

# Scaleformer : a scalable transformer with linear complexity and relative positional encoding

Benoit Favier<sup>1</sup> and Walter Dal’Maz Silva<sup>2</sup>

<sup>1</sup>Phealing

<sup>2</sup>Unaffiliated

October 9, 2021

## Abstract

To overcome the quadratic complexity of the vanilla transformer with sequence length, some previous works proposed a kernelized attention mechanism, which has a linear complexity. Other works proposed the change the way the position of each token is encoded so that the model depends on relative distance between tokens instead of absolute position. In this work we propose a model that combines kernelized attention with relative positional encoding while still scaling linearly in complexity. Furthermore we propose to dynamically chose the most favourable order of commutative operations. While giving the same results, different operation orders have memory usage and computational cost that depends differently on embedding dimension and sequence length.

## 1 Introduction

Vaswani et al. [11] introduced the transformers, a simpler network to model sequence-to-sequence translations tackling new heights of problem complexity. The improvements were due to the non sequential nature of its training

process which allowed better parallelization and thus fast training of big models, and due to improved long term dependencies thanks to the intrinsic good gradient flows of the attention mechanism. However the original transformer presented some drawbacks. It has a quadratic complexity with sequences length, what represents a major drawback to its scalability. The positional encoding induce a different initial embedding for the same group of words at different position in the sentence, which is detrimental for generalization. The absolute positional encoding is also detrimental for the extrapolation to sequences longer than training sequences - even for the initially proposed sinusoidal positional encoding.

Ever since, significant research effort have been made to alleviate these problems. Several strategies have been proposed to reduce the quadratic complexity with sequence length, and the relative positional encoding further improved the performances of the model. These improvements have opened the path to the application of transformers to image analysis where scalability and positional invariance are of utmost importance. However most proposed architecture we are aware of only alleviated parts of the issues, by being incompatible with relative positional encoding, having no scheme for masked attention - thus being restrained to encoder-only models, or being complex to use/implement - needing custom operations programmed in CUDA or introducing stochastic methods.

In the present work we propose a complete encoder-decoder transformer model that scales linearly with large sequence lengths, by putting together ideas developped across several works [2, 4, 6, 9–11]. It remains compatible with relative positional encoding, and can be implemented with usual functions of neural network frameworks without requiring custom CUDA code.

## 2 Background

### 2.1 The original multi-head attention mechanism

The scaled dot-product attention proposed by Vaswani et al. [11] transforms the vectorial embedding  $\vec{Y}_i$  of a token as a function of a sequence of other embedding  $\vec{X}_j$ . Where  $\vec{Y}_i$  and  $\vec{X}_j$  are all vectors of size  $D$ . A key  $\vec{K}_j$  and value  $\vec{V}_j$  are attributed to each vector  $\vec{X}_j$ , and query  $\vec{Q}_i$  is attributed to  $\vec{Y}_i$ . The vectors are obtained by linear projection from dimension  $D$  to  $d$  using three matrices of learnable parameters. The new vector  $\vec{y}_i$  is a weighted sum of the

$\vec{V}_j$ . The weights are scores of matching between the query  $\vec{Q}_i$  and the keys  $\vec{K}_j$ , calculated as the dot product between the two vectors. They are also *softmaxed* to sum up to 1. The transformation of  $L_Q$  vectors  $\vec{Y}_i$  as a function of  $L_K$  vectors  $\vec{X}_j$  can be efficiently computed with matrix multiplications:

$$A = \text{softmax} \left( \frac{Q \times K^T}{\sqrt{d}} \right) \times V \quad (1)$$

With  $Q$  a matrix of shape  $(L_Q, d)$ ,  $K$  a matrix of shape  $(L_K, d)$  and  $V$  a matrix of shape  $(L_K, d)$ . The  $\sqrt{d}$  at the denominator is a scaling factor used to avoid saturation in the exponential terms of the softmax function.

The multi-head attention performs  $h$  different projections into spaces of dimension  $d = D/h$ . The resulting vector  $\vec{Y}'_i$  is the concatenation of the  $h$  vectors  $\vec{y}'_i$  obtained. Thus the embedding dimension is preserved. Using multiple heads was found beneficial by the authors over using a single head of dimension  $d = D$ .

During training, the cross entropy of the  $n^{th}$  predicted token is calculated assuming all previous tokens have been generated correctly. This enables to parallelize training completely without need for recurrence. However as the  $n^{th}$  token should not depend of the following tokens, the cells in the upper right corner of the score matrix are set to  $-\infty$  such that after the softmax they are equal to 0, and the rows still sums up to 1.

## 2.2 Improving Transformer scalability with sequence length

The original attention mechanism requires the computation of a score matrix  $Q \times K^T$  of shape  $(L_Q, L_K)$ , with complexity  $O(L_Q d L_K)$ . If the query and key sequence lengths are multiplied by two, then the memory used and computation time are multiplied by 4. To improve the scalability of the transformer with sequence length, several axis of research have been explored.

Kitaev et al. [7] proposed the Reformer’s architecture, which uses an hash-bucketting algorithm to reduce the complexity of the original multi head attention operation from  $O(L^2)$  to  $O(L \log(L))$ .

Dai et al. [3] proposed the Transformer-XL’s architecture, which cuts the sequence in segments of length  $L$ . The model predicts each stage of the current segment as a function of the previous and current segment. All the segments are computed sequentially with a recurrence mechanism. The complexity

is linear with sequence length, but the computation cannot be completely parallelized due to the recurrence mechanism.

Other publications explored using a sparse attention matrix, such as the Longformer by Beltagy et al. [1] and the Big Bird model by Zaheer et al. [13]. As each token attends to a fixed number of all other tokens, the scalability is improved. These sparse attention models however require custom operations implemented in CUDA.

Some other works propose to modify the attention mechanism so that it's complexity scales linearly with sequence length. The Linformer by Wang et al. [12] projects the key and values onto a smaller sequence length dimension with matrix multiplication. It cannot however generalize to sequences longer than during training, as the weights of the projection for such tokens would be undefined.

Shen et al. [10] proposed to replace the softmax attention score.  $A = \text{softmax}\left(\frac{Q \times K^T}{\sqrt{d}}\right) \times V$  is changed into  $A = \rho(Q) \times \rho(K)^T \times V$ . With  $\rho$  the softmax function along the embedding dimension. Thanks to matrix multiplication commutativity, the order of the operations can be chosen. If  $Q$ ,  $K$  and  $V$  are of shape  $(L_Q, d)$ ,  $(L_K, d)$  and  $(L_V, d)$  respectively, the complexity of  $(\phi(Q) \times \phi(K)^T) \times V$  is  $O(L_Q \times d \times L_K)$  whereas the complexity of  $\phi(Q) \times (\phi(K)^T \times V)$  is  $O(\max(L_Q, L_K) \times d^2)$ . The right-side-first operation is linear in complexity with sequence length. The shape of the intermediate result matrix is also changed, allowing to scale better in memory requirements as well. The original *softmaxed* attention score matrix was giving rows of positive scores that sum to 1. With this change the elements of the score matrix remain positive as  $\phi(Q)$  and  $\phi(K)^T$  are matrices of positive values, but the rows of the score matrix does not sum up to 1. This work also does not give a linear complexity formulation for masked attention. If the right-side-first scheme is adopted, the attention score matrix  $\phi(Q) \times \phi(K)^T$  is never explicitly computed, and can't be masked.

Building on this idea of commutative attention function proposed by [10], Katharopoulos et al. [6] introduced their kernerlized attention function as:

$$A = \frac{\phi(Q) \times \phi(K)^T}{\sum_j (\phi(Q) \times \phi(K)^T)} \times V \quad (2)$$

The function  $\phi$  is applied element-wise and can be any positive function, for example  $\phi(x) = \text{elu}(x) + 1$ . This attention is row-wise normalized so that all rows of the score matrix are sets of positive weights adding up to one.

This preserves the objective of the original softmaxed attention scores, while allowing to perform operations in an optimal order.

The Performer by Choromanski et al. [2] exploits the same idea of a kernelized attention introduced by Katharopoulos et al. [6], with an algorithm that better approximates softmaxed attention. Most importantly they also give in annex a prefix sum algorithm to perform operations in the right-side-first order while giving the same result as masked left-side-first operation.

Although the author did not specify how to implement it, the only implementations we found of this operation requires custom CUDA code. In this work we will give an implementation of the right-side-first masked operation, with usual functions from neural network frameworks, that remains linear in complexity.

### 2.3 Alternatives to absolute positional encoding

The original multi-head attention operation introduced by Vaswani et al. [11] was intrinsically invariant by token order permutation. As token position was an important information for sequence to sequence models, they encoded the global position of each token in their embedding. Since then, some modified attention mechanisms, that depend on relative tokens position, have been proposed.

Shaw et al. [9] explored modifying the attention mechanism so that it depends on the relative distance between tokens. A second score matrix that is function of the query and the query/key relative distance is added to the original score matrix.  $A = \text{softmax} \left( \frac{Q \times K^T}{\sqrt{d}} \right) \times V$  becomes  $A = \left( \frac{Q \times K^T + S_{rel}}{\sqrt{d}} \right) \times V$  with  $S_{rel}$  of shape  $(L_Q, L_K)$  defined as  $S_{rel_{ij}} = \vec{Q}_i \cdot \vec{R}P_{clip(i-j, -k, k)}$ . Where  $k$  is the attention horizon length and  $\vec{R}P_n$  is one of  $2k + 1$  relative positional embedding, vectors of size  $d$ . Shaw et al. [9] and Huang et al. [5] observed that introducing this attention scheme improved performances. The naive calculation of this term however has a complexity of  $O(L_Q L_K d)$ . No algorithm was provided to linearize the complexity.

More recently Liutkus et al. [8] gives a stochastic positional encoding that is linear in complexity with regards to sequence length. However the implementation is complex and its stochastic nature requires that the operations be repeated several times in parallel.

Horn et al. [4] noted that the term  $S^{rel} \times V$  can be computed with linear complexity for the case where  $RP_{-k} = RP_k$ . However this is restraining

as the model can't make the difference between tokens before the attention horizon or after.

In this work we will show that the computation of  $S^{rel} \times V$  can also be done with linear complexity, without concession.

### 3 Linear complexity masked attention

In this section we will detail the prefix sum algorithm proposed by Choromanski et al. [2] for masked kernelized attention, and give an implementation with usual functions of neural network frameworks. The masked kernelized attention is defined as:

$$A^{masked} = \text{masked}(\phi(Q) \times \phi(K)^T) \times V \quad (3)$$

In this expression, masked being the operation that set all cells above the diagonal to 0 in a matrix. The naive implementation of this operation has complexity  $O(L_Q L_K d)$ .

We can change the complexity using the operation proposed by Choromanski et al. [2]. We will derive its formulation here. To start with it,  $A^{masked}$  is expressed as

$$A^{masked} = S^{masked} \times V \quad (4)$$

which in summation form is expressed as

$$A_{ij}^{masked} = \sum_k V_{kj} \times S_{ik}^{masked} \quad (5)$$

and the elements of S are defined as:

$$S_{ik}^{masked} = \begin{cases} k \leq i & \sum_l (\phi(Q)_{il} \times \phi(K)_{kl}) \\ \text{otherwise} & 0 \end{cases} \quad (6)$$

Putting these elements together leads to

$$A_{ij}^{masked} = \sum_{k=1}^i V_{kj} \times \sum_l (\phi(Q)_{il} \times \phi(K)_{kl}) \quad (7)$$

which can be reworked as

$$A_{ij}^{masked} = \sum_l \phi(Q)_{il} \times \sum_{k=1}^i (V_{kj} \times \phi(K)_{kl}) \quad (8)$$

In this work we make use of these ideas to implement the calculation of  $A^{masked}$  with complexity  $O(\max(L_Q, L_K) \times d^2)$  without custom GPU code as per the algorithm 1. With  $align(tensor, n)$  the function that truncates or repeats the last value so that  $tensor$  has length  $n$  along the first dimension,  $zeros(a, b, \dots)$  the 0-initialized tensor of shape  $(a, b, \dots)$ , and  $cumsum(tensor)$  the function that returns the cumulated sum along the first dimension.

**Algorithm 1:** calculation of  $A^{masked}$  with linear complexity

**Data:**  $\phi(Q)$ ,  $\phi(K)$ ,  $V$  tensors of shape  $(L_Q, d)$ ,  $(L_K, d)$ , and  $(L_K, d)$   
**Result:**  $A^{masked}$  tensor of shape  $(L_Q, d)$   
Unrolled  $:= zeros(L_K, d, d)$   
**for**  $i = 0$  **to**  $L_Q - 1$  **do**  
    **for**  $j = 0$  **to**  $d - 1$  **do**  
        Unrolled $_{kjl} = V_{kj} \times \phi(K)_{kl}$   
    **end**  
**end**  
Right  $:= cumsum(Unrolled, \dim = 0)$   
Aligned  $:= align(Right, L_Q)$   
 $A^{masked} := zeros(L_Q, d)$   
**for**  $i = 0$  **to**  $L_Q - 1$  **do**  
    **for**  $j = 0$  **to**  $d - 1$  **do**  
         $A_{ij}^{masked} = \sum_k (\phi(Q)_{ik} \times Aligned_{ikj})$   
    **end**  
**end**

## 4 Linear complexity RPE implementation

In this section we will detail how the relative positional encoding proposed by Shaw et al. [9] can in fact be computed with linear complexity. The  $S_{rel}$  matrix of scores is computed as  $S_{rel_{ij}} = \vec{Q}_i \cdot \vec{R}P_{clip(i-j, -k, k)}$ . In figure 1 the colors represent the index of the query and the number the index of the relative position.

The relative positional encoding's term of the attention is calculated as  $A_{rel} = S_{rel} \times V$ . Each row of the  $A_{rel}$  matrix is a weighted sum of the value

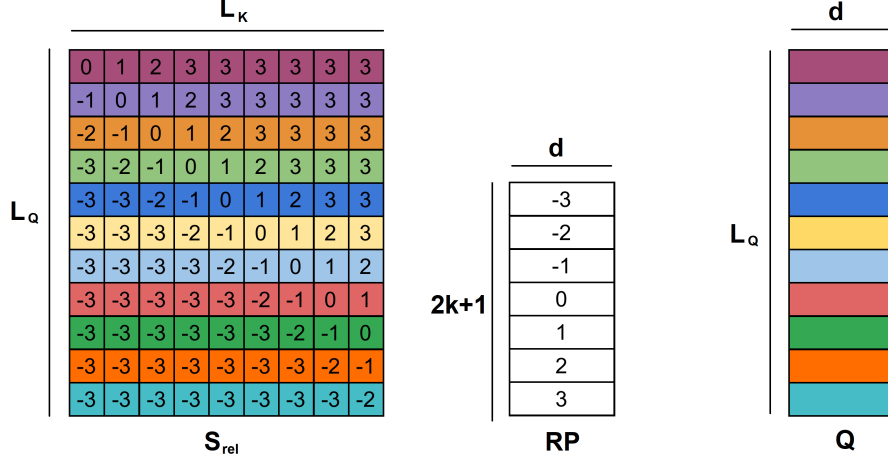


Figure 1:  $S_{rel}$  calculation

vectors  $V_i$ . Each row of  $S_{rel}$  is a set of weights.

One can observe in figure 2 that some weights are repeated several times in the  $S_{rel}$  matrix. Calculating the whole matrix can be avoided by instead calculating all possible weights only once, with complexity  $O(L_Q \times d \times (2k + 1))$ . As illustrated in figure 3, the matrix multiplication can then be replaced by a sum of three terms. The embedding dimension of size  $d$  is not represented here. Instead the horizontal axis in the figure represents all terms that must be summed. The grey squares represent some zero-padding. The first term is a cumulated sum of the value vectors that is weighted by the before-horizon set of weights (complexity  $O(\max(L_Q, L_K))$ ). The second term is a weighed sum of a moving window of the value vectors. Where the weights are the central diagonal of the  $S_{rel}$  matrix. (complexity  $O(L_Q \times (2k - 1) \times d)$ ). The last term is similar to the first term, but for the after-horizon set of weights (complexity  $O(\max(L_Q, L_K))$ ). The implementation is given in algorithm 2. The algorithm is given for bidirectional attention only, but it can easily be adapted for masked attention.



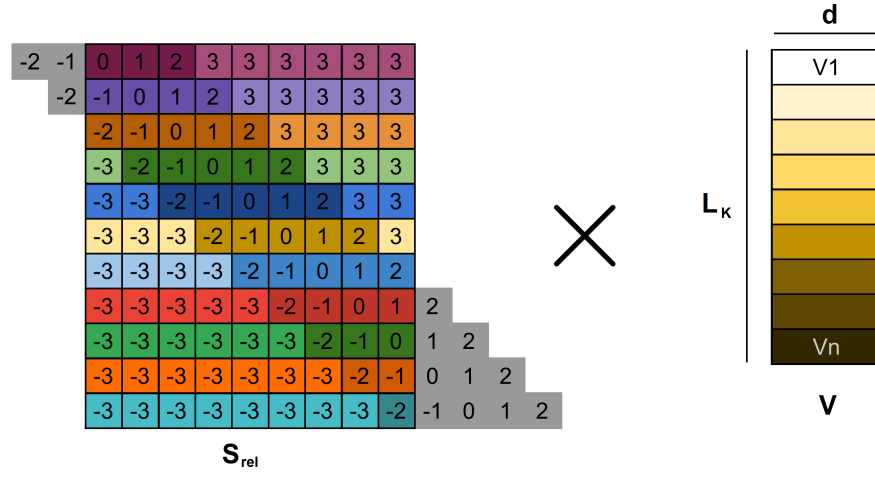


Figure 2:  $A_{rel}$  calculation

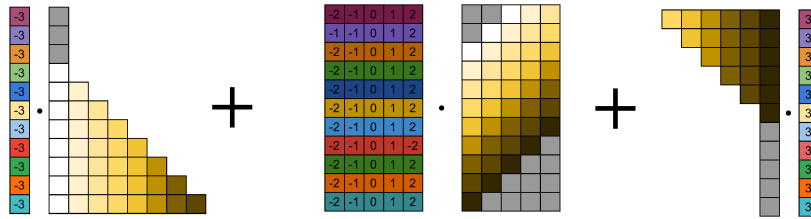


Figure 3:  $A_{rel}$  simplified calculation

**Algorithm 2:** calculation of  $A_{rel}$  with linear complexity

**Data:**  $\phi(Q)$ ,  $RP$ ,  $V$  tensors of shape  $(L_Q, d)$ ,  $(2k + 1, d)$ , and  $(L_K, d)$   
**Result:**  $A_{rel}$  tensor of shape  $(L_Q, d)$   
 $W := \phi(Q) \times RP^T$   
 $W_{before} := W_{[:,0]}$   
 $n_{before} := \min(\max(0, L_Q - k), L_K)$   
 $\text{Padding}_{before} := \text{zeros}(\min(k, L_Q), d)$   
 $\text{Padded}_{before} := \text{concatenate}(\text{Padding}_{before}, \text{cumsum}(V))$   
 $V_{before} := \text{aligned}(\text{Padded}_{before}, n_{before})$   
 $W_{horizon} := W_{[:,1:2k-1]}$   
 $V^{horizon} := \text{zeros}(L_Q, 2k - 1, d)$   
**for**  $i = 0$  **to**  $L_Q - 1$  **do**  
    **for**  $j = 0$  **to**  $2k - 2$  **do**  
        **for**  $l = 0$  **to**  $d - 1$  **do**  
             $V_{ijl}^{horizon} = \begin{cases} V_{[i-k+1+j,l]} & \text{if } 0 \leq (i - k + 1 + j) < L_K \\ 0 & \text{otherwise} \end{cases}$   
        **end**  
    **end**  
**end**  
 $W_{after} := W_{[:,2k]}$   
 $n_{after} := \min(L_Q + k, L_K)$   
 $\text{Summed} := \text{sum}(V[k - 1 : n_{after}, :], \text{dim}=0)$   
 $\text{rCumSum} := \text{Summed} - \text{cumsum}(V[k - 1 : n_{after} - 1, :])$   
 $\text{Padding}_{after} := \text{zeros}(\max(0, L_Q - \max(0, L_K - R)), D)$   
 $V_{after} := \text{concatenate}(\text{rCumSum}, \text{Padding}_{after})$   
 $A_{rel} := W_{before} \cdot V_{before} + \text{sum}(W_{horizon} \cdot V^{horizon}, \text{dim}=1) + W_{after} \cdot V_{after}$

## 5 Linear Scalable Transformer model

The proposed model replaces the scaled-dot-product-attention by a kernelized attention with RPE. Following the observations of Shaw et al. [9] that accumulating absolute positional encoding with relative positional encoding yield no benefits, the positional encoding is also removed. The algorithm used to calculate each term is chosen dynamically to occupy the least memory depending on the sequence lengths and embedding dimensions - as memory usage is easier to predict than execution time.

In this work we have chosen the following formulation, with  $\phi(x) = \text{elu}(x) + 1$ .

$$A = \frac{(\phi(Q) \times \phi(K^T) + S^{rel})}{\sum_j (\phi(Q) \times \phi(K^T) + S^{rel})} \times V \quad (9)$$

This is essentially a combination of two terms: the kernelized attention proposed by Katharopoulos et al. [6], and the relative positional encoding proposed by Shaw et al. [9]. The left term is the score matrix of shape  $(L_Q, L_K)$ , with a denominator which scales all rows so that they sum to 1. For the sake of the implementation, the multiplication must be distributed as:

$$A = \frac{(\phi(Q) \times \phi(K^T) \times V) + (S^{rel} \times V)}{\sum_j (\phi(Q) \times \phi(K^T)) + \sum_j (S^{rel})} \quad (10)$$

The denominator can be easily calculated by applying the linear complexity algorithms with  $V$  replaced by a matrix of shape  $(L_K, 1)$  full of 1, or by summing the rows of the score matrices for quadratic complexity algorithms.

For each case (masked/bidirectional) the algorithm is chosen between naive and linear complexity to occupy the least memory.

- for the masked  $\phi(Q) \times \phi(K)^T \times V$  term, the memory occupied by the naive algorithm is  $L_Q L_K$  while the linear complexity algorithm occupies  $d^2 \times \max(L_Q, L_K)$
- for the bidirectional  $\phi(Q) \times \phi(K)^T \times V$  term, the memory occupied by the naive algorithm is  $L_Q L_K$  while the linear complexity algorithm occupies  $d^2$
- for the  $S_{rel} \times V$  term (masked and bidirectional), the memory occupied by the naive algorithm is  $L_Q L_K$  while the linear complexity algorithm occupies  $L_Q \times (2k + 1 + 4)$

## 6 Results

The various algorithms have been timed on CPU and runtimes have been plotted against sequences length in log2-log2 scale. Algorithms that scale linearly with sequence length have a slope of 1, while algorithms that scale

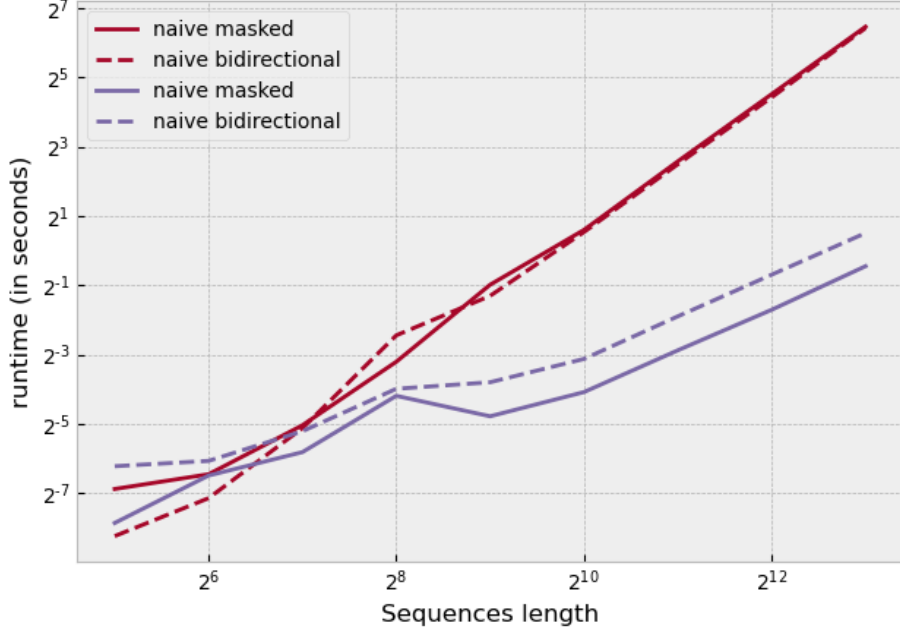


Figure 4:  $A_{rel}$  calculation runtimes on CPU for  $d = 256$  and  $k = 10$

quadratically with sequence length have a slope of 2. The Figure 4, shows the timings of the  $A_{rel}$  matrix calculation. On Figure 5 the kernelized attention algorithms have been timed. We can observe that although our algorithm for masked attention has a linear complexity, our linear algorithm is still often slower than the quadratic complexity algorithm. The bottleneck in our implementation was the cumulated sum operation over the *Unrolled* tensor, which was surprisingly slow. A CUDA implemented operation might still be beneficial for this particular algorithm, also there might still be room for optimization in our implementation. Finally Figure 6 shows the timings of the vanilla multi-head attention compared to our proposed multi-head kernelized attention calculated with dynamically chosen algorithms.

## 7 Conclusion

In the present work an algorithm to compute Shaw et al. [9] relative positional embedding with linear complexity has been presented. An implementation of Choromanski et al. [2] prefix sum algorithm that doesn't requires custom

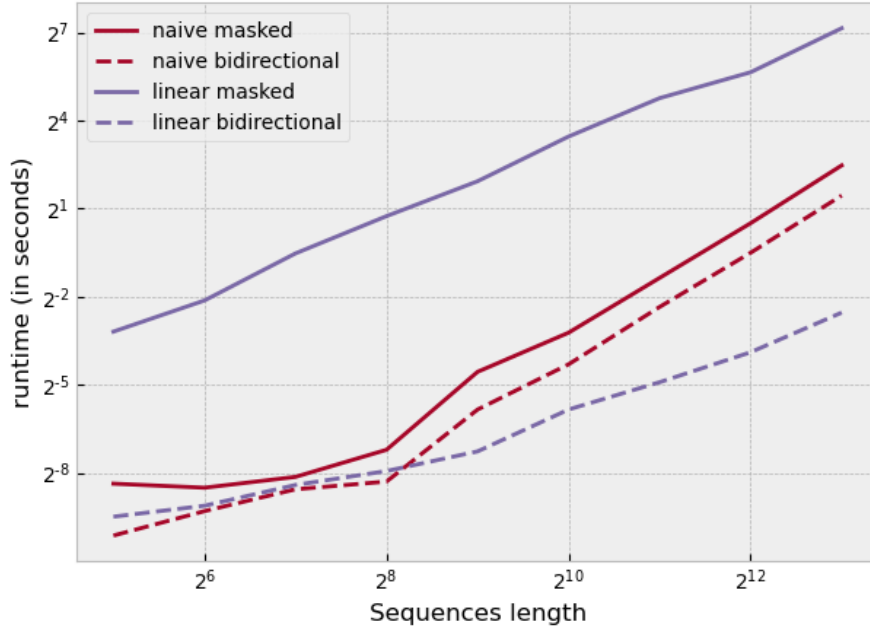


Figure 5: kernelized attention runtimes on CPU for  $d = 256$

CUDA code (while maintaining linear complexity) was also presented.

These two elements allowed to define a kernelized attention function with relative positional encoding, that can be computed with linear complexity with regards to sequence length. The proposed model presents linear scalability with high sequences lengths, can be implemented without CUDA code in neural network framework, and is adapted to sequence to sequence tasks instead of being restricted to encoders only models.

## References

- [1] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020. URL <https://arxiv.org/abs/2004.05150>.
- [2] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian

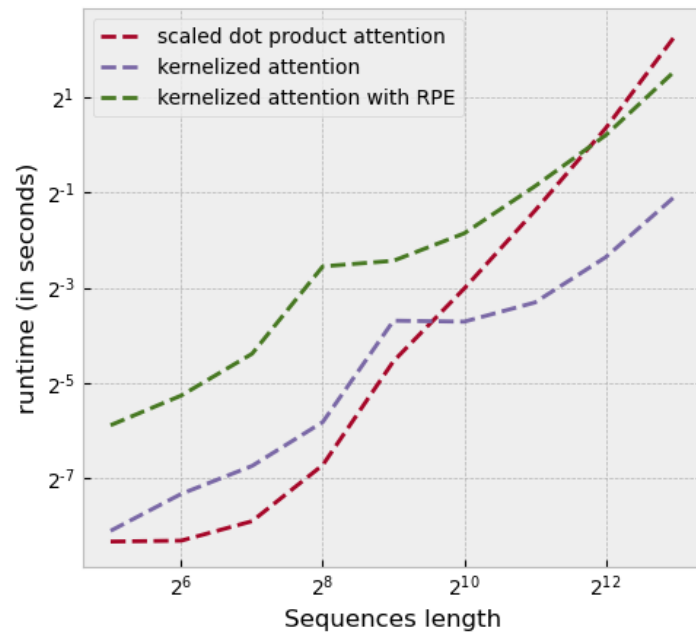


Figure 6: multi-head attention runtimes on CPU (bidirectional attentions only)

- Weller. Rethinking attention with performers, 2021. URL <https://arxiv.org/abs/2009.14794>.
- [3] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context, 2019. URL <https://arxiv.org/abs/1901.02860>.
  - [4] Max Horn, Kumar Shridhar, Elrich Groenewald, and Philipp F. M. Baumann. Translational equivariance in kernelizable attention, 2021. URL <https://arxiv.org/abs/2102.07680>.
  - [5] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinulescu, and Douglas Eck. Music transformer, 2018. URL <https://arxiv.org/abs/1809.04281>.
  - [6] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention, 2020. URL <https://arxiv.org/abs/2006.16236>.
  - [7] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer, 2020. URL <https://arxiv.org/abs/2001.04451>.
  - [8] Antoine Liutkus, Ondřej Cífka, Shih-Lun Wu, Umut Şimşekli, Yi-Hsuan Yang, and Gaël Richard. Relative positional encoding for transformers with linear complexity, 2021. URL <https://arxiv.org/abs/2105.08399>.
  - [9] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations, 2018. URL <https://arxiv.org/abs/1803.02155>.
  - [10] Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities, 2020. URL <https://arxiv.org/abs/1812.01243>.
  - [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017. URL <https://arxiv.org/abs/1706.03762>.

- [12] Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity, 2020. URL <https://arxiv.org/abs/2006.04768>.
- [13] Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transformers for longer sequences, 2021. URL <https://arxiv.org/abs/2007.14062>.