

SGFormer: Simplifying and Empowering Transformers for Large-Graph Representations

Advances in Neural Information Processing Systems (NeurIPS), 2023

Qitian Wu, Wentao Zhao, Chenxiao Yang, Hengrui Zhang, Fan Nie,
Haitian Jiang, Yatao Bian, Junchi Yan



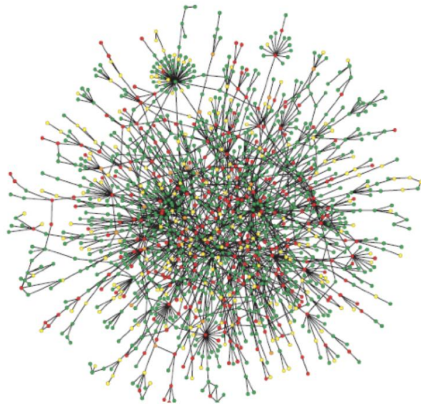
Tencent

- [1] Qitian Wu et al., NodeFormer: A Scalable Graph Structure Learning Transformer for Node Classification, NeurIPS 2022 (spotlight)
- [2] Qitian Wu et al., DIFFFormer: Scalable (Graph) Transformers Induced by Energy Constrained Diffusion, ICLR 2023 (spotlight)
- [3] Qitian Wu et al., SGFormer: Simplifying and Empowering Transformers for Large-Graph Representations, NeurIPS 2023

Data with Explicit Structures

- Real-world data involves inter-dependencies of observed structures

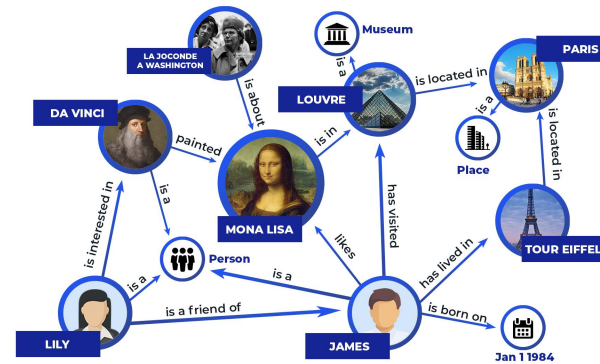
protein interactions



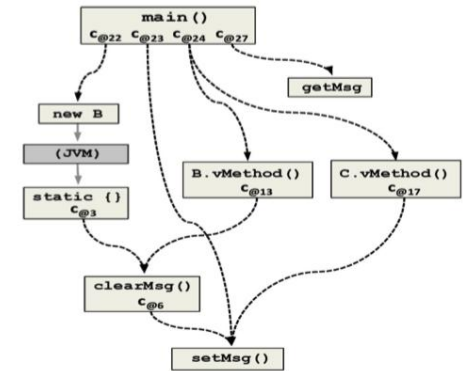
social networks



knowledge graphs



code



- Graph neural networks become a default class of model solutions

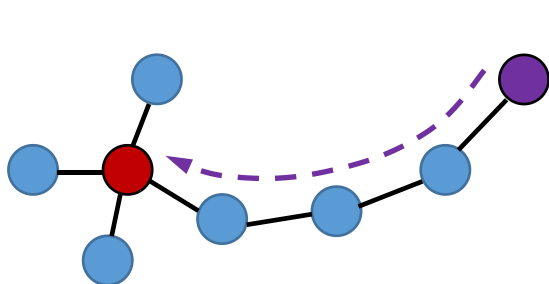
- 1) GNNs can handle variable-sized inputs with geometries
- 2) GNNs are expressive for modeling topological features

Pitfalls of Graph Neural Networks

❑ The designs of mainstream GNNs:

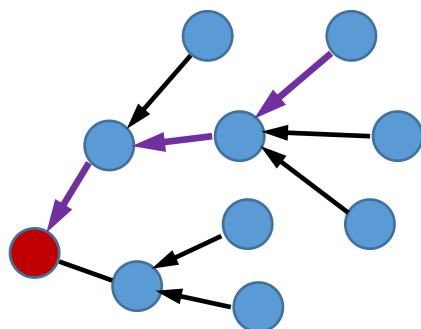
- Locally aggregate neighbored nodes' features in each layer
- Use neighbored nodes' embs for informative representation

❑ Common scenarios GNNs show deficient capability:



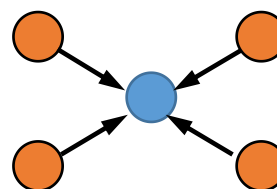
hard to capture long-range dependence
[Dai et al., 2018]

long-range reasoning



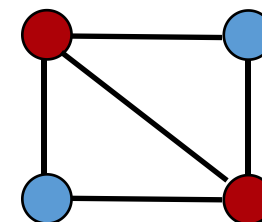
distant signals are overly squashed
[Alon et al., 2021]

over-squashing



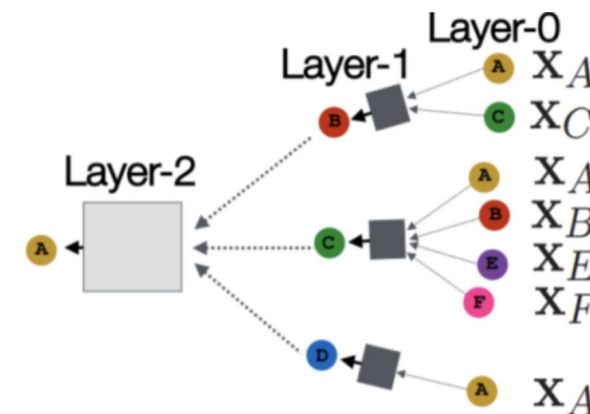
dissimilar linked nodes propagate wrong signals
[Zhu et al., 2020]

heterophily

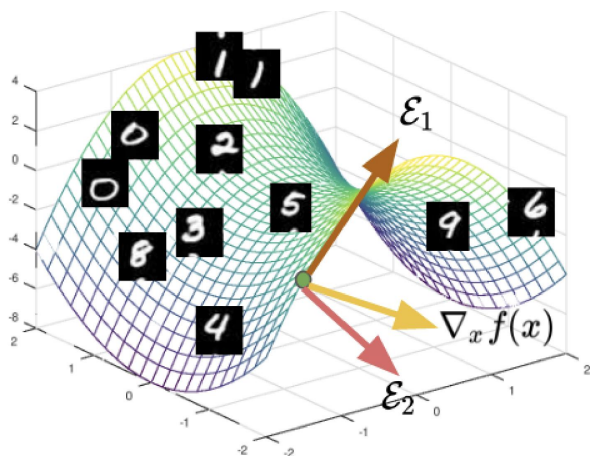


fail to distinguish two similar inputs
[Xu et al., 2019]

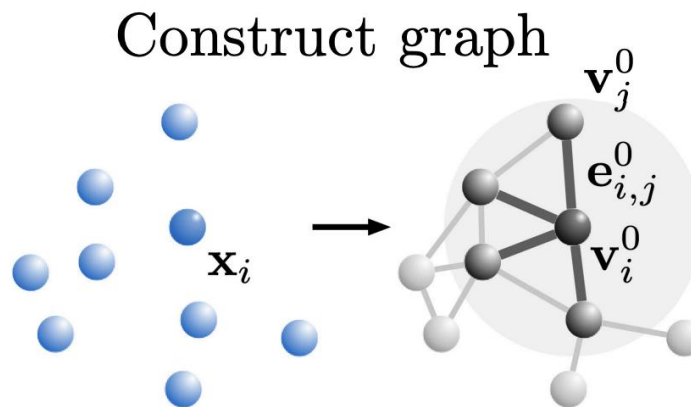
expressivity



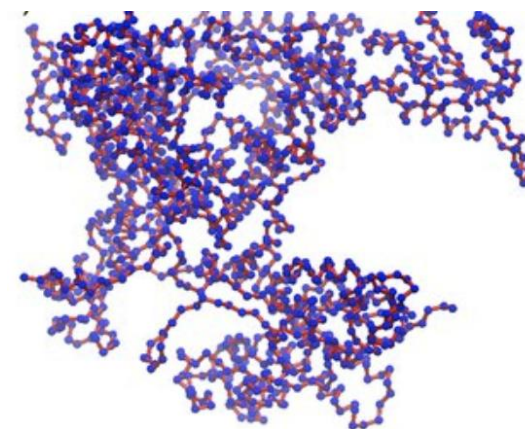
Inter-Dependent Data without Input Graphs



Observed data lies on low-dimensional manifold
[Sebastian et al., 2021]



Physical interactions affect data generation yet are not observed
[Alvaro et al., 2020]



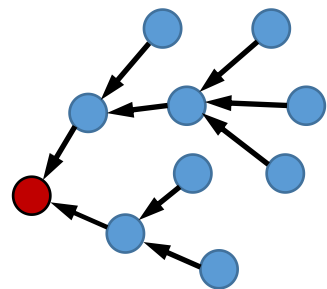
Complex hidden structures beyond observed geometry
[Xu et al., 2020]

□ GNNs require observed graphs as input:

- **Solution:** Pre-define a graph by some rules (e.g., k nearest neighbors)
- **Limitation:** the pre-defined graph is independent of downstream tasks

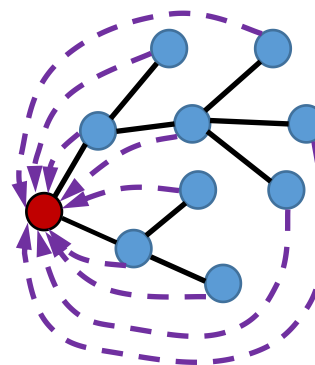
Message Passing Beyond Input Graphs

Graph Neural Networks



local message passing
defined over fixed input
topology

Transformers



all-pair message passing
on layer-specific latent
graphs

Graph Transformers (GT) have shown impressive power on graph learning tasks

Most of existing GT are designed for small-graph-based tasks (e.g., molecules)

Challenges of building GT for large-graph representations:

- The **quadratic** global attentions hinder the **scalability** for large graphs
- The **complicated** model architectures are prone for **over-fitting**

NodeFormer: All-Pair Attention with $O(N)$

□ Kernelized softmax message passing

$$\mathbf{z}_u^{(l+1)} = \sum_{v=1}^N \frac{\exp(\mathbf{q}_u^\top \mathbf{k}_v)}{\sum_{w=1}^N \exp(\mathbf{q}_u^\top \mathbf{k}_w)} \cdot \mathbf{v}_v \quad \text{where } \mathbf{q}_u = W_Q^{(l)} \mathbf{z}_u^{(l)}, \quad \mathbf{k}_u = W_K^{(l)} \mathbf{z}_u^{(l)}, \quad \mathbf{v}_u = W_V^{(l)} \mathbf{z}_u^{(l)}$$

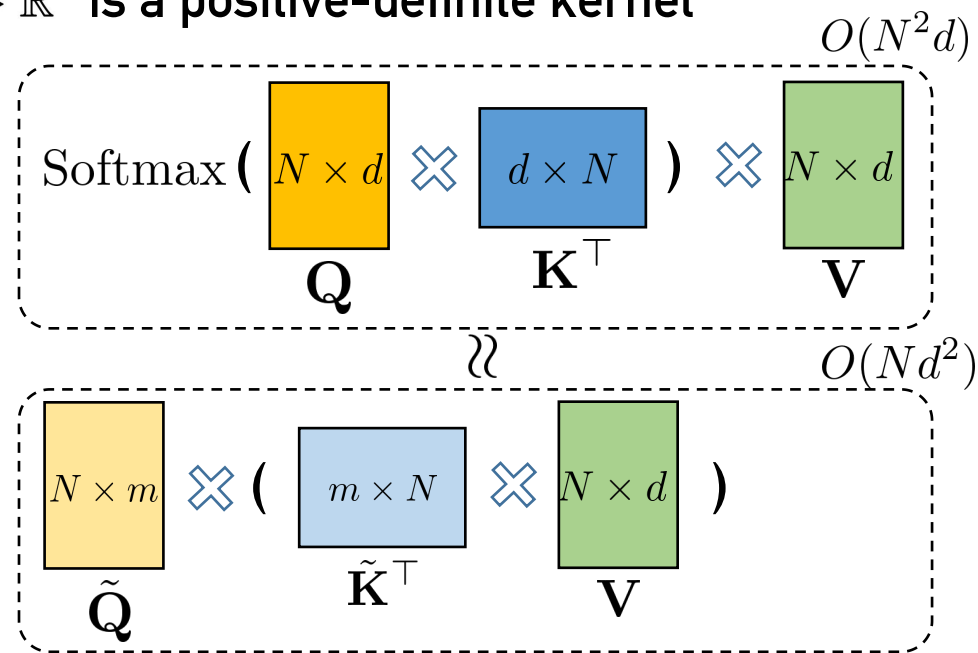
$$\mathbf{z}_u^{(l+1)} = \sum_{v=1}^N \frac{\kappa(\mathbf{q}_u, \mathbf{k}_v)}{\sum_{w=1}^N \kappa(\mathbf{q}_u, \mathbf{k}_w)} \cdot \mathbf{v}_v \quad \text{where } \kappa(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \text{ is a positive-definite kernel}$$

[Mercer's theorem] $\kappa(\mathbf{a}, \mathbf{b}) = \langle \Phi(\mathbf{a}), \Phi(\mathbf{b}) \rangle_{\mathcal{V}} \approx \phi(\mathbf{a})^\top \phi(\mathbf{b})$
 $\phi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is a random feature map

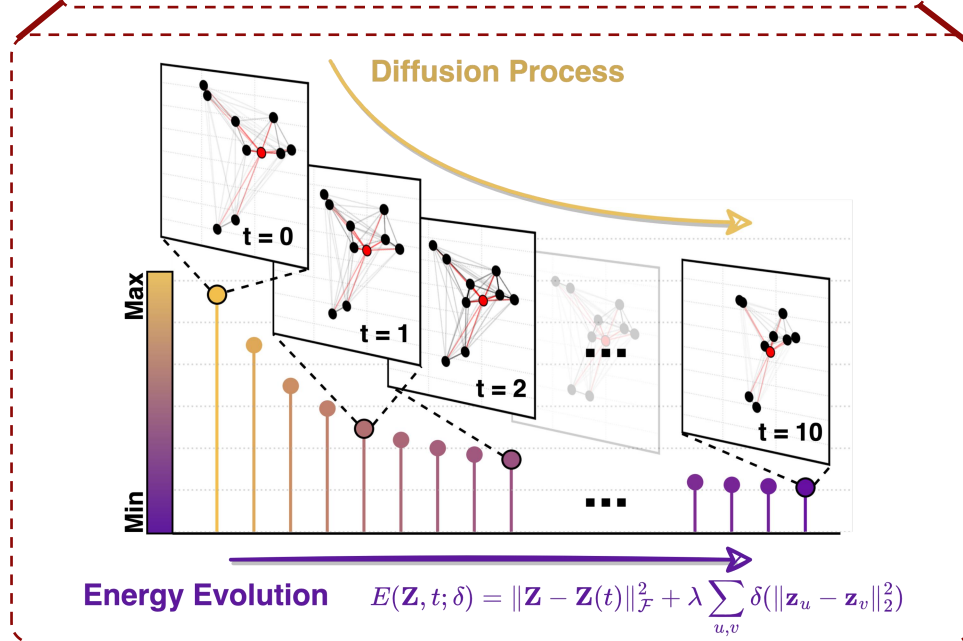
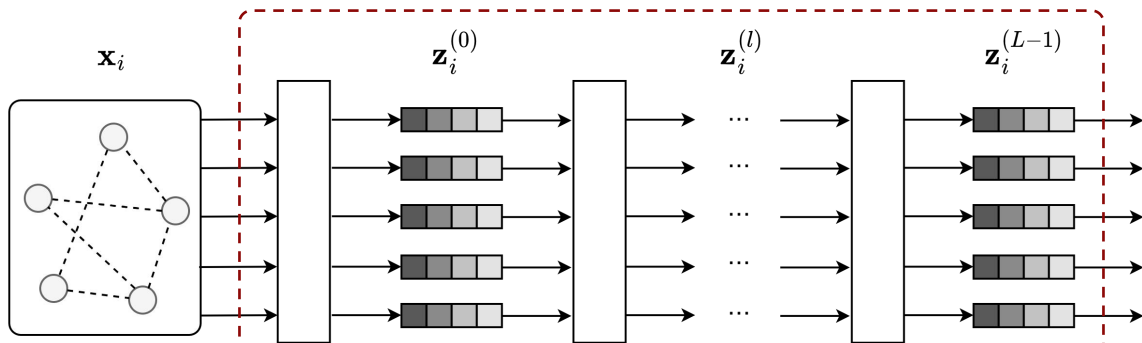
$$\mathbf{z}_u^{(l+1)} = \sum_{v=1}^N \frac{\phi(\mathbf{q}_u)^\top \phi(\mathbf{k}_v)}{\sum_{w=1}^N \phi(\mathbf{q}_u)^\top \phi(\mathbf{k}_w)} \cdot \mathbf{v}_v = \frac{\phi(\mathbf{q}_u)^\top \sum_{v=1}^N \phi(\mathbf{k}_v) \cdot \mathbf{v}_v^\top}{\phi(\mathbf{q}_u)^\top \sum_{w=1}^N \phi(\mathbf{k}_w)}$$

two summation are shared by all nodes (independent of u)
— only compute once

computation complexity $O(N) + N \cdot O(1) = O(N)$



DIFFormer: Transformers by Diffusion



$$\hat{\mathbf{S}}_{ij}^{(k)} = \frac{f(\|\mathbf{z}_i^{(k)} - \mathbf{z}_j^{(k)}\|_2^2)}{\sum_{l=1}^N f(\|\mathbf{z}_i^{(k)} - \mathbf{z}_l^{(k)}\|_2^2)}, \quad 1 \leq i, j \leq N$$

$$\mathbf{z}_i^{(k+1)} = \left(1 - \tau \sum_{j=1}^N \hat{\mathbf{S}}_{ij}^{(k)}\right) \mathbf{z}_i^{(k)} + \tau \sum_{j=1}^N \hat{\mathbf{S}}_{ij}^{(k)} \mathbf{z}_j^{(k)}, \quad 1 \leq i \leq N$$

Global attention inspired by diffusivity function

$$\omega_{ij}^{(k)} = f(\|\tilde{\mathbf{z}}_i^{(k)} - \tilde{\mathbf{z}}_j^{(k)}\|_2^2) = 1 + \left(\frac{\mathbf{z}_i^{(k)}}{\|\mathbf{z}_i^{(k)}\|_2}\right)^\top \left(\frac{\mathbf{z}_j^{(k)}}{\|\mathbf{z}_j^{(k)}\|_2}\right)$$

$$\sum_{j=1}^N \mathbf{S}_{ij}^{(k)} \mathbf{z}_j^{(k)} = \sum_{j=1}^N \frac{1 + (\tilde{\mathbf{z}}_i^{(k)})^\top \tilde{\mathbf{z}}_j^{(k)}}{\sum_{l=1}^N (1 + (\tilde{\mathbf{z}}_i^{(k)})^\top \tilde{\mathbf{z}}_l^{(k)})} \mathbf{z}_j^{(k)}$$

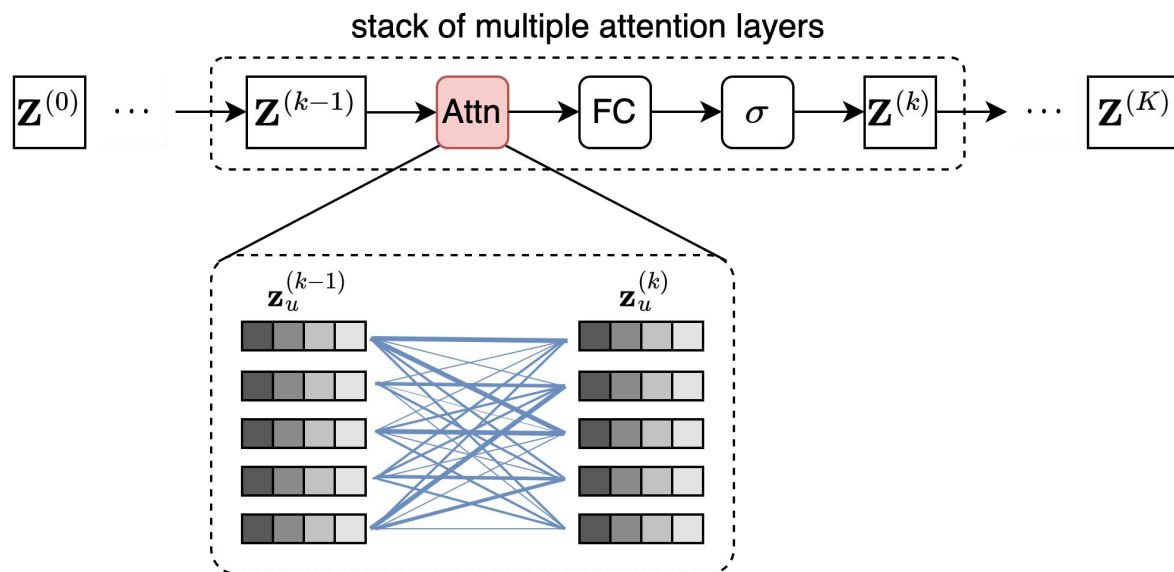
$$= \frac{\sum_{j=1}^N \mathbf{z}_j^{(k)} + \left(\sum_{j=1}^N \tilde{\mathbf{z}}_j^{(k)} \cdot (\mathbf{z}_j^{(k)})^\top\right) \cdot \tilde{\mathbf{z}}_i^{(k)}}{N + (\tilde{\mathbf{z}}_i^{(k)})^\top \sum_{l=1}^N \tilde{\mathbf{z}}_l^{(k)}}$$

$O(N)$

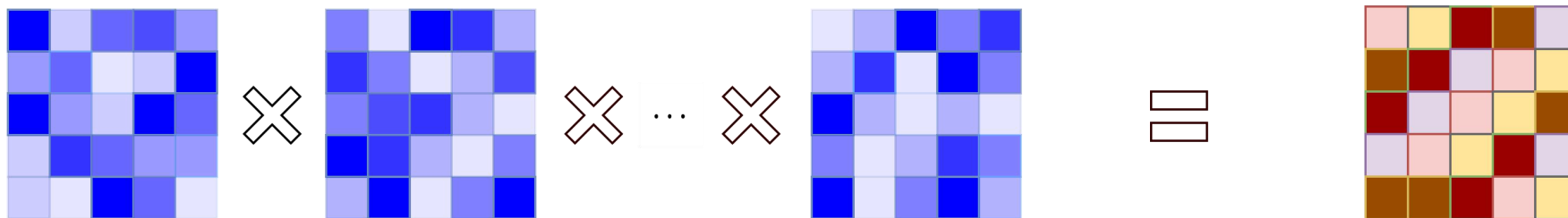
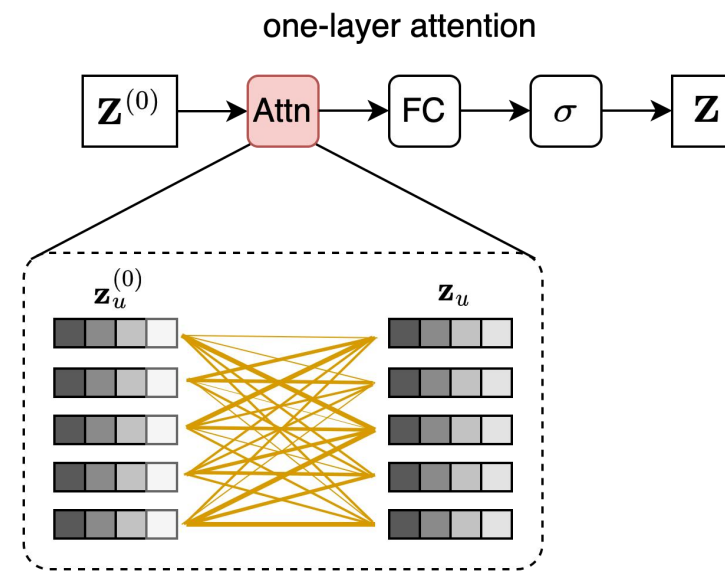
Qitian Wu et al., DIFFormer: Scalable (Graph) Transformers Induced by Energy Constrained Diffusion, ICLR 2023

Do We Really Need Deep Attention Layers?

Prior Art



Our Model



Qitian Wu et al., SGFormer: Simplifying and Empowering Transformers on Large-Graph Representations, NeurIPS 2023

SGFormer: Simplified Graph Transformers

Observation: one-layer all-pair attention is expressive enough for propagating global information among arbitrary node pairs

SGFormer: **one-layer single-head** global attention + auxiliary GNN

- Simple attention with linear complexity: $\mathbf{Z}^{(0)} = f_I(\mathbf{X})$

$$\mathbf{Q} = f_Q(\mathbf{Z}^{(0)}), \quad \tilde{\mathbf{Q}} = \frac{\mathbf{Q}}{\|\mathbf{Q}\|_{\mathcal{F}}}, \quad \mathbf{K} = f_K(\mathbf{Z}^{(0)}), \quad \tilde{\mathbf{K}} = \frac{\mathbf{K}}{\|\mathbf{K}\|_{\mathcal{F}}}, \quad \mathbf{V} = f_V(\mathbf{Z}^{(0)})$$

$$\mathbf{D} = \text{diag} \left(\mathbf{1} + \frac{1}{N} \tilde{\mathbf{Q}}(\tilde{\mathbf{K}}^\top \mathbf{1}) \right), \quad \mathbf{Z} = \beta \mathbf{D}^{-1} \left[\mathbf{V} + \frac{1}{N} \tilde{\mathbf{Q}}(\tilde{\mathbf{K}}^\top \mathbf{V}) \right] + (1 - \beta) \mathbf{Z}^{(0)}$$

- Add an auxiliary GNN at the **output layer**:

$$\mathbf{Z}_O = (1 - \alpha) \mathbf{Z} + \alpha \text{GN}(\mathbf{Z}^{(0)}, \mathbf{A}), \quad \hat{\mathbf{Y}} = f_O(\mathbf{Z}_O)$$

Qitian Wu et al., SGFormer: Simplifying and Empowering Transformers on Large-Graph Representations, NeurIPS 2023

Comparison of Existing Graph Transformers

	pos emb	multi-head	pre-processing	all-pair expressivity	complexity	largest demo of #nodes
GraphTransformer [Dwivedi et al. 2020]	R	R	R	yes	$O(N^2)$	0.2K
Graphormer [Ying et al. 2021]	R	R	R	yes	$O(N^2)$	0.3K
SAT [Chen et al. 2022]	R	R	R	yes	$O(N^2)$	0.2K
EGT [Hussain et al. 2022]	R	R	R	yes	$O(N^2)$	0.5K
GraphGPS [Ramp��se et al. 2022]	R	R	R	yes	$O(N^2)$	1.0K
NodeFormer [Wu et al. 2022]	R	R	-	yes	$O(N + E)$	2M
SGFormer	-	-	-	yes	$O(N + E)$	0.1B

Qitian Wu et al., SGFormer: Simplifying and Empowering Transformers on Large-Graph Representations, NeurIPS 2023

Experiment on Medium-Sized Graphs

Results on medium-sized node classification graphs

	homophilous graphs			heterophilic graphs			
Dataset	Cora	CiteSeer	PubMed	Actor	Squirrel	Chameleon	Deezer
# nodes	2,708	3,327	19,717	7,600	2,223	890	28,281
# edges	5,278	4,552	44,324	29,926	46,998	8,854	92,752
GCN	81.6 \pm 0.4	71.6 \pm 0.4	78.8 \pm 0.6	30.1 \pm 0.2	38.6 \pm 1.8	41.3 \pm 3.0	62.7 \pm 0.7
GAT	83.0 \pm 0.7	72.1 \pm 1.1	79.0 \pm 0.4	29.8 \pm 0.6	35.6 \pm 2.1	39.2 \pm 3.1	61.7 \pm 0.8
SGC	80.1 \pm 0.2	71.9 \pm 0.1	78.7 \pm 0.1	27.0 \pm 0.9	39.3 \pm 2.3	39.0 \pm 3.3	62.3 \pm 0.4
JKNet	81.8 \pm 0.5	70.7 \pm 0.7	78.8 \pm 0.7	30.8 \pm 0.7	39.4 \pm 1.6	39.4 \pm 3.8	61.5 \pm 0.4
APPNP	83.3 \pm 0.5	71.8 \pm 0.5	80.1 \pm 0.2	31.3 \pm 1.5	35.3 \pm 1.9	38.4 \pm 3.5	66.1 \pm 0.6
H ₂ GCN	82.5 \pm 0.8	71.4 \pm 0.7	79.4 \pm 0.4	34.4 \pm 1.7	35.1 \pm 1.2	38.1 \pm 4.0	66.2 \pm 0.8
SIGN	82.1 \pm 0.3	72.4 \pm 0.8	79.5 \pm 0.5	36.5 \pm 1.0	40.7 \pm 2.5	41.7 \pm 2.2	66.3 \pm 0.3
CPGNN	80.8 \pm 0.4	71.6 \pm 0.4	78.5 \pm 0.7	34.5 \pm 0.7	38.9 \pm 1.2	40.8 \pm 2.0	65.8 \pm 0.3
GloGNN	81.9 \pm 0.4	72.1 \pm 0.6	78.9 \pm 0.4	36.4 \pm 1.6	35.7 \pm 1.3	40.2 \pm 3.9	65.8 \pm 0.8
Graphormer _{SMALL}	OOM	OOM	OOM	OOM	OOM	OOM	OOM
Graphormer _{SMALLER}	75.8 \pm 1.1	65.6 \pm 0.6	OOM	OOM	40.9 \pm 2.5	41.9 \pm 2.8	OOM
Graphormer _{ULTRASmall}	74.2 \pm 0.9	63.6 \pm 1.0	OOM	33.9 \pm 1.4	39.9 \pm 2.4	41.3 \pm 2.8	OOM
GraphTrans _{SMALL}	80.7 \pm 0.9	69.5 \pm 0.7	OOM	32.6 \pm 0.7	41.0 \pm 2.8	42.8 \pm 3.3	OOM
GraphTrans _{ULTRASmall}	81.7 \pm 0.6	70.2 \pm 0.8	77.4 \pm 0.5	32.1 \pm 0.8	40.6 \pm 2.4	42.2 \pm 2.9	OOM
NodeFormer	82.2 \pm 0.9	72.5 \pm 1.1	79.9 \pm 1.0	36.9 \pm 1.0	38.5 \pm 1.5	34.7 \pm 4.1	66.4 \pm 0.7
SGFormer	84.5 \pm 0.8	72.6 \pm 0.2	80.3 \pm 0.6	37.9 \pm 1.1	41.8 \pm 2.2	44.9 \pm 3.9	67.1 \pm 1.1

Qitian Wu et al., SGFormer: Simplifying and Empowering Transformers on Large-Graph Representations, NeurIPS 2023

Experiment on Large-Sized Graphs

Results on large node classification graphs

Method	ogbn-proteins	Amazon2m	pokec	ogbn-arxiv	ogbn-papers100M
# nodes	132,534	2,449,029	1,632,803	169,343	111,059,956
# edges	39,561,252	61,859,140	30,622,564	1,166,243	1,615,685,872
MLP	72.04 \pm 0.48	63.46 \pm 0.10	60.15 \pm 0.03	55.50 \pm 0.23	47.24 \pm 0.31
GCN	72.51 \pm 0.35	83.90 \pm 0.10	62.31 \pm 1.13	71.74 \pm 0.29	OOM
SGC	70.31 \pm 0.23	81.21 \pm 0.12	52.03 \pm 0.84	67.79 \pm 0.27	63.29 \pm 0.19
GCN-NSampler	73.51 \pm 1.31	83.84 \pm 0.42	63.75 \pm 0.77	68.50 \pm 0.23	62.04 \pm 0.27
GAT-NSampler	74.63 \pm 1.24	85.17 \pm 0.32	62.32 \pm 0.65	67.63 \pm 0.23	63.47 \pm 0.39
SIGN	71.24 \pm 0.46	80.98 \pm 0.31	68.01 \pm 0.25	70.28 \pm 0.25	65.11 \pm 0.14
NodeFormer	77.45 \pm 1.15	87.85 \pm 0.24	70.32 \pm 0.45	59.90 \pm 0.42	-
SGFormer	79.53 \pm 0.38	89.09 \pm 0.10	73.76 \pm 0.24	72.63 \pm 0.13	66.01 \pm 0.37

SGFormer can be trained in full-graph manner on obgn-arxiv

Mini-batch training for proteins, Amazon2M, pokec with batch size 10K/100K

For Papers100M, using batch size **0.4M** only requires **3.5 hours** on a **24GB GPU**

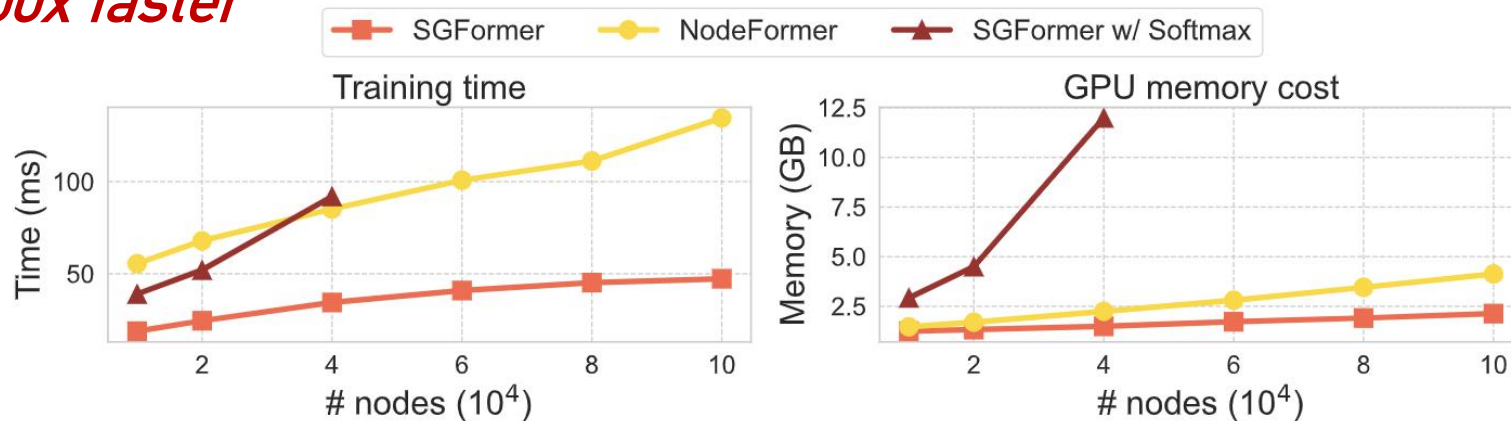
Qitian Wu et al., SGFormer: Simplifying and Empowering Transformers on Large-Graph Representations, NeurIPS 2023

Experiment Results

Comparison of training/inference time per epoch and memory cost

Method	Cora			PubMed			Amazon2M		
	Tr (ms)	Inf (ms)	Mem (GB)	Tr (ms)	Inf (ms)	Mem (GB)	Tr (ms)	Inf (ms)	Mem (GB)
Graphormer	563.5	537.1	5.0	-	-	-	-	-	-
GraphTrans	160.4	40.2	3.8	-	-	-	-	-	-
NodeFormer	68.5	30.2	1.2	321.4	135.5	2.9	5369.5	1410.0	4.6
SGFormer	15.0	3.8	0.9	15.4	4.4	1.0	2481.4	382.5	2.7

100x faster



Scalability test of training time/memory costs w.r.t. number of nodes

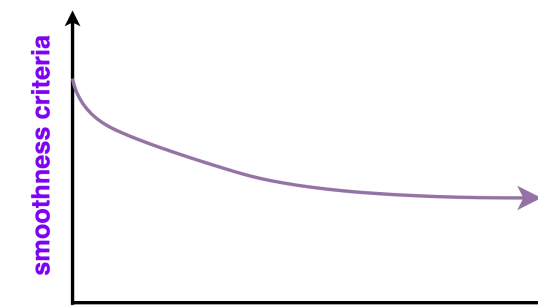
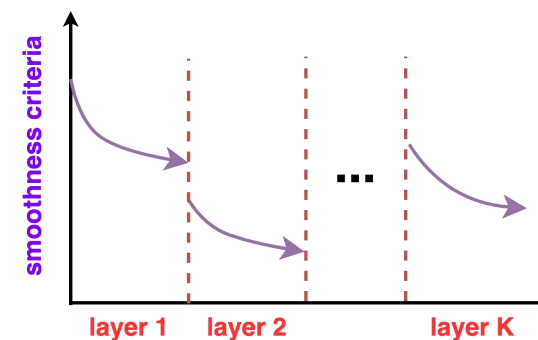
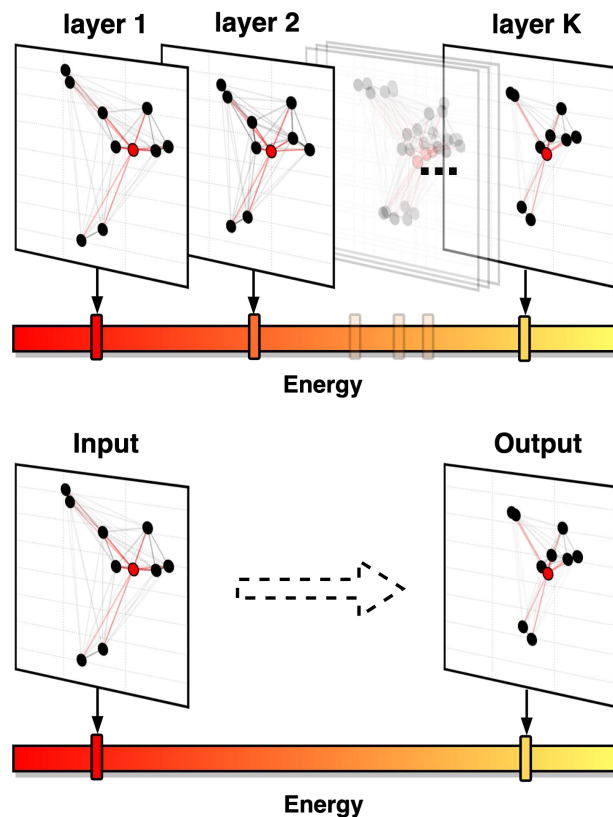
How Powerful Are One-Layer Attentions?

stack of multiple attention layers

$$\mathbf{z}_u^{(k)} = (1 - \tau)\mathbf{z}_u^{(k-1)} + \tau \sum_{v=1}^N c_{uv}^{(k)} \mathbf{z}_v^{(k-1)}$$

one-layer attention

$$\mathbf{z}_u = (1 - \tau)\mathbf{z}_u^{(0)} + \tau \sum_{v=1}^N c_{uv} \mathbf{z}_v^{(0)}$$



Theorem (Equivalence between Multi-Layer Attentions and One-Layer Attention)

For any K-layer attention, there exists a one-layer model that induces the same denoising effect.

Conclusions

Graph Transformers have become a popular research topic in ML community

Some open problems: 1) poor scalability (quadratic complexity)
2) lack of principled guidance for attention designs
3) inefficiency, complicated model

[1] NodeFormer: A Scalable Graph Structure Learning Transformer for Node Classification, in NeurIPS 2022

all-pair message passing with linear complexity scale to **2M** nodes handle no-graph tasks

Codes: <https://github.com/qitianwu/NodeFormer>

[2] DIFFormer: Scalable (Graph) Transformers Induced by Energy Constrained Diffusion, in ICLR 2023

principled global attention designs superiority for **low labeled rates**

Codes: <https://github.com/qitianwu/DIFFormer>

[3] SGFormerSimplifying and Empowering Transformers for Large-Graph Representations, in NeurIPS 2023

simple attention (one-layer single-head) **100x** inference speed-up scale to **0.1B** nodes

Codes: <https://github.com/qitianwu/SGFormer>