

## Background

Transformers are hard to train.

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

- **Quadratic** time and space complexity in the self-attention mechanism, which makes **long sequence** processing resource-consuming.
- Small perturbations in parameter updates tend to be amplified, resulting in significant disturbances in the model output.

**Kernel methods** may be the answer to both challenges. To justify the introduction of kernel methods, we first build up the 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.

More details are provided in the next block, in which we propose “Kernelized Attention”.

## Kernelized Attention

- Kernelized Attention replaces the softmax structure with a Gaussian kernel

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

where  $\kappa(\mathbf{q}_i, \mathbf{q}_j) := \exp(-\|\mathbf{q}_i - \mathbf{q}_j\|^2/2)$  is the standard Gaussian kernel function,

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

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

where  $\mathbf{D}_Q$  (resp.  $\mathbf{D}_K$ ) is a diagonal matrix with elements  $(\mathbf{D}_Q)_{ii} = \exp\left(\frac{\|\mathbf{q}_i\|^2}{\sqrt{p}}\right)$ ,  $\forall i \in [n]$ .

Comparing the two equations, we can check the un-normalized attention score matrix  $\mathbf{A}$  is involved, and  $\mathbf{D}_Q, \mathbf{D}_K$  perform the normalization similar to  $\mathbf{D}$  in self-attention.

## Attempts in past literature

The attention score matrix is known to exhibit a very fast rate of singular value decay, similar to that of an empirical kernel matrix. This near singular property motivates many low-rank attention approximation methods to skillfully leverage the computation techniques in kernel methods.

- Linformer [Wang et al., 2020] compresses the size of the key and value matrix with random projections based on the Johnson–Lindenstrauss transform;
- Reformer [Kitaev et al., 2020] applies locality-sensitive hashing (LSH) [Har-Peled et al., 2012] to simplify the computation of the attention score matrix, which is recently used in kernel density estimation [Charikar and Siminelakis, 2017, Backurs et al., 2019];
- Performer [Choromanski et al., 2020] projects both query and key matrix through random Fourier features [Rahimi et al., 2007], heavily exploiting Bochner Theorem for stationary kernels.

## Nyström methods

We write Nyström methods in a sketching form as follows:

- Approximate  $\mathbf{K}$  by  $\tilde{\mathbf{K}} = \mathbf{K} \mathbf{S} (\mathbf{S}^T \mathbf{K} \mathbf{S})^\dagger \mathbf{S}^T \mathbf{K}$ ;
- For classical Nyström methods,  $\mathbf{S}$  is an  $n$ -by- $d$  (uniform) sampling matrix;
- The complexity is reduced to  $\mathcal{O}(nd^2)$ , when  $d \ll n$ ;
- A fun fact: we can check Nyström methods by letting  $\mathbf{S} = \mathbf{I}$ .

An implicit requirement for the matrix  $\mathbf{K}$  to approximate is,  $\mathbf{K}$  should be positive semidefinite (PSD). However, in general  $\mathbf{A}$  is not PSD since  $\mathbf{A}$  is not even symmetric.

## Skyformer: a modified Nyström method

To resolve the issues mentioned above, we propose Skyformer:

- Complete the matrix into a PSD matrix  $\tilde{\mathbf{C}}$

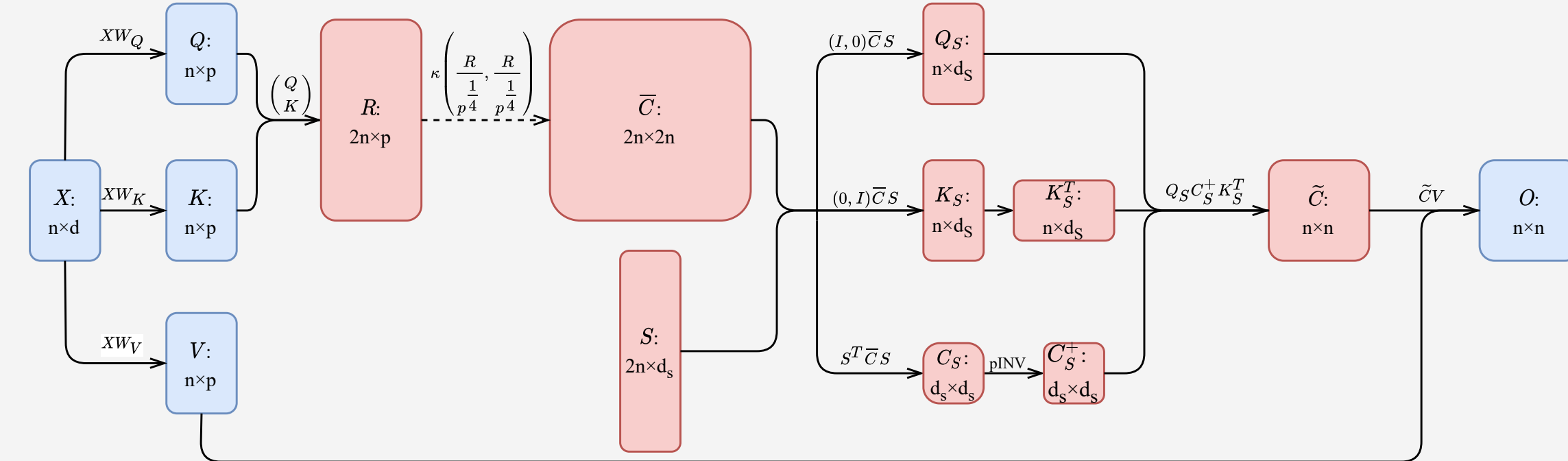
$$\tilde{\mathbf{C}} := \phi\left(\begin{pmatrix} \mathbf{Q} \\ \mathbf{K} \end{pmatrix}, \begin{pmatrix} \mathbf{Q} \\ \mathbf{K} \end{pmatrix}\right), \quad (\mathbf{C} := \phi(\mathbf{Q}, \mathbf{K}) = (\mathbf{I}, \mathbf{0}) \tilde{\mathbf{C}} (\mathbf{0}, \mathbf{I})^T) \quad (4)$$

- Approximate  $\tilde{\mathbf{C}}$  with  $\tilde{\tilde{\mathbf{C}}}$  through

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

- The final approximation will be given as

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



**Figure 1.** An overview of the proposed attention approximation method. This method approximates the original softmax matrix with a low dimensional randomized sketch of the complete symmetric matrix. Each box represents a matrix with the box size representing matrix size. The red boxes are matrices introduced by the proposed method besides the vanilla self-attention matrices in blue boxes. The dash line means that during computation, instead of the full  $\tilde{\mathbf{C}}$  we only compute the selected columns of  $\tilde{\mathbf{A}}$  with a sampling matrix  $\mathbf{S}$ .

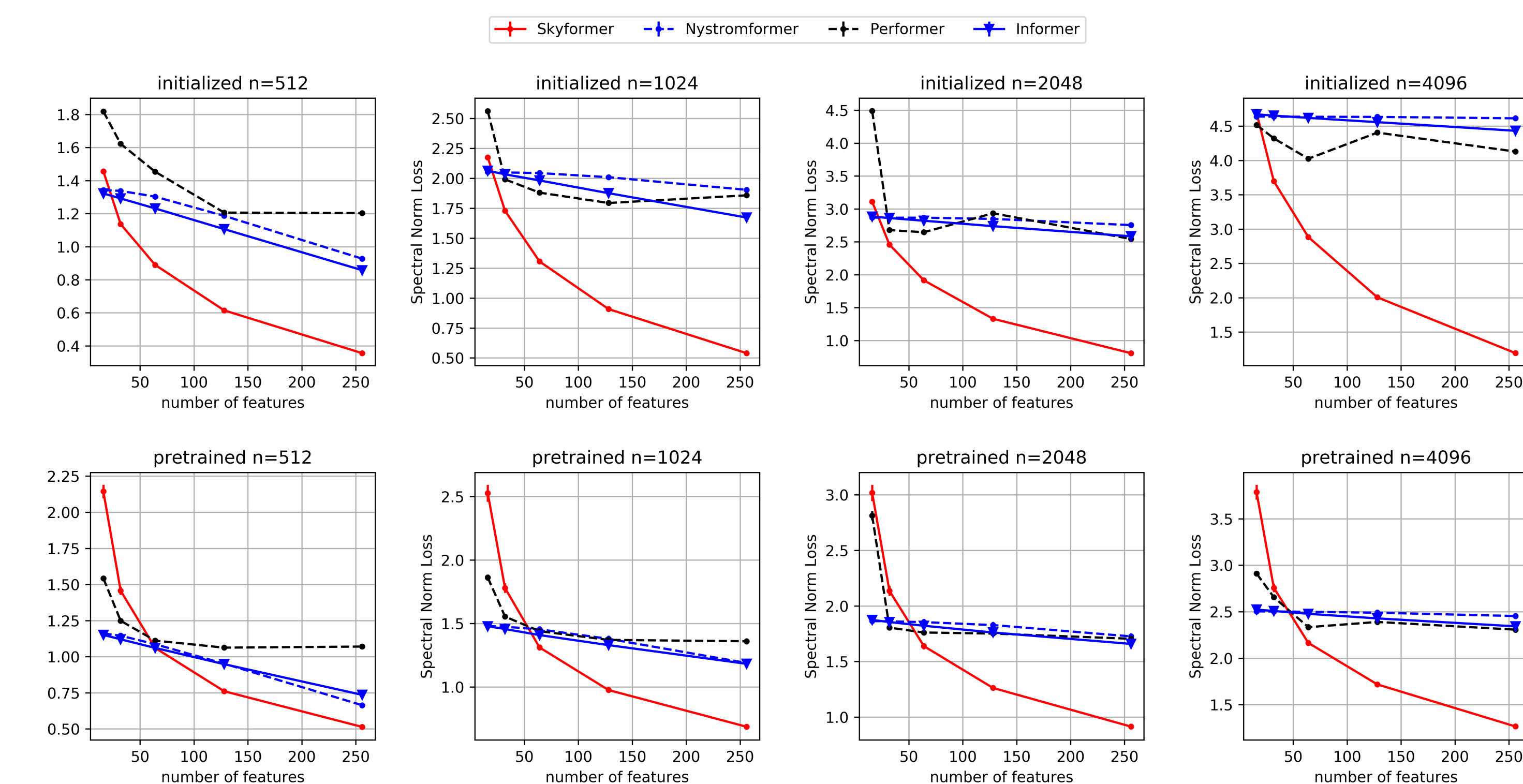
## Error analysis of Skyformer

An implicit advantage of using Kernelized Attention is that we can leverage the existing conclusions for kernel methods to analyze the theoretical properties of the model.

In the paper we state a high probability bound (Theorem 2) on the size  $d$  of the sub-sampling matrix used in Skyformer for the kernelized attention score matrix  $\mathbf{C}$  by the following theorem. This theorem implies the time and space complexity of our proposed approximation depends on the statistical dimension  $d_{stat} := \text{Tr}(\tilde{\mathbf{C}}(\tilde{\mathbf{C}} + \lambda \mathbf{I})^{-1})$ . If we directly use the conclusion from Gaussian kernels,  $d_{stat}$  should be  $\tilde{\mathcal{O}}(1)$  (complexity modulo poly-log term) [Yang et al., 2017] due to the exponential eigenvalue decay rate of Gaussian kernels, which is comparable to the complexity of most other efficient attention.

## Empirical results

### Empirical Approximation Evaluation



**Figure 2.** Spectral norm results with different sequence lengths under different  $W_Q, W_K, W_V$  settings, either from initialized or pretrained BERT models. All methods are approximating the original self-attention output. Y axis: Lower spectral loss means better approximation. X axis: Higher  $d$  (number of features) means visiting more elements in the original matrix and bringing more computation costs. The label “Skyformer” here means that we use the algorithm behind Skyformer, to approximate the raw attention score matrix  $\mathbf{A}$  in self-attention. In this experiment, “Skyformer” also needs to first approximate  $\mathbf{A}$ , and then approximate  $\mathbf{D}$ , as Performer does.

### 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>	<b>64.70</b>	<b>38.69</b>	<b>82.06</b>	<b>70.73</b>	<b>40.77</b>	<b>59.39</b>

## References and supplementary materials

Omitted due to space limit. Please refer to <https://openreview.net/forum?id=pZCYG7gjkKz>.

