

Fast Monte-Carlo Approximation of the Attention Mechanism

Hyunjun Kim and JeongGil Ko

School of Integrated Technology, Yonsei University
hyunjun.kim@yonsei.ac.kr, jeonggil.ko@yonsei.ac.kr

Abstract

We introduce Monte-Carlo Attention (MCA), a randomized approximation method for reducing the computational cost of self-attention mechanisms in Transformer architectures. MCA exploits the fact that the importance of each token in an input sequence vary with respect to their attention scores; thus, some degree of error can be tolerable when encoding tokens with low attention. Using approximate matrix multiplication, MCA applies different error bounds to encode input tokens such that those with low attention scores are computed with relaxed precision, whereas errors of salient elements are minimized. MCA can operate in parallel with other attention optimization schemes and does not require model modification. We study the theoretical error bounds and demonstrate that MCA reduces attention complexity (in FLOPS) for various Transformer models by up to $11\times$ in GLUE benchmarks without compromising model accuracy.

Introduction

Attention mechanisms have become a predominant modeling paradigm for state-of-the-art deep neural networks. Among many, the Transformer architecture (Vaswani et al. 2017) is a major contribution that drove this success, which demonstrated unprecedented performance in various domains such as neural machine translation (Chen et al. 2018), question answering (Raffel et al. 2020), image classification (Parmar et al. 2019; Dosovitskiy et al. 2021), and time-series forecasting (Lim et al. 2021).

A typical Transformer can be thought as a stack of self-attention layers. Self-attention can be viewed as an inductive bias with graph-like dependencies among input tokens, where important elements hold stronger influence on the final output. On an algorithmic perspective, this is a three step process: (i) computing bidirectional attention scores between input tokens, (ii) element-wise encoding of input token sequences, and (iii) computing the weighted sum of encoded sequences with respect to attention scores.

Recently, Transformer parameter sizes are ever-increasing with current trends of exploiting self-supervised pre-training using massive data. This has introduced significant computational cost when employing Transformers; thus, recent re-

search has moved on to designing more computationally efficient attention mechanisms.

To alleviate the computational cost of self-attention mechanisms, two mainstream approaches have been taken in previous work, where each optimize the process in different dimensions. First is to improve the quadratic time complexity of vanilla dot-product attentions relative to the input sequence length. Since computing attention scores for an input of n elements takes $\mathcal{O}(n^2)$ time and memory (typical attention being bidirectional), applying Transformers to long input sequences can be a challenge. A number of techniques have been proposed in this dimension of optimization, such as leveraging fixed or learnable patterns to sparsify the attention matrix (Child et al. 2019; Beltagy, Peters, and Cohan 2020) and approximating attention matrices with mathematical methods (Wang et al. 2020b; Choromanski et al. 2021).

Second, there have also been efforts to reduce cost by trimming the model itself to reduce the complexity of Transformers. Popular approaches include quantization (Shen et al. 2020) and compression via knowledge distillation or pruning (Sanh et al. 2019; Xu et al. 2020). Leveraging neural architecture search to remove redundant operations (Chen et al. 2020; Guo et al. 2019) is also a common strategy. The efficiency benefits of these schemes are orthogonal to the first approach; thus, can be applied in parallel.

We note that existing work on reducing attention complexity generally focuses on modifying weights or changing its structure, where the optimization is essentially “model driven”. Therefore, the performance and resource trade-off can be considered static and tightly coupled to the model or the optimization method themselves. In this work, we present a novel dimension of attention optimization, with no model weight or structure modifications, rather putting focus on the statistical characteristics of the attention matrix. Specifically, we present the *Monte-Carlo Attention (MCA)*, as a flexible approach to reduce the computational cost of self-attention. MCA is motivated by the observation that attention scores are highly divergent due to the softmax operation, which prompts the question: “*why should we use equal precision for encoding all input tokens?*” By suppressing computation for elements that minimally (if at all) affect the final output, we can significantly reduce the attention complexity. MCA optimizes attention using randomized matrix multiplication approximation by configuring different error

bounds for the element-wise encoding of input sequences with respect to their *a priori* attention scores. We do so such that elements with weak attention scores are computed with relaxed precision, whereas errors for encoding salient elements are minimized. We discuss an optimal strategy for curtailing operational complexity with respect to the attention matrix and evaluate its theoretical error bounds.

MCA can be applied orthogonally with already-proposed optimization approaches to complement the overall model efficiency, and holds strong advantages itself in that it (i) allows simple dynamic control of performance-resource trade-off - allowing easy adjustment to different platforms and resource availability, (ii) has predictable and parametrizable upper error bounds independent to the input length - making it applicable for various inputs, and (iii) can be used as a drop-in replacement for naive attention mechanisms by not requiring additional model preprocessing or training.

We implement MCA with custom CUDA kernels and conduct extensive experiments using BERT (Devlin et al. 2019) and GLUE benchmark datasets. We also evaluate MCA’s performance when applied to other common efficient Transformer variants, such as DistilBERT and Longformer to demonstrate its compatibility with other optimization schemes in parallel. Our experimental results show that MCA can effectively reduce FLOPS in attention by a large margin (up to $11\times$) with negligible accuracy drop.

Background and Preliminaries

We first recall to the basics of self-attention mechanisms for our problem formulation, and introduce the random sampling-based matrix approximation algorithm, which covers the background of our proposed scheme. Note that we use Python style notations for matrix indexing. For matrix M , its i -th row and column are denoted as $M[i]$, $M[:, i]$ respectively, and the element in row i and column j as $M[i, j]$.

Self-Attention Mechanisms

After the successful debut of Transformer architectures, various innovations in self-attention mechanisms have emerged. A typical self-attention mechanism in Transformers start by computing an attention matrix for the input tokens: for an input sequence $X \in \mathbb{R}^{n \times d}$ of length n and feature dimension d , it first computes the attention matrix $A \in \mathbb{R}^{n \times n}$ based on X . The vanilla Transformer (Vaswani et al. 2017) employs scaled dot-product operation to compute $A = \text{softmax}(aQK^T)$ where $Q = XW_q$, $K = XW_k$ and $a = 1/\sqrt{d}$ being the scaling factor. $W_q, W_k \in \mathbb{R}^{d \times d}$ are trainable linear projections. We note that this operation requires $\mathcal{O}(n^2)$ complexity, which induces significant overhead when long inputs are passed. Thus, researchers have tried to tackle this issue and we point interested readers to Tay et al. 2020 for a comprehensive overview. Once the attention matrix is available, the next operation performs element-wise encoding (i.e., linear projection) of the input tokens: $H = XW$, where $W \in \mathbb{R}^{d \times d}$. Finally, the output $Y \in \mathbb{R}^{n \times d}$ is computed as product of the corresponding attention and encoded values $Y = AH$. Putting it all together,

the i -th output element $Y[i]$ can be expressed with respect to the input element $X[i]$ as:

$$Y[i] = \sum_{j=1}^n A[i, j]X[j]W \quad (1)$$

Note that for simplicity, we omit the symbols for multi-head attentions, and point out that the same property holds. Asymptotically, this computation requires $\mathcal{O}(nd^2)$ time and memory. In this work we are interested in reducing the computational complexity with respect to the feature dimension d , which is different from approaches that optimize attention score computation for input length n . In many practical settings, d (e.g., 768 in BERT) is much larger than n and this tendency holds as Transformers scale to larger networks (e.g., $d=3072$ in Megatron-LM (Shoeybi et al. 2019)). Thus, we argue that optimizing for $\mathcal{O}(d^2)$ rather than $\mathcal{O}(n^2)$ can have more impact on model performance when dealing with datasets with modest sequence lengths.

Approximating Matrix Multiplication via Monte-Carlo Sampling

Matrix multiplication approximation schemes are designed to allow formidable optimizations on large matrices by permitting some uncertainties in the final output, leading to accelerated output computation. There are several known methods for approximating matrix multiplication each with distinct statistical and computational properties. One family of algorithms is based on the low-rank approximation using truncated singular value decomposition (Drineas and Kannan 2001; Denton et al. 2014; Osawa et al. 2017), which well-preserves the low-rank structures of the output matrix. Unfortunately, given that they require matrix factorization as prerequisite, it is challenging to apply them in iterative processes. Fast Fourier Transform is an alternative for approximating matrix multiplication by representing the input column-row outer product as a polynomial multiplication, which can be effective when input matrices are sparse (Pagh 2013). In addition, kernelization-based methods demonstrated their potential in optimizing square matrix multiplications (Drineas and Kannan 2001). Finally, randomized algorithms (Drineas, Kannan, and Mahoney 2006; Eriksson-Bique et al. 2011), use random sampling to estimate the original outputs with sub-sampled information. This approach does not constrain the matrix shape or contents, while offering computing acceleration.

Based on these observations, we exploit a random sampling-based approach proposed by Drineas, Kannan, and Mahoney 2006, which, despite its random nature, allows for the easy control of the output matrix error bound and can be efficiently implemented on modern GPUs. Given an input matrix $A \in \mathbb{R}^{m \times k}$ and $B \in \mathbb{R}^{k \times n}$, the key idea is to view matrix multiplication as a *sum of the outer products* of columns A and their corresponding rows B . With this, we apply random sampling to a subset of these column-row pairs to compute an approximated output as the following:

$$AB = \sum_{i=1}^k A[:, i]B[i] \approx \frac{1}{r} \sum_{i=1}^r \frac{1}{p(s[i])} A[:, s[i]]B[s[i]] \quad (2)$$

$r \geq 1$ is the number of samples and $p(i)$ is the sampling probability of the i -th column-row pair. The sampling sequence $s \in \mathbb{Z}^r$ is generated via a random process $\Pr[s[k] = i] = p(i)$ in i.i.d. trials with replacement. The approximated output is unbiased such that $\mathbb{E}(\tilde{A}\tilde{B}) = AB$ holds given sufficient samples. Since this algorithm relies on random sampling, it is important to assure that the error $\mathbb{E}(\tilde{A}\tilde{B})$, is kept low.

This algorithm allows for linear complexity reduction from $\mathcal{O}(mkn)$ to $\mathcal{O}(mrn)$ with respect to the number of samples r , which corresponds to the feature dimension d when applied to attention mechanisms. Later, we adjust r to the attention score as a way to achieve a dynamic balance between the computing overhead and output precision.

While any probability distribution can be used for p_i , Drineas et al. showed that applying non-uniform distribution and biasing towards higher-rank terms yields tighter upper error bounds in terms of Frobenius norm compared to a uniform distribution. Specifically, Drineas et al. proved that:

$$\|AB - \tilde{A}\tilde{B}\|_F = \mathcal{O}\left(\frac{\|A\|_F \|B\|_F}{\sqrt{s}}\right), \quad (3)$$

when the $p(i)$ term is proportional to the L2 norm of each column-row pair multiplied together:

$$p(i) = \frac{\|A[:, i]\|_2 \|B[i]\|_2}{\sum_{j=0}^k \|A[:, j]\|_2 \|B[j]\|_2}. \quad (4)$$

For details on randomized linear algebra and its algorithms, we point interested readers to Mahoney 2011.

Monte-Carlo Attention

We now present our proposed Monte-Carlo Attention (MCA), a novel approximation algorithm for reducing the complexity of self-attention mechanisms in Transformers. MCA optimizes the element-wise encoding of input tokens (c.f., Equation (1)) by replacing the matrix-vector product with an *approximated* operation derived from random sampling. The principle idea is to apply different sampling counts for each output $Y[i]$ with respect to their attention scores to suppress overhead induced by lightly-weighted elements. Specifically, MCA can be expressed as follows:

$$Y[i] \approx \sum_{j=1}^n \sum_{k=1}^{r_j} \frac{A[i, j]X[i, s_j[k]]}{r_j p(s_j[k])} W[s_j[k]] \quad (5)$$

r_i and s_i each denote the sampling count and sampled indices for input $X[i]$, respectively, where $r_i = \text{len}(s_i)$. While most terms in the equation are known, two need to be properly defined. The first is $p(n)$, the sampling probability for each element. Despite being able to acquire the optimal probability from Equation 4, we claim that this formulation is impractical given that it requires auxiliary computation for

all new input X . Instead, MCA disconnects dependence towards the input in computing $p(i)$, which allows the computing to be a one-time process. The second unknown term is r_i , the number of samples taken from taken from the attention matrix A with respect to each input element $X[i]$. For this, MCA identifies the maximum attention score of each element towards other elements, roughly representing the minimum level of precision required for its computation. A high maximum attention score suggests that the currently evaluated element should have high precision (i.e., use more samples for encoding element $X[i]$) as it can heavily affect the output. On the contrary, a small maximum attention would indicate that the element minimally contributes to the output; thus can set a small r_i . Based on the theoretical upper error bounds in approximating $Y[i]$, we can derive A - r pair relationships which guarantee robust error concentration invariant to the attention and input length.

We note that computing these parameters with respect to the error bounds is important given that MCA is essentially a randomized algorithm. The following subsections discuss how $p(i)$ and r_i can be configured, and show their impact on the overall error. In a nutshell, we show that substituting Transformer attention modules to MCA, despite the randomness, guarantees tight error bounds for the attention layers.

Formulation of Sampling Probability

The intuition behind exploiting non-uniform distributions for sampling column-row pairs is supported by the observation that column-row pairs with relatively small norms will have minimal impact to the output. Thus, configuring a sampling probability proportional to the norm can minimize output variance. The distribution presented in Equation (4) is known to be optimal in terms of minimizing the output Frobenius norm error (Drineas and Kannan 2001). Unfortunately, this formulation can be (practically) inefficient, as p must be computed freshly for all incoming inputs. Furthermore, as Transformers consist of multiple self-attention layers, computing p can induce an additional $\mathcal{O}(nd)$ complexity for each layer. Therefore, determining the sampling distribution only from the attention weights, and eliminating the use of input X , is a much more computationally-efficient strategy allowing a one-time p computation. MCA borrows findings from Drineas, Kannan, and Mahoney 2006 that suggests that the optimal probability can be approximated even when only one of two matrices is known. This, when applied to our domain, frees the dependency towards the input matrix. Specifically, we compute the sampling probability by taking $W^T W$ for approximation, which yields:

$$p(i) = \frac{\|W[i]\|_2^2}{\sum_{j=0}^d \|W[j]\|_2^2} = \frac{\|W[i]\|_2^2}{\|W\|_F^2} \quad (6)$$

While this asymptotically exhibits $\mathcal{O}(d^2)$ complexity, in practice, computed results can be embedded in the model or cached to reduce the overhead. With this probability, the maximum Frobenius norm error becomes proportional to the norm of the attention weights and input, scaled by the square of the sample size.

Lemma 1 Let $H[i] \in \mathbb{R}^d$ be an approximation to $X[i]W$ using $r_i \in \mathbb{Z}^+$ random samples with the probability distribution discussed in Equation (6). Then,

$$\mathbb{E}[\|H[i] - X[i]W\|_F] \leq \frac{1}{\sqrt{r_i}} \|X[i]\|_2 \|W\|_F. \quad (7)$$

Determining Sample Size from Attention

Self-attention mechanisms are designed to be bidirectional: meaning that each element in a sequence possess its own perspectives to all other elements. Hence, for any element, there exists n different attention scores to describe its importance within a sequence. To determine a proper sample size r_i , given an element and its outbound attention scores, we take the maximum of all available scores (for each element) in defining the minimum precision. This assures high precision (i.e., more samples) for highly weighted elements and low precision (i.e., less samples) for those with low weights. When done so, based on Lemma 1, the mean error bound can be expressed with respect to output $Y[i]$ via Equation (1):

$$\mathbb{E}[\|\tilde{Y}[i] - Y[i]\|_F] \leq \sum_{j=1}^n \frac{A[i, j]}{\sqrt{r_i}} \|X[i]\|_2 \|W\|_F. \quad (8)$$

When taking $\sqrt{r_j} = \max A[:, j]$, the attention term in Equation (8) can be eliminated from inequality as $\frac{A[i, j]}{\max A[:, j]} \leq 1$ for all i and j (assuming A as non-negative). Such a strategy of taking the maximum of all attention scores can be considered conservative. Potentially, one can define more aggressive schemes by taking the mean or median of the scores. Unfortunately, analyzing the theoretical error bounds for such cases can be challenging, as error becomes dependent on the statistical characteristics of the attention matrix. We leave such optimizations for future work.

The formulation above indicates that the output error, caused from the randomness of MCA, is dependent on input sequence length n ; thus, the error may seem to increase with long inputs. This is not desirable when bounding the error. To mitigate this, we consider an additional n term in r_i . As a result, r_i is calculated as:

$$\sqrt{r_i} = \frac{n \cdot \max A[:, i]}{\alpha} \quad (9)$$

Here, $\alpha \in (0, 1)$ is the attention error coefficient, a user-controllable parameter to linearly scale error bounds. $\alpha = 1$ indicates a performance-oriented configuration (i.e., reduced complexity with less samples), whereas lower α results in extra precision at the cost of additional computational cost. Still, we show in our evaluations that even a small α (e.g., 0.2) can result in significantly reduced computation.

Finally, we prove that this formulation allows robust concentration of output errors independent to the input sequence and the attention matrix statistics.

Theorem 2 Let $\tilde{Y}[i]$ be the estimate of $Y[i]$ by equation (5), and define β as mean Euclidean norm of $X[i]$ for all i (i.e., $\beta = \frac{1}{n} \sum_{i=1}^n \|X[i]\|_2$). Assume A is positive matrix for which $A[i, j] > 0$ for all i, j . Then

$$\mathbb{E}[\|\tilde{Y}[i] - Y[i]\|_F] \leq \alpha \beta \|W\|_F. \quad (10)$$

Furthermore, suppose $\delta \in (0, 1)$. Then with probability at least $1 - \delta$,

$$\|\tilde{Y}[i] - Y[i]\|_F \leq \frac{\alpha \beta}{\delta} \|W\|_F. \quad (11)$$

Therefore, we show that despite MCA’s randomness, with proper α , the overall error can still be tightly bounded; thus, achieve efficient encoding with minimal loss in accuracy.

Experiments

We implement MCA via custom CUDA kernels, and measure its computation complexity and accuracy as performance metrics. We replace the multi-head attention in BERT (Vaswani et al. 2017) with MCA, and test the performance with GLUE benchmark (Wang et al. 2018). We also show using DistilBERT and Longformer that MCA can co-exist with existing attention optimization methods.

Transformer Models

BERT is one of the earliest applications of the Transformer architecture that demonstrated unprecedented performance in a wide spectrum of NLP tasks. We use the BERT_{BASE} model, which encompasses 12 layers of multi-head attention with 12 heads and a feature dimension of 768.

DistilBERT (Sanh et al. 2019) is a variant of BERT which targets to reduce the computational cost by compression the Transformer’s model weights using knowledge distillation (Hinton, Vinyals, and Dean 2015). DistilBERT holds the same general structure as BERT (with $\frac{1}{2}$ of the attention layers), but the token type embedding and poolers are removed from the layers. We use this model to show that MCA can coalesce with existing model compression techniques.

Longformer (Beltagy et al. 2020) is a Transformer network suitable for lengthy documents by addressing the issue of complexity-blowup (i.e., $\mathcal{O}(n^2)$) with increasing input sequence lengths. Specifically, Longformers sparsify the attention matrix by employing a fixed-size window w attention surrounding each input element. This constrains the complexity to $\mathcal{O}(n \times w)$, scaling linearly with input length n .

Benchmark Datasets

GLUE Benchmark. To consider a broad range of input data, we adopt the General Language Understanding Evaluation (GLUE) benchmark. GLUE is commonly used in evaluating the performance of language models. It consists of nine individual task sets, with two tasks for single-sentence classification (CoLA and SST-2), three tasks for semantic similarity (MPRC, QQP, and STS-B), and four natural language inference tasks (MNLI, RTE, QNLI, and WNLI).

Document Classification. To evaluate MCA with the Longformer model, we use the arXiv Academic Paper dataset (AAPD) (Yang et al. 2016), the IMDB review classification dataset and the Hyperpartisan News Detection (HND) dataset (Kiesel et al. 2019). AAPD consists average of 167 tokens for each entry while IMDB has entries with 300 tokens. The HND dataset includes elements with the longest item size, with an average of 705 tokens for each article.

Task	Baseline		$\alpha=0.2$		$\alpha=0.4$		$\alpha=0.6$		$\alpha=1.0$	
	Metric	Result	Result	FLOPS	Result	FLOPS	Result	FLOPS	Result	FLOPS
CoLA	MC	53.74	53.74 \pm 0.1	11.44 \times	53.90 \pm 0.2	12.62 \times	53.78 \pm 0.4	14.28 \times	50.95 \pm 0.7	18.38 \times
SST-2	Acc.	92.43	92.26 \pm 0.0	5.58 \times	92.04 \pm 0.0	6.71 \times	91.22 \pm 0.0	8.13 \times	80.66 \pm 0.1	12.34 \times
MRPC	Acc.	84.55	84.36 \pm 0.1	3.04 \times	83.77 \pm 0.3	3.80 \times	82.28 \pm 0.4	4.62 \times	65.95 \pm 0.8	7.14 \times
	F1	89.41	88.96 \pm 0.0	3.04 \times	88.86 \pm 0.2	3.80 \times	87.69 \pm 0.3	4.62 \times	73.50 \pm 0.7	7.14 \times
STS-B	PC	88.04	87.84 \pm 0.3	4.80 \times	82.69 \pm 0.9	5.95 \times	67.90 \pm 1.3	7.35 \times	33.55 \pm 1.4	12.40 \times
	SC	87.63	87.21 \pm 0.3	4.80 \times	81.04 \pm 0.8	5.95 \times	67.54 \pm 1.2	7.35 \times	28.53 \pm 1.3	12.40 \times
QQP	Acc.	90.90	90.89 \pm 0.0	4.76 \times	90.74 \pm 0.1	5.73 \times	89.98 \pm 0.4	6.89 \times	78.20 \pm 0.6	10.30 \times
	F1	87.81	87.79 \pm 0.0	4.76 \times	87.63 \pm 0.1	5.73 \times	86.55 \pm 0.3	6.89 \times	73.61 \pm 0.6	10.30 \times
MNLI	Pos.	83.60	83.50 \pm 0.1	3.75 \times	82.60 \pm 0.2	4.76 \times	78.62 \pm 0.8	5.88 \times	71.43 \pm 1.1	9.61 \times
	Neg.	84.79	84.72 \pm 0.1	3.75 \times	83.83 \pm 0.2	4.76 \times	79.50 \pm 0.8	5.88 \times	72.05 \pm 1.2	9.61 \times
QNLI	Acc.	91.54	91.47 \pm 0.0	3.03 \times	90.24 \pm 0.0	3.83 \times	84.33 \pm 0.1	4.94 \times	56.43 \pm 0.1	10.64 \times
RTE	Acc.	72.56	71.68 \pm 0.2	2.50 \times	70.39 \pm 0.4	3.24 \times	64.72 \pm 0.8	4.55 \times	52.65 \pm 1.1	10.06 \times
WNLI	Acc.	56.33	56.33 \pm 0.0	4.08 \times	56.32 \pm 0.1	5.15 \times	54.96 \pm 0.6	6.49 \times	52.69 \pm 1.1	11.09 \times

Table 1: FLOPS reduction and its corresponding model accuracy (with 95% confidence intervals) when MCA-BERT is tested on the GLUE benchmark with different error bound coefficients α . MC, PC, and SC each denote Matthews, Pearson, and Spearman correlation coefficients.

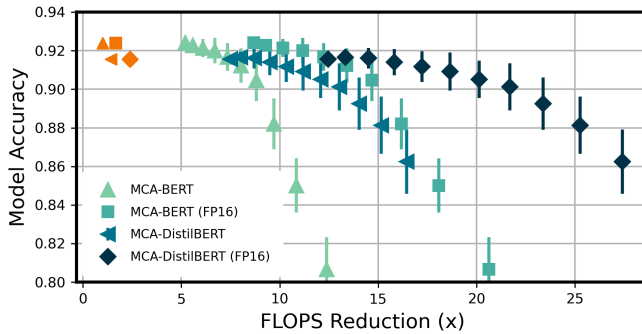


Figure 1: Mean accuracy and FLOPS trade-off for MCA-BERT and MCA-DistilBERT, along with FP16 quantized versions (SST-2 dataset). The orange plots on the top-left show the accuracy-FLOPS relationship for original BERT and DistilBERT (without MCA applied).

Implementation and Experimental Setup

MCA is implemented as a C++ PyTorch extension using custom CUDA kernels. Our implementation replaces the self-attention in Transformer models, and can be configured whether to run in regular mode (i.e., act as a normal self-attention) or approximation mode (i.e., activate MCA). We employ AWS p3.2xlarge instances for the experiment, which serves one NVIDIA V100 GPU and use the Huggingface library for implementations and download pre-trained weights for BERT, DistilBERT, and Longformer. Readers can refer to Appendix 1 for reproducibility details.

FLOPS Reduction and Accuracy Impact

Using the GLUE benchmark, we first evaluate the performance of MCA when applied to BERT. Specifically, we measure the Floating Operations Per Second (FLOPS) of

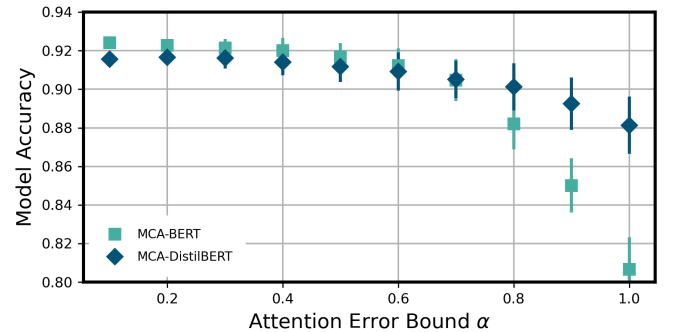


Figure 2: Impact of attention error bound α on model accuracy of MCA-BERT and MCA-DistilBERT (SST-2 dataset). The vertical bars represents 95% confidence interval.

the MCA-BERT attention operation and compare with when MCA is not applied. Note that we measure only the FLOPS for the attention (i.e., AXW) and exclude other Transformer operations (e.g., embedding layer, classification heads) as they will be consistent. We also report the full model accuracy to observe the complexity-accuracy trade-off. We compute 128 samples for each task (different random seeds) and report the mean with 95% confidence interval. Given that α is an adjustable parameter that determines the attention error bound tightness, we evaluate for $\alpha = 0.2, 0.4, 0.6$, and 1.0 .

Table 1 summarizes our results. Here, MCA shows a significant FLOPS reduction for all GLUE tasks. Using $\alpha = 0.2$, MCA exhibits 4.64 \times FLOPS reduction (mean of all tests), with negligible accuracy loss. With increasing α , we see further model complexity reduction. For example, with $\alpha = 0.4$, we observe 5.72 \times average FLOPS reduction with $\leq 1\%$ accuracy loss for all cases. A more generous $\alpha \geq 0.6$, offers even more drastic improvements in complexity reduction. However, this comes at the cost of accuracy drop,

some quite noticeable, suggesting that α should be configured with respect to the target task and dataset.

One important observation is that the FLOPS reduction, while maintaining consistent trends, varies with respect to the dataset. For example, at $\alpha = 0.2$, the CoLA task shows a $11.44\times$ FLOPS reduction compared to $2.50\times$ of RTE. Given that MCA reduces complexity by exploiting smaller sample sizes for elements with low attention scores, we conjecture that differences can rise from the fact that the attention matrix from CoLA is sparser than that of RTE. Generally, binary classification tasks (e.g., CoLA and SST-2) exhibit higher reduction rates while semantic similarity tasks and NLI tasks show modest reduction rates.

Another point to note is that previously proposed schemes that compress or prune the model require the tuning of many different hyperparameters. On the contrary we argue that the effort needed for per-task α adjustment in MCA is minimal, as this is the only hyperparameter to adjust when applying on different datasets or computing environments.

Performance-Accuracy Trade-off. Figure 1 visualizes the required FLOPS and model accuracy relationships of BERT, DistilBERT, MCA-BERT, and MCA-DistilBERT using the SST-2 dataset (top-left orange plots are cases with MCA not applied). We also present results for the corresponding 16-bit floating point (FP16) quantized models of each case, respectively. As the results show, MCA results in 40-60% FLOPS reduction with (nearly) the same accuracy and as much as 170% reduction with small accuracy loss. Nevertheless, if the FLOPS drops below a tipping point, MCA shows noticeable loss in accuracy, showing a logarithmic tradeoff relationship. Figure 1 also shows that performance trends of MCA holds for models with quantized weights as well.

Attention Error vs. Model Error. Figure 2 presents the impact of the attention error bound α on the model accuracy for BERT and DistilBERT with MCA applied when using the SST-2 dataset. As the plots show, with increasing α (higher errors at the attention), the overall model accuracy shows a gradually decreasing pattern. Nevertheless, configuring α to 0.4, or even 0.6 will only minimally impact the overall model accuracy. Comprehensively with the results reported in Table 1 we can see that despite the small loss in accuracy, the computational complexity reduction gain is significant.

Integration with Compressed Transformers

We next evaluate the impact of applying MCA on already compressed Transformer models. For this, we perform the same experiment as Table 1 on DistilBERT, which only possess 40% of BERT’s original parameter size. We present results in Table 2. Overall, results from MCA-DistilBERT mostly agree with those seen for MCA-BERT in terms of FLOPS reduction and accuracy. For more a detailed view, we look back on Figure 1 which plots a finer-grained FLOPS-accuracy relationship for both MCA-applied models (using SST-2), which shows that MCA effectively reduces the complexity even for already compressed models.

From Figures 1 and 2, we see that the accuracy drop with increasing α is more gradual for MCA-DistilBERT compared to MCA-BERT. Since, MCA is a randomized approach, while tightly bounded, a small amount of error accu-

mulates with more attention layers to process. With MCA-DistilBERT having only $\frac{1}{2}$ attention layers of MCA-BERT, we conjecture that MCA-DistilBERT benefits by having less errors accumulated.

Integration with Sparse Attention Patterns

To empirically validate the coalesced performance of MCA and Transformers with sparse attention mechanisms, we integrate MCA to the Longformer model and test its performance on the document classification task datasets. We configure the Longformer with a window size of 256, and apply global attention for the CLS token.

As results in Table 3 show, for all datasets, MCA achieved FLOPS reduction by a factor of at least $3\times$ with negligible error (for tight error bounds). The tendency of FLOPS-accuracy trade-off for the Longformer experiments also agree with those from prior observations. Comprehensively, these results suggest that MCA can be generally applied with different Transformer architecture variants.

Related Work

A number of previous work target to alleviate the computational cost of self-attention mechanisms, mostly in two orthogonal ways: (i) by optimizing for long input sequences, and (ii) by proposing new Transformer designed that are less resource demanding. We introduce examples below.

Efficient Attention for Long Sequences

The scaled dot-product attention in a vanilla Transformer induces significant overhead when dealing with lengthy data. Thus, a number of previous work have focused on optimizing the evaluation algorithms for attention scores as a way to efficiencize Transformers (Tay et al. 2021).

Fixed Sparse Patterns. Some early efforts exploited sparse attention with fixed patterns, which reduce the attention complexity proportional to a target sparsity. As an example, blocked attention (Qiu et al. 2020; Parmar et al. 2019) groups input elements in blocks for block-wise attention evaluation. More recently, strided-pattern attention has been proposed (Child et al. 2019; Beltagy, Peters, and Cohan 2020). These work limit attention evaluation to fixed-sized windows and evaluates only the neighboring elements.

Learnable Sparse Patterns. Methods that learn the access pattern have also been proposed. Reformer (Kitaev, Kaiser, and Levskaya 2020) replaces dot-product attention with locally-sensitive hashing, which clusters input tokens based on hash similarity and filters prominent tokens. Online K-means clustering (Roy et al. 2021) can be used as a similar concept by learning the input sparsity patterns.

Low-Rank Approximations. Some work directly approximate attention scores by assuming a low-rank structured attention matrix (Wang et al. 2020b). More recently, Performer (Choromanski et al. 2021) exploited kernelization for approximation, which does not rely on input sparsity patterns or low-rankness.

Low-cost Transformers

There have also been efforts to downscale the time complexity of Transformers themselves as we detail below.

Task	Baseline		$\alpha=0.2$		$\alpha=0.4$		$\alpha=0.6$		$\alpha=1.0$	
	Metric	Result	Result	FLOPS	Result	FLOPS	Result	FLOPS	Result	FLOPS
CoLA	MC	56.85	56.49 \pm 0.1	11.36 \times	55.69 \pm 0.3	12.33 \times	54.29 \pm 0.3	13.75 \times	50.08 \pm 0.6	17.30 \times
SST-2	Acc.	91.51	91.50 \pm 0.0	5.60 \times	91.41 \pm 0.0	6.64 \times	90.92 \pm 0.0	7.84 \times	88.20 \pm 0.1	10.58 \times
MRPC	Acc.	85.78	85.26 \pm 0.1	2.98 \times	83.83 \pm 0.2	3.73 \times	78.59 \pm 0.4	4.54 \times	64.76 \pm 0.5	7.35 \times
	F1	90.26	89.91 \pm 0.1	2.98 \times	88.99 \pm 0.2	3.73 \times	85.66 \pm 0.2	4.54 \times	76.04 \pm 0.5	7.35 \times
STS-B	PC	88.39	88.29 \pm 0.2	4.65 \times	79.73 \pm 0.3	5.91 \times	63.11 \pm 0.4	7.21 \times	30.12 \pm 0.7	12.33 \times
	SC	87.56	86.94 \pm 0.2	4.65 \times	78.24 \pm 0.3	5.91 \times	62.40 \pm 0.4	7.21 \times	29.02 \pm 0.6	12.33 \times
QQP	Acc.	88.24	88.11 \pm 0.0	4.75 \times	87.79 \pm 0.2	5.85 \times	86.77 \pm 0.5	7.14 \times	72.24 \pm 0.9	10.41 \times
	F1	85.48	84.99 \pm 0.0	4.75 \times	83.26 \pm 0.2	5.85 \times	82.95 \pm 0.8	7.14 \times	70.02 \pm 1.3	10.41 \times
MNLI	Pos.	82.50	82.37 \pm 0.1	3.62 \times	81.09 \pm 0.4	4.31 \times	77.70 \pm 0.6	5.09 \times	68.11 \pm 0.8	7.40 \times
	Neg.	83.83	83.66 \pm 0.1	3.62 \times	82.41 \pm 0.4	4.31 \times	79.26 \pm 0.6	5.09 \times	71.25 \pm 0.8	7.40 \times
QNLI	Acc.	88.48	88.31 \pm 0.0	2.86 \times	87.06 \pm 0.0	3.51 \times	83.19 \pm 0.1	4.28 \times	62.89 \pm 0.2	7.43 \times
RTE	Acc.	65.70	65.22 \pm 0.3	2.33 \times	64.43 \pm 0.6	2.91 \times	60.35 \pm 0.8	3.64 \times	52.79 \pm 0.9	6.47 \times
WNLI	Acc.	56.33	54.40 \pm 0.6	3.97 \times	52.55 \pm 1.4	5.10 \times	52.77 \pm 1.5	6.56 \times	52.24 \pm 1.9	10.55 \times

Table 2: Experiment results on MCA-DistilBERT. Experimental configurations are kept same as Table 1.

Task	Baseline		$\alpha=0.2$		$\alpha=0.4$		$\alpha=0.6$		$\alpha=1.0$	
	Metric	Result	Result	FLOPS	Result	FLOPS	Result	FLOPS	Result	FLOPS
AAPD	Acc.	74.81	74.70 \pm 0.1	3.44 \times	72.21 \pm 0.2	4.45 \times	69.68 \pm 0.2	6.22 \times	66.87 \pm 0.3	9.83 \times
	F1	71.39	71.02 \pm 0.1	3.44 \times	68.33 \pm 0.2	4.45 \times	65.35 \pm 0.3	6.22 \times	64.14 \pm 0.3	9.83 \times
HND	Acc.	80.33	80.18 \pm 0.0	5.32 \times	78.73 \pm 0.1	6.81 \times	74.23 \pm 0.3	7.93 \times	70.36 \pm 0.5	11.81 \times
	F1	76.57	76.32 \pm 0.0	5.32 \times	74.53 \pm 0.1	6.81 \times	72.94 \pm 0.3	7.93 \times	68.53 \pm 0.4	11.30 \times
IMDB	Acc.	89.11	89.05 \pm 0.0	3.67 \times	84.10 \pm 0.1	5.06 \times	81.17 \pm 0.3	7.50 \times	75.11 \pm 0.4	11.30 \times

Table 3: Experiment results on MCA-Longformer. A local attention with window size of 256 is used for the Longformer model.

Model Compression. Pruning redundant information in Transformer weights (Gordon, Duh, and Andrews 2020; Wang, Wohlwend, and Lei 2020) and less-important heads (Michel, Levy, and Neubig 2019) have shown good performance albeit its complicated post-training procedure. Another approach is to compress the model via knowledge distillation (Beltagy, Peters, and Cohan 2020), allowing flexible model sizes at the cost of (potentially heavy) re-training. **Weight Quantization.** Along with FP16 quantization, more aggressive quantization schemes, such as 8-bit integers (Zafir et al. 2019) or even down to 4- or 2-bits (Shen et al. 2020), have been applied to Transformers. While this line of research does not often involve heavy post-processing (e.g., additional model training), such low precision schemes, in practice, often require hardware accelerators to maximize the desired efficiency.

Neural Architecture Search (NAS). Employing NAS downsizes Transformer sizes by identifying more efficient model architectures. This can be very effective when targeting a specific hardware architecture (Wang et al. 2020a) or when meeting target latency goals (Guo et al. 2019).

Positioning MCA with Previous Work

Despite some limitations, we acknowledge that active research in both perspectives are meaningful. As discussed, MCA takes a different approach, by exploiting the statistical characteristics of attention matrices. Such an approach can

be considered to be orthogonal to the previously proposed improvements, and we showed through our evaluations that MCA can indeed operate in parallel with other computation acceleration techniques for Transformers.

Conclusion

This work presented the Monte-Carlo Attention (MCA) mechanism, a randomized approximation method for reducing the complexity of self-attention mechanisms in Transformers. MCA is designed under the philosophy that not all attention elements need to be treated with equal weight. Those with heavy attention should be computed with high precision, whereas the weaker, which give minimal impact to the final output, can use lower precision. For this MCA employs random sampling-based matrix-vector product operations to perform element-wise encoding of input tokens, where each element is allocated different amounts of samples with respect to their attention score. We evaluated the model complexity reduction performance (in terms of FLOPS) and the model accuracy using BERT (and its variations) with GLUE benchmark datasets. Our results suggest that MCA is an effective alternative to the attention operators used in today’s Transformer networks by reducing the FLOPS by up to 11 \times with negligible (near-zero) accuracy drop. Furthermore, we showed that MCA can be applied in parallel with widely-used approaches to efficiencize Transformers, while still maintaining its performance benefits.

Acknowledgements

This work was supported by the National Research Foundation of Korea (NRF) Grant funded by Ministry of Science and ICT (MSIT) Basic Science Research Program (2021R1A2C4002380), Information Technology Research Center (ITRC) Support Program supervised by the Institute for Information and Communications Technology Planning and Evaluation (IITP-2020-0-01461), and by the Ministry of Culture, Sports and Tourism and Korea Creative Content Agency (R2021040018).

References

- Beltagy, I.; Peters, M. E.; and Cohan, A. 2020. Longformer: The Long-Document Transformer. *CoRR*, abs/2004.05150.
- Chen, D.; Li, Y.; Qiu, M.; Wang, Z.; Li, B.; Ding, B.; Deng, H.; Huang, J.; Lin, W.; and Zhou, J. 2020. AdaBERT: Task-Adaptive BERT Compression with Differentiable Neural Architecture Search. In Bessiere, C., ed., *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, 2463–2469. ijcai.org.
- Chen, M. X.; Firat, O.; Bapna, A.; Johnson, M.; Macherey, W.; Foster, G. F.; Jones, L.; Schuster, M.; Shazeer, N.; Parmar, N.; Vaswani, A.; Uszkoreit, J.; Kaiser, L.; Chen, Z.; Wu, Y.; and Hughes, M. 2018. The Best of Both Worlds: Combining Recent Advances in Neural Machine Translation. In Gurevych, I.; and Miyao, Y., eds., *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, 76–86. Association for Computational Linguistics.
- Child, R.; Gray, S.; Radford, A.; and Sutskever, I. 2019. Generating Long Sequences with Sparse Transformers. *CoRR*, abs/1904.10509.
- Choromanski, K. M.; Likhoshesterov, V.; Dohan, D.; Song, X.; Gane, A.; Sarlós, T.; Hawkins, P.; Davis, J. Q.; Mohiuddin, A.; Kaiser, L.; Belanger, D. B.; Colwell, L. J.; and Weller, A. 2021. Rethinking Attention with Performers. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Denton, E. L.; Zaremba, W.; Bruna, J.; LeCun, Y.; and Fergus, R. 2014. Exploiting Linear Structure Within Convolutional Networks for Efficient Evaluation. In Ghahramani, Z.; Welling, M.; Cortes, C.; Lawrence, N. D.; and Weinberger, K. Q., eds., *Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, 1269–1277.
- Devlin, J.; Chang, M.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, 4171–4186. Association for Computational Linguistics.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; Uszkoreit, J.; and Houlsby, N. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Drineas, P.; and Kannan, R. 2001. Fast Monte-Carlo Algorithms for Approximate Matrix Multiplication. In *42nd Annual Symposium on Foundations of Computer Science, FOCS 2001, 14-17 October 2001, Las Vegas, Nevada, USA*, 452–459. IEEE Computer Society.
- Drineas, P.; Kannan, R.; and Mahoney, M. W. 2006. Fast Monte Carlo Algorithms for Matrices I: Approximating Matrix Multiplication. *SIAM J. Comput.*, 36(1): 132–157.
- Eriksson-Bique, S. D.; Solbrig, M.; Stefanelli, M.; Warkentin, S.; Abbey, R.; and Ipsen, I. C. F. 2011. Importance Sampling for a Monte Carlo Matrix Multiplication Algorithm, with Application to Information Retrieval. *SIAM J. Sci. Comput.*, 33(4): 1689–1706.
- Gordon, M. A.; Duh, K.; and Andrews, N. 2020. Compressing BERT: Studying the Effects of Weight Pruning on Transfer Learning. In Gella, S.; Welbl, J.; Rei, M.; Petroni, F.; Lewis, P. S. H.; Strubell, E.; Seo, M. J.; and Hajishirzi, H., eds., *Proceedings of the 5th Workshop on Representation Learning for NLP, RepL4NLP@ACL 2020, Online, July 9, 2020*, 143–155. Association for Computational Linguistics.
- Guo, Y.; Zheng, Y.; Tan, M.; Chen, Q.; Chen, J.; Zhao, P.; and Huang, J. 2019. NAT: Neural Architecture Transformer for Accurate and Compact Architectures. In Wallach, H. M.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox, E. B.; and Garnett, R., eds., *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 735–747.
- Hinton, G. E.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. *CoRR*, abs/1503.02531.
- Kiesel, J.; Mestre, M.; Shukla, R.; Vincent, E.; Adineh, P.; Corney, D. P. A.; Stein, B.; and Potthast, M. 2019. SemEval-2019 Task 4: Hyperpartisan News Detection. In May, J.; Shutova, E.; Herbelot, A.; Zhu, X.; Apidianaki, M.; and Mohammad, S. M., eds., *Proceedings of the 13th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2019, Minneapolis, MN, USA, June 6-7, 2019*, 829–839. Association for Computational Linguistics.
- Kitaev, N.; Kaiser, L.; and Levskaya, A. 2020. Reformer: The Efficient Transformer. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Lim, B.; Arik, S. Ö.; Loeff, N.; and Pfister, T. 2021. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*.
- Mahoney, M. W. 2011. Randomized algorithms for matrices and data. *arXiv preprint arXiv:1104.5557*.
- Michel, P.; Levy, O.; and Neubig, G. 2019. Are Sixteen Heads Really Better than One? In Wallach, H. M.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox,

- E. B.; and Garnett, R., eds., *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 14014–14024.
- Osawa, K.; Sekiya, A.; Naganuma, H.; and Yokota, R. 2017. Accelerating Matrix Multiplication in Deep Learning by Using Low-Rank Approximation. In *2017 International Conference on High Performance Computing & Simulation, HPCS 2017, Genoa, Italy, July 17-21, 2017*, 186–192. IEEE.
- Pagh, R. 2013. Compressed matrix multiplication. *ACM Trans. Comput. Theory*, 5(3): 9:1–9:17.
- Park, H.; Lee, Y.; and Ko, J. 2021. Enabling Real-time Sign Language Translation on Mobile Platforms with On-board Depth Cameras. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 5(2): 77:1–77:30.
- Park, N.; Lee, T.; and Kim, S. 2021. Vector Quantized Bayesian Neural Network Inference for Data Streams. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI , Virtual Event, February 2-9, 2021*, 9322–9330. AAAI Press.
- Parmar, N.; Ramachandran, P.; Vaswani, A.; Bello, I.; Levskaya, A.; and Shlens, J. 2019. Stand-Alone Self-Attention in Vision Models. In Wallach, H. M.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox, E. B.; and Garnett, R., eds., *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 68–80.
- Qiu, J.; Ma, H.; Levy, O.; Yih, W.; Wang, S.; and Tang, J. 2020. Blockwise Self-Attention for Long Document Understanding. In Cohn, T.; He, Y.; and Liu, Y., eds., *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, 2555–2565. Association for Computational Linguistics.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.*, 21: 140:1–140:67.
- Roy, A.; Saffar, M.; Vaswani, A.; and Grangier, D. 2021. Efficient Content-Based Sparse Attention with Routing Transformers. *Trans. Assoc. Comput. Linguistics*, 9: 53–68.
- Sanh, V.; Debut, L.; Chaumond, J.; and Wolf, T. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108.
- Shen, S.; Dong, Z.; Ye, J.; Ma, L.; Yao, Z.; Gholami, A.; Mahoney, M. W.; and Keutzer, K. 2020. Q-BERT: Hessian Based Ultra Low Precision Quantization of BERT. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI New York, NY, USA, February 7-12, 2020*, 8815–8821. AAAI Press.
- Shoeybi, M.; Patwary, M.; Puri, R.; LeGresley, P.; Casper, J.; and Catanzaro, B. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*.
- Tay, Y.; Bahri, D.; Metzler, D.; Juan, D.; Zhao, Z.; and Zheng, C. 2021. Synthesizer: Rethinking Self-Attention for Transformer Models. In Meila, M.; and Zhang, T., eds., *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, 10183–10192. PMLR.
- Tay, Y.; Dehghani, M.; Bahri, D.; and Metzler, D. 2020. Efficient Transformers: A Survey. *CoRR*, abs/2009.06732.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is All you Need. In Guyon, I.; von Luxburg, U.; Bengio, S.; Wallach, H. M.; Fergus, R.; Vishwanathan, S. V. N.; and Garnett, R., eds., *Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, 5998–6008.
- Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In Linzen, T.; Chrupala, G.; and Alishahi, A., eds., *Proceedings of the Workshop: Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2018, Brussels, Belgium, November 1, 2018*, 353–355. Association for Computational Linguistics.
- Wang, H.; Wu, Z.; Liu, Z.; Cai, H.; Zhu, L.; Gan, C.; and Han, S. 2020a. HAT: Hardware-Aware Transformers for Efficient Natural Language Processing. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J. R., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, 7675–7688. Association for Computational Linguistics.
- Wang, S.; Li, B. Z.; Khabsa, M.; Fang, H.; and Ma, H. 2020b. Linformer: Self-Attention with Linear Complexity. *CoRR*, abs/2006.04768.
- Wang, Z.; Wohlwend, J.; and Lei, T. 2020. Structured Pruning of Large Language Models. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, 6151–6162. Association for Computational Linguistics.
- Xu, C.; Zhou, W.; Ge, T.; Wei, F.; and Zhou, M. 2020. BERT-of-Theseus: Compressing BERT by Progressive Module Replacing. In Webber, B.; Cohn, T.; He, Y.; and Liu, Y., eds., *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, 7859–7869. Association for Computational Linguistics.
- Yang, Z.; Yang, D.; Dyer, C.; He, X.; Smola, A. J.; and Hovy, E. H. 2016. Hierarchical Attention Networks for Document Classification. In Knight, K.; Nenkova, A.; and Rambow, O., eds., *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, 1480–1489. The Association for Computational Linguistics.
- Zafir, O.; Boudoukh, G.; Izsak, P.; and Wasserblat, M. 2019. Q8BERT: Quantized 8Bit BERT. In *Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing - NeurIPS Edition, EMC2@NeurIPS 2019, Vancouver, Canada, December 13, 2019*, 36–39. IEEE.