# Learning interpretable positional encodings in transformers depends on initialization
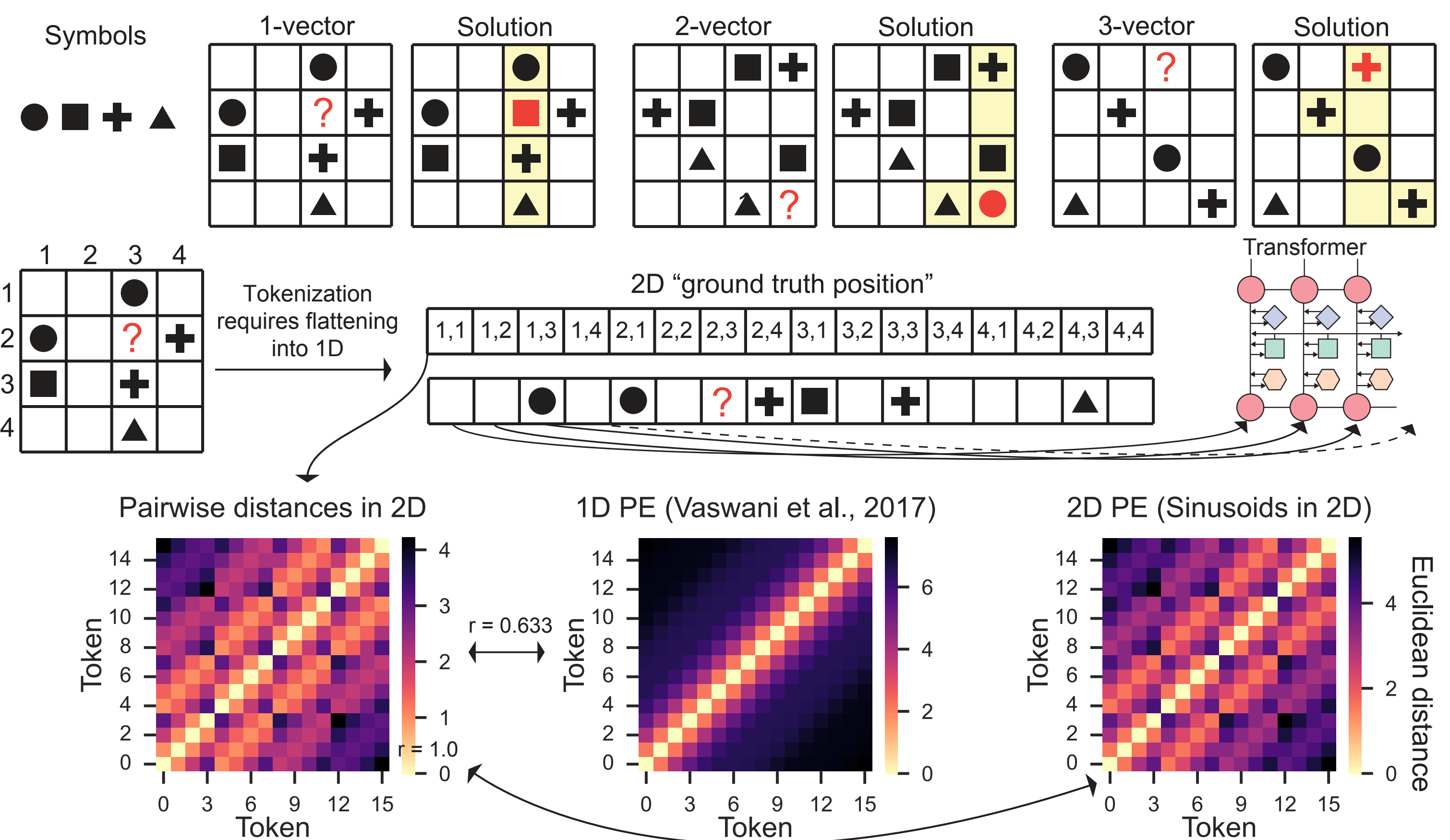
Takuya Ito[1], Luca Cocchi, Tim Klinger, Parikshit Ram, Murray Campbell, Luke J. Hearne

arXiv

[1]Mathematics of Computation, IBM Research, [2]QIMR Berghofer Medical Research Institute, QLD, Australia

E-mail: taku.ito1@gmail.com

## Motivation: The importance of positional encoding choice for transformer generalization

* Choice of positional encodings (PE) in transformers have been shown to be critical for learning and generalization.

* Most investigations into PE have been tailored towards 1D string-based tasks, such as arithmetic or context-free grammars using pre-specified PEs (e.g., ROPE, or absolute PEs)

* Here we investigate the importance on a suite of tasks with sequence data organized in higher dimensions (>greater than 1D sequences)

* Specifically, we study the conditions by which we can **learn interpretable positional encodings**, and study how they impact generalization

* Inspired by recent work on rich and lazy representation learning, we explored how initialization of a learnable PE parameter influences interpretability and generalization in transformers

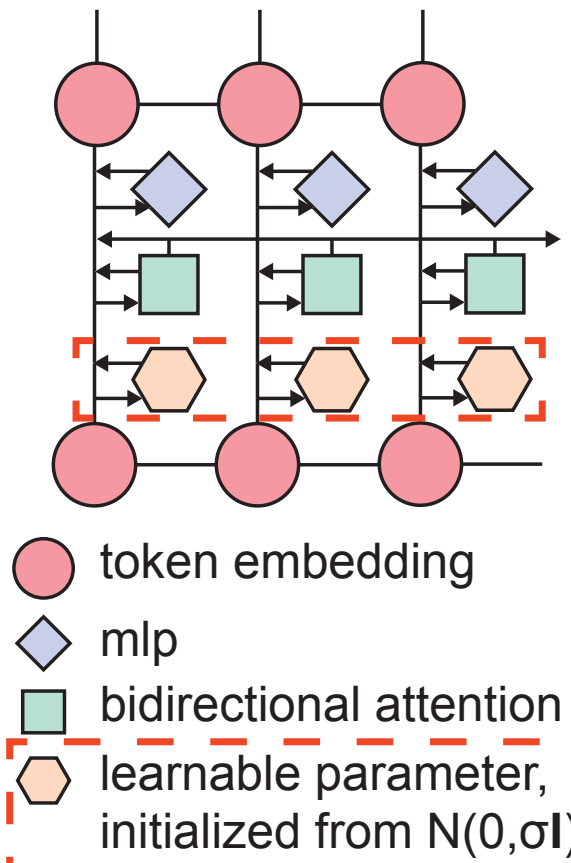Intuition from a simple 2D task: The Latin Squares Task (simplified Sudoku)



## Experiment 1: Learning interpretable positional encodings in the Latin Squares Task
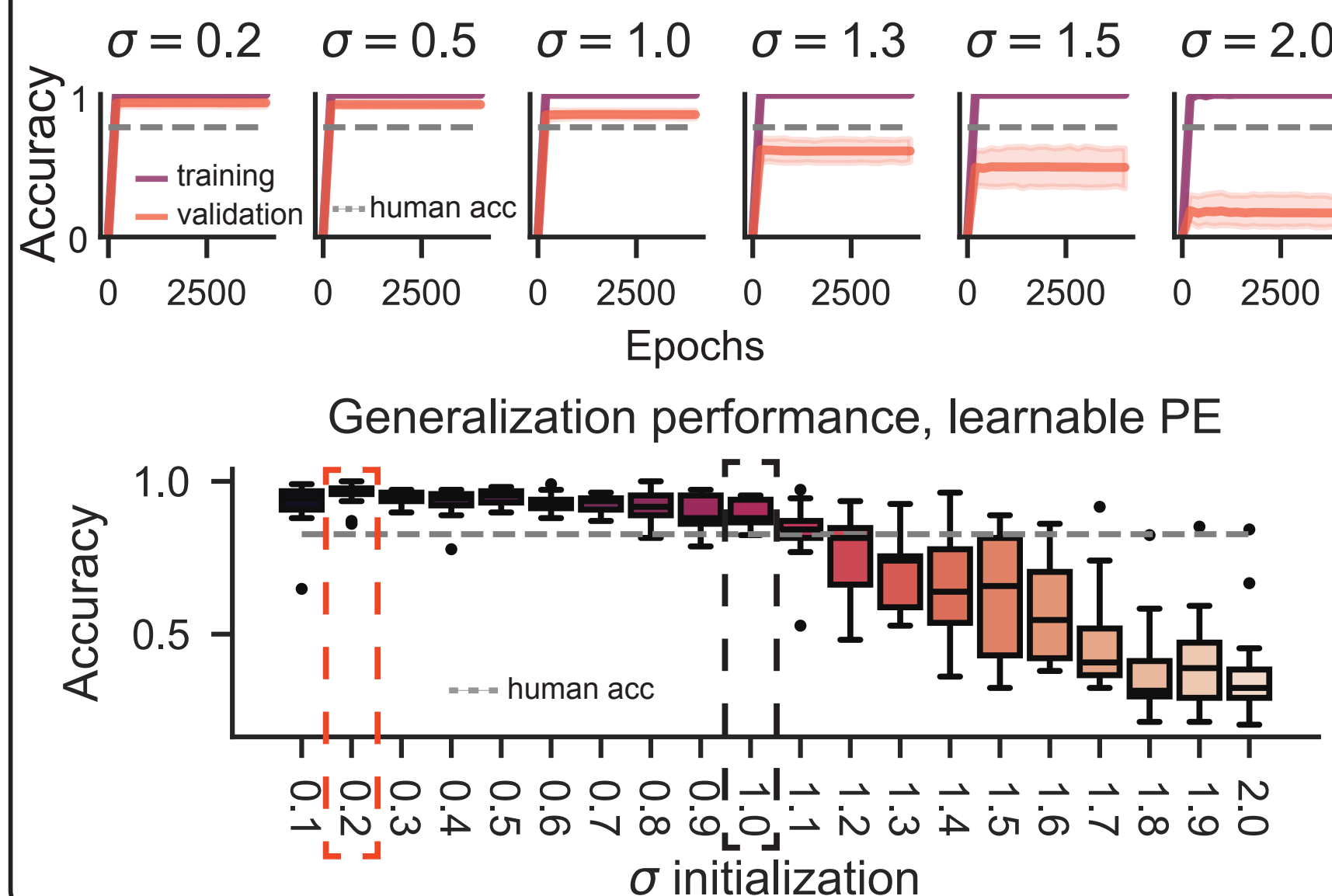
**Model manipulation**

Replace traditional PEs with a learnable parameter initialized from $N(0, \sigma I)$, and manipulate $\sigma$
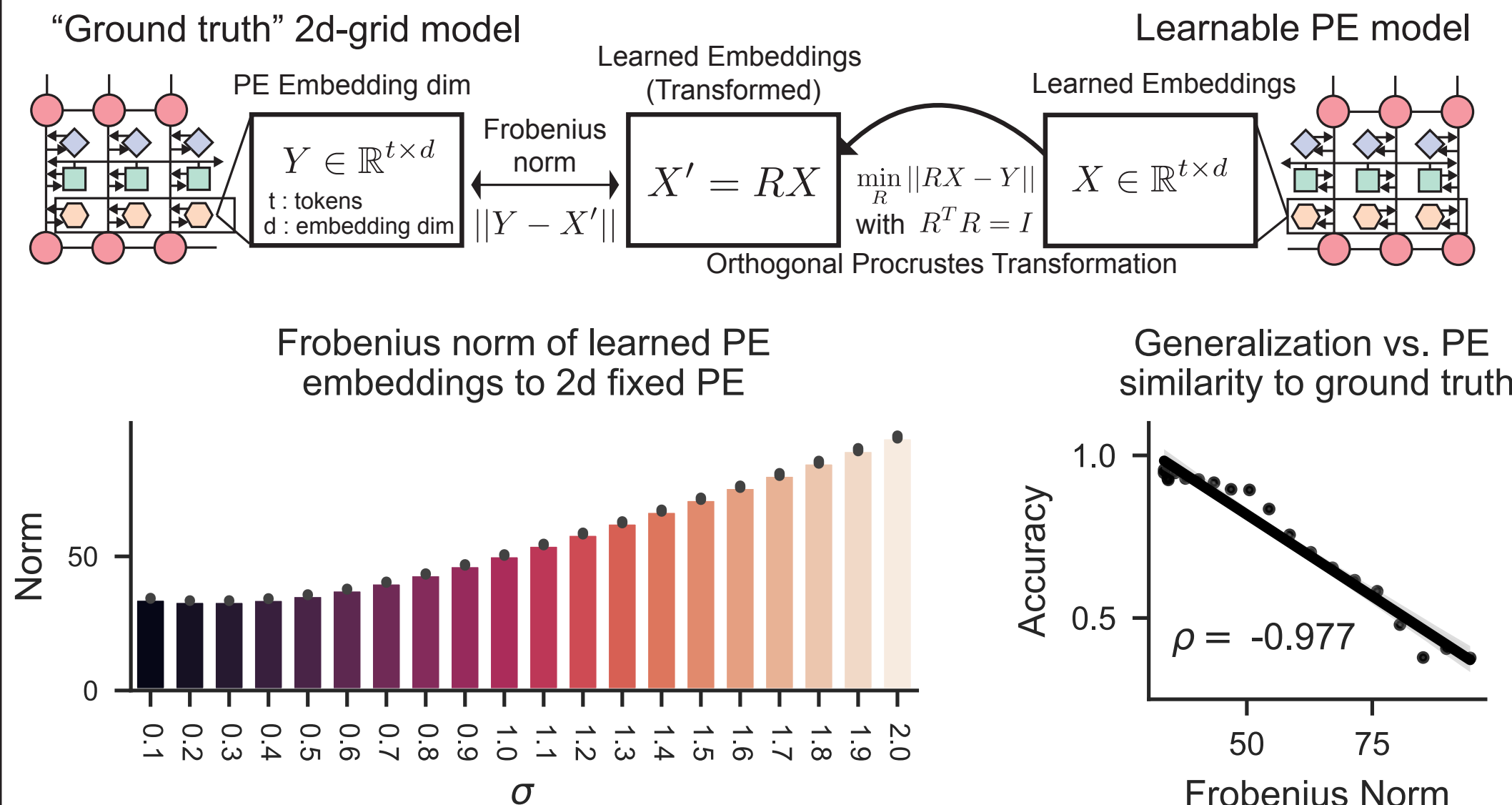
Transformer architecture

● token embedding
◆ mlp
■ bidirectional attention
⬡ learnable parameter, initialized from $N(0, \sigma I)$

**Generalization**

Enhanced test set generalization when training a learnable PE embedding initialized from small $\sigma$
Large $\sigma$: Lazy learning | Small $\sigma$: Rich learning

$\sigma = 0.2$  $\sigma = 0.5$  $\sigma = 1.0$  $\sigma = 1.3$  $\sigma = 1.5$  $\sigma = 2.0$

Generalization performance, learnable PE

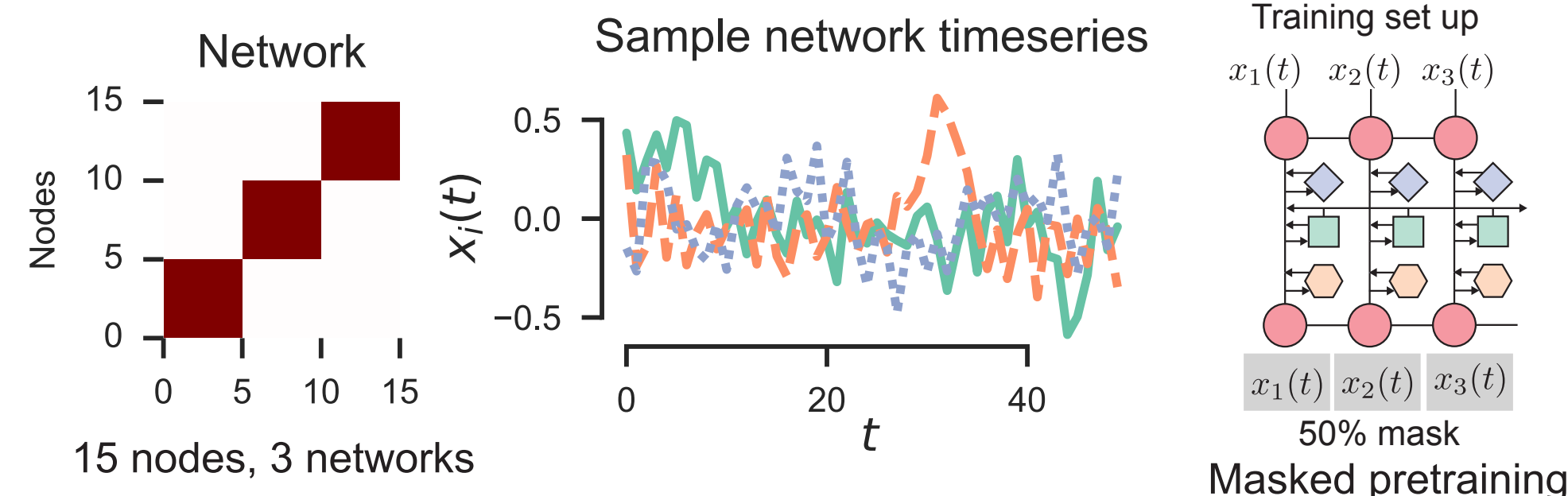$\sigma$ initialization

**Interpretability analyses**

Enhanced interpretability of PE embeddings when using a learnable PE embedding initialized with small $\sigma$
Learned PE embedding mimics a 2D grid structure, consistent with the LST input

"Ground truth" 2d-grid model

Learned Embeddings (Transformed)

$Y \in \mathbb{R}^{t \times d}$   t: tokens   d: embedding dim   Frobenius norm $||Y - X'||$

$X' = RX$  $\min_{R} ||RX - Y||$ with $R^T R = I$

$X \in \mathbb{R}^{t \times d}$  Learnable PE model

Orthogonal Procrustes Transformation

Frobenius norm of learned PE embeddings to 2d fixed PE

Generalization vs. PE similarity to ground truth

$\rho = -0.977$



## Experiment 2: Learnable PEs recover interpretable network clusters in network simulation
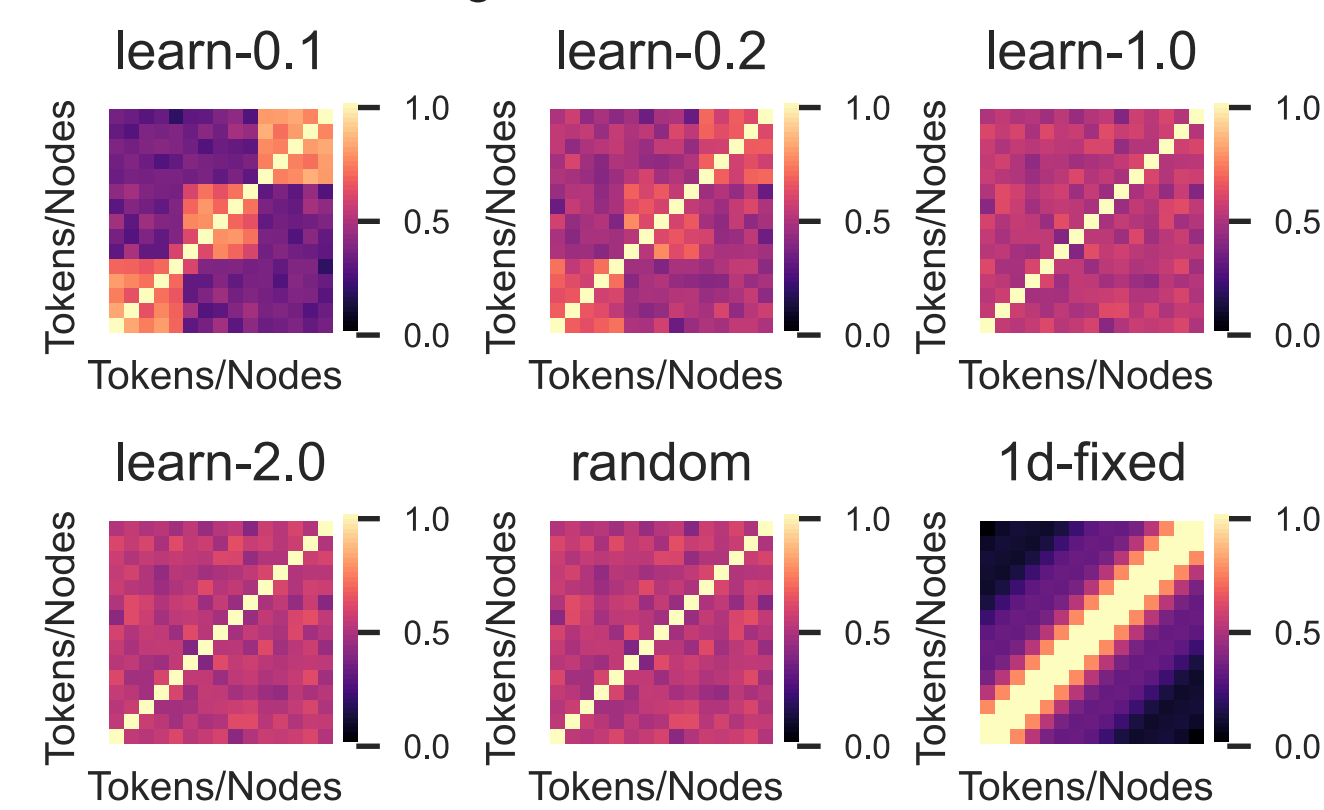
Experimental setup: Simulate nonlinear autoregressive network simulation

$$x_i(t) = \sum_{k=1}^{p} w_{i,k} \cdot x_i(t-k) + \sum_{j \in C_i, j \neq i} \lambda_{ij} \cdot f(x_j(t-1)) + \sum_{j \notin C_i} \eta_{ij} \cdot f(x_j(t-1)) + \epsilon_i(t)$$
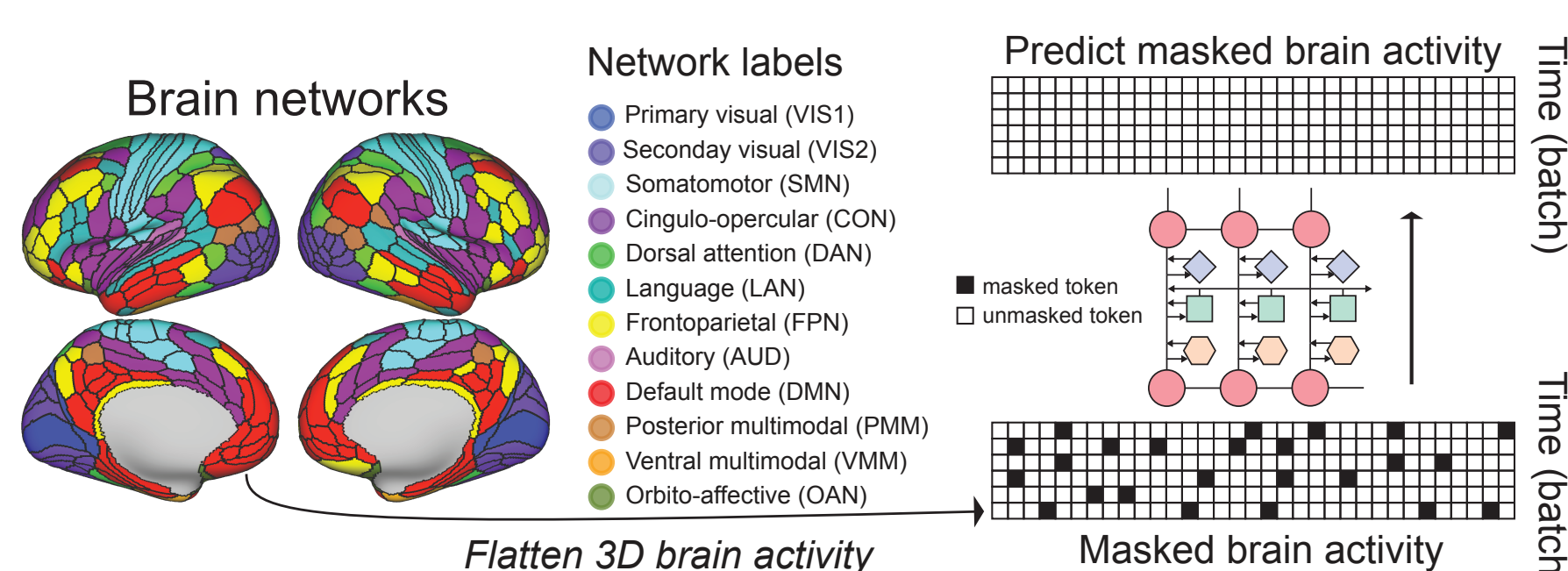
Network

15 nodes, 3 networks

Sample network timeseries

Training set up
$x_1(t)$ $x_2(t)$ $x_3(t)$

$x_1(t)$ $x_2(t)$ $x_3(t)$
50% mask
Masked pretraining

Interpretability analyses: Richly learned PE embeddings recover network clusters

* We measured the cosine distance of learned PE embeddings across a range of PE models

* Small $\sigma$ models (e.g., learn-0.1) accurately recovered the distance between PE embeddings

* Large $\sigma$ models (e.g., learn-1.0) did not

* Models with random and 1d-fixed PEs also did not exhibit network structure

learn-0.1  learn-0.2  learn-1.0
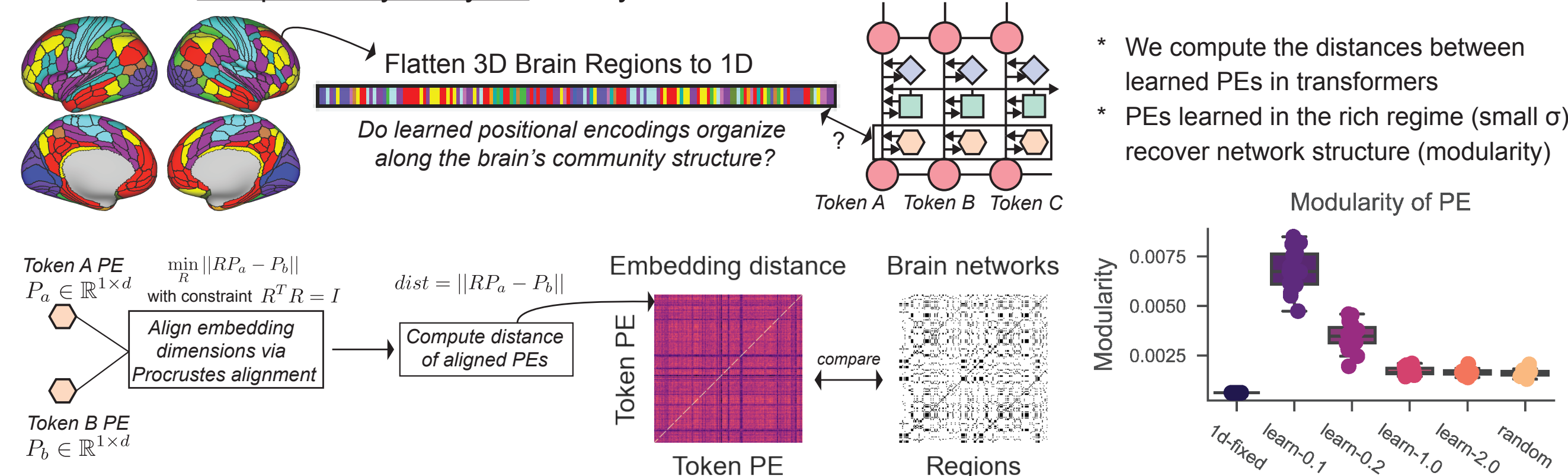learn-2.0  random  1d-fixed



## Experiment 3: Learnable PEs recover known network clusters in human brain fMRI data

Experimental setup: Predicting whole-brain activity from masked inputs

Brain networks

Network labels
● Primary visual (VIS1)
● Secondary visual (VIS2)
● Somatomotor (SMN)
● Cingulo-opercular (CON)
● Dorsal attention (DAN)
● Language (LAN)
● Frontoparietal (FPN)
● Auditory (AUD)
● Default mode (DMN)
● Posterior multimodal (PMM)
● Ventral multimodal (VMM)
● Orbito-affective (OAN)

Predict masked brain activity

Flatten 3D brain activity

Masked brain activity (tokens = brain regions)

Interpretability analyses: Richly learned PEs recover known functional network clusters

Flatten 3D Brain Regions to 1D

Do learned positional encodings organize along the brain's community structure?

* We compute the distances between learned PEs in transformers
* PEs learned in the rich regime (small $\sigma$) recover network structure (modularity)

Token A PE $P_a \in \mathbb{R}^{1 \times d}$   $\min_R ||RP_a - P_b||$ with constraint $R^T R = I$   Align embedding dimensions via Procrustes alignment   $dist = ||RP_a - P_b||$ Compute distance of aligned PEs

Token B PE $P_b \in \mathbb{R}^{1 \times d}$

Embedding distance   Brain networks

Modularity of PE



## Conclusion

* We extend prior transformer generalization studies from 1D sequences to n-dimensional sequences, which requires positional encoding schemes for higher dimensions.
* We demonstrate that **rich representation learning of positional encodings** – which is induced by initializing parameters with a small norm – **learns interpretable embeddings that also enhance generalization**