

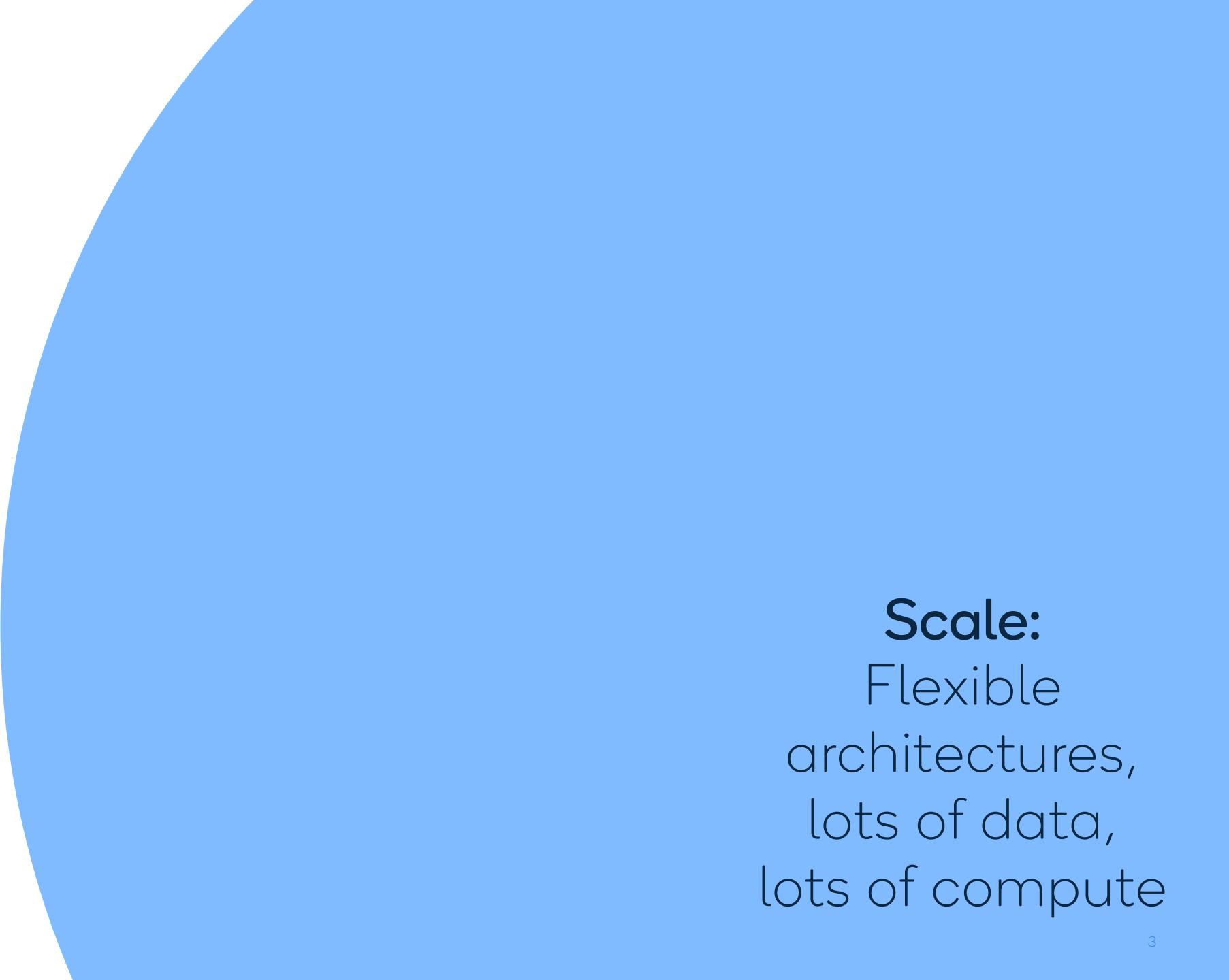
# Scalable equivariance

## with Geometric Algebra Transformers

**Johann Brehmer**  
Qualcomm AI Research

## **Structure:**

Problem-specific  
inductive biases in  
algorithms and  
architectures



**Scale:**  
Flexible  
architectures,  
lots of data,  
lots of compute

## **Structure:**

Problem-specific  
inductive biases in  
algorithms and  
architectures

## **Scale:**

Flexible  
architectures,  
lots of data,  
lots of compute

## **Structure:**

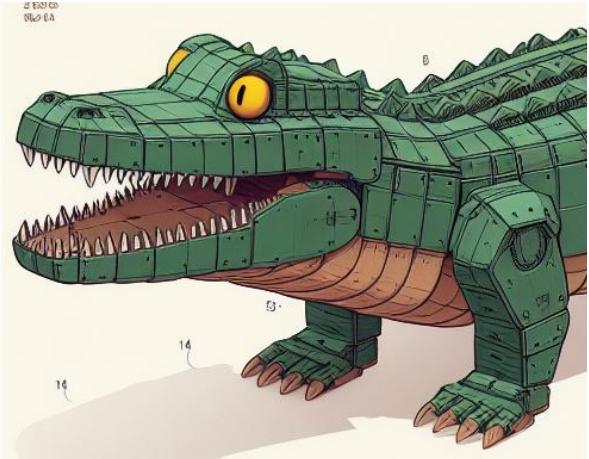
Problem-specific  
inductive biases in  
algorithms and  
architectures

## **Geometric Algebra Transformer:**

Our version of a  
versatile architecture  
for geometric  
problems

## **Scale:**

Flexible  
architectures,  
lots of data,  
lots of compute



**GATr 101**  
How to build an  
equivariant Transformer



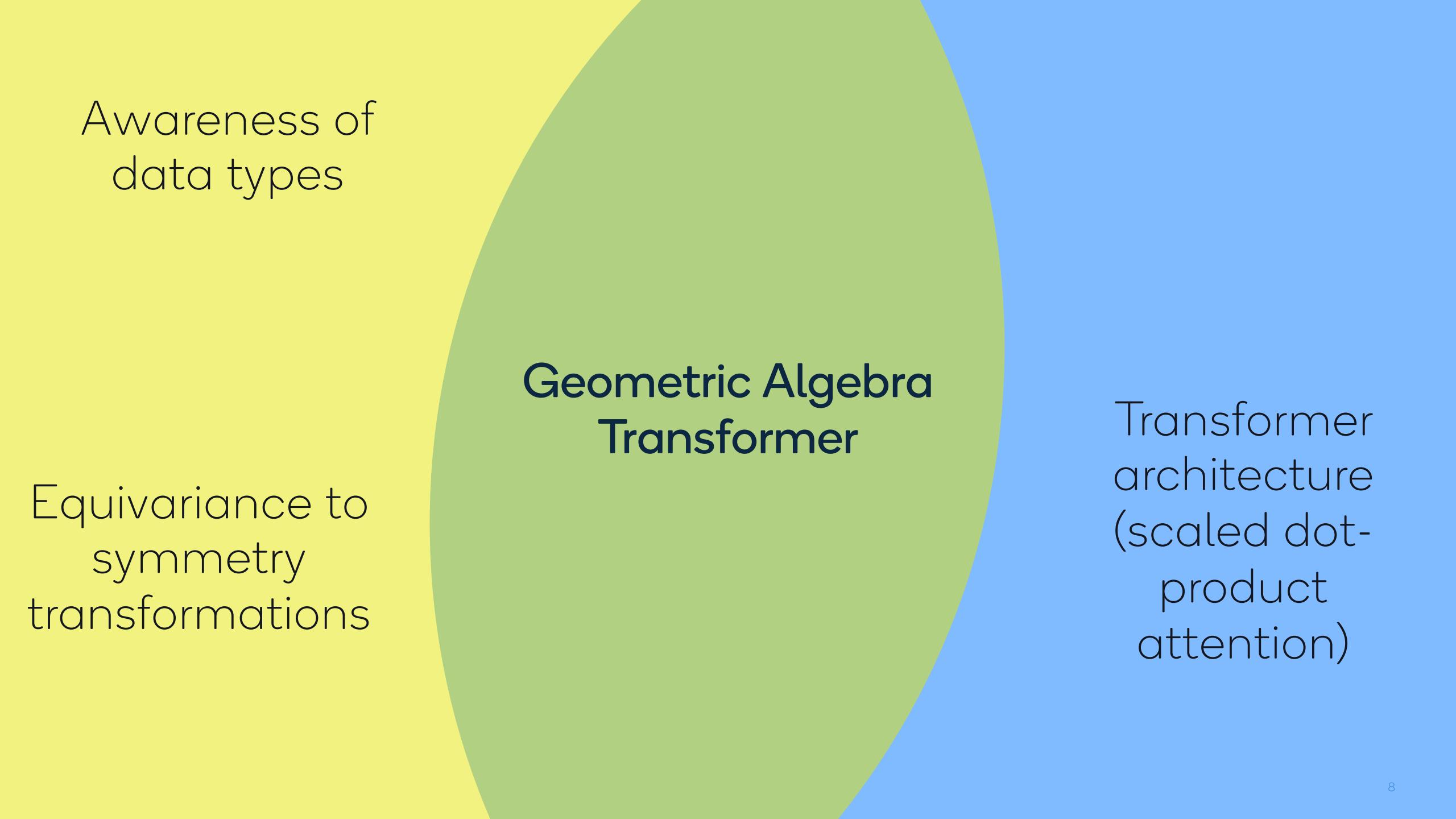
**Euclidean GATr**  
Applications far below  
the speed of light



**Lorentz-GATr**  
for particle physics

# GATr 101





Awareness of  
data types

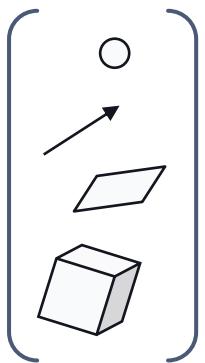
## Geometric Algebra Transformer

Transformer  
architecture  
(scaled dot-  
product  
attention)

Equivariance to  
symmetry  
transformations

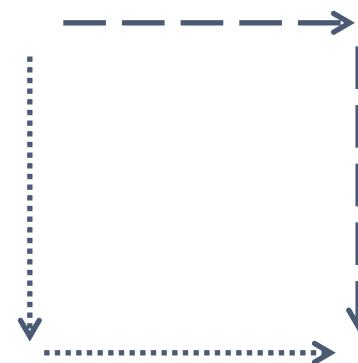
**Geometric Algebra  
Transformer**

=



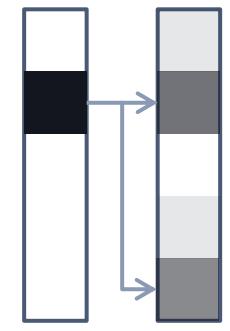
**Geometric algebra  
representations**

+



**Equivariant  
layers**

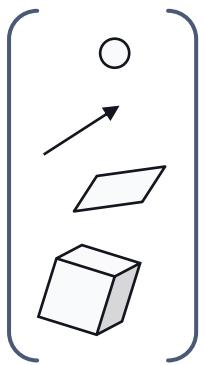
+



**Transformer  
architecture**

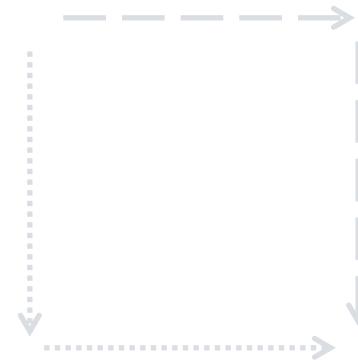
Geometric Algebra  
Transformer

=



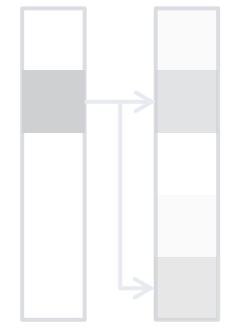
Geometric algebra  
representations

+



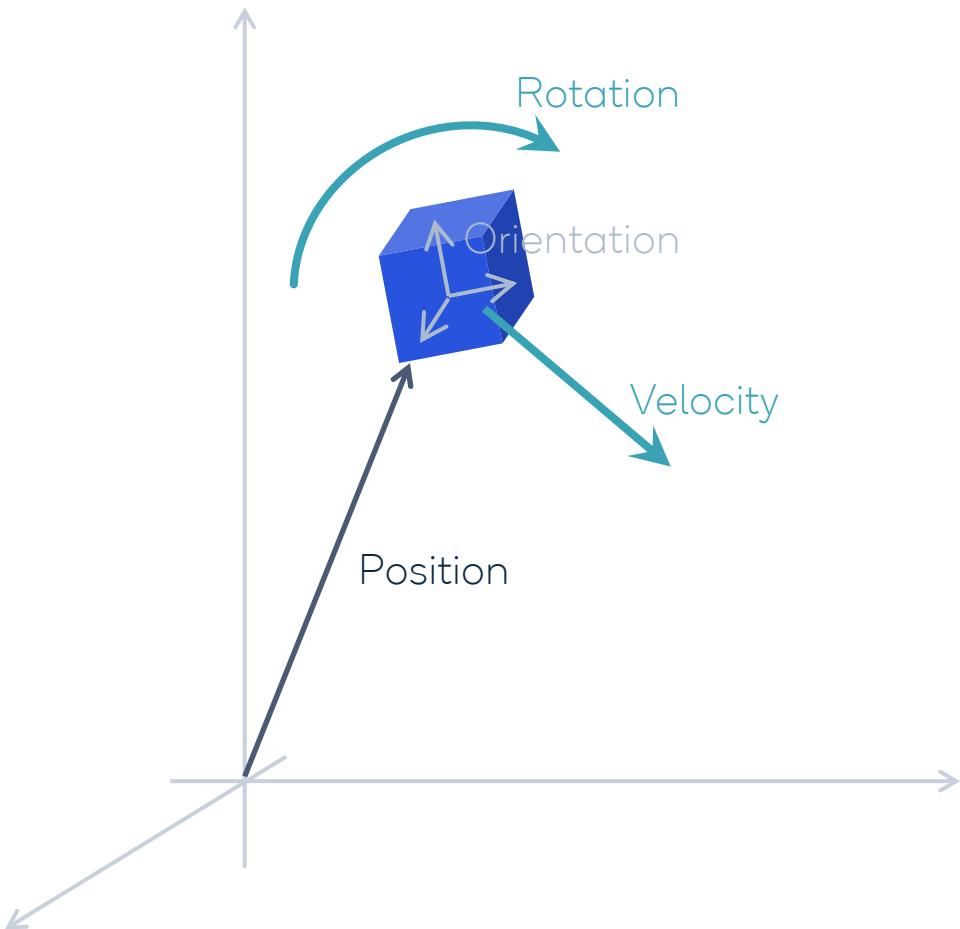
Equivariant  
layers

+



Transformer  
architecture

# Representing geometric data



How do you parameterize such a 3D object?

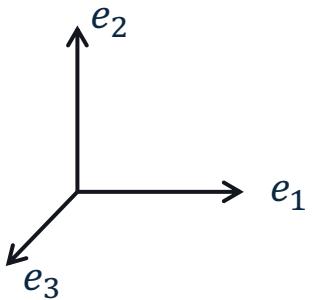
- Most deep learning: 14 numbers
- Previous **geometric deep learning**:  
e.g. 2 vectors, 2 rotation matrices
- **GATr**: 1 position, 3 directions, 1 translation, 1 rotation

Why use different types?

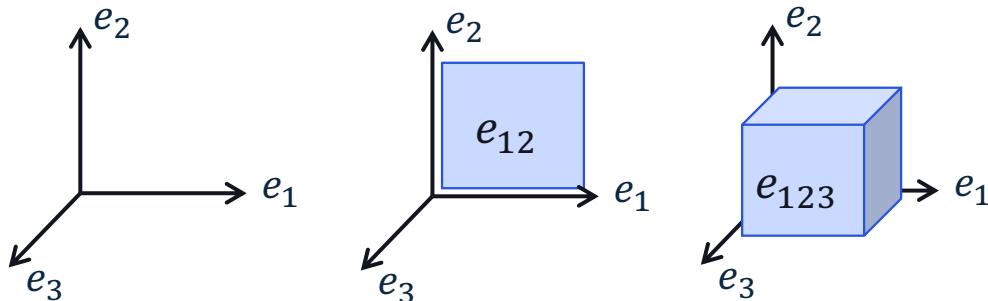
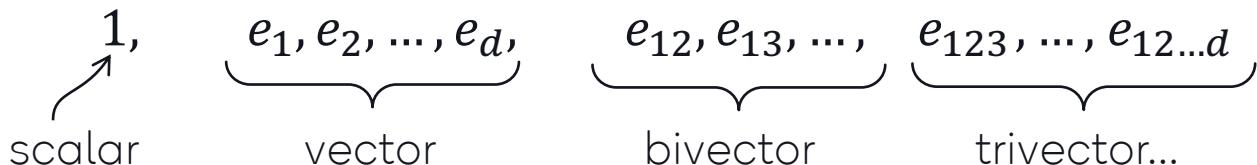
- Types have different **common patterns**  
(compute distances between positions, but not between direction vectors)
- Types differ in behaviour under **transformations**
- Types provide an inductive bias, potentially improving **sample efficiency** and **generalization**

# Geometric algebra

- Vector space  $V$  with  $d$  dimensions and inner products
  - Basis  $e_1, e_2, \dots, e_d$



- Geometric algebra  $\mathcal{G}(V)$  has  $2^d$  dimensions, basis



- Geometric product  $\mathcal{G}(V) \times \mathcal{G}(V) \rightarrow \mathcal{G}(V)$ 
  - Generalizes dot product and cross product



Graßmann



Clifford

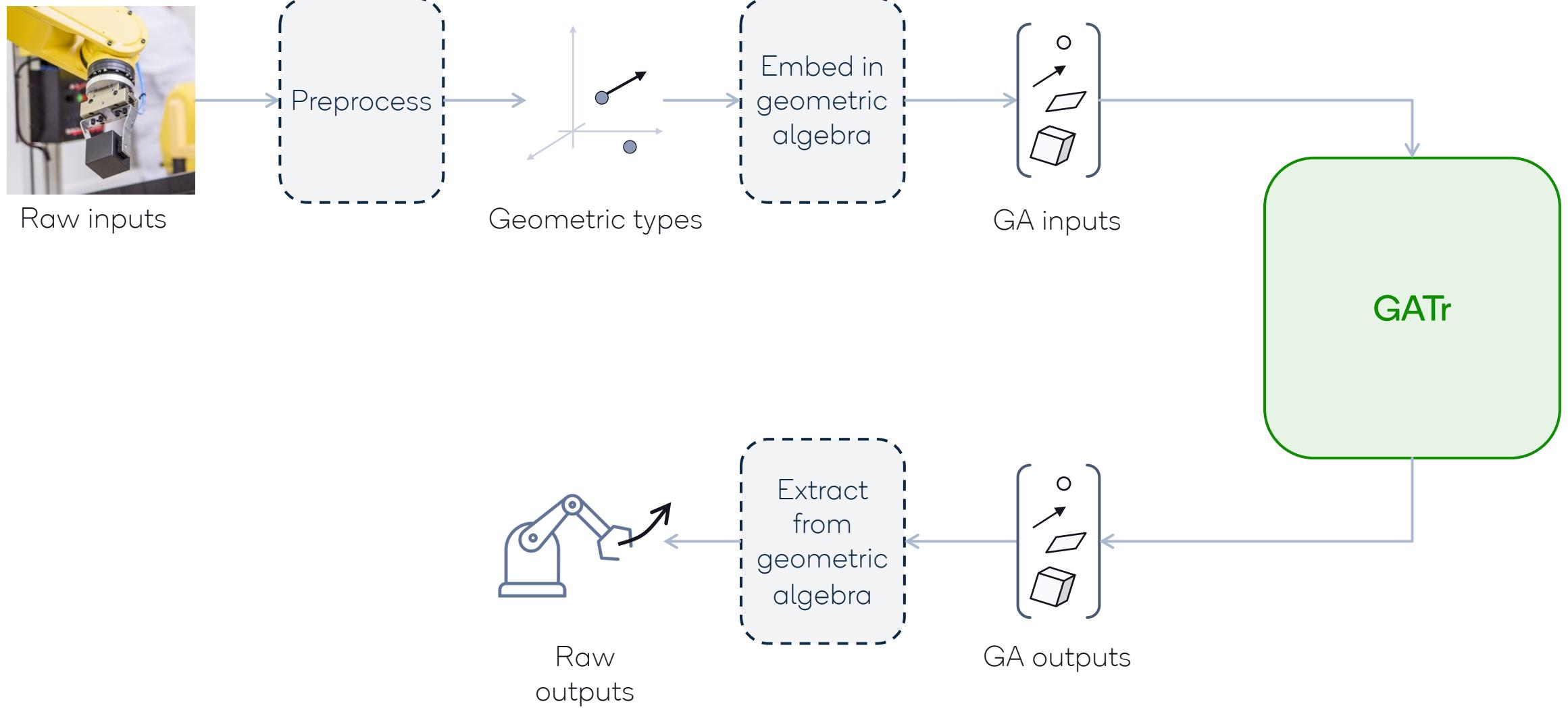
# Projective GA as representation for ML

- We use GA representations in addition to the usual unstructured vector space
- Offers  **$16n$ -dimensional representation** of 3D geometric data
  - “**Typing**”: a point is not a direction of movement is not the orientation of a plane
  - Established embeddings for **3D primitives and transformations**

| Object / operator   | Scalar    | Vector | Bivector | Trivector      | PS       |           |           |            |
|---|-----------|--------|----------|----------------|----------|-----------|-----------|------------|
|   | 1         | $e_0$  | $e_i$    | $e_{0i}$       | $e_{ij}$ | $e_{0ij}$ | $e_{123}$ | $e_{0123}$ |
| Scalar $\lambda \in \mathbb{R}$   | $\lambda$ | 0      | 0        | 0              | 0        | 0         | 0         | 0          |
| Plane w/ normal $n \in \mathbb{R}^3$ , origin shift $d \in \mathbb{R}$                    | 0         | $d$    | $n$      | 0              | 0        | 0         | 0         | 0          |
| Line w/ direction $n \in \mathbb{R}^3$ , orthogonal shift $s \in \mathbb{R}^3$            | 0         | 0      | 0        | $s$            | $n$      | 0         | 0         | 0          |
| Point $p \in \mathbb{R}^3$  | 0         | 0      | 0        | 0              | 0        | $p$       | 1         | 0          |
| Pseudoscalar $\mu \in \mathbb{R}$   | 0         | 0      | 0        | 0              | 0        | 0         | 0         | $\mu$      |
| Reflection through plane w/ normal $n \in \mathbb{R}^3$ , origin shift $d \in \mathbb{R}$ | 0         | $d$    | $n$      | 0              | 0        | 0         | 0         | 0          |
| Translation $t \in \mathbb{R}^3$  | 1         | 0      | 0        | $\frac{1}{2}t$ | 0        | 0         | 0         | 0          |
| Rotation expressed as quaternion $q \in \mathbb{R}^4$                                     | $q_0$     | 0      | 0        | 0              | $q_i$    | 0         | 0         | 0          |
| Point reflection through $p \in \mathbb{R}^3$   | 0         | 0      | 0        | 0              | 0        | $p$       | 1         | 0          |

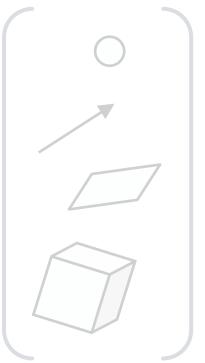
- Geometric product: **canonical operation** on these representations

# GA representations in practice



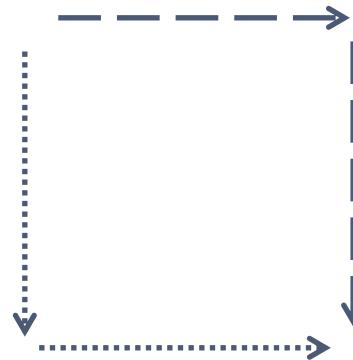
Geometric Algebra  
Transformer

=



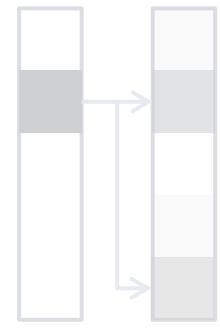
Geometric algebra  
representations

+



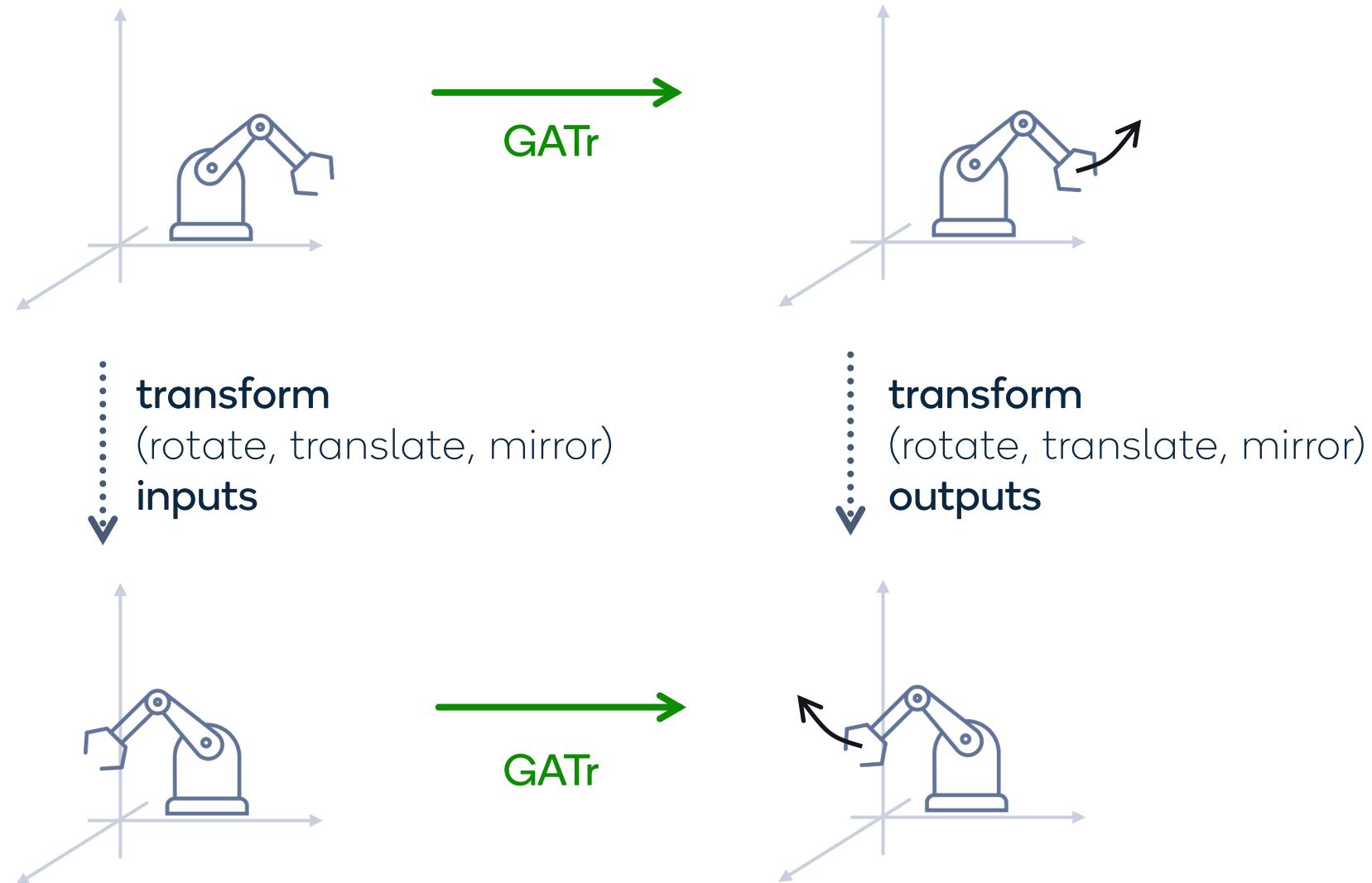
Equivariant  
layers

+



Transformer  
architecture

# $E(3)$ equivariance



# GATr: built out of new **E(3)-equivariant layers** between GAs

- We theoretically characterize how equivariance constrains **linear layers**:

**Proposition 1.** *Any linear map  $\phi : \mathbb{G}_{d,0,1} \rightarrow \mathbb{G}_{d,0,1}$  that is equivariant to  $\text{Pin}(d, 0, 1)$  is of the form*

$$\phi(x) = \sum_{k=0}^{d+1} w_k \langle x \rangle_k + \sum_{k=0}^d v_k e_0 \langle x \rangle_k \quad (4)$$

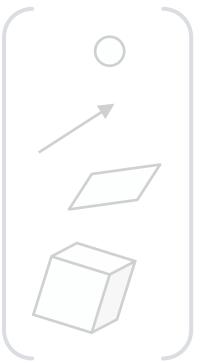
*for parameters  $w \in \mathbb{R}^{d+2}, v \in \mathbb{R}^{d+1}$ . Here  $\langle x \rangle_k$  is the blade projection of a multivector, which sets all non-grade- $k$  elements to zero.*

essentially E(d)

- Plus: equivariant attention, nonlinearities, normalization...

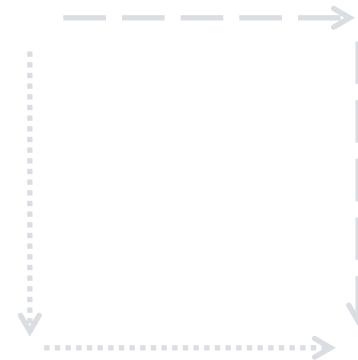
Geometric **A**lgebra  
**T**ransformer

=



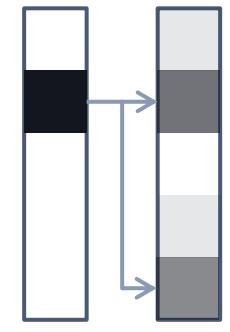
Geometric algebra  
representations

+



Equivariant  
layers

+



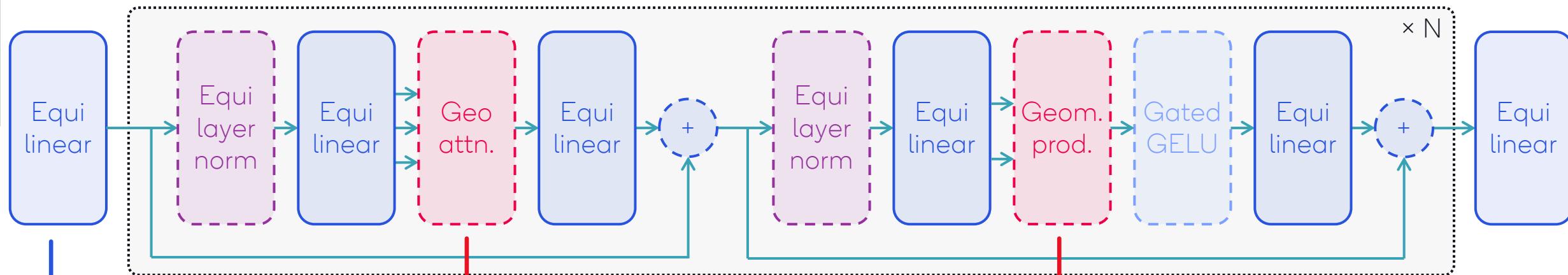
**Transformer**  
architecture

## 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

# GATr is powered by **dot-product attention**

Message-passing neural networks

$$m_{i \rightarrow j} = \Phi(x_i, x_j, e_{ij})$$

$$x'_j = \Psi\left(x_j, \sum_i m_{i \rightarrow j}\right)$$

GATr (and other transformers)

$$V_i, K_i, Q_i = \Phi(x_i)$$

$$x' = \sigma\left(QK^T / \sqrt{d}\right)V$$

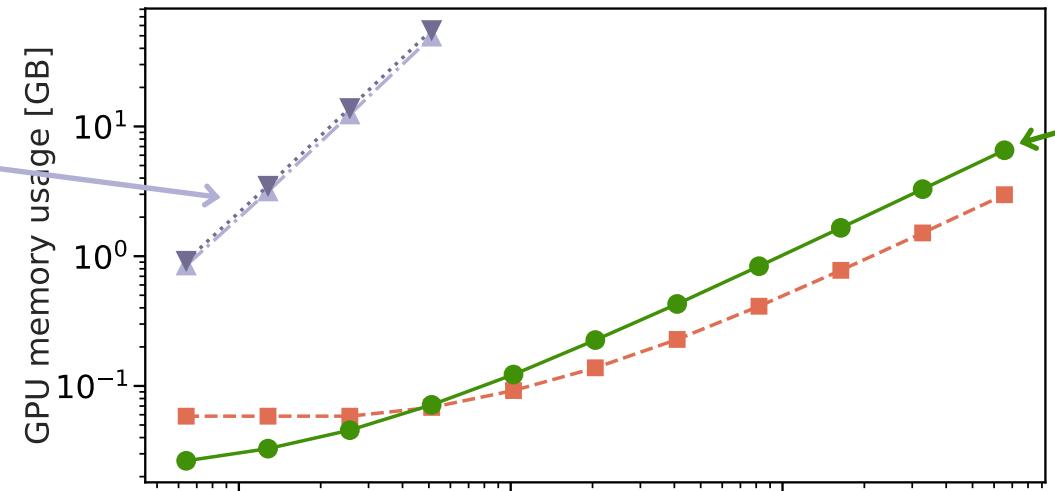
- Node and edge representations
- (Equivariant) network on each edge

- Only node representations
- Dot-product per edge,  
with highly optimized implementations  
(e.g. flash attention)

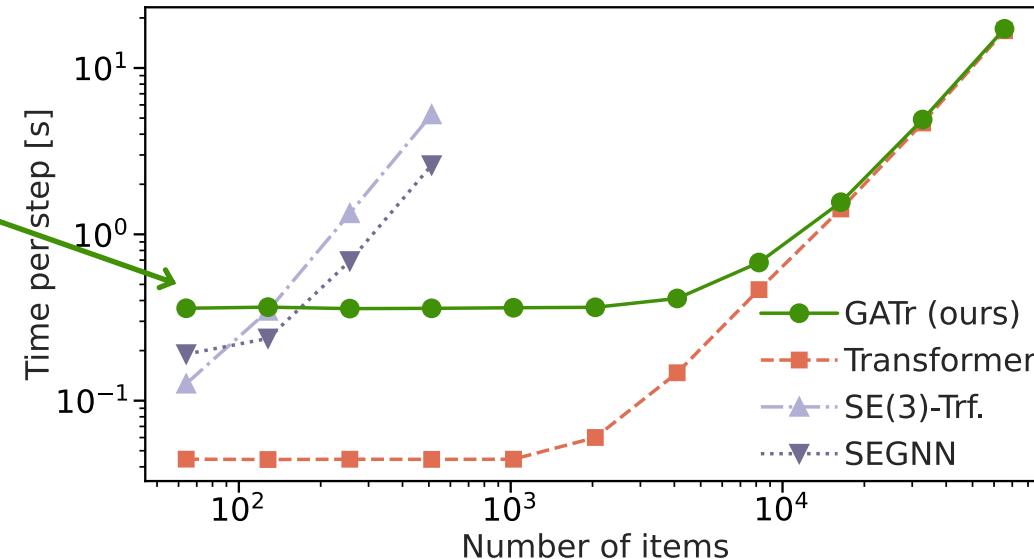
Transformers have same theoretical complexity, but dramatically more efficient in practice

# GATr is more scalable than GDL baselines

Heavily optimized  
Nvidia implementation  
of a classic GDL  
baseline



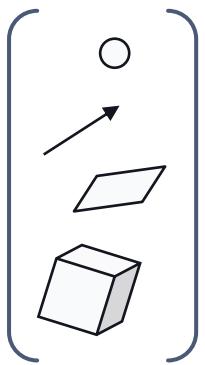
Overhead in small  
problems, but still  
room for optimization



GATr (with flash  
attention) scales like a  
transformer, to 10ks  
tokens!

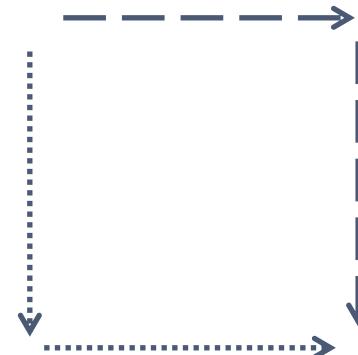
**Geometric Algebra  
Transformer**

=



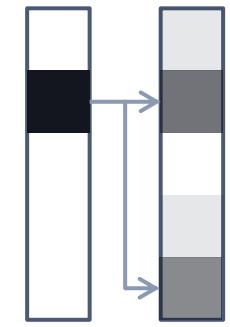
**Geometric algebra  
representations**

+



**Equivariant  
layers**

+

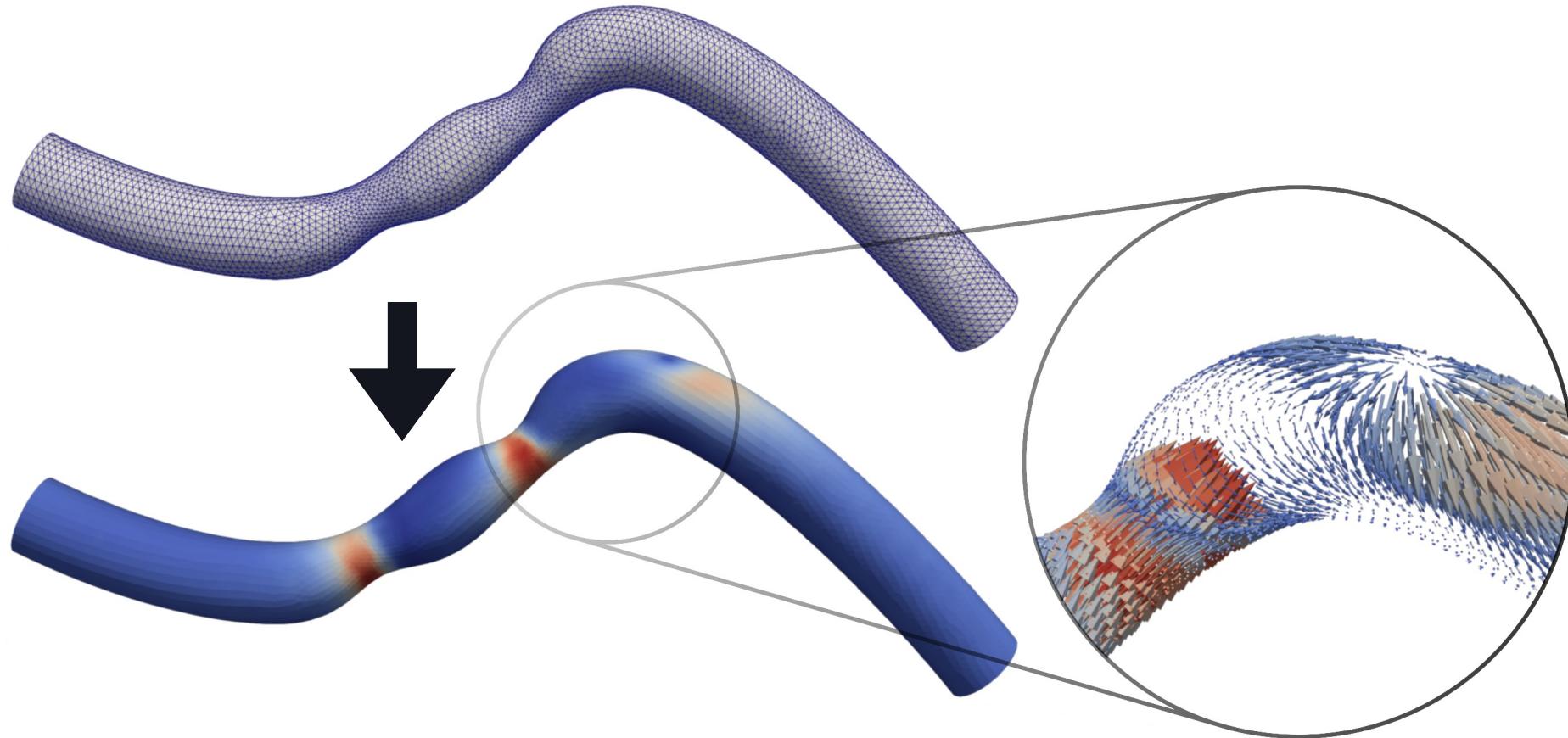


**Transformer  
architecture**

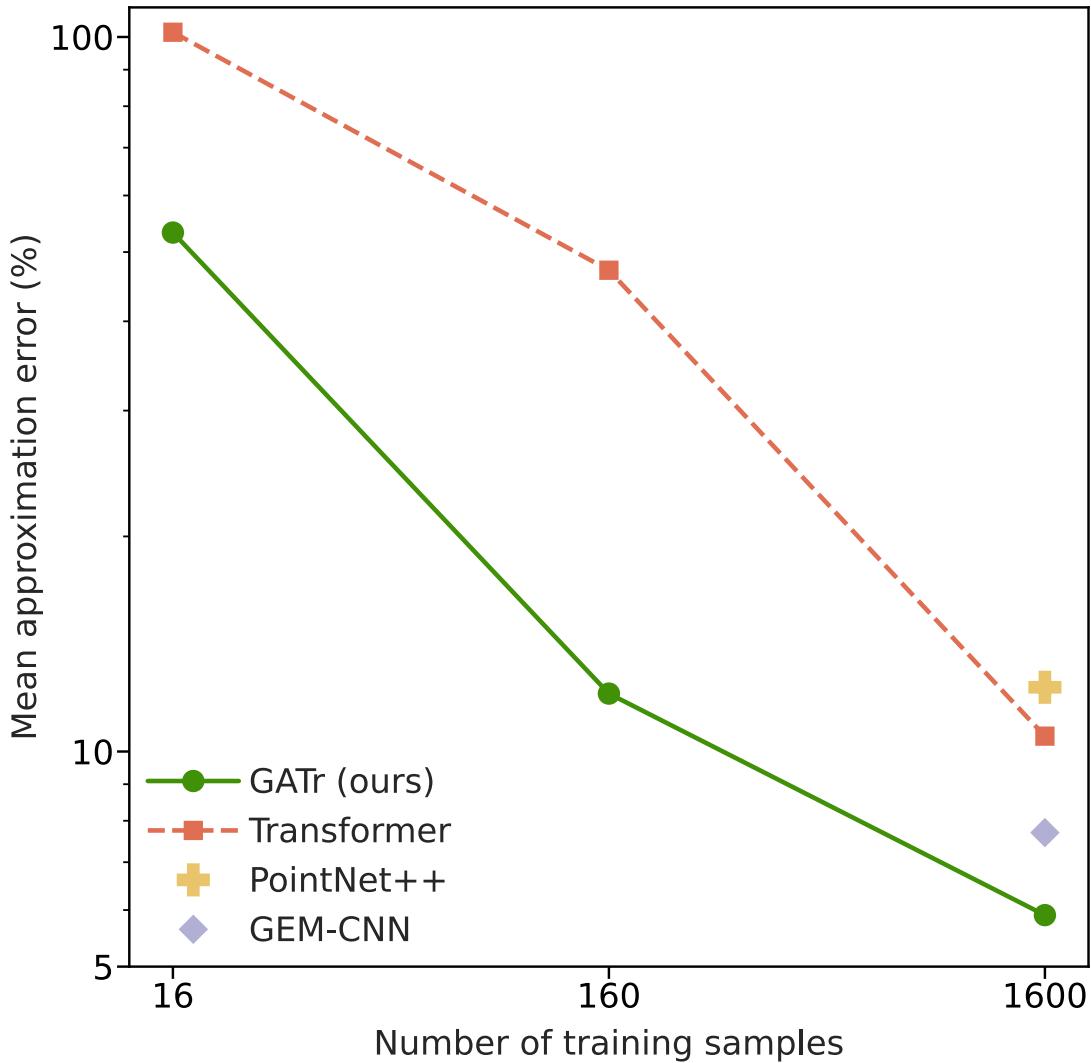
# Euclidean GATr



# Arterial wall-shear stress estimation with 7k tokens

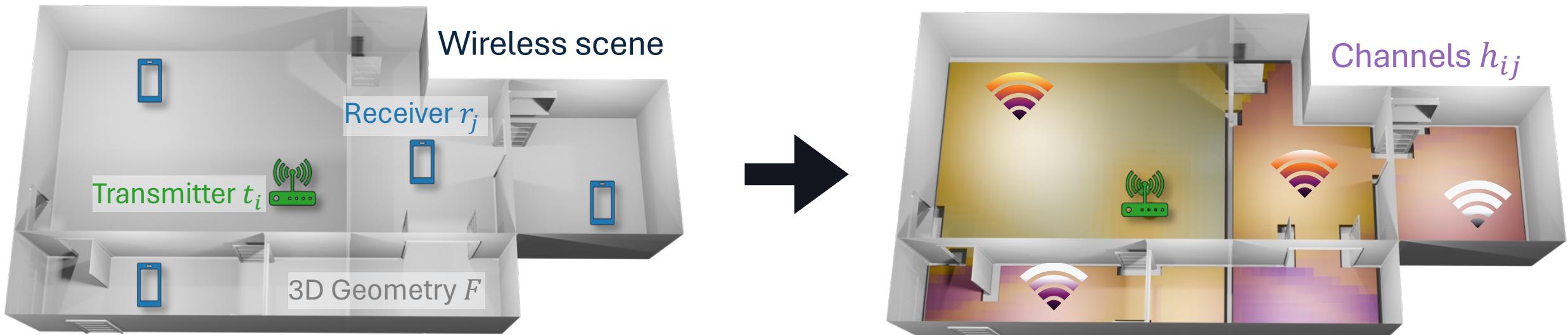


# Arterial wall-shear stress estimation with 7k tokens



**GATr** works well on low-data,  
high-complexity problems

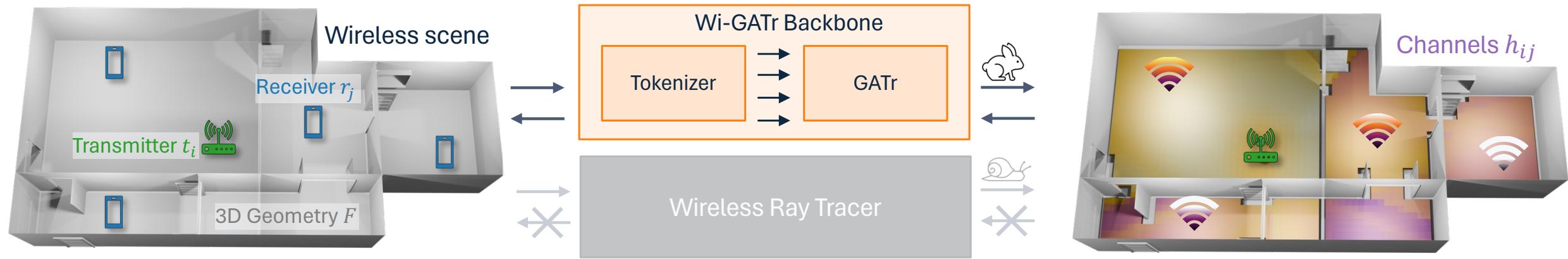
# Differentiable Wireless signal modelling



# Differentiable Wireless signal modelling

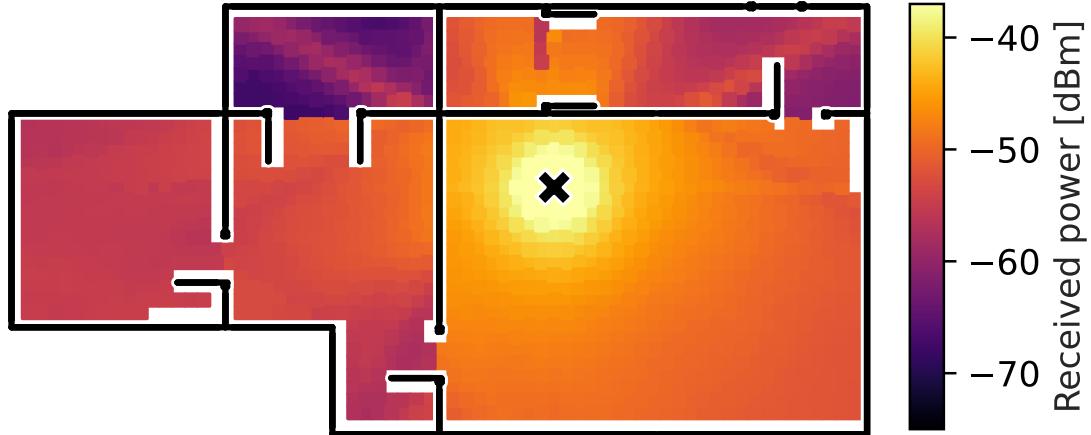


# Differentiable Wireless signal modelling

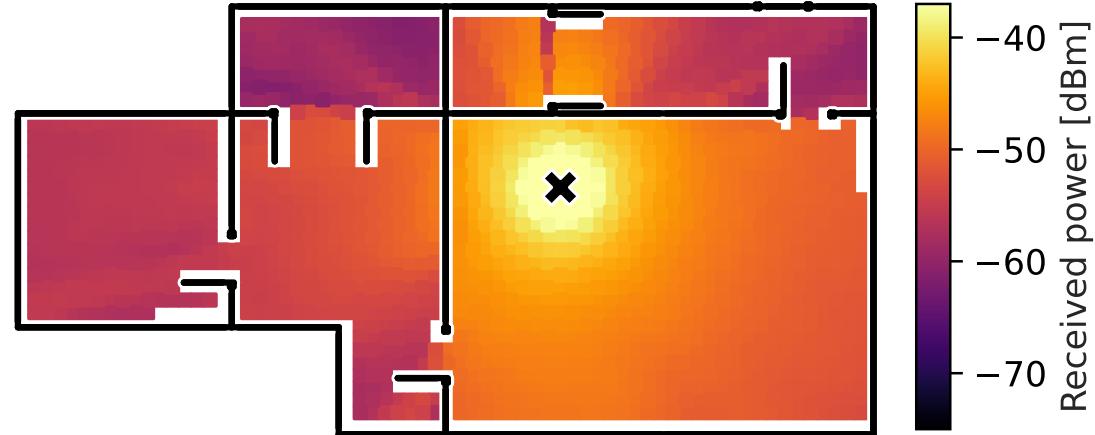


# Differentiable Wireless signal modelling from only 100 training samples

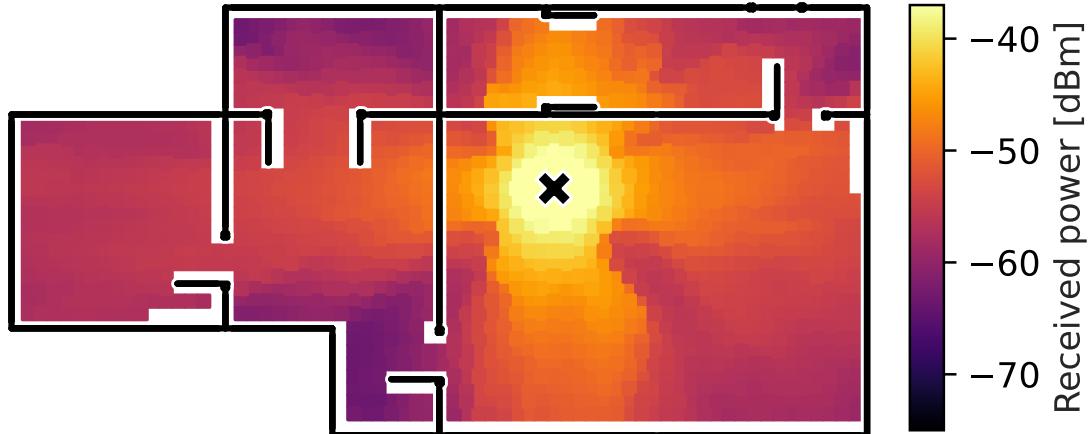
Ground truth



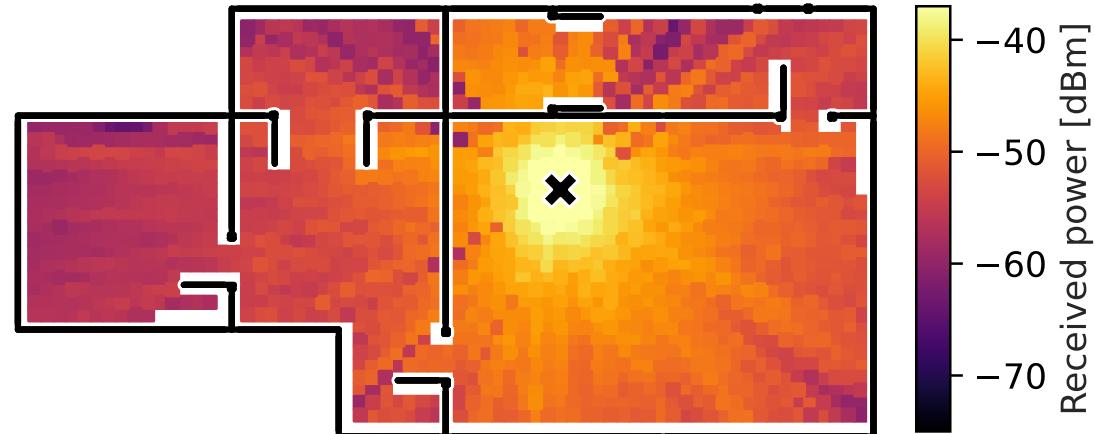
Wi-GATr



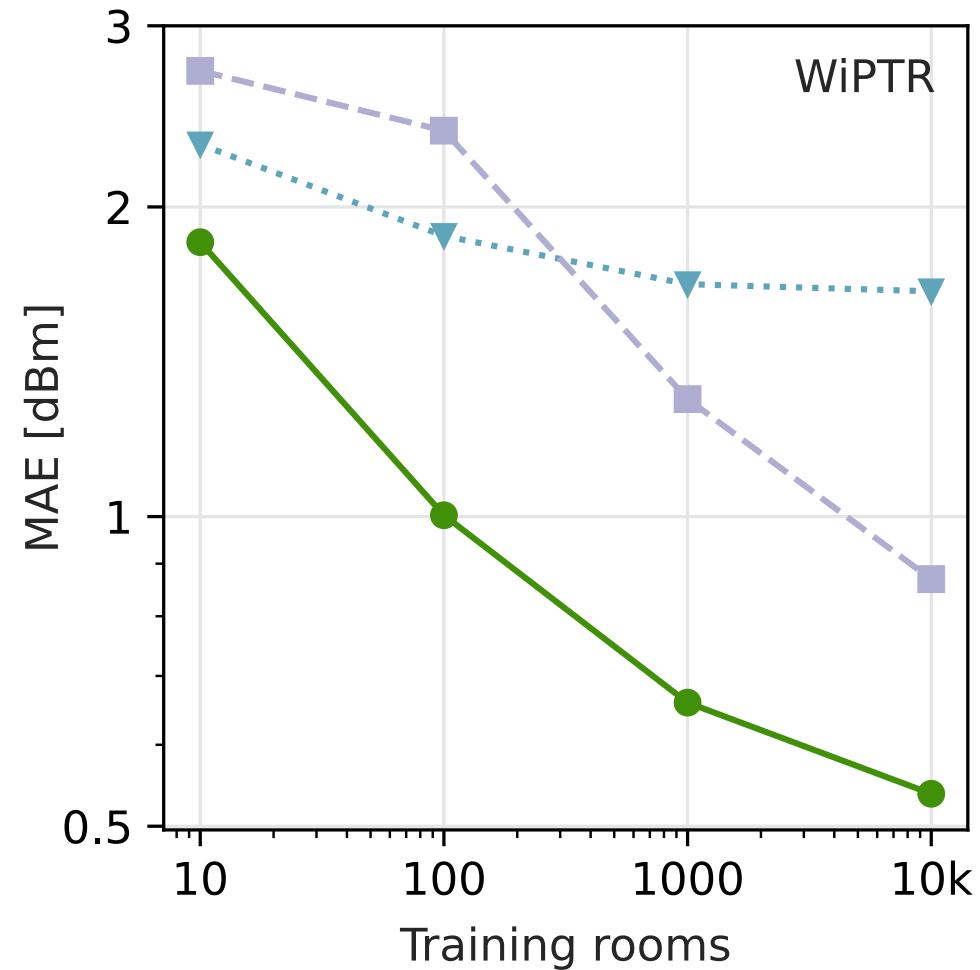
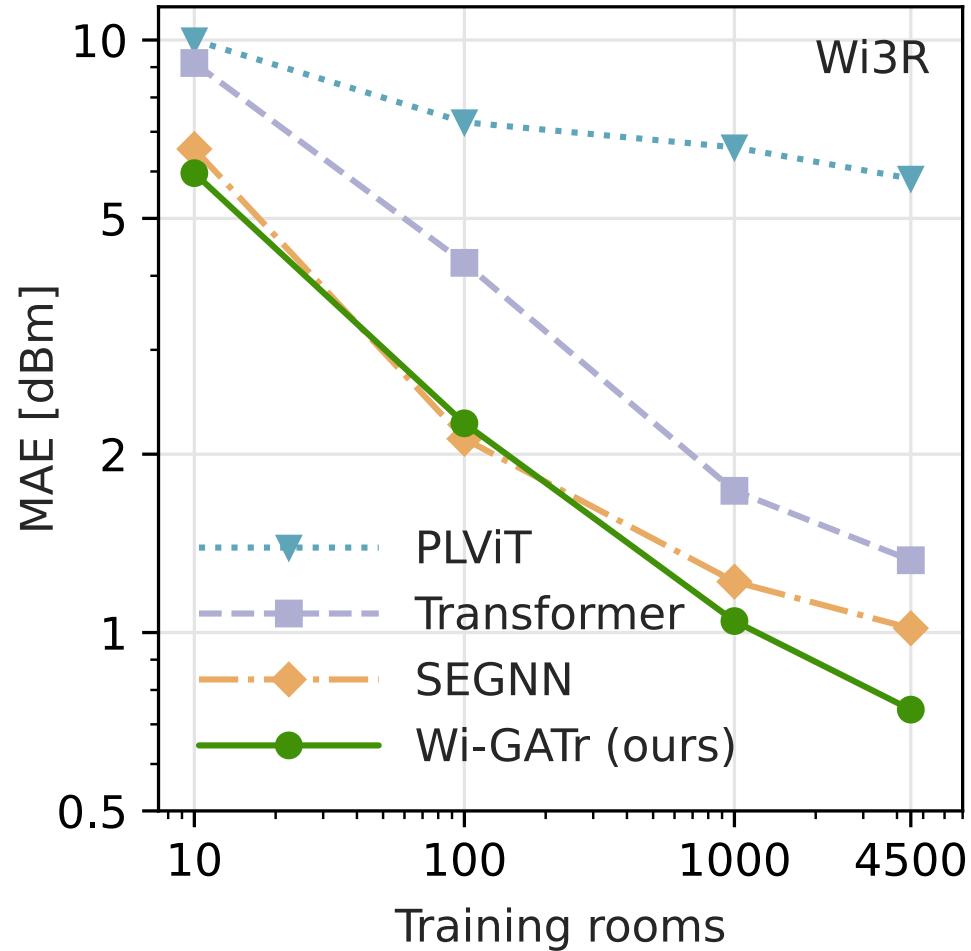
Transformer



ViT

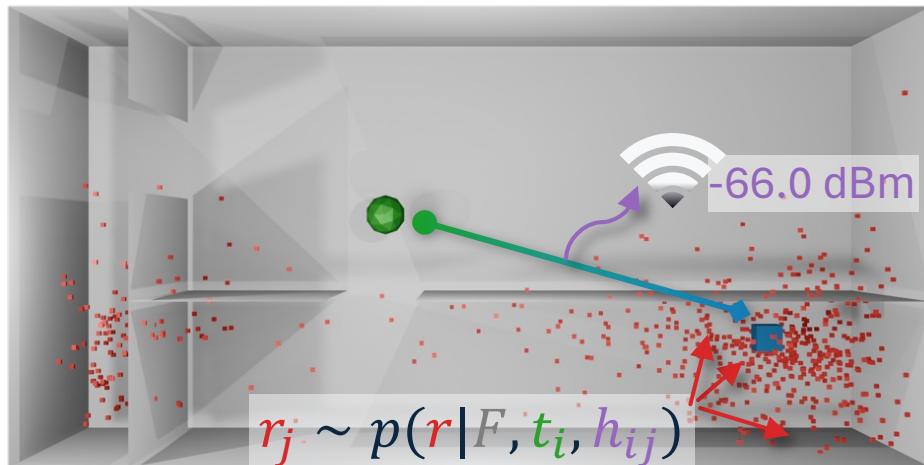


# Differentiable Wireless signal modelling

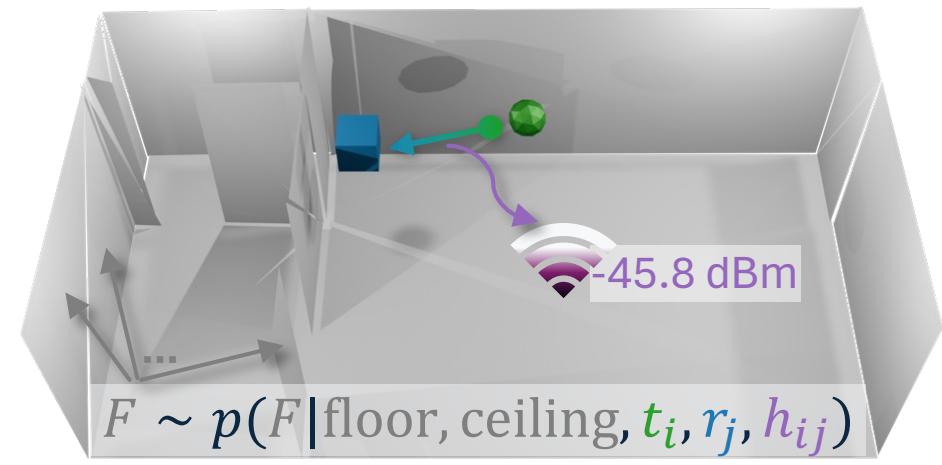


# Next level: Probabilistic wireless modelling with diffusion models

**GATr-based diffusion model** lets us solve various inference tasks through conditional sampling from a single joint density (a la “inpainting”)



Receiver localization  
from environment and  
wireless signal



Geometry reconstruction  
from wireless signal

# Lorentz- GATr



$$\frac{dL}{dA}$$

$$\frac{P_{EG}}{4\pi}$$

$$Na$$

$$\delta \vartheta = \frac{\pi}{2} \quad \epsilon \rightarrow 0$$

$$= C N a \left[ -\frac{\gamma}{2} + \gamma \right]$$

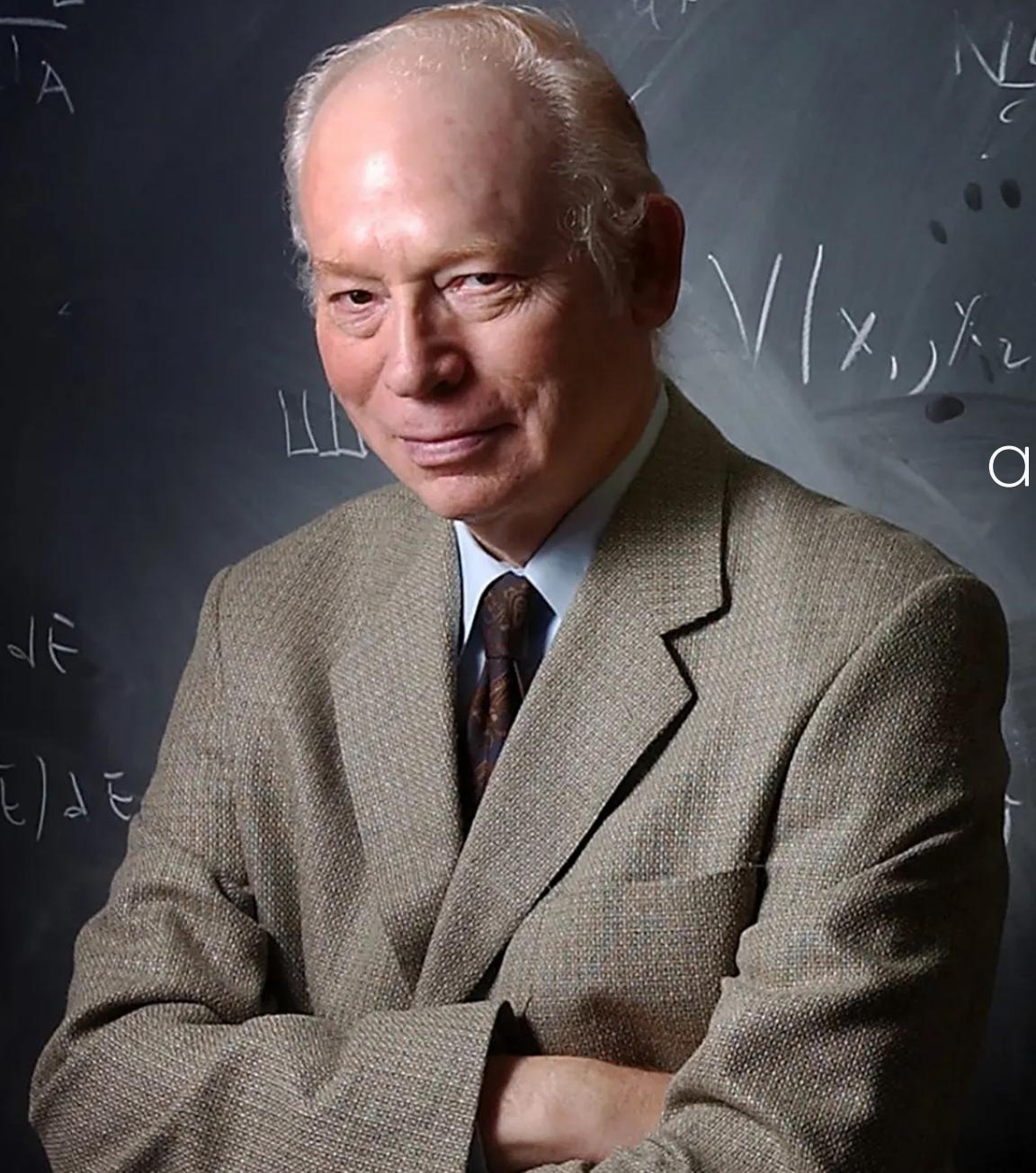
$$\nabla(x_1, x_2)$$

"The universe is  
an enormous direct product  
of representations  
of symmetry groups"

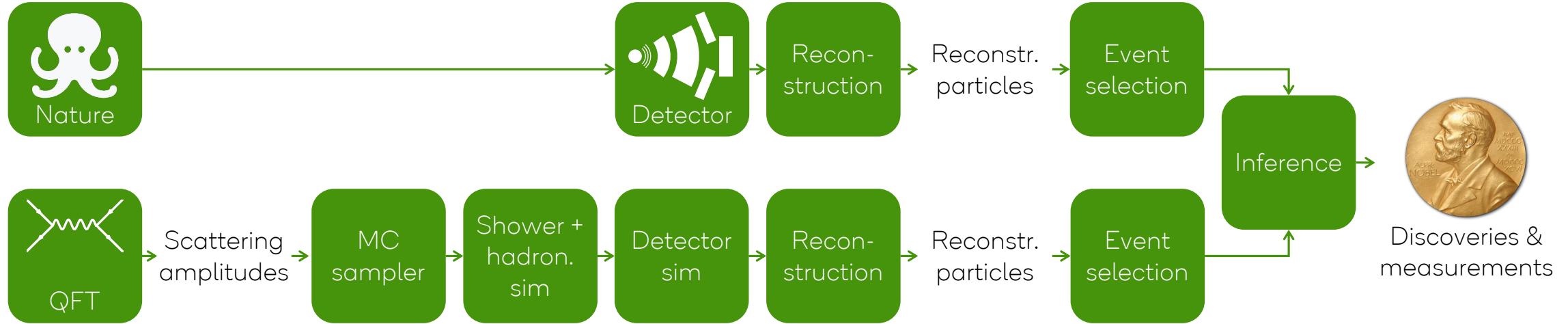
Steven Weinberg

$$F_{\alpha F}$$

$$Q(\bar{t}) \Delta \bar{t}$$

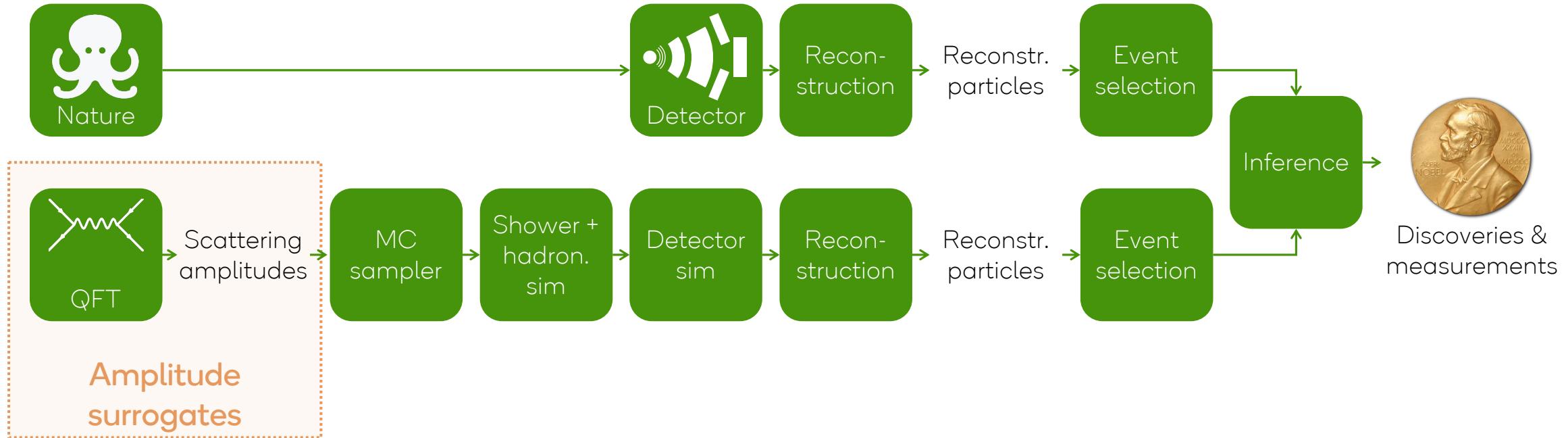


# hep data analysis pipeline



- Most steps have some degree of **Lorentz symmetry**
  - Detector / measurement process can partially break it or make it approximate
- Yet, equivariant architectures for high-energy physics seem understudied
  - Notable exceptions: LorentzNet, PELICAN, ...

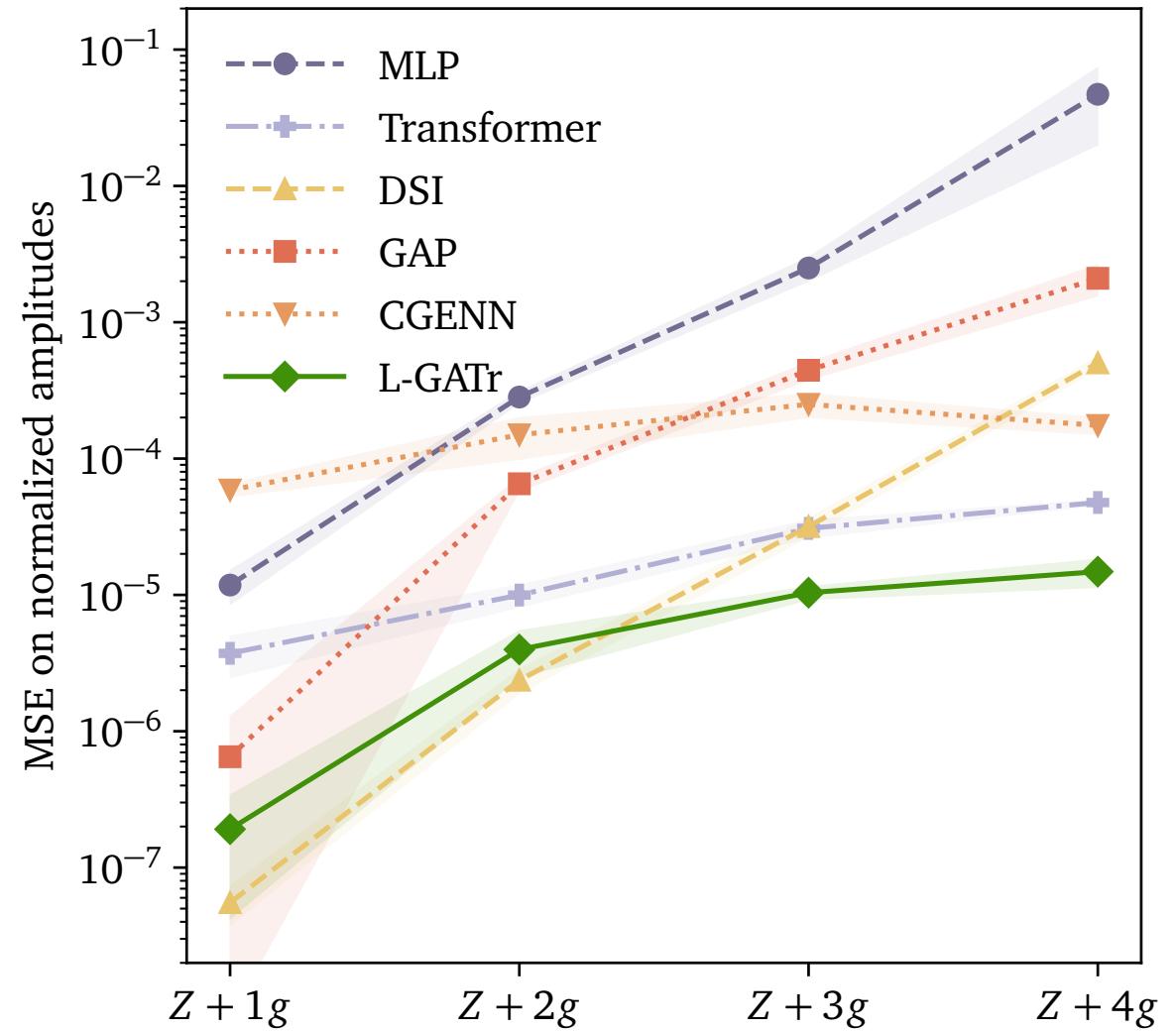
|                   | life at low speed  | high-energy physics  |
|-------------------|--|--|
| symmetry group    | $E(3)$   | Poincaré<br>but we only care about <b>Lorentz, <math>O^+(1,3)</math></b>   |
| data              | scalars<br>points<br>directions<br>...                         | scalars (e.g. particle ID)<br>four-vectors (e.g. four-momenta)<br>...  |
| GATr rep          | <b>projective GA <math>\mathcal{G}_{3,0,1}</math></b><br>(16D) | <b>space-time GA <math>\mathcal{G}_{1,3,0}</math></b><br>(also 16D, but different metric)                                      |
| GATr layers       | equi linear<br>attention<br>normalization<br>...               | equi linear (derived for different symmetry)<br>attention (based on new metric)<br>normalization (robust to neg. norms)<br>... |
| GATr architecture | transformer  | transformer  |



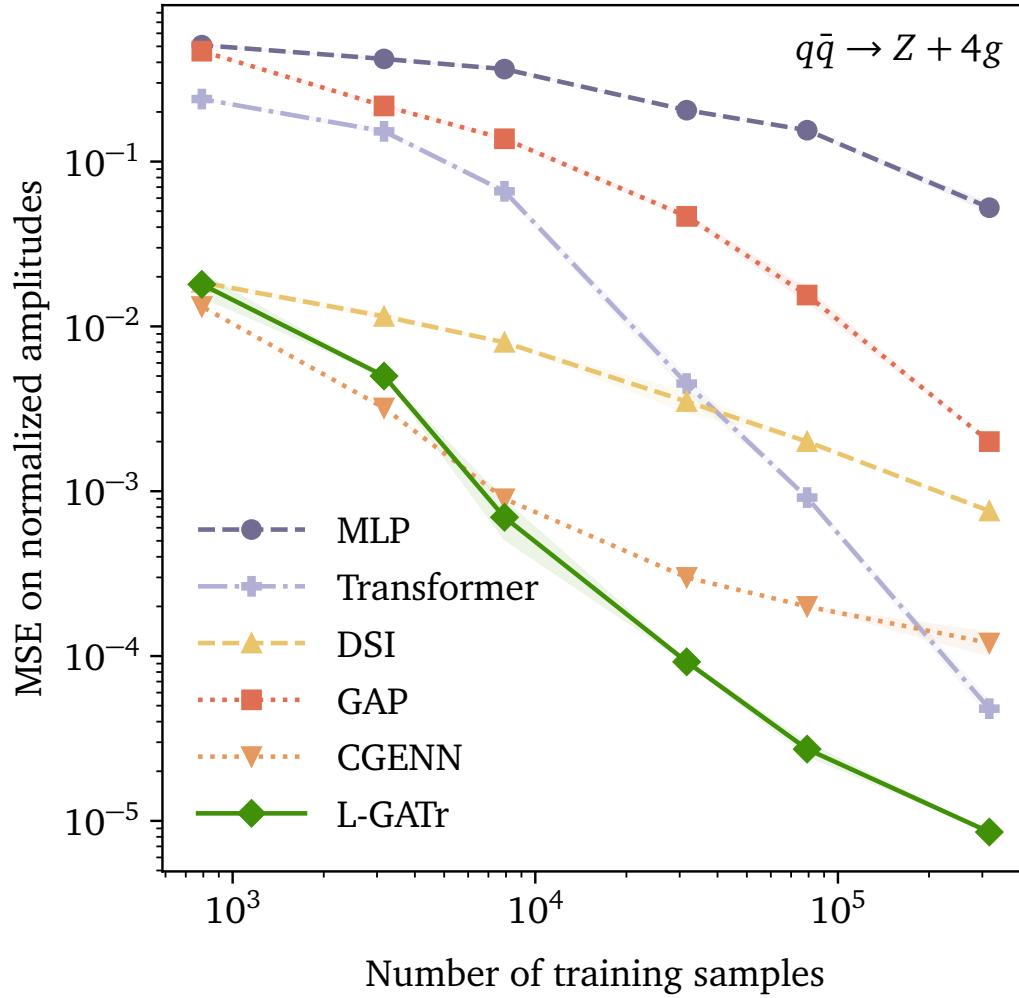
# Amplitude regression

Surrogate models for QFT amplitudes are most useful (and most challenging) at high multiplicity

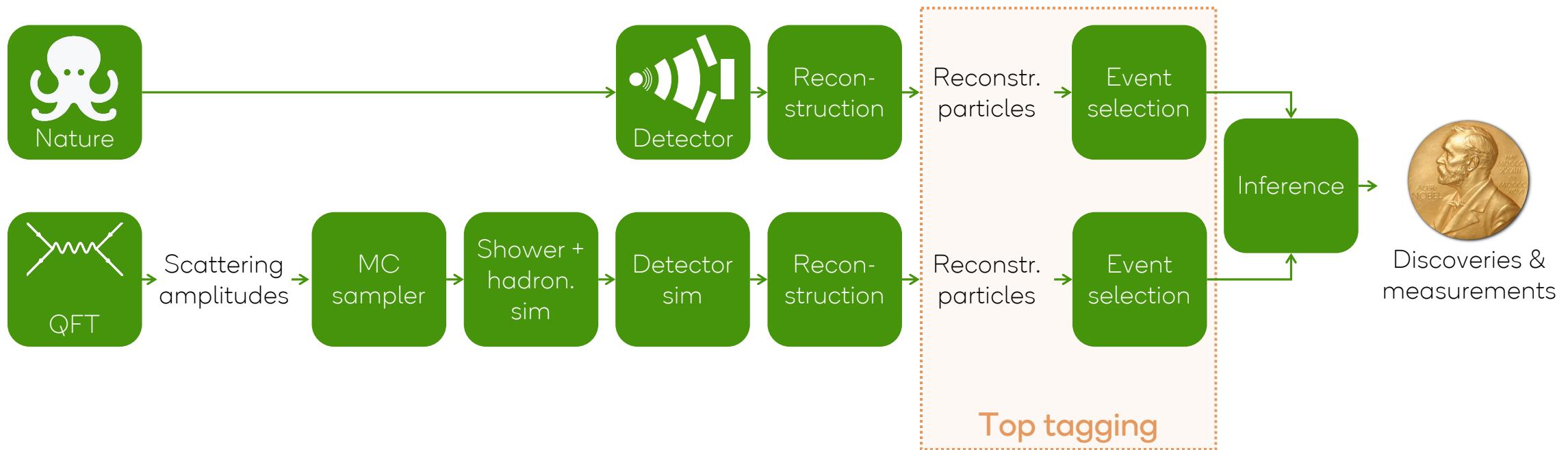
That's where **L-GATr** shines!



# Amplitude regression



**L-GATr** is both a **small-data** and a **big-data** architecture



# Top tagging on everyone's favorite dataset

| 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 [60] | 0.929    | 0.981  | –                                     | $867 \pm 15$                          |
| PFN [50]               | 0.932    | 0.9819 | $247 \pm 3$                           | $888 \pm 17$                          |
| TreeNiN [56]           | 0.933    | 0.982  | –                                     | $1025 \pm 11$                         |
| ParticleNet [62]       | 0.940    | 0.9858 | $397 \pm 7$                           | $1615 \pm 93$                         |
| ParT [63]              | 0.940    | 0.9858 | $413 \pm 16$                          | $1602 \pm 81$                         |

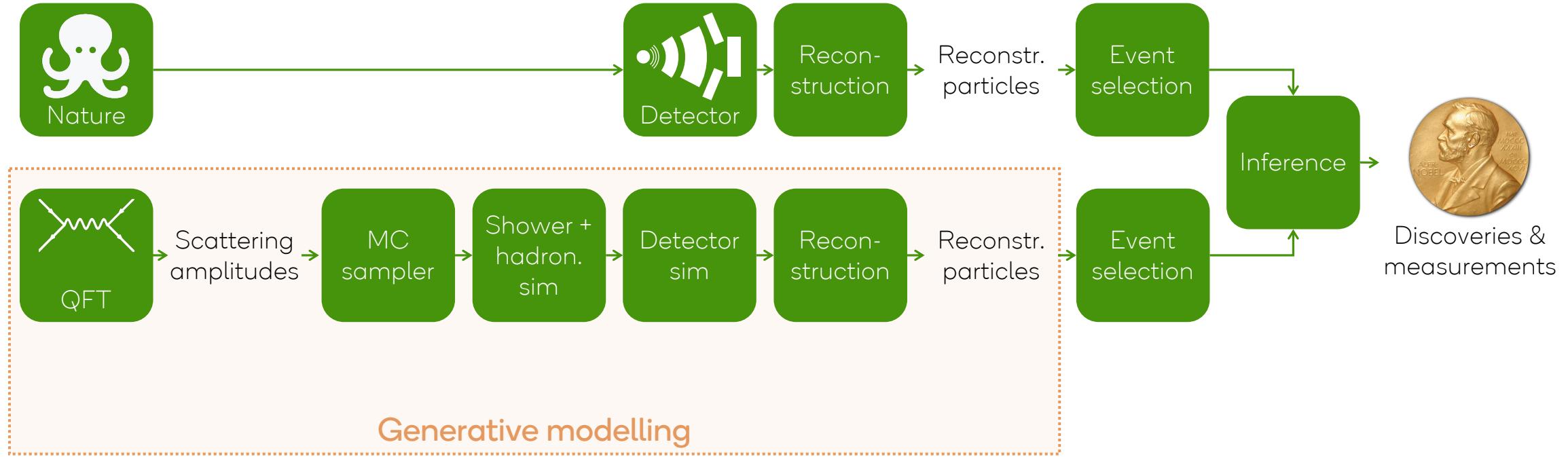
# Top tagging on everyone's favorite dataset

| 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 [60] | 0.929                                 | 0.981                                 | –                                     | $867 \pm 15$                          |
| PFN [50]               | 0.932                                 | 0.9819                                | $247 \pm 3$                           | $888 \pm 17$                          |
| TreeNiN [56]           | 0.933                                 | 0.982                                 | –                                     | $1025 \pm 11$                         |
| ParticleNet [62]       | 0.940                                 | 0.9858                                | $397 \pm 7$                           | $1615 \pm 93$                         |
| ParT [63]              | 0.940                                 | 0.9858                                | $413 \pm 16$                          | $1602 \pm 81$                         |
| LorentzNet* [41]       | 0.942                                 | 0.9868                                | $498 \pm 18$                          | $2195 \pm 173$                        |
| CGENN* [66]            | 0.942                                 | 0.9869                                | 500                                   | 2172                                  |
| PELICAN* [9]           | <b><math>0.9426 \pm 0.0002</math></b> | <b><math>0.9870 \pm 0.0001</math></b> | –                                     | <b><math>2250 \pm 75</math></b>       |

# Top tagging on everyone's favorite dataset

| Model                       | Accuracy               | AUC                    | $1/\epsilon_B$ ( $\epsilon_S = 0.5$ ) | $1/\epsilon_B$ ( $\epsilon_S = 0.3$ ) |
|-----------------------------|------------------------|------------------------|---------------------------------------|---------------------------------------|
| TopoDNN [48]                | 0.916                  | 0.972                  | –                                     | <b>295 ± 5</b>                        |
| LoLa [15]                   | 0.929                  | 0.980                  | –                                     | <b>722 ± 17</b>                       |
| P-CNN [1]                   | 0.930                  | 0.9803                 | <b>201 ± 4</b>                        | <b>759 ± 24</b>                       |
| <i>N</i> -subjettiness [60] | 0.929                  | 0.981                  | –                                     | <b>867 ± 15</b>                       |
| PFN [50]                    | 0.932                  | 0.9819                 | <b>247 ± 3</b>                        | <b>888 ± 17</b>                       |
| TreeNiN [56]                | 0.933                  | 0.982                  | –                                     | <b>1025 ± 11</b>                      |
| ParticleNet [62]            | 0.940                  | 0.9858                 | <b>397 ± 7</b>                        | <b>1615 ± 93</b>                      |
| ParT [63]                   | 0.940                  | 0.9858                 | <b>413 ± 16</b>                       | <b>1602 ± 81</b>                      |
| LorentzNet* [41]            | 0.942                  | 0.9868                 | <b>498 ± 18</b>                       | <b>2195 ± 173</b>                     |
| CGENN* [66]                 | 0.942                  | 0.9869                 | <b>500</b>                            | <b>2172</b>                           |
| PELICAN* [9]                | <b>0.9426 ± 0.0002</b> | <b>0.9870 ± 0.0001</b> | –                                     | <b>2250 ± 75</b>                      |
| <b>L-GATr (ours)*</b>       | <b>0.9417 ± 0.0002</b> | <b>0.9868 ± 0.0001</b> | <b>548 ± 26</b>                       | <b>2148 ± 106</b>                     |

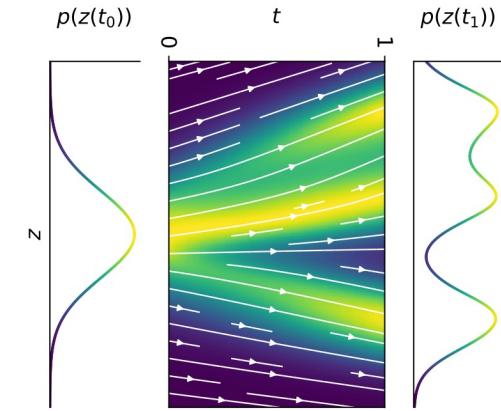
**L-GATr** is on par (but not better than) the best equivariant (\*) baselines



# Generative modelling with conditional flow matching

- **Continuous normalizing flows:**

- simple base density
- + continuous dynamics described by differential equation
- = generative model

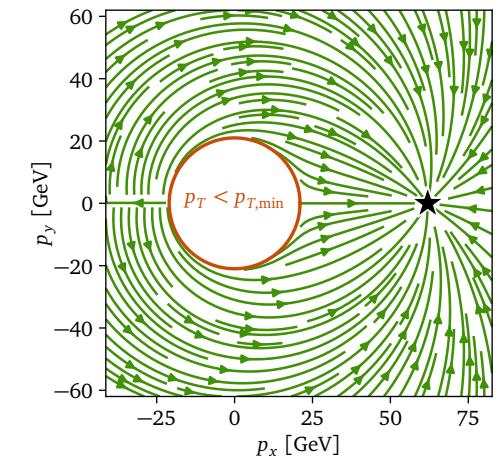


- **Conditional flow matching:**

- a brilliantly simple way to train them

- **Riemannian flow matching:**

- manifold-based approach that lets us incorporate boundaries



R. Chen et al, "Neural Ordinary Differential Equations", NeurIPS 2018

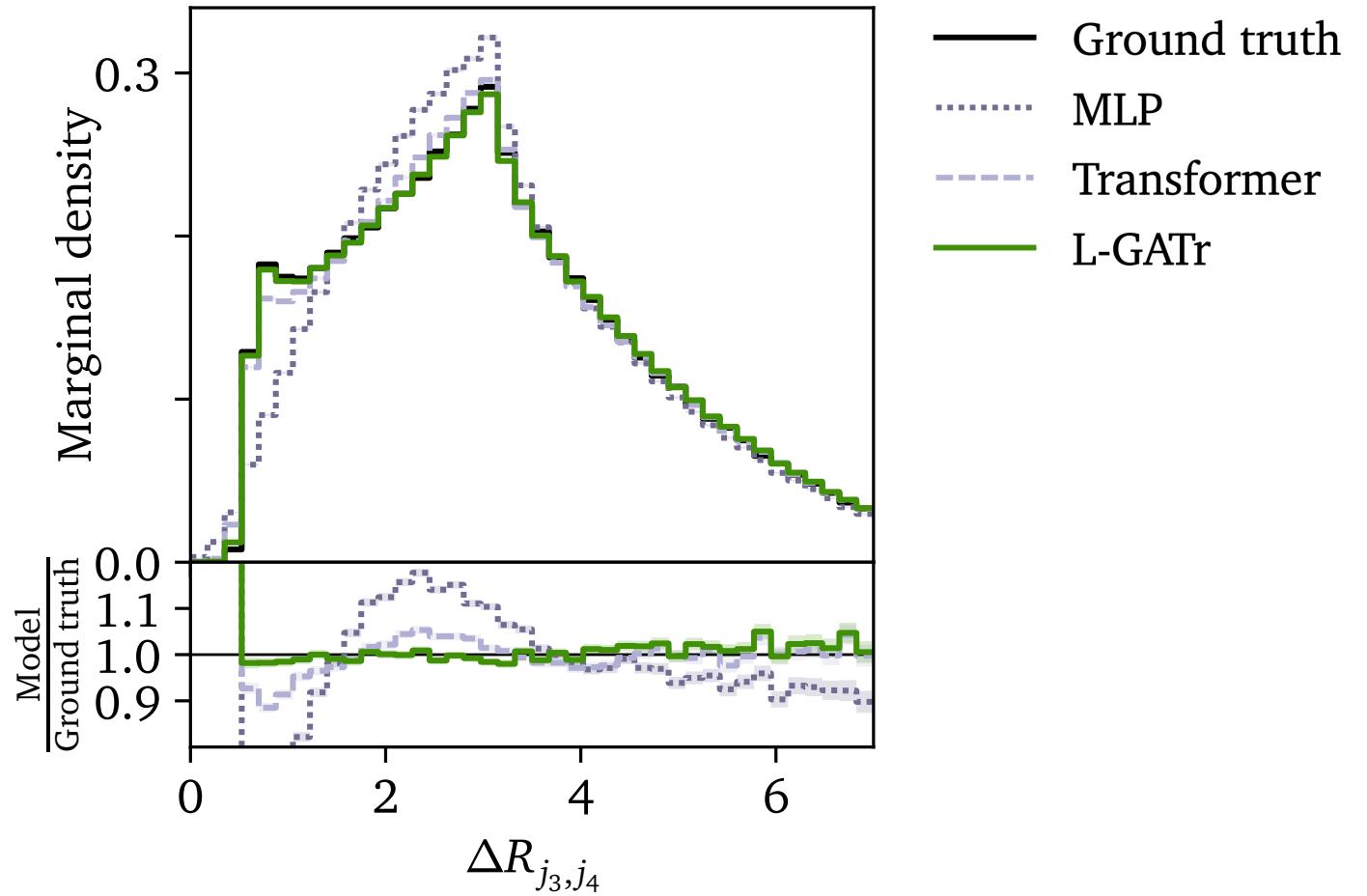
W. Grathwohl et al, "FFJORD: Free-form continuous dynamics for scalable reversible generative models", ICLR 2019

Y. Lipman et al, "Flow matching for generative modelling", ICLR 2023

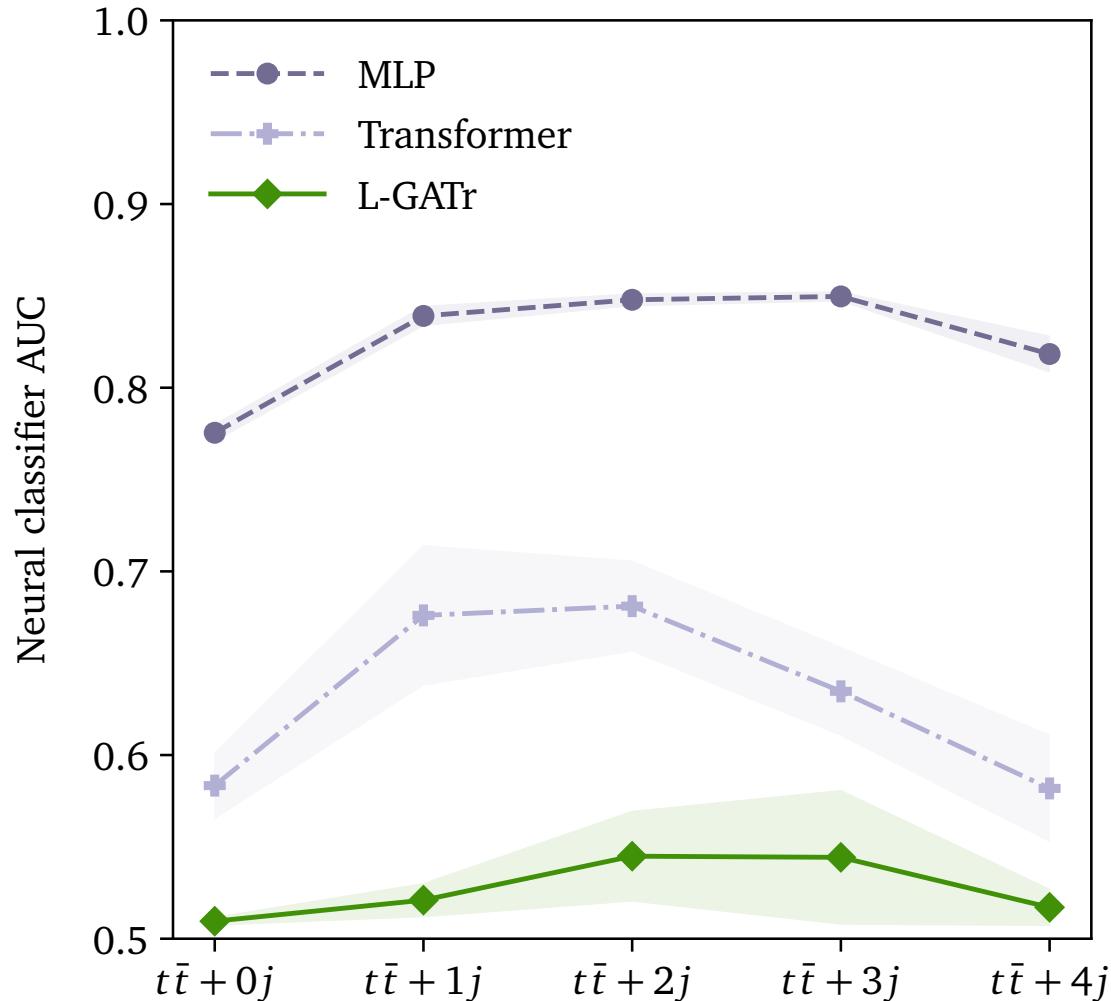
R. Chen et al, "Flow matching on general geometries", ICLR 2024

# Generative modelling with conditional flow matching

L-GATr leads to improvements  
on tricky kinematic features...



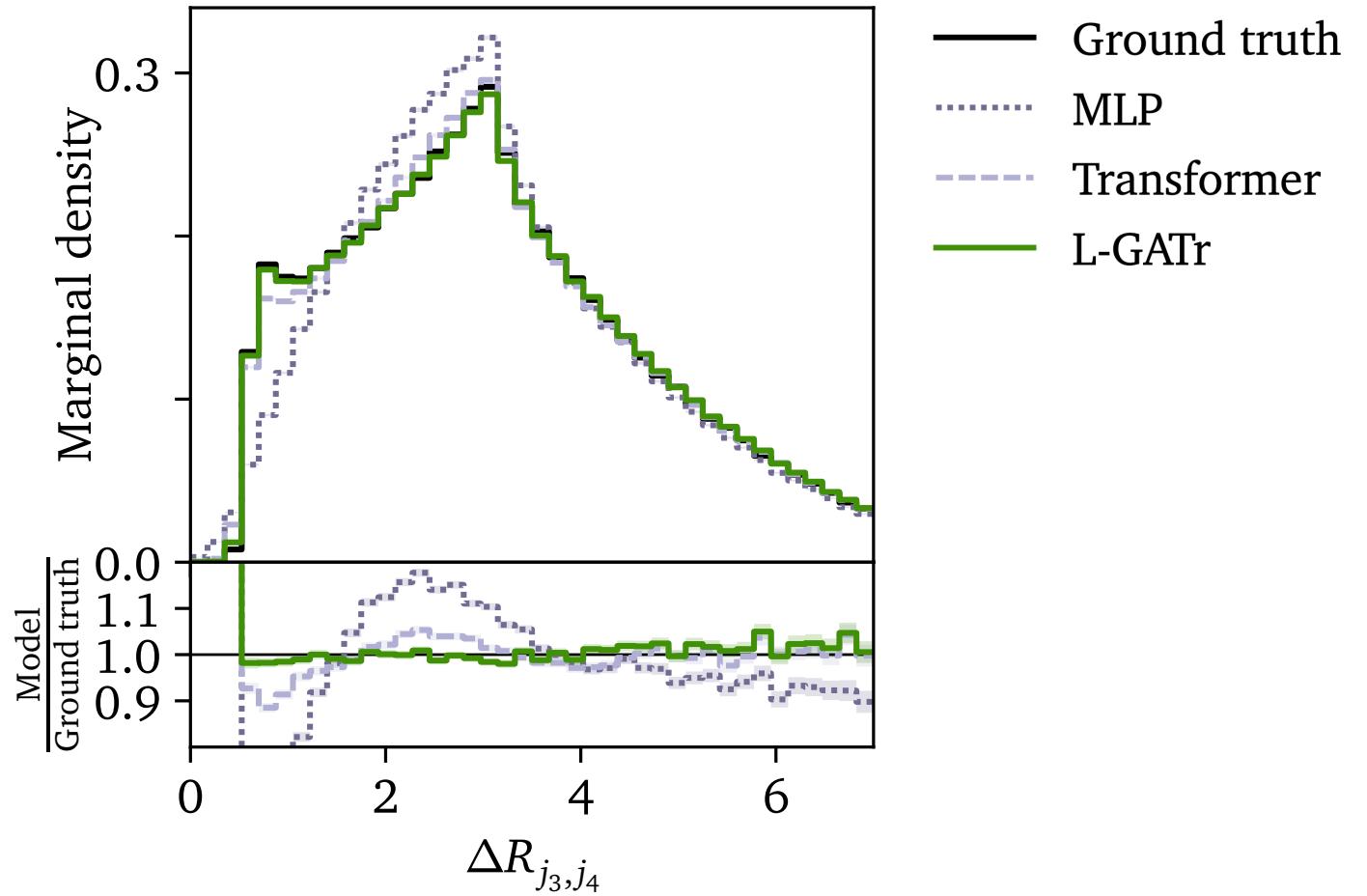
# Generative modelling with conditional flow matching

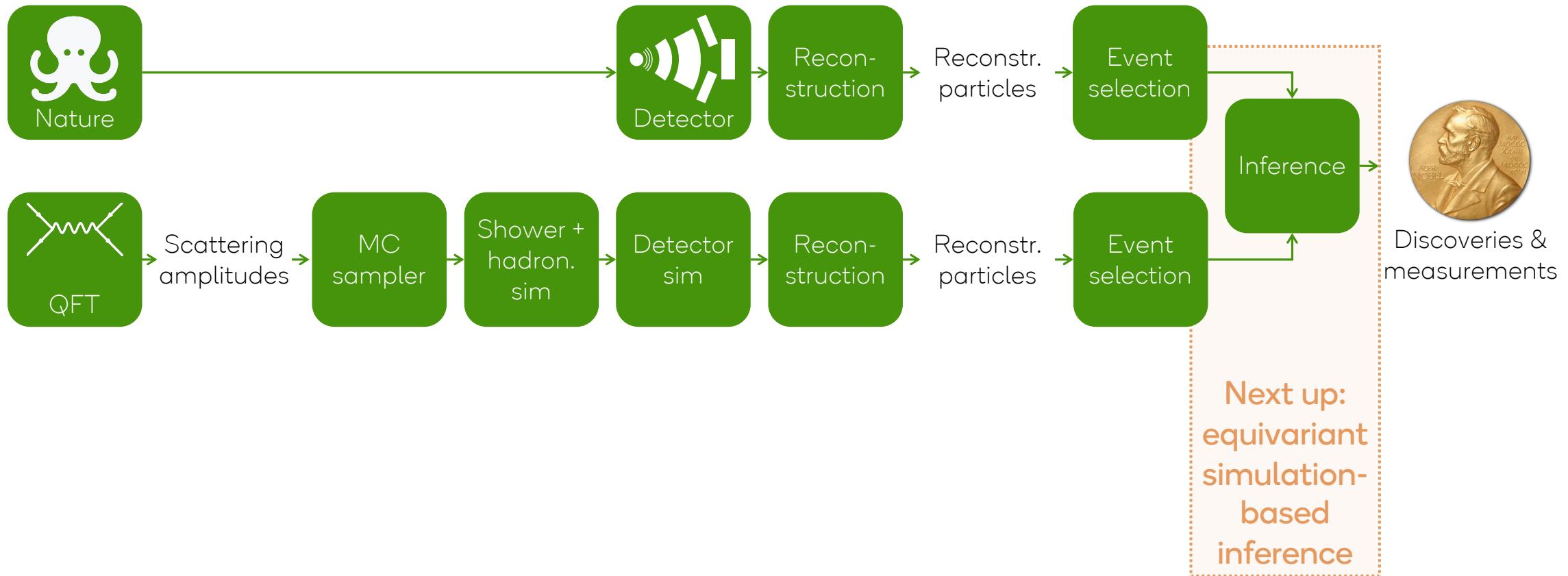


**L-GATr** generates samples that a classifier almost cannot distinguish from the ground truth...

# Generative modelling with conditional flow matching

...and gets tricky kinematic features right





## **Structure:**

Problem-specific  
inductive biases in  
algorithms and  
architectures

## **Geometric Algebra Transformer:**

Our version of a  
versatile architecture  
for geometric  
problems

## **Scale:**

Flexible  
architectures,  
lots of data,  
lots of compute

## Geometric Algebra Transformer

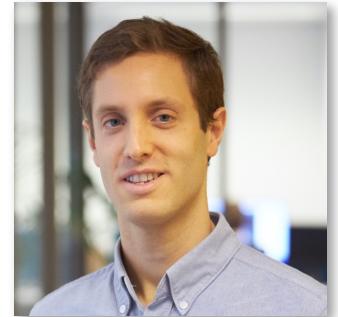
Johann Brehmer\*, Pim de Haan\*, Sönke Behrends, Taco Cohen  
NeurIPS 2023, [arXiv:2305.18415](https://arxiv.org/abs/2305.18415)



Pim de Haan



Sönke Behrends



Taco Cohen

## Euclidean, Projective, Conformal: Choosing a Geometric Algebra for your Equivariant Transformer

Pim de Haan, Taco Cohen, Johann Brehmer  
AISTATS 2024, [arXiv:2311.04744](https://arxiv.org/abs/2311.04744)

## Probabilistic and Differentiable Wireless Simulation with Geometric Transformers

Thomas Hehn, Markus Peschl, Tribhuvanesh Orekondy,  
Arash Behboodi, Johann Brehmer  
Under review



Thomas Hehn



Markus Peschl



Tribhuvanesh  
Orekondy



Arash  
Behboodi

## Lorentz-Equivariant Geometric Algebra Transformers for High-Energy Physics

Jonas Spinner\*, Victor Bresó\*, Pim de Haan, Tilman Plehn,  
Jesse Thaler, Johann Brehmer  
[arXiv:2405.14806](https://arxiv.org/abs/2405.14806)



Jonas Spinner



Victor Bresó



Jesse Thaler



Tilman Plehn

## Clifford group equivariant neural networks

David Ruhe, Johannes Brandstetter, Patrick Forré  
NeurIPS 2023, [arXiv:2305.11141](https://arxiv.org/abs/2305.11141)

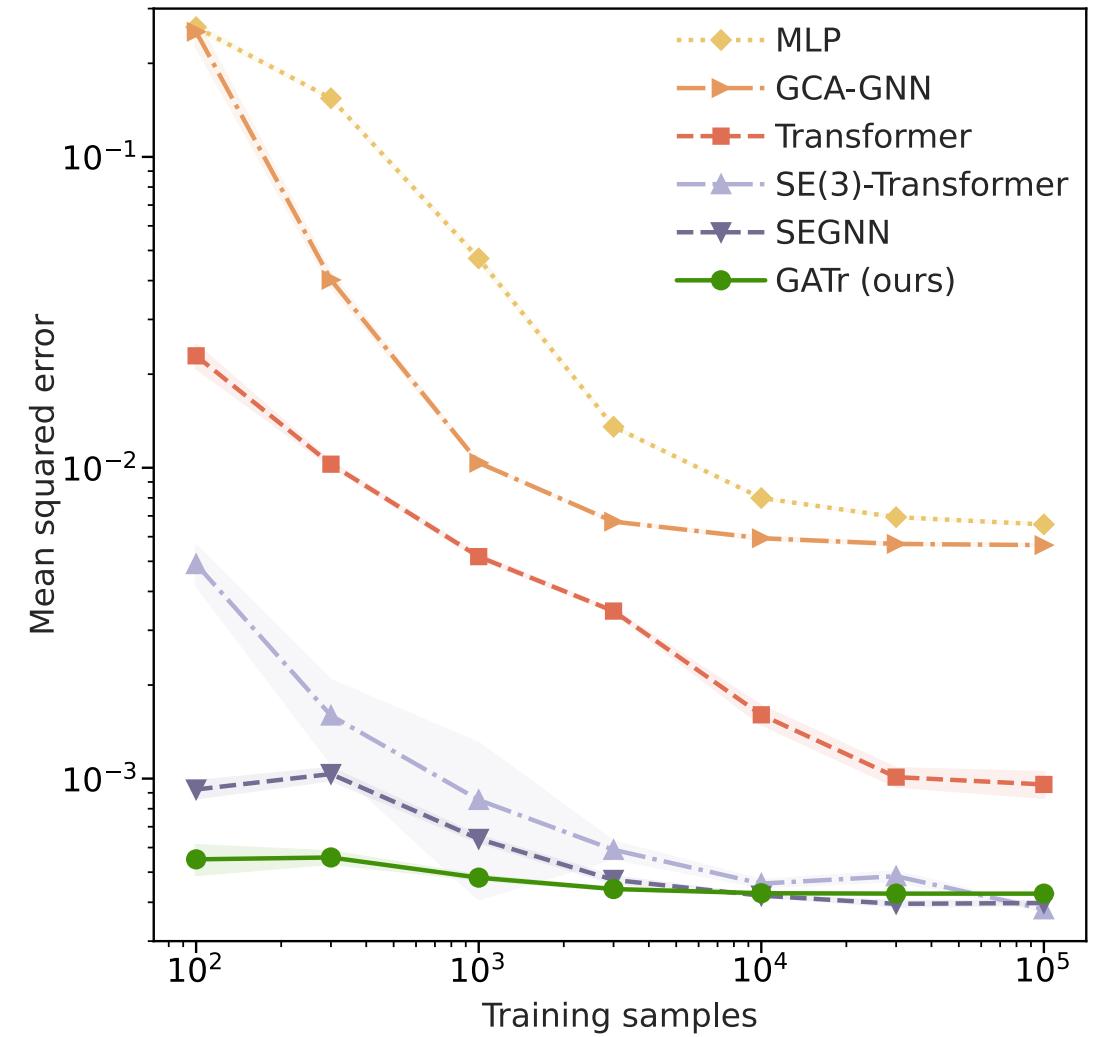
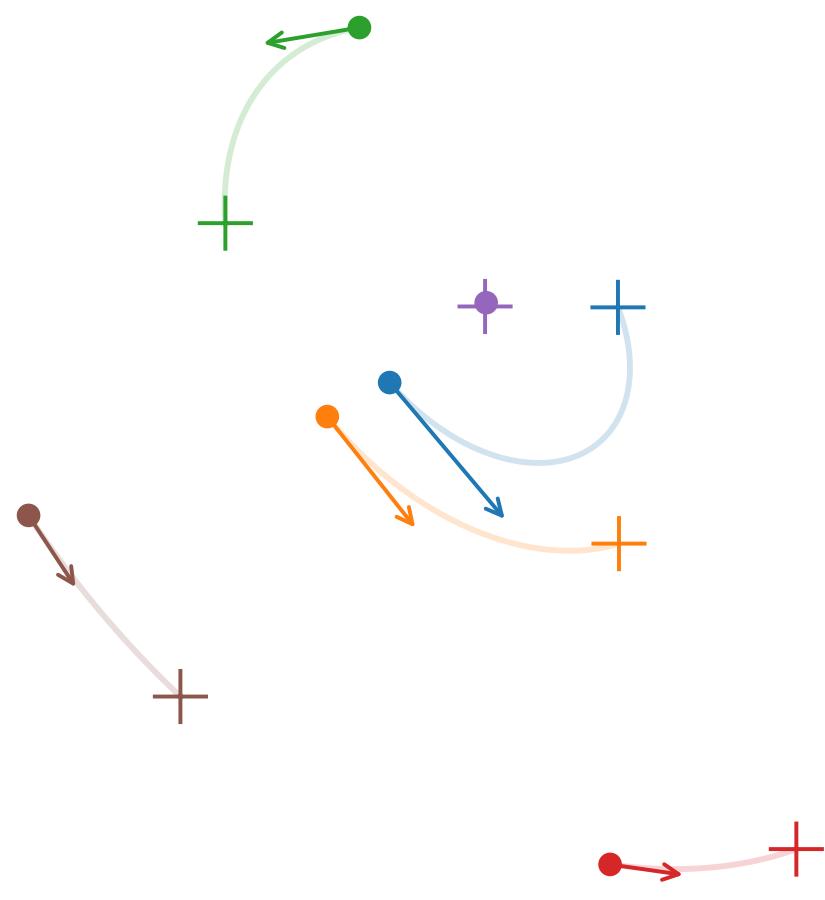
## LaB-GATr: geometric algebra transformers for large biomedical surface and volume meshes

Julian Suk, Baris Imre, Jelmer M. Wolterink  
[arXiv:2403.07536](https://arxiv.org/abs/2403.07536)

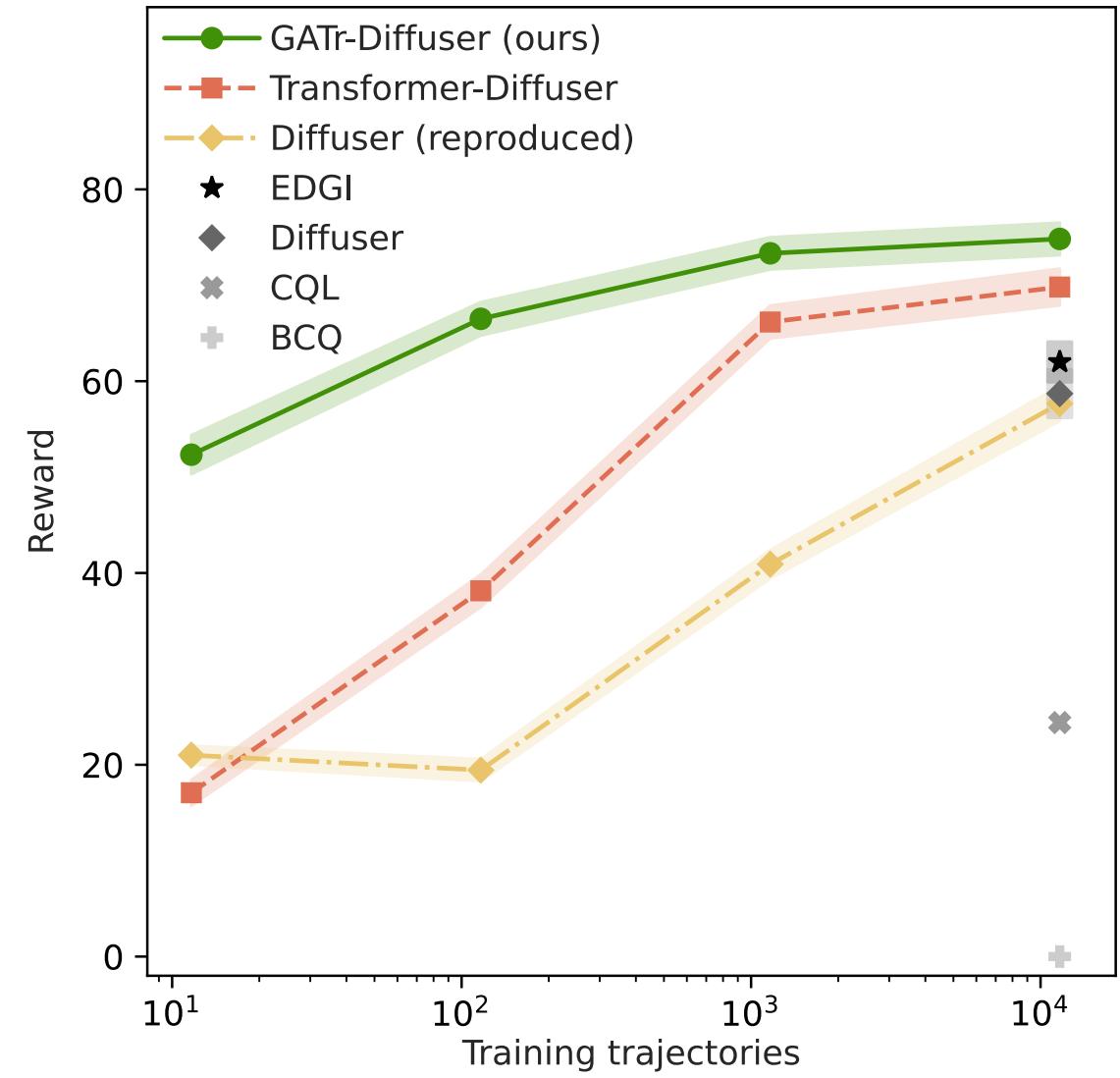
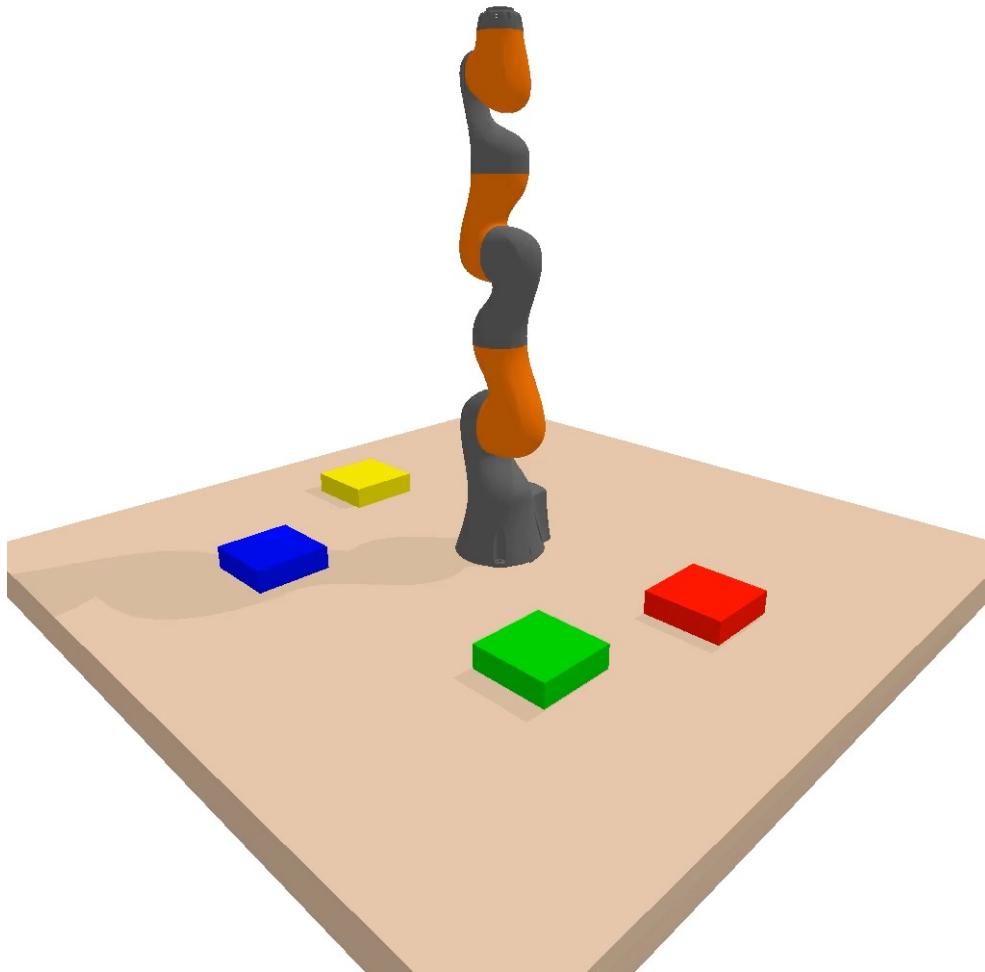
# Bonus material



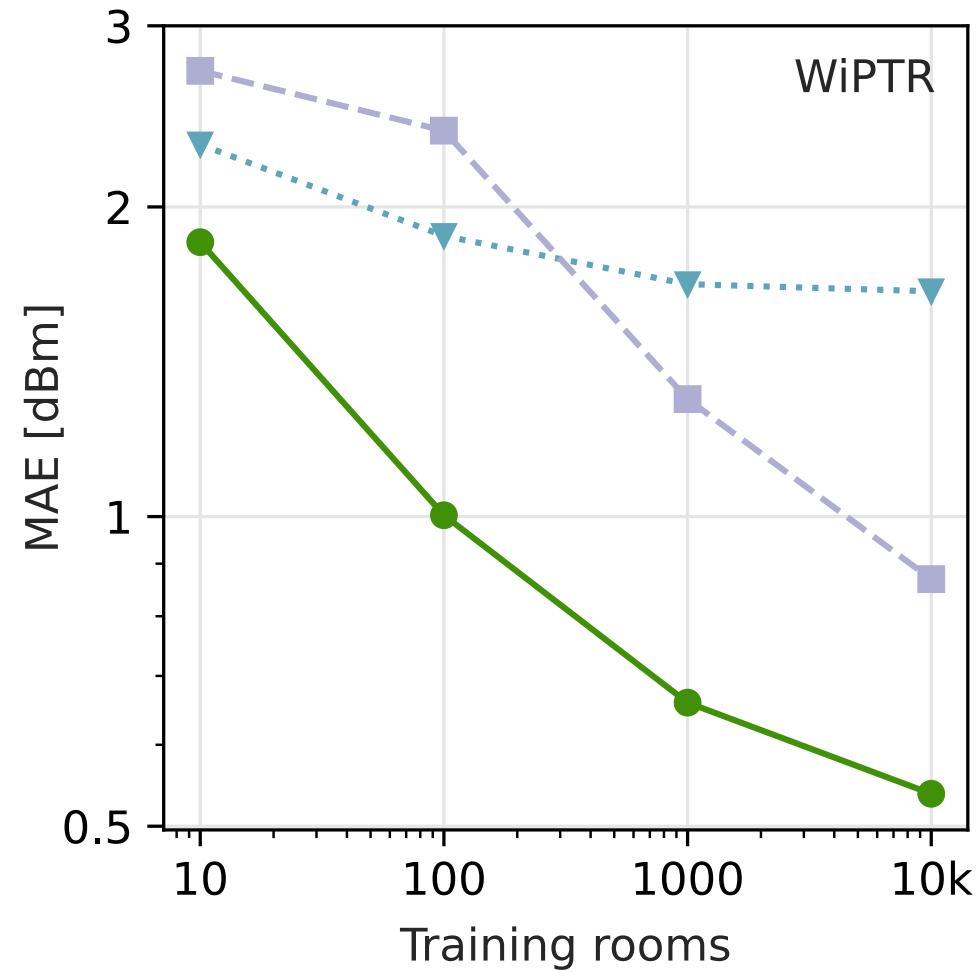
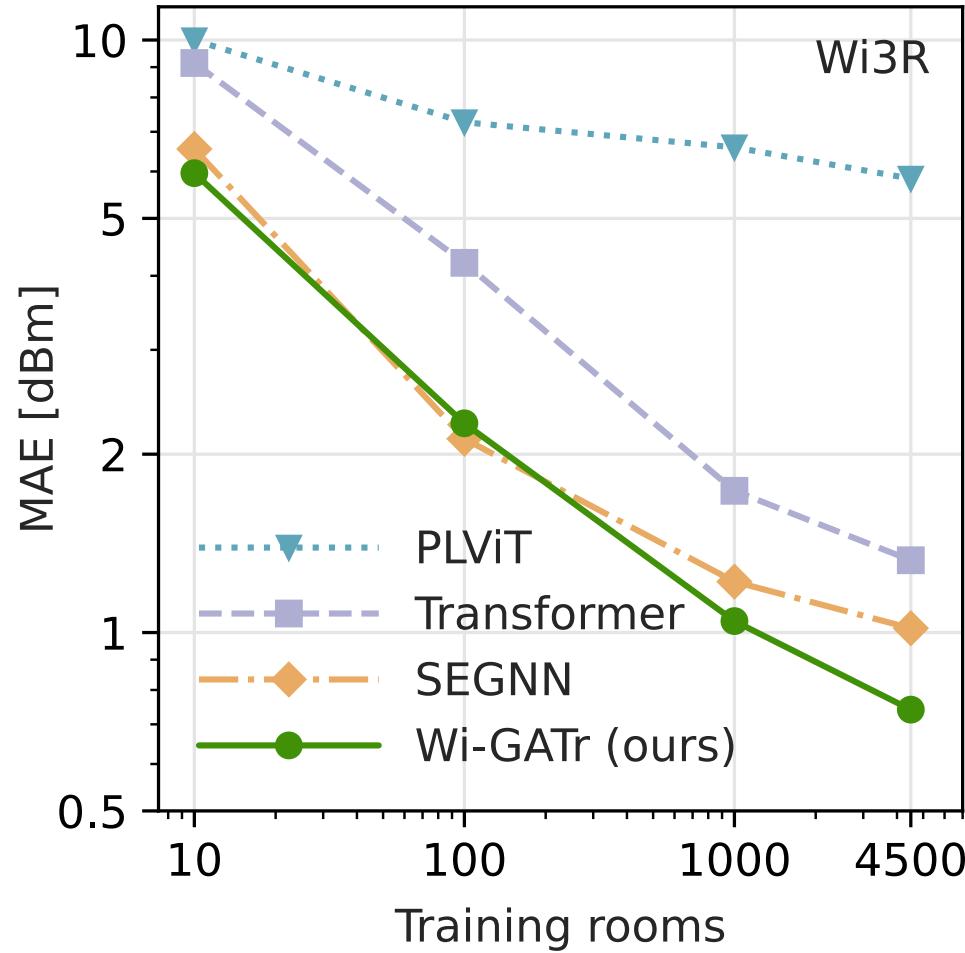
# Benchmarking: n-body modelling



# Diffusion-based planning: Robotic control



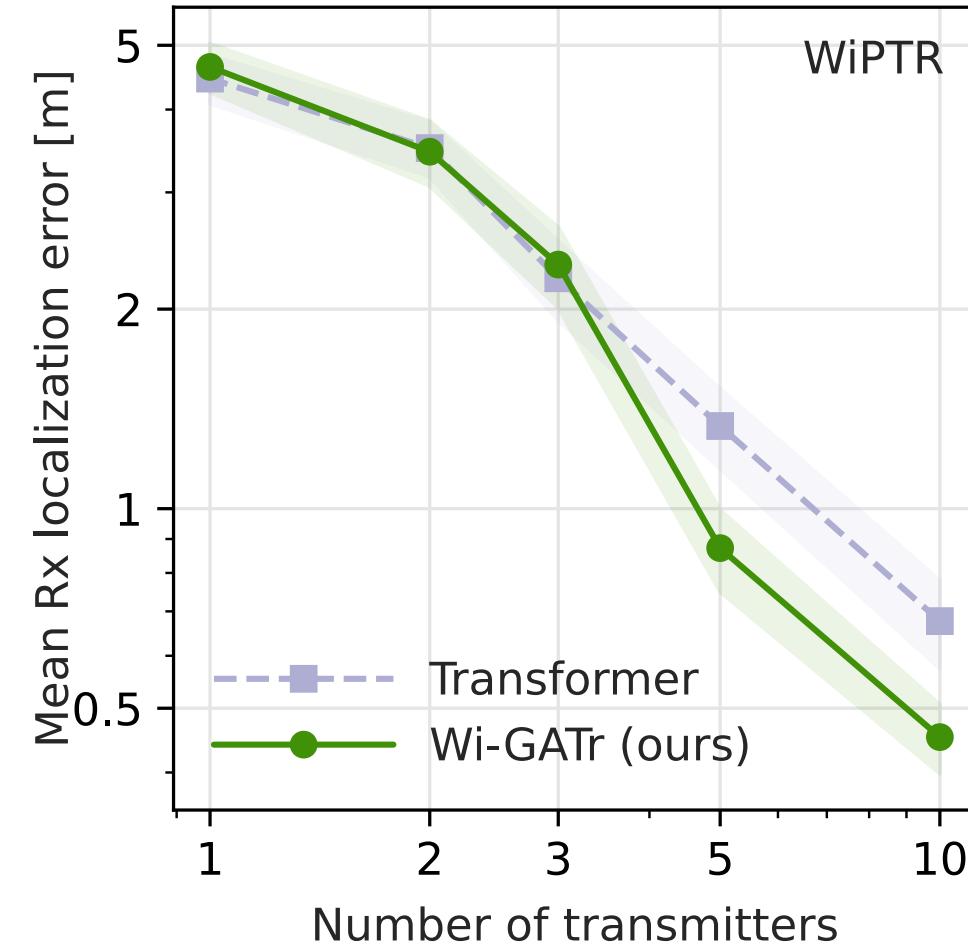
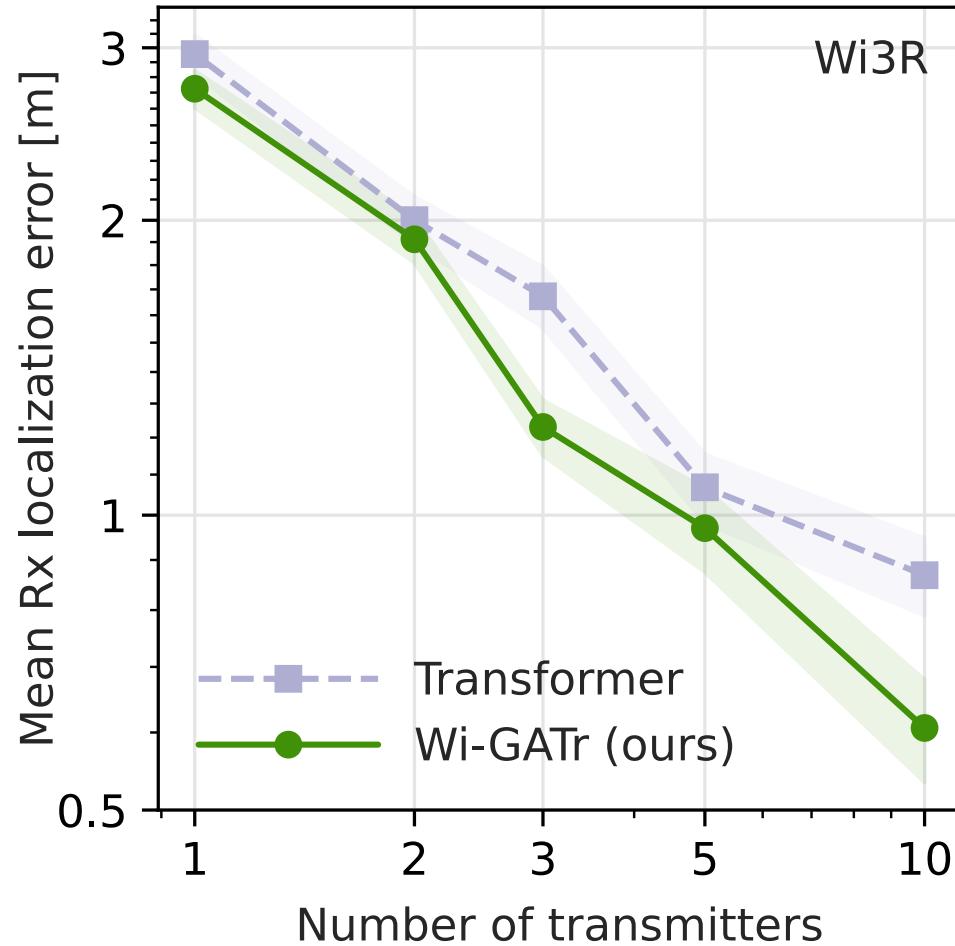
# Differentiable Wireless signal modelling



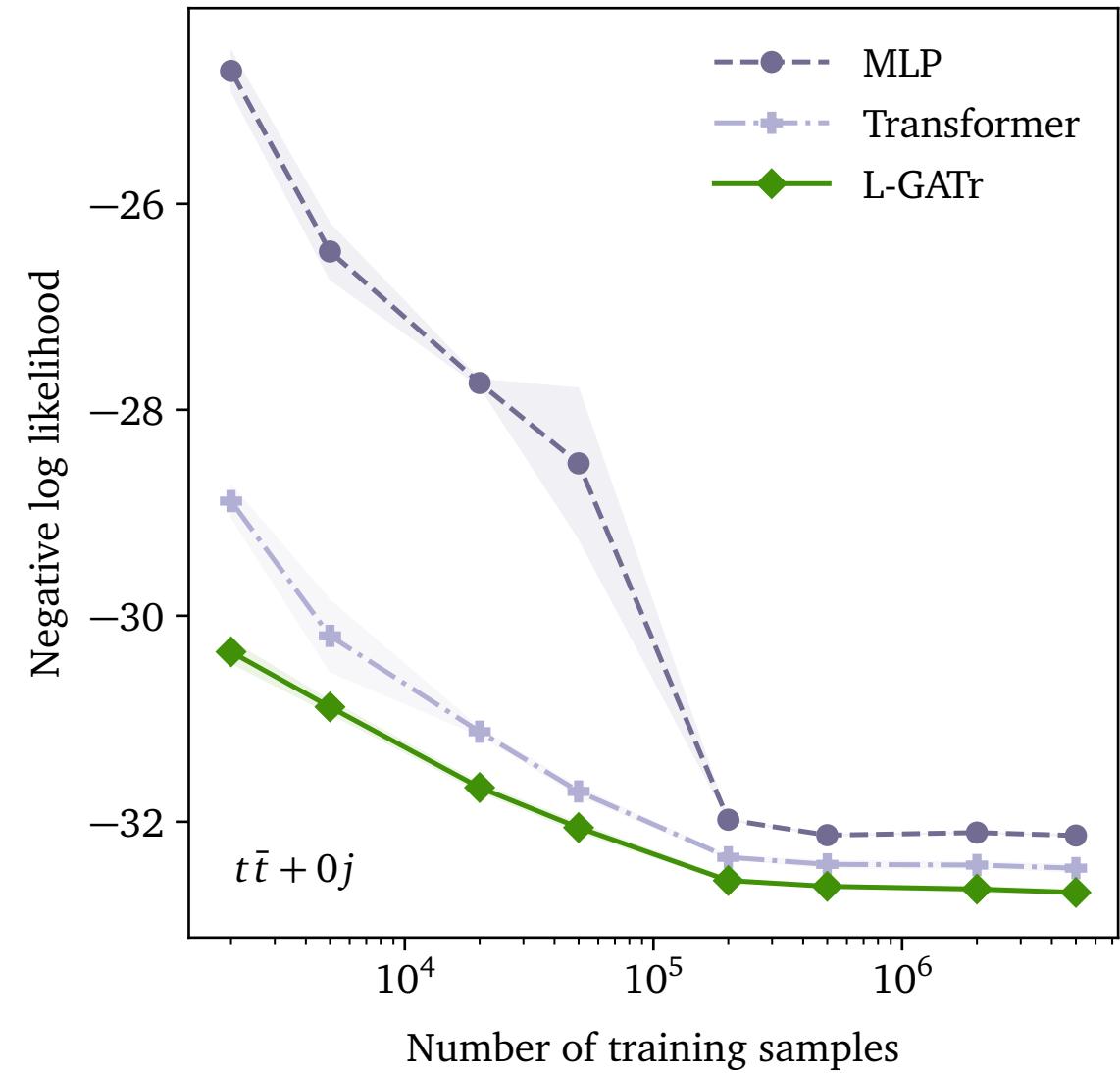
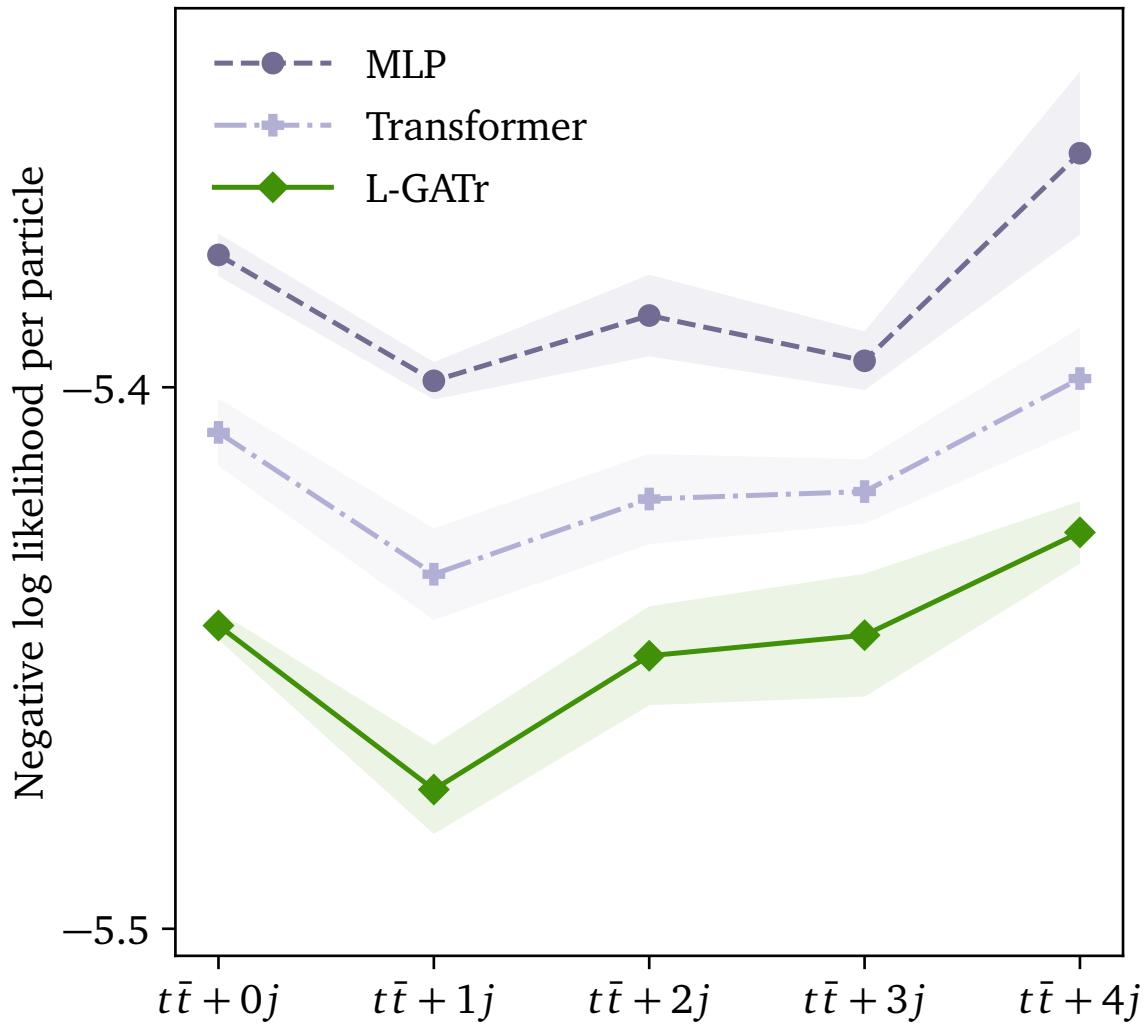
# Differentiable Wireless signal modelling

|                                 | Wi3R dataset      |         |             |       | WiPTR dataset     |         |       |
|---------------------------------|-------------------|---------|-------------|-------|-------------------|---------|-------|
|                                 | Wi-GATr<br>(ours) | Transf. | SEGNN       | PLViT | Wi-GATr<br>(ours) | Transf. | PLViT |
| <i>In distribution</i>          |                   |         |             |       |                   |         |       |
| Rx interpolation                | <b>0.63</b>       | 1.14    | 0.92        | 5.61  | <b>0.53</b>       | 0.84    | 1.67  |
| Unseen floor plans              | <b>0.74</b>       | 1.32    | 1.02        | 5.84  | <b>0.54</b>       | 0.87    | 1.66  |
| <i>Symmetry transformations</i> |                   |         |             |       |                   |         |       |
| Rotation                        | <b>0.74</b>       | 78.68   | 1.02        | 5.84  | <b>0.54</b>       | 28.17   | 1.66  |
| Translation                     | <b>0.74</b>       | 64.05   | 1.02        | 5.84  | <b>0.54</b>       | 4.04    | 1.66  |
| Permutation                     | <b>0.74</b>       | 1.32    | 1.02        | 5.84  | <b>0.54</b>       | 0.87    | 1.66  |
| Reciprocity                     | <b>0.74</b>       | 1.32    | 1.01        | 8.64  | <b>0.54</b>       | 0.87    | 1.65  |
| <i>Out of distribution</i>      |                   |         |             |       |                   |         |       |
| OOD layout                      | 9.24              | 14.06   | <b>2.34</b> | 7.00  | <b>0.54</b>       | 1.01    | 1.58  |

# Differentiable Wireless signal modelling



# Generative modelling for hep with conditional flow matching



# Structure vs scale

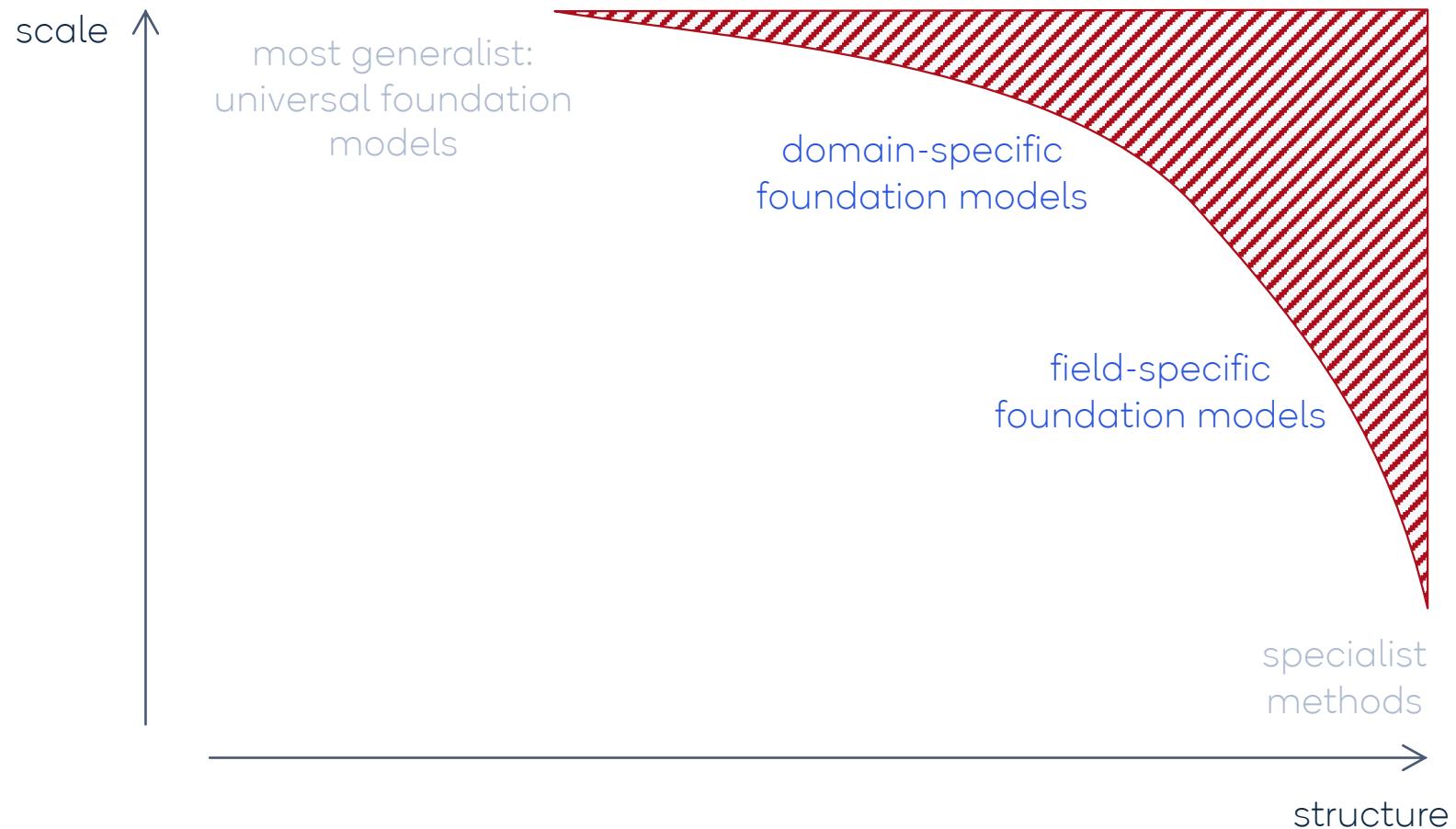
It is tempting to think of scale and structure as opposites

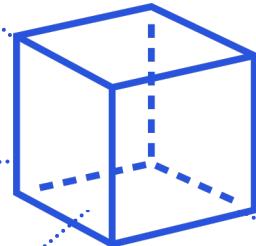
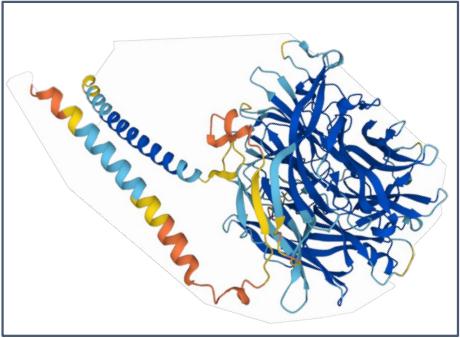


# Structure vs scale

That's a false dichotomy:

there are interesting, underexplored settings where **structure and scale can be complementary**

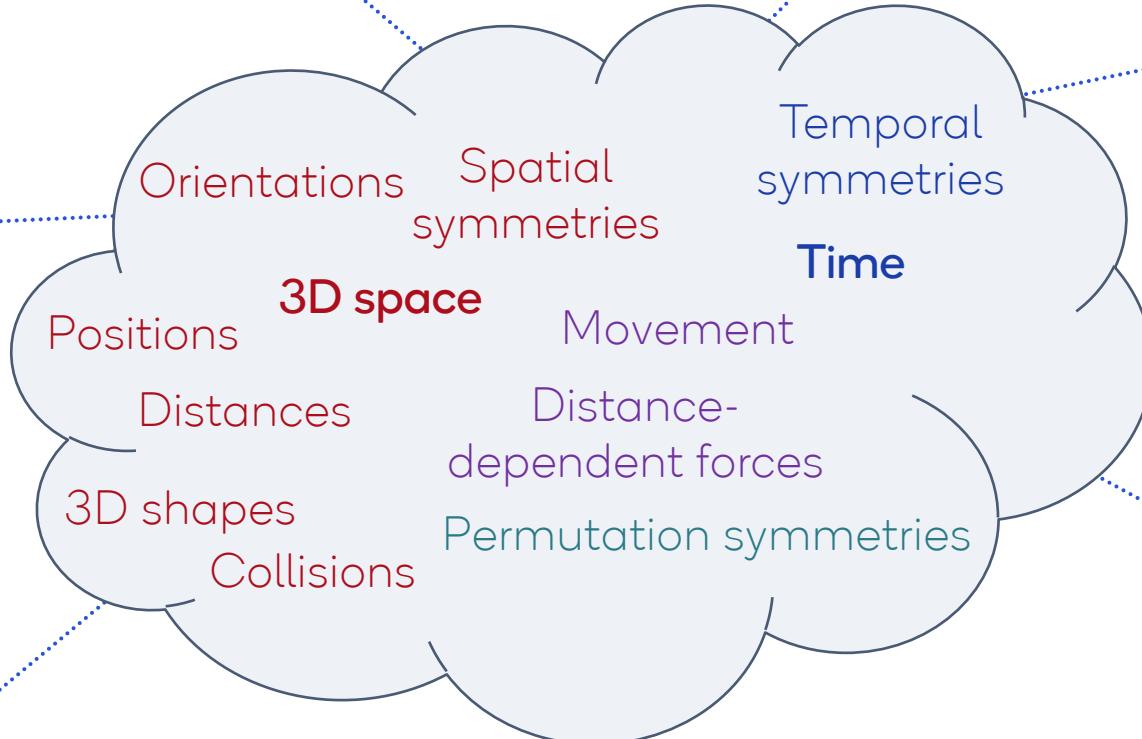
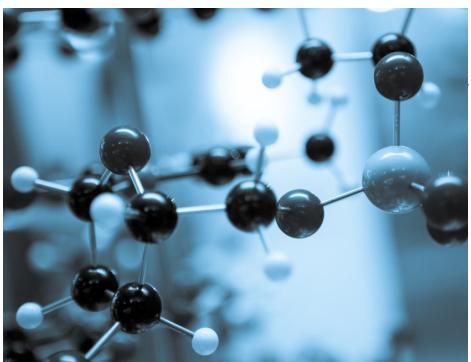
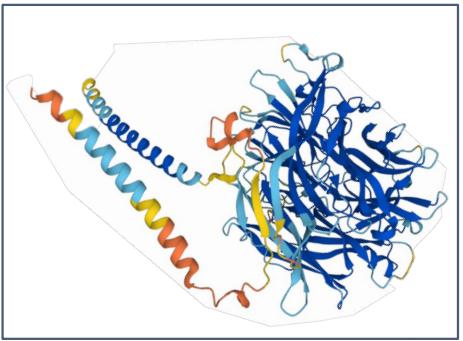




## Geometric foundation model

- Trained on (and finetuned to) various scientific and engineering problems
- Representations tailored to spatio-temporal data
- Architecture reflects symmetries of 3D space





## Foundation model checklist

- Lots of **data** the internet
- Universal **representations** tokenizers and embeddings
- A scalable, expressive **architecture** transformers
- Self-supervised **training** protocol max likelihood, ...
- Multiple **downstream tasks** with shared structure chatbots, ...
- Bonus points for **predictable scaling** neural scaling laws

## In language / vision foundation models

## Foundation model checklist

- Lots of **data**
- Universal **representations**
- A scalable, expressive **architecture**
- Self-supervised **training** protocol
- Multiple **downstream tasks** with shared structure
- Bonus points for **predictable scaling**

## Geometric foundation model



## Foundation model checklist

- Lots of **data**
- Universal **representations**
- A scalable, expressive **architecture**
- Self-supervised **training** protocol
- Multiple **downstream tasks** with shared structure
- Bonus points for **predictable scaling**

## Geometric foundation model

- ?
- could be geometric algebra–based**
- equivariant Transformers**
- ?
- GATr works on very different domains**
- ?

## Foundation model checklist

- Lots of **data**
- Universal **representations**
- A scalable, expressive **architecture**
- Self-supervised **training** protocol
- Multiple **downstream tasks** with shared structure
- Bonus points for **predictable scaling**

## Geometric foundation model

- data and simulators are out there?
- could be geometric algebra-based
- equivariant transformers
- predicting masked-out tokens?
- GATr works on very different domains
- ?

# Thank you



Follow us on: [in](#) [Twitter](#) [Instagram](#) [YouTube](#) [Facebook](#)

For more information, visit us at:

[qualcomm.com](http://qualcomm.com) & [qualcomm.com/blog](http://qualcomm.com/blog)

Nothing in these materials is an offer to sell any of the components or devices referenced herein.

©2018-2024 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

Qualcomm is a trademark or registered trademark of Qualcomm Incorporated. Other products and brand names may be trademarks or registered trademarks of their respective owners.

References in this presentation to "Qualcomm" may mean Qualcomm Incorporated, Qualcomm Technologies, Inc., and/or other subsidiaries or business units within the Qualcomm corporate structure, as applicable. Qualcomm Incorporated includes our licensing business, QTL, and the vast majority of our patent portfolio. Qualcomm Technologies, Inc., a subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of our engineering, research and development functions, and substantially all of our products and services businesses, including our QCT semiconductor business. Snapdragon and Qualcomm branded products are products of Qualcomm Technologies, Inc. and/or its subsidiaries. Qualcomm patented technologies are licensed by Qualcomm Incorporated.