

Graph Transformers

Dr. Leslie O'Bray
October 18, 2023



DBSSE

ETH zürich

 Isomorphic Labs

Outline of the lecture

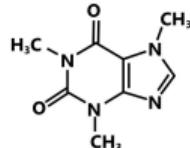
- 1** (Very brief) Context setting
- 2** Recap of GNNs & local methods
- 3** Limitations of modeling local information
- 4** Transformers & graph transformers
- 5** Structure-Aware Transformer

Graphs are widely applicable across domains

Social networks



Chemoinformatics



Product recommendations



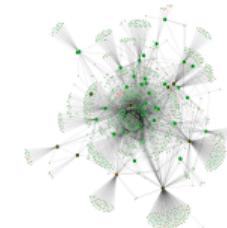
Traffic prediction



Weather forecasting

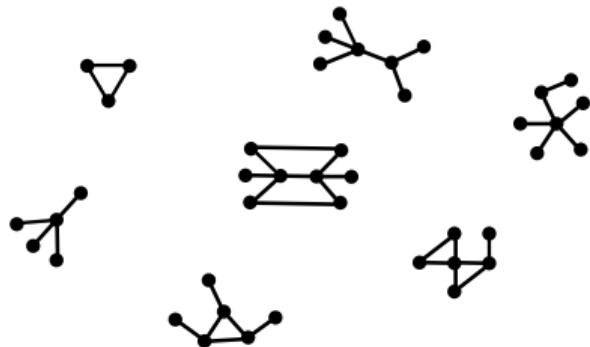


Protein-protein associations

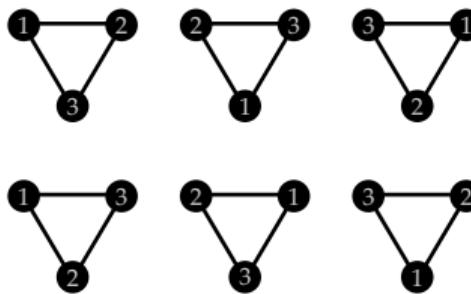


The flexibility of graphs presents unique challenges

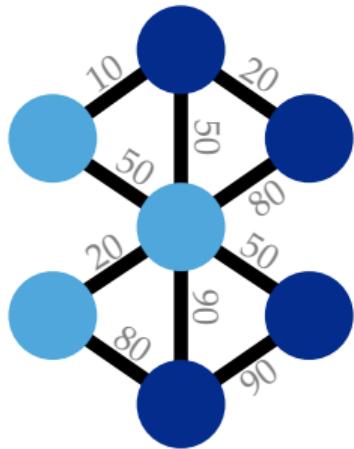
Variable sizes



No canonical ordering



Graph definitions and scope

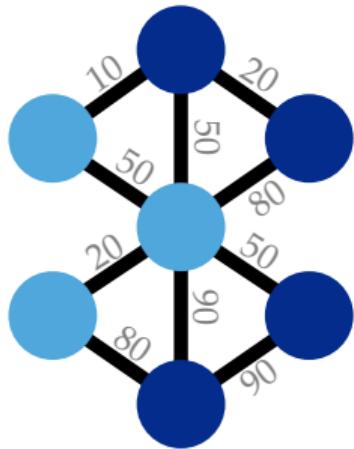


Today's lecture is concerned primarily with:

- * Graph classification
- * Simple graphs
- * Static graphs
- * Homogeneous graphs

- * $G = (V, E)$
- * Node labels or attributes X
- * (Optional) edge attributes or weights

Graph definitions and scope



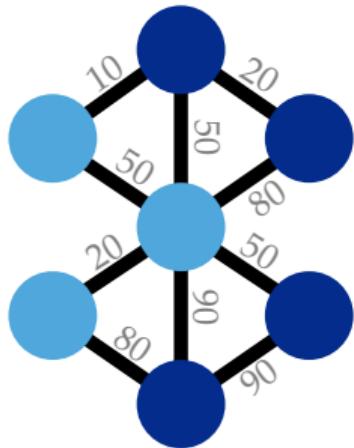
Today's lecture is concerned primarily with:

- * Graph classification
- * Simple graphs
- * Static graphs
- * Homogeneous graphs

- * $G = (V, E)$
- * Node labels or attributes X
- * (Optional) edge attributes or weights

...but can be extended to other types of graphs!

Graph definitions and scope



Today's lecture is concerned primarily with:

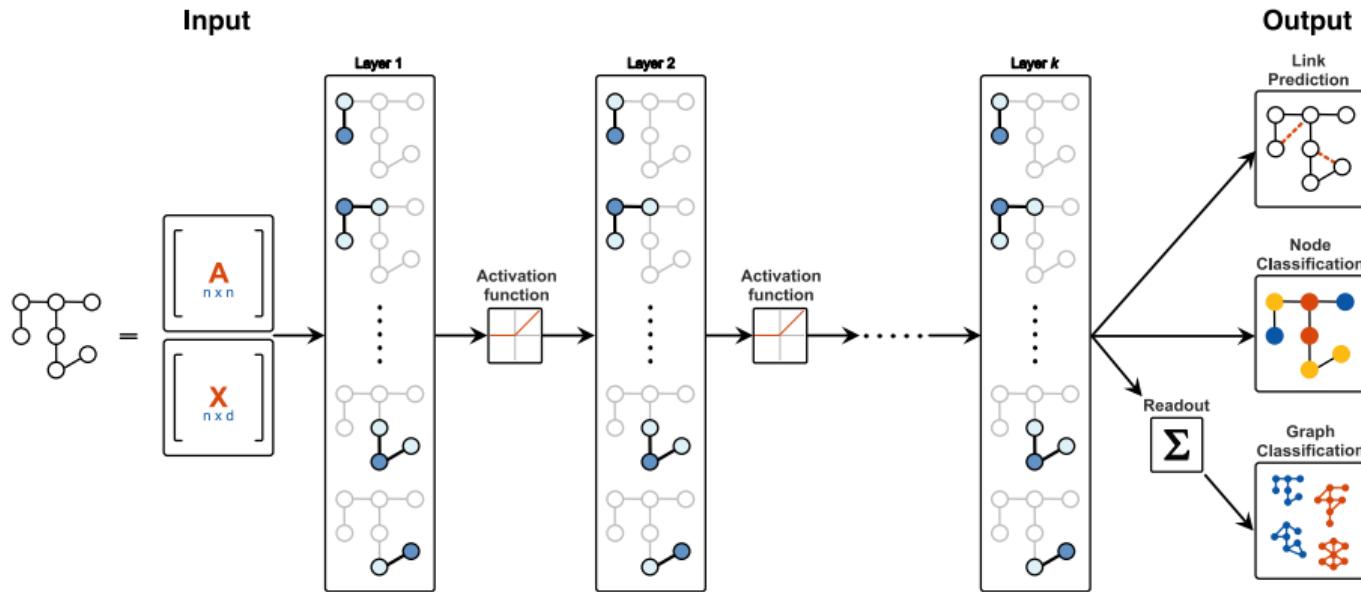
- * Graph classification
- * Simple graphs
- * Static graphs
- * Homogeneous graphs
- * $G = (V, E)$
- * Node labels or attributes X
- * (Optional) edge attributes or weights

...but can be extended to other types of graphs!

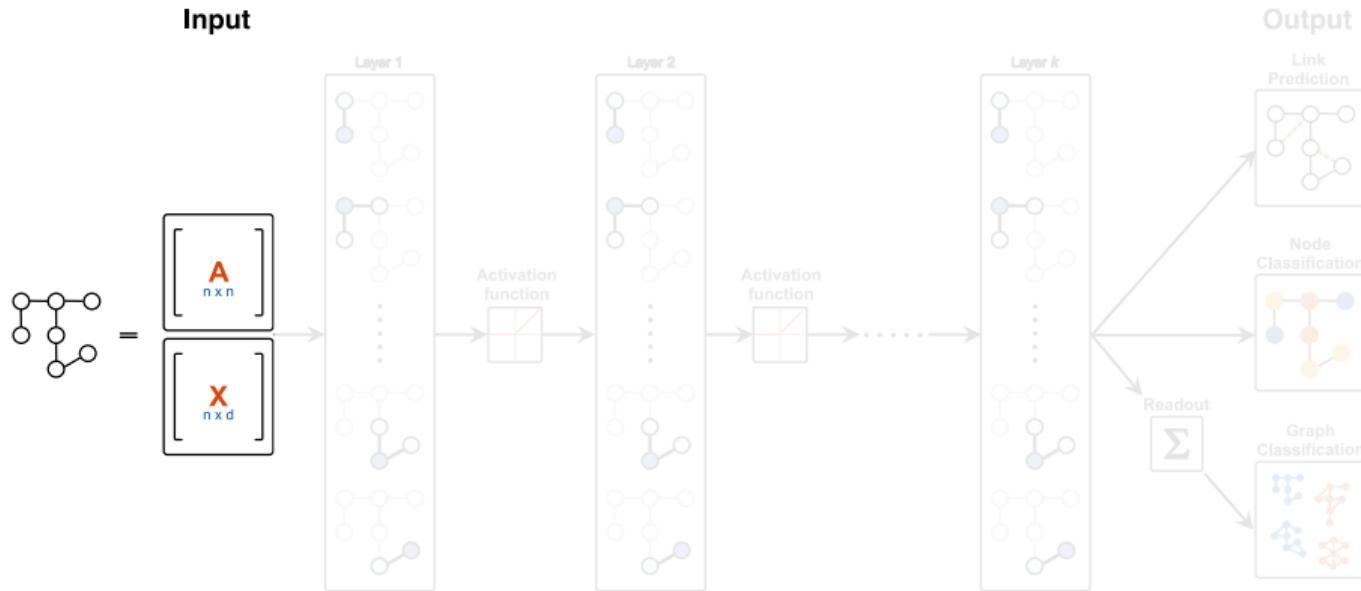
Graph representation learning: how to effectively incorporate the graph structure?

Recap: local message passing GNNs

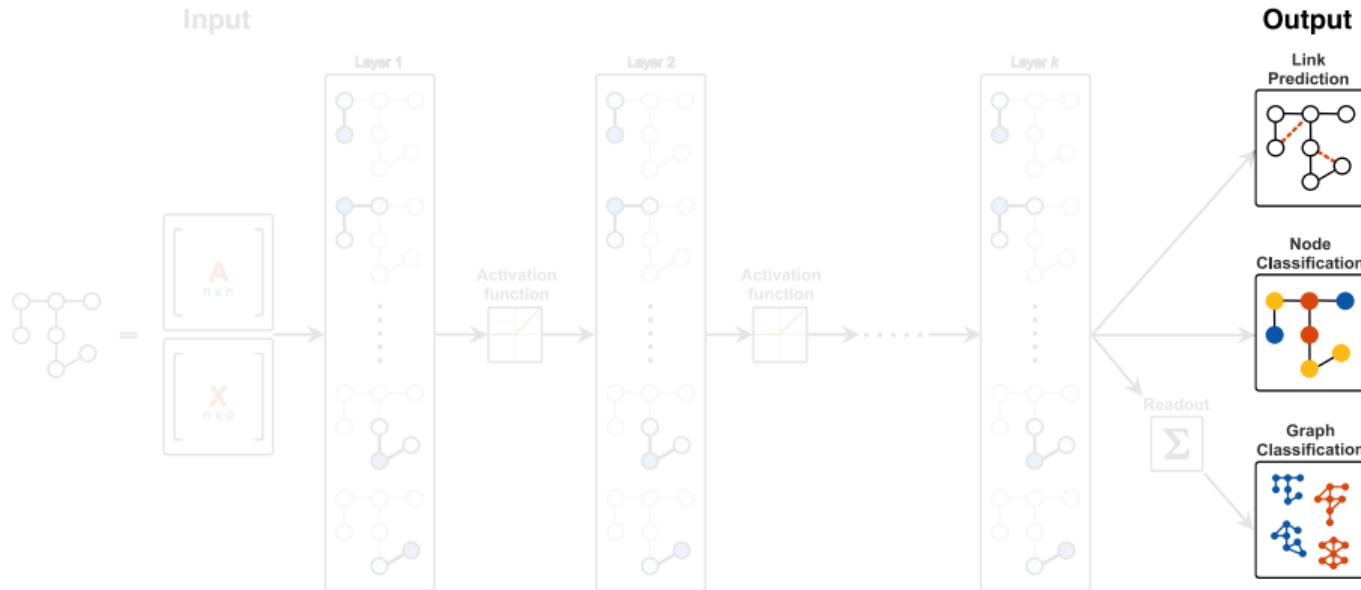
Current paradigm uses neighborhood aggregation



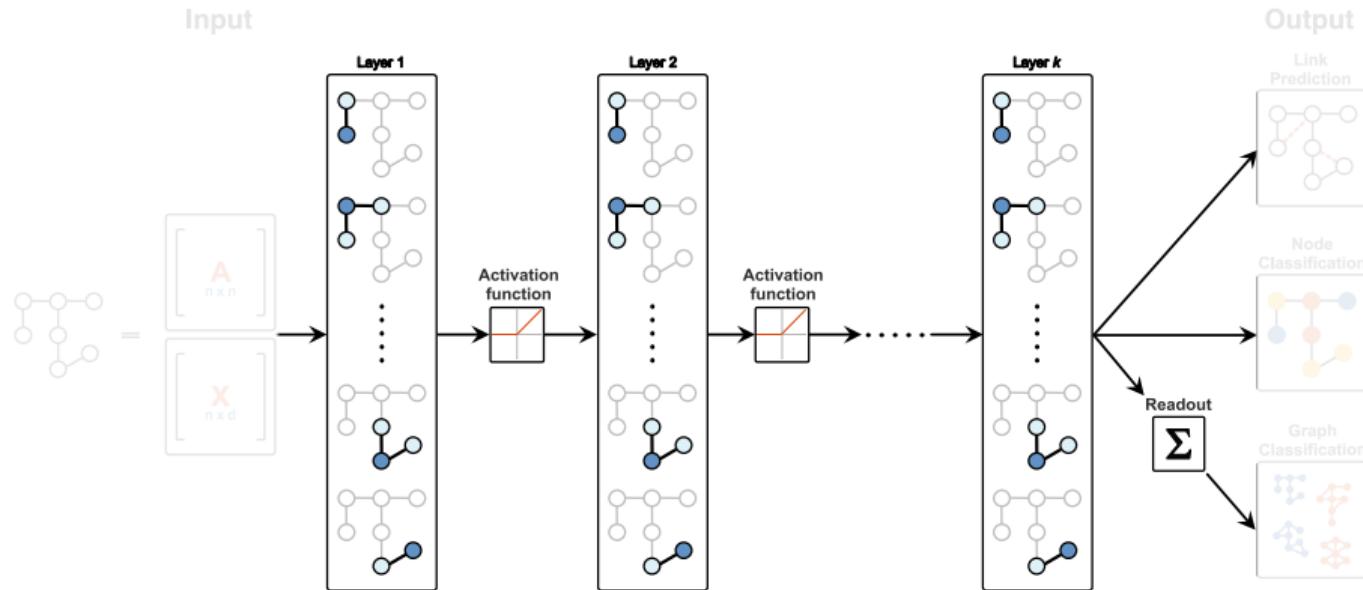
Current paradigm uses neighborhood aggregation



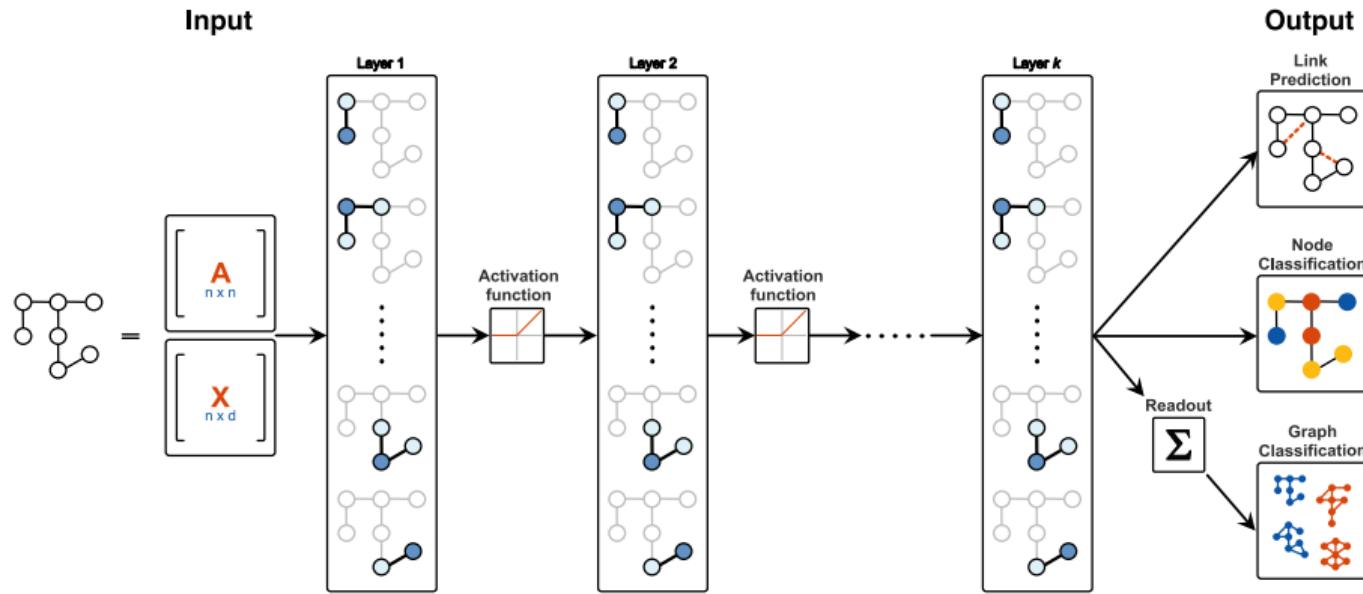
Current paradigm uses neighborhood aggregation



Current paradigm uses neighborhood aggregation



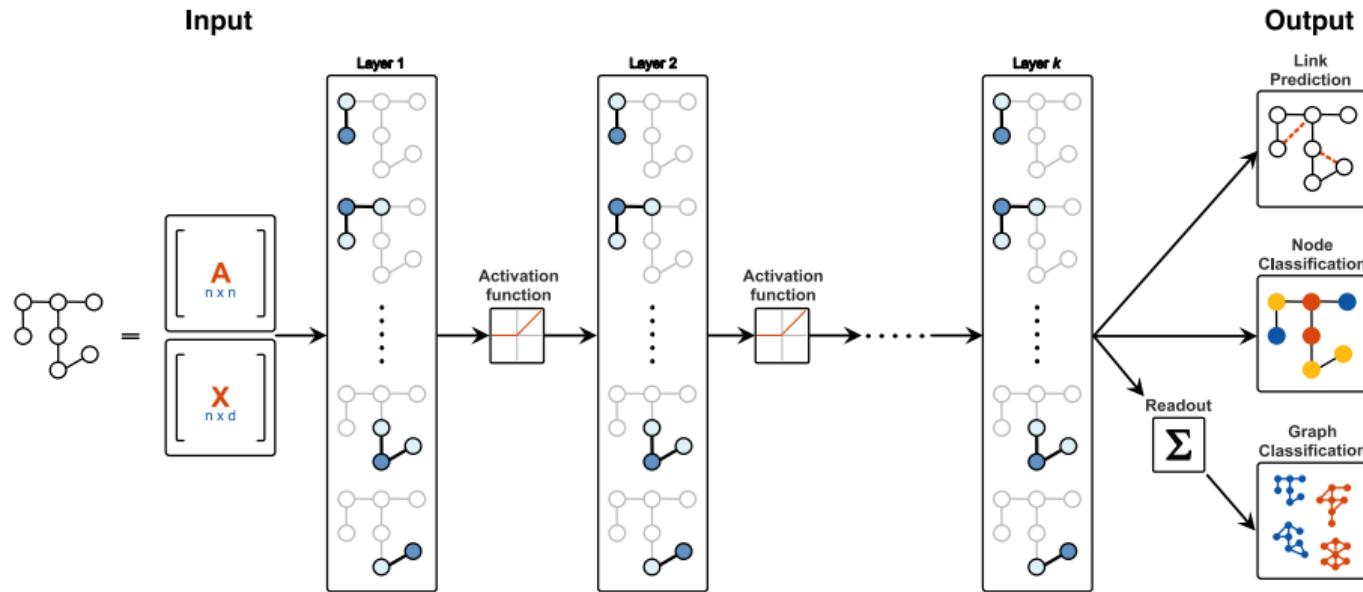
Current paradigm uses neighborhood aggregation



The representation of node u at layer l , h_u^l , is updated by:

$$a_u^l = \text{aggregate}\left(\left\{h_v^{l-1}\right\}_{v \in \mathcal{N}_1(u)}\right)$$

Current paradigm uses neighborhood aggregation

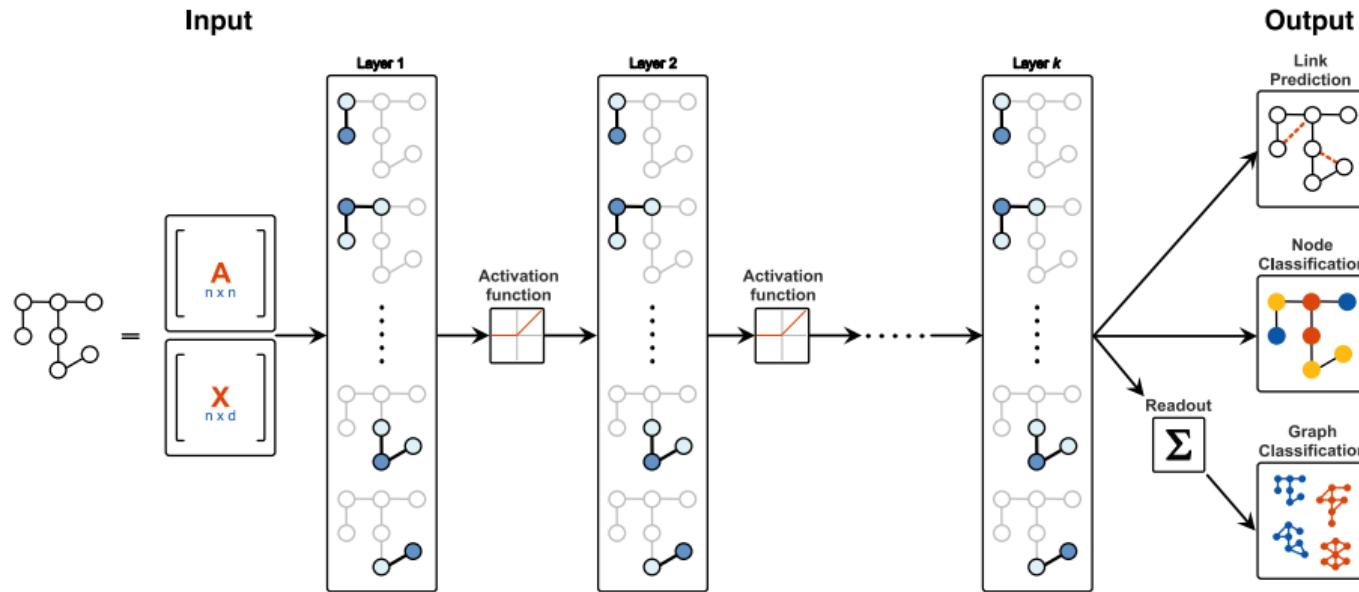


The representation of node u at layer l , h_u^l , is updated by:

$$a_u^l = \text{aggregate}\left(\{h_v^{l-1}\}_{v \in \mathcal{N}_1(u)}\right)$$

$$h_u^l = \text{combine}\left(a_u^l, h_u^{l-1}\right)$$

Current paradigm uses neighborhood aggregation



The representation of node u at layer l , h_u^l , is updated by:

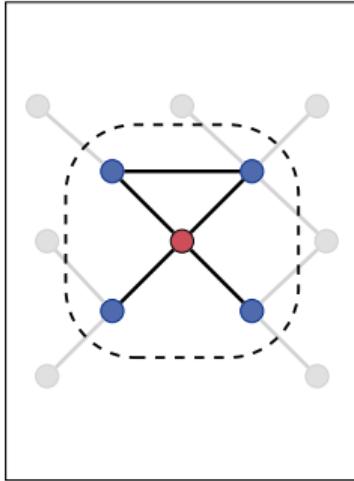
$$a_u^l = \text{aggregate}\left(\left\{h_v^{l-1}\right\}_{v \in \mathcal{N}_1(u)}\right)$$

$$h_u^l = \text{combine}\left(a_u^l, h_u^{l-1}\right)$$

$$h_G = \text{readout}\left(\left\{h_v^k\right\}_{v \in G}\right)$$

GNNs incorporate a local view of the graph, which can cause issues

Under-reaching

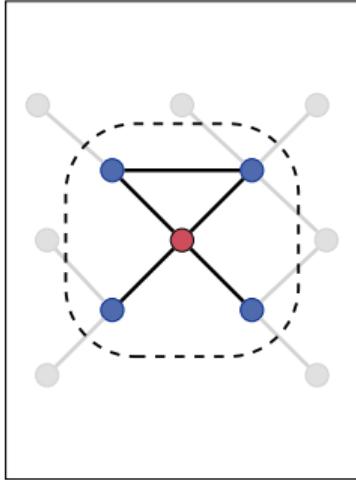


a node cannot see some nodes

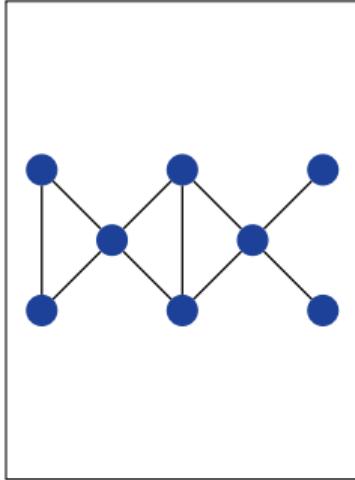
(Barceló et al., 2020)

GNNs incorporate a local view of the graph, which can cause issues

Under-reaching



Over-smoothing

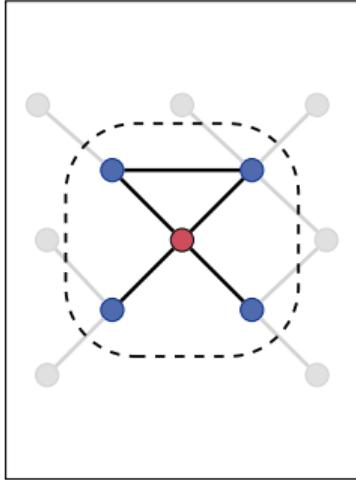


all nodes start to look *the same*

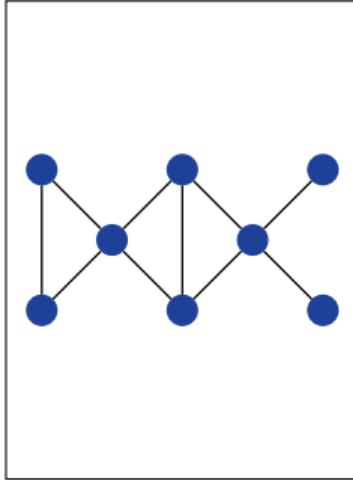
(Oono and Suzuki, 2020)

GNNs incorporate a local view of the graph, which can cause issues

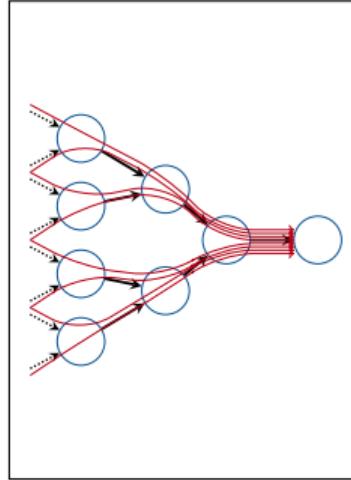
Under-reaching



Over-smoothing



Over-squashing

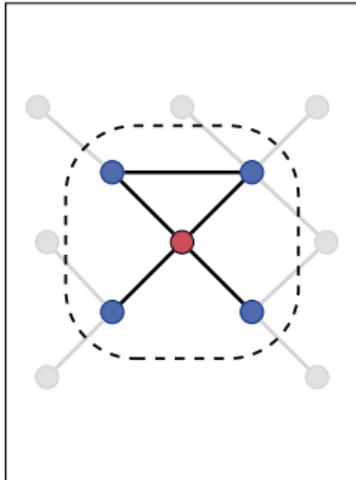


information gets “lost” in transit

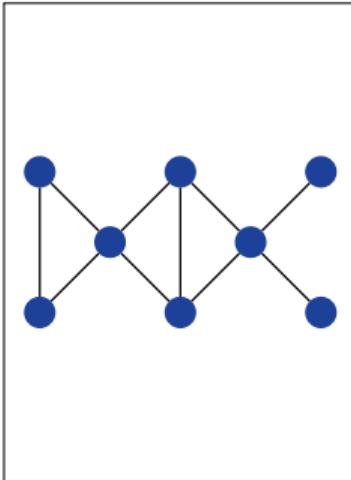
(Topping et al., 2022)

GNNs incorporate a local view of the graph, which can cause issues

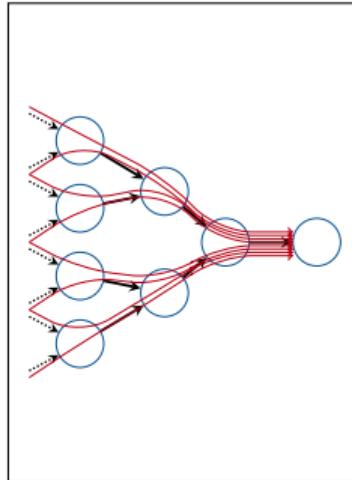
Under-reaching



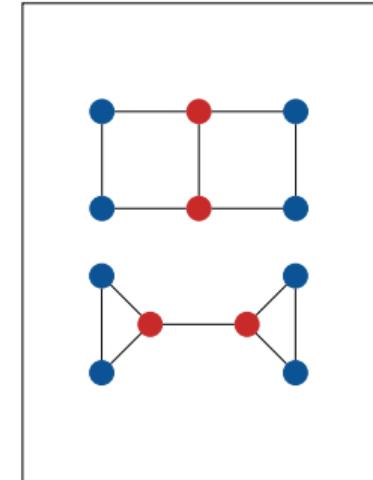
Over-smoothing



Over-squashing



Limited expressivity



at most as expressive as 1-WL

(Morris et al., 2019)

**Could we overcome these issues by incorporating
global information?**

Integrating global information

Ways of incorporating global information into graph learning

- * Virtual nodes
- * Unique node identifiers
- * Adding global information e.g. as node features
- * Transformers

Today's focus: integrating global information via transformers

Why use transformers?

- ⌘ Transformers have had success across a variety of domains
 - ⌘ Natural language processing (OpenAI, 2023)
 - ⌘ Computer vision (Dosovitskiy et al., 2020)
 - ⌘ Speech (Gulati et al., 2020)
 - ⌘ Biological sequence modeling (Rives et al., 2021)

Today's focus: integrating global information via transformers

Why use transformers?

- * Transformers have had success across a variety of domains
- * A single self-attention layer addresses the issues of message passing GNNs
 - * Models global information by design
 - * A transformer is a fully-connected graph
 - * Attention enables message passing between all nodes

Today's focus: integrating global information via transformers

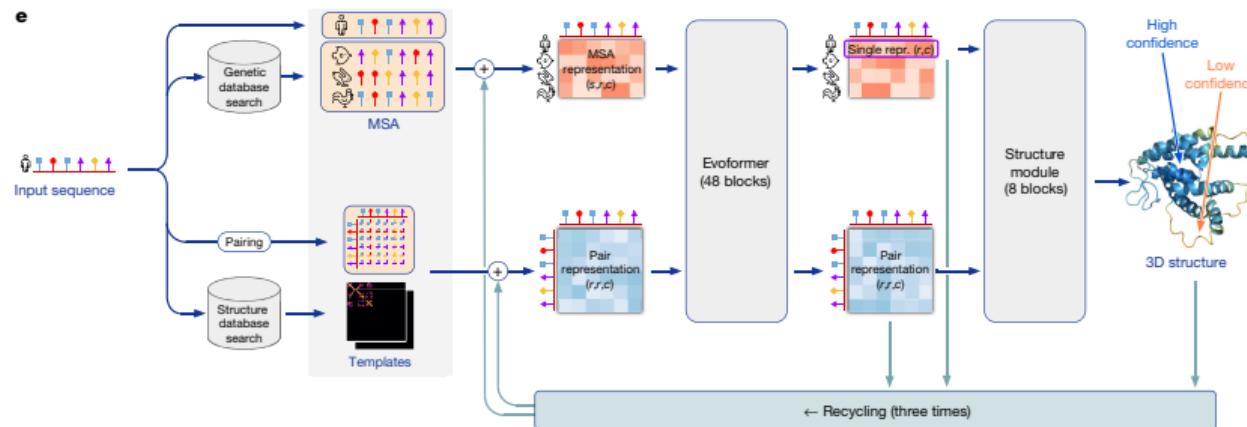
Why use transformers?

- * Transformers have had success across a variety of domains
- * A single self-attention layer addresses the issues of message passing GNNs
- * End-to-end trainable
 - * Doesn't require hand-crafted features

Today's focus: integrating global information via transformers

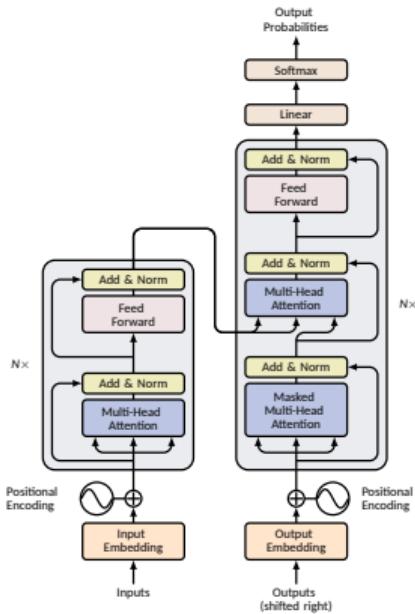
Why use transformers?

- * Transformers have had success across a variety of domains
- * A single self-attention layer addresses the issues of message passing GNNs
- * End-to-end trainable
- * Seen real world examples where global information is informative

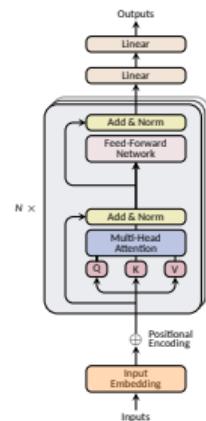


Types of transformer architectures

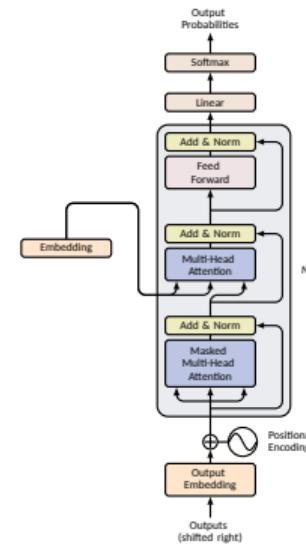
Encoder-Decoder



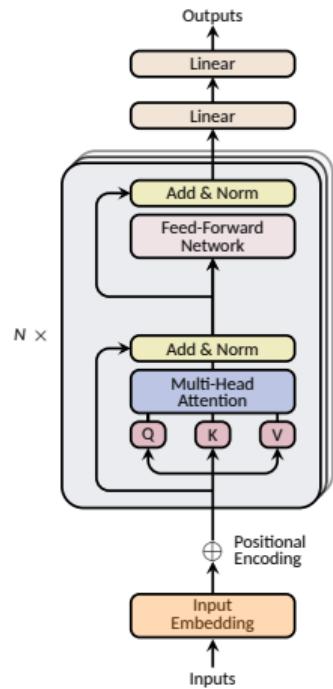
Encoder-only



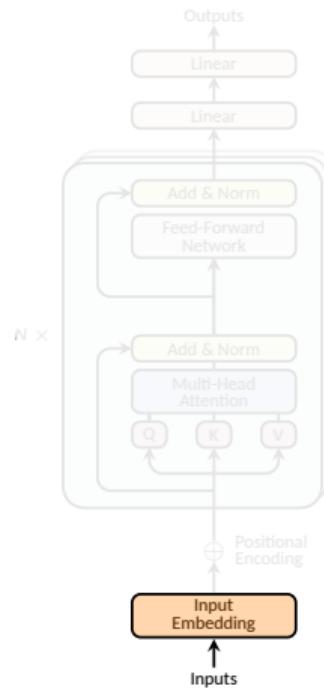
Decoder-only



Components of a transformer



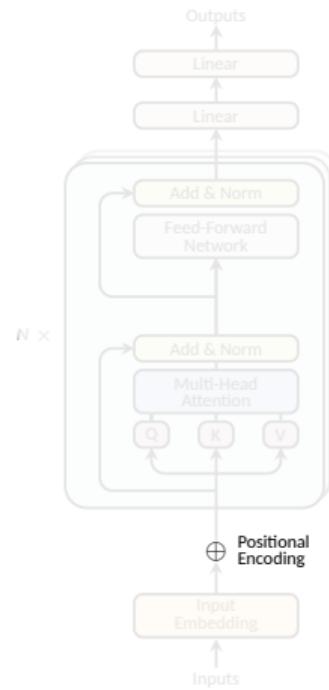
Components of a transformer



Architecture of a transformer

* Input tokenization

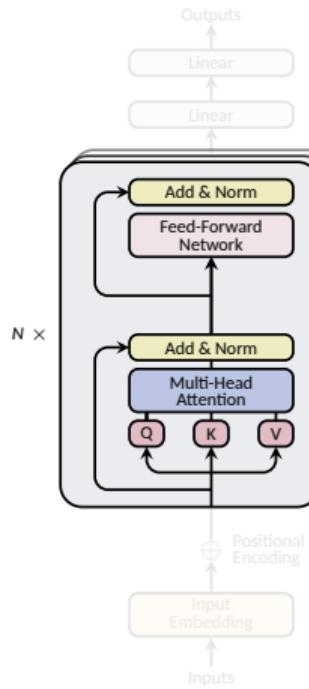
Components of a transformer



Architecture of a transformer

- * Input tokenization
- * Positional encoding

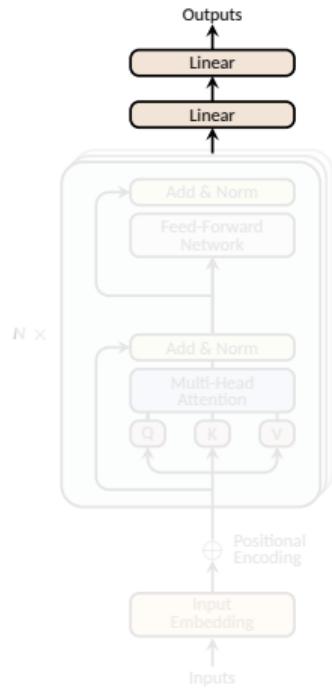
Components of a transformer



Architecture of a transformer

- ⌘ Input tokenization
- ⌘ Positional encoding
- ⌘ Transformer layer with **self-attention**

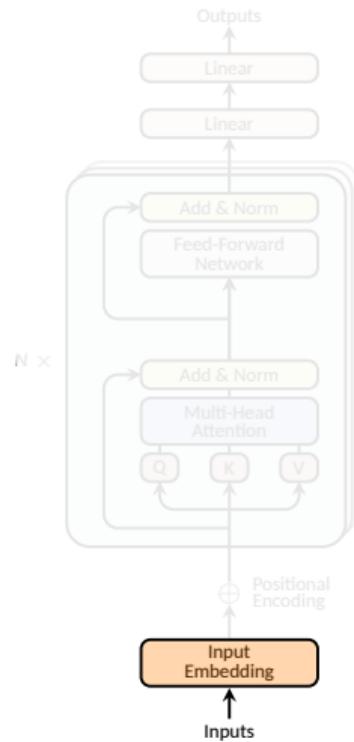
Components of a transformer



Architecture of a transformer

- ⌘ Input tokenization
- ⌘ Positional encoding
- ⌘ Transformer layer with **self-attention**

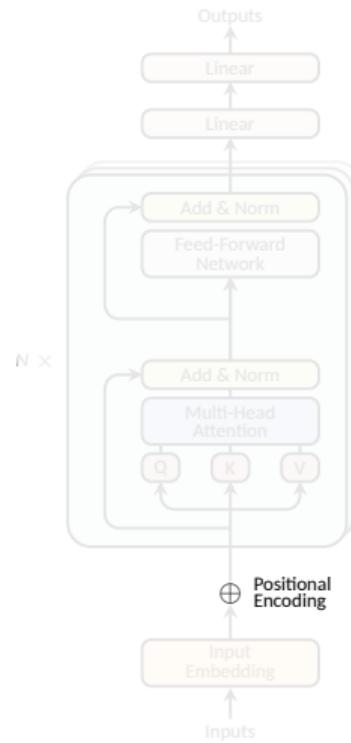
Components of a transformer: Tokenization



Tokenization & embedding

- ⌘ Uses a dictionary to convert a word to a token
- ⌘ “The dog sat by the door” → [5 129 87 902 5 398]
- ⌘ Then uses an embedding matrix to convert from tokens to a d -dim vector

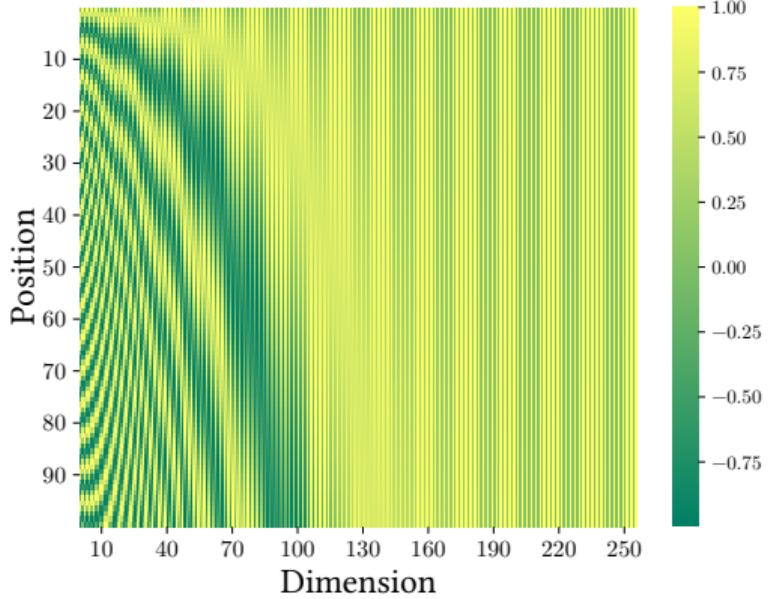
Components of a transformer: Positional encoding



Positional encoding

- * Captures a word's position in a text
- * Is added or concatenated to the input

Absolute positional encoding

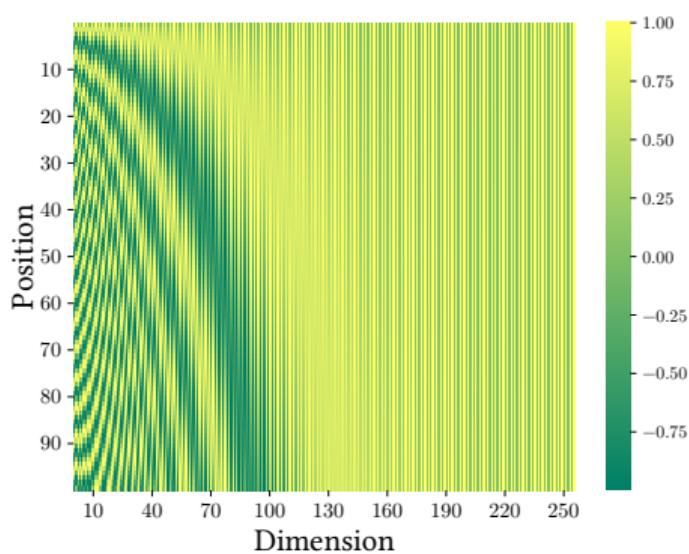


$$\text{PE}_{p,2i} = \sin\left(\frac{p}{10000^{\frac{2i}{d}}}\right)$$

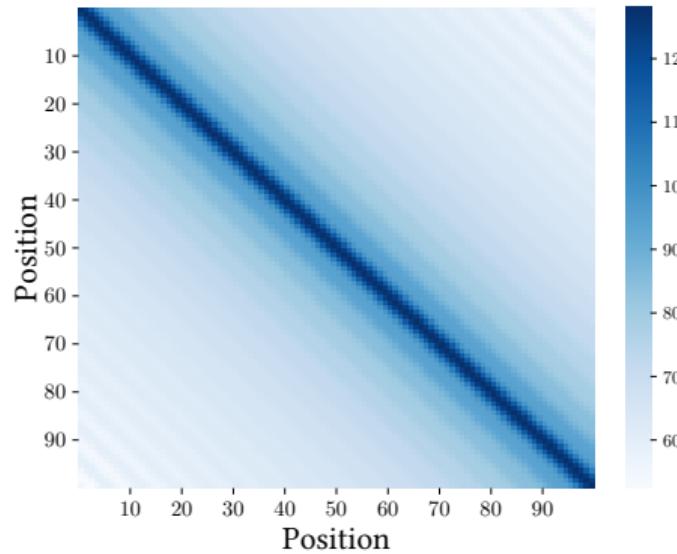
$$\text{PE}_{p,2i+1} = \cos\left(\frac{p}{10000^{\frac{2i}{d}}}\right)$$

- p : index position of the word
- i : dimension in PE vector
- d : dimension of the model

Absolute positional encoding

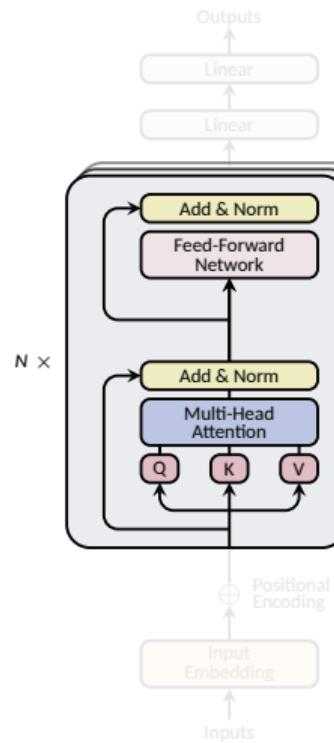


(a) APE



(b) Attention matrix

Components of a transformer: Transformer layer with self-attention



Transformer layer with **self-attention**

- ⌘ Input is projected to **Q**, **K**, and **V** matrices
- ⌘ Passes through (multi-head) attention
- ⌘ Residual connections + normalization
- ⌘ Passes through MLPs

What is self-attention: The components

Q: Query matrix

	d_1	d_2	d_3	d_4
w_1				
w_2				
w_3				
w_4				
w_5				

K: Key matrix

	w_1	w_2	w_3	w_4	w_5
d_1					
d_2					
d_3					
d_4					

V: Value matrix

	d_1	d_2	d_3	d_4
w_1				
w_2				
w_3				
w_4				
w_5				

$$\mathbf{Q} = \mathbf{XW}_Q \in \mathbb{R}^{n \times d}$$

$$\mathbf{K} = \mathbf{XW}_K \in \mathbb{R}^{n \times d}$$

$$\mathbf{V} = \mathbf{XW}_V \in \mathbb{R}^{n \times d}$$

where $|\# \text{ words}| = n$, $\mathbf{X} \in \mathbb{R}^{n \times d}$, and $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$

What is self-attention: The components

Q: Query matrix

	d_1	d_2	d_3	d_4
w_1				
w_2				
w_3				
w_4				
w_5				

K: Key matrix

	w_1	w_2	w_3	w_4	w_5
d_1					
d_2					
d_3					
d_4					

V: Value matrix

	d_1	d_2	d_3	d_4
w_1				
w_2				
w_3				
w_4				
w_5				

Interpretation of the different matrices:

To update an embedding x_i^l at layer l :

- * q_j : the embedding of a word when it receives information from the other nodes
- * k_j : the embedding of a word when it is sending information to update other nodes
- * v_j : the embedding actually used in the update of each word

What is self-attention: The intuition

$$\text{Attn}(\mathbf{X}) := \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

Attention matrix

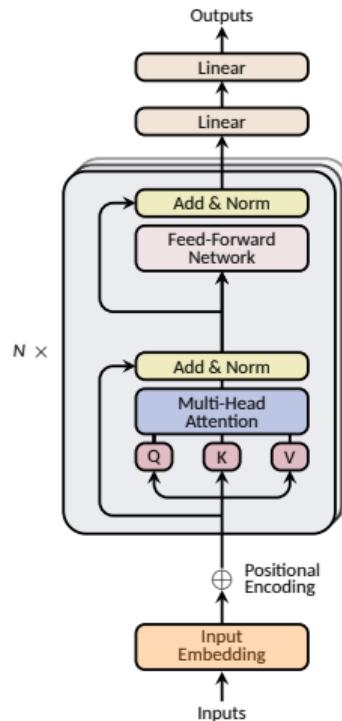
	w_1	w_2	w_3	w_4	w_5
w_1	Light Blue	Light Blue	Light Blue	Dark Blue	Light Blue
w_2	Light Blue	Dark Blue	Light Blue	Light Blue	Light Blue
w_3	Dark Blue	Light Blue	Light Blue	Light Blue	Light Blue
w_4	Dark Blue	Light Blue	Light Blue	Light Blue	Light Blue
w_5	Light Blue	Dark Blue	Light Blue	Light Blue	Light Blue

Value matrix

	d_1	d_2	d_3	d_4
w_1	Light Green	Light Green	Light Green	Light Green
w_2	Light Green	Light Green	Light Green	Light Green
w_3	Light Green	Light Green	Light Green	Light Green
w_4	Light Green	Light Green	Light Green	Light Green
w_5	Light Green	Light Green	Light Green	Light Green

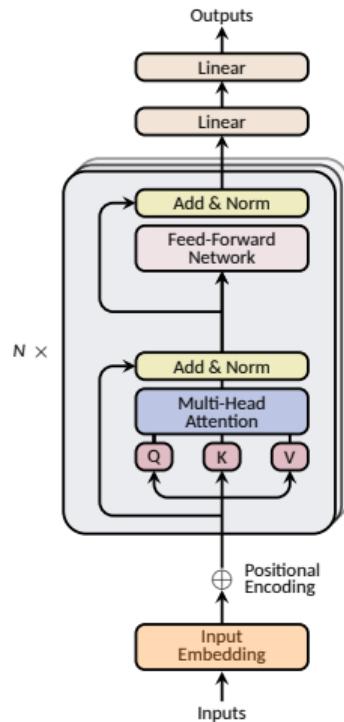
Adapting the transformer to graphs

How can we adapt the transformer to graph data?



Any ideas?

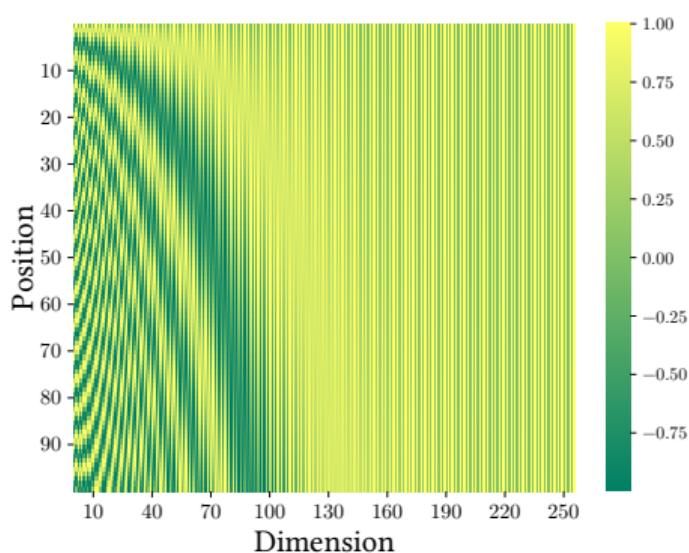
How can we adapt the transformer to graph data?



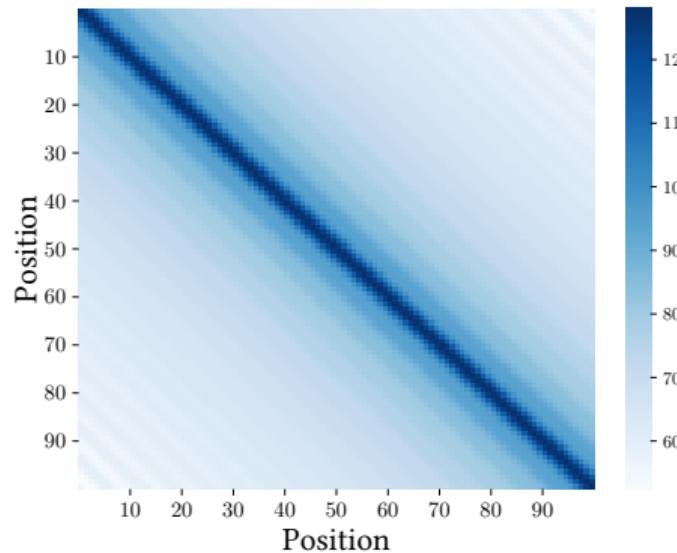
Adapting the transformer for graphs

- ✿ Input treats the graph as a bag of nodes
 - ✿ Encodes structure via positional encoding
- ⇒ Early work did this (Dwivedi and Bresson, 2021)

Absolute positional encoding



(a) APE



(b) Attention matrix

Positional encodings on graphs: LPE

Laplacian positional encoding

- Positioned as graph analog of the sine/cosine PEs
- \mathbf{U} is sorted by increasing order of the corresponding eigenvalue.
- Take the first $1, \dots, k$ eigenvectors of \mathbf{U}

Positional encodings on graphs: LPE

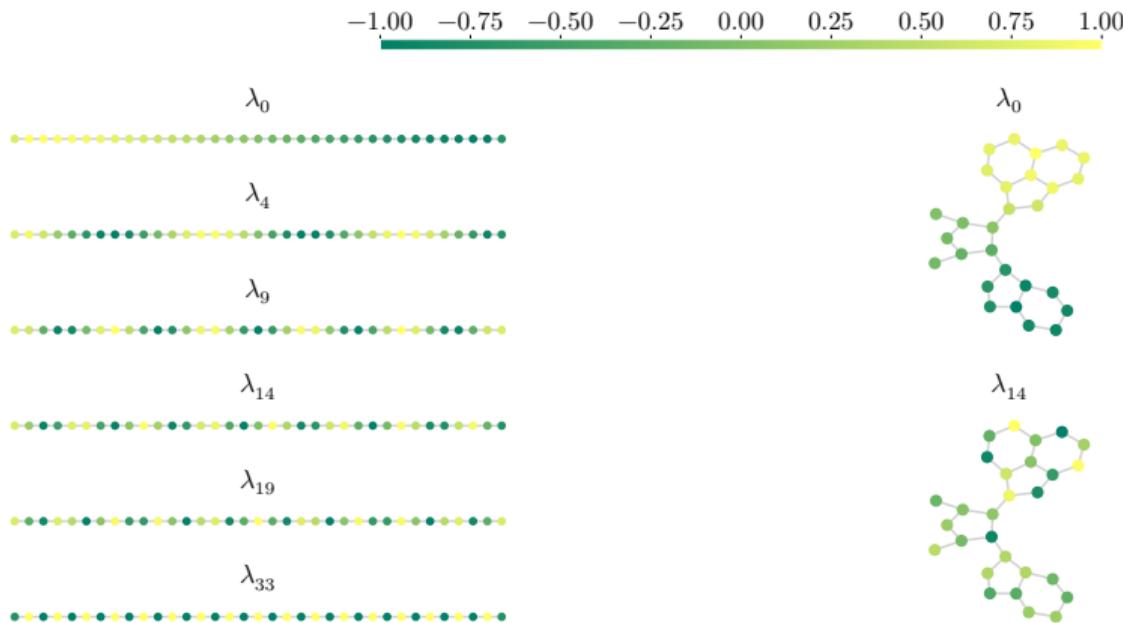
Laplacian positional encoding

- Positioned as graph analog of the sine/cosine PEs
- \mathbf{U} is sorted by increasing order of the corresponding eigenvalue.
- Take the first $1, \dots, k$ eigenvectors of \mathbf{U}
- Limitation: only determined up to sign!

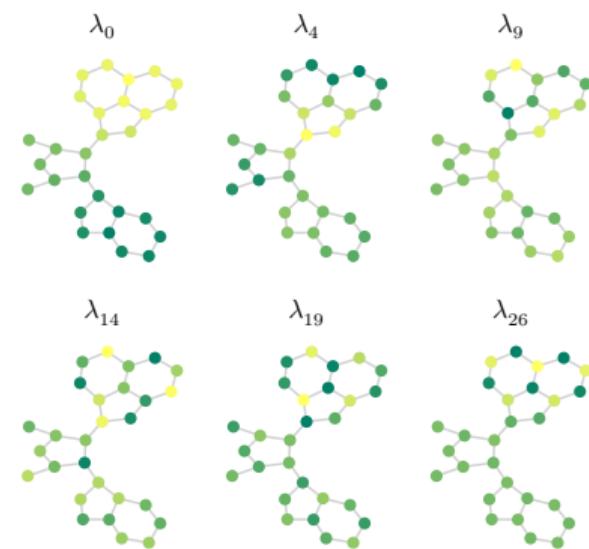
$$\Delta = I_n - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} = \mathbf{U}^T \Lambda \mathbf{U},$$

$$\text{PE}_{\text{LPE}} = \mathbf{U}_k = [\mathbf{U}_1, \dots, \mathbf{U}_k]$$

Laplacian positional encoding (LPE)



(a) LPE of line graph



(b) LPE of a molecule

Positional encodings on graphs: RWPE

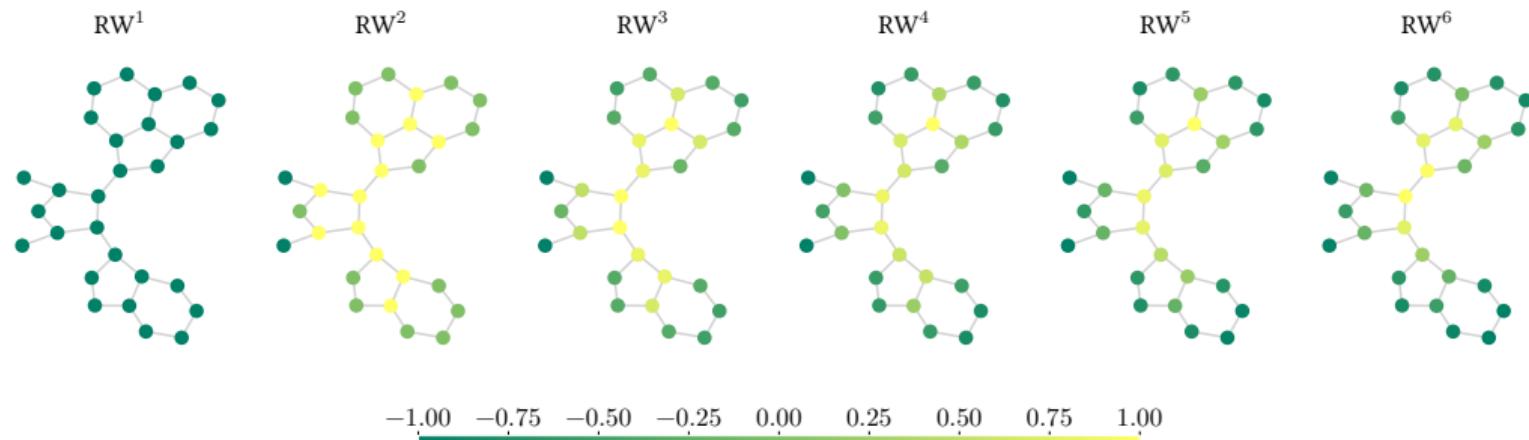
Random walk positional encoding (RWPE)

- Formed by taking powers of the random walk matrix, then taking the diagonal
- Captures some structural information in the graph

$$\text{PE}_{\text{RWPE}} = \begin{bmatrix} \text{RW}_{1,1}^1 & \dots & \text{RW}_{1,n}^k \\ \vdots & \ddots & \vdots \\ \text{RW}_{n,1}^1 & \dots & \text{RW}_{n,n}^k \end{bmatrix}$$

where $\text{RW} = \mathbf{D}^{-1} \mathbf{A}$.

Random walk positional encoding (RWPE)



Early work on graph transformers showed lackluster results

Early publications tried:

- ⌘ Different positional embeddings (LPE, RWPE, Weisfeler-Lehman, etc.) (Dwivedi and Bresson, 2021; Zhang et al., 2020)
- ⌘ Using a GNN as preprocessing before the transformer layer (Rong et al., 2020)
- ⌘ Limiting self-attention calculation to neighbors (Dwivedi and Bresson, 2021)

Do Transformers Really Perform Bad for Graph Representation?

Chengnan Ying¹, Tiansi Cai², Shengjie Lan²,
Shuxin Zhou¹, Mingkang Shi¹, Tie-Yan Liu¹
¹Dalton University of Technology, ²Peking University
Peking University Microsoft Research Asia
yingchengnan@gmail.com, tiansi.cai@princeton.edu, lanxj@tsinghua.edu.cn
(one) guanx, mika, tyliu@microsoft.com, alicebell@tsinghua.edu.cn

Abstract

The Transformer architecture has become a dominant choice in many domains, such as natural language processing and speech. It has also shown its promising performance on popular leaderboards of graph-level prediction compared to mainstream GNN variants. Therefore, it remains a mystery how Transformers could possibly outperform GNNs for graph representation learning. This paper reveals the mystery by presenting Graphmer, which is built upon the standard Transformer architecture, and could attain excellent results on a broad range of graph representation learning tasks, especially on the recent OGBl Large-Scale Challenge. Our key insight is that the Transformer architecture is well suited for graph representation because it can directly model the structural information of a graph into the model. To this end, we propose several simple yet effective structural encoding methods to help Graphmer better model graphs. Experimental results show that Graphmer can outperform state-of-the-art power of Graphmer and exhibit that with our ways of encoding the structural information of graphs, many popular GNN variants could be covered as the special cases of Graphmer. The code and models of Graphmer will be made publicly available at <https://github.com/Microsoft/Graphmer>.

1 Introduction

The Transformer [2] is well acknowledged as the most powerful neural network in modeling sequential data, such as natural language [10, 45, 2] and speech [18]. Model variants built upon the Transformer architecture have been applied to many other domains, such as machine translation [15, 16, 21], image captioning [10, 11, 12], and even the domain of natural language [13, 14, 15]. However, to the best of our knowledge, Transformer has still not been the de-facto standard on public graph representation leaderboards [20, 21, 22]. There are many attempts of leveraging Transformer into the graph domain, but the only effective way is replacing the multi-head attention mechanism with a graph-based attention mechanism [23, 24, 25, 26, 27, 28, 29, 30]. Therefore, it is still an open question whether Transformer architecture is suitable to model graphs and how to make it work in graph representation learning.

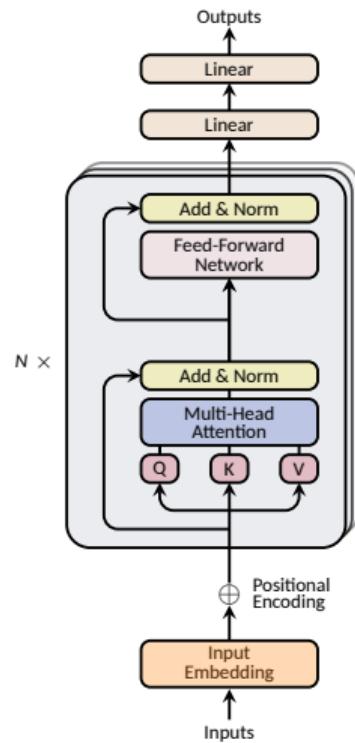
In this paper, we give an affirmative answer by developing Graphmer, which is directly built upon the Transformer architecture and achieves state-of-the-art results on a broad range of graph-level prediction tasks, including the very recent Open Graph Benchmark Large-Scale Challenge (OGBL-LSC) [31] and several popular leaderboards (e.g., OGBl [32], Benchmark-GCN [33]). The Transformer is originally designed for sequence modeling. To utilize its power in graphs, we believe

¹Intern at MSR,
²Corresponding authors.

35th Conference on Neural Information Processing Systems (NeurIPS 2021).

Can we do better?

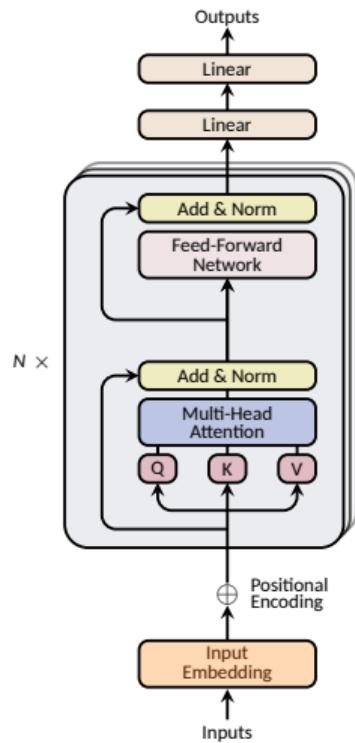
How can we adapt the transformer to work well on graph data?



Adapting the transformer for graphs

- * Input treats the graph as a bag of nodes
- * Encodes structure via positional encoding

How can we adapt the transformer to work well on graph data?



Adapting the transformer for graphs

- * Input treats the graph as a bag of nodes
- * Encodes structure via positional encoding

What if we *explicitly* incorporate structure into attention?

Combining global and local information using transformers

Dexiong Chen, * Leslie O'Bray* and Karsten Borgwardt. Structure-aware transformer for graph representation learning.
In *International Conference on Machine Learning (ICML)*, 2022.

* equal contribution

Self-attention in the context of graphs

Attention matrix

	n_1	n_2	n_3	n_4	n_5
n_1					
n_2					
n_3					
n_4					
n_5					

Value matrix

	d_1	d_2	d_3	d_4
n_1				
n_2				
n_3				
n_4				
n_5				

* $G = (V, E, \mathbf{X})$

* $|V| = n$

* $\mathbf{X} \in \mathbb{R}^{n \times d}$

* $\mathbf{Q} = \mathbf{XW}_Q \in \mathbb{R}^{n \times d}$

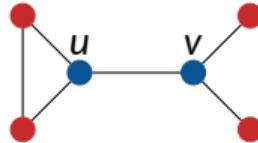
* $\mathbf{K} = \mathbf{XW}_K \in \mathbb{R}^{n \times d}$

* $\mathbf{V} = \mathbf{XW}_V \in \mathbb{R}^{n \times d}$

$$\text{Attn}(\mathbf{X}) := \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

Encoding structure into self-attention

Existing attention compares nodes on the basis of their node attributes



Instead compare based on their **structural and attribute similarity**

From attention to structure-aware attention

Reformulate the attention calculation as a **kernel smoother** (Tsai et al., 2019):

$$\text{Attn}(x_v) = \sum_{u \in V} \frac{\kappa_{\text{exp}}(x_v, x_u)}{\sum_{w \in V} \kappa_{\text{exp}}(x_v, x_w)} f(x_u), \quad \forall v \in V,$$

- * $f(u) = \mathbf{W}_V x_u$ is the linear value function
- * κ_{exp} is an exponential kernel on $\mathbb{R}^d \times \mathbb{R}^d$ parameterized by \mathbf{W}_Q and \mathbf{W}_K :

$$\kappa_{\text{exp}}(x, x') := \exp\left(\frac{\langle \mathbf{W}_Q x, \mathbf{W}_K x' \rangle}{\sqrt{d}}\right), \quad (1)$$

- * $\langle \cdot, \cdot \rangle$ is the dot product on \mathbb{R}^d .

From attention to structure-aware attention

Self-attention

$$\text{Attn}(x_v) = \sum_{u \in V} \frac{\kappa_{\text{exp}}(x_v, x_u)}{\sum_{w \in V} \kappa_{\text{exp}}(x_v, x_w)} f(x_u)$$

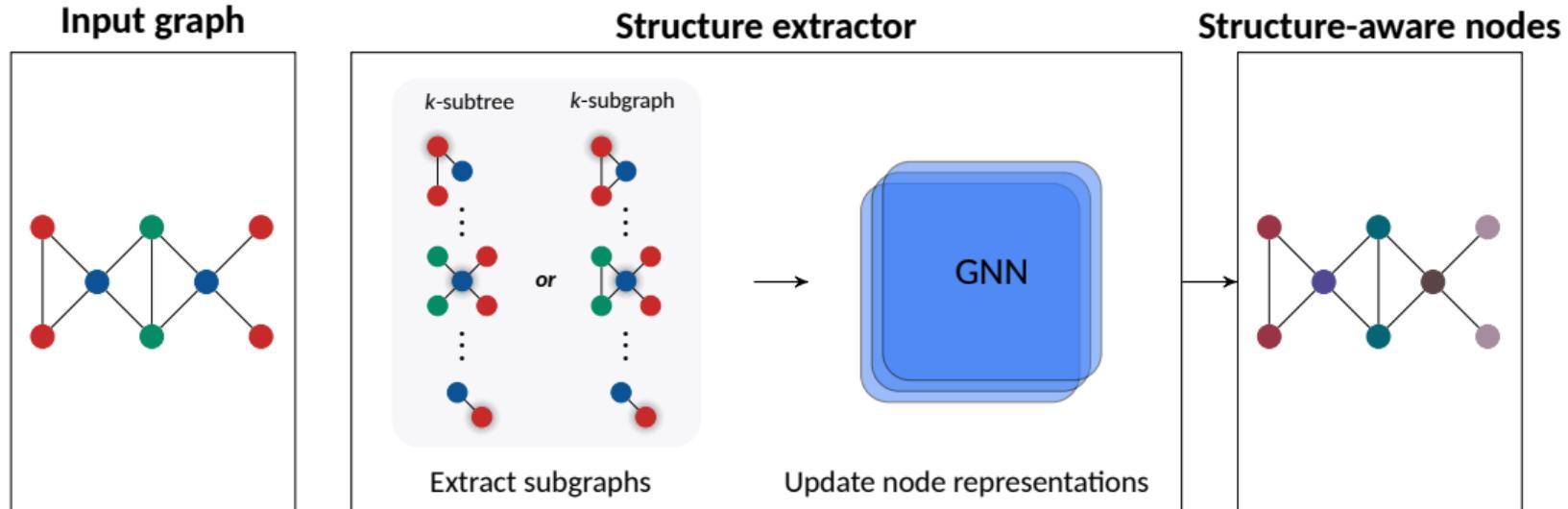
Structure-aware self-attention

$$\text{SA-attn}(v) = \sum_{u \in V} \frac{\kappa_{\text{graph}}(S_G(v), S_G(u))}{\sum_{w \in V} \kappa_{\text{graph}}(S_G(v), S_G(w))} f(x_u)$$

- * κ_{graph} measures similarity between subgraphs
- * $S_G(v)$ is a subgraph rooted at node v

$$\kappa_{\text{graph}}(S_G(v), S_G(u)) = \kappa_{\text{exp}}\left(\underbrace{\varphi(v, G)}_{\text{structure extractor}}, \varphi(u, G)\right)$$

Example structure extractors: k -subtree and k -subgraph

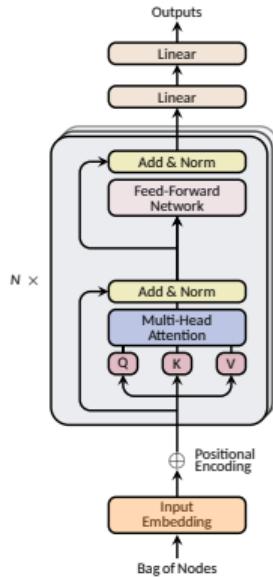


Example structure extractors:

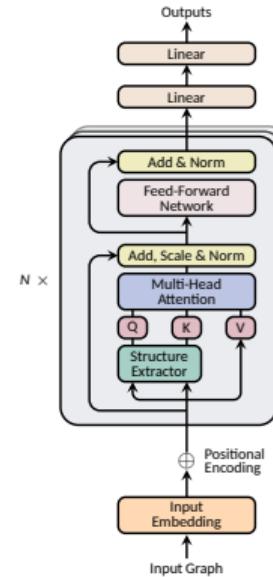
$$\text{k-subtree: } \varphi(u, G) = \text{GNN}_G^{(k)}(u)$$

$$\text{k-subgraph: } \varphi(u, G) = \sum_{v \in \mathcal{N}_k(u)} \text{GNN}_G^{(k)}(v)$$

Structure-Aware Transformer (SAT)

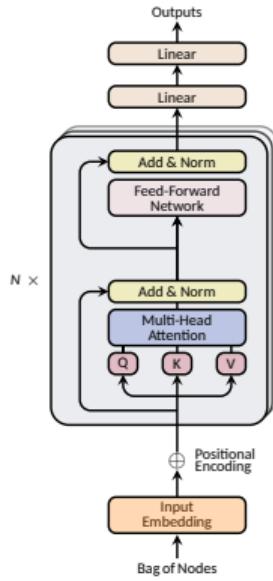


(a) Graph Transformer

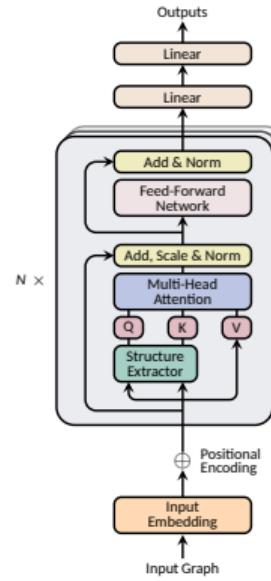


(b) SAT

Structure-Aware Transformer (SAT)



(a) Graph Transformer



(b) SAT

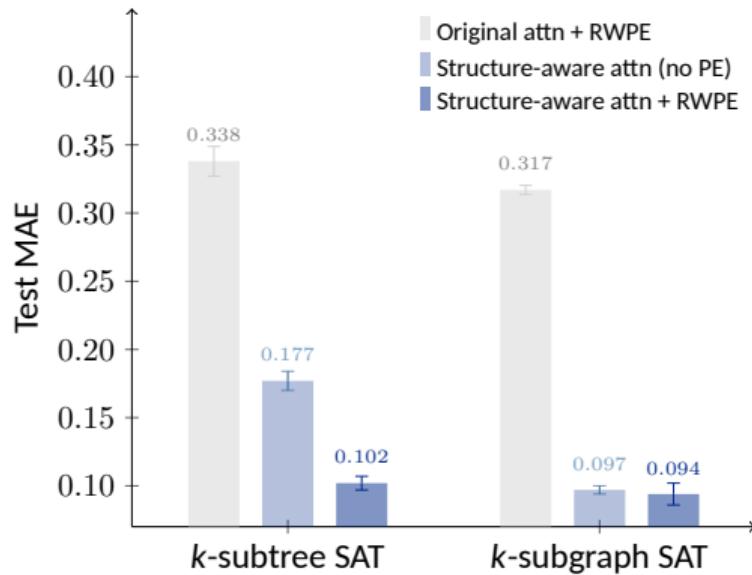
SAT is compatible with linear attention mechanisms, reducing the complexity from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)!$

Experiments: SAT achieves SOTA results on several datasets

	ZINC ↓	CLUSTER ↑	PATTERN ↑		OGBG-PPA ↑	OGBG-CODE2 ↑
# graphs	12,000	12,000	14,000	# graphs	158,100	452,741
Avg. # nodes	23.2	117.2	118.9	Avg. # nodes	243.4	125.2
Avg. # edges	49.8	4,303.9	6,098.9	Avg. # edges	2,266.1	124.2
Metric	MAE	Accuracy	Accuracy	Metric	Accuracy	F1 score
GNNs	GIN	0.387 ± 0.015	64.716 ± 1.553	85.590 ± 0.011	0.6857 ± 0.0061	0.1595 ± 0.0018
	GAT	0.384 ± 0.007	70.587 ± 0.447	78.271 ± 0.186		0.7037 ± 0.0107
	PNA	0.188 ± 0.004	67.077 ± 0.977	86.567 ± 0.075		0.1581 ± 0.0026
GTs	Transformer+RWPE	0.310 ± 0.005	29.622 ± 0.176	86.183 ± 0.019	Transformer	0.6454 ± 0.0033
	SAN	0.139 ± 0.006	76.691 ± 0.650	86.581 ± 0.037	GraphTrans	-
	Graphormer	0.122 ± 0.006	-	-	k-subtree SAT	0.7522 ± 0.0056
	k-subtree SAT	0.102 ± 0.005	77.751 ± 0.121	86.865 ± 0.043		0.1937 ± 0.0028
	k-subgraph SAT	0.094 ± 0.008	77.856 ± 0.104	86.848 ± 0.037		

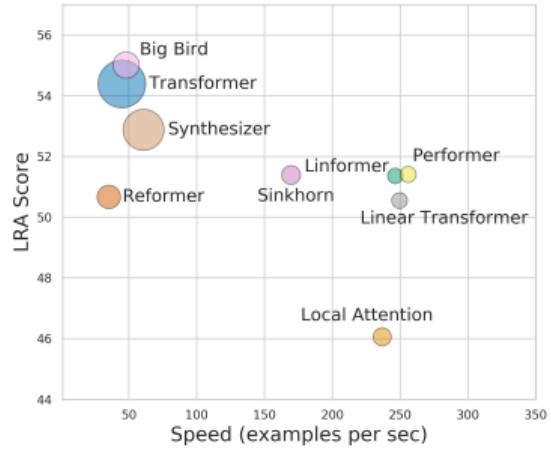
First | Second | Third

Benefit of the structure-aware attention on ZINC



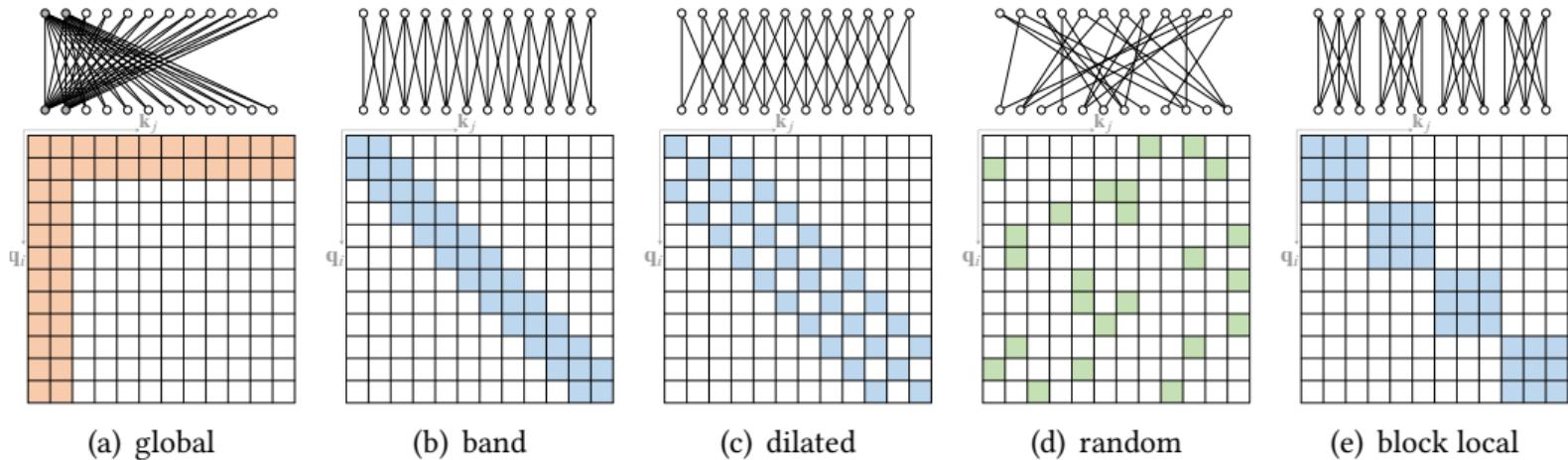
Limitations of transformers and graph transformers

- ⌘ Self-attention is $\mathcal{O}(N^2)$!
- ⌘ A lot of research in developing linear or sparse attention
- ⌘ Tradeoff between performance and speed



(Tay et al., 2021)

Example of different ideas for sparse attention



(Lin et al., 2022)

A lot of exciting recent work in graph transformers

- * Rethinking Graph Transformers with Spectral Attention (Kreuzer et al., 2021)
- * Do Transformers Really Perform Badly for Graph Representation? (Ying et al., 2021)
- * Recipe for a General, Powerful, Scalable Graph Transformer (Rampasek et al., 2022)
- * Pure Transformers are Powerful Graph Learners (Kim et al., 2022)

Summary & take-home

- * GNNs are SOTA but have known issues given that they only model local structural information
- * Transformers offer a way to counteract some of these issues by incorporating global information in the graph
- * Global information alone may not be enough, but the combination of local and global is very promising

References I

- P. Barceló, E. V. Kostylev, M. Monet, J. Pérez, J. Reutter, and J. P. Silva. The logical expressiveness of graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2020.
- V. P. Dwivedi and X. Bresson. A generalization of transformer networks to graphs. *AAAI Workshop on Deep Learning on Graphs: Methods and Applications*, 2021.
- A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. 2020.
- J. Kim, D. T. Nguyen, S. Min, S. Cho, M. Lee, H. Lee, and S. Hong. Pure transformers are powerful graph learners. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=um2BxfgkT2_.
- D. Kreuzer, D. Beaini, W. L. Hamilton, V. Létourneau, and P. Tossou. Rethinking graph transformers with spectral attention. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- T. Lin, Y. Wang, X. Liu, and X. Qiu. A survey of transformers. *AI Open*, 3:111–132, 2022. ISSN 2666-6510. doi: <https://doi.org/10.1016/j.aiopen.2022.10.001>.
- C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe. Weisfeiler and Leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.

References II

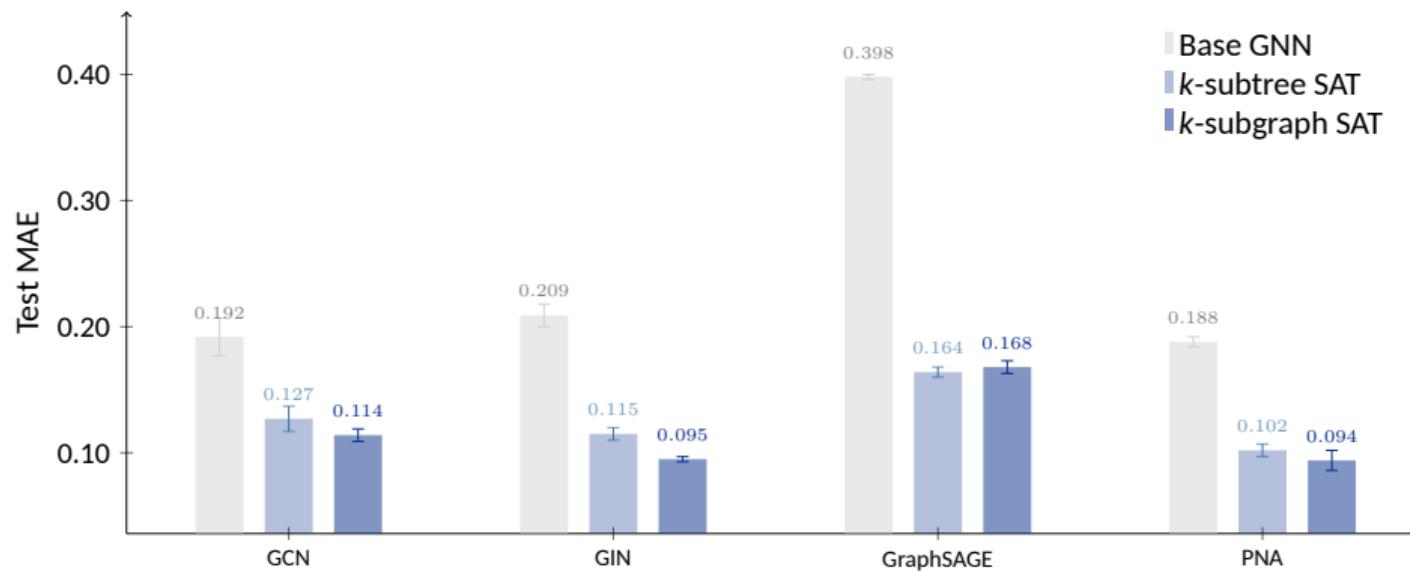
- K. Oono and T. Suzuki. Graph neural networks exponentially lose expressive power for node classification. In *International Conference on Learning Representations (ICLR)*, 2020.
- OpenAI. Gpt-4 technical report, 2023.
- L. Rampasek, M. Galkin, V. P. Dwivedi, A. T. Luu, G. Wolf, and D. Beaini. Recipe for a general, powerful, scalable graph transformer. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=1MMaNf6oxKM>.
- A. Rives, J. Meier, T. Sercu, S. Goyal, Z. Lin, J. Liu, D. Guo, M. Ott, C. L. Zitnick, J. Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), 2021.
- Y. Rong, Y. Bian, T. Xu, W. Xie, Y. Wei, W. Huang, and J. Huang. Self-supervised graph transformer on large-scale molecular data. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Y. Tay, M. Dehghani, S. Abnar, Y. Shen, D. Bahri, P. Pham, J. Rao, L. Yang, S. Ruder, and D. Metzler. Long range arena : A benchmark for efficient transformers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=qVyeW-grC2k>.
- J. Topping, F. D. Giovanni, B. P. Chamberlain, X. Dong, and M. M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature. In *International Conference on Learning Representations*, 2022.
- Y.-H. H. Tsai, S. Bai, M. Yamada, L.-P. Morency, and R. Salakhutdinov. Transformer dissection: A unified understanding of transformer's attention via the lens of kernel. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.

References III

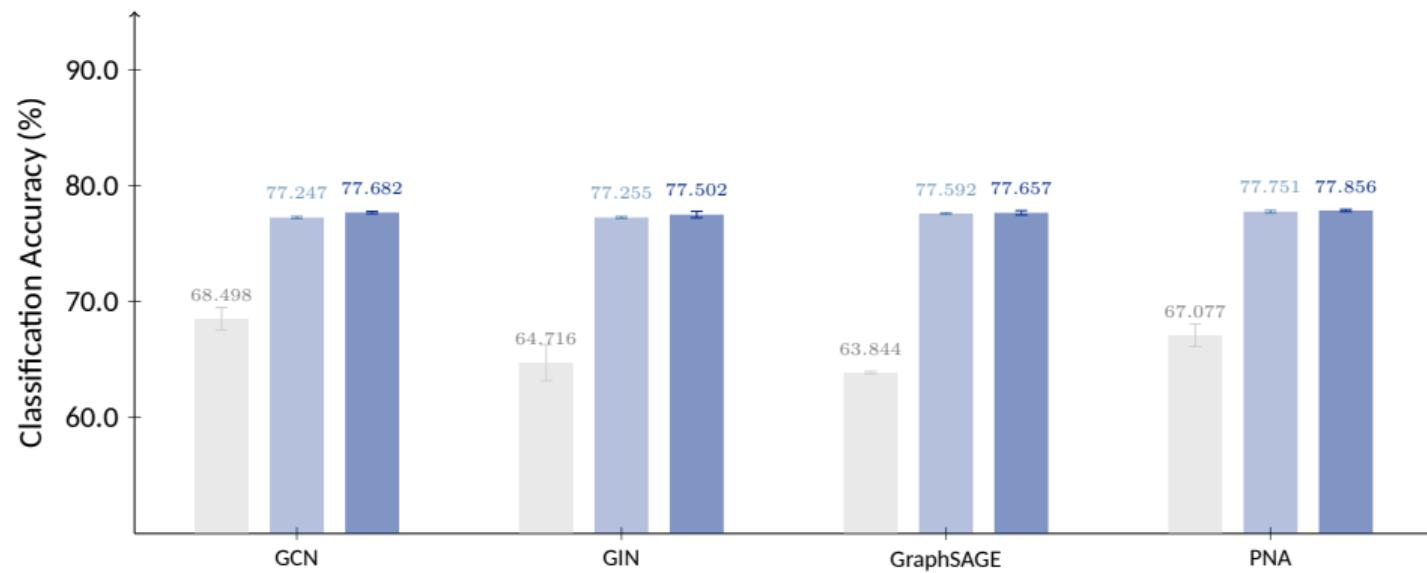
- C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, and T.-Y. Liu. Do transformers really perform badly for graph representation? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- J. Zhang, H. Zhang, C. Xia, and L. Sun. Graph-Bert: Only attention is needed for learning graph representations. 2020.

Backup slides: SAT

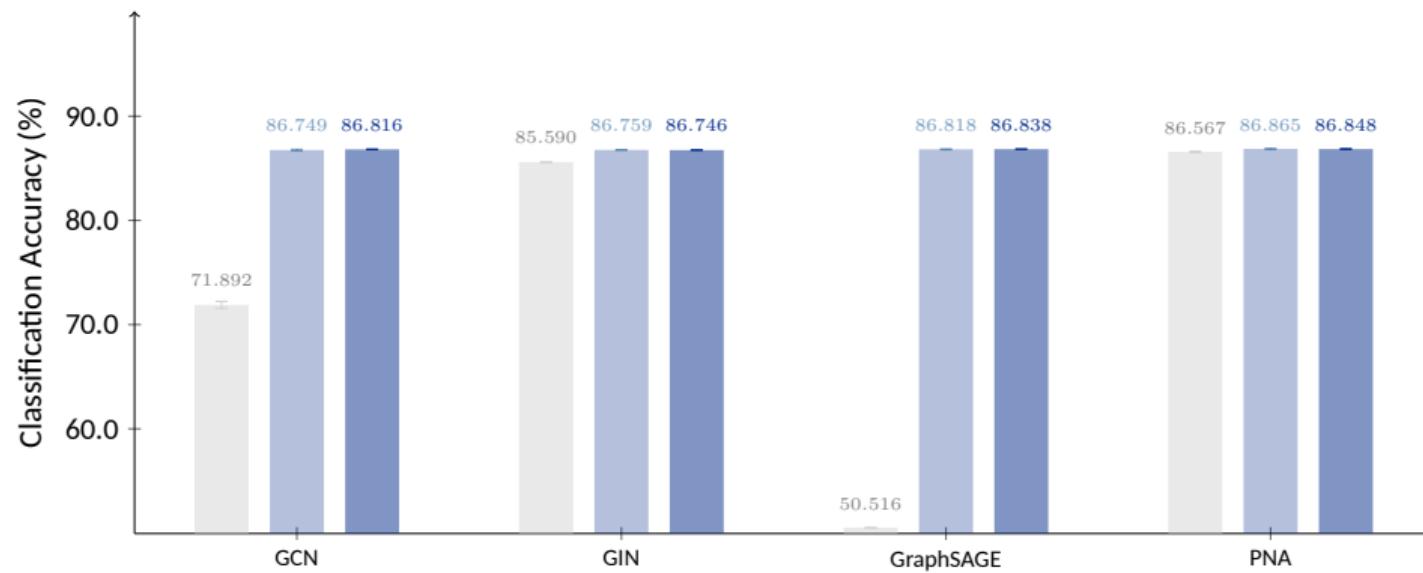
Uplift of integrating global information with local information on ZINC



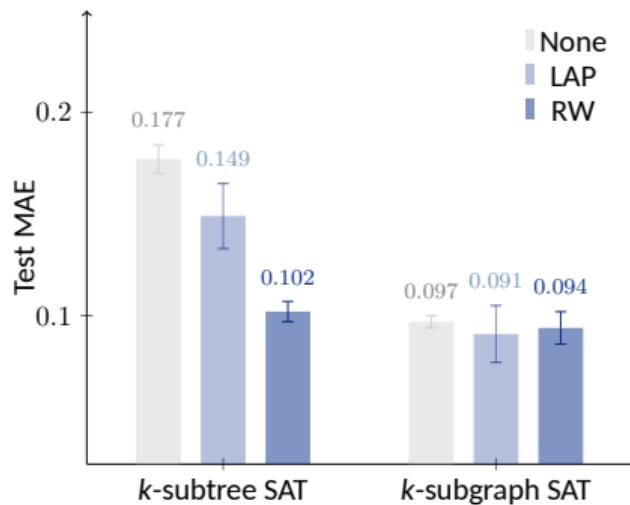
Uplift of integrating global information with local information, CLUSTER



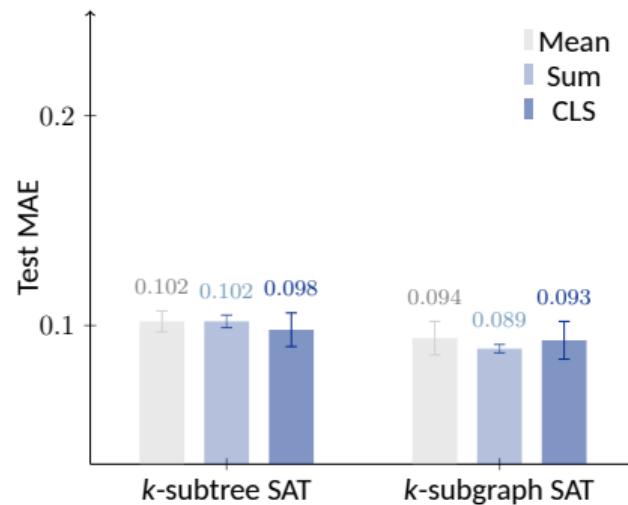
Uplift of integrating global information with local information, PATTERN



Effect of positional encoding and readout on ZINC

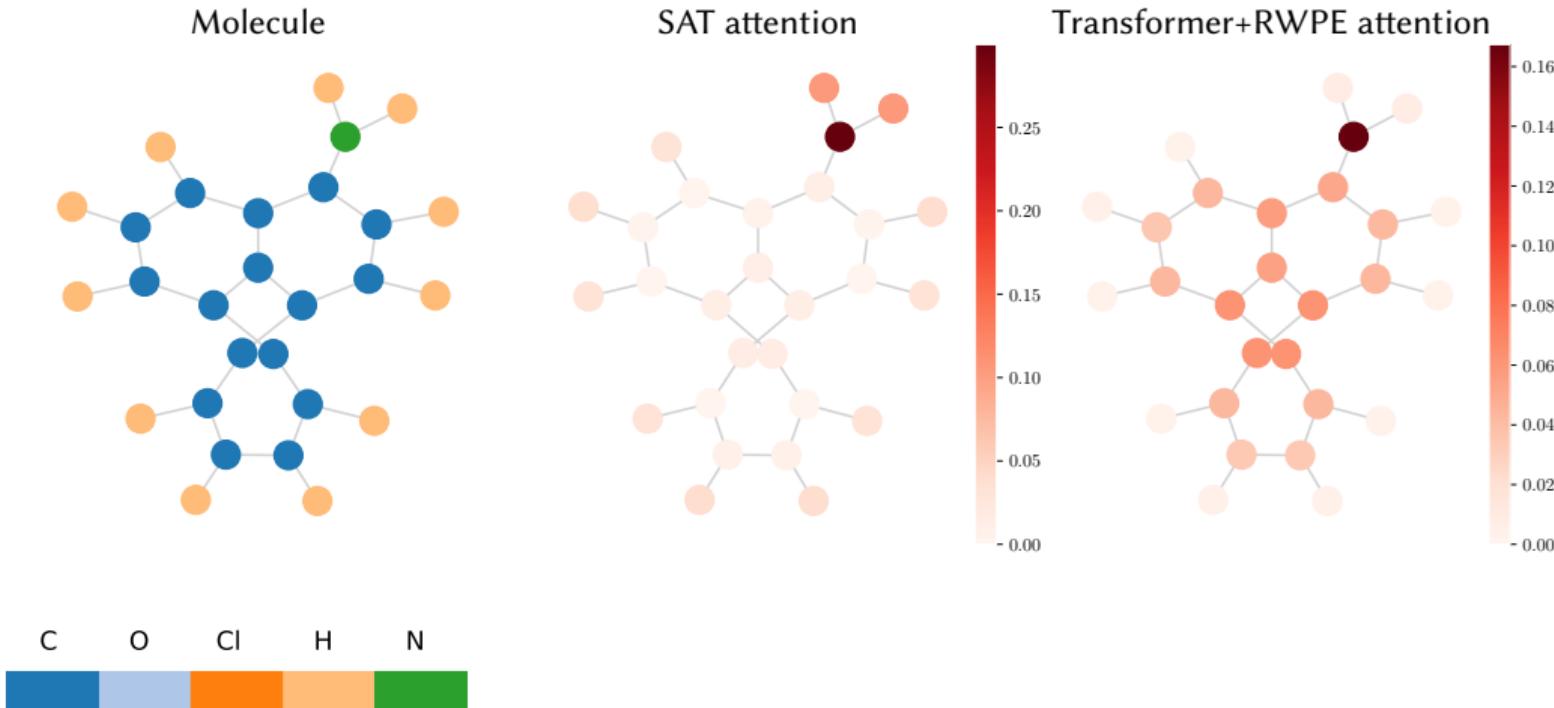


(a) Effect of positional encoding



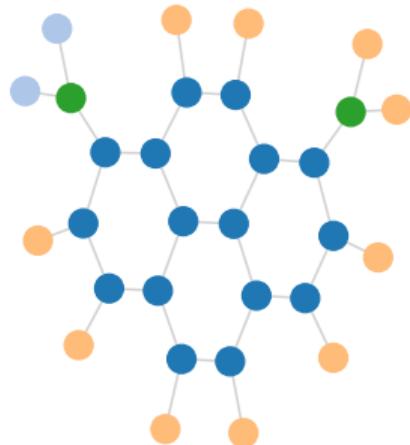
(b) Effect of readout

Structure-aware attention scores on Mutagenicity

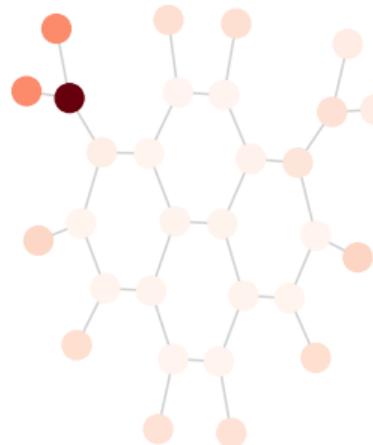


Structure-aware attention scores on Mutagenicity

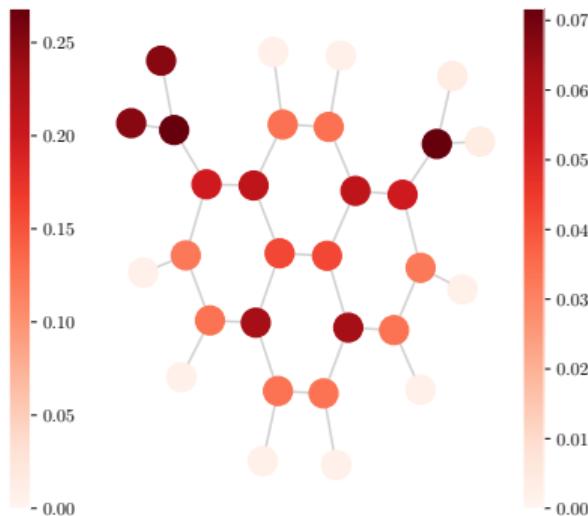
Molecule



SAT attention



Transformer+RWPE attention



Structure-aware attention scores on Mutagenicity

