

bert2BERT: Towards Reusable Pretrained Language Models

Cheng Chen^{1†}, Yichun Yin², Lifeng Shang², Xin Jiang², Yujia Qin¹,
Fengyu Wang¹, Zhi Wang^{3,4‡}, Xiao Chen², Zhiyuan Liu¹, Qun Liu²

¹Department of Computer Science and Technology, Tsinghua University

²Huawei Noah’s Ark Lab, ³Tsinghua Shenzhen International Graduate School

⁴Peng Cheng Laboratory

{c-chen19, qyj20, wangfy20}@mails.tsinghua.edu.cn

{yinyichun, shang.lifeng, jiang.xin, chen.xiao2, qun.liu}@huawei.com

wangzhi@sz.tsinghua.edu.cn, liuzy@tsinghua.edu.cn

Abstract

In recent years, researchers tend to pre-train ever-larger language models to explore the upper limit of deep models. However, large language model pre-training costs intensive computational resources, and most of the models are trained from scratch without reusing the existing pre-trained models, which is wasteful. In this paper, we propose bert2BERT¹, which can effectively transfer the knowledge of an existing smaller pre-trained model to a large model through parameter initialization and significantly improve the pre-training efficiency of the large model. Specifically, we extend the previous *function-preserving* (Chen et al., 2016) method proposed in computer vision on the Transformer-based language model, and further improve it by proposing a novel method, *advanced knowledge* for the large model’s initialization. In addition, a two-stage learning method is proposed to further accelerate the pre-training. We conduct extensive experiments on representative PLMs (e.g., BERT and GPT) and demonstrate that (1) our method can save a significant amount of training cost compared with baselines including learning from scratch, StackBERT (Gong et al., 2019) and MSLT (Yang et al., 2020); (2) our method is generic and applicable to different types of pre-trained models. In particular, bert2BERT saves about 45% and 47% computational cost of pre-training BERT_{BASE} and GPT_{BASE} by reusing the models of almost their half sizes.

1 Introduction

Pre-trained language models (PLMs), such as BERT (Devlin et al., 2019), GPT (Radford et al., 2018, 2019; Brown et al., 2020), ELECTRA (Clark et al., 2020), XLNet (Yang et al., 2019) and RoBERTa (Liu et al., 2019), have achieved great

[†] This work is done when Cheng Chen is an intern at Huawei Noah’s Ark Lab.

[‡] Corresponding author.

¹Our code is available at <https://github.com/huawei-noah/Pretrained-Language-Model>.

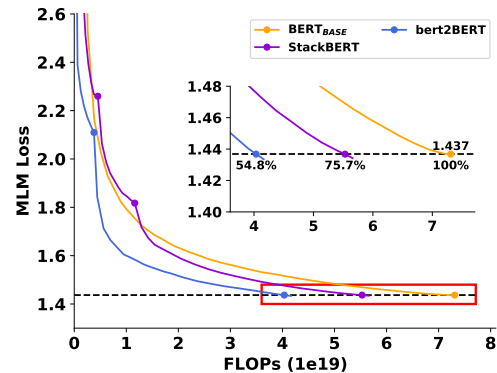


Figure 1: Loss curves of bert2BERT and baselines. StackBERT (Gong et al., 2019) is based on the progressive training setting. More details are shown in Table 2.

success in natural language processing (NLP). However, the pre-training process of large PLMs can be extremely computationally expensive and produces huge carbon footprints. For example, GPT-3 uses 3.1E+6 GPU hours for training, at an estimated cost of \$4.6 million², consuming a lot of computing resources. Therefore, how to reduce the training cost of PLM is of great importance to Green AI (Schwartz et al., 2020).

Recently, there is a trend of training extremely large models to explore the upper limits of PLMs. For example, large pre-trained models, including GPT-3 (Brown et al., 2020) (175B), PanGu- α (Zeng et al., 2021) (200B) and Switch Transformers (Fedus et al., 2021) (1571B), have been proved promising in language understanding and generation. However, these models are all pre-trained from scratch independently without utilizing the knowledge of smaller ones that have already been trained. On the other hand, our empirical studies show that the pre-trained models of different scales could share similar knowledge, for example in Figure 2, the attention patterns of the two PLMs with different sizes are similar.

To save the training cost of large models, we

²<https://lambdalabs.com/blog/demystifying-gpt-3/>

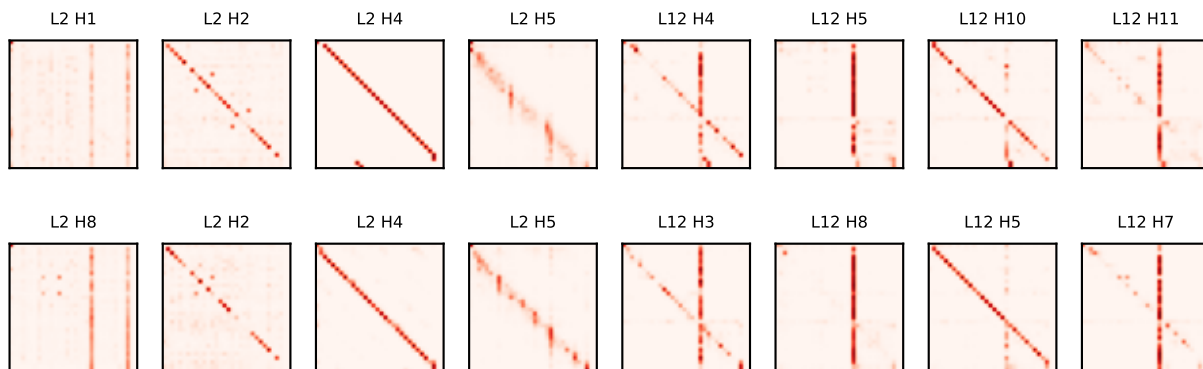


Figure 2: The comparisons of attention patterns between small and large PLMs. The upper ones are the attention patterns of BERT_{BASE} model whose architecture is $\{L=12, D=768\}$, and the lower ones are the attention patterns of one small BERT model whose architecture is $\{L=12, D=512\}$. We find that there are a large number of similar attention patterns in the same layer of the two models, indicating the possibility of reusing parameters of trained small PLMs to speed up the pre-training of large PLMs. The attention maps of PLMs with different layers are also similar, which is visualized in previous work (Gong et al., 2019; Yang et al., 2020).

propose the bert2BERT method, which can efficiently transfer the learned knowledge of the smaller model to the large model. bert2BERT consists of two components: (1) For parameter initialization, we first extend the function preserving training (Chen et al., 2016) to PLMs by duplicating and stacking the parameters of the existing smaller PLM, which we call function-preserving initialization (FPI). FPI ensures that the initialized large model has almost the same behavior as the small model, so that the large model has a good starting point for later optimization. We also find that duplicating the weights of the upper layer to the current layer can further accelerate the convergence of the large model, which we call advanced knowledge initialization (AKI). Although the AKI somewhat violates the principle of function preserving, we find that empirically it also has a good starting point as shown in Table 1, which leads to a faster convergence rate and achieves higher training efficiency. (2) Secondly, a two-stage training strategy is further applied to the large model to accelerate the training process.

To demonstrate the superiority of our method, we conduct extensive experiments on two representative PLMs: BERT and GPT, with different source model sizes. The results show that: (1) our method can save a significant amount of computation in pre-training compared to the traditional way of learning from scratch and progressive stacking methods such as StackBERT (Gong et al., 2019) and MSLT (Yang et al., 2020); (2) our method is model-agnostic, which can be applied on a wide range of Transformer-based PLMs. One typical example is that, when using a small pre-trained

model with half the size of BERT_{BASE} for initialization, bert2BERT saves **45%** computation cost of the original BERT_{BASE} pre-training.

In general, our contributions are summarized as follows: (1) We explore a new direction for the efficient pre-training by reusing the trained parameters of small models to initialize the large model; (2) We successfully extend function preserving method (Chen et al., 2016) on BERT and further propose advanced knowledge initialization, which can effectively transfer the knowledge of the trained small model to the big model and improve the pre-training efficiency; (3) The proposed method outperforms other training methods and achieves 45% computation reduction on BERT_{BASE}; (4) Our method is generic, effective for both the BERT and GPT models, and have great potential to become an energy-efficient solution for pre-training super large-scale language models.

2 Related Work

Efficient Pre-training in NLP. The efficiency of pre-training has been explored by previous work. Some works (Gong et al., 2019; Yang et al., 2020; Gu et al., 2021) propose progressive learning to accelerate the pre-training, which are motivated by the fact that different layers have some similar knowledge (e.g., attention patterns). They start pre-training a small model with fewer Transformer layers, and then iteratively expand the model by stacking the already trained layers on the top. Another line of work proposes to “back distill” the knowledge of the small models into large models, which is termed as knowledge inheritance (Qin et al., 2021). Some works focus on the data effi-

ciency (Wu et al., 2021) and take notes for rare words during the pre-training process to help the model understand them when they occur next. ELECTRA (Clark et al., 2020) proposes a task of replaced token detection to predict whether each token in the input was replaced or not, which improves the pre-training efficiency. Our method is orthogonal to this kind of work and the combination of ELECTRA and bert2BERT could achieve better efficiency. In addition, there are several other orthogonal techniques for efficient pre-training: mixed-precision training (Shoeybi et al., 2019), large batch optimization (You et al., 2020), model architecture innovation (Lan et al., 2020), layer dropping technique (Zhang and He, 2020), etc.

Reusable Neural Network. Reusable neural network, a topic related to transfer learning (Pan and Yang, 2010), is introduced to accelerate the model training in computer vision. One classical work is Net2Net (Chen et al., 2016), which first proposes the concept of the function-preserving transformation to make neural networks reusable. However, Net2Net randomly selects the neurons to be split. To handle this problem, some works (Wu et al., 2019, 2020b; Wang et al., 2019b; Wu et al., 2020a) leverage a functional steepest descent idea to decide the optimal subset of neurons to be split. The pruning technique (Han et al., 2015) is also introduced for reusable neural networks (Feng and Panda, 2020). In this paper, we study the reusable pre-trained language model and propose a new method, bert2BERT, to accelerate the pre-training of BERT and GPT.

3 Preliminary

BERT consists of one embedding layer and multiple Transformer (Vaswani et al., 2017) layers.

3.1 Embedding Layer

The embedding layer first maps the tokens in a sentence into vectors with an embedding matrix \mathbf{W}^E . Then one normalization layer is employed to produce the initial hidden states \mathbf{H}_0 .

3.2 Transformer Layer

The hidden states are iteratively processed by multiple Transformer layers as follows:

$$\mathbf{H}_l = \text{Transformer}_l(\mathbf{H}_{l-1}), l \in [1, L] \quad (1)$$

where L denotes the number of Transformer layers, each including a multi-head attention (MHA) and a feed-forward network (FFN).

MHA. It is composed of multiple parallel self-attention heads. The hidden states of the previous layer are fed into each head and then the outputs of all heads are summed to obtain the final output as follows:

$$\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i = \mathbf{H}_{l-1} \mathbf{W}_{l,i}^Q, \mathbf{H}_{l-1} \mathbf{W}_{l,i}^K, \mathbf{H}_{l-1} \mathbf{W}_{l,i}^V,$$

$$\mathbf{H}_{l,i}^{\text{HEAD}} = \text{softmax}\left(\frac{\mathbf{Q}_i \mathbf{K}_i^T}{\sqrt{d_k}}\right) \mathbf{V}_i \mathbf{W}_{l,i}^O,$$

$$\text{MHA}(\mathbf{H}_{l-1}) = \sum_{i=1}^a \mathbf{H}_{l,i}^{\text{HEAD}},$$

$$\mathbf{H}_l^{\text{MHA}} = \text{LayerNorm}(\mathbf{H}_{l-1} + \text{MHA}(\mathbf{H}_{l-1})). \quad (2)$$

\mathbf{H}_{l-1} is linearly projected to queries (\mathbf{Q}_i), keys (\mathbf{K}_i) and values (\mathbf{V}_i) using $\mathbf{W}_{l,i}^Q$, $\mathbf{W}_{l,i}^K$, $\mathbf{W}_{l,i}^V$ respectively. $\mathbf{H}_{l,i}^{\text{HEAD}}$ indicates the context-aware vector which is obtained by the scaled dot-product of queries and keys in the i -th attention head. a represents the number of self-attention heads. d_k is the head dimension acting as the scaling factor.

FFN. It consists of two linear layers and one GeLU activation function (Hendrycks and Gimpel, 2016), that is:

$$\mathbf{H}_l^{\text{FFN}} = \text{GeLU}(\mathbf{H}_l^{\text{MHA}} \mathbf{W}_l^1 + \mathbf{b}_l^1) \mathbf{W}_l^2 + \mathbf{b}_l^2,$$

$$\mathbf{H}_l = \text{LayerNorm}(\mathbf{H}_l^{\text{MHA}} + \mathbf{H}_l^{\text{FFN}}). \quad (3)$$

Layer Normalization. Both the modules of MHA and FFN have one layer normalization (Ba et al., 2016) that stabilizes the dynamics of the hidden state in the Transformer. Formally, it is written as:

$$\text{LayerNorm}(\mathbf{H}) = \left(\frac{\mathbf{H} - \mu_{\mathbf{H}}}{\sigma_{\mathbf{H}}}\right) \odot \mathbf{W}^{\text{LN}} + \mathbf{b}^{\text{LN}}, \quad (4)$$

where \odot means the element-wise multiplication. The statistics of $\mu_{\mathbf{H}}$ and $\sigma_{\mathbf{H}}$ are the mean and variance of hidden states \mathbf{H} respectively.

4 Methodology

4.1 Problem Statement

We aim to accelerate the pre-training of target model $\mathcal{T}(L^t, D^t)$ by transferring the knowledge of an existing pre-trained source model $\mathcal{S}(L^s, D^s)$, where $L^{s|t}$ means the numbers of Transformer layer and $D^{s|t}$ means the model width (i.e., hidden size), satisfying $L^s \leq L^t$ and $D^s \leq D^t$. Formally, our problem is two-fold: (1) how to perform an effective parameter initialization for \mathcal{T} by reusing the trained parameters of \mathcal{S} , and (2) how to efficiently

train the initialized \mathcal{T} , so that \mathcal{T} can have a faster convergence rate in pre-training.

4.2 Overview

Targeting the above problems, bert2BERT first initializes the target model \mathcal{T} with the parameters of the existing model \mathcal{S} by the width-wise expansion ($D^s \rightarrow D^t$) and depth-wise expansion ($L^s \rightarrow L^t$). Through this expansion, the knowledge contained in the parameters of the source model is directly transferred to the target model. Then we further pre-train the initialized target model with a two-stage pre-training method. The overall workflow is illustrated in Section 4.5.

Essentially, the width-wise expansion can be decomposed into expansions of parameter matrices (or vectors³). As illustrated in Figure 3, the matrix expansion enlarges $\mathbf{W} \in \mathbb{R}^{d_{in}^w \times d_{out}^w}$ of \mathcal{S} to $\mathbf{U} \in \mathbb{R}^{d_{in}^u \times d_{out}^u}$ of \mathcal{T} by two kinds of operations: in-dimension and out-dimension expansion.

In the following sections, we first introduce two strategies of width-wise expansion: function-preserving and advanced knowledge initialization. Then, we introduce the depth-wise expansion and detail the two-stage pre-training process.

4.3 Width-wise Expansion

For the paper clarity, we introduce two index mapping functions: g_{in} and g_{out} , where $g_{in}(i)$ means the i -th in-dimension of \mathbf{U} reuses the $g_{in}(i)$ -th in-dimension parameters of \mathbf{W} , $g_{out}(j)$ means the j -th out-dimension of \mathbf{U} reuses the $g_{out}(j)$ -th out-dimension parameters of \mathbf{W} . Both our two methods are defined with these two mapping functions. $\mathbf{W}_{(i,j)}$ means the parameter element, i and j refer to the i -th in-dimension index and j -th out-dimension index respectively. As shown in Figure 3, the i -th in-dimension parameters of \mathbf{W} are the parameters of the i -th input neuron of \mathbf{W} or the i -th column of \mathbf{W} .

4.3.1 Function Preserving Initialization

Function preserving initialization (FPI) (Chen et al., 2016) aims to make the initialized target model have the same function as the source model, which means that given the same input, the initialized target model has the same output as the source model. In this paper, we extend FPI on a different architecture, Transformer-based pre-trained language model. We give an example in Figure 3 to illustrate

³We omit the expansion of bias (vector) for simplicity. It follows a similar process as the matrix expansion.

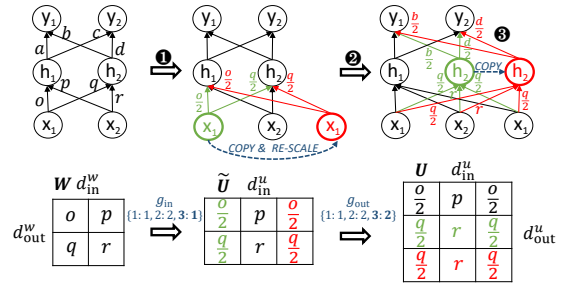


Figure 3: Overview of the function preserving initialization (FPI). Given the same input $\{x_1, x_2\}$, FPI ensures the initialized target model has the same output $\{y_1, y_2\}$ with the source model. The first and the second steps are expanding the in-dimension and out-dimension of the parameter matrix according to mapping functions g_{in} and g_{out} respectively. After we expand the matrix \mathbf{W} into \mathbf{U} , we use the in-dimension expansion on the upper parameter matrix again to ensure the output $\{y_1, y_2\}$ same as the original one. From the view of neurons, FPI copies the corresponding input and output neurons to expand the neural network.

FPI. Formally, the mapping functions are defined as follows:

$$g_{in}(i) = \begin{cases} i & i \in [1, d_{in}^w] \\ f(\{1, 2, \dots, d_{in}^w\}) & i \in (d_{in}^w, d_{in}^u], \end{cases} \quad (5)$$

$$g_{out}(j) = \begin{cases} j & j \in [1, d_{out}^w] \\ f(\{1, 2, \dots, d_{out}^w\}) & j \in (d_{out}^w, d_{out}^u], \end{cases} \quad (6)$$

where $f(\cdot)$ is uniform sampling. We denote the weight expansion as $\mathbf{U} = \text{EXPN}(\mathbf{W}; g_{in}, g_{out})$, which includes in-dimension expansion (Eq. 7) and out-dimension expansion (Eq. 8):

$$C_{g_{in}(i)} = \sum_{i'=1}^{d_{in}^u} \mathbb{I}(g_{in}(i') = g_{in}(i)) \quad (7)$$

$$\begin{aligned} \tilde{\mathbf{U}}_{(i,*)} &= \frac{1}{C_{g_{in}(i)}} \mathbf{W}_{(g_{in}(i),*)}, \\ \mathbf{U}_{(*,j)} &= \tilde{\mathbf{U}}_{(*,g_{out}(j))}, \end{aligned} \quad (8)$$

where $\mathbb{I}(\cdot)$ is an indicator function, and $C_{g_{in}(i)}$ is the count of $g_{in}(i)$ in the values of $g_{in}(\cdot)$, which is used to re-scale the original parameters to keep the function preserving property.

Expansion for All Modules. We apply FPI for all modules of BERT via matrix expansion $\text{EXPN}(\cdot)$. Specifically, for the embedding matrix \mathbf{W}^E , we only conduct the out-dimension expansion:

$$\mathbf{U}_{(*,j)}^E = \mathbf{W}_{(*,g_{out}^E(j))}^E. \quad (9)$$

MHA module can be decomposed into multiple parallel self-attention heads and we conduct the head-wise expansion for this module, which means

increasing the number of attention heads. The head-wise expansion is formulated as:

$$\mathbf{U}^{Q|K|V|O} = \text{EXPN}(\mathbf{W}^{Q|K|V|O}; g_{\text{in}}^{q|k|v|o}, g_{\text{out}}^{q|k|v|o}). \quad (10)$$

Specifically, the head-wise expansion means that we reuse the head group parameters to construct the new matrices. The i -th head group in l -th layer contains $\mathbf{W}_{l,i}^Q|\mathbf{W}_{l,i}^K|\mathbf{W}_{l,i}^V|\mathbf{W}_{l,i}^O$ in Eq. 2 and the out-dimension expansion for $\mathbf{W}_{l,i}^Q|\mathbf{W}_{l,i}^K|\mathbf{W}_{l,i}^V$ is:

$$g_{\text{out}}^{q|k|v}(j) = \begin{cases} j & j \in [1, a^s] \\ f(\{1, 2, \dots, a^s\}) & j \in (a^s, a^t], \end{cases} \quad (11)$$

where j is the head index and $a^{s|t}$ mean the head numbers of source model and target model respectively. The module has three constraints: $\{g_{\text{out}}^e = g_{\text{in}}^{q|k|v}; g_{\text{out}}^{q|k|v} = g_{\text{in}}^o; g_{\text{in}}^{q|k|v} = g_{\text{out}}^o\}$, with the first two constraints for hidden dimension consistency (Wen et al., 2018; Chen et al., 2021) and the third one for residual connection (Eq. 2).

For the FFN module, we perform the expansion on the parameter matrices $\mathbf{W}^{1|2}$ (Eq. 3) as follows:

$$\mathbf{U}^{1|2} = \text{EXPN}(\mathbf{W}^{1|2}; g_{\text{in}}^{1|2}, g_{\text{out}}^{1|2}). \quad (12)$$

Similar to the MHA module, the mapping functions of FFN also have three constraints: $\{g_{\text{out}}^o = g_{\text{in}}^1; g_{\text{out}}^1 = g_{\text{in}}^2; g_{\text{in}}^1 = g_{\text{out}}^2\}$.

For the layer normalization, we take the layer normalization of FFN as an example, its expansion is formulated as:

$$\mathbf{U}_j^{LN} = \mathbf{W}_{g_{\text{out}}^2(j)}^{LN}. \quad (13)$$

Note that in layer normalization (Eq. 4), the mean μ and variance σ are calculated based on the hidden representations \mathbf{H} . Thus, the expansion of this parameter inevitably induces a gap and prevents the target model from strictly following the function preserving principle. However, we empirically find that the gap is so small that it can hardly affect the initialization and convergence of the target model. Thus we ignore this discrepancy.

We have validated the effectiveness of the adapted FPI in different settings in Table 1. The results show that the initialized model \mathcal{T} achieves almost the same loss as \mathcal{S} , demonstrating that FPI successfully retains the knowledge of the small model when performing parameter expansion.

4.3.2 Advanced Knowledge Initialization

To further improve the convergence rate of the pre-training target model, we propose the advanced knowledge initialization (AKI), which expands new

Method	$\mathcal{S}(12, 384)$	$\mathcal{S}(12, 512)$
Original	1.89	1.67
Rand	10.40	10.42
DirectCopy	9.05	6.45
FPI	1.89	1.70
AKI	2.08	1.96

Table 1: The comparison of MLM losses between FPI and baselines. ‘‘Original’’ refers to the MLM losses of source pre-trained models \mathcal{S} . ‘‘Rand’’ refers to the MLM losses of randomly initialized target models. ‘‘Direct-Copy’’ refers to a naive method that directly copies the source model to the target model and the unfilled part is randomly initialized, ‘‘FPI’’ represents the function preserving method. We expand both models to the target model $\mathcal{T}(12, 768)$ and find that FPI can make the target model have similar losses with these trained source models. The loss gap between FPI and Original is brought by layer normalization. ‘‘AKI’’ represents the advanced knowledge initialization method.

matrices based on not only the parameters of the same layer but also the parameters of the upper layer in the source model. The intuition is based on previous findings (Jawahar et al., 2019; Clark et al., 2019) that adjacent Transformer layers have similar functionality, which ensures that it will not damage the knowledge contained in the parameters of the current layer. Moreover, the knowledge that comes from adjacent layers can break the symmetry (Chen et al., 2016) appeared in FPI, which has been demonstrated beneficial. We give an illustrative example in Figure 4 and formulate AKI as:

$$\mathbf{U}^l = \text{EXPN}(\mathbf{W}^l, \mathbf{W}^{l+1}; g_{\text{in}}^{l|l+1}, g_{\text{out}}^l). \quad (14)$$

Specifically, we first do the in-dimension expansion for $\mathbf{W}^{l|l+1}$. Here we take \mathbf{W}^l as an example:

$$C_{g_{\text{in}}^l(i)} = \sum_{i'=1}^{d_{\text{in}}^u} \mathbb{I}(g_{\text{in}}^l(i') = g_{\text{in}}^l(i)) \quad (15)$$

$$\tilde{\mathbf{U}}_{(i,*)}^l = \frac{1}{C_{g_{\text{in}}^l(i)}} \mathbf{W}_{(g_{\text{in}}^l(i),*)}^l.$$

It is similar with Eq. 7. Then we stack the expanded matrices of $\tilde{\mathbf{U}}^l$ and $\tilde{\mathbf{U}}^{l+1}$ to construct the final matrix:

$$\mathbf{U}_{(*,j)}^l = \begin{cases} \tilde{\mathbf{U}}_{(*,j)}^l & j \in [1, d_{\text{out}}^w] \\ \tilde{\mathbf{U}}_{(*,g_{\text{out}}^l(j)}^{l+1} & j \in (d_{\text{out}}^w, d_{\text{out}}^u]. \end{cases} \quad (16)$$

We directly copy the expanded $\tilde{\mathbf{U}}^l$ as the top part of the new matrix and place the sampled parameters from $\tilde{\mathbf{U}}^{l+1}$ on the bottom of the new matrix.

We aggregate upper-layer information into a new matrix for two intuitions: (1) it breaks the FPI symmetry that hinders model convergence (Chen et al.,

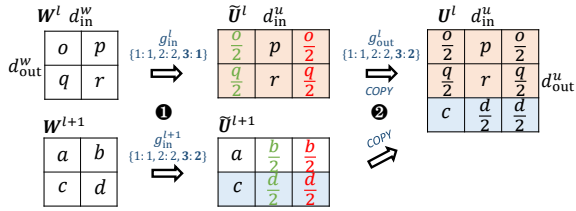


Figure 4: Overview of AKI. It first performs the in-dimension expansion on both the matrixes of current and upper layers. Then it uses the widened matrix of the current layer as the top part of the new matrix and samples the row of the widened matrix of the upper layer as the bottom part of the new matrix.

2016). For example, FPI makes the attention patterns in the same layer repeated, which is redundant and called symmetry; (2) upper-layer information can be used as similar but high-level knowledge to guide the model to converge faster. We display the attention patterns of the target model initialized by AKI in Appendix E and find that the target model can maintain the attention patterns of both current and upper layers very well.

Expansion for All Modules. For embedding matrix, we only do the out-dimension expansion as Eq. 9 in the FPI. Both the modules of MHA and FFN do the matrix expansion by following the defined operation in Eq. 15 and Eq. 16. The constraints of mapping functions follow the setting of FPI.

Empirically, we find that the AKI method outperforms FPI, while the performance is worse if we build a new matrix based on the matrix of the lower layer (or low-level knowledge). How to construct the optimal initialization for the target model with the parameters of different layers remains an open question and we leave it as future work.

For more details, we give a clear illustration of the FPI and AKI process in Appendix F.

4.4 Depth-wise Expansion

After the width-wise expansion, we obtain a widened model with the same width as the target model. To bridge the depth gap, we perform depth-wise expansion to increase model depth to the depth of the target model. We illustrate this process in Algorithm 1 and the main idea is to iteratively stack the widened model until its depth is equal to the target model (Gong et al., 2019).

4.5 Two-stage Pre-training

To further improve the pre-training efficiency of initialized target model, we propose a two-stage training method: (1) train sub-models with different

Algorithm 1 Target Model Initialization

Input: the target model $\mathcal{T}(L^t, D^t)$ and the source model $\mathcal{S}(L^s, D^s)$.

- 1: $\mathcal{T}_1(L^s, D^t) \leftarrow$ do AKI or FPI with $\mathcal{S}(L^s, D^s)$
- 2: $k \leftarrow \lfloor L^t/L^s \rfloor$
- 3: **for** $t = 2 \rightarrow k$ **do**
- 4: $\mathcal{T}_t(L^s \cdot t, D^t) \leftarrow$ stack \mathcal{T}_1 on top of \mathcal{T}_{t-1}
- 5: **end for**
- 6: $\mathcal{T} \leftarrow$ stack top $L^t - L^s \cdot k$ layers of \mathcal{T}_1 .

Output: the initialized model $\mathcal{T}(L^t, D^t)$

Algorithm 2 Two-stage Pre-training

Input: the initialized model \mathcal{T} , large-scale unsupervised dataset \mathcal{D} , the epoch number of sub-model training E_b and the epoch number of whole training process E , the layer number l_b .

- 1: Construct sub-models and these models have the layer numbers of $\{l_b, 2 \cdot l_b, \dots, L^t\}$.
- 2: **for** $e = 1 \rightarrow E_b$ **do**
- 3: **for** $batch$ in D **do**
- 4: $\mathcal{T}' \leftarrow$ sample one sub-model.
- 5: Perform forward and backward of \mathcal{T}' .
- 6: Update only top l_b layers of \mathcal{T}' .
- 7: **end for**
- 8: **end for**
- 9: **for** $e = E_b \rightarrow E$ **do**
- 10: **for** $batch$ in D **do**
- 11: Perform forward and backward of \mathcal{T} .
- 12: Update whole model \mathcal{T} .
- 13: **end for**
- 14: **end for**

Output: the pre-trained model \mathcal{T}

layers in a random manner to make the complete model converge at a low cost. These sub-models are built with bottom Transformer layers of the initialized target model and share one classification layer. At each optimization step, we randomly sample one sub-model and only update its top Transformer layers and the shared classification layer. (2) After the sub-structure training, we further perform the traditional full-model training. The details of our method are displayed in Algorithm 2.

5 Experiment

5.1 Experimental Setup

Pre-training Details. We use the English Wikipedia and Toronto Book Corpus (Zhu et al., 2015) as the pre-training data. The settings of pre-training are: peak learning rate of 1e-4, warmup

Model	FLOPs ($\times 1e19$)	Ratio (Saving)	Loss (MLM)	SQuADv1.1 (F1)	SST-2 (Acc)	MNLI (Acc)	MRPC (Acc)	CoLA (Mcc)	QNLI (Acc)	QQP (Acc)	STS-B (Acc)	Avg.
BERT _{BASE} (Google)	-	-	-	88.4(0.1)	93.6(0.2)	84.7(0.1)	87.9(0.9)	59.6(1.5)	91.6(0.1)	91.4(0.1)	89.6(0.5)	85.8(0.1)
BERT _{BASE} † (Ours)	7.3	0%	1.437	89.6(0.1)	92.7(0.2)	84.6(0.2)	88.6(0.5)	57.3(4.0)	90.6(0.7)	90.6(0.1)	89.9(0.3)	85.5(0.5)
<i>Progressive Training</i>												
MSLT†	6.5	10.7%	1.436	90.4(0.2)	92.9(0.2)	85.1(0.2)	87.9(2.1)	55.6(4.1)	90.7(0.2)	90.6(0.2)	88.2(0.6)	85.2(0.7)
StackBERT†	5.5	24.3%	1.433	90.4(0.2)	92.6(0.4)	85.3(0.1)	88.2(1.0)	63.2(0.9)	91.0(0.4)	91.0(0.1)	86.7(0.7)	86.0(0.2)
<i>bert2BERT : $\mathcal{S}(12, 512) \rightarrow \mathcal{T}(12, 768)$</i>												
DirectCopy	6.4	12.2%	1.436	89.8(0.2)	92.9(0.3)	84.7(0.2)	86.2(0.6)	62.2(0.7)	90.2(0.6)	90.4(0.1)	89.2(0.1)	85.7(0.1)
FPI	5.1	30.4%	1.436	90.0(0.2)	92.6(0.4)	85.2(0.1)	87.1(0.5)	61.5(0.9)	90.9(0.6)	90.8(0.2)	89.7(0.2)	86.0(0.1)
AKI	4.5	38.4%	1.434	90.4(0.1)	92.5(0.4)	85.3(0.4)	87.8(0.9)	61.0(1.4)	91.2(0.2)	90.5(0.1)	89.5(0.2)	86.0(0.2)
bert2BERT	4.0	45.2%	1.433	90.0(0.2)	92.9(0.1)	85.1(0.1)	87.7(0.7)	60.0(1.2)	90.5(0.8)	90.4(0.1)	89.2(0.2)	85.7(0.4)

Table 2: Comparison between bert2BERT and baselines. We report mean (and standard deviation) performance over 3 runs on the dev set. bert2BERT means the combination of AKI and two-stage pre-training here. FPI and AKI mean that the function preserving initialization, advanced knowledge initialization respectively. † means the re-implemented results, where the BERT_{BASE} and StackBERT achieve similar results with the original paper, and the MSLT result is different from the original paper may be due to the different training settings (e.g., in the original paper, it uses the LAMB optimizer (You et al., 2020) and only trains the corpus with a max sequence length of 128).

steps of 10k, training epochs of $E=40$, batch size of 512, sub-model training epochs of $E_b=5$, layer number of $l_b=3$. Unless otherwise noted, all methods including bert2BERT and baselines use the same pre-training settings for fair comparisons. In the settings of bert2BERT, the target model has a BERT_{BASE} architecture of $\mathcal{T}(12, 768)$ and the source model has an architecture of $\mathcal{S}(12, 512)$.

Fine-tuning Details. For the evaluation, we use tasks from GLUE benchmark (Wang et al., 2019a) and SQuADv1.1 (Rajpurkar et al., 2016). We report F1 for SQuADv1.1, Matthews correlation coefficient (Mcc) for CoLA (Warstadt et al., 2019) and accuracy (Acc) for other tasks. For the GLUE tasks fine-tuning, we set the batch size to 32, choose the learning rate from $\{5e-6, 1e-5, 2e-5, 3e-5\}$ and epochs from $\{4, 5, 10\}$. For the SQuADv1.1 fine-tuning, we set the batch size to 16, the learning rate to $3e-5$, and the number of training epochs to 4. All results are the average of 3 runs on the dev set.

Baselines. We first introduce a naive bert2BERT baseline named DirectCopy, which directly copies the small model to the target model and randomly initializes the unfilled parameters. StackBERT (Gong et al., 2019) and MSLT (Yang et al., 2020) are also included as the baselines. Both of them are trained in a progressive manner. Following the original setting, for the StackBERT, we first train the 3-layer BERT for 5 epochs, stack it twice into a 6-layer BERT and then train it for 7 epochs. In the final step, we stack the 6-layer model into BERT_{BASE} and further train it with 28 epochs. For MSLT, we first perform 4-stage training. In each stage, we add the top 3 layers of the model already trained to the top of the model and then pre-train the new model by partially updating

the top 3 layers. Each stage of the partial training process has 8 epochs. Finally, we further perform 20 full-model training epochs⁴ to achieve the same loss as BERT_{BASE} trained from scratch. The baselines are trained using the same optimizer, training steps, and warmup steps as the bert2BERT.

5.2 Results and Analysis

We demonstrate the effectiveness of the proposed method on the SQuAD and GLUE benchmark. The results are shown in Table 2. We also represent the loss curves in Figure 1 and Appendix A. The results show that: (1) DirectCopy only saves 12.2% computational costs, which indicates this naive method of directly copying the trained parameters of the source model to the target model is not effective; (2) our proposed methods, FPI and AKI, achieve better performances than the baselines. Although AKI does not follow the function preserving, it has a bigger loss than FPI at the start of training, AKI achieves a faster convergence rate by using the advanced knowledge and breaking the symmetry; (3) by performing the two-stage pre-training on the target model initialized by AKI, we can save 45.2% computational costs. Note that the total parameters of the source model are half of those of the target model (54M vs. 110M). The loss of bert2BERT in Figure 1 is high at the stage of sub-model training because it represents the average loss of all sub-models. We also compare the attention patterns of the target models initialized by DirectCopy, FPI, and AKI. The attention patterns and their discussions are displayed in Appendix E.

⁴We have tried the same setting as the original paper with 8 epoch full-model running but it does not achieve the same loss with BERT_{BASE} (1.511 vs. 1.437).

bert2BERT with Smaller Source Model. We also evaluate bert2BERT on different settings, where the source model $\mathcal{S}(6, 512)$, $\mathcal{S}(8, 512)$, $\mathcal{S}(10, 512)$ are significantly smaller than the target model (35M | 42M | 48M vs. 110M). The results are shown in Table 3 and loss curves are displayed in Appendix B. We observe that DirectCopy for $\mathcal{S}(6, 512)$ achieves no efficiency improvement over the original pre-training, which indicates that the significant size gap between the source and target model greatly reduces the benefit of DirectCopy methods. Compared with DirectCopy, our proposed method reduces the computation cost by 23.3%, which again demonstrates the effectiveness of bert2BERT. The results show that the smaller the size gap between the source model and target model, the greater the cost savings of bert2BERT. We also note that it is more challenging to speed up the target model with a small source model $\mathcal{S}(6, 512)$. We encourage future work to explore to transfer the knowledge from smaller source models to improve the pre-training efficiency of the target model.

Settings	Model	FLOPs ($\times 1e19$)	Ratio (Saving)	Loss (MLM)	Avg.
$\mathcal{S}(6, 512)$	DirectCopy	7.3	0%	1.440	89.1
	bert2BERT	5.6	23.3%	1.435	89.3
$\mathcal{S}(8, 512)$	bert2BERT	4.6	36.8%	1.435	89.2
$\mathcal{S}(10, 512)$	bert2BERT	4.2	42.7%	1.434	89.1

Table 3: bert2BERT with smaller source model. Avg means the average score of SST-2/MNLI/SQuADv1.1.

Effect of Sub-model Training Epochs. Our training procedure includes two stages: sub-model training and full-model training. Here, we study the effect of the number of sub-model training epochs by performing bert2BERT on the different settings of $E_b=\{0, 5, 10, 20\}$. The results are presented in Table 4 and the loss curves are displayed in Appendix C. We observe that our method achieves the best efficiency when the epoch number is set to 5, while a larger or smaller epoch number will bring a negative impact.

Model	FLOPs ($\times 1e19$)	Ratio (Saving)	Loss (MLM)	Avg.
<i>bert2BERT : $\mathcal{S}(12, 512) \rightarrow \mathcal{T}(12, 768)$</i>				
bert2BERT ($E_b = 0$)	4.5	38.4%	1.434	89.4
bert2BERT ($E_b = 5$)	4.0	45.2%	1.433	89.3
bert2BERT ($E_b = 10$)	4.1	43.9%	1.436	89.3
bert2BERT ($E_b = 20$)	5.4	25.4%	1.448	89.1

Table 4: Effect of sub-model training epochs. Avg means the average score of SST-2/MNLI/SQuADv1.1.

5.3 Application on GPT

Datasets. To demonstrate that our method is generic, following the BERT setting, we also use the English Wikipedia and Book Corpus in the GPT-training. For the evaluation, we use the datasets of WikiText-2, PTB, and WikiText103 and evaluate these models under the zero-shot setting without fine-tuning on the training set.

Implementation Details. We use the architecture of $\{L=12, D=768\}$ for the GPT target model, and pre-train it with the learning rate of $1e-4$, training epochs of 20. For bert2BERT, we use the source model with an architecture of $\{L=12, D=512\}$, initialize the target model with AKI, and pre-train it by the full-model training ($E_b=0$).

Results and Analysis. We compare the original pre-training method and bert2BERT, the results are shown in Table 5 and Appendix D. We observe that the proposed method saves **47%** computation cost of GPT pre-training, exhibiting a similar trend to BERT pre-training. Although GPT and BERT have different architectures (e.g., post-LN and pre-LN (Xiong et al., 2020)) and are pre-trained with different tasks, bert2BERT saves a significant amount of training cost on both these two models, which shows that the proposed method is generic and is effective for different kinds of PLMs.

Model	FLOPs ($\times 1e19$)	PTB (w/o FT)	WikiText-2 (w/o FT)	WikiText103 (w/o FT)
<i>bert2BERT : $\mathcal{S}(12, 512) \rightarrow \mathcal{T}(12, 768)$</i>				
GPT	4.9	133.8	47.0	53.5
bert2BERT	2.6 (47% ↓)	132.1	47.9	53.0

Table 5: Experiments on GPT. We report the perplexity for these tasks. “w/o FT” means that the pre-trained model is directly evaluated on the test set without fine-tuning on the train set.

5.4 Application on T5

Datasets. To demonstrate that our method can be used to train larger models, we use the Baidu Wikipedia, Sougou Wikipedia, and Zhihu to train the T5 model (Raffel et al., 2020). For the evaluation, we use the dataset of the original Chinese natural language inference task (OCNLI) (Hu et al., 2020).

Implementation Details. Since the bert2BERT method is suitable for BERT and GPT, it can also be used for the T5 model, which consists of an encoder and a decoder. The target T5 model’s architecture is $\{L_e=12, L_d=12, D=1024, A=16\}$, where

L_e and L_d means the numbers of encoder and decoder Transformer layers respectively, D means the hidden size, A means the number of attention heads. We pre-train it with the learning rate of $1e-4$, batch size of 1024. For bert2BERT, we use the source model with an architecture of $\{L_e=12, L_d=12, D=256, A=4\}$, initialize the target model with FPI, and pre-train it by the full-model training ($E_b=0$). Note that the scale gap between the source model and the target model is over 10 times (31M vs. 360M), which is a challenging setting.

Results and Analysis. We compare the original pre-training method and bert2BERT method on the T5 model, the results are shown in Table 6. We observe that the proposed method saves at least 25% computation cost of T5 pre-training. It demonstrates the effectiveness of the method on larger models.

Model	FLOPs ($\times 1e20$)	Loss (MLM)	OCNLI (Acc)
<i>bert2BERT</i> : $\mathcal{S}(12, 12, 256, 4) \rightarrow \mathcal{T}(12, 12, 1024, 16)$			
T5	1.6	1.90	72.03
bert2BERT	1.2 (25% ↓)	1.90	72.75

Table 6: Experiments on the T5 model.

6 Conclusion and Future Work

This paper proposes an efficient pre-training method, bert2BERT, which reuses the parameters of the small trained model as the initialization parameters of the large model. We employ the proposed method in BERT and GPT under different settings of model sizes. The extensive results show that bert2BERT is generic to Transformer-based models and saves a significant amount of computation cost. Moreover, the detailed analysis shows that our techniques, function-preserving, advanced knowledge initialization, and two-stage pre-training, are all effective. In the future, we will apply bert2BERT on training super large-scale language models (e.g., use the 10B source model to train the 100B target model) and extends its scope to other PLMs such as ELECTRA and BART (Lewis et al., 2020).

Acknowledgements

This work is supported in part by NSFC (Grant No. 61872215), and Shenzhen Science and Technology Program (Grant No. RCYX20200714114523079). We would like to thank Yifeng Liu, Binbin Deng, Ziliang Yang, Jiabin Shi for their support of this work.

References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#). *ArXiv preprint, abs/1607.06450*.
- 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. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Cheng Chen, Yichun Yin, Lifeng Shang, Zhi Wang, Xin Jiang, Xiao Chen, and Qun Liu. 2021. [Extract then distill: Efficient and effective task-agnostic BERT distillation](#). In *Artificial Neural Networks and Machine Learning - ICANN 2021 - 30th International Conference on Artificial Neural Networks, Bratislava, Slovakia, September 14-17, 2021, Proceedings, Part III*, volume 12893 of *Lecture Notes in Computer Science*, pages 570–581. Springer.
- Tianqi Chen, Ian J. Goodfellow, and Jonathon Shlens. 2016. [Net2net: Accelerating learning via knowledge transfer](#). In *ICLR*.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: pre-training text encoders as discriminators rather than generators](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- William Fedus, Barret Zoph, and Noam Shazeer. 2021. [Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity](#). *CoRR*.

- A. Feng and P. Panda. 2020. [Energy-efficient and robust cumulative training with net2net transformation](#). In *IJCNN*.
- Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2019. [Efficient training of BERT by progressively stacking](#). In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2337–2346. PMLR.
- Xiaotao Gu, Liyuan Liu, Hongkun Yu, Jing Li, Chen Chen, and Jiawei Han. 2021. [On the transformer growth for progressive BERT training](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5174–5180, Online. Association for Computational Linguistics.
- Song Han, Jeff Pool, John Tran, and William Dally. 2015. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28.
- Dan Hendrycks and Kevin Gimpel. 2016. [Gaussian error linear units \(gelus\)](#). *ArXiv preprint*, abs/1606.08415.
- Hai Hu, Kyle Richardson, Liang Xu, Lu Li, Sandra Kübler, and Lawrence S. Moss. 2020. [OCNLI: original chinese natural language inference](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 3512–3526. Association for Computational Linguistics.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. [What does BERT learn about the structure of language?](#) In *ACL*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [ALBERT: A lite BERT for self-supervised learning of language representations](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *ACL*, pages 7871–7880, Online. Association for Computational Linguistics.
- 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](#). *CoRR*.
- Sinno Jialin Pan and Qiang Yang. 2010. [A survey on transfer learning](#). *TKDE*.
- Yujia Qin, Yankai Lin, Jing Yi, Jiajie Zhang, Xu Han, Zhengyan Zhang, YuSheng Su, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2021. [Knowledge inheritance for pre-trained language models](#). *CoRR*.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. [Improving language understanding by generative pre-training](#).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. 2020. Green ai. *Communications of the ACM*, 63(12):54–63.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. [Megatron-lm: Training multi-billion parameter language models using model parallelism](#). *CoRR*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019a. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Dilin Wang, Meng Li, Lemeng Wu, Vikas Chandra, and Qiang Liu. 2019b. [Energy-aware neural architecture optimization with fast splitting steepest descent](#). *CoRR*.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. [Neural network acceptability judgments](#). *Transactions of the Association for Computational Linguistics*, 7:625–641.

- Wei Wen, Yuxiong He, Samyam Rajbhandari, Minjia Zhang, Wenhan Wang, Fang Liu, Bin Hu, Yiran Chen, and Hai Li. 2018. [Learning intrinsic sparse structures within long short-term memory](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- Lemeng Wu, Bo Liu, Peter Stone, and Qiang Liu. 2020a. [Firefly neural architecture descent: a general approach for growing neural networks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Lemeng Wu, Dilin Wang, and Qiang Liu. 2019. [Splitting steepest descent for growing neural architectures](#). In *NeurIPS*, pages 10655–10665.
- Lemeng Wu, Mao Ye, Qi Lei, Jason D. Lee, and Qiang Liu. 2020b. [Steepest descent neural architecture optimization: Escaping local optimum with signed neural splitting](#). *CoRR*.
- Qiyu Wu, Chen Xing, Yatao Li, Guolin Ke, Di He, and Tie-Yan Liu. 2021. [Taking notes on the fly helps language pre-training](#). In *ICLR*.
- Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. 2020. [On layer normalization in the transformer architecture](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 10524–10533. PMLR.
- Cheng Yang, Shengnan Wang, Chao Yang, Yuechuan Li, Ru He, and Jingqiao Zhang. 2020. [Progressively stacking 2.0: A multi-stage layerwise training method for bert training speedup](#). *arXiv preprint arXiv:2011.13635*.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. [Xlnet: Generalized autoregressive pretraining for language understanding](#). In *NeurIPS*, pages 5754–5764.
- Yang You, Jing Li, Sashank J. Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. 2020. [Large batch optimization for deep learning: Training BERT in 76 minutes](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Wei Zeng, Xiaozhe Ren, Teng Su, Hui Wang, Yi Liao, Zhiwei Wang, Xin Jiang, ZhenZhang Yang, Kaisheng Wang, Xiaoda Zhang, Chen Li, Ziyang Gong, Yifan Yao, Xinjing Huang, Jun Wang, Jianfeng Yu, Qi Guo, Yue Yu, Yan Zhang, Jin Wang, Hengtao Tao, Dasen Yan, Zexuan Yi, Fang Peng, Fangqing Jiang, Han Zhang, Lingfeng Deng, Yehong Zhang, Zhe Lin, Chao Zhang, Shaojie Zhang, Mingyue Guo, Shanzhi Gu, Gaojun Fan, Yaowei Wang, Xuefeng Jin, Qun Liu, and Yonghong Tian. 2021. [Pangu- \$\alpha\$: Large-scale autoregressive pretrained chinese language models with auto-parallel computation](#).
- Minjia Zhang and Yuxiong He. 2020. [Accelerating training of transformer-based language models with progressive layer dropping](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Aligning books and movies: Towards story-like visual explanations by watching movies and reading books](#). In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 19–27. IEEE Computer Society.

A Ablation Study of bert2BERT

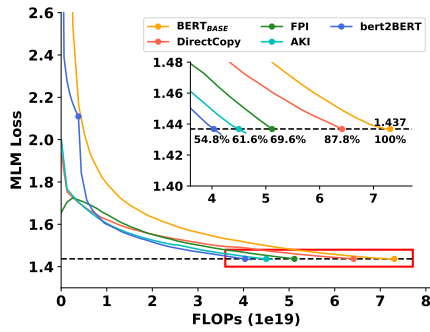


Figure 5: Ablation study of bert2BERT. bert2BERT means the combination of AKI and two-stage pre-training.

The ablation study of bert2BERT is displayed in Table 5. From the table, we observe that: (1) all the proposed methods is better than the original pre-training method and DirectCopy; (2) although AKI has a worse initialization than FPI, it achieves faster convergence rate than FPI; (3) the two-stage pre-training further reduce the cost from 61.6% to 54.8%; (4) the FPI curve has an upward trend at the beginning. We conjecture that it is due to the symmetry brought by FPI and the model needs some optimization time to break this symmetry.

B bert2BERT with smaller source model

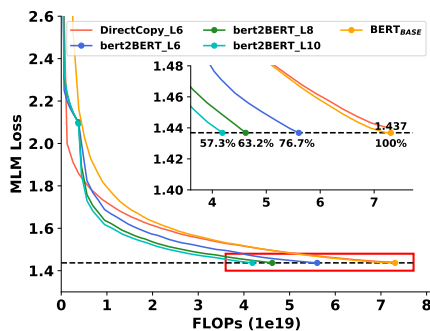


Figure 6: Loss curves of bert2BERT and baselines with smaller source models.

We test bert2BERT with different source models and the loss curves are represented in Figure 6.

C Effect of sub-model training epochs

We study the effect of sub-model training epochs on the pre-training efficiency. The loss curves are represented in Figure 7. Note that the setting $E_b = 20$ has not achieved the same loss (1.437) as the baseline BERT_{BASE} in the 40 training epochs.

D Application on GPT

The loss curve of our method on GPT application is displayed in Figure 8.

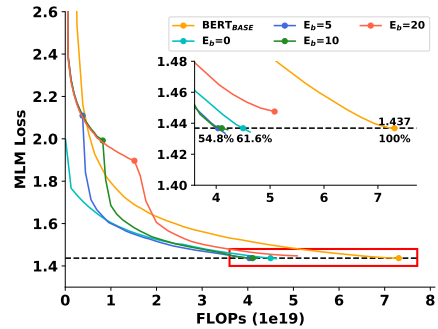


Figure 7: Loss curves of bert2BERT with different sub-model training epochs.

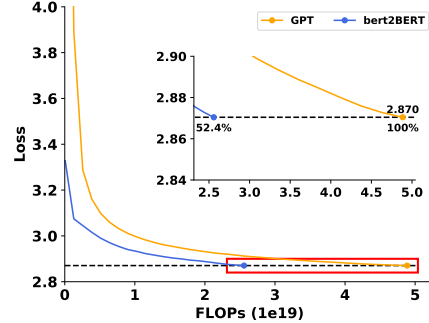


Figure 8: Pre-training loss curves of GPT.

E Comparisons of Attention Patterns

We take the source model $\mathcal{S}(4, 256)$ and target model $\mathcal{T}(4, 512)$ as an example to analyze the attention patterns of DirectCopy in Figure 10, FPI in Figure 11 and AKI in Figure 12.

We display the attention patterns of the source model $\mathcal{S}(4, 256)$ in Figure 9. Compared with the source model, we observe that the newly added attention patterns of DirectCopy are messy, and the randomly initialized parameters destroy the attention patterns of the source model. The proposed FPI method makes the new model have the same attention patterns as the source model, thus the knowledge of the source model is preserved. However, FPI always induces symmetrical attention patterns in the same layer. This symmetry will hinder the convergence. To handle this problem, we use AKI method to reuse the parameters of the upper layer (advanced knowledge) to break the symmetry, and meanwhile make the knowledge in the same layer richer. Through the AKI method, the attention patterns of the upper layer can be also maintained well in the target model. For example, as shown in Figure 12, the newly added attention patterns of the 1st layer in the target model are similar to the ones of the 2nd layer in the source model.

F Illustration of FPI and AKI process

We illustrate the process of FPI and AKI in Figure 13 and 14 respectively.

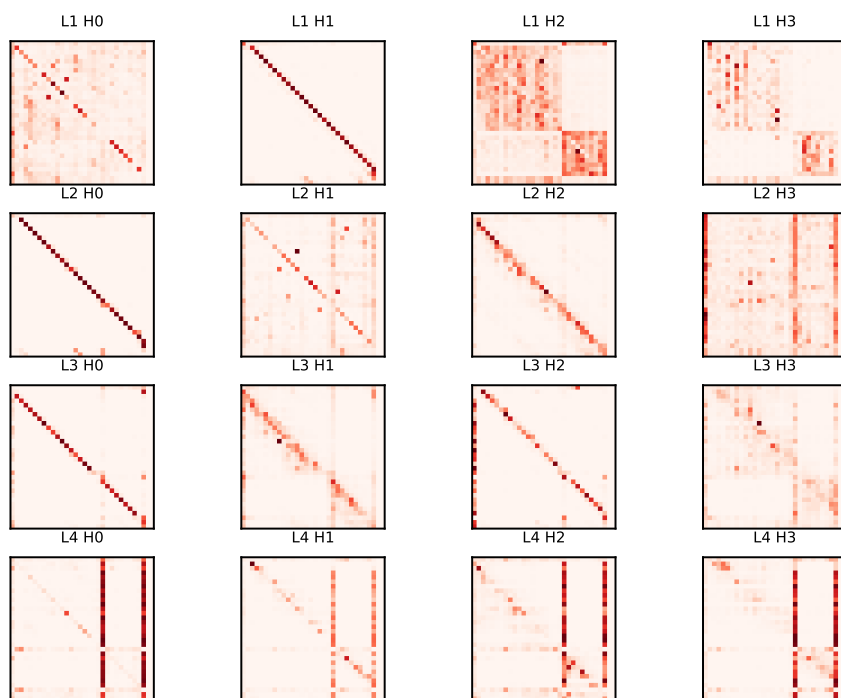


Figure 9: Attention patterns of the source model $\mathcal{S}(4, 256)$, which has 4 attention heads in each layer.

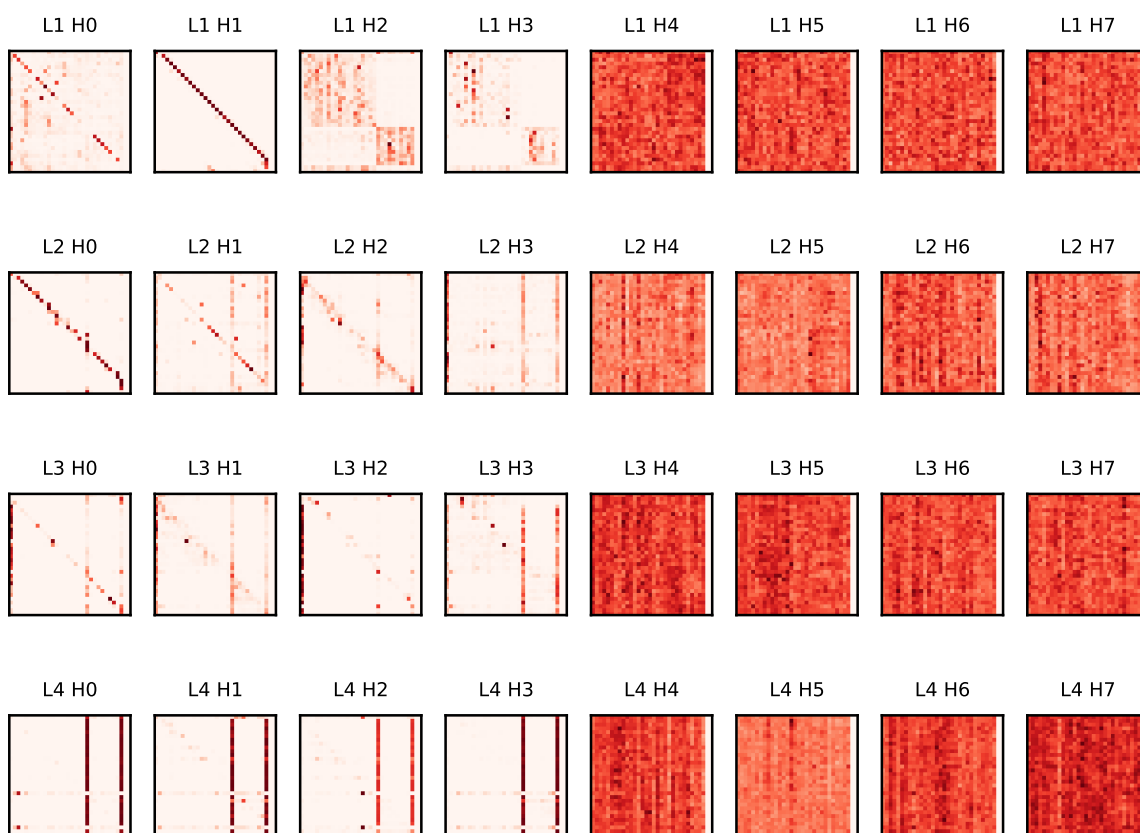


Figure 10: Attention patterns of the target model $\mathcal{T}(4, 512)$ based on the baseline DirectCopy method. The first 4 attention patterns (H0-H3) in each row correspond to the source model's attention patterns, and the last 4 attention patterns (H4-H7) are newly added.

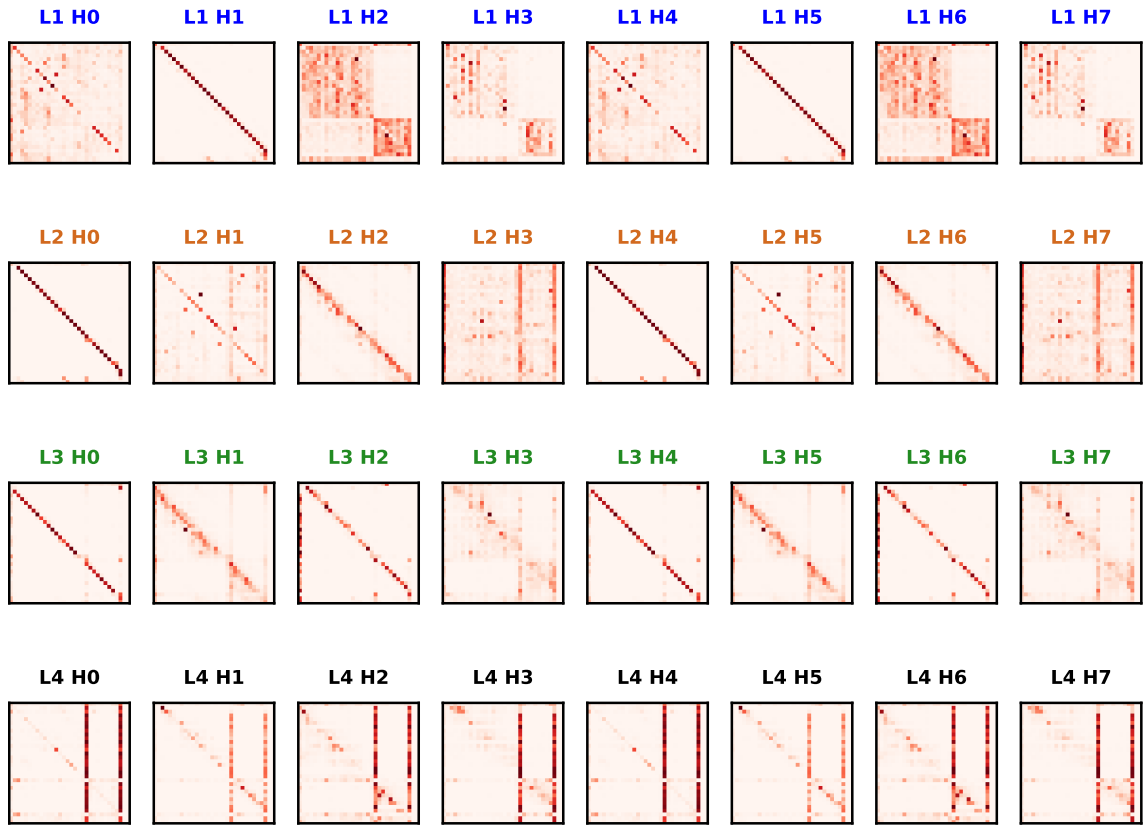


Figure 11: Attention patterns of the target model $\mathcal{T}(4, 512)$ based on our FPI method. The last 4 attention patterns (H4-H7) in each row are obtained by FPI expansion.

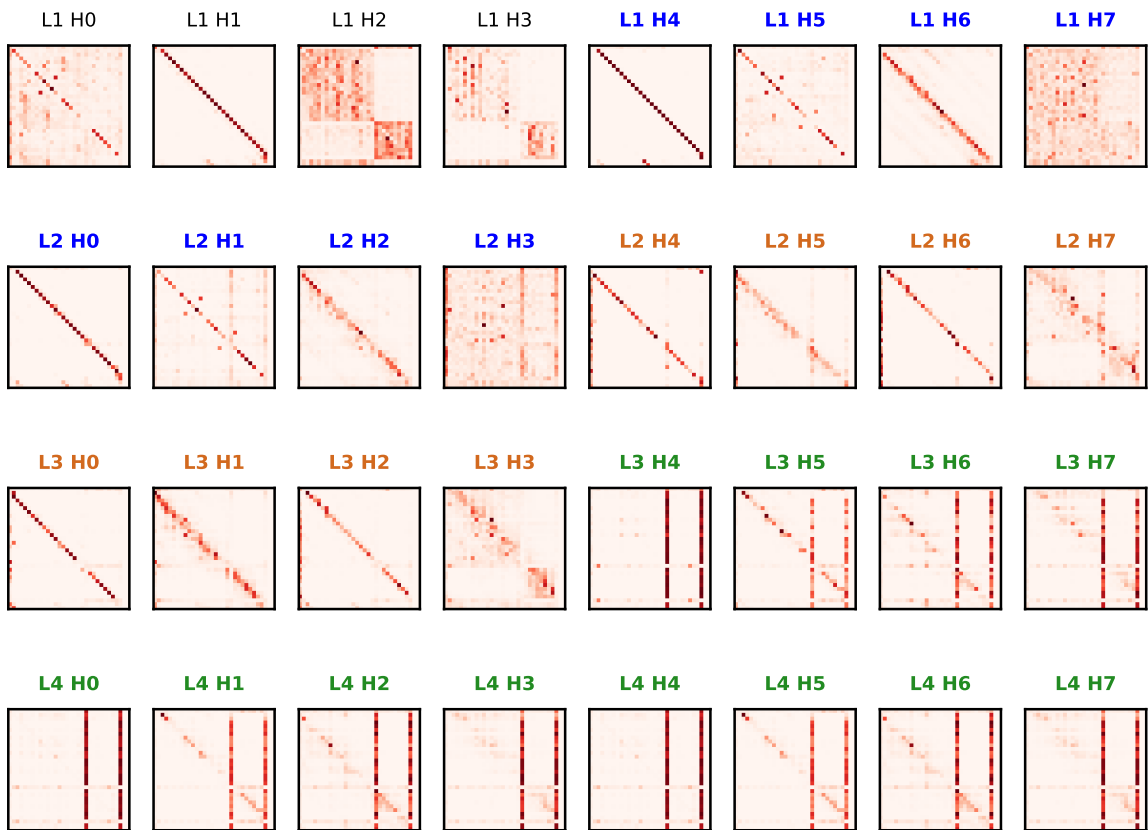


Figure 12: Attention patterns of the target model $\mathcal{T}(4, 512)$ based on our AKI method. The last 4 attention patterns (H4-H7) in each row are obtained by AKI expansion.

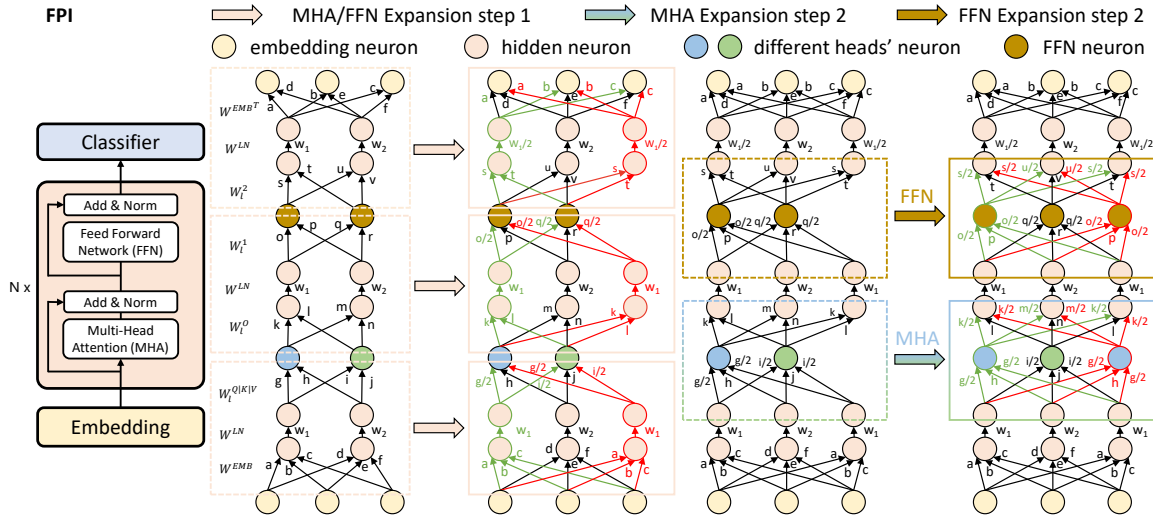


Figure 13: FPI process. We use the FPI method to widen the source model with a width of 2 into a target model with a width of 3. In the example, the source model and the target model have 2 and 3 attention heads respectively. And the head dimension is 1. To facilitate the illustration, we reduce the number of neurons in the FFN layer. We also note that since the MLM classifier of BERT is a transposition of the Embedding layer, they share a parameter matrix. Therefore, in step 1, we expand the MLM classifier by re-scaling the parameter values of the LN layer below the MLM classifier instead of following formula 7.

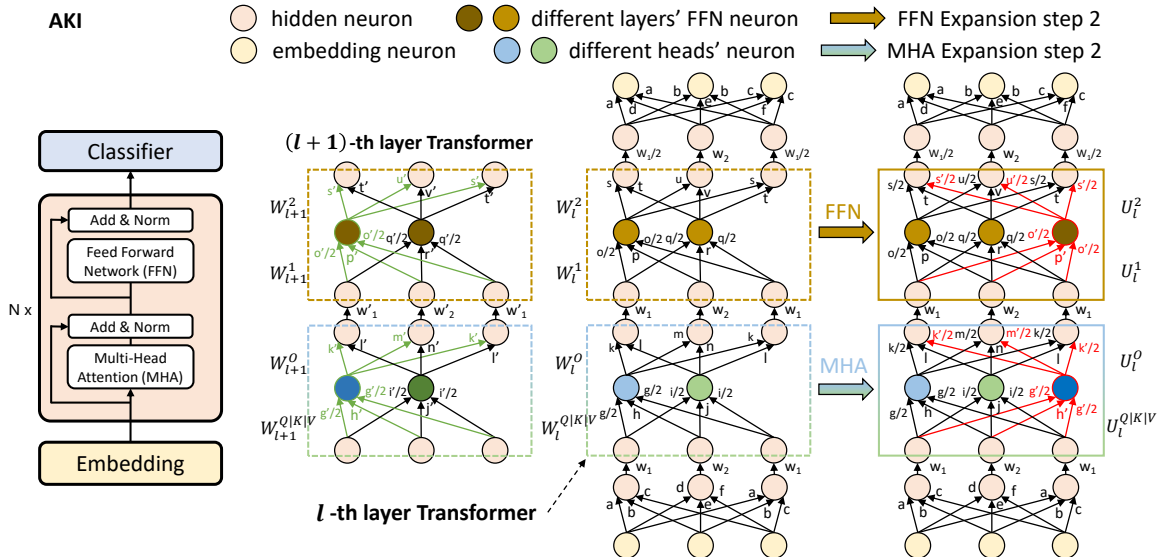


Figure 14: AKI process. We ignore its first step because it is the same as FPI's first step. The main difference between AKI and FPI is that in step 2, AKI copies some attention heads in MHA and some parameters in FFN of the $(l+1)$ -th layer instead of only copying the l -th layer to construct the new l -th layer Transformer.