
Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention

Angelos Katharopoulos^{1,2} Apoorv Vyas^{1,2} Nikolaos Pappas³ François Fleuret^{2,4,*}

Abstract

Transformers achieve remarkable performance in several tasks but due to their quadratic complexity, with respect to the input's length, they are prohibitively slow for very long sequences. To address this limitation, we express the self-attention as a linear dot-product of kernel feature maps and make use of the associativity property of matrix products to reduce the complexity from $\mathcal{O}(N^2)$ to $\mathcal{O}(N)$, where N is the sequence length. We show that this formulation permits an iterative implementation that dramatically accelerates autoregressive transformers and reveals their relationship to recurrent neural networks. Our *linear transformers* achieve similar performance to vanilla transformers and they are up to 4000x faster on autoregressive prediction of very long sequences.

1. Introduction

Transformer models were originally introduced by Vaswani et al. (2017) in the context of neural machine translation (Sutskever et al., 2014; Bahdanau et al., 2015) and have demonstrated impressive results on a variety of tasks dealing with natural language (Devlin et al., 2019), audio (Sperber et al., 2018), and images (Parmar et al., 2019). Apart from tasks with ample supervision, transformers are also effective in transferring knowledge to tasks with limited or no supervision when they are pretrained with autoregressive (Radford et al., 2018; 2019) or masked language modeling objectives (Devlin et al., 2019; Yang et al., 2019; Song et al., 2019; Liu et al., 2020).

However, these benefits often come with a very high computational and memory cost. The bottleneck is mainly caused

¹Idiap Research Institute, Switzerland ²EPFL, Switzerland ³University of Washington, Seattle, USA ⁴University of Geneva, Switzerland. *Work done at Idiap. Correspondence to: Angelos Katharopoulos <firstname.lastname@idiap.ch>.

Proceedings of the 37th International Conference on Machine Learning, Online, PMLR 119, 2020. Copyright 2020 by the author(s).

by the global receptive field of self-attention, which processes contexts of N inputs with a quadratic memory and time complexity $\mathcal{O}(N^2)$. As a result, in practice transformers are slow to train and their context is *limited*. This disrupts temporal coherence and hinders the capturing of long-term dependencies. Dai et al. (2019) addressed the latter by attending to memories from previous contexts albeit at the expense of computational efficiency.

Lately, researchers shifted their attention to approaches that increase the context length without sacrificing efficiency. Towards this end, Child et al. (2019) introduced sparse factorizations of the attention matrix to reduce the self-attention complexity to $\mathcal{O}(N\sqrt{N})$. Kitaev et al. (2020) further reduced the complexity to $\mathcal{O}(N \log N)$ using locality-sensitive hashing. This made scaling to long sequences possible. Even though the aforementioned models can be efficiently trained on large sequences, they do not speed-up autoregressive inference.

In this paper, we introduce the *linear transformer* model that significantly reduces the memory footprint and scales linearly with respect to the context length. We achieve this by using a kernel-based formulation of self-attention and the associative property of matrix products to calculate the self-attention weights (§ 3.2). Using our linear formulation, we also express causal masking with linear complexity and constant memory (§ 3.3). This reveals the relation between transformers and RNNs, which enables us to perform autoregressive inference orders of magnitude faster (§ 3.4).

Our evaluation on image generation and automatic speech recognition demonstrates that *linear transformer* can reach the performance levels of transformer, while being up to three orders of magnitude faster during inference.

2. Related Work

In this section, we provide an overview of the most relevant works that seek to address the large memory and computational requirements of transformers. Furthermore, we discuss methods that theoretically analyze the core component of the transformer model, namely self-attention. Finally, we present another line of work that seeks to alleviate the softmax bottleneck in the attention computation.

2.1. Efficient Transformers

Existing works seek to improve memory efficiency in transformers through weight pruning (Michel et al., 2019), weight factorization (Lan et al., 2020), weight quantization (Zafir et al., 2019) or knowledge distillation. Clark et al. (2020) proposed a new pretraining objective called replaced token detection that is more sample efficient and reduces the overall computation. Lample et al. (2019) used product-key attention to increase the capacity of any layer with negligible computational overhead.

Reducing the memory or computational requirements with these methods leads to training or inference time speedups, but, fundamentally, the time complexity is still quadratic with respect to the sequence length which hinders scaling to long sequences. In contrast, we show that our method reduces both memory and time complexity of transformers both theoretically (§ 3.2) and empirically (§ 4.1).

Another line of research aims at increasing the “context” of self-attention in transformers. Context refers to the maximum part of the sequence that is used for computing self-attention. Dai et al. (2019) introduced Transformer-XL which achieves state-of-the-art in language modeling by learning dependencies beyond a fixed length context without disrupting the temporal coherence. However, maintaining previous contexts in memory introduces significant additional computational cost. In contrast, Sukhbaatar et al. (2019) extended the context length significantly by learning the optimal attention span per attention head, while maintaining control over the memory footprint and computation time. Note that both approaches have the same asymptotic complexity as the vanilla model. In contrast, we improve the asymptotic complexity of the self-attention, which allows us to use significantly larger context.

More related to our model are the works of Child et al. (2019) and Kitaev et al. (2020). The former (Child et al., 2019) introduced sparse factorizations of the attention matrix reducing the overall complexity from quadratic to $\mathcal{O}(N\sqrt{N})$ for generative modeling of long sequences. More recently, Kitaev et al. (2020) proposed Reformer. This method further reduces complexity to $\mathcal{O}(N \log N)$ by using locality-sensitive hashing (LSH) to perform fewer dot products. Note that in order to be able to use LSH, Reformer constrains the keys, for the attention, to be identical to the queries. As a result this method cannot be used for decoding tasks where the keys need to be different from the queries. In comparison, *linear transformers* impose no constraints on the queries and keys and scale linearly with respect to the sequence length. Furthermore, they can be used to perform inference in autoregressive tasks three orders of magnitude faster, achieving comparable performance in terms of validation perplexity.

2.2. Understanding Self-Attention

There have been few efforts to better understand self-attention from a theoretical perspective. Tsai et al. (2019) proposed a kernel-based formulation of attention in transformers which considers attention as applying a kernel smoother over the inputs with the kernel scores being the similarity between inputs. This formulation provides a better way to understand attention components and integrate the positional embedding. In contrast, we use the kernel formulation to speed up the calculation of self-attention and lower its computational complexity. Also, we observe that if a kernel with positive similarity scores is applied on the queries and keys, linear attention converges normally.

More recently, Cordonnier et al. (2020) provided theoretical proofs and empirical evidence that a multi-head self-attention with sufficient number of heads can express any convolutional layer. Here, we instead show that a self-attention layer trained with an autoregressive objective can be seen as a recurrent neural network and this observation can be used to significantly speed up inference time of autoregressive transformer models.

2.3. Linearized softmax

For many years, softmax has been the bottleneck for training classification models with a large number of categories (Goodman, 2001; Morin & Bengio, 2005; Mnih & Hinton, 2009). Recent works (Blanc & Rendle, 2017; Rawat et al., 2019), have approximated softmax with a linear dot product of feature maps to speed up the training through sampling. Inspired from these works, we linearize the softmax attention in transformers. Concurrently with this work, Shen et al. (2020) explored the use of linearized attention for the task of object detection in images. In comparison, we do not only linearize the attention computation, but also develop an autoregressive transformer model with linear complexity and constant memory for both inference and training. Moreover, we show that through the lens of kernels, every transformer can be seen as a recurrent neural network.

3. Linear Transformers

In this section, we formalize our proposed *linear transformer*. We present that changing the attention from the traditional *softmax* attention to a feature map based dot product attention results in better time and memory complexity as well as a causal model that can perform sequence generation in linear time, similar to a recurrent neural network.

Initially, in § 3.1, we introduce a formulation for the transformer architecture introduced in (Vaswani et al., 2017). Subsequently, in § 3.2 and § 3.3 we present our proposed *linear transformer* and finally, in § 3.4 we rewrite the transformer as a recurrent neural network.

3.1. Transformers

Let $x \in \mathbb{R}^{N \times F}$ denote a sequence of N feature vectors of dimensions F . A transformer is a function $T : \mathbb{R}^{N \times F} \rightarrow \mathbb{R}^{N \times F}$ defined by the composition of L transformer layers $T_1(\cdot), \dots, T_L(\cdot)$ as follows,

$$T_l(x) = f_l(A_l(x) + x). \quad (1)$$

The function $f_l(\cdot)$ transforms each feature independently of the others and is usually implemented with a small two-layer feedforward network. $A_l(\cdot)$ is the self attention function and is the only part of the transformer that acts across sequences.

The self attention function $A_l(\cdot)$ computes, for every position, a weighted average of the feature representations of all other positions with a weight proportional to a similarity score between the representations. Formally, the input sequence x is projected by three matrices $W_Q \in \mathbb{R}^{F \times D}$, $W_K \in \mathbb{R}^{F \times D}$ and $W_V \in \mathbb{R}^{F \times M}$ to corresponding representations Q , K and V . The output for all positions, $A_l(x) = V'$, is computed as follows,

$$\begin{aligned} Q &= xW_Q, \\ K &= xW_K, \\ V &= xW_V, \end{aligned} \quad (2)$$

$$A_l(x) = V' = \text{softmax} \left(\frac{QK^T}{\sqrt{D}} \right) V.$$

Note that in the previous equation, the softmax function is applied rowwise to QK^T . Following common terminology, the Q , K and V are referred to as the “queries”, “keys” and “values” respectively.

Equation 2 implements a specific form of self-attention called softmax attention where the similarity score is the exponential of the dot product between a query and a key. Given that subscripting a matrix with i returns the i -th row as a vector, we can write a generalized attention equation for any similarity function as follows,

$$V'_i = \frac{\sum_{j=1}^N \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^N \text{sim}(Q_i, K_j)}. \quad (3)$$

Equation 3 is equivalent to equation 2 if we substitute the similarity function with $\text{sim}(q, k) = \exp\left(\frac{q^T k}{\sqrt{D}}\right)$.

3.2. Linearized Attention

The definition of attention in equation 2 is generic and can be used to define several other attention implementations such as polynomial attention or RBF kernel attention (Tsai et al., 2019). Note that the only constraint we need to impose to $\text{sim}(\cdot)$, in order for equation 3 to define an attention function, is to be non-negative. This includes all kernels $k(x, y) : \mathbb{R}^{2 \times F} \rightarrow \mathbb{R}_+$.

Given such a kernel with a feature representation $\phi(x)$ we can rewrite equation 2 as follows,

$$V'_i = \frac{\sum_{j=1}^N \phi(Q_i)^T \phi(K_j) V_j}{\sum_{j=1}^N \phi(Q_i)^T \phi(K_j)}, \quad (4)$$

and then further simplify it by making use of the associative property of matrix multiplication to

$$V'_i = \frac{\phi(Q_i)^T \sum_{j=1}^N \phi(K_j) V_j^T}{\phi(Q_i)^T \sum_{j=1}^N \phi(K_j)}. \quad (5)$$

The above equation is simpler to follow when the numerator is written in vectorized form as follows,

$$\left(\phi(Q) \phi(K)^T \right) V = \phi(Q) \left(\phi(K)^T V \right). \quad (6)$$

Note that the feature map $\phi(\cdot)$ is applied rowwise to the matrices Q and K .

From equation 2, it is evident that the computational cost of softmax attention scales with $\mathcal{O}(N^2)$, where N represents the sequence length. The same is true for the memory requirements because the full attention matrix must be stored to compute the gradients with respect to the queries, keys and values. In contrast, our proposed *linear transformer* from equation 5 has time and memory complexity $\mathcal{O}(N)$ because we can compute $\sum_{j=1}^N \phi(K_j) V_j^T$ and $\sum_{j=1}^N \phi(K_j)$ once and reuse them for every query.

3.2.1. FEATURE MAPS AND COMPUTATIONAL COST

For softmax attention, the total cost in terms of multiplications and additions scales as $\mathcal{O}(N^2 \max(D, M))$, where D is the dimensionality of the queries and keys and M is the dimensionality of the values. On the contrary, for linear attention, we first compute the feature maps of dimensionality C . Subsequently, computing the new values requires $\mathcal{O}(NCM)$ additions and multiplications.

The previous analysis does not take into account the choice of kernel and feature function. Note that the feature function that corresponds to the exponential kernel is infinite dimensional, which makes the linearization of exact softmax attention infeasible. On the other hand, the polynomial kernel, for example, has an exact finite dimensional feature map and has been shown to work equally well with the exponential or RBF kernel (Tsai et al., 2019). The computational cost for a linearized polynomial transformer of degree 2 is $\mathcal{O}(ND^2M)$. This makes the computational complexity favorable when $N > D^2$. Note that this is true in practice since we want to be able to process sequences with tens of thousands of elements.

For our experiments, that deal with smaller sequences, we employ a feature map that results in a positive similarity

function as defined below,

$$\phi(x) = \text{elu}(x) + 1, \quad (7)$$

where $\text{elu}(\cdot)$ denotes the exponential linear unit (Clevert et al., 2015) activation function. We prefer $\text{elu}(\cdot)$ over $\text{relu}(\cdot)$ to avoid setting the gradients to 0 when x is negative. This feature map results in an attention function that requires $\mathcal{O}(NDM)$ multiplications and additions. In our experimental section, we show that the feature map of equation 7 performs on par to the full transformer, while significantly reducing the computational and memory requirements.

3.3. Causal Masking

The transformer architecture can be used to efficiently train autoregressive models by masking the attention computation such that the i -th position can only be influenced by a position j if and only if $j \leq i$, namely a position cannot be influenced by the subsequent positions. Formally, this causal masking changes equation 3 as follows,

$$V'_i = \frac{\sum_{j=1}^i \text{sim}(Q_i, K_j) V_j}{\sum_{j=1}^i \text{sim}(Q_i, K_j)}. \quad (8)$$

Following the reasoning of § 3.2, we linearize the masked attention as described below,

$$V'_i = \frac{\phi(Q_i)^T \sum_{j=1}^i \phi(K_j) V_j^T}{\phi(Q_i)^T \sum_{j=1}^i \phi(K_j)}. \quad (9)$$

By introducing S_i and Z_i as follows,

$$S_i = \sum_{j=1}^i \phi(K_j) V_j^T, \quad (10)$$

$$Z_i = \sum_{j=1}^i \phi(K_j), \quad (11)$$

we can simplify equation 9 to

$$V'_i = \frac{\phi(Q_i)^T S_i}{\phi(Q_i)^T Z_i}. \quad (12)$$

Note that, S_i and Z_i can be computed from S_{i-1} and Z_{i-1} in constant time hence making the computational complexity of linear transformers with causal masking linear with respect to the sequence length.

3.3.1. GRADIENT COMPUTATION

A naive implementation of equation 12, in any deep learning framework, requires storing all intermediate values S_i in order to compute the gradients. This increases the memory consumption by $\max(D, M)$ times; thus hindering the

applicability of causal linear attention to longer sequences or deeper models. To address this, we derive the gradients of the numerator in equation 9 as cumulative sums. This allows us to compute both the forward and backward pass of causal linear attention in **linear time** and **constant memory**. A detailed derivation is provided in the supplementary material.

Given the numerator \bar{V}_i and the gradient of a scalar loss function with respect to the numerator $\nabla_{\bar{V}_i} \mathcal{L}$, we derive $\nabla_{\phi(Q_i)} \mathcal{L}$, $\nabla_{\phi(K_i)} \mathcal{L}$ and $\nabla_{V_i} \mathcal{L}$ as follows,

$$\nabla_{\phi(Q_i)} \mathcal{L} = \nabla_{\bar{V}_i} \mathcal{L} \left(\sum_{j=1}^i \phi(K_j) V_j^T \right)^T, \quad (13)$$

$$\nabla_{\phi(K_i)} \mathcal{L} = \left(\sum_{j=i}^N \phi(Q_j) \left(\nabla_{\bar{V}_j} \mathcal{L} \right)^T \right) V_i, \quad (14)$$

$$\nabla_{V_i} \mathcal{L} = \left(\sum_{j=i}^N \phi(Q_j) \left(\nabla_{\bar{V}_j} \mathcal{L} \right)^T \right) \phi(K_i). \quad (15)$$

The cumulative sum terms in equations 9, 13-15 are computed in linear time and require constant memory with respect to the sequence length. This results in an algorithm with computational complexity $\mathcal{O}(NCM)$ and memory $\mathcal{O}(N \max(C, M))$ for a given feature map of C dimensions. A pseudocode implementation of the forward and backward pass of the numerator is given in algorithm 1.

3.3.2. TRAINING AND INFERENCE

When training an autoregressive transformer model the full ground truth sequence is available. This makes layerwise parallelism possible both for $f_i(\cdot)$ of equation 1 and the attention computation. As a result, transformers are more efficient to train than recurrent neural networks. On the other hand, during inference the output for timestep i is the input for timestep $i + 1$. This makes autoregressive models impossible to parallelize. Moreover, the cost per timestep for transformers is not constant; instead, it scales with the square of the current sequence length because attention must be computed for all previous timesteps.

Our proposed *linear transformer* model *combines the best of both worlds*. When it comes to training, the computations can be parallelized and take full advantage of GPUs or other accelerators. When it comes to inference, the cost per time and memory for one prediction is constant for our model. This means we can simply store the $\phi(K_j) V_j^T$ matrix as an internal state and update it at every time step like a recurrent neural network. This results in inference **thousands of times faster** than other transformer models.

3.4. Transformers are RNNs

In literature, transformer models are considered to be a fundamentally different approach to recurrent neural networks. However, from the causal masking formulation in § 3.3 and the discussion in the previous section, it becomes evident that any transformer layer with causal masking can be written as a model that, given an input, modifies an internal state and then predicts an output, namely a Recurrent Neural Network (RNN). Note that, in contrast to Universal Transformers (Dehghani et al., 2018), we consider the recurrence with respect to time and not depth.

In the following equations, we formalize the transformer layer of equation 1 as a recurrent neural network. The resulting RNN has two hidden states, namely the attention memory s and the normalizer memory z . We use subscripts to denote the timestep in the recurrence.

$$s_0 = 0, \quad (16)$$

$$z_0 = 0, \quad (17)$$

$$s_i = s_{i-1} + \phi(x_i W_K)(x_i W_V)^T, \quad (18)$$

$$z_i = z_{i-1} + \phi(x_i W_K), \quad (19)$$

$$y_i = f_l \left(\frac{\phi(x_i W_Q)^T s_i}{\phi(x_i W_Q)^T z_i} + x_i \right). \quad (20)$$

In the above equations, x_i denotes the i -th input and y_i the i -th output for a specific transformer layer. Note that our formulation does not impose any constraint on the feature function and it can be used for representing *any transformer* model, in theory even the ones using softmax attention. This formulation is a first step towards better understanding the relationship between transformers and popular recurrent networks (Hochreiter & Schmidhuber, 1997) and the processes used for storing and retrieving information.

4. Experiments

In this section, we analyze experimentally the performance of the proposed *linear transformer*. Initially, in § 4.1, we evaluate the linearized attention in terms of computational cost, memory consumption and convergence on synthetic data. To further showcase the effectiveness of *linear transformers*, we evaluate our model on two real-world applications, image generation in § 4.2 and automatic speech recognition in § 4.3. We show that our model achieves competitive performance with respect to the state-of-the-art transformer architectures, while requiring significantly less GPU memory and computation.

Throughout our experiments, we compare our model with two baselines, the full transformer with softmax attention and the Reformer (Kitaev et al., 2020), the latter being a state-of-the-art accelerated transformer architecture. For the Reformer, we use a PyTorch reimplementation of the pub-

Algorithm 1 Linear transformers with causal masking

```

function forward( $\phi(Q), \phi(K), V$ ):
   $V' \leftarrow 0, S \leftarrow 0$ 
  for  $i = 1, \dots, N$  do
     $S \leftarrow S + \phi(K_i) V_i^T$  equation 10
     $\bar{V}_i \leftarrow \phi(Q_i) S$ 
  end
  return  $\bar{V}$ 
end

function backward( $\phi(Q), \phi(K), V, G$ ):
  /*  $G$  is the gradient of the loss
     with respect to the output of
     forward */
   $S \leftarrow 0, \nabla_{\phi(Q)} \mathcal{L} \leftarrow 0$ 
  for  $i = 1, \dots, N$  do
     $S \leftarrow S + \phi(K_i) V_i^T$ 
     $\nabla_{\phi(Q_i)} \mathcal{L} \leftarrow G_i S^T$  equation 13
  end
   $S \leftarrow 0, \nabla_{\phi(K)} \mathcal{L} \leftarrow 0, \nabla_V \mathcal{L} \leftarrow 0$ 
  for  $i = N, \dots, 1$  do
     $S \leftarrow S + \phi(Q_i) G_i^T$ 
     $\nabla_{V_i} \mathcal{L} \leftarrow S^T \phi(K_i)$  equation 15
     $\nabla_{\phi(K_i)} \mathcal{L} \leftarrow S V_i$  equation 14
  end
  return  $\nabla_{\phi(Q)} \mathcal{L}, \nabla_{\phi(K)} \mathcal{L}, \nabla_V \mathcal{L}$ 
end

```

lished code and for the full transformer we use the default PyTorch implementation. Note that for Reformer, we do not use the reversible layers, however, this does not affect the results as we only measure the memory consumption with respect to the self attention layer. In all experiments, we use **softmax** (Vaswani et al., 2017) to refer to the standard transformer architecture, **linear** for our proposed *linear transformers* and **lsh-X** for Reformer (Kitaev et al., 2020), where X denotes the hashing rounds.

For training the *linear transformers*, we use the feature map of equation 7. Our PyTorch (Paszke et al., 2019) code with documentation and examples can be found at <https://linear-transformers.com/>. The constant memory gradient computation of equations 13-15 is implemented in approximately 200 lines of CUDA code.

4.1. Synthetic Tasks

4.1.1. CONVERGENCE ANALYSIS

To examine the convergence properties of *linear transformers* we train on an artificial copy task with causal masking. Namely, the transformers have to copy a series of symbols similar to the sequence duplication task of Kitaev et al. (2020). We use a sequence of maximum length 128 with 10

Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

Real-Time Litigation Alerts



Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

Advanced Docket Research



With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

Analytics At Your Fingertips



Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.