
**Lorentz-Equivariant
Geometric **Algebra**
Transformer**



Strong performance
on diverse problems

Experiments

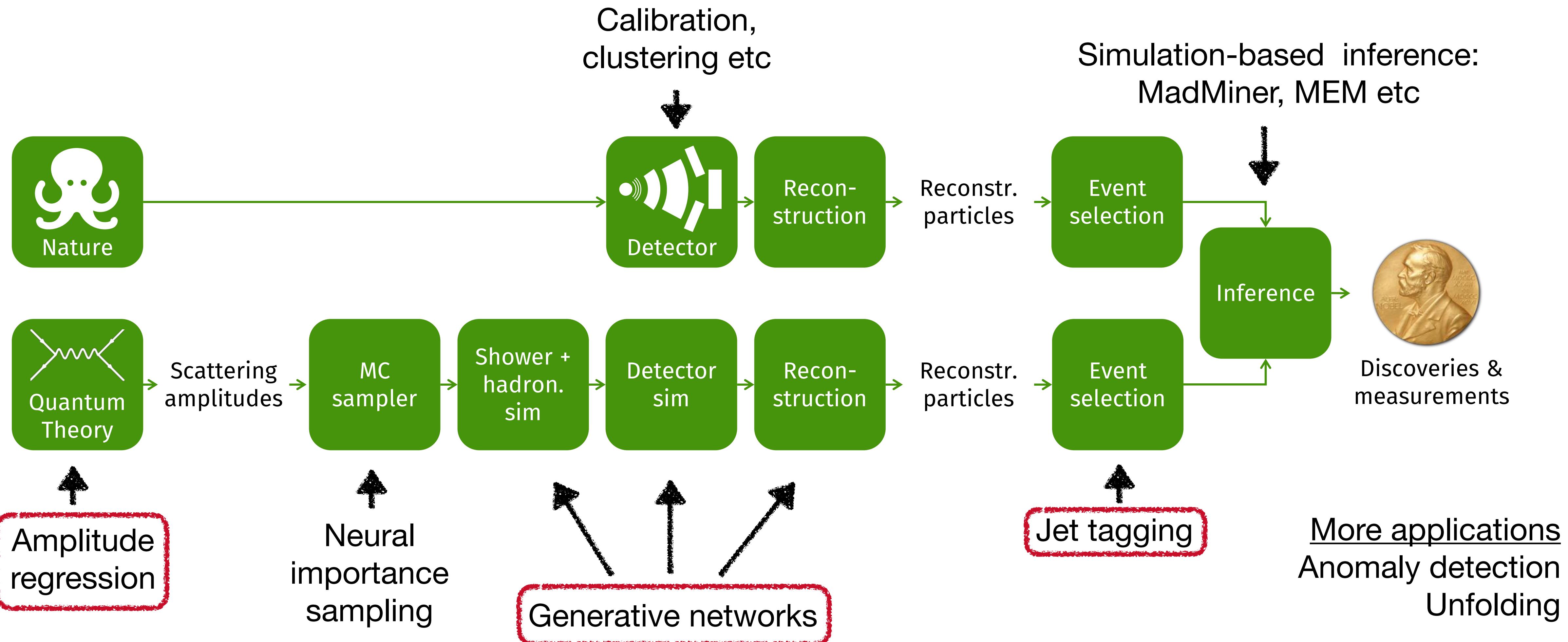
LHC simulation chain



Experiments

LHC simulation chain meets ML

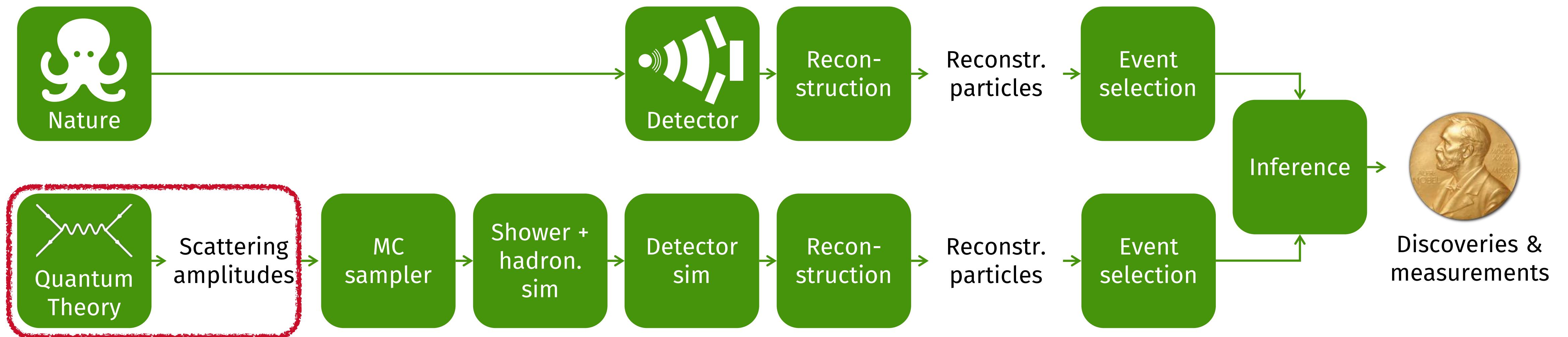
A diagram illustrating a mathematical operation. On the left is a 4x4 grid of gray and black squares. To its right is a plus sign. Next is a matrix with columns labeled 1 , γ^μ , $\sigma^{\mu\nu}$, $\gamma^\mu\gamma_5$, and γ_5 . An equals sign follows. To the right of the equals sign is a circled green arrow pointing to a plot. The plot shows 'MSB on normalized amplitudes' on a logarithmic y-axis (from 10^{-7} to 10^{-1}) versus 'Z+ig' through 'Z+4g' on the x-axis. Several data series are shown: MLP (blue circles), Transformer (light blue circles), GAP (orange triangles), CGNN (green diamonds), L-GAT (green squares), and a theoretical curve (red dashed line).



Experiments

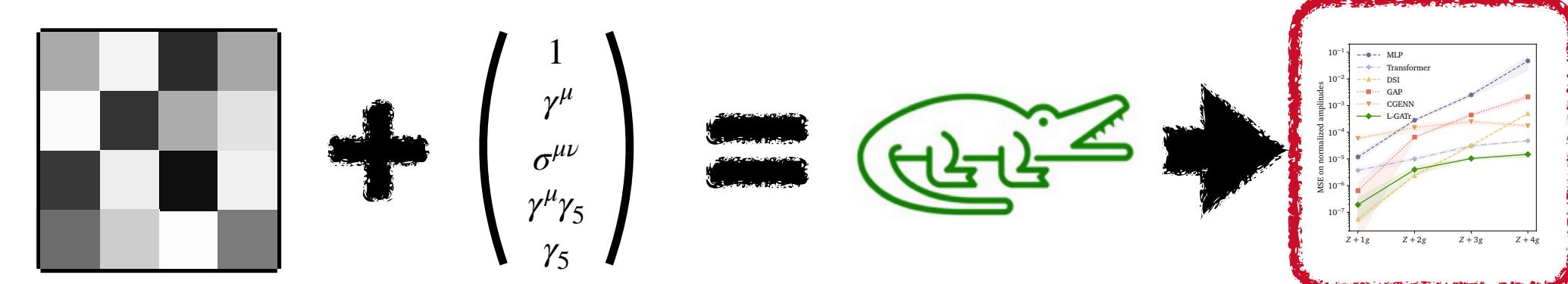
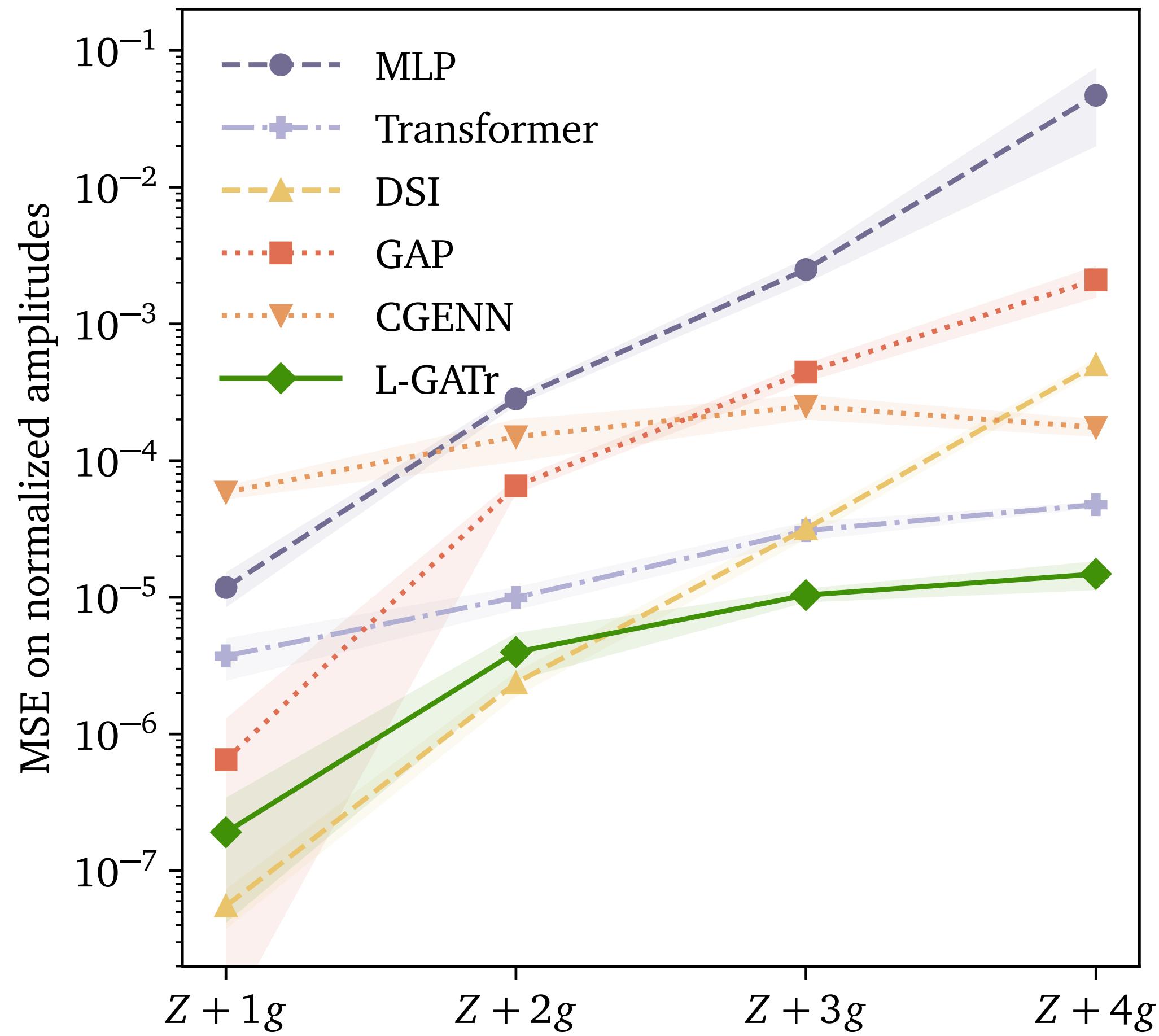
Amplitude regression

$$\begin{matrix} & + \\ \begin{matrix} \text{grid} \\ \text{+} \end{matrix} & \left(\begin{array}{c} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu \gamma_5 \\ \gamma_5 \end{array} \right) \end{matrix} = \boxed{\text{green checkmark}}$$



Experiments

Amplitude regression

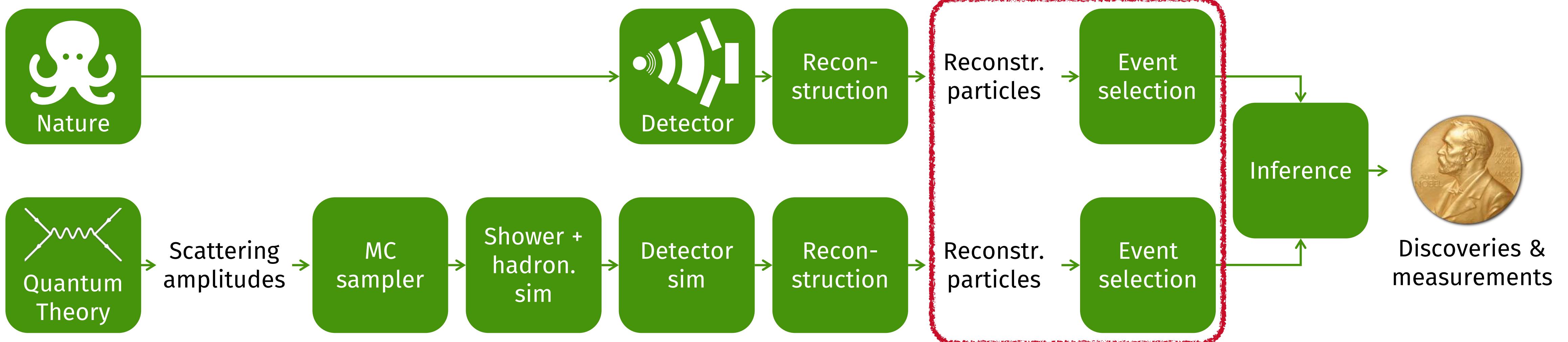


L-GATr scales best to **high multiplicity**, where amplitude surrogates are most useful

Experiments

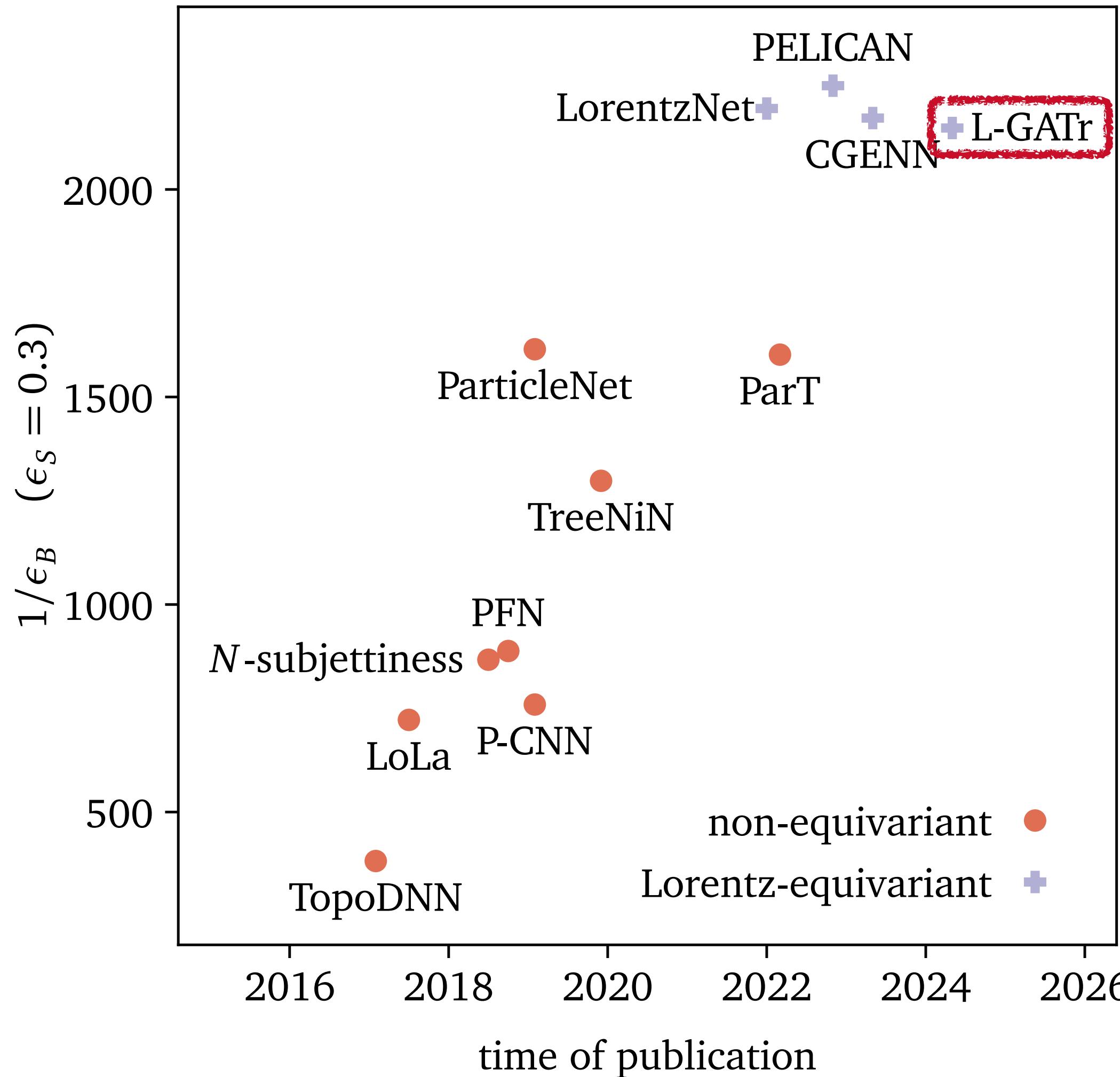
Top tagging

$$\text{Measurement} + \begin{pmatrix} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu \gamma_5 \\ \gamma_5 \end{pmatrix} = \text{True Value}$$



Experiments

Top tagging

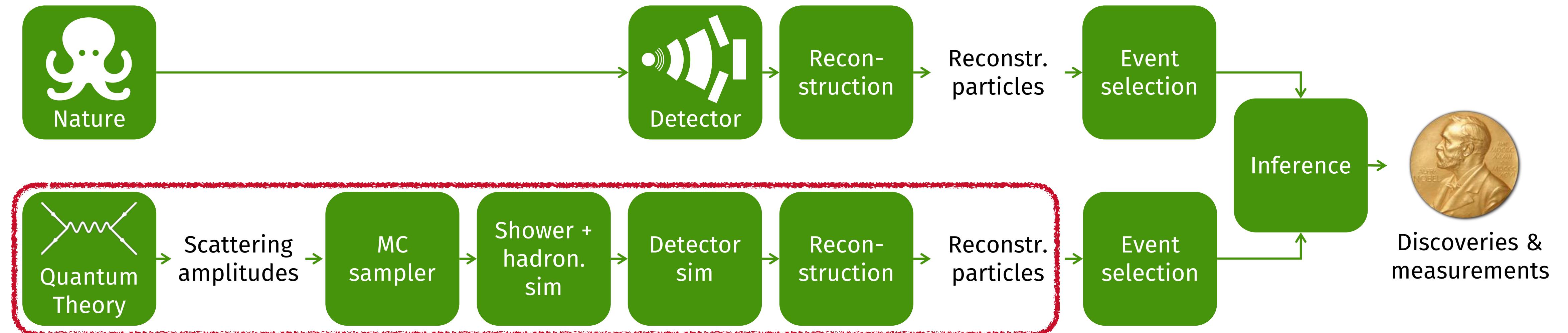
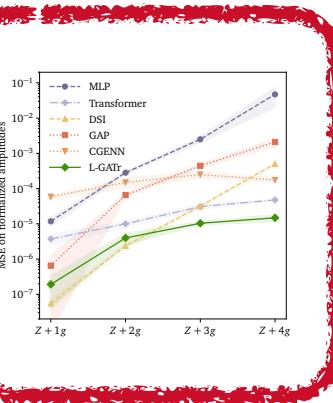


L-GATr is on par with the best equivariant (*) baselines

Experiments

Event generation

$$\begin{pmatrix} 1 \\ \gamma^\mu \\ \sigma^{\mu\nu} \\ \gamma^\mu\gamma_5 \\ \gamma_5 \end{pmatrix} = \text{checkmark}$$



Experiments

Event generation



Continuous normalising flows (CNF)

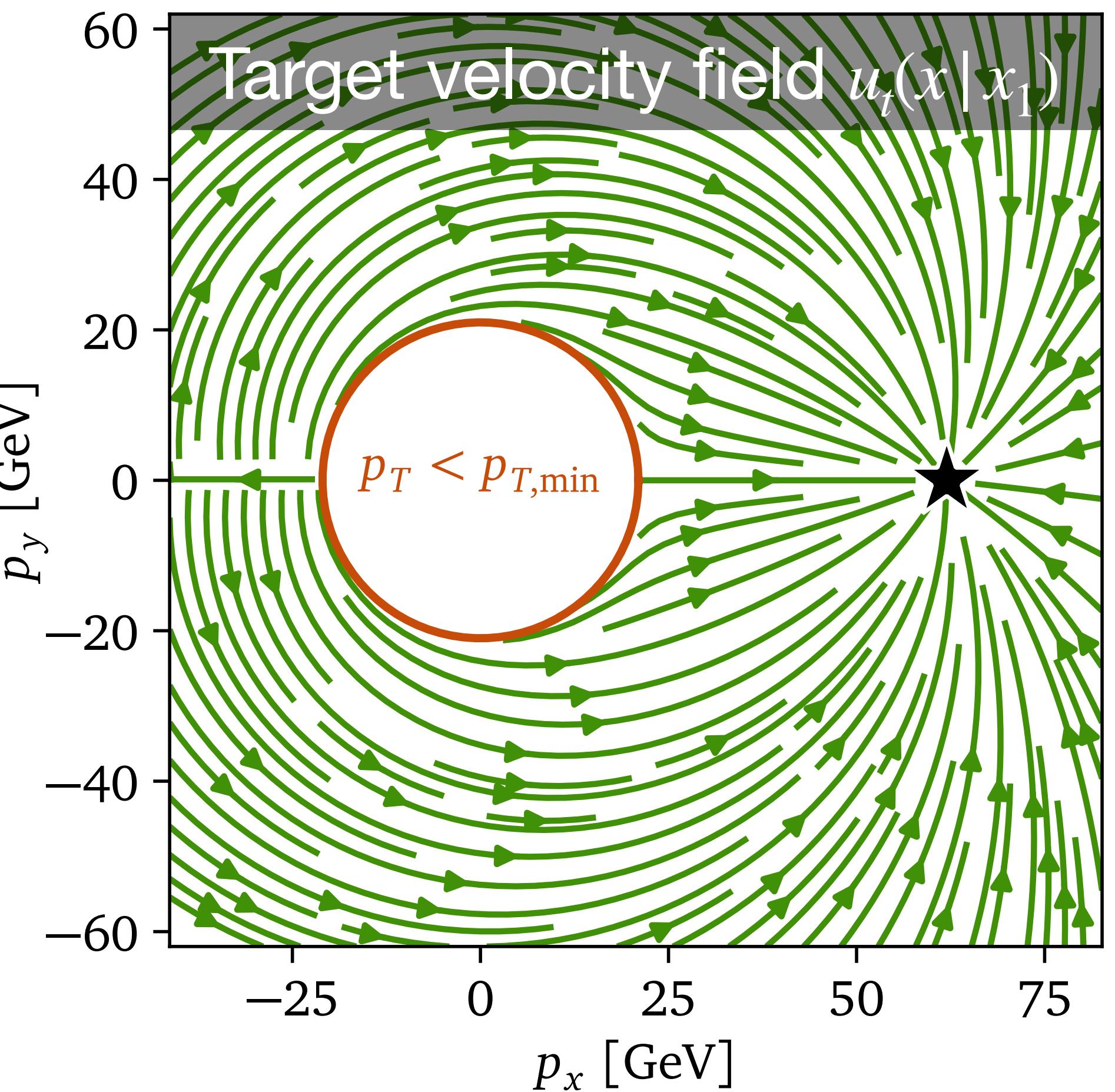
connect a simple base density
to a complex target density
through a neural differential equation

$$\frac{d}{dt}x = v_t(x)$$

Conditional flow matching (CFM)
is a simple way to train CNFs
by comparing the learned velocity $v_t(x)$
to a conditional target velocity $u_t(x | x_1)$

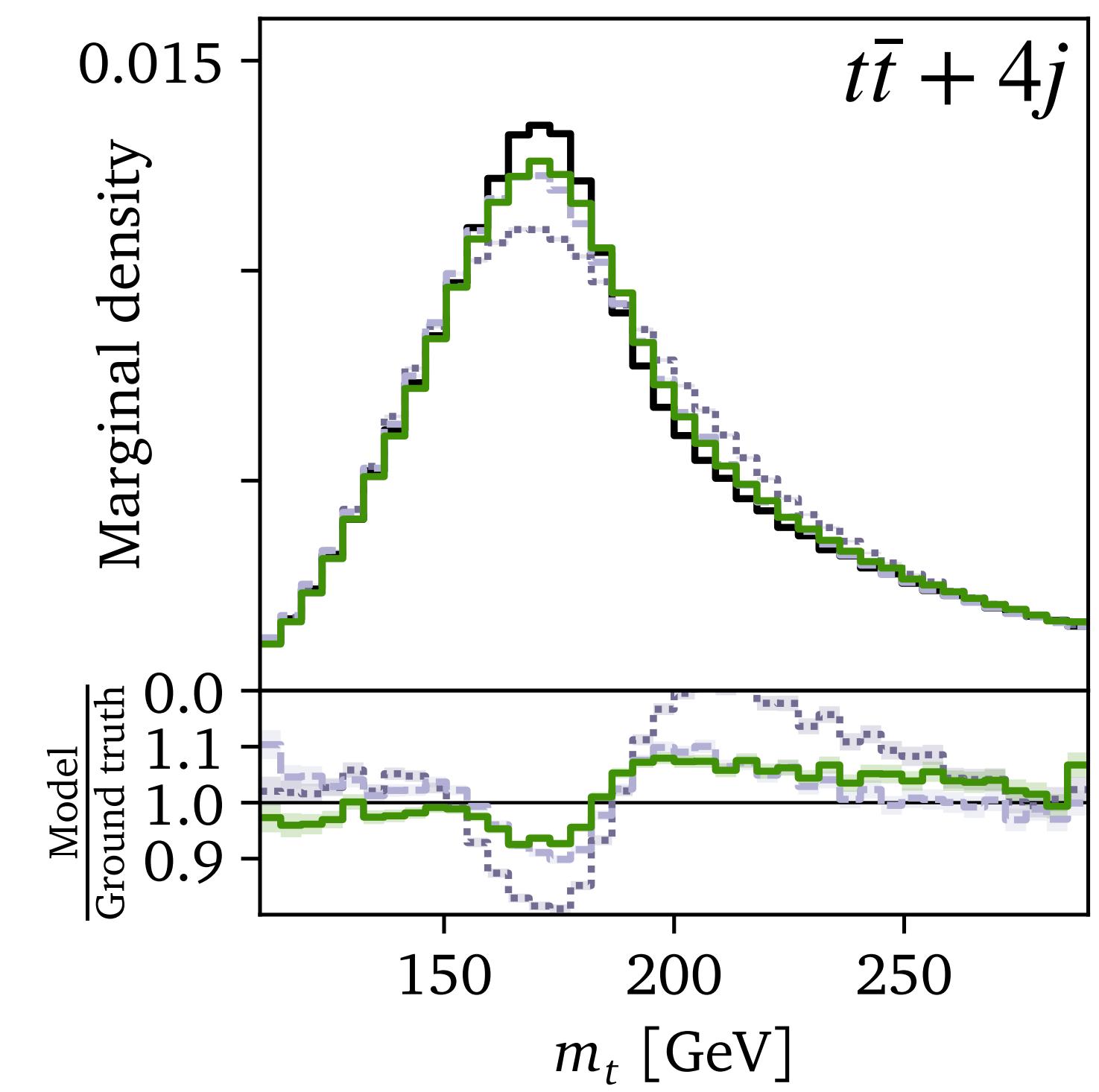
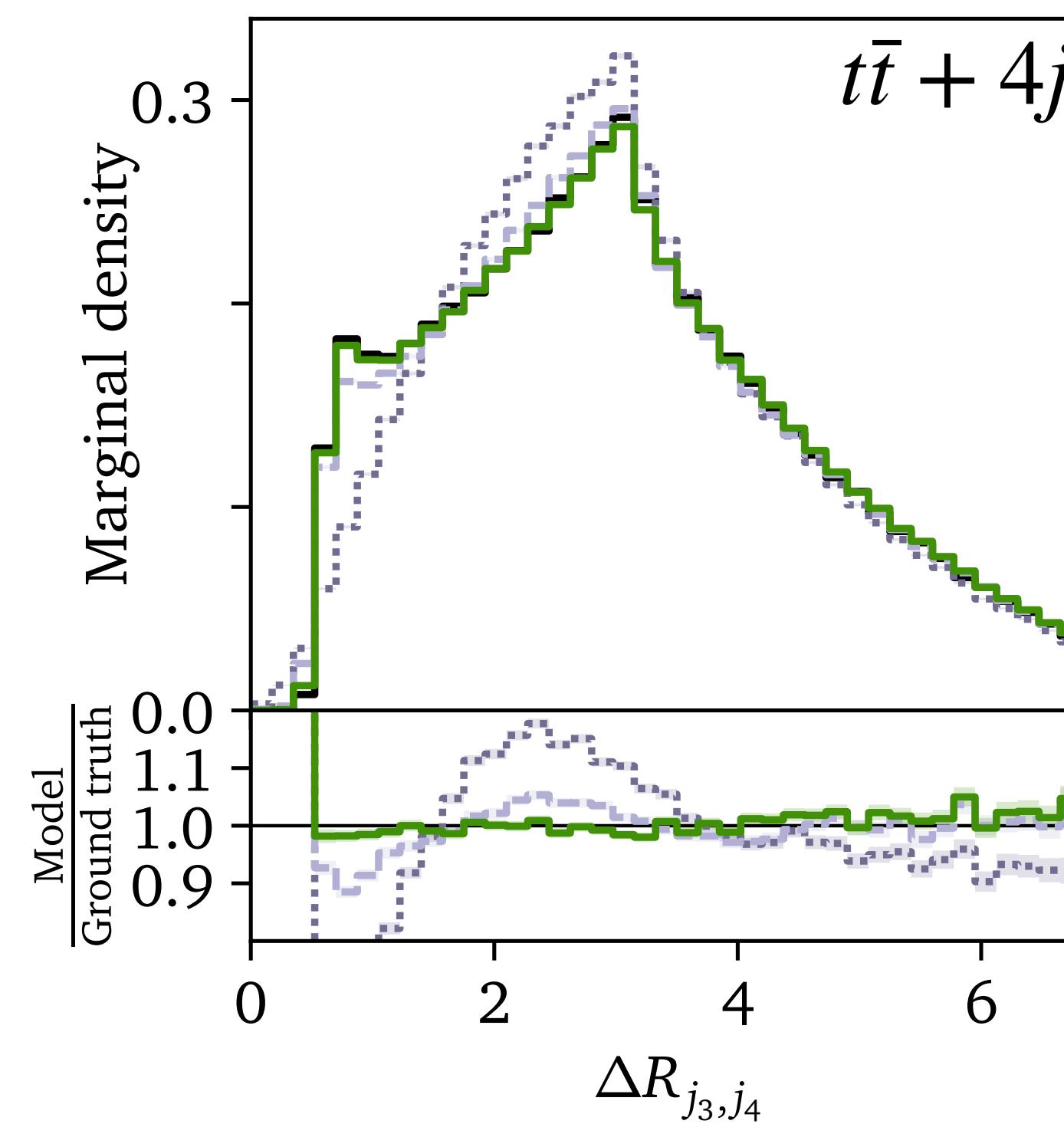
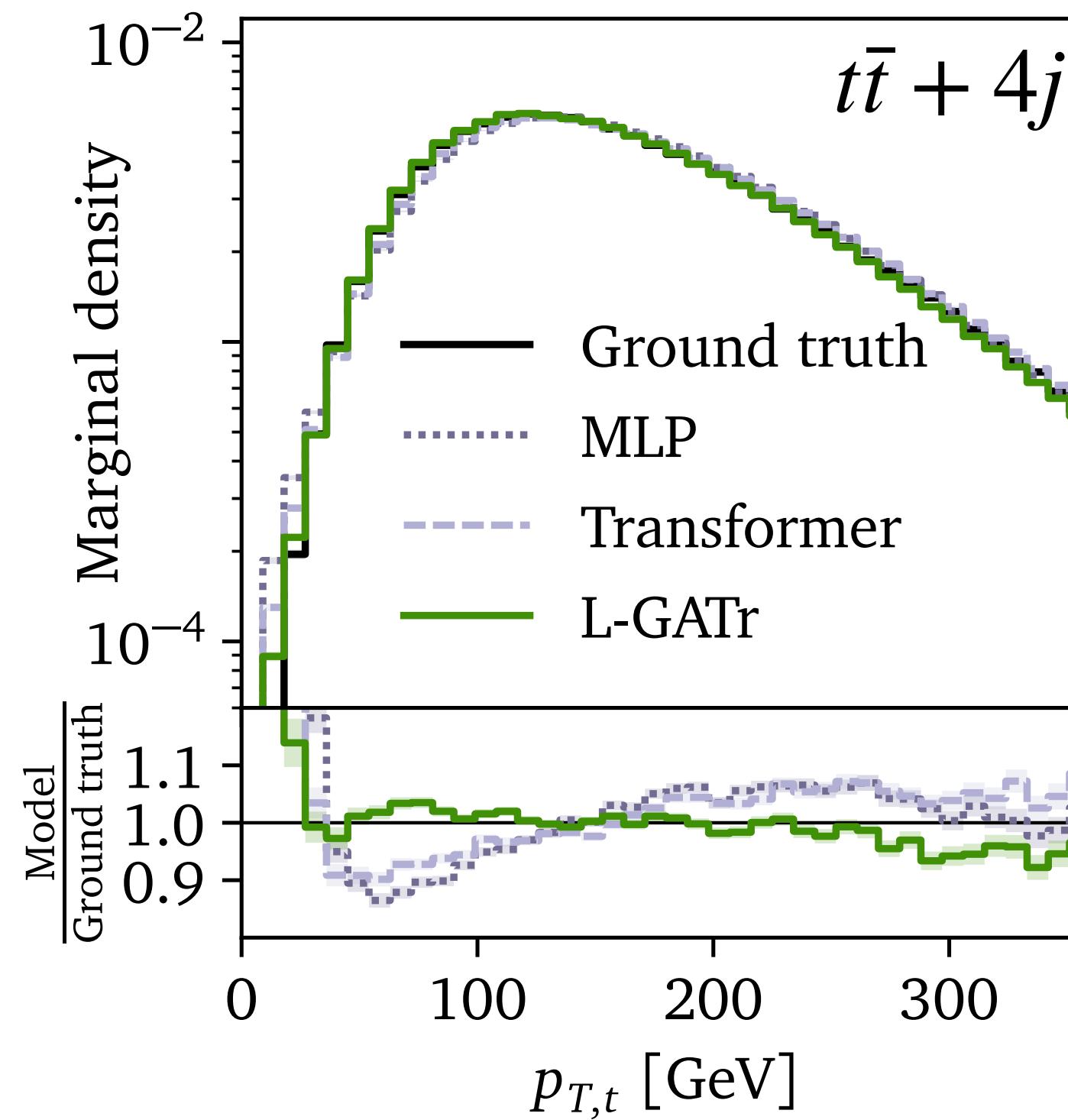
Continuous normalising flows
arXiv:1806.07366

Conditional flow matching
arXiv:2210.02747



Experiments

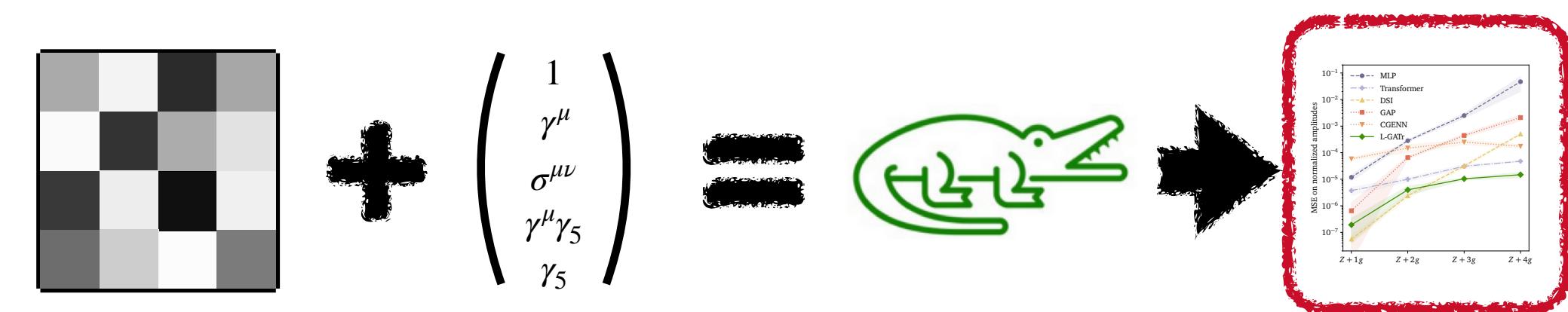
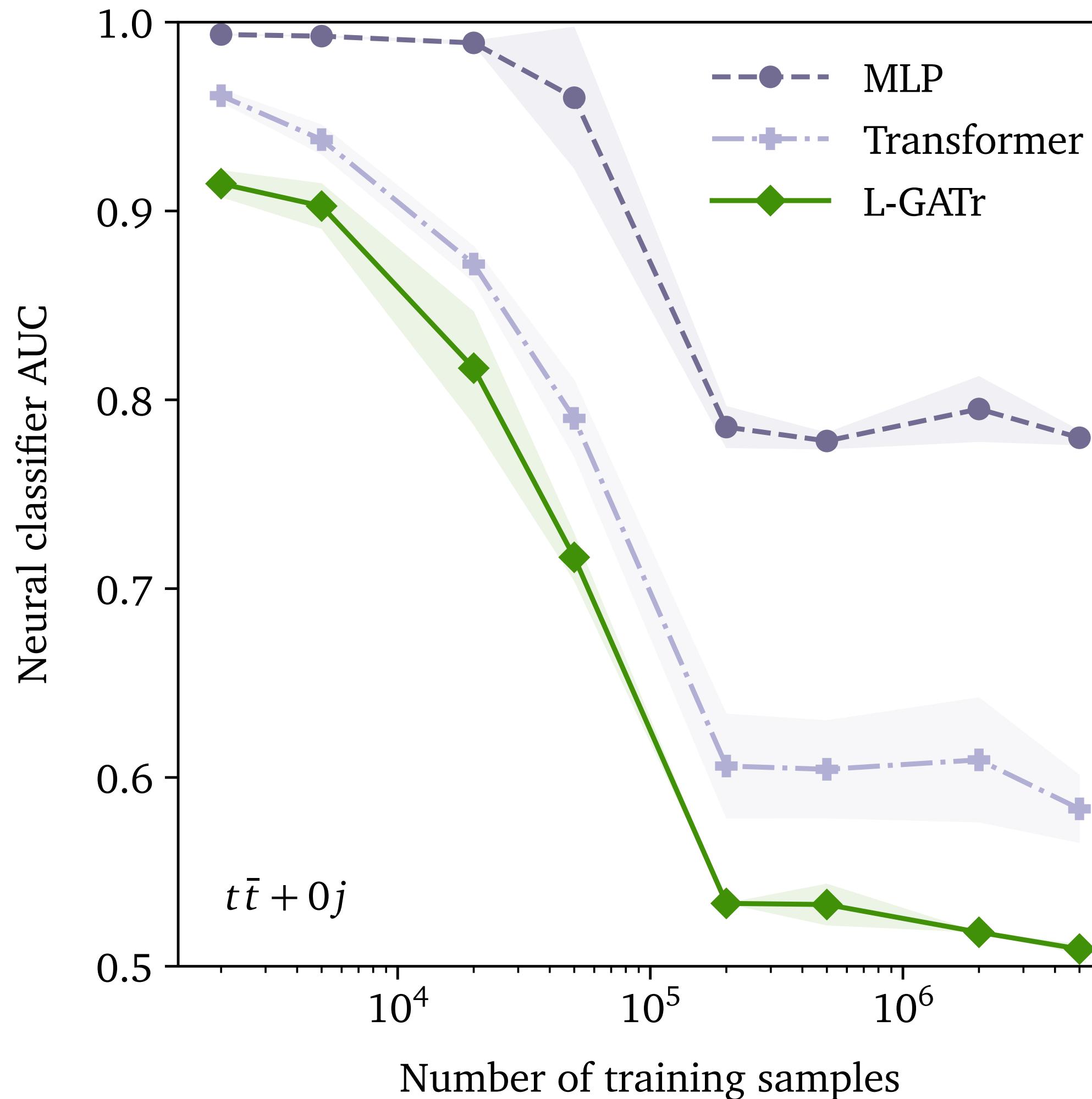
Event generation



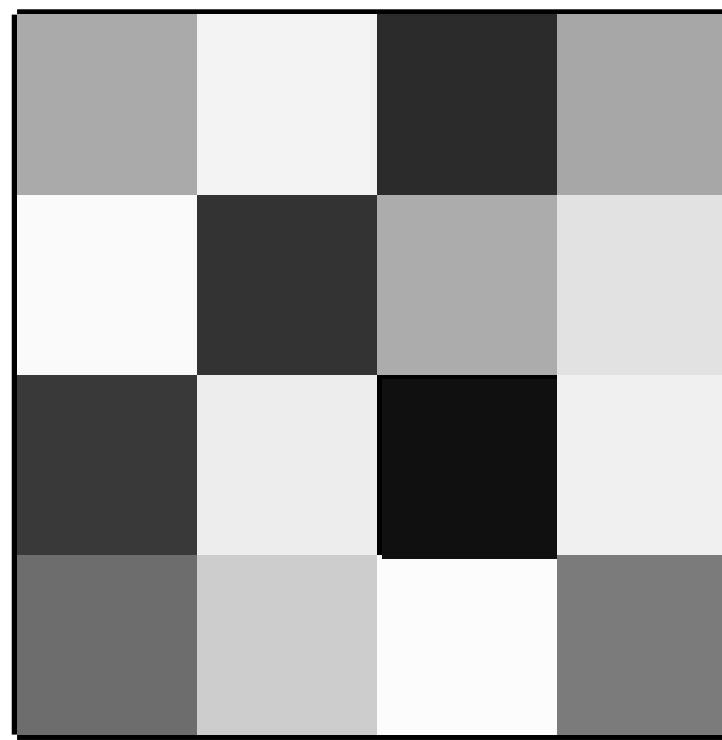
L-GATr helps with tricky kinematic features

Experiments

Event generation



L-GATr generates samples that a classifier can almost not distinguish from the ground truth

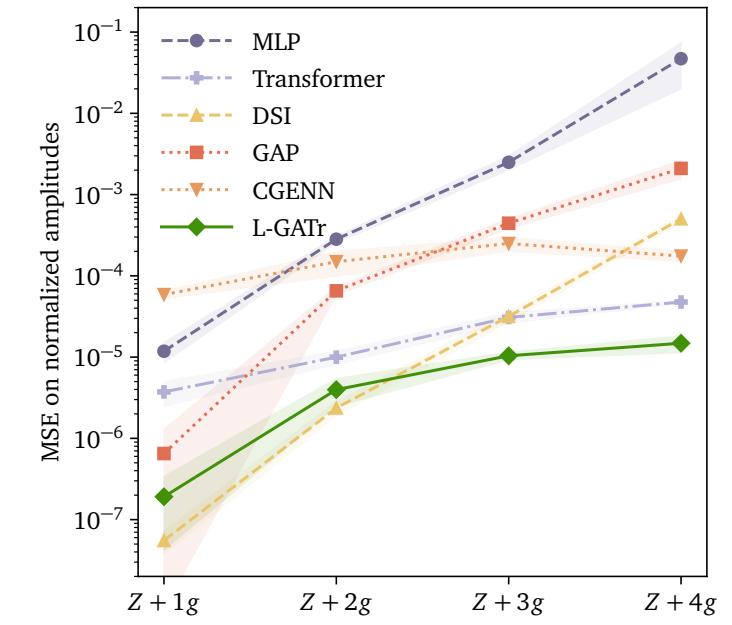


Transformer
architecture



Geometric algebra
representations

**Lorentz-Equivariant
Geometric **A**lgebra
Transformer**



Strong performance
on diverse problems

L-GATr combines **equivariance** and **scalability**



Victor Bresó



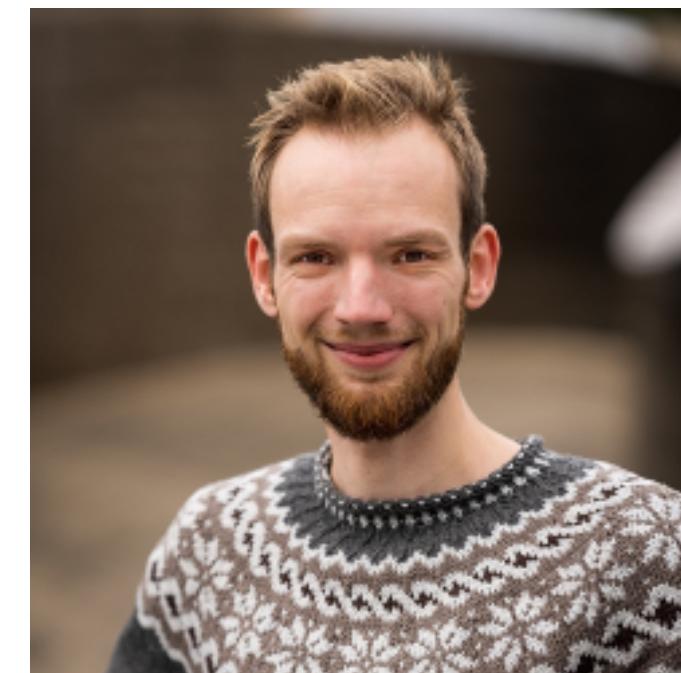
Pim de Haan



Tilman Plehn



Jesse Thaler



Johann Brehmer

Geometric Algebra Transformer

E(3)-equivariant version

Johann Brehmer*, Pim de Haan*, Sönke Behrends, Taco Cohen
NeurIPS 2023, arXiv:2305.18415



E(3)-GATr paper



E(3)-GATr code

Lorentz-Equivariant Geometric Algebra Transformer for High-Energy Physics

Jonas Spinner*, Victor Bresó*, Pim de Haan,
Tilman Plehn, Jesse Thaler, Johann Brehmer
NeurIPS 2024, arXiv:2405.14806



L-GATr paper

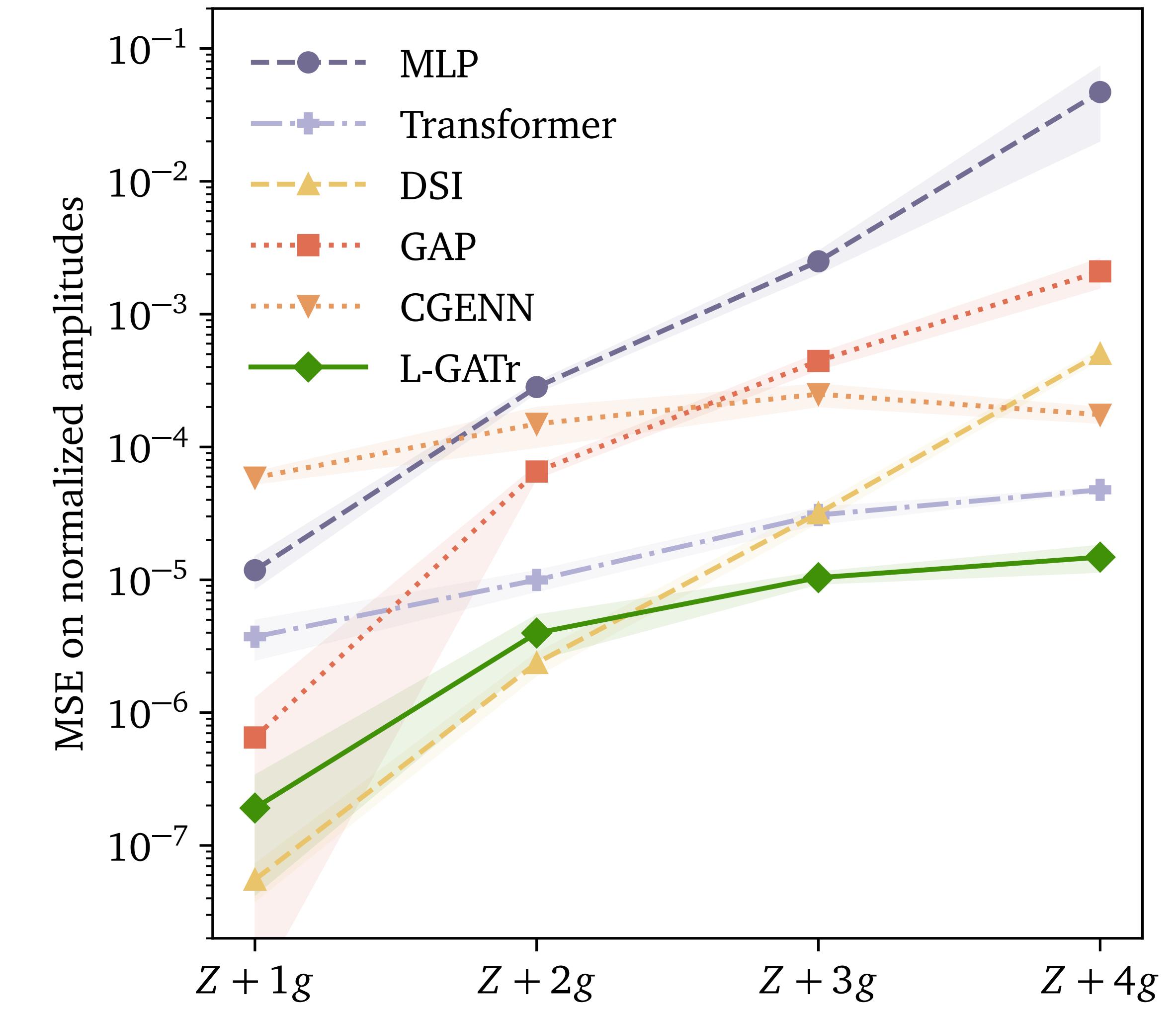
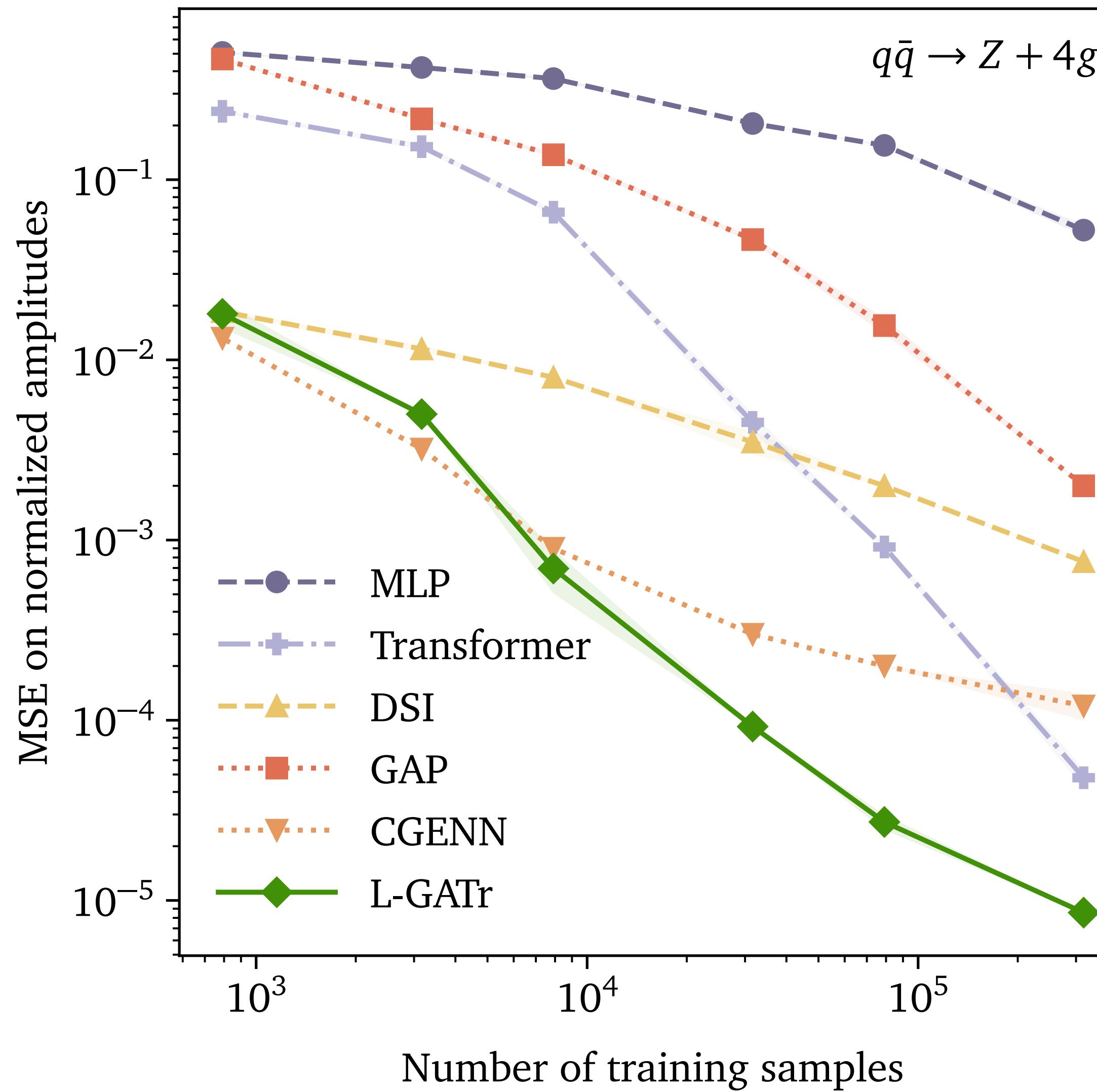


L-GATr code

What would **you** use L-GATr for?

Bonus material

Amplitude regression



Experiments

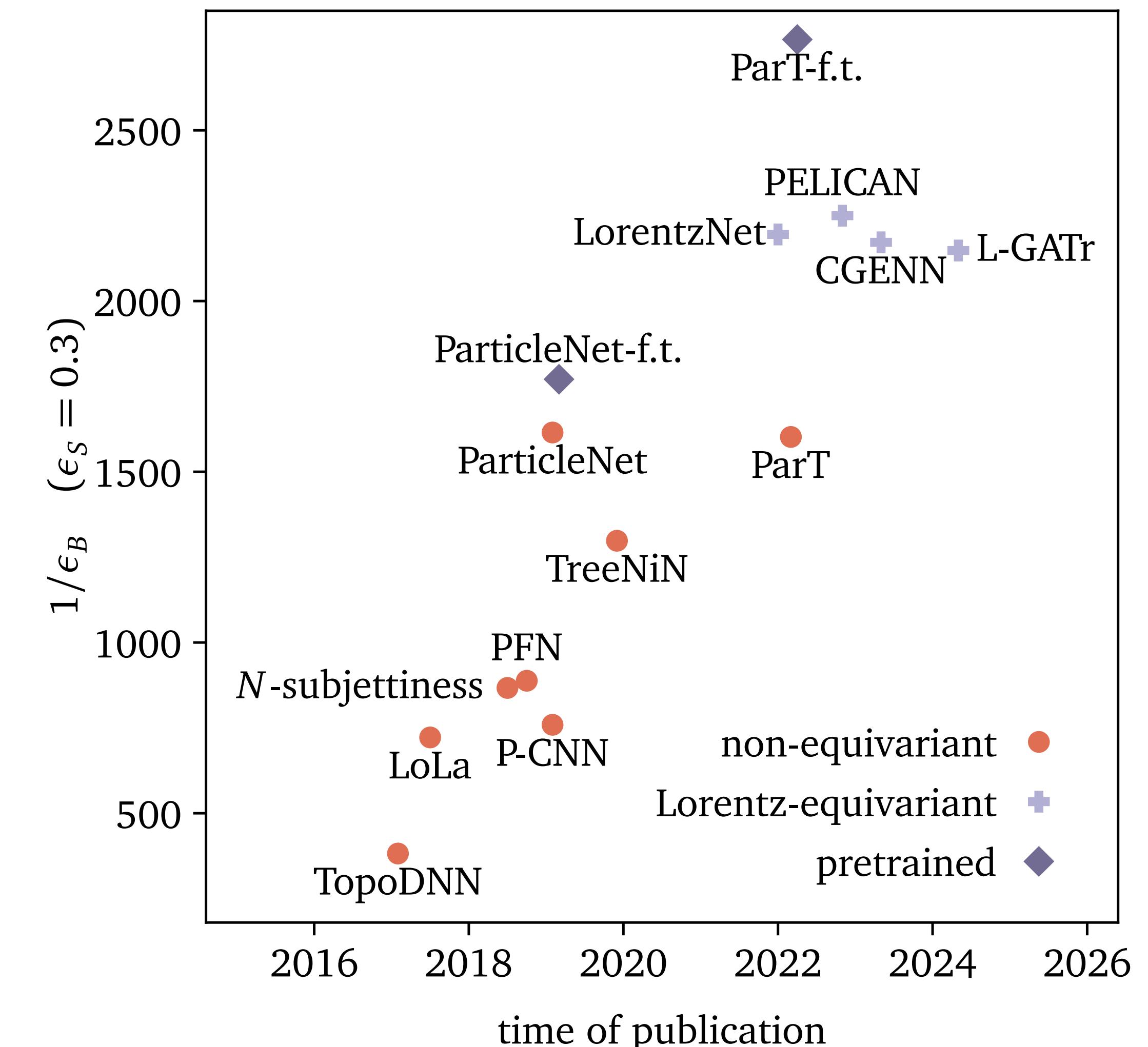
Top tagging

Model	Accuracy	AUC	$1/\epsilon_B$ ($\epsilon_S = 0.5$)	$1/\epsilon_B$ ($\epsilon_S = 0.3$)
TopoDNN [48]	0.916	0.972	–	295 \pm 5
LoLa [15]	0.929	0.980	–	722 \pm 17
P-CNN [1]	0.930	0.9803	201 \pm 4	759 \pm 24
N -subjettiness [61]	0.929	0.981	–	867 \pm 15
PFN [50]	0.932	0.9819	247 \pm 3	888 \pm 17
TreeNiN [57]	0.933	0.982	–	1025 \pm 11
ParticleNet [63]	0.940	0.9858	397 \pm 7	1615 \pm 93
ParT [64]	0.940	0.9858	413 \pm 16	1602 \pm 81
LorentzNet* [41]	0.942	0.9868	498 \pm 18	2195 \pm 173
CGENN* [67]	0.942	0.9869	500	2172
PELICAN* [9]	0.9426 \pm 0.0002	0.9870 \pm 0.0001	–	2250 \pm 75
L-GATr (ours)*	0.9417 \pm 0.0002	0.9868 \pm 0.0001	548 \pm 26	2148 \pm 106

Experiments

Top tagging

- New paradigm: **Transfer learning**
Pretrain model on large dataset, then fine-tune on target dataset
- Transformers transfer better than graph networks



Experiments

Conditional Flow Matching

Continuous normalising flows (CNF)

connect a simple base density
to a complex target density
through a neural differential equation

$$\frac{d}{dt}x = v_t(x)$$

Conditional flow matching (CFM)

is a simple way to train CNFs
by comparing the learned velocity $v_t(x)$
to a conditional target velocity $u_t(x | x_1)$

$$\mathcal{L} = \mathbb{E}_{t,x,x_1} \|v_t(x) - u_t(x | x_1)\|^2$$

Continuous normalising flows
arXiv:1806.07366

Conditional flow matching
arXiv:2210.02747

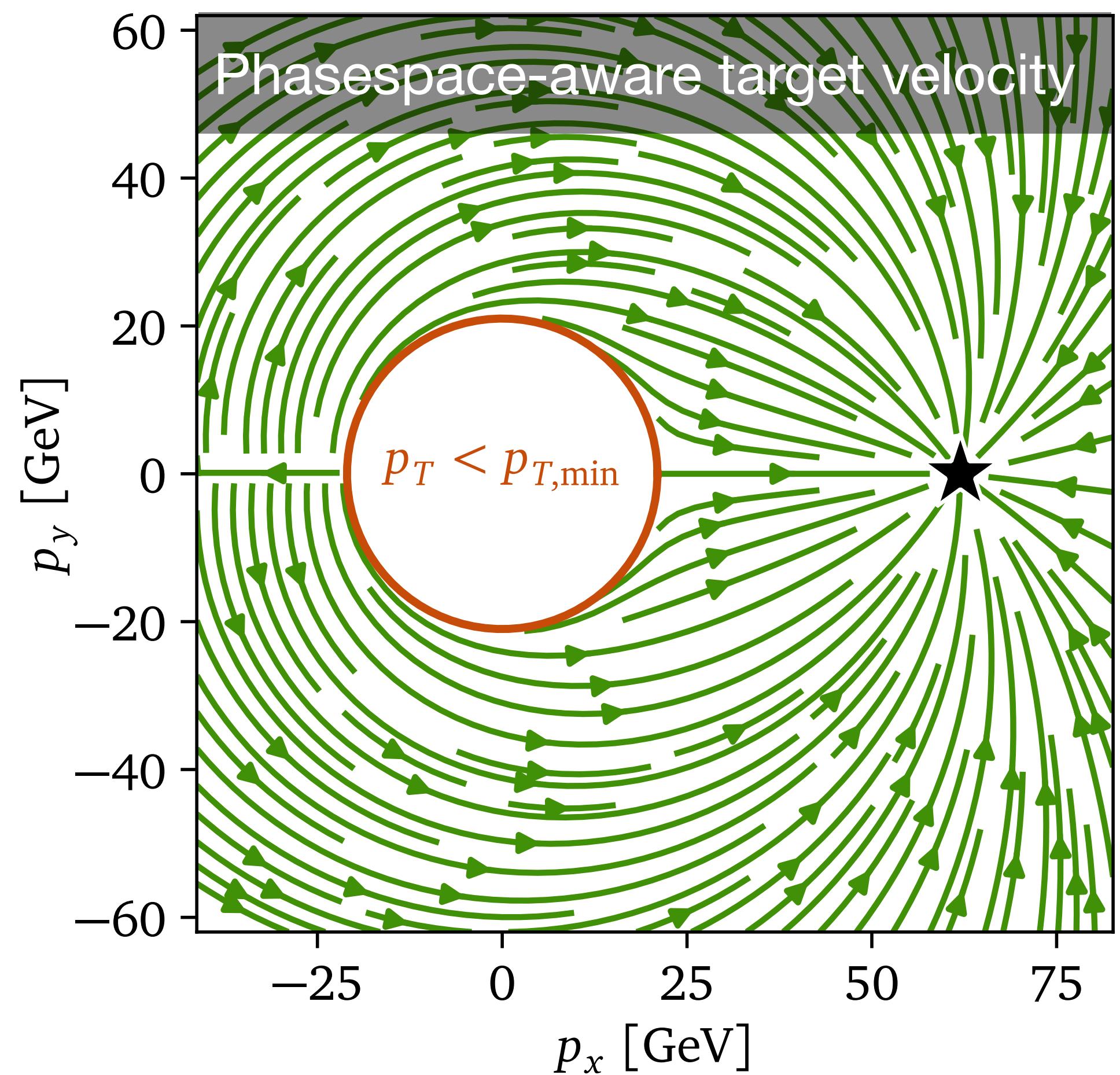
Experiments

Target velocities for CFM

In conditional flow matching (CFM),
the **choice of target velocity** can be
more important than the architecture

Target velocity	Architecture	AUC
Euclidean	L-GATr	0.99
Phasespace-aware	MLP	0.78
Phasespace-aware	L-GATr	0.51

Riemannian Flow Matching
arXiv:2302.03660



Event generation

Target velocities for CFM

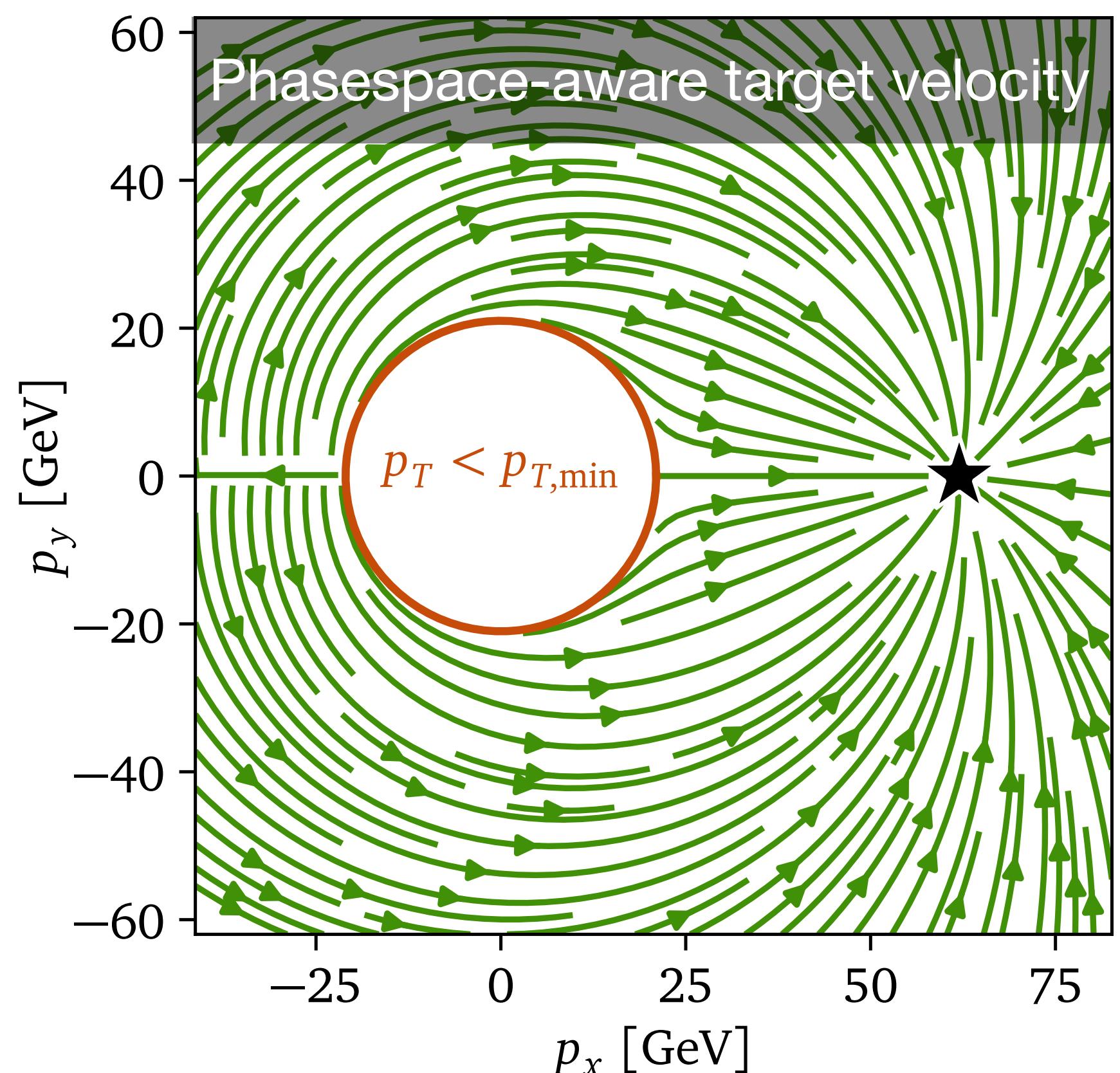
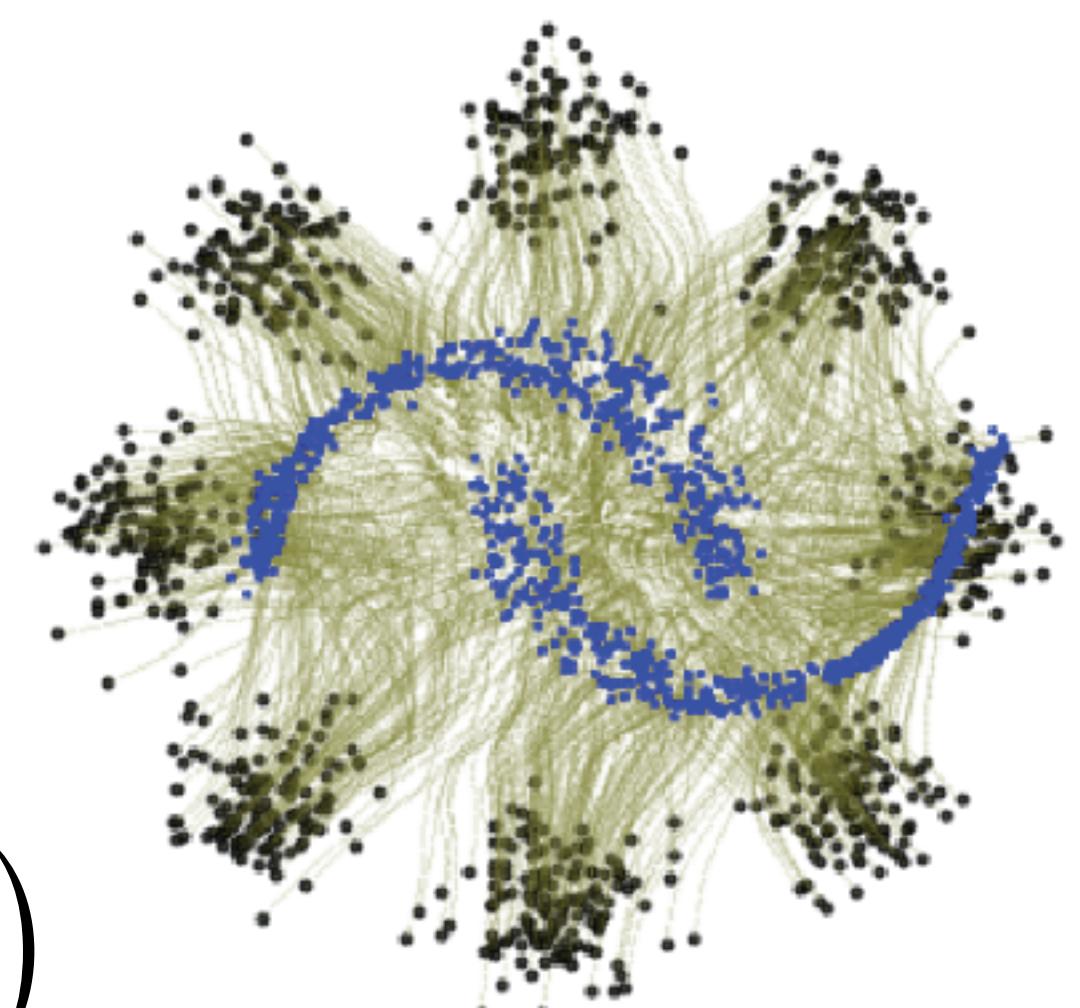
$$p = (E, p_x, p_y, p_z) = f(y) = \left(\sqrt{m^2 + p_T^2 \cosh^2 \eta}, p_T \cos \phi, p_T \sin \phi, p_T \sinh \eta \right)$$

$$y = (y_m, y_p, \phi, \eta), \quad m^2 = \exp(y_m), \quad p_T = p_{T,\min} + \exp(y_p)$$

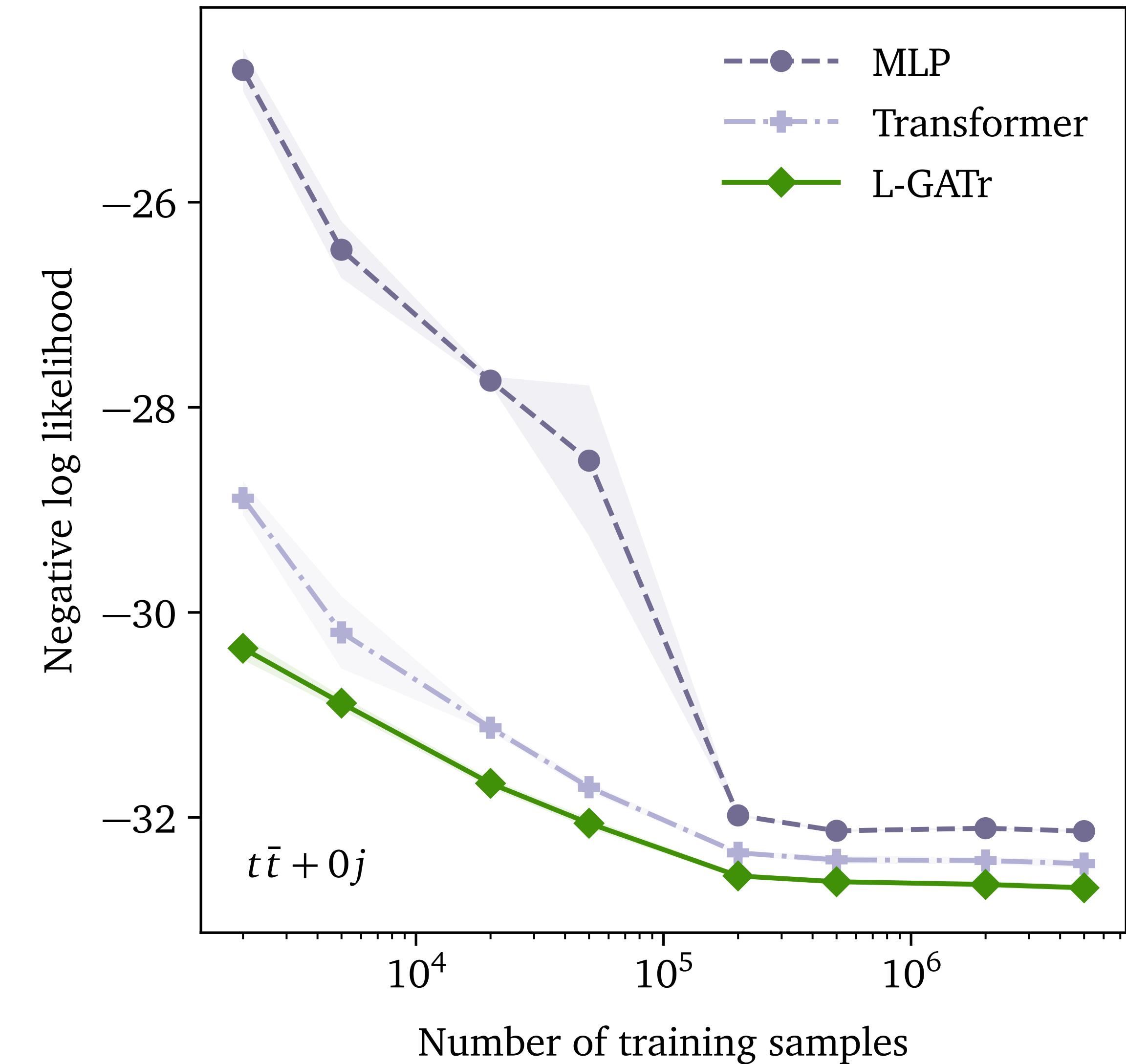
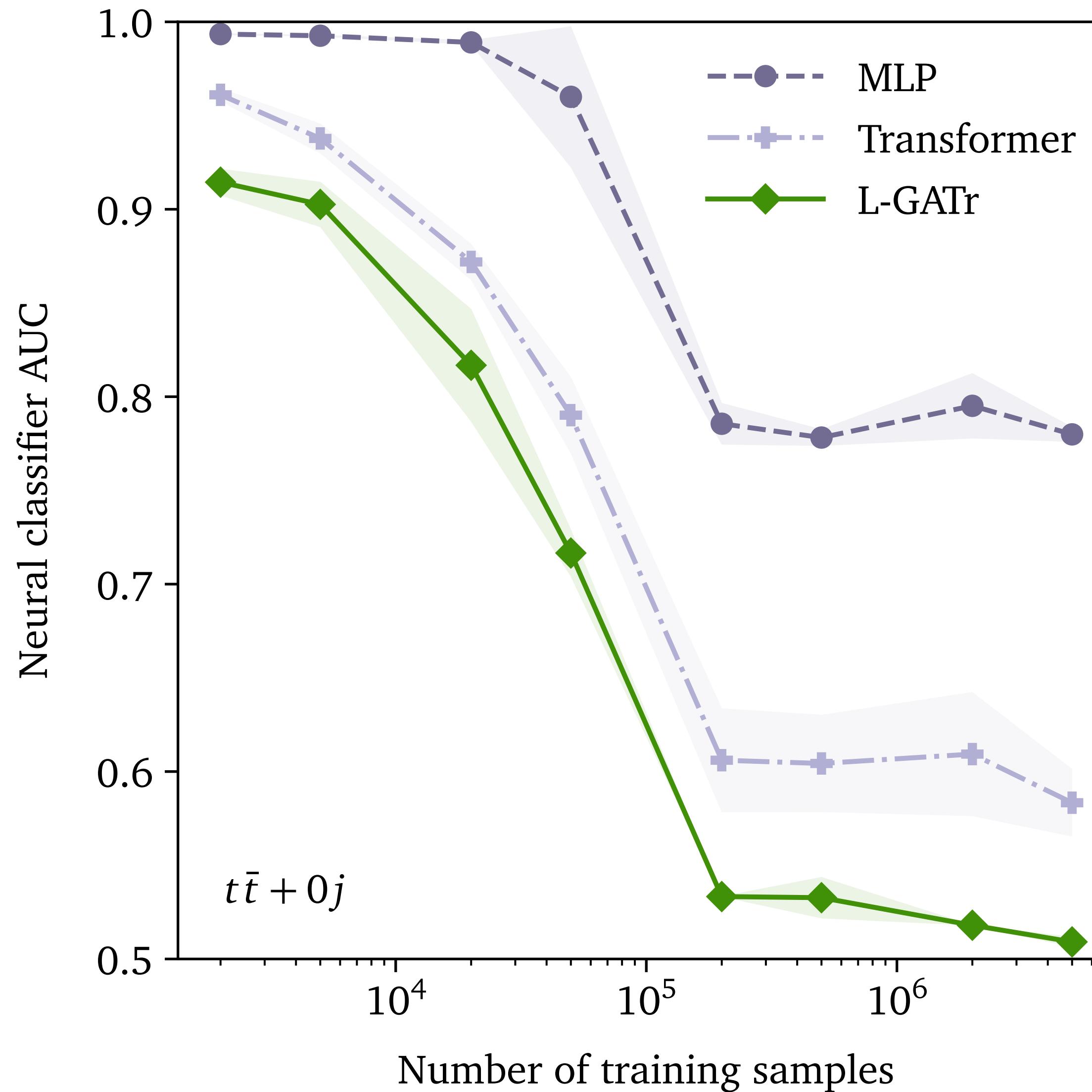
Target velocities can be

constant in $p = (E, p_x, p_y, p_z)$ ('euclidean')

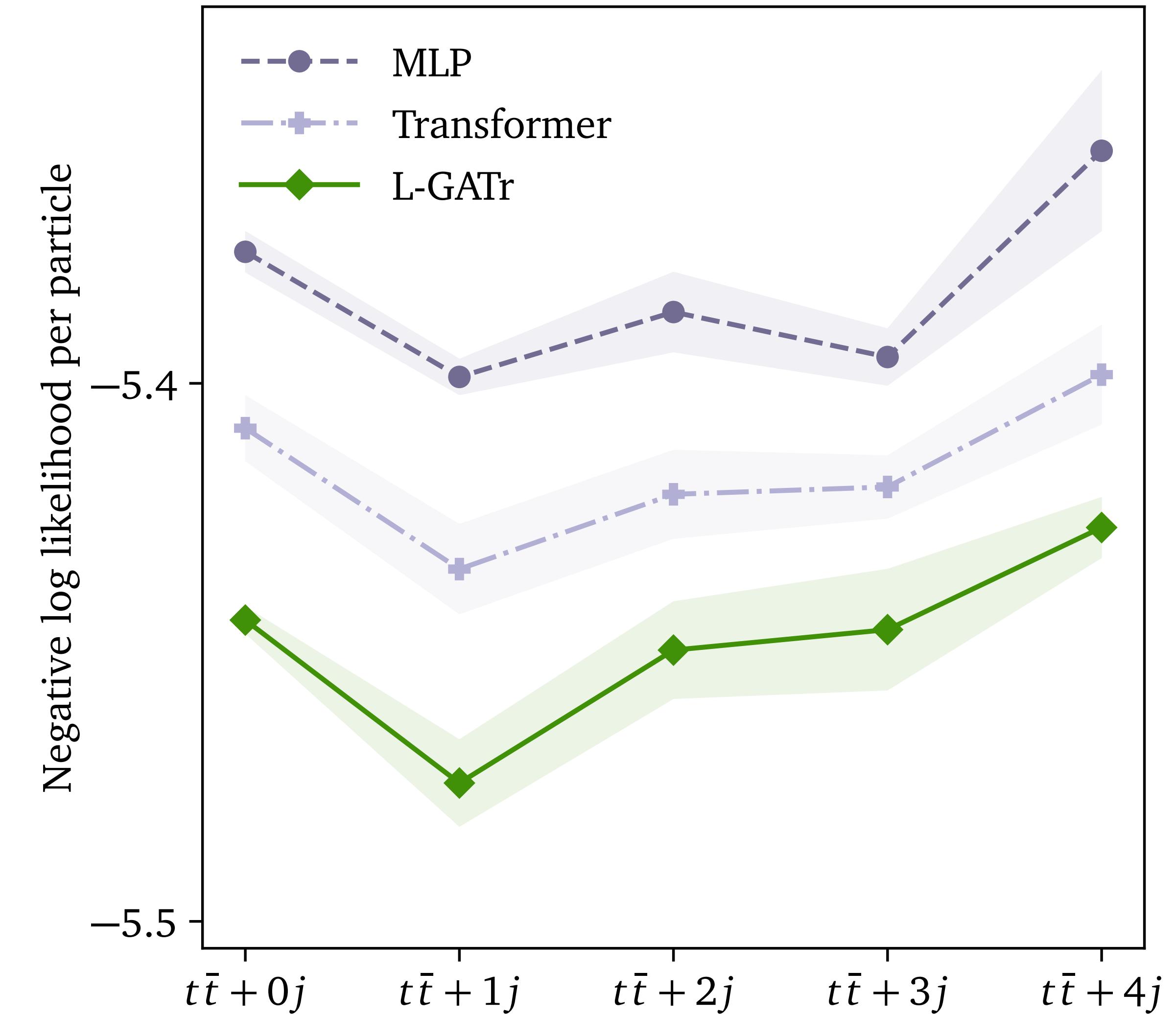
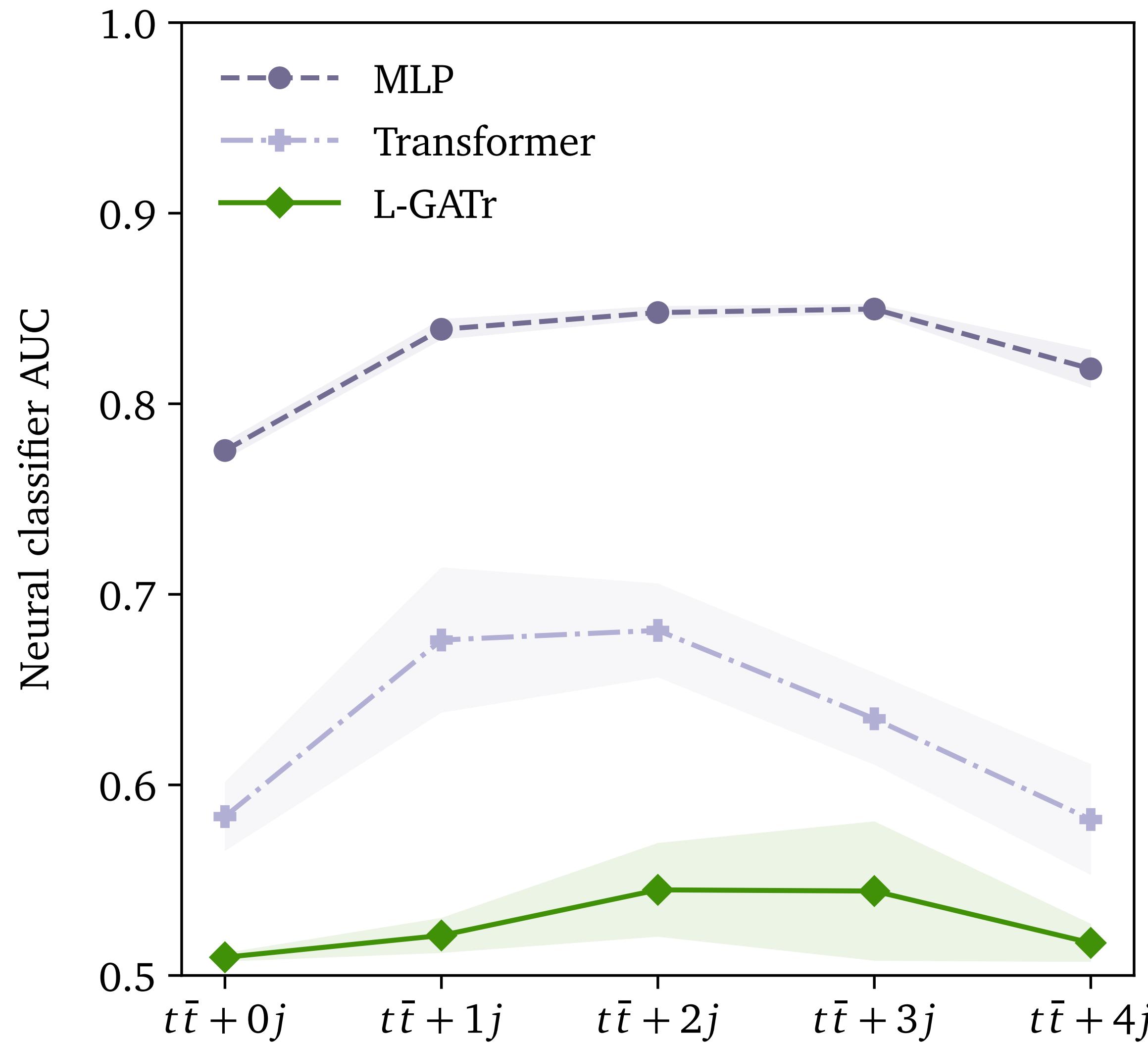
constant in $y = (y_m, y_p, \phi, \eta)$ ('phasespace-aware')



Event generation



Event generation



Symmetry breaking with spurious

Sources of symmetry breaking

- Real world: Beam direction, detector geometry...
Symmetry-breaking object: Beam direction spurion
- Generation: Have to break $SO(1,3) \rightarrow SO(3)$ because generative networks can only be defined on compact groups
Symmetry-breaking object: Time direction spurn

We break the symmetry by adding the spurious as extra token or as extra channel for each token