

# Primphormer: Efficient Graph Transformers with Primal Representations

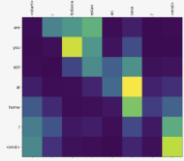
Mingzhen He · Ruihai Yang · Hanling Tian · Youmei Qiu · Xiaolin Huang\*



上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY

## Attention mechanism

$$\begin{cases} K(x_i, x_j) = \sigma(\langle q(x_i), k(x_j) \rangle) \\ o(x_i) = \sum_{j=1}^N v(x_j) K(x_i, x_j) \end{cases}$$



## Asymmetry

Where does asymmetry come from?

$$\begin{matrix} q(x_i) & k(x_j) \\ \langle W_q x_i, W_k x_j \rangle & \neq \langle W_q x_j, W_k x_i \rangle \end{matrix}$$

## Computational overhead

Large-scale problem (huge  $N$ )  $\xrightarrow{\mathcal{O}(N^2)}$  Sometime infeasible

## Primphormer

$$\begin{aligned} \min_{\mathbf{W}_e, \mathbf{W}_r} \quad & J = \frac{1}{2} \sum_i \mathbf{e}_i^T \Lambda \mathbf{e}_i + \frac{1}{2} \sum_j \mathbf{r}_j^T \Lambda \mathbf{r}_j - \text{Tr}(\mathbf{W}_e^T \mathbf{W}_r), \\ \text{s.t.} \quad & \mathbf{e}_i = \mathbf{f}_X \mathbf{W}_e \phi_q(x_i) \\ & \mathbf{r}_j = \mathbf{f}_X \mathbf{W}_r \phi_k(x_j). \end{aligned}$$

**Key idea:** Using primal-dual relationship to remodel the attention mechanism on graphs

Primal  $\begin{cases} \mathbf{e}(x) = \mathbf{f}_X \mathbf{W}_e \phi_q(x) \\ \mathbf{r}(x) = \mathbf{f}_X \mathbf{W}_e \phi_k(x). \end{cases}$

Dual  $\begin{cases} \mathbf{e}(x) = \sum_{j=1}^N \mathbf{F}_X \mathbf{h}_{r_j} \phi_k(x_j)^T \phi_q(x) := \sum_{j=1}^N \tilde{\mathbf{h}}_{r_j} K(x, x_j), \\ \mathbf{r}(x) = \sum_{i=1}^N \mathbf{F}_X \mathbf{h}_{e_i} \phi_q(x_i)^T \phi_k(x) := \sum_{i=1}^N \tilde{\mathbf{h}}_{e_i} K(x_i, x). \end{cases}$

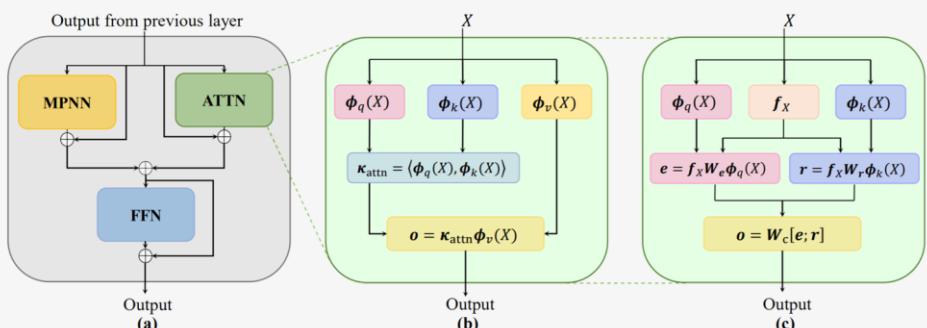


Figure 1 Illustrations of the architectures in one layer. (a) The GPS architecture. (b) The standard self-attention architecture. The attention score  $\kappa_{\text{attn}}$  involves pair-wise computations. (c) Primphormer eliminates the need for pair-wise computations by introducing the primal representation, resulting in a new computationally efficient GT.

## Expressivity and UAT

**Theorem 3.3.** For any continuous function  $f : [0, 1]^{d \times N} \rightarrow \mathbb{R}^{d \times N}$  and for each  $\epsilon > 0$ , there exists a Primphormer  $\mathcal{T}_{\text{PE}}$  with the positional encoding  $\mathbf{E}$  such that

$$\sup_{\mathbf{X} \in \mathcal{X}^N} \|f(\mathbf{X}) - \mathcal{T}_{\text{PE}}(\mathbf{X})\|_\infty < \epsilon. \quad (3.2)$$

**Theorem 3.4.** Let  $G = (V, E, \ell)$  be a labeled graph with  $N$  nodes, and node feature matrix  $\mathbf{X}^{(0)} := \mathbf{H} \in \mathbb{R}^{d \times N}$  consistent with the label  $\ell$ . Then, for all iterations  $t \geq 0$ , there exists a parameterization of Primphormer such that

$$C_t^1(v) = C_t^1(w) \iff \mathbf{X}^{(t)}(v) = \mathbf{X}^{(t)}(w), \quad (3.3)$$

for all nodes  $v, w \in V$ , where  $C_t^1 : V \rightarrow \mathbb{N}$  is the coloring function of the 1-WL test at  $t$ -th iteration.

**Corollary 3.5.** Primphormer is as powerful as Transformer in terms of distinguishing non-isomorphic graphs.

Table 3 Comparison of attentions in GPS. Best results are colored in first, second, third. OOM means out of memory.

MODEL	CIFAR10 GPS ACCURACY↑	MALNET-TINY ACCURACY↑	PASCALVOC-SP F1↑	PEPTIDES-FUNC AP↑	OGBN-PRODUCTS ACCURACY↑
MPNN-ONLY	$69.95 \pm 0.499$	$92.23 \pm 0.650$	$0.3016 \pm 0.0031$	$0.6159 \pm 0.0048$	$74.25 \pm 0.2148$
+TRANSFORMER	$72.31 \pm 0.344$	$93.50 \pm 0.410$	$0.3748 \pm 0.0109$	$0.6535 \pm 0.0041$	OOM
+BIGBIRD	$70.48 \pm 0.106$	$92.34 \pm 0.340$	$0.2762 \pm 0.0069$	$0.5854 \pm 0.0079$	$73.82 \pm 0.412$
+PERFORMER	$70.67 \pm 0.338$	$92.64 \pm 0.780$	$0.3724 \pm 0.0131$	$0.6475 \pm 0.0056$	$74.30 \pm 0.211$
+PRIM-ATTEN	$71.57 \pm 0.256$	$92.97 \pm 0.228$	$0.3173 \pm 0.0055$	$0.6447 \pm 0.0046$	$74.47 \pm 0.134$
+EXPHORMER	$74.69 \pm 0.125$	$94.02 \pm 0.209$	$0.3975 \pm 0.0037$	$0.6527 \pm 0.0043$	$74.67 \pm 0.179$
+PRIMPHORMER	$74.13 \pm 0.241$	$93.62 \pm 0.242$	$0.4602 \pm 0.0077$	$0.6612 \pm 0.0065$	$74.89 \pm 0.281$

Table 4 Efficiency comparisons on running time and peak memory consumption.

MODEL	TIME (S/EPOCH)					PEAK MEMORY USAGE (GB)				
	CIFAR.	MALNET.	PASCAL.	FUNC.	PROD.	CIFAR.	MALNET.	PASCAL.	FUNC.	PROD.
MPNN-ONLY	20.3	24.5	15.7	4.8	21.1	2.31	1.92	4.18	2.45	11.97
+TRANSFORMER	28.0	232.4	35.6	12.8	-	3.81	35.32	7.82	8.46	OOM
+BIGBIRD	55.2	325.6	52.3	51.9	93.9	2.81	2.71	4.99	17.29	
+PERFORMER	50.8	73.5	49.7	21.7	22.7	10.5	11.59	6.14	7.71	16.14
+PRIM-ATTEN	<b>32.1</b>	<b>62.5</b>	<b>25.7</b>	<b>7.9</b>	<b>22.6</b>	<b>2.74</b>	<b>2.58</b>	<b>4.74</b>	<b>3.38</b>	<b>13.63</b>
+EXPHORMER	44.5	62.1	35.2	7.6	25.4	5.54	10.38	7.35	4.81	31.09
+PRIMPHORMER	<b>32.6</b>	<b>61.9</b>	<b>25.3</b>	<b>7.7</b>	<b>22.1</b>	<b>2.74</b>	<b>2.86</b>	<b>4.72</b>	<b>3.41</b>	<b>13.35</b>

Table 5 Results on the BREC benchmark. Basic, Regular, Extension, and CFI are subsets of the BREC benchmark. Experiments are averaged over 5 runs.

MODEL	PE	BAS.↑	REG.↑	EXT.↑	CFI↑	ALL↑
GRAPHORMER		16	12	41	10	79
APE-GT		50.6	31.3	62.4	1	145.3
3-WL	<b>60</b>	<b>50</b>	<b>100</b>	<b>60</b>	<b>270</b>	
TRANSFORMER		47.2	39	65.2	<b>3</b>	154.4
PRIM-ATTEN	LAP	12.8	19	13.6	0.6	46
PRIMPHORMER		<b>51.6</b>	<b>42</b>	<b>72.4</b>	<b>3</b>	<b>169</b>
TRANSFORMER		59.8	49.4	98.6	5.2	213
PRIM-ATTEN	SPE	46.4	49	73	3	171.4
PRIMPHORMER		<b>60</b>	<b>50</b>	<b>100</b>	<b>9.4</b>	<b>219.4</b>

## Takes away

- Primal representation can be an efficient module
- Primphormer preserves good theoretical property

