

The background is a gradient from dark purple to blue, speckled with white dots. On the left side, there are several concentric circles and a large arc with a frequency scale ranging from 140 to 260. Smaller circular patterns with arrows are scattered across the left half of the image.

# SPECTRAL ATTENTION NETWORKS

# Meet the authors!



Devin Kreuzer



Dominique Beaini

Prof. William L. Hamilton

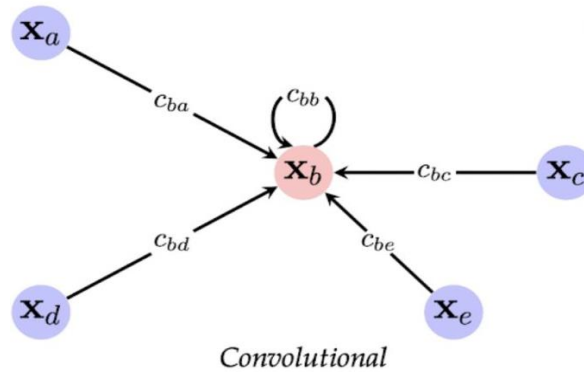
Vincent Letourneau

Prudencio Tossou

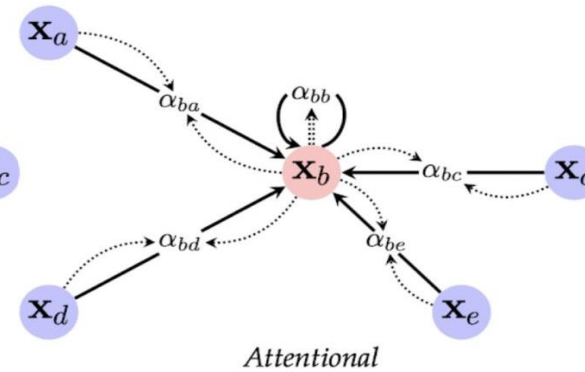


# Preliminaries

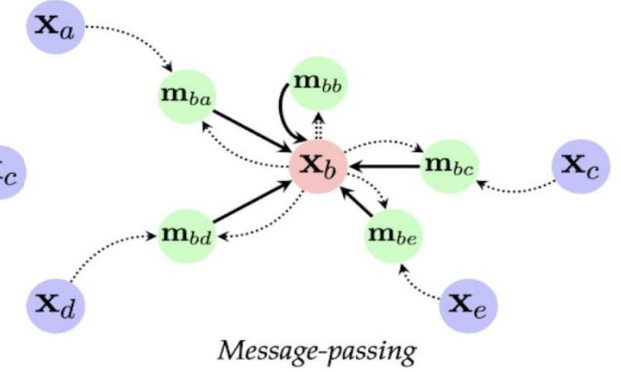
Current GNNs strategies:



$$\mathbf{h}_i = \phi \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} c_{ij} \psi(\mathbf{x}_j) \right)$$



$$\mathbf{h}_i = \phi \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} a(\mathbf{x}_i, \mathbf{x}_j) \psi(\mathbf{x}_j) \right)$$

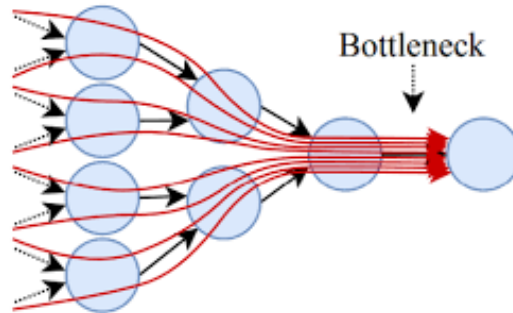


$$\mathbf{h}_i = \phi \left( \mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j) \right)$$



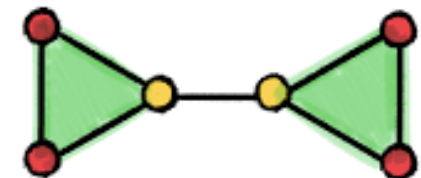
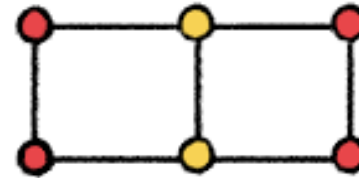
Image from Petar Velickovic, DeepMind

Limitations:



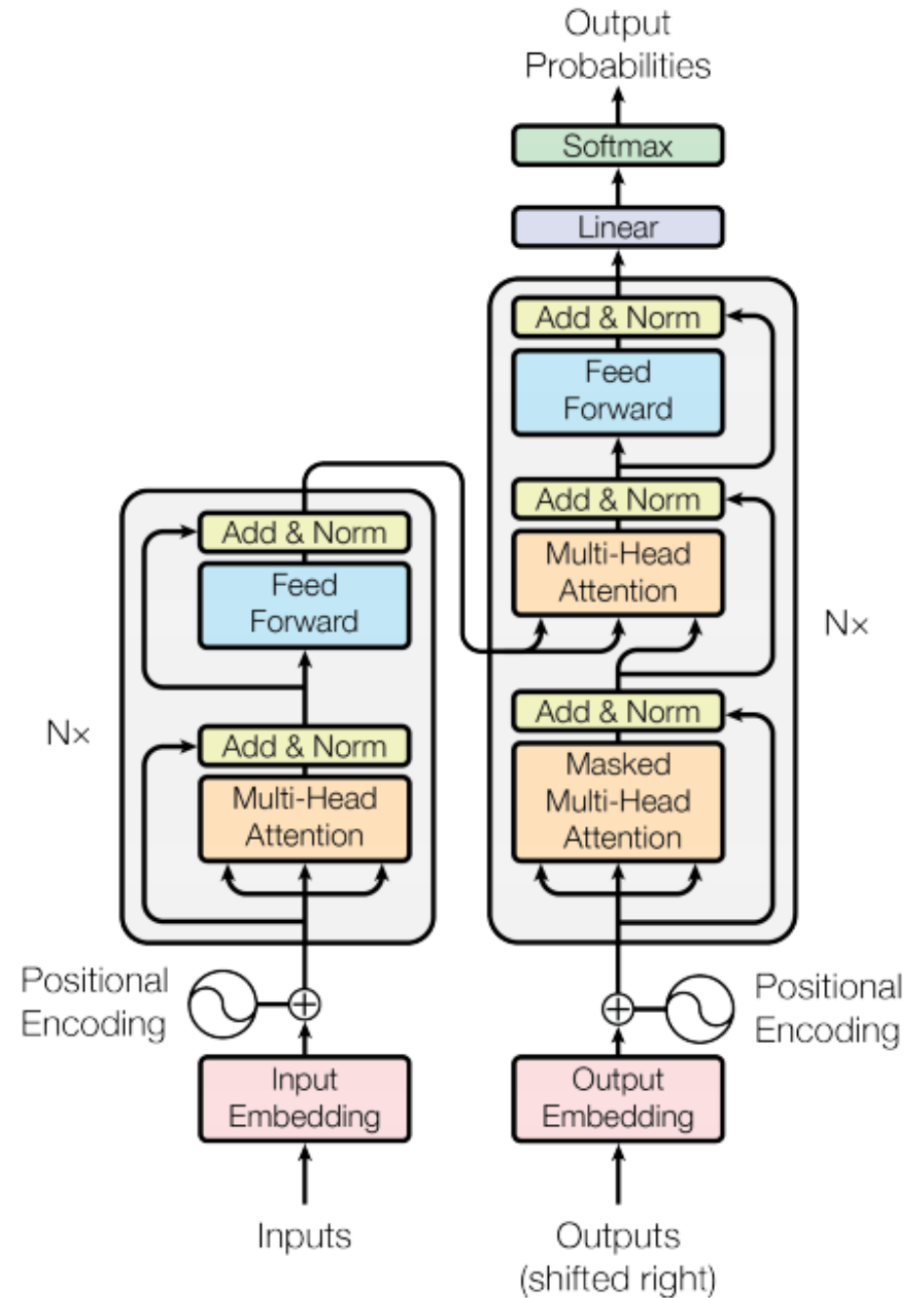
(b) The bottleneck of graph neural networks

Image from Uri Alon and Eran Yahav, Technion



# Transformer

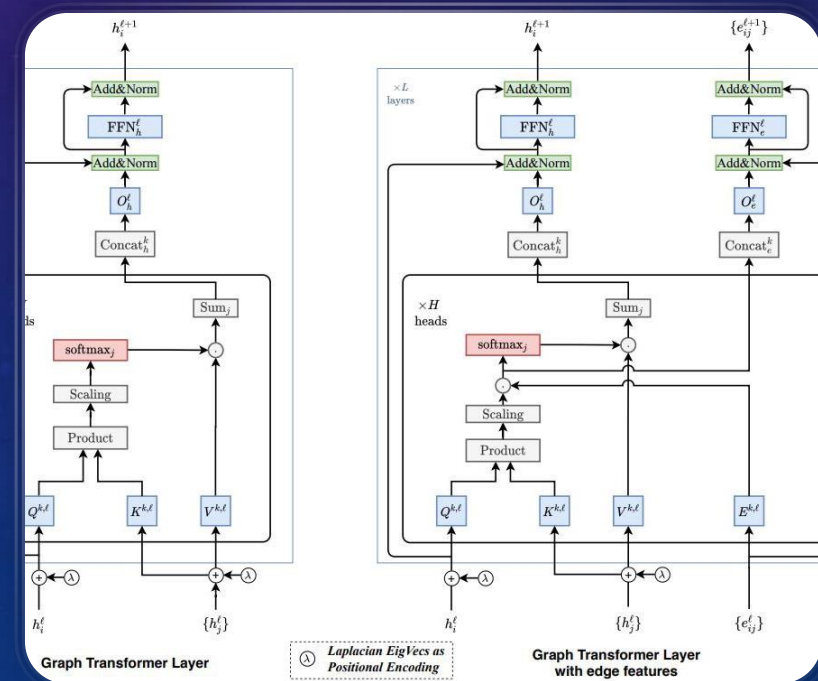
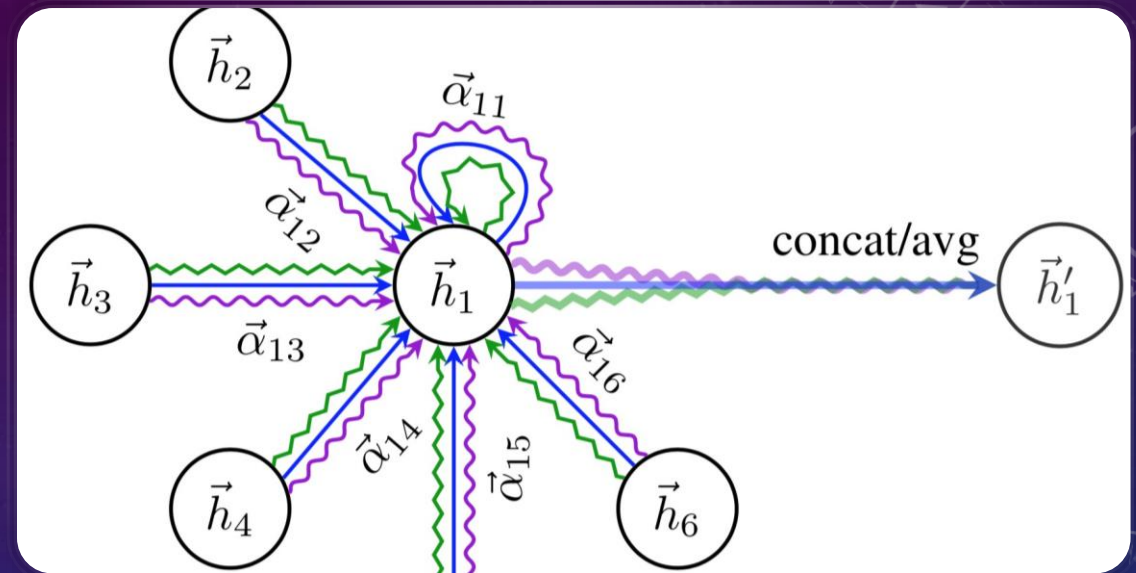
- Revolutionized NLP by removing many inductive biases
- **Main challenge: Positional Encodings for arbitrary graphs**
- Answer: graph Laplacian





# PREVIOUS WORK

- We expand primarily on:
  1. GAT [1]
  2. GT [2]

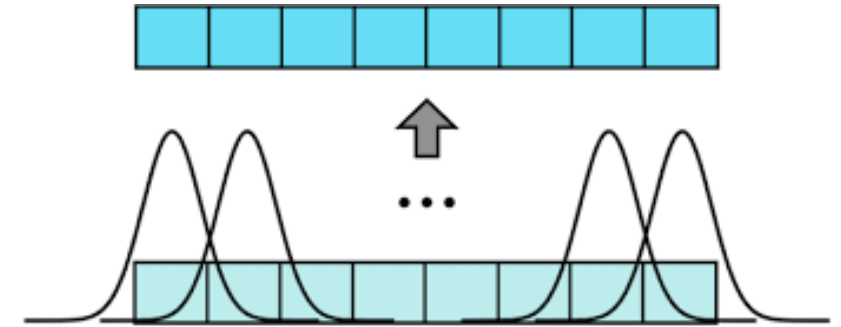


[1] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. arXiv preprint arXiv:1710.10903, 2017

[2] Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs, 2020

# Positional Encodings

- Original Transformer PEs
  - Frequencies are pre-determined
- In arbitrary graphs, this process isn't feasible
  - No natural 'pos' value
  - Equivalent given by eigenvectors  $\phi$  of the Laplacian  $L$
  - $L = D - A$ ,  $L_{sym} = D^{-1/2} L D^{-1/2}$



$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

# Some basic Mathematical intuition

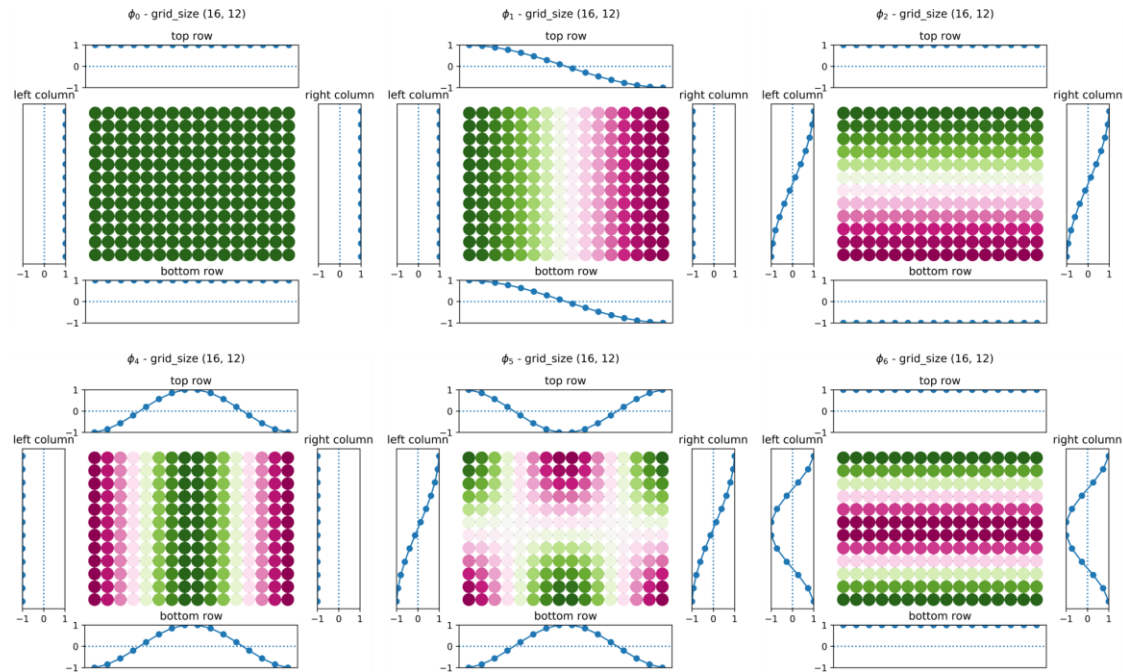
- $\phi$  : function over domain of nodes
- $N_i$ : neighborhood of node  $i$
- $L$ : operator on functions over nodes

$$(L\phi)_i = |N_i| \phi_i - \sum_{j \in N_i} \phi_j$$
$$\phi^T L \phi = \left[ |N_1| \phi_1^2 - \sum_{j \in N_1} \phi_j^2 \right] - \left[ |N_2| \phi_2^2 - \sum_{j \in N_2} \phi_j^2 \right] - \dots = \lambda$$

- The eigenvectors  $\phi$  assign values to each node which minimize the quantity  $\lambda$ 
  - $\phi$  are sine/cosine functions traversing the graph with frequency  $\lambda$
  - Can as positional encodings!

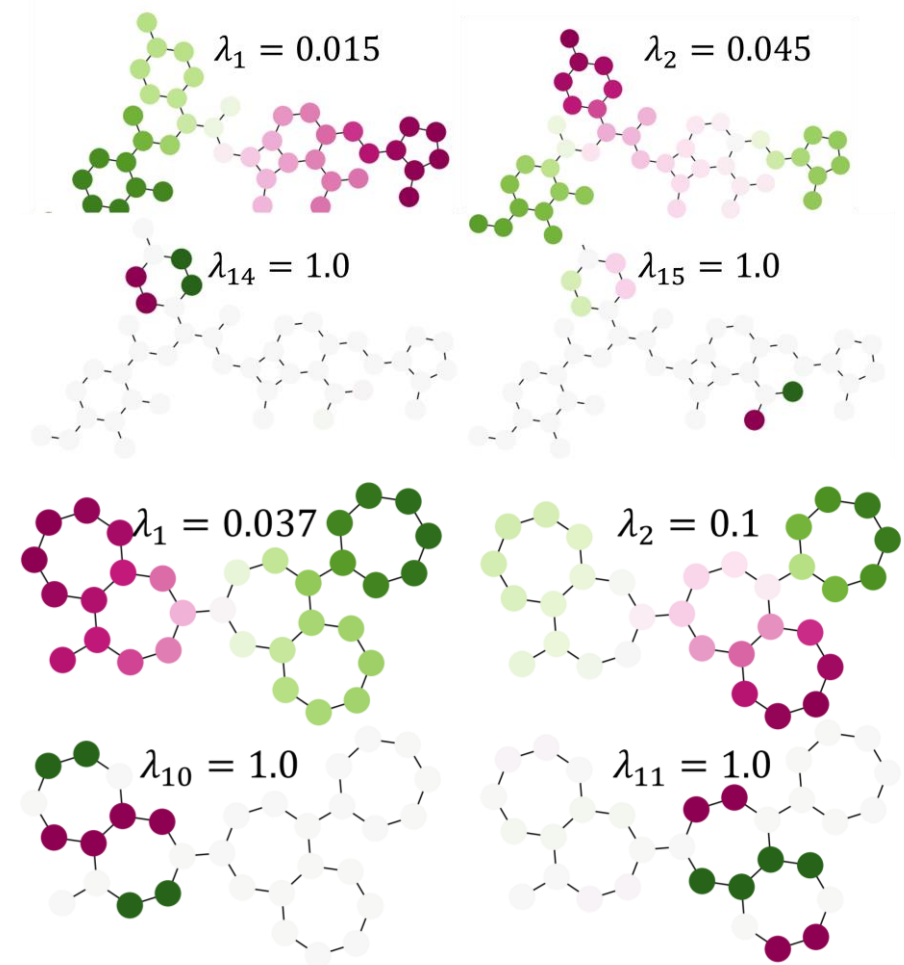
# Laplace eigenvectors are sine-wave equivalents

## Eigenvectors of a grid



$\lambda$ : Eigenvalue, represents the frequency of the wave  
 $\phi$ : Eigenvector, represents the sine/cosine function

## Eigenvectors of a molecule





The background features a dark blue gradient with faint, light blue geometric patterns. On the left side, there are several concentric circles and arcs, some of which are marked with degree values ranging from 40 to 260. These markings are arranged in a way that suggests a circular scale or a series of orbits. The overall aesthetic is technical and scientific.

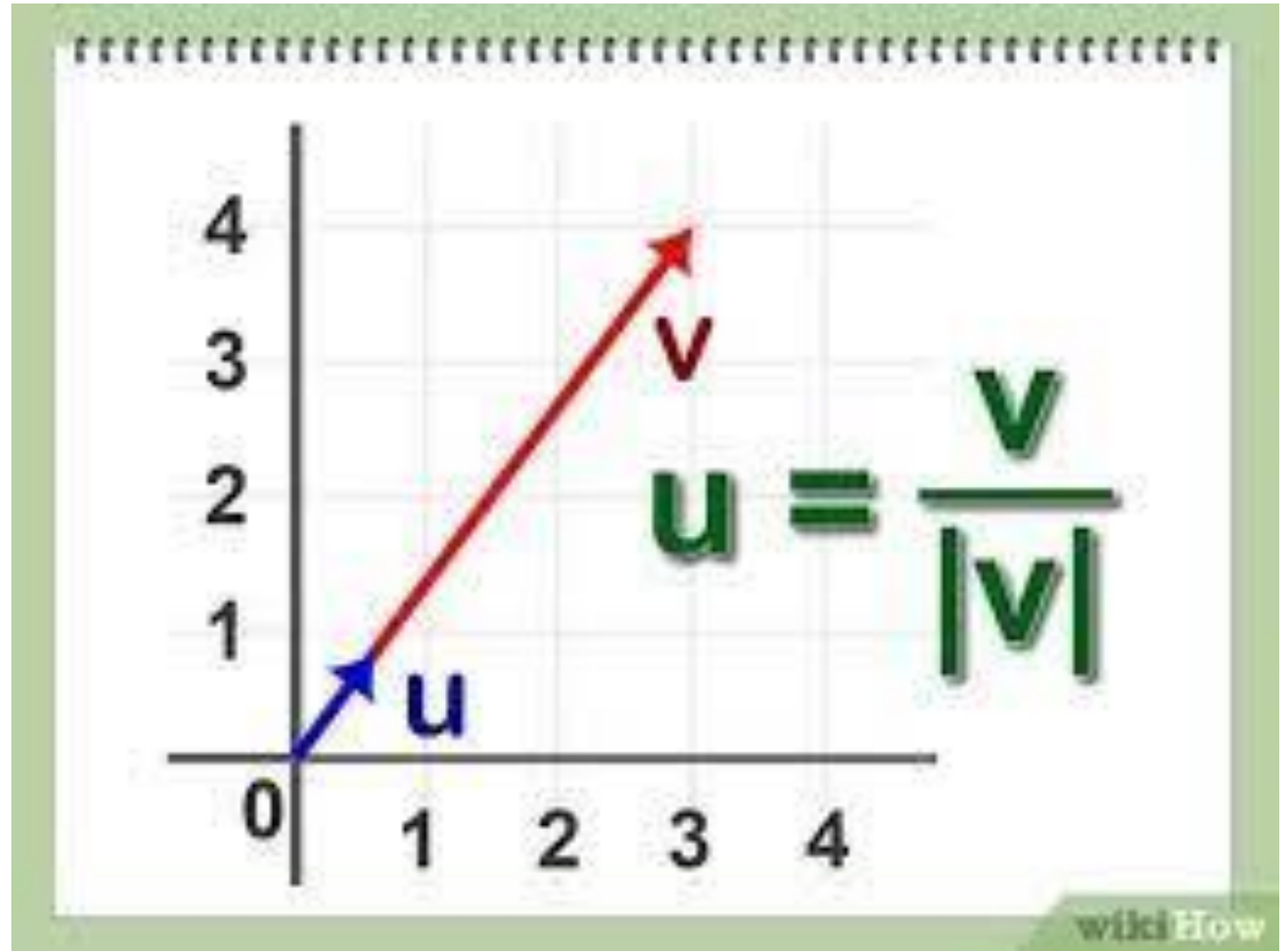
# LAPLACE EIGENFUNCTIONS

## *ETIQUETTE*

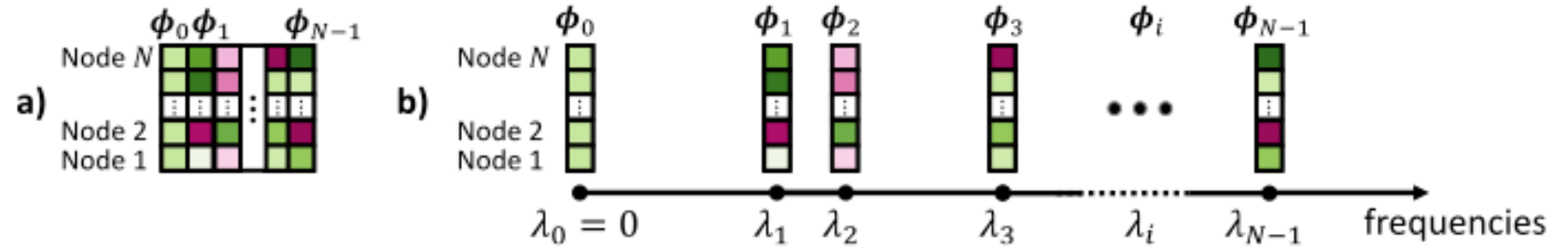
KEY CHALLENGES IN DIRECTLY CONCATENATING/ADDING  $\phi$  AS PE

## *Etiquette #1:* Normalization

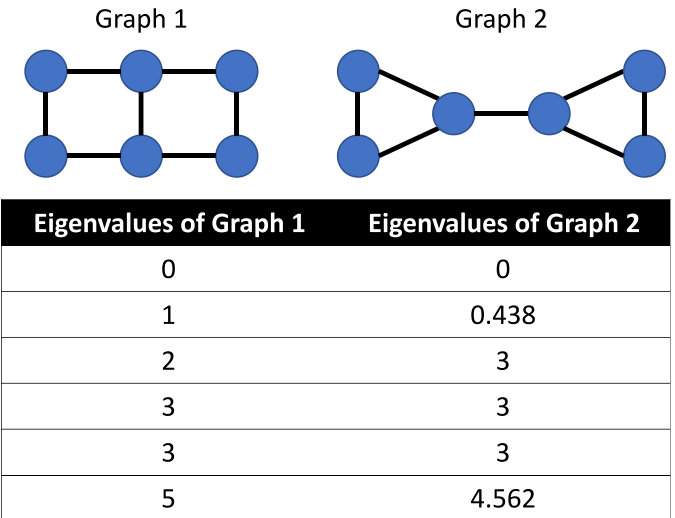
- Given  $\phi$ ,  $c\phi$  is also a valid option
- Must pick one vector in the span of  $\phi$
- Select normalized eigenvectors:  $\langle \phi, \phi \rangle = 1$



## Etiquette #2: Eigenvalues



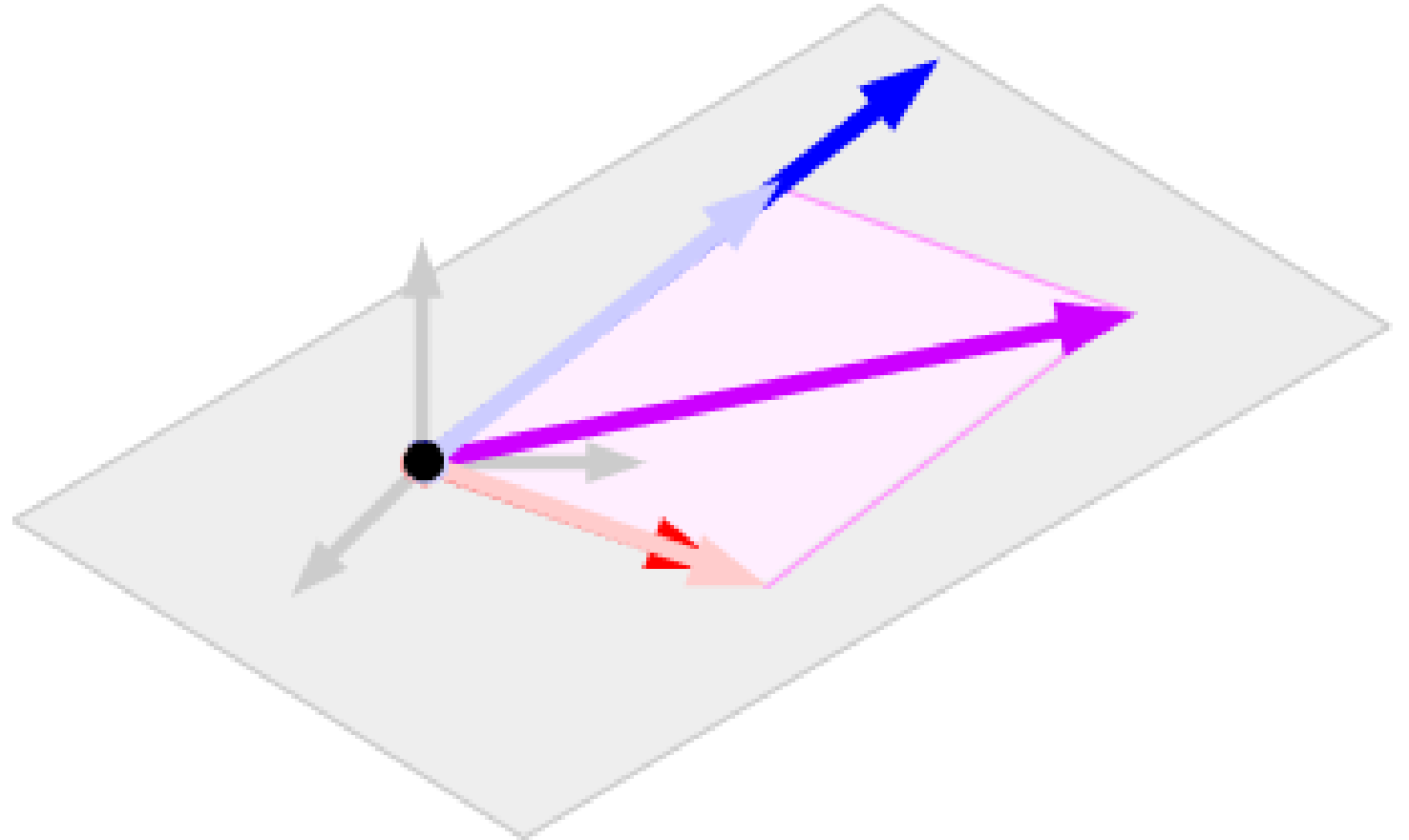
- An eigenvalue supplies valuable information about its eigenvector
- Eigenvalue spectrum provides important information about graph structure
- Important physical properties depend on the eigenvalues:



$$d_D^2(j_1, j_2) = \sum_{k>0} e^{-2t\lambda_i} (\phi_{i,j_1} - \phi_{i,j_2})^2 \quad , \quad d_B^2(j_1, j_2) = \sum_{i>0} \frac{(\phi_{i,j_1} - \phi_{i,j_2})^2}{\lambda_i^2}$$

## *Etiquette #3:* Multiplicities

- Two eigenvectors can have the same eigenvalue
- Need to pick a vector in the space spanned by both!
  - Even more challenging!



## *Etiquette #4: Variable number of eigenvectors*

- A graph  $G_i$  can have at most  $N_i$  linearly independent eigenvectors
- Prior elected to select a fixed number  $k$  eigenvectors for each graph, where  $k \leq N_i, \forall i$
- Major bottleneck when the smallest graphs  $\ll$  largest graphs in the dataset  $s$
- Need for model that uses fixed dimension PE that captures entire spectrum

## *Etiquette #5: Sign ambiguity*

- Given  $\phi$ ,  $-\phi$  is also a valid
- Left with  $2^k$  possible combination of valid eigenvectors
- Previous work proposed data augmentation: random sign flipping



The background is a gradient of deep blue and purple, speckled with white dots resembling a starry sky. Overlaid on this are several faint, white, circular and semi-circular patterns. Some of these patterns have tick marks and numbers, suggesting a circular scale or a compass. There are also curved arrows indicating a direction of movement or flow. The overall aesthetic is technical and futuristic.

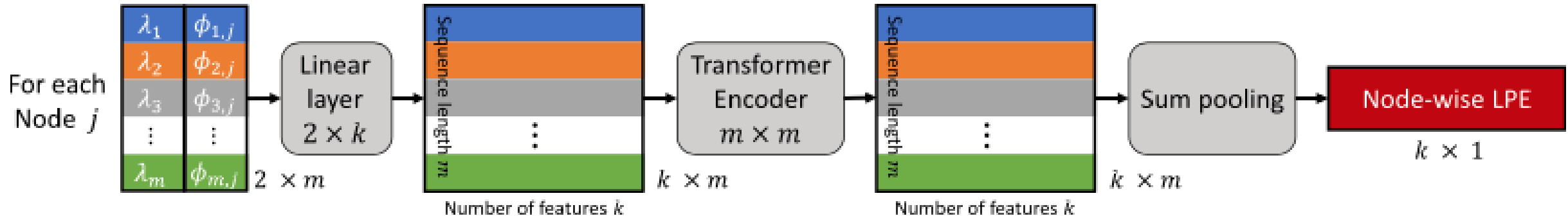
# MODEL ARCHITECTURE

# OVERVIEW

1. Pre-compute first  $m$  eigenvectors and fully-connect each graph (noting the added edges)
2. Learn positional encodings that address the *etiquettes* previously described
3. Concatenate the learned positional encodings (LPE) to the input node features
4. Apply the Main Transformer

# LPE over Nodes

- Treat Laplace spectrum as a variable-length sequence where  $\lambda$  is the position of  $\phi$
- A Transformer is naturally suited to learn PEs

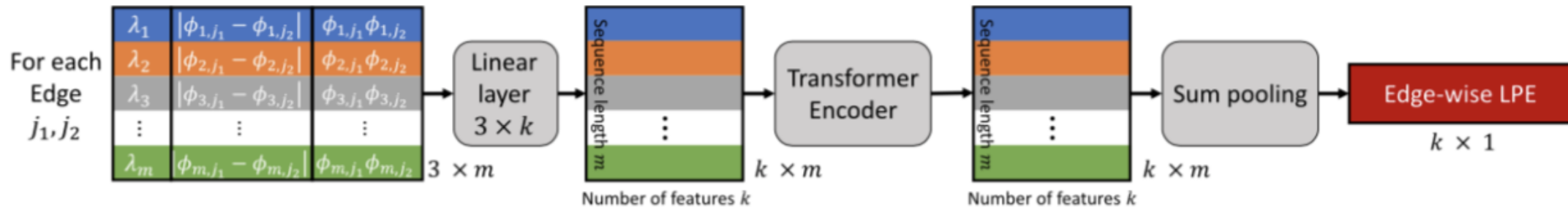


# LPE over Nodes

- This technique addresses *etiquettes #1-4*
  - L2 normalized eigenvectors
  - Eigenvalues paired with eigenvectors
  - Number of eigenvectors treated as a variable
  - Aware of multiplicities
- *Etiquette #5* is unaddressed
  - Randomly flip sign of eigenvectors during training
  - Resolved in the LPE Over Edges

# LPE Transformer over edges

- Very important notion:
  - $\phi$  on its own reveals no information at all
  - $\Delta\phi$  tells us everything!
  - Just like potential energy
- Laplace PEs are more naturally represented as edge features



- Downside: Huge computational burden
  - Extra factor of  $N$  if considering fully connected graphs

# Main Graph Transformer

- Based on previous work
  - Multi-head attention:

$$\hat{\mathbf{h}}_i^{l+1} = \mathbf{O}_h^l \parallel \left( \sum_{k=1}^H \sum_{j \in V} w_{ij}^{k,l} \mathbf{V}^{k,l} \mathbf{h}_j^l \right)$$

$$\hat{\mathbf{h}}^{l+1} = \text{Norm}(\mathbf{h}_i^l + \hat{\mathbf{h}}_i^{l+1}), \quad \hat{\hat{\mathbf{h}}}_i^{l+1} = \mathbf{W}_2^l \text{ReLU}(\mathbf{W}_1^l \hat{\mathbf{h}}_i^{l+1}), \quad \mathbf{h}_i^{l+1} = \text{Norm}(\hat{\mathbf{h}}^{l+1} + \hat{\hat{\mathbf{h}}}_i^{l+1})$$

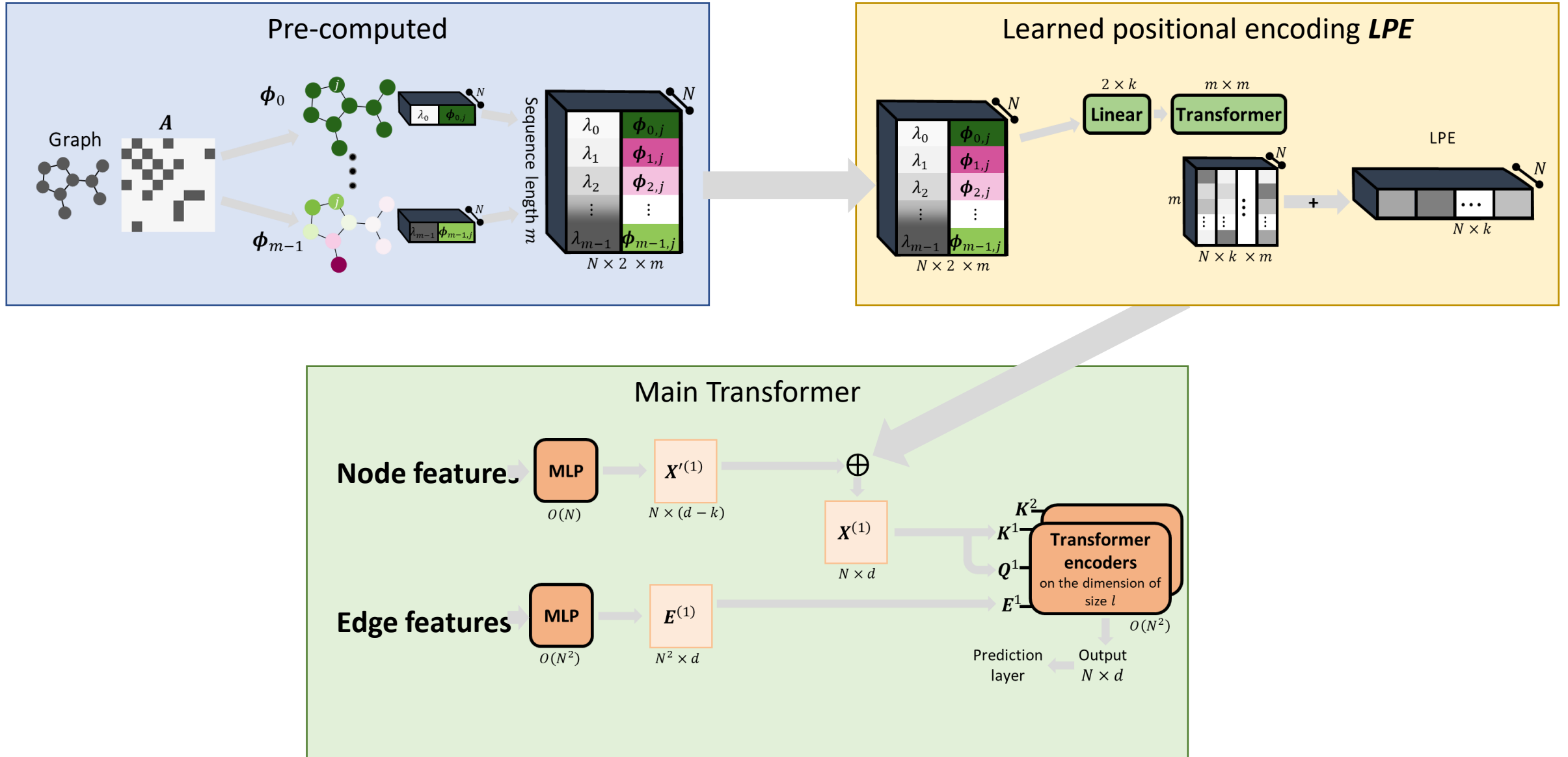
- Two key additions
  - Separate** weights for real and added edges
  - Hyperparameter that tunes the weight between local and global attention

$$\hat{w}_{ij}^{k,l} = \left\{ \begin{array}{ll} \frac{Q^{1,k,l} \mathbf{h}_i^l \circ \mathbf{K}^{1,k,l} \mathbf{h}_j^l \circ \mathbf{E}^{1,k,l} \mathbf{e}_{ij}}{\sqrt{d_k}} & \text{if } i \text{ and } j \text{ are connected in sparse graph} \\ \frac{Q^{2,k,l} \mathbf{h}_i^l \circ \mathbf{K}^{2,k,l} \mathbf{h}_j^l \circ \mathbf{E}^{2,k,l} \mathbf{e}_{ij}}{\sqrt{d_k}} & \text{otherwise} \end{array} \right\}$$

$$w_{ij}^{k,l} = \left\{ \begin{array}{ll} \frac{1}{1+\gamma} \cdot \text{softmax}(\sum_{d_k} \hat{w}_{ij}^{k,l}) & \text{if } i \text{ and } j \text{ are connected in sparse graph} \\ \frac{\gamma}{1+\gamma} \cdot \text{softmax}(\sum_{d_k} \hat{w}_{ij}^{k,l}) & \text{otherwise} \end{array} \right\}$$



# Overview



# LIMITATIONS

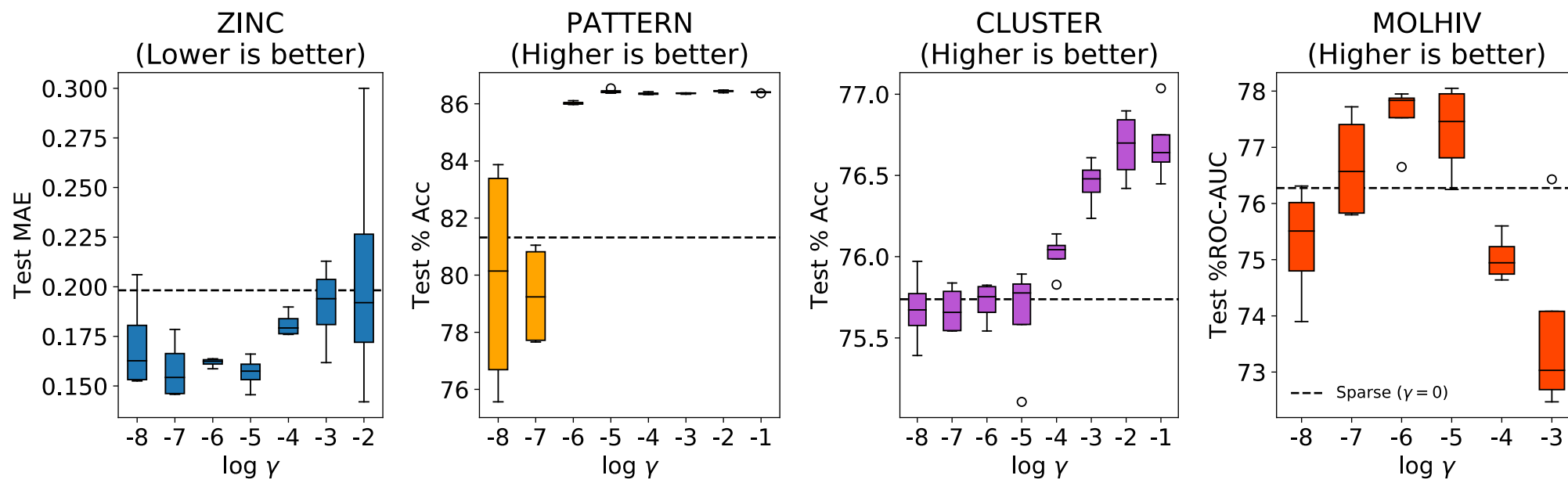
- Node LPE does not address sign invariance issue
- Computational complexity
  - Node LPE:  $O(N^3)$  if considering all eigenfunctions
  - Edge LPE:  $O(N^4)$  if considering all eigenfunctions and fully-connected graphs
  - Main Graph Transformer:  $O(N^2)$

The background is a gradient of dark blue and purple, speckled with small white dots resembling stars. Overlaid on this are several faint, white, circular and semi-circular patterns. Some of these patterns have tick marks and numbers, suggesting a circular scale or a clock face. The numbers visible include 40, 150, 160, 170, 180, 190, 200, 210, 220, 230, 240, 250, and 260. There are also curved arrows indicating a direction of movement or rotation.

# EXPERIMENTS

# Ablation studies


- Study the effect of hyperparameter  $\gamma$



# Ablation studies

- Compare sparse to full attention

| Model details |      | ZINC              | PATTERN            | CLUSTER            | MOLHIV           |
|---------------|------|-------------------|--------------------|--------------------|------------------|
| Attention     | LPE  | MAE               | % ACC              | % ACC              | % ROC-AUC        |
| Sparse        | -    | $0.267 \pm 0.032$ | $83.613 \pm 0.663$ | $75.683 \pm 0.098$ | $73.46 \pm 0.71$ |
| Sparse        | Node | $0.198 \pm 0.004$ | $81.329 \pm 2.150$ | $75.738 \pm 0.106$ | $76.61 \pm 0.62$ |
| Full          | -    | $0.392 \pm 0.055$ | $86.322 \pm 0.049$ | $76.447 \pm 0.177$ | $73.84 \pm 1.80$ |
| Full          | Node | $0.157 \pm 0.006$ | $86.441 \pm 0.040$ | $76.691 \pm 0.247$ | $77.57 \pm 0.61$ |

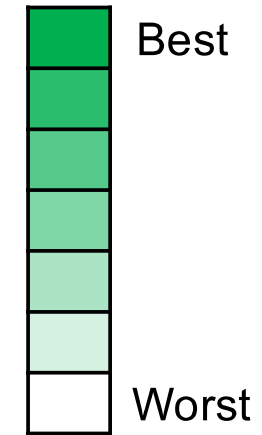


Best

Worst

# Comparison to SOTA

|                 | ZINC                                | PATTERN                              | CLUSTER                              | MOLHIV                             |
|-----------------|-------------------------------------|--------------------------------------|--------------------------------------|------------------------------------|
| Model           | MAE                                 | % ACC                                | % ACC                                | % ROC-AUC                          |
| GCN             | 0.367 $\pm$ 0.011                   | 71.892 $\pm$ 0.334                   | 68.498 $\pm$ 0.976                   | 76.06 $\pm$ 0.97                   |
| GraphSage       | 0.398 $\pm$ 0.002                   | 50.492 $\pm$ 0.001                   | 63.844 $\pm$ 0.110                   | -                                  |
| GatedGCN        | 0.282 $\pm$ 0.015                   | 85.568 $\pm$ 0.088                   | 73.840 $\pm$ 0.326                   | -                                  |
| GatedGCN-PE     | 0.214 $\pm$ 0.013                   | 86.508 $\pm$ 0.085                   | 76.082 $\pm$ 0.196                   |                                    |
| GIN             | 0.526 $\pm$ 0.051                   | 85.387 $\pm$ 0.136                   | 64.716 $\pm$ 1.553                   | 75.58 $\pm$ 1.40                   |
| PNA             | 0.142 $\pm$ 0.010                   | -                                    | -                                    | 79.05 $\pm$ 1.32                   |
| DGN             | -                                   | -                                    | -                                    | <b>79.70 <math>\pm</math> 0.97</b> |
| Attention-based |                                     |                                      |                                      |                                    |
| GAT             | 0.384 $\pm$ 0.007                   | 78.271 $\pm$ 0.186                   | 70.587 $\pm$ 0.447                   | -                                  |
| GT (sparse)     | 0.226 $\pm$ 0.014                   | 84.808 $\pm$ 0.068                   | 73.169 $\pm$ 0.662                   | -                                  |
| GT (full)       | 0.598 $\pm$ 0.049                   | 56.482 $\pm$ 3.549                   | 27.121 $\pm$ 8.471                   | -                                  |
| SAN (ours)      | <b>0.139 <math>\pm</math> 0.006</b> | <b>86.581 <math>\pm</math> 0.037</b> | <b>76.691 <math>\pm</math> 0.247</b> | 77.85 $\pm$ 0.65                   |





# CONCLUSION

- First successful fully-connected GNN model to perform well
- Outperforms other Attention-based models by wide margin
- Large computational bottleneck
  - Need to implement variations that scale linearly or logarithmically to enable Edge LPE



# THANK YOU!!!

Follow us on Twitter!

@dom\_beaini

@KreuzerDevin

@valence\_ai