

# Primal-Attention: Self-attention through Asymmetric Kernel SVD in Primal Representation

Yingyi Chen\*, Qinghua Tao\*, Francesco Tonin, Johan A.K. Suykens

ESAT-STADIUS, KU Leuven, Belgium \*Equal contribution



European Research Council  
Established by the European Commission



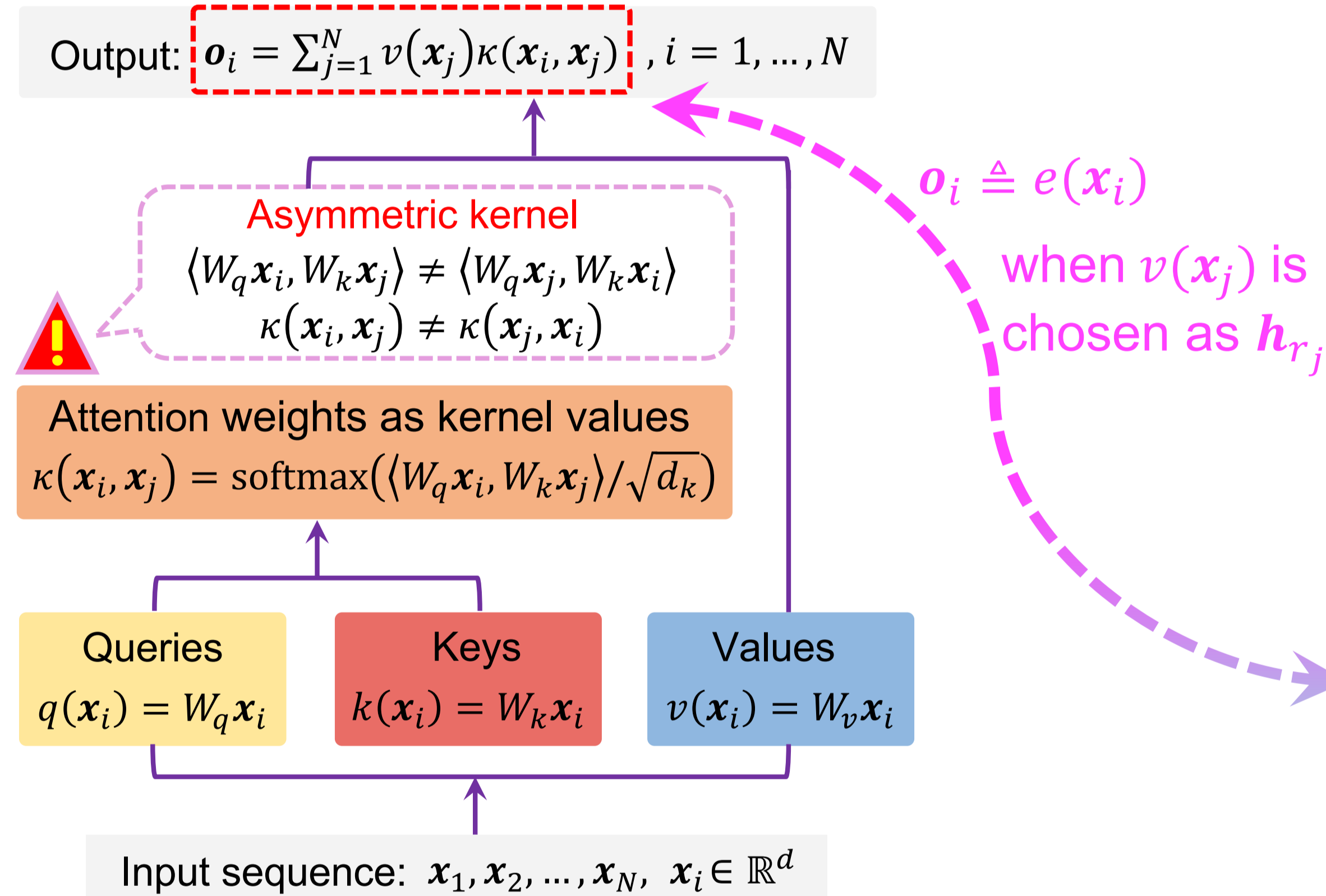
## Summary

Represent and optimize self-attention through *asymmetric Kernel Singular Value Decomposition* (KSVD):

- A **primal-dual representation** of self-attention in Transformers is formulated;
- A new attention mechanism named **Primal-Attention** based on primal representation of KSVD is proposed;
- A **KSVD optimization** designed for Primal-Attention is implemented.

## Canonical self-attention is with Asymmetric Kernel

- Attention matrix in self-attention is asymmetric



- For *asymmetric kernel*, the kernel trick from *Reproducing Kernel Banach Spaces* (RKBS) with  $\kappa(\cdot, \cdot): \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$  can be defined by the inner product of two feature maps from  $\mathcal{B}_X, \mathcal{B}_Z$ :

$$\kappa(x, z) = \langle \phi_x(x), \phi_z(z) \rangle, \forall x \in \mathcal{X}, \phi_x \in \mathcal{B}_X, z \in \mathcal{Z}, \phi_z \in \mathcal{B}_Z.$$

## SVD and Shifted Eigenvalue Problem

- SVD factorizes a given matrix  $A \in \mathbb{R}^{N \times M}$  by two sets of orthonormal eigenbases:

$$A = U \Sigma V^T, \Sigma = \text{diag}\{\sigma_1, \dots, \sigma_s\}, U \in \mathbb{R}^{N \times s}, V \in \mathbb{R}^{M \times s}$$

left singular vectors    right singular vectors

- Decomposition theorem (Lanczos, 1958): Any non-zero matrix  $A \in \mathbb{R}^{N \times M}$  can be written as  $A = \tilde{U} \tilde{\Sigma} \tilde{V}^T$ , where  $\tilde{U} \in \mathbb{R}^{N \times s}$ ,  $\tilde{V} \in \mathbb{R}^{M \times s}$ ,  $\tilde{\Sigma} \in \mathbb{R}^{s \times s}$  are defined by the *shifted eigenvalue problem*:

$$A \tilde{V} = \tilde{U} \tilde{\Sigma},$$

$$A^T \tilde{U} = \tilde{V} \tilde{\Sigma},$$

$$\tilde{U}^T \tilde{U} = I_s, \tilde{V}^T \tilde{V} = I_s, \tilde{\Sigma} \text{ is diagonal with positive numbers.}$$

## Primal-dual Representation of Self-attention with KSVD

- Primal problem** with KSVD for self-attention: we extend SVD under *Least Squares Support Vector Machines* (Suykens et al., 2002) framework (Suykens, 2016) to a nonlinear version

$$\max_{W_e, W_r, e_i, r_j} J = \frac{1}{2} \sum_{i=1}^N e_i^T \Lambda e_i + \frac{1}{2} \sum_{j=1}^N r_j^T \Lambda r_j - \text{Tr}(W_e^T W_r)$$

s.t.  $e_i = (f(X)^T W_e)^T \phi_q(x_i), i = 1, \dots, N,$   
 $r_j = (f(X)^T W_r)^T \phi_k(x_j), j = 1, \dots, N,$

Data-dependent projection weights	Feature maps related to queries and keys	Asymmetric attention kernel / Regu. coeff.	Projection scores w.r.t. queries, keys
$f(X)^T W_e =: W_{e X} \in \mathbb{R}^{p \times s}$	$\phi_q(\cdot): \mathbb{R}^d \rightarrow \mathbb{R}^p$	$K := \{(\phi_q(x_i), \phi_k(x_j))\}$	$e_i := W_{e X}^T \phi_q(x_i)$
$f(X)^T W_r =: W_{r X} \in \mathbb{R}^{p \times s}$	$\phi_k(\cdot): \mathbb{R}^d \rightarrow \mathbb{R}^p$	$\Lambda \in \mathbb{R}^{s \times s} > 0$ diagonal	$r_j := W_{r X}^T \phi_k(x_j)$

- Dual problem** with KSVD for self-attention: with Lagrangian duality and KKT conditions, the dual problem of above leads to the shifted eigenvalue problem

$$K H_r = H_e \Sigma,$$

$$K^T H_e = H_r \Sigma,$$

$H_e = [h_{e_1}, \dots, h_{e_N}]^T \in \mathbb{R}^{N \times s}$ ,  $H_r = [h_{r_1}, \dots, h_{r_N}]^T \in \mathbb{R}^{N \times s}$  are dual variables serving as *left* and *right singular vectors*.

- Primal-dual representation** of KSVD in self-attention:

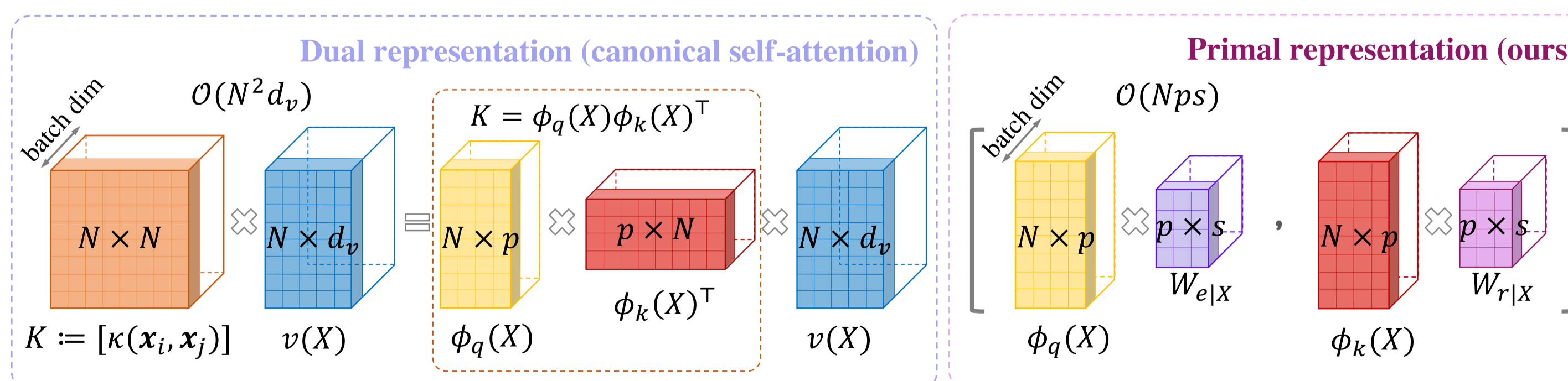
$$\text{Primal: } \begin{cases} e(x) = W_{e|X}^T \phi_q(x) \\ r(x) = W_{r|X}^T \phi_k(x) \end{cases}, \text{ Dual: } \begin{cases} e(x) = \sum_{j=1}^N h_{r_j} \kappa(x, x_j) \\ r(x) = \sum_{i=1}^N h_{e_i} \kappa(x_i, x) \end{cases}$$

## Primal-Attention

- Modelling**: leveraging primal representation of KSVD with  $\phi_q, \phi_k$

$$o_i := [e_i; r_i] = [W_{e|X}^T \phi_q(x_i); W_{r|X}^T \phi_k(x_i)],$$

experimentally,  $\phi_q(x) := q(x) / \|q(x)\|_2, \phi_k(x) := k(x) / \|k(x)\|_2$ , reducing time complexity from  $\mathcal{O}(N^2 d_v)$  to  $\mathcal{O}(Nps)$ .



- Optimization**: stationary solutions of KSVD for each head can be obtained by a zero-value of the primal objective  $J = 0$

$$J(W_e, W_r, \Lambda) = \frac{1}{2} \sum_{i=1}^N \left\| (W_{e|X} \Lambda^{1/2})^T \phi_q(x_i) \right\|_2^2 + \frac{1}{2} \sum_{j=1}^N \left\| (W_{r|X} \Lambda^{1/2})^T \phi_k(x_j) \right\|_2^2 - \text{Tr}(W_e^T W_r).$$

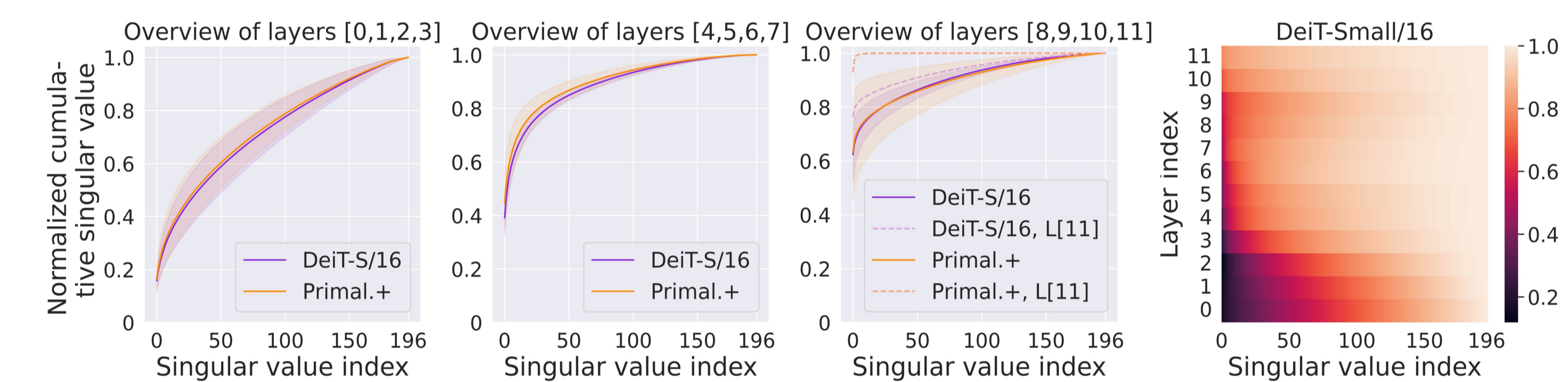
The *objective* of Primal-Attention is

$$\min L_{CE} + \eta \sum_l J_l^2,$$

$\sum_l J_l$  denotes additive objectives  $J_l$  of all Primal-attention blocks, where  $J_l$  is implemented as the mean of all heads.

## Numerical Experiments

- Enhanced low-rank property**: spectrum analysis of the self-attention matrix on ImageNet-1K



- Enhanced accuracy & efficiency**: Long-Range Arena Benchmark

Dataset (seq. length)	Transformer	Re-former	Per-former	Lin-former	Nyström-former	Long-former	YOSO-E	Primal.	Primal.+Trans.
ListOps (2K)	37.1	19.1	18.8	37.3	37.2	37.2	37.3	37.3	37.3
Text (4K)	65.0	64.9	63.8	55.9	65.5	64.6	64.7	61.2	65.4
Retrieval (4K)	79.4	78.6	78.6	79.4	79.6	81.0	81.2	77.8	81.0
Image (1K)	38.2	43.3	37.1	37.8	41.6	39.1	39.8	43.0	43.9
Pathfinder (1K)	74.2	69.4	69.9	67.6	70.9	73.0	72.9	68.3	74.3
Avg. Acc. (%)	58.8	55.1	53.6	55.6	59.0	59.0	59.2	57.5	60.4

Model	Time (s/1K-steps)					Memory (GB)				
	ListOps	Text	Retrieval	Image	Pathfinder	ListOps	Text	Retrieval	Image	Pathfinder
Transformer	194.5 (1x)	694.8 (1x)	1333.7 (1x)	334.5 (1x)	405.5 (1x)	5.50 (1x)	21.24 (1x)	18.72 (1x)	5.88 (1x)	5.88 (1x)
Nyströmformer	68.4 (2.8x)	120.9 (5.7x)	235.5 (5.7x)	179.5 (1.9x)	221.2 (1.8x)	0.89 (6.2x)	12.6x (12.6x)	3.29 (5.7x)	1.93 (3.0x)	1.93 (3.0x)
Linformer	63.4 (3.1x)	116.5 (6.0x)	226.2 (5.9x)	158.5 (2.1x)	204.0 (2.0x)	1.73 (3.2x)	3.45 (6.2x)	6.33 (3.0x)	3.45 (1.7x)	3.45 (1.7x)
Performer	83.8 (2.3x)	157.5 (4.4x)	320.6 (4.2x)	211.4 (1.6x)	278.1 (1.5x)	1.67 (3.3x)	3.34 (6.4x)	6.28 (3.0x)	3.34 (1.8x)	3.34 (1.8x)
Reformer	87.0 (2.2x)	168.5 (4.1x)	339.9 (3.9x)	223.7 (1.5x)	286.7 (1.4x)	1.64 (3.3x)	3.29 (6.5x)	6.09 (3.1x)	3.29 (1.8x)	3.29 (1.8x)
Primal.+Trans.	113.4 (1.7x)	367.6 (1.9x)	546.5 (2.4x)	212.1 (1.6x)	263.2 (1.5x)	5.24 (1.1x)	20.7 (1.0x)	18.59 (1.0x)	5.35 (1.1x)	5.35 (1.1x)
Primal.	56.5 (3.4x)	93.6 (7.4x)	185.3 (7.2x)	142.9 (2.3x)	180.0 (2.3x)	0.69 (7.9x)	1.37 (15.5x)	2.99 (6.3x)	1.39 (4.2x)	1.52 (3.9x)

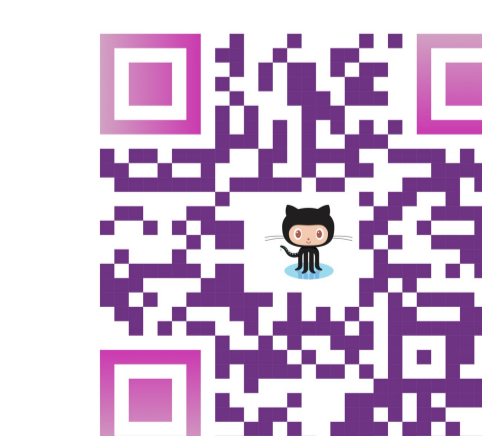
- Other benchmarks including *UEA time series classification*, *D4RL reinforcement learning*, *ImageNet-100*, *ImageNet-1K*, *WikiText-103* and more *ablation studies* can be found in the paper.

Paper:



arXiv:2305.19798

Code:



github.com/yingyiche n-cyy/PrimalAttention

References:

- Lanczos. "Linear systems in self-adjoint form." *The American Mathematical Monthly*, 1958.
- Suykens et al. *Least Squares Support Vector Machines*. World scientific, 2002.
- Suykens. "SVD revisited: A new variational principle, compatible feature maps and nonlinear extensions." *Applied and Computational Harmonic Analysis*, 2016.