

Graph Transformers

Prof. O-Joun Lee

Dept. of Artificial Intelligence,
The Catholic University of Korea
ojlee@catholic.ac.kr



네트워크 과학연구실
NETWORK SCIENCE LAB



가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA

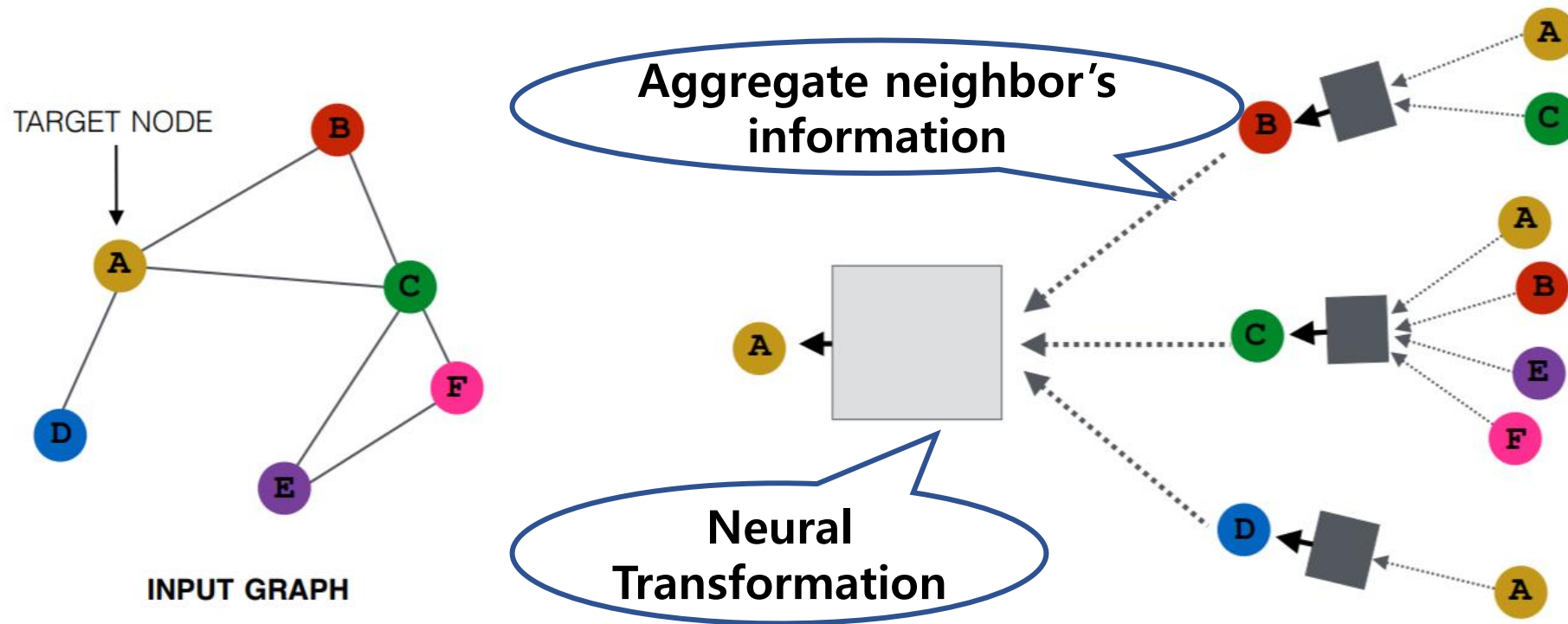


Contents

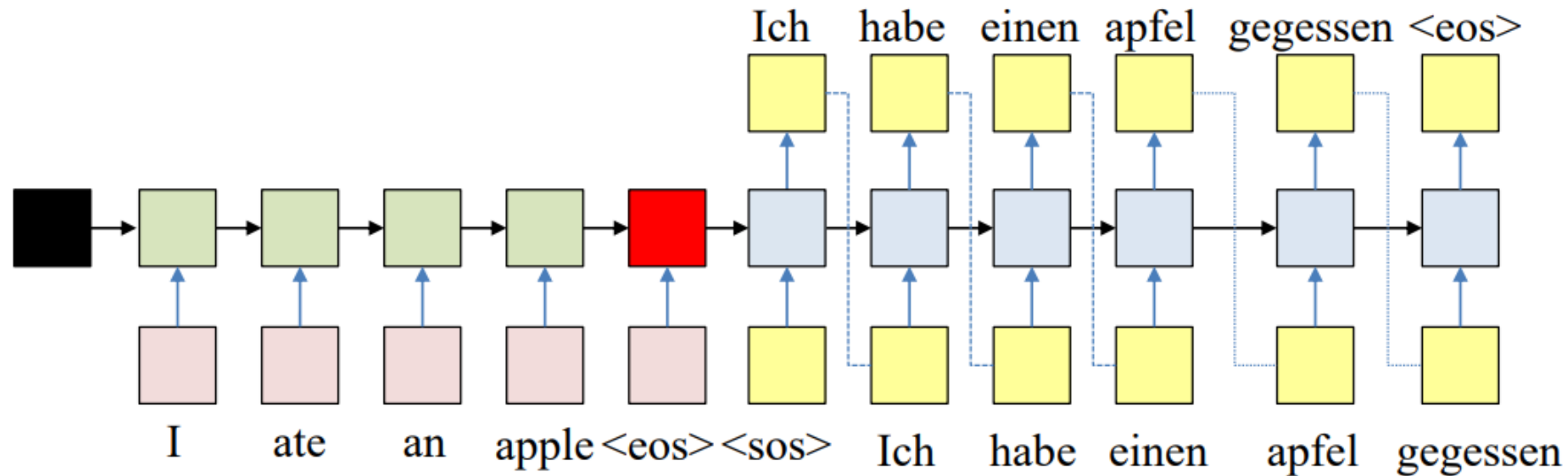


- From Message Passing to Self-attention
- Positional Encoding in Graphs
- Transformers are Graph Neural Networks
- Representative Transformer models

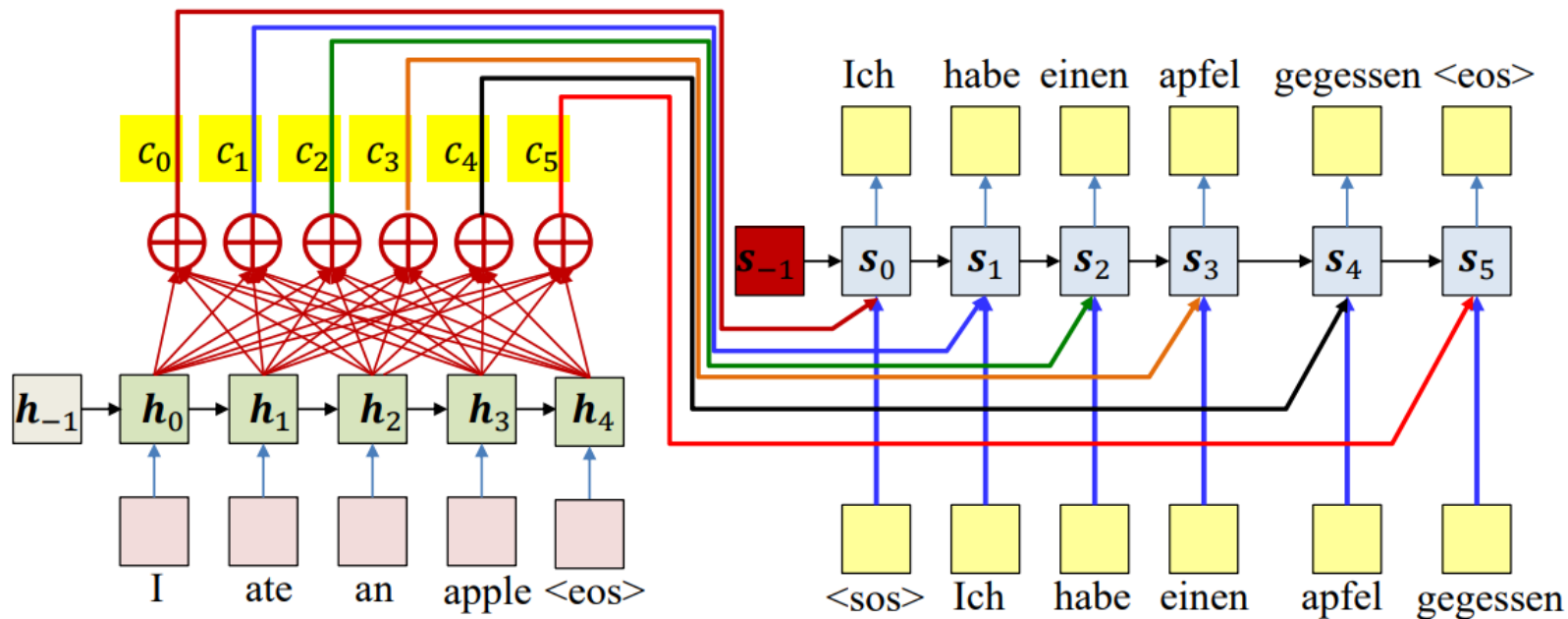
- **Key Idea:** Each node aggregates messages from its neighborhood to get contextualized node embedding.
- **Limitation:** Most GNNs focus on homogeneous graph.



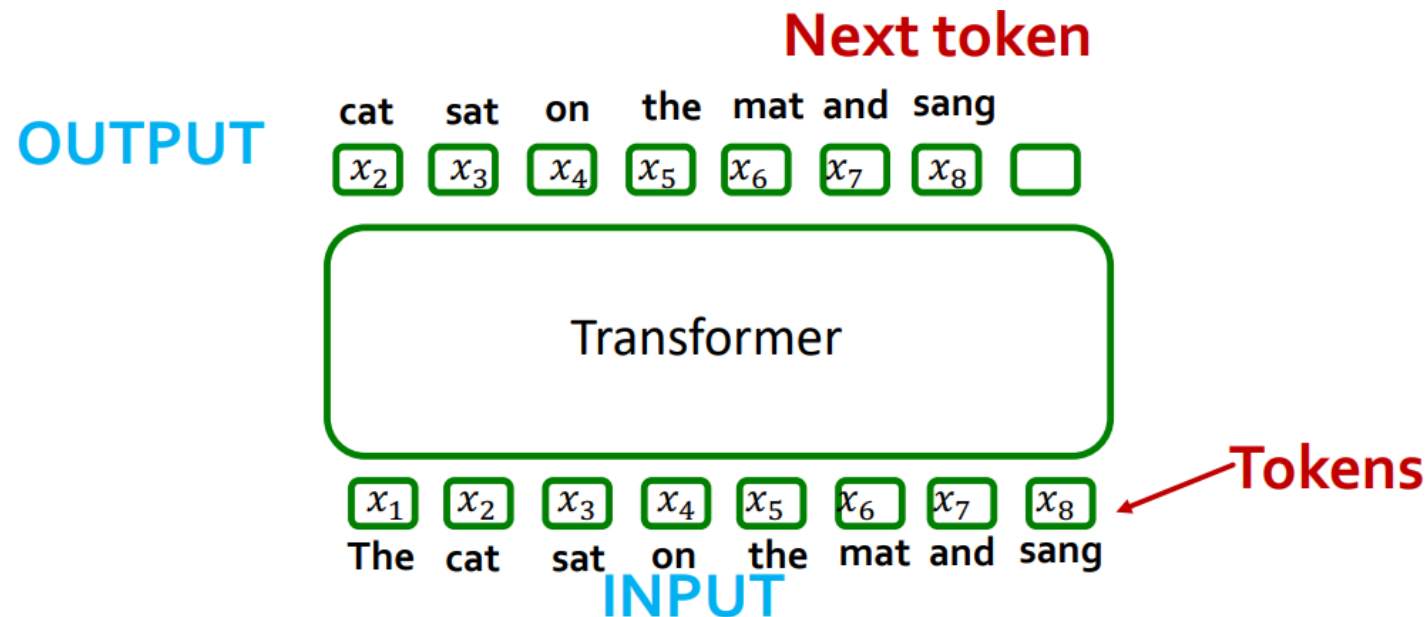
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence
- Subsequently a second RNN uses the hidden activation as initial state to produce a sequence of outputs



- Encoder recurrently produces hidden representations of input word sequence
- Decoder recurrently generates output word sequence
 - For each output word the decoder uses a weighted average of the hidden input representations as input “context”, along with the recurrent hidden state and the previous output word

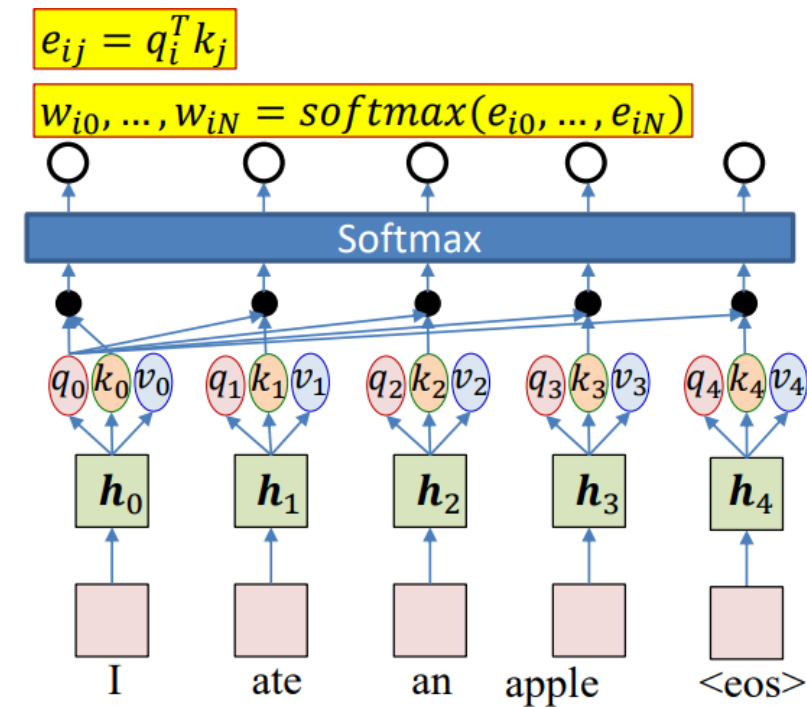


- Transformer ingest **TOKENS**
- Transformers map 1D sequences of vectors to 1D sequences of vectors known as tokens
 - Tokens describe a "piece" of data – e.g., a word
- What output sequence?
 - Option 1: next token => GPT

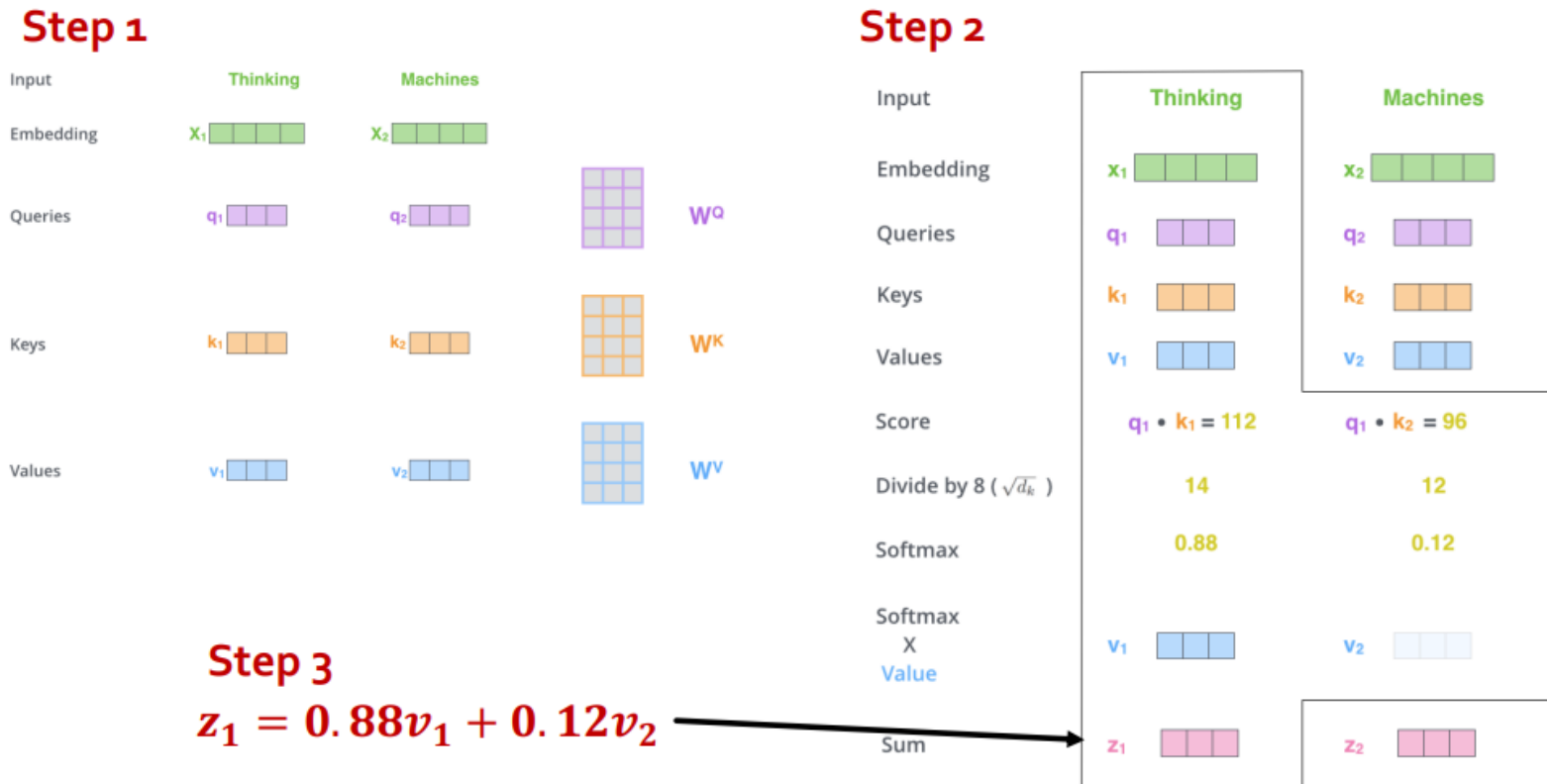


- First, for every word in the input sequence we compute an initial representation
 - E.g. using a single MLP layer
- Then, from each of the hidden representations, we compute a query, a key, and a value.
 - Using separate linear transforms
 - The weight matrices W_q , W_k and W_v are learnable parameters
- The updated representation for the word is the attention-weighted sum of the values for all words (Including itself)

$$\begin{aligned} q_i &= W_q h_i \\ k_i &= W_k h_i \\ v_i &= W_v h_i \\ w_{ij} &= \text{attn}(q_i, k_{0:N}) \end{aligned}$$

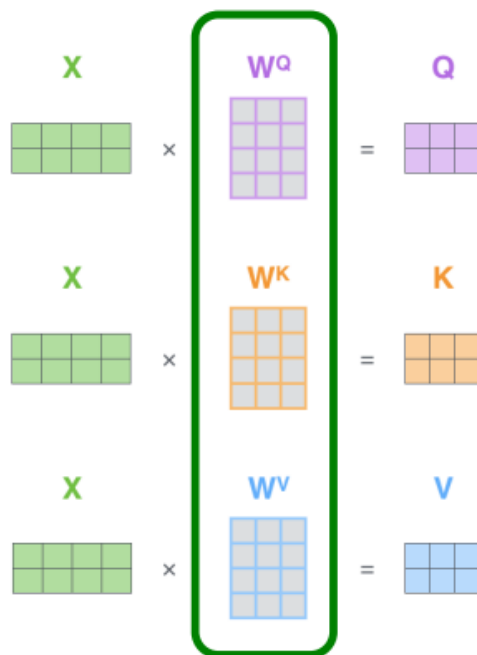


- Step 1: compute “key, value, query” for each input
- Step 2 (just for x_1): compute scores between pairs, turn into probabilities (same for x_2)
- Step 3: get new embedding z_1 by weighted sum of v_1, v_2

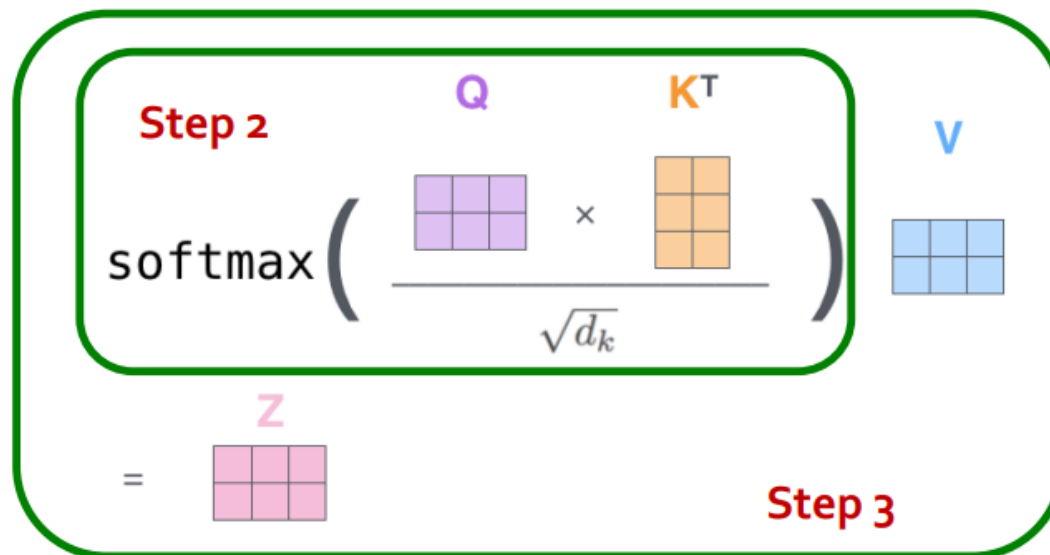


- Same calculation in matrix form

Step 1



Model
parameters



- We can have multiple such attention “heads”
 - Each will have an independent set of queries, keys and values
 - Each will obtain an independent set of attention weights
 - Potentially focusing on a different aspect of the input than other heads
 - Each computes an independent output
- The final output is the concatenation of the outputs of these attention heads
- “MULTI-HEAD ATTENTION”
(actually Multi-head self attention)

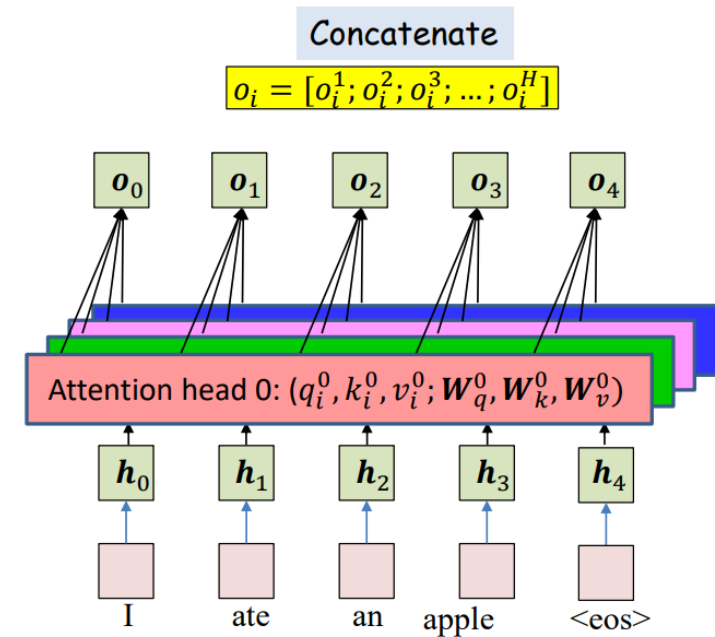
$$q_i^a = W_q^a h_i$$

$$k_i^a = W_k^a h_i$$

$$v_i^a = W_v^a h_i$$

$$w_{ij}^a = \text{attn}(q_i^a, k_{0:N}^a)$$

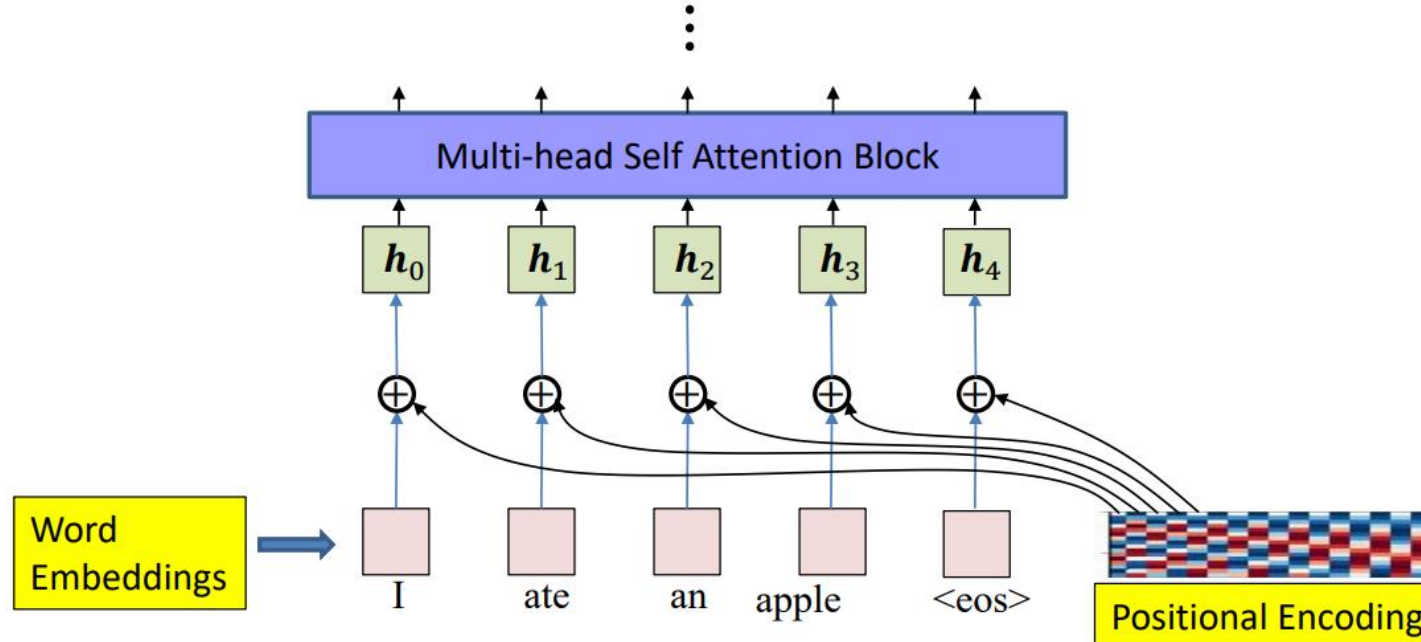
$$o_i^a = \sum_j w_{ij}^a v_j^a$$



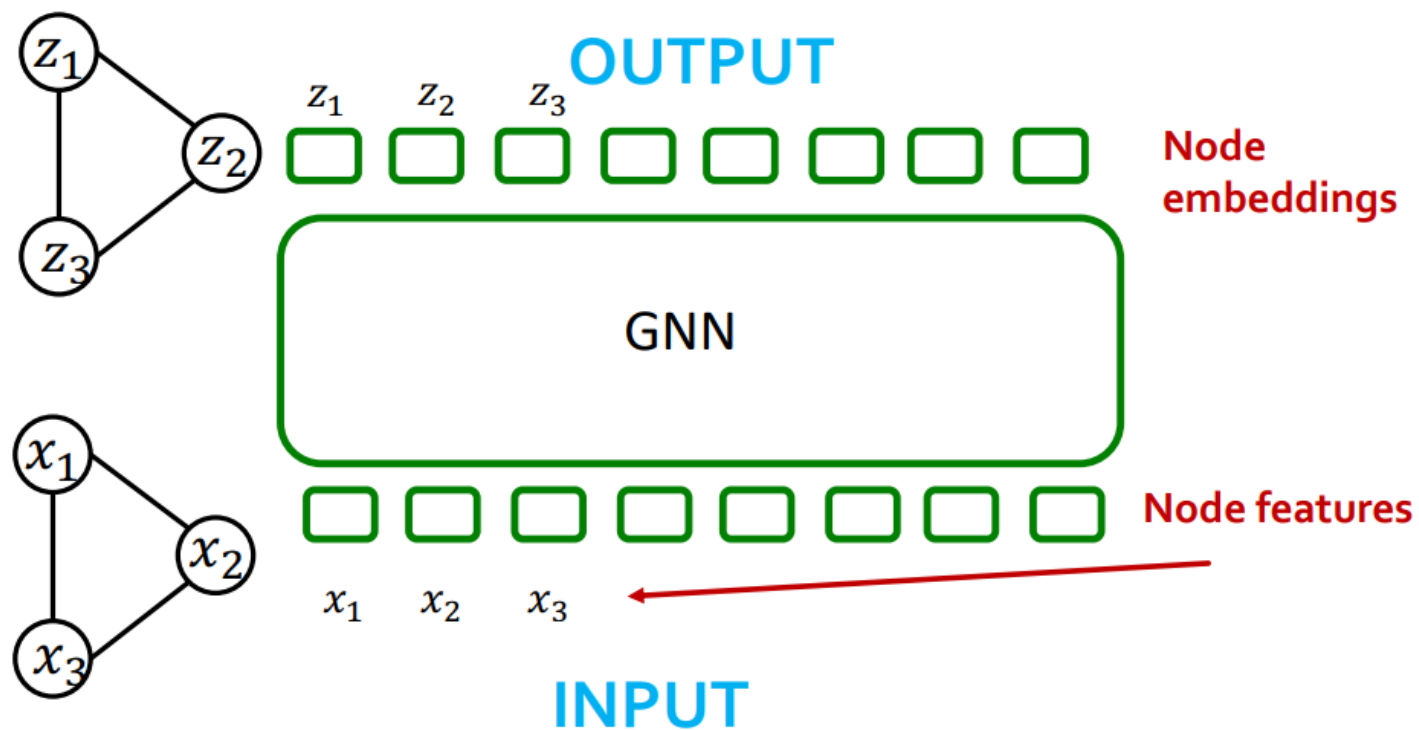
- Positional Encoding: A sequence of vectors P_0, P_1, \dots, P_N to encode position
 - Every vector is unique (and uniquely represents time)
 - Relationship between P_t and P_{t+k} only depends on the distance between them:

$$P_{t+k} = M_k P_t$$

- The linear relationship between P_t and P_{t+k} enables the net to learn shiftinvariant “gap” dependent relationships

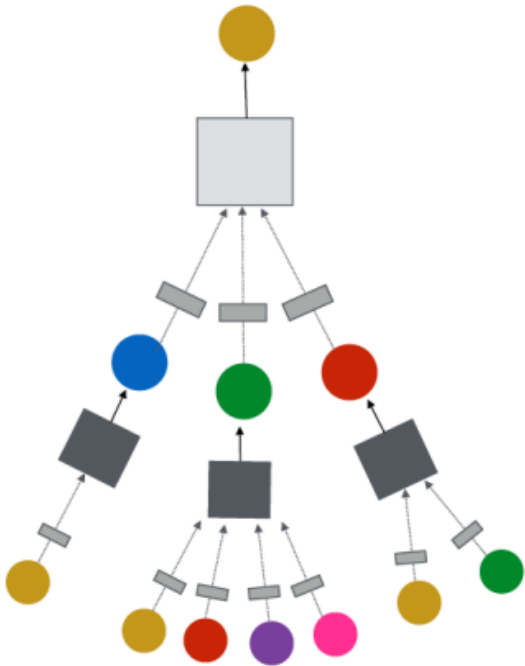


- **Similarity:** GNNs also take in a sequence of vectors (in no particular order) and output a sequence of embeddings
- **Difference:** GNNs use message passing, Transformer uses self-attention



- **Difference:** GNNs use message passing, Transformer uses self-attention
- Are self-attention and message passing really different?

Message Passing

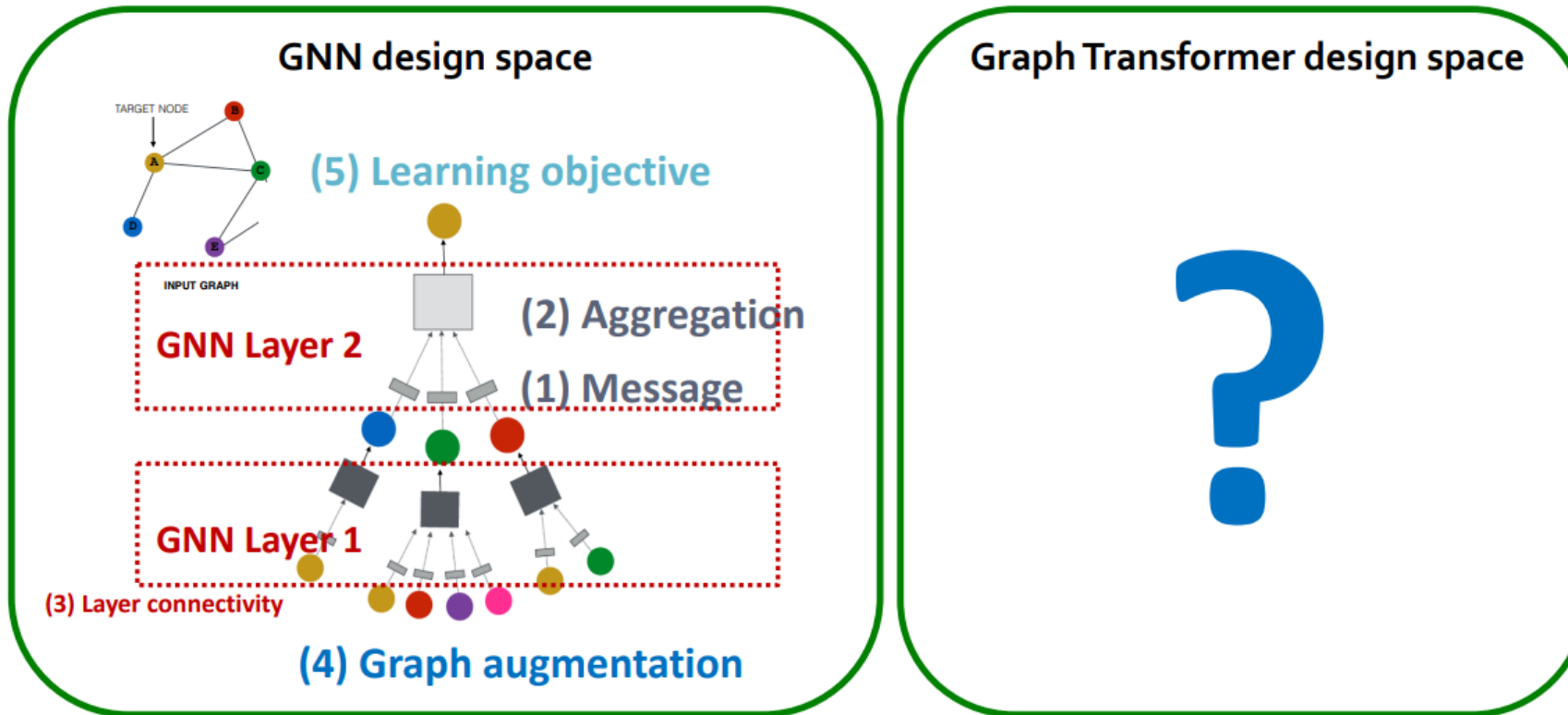


Vs.

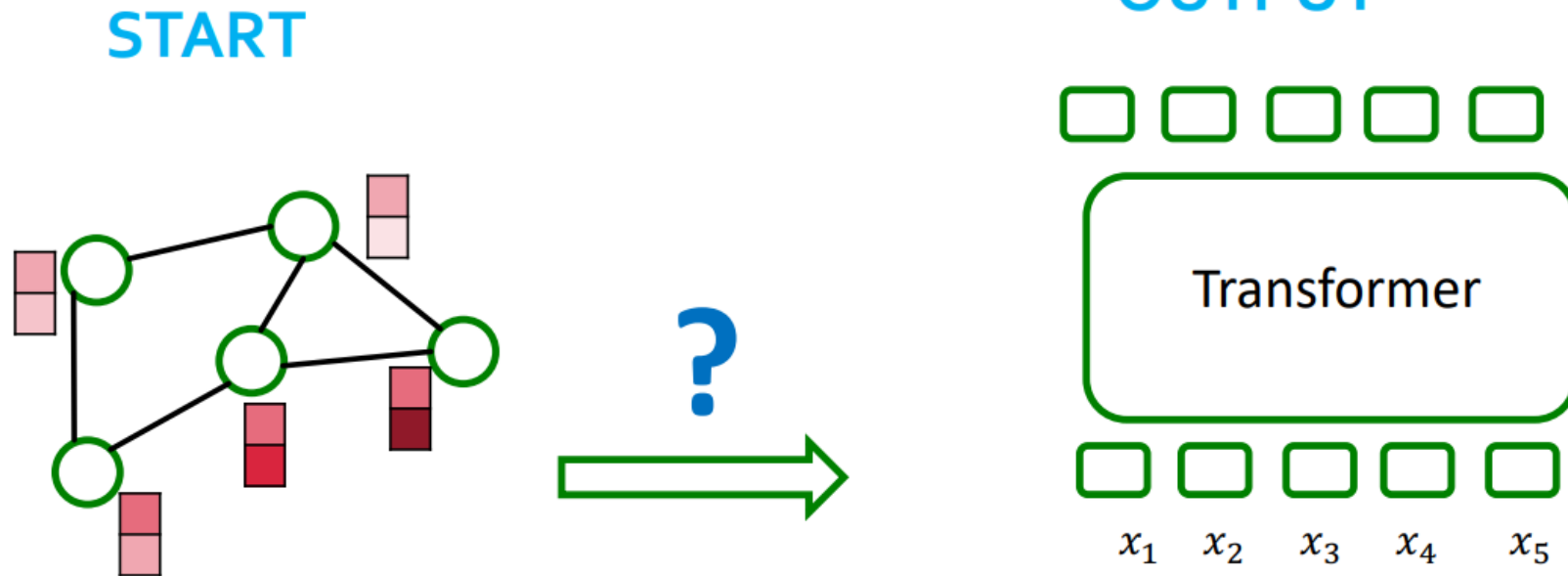
Self-attention

$$\begin{aligned} X &\times W^Q = Q \\ X &\times W^K = K \\ X &\times W^V = V \\ \text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) &\times V = Z \end{aligned}$$

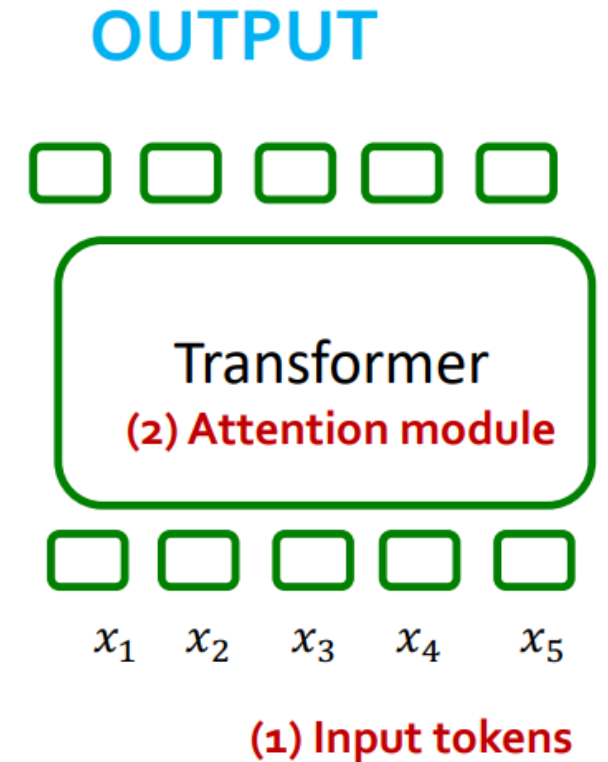
- We know a lot about the design space of GNNs
- What does the corresponding design space for Graph Transformers look like?



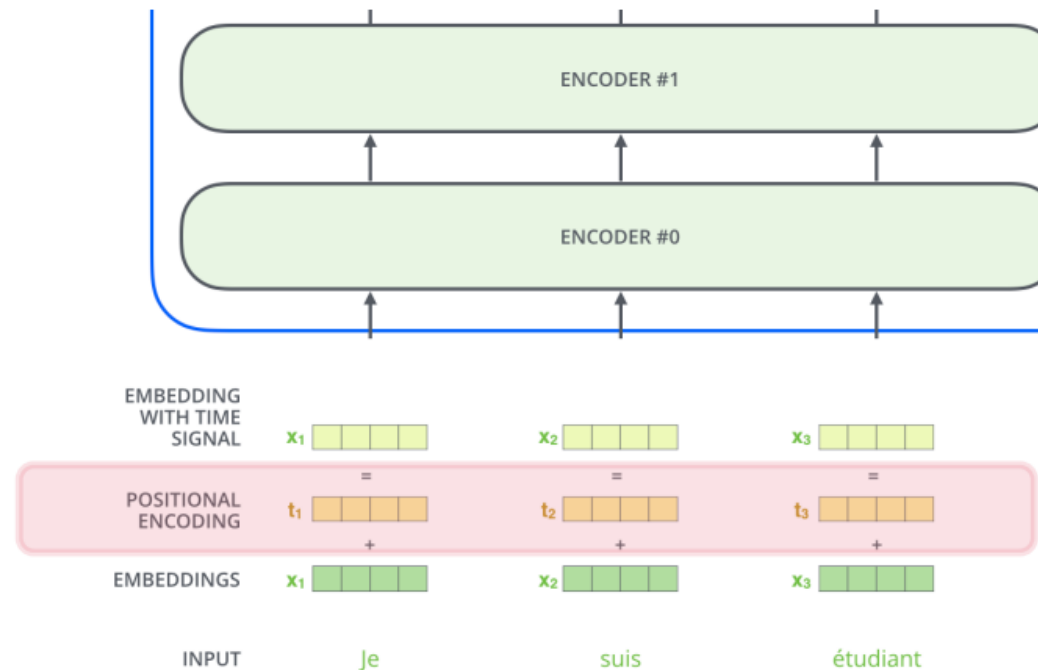
- We start with graph(s)
- How to input a graph into a Transformer?



- To understand how to process graphs with Transformers, we must:
 - Understand the key components of the Transformer. Seen already:
 - 1) tokenizing,
 - 2) self-attention
 - Decide how to make suitable graph versions of each



- Transformer doesn't know order of inputs
- Extra positional features needed so it knows that
 - Je = word 1
 - suis = word 2
 - etc.
- For NLP, positional encoding vectors are learnable parameters

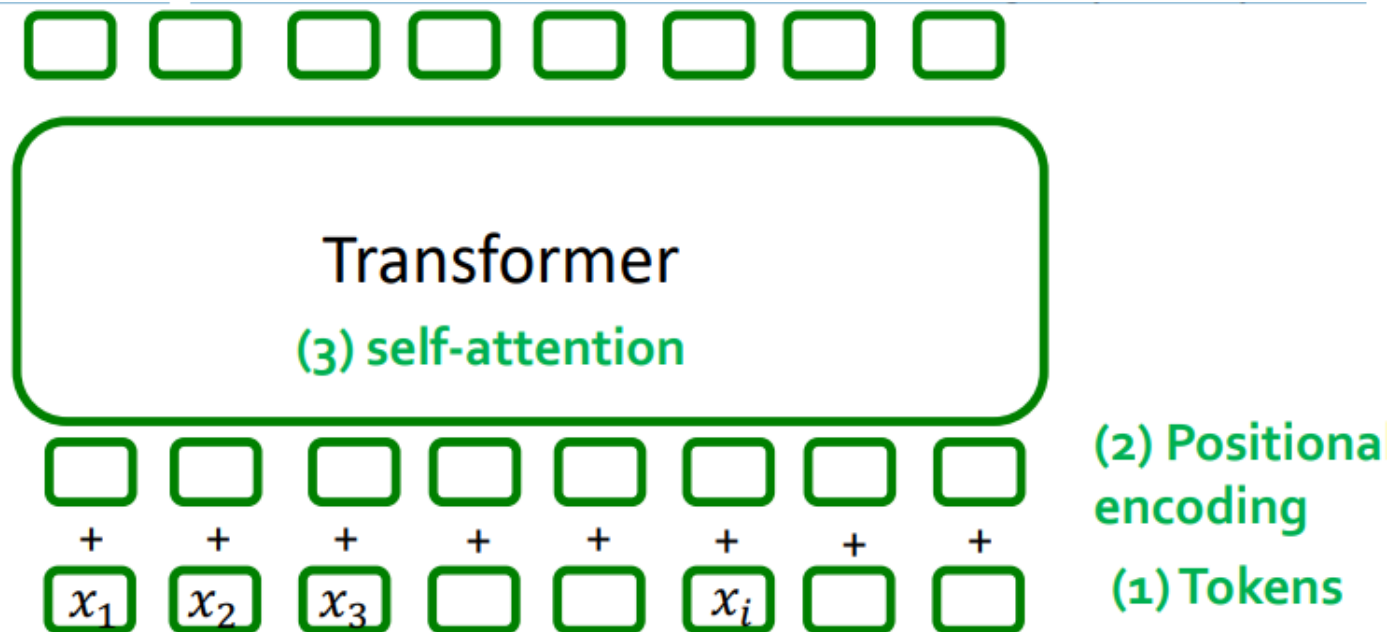


➤ Key components of Transformer:

- (1) tokenizing
- (2) positional encoding
- (3) self-attention

How to chose these for graph data?

➤ Key question: What should these be for a graph input?

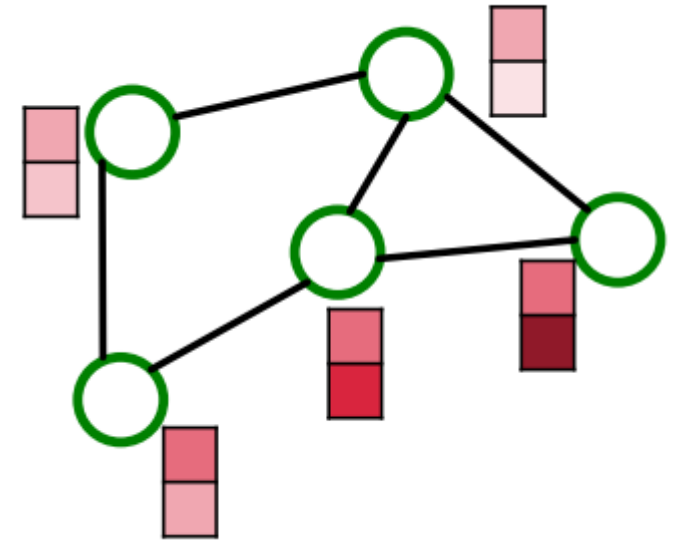


- A graph Transformer must take the following inputs:
 - (1) Node features?
 - (2) Adjacency information?
 - (3) Edge features (if any)

- Key components of Transformer:
 - (a) tokenizing
 - (b) positional encoding
 - (c) self-attention

SOLUTIONS:

- There are many ways to do this:
- Different approaches correspond to different “matchings” between graph inputs (1), (2), (3) transformer components (a), (b), (c)



➤ A graph Transformer must take the following inputs:

➤ (1) Node features?



➤ (2) Adjacency information?



➤ (3) Edge features (if any)



➤ Key components of Transformer:

➤ (a) tokenizing

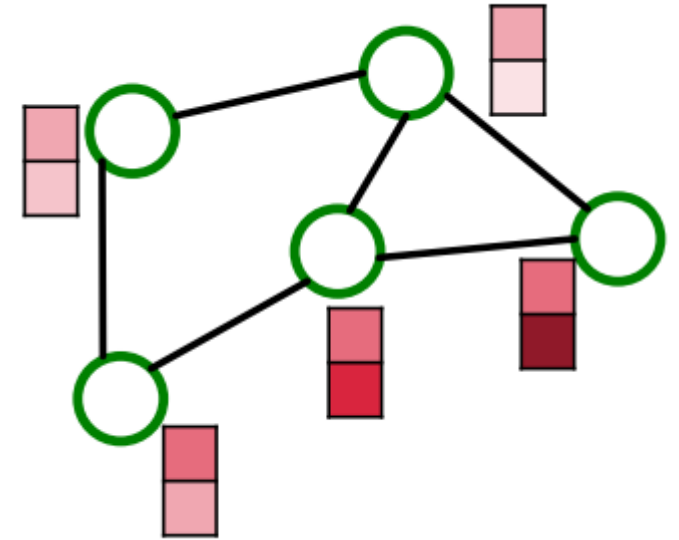
➤ (b) positional encoding

➤ (c) self-attention

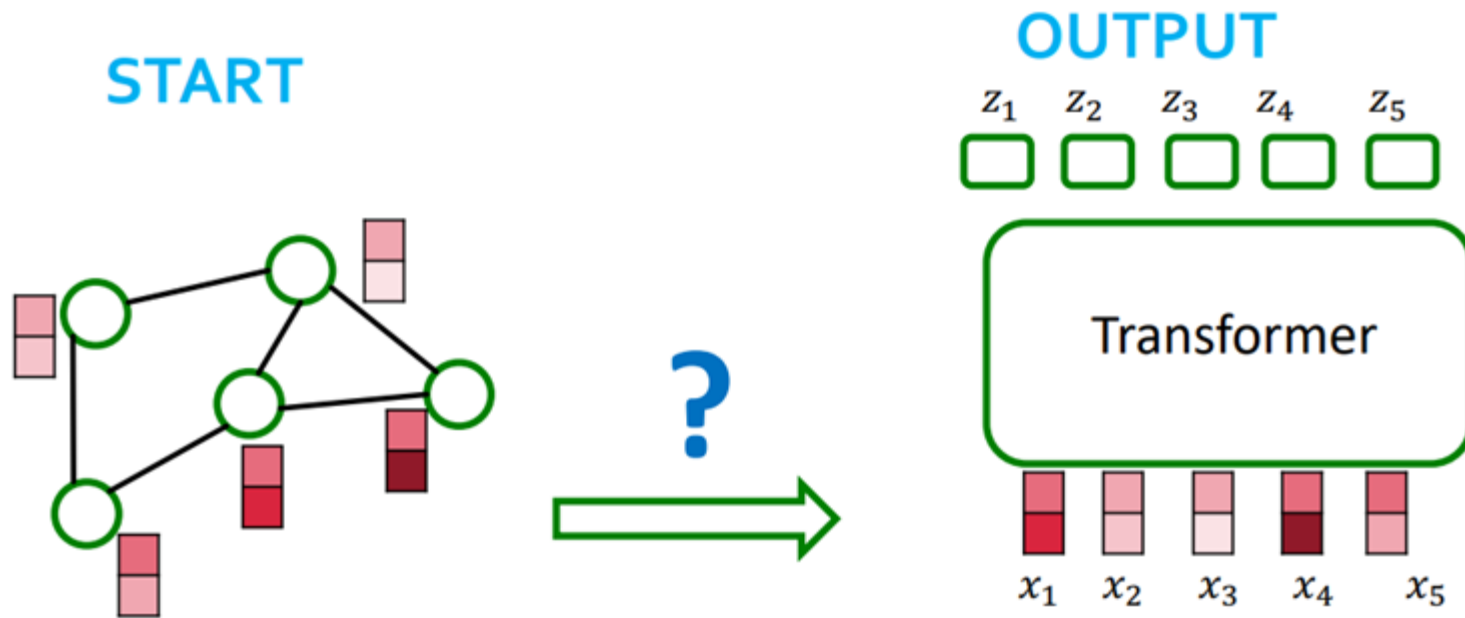
SOLUTIONS:

➤ There are many ways to do this:

➤ Different approaches correspond to different “matchings” between graph inputs (1), (2), (3) transformer components (a), (b), (c)

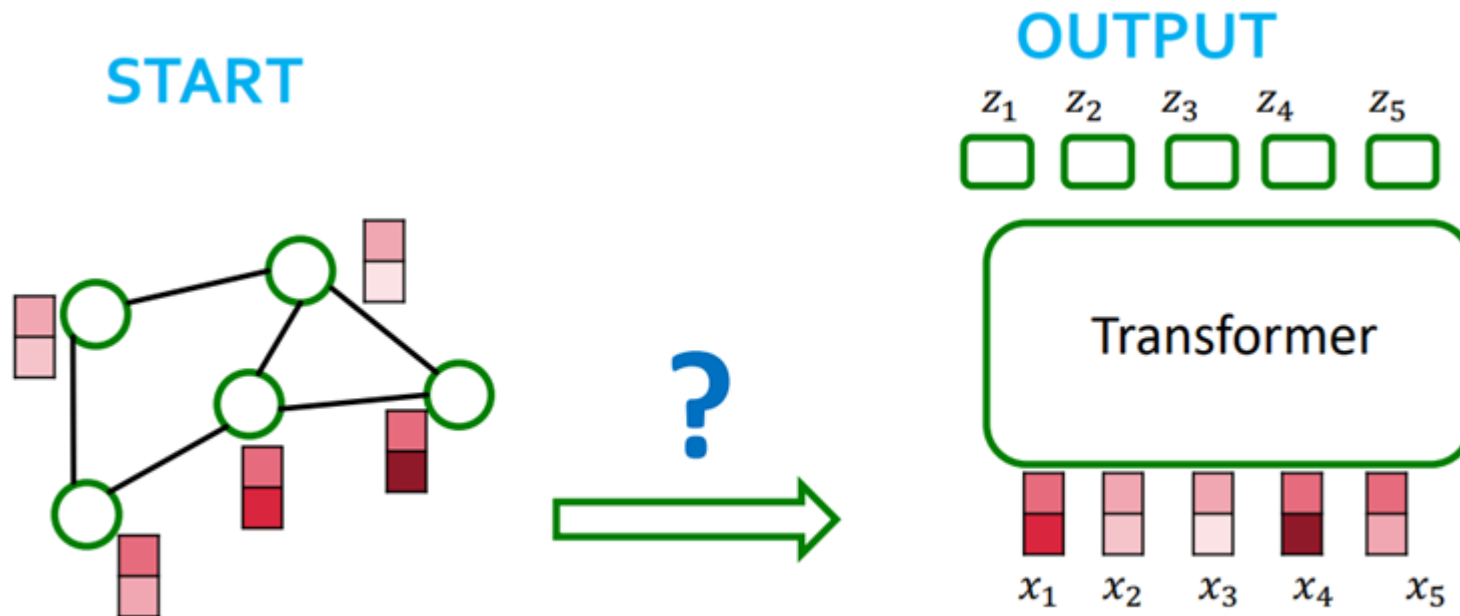


- **Q1:** what should our tokens be?
- **Sensible Idea:** node features = input tokens
- This matches the setting for the “attention is message passing on the fully connected graph” observation



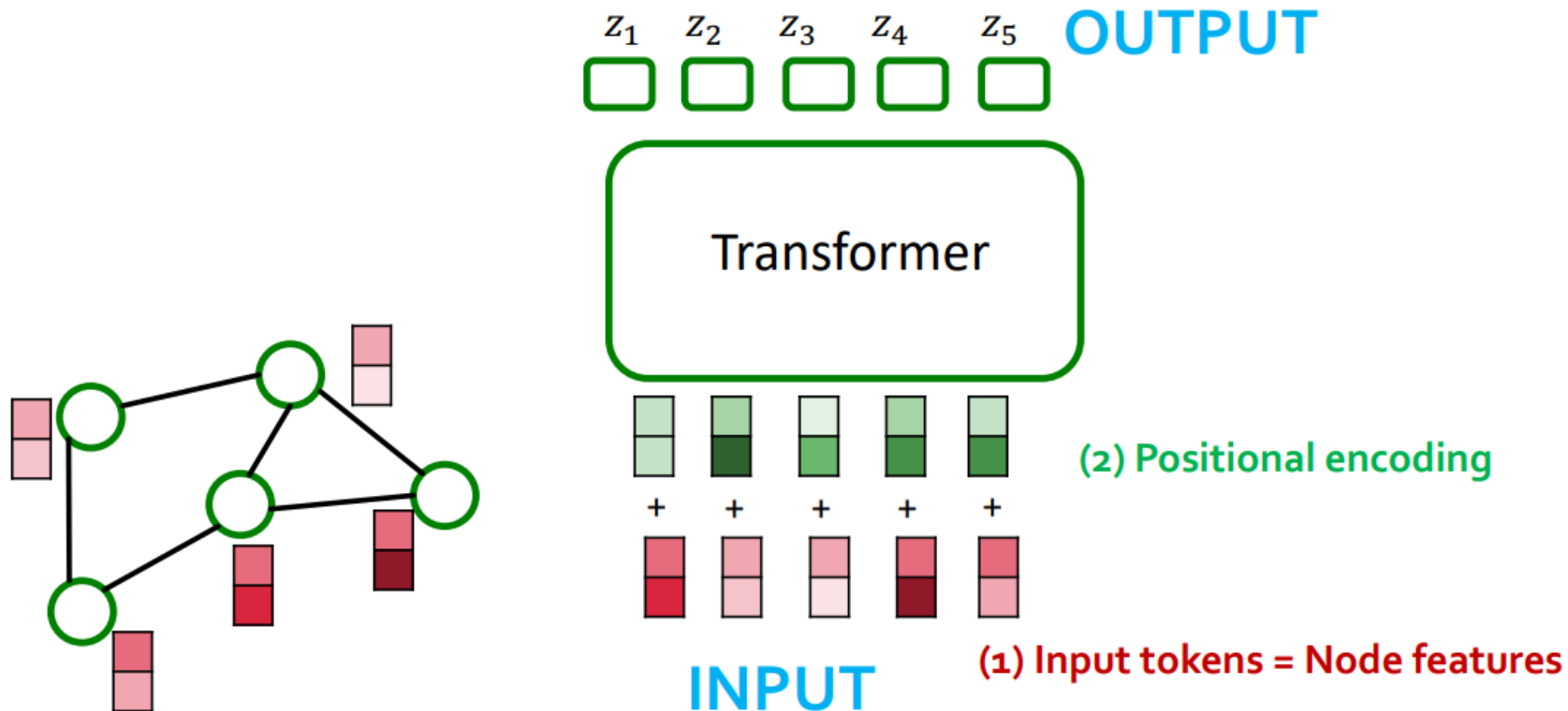
(1) Input tokens = Node features

- Q1: what should our tokens be?
- Sensible Idea: node features = input tokens
- This matches the setting for the “attention is message passing on the fully connected graph” observation
- **Problem?** We completely lose adjacency info!
- **How to also inject adjacency information?**



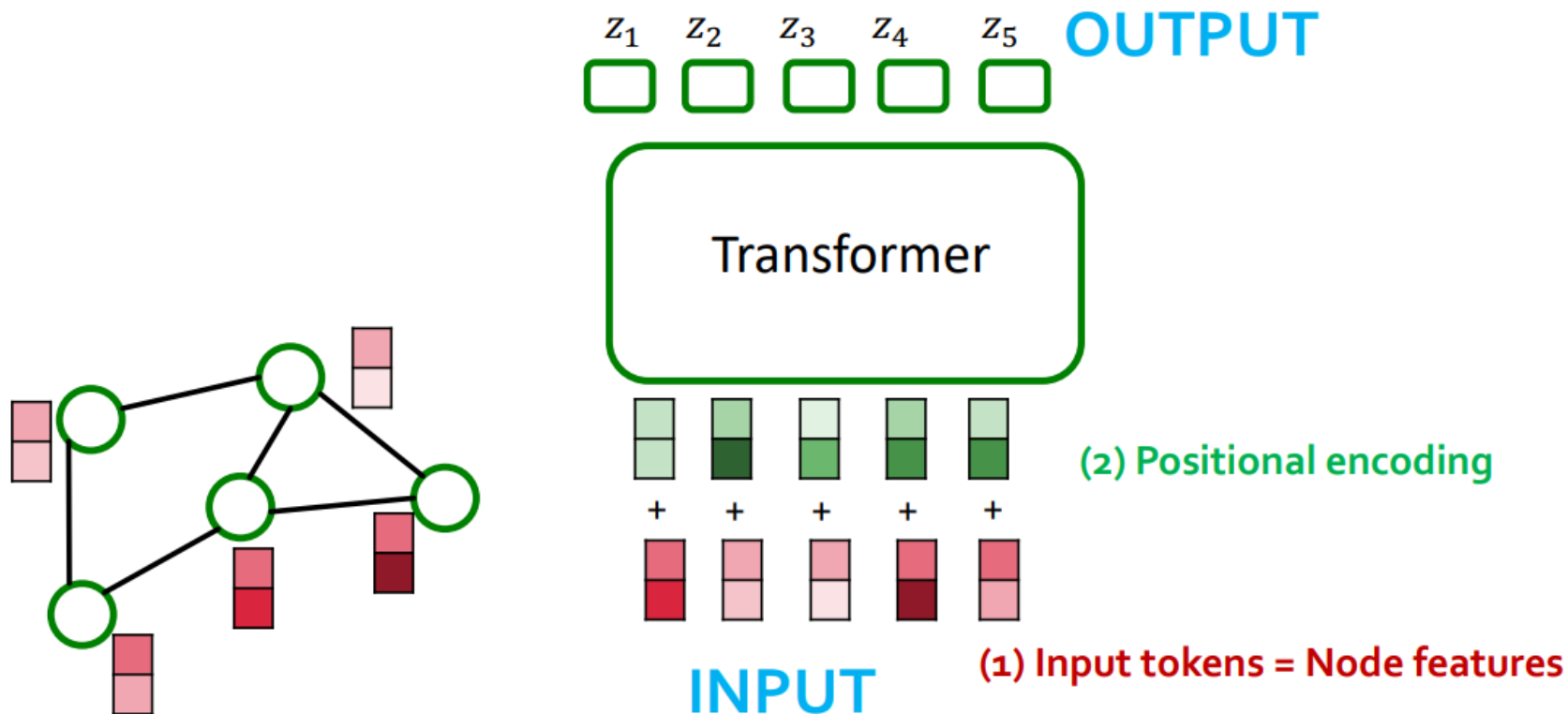
(1) Input tokens = Node features

- **Problem?** We completely lose adjacency info!
- **How to also inject adjacency information?**
- **Idea:** Encode adjacency info in the positional encoding for each node
- Positional encoding describes where a node is in the graph

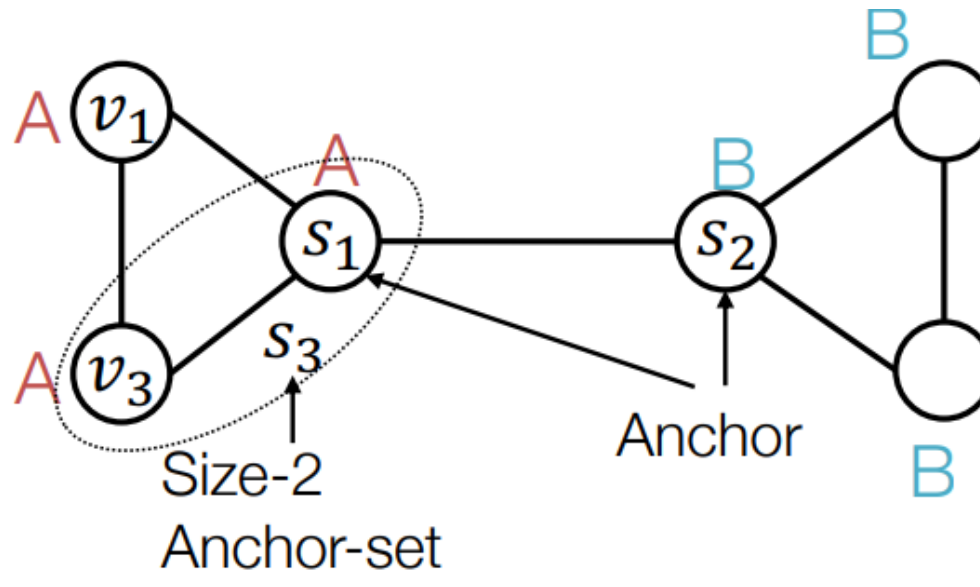


➤ **Q2: How to design a good positional encoding?**

- Option 1: relative distance
- Option 2: Laplacian Eigenvector PE



- Similar methods based on random walks
- This is a good idea. It works well in many cases
- Especially strong for tasks that require counting cycles



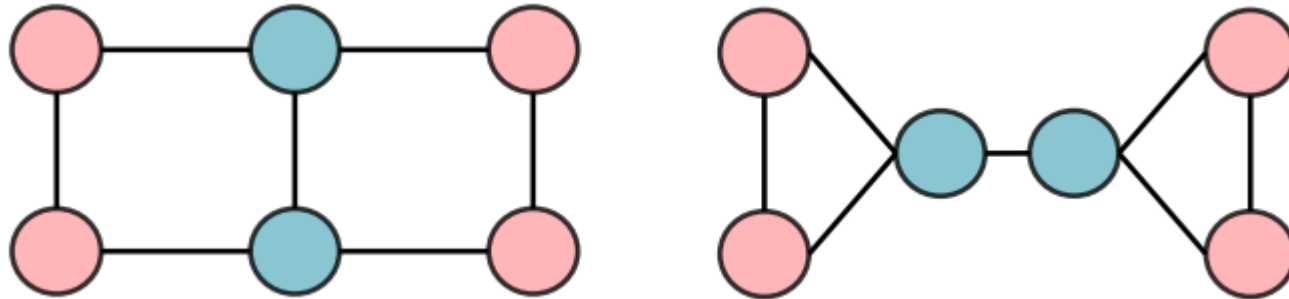
Positional encoding for node v_1

Relative Distances

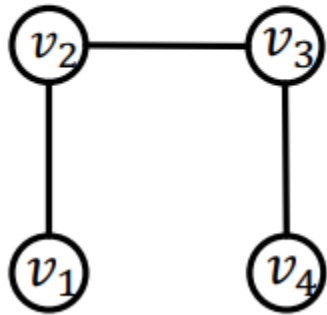
	s_1	s_2	s_3
v_1	1	2	1
v_3	1	2	0

Anchor s_1, s_2 cannot differentiate node v_1, v_3 , but anchor-set s_3 can

- Relative distances useful for position-aware task
- SPD can be used to improve WL-Test:
 - These two graphs cannot be distinguished by 1-WL-test.
 - But the SPD sets, i.e., the SPD from each node to others, are different:
 - The two types of nodes in the left graph have SPD sets $\{0, 1, 1, 2, 2, 3\}$, $\{0, 1, 1, 1, 2, 2\}$ while the nodes in the right graph have SPD sets $\{0, 1, 1, 2, 3, 3\}$, $\{0, 1, 1, 1, 2, 2\}$.



- Draw on knowledge of Graph Theory (many useful and powerful tools)
- Key object: Laplacian Matrix $L = \text{Degrees} - \text{Adjacency}$
 - Each graph has its own Laplacian matrix
 - Laplacian encodes the graph structure
 - Several Laplacian variants that add degree information differently



$L =$

1	0	0	0
0	2	0	0
0	0	2	0
0	0	0	1

Degree of each node

$-$

0	1	0	0
1	0	1	0
0	1	0	1
0	0	1	0

Adjacency

- Laplacian matrix captures graph structure
- Its eigenvectors inherit this structure
- This is important because eigenvectors are vectors and so can be fed into a Transformer
- Eigenvectors with small eigenvalue = local structure, large eigenvalue = global symmetries

Refresher

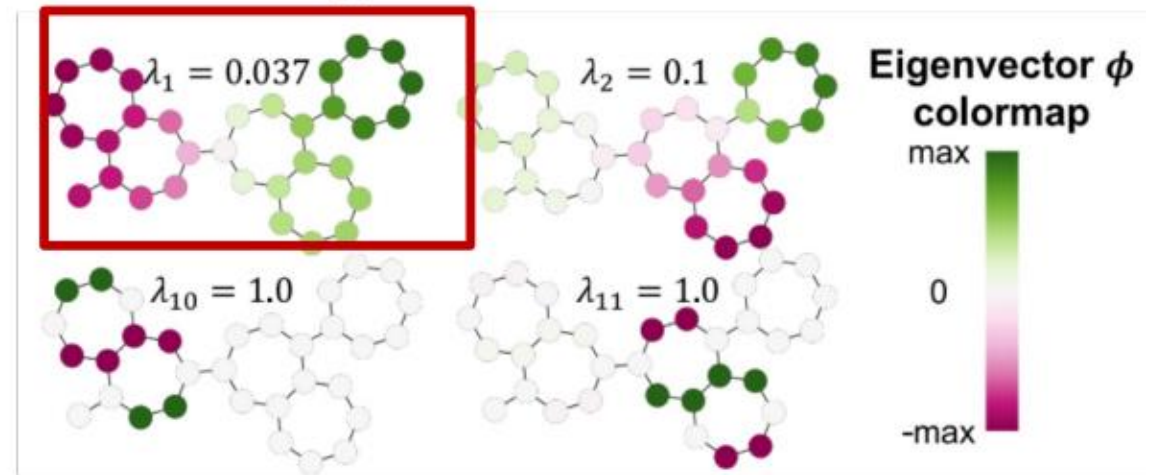
Eigenvector: v such that $Lv = \lambda v$

$L: n \times n$ matrix

$v: n$ dimensional vector

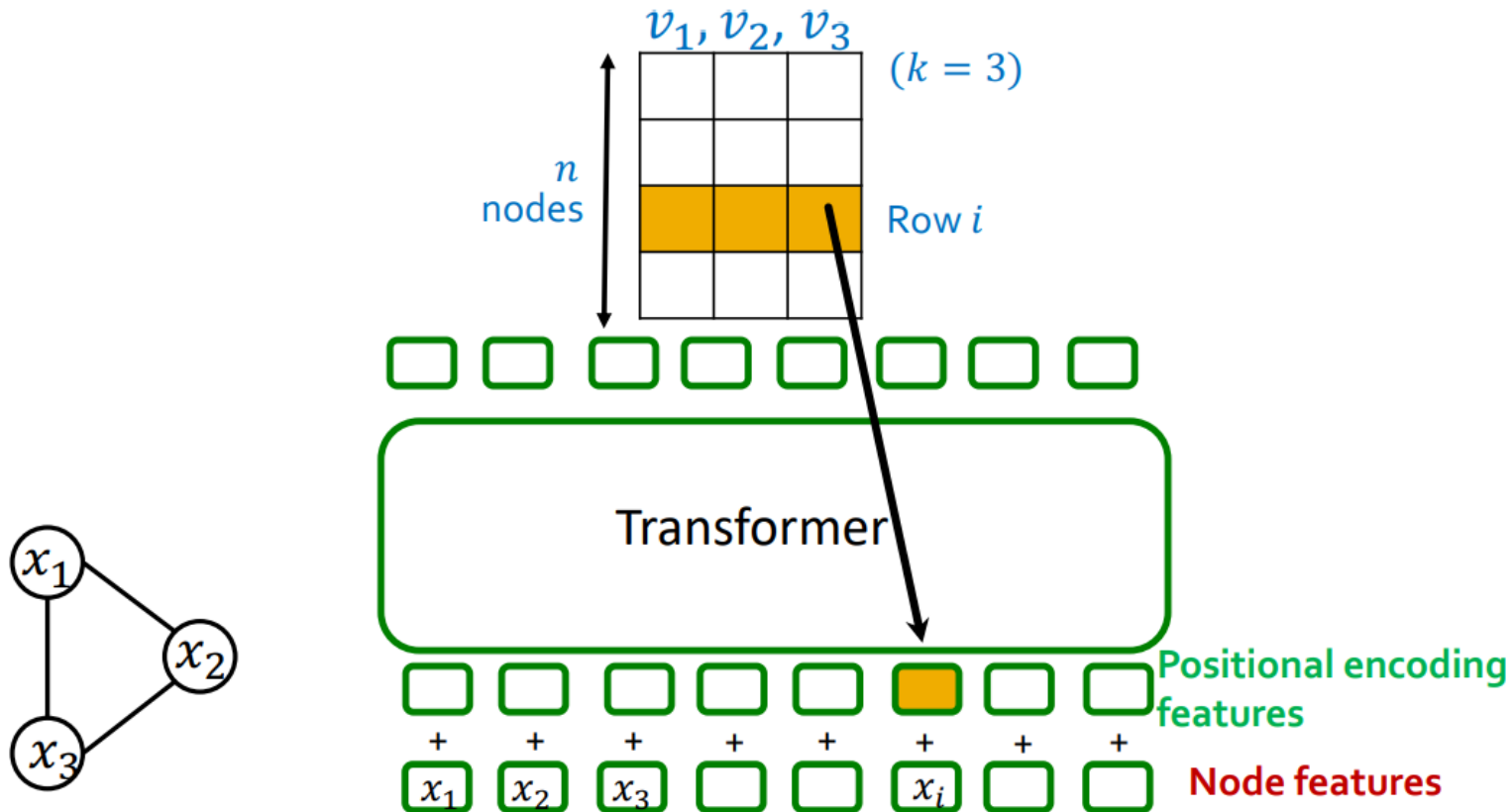
λ : Scalar eigenvalue

Visualize one eigenvector

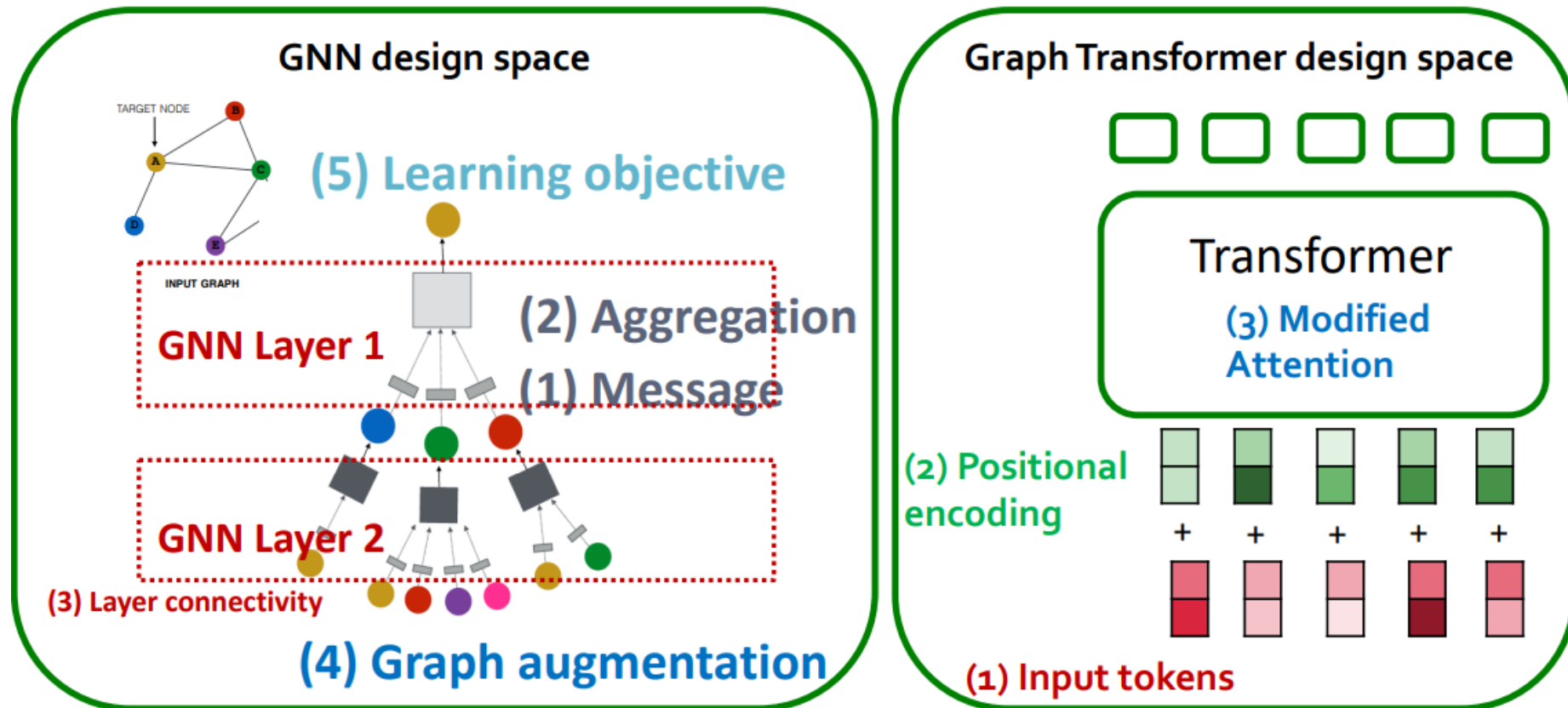


(Figure from Kreuzer* and Beaini* et al. 2021)

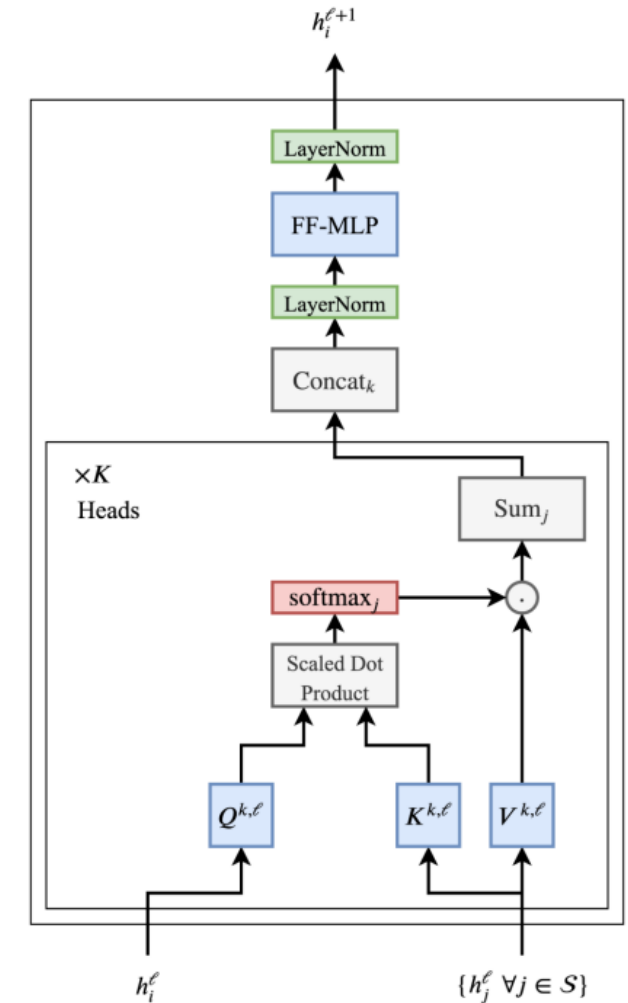
- Positional encoding steps:
 - 1. compute k eigenvectors
 - 2. Stack into matrix:
 - i -th row is positional encoding for node i



- Transformer are GNNs



- **Breaking down the Transformer:** Update each node's features through Multi-head Attention mechanism as a weighted sum of features of other words in the sentence.
 - Scaling dot product attention
 - Normalization layers
 - Residual links

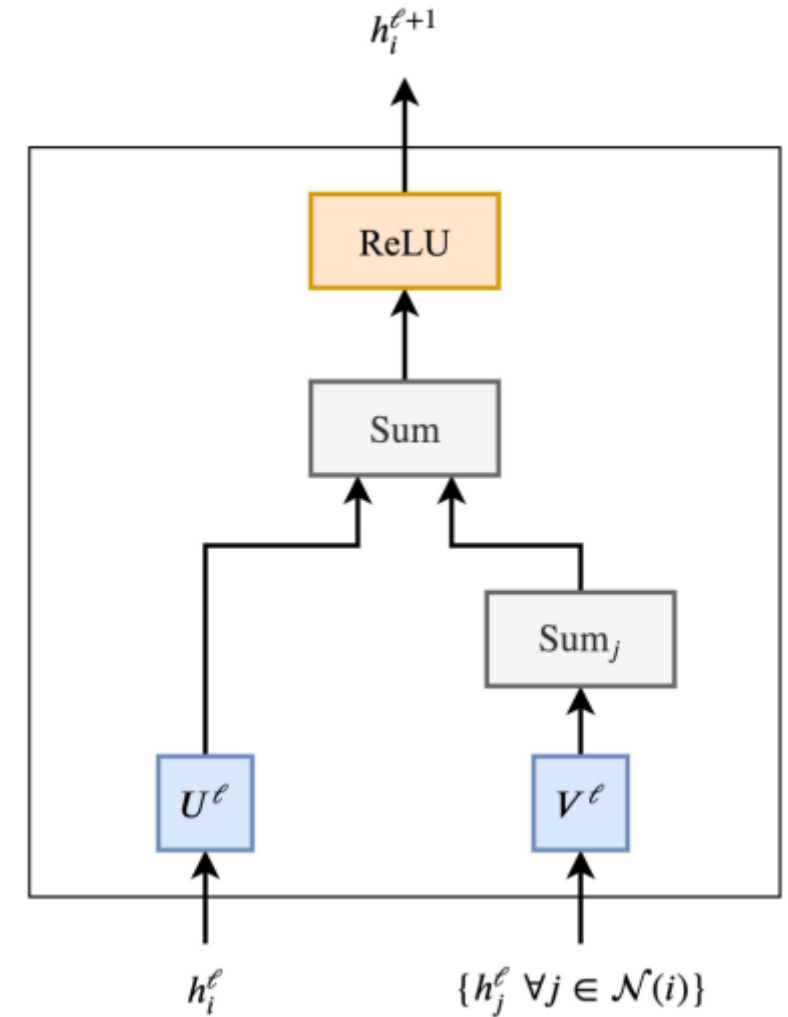


Transformer

- **Breaking down the GNNs:** GNNs update the hidden features h of node i at layer l via a non-linear transformation of the node's own features added to the aggregation of features from each neighbouring node $j \in N(i)$:

$$h_i^{\ell+1} = \sigma \left(U^\ell h_i^\ell + \sum_{j \in \mathcal{N}(i)} (V^\ell h_j^\ell) \right),$$

- where U, V are learnable weight matrices of the GNN layer and σ is a non-linearity.



GNNs

➤ Breaking down the Transformer and GNNs:

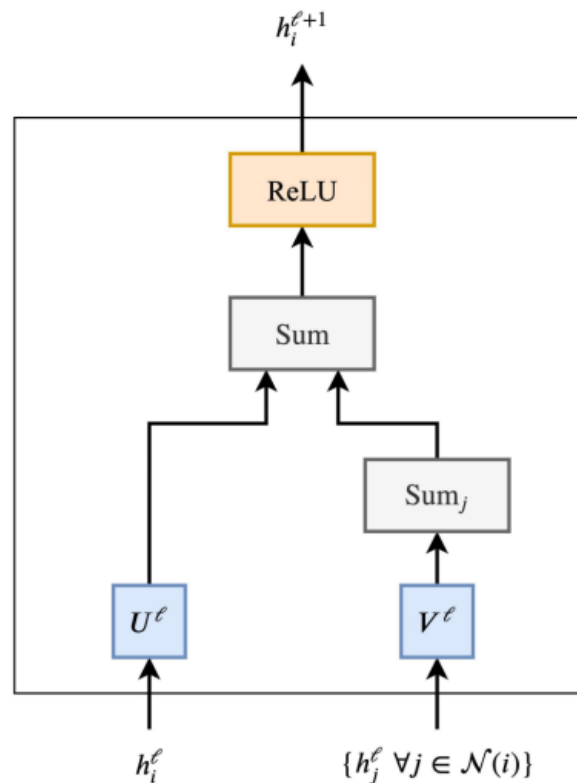
➤ GNNs:

$$h_i^{\ell+1} = \sigma \left(U^\ell h_i^\ell + \sum_{j \in \mathcal{N}(i)} (V^\ell h_j^\ell) \right),$$

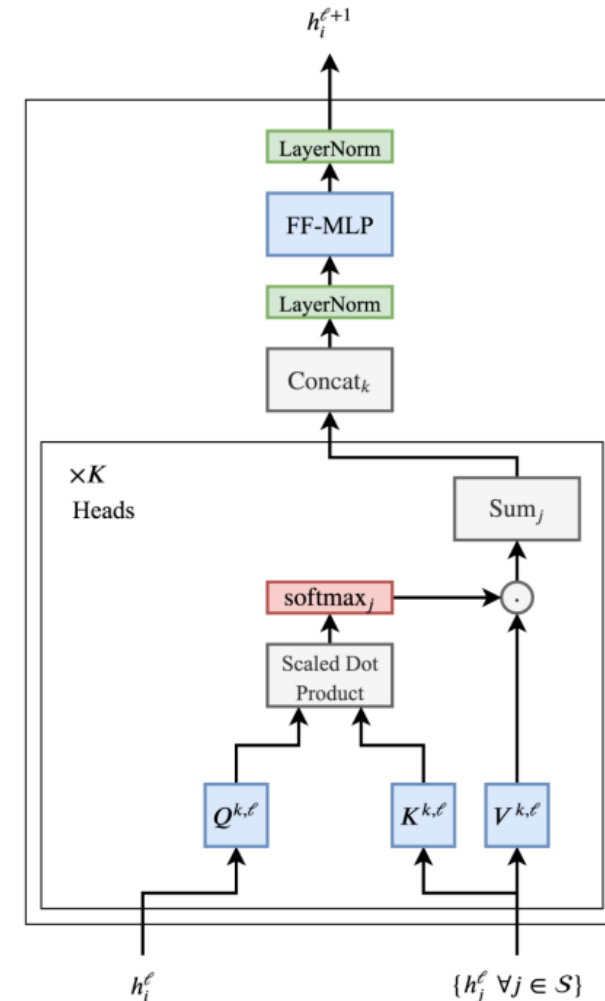
➤ Transformers:

$$i. e., h_i^{\ell+1} = \sum_{j \in \mathcal{S}} w_{ij} (V^\ell h_j^\ell),$$

where $w_{ij} = \text{softmax}_j(Q^\ell h_i^\ell \cdot K^\ell h_j^\ell)$,



GNNs



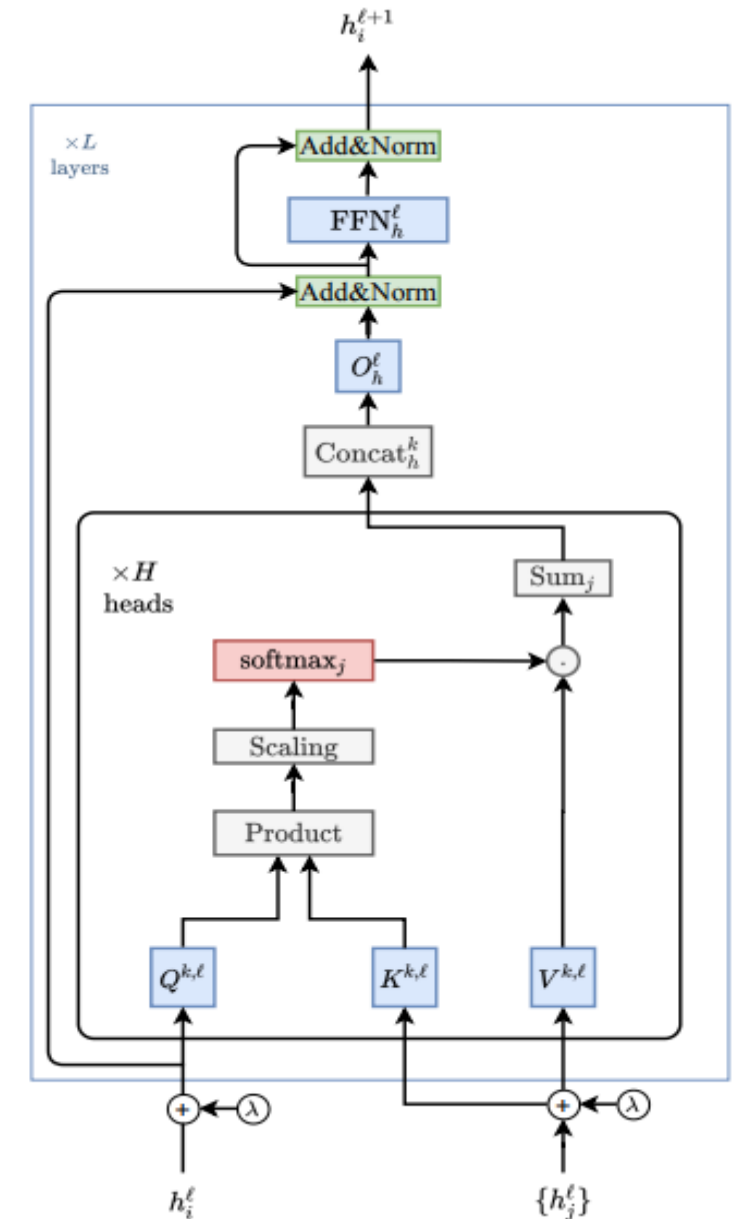
Transformer

GT (Graph Transformers *)

- Using Laplacian Eigenvectors (λ) used as positional encoding (LapPE).
- Graph Transformer Layer:

$$\hat{h}_i^{\ell+1} = O_h^\ell \parallel \left(\sum_{k=1}^H w_{ij}^{k,\ell} V^{k,\ell} h_j^\ell \right),$$

$$\text{where, } w_{ij}^{k,\ell} = \text{softmax}_j \left(\frac{Q^{k,\ell} h_i^\ell \cdot K^{k,\ell} h_j^\ell}{\sqrt{d_k}} \right)$$



Graphormer (*)

➤ Centrality Encoding:

$$h_i^{(0)} = x_i + z_{\text{deg}^-(v_i)}^- + z_{\text{deg}^+(v_i)}^+,$$

(learnable indegree z^- , and outdegree z^+)

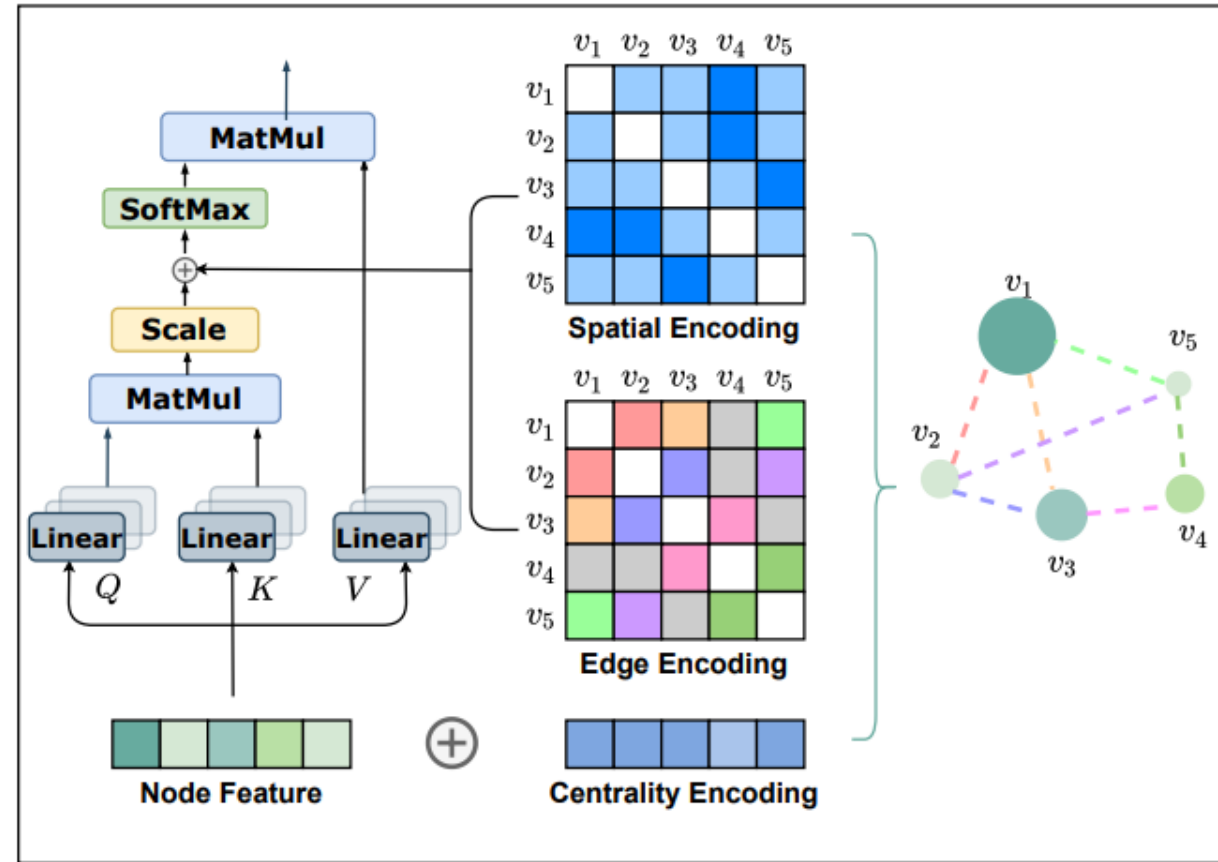
➤ Self-attention bias:

$$A_{ij} = \frac{(h_i W_Q)(h_j W_K)^T}{\sqrt{d}} + b_{\phi(v_i, v_j)} + c_{ij}$$

$$\text{where } c_{ij} = \frac{1}{N} \sum_{n=1}^N x_{e_n} (w_n^E)^T$$

presents the path between two nodes i and j
via edge feature path: $\text{SP}_{ij} = (e_1, e_2, \dots, e_N)$

$b_{\phi(v_i, v_j)}$: the distance of the shortest path (SPD) between two nodes i and j

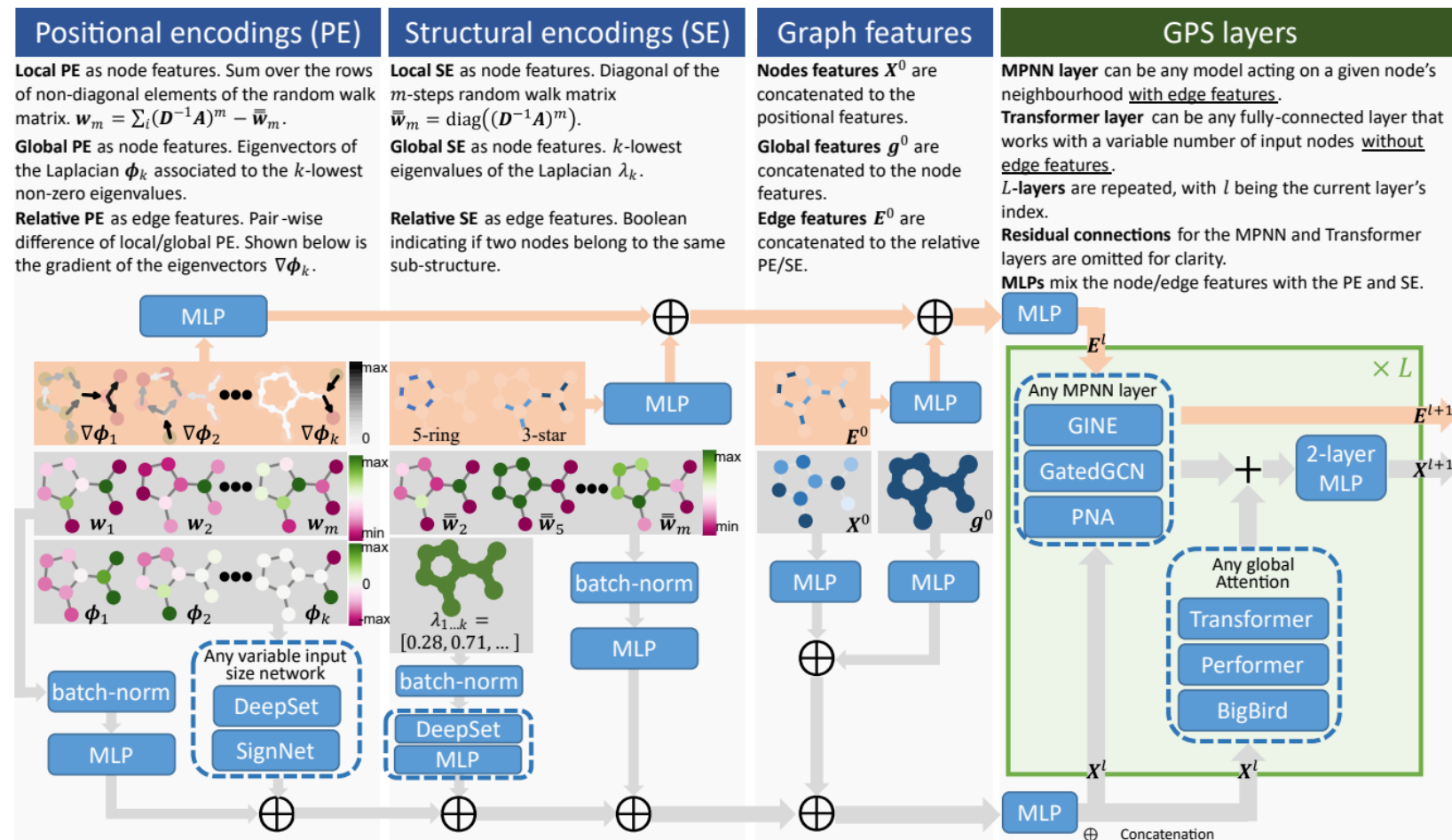


➤ GPS uses:

- Randomwalk PE
- GPS layers:
 - An MPNN+
 - Transformer hybrid

computed as

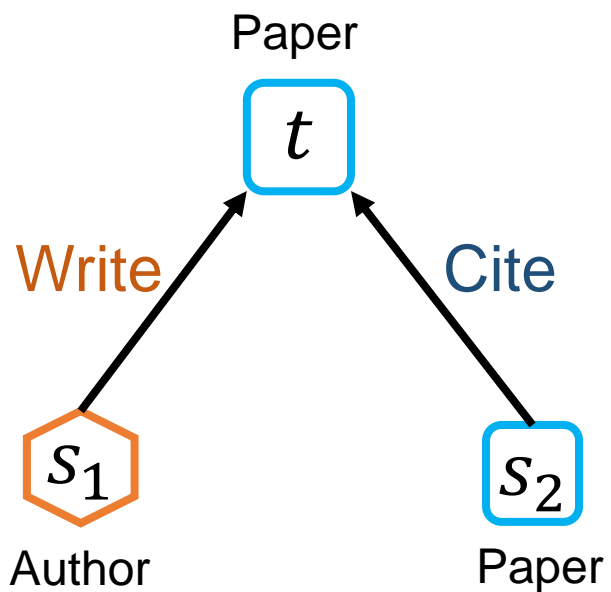
$$\begin{aligned} \mathbf{X}^{\ell+1}, \mathbf{E}^{\ell+1} &= \text{GPS}^{\ell}(\mathbf{X}^{\ell}, \mathbf{E}^{\ell}, \mathbf{A}) \\ \mathbf{X}_M^{\ell+1}, \mathbf{E}^{\ell+1} &= \text{MPNN}_e^{\ell}(\mathbf{X}^{\ell}, \mathbf{E}^{\ell}, \mathbf{A}), \\ \mathbf{X}_T^{\ell+1} &= \text{GlobalAttn}^{\ell}(\mathbf{X}^{\ell}), \\ \mathbf{X}^{\ell+1} &= \text{MLP}^{\ell}(\mathbf{X}_M^{\ell+1} + \mathbf{X}_T^{\ell+1}), \end{aligned}$$



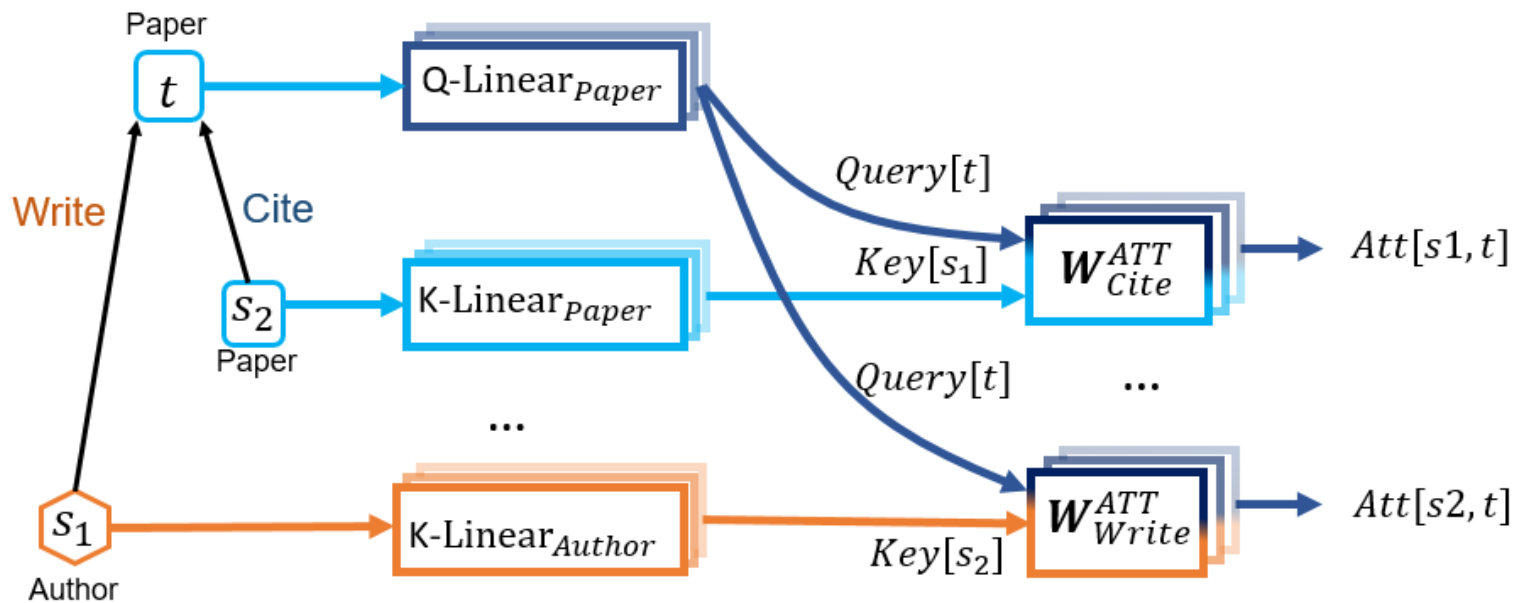
GPS: a General, Powerful, Scalable graph Transformer

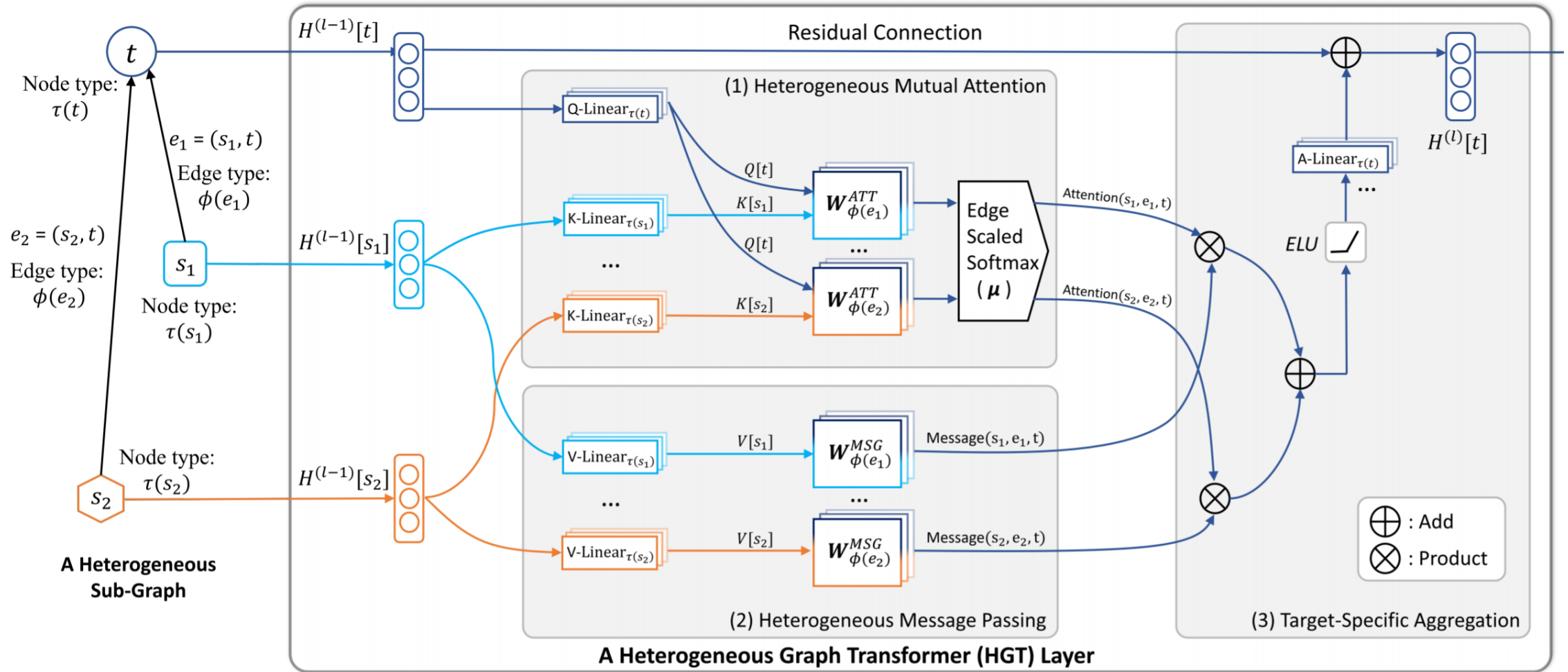
➤ Heterogeneous Mutual Attention in heterogeneous Graphs

$$W_{\langle \text{Author}, \text{Write}, \text{Paper} \rangle} = W_{\text{Author}} W_{\text{Write}} W_{\text{Paper}}$$



$$W_{\langle \text{Paper}, \text{Cite}, \text{Paper} \rangle} = W_{\text{Paper}} W_{\text{Cite}} W_{\text{Paper}}$$







네트워크 과학연구실
NETWORK SCIENCE LAB



가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA

