

# Efficient Fine-Tuning of BERT Models on the Edge

Danilo Vucetic\*, Mohammadreza Tayaranian\*, Maryam Ziaeeefard,  
James J. Clark, Brett H. Meyer and Warren J. Gross  
*Department of Electrical and Computer Engineering,  
McGill University  
Montreal, Canada*

{danilo.vucetic, mohammadreza.tayaranian}@mail.mcgill.ca,  
{maryam.ziaeeefard, james.clark1, brett.meyer, warren.gross}@mcgill.ca

**Abstract**—Resource-constrained devices are increasingly the deployment targets of machine learning applications. Static models, however, do not always suffice for dynamic environments. On-device training of models allows for quick adaptability to new scenarios. With the increasing size of deep neural networks, as noted with the likes of BERT and other natural language processing models, comes increased resource requirements, namely memory, computation, energy, and time. Furthermore, training is far more resource intensive than inference. Resource-constrained on-device learning is thus doubly difficult, especially with large BERT-like models. By reducing the memory usage of fine-tuning, pre-trained BERT models can become efficient enough to fine-tune on resource-constrained devices. We propose Freeze And Reconfigure (FAR), a memory-efficient training regime for BERT-like models that reduces the memory usage of activation maps during fine-tuning by avoiding unnecessary parameter updates. FAR reduces fine-tuning time on the DistilBERT model and CoLA dataset by 30%, and time spent on memory operations by 47%. More broadly, reductions in metric performance on the GLUE and SQuAD datasets are around 1% on average.

**Index Terms**—Transformers, BERT, DistilBERT, NLP, Language Models, Efficient Transfer Learning, Efficient Fine-Tuning, Memory Efficiency, Time Efficiency, Edge Machine Learning

## I. INTRODUCTION

Large language models employing attention-based architectures have revolutionized Natural Language Processing (NLP) with the likes of BERT and GPT-3, achieving numerous state-of-the-art scores on well-studied tasks [1], [2]. As language models move from the cloud to resource-constrained devices (i.e., edge devices), their deployment is made difficult by constraints in available resources such as memory, computation, energy, and time [3]–[5]. Solely considering inference for edge learning is insufficient as the operating environment may be dynamic, and input data may change over time [5]. Learning on edge devices allows models to continually adapt to their environments and to the requirements of their users, without revealing data to the cloud [4]–[6]. However, training requires about three times more memory operations

than inference<sup>1</sup> and more resource utilization generally. Tools like model compression are not enough to bridge the gap (i.e., a model 1/3 the size with similar performance is hard to achieve), and compressed models may not be training-efficient [7]. The problem specifically lies with the greater memory requirements of training compared to inference, since: 1) accesses to main memory are two orders of magnitude more energy intensive than computation; and 2) large models can not fit in cache, meaning that higher latency accesses to main memory are required [7]–[9].

Enabling on-device learning for large language models necessitates two conditions: 1) using a compressed model to reduce baseline memory requirements, and 2) training in a resource-efficient manner. To achieve the latter condition we propose *Freeze And Reconfigure (FAR)*, a memory-efficient fine-tuning regime for BERT-like models. We focus first and foremost on reducing memory usage during the backward pass by reducing the largest contributors to said usage: activations and gradients. We introduce a methodology that identifies which parameters to *freeze* and which to continue training. *Freeze* here refers to disabling parameter updates during back-propagation. We also propose a method of dynamic architecture reconfiguration such that frozen and non-frozen parameters are grouped separately. Reconfiguration reduces poorly-structured memory accesses which, for example, encumber unstructured pruning methods. The contributions of this paper are summarized as follows:

- We introduce a novel learning metric that tracks the size of weight updates over contiguous groups of parameters during fine-tuning using the  $L_1$  norm. According to the metric, the best-learning subsets of nodes are selected.
- We introduce a novel parameter freezing scheme that exploits architectural reconfiguration, grouping frozen and nonfrozen parameters, reducing training-time resource consumption.

FAR achieves sizeable reductions in the resource consumption of fine-tuning. Our results show that when fine-tuning DistilBERT on CoLA with 60% of the model parameters frozen, fine-tuning time is reduced by 30% and memory

\* Equal contribution.

<sup>1</sup>During inference only the parameters and data are retrieved from memory (2 operations). During training, the parameters must be retrieved at least twice and stored once after being updated, activations are stored and retrieved, and data is retrieved (6 operations).

access time is reduced by 47%, while metric performance is maintained. Similar results are achieved for other GLUE and SQuAD tasks. Evidently, FAR makes training more efficient while maintaining metric performance.

## II. BACKGROUND AND RELATED WORK

### A. Transformers and BERT Models: Training and Efficiency

Transformers and BERT-like models make use of Attention and Feed-Forward Network (FFN) sublayers to extract linguistic knowledge and achieve state-of-the-art scores [10]. BERT in particular uses the Transformer Encoder architecture to model languages by pre-training on large datasets of relatively simple language tasks. To deploy BERT-like models to a new downstream task, the model must be fine-tuned [1]. This allows the use of a single set of weights, the pre-trained BERT weights, as the starting point for various, faster converging, fine-tuning tasks, avoiding the large overhead of pre-training the full BERT model<sup>2</sup>. Various methods have been proposed to compress BERT including DistilBERT [11], MobileBERT [12], and TinyBERT [13]. A detailed discussion of compressing large language models is found in [14], including analysis of the major bottleneck of Transformer-based models: the FFN sublayers, which contribute more parameters, runtime memory consumption, and inference latency than all other sublayers. Model compression is useful for enabling on-device learning, but greater efficiency during training is desirable to further enable the realization of large models on edge devices.

### B. Efficient Learning

On-device learning has been detailed in [5], including methods of data compression and theoretical approaches to efficiency. When considering on-device learning, a fundamental trade-off must be kept in mind to train in a reasonable amount of time: either the training-time resource requirements of a model must be decreased, or the hardware resources must be increased. The former approach is usually studied in terms of model compression such as in [11]–[13]. Others take the problem as reducing the real requirements of training without altering model architecture. One such approach, [15], trains only the best-learning weights of a feed-forward neural network, while keeping the others at their initialization values. Trained weights are selected by tracking the largest accumulated weight updates over some initial training steps. The authors also further reduce memory requirements by regenerating frozen weights from a pseudorandom number generator, thereby trading memory for extra computation. This approach does not work for BERT-like models as the weights of the pre-trained model are not randomly generated. Additionally, the unstructured selection of weights leads to inefficient memory accesses, an issue which we address with structured selection in FAR. A similar approach, BitFit [16], freezes all BERT weights during fine-tuning and only applies updates to bias terms. The authors report near-baseline performance

<sup>2</sup>As reported in [1], 4 cloud tensor processing units running for 4 days are required to pre-train BERT’s 110 million parameters.

despite the large drop in trained parameters. While BitFit works on uncompressed (i.e., highly overparameterized) BERT models, on the more compressed variants like DistilBERT, it is expected to fail due to the smaller network size and lower capacity. We address this issue by fine-tuning a greater proportion of the parameters.

## III. FREEZE AND RECONFIGURE

On-device learning of large language models requires compressed models and efficient training procedures. The reduction in memory usage provided by these requirements has knock-on effects: higher energy efficiency, reduced computation, and reduced memory access times. In the case of memory access time, reducing the number of parameters decreases the amount of information transferred across a system.

FAR is specifically designed for use with compressed language models, specifically DistilBERT, whose capacity and adaptability are diminished from their uncompressed counterparts [17]. This is accomplished by fine-tuning a subset of linear layer weights (called FFNs in DistilBERT) as well as their biases, in contradistinction to BitFit, which trains only the biases of BERT.

In FAR, the sets of fine-tuned weights in the linear layers are selected in a structured manner to avoid sparse memory accesses which would delay memory operations [18]. Each node in an FFN layer is considered a single group of parameters. These nodes are classified as *learner* and *nonlearner* nodes depending on how well a node performs during fine-tuning. The parameters of the nonlearner nodes are frozen while learner nodes are fine-tuned. Frozen nodes no longer require gradient calculations during the backward pass, storage of activations, or memory accesses during the backward pass. Reconfiguration of the FFN layers is completed during fine-tuning by separating learner and nonlearner nodes into distinct, newly created, FFN sublayers.

Each Transformer block in DistilBERT is made up of a multi-head attention layer which feeds an FFN. The FFN consists of two dense layers which contain 3072 and 768 nodes respectively. The weights of the FFNs make up more than 66% of the parameters of the whole network. Thus, freezing a subset of these weights has a considerable effect on the overall number of updated parameters of the model. Past work has shown that freezing up to 80% of feed-forward neural network parameters results in minimal changes in accuracy for smaller models and simpler tasks [15].

### A. Selection of Learner Nodes

The set of learner nodes is decided upon after an initial set of fine-tuning iterations during a process we call *priming* (cf. Figure 1, green nodes in left figure are learner nodes). During priming, a copy of the pre-trained FFN weights is stored and the network is fine-tuned for some percentage,  $p$ , of the total number of optimization steps. The priming percentage,  $p$ , is kept small so as to avoid the additional burden of fully fine-tuning a large model for a greater proportion of the optimization steps. It is expected, however, that more

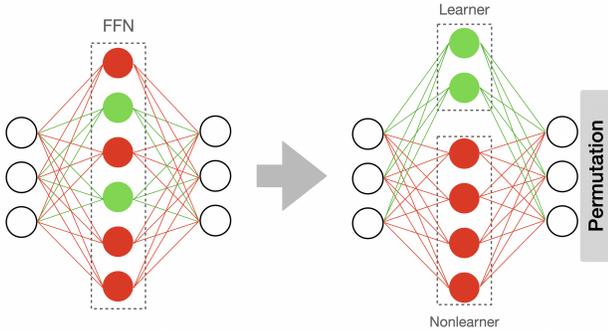


Fig. 1. Overview of priming steps and the dynamic reconfiguration of nodes in FFN. In the left, the learner (green) and nonlearner (red) nodes are chosen after the priming steps. In the right, the nodes are reconfigured to disable gradient computation for red weights. The output of the FFN is permuted after reconfiguration.

priming should result in higher metric performance. After fine-tuning the model during priming, the initial FFN weight vectors,  $\mathbf{w}_n^{e,i}$ , are subtracted from the fine-tuned FFN weight vectors,  $\phi_n^{e,i}$ , for each Transformer Encoder,  $e$ , FFN sublayer,  $i$ , and node,  $n$ . We employ the  $L_1$  norm to compute the total amount learned by each node in the FFN. This is in distinction to [15], where each weight is chosen individually based on the difference in its initial and trained value. Usage of the  $L_1$  norm allows for the weights of the FFN node to be grouped so that architectural modifications can be made later for greater memory savings. The learning metric is computed as:

$$m_n^{e,i} = \|\phi_n^{e,i} - \mathbf{w}_n^{e,i}\|_1 \quad (1)$$

Finally, to compute the set of learner nodes, a retention percentage,  $r$ , is defined as the ratio of learner nodes to the total number of nodes. After sorting the learning metric for each FFN, the  $r\%$  of nodes with the largest metric are classified as learners; the rest as nonlearners. Note that decreasing the node retention results in more nonlearner nodes, which in turn reduces memory utilization during fine-tuning. Using higher retention values is expected to have a positive impact on metric performance as it increases the model’s ability to adapt to the down-stream task. In addition, the priming and retention hyperparameters are defined as percentages to allow them to scale to larger datasets and varying models. For example, a more compressed model may require higher node retention to avoid substantial losses in metric performance.

### B. Dynamic FFN Reconfiguration

After assigning the learner and nonlearner nodes, the FFNs of each Transformer block are reconfigured to separate the nonlearner nodes from learner nodes. As shown in Figure 1, FFN nodes are reconfigured into two parallel sub-modules. Such reconfiguration ensures that the memory of learner nodes is grouped together, allowing only the fine-tuned parameters to be retrieved during fine-tuning. Finally, the reconfigured FFNs need to have their output order restored such that their outputs are coherent with the parameters of the model, so a

TABLE I  
EFFECTIVENESS OF PRIMING AND RETENTION WITH 5 AVERAGED RUNS ON DISTILBERT FOR EACH DATASET

		MNLI			CoLA		
$p$	$r$	10%	25%	40%	10%	25%	40%
1%		81.8	81.9	82.1	52.8	53.1	53.3
5%		81.8	82.1	81.9	51.9	52.7	53.7
10%		81.9	82.0	<b>82.2</b>	53.4	52.5	<b>54.0</b>

permutation is applied to the output vectors. The combination of reconfiguration and permutation trades memory accesses for additional time spent on computation and data movement, which is far less costly [8].

Using nodes as parameter groupings and reconfiguring the model this way allows us to avoid memory access issues that arise in methods like unstructured pruning. Since individual weights are never solely accessed, we avoid sparse memory accesses. Reconfiguration is also beneficial for practical implementation since in PyTorch, automatic gradient computation can not be disabled for singular parameters, but instead only for whole layers or sublayers. Thus, unlike unstructured pruning, no special hardware or libraries are needed to implement FAR.

## IV. EXPERIMENTS

We test FAR using the compressed BERT-based model DistilBERT [11]. DistilBERT is a pre-trained language model and can be fine-tuned to downstream tasks. DistilBERT has six Transformer encoder layers, each with two FFN sublayers (i.e.,  $e \in \{1, \dots, 6\}$  and  $i \in \{1, 2\}$  in Equation 1). We fine-tune the model using a learning rate of  $2e^{-5}$ , with a linear learning rate schedule, batch size of 16, using at most 5 epochs.

Experiments are run on *NVIDIA Tesla V100 Volta* GPU accelerators to gather metric results for all datasets. Whenever timing is reported, these results have been gathered on the *NVIDIA Jetson Xavier NX* edge device. Memory access time is computed using the NVIDIA’s profiling tool, *nvprof*. Memory access time is the summation of the execution time for memory-related API calls, namely *cudaStreamSynchronize*, *cudaMemcpyAsync* and *cudaMemsetAsync*.

FAR is tested against conventional and efficient fine-tuning approaches using standard NLP datasets from the GLUE benchmark and SQuAD 2.0 [19], [20]. All the GLUE tasks are included in the experiments except for RTE, MRPC, STS-B because of their small dataset sizes and WNLI because of its confounding results [21]. This leaves MNLI, QQP, QNLI, SST-2, and CoLA. CoLA (9594 training examples) and MNLI (433k training examples) are used to find effective priming and retention percentages in grid searches (Cf. Table I). MNLI and CoLA are believed to represent small and large datasets, respectively [19]. Each reported value is the average of 5 fine-tuning attempts with random seeds to ensure consistency.

In order to test the hypothesis that higher priming and higher retention percentages produce better metric performance, we complete a grid-search style test on CoLA and MNLI. The

TABLE II

METRIC PERFORMANCE OF DISTILBERT WITH AND WITHOUT FAR ON GLUE AND SQUAD 2.0 - 1% PRIMING AND 10% AND 40% RETENTION. EACH FAR RUN IS ACCOMPANIED WITH ITS PERCENTAGE DROP COMPARED TO THE BASELINE. THE LAST COLUMN SHOWS THE AVERAGE SCORE DROP PERCENTAGE ACROSS ALL TASKS. † SHOWS OUR IMPLEMENTATION.

	CoLA	MNLI	QNLI	QQP	SST2	SQuAD 2.0		Avg
	Mathews Corr.	Acc.	Acc.	Acc.	Acc.	EM	F1	
Baseline †	<b>53.6</b>	<b>82.1</b>	<b>89.1</b>	<b>90.2</b>	91.1	<b>65.0</b>	<b>67.9</b>	
FAR <sub>10</sub>	52.8	81.9	88.2	90.0	91.1	63.3	66.4	
	-1.58%	-0.24%	-1.01%	-0.22%	0.00%	-2.59%	-2.27%	-1.13%
FAR <sub>40</sub>	53.3	<b>82.1</b>	88.5	<b>90.2</b>	<b>91.3</b>	63.6	66.6	
	-0.55%	0.00%	-0.67%	0.00%	0.22%	-2.14%	-2.01%	-0.74%

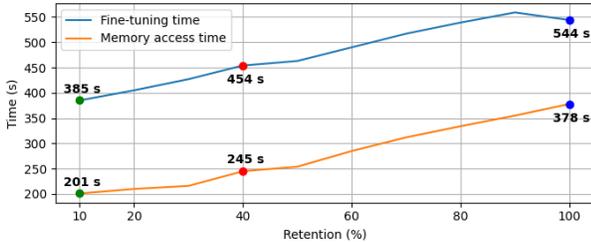


Fig. 2. Overall time and total memory access time for fine-tuning DistilBERT on CoLA using FAR with various retention percentages on the NVIDIA Jetson Xavier NX. Priming percentage is set to 1% across all runs. Values for FAR<sub>10</sub>, FAR<sub>40</sub> and the baseline (FAR<sub>100</sub>) are indicated with green, red and blue dots, respectively.

results are listed in Table I. The table shows that with higher retention percentage and higher priming percentage, fine-tuning produces a higher metric score, with the highest score occurring with 10% priming and 40% retention. Moreover, the difference between 1% and 10% priming is small enough to warrant the use of 1% priming for FAR. This also means that the training overhead is made minimal since extra priming does not necessarily correlate to better performance. Henceforth, FAR<sub>x</sub> will represent FAR using 1% priming and x% node retention. FAR<sub>10</sub> and FAR<sub>40</sub> are used to demonstrate two trade-off points between metric performance and resource utilization. The former has lower resource utilization while also producing lower metric performance and vice-versa for the latter. The results of this comparison are listed in Table II alongside with the baseline implementation of DistilBERT using the Hugging Face python library [22], without freezing any weights during fine-tuning. We observe higher metric scores for FAR<sub>40</sub> compared to FAR<sub>10</sub> and a lower drop in performance compared to the baseline at just 0.74% compared to 1.13% respectively. The largest drops in performance occur on SQuAD 2.0, which is the most complex dataset tested.

Figure 2 shows the effect of changing the retention percentage on training time and memory access time. A similar trend is observed in both figures, showing that using setups with lower retention values reduces both fine-tuning and memory access time. These results indicate that FAR reduces the resource consumption of fine-tuning, and as a result also reduces training time. FAR<sub>10</sub> and FAR<sub>40</sub> are also indicated with green and red dots, respectively. In terms of overall fine-tuning time, FAR<sub>10</sub> and FAR<sub>40</sub> each respectively show 30%

TABLE III

COMPARISON OF FAR<sub>10</sub> WITH RANDOM NODE SELECTION AND BITFIT

	CoLA	MNLI	SQuAD 2.0	
	Mathews Corr.	Acc.	EM	F1
FAR <sub>10</sub>	<b>52.8</b>	<b>81.9</b>	<b>63.3</b>	<b>66.4</b>
Random Selection <sub>10</sub>	52.0	81.8	62.8	65.7
BitFit	51.7	81.5	49.8	49.8

and 17% reductions compared to the baseline. For memory access time, 47% and 36% are the time reductions associated to FAR<sub>10</sub> and FAR<sub>40</sub>.

Finally, we carry out an ablation study to measure the effectiveness of our method against related works and to show that FAR is in fact necessary for efficient fine-tuning. First, to show the effect of using priming, we perform an experiment in which the learner and nonlearner nodes are randomly selected. Table III shows that for the same retention value of 10%, FAR, which uses priming to select nodes, performs better than random node selection, with the gap being wider in the more complex SQuAD task. Comparing FAR with BitFit, which freezes all the weights of all dense layers across the network, we see a drop in performance of less than 2% for CoLA and MNLI. However, BitFit performs poorly on SQuAD 2.0, proving its failure to cope with complex tasks on compressed models like DistilBERT. The memory access and training times of fine-tuning on DistilBERT using BitFit on CoLA are 237 and 384 seconds. Clearly, BitFit spends more time accessing memory, but around the same time fine-tuning as FAR<sub>10</sub> (Figure 2).

## V. CONCLUSIONS AND FUTURE WORK

In this paper we propose FAR, a novel method to reduce fine-tuning memory consumption by exploiting the overparameterization of large language models. We show that for a compressed BERT-based model, it is not enough to simply freeze large sections of the model during fine-tuning, but that freezing must be selective and structured. Despite large decreases in the number of fine-tuned parameters, metric performance remains near baseline levels, while fine-tuning time and memory usage are significantly reduced. Future work involves investigations into new domains such as Transformers on NLP and vision tasks. We will also apply FAR to other backbone models such as MobileBERT to demonstrate the generalizability of the method.

## ACKNOWLEDGMENT

We thank Huawei Canada for their sponsorship of this work and Compute Canada for providing computational resources.

## REFERENCES

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186. [Online]. Available: <https://aclanthology.org/N19-1423>
- [2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. [Online]. Available: <https://proceedings.neurips.cc>
- [3] E. Strubell, A. Ganesh, and A. McCallum, "Energy and policy considerations for deep learning in nlp," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3645–3650. [Online]. Available: <https://aclanthology.org/P19-1355.pdf>
- [4] R. Sharma, S. Biookaghazadeh, B. Li, and M. Zhao, "Are existing knowledge transfer techniques effective for deep learning with edge devices?" in *2018 IEEE International Conference on Edge Computing (EDGE)*, 2018, pp. 42–49.
- [5] S. Dhar, J. Guo, J. J. Liu, S. Tripathi, U. Kurup, and M. Shah, "A survey of on-device machine learning: An algorithms and learning theory perspective," *ACM Trans. Internet Things*, vol. 2, no. 3, Jul. 2021. [Online]. Available: <https://doi.org/10.1145/3450494>
- [6] M. Satyanarayanan, "The emergence of edge computing," *Computer*, vol. 50, no. 1, pp. 30–39, 2017.
- [7] H. Cai, C. Gan, L. Zhu, and S. Han, "Tinytl: Reduce memory, not parameters for efficient on-device learning," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. [Online]. Available: <https://proceedings.neurips.cc>
- [8] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally, "Eie: Efficient inference engine on compressed deep neural network," in *Proceedings of the 43rd International Symposium on Computer Architecture*, ser. ISCA '16. IEEE Press, 2016, p. 243–254. [Online]. Available: <https://doi.org/10.1109/ISCA.2016.30>
- [9] J. L. Hennessy and D. A. Patterson, *Computer architecture: a quantitative approach*, 5th ed. Elsevier, 2012, pp. B–3.
- [10] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [11] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter," 2019. [Online]. Available: <http://arxiv.org/abs/1910.01108>[unpublished]
- [12] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou, "MobileBERT: a compact task-agnostic BERT for resource-limited devices," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp. 2158–2170. [Online]. Available: <https://www.aclweb.org/anthology/2020.acl-main.195>
- [13] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu, "TinyBERT: Distilling BERT for natural language understanding," in *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association for Computational Linguistics, Nov. 2020, pp. 4163–4174. [Online]. Available: <https://aclanthology.org/2020.findings-emnlp.372>
- [14] P. Ganesh, Y. Chen, X. Lou, M. A. Khan, Y. Yang, D. Chen, M. Winslett, H. Sajjad, and P. Nakov, "Compressing large-scale transformer-based models: A case study on bert," *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 1061–1080, 2021.
- [15] M. Golub, G. Lemieux, and M. Lis, "Full deep neural network training on A pruned weight budget," in *Proceedings of Machine Learning and Systems 2019, MLSys 2019, Stanford, CA, USA, March 31 - April 2, 2019*, A. Talwalkar, V. Smith, and M. Zaharia, Eds. mlsys.org, 2019. [Online]. Available: <https://proceedings.mlsys.org/book/260.pdf>
- [16] E. B. Zaken, S. Ravfogel, and Y. Goldberg, "Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models," *CoRR*, vol. abs/2106.10199, 2021. [Online]. Available: <https://arxiv.org/abs/2106.10199>[unpublished]
- [17] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh, "Improved knowledge distillation via teacher assistant," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5191–5198.
- [18] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, "Learning structured sparsity in deep neural networks," *Advances in neural information processing systems*, vol. 29, pp. 2074–2082, 2016.
- [19] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, "Glue: A multi-task benchmark and analysis platform for natural language understanding," in *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 2018, pp. 353–355.
- [20] P. Rajpurkar, R. Jia, and P. Liang, "Know what you don't know: Unanswerable questions for squad," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 2018, pp. 784–789.
- [21] M. A. Gordon, K. Duh, and N. Andrews, "Compressing bert: Studying the effects of weight pruning on transfer learning," *ACL 2020*, p. 143, 2020.
- [22] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. Rush, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45. [Online]. Available: <https://aclanthology.org/2020.emnlp-demos.6>