EXPHORMER SPARSE TRANSFORMERS FOR GRPAHS

Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J. Sutherland , Ali Kemal Sinop

> Seminar in Deep Neural Networks Presenter: Johannes Herter May 21, 2023



GRAPH LEARNING: GRAPHS ARE WIDELY APPLICABLE ACROSS DOMAINS

Social networks



Chemoinformatics



Product recommendations



Traffic prediction



Weather forecasting



Protein-protein associations











Under-reaching



a node cannot see some nodes (Barceló et al., 2020)

Under-reaching



Over-smoothing



all nodes start to look *the same* (Oono and Suzuki, 2020)



Over-smoothing



Over-squashing

information gets "*lost*" in transit (Topping et al., 2022)



Over-smoothing



Limited Expressivity





Over-smoothing



Limited Expressivity





Over-smoothing



Limited Expressivity



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

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 modelling (Rives et al., 2021)

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

TYPES OF TRANSFORMER ARCHITECTURES

(shifted right)



Decoder-only

Output

Probabilities

Softmax

Linear

Add & Norm

Feed

Forward

Add & Norm

Multi-Head

Attention

Add & Norm

Masked

Multi-Head

Attention

⊕-

Output

Embedding

Outputs

(shifted right)

V

N×

Positional

Encoding

TYPES OF TRANSFORMER ARCHITECTURES



Outputs (shifted right) $N \times$

Positional

Encoding





NLP

Input tokenization

Graph Learning



NLP

Input tokenization

Graph Learning





Graph Learning





NLP

Graph Learning

- Input tokenization
- Positional encoding
 - Captures a word's position in a text
 - Is added/ concatenated to the input



NLP

- Input tokenization
- Positional encoding

Graph Learning

- Input treats the graph as bag of nodes
- Encode structure via positional encoding



NLP

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

Graph Learning

- Input treats the graph as bag of nodes
- Encode structure via positional encoding



NLP

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

Graph Learning

- Input treats the graph as bag of nodes
- Encode structure via positional encoding
- Transformer layer
 with self-attention

ABSOLUTE POSITIONAL ENCODING (APE)



PE ON GRAPHS ?

PE ON GRAPHS: LAPLACIAN POSITIONAL ENCODING

- Graph Laplacian: How does a node relate to its neighbours
- \bullet U is sorted by increasing order of the corresponding eigenvalue.
- ↔ Take the first 1, ..., k eigenvectors of Δ

 $\Delta = I_n - D^{-\frac{1}{2}}AD^{\frac{1}{2}}$ $PE_{LPE} = [U_1, \dots, U_k]$

PE ON GRAPHS: LAPLACIAN POSITIONAL ENCODING

Graph Laplacian: How does a node relate to its neighbours

DEMO

ABSOLUTE POSITIONAL ENCODING VS LAPLACIAN POSITIONAL ENCODING





APE of a sequence

LPE of a graph



LIMITATIONS OF (GRAPH) TRANSFORMERS

- Loss of inductive bias (locality)
- Self-Attention is $\mathcal{O}(N^2)$!
- ✤ A lot of research in developing sparse attention
- Tradeoff between performance and speed





WHAT WE HAVE SEEN SO FAR

- Graph Representation Learning
- Adapting Transformers for Graph Data
- Limitations of Graph Transformers
- Now: Graph-GPS

GRAPH-GPS: GENERAL POWERFUL SCALABLE GRAPH TRANSFORMER

Positional & Structural Features



Combine MPNN and Transformer



(Rampášek et al.,2023)

GRAPH-GPS: GENERAL POWERFUL SCALABLE GRAPH TRANSFORMER



Why not using *sparse attention* mechanisms more tailored *for graphs*?

Original Graph



Full Attention



Too memory intensitive

(Schirzad et al.,2023)

Original Graph



 Graphs carry much more topological structure than sequences



Original Graph



 Graphs carry much more topological structure than sequences
 Weak long-range information flow

Local Neighbourhood Attention

Original Graph







Diameter 2Information bottleneck



Global Attention



Expander Graph Attention



Approximate complete graphs

- Spectral properties
- Mixing properties
- Diameter is $\mathcal{O}(\log N)$



- Spectral properties:
 - A d-regular expander graph on n vertices approximates the complete graph K_n



- Spectral properties:
 - A *d*-regular expander graph on *n* vertices approximates the complete graph *K*_n



- Spectral properties:
 - A d-regular expander graph on n vertices approximates the complete graph K_n
 - There exist simple algorithms to create *d*-regular expander graphs
- **\bigstar** Compared to K_n it has only O(n) edges



Mixing properties:

• For a *d*-regular expander graph a random walk mixes well



Mixing properties:

- For a *d*-regular expander graph a random walk mixes well
- The diameter is $\mathcal{O}_{d,\epsilon}(\log N)$





Global Attention





Expander Graph Attention



 $\mathcal{O}(|V| + |E|)$ edges in the attention graph



REMEMBER GRAPH-GPS ?

Positional & Structural Features



Combine MPNN and Transformer



EXPHORMER: GRAPH-GPS WITH EXPANDER ATTENTION MODULE

Positional & Structural Features



Combine MPNN and Transformer



EXPHORMER models have universal approximability !

EXPHORMER: EXPERIMENTS & RESULTS - BASELINES

Model	CIFAR10 Accuracy ↑	MalNet-Tiny Accuracy ↑	MNIST Accuracy ↑	CLUSTER Accuracy ↑	PATTERN Accuracy ↑
GCN (Kipf & Welling, 2017) GIN (Xu et al., 2018) GAT (Veličković et al., 2018)	55.71 ± 0.381 55.26 ± 1.527 64.22 ± 0.455	81.0 88.98 ± 0.557 92.1 ± 0.242	90.71 ± 0.218 96.49 ± 0.252 95.54 ± 0.205	$\begin{array}{c} 68.50 \pm 0.976 \\ 64.72 \pm 1.553 \\ 70.59 \pm 0.447 \end{array}$	$71.89 \pm 0.334 \\ 85.39 \pm 0.136 \\ 78.27 \pm 0.186$
GatedGCN (Bresson & Laurent, 2017; Dwivedi et al., 2020)	67.31±0.311	92.23±0.65	97.34±0.143	73.84 ± 0.326	85.57 ± 0.088
PNA (Corso et al., 2020)	70.35 ± 0.63	_	$97.94{\pm}0.12$	_	_
DGN (Beaini et al., 2021)	72.84±0.417	-	-	-	$86.68 {\pm} 0.034$
CRaWl (Toenshoff et al., 2021)	69.01±0.259	_	97.94±0.050	_	_
GIN-AK+ (Zhao et al., 2022b)	72.19 ± 0.13	—	—	—	86.85±0.057
SAN (Kreuzer et al., 2021)	_	_	_	$76.69 {\pm} 0.65$	$86.58 {\pm} 0.037$
K-Subgraph SAT (Chen et al., 2022a)	_	-	-	$77.86 {\pm} 0.104$	86.85±0.037
EGT (Hussain et al., 2021)	68.70 ± 0.409		98.17±0.087	79.23±0.348	86.82±0.020
GraphGPS (Rampásek et al., 2022)	72.30±0.356	93.50±0.41	98.05±0.126	$78.02{\pm}0.180$	86.69±0.059
EXPHORMER (ours)	74.69±0.125	$\textbf{94.02} \pm \textbf{0.209}$	98.55 ± 0.039	$\textbf{78.07} \pm \textbf{0.037}$	86.74±0.015

EXPHORMER: PROMISING RESULTS ON LONG-RANGE GRAPH BENCHMARKS

Model	PascalVOC-SP F1 score ↑	COCO-SP F1 score ↑	Peptides-Func AP ↑	Peptides-Struct MAE↓	PCQM-Contact MRR ↑	 first second
GCN GINE GatedGCN GatedGCN+RWSE	$\begin{array}{c} 0.1268 \pm 0.0060 \\ 0.1265 \pm 0.0076 \\ 0.2873 \pm 0.0219 \\ 0.2860 \pm 0.0085 \end{array}$	$\begin{array}{c} 0.0841 \pm 0.0010 \\ 0.1339 \pm 0.0044 \\ \textbf{0.2641} \pm \textbf{0.0045} \\ 0.2574 \pm 0.0034 \end{array}$	$\begin{array}{c} 0.5930 \pm 0.0023 \\ 0.5498 \pm 0.0079 \\ 0.5864 \pm 0.0077 \\ 0.6069 \pm 0.0035 \end{array}$	$\begin{array}{c} 0.3496 \pm 0.0013 \\ 0.3547 \pm 0.0045 \\ 0.3420 \pm 0.0013 \\ 0.3357 \pm 0.0006 \end{array}$	$\begin{array}{c} 0.3234 \pm 0.0006 \\ 0.3180 \pm 0.0027 \\ 0.3218 \pm 0.0011 \\ 0.3242 \pm 0.0008 \end{array}$	third
Transformer+LapPE SAN+LapPE SAN+RWSE GraphGPS	$\begin{array}{c} 0.2694 \pm 0.0098 \\ \textbf{0.3230} \pm \textbf{0.0039} \\ 0.3216 \pm 0.0027 \\ \textbf{0.3748} \pm \textbf{0.0109} \end{array}$	$\begin{array}{c} 0.2618 \pm 0.0031 \\ 0.2592 \pm 0.0158 * \\ 0.2434 \pm 0.0156 * \\ \textbf{0.3412} \pm \textbf{0.0044} \end{array}$	$\begin{array}{c} 0.6326 \pm 0.0126 \\ 0.6384 \pm 0.0121 \\ \textbf{0.6439} \pm \textbf{0.0075} \\ \textbf{0.6535} \pm \textbf{0.0041} \end{array}$	$\begin{array}{c} \textbf{0.2529} \pm \textbf{0.0016} \\ 0.2683 \pm 0.0043 \\ 0.2545 \pm 0.0012 \\ \textbf{0.2500} \pm \textbf{0.0005} \end{array}$	$\begin{array}{c} 0.3174 \pm 0.0020 \\ \textbf{0.3350} \pm \textbf{0.0003} \\ \textbf{0.3341} \pm \textbf{0.0006} \\ 0.3337 \pm 0.0006 \end{array}$	
Exphormer (ours)	$\textbf{0.3975} \pm \textbf{0.0037}$	$\textbf{0.3455} \pm \textbf{0.0009}$	$\textbf{0.6527} \pm \textbf{0.0043}$	$\textbf{0.2481} \pm \textbf{0.0007}$	$\textbf{0.3637} \pm \textbf{0.0020}$	

EXPHORMER: NO ONE-FITS-ALL SOLUTION!

Dataset	No Local Edges	No Expander Edges	No Global Nodes	All Components
Cifar10	74.62 ± 0.12	74.53 ± 0.19	74.68 ± 0.19	74.69 ± 0.13
Malnet-Tiny	92.64 ± 0.55	94.02 ± 0.21	92.48 ± 0.33	92.06 ± 0.18
Pattern	86.59 ± 0.03	86.70 ± 0.02	86.14 ± 0.08	86.74 ± 0.02
PascalVOC-SP	0.3708 ± 0.0039	0.3588 ± 0.0013	0.3975 ± 0.0037	0.3682 ± 0.0042
Peptides-Struct	0.2631 ± 0.0007	$\boldsymbol{0.2481 \pm 0.0007}$	0.2655 ± 0.0003	0.2643 ± 0.0008
Computer	90.34 ± 0.45	91.48 ± 0.41	91.59 ± 0.31	91.43 ± 0.53

EXPHORMER: SOME IMPROVEMENTS BUT NO ONE-FITS-ALL SOLUTION!



node cannot see some nodes (Barceló et al., 2020)



all nodes start to look *the same* (Oono and Suzuki, 2020)



information gets "*lost*" in transit (Topping et al., 2022) Limited Expressivity



EXPHORMER: SOME IMPROVEMENTS BUT NO ONE-FITS-ALL SOLUTION!

- Idea to use expanders as attention graphs is intuitive
- * Exphormer promotes its ideas well, but in the end is just a small extension to

Graph-GPS

- Need for careful tuning of the components !
- Benchmarking is hard to compare







THANK YOU FOR YOUR ATTENTION !

Johannes Herter

johannes.herter@inf.ethz.ch

BACKUP: EXPANDER GRAPH CONSTRUCTION IS SIMPLE AND FAST!

Generating a Random Regular Expander We now describe how we generate a random regular expander. Let G = (V, E) be the original graph, where $V = \{1, 2, ..., n\}$. For the purposes of experimentation (in Tables 1 to 5), we use the random graph process analyzed in Friedman (2003) (see Theorem C.2 in Appendix C) to generate a random d-regular graph G' = (V, E') on the same node set V:

- Pick d/2 permutations π₁, π₂,..., π_{d/2} on V, each π_i chosen independently and uniformly among all possible permutations of n elements.
- Then, letting [k] denote $\{1, 2, \ldots, k\}$, choose

 $E' = \left\{ (i, \pi_j(i)), (i, \pi_j^{-1}(i)) : j \in [d/2], i \in [n] \right\}.$

Theorem C.2. (Friedman, 2003, Theorem 1.1) Fix $\epsilon > 0$ and an even integer d > 2. Then, suppose G = (V, E) is a random graph generated by taking d/2 independent uniformly random permutations $\pi_1, \pi_2, \ldots, \pi_{d/2}$ on $V = \{1, 2, \ldots, n\}$ and then choosing the edge set as

$$E = \left\{ (i, \pi_j(i)), (i, \pi_j^{-1}(i)) : 1 \le j \le d, 1 \le i \le n \right\}.$$

Then, with probability $1 - O(n^{-\Omega(\sqrt{d})})$, G satisfies $\lambda_j(G) \le 2\sqrt{d-1} + \epsilon$ for j = 2, ..., n, where $d = \lambda_1(G) \ge \lambda_2(G) \ge \cdots \ge \lambda_n(G) \ge -d$ are the eigenvalues of the adjacency matrix of G.

BACKUP: EXPANDER HAVE NICE SPECTRAL AND MIXING PROPERTIES

Theorem 4.1. (Spielman, 2019, Section 27.2) A d-regular ϵ -expander G on n vertices spectrally approximates the complete graph K_n on n vertices:²

$$(1-\epsilon)\frac{1}{n}L_K \preceq \frac{1}{d}L_G \preceq (1+\epsilon)\frac{1}{n}L_K.$$

Lemma 4.2. (*Hoory et al.*, 2006, *Theorem 3.2*) Let G = (V, E) be a d-regular ϵ -expander graph on n = |V| nodes. For any initial distribution $\pi^{(0)} : V \to \mathbb{R}^+$ and any $\delta > 0$, $\pi^{(t)}$ satisfies $\|\pi^{(t)} - \frac{1}{n}\|_1 \leq \delta$ as long as $t \geq \frac{1}{2(1-\varepsilon)} \log(n/\delta^2)$.

BACKUP: RESULTS GRAPHGPS VS EXPHORMER

Model/Dataset	Cifar10 Accuracy ↑	MalNet-Tiny Accuracy ↑	PascalVOC-SP F1 score ↑	Peptides-Func AP↑
GPS (MPNN-only)	$ 69.948 \pm 0.499$	92.23 ± 0.65	0.3016 ± 0.0031	0.6159 ± 0.0048
GPS-BigBird GPS-Performer GPS-Transformer	$ \begin{array}{c} 70.480 \pm 0.106 \\ \textbf{70.670} \pm \textbf{0.338} \\ \textbf{72.305} \pm \textbf{0.344} \end{array} $	$\begin{array}{c} 92.34 \pm 0.34 \\ \textbf{92.64} \pm \textbf{0.78} \\ \textbf{93.50} \pm \textbf{0.41} \end{array}$	$\begin{array}{c} 0.2762 \pm 0.0069 \\ \textbf{0.3724} \pm \textbf{0.0131} \\ \textbf{0.3736} \pm \textbf{0.0158} \end{array}$	$\begin{array}{c} 0.5854 \pm 0.0079 \\ \textbf{0.6475} \pm \textbf{0.0056} \\ \textbf{0.6535} \pm \textbf{0.0041} \end{array}$
EXPHORMER	74.69±0.125	$\textbf{94.02} \pm \textbf{0.21}$	$\textbf{0.3975} \pm \textbf{0.0037}$	$\textbf{0.6527} \pm \textbf{0.0043}$

BACKUP: NO COMPARISON TO GRIT ??

Model	ZINC	MNIST	CIFAR10	PATTERN	CLUSTER
	MAE↓	Accuracy ↑	Accuracy	Accuracy ↑	Accuracy ↑
GCN	0.367 ± 0.011	90.705 ± 0.218	55.710 ± 0.381	71.892 ± 0.334	68.498 ± 0.976
GIN	0.526 ± 0.051	96.485 ± 0.252	55.255 ± 1.527	85.387 ± 0.136	64.716 ± 1.553
GAT	0.384 ± 0.007	95.535 ± 0.205	64.223 ± 0.455	78.271 ± 0.186	70.587 ± 0.447
GatedGCN	0.282 ± 0.015	97.340 ± 0.143	67.312 ± 0.311	85.568 ± 0.088	73.840 ± 0.326
GatedGCN-LSPE	0.090 ± 0.001	-	-	-	_
PNA	0.188 ± 0.004	97.94 ± 0.12	70.35 ± 0.63	-	_
DGN	0.168 ± 0.003	—	72.838 ± 0.417	86.680 ± 0.034	—
GSN	0.101 ± 0.010	_	_	_	_
CIN	0.079 ± 0.006	_	_	_	_
CRaW1	0.085 ± 0.004	97.944 ± 0.050	69.013 ± 0.259	_	_
GIN-AK+	0.080 ± 0.001	_	72.19 ± 0.13	86.850 ± 0.057	_
SAN	0.139 ± 0.006	_	_	86.581 ± 0.037	76.691 ± 0.65
Graphormer	0.122 ± 0.006	_	-	_	—
K-Subgraph SAT	0.094 ± 0.008	_	-	86.848 ± 0.037	77.856 ± 0.104
EGT	0.108 ± 0.009	98.173 ± 0.087	68.702 ± 0.409	86.821 ± 0.020	79.232 ± 0.348
Graphormer-URPE	0.086 ± 0.007	-	_	—	—
Graphormer-GD	0.081 ± 0.009	_	-	_	_
GPS	0.070 ± 0.004	98.051 ± 0.126	72.298 ± 0.356	86.685 ± 0.059	78.016 ± 0.180
GRIT (ours)	$\boldsymbol{0.059 \pm 0.002^*}$	98.108 ± 0.111	${\bf 76.468 \pm 0.881^{*}}$	${\bf 87.196 \pm 0.076^{*}}$	${\bf 80.026 \pm 0.277^*}$

Model	CIFAR10 Accuracy ↑	MalNet-Tiny Accuracy ↑	MNIST Accuracy ↑	CLUSTER Accuracy ↑	PATTERN Accuracy ↑
GCN (Kipf & Welling, 2017) GIN (Xu et al., 2018) GAT (Veličković et al., 2018) GatedGCN (Bresson & Laurent, 2017; Dwivedi et al., 2020) PNA (Corso et al., 2020) DGN (Beaini et al., 2021)	55.71 ± 0.381 55.26 ± 1.527 64.22 ± 0.455 67.31 ± 0.311 70.35 ± 0.63 72.84 ± 0.417	81.0 88.98±0.557 92.1 ±0.242 92.23±0.65	90.71 \pm 0.218 96.49 \pm 0.252 95.54 \pm 0.205 97.34 \pm 0.143 97.94 \pm 0.12	$\begin{array}{c} 68.50 \pm 0.976 \\ 64.72 \pm 1.553 \\ 70.59 \pm 0.447 \\ 73.84 \pm 0.326 \\ \end{array}$	$71.89 \pm 0.334 \\ 85.39 \pm 0.136 \\ 78.27 \pm 0.186 \\ 85.57 \pm 0.088 \\ - \\ 86.68 \pm 0.034$
CRaWl (Toenshoff et al., 2021) GIN-AK+ (Zhao et al., 2022b)	69.01±0.259 72.19±0.13		97.94±0.050 _		_ 86.85±0.057
SAN (Kreuzer et al., 2021) K-Subgraph SAT (Chen et al., 2022a) EGT (Hussain et al., 2021) GraphGPS (Rampásek et al., 2022)	 68.70±0.409 72.30±0.356	_ _ 93.50±0.41	- 98.17±0.087 98.05±0.126	$\begin{array}{c} 76.69 {\pm} 0.65 \\ 77.86 {\pm} 0.104 \\ \textbf{79.23} {\pm} \textbf{0.348} \\ \textbf{78.02 {\pm} 0.180} \end{array}$	86.58±0.037 86.85±0.037 86.82±0.020 86.69±0.059
EXPHORMER (ours)	74.69±0.125	$\textbf{94.02} \pm \textbf{0.209}$	$\textbf{98.55} \pm \textbf{0.039}$	$\textbf{78.07} \pm \textbf{0.037}$	86.74±0.015