

Compressing Large-Scale Transformer-Based Models: A Case Study on BERT

Prakhar Ganesh¹, Yao Chen¹, Xin Lou¹, Mohammad Ali Khan¹, Yin Yang²,
Deming Chen³, Marianne Winslett³, Hassan Sajjad^{4,2} and Preslav Nakov^{4,2}

¹Advanced Digital Sciences Center

²Hamad Bin Khalifa University

³University of Illinois at Urbana-Champaign

⁴Qatar Computing Research Institute

{prakhar.g, yao.chen, lou.xin, mohammad.k}@adsc-create.edu.sg,
{yyang, hsajjad, pnakov}@hbku.edu.qa, {dchen, winslett}@illinois.edu

Abstract

Transformer-based models pre-trained on large-scale corpora achieve state-of-the-art accuracy for natural language processing tasks, but are too resource-hungry and compute-intensive to suit low-capability devices or applications with strict latency requirements. One potential remedy is model compression, which has attracted extensive attention. This paper summarizes the branches of research on compressing Transformers, focusing on the especially popular BERT model. BERT’s complex architecture means that a compression technique that is highly effective on one part of the model, *e.g.*, attention layers, may be less successful on another part, *e.g.*, fully connected layers. In this systematic study, we identify the state of the art in compression for each part of BERT, clarify current best practices for compressing large-scale Transformer models, and provide insights into the inner workings of various methods. Our categorization and analysis also shed light on promising future research directions for achieving a lightweight, accurate, and generic natural language processing model.

1 Introduction

Sentiment analysis, paraphrase detection, machine reading comprehension, text summarization, question answering: all these natural language processing (NLP) tasks benefit from pre-training a large-scale generic model on an enormous corpus such as a Wikipedia dump and/or a book collection, and then fine-tuning the resulting model for these specific downstream tasks. Earlier solutions following this methodology used recurrent neural networks (RNNs) as the generic base model, *e.g.*, ULMFiT [Howard and Ruder, 2018] and ELMo [Peters *et al.*, 2018]. More recent methods mostly employ the Transformer architecture [Vaswani *et al.*, 2017] as their base model design, which does not process sequentially like RNNs and instead relies entirely on the attention mechanism, *e.g.*, BERT [Devlin *et al.*, 2019], GPT-2 [Radford *et al.*, 2019], XLNet [Yang *et al.*, 2019], MegatronLM [Shoeybi *et al.*, 2019], and Turing-NLG [Microsoft, 2020].

These Transformers are powerful: for instance, BERT, when first released, improved the state of the art for eleven different NLP tasks by significant margins [Devlin *et al.*, 2019]. However, Transformers are also bulky and resource-hungry: for instance, the most recent large-scale Transformer, Turing-NLG [Microsoft, 2020], has over 17 billion parameters. To quote the Microsoft blog, the general trend is that “larger natural language models lead to better results”. Models of this size incur high memory consumption, computational overhead, and energy cost. To put things in perspective, a model with over 1.3 billion parameters cannot fit into a GPU with 32GB video memory, which is the capacity of Tesla V100, the most advanced data center GPU today. This problem is exacerbated when we consider devices with lower capacity, *e.g.*, smartphones, and applications with strict latency constraints, *e.g.*, interactive chatbots.

One way to address this problem is through model compression, an intricate part of deep learning that has attracted much attention from both researchers and practitioners. Although most methods in this area were originally proposed for convolutional neural networks (CNNs), *e.g.*, [Cheng *et al.*, 2017], many ideas are directly applicable to Transformers. Common techniques include (i) pruning, which removes less important parameters, (ii) weight quantization, which uses fewer bits to represent the parameters, (iii) parameter sharing across similar model units, and (iv) knowledge distillation, which trains a smaller student model that learns from intermediate outputs from the original model. Unlike CNNs, a Transformer model has a complex architecture consisting of multiple parts such as embedding layers, self-attention layers, and feed forward layers (see Section 2). The relative effectiveness of different compression methods can vary when applied to different parts of a Transformer model (see Section 3). There are also methods designed specifically for Transformers, *e.g.*, attention head pruning (see Section 3.2), and replacing the entire Transformer with an RNN or a CNN (see Section 3.3).

Several recent surveys focus on pre-trained representations and large-scale Transformer-based models *e.g.*, [Storks *et al.*, 2019; Jing and Xu, 2019; Wang *et al.*, 2019b]. However, to the best of our knowledge, no comprehensive, systematic study compares the effectiveness of different model compression

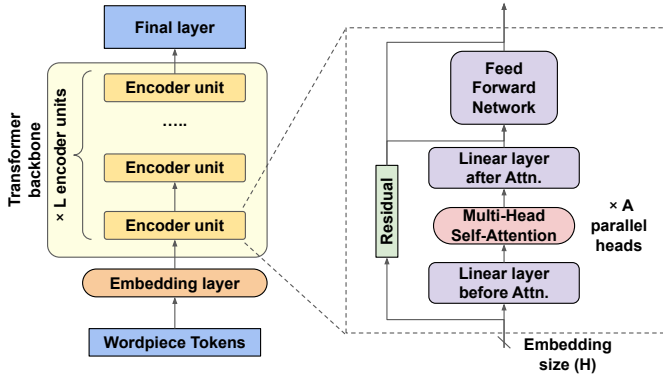


Figure 1: BERT model flowchart.

sion techniques on Transformer-based large-scale NLP models, even though a variety of approaches to compressing such models have been proposed.

Motivated by this, here we offer a thorough and in-depth comparative study on compressing Transformer-based natural language models, with a special focus on the widely used BERT [Devlin *et al.*, 2019] model for NLP. This topic is timely, since (i) the use of Transformer-based BERT-like models has grown dramatically, as demonstrated by recent EMNLP 2019 workshop winners in the FLC task [Yoosuf and Yang, 2019], the DGT task [Saleh *et al.*, 2019], the MRC task [Li *et al.*, 2019], *etc.*, as well as the current leaders of the GLUE benchmark [Wang *et al.*, 2019a] and the SQuAD dataset [Rajpurkar *et al.*, 2016]; (ii) many researchers are being left behind as they do not have expensive GPUs (or a multi-GPU setup) with a large amount of video memory, and thus cannot use the large BERT model; and (iii) AI-powered devices such as smartphones and smart speakers can benefit tremendously from an on-board BERT-like model, but do not have the capability to run it. In addition to summarizing existing techniques and best practices for BERT compression, we point out several promising future research directions for compressing large-scale Transformer-based models.

2 A Breakdown & Analysis of BERT

Bidirectional Encoder Representations from Transformers, or BERT [Devlin *et al.*, 2019], is a Transformer model [Vaswani *et al.*, 2017] pre-trained on large corpora from Wikipedia and the Bookcorpus [Zhu *et al.*, 2015] using two training objectives: (i) Masked Language Model (MLM), which helps it learn the context in a sentence, and (ii) Next Sentence Prediction (NSP), from which it learns the relationship between two sentences. A subsequent Transformer called ALBERT [Lan *et al.*, 2019] replaces NSP with the more challenging Sentence Order Prediction (SOP) task, but here we will focus on the original BERT model. We will mainly consider model size reduction and inference time speedup, rather than the training time.

BERT decomposes its input sentence(s) into WordPiece tokens [Wu *et al.*, 2016], through a tokenizer trained on the training corpora. The use of WordPieces enables BERT to reduce the size of the vocabulary, and at the same time to

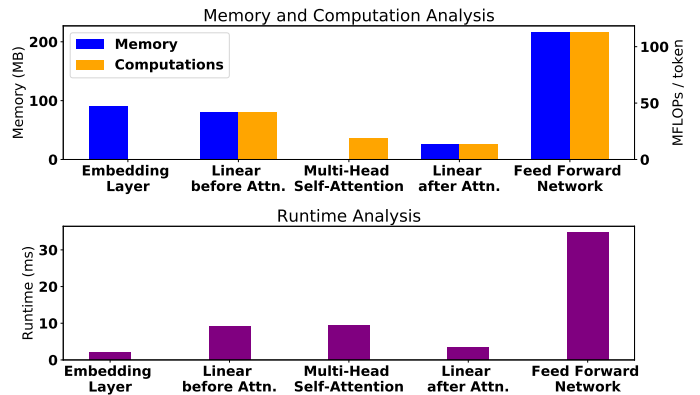


Figure 2: Memory and theoretical computational cost breakdown (top) and runtime cost analysis (bottom) for BERT_{BASE}. All measurements are the sums over all encoder units.

become more robust in the presence of out-of-vocabulary (OOV) words. A special classification token ([CLS]) is inserted before the input, and the final output corresponding to this token is used for classification tasks. For sentence pair tasks, the two sentences are packed together by inserting a separator token ([SEP]) between them. BERT represents each WordPiece token with three vectors, called its token, segment, and position embedding, respectively. These embeddings are summed together and then passed through the main body of the model, *i.e.*, the Transformer backbone, which produces the output representations that are fed to the final, application-dependent layer, *e.g.*, a classifier for sentiment analysis.

As shown in Fig. 1, the Transformer layers have multiple stacked encoder units, each with two major sub-units: a self-attention sub-unit and a feed forward network (FFN) sub-unit [Vaswani *et al.*, 2017], both with residual connections. Each self-attention sub-unit consists of a multi-head self-attention layer, and fully-connected layers before and after it. An FFN sub-unit exclusively contains fully-connected layers.

The size of a BERT model can be specified using three hyper-parameters: the number of encoder units (L), the size of each embedding vector (H), and the number of attention heads in each self-attention layer (A). L and H determine the depth and the width of the model, whereas A is an internal hyper-parameter of the Transformer that affects the number of contextual relations that each encoder can focus on. The authors of BERT provided two pre-trained models: BERT_{BASE} ($L = 12$; $H = 768$; $A = 12$) and BERT_{LARGE} ($L = 24$; $H = 1024$; $A = 16$) [Devlin *et al.*, 2019].

Fig. 2 analyzes the memory and the theoretical computational requirements (measured in millions of FLOPs) of different parts of the BERT_{BASE} model. Clearly, the parts consuming the most memory and executing the highest number of FLOPs are the FFN sub-units. The embedding layer also requires substantial memory, due to the large vector size (H) used to represent each embedding vector. Note that the embedding layer has zero FLOPs, since it is simply a lookup table that involves no arithmetic computations at inference time.

For the self-attention sub-units, we further break down the costs into multi-head self-attention layers and the linear

(*i.e.*, fully-connected) layers before and after them. The self-attention layers consume zero memory since they do not have any learnable parameters; however, their computational cost is non-zero due to the softmax operations. The linear layers surrounding each attention layer incur additional memory and computational overhead, though it is relatively small compared to the FFN sub-units. Note that the input to the attention layer is divided among various heads, and thus each head operates in a lower-dimensional space (H/A). The linear layer before each attention is roughly three times the size of that after the attention, since each attention has three inputs (key, value, and query) and only one output.

The theoretical computational overhead may differ from the actual inference cost at run-time, which depends on the hardware where the model runs. To measure the real run-time costs, we conducted experiments using an Nvidia Titan X GPU with 12GB of video RAM, which is popular among researchers due to its high performance/price ratio. We made all measurements using the TensorFlow profiling tool. The most notable difference between the theoretical analysis and the run-time measurements is that the multi-head self-attention layers are significantly more costly in practice than in theory. This is because the softmax operations in these layers are rather complex, and are implemented as several matrix transformations followed by a matrix multiplication. Consequently, each FLOP in such a layer takes longer compared to other layers. Meanwhile, the execution time of the embedding layers is non-zero (due to memory access costs), though still relatively small. The FFN sub-units are the bottleneck of the whole model, which is consistent with the results from the theoretical analysis.

3 Model Compression Methods

In this section we review the BERT model compression methods in the literature. Due to BERT’s complex architecture, no existing method takes the entire model as the compression target. Instead, each compression technique applies to certain components of BERT, as summarized in Table 1. The table also classifies existing methods based on whether they require access to the original BERT training corpus (marked with \checkmark) or are compressed directly on the task-specific dataset (marked with \triangle). Below, we review each compression technique individually.

3.1 Data Quantization (DQ)

DQ refers to representing each model weight using fewer bits, which reduces the memory footprint of the model and lowers the precision of its numerical calculations. DQ may improve inference speed when the underlying computational device is optimized to process lower-precision numbers faster, *e.g.*, tensor cores in newer generations of Nvidia GPUs. Moreover, programmable hardware such as FPGAs can be specifically optimized for any bitwidth representation.

DQ is generally applicable to all model weights as BERT weights reside in fully-connected layers (*i.e.*, the embedding layer, linear layers, and FFN sub-units), which have been shown to be quantization-friendly [Hubara *et al.*, 2018]. The original BERT model provided by Google represents

Table 1: Compression methods and their corresponding target architectural components in BERT.

Component	DQ	PR		KD-			AIC		
		EP	SP	EO	OL	AM	PS	EMC	AD
Embed. Layer	\checkmark			\triangle	\triangle			\triangle	
Self-Attn.	\checkmark	\checkmark	\triangle	\triangle	\triangle	\triangle			\checkmark
FFN	\checkmark	\checkmark		\triangle	\triangle				
Transformer				\triangle	\triangle		\checkmark		

\checkmark indicates the requirement of the original BERT training corpus.

\triangle indicates the requirement of only task-specific datasets.

each weight by a 32-bit floating number. A naïve approach is to simply truncate each 32-bit weight to the target bitwidth, which often produces a sizable drop in accuracy. Instead, the Quantization Aware Training (QAT) approach is more effective at retaining the model accuracy. Unlike the naïve approach, QAT involves additional training steps to adjust the quantized weights. QAT for BERT is used in the fake nodes [Zafir *et al.*, 2019; Shen *et al.*, 2020], 8-bit integer [Zafir *et al.*, 2019], and Hessian-based mixed-precision [Shen *et al.*, 2020] methods. Interestingly, the embedding layer is more sensitive to quantization than other layers and requires more bits to maintain its accuracy [Shen *et al.*, 2020].

3.2 Pruning (PR)

PR refers to identifying and removing redundant or less important weights and/or components, which sometimes even makes the model more robust and better-performing. PR methods for BERT largely fall into two categories.

Elementwise Pruning (EP). EP, also known as sparse pruning, focuses on PR individual weights, and focuses on locating the set of least important weights in the model. The importance of weights can be judged by their absolute values, gradients, or another measurement defined by the designer [Gordon *et al.*, 2020; Guo *et al.*, 2019]. PR could be effective for BERT, given BERT’s massive fully-connected layers. However, existing EP methods have only experimented with the Transformer backbone (*i.e.*, the self-attention and the FFN sub-units) and have left the embedding layers unexplored. PR methods include magnitude weight pruning [Gordon *et al.*, 2020], which simply removes weights close to zero; and reweighted proximal pruning (RPP) [Guo *et al.*, 2019], which uses iteratively reweighted ℓ_1 minimization followed by the proximal algorithm for decoupling PR and error back-propagation. Since EP considers each weight individually, the set of pruned weights can be arbitrary and is usually irregular.

Structured Pruning (SP). Unlike EP, SP focuses on pruning the architectural components by reducing and simplifying the numerical component modules in the BERT model:

- *Attention head pruning.* As demonstrated in Section 2, the self-attention layer incurs considerable computational overhead at inference time and its importance has often been questioned [Kovaleva *et al.*, 2019]. Recently, it was also shown that most attention heads focus only on trivial positional relations and can be replaced with fixed attention patterns [Raganato *et al.*, 2020]. Therefore, reducing the number of heads of the self-attention layer can signif-

icantly reduce the execution time of the model. In fact, high accuracy is possible with only 1–2 attention heads per encoder unit, though the original model has 16 attention heads [Michel *et al.*, 2019].

- *Encoder unit pruning.* Instead of training a smaller student model from scratch, we could reduce the number of encoder units (L) by pruning less important layers. Layer dropout has been proposed which randomly drops encoder units during training. This makes extracting a smaller model of any desired depth possible during inference, since the original model has been trained to be robust to such pruning [Fan *et al.*, 2019].

3.3 Knowledge Distillation (KD)

KD refers to training a smaller model (called the *student*) using outputs (usually from intermediate layers) of one or more larger pre-trained models (called the *teachers*). In the BERT model, there are multiple intermediate results that the student can learn from, such as the outputs of the encoder units, the logits in the final layer and the attention maps. Based on what the student learns from the teacher, we categorize existing methods as follows.

Distillation on encoder outputs (EO). Every encoder unit is responsible for providing some semantic and contextual relationship between words in the input sentence(s) and progressively improving the pre-trained representation. Thus we can create a smaller BERT model by reducing the size of the encoder unit (H), the number of encoder units (L), or both.

- Reducing H leads to more compact representations in the student [Zhao *et al.*, 2019; Sun *et al.*, 2020; Jiao *et al.*, 2019]. One challenge is that the student cannot directly learn from the teacher’s outputs, due to different sizes. To overcome this, the student also learns a transformation, which can be implemented as either down-projecting the teacher’s outputs to a lower dimensionality, or up-projecting the student’s outputs to the original dimensionality; the performance is similar [Zhao *et al.*, 2019].
- Reducing L , *i.e.*, the number of encoder units, forces each encoder unit in the student to learn from the behavior of a sequence of multiple encoder units in the teacher. Distilling knowledge from the outputs of equally spaced encoder units in the teacher captures relatively more homogeneous information, which leads to a better-performing student model compared to only distilling knowledge from the encoder units towards the end [Sun *et al.*, 2019; Sanh *et al.*, 2019; Sun *et al.*, 2020; Jiao *et al.*, 2019; Zhao *et al.*, 2019].

Distillation on output logits (OL). Similarly to knowledge distillation for CNNs [Cheng *et al.*, 2017], the student can directly learn from the logits (*i.e.*, from soft labels) of the final softmax layer. A common hyper-parameter in softmax calculation is the *temperature*, which controls the smoothness of the output by scaling the logits before applying softmax. Thus, this hyper-parameter determines the extent to which the student relies on the soft labels provided by the teacher [Hinton *et al.*, 2015]. Depending on the temperature setting, existing solutions can be further divided into two sub-categories: (i) methods that experiment with different temperature values [Sun *et al.*, 2019;

Sanh *et al.*, 2019; Chen *et al.*, 2020], and (ii) methods that fix the temperature value to 1 [Jiao *et al.*, 2019; Zhao *et al.*, 2019; Tang *et al.*, 2019b; Liu *et al.*, 2019; Tang *et al.*, 2019a; Turc *et al.*, 2019; Cao *et al.*, 2020].

Note that here the student does not need to be a smaller version of BERT or even a Transformer, and can follow a completely different architecture. This kind of replacement can help preserve the nature of the model while compressing it significantly. The two commonly used replacements for the Transformer architecture are the following.

- *Replace the Transformer with a BiLSTM*, to create a lighter backbone. Recurrent models such as BiLSTMs process words sequentially instead of simultaneously attending to each word in the sentence as in Transformers. Both can create bidirectional representations, and thus BiLSTMs can be considered a lighter alternative to Transformers. Most of the work in compressing a BERT model into a BiLSTM is done downstream directly on a specific NLP task. However, due to the smaller size of task-specific datasets, various methods have been proposed to create additional synthetic training data using rule-based data augmentation techniques [Tang *et al.*, 2019b; Tang *et al.*, 2019a] or to collect data from multiple tasks together to train a single model [Liu *et al.*, 2019].
- *Replace the Transformer with a CNN*, so that the CNN’s parallelism will improve run-time performance during inference [Chia *et al.*, 2018]. Like Transformers, CNNs can process in all directions at the same time. However, due to the smaller size of its kernel, a single CNN convolution layer can only focus on local context, while a single encoder unit in a Transformer focuses on the complete global context. Thus, the Transformer backbone can be replaced with a deep CNN network, obtained using a differentiable Neural Architecture Search (NAS) algorithm [Chen *et al.*, 2020] in the student network to obtain a lighter and faster model.

Distillation on attention maps (AM). An AM refers to the softmax distribution output by the self-attention layers that indicates the contextual dependence between various input tokens. An AM in BERT shows distinguishable linguistic relations such as identical words across sentences, verbs and the corresponding objects, or possessive pronouns and the corresponding nouns [Clark *et al.*, 2019]. To capture such information, the student can be directed to learn these linguistic relations [Sun *et al.*, 2020; Jiao *et al.*, 2019].

3.4 Architecture-Invariant Compression (AIC)

Architecture-invariant methods compress the input model without changing its architecture. DQ (Section 3.1) is one such method. Here, we examine other types of AIC.

Parameter sharing (PS). ALBERT [Lan *et al.*, 2019] follows the same architecture as BERT, but shares weights across all encoder units, which yields a significantly reduced memory footprint. Moreover, ALBERT has been shown to enable training larger and deeper models. For example, while BERT’s performance peaks at BERT_{LARGE} (performance of BERT_{XLARGE} drops significantly), the performance of ALBERT keeps improving until the far larger ALBERT_{XXLARGE}

Table 2: Evaluation comparison of various compression methods.

	Model Size	Speedup	Accuracy/F1				$\frac{\Delta MS}{\Delta A}$ †
			MNLI	QQP	SST-2	SQuAD1.1	
BERT _{BASE}	100%	1x	84.6	89.2/71.2	93.5	88.5	–
DQ	25%	4x	–	87.9/–	92.2	88.5	–
	11.5%	–	81.7	–	92.1	87.0	30.52
PR	53.6%	–	86.1	91.2/–	92.4	90.2	–
	30.9%	–	81.4	89.2/–	90.5	81.7	21.59
KD-OL	11.7%	–	78.6	88.6/70.7	91.0	–	14.72
	4.4%	17.7x	72.6	86.2/64.4	89.8	–	7.96
KD-EO+OL	60.5%	1.63x	–	89.2/–	92.7	86.9	–
	1.6%	–	71.0	–	82.8	–	7.23
KD-EO+AM	23.2%	2.36x	84.3	88.3/70.5	92.6	90.2	256
KD-EO+OL+AM	13.3%	9.4x	82.5	89.2/71.3	92.6	–	41.28
AIC	100%	2.7x*	82.6	90.3/–	–	87.1	0
	11.1%	–	81.6	–	90.3	89.3	44.45
	1.6%	–	71.0	–	82.8	–	7.23

*For methods that allow varying model size for different tasks, we used the average to represent the model size and the speedup.

† ΔMS is the drop in model size, and ΔA is the drop in accuracy for the MNLI task.

($L = 12$; $H = 4096$; $A = 64$) model [Lan *et al.*, 2019].

Embedding matrix compression (EMC). The embedding matrix is the lookup table for the embedding layer, which is about 21% of the size of the complete BERT model. To reduce the size of the embedding matrix, one possibility is to reduce the vocabulary size V , which is about 30k in the original BERT. Recall from Section 2 that the vocabulary of BERT is learned from the training data using a WordPiece tokenizer. When $|V|$ is set to 5k, about 94% of the tokens in the new vocabulary that BERT learns are also present in the vocabulary of the original BERT [Zhao *et al.*, 2019]. This suggests that most of the tokens in the original BERT vocabulary might have been redundant. Another approach applies matrix factorization, which replaces the original $V \times H$ embedding matrix with the product of two smaller matrices ($V \times E$ and $E \times H$). The reduction in memory usage is sizable if $E \ll H$ [Lan *et al.*, 2019].

Attention layer decomposition (AD). The lower multi-head self-attention layers learn local context, whereas upper ones learn global context [Clark *et al.*, 2019]. For layers that focus on the local context, computing self-attention across the complete input can be unnecessarily expensive. Thus, for tasks involving a pair of input sentences, we can decompose the self-attention layer in the lower layers of BERT by calculating self-attention individually for each sentence in the pair [Cao *et al.*, 2020]. Because the self-attention layer contains no weights, this method only reduces computational costs, and the model size remains unchanged. Moreover, since the local context outputs for both sentences are calculated independently of each other, this method also enables a higher degree of parallel processing and caching during inference.

4 Effectiveness of the Compression

In this section, we compare the performance of BERT compression techniques based on their final model size and runtime speedup, as well as their accuracy/F1 scores on various NLP tasks.

4.1 Datasets and Metrics

From the General Language Understanding Evaluation (GLUE) benchmark [Wang *et al.*, 2019a] and the Stanford

Question Answering Dataset (SQuAD) [Rajpurkar *et al.*, 2016], we rely here on the following most commonly reported tasks: MNLI and QQP for sentence pair classification, SST-2 for single sentence classification and SQuAD v1.1 for comprehension-based question answering. As in the GLUE leaderboard, we report accuracy comparison for MNLI and SST-2, and both accuracy and F1 for QQP. For SQuAD v1.1 we report F1, which is the official evaluation measure on the SQuAD leaderboard. We also use the ratio between the drop in the model size and the drop in accuracy on MNLI before and after compression to evaluate the methods on a common scale, since this is an important trade-off.

4.2 Comparison and Analysis

Table 2 compares the effectiveness of BERT compression methods. Note that some methods focused on compression for only part of the model. Further, some methods are trained on BERT_{BASE} and others on BERT_{LARGE} as the teacher model. For uniformity, all model sizes and speedups reported here are for the final complete model after compression and are compared against BERT_{BASE}, even if the method was originally applied to BERT_{LARGE}. However, methods trained with BERT_{LARGE} as a teacher still have an advantage over other methods [Mukherjee and Awadallah, 2019]. We also compare various methods based on their absolute drop in accuracy or F1 scores across different tasks relative to BERT_{BASE}. Next we describe multiple interesting trends present in the table.

Accuracy vs. Model size. The model size drop to accuracy drop ratio ($\frac{\Delta MS}{\Delta A}$) for extreme compression using KD-EO and KD-OL is quite low because of the noticeable drop in accuracy to achieve the required compression. In contrast, DQ and PR achieve respectable ratios since they mainly focus on reducing the redundancies present in the model. However, the ratios involving KD-AM are high (with KD-EO+AM the highest), as KD-AM reduces the size of the model with almost no loss in accuracy.

Pruning and Quantization. Data quantization (DQ) and element-wise pruning (EP) are well-suited for BERT and have shown performance that is on par with other methods. As is shown in Table 2, pruning is shown effective and reduces the original BERT model to 53.6% of its size [Guo *et al.*, 2019], while quantization into eight bits reduces the original size of BERT to 25% [Zafriir *et al.*, 2019], both with only a 1-2% drop in accuracy/F1 for all the tasks. However, when forced to provide higher compression ratios, such as pruning the model to almost 30% of its size [Guo *et al.*, 2019] or quantization to 2-3 bits (11% of the original size) [Shen *et al.*, 2020], these methods experience a sizable drop in accuracy/F1 (upto 4-7% in certain tasks).

Specific distillation. There are multiple ways to distill information while trying to compress BERT. The large variation of final compression ratios (1.6% to 60% of the original model size) for methods using KD on encoder unit outputs (KD-EO) and output logits (KD-OL) shows their complementary nature with respect to a wide range of solutions. We can create an extremely small BERT student model (1.6% of the original size) by reducing the hidden size (H) and the number of encoder units (L) [Zhao *et al.*, 2019], and then

train it to provide respectable performance using distillation from the original BERT model. We also find that KD on attention maps (KD-AM) yields smaller yet better models than KD-EO or KD-OL, which shows the importance of replicating the attention distribution in student networks.

Compression with minimal accuracy loss. An important research direction is to find ways to maximize model compression without hurting accuracy/F1. These compression methods rely on finding redundancies in the original model. We observe that methods such as DQ, PR and multiple forms of KD [Sun *et al.*, 2020; Sanh *et al.*, 2019; Guo *et al.*, 2019; Zafrir *et al.*, 2019] can preserve most of the model’s accuracy/F1 (less than 1-2% drop for all tasks) and are good examples of model compression with minimal accuracy loss. However, as expected, the reduction in model size for these methods is also limited (50-60% of the original model size in most cases), compared to other methods. Specifically, out of all BERT compression approaches, EP provides the best accuracy/F1 (almost no loss) while compressing the model to 53.6% of its original size [Guo *et al.*, 2019].

Extreme compression methods. Sometimes the model must be compressed for deployment on a memory- or latency-bound edge device, while retaining the the best possible accuracy. Most work in this direction relies on replacing the Transformer backbone by a lighter alternative and then using KD-OL to train the model [Tang *et al.*, 2019b; Tang *et al.*, 2019a; Wu *et al.*, 2019; Chen *et al.*, 2020; Liu *et al.*, 2019]. These methods are capable of achieving the highest compression ratio among the existing methods (4-11% of the original model size). Again, as expected, the accuracy/F1 drop for these methods is sizable (upto 6-12% in certain tasks like MNLI) because the focus is on obtaining extremely small models. Out of all the different BERT compression methods, [Zhao *et al.*, 2019] yields the smallest model size, 1.6% of the original model, by reducing both the hidden size (H) and the number of encoder units (L) in the model, but suffers a huge loss in accuracy/F1 (around 11-13% in all tasks).

5 Open Issues & Research Directions

With our analysis and comparison of different BERT compression methods present in literature, we can see that traditional model compression methods such as DQ and pruning do show some benefit for BERT, but the most effective techniques are specific to BERT, including variants of KD and methods to reduce the size of the encoder units and the number of encoder units. These methods also offer insights into the way BERT works and the importance of the different parts of its complicated architecture.

Major Issues. We see the following open issues.

1. A very prominent feature of BERT compression methods is the coupled nature of encoder units as well as the inner architecture of the encoder itself. However, as generally noted in deep learning, some layers in the model might be capable of handling more compression than others, and thus the methods treating different layers equally for compression are not optimal.
2. The very nature of the Transformer backbone that forces

the model to be parameter-heavy makes model compression for BERT more challenging. Existing work in replacing the Transformer backbone by BiLSTMs and CNNs has yielded extraordinary compression ratios, but with a sizable accuracy drop.

3. The accuracy of the Transformer model could be improved by increasing the number of stacked encoder units, *e.g.*, MegatronLM, T-NLG. However, increasing the number of these units will hurt execution time.

Research Directions. BERT compression is still in its early stages; we see multiple interesting avenues for future research.

1. There are interesting structured patterns even in sparsely pruned weight matrices for BERT [Guo *et al.*, 2019]. How to use such patterns to achieve better compression is a promising research direction.
2. The existing idea of compressing self-attention layers for the lower encoder units only [Cao *et al.*, 2020] gives an incentive to further explore decoupling the encoder units. Compressing each encoder unit based on its individual importance, *e.g.*, working with different numbers of attention heads or varied hidden size across encoder units, can be an important step.
3. Decoupling and then individually studying the self-attention and the FFN sub-units is also important for understanding the progressive improvement they provide. As explored elsewhere [Guo *et al.*, 2019; Clark *et al.*, 2019], their importance can vary with the depth. This means that certain FFN sub-units might be more amenable to compression, even if the corresponding self-attention sub-unit is not.
4. The high compression ratios of the approaches that replace the Transformer backbone with a lighter-weight architecture suggest further exploration of more complicated variations of these models and hybrid BiLSTM/CNN/Transformer models (still less parameter-extensive than the Transformer) to limit the accuracy drop [Tang *et al.*, 2019a].
5. Other ideas for improving CNN model performance can be applied to BERT. For example, the parallel convolutions adopted in the Inception network inspires the use of parallel encoder units instead of stacking them together to achieve better accuracy. Similarly, ideas like cross-layer shortcut connections, hourglass architecture, and use of NAS, *etc.*, can also be applied to the backbone architecture of BERT.
6. Many existing methods for BERT compression work on specific parts of the model; we could combine such complementary methods to achieve better overall model compression, *e.g.*, one could combine attention head pruning with quantization.

References

- [Cao *et al.*, 2020] Q. Cao, et al. Faster and just as accurate: A simple decomposition for transformer models. *ICLR Openreview*, 2020.
- [Chen *et al.*, 2020] D. Chen, et al. AdaBERT: Task-adaptive

- BERT compression with differentiable neural architecture search. *Arxiv*, 2020.
- [Cheng *et al.*, 2017] Y. Cheng, et al. A survey of model compression and acceleration for deep neural networks. *Arxiv*, 2017.
- [Chia *et al.*, 2018] Y. K. Chia, et al. Transformer to CNN: Label-scarce distillation for efficient text classification. *NeurIPS-CDNNRIA*, 2018.
- [Clark *et al.*, 2019] K. Clark, et al. What does BERT look at? an analysis of BERT’s attention. *ACL Workshops*, 2019.
- [Devlin *et al.*, 2019] J. Devlin, et al. BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL-HLT*, 2019.
- [Fan *et al.*, 2019] A. Fan, et al. Reducing transformer depth on demand with structured dropout. *ICLR*, 2019.
- [Gordon *et al.*, 2020] M. Gordon, et al. Compressing BERT: Studying the effects of weight pruning on transfer learning. *ICLR Openreview*, 2020.
- [Guo *et al.*, 2019] F. Guo, et al. Reweighted proximal pruning for large-scale language representation. *Arxiv*, 2019.
- [Hinton *et al.*, 2015] G. Hinton, et al. Distilling the knowledge in a neural network. *NIPS Workshops*, 2015.
- [Howard and Ruder, 2018] J. Howard et al. Universal language model fine-tuning for text classification. *ACL*, 2018.
- [Hubara *et al.*, 2018] I. Hubara, et al. Quantized neural networks: Training neural networks with low precision weights and activations. *JMLR*, 2018.
- [Jiao *et al.*, 2019] X. Jiao, et al. Tinybert: Distilling bert for natural language understanding. *Arxiv*, 2019.
- [Jing and Xu, 2019] K. Jing et al. A survey on neural network language models. *Arxiv*, 2019.
- [Kovaleva *et al.*, 2019] O. Kovaleva, et al. Revealing the dark secrets of BERT. *EMNLP*, 2019.
- [Lan *et al.*, 2019] Z. Lan, et al. ALBERT: A lite BERT for self-supervised learning of language representations. *ICLR*, 2019.
- [Li *et al.*, 2019] H. Li, et al. D-net: A simple framework for improving the generalization of machine reading comprehension. *EMNLP-MRQA*, 2019.
- [Liu *et al.*, 2019] L. Liu, et al. Attentive student meets multi-task teacher: Improved knowledge distillation for pre-trained models. *Arxiv*, 2019.
- [Michel *et al.*, 2019] P. Michel, et al. Are sixteen heads really better than one? *NeurIPS*, 2019.
- [Microsoft, 2020] Microsoft. Turing-nlg : A 17-billion parameter language model by microsoft, Feb 2020.
- [Mukherjee and Awadallah, 2019] S. Mukherjee et al. Distilling transformers into simple neural networks with unlabeled transfer data. *Arxiv*, 2019.
- [Peters *et al.*, 2018] M. E. Peters, et al. Deep contextualized word representations. *NAACL-HLT*, 2018.
- [Radford *et al.*, 2019] A. Radford, et al. Language models are unsupervised multitask learners. *OpenAI Blog*, 2019.
- [Raganato *et al.*, 2020] A. Raganato, et al. Fixed encoder self-attention patterns in transformer-based machine translation. *Arxiv*, 2020.
- [Rajpurkar *et al.*, 2016] P. Rajpurkar, et al. Squad: 100,000+ questions for machine comprehension of text. *EMNLP*, 2016.
- [Saleh *et al.*, 2019] F. Saleh, et al. Naver labs europe’s systems for the document-level generation and translation task at wngt 2019. *EMNLP-WNGT*, 2019.
- [Sanh *et al.*, 2019] V. Sanh, et al. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *Arxiv*, 2019.
- [Shen *et al.*, 2020] S. Shen, et al. Q-bert: Hessian based ultra low precision quantization of bert. *AAAI*, 2020.
- [Shoeybi *et al.*, 2019] M. Shoeybi, et al. Megatron-lm: Training multi-billion parameter language models using gpu model parallelism. *Arxiv*, 2019.
- [Storks *et al.*, 2019] S. Storks, et al. Recent advances in natural language inference: A survey of benchmarks, resources, and approaches. *Arxiv*, 2019.
- [Sun *et al.*, 2019] S. Sun, et al. Patient knowledge distillation for BERT model compression. *EMNLP*, 2019.
- [Sun *et al.*, 2020] Z. Sun, et al. MobileBERT: Task-agnostic compression of BERT for resource limited devices. *ICLR Openreview*, 2020.
- [Tang *et al.*, 2019a] R. Tang, et al. Natural language generation for effective knowledge distillation. *DeepLo*, 2019.
- [Tang *et al.*, 2019b] R. Tang, et al. Distilling task-specific knowledge from BERT into simple neural networks. *Arxiv*, 2019.
- [Turc *et al.*, 2019] I. Turc, et al. Well-read students learn better: The impact of student initialization on knowledge distillation. *Arxiv*, 2019.
- [Vaswani *et al.*, 2017] A. Vaswani, et al. Attention is all you need. *NeurIPS*, 2017.
- [Wang *et al.*, 2019a] A. Wang, et al. Glue: A multi-task benchmark and analysis platform for natural language understanding. *ICLR*, 2019.
- [Wang *et al.*, 2019b] S. Wang, et al. A survey of word embeddings based on deep learning. *Computing*, 2019.
- [Wu *et al.*, 2016] Y. Wu, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *Arxiv*, 2016.
- [Wu *et al.*, 2019] B. Wu, et al. Towards non-task-specific distillation of bert via sentence representation approximation. *10.13140/RG.2.2.28187.62246*, 2019.
- [Yang *et al.*, 2019] Z. Yang, et al. Xlnet: Generalized autoregressive pretraining for language understanding. *NeurIPS*, 2019.
- [Yoosuf and Yang, 2019] S. Yoosuf et al. Fine-grained propaganda detection with fine-tuned BERT. *EMNLP*, 2019.
- [Zafir *et al.*, 2019] O. Zafir, et al. Q8bert: Quantized 8bit bert. *NeurIPS-EEMLCC*, 2019.
- [Zhao *et al.*, 2019] S. Zhao, et al. Extreme language model compression with optimal subwords and shared projections. *Arxiv*, 2019.
- [Zhu *et al.*, 2015] Y. Zhu, et al. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. *IEEE ICCV*, 2015.