Structure-Aware Transformer for Graph Representation Learning Dexiong Chen^{\dagger ,1,2}, Leslie O'Bray^{\dagger ,1,2} and Karsten Borgwardt^{1,2}

[†]: These authors contributed equally. ¹: Department of Biosystems Science and Engineering, ETH Zürich, Switzerland. ²: SIB Swiss Institute of Bioinformatics, Switzerland.

Motivation and contribution

- GNNs struggle from problems of under-reaching, oversmoothing, and over-squashing.
- Transformers can address those issues, but have had limited performance to date on graphs because the existing self-attention calculation uses *attributed similarity* between nodes, as opposed to *structural similarity*.
- Information about the graph is typically only incorporated via a positional encoding.
- Nodes can have the same position while being structurally different (see nodes *u* and *v*).



Structure-Aware Transformer (SAT)

- We extend self-[™]attention to account for local structures by extracting a *subgraph representation* rooted at each node.
- We propose several methods for automatically generating subgraph representations.
- The resulting framework can leverage *any* GNN to extract subgraph representations and empirically outperforms the base GNN \Rightarrow SAT is an effortless enhancer of any GNN.
- The representations from the structure-aware attention are *at least as expressive as* the representations from the structure extractor.
- We provide a bound on the distance between the structure-aware representations.

SAT enhances sparse GNNs (ZINC)



- SAT uses a GNN (i.e. a "base GNN") to create structureaware node representations in the graph Transformer.
- Empirically, SAT always improves upon the performance of the base GNN it uses.

ETHzürich

A structure-aware self attention

The original self-attention calculation is:

Attn(**X**) := softmax(
$$\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_{out}}}$$
) **V** $\in^{n \times d_{out}}$
where **Q** = **XW**_{**Q**}, **K** = **XW**_{**K**} and **V** = **XW**_{**V**}.

We can rewrite this as a kernel smoother:

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

where

$$\kappa_{\exp}(x, x') := \exp\left(\langle \mathbf{W}_{\mathbf{Q}} x, \mathbf{W}_{\mathbf{K}} x' \rangle / \sqrt{d_{out}}\right) f(x) = \mathbf{W}_{\mathbf{V}} x \in \mathbb{R}^{d_{out}}$$

We define our **structure-aware attention** as:

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

where

$$\kappa_{\mathsf{graph}}(S_G(v), S_G(u)) = \kappa_{\exp}(\varphi(v, G), \varphi(u, G)),$$

and $\varphi(v, G)$ is a structure extractor that extracts vector representations of some subgraph centered at u with node features **X**, e.g. k-subtree SAT and k-subgraph SAT:

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

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

Effect of k in SAT (ZINC \downarrow)



- The effect when considering different values of k(which determines the k-hop subgraph).
- k = 0 is equivalent to a vanilla Transformer. $k \ge 1$ incorporates structures into the node embeddings.

DBSSE

An instance of a Structure-Aware Transformer: k-subgraph SAT



Comparison to state-of-the-art GNNs and graph Transformers

	ZINC ↓	CLUSTER ↑	PATTERN ↑		OGBG-PPA ↑	OGBG-CODE2 ↑
# graphs Avg. # nodes Avg. # edges Metric	12,000 23.2 49.8 MAE	12,000 117.2 4,303.9 Accuracy	14,000 118.9 6,098.9 Accuracy	# graphs Avg. # nodes Avg. # edges Metric	158,100 243.4 2,266.1 Accuracy	452,741 125.2 124.2 F1 score
GIN GAT PNA	0.387 ± 0.015 0.384 ± 0.007 0.188 ± 0.004	64.716 ± 1.553 70.587 ± 0.447 67.077 ± 0.977	85.590 ± 0.011 78.271 ± 0.186 86.567 ± 0.075	GCN GCN-Virtual Node GIN	0.6839 ± 0.0084 0.6857 ± 0.0061 0.6892 ± 0.0100	0.1507 ± 0.0018 0.1595 ± 0.0018 0.1495 ± 0.0023
Transformer+RWPE Graph Transformer SAN	0.310 ± 0.005 0.226 ± 0.014 0.139 ± 0.006	29.622 ± 0.176 73.169 ± 0.622 76.691 ± 0.650	86.183 ± 0.019 84.808 ± 0.068 86.581 ± 0.037	GIN-Virtual Node DeeperGCN ExpC Transformer GraphTrans	0.7037±0.0107 0.7712±0.0071 0.7976 ± 0.0072	0.1581±0.0026 _ _
Graphormer k-subtree SAT	0.122±0.006 0.102±0.005	- 77.751±0.121	- 86.865±0.043		0.6454±0.0033 -	0.1670±0.0015 0.1830±0.0024
K-SUDGIAPH SAT	0.074±0.008	//.030±0.104	00.040±0.037	k-subtree SAT	0.7522 ± 0.0056	$0.1937{\pm}0.0028$

Effect of absolute encoding (ZINC \downarrow)



• Adding an absolute encoding improves performance. • The performance gain from including an absolute encoding is much smaller compared to the performance gain from using the structure-aware node embeddings.







Interpretability of SAT

• When molecules were screened for mutagenicity, SAT was better able to identify motifs that are known to be mutagenic than a vanilla Transformer with RWPE.

• SAT puts more attention weights on the known mutagenic motifs.

