

# MoEfication: Conditional Computation of Transformer Models for Efficient Inference

Anonymous ACL submission

## Abstract

Transformer-based pre-trained language models achieve superior performance on most NLP tasks due to large parameter capacity, but also lead to huge computation cost. Fortunately, we observe that most inputs only activate a tiny ratio of neurons of large Transformer-based pre-trained models during inference. Hence, we propose to convert a model into its mixture-of-experts (MoE) version with the same parameters, namely MoEfication, which accelerates large-model inference by conditional computation based on the sparse activation phenomenon. Specifically, MoEfication consists of two phases: (1) splitting the parameters of feed-forward neural networks (FFNs) into multiple parts as experts, and (2) building expert routers to decide which experts will be used for each input. Experimental results show that MoEfication can save 80% computation cost of FFNs while maintaining over 95% original performance for different models, including models with different sizes (up to 3 billion parameters) and distilled models, on various downstream tasks. Moreover, we find that the MoEfied model achieves better performance than the MoE model pre-trained from scratch with the same model size. We will release all the code and models of this paper.

## 1 Introduction

Recent years have witnessed an exponential increase in the size of Transformer-based pre-trained language models (PLMs) (Han et al., 2021). From BERT (Devlin et al., 2019) in 2018 to GPT-3 (Brown et al., 2021) in 2020, the number of parameters has already increased by nearly 600 times. Moreover, the exploration of larger models is continuing. The increasing model size significantly improves the model performance on a variety of downstream NLP tasks (Raffel et al., 2020; He et al., 2021b), but also comes with huge computation cost, which limits the potential applications of

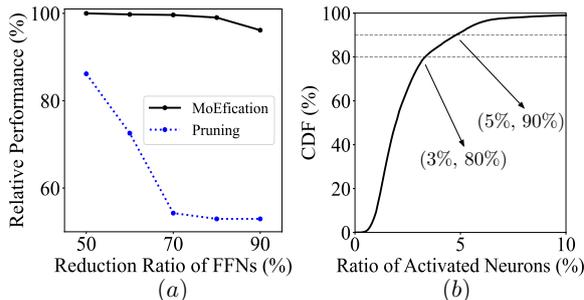


Figure 1: Results of a fine-tuned T5-Large (Raffel et al., 2020) on SST-2. (a) Relative performance compared to the original performance with different reduction ratios of FFNs. Large pruning ratios significantly degrade the performance. (b) Cumulative distribution function (CDF) of the ratio of activated neurons for each input. SST-2’s training set is used as inputs. 90% inputs only activate less than 5% neurons of FFNs.

large-scale PLMs. Hence, it is essential to explore novel techniques to make PLMs more efficient.

The computation of Transformer mainly consists of two parts: attention networks and feed-forward neural networks (FFNs). Much effort has been made to reduce the cost of attention networks (Beltagy et al., 2020; Kitaev et al., 2020; Tay et al., 2020) while little has been made for FFNs. Previous work on the acceleration of FFNs usually uses general pruning algorithms and ignores the characteristics of FFNs (Li et al., 2020; Xu et al., 2021). Hence, large pruning ratios will lead to poor results as shown in Figure 1. In this work, we explore to further accelerate FFNs beyond model pruning.

Fortunately, according to our observation on FFNs in Transformer models, we find a phenomenon of **sparse activation**, i.e., only a tiny fraction of neurons are activated for a single input. As shown in Figure 1, when we perform inference on a fine-tuned T5-Large model with 700-million parameters, 90% inputs only activate less than 5% neurons<sup>1</sup>. Hence, we can omit the computation of

<sup>1</sup>T5 uses ReLU as the activation function. We treat the neurons having positive outputs as activated neurons.

064 inactive neurons to reduce the cost. Meanwhile,  
065 most neurons will be eventually activated by some  
066 inputs. As a result, model pruning is not applicable  
067 and will significantly degrade the performance. In-  
068 stead of model pruning, we explore efficient FFNs  
069 based on conditional computation (Bengio, 2013),  
070 which selectively activates parts of the network ac-  
071 cording to input. This mechanism **naturally** exists  
072 in the FFNs of pre-trained Transformers.

073 Inspired by the sparse activation phenomenon,  
074 we propose to convert a large-scale PLM into its  
075 mixture-of-experts (MoE) version with the same  
076 parameters for efficient conditional computation  
077 in inference, namely *MoEfication*. Different from  
078 previous work on MoE Transformers that typically  
079 *breeds* models into multiple experts (Lepikhin et al.,  
080 2021; Fedus et al., 2021), *MoEfication* aims to  
081 *split* existing models into multiple experts while  
082 keeping the model size unchanged. We expect an  
083 *MoEfied* model will improve the model efficiency  
084 and maintain the performance of the original model  
085 by dynamically selecting experts.

086 *MoEfication* consists of two phrases. (1) **Expert**  
087 **Construction**: Split a whole feed-forward layer  
088 into multiple experts. The goal is to group those  
089 neurons that are often activated simultaneously into  
090 the same expert network. To achieve this goal, we  
091 build a co-activation graph based on the activation  
092 results and divide this graph into several subgraphs  
093 as experts by graph partition. (2) **Expert Selection**:  
094 Select those experts that contain as many activated  
095 neurons as possible for each input to approximate  
096 to the original results. To reach this target, we  
097 first find the best selections based on the activation  
098 results and then use them to train shallow neural  
099 networks as expert routers.

100 In the experiments, we validate the effectiveness  
101 of *MoEfication* on two typical kinds of downstream  
102 tasks, including GLUE and QA benchmarks (Wang  
103 et al., 2019; Rajpurkar et al., 2016; Lai et al., 2017),  
104 using T5 with different sizes (Raffel et al., 2020).  
105 Experimental results show that *MoEfication* can  
106 save 80% computation cost of FFNs while main-  
107 taining over 95% original performance for both  
108 conventional models (the number of parameters  
109 varies from 60 millions to 3 billions) and distilled  
110 models. Besides, we find that the *MoEfied* model  
111 achieves better performance than the MoE model  
112 pre-trained from scratch with the same model size.  
113 Then, we study the routing patterns of *MoEfied*  
114 models and hope these findings can help future

work on the design and training of MoE models. 115

## 2 Related Work 116

**Model Acceleration for PLMs.** Model accelera- 117  
tion aims to reduce the time and space complex- 118  
ity of PLMs for faster inference and deployment 119  
on resource-constrained devices. There are sev- 120  
eral techniques for model acceleration, including 121  
knowledge distillation (Sanh et al., 2019; Sun et al., 122  
2019; Jiao et al., 2020), model pruning (Voita et al., 123  
2019; Michel et al., 2019; Zhang et al., 2021), 124  
model quantization (Zafir et al., 2019; Zhang et al., 125  
2020), and dynamic inference (Xin et al., 2020; 126  
Li et al., 2021). Among these techniques, model 127  
pruning and dynamic inference explore to omit un- 128  
necessary computation for acceleration, which is 129  
similar to the target of *MoEfication*. Different from 130  
model pruning, which omits redundant param- 131  
eters, *MoEfication* keeps the original model size and 132  
dynamically selects parts of parameters at a time. 133  
For dynamic inference, previous work focuses on 134  
how to dynamically drop layers to accelerate infer- 135  
ence (Huang et al., 2018; Wu et al., 2020; Li et al., 136  
2021). In this manner, the output of each layer is ex- 137  
pected to be able to predict labels, and hence it will 138  
introduce additional training objectives and predic- 139  
tion strategies. In contrast, *MoEfication* simplifies 140  
models in a finer granularity, and does not change 141  
the process of training and inference. In summary, 142  
*MoEfication* can be regarded as a novel direction 143  
diagonal with the above-mentioned approaches. 144

**Large-scale PLMs with MoE.** Jacobs et al. 145  
(1991) propose mixture-of-experts to build a sys- 146  
tem composed of many separate networks, which 147  
learn to handle a subset of the training examples in- 148  
dependently. When deep neural networks achieve 149  
great success (Hinton et al., 2012; Krizhevsky et al., 150  
2012; Goodfellow et al., 2013), Bengio (2013) 151  
thinks the model size is a key factor and MoE 152  
is an important technique to scaling model com- 153  
putation and proposes the idea of “conditional 154  
computation”. The first large-scale MoE lan- 155  
guage model is proposed by Shazeer et al. (2017), 156  
which adds an MoE layer between two LSTM lay- 157  
ers and independently assigns tokens to combi- 158  
nations of experts. Recently, GShard (Lepikhin 159  
et al., 2021), Switch-Transformer (Fedus et al., 160  
2021), BASELayer (Lewis et al., 2021), and Hash- 161  
Layer (Roller et al., 2021) study how to build large- 162  
scale Transformer-based models with MoE and op- 163  
timal training strategies, which can fully utilize the 164

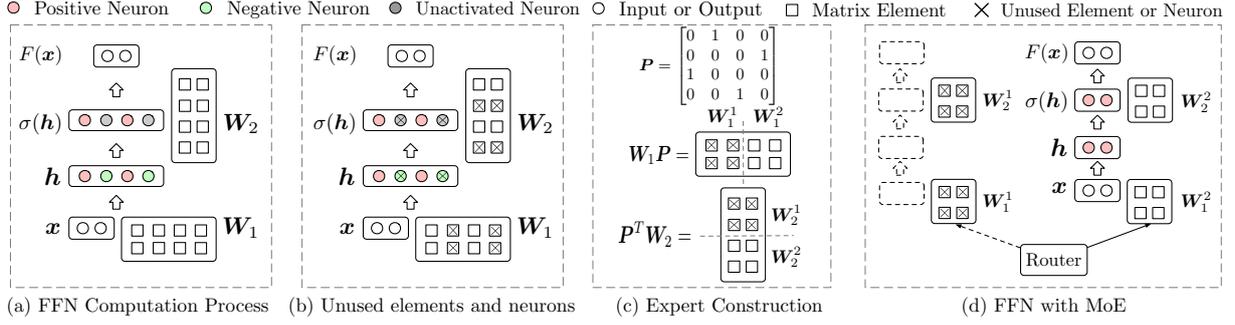


Figure 2: An example of the sparse activation phenomenon and MoEfication. (a) shows the computation process of an FFN for a given input. (b) shows the unused elements and neurons for this input. (c) shows how to construct experts. (d) shows how the MoEfied model handles this input efficiently.

model capacity. Different from them, we utilize the naturally-existing sparse activation phenomenon to convert a model into its MoE version for better efficiency during inference.

### 3 Method

In this section, we will introduce the general idea of MoEfication and divide it into two phases: expert construction and expert selection.

#### 3.1 Overall Framework

MoEfication aims to utilize the sparse activation phenomenon in the FFNs of Transformers to reduce the computation cost.

We first formally describe the sparse activation phenomenon. The FFNs of Transformers are two-layer fully connected networks, which process an input representation  $x \in \mathbb{R}^{d_{model}}$  by

$$\begin{aligned} h &= xW_1 + b_1, \\ F(x) &= \sigma(h)W_2 + b_2, \end{aligned} \quad (1)$$

where  $W_1 \in \mathbb{R}^{d_{model} \times d_{ff}}$  and  $W_2 \in \mathbb{R}^{d_{ff} \times d_{model}}$  are the weight matrices,  $b_1 \in \mathbb{R}^{d_{ff}}$  and  $b_2 \in \mathbb{R}^{d_{model}}$  are the bias vectors, and  $\sigma(\cdot)$  is a non-linear activation function, which prefers to retain positive values and discard negative ones. In this work, we study the activation function ReLU (Nair and Hinton, 2010), which is used by the original Transformer (Vaswani et al., 2017) and some widely-used Transformer-based PLMs (Sun et al., 2020; Raffel et al., 2020).

As shown in Figure 1, there are many inactive (zero) values in the intermediate output  $\sigma(h)$ . The computation of these values can be omitted for acceleration. Meanwhile, different inputs will activate different neurons. Hence, we explore to select the possibly-activated neurons of  $h$  before the FFN computation instead of model pruning.

We show an example in Figure 2. In this FFN,  $d_{model}$  is 2,  $d_{ff}$  is 4, and the bias vectors are omitted for simplification. For a given input representation  $x$ , there are two positive values in  $h$ . Hence, we only need to compute part of the FFN, i.e., a  $2 \times 2$  submatrix of  $W_1$  and a  $2 \times 2$  submatrix of  $W_2$ , to obtain the same output  $F(x)$ . Correspondingly, we can MoEfy the original FFN to have an MoE layer with two experts and select the one on the right-hand side for this input  $x$ .

For MoEfication, we first split the FFN into several independent parts, namely expert construction, and then design a router to select suitable experts for each input, namely expert selection.

#### 3.2 Expert Construction

In this subsection, we introduce how to split an FFN into several parts. The core idea is to group together the neurons that are often activated simultaneously. In this way, for each input, we can select a small number of experts to cover all its activated neurons. To achieve better parallel computation performance, we set the size of each expert to be the same. If the number of experts is  $k$ , the input and output dimension of experts is still  $d_{model}$  and their intermediate dimension is  $d_e = \frac{d_{ff}}{k}$ . Then, the parameters of  $i$ -th expert are denoted by

$$W_1^i \in \mathbb{R}^{d_{model} \times d_e}, b_1^i \in \mathbb{R}^{d_e}, W_2^i \in \mathbb{R}^{d_e \times d_{model}}. \quad (2)$$

Given the result of splitting, we construct the corresponding permutation of intermediate neurons by  $\begin{pmatrix} 1 & 2 & \dots & d_{ff} \\ f(1) & f(2) & \dots & f(d_{ff}) \end{pmatrix}$ , where  $f(n)$  is the mapping function from the original neuron index to the permuted neuron index. We compute  $f(n)$  by

$$f(n) = (e(n) - 1)d_e + |\{m|m \leq n, e(m) = e(n)\}|, \quad (3)$$

where  $e(n)$  is the expert index of the  $n$ -th neuron, which varies from 1 to  $k$ , and  $|\{m|m \leq n, e(m) =$

$e(n)$  is the index of the  $n$ -th neuron in the expert. Then, we use its permutation matrix  $\mathbf{P} \in \mathbb{R}^{d_{ff} \times d_{ff}}$  to permute the rows or columns of parameters and have the following split:

$$\begin{aligned} [\mathbf{W}_1^1, \mathbf{W}_1^2, \dots, \mathbf{W}_1^k] &= \mathbf{W}_1 \mathbf{P}, \\ \mathbf{b}_1^1 \oplus \mathbf{b}_1^2 \oplus \dots \oplus \mathbf{b}_1^k &= \mathbf{b}_1 \mathbf{P}, \\ [(\mathbf{W}_2^1)^T, (\mathbf{W}_2^2)^T, \dots, (\mathbf{W}_2^k)^T] &= (\mathbf{P}^T \mathbf{W}_2)^T, \end{aligned} \quad (4)$$

where  $\oplus$  represents the vertical concatenation. Note that the permutation will not influence the output representation:

$$\begin{aligned} \sigma(\mathbf{h})\mathbf{W}_2 + \mathbf{b}_2 &= \sigma(\mathbf{h})\mathbf{P}\mathbf{P}^T\mathbf{W}_2 + \mathbf{b}_2, \\ &= \sigma(\mathbf{h}\mathbf{P})\mathbf{P}^T\mathbf{W}_2 + \mathbf{b}_2, \\ &= \sigma(\mathbf{x}\mathbf{W}_1\mathbf{P} + \mathbf{b}_1\mathbf{P})\mathbf{P}^T\mathbf{W}_2 + \mathbf{b}_2. \end{aligned} \quad (5)$$

In this work, we propose two methods to split an FFN into  $k$  parts.

**Parameter Clustering Split.** To take the parameter information into consideration, we treat the columns of  $\mathbf{W}_1$  as a collection of vectors with  $d_{model}$  dimension. Based on the intuition that the neurons with similar vectors will be activated simultaneously, we apply balanced K-Means (Malinen and Fränti, 2014) to the vector collection to obtain  $k$  clusters to construct the mapping function.

**Co-Activation Graph Split.** To directly use the information of co-activation, we construct a co-activation graph by counting co-activations of PLMs for the samples of the training set. Each neuron will be represented by a node in the graph, and the edge weight between two nodes are their co-activation values. The co-activation value is computed by

$$\text{co-activation}(n, m) = \sum_{\mathbf{x}} \mathbf{h}_n^{(\mathbf{x})} \mathbf{h}_m^{(\mathbf{x})} \mathbb{1}_{\mathbf{h}_n^{(\mathbf{x})} > 0, \mathbf{h}_m^{(\mathbf{x})} > 0}, \quad (6)$$

where  $\mathbf{h}_n^{(\mathbf{x})}$ ,  $\mathbf{h}_m^{(\mathbf{x})}$  are the  $n$ -th and the  $m$ -th neurons of  $\mathbf{h}$  for the input  $\mathbf{x}$  and  $\mathbb{1}_{\mathbf{h}_n^{(\mathbf{x})} > 0, \mathbf{h}_m^{(\mathbf{x})} > 0}$  indicates  $\mathbf{h}_n^{(\mathbf{x})}$  and  $\mathbf{h}_m^{(\mathbf{x})}$  are activated simultaneously. Then, we apply graph partitioning algorithms (Karypis and Kumar, 1998) to the co-activation graph to obtain the split, where the internal connections for each group will be strong. It means that the neurons splitted into the same group are often activated simultaneously for the training samples.

### 3.3 Expert Selection

In this subsection, we introduce how to create a router for expert selection. An MoEified FFN processed an input  $\mathbf{x}$  by

$$F_m(\mathbf{x}) = \sum_{i \in S} \sigma(\mathbf{x}\mathbf{W}_1^i + \mathbf{b}_1^i)\mathbf{W}_2^i + \mathbf{b}_2, \quad (7)$$

where  $S$  is the set of the selected experts. If all experts are selected, we have  $F_m(\mathbf{x}) = F(\mathbf{x})$ . Considering that  $\sigma(\mathbf{x}\mathbf{W}_1^i + \mathbf{b}_1^i)\mathbf{W}_2^i$  equals to  $\mathbf{0}$  for most experts, we try to select  $n$  experts, where  $n < k$ , minimize  $\|F_m(\mathbf{x}) - F(\mathbf{x})\|_2$ . The selection methods will assign a score  $s_i$  to each expert for the given input  $\mathbf{x}$  and select the experts with the  $n$  highest scores by

$$S = \underset{A \subset \{1, 2, \dots, k\}, |A|=n}{\text{arg max}} \sum_{i \in A} s_i. \quad (8)$$

**Groundtruth Selection** for the intermediate output  $\sigma(\mathbf{h})$ . We can obtain the groundtruth selection, which minimizes  $\|\sum_{i \in S} \sigma(\mathbf{x}\mathbf{W}_1^i + \mathbf{b}_1^i) - \sigma(\mathbf{h})\|_2$ , by a greedy algorithm. We calculate the sum of positive values in each expert as  $s_i$  and select experts using Equation 8. This selection should approximate to the lower bound of  $\|F_m(\mathbf{x}) - F(\mathbf{x})\|_2$ . Correspondingly, its performance will approximate to the ideal performance of an MoEified model. Meanwhile, it is intractable to directly optimize  $\|F_m(\mathbf{x}) - F(\mathbf{x})\|_2$  because there are too many possible combinations of experts.

**Similarity Selection.** To utilize the parameter information, we average all columns of  $\mathbf{W}_1^i$  and use it as the expert representation. Given an input  $\mathbf{x}$ , we calculate the cosine similarity between the expert representation and  $\mathbf{x}$  as  $s_i$ .

**MLP Selection.** We train a multi-layer perceptron (MLP), which takes the  $\mathbf{x}$  as input and predicts the sum of positive values in each expert. Then, we use the prediction as  $s_i$ . This method tries to approximate to the performance of groundtruth selection.

## 4 Experiment

### 4.1 Experimental Setups

**Models and Hyperparameters** We use four variants of T5 (Raffel et al., 2020), which are the 60-million-parameter T5-Small, the 200-million-parameter T5-Base, the 700-million-parameter T5-Large, and the 3-billion-parameter T5-XLarge. The non-linear activation function is ReLU (Nair and Hinton, 2010). We use Adam as the optimizer and a learning rate of  $10^{-6}$  for fine-tuning on downstream tasks. The batch size is set to 64 and the number of epochs is set to 3.

**Datasets.** We use several natural language understanding datasets to evaluate our models. For text classification, we use GLUE benchmark (Wang et al., 2019), including MNLI-matched (Williams

et al., 2018), QNLI (Rajpurkar et al., 2016), QQP<sup>2</sup>, RTE (Dagan et al., 2006), SST-2 (Socher et al., 2013), MRPC (Dolan and Brockett, 2005), CoLA (Warstadt et al., 2019), and STS-B (Giampiccolo et al., 2007). For reading comprehension, we use SQuAD (Rajpurkar et al., 2016) and RACE (Lai et al., 2017), which are the representative datasets for span extraction and multi-choice QA, respectively. We report the results on their development sets. For MNLI, QNLI, QQP, RTE, SST-2, MRPC, RACE, we use accuracy as the metric. For CoLA, we use matthews correlation coefficient as the metric. For STS-B, we use pearson and spearman correlation as the metrics. For SQuAD, we use F1 score as the metric.

**Expert Construction.** For balanced K-Means, we use an open-source implementation<sup>3</sup>. Besides Parameter Clustering Split and Co-activation Graph Split, we also implement Random Split as a naive baseline, which uses an identity matrix as  $P$ . We set the number of neurons in each expert to 32. Correspondingly, the number of experts varies from 64 to 512 for different T5 variants. With the same expert size, the relative computation cost of routing is the same as shown in Appendix.

**Expert Selection.** Besides Similarity Selection and MLP Selection, we also implement Random Selection, where we treat each expert as a collection of vectors with  $d_{model}$  dimension and randomly select one of them as the expert representation. For Random Selection and Similarity Selection, the computation complexity for routing is  $O(kd_{model})$ . For MLP Selection, we use a two-layer feed-forward network as the architecture. The input dimension is  $d_{model}$ , the intermediate dimension is  $k$ , and the output dimension is  $k$ . The non-linear activation function is  $\tanh(\cdot)$ . Its computation complexity is  $O(kd_{model} + k^2)$ . Compared to the computation complexity of FFNs of the original model,  $O(d_{model} \cdot d_{ff})$ , the computation cost of routers is ignorable because  $k$  is much smaller than  $d_{ff}$ . For example,  $k$  is 128 and  $d_{ff}$  is 4096 for T5-Large. For the training of our MLP routers, we adopt cross-entropy as the training objective and use the Adam optimizer with the learning rate of  $10^{-2}$ . The batch size is set to 512 and the number of epochs is set to 10. We sample nearly 500 thousand input representations from the training data and split them into the training and develop-

<sup>2</sup><https://data.quora.com>

<sup>3</sup><https://github.com/ndanielsen/Same-Size-K-Means>

Model	SST-2	MNLI	RACE
Small	90.9	82.4	44.7
Small-Distill	91.9	82.6	50.6
Base	94.0	86.4	71.7
Large	96.2	89.5	81.3
XLarge	96.9	90.5	85.6

Table 1: Original Performance of different models on three downstream tasks. The model architecture is T5.

ment sets with the ratio of 9 : 1. Note that we only use the activation information as supervision. The training time of each FFN is about several minutes on a single GPU.

## 4.2 MoEfication with Different Models

In this subsection, we evaluate MoEfication on different PLMs. We consider two factors: the model size and whether the model is compressed. For the model size, we use five variants of T5 (Raffel et al., 2020), from T5-Small to T5-XLarge. For convenience, we directly use the scale names as the abbreviations. To investigate the influence of model compression, we compress T5-Large to T5-Small by classic knowledge distillation (Hinton et al., 2015). Specifically, the teacher model is a fine-tuned T5-Large and the student model is a pre-trained T5-Small. The distilled model is denoted by T5-Small-Distill. The expert construction and selection methods used here are Co-activation Graph Split and MLP Selection, which are proved to be the best combination in Section 4.4.

We report the performance of these models on three datasets, SST-2, MNLI, and RACE, in Table 1. They are the representative datasets for single-sentence classification, sentence-pair classification, and reading compression, respectively. The original performance of PLMs grows as the model size grows, and knowledge distillation improves the performance of T5-small.

We first calculate the activation statistics of different models by inputting the training data of each dataset. The results are shown in Figure 3. From the figure, we have three observations. (1) The activations of these models are sparse. Different from the previous study on models trained with smaller datasets, where the activation ratios are range from 10% to 50% (Geva et al., 2021)<sup>4</sup>, we find most inputs activate less than 10% of the neurons. (2) The activations of larger models are sparser than

<sup>4</sup>Since the activation ratios of a randomly-initialized model are around 50%, we guess these models do not make full use of their parameters.

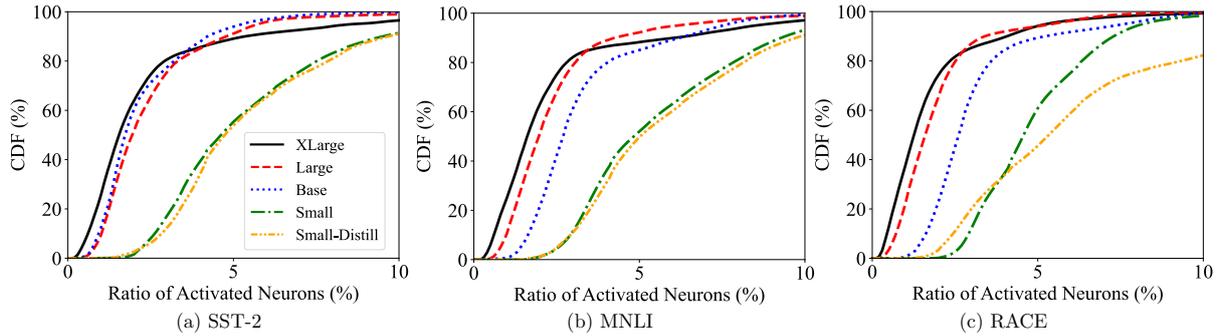


Figure 3: CDF of the ratio of activated neurons for each input with different models on three datasets.

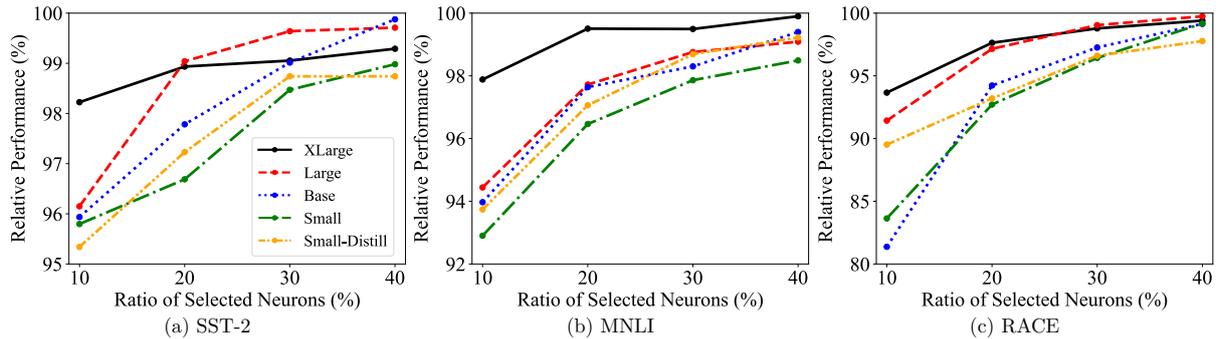


Figure 4: Relative performance of MoEified models with different sizes on three datasets. Dynamically selecting 10% to 20% neurons can recover nearly 98% original performance for large models such as T5-XLarge.

those of smaller models. For example, 80% inputs only activate less than 3% neurons in T5-XLarge while 40% inputs activate more than 3% neurons in T5-Small. (3) The sparsity is less related to distillation than the model size. The CDF curve of T5-Small-Distill is close to that of T5-Small.

Then, we compare the performance of MoEified models with different sizes and ratios of selected neurons and report the results in Figure 4. To measure the performance of MoEification, we calculate the relative performance of the MoEified model to the original model. From the figure, we have four observations. (1) MoEification works well with all models on all three datasets. MoEified models save 80% computation cost of FFNs while maintaining over 95% original performance. (2) The larger models can use fewer neurons to recover the original performance. For example, T5-XLarge achieves nearly 98% relative performance on SST-2 and MNLI with 10% neurons while T5-Small achieves the same results with 30% to 40% neurons. This result is consistent with the activation statistics, that is, larger models are sparser. We can expect that MoEification can work better with super large models. (3) Difficult tasks require models to select more experts to maintain the perfor-

mance. From Table 1, we can see that the accuracy of RACE is much lower than the other two tasks, and hence we think RACE is more difficult. Correspondingly, the relative performance with 10% neurons on RACE is also lower than those on the other tasks. (4) MoEification works similarly on T5-Small and T5-Small-Distill, which indicates that MoEification can work with knowledge distillation for more efficient inference.

### 4.3 Parameter Calibration

In practice, there is still a gap between the performance of MoEified models and that of original models because selected experts cannot cover all positive neurons with a limited computation budget. Hence, the outputs of MoEified models will be slightly different from those of original models. To calibrate MoEified models, we further fine-tune the models on the training set, namely parameter calibration. Considering that current routers are based on the first layers of FFNs ( $\mathbf{W}_1$  and  $\mathbf{b}_1$ ), we only optimize the second layers of FFNs ( $\mathbf{W}_2$  and  $\mathbf{b}_2$ ) to ensure routers can also work well after fine-tuning.

We evaluate this method on several downstream natural language understanding tasks with T5-Large. The ratio of selected neurons is set to 20%,

	MNLI	QNLI	QQP	RTE	SST-2	MRPC	CoLA	STS-B	RACE	SQuAD 1.1	Avg.
Original	89.5	94.4	91.7	87.1	96.2	88.0	59.4	91.2/90.9	81.3	93.2	87.2
MoEified	87.5	93.2	90.2	86.4	95.4	87.5	55.5	90.6/90.3	79.0	92.2	85.7 (-1.5)
+GT	89.1	94.1	91.4	86.4	96.3	88.3	58.8	90.9/90.8	80.8	93.2	86.9 (-0.3)
+Calib	88.7	93.6	91.3	87.5	96.2	89.3	59.4	91.0/90.6	79.9	92.3	86.9 (-0.3)

Table 2: Results of T5-Large on GLUE benchmark and two QA datasets. The last row reports the differences between the original model and MoE+Calib. MoEified models with parameter calibration achieve comparable performance to original models.

Construction	Selection	SST-2	MNLI	RACE
-	-	96.2	89.5	81.3
Random	Groundtruth	95.9	87.3	80.0
	Random	65.9	36.3	29.2
	Similarity	90.3	75.9	56.7
	MLP	94.1	84.1	75.0
Parameter Clustering	Groundtruth	95.5	88.8	80.9
	Random	70.6	36.4	41.8
	Similarity	86.7	66.3	63.6
	MLP	<b>95.9</b>	<b>87.5</b>	<b>78.7</b>
Co-Activation Graph	Groundtruth	96.3	89.1	80.8
	Random	85.3	68.5	54.7
	Similarity	92.2	81.4	71.0
	MLP	95.4	<b>87.5</b>	<b>79.0</b>

Table 3: Comparisons of different combinations of expert construction and selection methods using T5-Large. The first row is the original performance. The best results in each group are underlined and the best results on each dataset are in **boldface**.

which is sufficient for T5-Large as show in Figure 3. We use a small learning rate of  $10^{-7}$  for calibration. The other hyper-parameters remain the same as fine-tuning. The results are shown in Table 2. MoEified refers to the combination of Co-activation Graph Split and MLP Selection. MoEified+GT refers to the combination of Co-activation Graph Split and Groundtruth Selection. MoEified+Calib is the calibrated version of MoEified.

We observe that MoEfication introduces small performance loss (about 1.5% on average) with an 80% reduction of the computation cost in FFNs. Meanwhile, calibration can effectively deal with the issue of the precision errors brought by MoEfication. For example, MoEified+Calib improves MoEified by nearly 4% on CoLA and achieves the same average performance as MoEified+GT.

#### 4.4 Comparisons of MoEfication Strategies

To find the most effective MoEfication strategy, we evaluate different combinations of expert construction and selection methods. We use T5-Large and also set the ratio of selected neurons to 20%. The

Model	MLM Loss
MoE Pre-training	3.09
Standard Pre-training	2.88 (-0.21)
+MoEfication	3.02 (-0.07)
+GT	2.95 (-0.14)

Table 4: Comparisons of MoE models pre-trained from scratch and modified by MoEfication. We report the MLM loss on the validation set. Standard pre-training with MoEfication is better than pre-training a MoE model from scratch.

results are shown in Table 3. From the table, we have two observations:

(1) For expert construction, Co-activation Graph Split is the best method according to the overall performance. Compared to the other two methods, Co-activation Graph Split directly uses the co-activation information to group the neurons activating simultaneously into the same expert.

(2) For expert selection, the performance of Groundtruth Selection is close to that of the original model, which indicates that 20% parameters of FFNs are sufficient to achieve good performance on T5-Large. Meanwhile, MLP Selection is the best expert selection method and can work well with both Parameter Clustering Split and Co-activation Graph Split.

#### 4.5 MoEfication vs. MoE pre-training

In this subsection, we compare the performance of two kinds of MoE models. The first one is pre-trained from scratch. The second one is transformed from a standard model by MoEfication. For fair comparisons, we pre-train one MoE model and one standard model with the same model size from scratch using WikiText-103 (Merity et al., 2017). The pre-training objective is masked language modeling (MLM). The model architecture is the same as T5-Small. For pre-training, we use the batch size of 4096, the learning rate of 0.01, the maximum sequence length of 512, and the Adam optimizer. The number of experts is set to 64 and the router

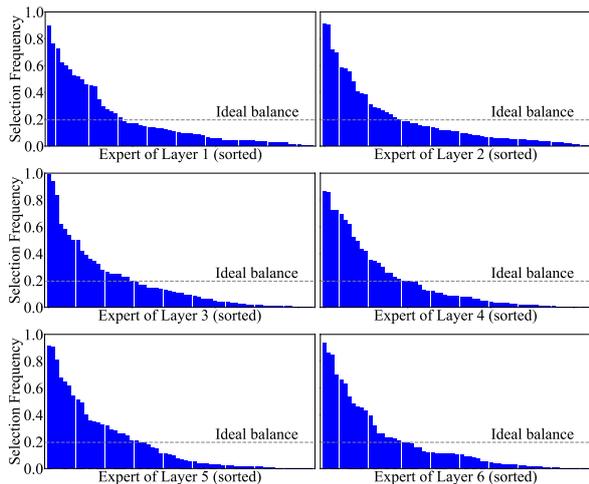


Figure 5: Selection Frequency of 64 experts in each encoder layer of MoEified T5-Small. The frequency of ideal balance selection is 0.2 while the distribution is much unbalanced.

will select 32 of them for a single input.

We report the MLM loss on the validation set in Table 4. From the table, we have two observations. (1) The loss of the standard pre-trained model is lower than that of the pre-trained MoE model. We guess that the optimization of MoE models is difficult than that of the standard models because of the restricted selection of MoE models. (2) MoEified models achieve better performance than the pre-trained MoE model. It indicates that pre-training a standard model then conducting MoEification can be a better option than pre-training an MoE model from scratch.

## 5 Analysis

In this section, we investigate the routing patterns of MoEified models and validate whether they are consistent with those of MoE models trained from scratch.

First, we count the selection frequency of each expert. Previous work introduces training objectives to ensure balance selection to make full use of model parameters (Lepikhin et al., 2021; Fedus et al., 2021). We report the results of the MoEified T5-Small with 20% experts on SST-2 in Figure 5. From the figure, we observe that the frequency distribution of expert selection is much unbalanced<sup>5</sup>. There are some commonly-used experts, whose frequencies are higher than 80%. Meanwhile, there

<sup>5</sup>Unbalanced selection will not influence the computation efficiency with current MoE implementations such as Fast-MoE (He et al., 2021a).

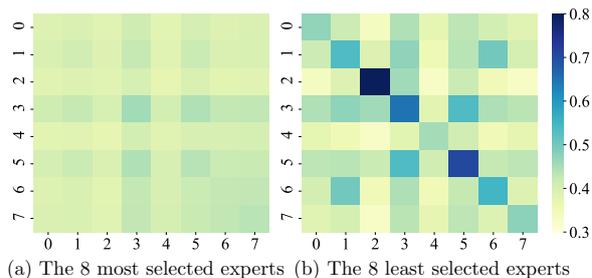


Figure 6: Input similarities between experts in the last encoder layer of MoEified T5-Small. For the most selected experts, both the self-similarities and inter-similarities are low. For the least selected experts, the self-similarities are much higher than inter-similarities.

are also some long-tail experts whose frequencies are lower than 10%.

Then, we calculate the self-similarities and inter-similarities of inputs between experts by sampling 10,000 inputs for each expert. We report the results of the last layer in Figure 6. For the most selected experts, which are selected by most inputs, the self-similarities are close to the inter-similarities. For the least selected experts, the self-similarities are much higher than the inter-similarities, which suggests that the inputs of each expert have obvious cluster structure.

From these results, we can conclude the routing patterns of MoEified models: there are some general experts, which can work for most inputs, and some input-specific experts, which are seldom used and may work in specific domains or tasks. This observation may inspire future work on training MoE models from scratch.

## 6 Conclusion

In this work, we propose MoEification, a new model acceleration technique, for large-scale Transformer models. MoEification utilizes the sparse activation phenomenon in FFNs of Transformer to convert a normal model to its MoE version with the same parameters. Experimental results show that large MoEified models can achieve comparable performance to the original models using only 10% to 20% computation cost of FFNs. By studying the routing patterns of MoEified models, we find that there are general and input-specific experts, which may inspire future work on training MoE models. In the future, we plan to extend MoEification to other Transformer models, such as BERT and GPT, and design better strategies for MoEification. We hope MoEification can benefit the real-world applications of large PLMs with better efficiency.

580  
581  
582  
583  
  
584  
585  
586  
  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
  
600  
601  
602  
603  
  
604  
605  
606  
607  
608  
  
609  
610  
611  
  
612  
613  
614  
615  
  
616  
617  
618  
  
619  
620  
621  
622  
  
623  
624  
625  
626  
  
627  
628  
629  
630  
631  
632  
633  
634

## References

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. [Longformer: The Long-Document transformer](#). *arXiv preprint 2004.05150*.

Yoshua Bengio. 2013. [Deep learning of representations: Looking forward](#). In *Proceedings of SLSP*, pages 1–37.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2021. [Language models are Few-Shot learners](#). In *Proceedings of NeurIPS*, pages 1877–1901.

Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. [The PASCAL recognising textual entailment challenge](#). In *Machine learning challenges.*, pages 177–190.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of NAACL*, pages 4171–4186.

William B Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of IWP*, pages 9–16.

William Fedus, Barret Zoph, and Noam Shazeer. 2021. [Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity](#). *arXiv preprint 2101.03961*.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer Feed-Forward layers are Key-Value memories](#). *arXiv preprint 2012.1491*.

Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. 2007. [The third PASCAL recognizing textual entailment challenge](#). In *Proceedings of TEP*, pages 1–9.

Ian Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. 2013. [Max-out networks](#). In *Proceedings of ICML*, pages 1319–1327.

Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao, Ao Zhang, Liang Zhang, Wentao Han, Minlie Huang, Qin Jin, Yanyan Lan, Yang Liu, Zhiyuan Liu, Zhiwu Lu, Xipeng Qiu, Ruihua Song, Jie Tang, Ji-Rong Wen, Jinhui Yuan, Wayne Xin Zhao, and Jun Zhu. 2021. [Pre-Trained models: Past, present and future](#). *arXiv preprint 2106.07139*.

Jiaao He, Jiezhong Qiu, Aohan Zeng, Zhilin Yang, Jidong Zhai, and Jie Tang. 2021a. [FastMoE: A fast Mixture-of-Expert training system](#). *arXiv preprint 2103.13262*.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021b. [DeBERTa: decoding-enhanced bert with disentangled attention](#). In *Proceedings of ICLR*.

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the knowledge in a neural network](#). *arXiv preprint arXiv:1503.02531*.

Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. 2012. [Improving neural networks by preventing co-adaptation of feature detectors](#). *arXiv preprint 1207.0580*.

Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Weinberger. 2018. [Multi-Scale dense networks for resource efficient image classification](#). In *Proceedings of ICLR*.

Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. 1991. [Adaptive mixtures of local experts](#). *Neural Comput.*, 3(1):79–87.

Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. [TinyBERT: Distilling BERT for natural language understanding](#). In *Findings of EMNLP*, pages 4163–4174.

George Karypis and Vipin Kumar. 1998. [A fast and high quality multilevel scheme for partitioning irregular graphs](#). *SIAM J. Sci. Comput.*, 20(1):359–392.

Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. [Reformer: The efficient transformer](#). In *Proceedings of ICLR*.

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. [ImageNet classification with deep convolutional neural networks](#). In *Proceedings of NeurIPS*, pages 1106–1114.

Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale Reading comprehension dataset from examinations](#). In *Proceedings of EMNLP*, pages 785–794.

Dmitry Lepikhin, Hyoukjoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2021. [GShard: Scaling giant models with conditional computation and automatic sharding](#). In *Proceedings of ICLR*.

Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. 2021. [BASE layers: Simplifying training of large, sparse models](#). *arXiv preprint 2103.16716*.

687	Bingbing Li, Zhenglun Kong, Tianyun Zhang, Ji Li,	Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie	739
688	Zhengang Li, Hang Liu, and Caiwen Ding. 2020. Ef-	Liu, Yiming Yang, and Denny Zhou. 2020. <b>Mo-</b>	740
689	efficient transformer-based large scale language repre-	mobileBERT: a compact Task-Agnostic BERT for	741
690	sentations using hardware-friendly block structured	Resource-Limited devices. In <i>Proceedings of ACL</i> ,	742
691	pruning. In <i>Findings of EMNLP</i> , pages 3187–3199.	pages 2158–2170.	743
692	Lei Li, Yankai Lin, Deli Chen, Shuhuai Ren, Peng Li,	Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald	744
693	Jie Zhou, and Xu Sun. 2021. <b>CascadeBERT: Accel-</b>	Metzler. 2020. <b>Efficient transformers: A survey.</b>	745
694	erating inference of pre-trained language models via	<i>arXiv preprint arXiv:2009.06732</i> .	746
695	calibrated complete models cascade. In <i>Findings of</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, and	747
696	<i>EMNLP</i> , pages 475–486.	Jakob Uszkoreit. 2017. <b>Attention is all you need.</b> In	748
697	Mikko I. Malinen and Pasi Fränti. 2014. <b>Balanced k-</b>	<i>Proceedings of NeurIPS</i> , pages 5998–6008.	749
698	means for clustering. In <i>Proceedings of SSSPR</i> , vol-	Elena Voita, David Talbot, Fedor Moiseev, Rico Sen-	750
699	ume 8621, pages 32–41.	nrich, and Ivan Titov. 2019. <b>Analyzing Multi-Head</b>	751
700	Stephen Merity, Caiming Xiong, James Bradbury, and	<b>Self-Attention: Specialized heads do the heavy lift-</b>	752
701	Richard Socher. 2017. <b>Pointer sentinel mixture mod-</b>	<b>ing, the rest can be pruned.</b> In <i>Proceedings of ACL</i> ,	753
702	els. In <i>Proceedings of ICLR</i> .	pages 5797–5808.	754
703	Paul Michel, Omer Levy, and Graham Neubig. 2019.	Alex Wang, Amanpreet Singh, Julian Michael, Felix	755
704	<b>Are sixteen heads really better than one?</b> In <i>Pro-</i>	Hill, Omer Levy, and Samuel R Bowman. 2019.	756
705	<i>ceedings of NeurIPS</i> , pages 14014–14024.	<b>GLUE: A multi-task benchmark and analysis plat-</b>	757
706	Vinod Nair and Geoffrey E. Hinton. 2010. <b>Rectified</b>	<b>form for natural language understanding.</b> In <i>Pro-</i>	758
707	<b>linear units improve restricted boltzmann machines.</b>	<i>ceedings of ICLR</i> .	759
708	In <i>Proceedings of ICML</i> , pages 807–814.	Alex Warstadt, Amanpreet Singh, and Samuel R. Bow-	760
709	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	man. 2019. <b>Neural network acceptability judgments.</b>	761
710	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	<i>TACL</i> , 7:625–641.	762
711	Wei Li, and Peter J Liu. 2020. <b>Exploring the limits</b>	Adina Williams, Nikita Nangia, and Samuel R. Bow-	763
712	<b>of transfer learning with a unified Text-to-Text trans-</b>	man. 2018. <b>A broad-coverage challenge corpus for</b>	764
713	<b>former.</b> <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.	<b>sentence understanding through inference.</b> In <i>Pro-</i>	765
714	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	<i>ceedings of NAACL-HLT</i> , pages 1112–1122.	766
715	Percy Liang. 2016. <b>SQuAD: 100,000+ questions for</b>	Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen,	767
716	<b>machine comprehension of text.</b> In <i>Proceedings of</i>	Yi Yang, and Shilei Wen. 2020. <b>Dynamic inference:</b>	768
717	<i>EMNLP</i> , pages 2383–2392.	<b>A new approach toward efficient video action recog-</b>	769
718	Stephen Roller, Sainbayar Sukhbaatar, Arthur Szlam,	<b>niton.</b> <i>arXiv preprint 2002.03342</i> .	770
719	and Jason Weston. 2021. <b>Hash layers for large</b>	Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and	771
720	<b>sparse models.</b> <i>arXiv preprint arXiv:2106.04426</i> .	Jimmy Lin. 2020. <b>DeeBERT: Dynamic early exiting</b>	772
721	Victor Sanh, Lysandre Debut, Julien Chaumond, and	<b>for accelerating BERT inference.</b> In <i>Proceedings of</i>	773
722	Thomas Wolf. 2019. <b>DistilBERT, a distilled version</b>	<i>ACL</i> , pages 2246–2251.	774
723	<b>of BERT: smaller, faster, cheaper and lighter.</b> <i>arXiv</i>	Dongkuan Xu, Ian En-Hsu Yen, Jinxi Zhao, and Zhibin	775
724	<i>preprint 1910.01108</i> .	Xiao. 2021. <b>Rethinking network pruning – under the</b>	776
725	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz,	<b>pre-train and fine-tune paradigm.</b> In <i>Proceedings of</i>	777
726	Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff	<i>NAACL-HLT</i> , pages 2376–2382.	778
727	Dean. 2017. <b>Outrageously large neural networks:</b>	Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe	779
728	<b>The Sparsely-Gated Mixture-of-Experts layer.</b> In	Wasserblat. 2019. <b>Q8BERT: Quantized 8bit BERT.</b>	780
729	<i>Proceedings of ICLR</i> .	<i>arXiv preprint 1910.06188</i> .	781
730	Richard Socher, Alex Perelygin, Jean Wu, Jason	Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao	782
731	Chuang, Christopher D Manning, Andrew Ng, and	Chen, Xin Jiang, and Qun Liu. 2020. <b>TernaryBERT:</b>	783
732	Christopher Potts. 2013. <b>Recursive deep models</b>	<b>Distillation-aware ultra-low bit BERT.</b> In <i>Proceed-</i>	784
733	<b>for semantic compositionality over a sentiment tree-</b>	<i>ings of EMNLP</i> , pages 509–521.	785
734	<b>bank.</b> In <i>Proceedings of EMNLP</i> , pages 1631–1642.	Zhengyan Zhang, Fanchao Qi, Zhiyuan Liu, Qun Liu,	786
735	Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. 2019.	and Maosong Sun. 2021. <b>Know what you don’t</b>	787
736	<b>Patient knowledge distillation for BERT model com-</b>	<b>need: Single-Shot Meta-Pruning for attention heads.</b>	788
737	<b>pression.</b> In <i>Proceedings of EMNLP</i> , pages 4323–	<i>AI Open</i> , 2:36–42.	789
738	4332.		

## A Activation Statistics before Fine-tuning

We count the activation statistics of PLMs before fine-tuning on the pre-training data containing about 50,000 input tokens. The results are shown in Figure 7. We observe that PLMs before fine-tuning also have the sparse activation phenomenon and fine-tuning brings little change.

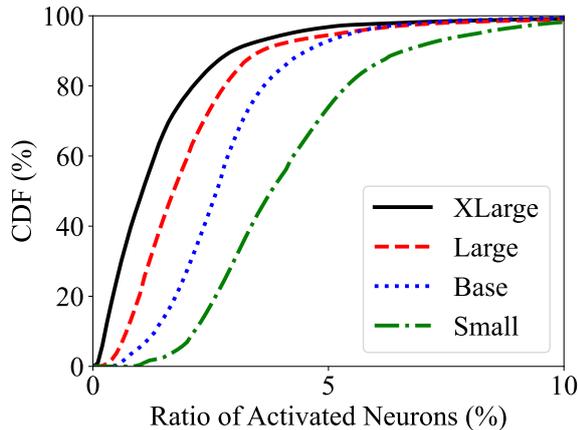


Figure 7: CDF of the ratios of activated neurons for each input with different models before fine-tuning.

Then, we compare the activations of pre-trained models and those of fine-tuned models. We use the average ratio of activated neurons as the index. The results are shown in Table 5. We observe that fine-tuning increases the average activation ratio for most models. The reason may be that different neurons start to learn the same task-specific patterns during fine-tuning. Interestingly, the increase on RACE is smaller than that on the other datasets. Since RACE is more difficult than the other datasets, there should be more task-specific patterns in RACE and less neurons learn the same patterns. Moreover, the pre-training task MLM requires more patterns than RACE so the ratios of MLM are lowest.

	Small	Base	Large	XLarge
MLM	4.18	2.85	2.17	1.52
SST-2	5.53	2.24	2.50	2.46
MNLI	5.59	3.25	2.44	2.45
RACE	4.94	3.08	1.98	1.79

Table 5: Average ratio of activated neurons for each input. MLM represents the pre-trained models with masked language modeling. SST-2, MNLI, RACE represent the fine-tuned models on each dataset.

## B Results of Graph Partition

Co-activation Graph Split achieves good performance in expert construction. Here, we study whether the co-activation graph is suitable for partitioning. We report the results of graph partition of T5-Large on SST-2 in Figure 8. Smaller ratios of edgecuts, which straddle partitions, mean that more co-activation pairs are included in experts. We only report the results of encoder layers because all ratios of decoder layers are smaller than 0.001. From this figure, we can see that the overall ratio is small and these graphs are suitable for partitioning.

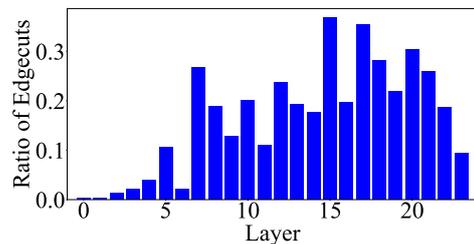


Figure 8: Ratio of edgecuts in different layers.

## C Accuracy of MLP Selection

MLP selection trains MLPs to fit the groundtruth selection. In this part, we report the accuracy of MLPs in T5-Large fine-tuned on SST-2. The results are shown in Figure 9 and 10. The overall accuracy of the encoder is about 0.8 and the overall accuracy of the decoder is about 0.7.

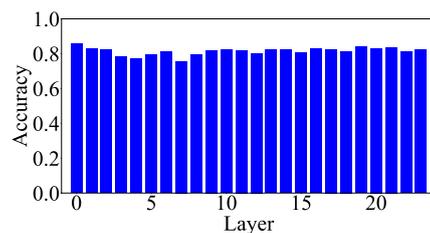


Figure 9: Accuracy of MLPs of encoder layers.

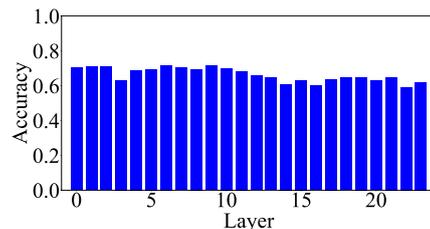


Figure 10: Accuracy of MLPs of decoder layers.

## D Relative Cost of Routing

In this work, we set the number of neurons in each expert to 32. Then, the number of experts in each layer  $k$  is  $\frac{d_{ff}}{32}$ . In most Transformer models,  $d_{ff} = 4d_{model}$ . The computation complexity of Similarity Selection for each input is

$$O(kd_{model}) = O\left(\frac{d_{model}^2}{8}\right). \quad (9)$$

The computation complexity of FFNs for each input is

$$O(d_{model} \cdot d_{ff}) = O(4d_{model}^2). \quad (10)$$

Then, the relative cost of routing to that of FFNs is constant for different models. It is also similar to MLP Selection.