

Graph Convolutions Enrich the Self-Attention in Transformers!

¹Jeongwhan Choi ²Hyowon Wi ²Jayoung Kim ²Yehjin Shin
³Kookjin Lee ⁴Nathaniel Trask ²Noseong Park

¹Yonsei University

²Korea Advanced Institute of Science and Technology (KAIST)

³Arizona State University

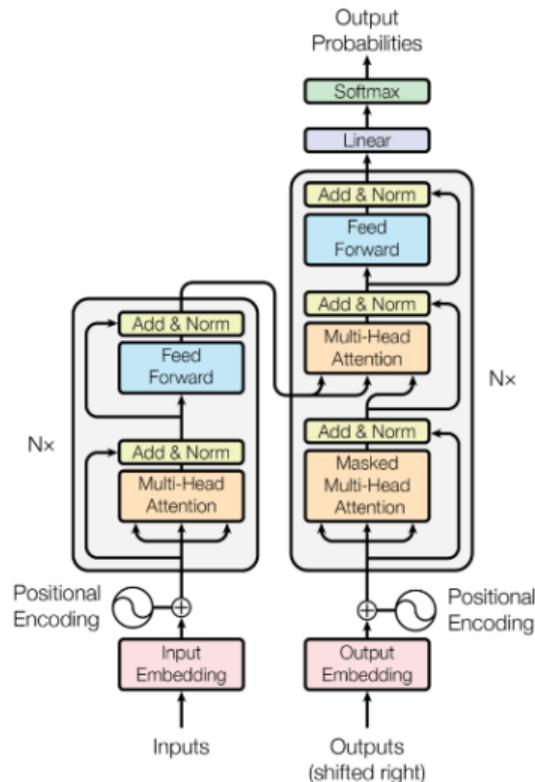
⁴University of Pennsylvania



Self-Attention is the Heart of Transformers

- The self-attention mechanism, denoted as SA, can be expressed as follows:

$$\begin{aligned} \text{SA}(\mathbf{X}) &= \text{softmax}\left(\frac{\mathbf{X}\mathbf{W}_{\text{key}}(\mathbf{X}\mathbf{W}_{\text{qry}})^{\top}}{\sqrt{d}}\right)\mathbf{X}\mathbf{W}_{\text{val}} \\ &= \bar{\mathbf{A}}\mathbf{X}\mathbf{W}_{\text{val}}. \end{aligned}$$



Graph Signal Processing

Discrete signal processing (DSP)

- Signal \mathbf{x} \rightarrow Apply Filter \mathbf{g} \rightarrow Output \mathbf{y}

$$y_i = \sum_{j=1}^n x_j g_{i-j}. \quad (1)$$

Graph signal processing (GSP) – generalization of DSP to graph domain

- The graph filter \mathbf{H} can be written with a shift operator \mathbf{S} (i.e., adjacency matrix \mathbf{A}):

$$\mathbf{y} = \mathbf{H}\mathbf{x} = \sum_{k=0}^K w_k \mathbf{S}^k \mathbf{x}, \quad (2)$$

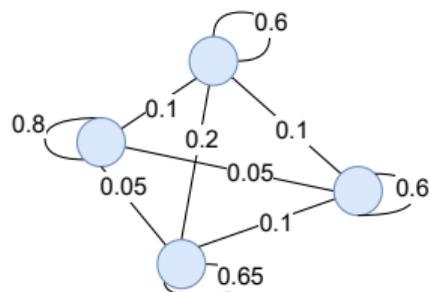
where K is the maximum order of polynomial, and w_k is a coefficient.

Self-Attention as a Graph Filter

- Self-attention matrix: $\bar{\mathbf{A}} = \mathbf{D}^{-1}\mathbf{A}$,
 - \mathbf{A} is an adjacency matrix
 - \mathbf{D} is a degree matrix.
- Self-attention can be considered as a simple graph filter ($\mathbf{H} = \bar{\mathbf{A}}$)

0.8	0.1	0.05	0.05
0.1	0.6	0.2	0.1
0.05	0.2	0.65	0.1
0.2	0.1	0.1	0.6

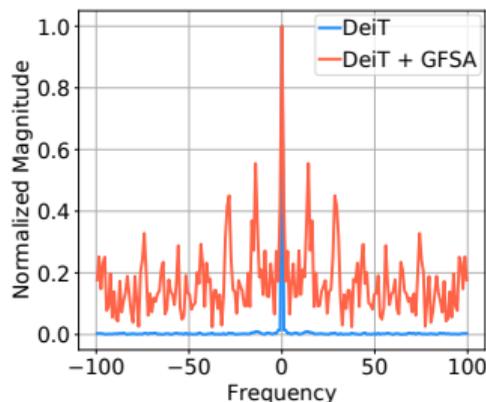
Self-attention $\bar{\mathbf{A}}$
(Probability score matrix)



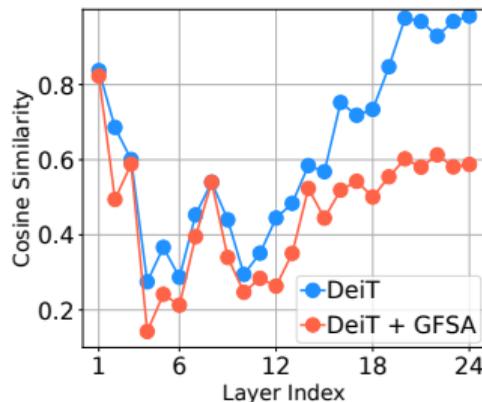
Self-Attention as
a weighted graph

- Self-attention is a weighted graph
 - Nodes \Leftrightarrow Tokens
 - Edge weights \Leftrightarrow Attention scores

Oversmoothing Problem in Transformers



(a) Filter response



(b) Cosine similarity

Figure: Filter frequency response and cosine similarity

- Oversmoothing problem:
 - As self-attention is a “low-pass filter”, high-frequency information is attenuated.
 - Latent representations tend to become similar to each other.

Graph Filter-based Self-Attention (GFSA)

- We redesign self-attention from a graph signal processing perspective
- Our proposed GFSA is defined with the graph filter $\tilde{\mathbf{H}}_{\text{GFSA}}$:

$$\text{GFSA}(\mathbf{X}) := \tilde{\mathbf{H}}_{\text{GFSA}} \mathbf{X} \mathbf{W}_{\text{val}}, \quad (3)$$

$$\tilde{\mathbf{H}}_{\text{GFSA}} = w_0 \mathbf{I} + w_1 \bar{\mathbf{A}} + w_K \underbrace{(\bar{\mathbf{A}} + (K - 1)(\bar{\mathbf{A}}^2 - \bar{\mathbf{A}}))}_{\simeq \bar{\mathbf{A}}^K}, \quad (4)$$

- We approximate $\bar{\mathbf{A}}^K$ with the first-order Taylor approximation:

$$\bar{\mathbf{A}}^K \simeq \bar{\mathbf{A}} + (K - 1)(\bar{\mathbf{A}}^2 - \bar{\mathbf{A}}). \quad (5)$$

- GFSA learns the appropriate coefficients for downstream tasks, so it can be reduced to a low-pass-only, high-pass-only, or combined filter.

Analysis of Frequency Responses with Visualization

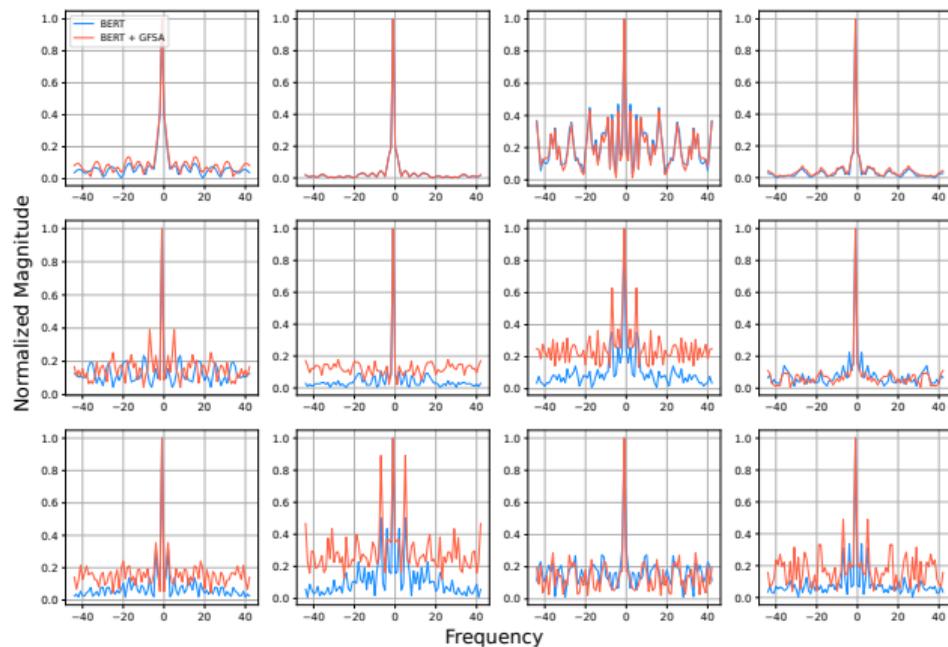


Figure: Visualization of the frequency responses for all 12 layers of BERT trained on STS-B dataset. The top-left figure corresponds to the first layer, and the bottom-right figure corresponds to the last layer.

GFSA Improves the Performance of Transformers!

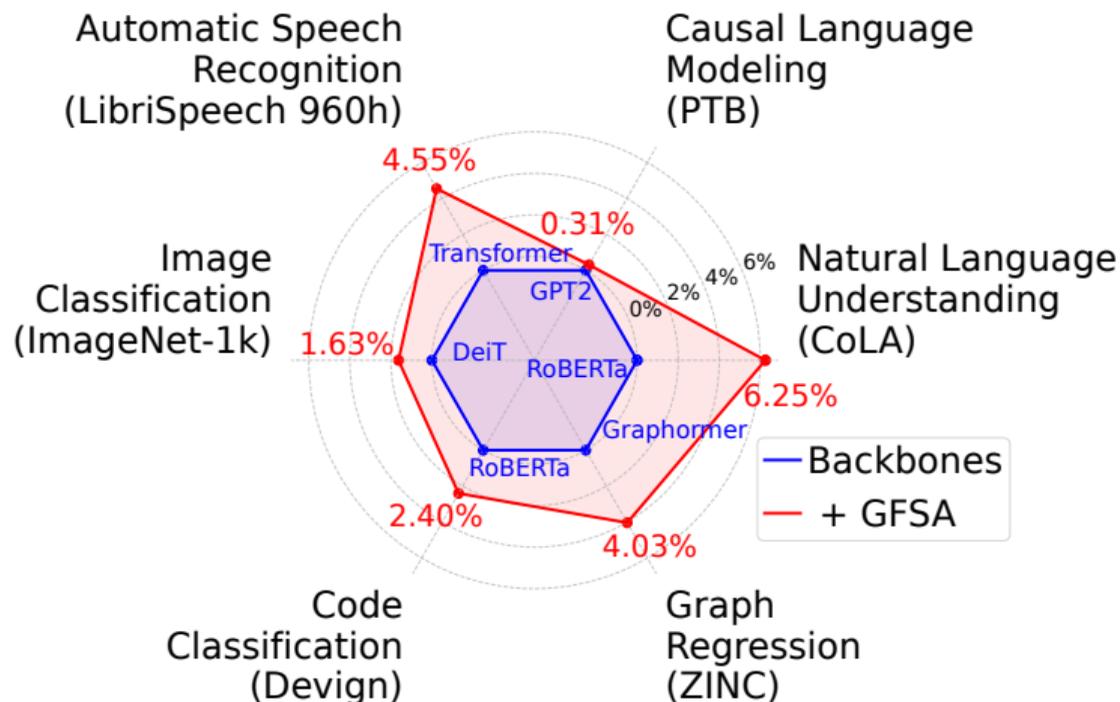


Figure: Performance improvements (%) of our GFSA when integrated with different Transformer backbones in various domains

GFSA with Efficient Design Strategies

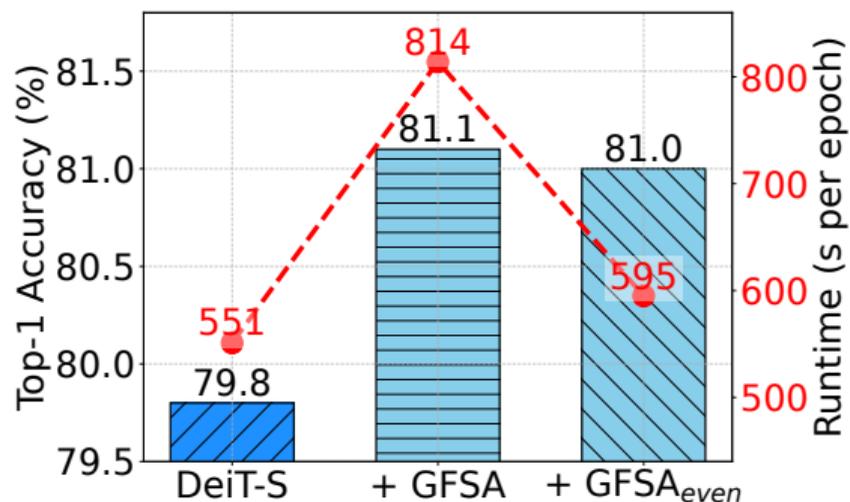


Figure: GFSA in selected layer. Effectiveness of our selective layer strategy on ImageNet-1k

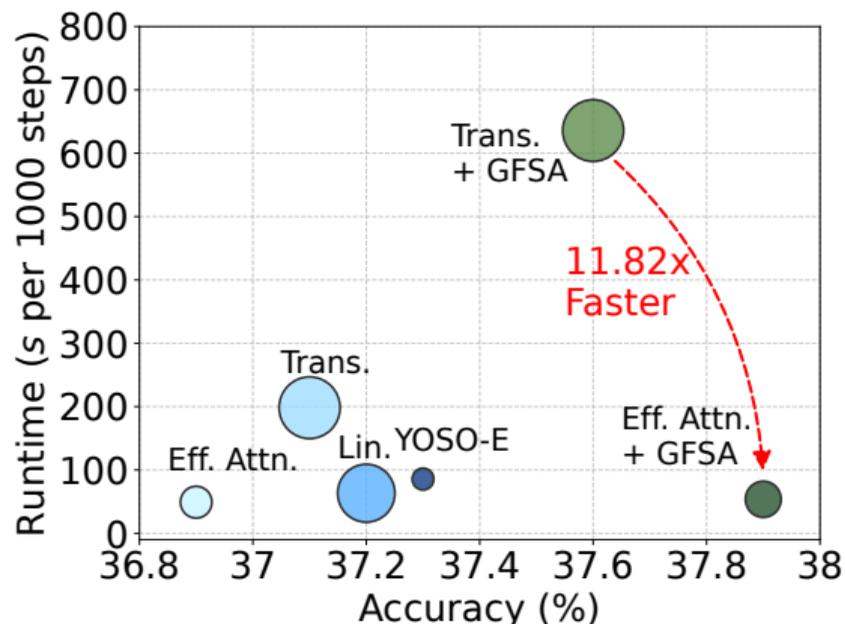
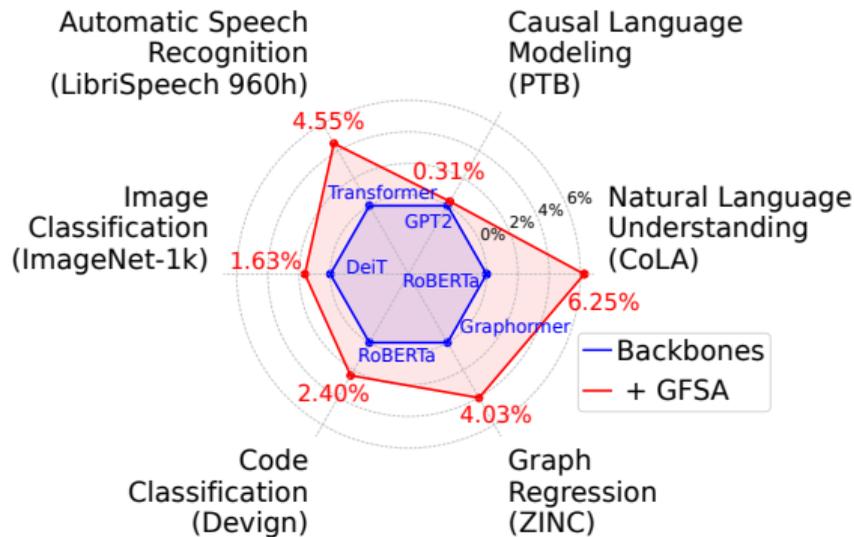
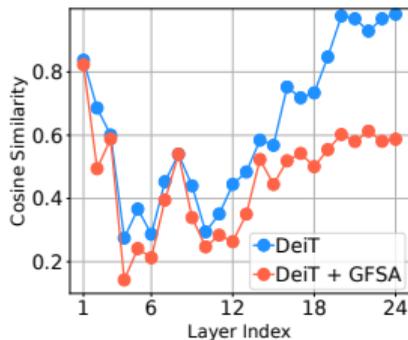
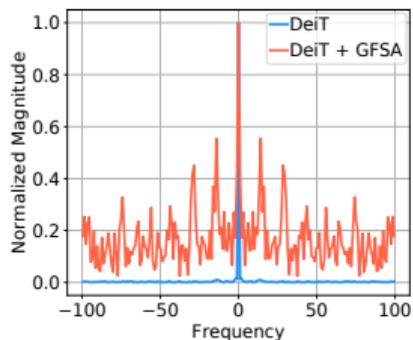


Figure: GFSA in linear Transformers. Performance, runtime, and GPU usage (circle sizes) of models on ListOps (2K) from Long Range Arena benchmark

Conclusion



- Considering the ongoing advancements in large language models, we hope that our approach may offer new insights for enhancing their performance and efficiency.

Thank You!

Email: jeongwhan.choi@yonsei.ac.kr

Paper



GitHub



KAIST

ASU



Penn

