# GraphGPS : General Powerful Scalable Graph Transformers

One-time paper review

ii tae jeong
[jeongiitae6@gmail.com](mailto:jeongiitae6@gmail.com)
seoul , south korea

# Outline

## Preliminary

- ➢ Message passing GNN vs. Graph Transformers
- ➢ isomorphism test
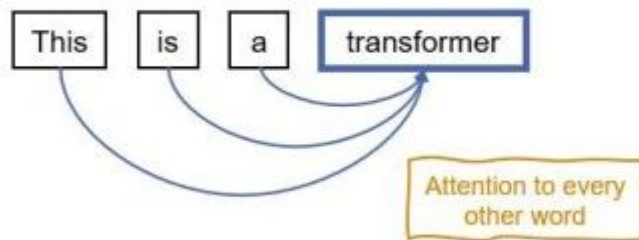- ➢ Linear Transformers

## Novelty

## Recipes

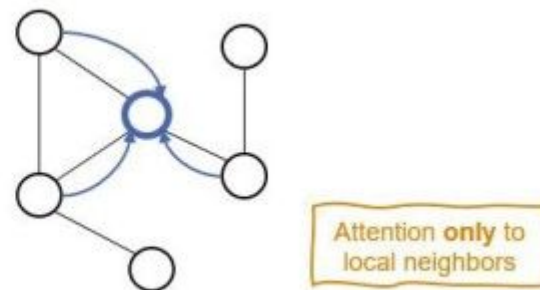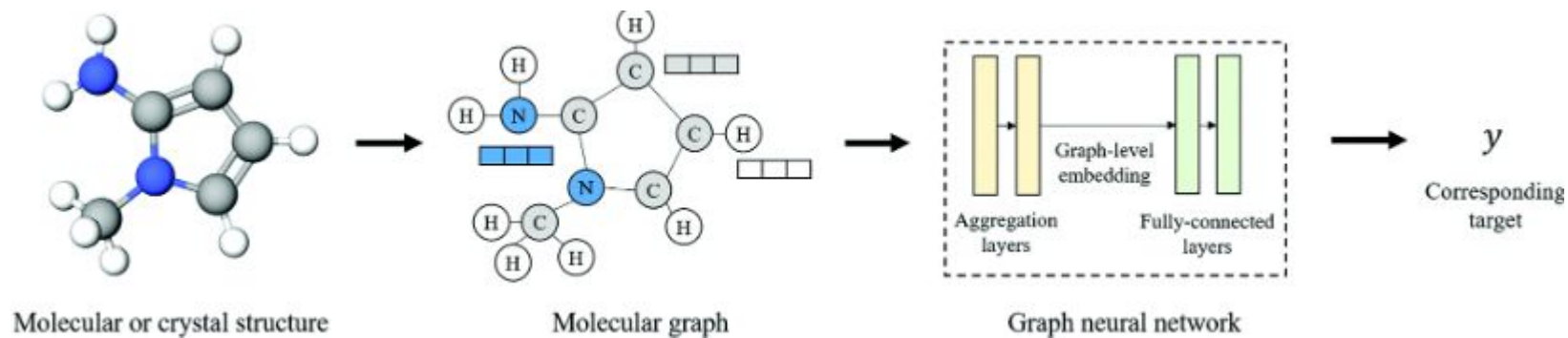## Open Discussion

# Preliminary

## Graph Transformer



<Sequence>
This is a transformer

v/s

<Multi-set>
{This, is, a, transformer}

why important these PE&SE information?



| Molecular or crystal structure | Molecular graph | Graph neural network |

Graph-level embedding

Aggregation layers

Fully-connected layers
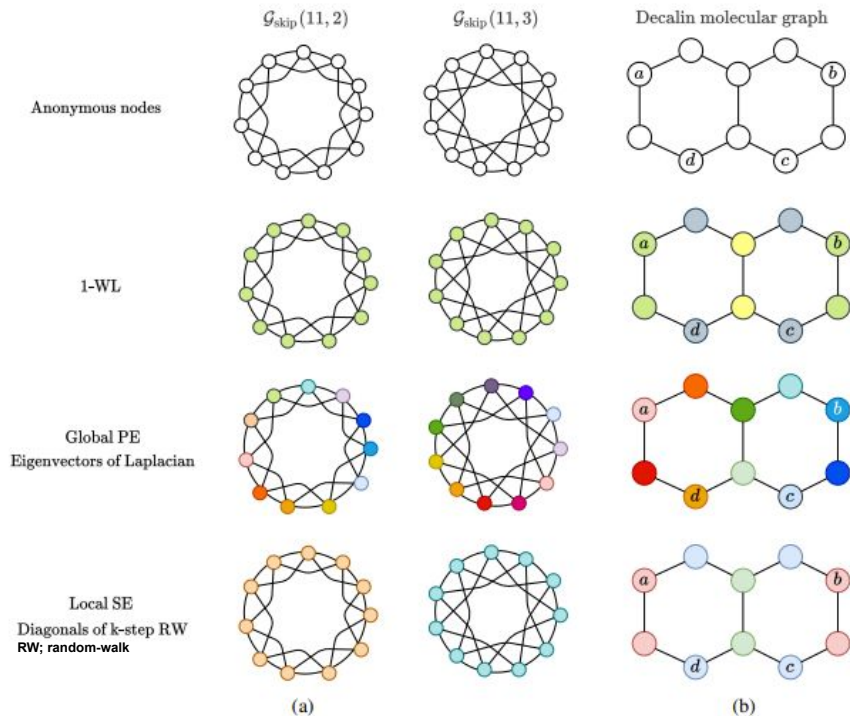
$y$

Corresponding target

MPNN vs. Graph Transformer



- **over-smoothing** (increasing the number of GNN layers, the features tend to converge to the same value),
- **over-squashing** (losing information when trying to aggregate messages from many neighbors into a single vector)
- poor capturing of long-range dependencies which is noticeable already on small but sparse molecular graphs.

**Keypoint**

Transformer architecture had been mainly used to text data at NLP industry. It was able to distinguish by specifying the data. **But GNN couldn't it since the node ordering.** so we need to that tools for specifying what the data identifiability.



$\mathcal{G}_{skip}(11,2)$      $\mathcal{G}_{skip}(11,3)$      Decalin molecular graph

Anonymous nodes

1-WL

Global PE
Eigenvectors of Laplacian

Local SE
Diagonals of k-step RW
**RW; random-walk**

(a)             (b)

❏    node) CSL(Circular Skip Link) graph for isomorphic task which is capture differentiate between two potential links
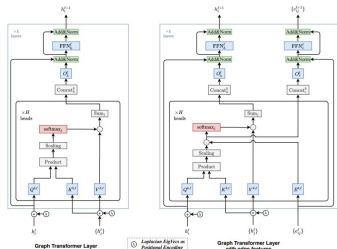
task → the feature embeddings of the two graphs **which are the hash function outputs of the collection of node colors are different**, thus making the task to distinguish the graphs successful

❏    edge) Decalin molecular graph, the node a is isomorphic to node b, and so is the node c to node d.

task → Identifying a potential link between the node-set(a,d) and (b,d), **the combination of the node colors of the node-sets will produce the same embedding for the two links.**

# Graph Transformer roadmap



Laplacian eigenvector
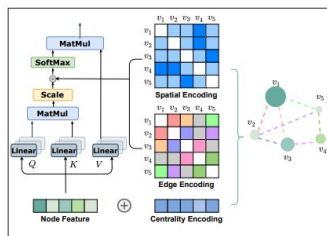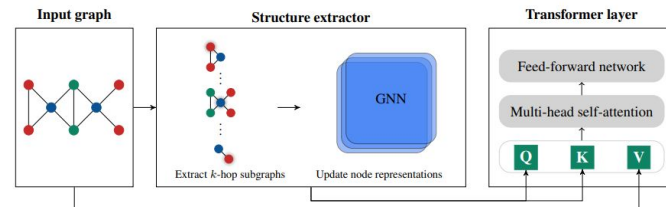
shortest path distance

aggregating a k-hop subgraph around each node
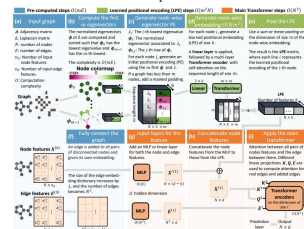
**2021, Jun**

**2022, Jan**

**2020, dec**

**2021, Nov**

**2022, Jun**

Laplacian eigenvalues to re-weight attention
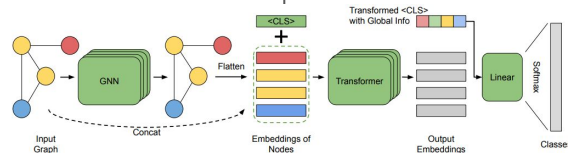
Run a GT after passing a graph through a GNN

# Preliminary

## Bigbird



(a) Random attention      (b) Window attention      (c) Global Attention      (d) BIGBIRD

Figure 1: Building blocks of the attention mechanism used in BIGBIRD. White color indicates absence of attention. (a) random attention with $r = 2$, (b) sliding window attention with $w = 3$ (c) global attention with $g = 2$. (d) the combined BIGBIRD model.

## Performer



Figure 1: Approximation of the regular attention mechanism $\mathbf{AV}$ (before $\mathbf{D}^{-1}$-renormalization) via (random) feature maps. Dashed-blocks indicate order of computation with corresponding time complexities attached.

- A set of g global tokens attending on all parts of the sequence.
- All tokens attending to a set of w local neighboring tokens.
- All tokens attending to a set of r random tokens.

$$\text{ATTN}_D(\boldsymbol{X})_i = \boldsymbol{x}_i + \sum_{h=1}^{H} \sigma \left( Q_h(\boldsymbol{x}_i) K_h(\boldsymbol{X}_{N(i)})^T \right) \cdot V_h(\boldsymbol{X}_{N(i)})$$

| Model | MLM | SQuAD | MNLI |
|---|---|---|---|
| BERT-base | 64.2 | 88.5 | 83.4 |
| Random (R) | 60.1 | 83.0 | 80.2 |
| Window (W) | 58.3 | 76.4 | 73.1 |
| R + W | 62.7 | 85.1 | 80.5 |

**Let N(i) denote the out-neighbors set of node i in D, then the i^th output vector of the generalized attention mechanism**

Table 1: Building block comparison @512

softmax kernel approximation

# Novelty

1. Scalable (linear global attention)
   → efficient implementation at graph transformer speed ) 400% faster without explicit edge features within the attention module.

2. Generalization (isomorphism)
   → using the positional , structural encoding tricks and ablation study

3. Powerful
   →    SOTA

# Recipes

**Positional encodings (PE)**

Local PE as node features. Sum over the rows of non-diagonal elements of the random walk matrix. $w_m = \sum_i (D^{-1}A)^m - \bar{w}_m$.
Global PE as node features. Eigenvectors of the Laplacian $\phi_k$ associated to the $k$-lowest non-zero eigenvalues.
Relative PE as edge features. Pair-wise difference of local/global PE. Shown below is the gradient of the eigenvectors $\nabla\phi_k$.
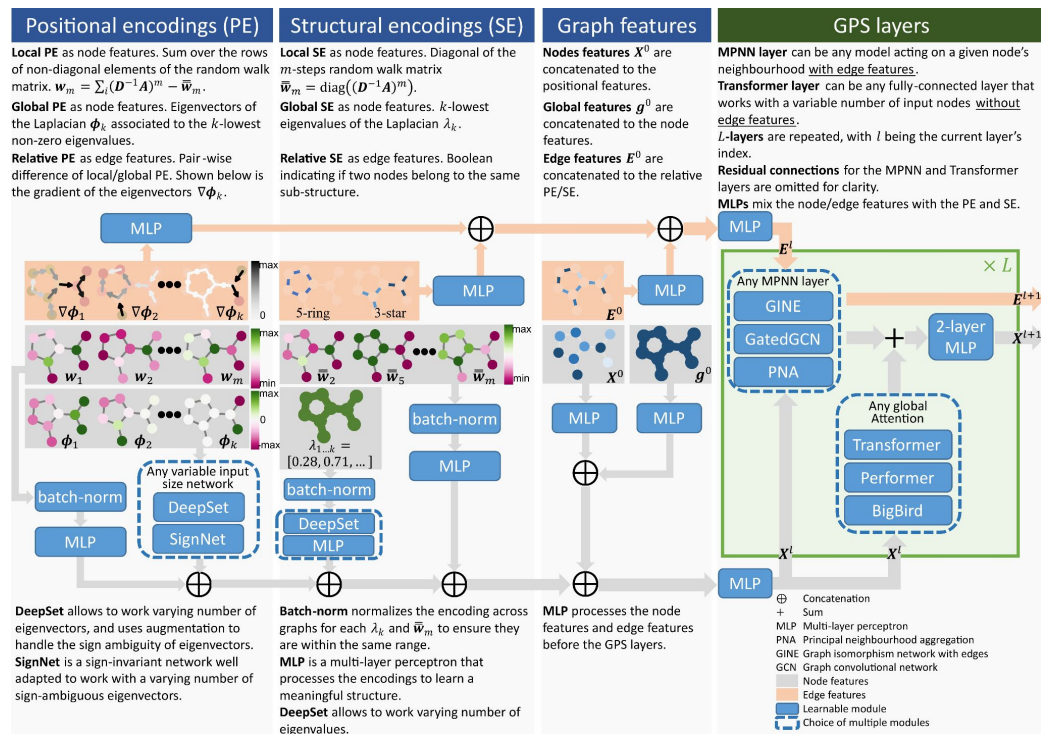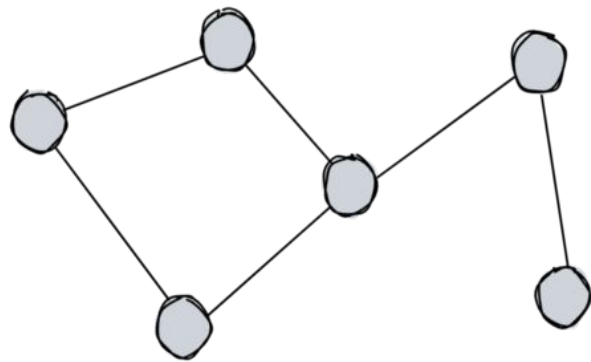
MLP

$\nabla\phi_1$ $\nabla\phi_2$ $\nabla\phi_k$

$w_1$ $w_2$ $w_m$

$\phi_1$ $\phi_2$ $\phi_k$

Any variable input size network
DeepSet
SignNet

batch-norm

MLP

DeepSet allows to work varying number of eigenvectors, and uses augmentation to handle the sign ambiguity of eigenvectors.
SignNet is a sign-invariant network well adapted to work with a varying number of sign-ambiguous eigenvectors.

**Structural encodings (SE)**

Local SE as node features. Diagonal of the $m$-steps random walk matrix $\bar{w}_m = \mathrm{diag}((D^{-1}A)^m)$.
Global SE as node features. $k$-lowest eigenvalues of the Laplacian $\lambda_k$.

Relative SE as edge features. Boolean indicating if two nodes belong to the same sub-structure.

MLP

5-ring  3-star

$\bar{w}_2$ $\bar{w}_5$ $\bar{w}_m$

$\lambda_{1...k} = [0.28, 0.71, ...]$

batch-norm

MLP

batch-norm

DeepSet
MLP

Batch-norm normalizes the encoding across graphs for each $\lambda_k$ and $\bar{w}_m$ to ensure they are within the same range.
MLP is a multi-layer perceptron that processes the encodings to learn a meaningful structure.
DeepSet allows to work varying number of eigenvalues.

**Graph features**

Nodes features $X^0$ are concatenated to the positional features.
Global features $g^0$ are concatenated to the node features.
Edge features $E^0$ are concatenated to the relative PE/SE.

MLP

$E^0$

$X^0$ $g^0$

MLP MLP

MLP processes the node features and edge features before the GPS layers.

**GPS layers**

MPNN layer can be any model acting on a given node's neighbourhood with edge features.
Transformer layer can be any fully-connected layer that works with a variable number of input nodes without edge features.
$L$-layers are repeated, with $l$ being the current layer's index.
Residual connections for the MPNN and Transformer layers are omitted for clarity.
MLPs mix the node/edge features with the PE and SE.

MLP

$E^l$

$\times L$

Any MPNN layer
GINE
GatedGCN
PNA

+

2-layer MLP

$E^{l+1}$

$X^{l+1}$

Any global Attention
Transformer
Performer
BigBird

$X^l$  $X^l$

MLP

$\oplus$ Concatenation
+ Sum
MLP Multi-layer perceptron
PNA Principal neighbourhood aggregation
GINE Graph isomorphism network with edges
GCN Graph convolutional network
Node features
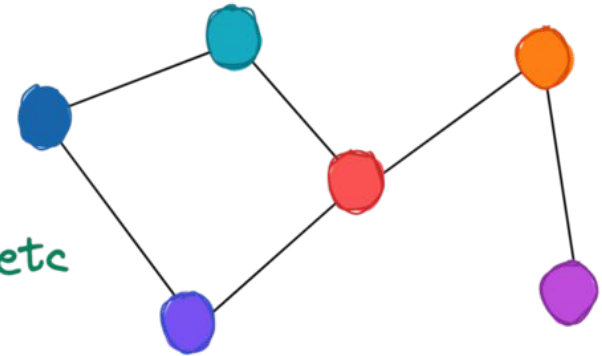Edge features
Learnable module
Choice of multiple modules

1. Node identification through positional and structural encodings. After analyzing many recently published methods for adding positionality in graphs, we found they can be broadly grouped into 3 buckets: local, global, and relative. Such features are provably powerful and help to overcome the notorious 1-WL limitation.

2. Aggregation of node identities with original graph features — those are your input node, edge, and graph features.

3. Processing layers (GPS layers) — how we actually process the graphs with constructed features, here combine both local message passing (any MPNNs) and global attention models (any graph transformer)
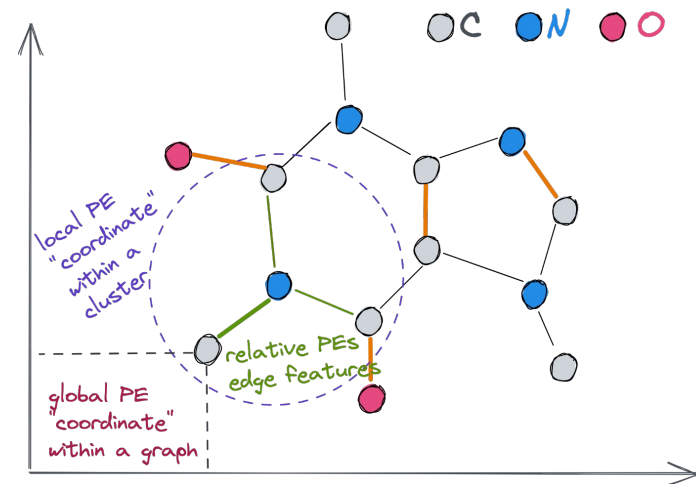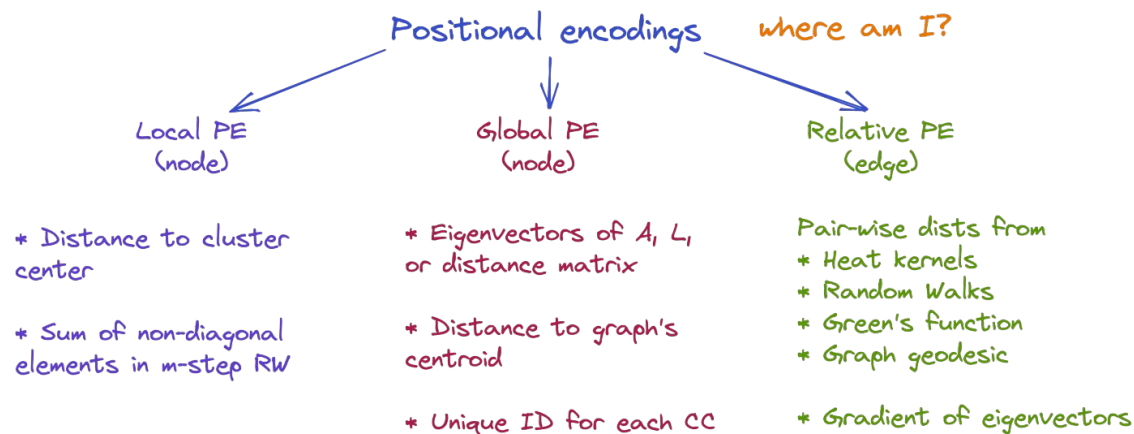
component 1. SE & PE encoding



structural & positional features

RWSE, Laplacian, SignNet, etc

Positional encodings        where am I?

Local PE
(node)

Global PE
(node)

Relative PE
(edge)

* Distance to cluster
center

* Sum of non-diagonal
elements in m-step RW

* Eigenvectors of $A$, $L$,
or distance matrix

* Distance to graph's
centroid

* Unique ID for each CC

Pair-wise dists from
* Heat kernels
* Random Walks
* Green's function
* Graph geodesic

* Gradient of eigenvectors



local PE
"coordinate"
within a
cluster

relative PEs
edge features

global PE
"coordinate"
within a graph

C    N    O

Structural encodings — what does my neighborhood look like?

Local SE (node)

Global SE (graph)

Relative SE (edge)

* Node degree

* RW diagonals

* Ricci curvature

* Enumerate substructures (triangles, rings)

* Eigenvalues of $A$, $L$

* Graph diameter, girth, degree, #CC

* Gradient of any Local SE

* Gradient of sub-structure enumeration

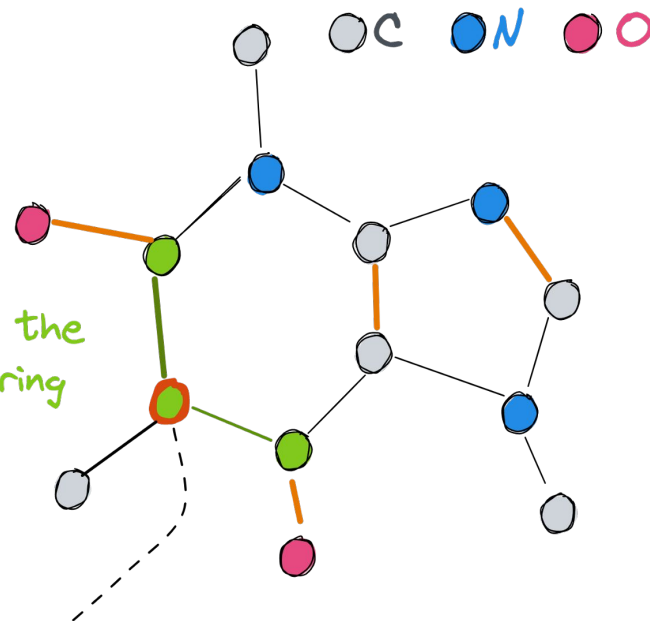are in the same ring



$Deg(\bullet) = 3$     Diameter $(G) = 6$

Table 1: The proposed categorization of positional encodings (PE) and structural encodings (SE). Some encodings are assigned to multiple categories in order to show their multiple expected roles.

| Encoding type | Description | Examples |
|---|---|---|
| **Local PE** node features | Allow a node to know its position and role within a local cluster of nodes. *Within a **cluster**, the closer two nodes are to each other, the closer their local PE will be, such as the position of a word in a sentence (not in the text).* | • Sum each column of non-diagonal elements of the $m$-steps random walk matrix.<br>• Distance between a node and the centroid of a cluster containing the node. |
| **Global PE** node features | Allow a node to know its global position within the graph. *Within a **graph**, the closer two nodes are, the closer their global PE will be, such as the position of a word in a text.* | • Eigenvectors of the Adjacency, Laplacian [14, 34] or distance matrices.<br>• Distance from the graph's centroid.<br>• Unique identifier for each connected component of the graph. |
| **Relative PE** edge features | Allow two nodes to understand their distances or directional relationships. *Edge embedding that is correlated to the distance given by any global or local PE, such as the distance between two words.* | • Pair-wise node distances from heat kernels, random-walks, Green's function, graph geodesic [3, 34, 41], or any local/global PE.<br>• Gradient of eigenvectors [3, 34] or any local/global PE.<br>• Boolean indicating if two nodes are in the same cluster. |
| **Local SE** node features | Allow a node to understand what substructures it is a part of. *Given an SE of radius $m$, the more similar the $m$-hop subgraphs around two nodes are, the closer their local SE will be.* | • Degree of a node [59].<br>• Diagonal of the $m$-steps random-walk matrix [15].<br>• Time-derivative of the heat-kernel diagonal (gives the degree at $t = 0$).<br>• Enumerate or count predefined structures such as triangles, rings, etc. [6, 64].<br>• Ricci curvature [51]. |
| **Global SE** graph features | Provide the network with information about the global structure of the graph. *The more similar two graphs are, the closer their global SE will be.* | • Eigenvalues of the Adjacency or Laplacian matrices [34].<br>• Graph properties: diameter, girth, number of connected components, number of nodes, number of edges, nodes-to-edges ratio. |
| **Relative SE** edge features | Allow two nodes to understand how much their structures differ. *Edge embedding that is correlated to the difference between any local SE.* | • Gradient of any local SE.<br>• Boolean indicating if two nodes are in the same substructure [5] (similar to the gradient of sub-structure enumeration). |

## positional/structural encoding

**LapPE,**
→ SAN(Spectral Attention Network) extended version, add to the node features of the graph and passed to fully-connected Transformer

**RWSE,**
→ Using the LSPE(Learnable STructural and Positional Encodings.) random-walk diffusion based positional encoding scheme.

**SignNet,**
→ general basis symmetries (fresh eigenvector selection)

**EquivStableLapPE**
→ Separating channel for update the original node features and positional features. and utilize extra positional features 2nd and 3rd smallest eigenvalues and Rotation equivariance fashion.

## local message-passing

**GatedGCN,**
→ gate will close to let the information flow from neighbor j to vertex i, or it will open to stop it.

$$h_i^{\ell+1} = f_{\text{G-GCNN}}^\ell \left( h_i^\ell, \{ h_j^\ell : j \to i \} \right) = \text{ReLU} \left( U^\ell h_i^\ell + \sum_{j \to i} \eta_{ij} \odot V^\ell h_j^\ell \right)$$

**GINE,**
→

$$\mathbf{x}_i' = h_\Theta \left( (1 + \epsilon) \cdot \mathbf{x}_i + \sum_{j \in \mathcal{N}(i)} \text{ReLU}(\mathbf{x}_j + \mathbf{e}_{j,i}) \right)$$

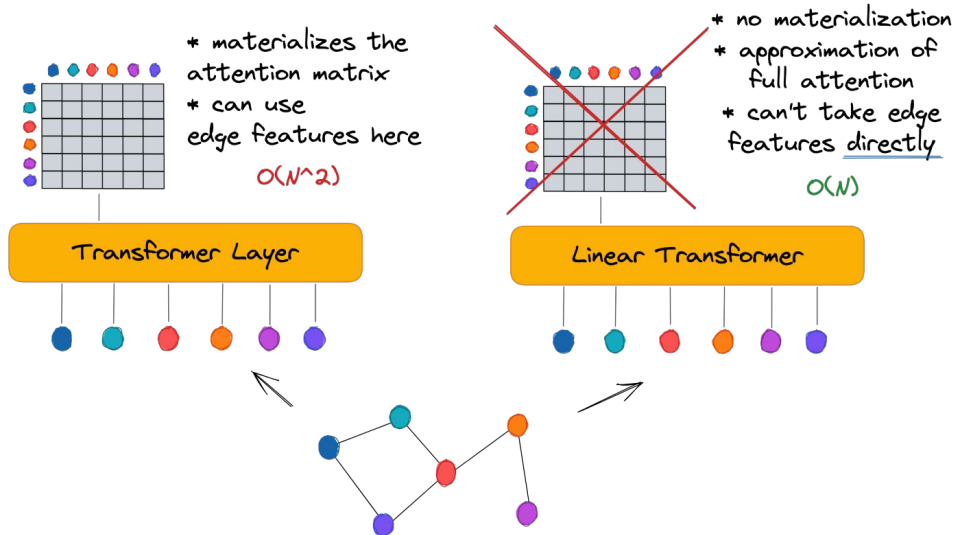that is able to incorporate edge features $\mathbf{e}_{j,i}$ into the aggregation procedure.

**PNA**
→ combining multiple aggregators with degree-scalers

## global attention

Transformer,
Performer,
BigBird

component 2. Linear Transformer



* materializes the attention matrix
* can use edge features here
$O(N^2)$

Transformer Layer

* no materialization
* approximation of full attention
* can't take edge features directly
$O(N)$

Linear Transformer

❏ interleaving tricks for utilizing edge information at transformer layer

$$h_u^{l+1} = \sum_{v \in \mathcal{N}_u} f(h_u^l, h_v^l, e_{uv}),$$

An example of such function $\mu_{uv}$ is the tensor product $\otimes$ of a one-hot encoding unique for each edge $o_{uv}$ and the edge features $e_{uv}$. For example, if $e_{uv} = [e_1, e_2, e_3]$ and the edge is represented with $o_{uv} = [0, 1, 0, 0]$, then $\mu_{uv} = o_{uv} \otimes e_{uv} = [0, 0, 0, e_1, e_2, e_3, 0, 0, 0, 0, 0, 0]$ satisfies all the above conditions. Although this function requires an exponential increase in the hidden dimension, this is also the case for the Lemma 5 in Xu et al. [57].

To model injective multiset functions for the neighbor aggregation, we develop a theory of "deep multisets", *i.e.*, parameterizing universal multiset functions with neural networks. Our next lemma states that sum aggregators can represent injective, in fact, *universal* functions over multisets.

**Lemma 5.** *Assume $\mathcal{X}$ is countable. There exists a function $f : \mathcal{X} \to \mathbb{R}^n$ so that $h(X) = \sum_{x \in X} f(x)$ is unique for each multiset $X \subset \mathcal{X}$ of bounded size. Moreover, any multiset function $g$ can be decomposed as $g(X) = \phi\left(\sum_{x \in X} f(x)\right)$ for some function $\phi$.*
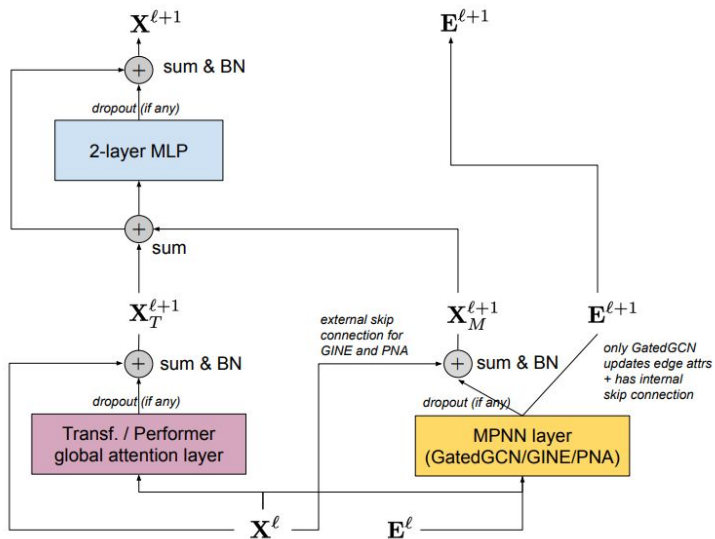
Figure D.1: Modular GPS layer that combines local MPNN and global attention blocks. Local MPNN encodes real edge features into the node-level hidden representations, while global attention mechanism can implicitly make use of this information together with PE/SE to infer relation between two nodes without explicit edge features. After each functional block (an MPNN layer, a global attention layer, an MLP) we apply residual connections followed by batch normalization (BN) [28]. In the 2-layer MLP block we use ReLU activations and its inner hidden dimension is twice the layer-input feature dimensionality $d_\ell$. Note, similarly to Transformer, the input and output dimensionality of the GPS-layer as a whole is the same.

$$\mathbf{X}^{\ell+1}, \mathbf{E}^{\ell+1} = \text{GPS}^\ell \left( \mathbf{X}^\ell, \mathbf{E}^\ell, \mathbf{A} \right) \tag{6}$$

$$\hat{\mathbf{X}}_M^{\ell+1}, \mathbf{E}^{\ell+1} = \text{MPNN}_e^\ell \left( \mathbf{X}^\ell, \mathbf{E}^\ell, \mathbf{A} \right), \tag{7}$$

$$\hat{\mathbf{X}}_T^{\ell+1} = \text{GlobalAttn}^\ell \left( \mathbf{X}^\ell \right), \tag{8}$$

$$\mathbf{X}_M^{\ell+1} = \text{BatchNorm} \left( \text{Dropout} \left( \hat{\mathbf{X}}_M^{\ell+1} \right) + \mathbf{X}^\ell \right), \tag{9}$$

$$\mathbf{X}_T^{\ell+1} = \text{BatchNorm} \left( \text{Dropout} \left( \hat{\mathbf{X}}_T^{\ell+1} \right) + \mathbf{X}^\ell \right), \tag{10}$$

$$\mathbf{X}^{\ell+1} = \text{MLP}^\ell \left( \mathbf{X}_M^{\ell+1} + \mathbf{X}_T^{\ell+1} \right) \tag{11}$$

# Open Discussion

(a) Distances in Poincaré disk.  (b) Distances projection.  (c) Distances with different curvature.
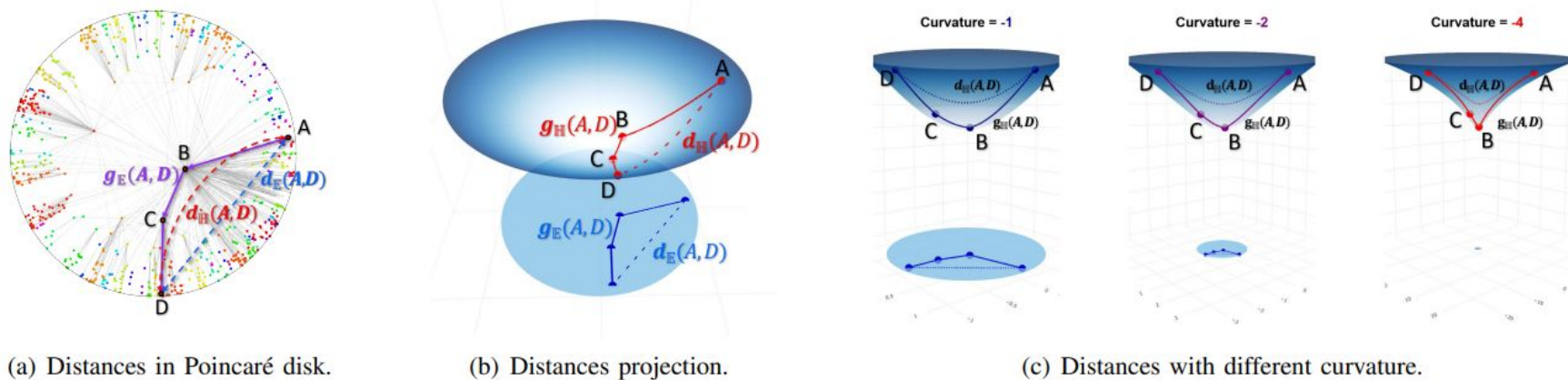
Figure 2. An illustration of different distance metrics in hyperbolic spaces with different curvature. (a) Graph distance (purple solid lines), Euclidean distance (blue dashed line) and hyperbolic geodesics (red dashed curve) on a tree-like graph in Poincaré disk. (b) Graph distance and embedded distance in Poincaré disk (Euclidean projection, blue solid and dashed lines) and hyperboloid (curvature $K=-1$, red solid and dashed curves). (c) Graph distance and hyperbolic distance on the hyperboloid of different curvature.

# References & ad

- https://towardsdatascience.com/graphgps-navigating-graph-transformers-c2cc223a051c

- Rampášek, L., Galkin, M., Dwivedi, V. P., Luu, A. T., Wolf, G., & Beaini, D. (2022). Recipe for a General, Powerful, Scalable Graph Transformer. *arXiv preprint arXiv:2205.12454*.

- Fu, X., Li, J., Wu, J., Sun, Q., Ji, C., Wang, S., ... & Philip, S. Y. (2021, December). ACE-HGNN: Adaptive curvature exploration hyperbolic graph neural network. In *2021 IEEE International Conference on Data Mining (ICDM)* (pp. 111-120). IEEE.

Next Pseudo Lab study group builder
**[Application] Value extraction from real graph data using Network theory & Graph neural network**
if u want to upgrade your network analysis & prediction skill , contact me !