



# Skyformer: Remodel Self-Attention with Gaussian Kernel and Nyström Method

Yifan Chen\*, Qi Zeng\*, Heng Ji, Yun Yang

University of Illinois at Urbana-Champaign

October 18, 2021

1 Background

2 Method

3 Results

4 Conclusion

# Background

## Transformers are expensive to train

### Quadratic Complexity

- Quadratic time and space complexity in the self-attention mechanism
- Transformers cannot support long sequence processing and large batch size

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{p}}\right) \mathbf{V} := \mathbf{D}^{-1}\mathbf{A}\mathbf{V} \quad (1)$$

### Training Instability

- Small perturbations in parameter updates tend to be amplified, resulting in significant disturbances in the model output

# Background

Kernel methods may be the answer to both challenges

Connections between Self-attention and Gaussian Kernels:

- The un-normalized attention score matrix can be formed via basic matrix operations on an empirical Gaussian kernel matrix
- The form of Gaussian kernels has the natural interpretation of assigning “attention” to different tokens
- Gaussian kernels automatically perform the normalization similar to how softmax does

# Method

## Kernelized Attention

- Kernelized Attention replaces the softmax structure with a Gaussian kernel

$$\text{Kernelized-Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{C}\mathbf{V} := \kappa\left(\frac{\mathbf{Q}}{p^{1/4}}, \frac{\mathbf{K}}{p^{1/4}}\right) \mathbf{V} \quad (2)$$

- It can be rewritten in terms of the un-normalized attention score matrix as

$$\text{Kernelized-Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{D}_Q^{-1/2} \cdot \mathbf{A} \cdot \mathbf{D}_K^{-1/2} \quad (3)$$

# Method

## Skyformer: a modified Nyström method

- Complete the matrix into a PSD matrix  $\bar{\mathbf{B}}$

$$\bar{\mathbf{B}} := \phi \left( \begin{pmatrix} \mathbf{Q} \\ \mathbf{K} \end{pmatrix}, \begin{pmatrix} \mathbf{Q} \\ \mathbf{K} \end{pmatrix} \right), \quad (\mathbf{B} := (\mathbf{I}, \mathbf{0}) \bar{\mathbf{B}} (\mathbf{0}, \mathbf{I})^T) \quad (4)$$

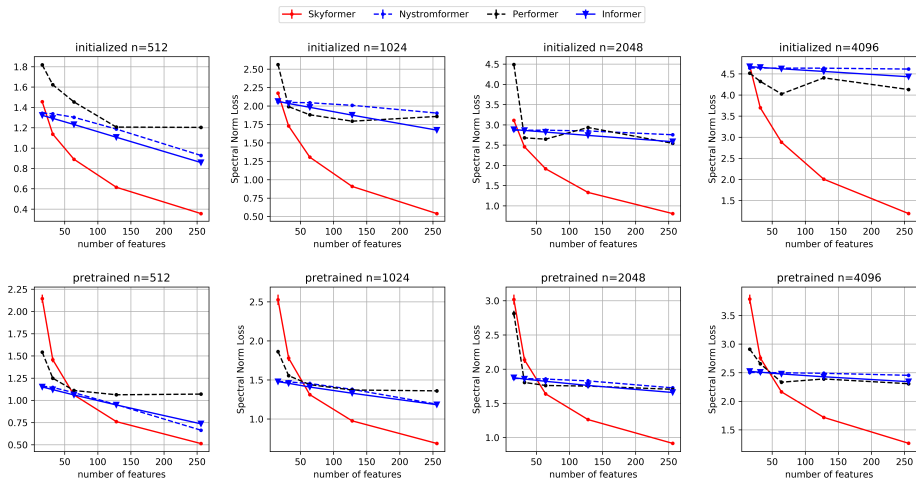
- Approximate  $\bar{\mathbf{B}}$  with  $\tilde{\tilde{\mathbf{B}}}$  through

$$\tilde{\tilde{\mathbf{B}}} = \bar{\mathbf{B}} \mathbf{S} (\mathbf{S}^T \bar{\mathbf{B}} \mathbf{S})^\dagger \mathbf{S}^T \bar{\mathbf{B}} \quad (5)$$

- The final approximation will be given as

$$\tilde{\mathbf{B}} := (\mathbf{I}, \mathbf{0}) \tilde{\tilde{\mathbf{B}}} (\mathbf{0}, \mathbf{I})^T \quad (6)$$

# Results: Empirical Approximation Evaluation



## Results: Classification accuracy on LRA benchmark

Model	Text	ListOps	Retrieval	Pathfinder	Image	AVG.
Self-Attention	61.95	38.37	80.69	65.26	40.57	57.37
Kernelized Attention	60.22	38.78	81.77	70.73	41.29	58.56
Nystromformer	64.83	38.51	80.52	69.48	41.30	58.93
Linformer	58.93	37.45	78.19	60.93	37.96	54.69
Informer	62.64	32.53	77.57	57.83	38.10	53.73
Performer	64.19	38.02	80.04	66.30	41.43	58.00
Reformer	62.93	37.68	78.99	66.49	48.87	58.99
BigBird	63.86	39.25	80.28	68.72	43.16	59.05
<b>Skyformer</b>	64.70	38.69	82.06	70.73	40.77	<b>59.39</b>



# Conclusion

- We revisit the intrinsic connection between self-attention and kernel methods, and explore a new kernel-based structure to stabilize the training of Transformers
- We approximate the Kernelized Attention via low dimensional randomized sketches by adapting the Nyström method to a non-PSD matrix.
- We conduct extensive experiments showing that Skyformer achieves comparable performance to the original self-attention with fewer computational costs.