

FastFormers: Highly Efficient Transformer Models for Natural Language Understanding

Young Jin Kim

Microsoft

One Microsoft Way

Redmond, WA 98052, USA

youki@microsoft.com

Hany Hassan Awadalla

Microsoft

One Microsoft Way

Redmond, WA 98052, USA

hanyh@microsoft.com

Abstract

Transformer-based models are the state-of-the-art for Natural Language Understanding (NLU) applications. Models are getting bigger and better on various tasks. However, Transformer models remain computationally challenging since they are not efficient at inference-time compared to traditional approaches. In this paper, we present *FastFormers*, a set of recipes to achieve efficient inference-time performance for Transformer-based models on various NLU tasks. We show how carefully utilizing knowledge distillation, structured pruning and numerical optimization can lead to drastic improvements on inference efficiency. We provide effective recipes that can guide practitioners to choose the best settings for various NLU tasks and pretrained models. Applying the proposed recipes to the SuperGLUE benchmark, we achieve from 9.8x up to 233.9x speed-up compared to out-of-the-box models on CPU. On GPU, we also achieve up to 12.4x speed-up with the presented methods. We show that *FastFormers* can drastically reduce cost of serving 100 million requests from 4,223 USD to just 18 USD on an *Azure F16s_v2*¹ instance. This translates to a sustainable runtime by reducing energy consumption 6.9x - 125.8x according to the metrics used in the SustaiNLP 2020 shared task.

1 Introduction

Since *BERT* (Devlin et al., 2018) has been introduced, Transformer (Vaswani et al., 2017)-based pretrained language models have dominated the Natural Language Understanding (NLU) field. Transformer models have provided unprecedented accuracy improvement compared to traditional models (Devlin et al., 2018; Liu et al., 2019). However, the models’ computational cost at inference time is prohibitively challenging to be widely

adopted in real world production scenarios which requires low latency, fast inference and low serving costs. In this work, we present *FastFormers*, a set of methods and recipes that provides highly efficient inference for Transformer models which enables deployment in large scale production scenarios. We specifically focus on the inference time efficiency since it mostly dominates the cost of production deployment.

Mainly, we utilize three methods: Knowledge Distillation, Structured Pruning and Model Quantization. First, we investigate the efficacy of various Knowledge Distillation techniques to significantly reduce the size of the models with respect to the depth and hidden state sizes while preserving the accuracy. Second, we explore Structured Pruning that further reduces the size of the models by reducing the number of self-attention heads and the number of intermediate hidden states in the feed-forward layers to achieve more efficiency while trying to preserve the accuracy as well. Finally, we explore Model Quantization which enables faster model executions by optimally utilizing hardware acceleration capabilities. On CPU, 8-bit integer quantization method is applied to utilize the most efficient CPU instructions available: *Vector Neural Network Instructions (VNNI)*. On GPU, all the model parameters are converted into 16-bit floating point data type to maximally utilize efficient *Tensor Cores*. Furthermore, computational graph optimizations which fuse multiple graph nodes are performed by utilizing *onnxruntime*² library. In addition, we explore optimal settings of allocating CPU and GPU resources to achieve better utilization.

The proposed methods are optimized and evaluated on both CPUs and GPUs which are the most commonly available hardware platforms. Per-

¹<https://docs.microsoft.com/en-us/azure/virtual-machines/fsv2-series>

²<https://github.com/microsoft/onnxruntime>

formance evaluations are conducted on SuperGLUE (Wang et al., 2019) which is one of the general purpose open domain NLU benchmarks. For the efficiency measurement, wall clock times and energy efficiency are measured while performing inferences on the test sets of *BoolQ*, *CB*, *COPA*, *MultiRC*, *ReCoRD*, *RTE* and *WiC* tasks from SuperGLUE. The energy efficiency is measured by an open source python library called *experiment-impact-tracker* proposed in (Henderson et al., 2020). To make sure the optimized models preserve similar accuracy, the accuracy of all the models is measured together.

The contributions of this paper are presenting a set of recipes for efficient inference of Transformer NLU models, analyzing the effect of various optimization techniques and finally making the code publicly available³ to facilitate utilizing *FastFormers* for efficient inference of Transformers models.

The rest of the paper is organized as follows: Section 2 presents Knowledge Distillation techniques, Section 3 presents Structured Pruning techniques, Section 4 discusses Model Quantization approaches, Section 5 presents runtime optimization techniques, Section 6 presents results on various tasks and finally Section 7 concludes the findings and future directions.

2 Knowledge Distillation

Knowledge distillation (Hinton et al., 2015) is a well known model compression technique where the large model is used as a teacher for a smaller student model. The knowledge distillation is the process of training the student model to mimic the behaviour of the larger teacher model. Knowledge distillation has been shown to improve the efficiency of the Transformer-based architectures for the NLU tasks (Sanh et al., 2019; Jiao et al., 2019) as well as natural language generation tasks such as machine translation (Kim et al., 2019).

Knowledge distillation methods: We utilize two different distillation approaches, namely *task-specific* and *task-agnostic* distillation. In the task-specific distillation, we distill fine-tuned teacher models into smaller student architectures following the procedure proposed by *TinyBERT* (Jiao et al., 2019). In the task-agnostic distillation approach, we directly apply fine-tuning on general distilled

models to tune for a specific task. These two methods are illustrated in Figure 1. In the illustration, the preceding number attached to each arrow means the order of distillation steps. We use soft cross-entropy function as the knowledge distillation loss function in all our experiments as used in (Sanh et al., 2019; Jiao et al., 2019).

As teacher models, we choose 12 stacked layer *BERT* (Devlin et al., 2018) and *RoBERTa* (Liu et al., 2019) models with 768 hidden state dimension and 12 self-attention heads. This is referred as *Base* size from the original BERT paper. In particular, we use HuggingFace’s pretrained BERT and RoBERTa⁴ models.

In *task-specific* scenario, we observe that the initialization of the student models affects the final accuracy of the distilled models. We utilize open domain pre-distilled models, namely *distilroberta-base*⁵ and *TinyBERT*⁶ as the initializers for the corresponding student models. Those models have been distilled with one of the original BERT model’s pretraining objective which is masked language model. So, they are not specific to any task and can be considered as smaller generic pretrained models.

Knowledge distillation results: The main goal of the knowledge distillation process is to acquire the smallest possible student model while preserving the accuracy of the teacher model. Since we are experimenting with various NLU tasks, the capacity of the optimal student model that preserves accuracy may vary with varying level of task’s difficulty. Therefore, we experiment with distilling various sized student models; then, we pick the smaller model among the distilled models that can offer higher accuracy than the original BERT model for each task.

In our experiments, we have observed that distilled models do not work well when distilled to a different model type. Therefore, we restricted our setup to avoid distilling RoBERTa model to BERT or vice versa. The major difference between the two model groups is the input token (sub-word) embedding. We think that different input embedding spaces result in different output embedding spaces,

³<https://github.com/microsoft/fastformers>

⁴<https://github.com/huggingface/Transformers>

⁵<https://huggingface.co/distilroberta-base>

⁶<https://github.com/huawei-noah/Pretrained-Language-Model/tree/master/TinyBERT>

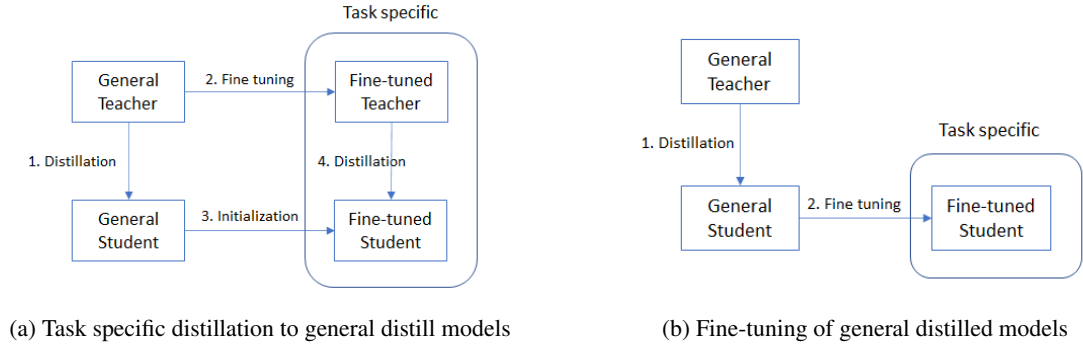


Figure 1: Knowledge distillation methods

Model	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC
BERT [†] (Reference, 12L, 768)	72.7	80.7	57.0	41.8	54.9	65.7	65.6
BERT (Teacher, 12L, 768)	75.99	87.96	64.00	42.00	64.81	69.68	72.41
BERT (Student, 6L, 768)	76.06	90.12	69.00	37.47	64.47	68.59	71.79
BERT (Student, 4L, 312)	72.63	87.63	64.00	36.71	33.68	64.62	65.20
RoBERTa (Teacher, 12L, 768)	81.59	89.34	56.00	50.30	79.66	79.06	71.63
RoBERTa (Student, 6L, 768)	75.19	90.68	57.00	42.90	67.33	66.43	65.83

Table 1: Accuracy of teacher and student models on the validation data set for each task of SuperGLUE benchmark with knowledge distillation. Model marked with [†] represents accuracy numbers on the test set provided by SustaiNLP organizers.

and knowledge transfer with different spaces does not work well.

The result of the knowledge distillation on the tasks are summarized in Table 1 together with the teacher models’ accuracy numbers on the validation data sets. It also includes the accuracy values on test data set for BERT model presented by SustaiNLP 2020 organizers⁷. We train both *cased* and *uncased* models using both task-specific and task-agnostic approaches, and present the model with higher accuracy values. For the more challenging tasks such as MultiRC and ReCoRD, we observe that RoBERTa based models provide better accuracy than BERT based models.

3 Structured Pruning

There has been a significant amount of research on model’s weights pruning approaches inspired by *The Lottery Ticket Hypothesis* (Frankle and Carbin, 2018). Furthermore, there are various papers published to apply this pruning strategy to Transformer models including: (Yu et al., 2019; Sanh et al., 2020; Gordon et al., 2020). Most of such approaches focused on random pruning to reduce the number of parameters following the lottery ticket

hypothesis. While this can reduce the size of the model on the computer storage, it may not improve the inference performance since it is not focusing on better utilization of the computing resources. Since our main focus in *FastFormers* is to improve inference efficiency, randomly pruning a subset of the model’s parameters may not improve performance. In this work, we focus on structured pruning which directly reduces the computation requirements.

Voita et al. (2019); Michel et al. (2019); Hou et al. (2020) proposed methods to prune some of the *heads* of Multi-Head Attention (MHA) in Transformer architecture. *DynaBERT* (Hou et al., 2020) additionally proposed pruning intermediate hidden states in feed-forward layer of Transformer architecture together with rewiring of these pruned attention module and feed-forward layers. In the paper, we define a target model size in terms of the number of heads and the hidden state size of feed-forward network, and use the pruning/distillation approach proposed in DynaBERT as a model compression method. This effectively reduces the dimensions of Transformer models.

Structured pruning methods: The first step of our structured pruning method is to identify the least important *heads* in MHA and the least im-

⁷<https://sites.google.com/view/sustainlp2020/shared-task?authuser=0>

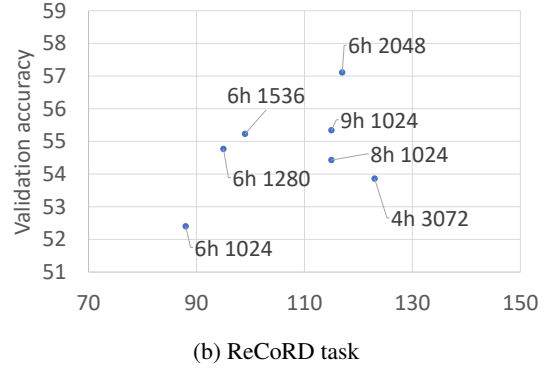
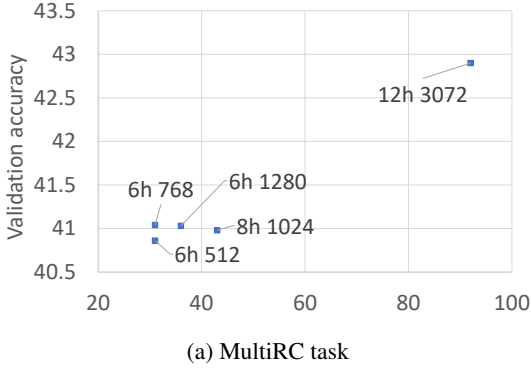


Figure 2: Structured pruning - inference time versus accuracy. Each data point indicates a pruned model with the number of remaining heads and the number of remaining intermediate hidden states of feed-forward layers.

portant hidden states in the feed-forward layers. We use a first order method for computing the importance score, which utilizes the first order gradient information proposed by Michel et al. (2019); Sanh et al. (2020); Hou et al. (2020) instead of using magnitude based pruning (Frankle and Carbin, 2018; Gordon et al., 2020). Before doing the importance score computation, we add a mask variable to each attention head for the gradient computation of the heads. Next, we run forward and backward passes of the model on the entire validation data set, then the absolute values of the gradients are accumulated. These accumulated values are used as importance scores which we use to sort the importance of the heads and the intermediate hidden states. Based on the target model size, we select a given number of top heads and top hidden states from the network. Once the sorting and selection steps are done, we re-group and reconnect the remaining heads and hidden states which result in a smaller sized model. When heads and hidden states are pruned, we use the same pruning ratio across different layers. This enables further optimizations to work seamlessly with the pruned models. In our experiments, we observed that the pruned model can get better accuracy when it goes through another round of knowledge distillation; this has also been noted in Hou et al. (2020). Therefore, we do one more knowledge distillation by using the non-pruned model as a teacher model and the pruned model as an initializer of student model.

Structured pruning results: We apply the structured pruning method mainly to MultiRC and ReCoRD tasks. For the other tasks, the test sets are not that big; therefore, knowledge distillation and other optimizations could make the models quite efficient. In both MultiRC and ReCoRD tasks, the

base model is RoBERTa based distilRoberta-base⁸ which has 6 stacked layers with 768 hidden states, 12 self-attention heads and 3072 hidden states in feed-forward layer. The structured pruning method brings better efficiency by trading off the accuracy. So, it is important to pick a reasonable model size that doesn't compromise the accuracy of the task. Figure 2 presents trade-off between inference speed and accuracy on MultiRC and ReCoRD validation data sets. For the MultiRC task, we could get 2.97x speed-up while losing 1.9 point of accuracy by pruning 50% of heads and 75% of intermediate hidden states from 12 heads and 3072 hidden sizes. For the ReCoRD task, we got 1.95x speed-up while trading off 12.1 point of accuracy by pruning 50% of heads and 50% of hidden states in feed-forward layer. In both cases, they are still exceeding the teacher sized BERT model's accuracy.

4 Low Precision Inference

After knowledge distillation and structured pruning compression, models can benefit from more efficient numerical computations by quantizing the model parameters on CPU and GPU (Rodriguez et al., 2018). It has been shown that the accuracy of Transformer models doesn't get severely compromised when utilizing 8-bit or 16-bit lower precision arithmetic (Devlin, 2017; Kim et al., 2019; Bhandare et al., 2019; Zafrir et al., 2019; Shen et al., 2020; Fan et al., 2020; Aji and Heafield, 2020). Modern CPUs and GPU accelerators are capable of computing those lower numerical arithmetic efficiently. For example, Cascade Lake CPUs have a special 8-bit vector instruction set called AVX (Advanced Vector eXtensions)-512-VNNI and V100

⁸<https://huggingface.co/distilRoberta-base>

GPUs can utilize its efficient Tensor Cores with 16-bit floating point data.

8-bit quantized matrix multiplications on the CPU: 8-bit quantized matrix multiplication brings a significant amount of speed-up compared to 32-bit floating point arithmetic, thanks to the relieved memory bandwidth bottleneck and reduced number of CPU instructions. Efficient utilization of the vector registers and CPU cache memories requires that the parameter weight matrix to be tiled and transposed in a cache efficient layout. This is referred as *packing* operation in the matrix multiplication. Packing itself is a non-negligible operation and repeated packing operation could cancel out all the benefits from the quantized matrix multiplications. Therefore, the result of the packing operation needs to be properly cached to get efficient performance. Moreover, some of the matrix products should stay 32-bit floating point to avoid repeated packing operations. Therefore, we do not use 8-bit matrix product for the Q, K inner product because both matrices are not constant. All the other matrix products have constant weight matrix, so we utilize 8-bit matrix products for them with cached weight packing. For the 8-bit quantized matrix product API, we utilize an open source quantized matrix multiplication library FBGEMM⁹ which we have integrated into the onnxruntime framework. It provides various quantization methods and a separate matrix packing functionality explicitly. We use dynamic quantization method which decides the quantization range of input matrix dynamically every time¹⁰. This enables the quantized values to effectively represent all the values in the input matrix. The weight matrix' quantization range is selected for each column separately and the quantization range for the input matrix is selected for entire input tensor. In our experiments, this 8-bit quantization brings up to around 3.0x speed-up on Cascade Lake CPUs for the Transformer models by trading off small amount of accuracy loss.

16-bit model conversion for the GPU: V100 GPU supports full 16-bit operations for the Transformer architecture. Also, 16-bit floating point operations do not require special handling of inputs and outputs except for having smaller value ranges. The impact of the numerical overflow due

to the smaller range in 16-bit float points is minimal at inference time, so we have not observed any differences in accuracy. Therefore, the model can be fully converted into 16-bit floating point data type before the model is utilized for inference. This 16-bit model conversion brings quite significant speed gain, since the Transformer models are memory bandwidth bound workload. We observe up to 3.53x speed-up depending on the model settings. V100 GPU also supports 8-bit quantized arithmetic, but it is not supported with its efficient Tensor cores; hence we do not utilize 8-bit quantization on GPUs.

5 Runtime Optimization

On top of the structural and numerical optimizations applied, we can utilize various ways to further optimize the computations. In particular, we focus on multi-processing optimizations and computational graph optimization.

Multi-processing optimization: The evaluation machine for the SustaiNLP 2020 shared task has two Cascade Lake 6248 CPUs with 40 physical cores. **The default execution engines for HuggingFace's transformers including PyTorch¹¹ and TensorFlow¹² usually use all available CPU cores for a single operator. This is not the optimal way of utilizing available CPU cores for several reasons. The operators in compressed Transformer architectures are not big enough to fully utilize the parallelism of 40 CPU cores. Therefore, the overheads of parallelizing the operation significantly overshadow the actual gains from the parallelism. Another important factor is that the parallelization to all cores reduces the cache locality and degrades the overall efficiency of CPU utilization.** Therefore, we implement a multiple instance inference by leveraging multiprocessing module of python¹³ programming language. It is preferable to use multi-threading instead of multi-processing to avoid additional copy of program memory, but python's multi-threading cannot really utilize multiple threads due to the Global Interpreter Lock (GIL)¹⁴. Even with this limitation, multi-instance inference with multiprocessing brings much more efficient computation.

¹¹<https://github.com/pytorch/pytorch>

¹²<https://github.com/tensorflow/tensorflow>

¹³<https://docs.python.org/3/library/multiprocessing.html>

¹⁴<https://wiki.python.org/moin/GlobalInterpreterLock>

⁹<https://github.com/pytorch/FBGEMM>

¹⁰https://pytorch.org/tutorials/recipes/recipes/dynamic_quantization.html

Number of inference instances	Time (sec)	Speed-up
Baseline (no thread control)	433	1.00x
1 instance (20 threads/instance)	319	1.36x
2 instances (10 threads/instance)	243	1.78x
4 instance (5 threads/instance)	247	1.75x
5 instance (4 threads/instance)	255	1.70x
10 instance (2 threads/instance)	300	1.44x
20 instance (1 thread/instance)	351	1.23x

Table 2: Speed comparison of different number of inference instances with thread control - time to perform inference on 1,000 ReCoRD validation data samples.

To maximize the cache locality, each inference instance is pinned to specific physical cores using *Linux*’ *taskset* command. Utilizing hyper-threading harms the cache utilization, so we always keep the total number of utilized threads (the number of threads per one inference instance multiplied by the number of instances) not exceeding the number of physical cores in the machine. The optimal number of multiple processes for the best efficiency varies by the model, hardware settings and the data set. We conduct experiment with all target tasks and investigate the best setting for each task. Table 2 shows one example of the speed-up achieved by optimized multi-instance inference on 1,000 samples of validation data set of ReCoRD task.

Computational graph optimizations: Computational graph optimization can further improve the efficiency of the neural network inference by pruning unused graph nodes and fusing multiple operations together. In particular, graph node fusion reduces additional memory allocation and copy which potentially degrade the efficiency. Also, graph node fusion potentially improves parallelism by increasing the size of individual operators. We replace Gaussian Error Linear Units (GELU) with Rectified Linear Units (ReLU) for the computational efficiency while model is distilled without losing any accuracy. We fuse ReLU and bias addition operations followed by matrix multiplications into FBGEMM’s fused post processing operation. Moreover, we utilize multi-head attention node fusion provided by onnxruntime. We use a customized onnxruntime based on v1.3.1.

6 Results

6.1 Combined results

Table 3 and Figure 3 present how the proposed methods work together on BoolQ task using CPU.

For the ablation study, an Azure F16s_v2 instance which has Intel(R) Xeon(R) Platinum 8168 CPUs (8 physical cores) is utilized. When all the optimizations are applied, it could achieve around 233x speed-up while only losing 1.2 of accuracy. **First, it is worth noting that the batch generation code in many frameworks including HuggingFace’s transformers uses fixed sequence length for the input.**

Whenever there comes a shorter sentence than the fixed length for the inference, additional zeros are padded at the end of each input sentence. This degrades CPU performance greatly which has relatively small parallelism than GPU. We modify the batch generation to support dynamic sequence length for each batch to avoid such redundant computation. By doing this dynamic shape batching, we could get around 3.51x speed-up on the test CPU. On top of the dynamic shape batching, all the optimization techniques introduced are applied. As described in Section 2, knowledge distillation brings a significant speed-up which is more than 9 times faster than the original teacher model while preserving the accuracy. One thing to note is that the compressible student model size varies with the task. In this BoolQ example, 4 stacked layers with smaller hidden size (312) model could learn original teacher model’s knowledge without losing any accuracy score on the validation data set. 8-bit quantization together with onnxruntime graph optimization brings around 2.26x speed-up on the 4 layer distilled model. This is lower than 3.0x speed-up acquired when the same methods is applied to 12 layer base size models mentioned in Section 4, because the system is already compressed into a smaller size. By using 8 concurrent and independent instances for inference instead of one inference instance with 8 threads, more than 1.7x additional speed-up could be achieved. Finally, by pruning heads in multi-head attentions and intermediate hidden states in feed-forward layers, another speed-up from 1.38x to 1.81x could be accomplished while trading off the accuracy. One interesting observation is that the structured pruning approach brings better speed-up when it is combined with multi-instance inference which is up to 1.81x. The same pruned model brings at most 1.26x speed-up when it is used with one inference instance. This indicates that multi-instance inference gets more performance benefits when each individual model size gets smaller.

Optimization methods added	Time (sec)	Cumulative speed-up	Speed-up	Accuracy	USD for 100 M queries
Baseline (PyTorch out-of-the-box, 12L, 768)	734.35	1.00x	-	74.01	\$4,223
+ dynamic sequence length	209.29	3.51x	3.51x	74.01	\$1,204
+ knowledge distillation (4L, 312)	22.5	32.64x	9.30x	74.04	\$129
+ 8-bit quantization + graph optimization	9.97	73.66x	2.26x	73.43	\$57
+ multi-instance inference	5.68	129.29x	1.76x	73.43	\$33
+ structured pruning					
25% heads and 25% hidden states pruned	4.11	178.67x	1.38x	73.36	\$24
33% heads and 50% hidden states pruned	3.14	233.87x	1.81x	72.81	\$18

Table 3: An ablation study of CPU inference speed-up on BoolQ validation data set with batch size 1. All performance numbers are measured on an Azure F16s-v2 instance.

6.2 Shared task submissions

We submit our optimized *FastFormers* systems to SustaiNLP 2020 shared task; four systems for track 1 and two systems for track 3. For track 1, we have GPU-only submissions and hybrid submissions which utilize both CPUs and GPUs. We observe that the inference time for ReCoRD only exceeds the inference time of the other tasks all together. Therefore, for the hybrid systems, we use GPUs for ReCoRD task inference and CPUs for all the other tasks. Those CPU and GPU inferences can be executed in parallel for the best throughput. For track 3, we only use CPUs for all tasks as indicated in the shared task description. For the GPU inference, we always limit the number of GPUs we utilize to one. Our highly optimized and compressed models can be executed on a single GPU fast enough. And, the scaling of multiple GPUs is sub-linear (3.0x with 4 GPUs and 1.8x with 2 GPUs) which indicates a single GPU inference is most energy efficient. We apply different types of optimization depending on the time consumed to run the task. For most time consuming tasks, MultiRC and ReCoRD, all compression and optimization techniques mentioned above are applied. All GPU inference uses batch size of 256 which gives the highest throughput and the best efficiency. On the other hand, a single batch works better for most of the cases on CPUs. We use batch size of 1 for all tasks except for CB (batch size of 4) and COPA (batch size of 8).

Table 4 summarizes the accuracy, speed-up numbers and energy savings of *FastFormers* systems prepared for the shared task. For GPUs, CB, COPA and WiC data sets are quite small, so the initial performance without any optimization took 1 or 2 seconds close to the resolution (1 second) of the wall clock time measurement. Therefore, it was

hard to observe big speed improvements for those data sets. For the other tasks, we acquire a good amount of speed-up ranging from 5.8x to 12.4x by our optimization. On CPUs, model compression gives more benefits all across the board in terms of speed. The smallest gain is 9.8x for WiC data set which could not utilize onnxruntime optimization due to the control ‘for’ loop in the output layer. This is currently not supported in onnxruntime. For the other tasks, we observe up to 40.3X speed-up. On the other hand, the energy savings while preserving BERT model accuracy measured by the shared task organizers range from 6.9X up to 125.8X. When combining all the tasks, our best systems save 22.1x energy on CPUs and 21.6x energy on GPUs.

7 Conclusion

In this paper, we have introduced *FastFormers*, which achieves efficient inference-time performance for Transformer-based models on various NLU tasks. We showed that utilizing knowledge distillation, structured pruning and numerical optimization can lead to drastic improvements on inference efficiency. We showed that the improvements can be up to 200X speed-up and results in more than 200X inference cost saving with 22X energy saving. We open source *FastFormers* for the community hoping it can drive more sustainable optimizations for Transformers models. For future directions, some other methods such as early exiting (Xin et al., 2020; Liu et al., 2020; Schwartz et al., 2020; Zhou et al., 2020) and linear time complexity self-attention models (Shen et al., 2018; Wang et al., 2020) could be added to the proposed recipes and possibly improve the efficiency further.

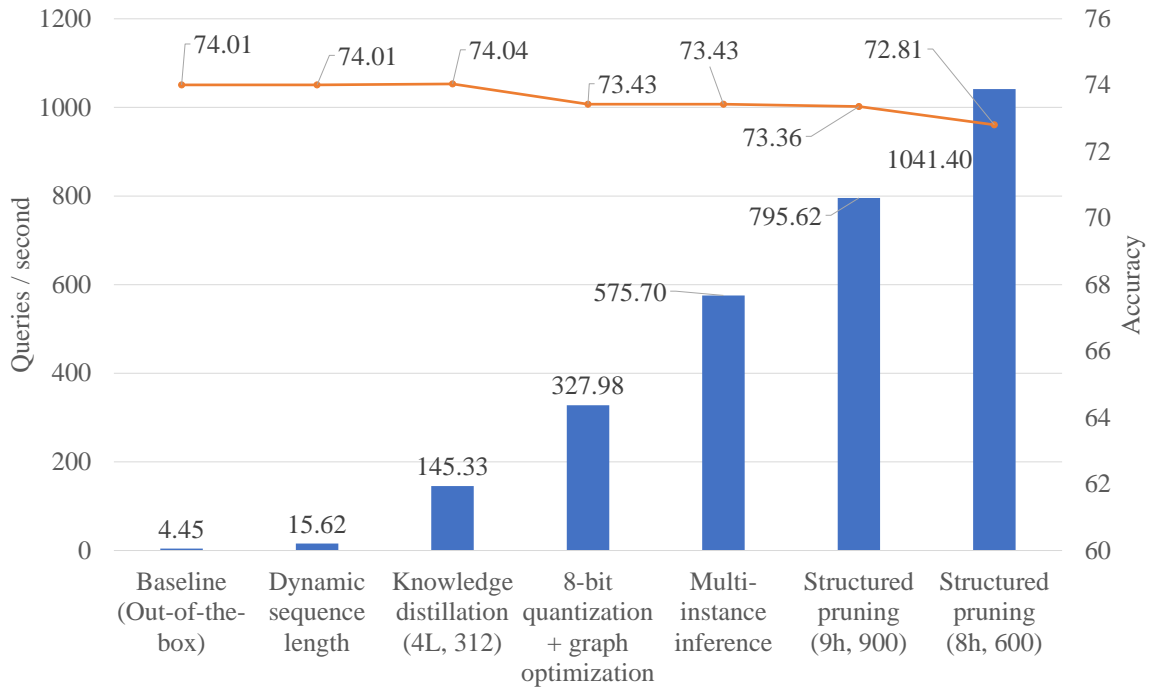


Figure 3: Accuracy versus queries per second with various optimizations on CPU.

Submitted systems	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	Overall
BERT [†] (Reference, 12L, 768)	72.7	80.7	57.0	41.8	54.9	65.7	65.6	62.6
Inference time speed-up	1.0x	1.0x	1.0x	1.0x	1.0x	1.0x	1.0x	1.0x
Energy savings	1.0x	1.0x	1.0x	1.0x	1.0x	1.0x	1.0x	1.0x
System 1 (GPU only)	74.0	82.7	58.0	41.8	56.2	66.4	66.0	63.6
Inference time speed-up	5.8x	2.0x	2.0x	9.1x	12.4x	9.0x	2.0x	11.3x
Energy savings	11.0x	57.9x	125.8x	6.9x	23.9x	20.5x	28.7x	20.3x
System 2 (GPU only)	74.0	82.7	58.0	43.2	56.2	66.4	66.0	63.8
Inference time speed-up	5.8x	2.0x	2.0x	9.1x	12.4x	9.0x	2.0x	11.3x
Energy savings	11.9x	58.3x	80.1x	8.1x	25.0x	15.9x	22.4x	21.6x
System 3 (CPU/GPU hybrid)	73.7	82.7	58.0	43.0	56.2	66.9	65.9	63.8
Inference time speed-up	40.3x	22.0x	38.0x	25.0x	12.4x	22.4x	9.8x	14.5x
Energy savings	13.4x	49.4x	92.8x	13.8x	23.0x	32.7x	30.4x	16.5x
System 4 (CPU/GPU hybrid)	73.7	82.7	58.0	43.1	56.2	66.9	65.9	63.8
Inference time speed-up	40.3x	22.0x	38.0x	17.5x	12.4x	22.4x	9.8x	14.5x
Energy savings	13.2x	45.8x	75.4x	11.8x	23.5x	34.0x	28.4x	16.1x
System 5 (CPU only)	73.7	82.7	58.0	43.0	56.6	66.9	65.9	63.8
Inference time speed-up	40.3x	22.0x	38.0x	25.0x	15.5x	22.4x	9.8x	16.6x
Energy savings	14.0x	41.1x	75.0x	15.3x	22.5x	34.7x	26.2x	22.1x
System 6 (CPU only)	73.7	82.7	58.0	43.1	56.6	66.9	65.9	63.8
Inference time speed-up	40.3x	22.0x	38.0x	17.5x	15.5x	22.4x	9.8x	16.1x
Energy savings	11.4x	22.4x	27.2x	9.2x	16.0x	85.5x	13.5x	15.6x

Table 4: Accuracy numbers with speed-up and energy savings on the test data set of submitted systems to the SustainNLP 2020 shared task. Model marked with [†] was accuracy numbers on the test set provided by organizers.

Acknowledgments

We would like to thank the organizers of SustainNLP 2020 for their various efforts and the reviewers for the thoughtful and helpful suggestions.

References

- Alham Fikri Aji and Kenneth Heafield. 2020. Compressing neural machine translation models with 4-bit precision. In *Proceedings of the Fourth Workshop on Neural Generation and Translation*, pages 35–42.
- Aishwarya Bhandare, Vamsi Sripathi, Deepthi Karkada, Vivek Menon, Sun Choi, Kushal Datta, and Vikram Salelore. 2019. Efficient 8-bit quantization of transformer neural machine language translation model. *arXiv preprint arXiv:1906.00532*.
- Jacob Devlin. 2017. Sharp models on dull hardware: Fast and accurate neural machine translation decoding on the cpu. *arXiv preprint arXiv:1705.01991*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Rémi Gribonval, Hervé Jégou, and Armand Joulin. 2020. Training with quantization noise for extreme model compression. *arXiv Preprint arXiv:2004.7320:1–18*.
- Jonathan Frankle and Michael Carbin. 2018. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*.
- Mitchell A Gordon, Kevin Duh, and Nicholas Andrews. 2020. Compressing bert: Studying the effects of weight pruning on transfer learning. *arXiv preprint arXiv:2002.08307*.
- Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. 2020. Towards the systematic reporting of the energy and carbon footprints of machine learning. *arXiv preprint arXiv:2002.05651*.
- Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. *Distilling the knowledge in a neural network*. In *NIPS Deep Learning and Representation Learning Workshop*.
- Lu Hou, Lifeng Shang, Xin Jiang, and Qun Liu. 2020. Dynabert: Dynamic bert with adaptive width and depth. *arXiv preprint arXiv:2004.04037*.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2019. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*.
- Young Jin Kim, Marcin Junczys-Dowmunt, Hany Hassan, Alham Fikri Aji, Kenneth Heafield, Roman Grundkiewicz, and Nikolay Bogoychev. 2019. From research to production and back: Ludicrously fast neural machine translation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 280–288.
- Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Haotang Deng, and Qi Ju. 2020. Fastbert: a self-distilling bert with adaptive inference time. *arXiv preprint arXiv:2004.02178*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In *Advances in Neural Information Processing Systems*, pages 14014–14024.
- Andres Rodriguez, Eden Segal, Etay Meiri, Evarist Fomenko, Y Jim Kim, Haihao Shen, and Barukh Ziv. 2018. Lower numerical precision deep learning inference and training. *Intel White Paper*, 3.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Victor Sanh, Thomas Wolf, and Alexander M Rush. 2020. Movement pruning: Adaptive sparsity by fine-tuning. *arXiv preprint arXiv:2005.07683*.
- Roy Schwartz, Gabi Stanovsky, Swabha Swayamdipta, Jesse Dodge, and Noah A Smith. 2020. The right tool for the job: Matching model and instance complexities. *arXiv preprint arXiv:2004.07453*.
- Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2020. Q-bert: Hessian based ultra low precision quantization of bert. In *AAAI*, pages 8815–8821.
- Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. 2018. Efficient attention: Attention with linear complexities. *arXiv preprint arXiv:1812.01243*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Elena Voita, David Talbot, Fedor Moiseev, Rico Senrich, and Ivan Titov. 2019. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint arXiv:1905.09418*.

- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. In *Advances in Neural Information Processing Systems*, pages 3266–3280.
- Sinong Wang, Belinda Li, Madian Khabisa, Han Fang, and Hao Ma. 2020. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*.
- Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. 2020. Deebert: Dynamic early exiting for accelerating bert inference. *arXiv preprint arXiv:2004.12993*.
- Haonan Yu, Sergey Edunov, Yuandong Tian, and Ari S Morcos. 2019. Playing the lottery with rewards and multiple languages: lottery tickets in rl and nlp. *arXiv preprint arXiv:1906.02768*.
- Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. 2019. Q8bert: Quantized 8bit bert. *arXiv preprint arXiv:1910.06188*.
- Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. Bert loses patience: Fast and robust inference with early exit. *arXiv preprint arXiv:2006.04152*.