

# Pure Transformers are Powerful Graph Learners

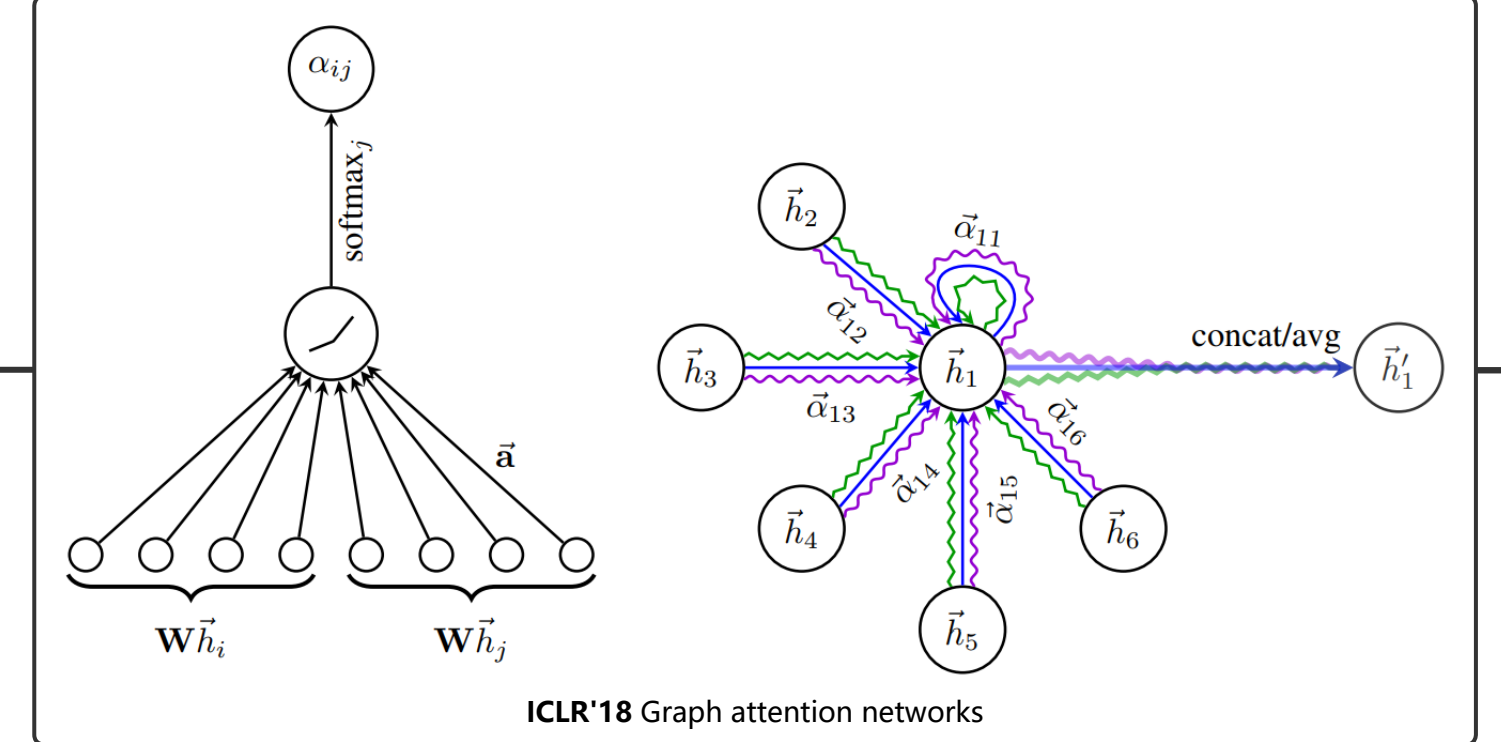
Authors

Jinwoo Kim<sup>1\*</sup>, Tien Dat Nguyen<sup>1</sup>, Seunwoo Min<sup>2</sup>, Sungjun Cho<sup>2</sup>,  
 Montae Lee<sup>2,3</sup>, Honglak Lee<sup>2,3</sup>, Sungchun Hong<sup>2,3</sup>  
<sup>1</sup>KAIST, <sup>2</sup>LG AI Research, <sup>3</sup>University of Illinois Chicago, <sup>4</sup>University of Michigan

## Background

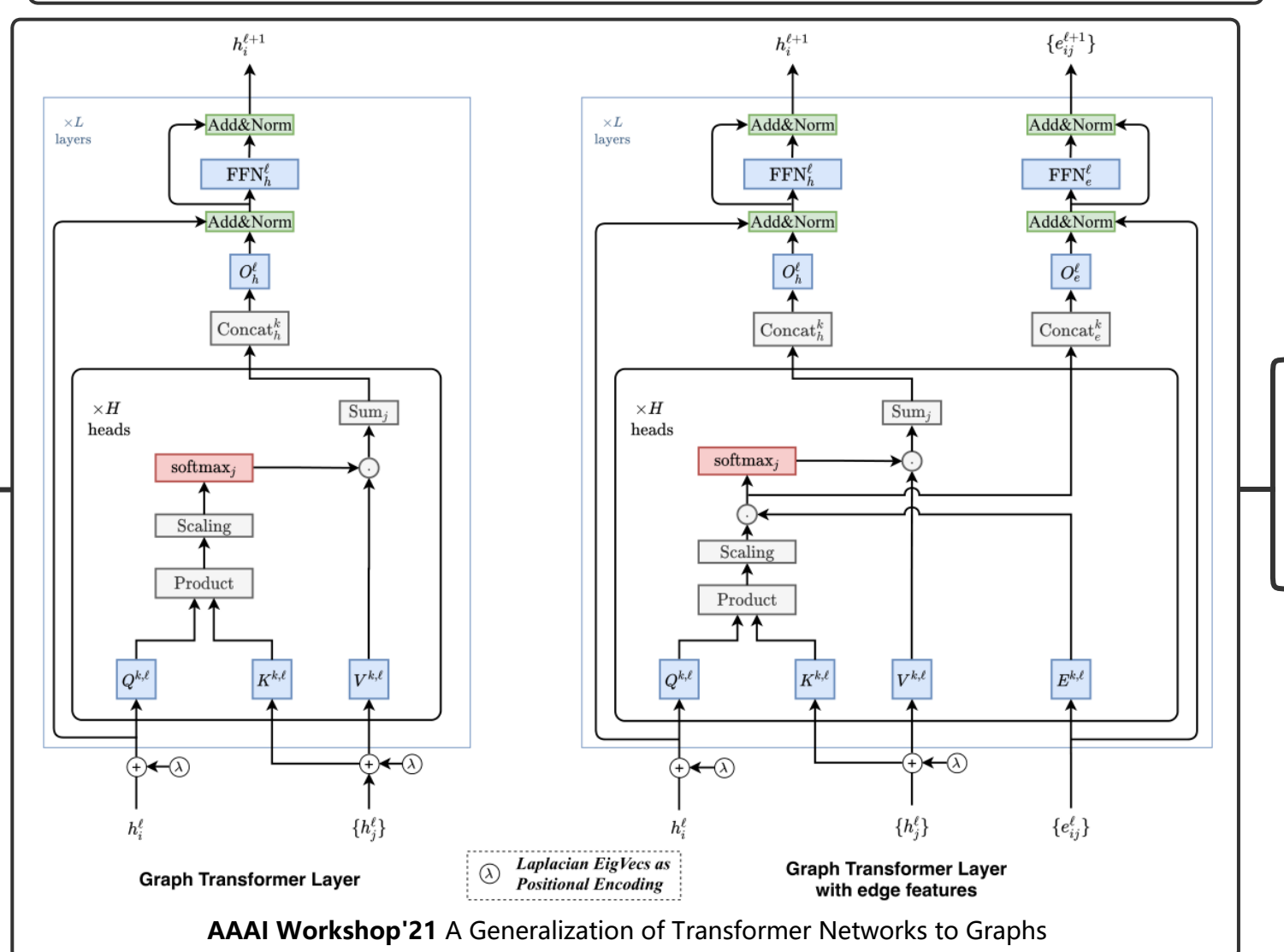
Different types of Graph Transformer

Restricting self-attention to local neighborhoods  
 Using global self-attention in conjunction with message-passing GNN  
 Injecting edge information into global self-attention via attention bias



$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(a^T(W_{v_i}||W_{v_j})))}{\sum_{k \in \mathcal{N}(v_i)} \exp(\text{LeakyReLU}(a^T(W_{v_i}||W_{v_k})))}$$

$$R_i^k = \sigma \left( \frac{1}{K} \sum_{j \in \mathcal{N}(v_i)} \alpha_{ij} W_{v_j}^k \right)$$



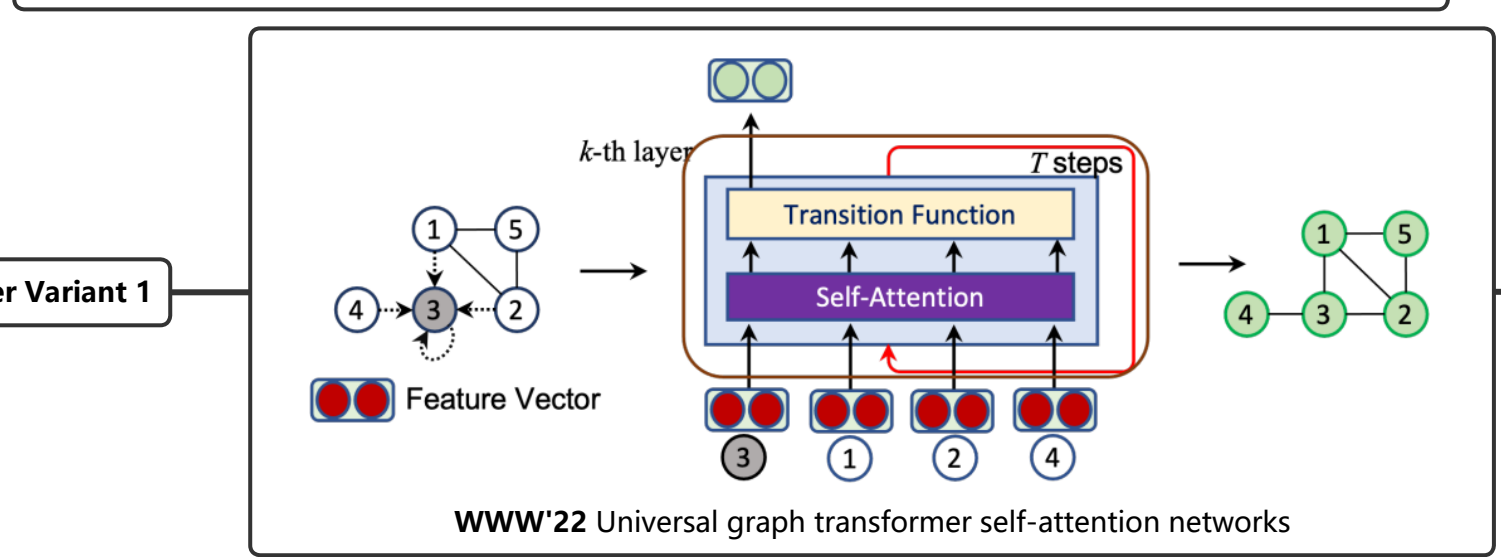
$$h_i^{k+1} = \text{LN} \left( h_i^k + \sum_{j \in \mathcal{N}(v_i)} \alpha_{ij} h_j^k \right) \quad (4)$$

where,  $\alpha_{ij}^k = \text{softmax} \left( \frac{e^{a^T(W_{v_i}^k || W_{v_j}^k)}}{\sum_{k \in \mathcal{N}(v_i)} e^{a^T(W_{v_i}^k || W_{v_k}^k)}} \right)$  (5)

$$\tilde{h}_i^{k+1} = \text{Norm}(h_i^k + h_i^{k+1}) \quad (6)$$

$$h_i^{k+1} = W_{\text{out}}^k \text{ReLU}(W_{\text{in}}^k \tilde{h}_i^{k+1}) \quad (7)$$

$$M_i^{k+1} = \text{Norm}(h_i^{k+1} + \tilde{h}_i^{k+1}) \quad (8)$$



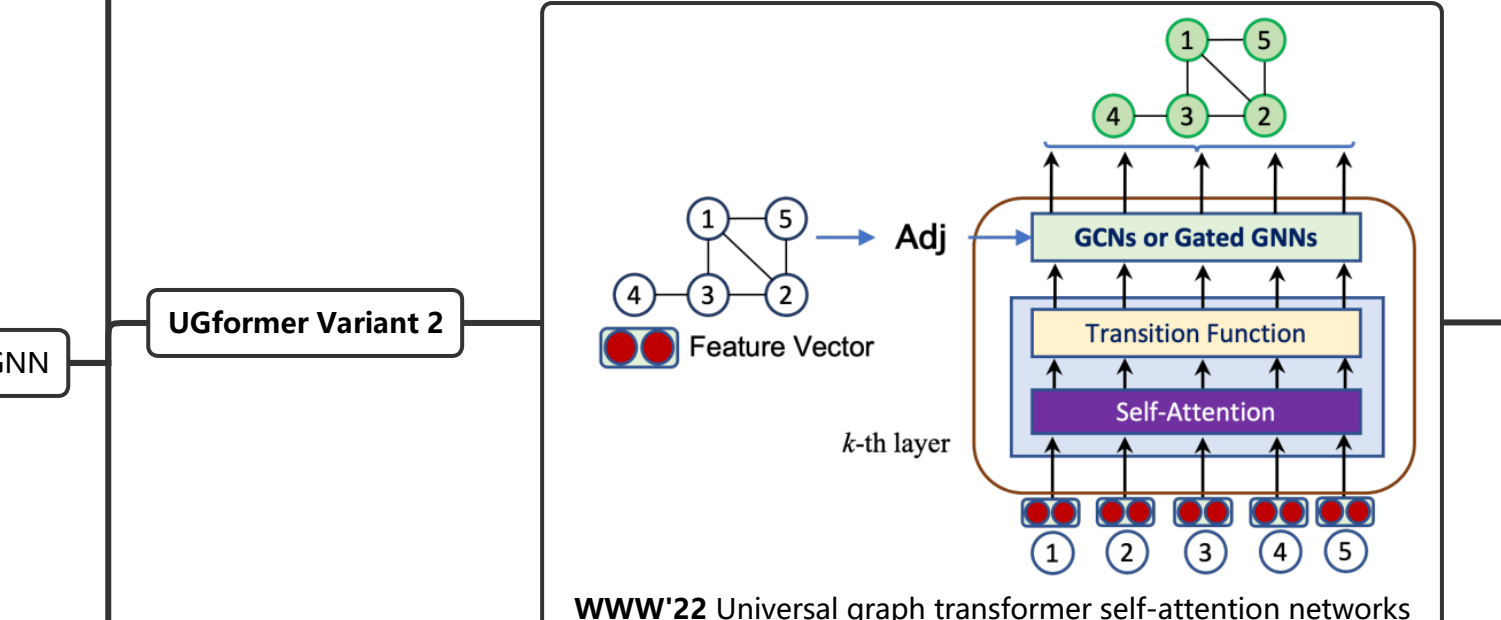
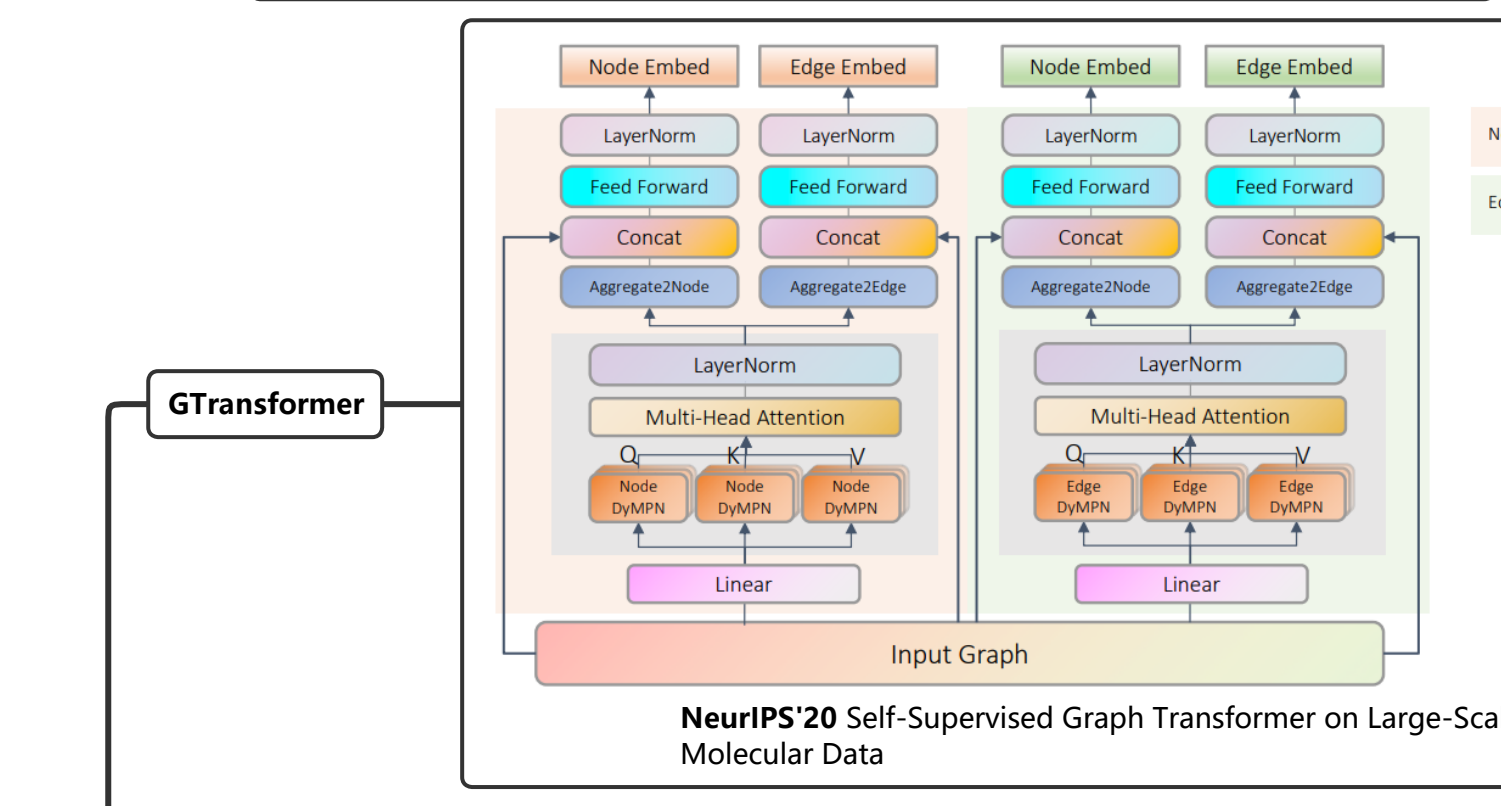
$$h_{v,u}^{(k)} = \text{LAYERNORMALIZATION}(h_{v,u}^{(k-1)} + \text{ATT}(h_{v,u}^{(k-1)})) \quad (1)$$

$$h_{v,u}^{(k)} = \text{LAYERNORMALIZATION}(h_{v,u}^{(k)} + \text{TRANS}(h_{v,u}^{(k)})) \quad (2)$$

$$\text{ATT}(h_{v,u}^{(k)}) = \sum_{v' \in \mathcal{N}(v)} \alpha_{v'u}^{(k)} (v^{(k)} h_{v',u}^{(k)})$$

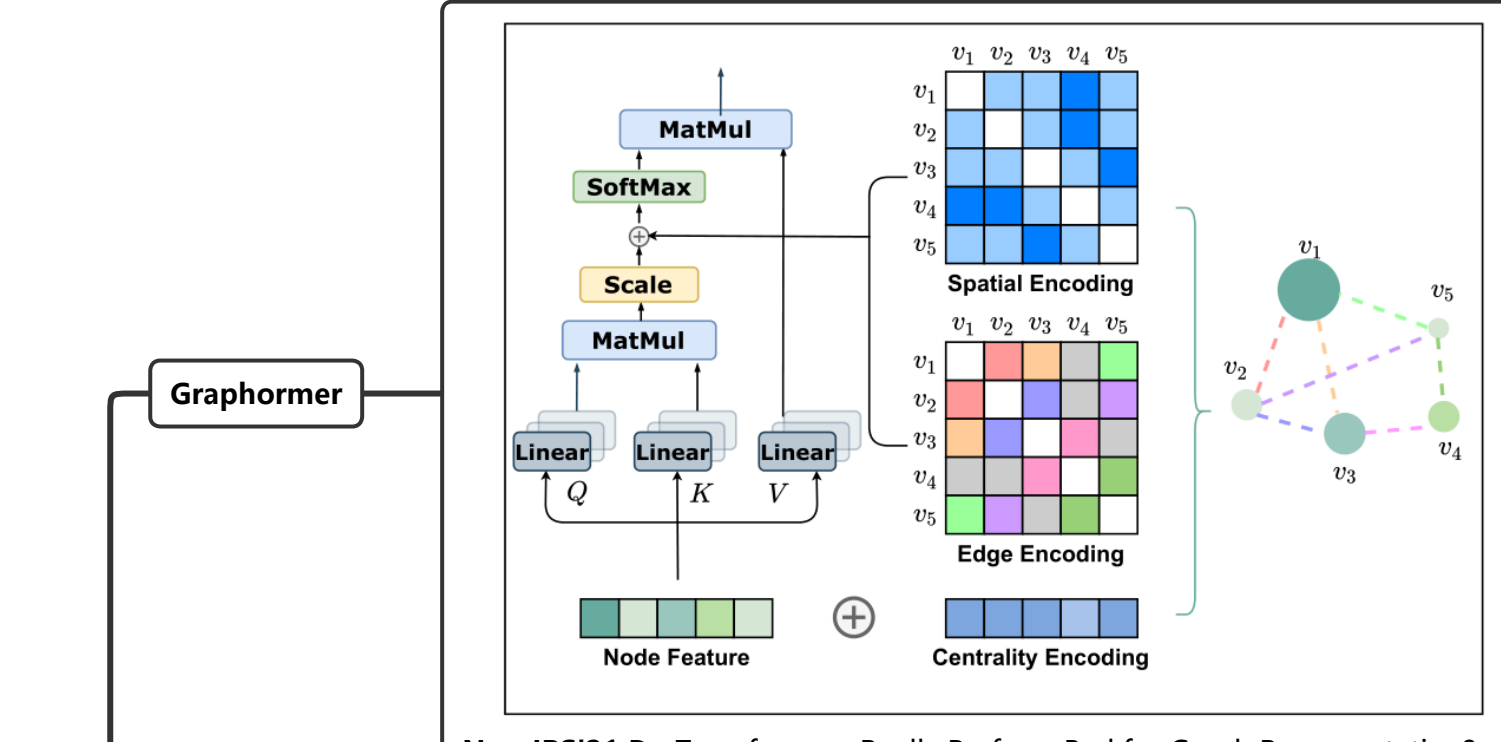
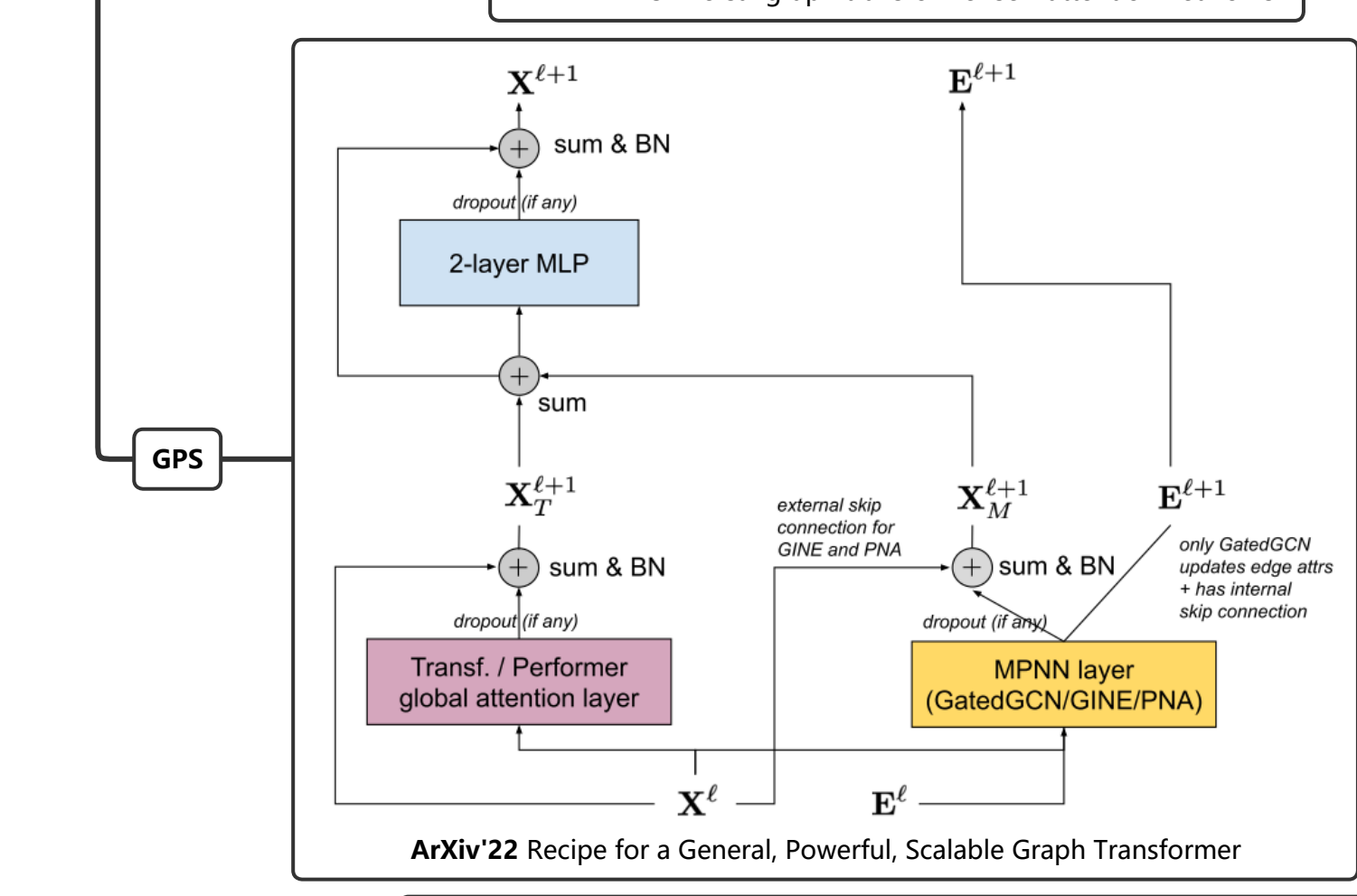
$$\alpha_{v'u}^{(k)} = \text{softmax} \left( \frac{Q^{(k)} h_{v,u}^{(k-1)} \cdot K^{(k)} h_{v',u}^{(k-1)}}{\sqrt{d}} \right)$$

$$Q^{(k)} \in \mathbb{R}^{d \times d} \text{ and } K^{(k)} \in \mathbb{R}^{d \times d}$$



$$H^{(k)} = \text{ATTENTION}_{\text{V}}(H^{(k-1)}, H^{(k-1)}, H^{(k-1)}, K^{(k)}, V^{(k)}, V^{(k)})$$

$$H^{(k+1)} = \text{GNN}(A, H^{(k)}) \quad (8)$$



$$h_i^{(0)} = d_i + \sum_{j \in \mathcal{N}(v_i)} \text{deg}^-(v_j) + \sum_{j \in \mathcal{N}(v_i)} \text{deg}^+(v_j)$$

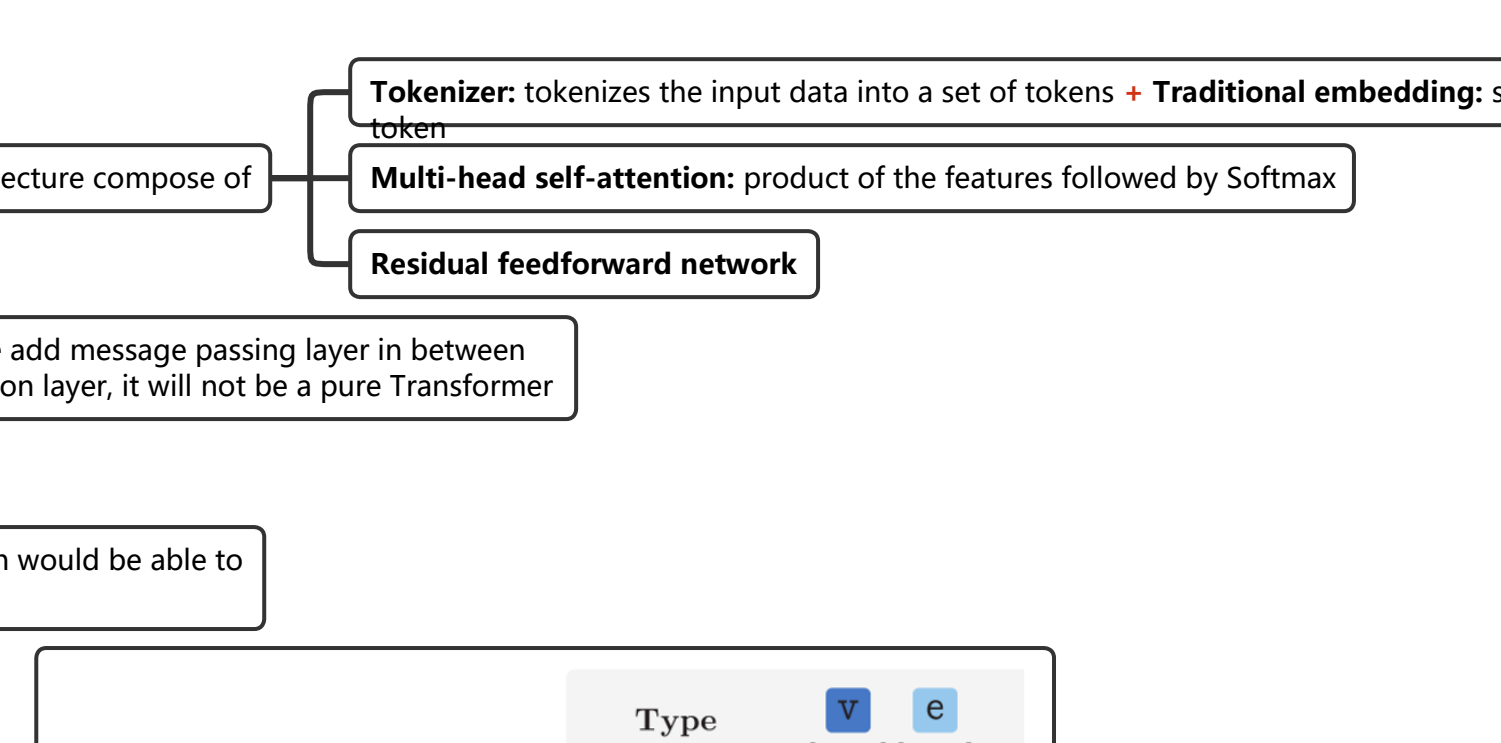
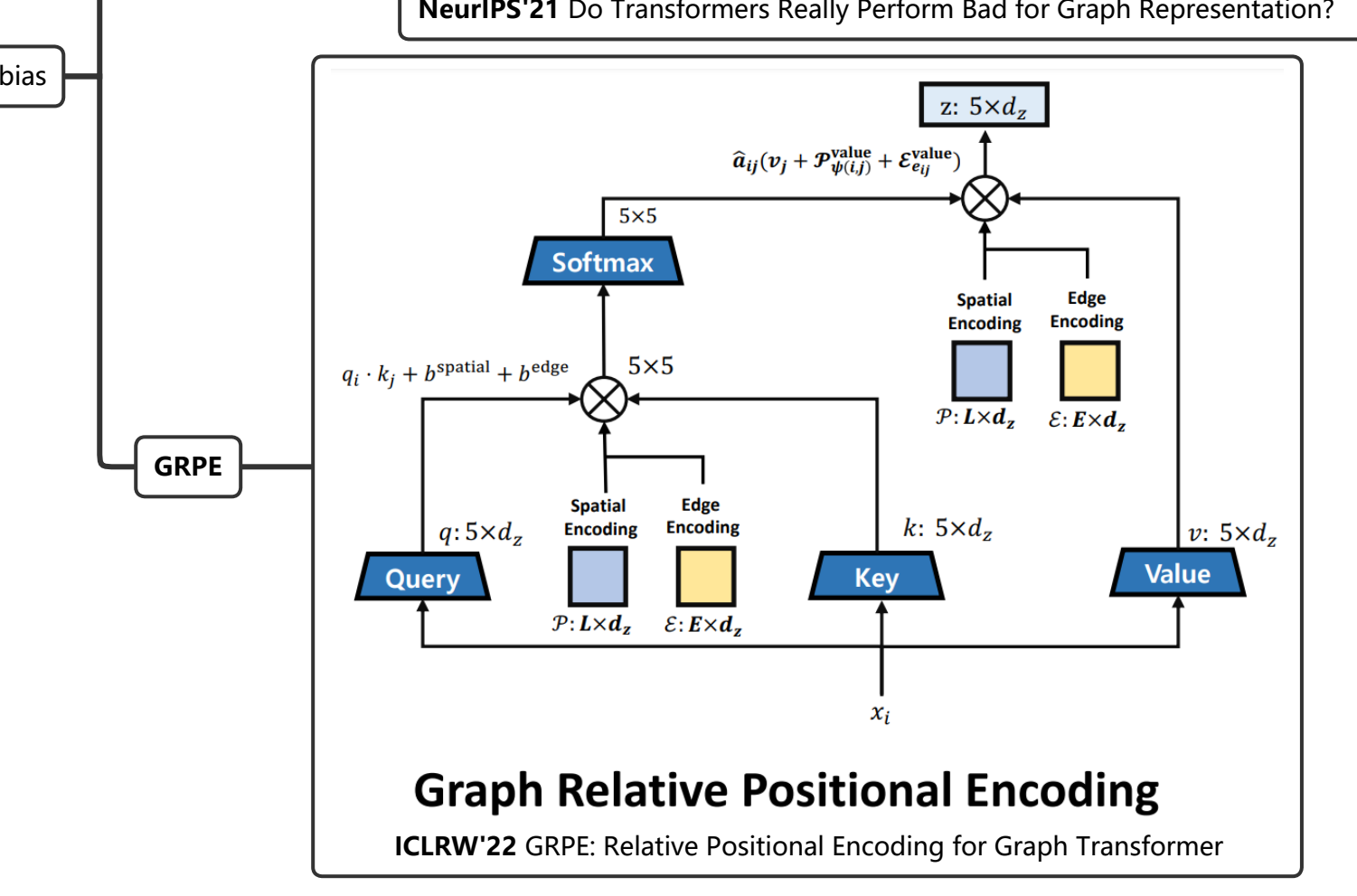
centrality encoding

$$A_{ij} = \frac{(A_{ij} + |A_{ij}|) |A_{ij}|^2 + \delta_{ij} c_{ij}}{\sqrt{d_i d_j}}$$

Spatial encoding,  $\delta$  can be connectivity between the nodes in the graph

$$c_{ij} = \frac{1}{N} \sum_{v \in \mathcal{V}} \sum_{v' \in \mathcal{V}} (e_{ij}^{v,v'})^T$$

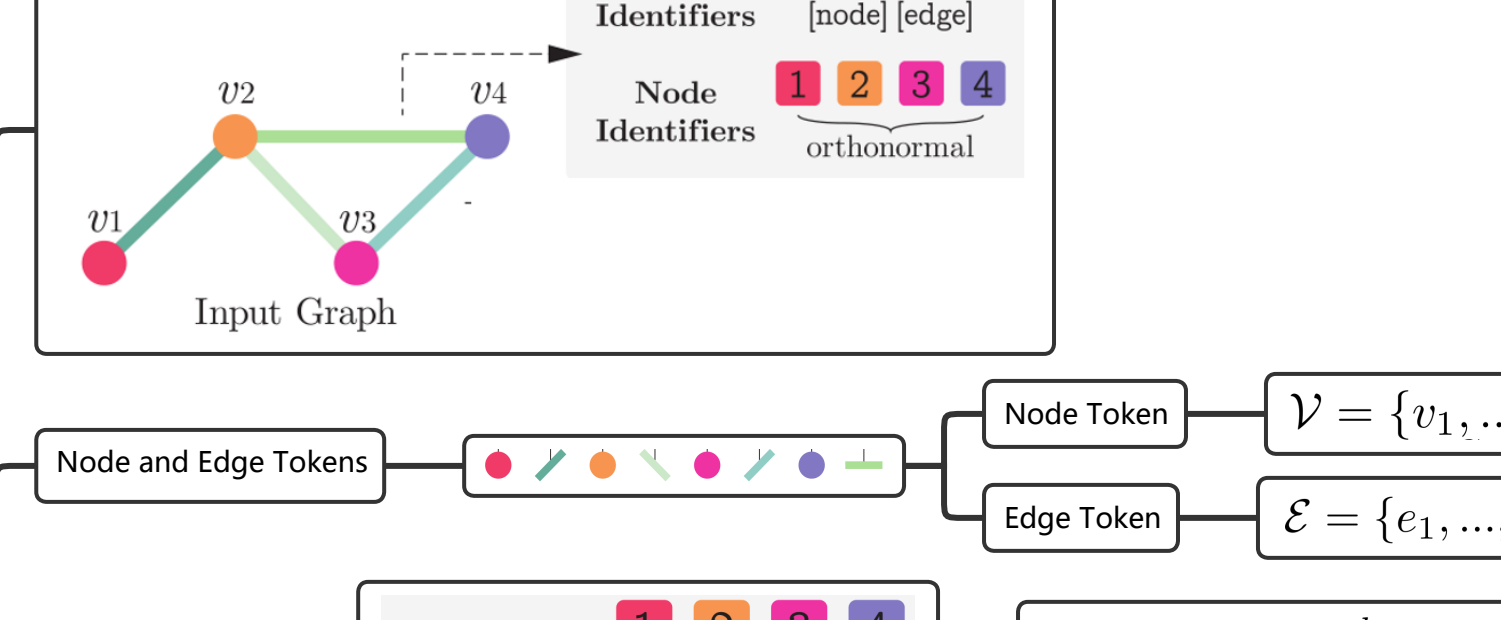
Edge encoding



## Tokenized Graph Transformer (TokenGT)

What is the definition of pure Transformer?  
 Network architecture compose of:  
 - Multi-head self-attention: product of the features followed by Softmax  
 - Residual feedforward network  
 Example: if we add message passing layer in between the self-attention layer, it will not be a pure Transformer

Intuition:  
 - Treat both node and edge as token  
 - With appropriate token-wise information, self-attention would be able to properly recognize the graph structure

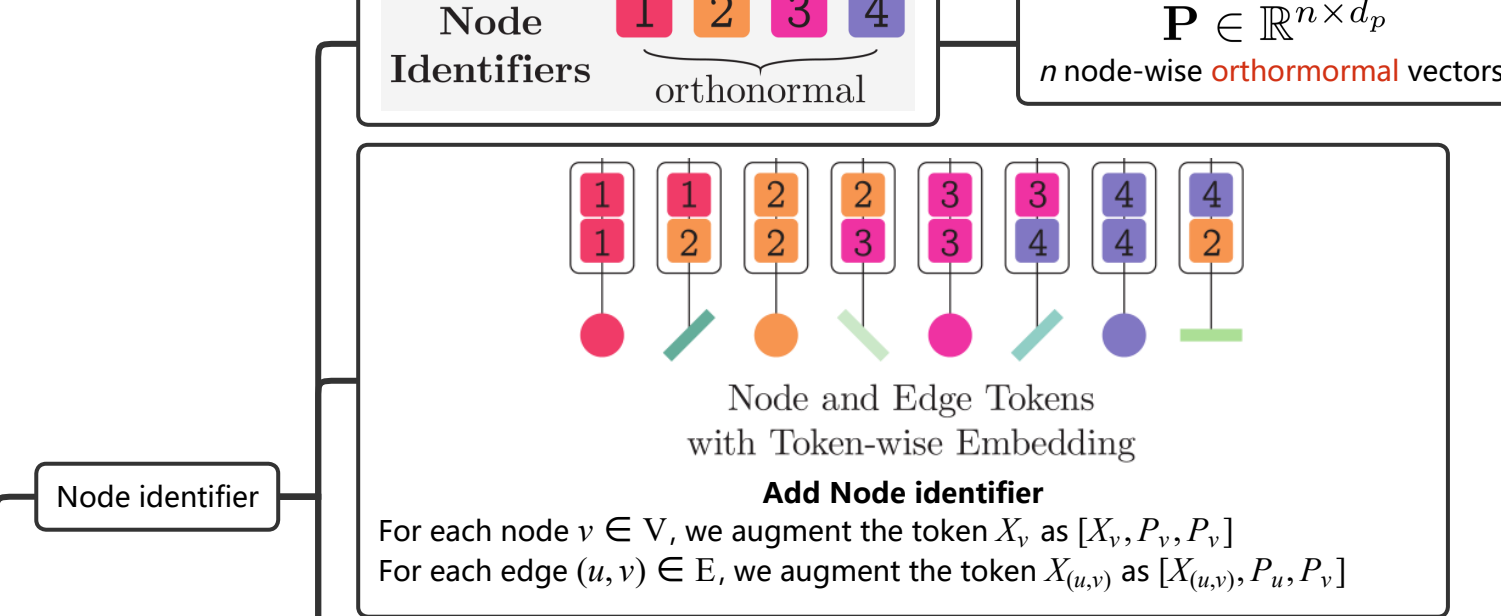


$$V = \{v_1, \dots, v_n\} \quad \mathbf{X}^V \in \mathbb{R}^{n \times C}$$

$$E = \{e_1, \dots, e_m\} \subseteq \mathcal{V}^2 \quad \mathbf{X}^E \in \mathbb{R}^{m \times C}$$

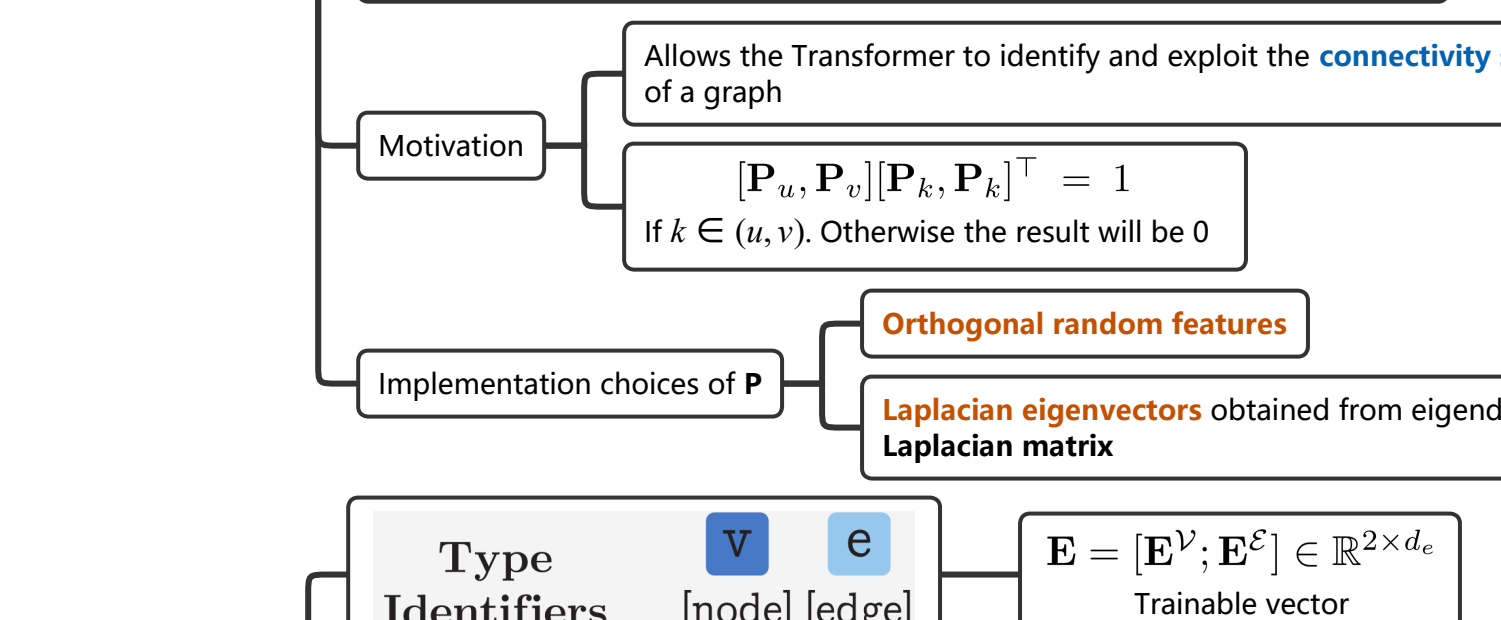
$$\mathbf{X} = [\mathbf{X}^V; \mathbf{X}^E] \in \mathbb{R}^{(n+m) \times C}$$

## TokenGT architecture

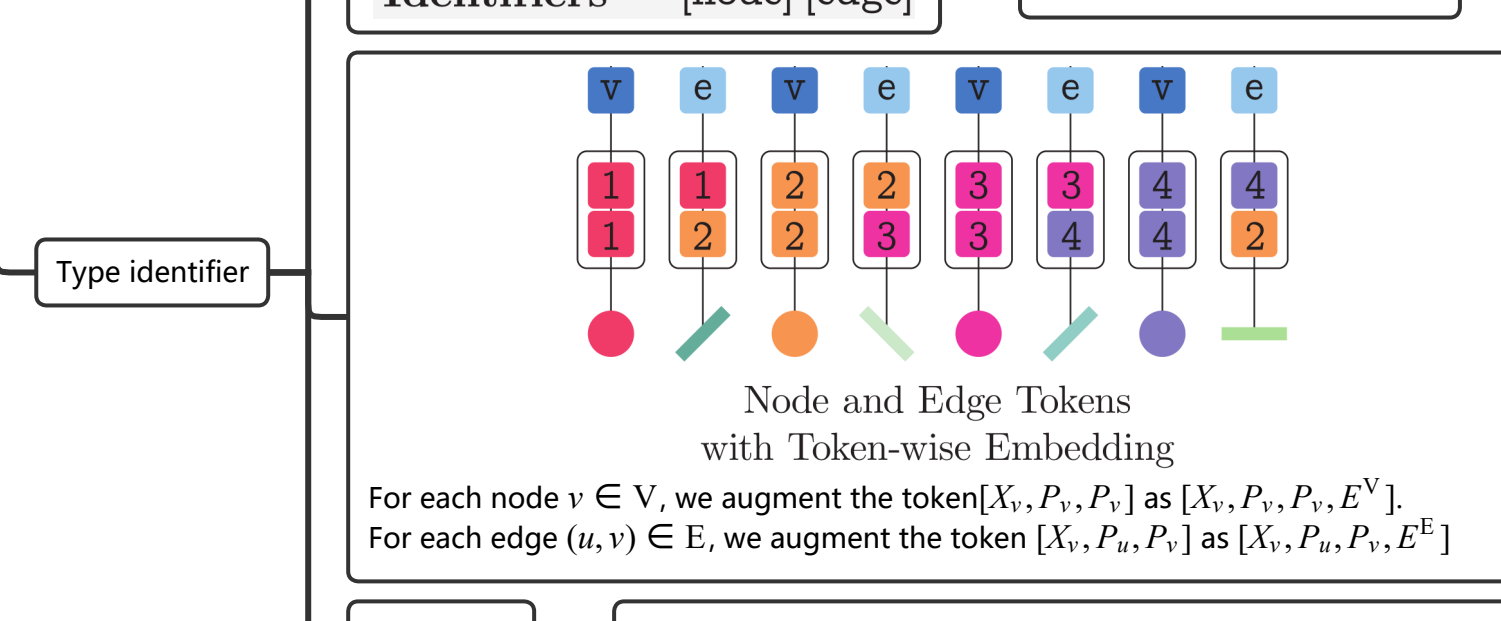


Motivation:  
 $\|P_{v_i} P_{v_j}\| \|P_{v_i} P_{v_j}\|^T = 1$   
 if  $k \in (u, v)$ . Otherwise the result will be 0

Implementation choices of P:  
 - Orthogonal random features  
 - Laplacian eigenvectors obtained from eigendecomposition of graph Laplacian matrix



Motivation:  
 Allow Transformer to identify whether a given token is a node or an edge



$$\mathbf{X}^n \in \mathbb{R}^{(n+m) \times (C+2d_p+d_e)}$$

$$W^{in} \in \mathbb{R}^{(C+2d_p+d_e) \times d}$$

$$\mathbf{X}^{in} = W^{in} \mathbf{X}^n$$

$$\mathbf{Z}^{(0)} = [\mathbf{X}_{\text{graph}}; \mathbf{X}^{in}] \in \mathbb{R}^{(1+n+m) \times d}$$

Final input to the main encoder

## Experimental Results

method	# parameters	validate MAE	asymptotics
<b>Message-passing GNNs</b>			
GCN [25]	2.0M	0.1379	$\mathcal{O}(n+m)$
GIN [26]	3.8M	0.1195	$\mathcal{O}(n+m)$
GAT	6.7M	0.1302	$\mathcal{O}(n+m)$
GCNv2 [25]	4.9M	0.1153	$\mathcal{O}(n+m)$
GINv2 [26]	6.7M	0.1183	$\mathcal{O}(n+m)$
GATv2	6.7M	0.1192	$\mathcal{O}(n+m)$
GCNv2 (large)	55.2M	0.1361	$\mathcal{O}(n+m)$
<b>Transformers with strong graph-specific modifications</b>			
Graphormer [23]	48.5M	0.0864	$\mathcal{O}(n^2)$
EGT [23]	89.3M	0.0869	$\mathcal{O}(n^2)$
GRPE [23]	46.2M	0.0890	$\mathcal{O}(n^2)$
<b>Pure Transformers</b>			
Transformer	48.5M	0.2340	$\mathcal{O}(n+m)^2$
TokenGT (OPE)	48.8M	0.0962	$\mathcal{O}(n+m)^2$
TokenGT (Lap)	48.5M	0.0910	$\mathcal{O}(n+m)^2$
TokenGT (Lap) + Performer	48.5M	0.0935	$\mathcal{O}(n+m)$