



# LaplacianFormer: Rethinking Linear Attention with Laplacian Kernel

Zhe Feng<sup>12\*</sup> · Sen Lian<sup>3\*</sup> · Changwei Wang<sup>56†</sup> · Muyang Zhang<sup>21</sup> ·  
Tianlong Tan<sup>4</sup> · Rongtao Xu<sup>7</sup> · Weiliang Meng<sup>12†</sup> · Xiaopeng Zhang<sup>12</sup>

<sup>1</sup> MAIS, Institute of Automation, CAS

<sup>2</sup> University of Chinese Academy of Sciences

<sup>3</sup> China Electronics Data Corporation

<sup>4</sup> Institute of Computing Technology, CAS

<sup>5</sup> Qilu University of Technology

<sup>6</sup> Shandong Fundamental Research Center

<sup>7</sup> Spatiotemporal AI

\* Equal contribution

† Corresponding author





ICLR



中国科学院大学  
University of Chinese Academy of Sciences



中国科学院自动化研究所  
INSTITUTE OF AUTOMATION  
CHINESE ACADEMY OF SCIENCES



中国科学院计算技术研究所  
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

CEC

中国电子

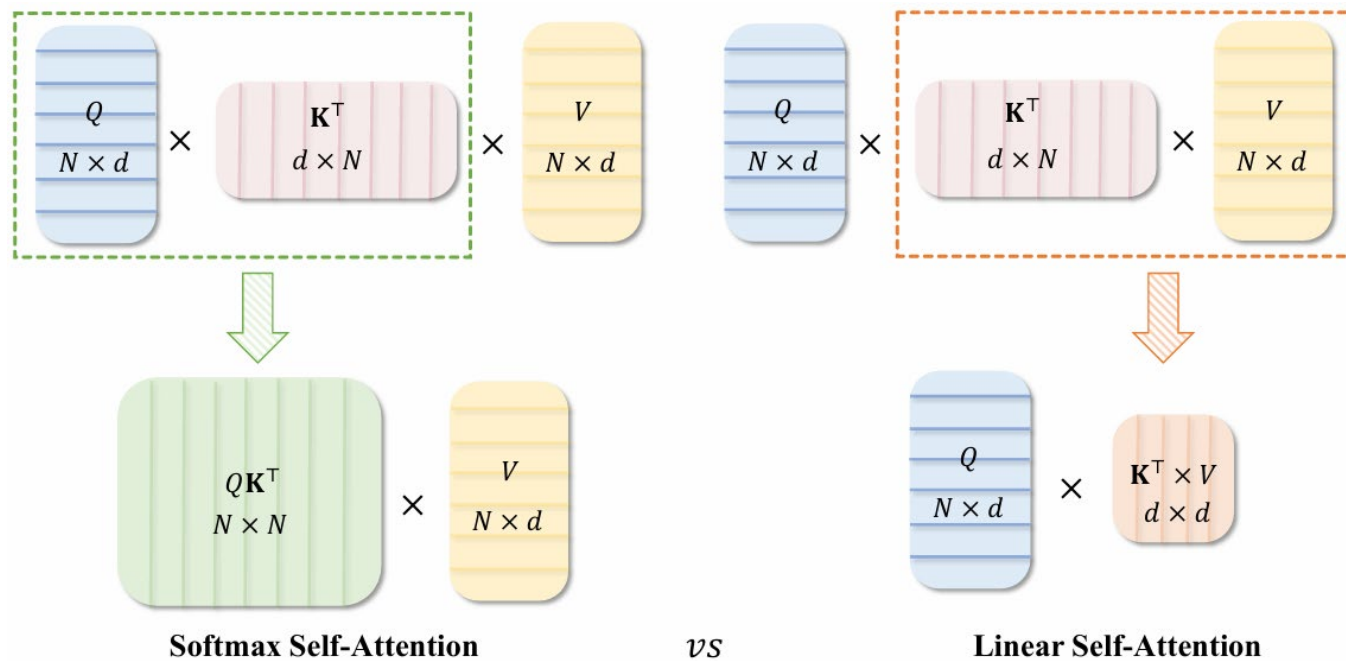
# Contents

- Introduction
- Motivation
- Problem Statement
- Methodology
- Results
- Conclusion



# Introduction

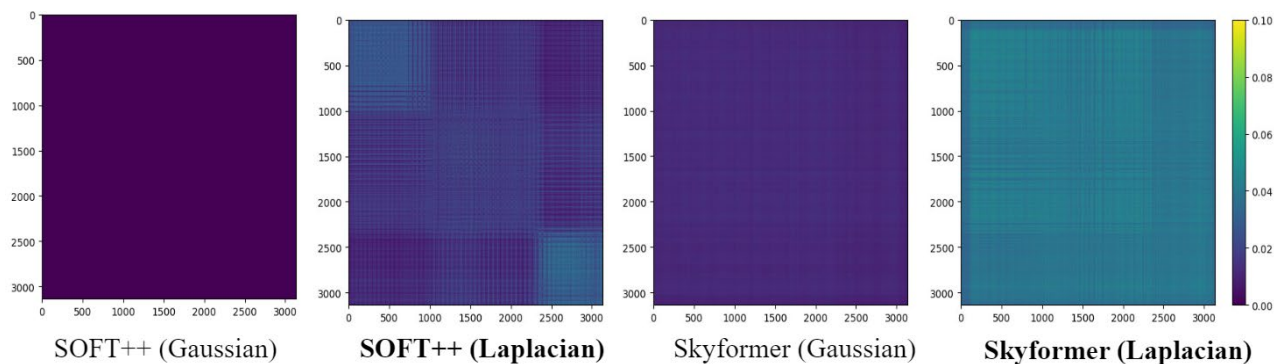
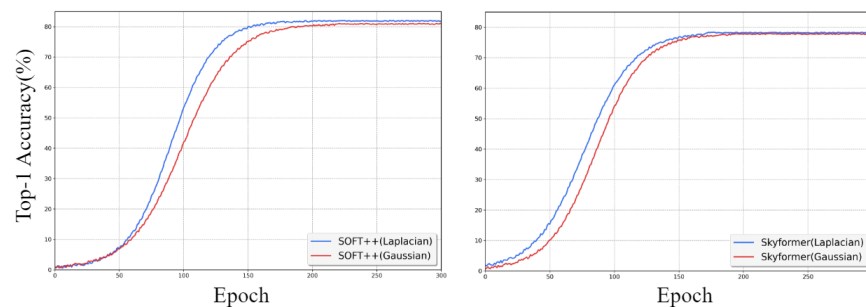
- Transformer is powerful for vision.
- Softmax attention has  $\mathcal{O}(N^2)$  complexity.
- Linear attention improves scalability.
- But kernel choice remains underexplored.





# Motivation

- Existing linear attention mostly uses **Gaussian kernels**.
- Empirical Q/K distances are **heavy-tailed**.
- Gaussian decays too aggressively and suppresses useful mid-range interactions.
- Laplacian kernel better matches sparse and heavy-tailed token relations.





ICLR



中国科学院大学  
University of Chinese Academy of Sciences



中国科学院自动化研究所  
INSTITUTE OF AUTOMATION  
CHINESE ACADEMY OF SCIENCES



中国科学院计算技术研究所  
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

CEC

中国电子

# Problem Statement

- How to design a better kernel for linear attention?
- How to preserve expressiveness under **low-rank approximation**?
- How to implement it **efficiently on GPU**?



# Methodology

- **Laplacian** kernel attention
- Injective normalized feature map
- **Nyström approximation** for low-rank efficiency
- **Newton–Schulz** iteration for fast inverse approximation
- CUDA implementation for practical acceleration

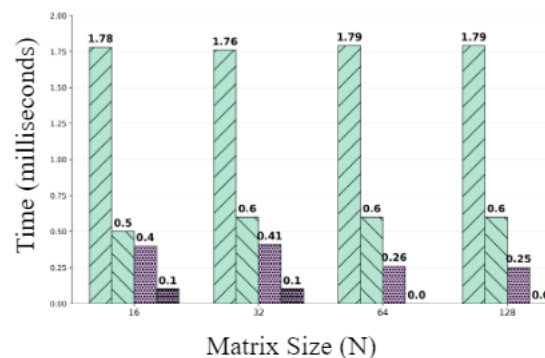
## Algorithm 1 Laplacian Kernel with Nyström Approximation

- 1: **Input:** Queries  $Q \in \mathbb{R}^{N \times d}$ , Keys  $K \in \mathbb{R}^{N \times d}$ , Nyström sampling function  $f_s$
- 2: **Sampling:**  $\tilde{Q}, \tilde{K} \leftarrow f_s(Q), f_s(K)$  ▷ Select  $m \ll n$  landmark points
- 3:  $W \leftarrow \exp\left(-\frac{1}{\lambda} \|\tilde{Q} \circ \tilde{K}\|_1\right)$  ▷ Kernel matrix on sampled points
- 4:  $C \leftarrow \exp\left(-\frac{1}{\lambda} \|Q \circ \tilde{K}\|_1\right)$  ▷ Cross-kernel between all queries and landmarks
- 5:  $\hat{G} \leftarrow CW^\dagger C^\top$  ▷ Low-rank approximation of full kernel matrix
- 6: **Output:**  $\hat{G}$

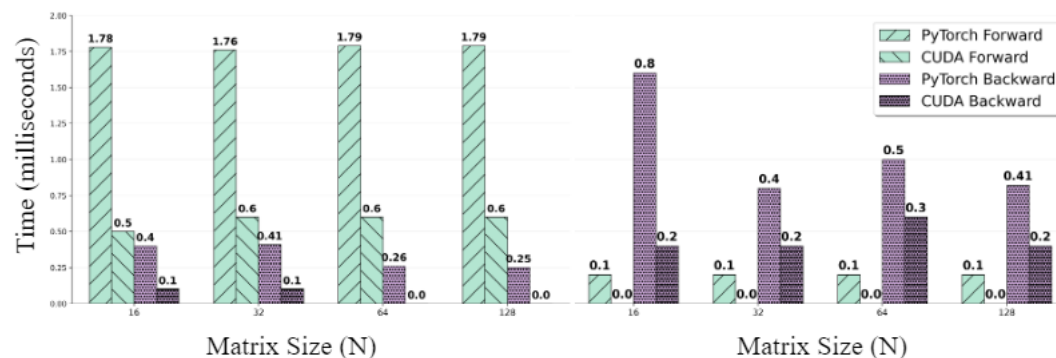
## Algorithm 2 Newton–Schulz Iteration for Approximating $W^\dagger$

- 1: **Input:** Landmark kernel matrix  $W \in \mathbb{R}^{m \times m}$ , number of iterations  $\mathcal{T} \in \mathbb{Z}^+$
- 2: Add small perturbation:  $W \leftarrow W + \epsilon I$ , where  $\epsilon > 0$
- 3: Initialize scaling factor:  $\alpha \leftarrow \frac{2}{\|W\|_2}$
- 4: Initialize:  $X_0 \leftarrow \alpha W^\top$
- 5: **for**  $k = 1$  to  $\mathcal{T}$  **do**
- 6:      $X_k \leftarrow X_{k-1}(2I - WX_{k-1})$
- 7: **end for**
- 8: **Output:** Approximate pseudoinverse  $X_{\mathcal{T}} \approx W^\dagger$

### Newton-Schulz Kernel Execution Time



### Laplace Kernel Execution Time





# Results

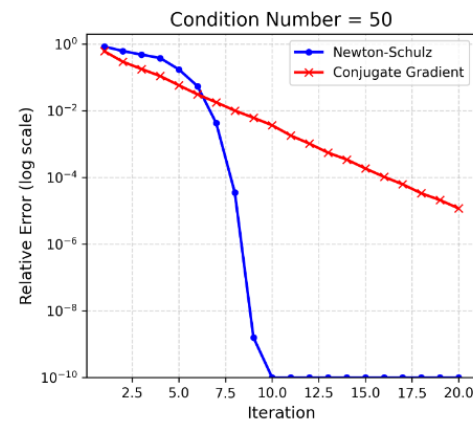
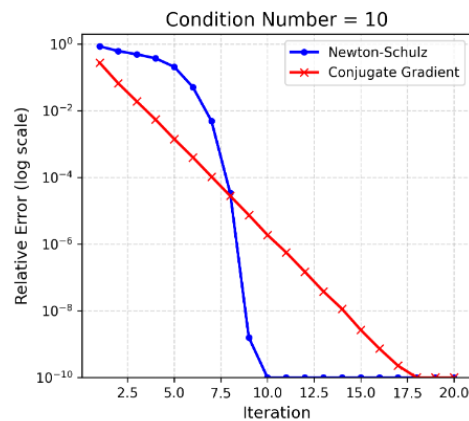
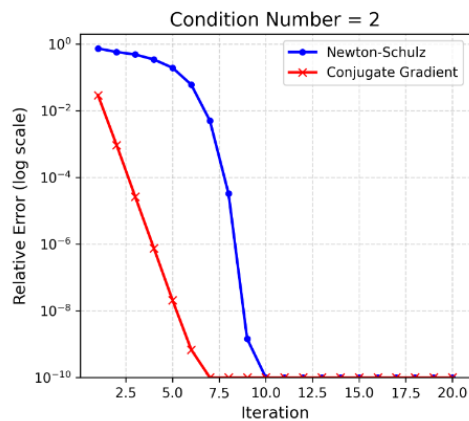
- Better accuracy-efficiency trade-off on ImageNet
- Linear memory scaling with sequence length
- Faster convergence than Gaussian-kernel counterparts
- Competitive downstream performance on detection/segmentation

FLOPs range	Model	Params	FLOPs	Top-1 %↑	Image Size
< 3G	Agent-DeiT-T Han et al. (2024c)	6.0M	1.2G	74.9	224
	VRWKV-T Duan et al. (2025)	6.2M	1.2G	75.1	256
	PVT-T-PolaFormer Meng et al. (2025)	12M	2.0G	78.8	224
	FL-PVTv2-B1 Han et al. (2023)	13M	2.2G	79.5	224
	BiFormer-T Zhu et al. (2023)	13.1M	2.2G	81.4	224
	<b>LaplacianFormer-Tiny</b>	<b>12.1M</b>	<b>2.1G</b>	<b>81.4</b>	<b>224</b>
3~8G	InLine-CSwin-S Han et al. (2024a)	33M	6.8G	83.8	224
	SViT-S Huang et al. (2023)	25M	4.4G	83.6	224
	BiFormer-S Zhu et al. (2023)	25.5M	4.5G	83.8	224
	HiViT-T Zhang et al. (2023b)	19M	4.6G	82.1	224
	Agent-PVT-S Han et al. (2024c)	20.6M	4.0G	82.2	224
	<b>LaplacianFormer-Small</b>	<b>25.7M</b>	<b>4.8G</b>	<b>83.8</b>	<b>224</b>
8~10G	SViT-B Huang et al. (2023)	52M	9.9G	84.8	224
	SOFT++-Medium Lu et al. (2024)	45M	7.2G	83.7	224
	BiFormer-B Zhu et al. (2023)	56.8M	9.8G	84.3	224
	Swin-S-PolaFormer Meng et al. (2025)	50M	8.7G	83.6	224
	SLAB-Swin-S Guo et al. (2024)	-	8.7G	81.8	224
	<b>LaplacianFormer-Medium</b>	<b>46.3M</b>	<b>7.43G</b>	<b>85.3</b>	<b>224</b>
10~14G	StructViT-B-8-1 Kim et al. (2024)	52M	12G	84.3	224
	SOFT++-Large Lu et al. (2024)	64M	11G	84.1	224
	NAT-B Hassani et al. (2023)	90M	13.7G	84.3	224
	MogaNet-L Li et al. (2024)	82.5M	15.9G	84.7	224
	FLatten-CSwin-S Han et al. (2023)	35M	6.9G	83.6	224
	<b>LaplacianFormer-Large</b>	<b>63.1M</b>	<b>11.2G</b>	<b>85.6</b>	<b>224</b>
> 14G	VRWKV-B Duan et al. (2025)	93.7M	18.2G	82.0	224
	SViT-L Huang et al. (2023)	95M	15.6G	85.3	224
	MLLA-B Han et al. (2024b)	96M	16.2G	85.3	224
	HiViT-B Zhang et al. (2023b)	66M	15.9G	83.8	224
	<b>LaplacianFormer-Huge</b>	<b>78.5M</b>	<b>15.5G</b>	<b>85.8</b>	<b>224</b>

Backbone	Mask R-CNN 1x						RetinaNet 1x					
	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP^m$	$AP_{50}^m$	$AP_{75}^m$	$AP^b$	$AP_{50}^b$	$AP_{75}^b$	$AP_S^b$	$AP_M^b$	$AP_L^b$
Swin-T-PRRepBN Guo et al. (2024)	42.9	65.8	46.8	39.3	62.6	41.9	-	-	-	-	-	-
FL-PVT-T Han et al. (2023)	38.2	61.6	41.9	37.0	57.6	39.0	-	-	-	-	-	-
SOFT++-Tiny Lu et al. (2024)	41.2	63.7	44.7	38.2	61.0	41.0	41.9	62.7	44.7	27.8	45.4	55.6
<b>LaplacianFormer-Tiny</b>	<b>43.2</b>	<b>66.1</b>	<b>47.2</b>	<b>40.3</b>	<b>63.0</b>	<b>42.9</b>	<b>42.5</b>	<b>64.1</b>	<b>46.4</b>	<b>29.1</b>	<b>46.9</b>	<b>57.8</b>
PVT-S-PolaFormer Meng et al. (2025)	43.9	66.1	47.9	40.2	63.1	43.0	43.2	64.1	46.4	-	-	-
InLine-PVT-S Han et al. (2024a)	43.4	66.4	47.1	40.1	63.1	43.3	-	-	-	-	-	-
SOFT++-Small Lu et al. (2024)	43.8	66.0	47.5	40.1	63.0	43.0	43.7	64.9	46.8	28.7	47.4	57.6
<b>LaplacianFormer-Small</b>	<b>45.8</b>	<b>68.2</b>	<b>49.8</b>	<b>42.0</b>	<b>65.1</b>	<b>45.2</b>	<b>45.5</b>	<b>66.8</b>	<b>49.1</b>	<b>30.7</b>	<b>51.8</b>	<b>59.5</b>
Agent-PVT-M Han et al. (2024c)	45.9	67.8	50.4	42.0	65.0	45.4	-	-	-	-	-	-
FL-Swin-M Han et al. (2023)	44.0	66.4	48.0	40.3	63.4	43.5	-	-	-	-	-	-
SOFT++-Medium Lu et al. (2024)	46.6	67.8	51.2	42.0	64.8	45.2	44.3	64.7	47.4	29.0	48.2	59.9
<b>LaplacianFormer-Medium</b>	<b>48.0</b>	<b>70.3</b>	<b>52.5</b>	<b>43.5</b>	<b>65.8</b>	<b>46.5</b>	<b>47.2</b>	<b>68.5</b>	<b>51.5</b>	<b>31.8</b>	<b>53.0</b>	<b>61.4</b>
Swin-T-PolaFormer Meng et al. (2025)	44.8	67.6	49.1	40.5	64.1	43.5	-	-	-	-	-	-
Agent-PVT-L Han et al. (2024c)	46.9	69.2	51.4	42.8	66.2	46.2	-	-	-	-	-	-
SOFT++-Large Lu et al. (2024)	47.0	68.3	51.7	42.2	65.2	45.4	47.0	67.8	50.4	30.2	50.9	62.0
<b>LaplacianFormer-Large</b>	<b>48.2</b>	<b>70.5</b>	<b>53.0</b>	<b>43.8</b>	<b>67.1</b>	<b>47.4</b>	<b>48.5</b>	<b>69.3</b>	<b>52.4</b>	<b>32.6</b>	<b>52.3</b>	<b>63.8</b>



# Results



Model	CG (%)	NS (%)
LaplacianFormer-Tiny	79.2	<b>81.4</b>
LaplacianFormer-Small	81.4	<b>83.8</b>

$\lambda$	0.5	1	2	4	8
Top-1 Acc (%) $\uparrow$	79.4	79.6	80.1	<b>81.4</b>	78.5



ICLR



中国科学院大学  
University of Chinese Academy of Sciences



中国科学院自动化研究所  
INSTITUTE OF AUTOMATION  
CHINESE ACADEMY OF SCIENCES



中国科学院计算技术研究所  
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

CEC

中国电子

# Conclusion

- Laplacian is a **better** kernel choice than Gaussian for linear attention
- **Injective representation** improves expressiveness
- Efficient approximation and **CUDA implementation** make deployment practical
- Limitations & Future Works:
  - ❑ Developing learnable or adaptive kernel functions.
  - ❑ Exploring task-specific kernel parameterization.



ICLR



中国科学院大学  
University of Chinese Academy of Sciences



中国科学院自动化研究所  
INSTITUTE OF AUTOMATION  
CHINESE ACADEMY OF SCIENCES



中国科学院计算技术研究所  
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

CEC

中国电子

# Thank you!

*The Fourteenth International Conference on Learning Representations  
Rio de Janeiro, Brazil*