

# Linformer

Monday, June 28, 2021 3:30 AM

- Paper link <https://arxiv.org/pdf/2006.04768.pdf>

## - Abstract

- $\Rightarrow$  Self attention mechanism in transformers take  $O(n^2)$  time and space.
- $\Rightarrow$  It can, however, be approximated by a low-rank matrix.
- $\Rightarrow$  This reduces the complexity from  $O(n^2)$  to  $O(n)$ .
- $\Rightarrow$  The resulting model that uses the mechanism above, called **Linformer**, performs on par with standard transformer models.

## - Intro

- $\Rightarrow$  Question: How to avoid quadratic time complexity of self-attention mechanism?
- $\Rightarrow$  Approaches
  - Child et al. 2019  $\rightarrow O(n\sqrt{n})$  [\[Link\]](#)
    - $\hookrightarrow$  large performance drop (2%)  
with limited speedup (20%)
  - Kitaev et al. 2020  $\rightarrow O(n \log n)$  [\[Link\]](#)
    - $\hookrightarrow$  uses locality sensitive hashing
    - $\hookrightarrow$  gains only appear with length  $> 2048$
    - $\hookrightarrow$  multi-round hashing  $\rightarrow$  increased sequential operations
- $\Rightarrow$  Insight: self-attention is low rank
  - The stochastic matrix formed by self-attention can be approximated by a low-rank matrix through the

approximated by a low-rank matrix through the Johnson-Lindenstrauss lemma.

- Turn the scaled dot product attention to a number of smaller attentions through linear projections.

$\Rightarrow$  Linformer models performed well compared to standard transformers

- pretraining performance
- fine-tuned performance

$\Rightarrow O(n)$  complexity leads to significant training and inference speedup.

### - Background

$\Rightarrow$  Transformers are made of multi-head attention (MHA) blocks.

$$\text{MHA}(Q, V, K) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where  $Q, V, K \in \mathbb{R}^{n \times d_m}$

$d_m = \text{embedding dimension}$

$n = \text{sequence length.}$

$$W^O \in \mathbb{R}^{h d_v \times d_m}$$

$d_v = \text{dimension of the projection subspace of each head.}$

$\Rightarrow$  For each head,

$$\begin{aligned} \text{head}_i &= \text{Attention}(Q W_i^Q, K W_i^K, V W_i^V) \\ &= \text{softmax}\left(\frac{Q W_i^Q (K W_i^K)^T}{\sqrt{d_k}}\right) V W_i^V \end{aligned}$$

where  $W_i^Q, W_i^K \in \mathbb{R}^{d_m \times d_k}$   
 $W_i^V \in \mathbb{R}^{d_m \times d_v}$

$$W_i^v \in \mathbb{R}^{d_m \times d_v}$$

$d_k$  = dimension of the embeddings of the keys and queries.

$$\Rightarrow \text{Define } P = \text{softmax} \left( \frac{QW_i^Q (KW_i^K)^T}{\sqrt{d_k}} \right) \quad \leftarrow \text{"context mapping matrix"}$$

$\Rightarrow$  Computing  $P \in \mathbb{R}^{n \times n}$  involves multiplying  $QW_i^Q \in \mathbb{R}^{n \times d_k}$  and  $KW_i^K \in \mathbb{R}^{n \times d_k}$

which takes  $O(n^2 d_k)$

- Self-attention is low-rank.

$\Rightarrow$  The context mapping matrix  $P$  is low-rank.

- The paper took two pretrained transformer models

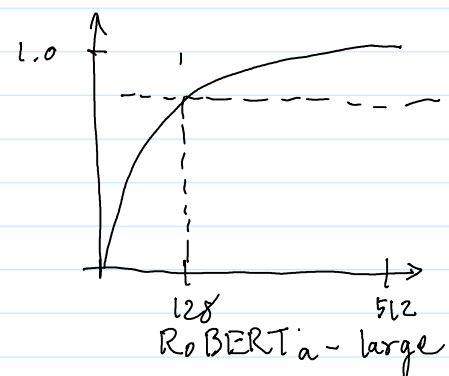
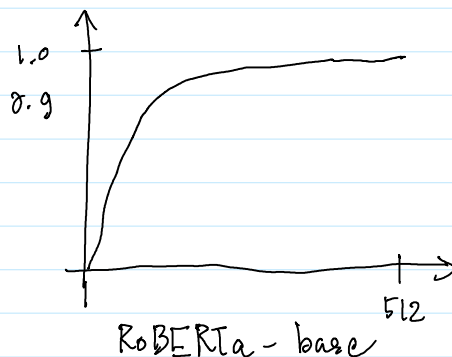
$\Rightarrow$  RoBERTa-base (12 layers)

$\Rightarrow$  RoBERTa-large (24 layers)

- They computed  $P$  over different layers, different heads, over 10k sentences.

- They find the singular values of the matrices and plotted the cumulative normalized singular value graphs. The graphs are then averaged.

- The graphs look kind of like this.



<sup>512</sup>  
RoBERTa - base

<sup>128</sup> <sup>512</sup>  
RoBERTa - large

- So, it seems that 90% of the mass of the singular values are located in the first 128.
- In higher layers, the mass concentrates on the large singular values even more.

⇒ Theoretically, we have the following result.

Thm For any  $Q, K, V \in \mathbb{R}^{n \times d}$  and  $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d \times d}$  for any column vector  $w \in \mathbb{R}^n$  of matrix  $VW_i^V$ , there exists a random matrix  $\tilde{P} \in \mathbb{R}^n \times \mathbb{R}^n$  of rank  $\mathcal{O}(\log n)$  s.t.  
$$\Pr(\|\tilde{P}w^T - Pw^T\| < \epsilon \|Pw^T\|) > 1 - o(1)$$

⇒ The proof uses the distributional Janson-Lindenstrauss lemma.

⇒ In fact,  $\tilde{P} = PR^T R$  where  $R$  is a  $k \times n$  matrix whose entries are sampled from the distribution  $\frac{1}{\sqrt{k}} \mathcal{N}(0, 1)$ .

- Model

⇒ First, let us say that  $d_k = d_v = d$  to simplify stuffs.

⇒ We introduce linear mappings  $E_i, F_i \in \mathbb{R}^{k \times n}$  when computing the keys and values. So,

$KW_i^K \rightarrow E_i KW_i^K$ ( $n \times d$ )                      ( $k \times d$ )	$VW_i^V \rightarrow F_i VW_i^V$ ( $n \times d$ )                      ( $k \times d$ )
---	---

⇒ So, the new head function becomes

$$\begin{aligned} \text{head}_i &= \text{Attention}(QW_i^Q, E_i KW_i^K, F_i VW_i^V) \\ &= \text{softmax} \left( \frac{\overbrace{QW_i^Q (E_i KW_i^K)^T}^{n \times d \quad k \times d}}{\underbrace{\quad}_{\sqrt{d}}} \overbrace{F_i VW_i^V}^{k \times d} \right) \end{aligned}$$

$$\tilde{P} \in \mathbb{R}^{n \times k}$$

Note that the calculation above can be carried out in  $O(ndk)$  time, which is linear in  $n$  if  $k$  is not a function of  $n$ .

$\Rightarrow$  Thm Let  $k = c \min \left\{ \frac{q \log d}{\epsilon^2}, \frac{5 \log n}{\epsilon^2} \right\}$ . There exists a way to sample  $E_i, F_i \in \mathbb{R}^{k \times n}$  such that, for any row vector  $w$  of  $QW_i^Q (KW_i^K)^T / \sqrt{d}$ , we have that

$$\Pr(\| \text{softmax}(wE_i^T) F_i V W_i^V - \text{softmax}(w) V W_i^V \| \leq \epsilon \| \text{softmax}(w) V W_i^V \|) < 1 - O(1).$$

$\Rightarrow$  The forms of  $E_i$  and  $F_i$  are

$$E_i = SR, \quad F_i = e^{-\delta} R$$

where  $R \in \mathbb{R}^{k \times n}$  is a matrix whose entries are sampled from the distribution  $\frac{1}{\sqrt{k}} \mathcal{N}(\sigma, 1)$ .

$\Rightarrow$  Note that the projection matrices can be shared between heads and layers. Moreover, we can use the same matrix for the keys and values.

### - Experiments

$\Rightarrow$  Pretraining performance

- Measured with perplexity curves.

- Linformer with  $k = 128, n = 512$   
 $k = 256, n = 1024$  } are already on par with standard transformer

- The paper also evaluated projection matrix sharing schemes, and it found that performances were similar.

- If we fix  $k$  and increase  $n$ , we find that performance

degrades when the model is not trained to convergence. However, performances become about the same after convergence.

⇒ Finetuning performance

- Linformer with  $n=512$ ,  $k=128$  performed on par with RoBERTa
- $k=256$  actually performed better.
- Best results were obtained by sharing one projection matrix across the whole model.