

Rotary Position Encodings for Graphs

ICML 2026 Spotlight

Isaac Reid*^{1,2} Arijit Sehanobish*³ Cederik Höfs*¹
Bruno Mlodozieniec^{1,4} Leonhard Vulpius¹ Federico Barbero^{5,2}
Adrian Weller^{1,6} Krzysztof Choromanski²
Richard E. Turner^{1,6} Petar Veličković^{2,1}

*Core contributors. ¹Cambridge ²Google DeepMind ³Independent
⁴MPI for Intelligent Systems ⁵Oxford ⁶Alan Turing Institute

arXiv:2509.22259 github.com/cederikhoefs/Graph-RoPE

Position encodings: same problem, three domains

- ▶ Sentences: **order of words**
- ▶ Images: **layout of patches**
- ▶ Graphs: **node connectivity**

Transformers are permutation invariant.
Position encodings inject structure into attention.

Rotary Position Encodings

Rotate $\mathbf{q}_i, \mathbf{k}_i$ per token by position-dependent angles:

$$\{\mathbf{q}_i, \mathbf{k}_i\} \mapsto \bigoplus_{n=1}^{d/2} \begin{pmatrix} \cos \theta_n & -\sin \theta_n \\ \sin \theta_n & \cos \theta_n \end{pmatrix} \{\mathbf{q}_i, \mathbf{k}_i\}, \quad \theta_n = \boldsymbol{\omega}_n^\top \mathbf{r}_i$$

Logit depends only on the *relative* position:

$$\mathbf{q}_i^\top \mathbf{k}_j \rightarrow \mathbf{q}_i^\top \text{RoPE}(\mathbf{r}_j - \mathbf{r}_i) \mathbf{k}_j$$

- ▶ Powers Llama, Gemma, modern ViTs

Can we apply RoPE to graphs?

Existing graph positional encodings

Two existing families of graph positional encodings to draw on:

- ▶ **Laplacian PE** (SAN, GraphGPS): *added* to each token

$$\mathbf{x}_i \rightarrow [\mathbf{x}_i, f(\{\lambda_k, \psi_k[i]\}_k)]$$

Spectral features carry rich structure; shown to help on long-range tasks

- ▶ **Structural-distance bias** (Graphormer, GraphiT): *added* to the logit

$$\mathbf{q}_i^\top \mathbf{k}_j + b(\text{SPD}, R, \dots)$$

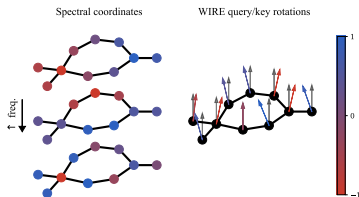
Strong inductive bias, but instantiates $N \times N$ attention \Rightarrow incompatible with linear attention

Idea: feed Laplacian features into RoPE \Rightarrow structure as rotation, $\mathcal{O}(N)$ scaling.

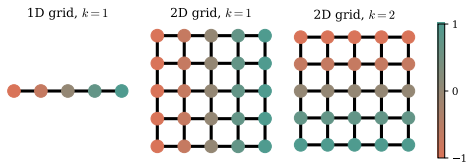
WIRE: Wave-Induced Rotary Encodings

Algorithm

1. Diagonalise Laplacian $L = D - A = U\Lambda U^\top$
2. Spectral coordinates from lowest m eigenvectors: $\mathbf{r}_i = [\mathbf{u}_k[i]]_{k=0}^{m-1} \in \mathbb{R}^m$
3. Rotate $\mathbf{q}_i, \mathbf{k}_i$ by $\theta_n = \omega_n^\top \mathbf{r}_i$; learnable $\{\omega_n\}_{n=1}^{d/2} \subset \mathbb{R}^m$



Theory I: RoPE \subset WIRE on grids



- ▶ Path graph P_N : first nontrivial eigenvector

$$\mathbf{u}_1[i] = -\cos\left(\frac{\pi}{N}\left(i + \frac{1}{2}\right)\right)$$

is monotone in token index

- ▶ $\omega_n = [0, \omega_n, 0, \dots] \Rightarrow$ **classical RoPE** (LLMs)
- ▶ 2D grids factorise \Rightarrow **ViT RoPE**

Theory II: WIRE depends on effective resistance

Normalised features $\mathbf{r}_i = [\mathbf{u}_k[i]/\sqrt{\lambda_k}]$, random $\omega_n \sim \mathcal{N}(0, \omega \mathbf{I})$:

$$\mathbb{E} \left[(\text{RoPE}(\mathbf{r}_i) \mathbf{q}_i)^\top \text{RoPE}(\mathbf{r}_j) \mathbf{k}_j \right] = \mathbf{q}_i^\top \mathbf{k}_j \left(1 - \frac{\omega^2}{2} R(i, j) \right) + \mathcal{O}(\omega^4).$$

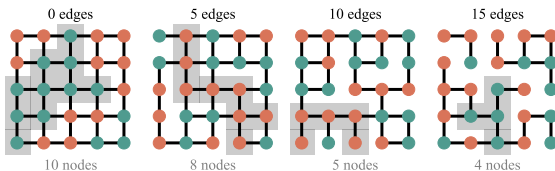
$R(i, j) = \mathbf{L}_{ii}^\dagger + \mathbf{L}_{jj}^\dagger - 2\mathbf{L}_{ij}^\dagger$: effective resistance, a metric on \mathcal{N} that lower-bounds SPD.

Why this is remarkable

- ▶ Topological masking by structural distance
- ▶ **No $N \times N$ attention matrix required**
- ▶ Limit is **gauge invariant**

Synthetic: monochromatic subgraphs

Predict size of largest monochromatic connected subgraph on a coloured 5×5 grid with random edges deleted.



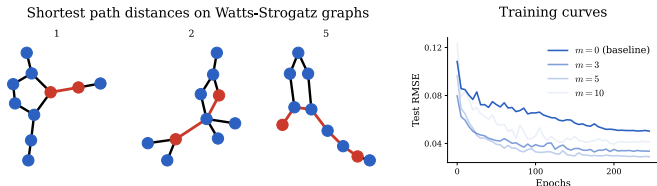
Edge deletion **interpolates grid** \rightarrow **irregular topology**; stress-tests WIRE away from the ViT-RoPE limit.

m	Test RMSE \downarrow (deleted edges)			
	0	5	10	15
0 (no WIRE)	0.060(1)	0.087(1)	0.081(1)	0.068(2)
3	0.053(2)	0.075(2)	0.072(3)	0.064(3)
5	0.057(2)	0.075(1)	0.070(2)	0.056(4)
10	0.055(2)	0.068(5)	0.063(2)	0.058(2)

Synthetic: shortest-path distance

Watts–Strogatz, $N = 10$, $k = 2$, $p = 0.6$; predict SPD between marked nodes.

Purely **topological** target; isolates WIRE's resistive-distance bias (Theorem II).



m	0 (base)	3	5	10
Test RMSE ↓	0.065(5)	0.048(6)	0.038(6)	0.045(4)

WIRE nearly **halves** the error.

Benchmarks: WIRE Performers close the gap

- ▶ Drop-in for GraphGPS + linear (Performer) attention
- ▶ Often matches $\mathcal{O}(N^2)$ softmax baselines

Dataset	Performer $\mathcal{O}(N)$		Transformer $\mathcal{O}(N^2)$
	Baseline	+ WIRE	Baseline
MNIST \uparrow	97.56(2)	98.10(1)	98.05(4)
CIFAR10 \uparrow	70.61(4)	71.15(3)	72.3(1)
PATTERN \uparrow	85.71(3)	86.63(6)	86.69(2)
CLUSTER \uparrow	76.90(3)	77.53(3)	78.02(6)
ogbg-molhiv \uparrow	0.776(2)	0.785(2)	0.788(1)
ogbg-molpcba \uparrow	0.238(3)	0.264(1)	0.291(3)
Peptides-struct \downarrow	0.2616(4)	0.2566(4)	0.2500(5)
MalNet-Tiny \uparrow	92.81(5)	93.46(2)	93.36(6)

MalNet-Tiny: WIRE Performer matches the quadratic Transformer; trains on a single T4 12 GB GPU.

Takeaway

WIRE: RoPE for graphs via Laplacian spectra

- ▶ Recovers standard RoPE on grids (LLMs, ViTs)
- ▶ Asymptotically downweights attention by **effective resistance**
- ▶ Stays $\mathcal{O}(N)$; fully compatible with linear attention
- ▶ $< 1\%$ of parameters; consistent gains across > 200 trained models

Paper: [arXiv:2509.22259](https://arxiv.org/abs/2509.22259) Code: github.com/cederikhoefs/Graph-RoPE

Come find us at the poster!