# TRANSPA: TOWARDS EFFICIENT STRUCTURED SPARSE TRAINING FOR TRANSFORMERS

Anonymous authors

Paper under double-blind review

#### Abstract

Transformers have emerged as the backbone neural network architecture in today's AI applications. Due to their high complexity, sparsifying transformers, at both pre-training and fine-tuning stages, is very attractive for lower the training and inference costs. In this paper, we propose TranSpa, an efficient structured sparse training approach for language and vision transformers. Unlike prior works focusing on individual building blocks, TranSpa fully considers the correlation between the weight matrices and their component rows/columns, and performs the coupled estimation and coupled sparsification. To achieve that, TranSpa introduces the use of new granularity when calibrating the importance of structural components in the transformer and removing the insignificant parts. Evaluations across different models, in both pre-training and fine-tuning scenarios, demonstrate the effectiveness of the proposed approach. TranSpa can bring  $1.6 \times$  size reduction with 0.6 lower perplexity when training GPT-2 model from scratch. It also enables  $1.6 \times$ training speedup over the existing sparse pre-training method. For training sparse LLaMA-1B from scratch, our approach reduces GPU memory usage by 50%, decreases training time by 21%, and achieves a  $1.6 \times$  speedup in inference throughput while maintaining model performance. Experiments of applying TranSpa for fine-tuning tasks also show significant performance improvement with respect to model accuracy and pruning cost reduction.

028 029

031

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

#### 1 INTRODUCTION

Thanks to their excellent performance for sequence modeling and scalability, transformers Vaswani et al. (2017) have served as the backbone neural network architecture across various important domains, such as natural language processing (NLP) Vaswani et al. (2017); Devlin et al. (2019); Wang et al. (2019a); OpenAI et al. (2023); Touvron et al. (2023), computer vision Dosovitskiy et al. (2020); Liu et al. (2021); Carion et al. (2020) and speech processing Dong et al. (2018); Gulati et al. (2020); Hsu et al. (2021). However, despite their unprecedented popularity, modern transformers suffer from high model complexity, causing expensive costs in both training and inference phases.

To alleviate these challenging issues, **sparse training**, a strategy that originated for optimizing convolutional neural networks (CNNs), has emerged as an attractive solution for efficient transformers. In general, given an input model  $\theta$ , sparse training aims to output a sparse model  $\theta_s$  with high task performance. As illustrated in Fig. 1, when  $\theta$  is randomly initialized, sparse training serves as an efficient *pre-training* method that can reduce the computational and memory costs of the expensive training-from-scratch process. When  $\theta$  is a pre-trained high-performance dense model, sparse training is essentially the well-known pruning technique, which imposes the sparsity on  $\theta$  in the *fine-tuning* process and trims down the inference cost.

Although sparse training has been well studied for slimming CNN models, the fundamentally different mechanism of transformers, *e.g.*, multi-head self-attention, make the existing CNN-oriented solutions not suitable for transformers. To date, the efficient sparse transformer training is still under-explored. More specifically, 1) for applying sparse training at the pre-training stage, Frankle & Carbin (2018), Chen et al. (2021a;b) propose to train sparse vision and NLP transformers from scratch, with focus on exploring head-level sparsity. Also Dao et al. (2022) proposes to use specially structured matrix for sparse transformer training; 2) for applying sparse training at the fine-tuning stage, Ma et al. (2023); Chen et al. (2021b); Yu et al. (2022) explore the structured sparsity at head



Figure 1: Our proposed TranSpa is a structured sparse training approach, applicable for both pretraining and fine-tuning stages.

levels. Observing the potential layer-wise redundancy, Chen et al. (2024); Men et al. (2024) propose
to remove unimportant layers in the transformer architecture. In addition, Ashkboos et al. (2024);
van der Ouderaa et al. (2023) use finer granularity at row and column level for structured pruning,
achieving good compression performance.

070 **Technical Contributions.** In this paper, we propose to perform structured sparse transformer training 071 <sup>1</sup> from a new perspective. Unlike prior works focusing on evaluating the importance of individual 072 components of transformer architecture, e.g., a weight matrix or its component row/column, we 073 propose to fully consider the inter-matrix and inter-row/column correlation, leading to new mechanism for importance assessment and structural removal. More specifically, at the *coupled estimation* step, 074 we introduce a new granularity, namely, coupled weight matrices, to calibrate the architectural 075 importance of transformer model, providing more fine-grained measurement via exploring the 076 inherent coupling effects of weight matrices. Then, at the *coupled sparsification* step, for the row 077 and column within those unimportant coupled matrices, we propose to rank and remove them in a coupled way, maximally preserving the inherent inter-row/column correlation in the weight matrix. 079

We apply our approach, namely TranSpa, in both pre-training and fine-tuning scenarios. Evaluation 080 across different model types (e.g., vision transformers, and large language model (LLMs)) show that 081 TranSpa brings significant performance improvement with respect to training speed, model accuracy, and cost reduction. For instance, when training the GPT-2 model from scratch, our approach can bring 083  $1.6 \times$  size reduction with 0.6 lower perplexity (PPL). It also brings  $1.6 \times$  training speedup compared 084 to the existing sparse pre-training methods. For pre-training sparse LLaMA-1B from scratch, our 085 approach cuts GPU memory usage by half, shortens training time by 21%, and boosts inference 086 throughput by  $1.6 \times$ . For training sparse DeiT from scratch, TranSpa brings 0.8% accuracy increase 087 over the state-of-the-art solutions. Experiments on pruning LLMs also demonstrate the effectiveness 088 of our approach.

090 091

092 093

094

099

100

101

102 103

104

065

### 2 BACKGROUND

2.1 PRELIMINARIES

**Notation.** We represent tensors using boldface calligraphic script, denoted as  $\mathcal{X}$ . Matrices and vectors are indicated with boldface capital and lowercase letters, such as  $\mathbf{X}$  and  $\mathbf{x}$ , respectively. Furthermore, non-boldface letters with indices, *e.g.*,  $\mathcal{X}(i_1, \dots, i_d)$ , X(i, j), and x(i), denote the entries for a *d*-dimensional tensor  $\mathcal{X}$ , a matrix  $\mathbf{X}$ , and a vector  $\mathbf{x}$ , respectively.

**Transformer.** For Transformer-based models, the key components include the Multi-Head Attention (MHA) and Feed-Forward Network (FFN). More specifcally, the MHA operation is defined as follows:

$$MHA(\boldsymbol{X}_Q, \boldsymbol{X}_K, \boldsymbol{X}_V) = Concat(head_1, \dots, head_h) \boldsymbol{W}^O,$$
(1)

 <sup>&</sup>lt;sup>1</sup>In this paper we focus on structured sparsity, which can bring measured speedup on off-the-shelf hardware;
 while the unstructured sparsity typically cannot. Though some Nvidia GPUs now provide hardware support for
 semi-structured 2:4 sparsity, this feature is only available for the high-end GPUs such as A100 and up, limiting its practice in many budget-limited resource-limited applications.

where  $X_Q, X_K, X_V \in \mathbb{R}^{l \times e_m}$  are the length-*l* input sequences,  $e_m$  is the embedding dimension, and *h* is the number of attention heads. Here each attention head  $head_i$  operates as below:

110 111

120

121 122 123

124 125

 $head_{i} = \text{Attention}(\boldsymbol{X}_{Q}\boldsymbol{W}_{i}^{Q}, \boldsymbol{X}_{K}\boldsymbol{W}_{i}^{K}, \boldsymbol{X}_{V}\boldsymbol{W}_{i}^{V}) = \text{softmax}\left(\frac{\boldsymbol{X}_{Q}\boldsymbol{W}_{i}^{Q}(\boldsymbol{X}_{K}\boldsymbol{W}_{i}^{K})^{\top}}{\sqrt{e}}\right)\boldsymbol{X}_{V}\boldsymbol{W}_{i}^{V},$ (2)

where  $W_i^Q$ ,  $W_i^K$ ,  $W_i^V \in \mathbb{R}^{e_m \times e}$ ,  $W^O \in \mathbb{R}^{e_m \times e_m}$ , and  $e_m = e \times h$ . The FFN consists of two fully connected layers with a Gaussian Error Linear Unit (GELU) activation function applied in between. Let X represent the input embeddings, and  $W_{in} \in \mathbb{R}^{e_m \times d}$ ,  $W_{out} \in \mathbb{R}^{d \times e_m}$ ,  $b_{in}$ ,  $b_{out}$  denote the weight matrices and bias vectors, respectively. The operation of FFN can be then defined as follows:

$$FFN(\boldsymbol{X}) = GELU(\boldsymbol{X}\boldsymbol{W}_{in} + \boldsymbol{b}_{in})\boldsymbol{W}_{out} + \boldsymbol{b}_{out}.$$
(3)

#### 2.2 RELATED WORKS

**Sparse Training for Transformer Pre-training.** Applying sparse training at the pre-training stage 126 is very attractive for reducing the high complexity of the costly train-from-scratch process. However, 127 unlike the well-studied sparse CNN pre-training, to date sparse transformer pre-training is still 128 under-explored. Chen et al. (2021a) dynamically extracts and trains sparse sub-networks, either 129 in unstructured or structured ways, thereby alleviating the training memory bottleneck for Vision 130 Transformers. Inspired by the Lottery Ticket Hypothesis originated in CNN, Chen et al. (2021b) 131 performs early detection of the structured winning lottery tickets for transformers, improving the pre-132 training efficiency. In Dao et al. (2022), a class of structured matrices, which can closely approximate 133 the dense weight matrices, are proposed for hardware-efficient sparse pre-training, accelerating the 134 overall training process. Notice that the methods developed in Chen et al. (2021b;a); Dao et al. (2022) 135 can also be applied at the fine-tuning stage for transformer pruning.

136 **Sparse Training for LLM Pruning.** Due to the high costs of pre-trained LLM models, using sparse 137 training to trim down LLM size is a promising solution for efficient deployment. In general, pruning 138 LLM can be performed in either unstructured or structured way. Unstructured pruning Frantar & 139 Alistarh (2023); Sun et al. (2023); Zhang et al. (2023); Sun et al. (2023); Xia et al. (2023); Malla et al. 140 (2024) focus on removing unimportant individual weights, achieving high compression performance 141 but limited computational efficiency due to the unstructured sparsity pattern. Structured pruning 142 Ma et al. (2023); Xia et al. (2023); Ashkboos et al. (2024); An et al. (2024); Chen et al. (2023b;a); Zhao et al. (2024a); Yang et al. (2024) aims to sparsify some building components of transformer 143 architecture, e.g., head, row, layer, etc., offering wall-clock time inference speedup. To guide the 144 guide the selection of elements to be removed, typical pruning metrics include magnitude-based 145 Sun et al. (2023); An et al. (2024) and loss-based Ma et al. (2023); van der Ouderaa et al. (2023). 146 Magnitude-based methods use the absolute values of weights, whereas loss-based approaches assess 147 the impact of pruning on model loss, often using the gradient information obtained from Taylor 148 expansion. 149

Sparse Training for CNN Pre-training/Pruning. Sparse training for CNN, at pre-training and 150 fine-tuning stages, has been well-studied in the literature. For efficient sparse CNN pre-training, 151 Mocanu et al. (2018) explores to prune and grow the same amount of weights periodically and 152 iteratively. Mostafa & Wang (2019) proposes to automatically adjust the sparsity levels, achieving 153 good scalability and high computational efficiency. Evci et al. (2020) uses the gradient-based criteria 154 to grow the weights with Erdos-Renyi Kernel initialization, and a memory-economic scheme is 155 proposed in Yuan et al. (2021) for training on the edge devices. Recently, Chen et al. (2023c) 156 proposes to automatically sparsify the CNN model from scratch, obtaining a compact model without 157 iterative fine-tuning. On the other hand, CNN pruning can be roughly categorized to unstructured 158 pruning (for weights) Han et al. (2015a;b) and structured pruning (for channels/filters) He et al. 159 (2019); Wang et al. (2019b); Yvinec et al. (2021); Hou et al. (2022); Lin et al. (2020); Dubey et al. (2018); Sui et al. (2021); Tan & Motani (2020); Luo et al. (2017). Considering the importance of 160 structured sparsity for inference speedup, most of the state-of-the-art CNN pruning works focus on 161 structured pruning.



Figure 2: Key steps of TranSpa: 1) Estimate the importance of coupled weight matrix  $W_{couple}$ ; and 2) Remove the coupled row/column pair in unimportant  $W_{couple}$ .

#### 3 Method

Fig. 2 shows the overall procedure of TranSpa, which consists of two key steps. (1) Coupled Estimation (Section 3.1). This step assesses the importance of the building components of transformers, using our proposed new calibration granularity at the coupled weight matrix level. (2) Coupled Sparsification (Section 3.2). Once the unimportant coupled weight matrices are identified, TranSpa further gauges the fine-grained importance of the coupled row/column pairs in the component matrices and discards the insignificant ones.

185

173

174

175 176 177

178

186 187

#### 3.1 COUPLED ESTIMATION: COUPLED WEIGHT MATRIX-WISE IMPORTANCE CALIBRATION

As extensively studied in the literature Frantar & Alistarh (2023); Sun et al. (2023); van der Ouderaa et al. (2023), estimating the unimportant components of neural networks plays a crucial role for model sparsification. In general, importance estimation can be performed at different granularity levels, such as weight, neuron, layer, and head. More specifically, when aiming for obtaining the structured sparse transformers, identifying the insignificant heads and layers is the most common practice Chen et al. (2021a;b); Men et al. (2024); Ma et al. (2023).

Unlike existing works, we propose to assess the importance of structural components within transformers using a new granularity, namely *coupled weight matrix*. This idea is motivated by the observation that the transformer has its unique computing pattern and model topology, and hence coupling the weight matrices can provide rich information for importance estimation. Next, we describe the details of our proposal.

**Coupled Weight Matrix in MHA.** Recall that Eq. 1 and Eq. 2 depict the computations of MHA, which consists of four types of weight matrices  $W^Q$ ,  $W^K$ ,  $W^V$  and  $W^O$ . We propose to estimate the structural importance of MHA by analyzing the combined effect of these matrices. More specifically, consider the following mathematical reformulation of Eq. 2:

203 204 205

206

$$MHA(\boldsymbol{X}_{Q}, \boldsymbol{X}_{K}, \boldsymbol{X}_{V}) = \sum_{i=1}^{h} head_{i} \boldsymbol{W}_{i}^{O} = \sum_{i=1}^{h} \operatorname{softmax}(\frac{\boldsymbol{X}_{Q}(\boldsymbol{W}_{i}^{Q} \boldsymbol{W}_{i}^{K^{\top}}) \boldsymbol{X}_{K}^{\top}}{\sqrt{e}} \boldsymbol{X}_{V}(\boldsymbol{W}_{i}^{V} \boldsymbol{W}_{i}^{O}),$$
(4)

where  $W_i^O \in \mathbb{R}^{e \times e_m}$  and  $W^O = \text{Concat}(W_1^O, \dots, W_h^O)$ . It is seen that  $W_i^{QK} = W_i^Q W_i^{K^\top}$ and  $W_i^{VO} = W_i^V W_i^O$ , as the combination of two weight matrices, can serve as the structural components in MHA. Following this perspective, we can then set the granularity of importance estimation at the level of those coupled weight matrices. We believe this strategy brings two benefits: *i*) it provides more fine-grained measurement than the commonly adopted head-level calibration; and *ii*) it meanwhile naturally explores the inter-matrix correlation within the attention heads, avoiding the limitations if only focusing on individual weight matrices.

**Coupled Weight Matrix in FFN.** Following the same philosophy, we also evaluate the importance of the building blocks in FFN from the lens of the coupled weight matrix. More specifically, consider

216 the FFN computation (Eq. 3) can be reformulated as follows: 217

218 
$$FFN(X) = 0.5Y_1 \odot (1 + erf(Y_1/\sqrt{2}))$$

219 220

221

222 223

224

225

231 232 233

237

241

 $W_{out}$ (5) $= 0.5 \mathbf{X} \mathbf{W}_{in} \mathbf{W}_{out} + 0.5 \mathrm{erf}(\mathbf{Y}_1 / \sqrt{2}) \odot \mathbf{X} \mathbf{W}_{in} \mathbf{W}_{out},$ 

where  $Y_1 = XW_{in}$ ,  $\odot$  is the Hadamard product, and  $erf(\cdot)$  denotes the error function as erf(x) = $\frac{2}{\sqrt{\pi}}\int_0^x \exp(-t^2)dt$  bounded by [-1,1]. From Eq. 5 it is seen that the coupled weight matrix  $\dot{W}_{12} = W_{in}W_{out}$  contributes a significant portion of the overall computation. Therefore, we propose to calibrate the importance of this combined matrix, as a good proxy, to identify the significance of the different structural components of FFN.

226 Coupled Weight Matrix-wise Importance. After setting the resolution of the importance estimator 227 with the proposed granularity, we can gauge the structural criticality of MHA and FFN accordingly. 228 To that end, we calibrate the empirical Fisher information Kunstner et al. (2019) in a coupled weight 229 matrix-wise way as follows: 230

$$\hat{I}(\boldsymbol{W}_{couple}) = \frac{1}{2} \sum_{k=1}^{2} \frac{1}{|\boldsymbol{W}_k|} \frac{1}{|\boldsymbol{\mathcal{D}}|} \sum_{i,j} \left( \frac{\partial \mathcal{L}(\boldsymbol{W}_k; \boldsymbol{\mathcal{D}})}{\partial \boldsymbol{W}_k} \right)_{i,j}^2, \text{ where } \boldsymbol{W}_{couple} = \boldsymbol{W}_1 \boldsymbol{W}_2.$$
(6)

234 Here for MHA,  $\{W_1, W_2\}$  is  $\{W_i^Q, W_i^{K^{\top}}\}$  and  $\{W_i^V, W_i^O\}$ ; while for FFN,  $\{W_1, W_2\}$  is  $\{W_{in}, W_{out}\}$ . Also,  $\mathcal{D}$  is the training data and  $|\cdot|$  returns the size of the operand. Notice that 235 236 different from prior works Hsu et al. (2022); Sung et al. (2021), the empirical Fisher information is calculated for the entire coupled weight matrix  $W_{couple}$  instead of its entries. Smaller  $I(W_{couple})$ 238 indicates the lower importance of  $W_{couple}$  when sparsifying the model. 239

#### 240 3.2 COUPLED SPARSIFICATION: REMOVING THE COUPLED ROW/COLUMN PAIR

242 Upon obtaining the importance information of all the coupled weight matrices, the next step is to sparsify  $W_1$  and  $W_2$  belonging to those less significant  $W_{couple}$ . To that end, we perform 243 row/column-wise sparsification, a strategy with finer granularity than removing the entire  $W_i$ , 244 towards minimizing performance loss and preserving model structuredness. Notice that though 245 removing the rows/columns in the weight matrices of transformers has been reported in Chen et al. 246 (2021b;a); van der Ouderaa et al. (2023), we believe those individual weight matrix-oriented methods 247 are not the best suited in the scenario involved with coupled weight matrices. In other words, the 248 coupling effect between  $W_1$  and  $W_2$  should be fully considered and leveraged in the sparsification 249 process. Aiming at that, we propose coupled sparsification, a solution that removes the coupled 250 row/column pairs in  $W_1$  and  $W_2$ , with details described as below. 251

Inspiration from tSVD. The key idea of coupled sparsification is inspired by the philosophy of 252 truncated singular value decomposition (tSVD). More specifically, recall that a matrix  $M \in \mathbb{R}^{m \times n}$ 253 can be exactly factorized to two matrices via SVD as: 254

255 256

257

258 259

260

261

262

263 264 265

$$\boldsymbol{M} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top} = \begin{pmatrix} \boldsymbol{U}_{m,r} & \boldsymbol{U}_{m,(m-r)} \end{pmatrix} \begin{pmatrix} \boldsymbol{\Sigma}_{r,r} & \\ \boldsymbol{\Sigma}_{(m-r),(n-r)} \end{pmatrix} \begin{pmatrix} \boldsymbol{V}_{r,n} \\ \boldsymbol{V}_{(n-r),n}^{\top} \end{pmatrix}$$
$$= \begin{pmatrix} \boldsymbol{U}_{m,r}\boldsymbol{\Sigma}_{r,r}^{1/2} & \boldsymbol{U}_{m,(m-r)}\boldsymbol{\Sigma}_{(m-r),(n-r)}^{1/2} \end{pmatrix} \begin{pmatrix} \boldsymbol{\Sigma}_{r,r}^{1/2}\boldsymbol{V}_{r,n}^{\top} \\ \boldsymbol{\Sigma}_{(m-r),(n-r)}^{1/2} \boldsymbol{V}_{(n-r),n}^{\top} \end{pmatrix} = \boldsymbol{M}_{1}\boldsymbol{M}_{2},$$
(7)

\_ \_ \_

where  $\Sigma_{r,r} \in \mathbb{R}^{r \times r}$  is a diagonal matrix containing r largest singular values  $\sigma_i$ 's and  $\Sigma_{(m-r),(n-r)} \in \mathbb{R}^{r \times r}$  $\mathbb{R}^{(m-r)\times(n-r)}$  is a diagonal rectangular matrix containing the smallest  $(\min(m, n) - r) \sigma_i$ 's. Then, consider a rank-r tSVD Wall et al. (2003) for approximating M is performed as:

$$\boldsymbol{M} \approx (\boldsymbol{U}_{m,r}\boldsymbol{\Sigma}_{r,r}^{1/2})(\boldsymbol{\Sigma}_{r,r}^{1/2}\boldsymbol{V}_{r,n}^{\top}) = (\boldsymbol{M}_1 \odot \boldsymbol{S}_{mask1})(\boldsymbol{M}_2 \odot \boldsymbol{S}_{mask2}).$$
(8)

266 Notice that here  $S_{maski}$  is the binary matrix that essentially removes a sets of rows/columns of 267  $M_i$ . Meanwhile, it is theoretically proven that tSVD provides the *optimal* rank-r approximation for M Wall et al. (2003). Therefore, the row/column-wise sparsification policy for  $M_1$  and  $M_2$  in 268 tSVD, by its nature, brings important insights for sparsifying two component matrices with closely 269 approximating their combination. More specifically, comparing Eq. 7 and Eq. 8, we can see that the

removed rows/columns are  $U_{m,(m-r)}\Sigma_{(m-r),(n-r)}^{1/2}$  and  $\Sigma_{(m-r),(n-r)}^{1/2}V_{(n-r),n}^{\top}$ , which exhibit the following characteristics:

273 <u>**Observation #1**</u> (Small  $\ell_1$  Norm) The removed rows/columns of  $M_1$  and  $M_2$  typically have smaller 274  $\ell_1$  norm, as  $\Sigma_{(m-r),(n-r)}$  only contains very least significant  $\sigma_i$ 's.

275 <u>Observation #2</u> (Aligned Removal) The removed rows in  $M_1$  and columns in  $M_2$  always have the 276 same indices, exhibiting the positional alignment of the sparsification.

**Sparsifying Coupled Weight Matrices.** Considering the mathematical format of Eq. 8 is highly similar to our desired task of sparsifying the coupled  $W_1$  and  $W_2$ , a natural idea is directly performing tSVD (Eq. 8) on  $W = W_1W_2$  to obtain sparse  $W_1$  and  $W_2$ . However, this strategy is mathematically infeasible because Eq. 7 is generally not invertible, *e.g.*, we cannot use SVD( $W_1W_2$ ) to reconstruct  $W_1$  or  $W_2$ . Fortunately, we can still leverage the observations extracted from tSVD process to design the sparsification policy for  $W_1$  and  $W_2$ . More specifically, we propose to calculate the importance scores of the row/column vectors of  $W_i$ :

$$\boldsymbol{v} = [v_1, v_2, \dots, v_d] = \frac{\|\operatorname{vec}(\boldsymbol{W}_1)\|_{1, col} \odot \|\operatorname{vec}(\boldsymbol{W}_2)\|_{1, row}}{\|\|\operatorname{vec}(\boldsymbol{W}_1)\|_{1, col} \odot \|\operatorname{vec}(\boldsymbol{W}_2)\|_{1, row}\|_1},\tag{9}$$

where  $\|\operatorname{vec}(\cdot)\|_{1,col}$  and  $\|\operatorname{vec}(\cdot)\|_{1,row}$  represent the column-wise and row-wise  $\ell_1$  norm of a matrix, respectively, *e.g.*,  $\|\operatorname{vec}(W_1)\|_{1,col} = \{\sum_{i=1}^m |W_1(i,j)|, (j = 1, 2, ..., d)\}$ . And  $\|\operatorname{vec}(\cdot)\|_1$  calculates the  $\ell_1$  norm of a vector. From Eq. 9 it is seen that  $v_i$  essentially calculates the normalized combined  $\ell_1$  for the *i*-th column of  $W_1$  and the *i*-th row of  $W_2$ , reflecting the two key observations we obtain from tSVD. Hence we can rank  $v_i$ 's to determine the important coupled row/column in  $W_1$  and  $W_2$ , and then remove those insignificant pairs.

**Overall Sparse Training Procedure.** Algorithm 1 describes the processing scheme of TranSpa. Here in each epoch, after identifying top-K least significant  $W_{couple}$ , the importance scores of the rows/columns in those component weight matrices are calculated. Once the cumulative important scores exceed the threshold  $\theta$ , the pair of rows and columns corresponding to the smallest score are removed. This gradual sparsification process continues till the overall model reaches the target budget size.

Algorithm 1 Processing Scheme of TranSpa	
<b>Input:</b> Random Initialized/Pre-trained model $\mathcal{W}$ , Dataset $\mathcal{D}$ ,	
top-K ratio, Target model size c, Cumulative threshold $\theta$	
<b>Output:</b> Sparse Model $\hat{W}$	
for $\hat{t}$ in $[1, 2, \cdots, T]$ do	
<b>Compute:</b> Loss of $\mathcal{W}$ on dataset $\mathcal{D}$	
Update: $W_t = W_{t-1} + \delta W$	
if $Param(\mathcal{W}_t) > c$ then	
Compute: $\hat{I}(\boldsymbol{W}_{couple}), \forall \boldsymbol{W}_{couple} \in \boldsymbol{\mathcal{W}}$	⊳ Eq. 6
<b>Select:</b> top- $K W_{couple}$ with smallest $\hat{I} \to \{W_{couple}\}^{top-K}$	
for $W_{couple}$ in $\{W_{couple}\}^{top-K}$ do	
<b>Compute:</b> $v$ of $W_{couple}$	⊳ Eq. 9
while $\sum v_i > \theta$ do	
<b>Remove</b> smallest $v_i$ and <i>i</i> -th row/col. of $W_{couple}$	
<b>Remove</b> <i>i</i> -th row/col. optimizer moments of $W_{couple}$	

314 315

284 285 286

316 317

318 319

320

4 **EXPERIMENTS** 

4.1 EXPERIMENTS ON PRE-TRAINING (TRAIN FROM SCRATCH)

NLP Transformer (GPT-2): We pre-train sparse NLP transformers by following the training settings
 in Zhao et al. (2023). Specifically, we train sparse GPT-2 Radford et al. (2019) models from scratch
 on the WikiText-103 dataset Merity et al. (2016), using the AdamW optimizer with an initial learning
 rate of 0.001 and a batch size of 512.

340 341

342

343

354

355

361

362

Table 1: Evaluation of training sparse/low-rank GPT-2 models from scratch on WikiText-103. Validation perplexity is provided, along with memory estimates for total parameters and optimizer states in FP16 format.

328	Method	PPL	Training Time	Training Epochs	# Params. (M)	# Mem. (MB)	Inference Throughput (samples/s)
220	GPT2-Small (Baseline)	18.5	24.5h, 8×V100	100	124.4	711.8	42.64
330	In-Rank Zhao et al. (2023)	18.9	22.2h, 8×V100	100	91.2	521.9	-
331	Monarch Dao et al. (2022)	20.7	-	100	72	412.0	-
332	Galore Zhao et al. (2024b)	18.3	22.6h, 4×V100	50	124.4	504.9	42.64
333	TranSpa	18.1	19.6h, 4×V100	50	88.2	504.7	54.79
334	TranSpa	20.7	21.9h, 4×V100	60	71.3	408.0	61.16
335	GPT2-Medium (Baseline)	19.5	60.5h, 8×V100	100	354.8	2030.2	19.74
000	In-Rank Zhao et al. (2023)	19.9	48.6h, 8×V100	100	223.0	1276.0	-
336	Monarch Dao et al. (2022)	20.3	-	100	165	994.1	-
337	Galore Zhao et al. (2024b)	19.6	60.4h, 4×V100	50	354.8	1260.9	19.74
338	TranSpa	19.0	52.1h, 4×V100	50	219.9	1258.3	24.48
339	TranSpa	19.8	56.8h, 4×V100	60	160.3	917.2	28.46

Table 2: Evaluation of training LLaMA-1B models from scratch on C4 dataset for 100K steps. Validation perplexity is provided, along with memory estimates for total parameters and optimizer states in BF16 format. The results for LoRA and ReLoRA are sourced from Zhao et al. (2024b).

Method	PPL	Training Time 8×A100	# Params. (M)	# Mem. (GB)	Inference Throughput (tokens/s)
Baseline	15.56	51.1h	1339.08	7.80	21786.49
LoRA Hu et al.	19.21	-	1339.08	6.17	21786.49
ReLoRA Lialin et al. (2023)	18.33	-	1339.08	6.17	21786.49
Galore Zhao et al. (2024b)	15.64	60.6h	1339.08	4.38	21786.49
TranSpa	15.60	40.3h	933.94	3.88	35211.27

**Large Language Model (LLaMA-1B):** For large language models, we follow the training settings in Zhao et al. (2024b) to train sparse LLaMA-1B from scratch on the C4 dataset. The training employs the Galore optimizer with an initial learning rate of 0.01 and a batch size of 512.

Vision Transformer (DeiT): We apply our approach to train sparse DeiT models Touvron et al. (2021)
from scratch on the ImageNet-1K dataset Deng et al. (2009). We use the same hyper-parameters as in
Chen et al. (2021a), which include the AdamW optimizer with an initial learning rate of 0.0005 and a
batch size of 512.

4.2 COMPARISON RESULTS

Table 1 compares our approach with the existing sparse and low-rank pre-training works. It is seen that 364 TranSpa achieves better performance than baseline and prior efforts using only half GPU resources 365 and less training time. Specifically, our approach can train sparse GPT2-Small and GPT2-Medium 366 models from scratch, with  $1.4 \times$  and  $1.6 \times$  model size reduction, respectively; and meanwhile, it 367 brings 0.4 and 0.6 lower perplexity over the baseline models. Compared with Dao et al. (2022), 368 TranSpa achieves the same or lower perplexity with  $1.6 \times$  training speedup (measured using the 369 number of epochs). Compared to Galore Zhao et al. (2024b), a method designed for low-rank compression of optimizer memory, TranSpa delivers better results with faster training speeds under 370 the same memory consumption and provides  $1.2 \times$  training speedup and  $1.4 \times$  inference acceleration. 371

Table 2 presents the pre-training results for LLaMA-1B. Our approach reduces GPU memory usage by 50%, decreases training time by 21%, and achieves a  $1.6 \times$  speedup in inference throughput without any performance loss. Compared to Galore, our approach reduces memory usage by 10%, increases training speed by  $1.5 \times$ , and boosts inference speed by  $1.6 \times$ . The results for LoRA Hu et al. and ReLoRA Lialin et al. (2023) are sourced from Zhao et al. (2024b). For low-rank methods, LoRA Hu et al. (2021) fine-tunes pre-trained models using low-rank adaptors:  $W = W_0 + BA$ , where  $W_0$  is the fixed initial weights and BA is a learnable low-rank adaptor. For pre-training,  $W_0$  is the full-rank

7

379	Table 3: Results of training sparse DeiT models from scratch on ImageNet-1k. Results for SSP-Tiny
380	and SSP-Small are sourced from Chen et al. (2021a). Inference speedup is measured on Nvidia RTX
381	30 <u>9</u> 0 GPU.

Method	# Params. (M)	# Mem. (MB)	FLOPs Saving (%)	Top-1 Acc. (%)	Inference Speedup
DeiT-Tiny (Baseline)	5.72	32.73	-	71.80	-
SSP-Tiny Bartoldson et al. (2020)	4.21	24.09	23.69	68.59	1.12×
S <sup>2</sup> ViTE-Tiny Chen et al. (2021a)	4.21	24.09	23.69	70.12	1.12×
TranSpa	3.95	22.60	32.01	70.92	1.20×
DeiT-Small (Baseline)	22.10	126.46	-	79.78	-
SSP-Small Bartoldson et al. (2020)	14.60	83.54	33.13	77.74	1.29×
S <sup>2</sup> ViTE-Small Chen et al. (2021a)	14.60	83.54	33.13	79.22	1.29×
TranSpa	13.98	79.99	37.30	79.57	<b>1.29</b> ×

Table 4: Perplexity of compressed LLaMA2-7B on WikiText-2 with different target model sizes.
SVD-LLM Wang et al. (2024) and SliceGPT Ashkboos et al. (2024) are low-rank based methods.
LLM Surgeon van der Ouderaa et al. (2023) is a pruning method, K-OBD van der Ouderaa et al. (2023), as a baseline comparison method, uses Kronecker-factored curvature and only prunes without updating the remaining weights.

Method		K-OBD	SVD-LLM	LLM Surgeon	SliceGPT	TranSpa
Training Time		16h58m, H100	15m, A100	17h08m, H100	1h07m, H100	1h41m, A100
PPL @ Target Size	80% 70% 60% 50%	9.14 15.43 28.03 46.64	7.94 9.56 13.11 23.97	6.18 7.83 10.39 15.38	6.64 8.12 -	6.36 <b>7.66</b> <b>10.24</b> <b>14.02</b>

initialization matrix. ReLoRA Lialin et al. (2023) is a variant of LoRA designed for pre-training. It periodically merges BA into W and reinitializes BA with a reset on optimizer states and learning rate. TranSpa surpasses these low-rank methods, reducing PPL by 3.6 and 2.7, respectively, while decreasing memory usage by 37% and offering  $1.6 \times$  acceleration in inference throughput.

Table 3 lists the performance results of various sparse vision transformer pre-training methods.
TranSpa achieves a 0.8% increase in top-1 accuracy for DeiT-Tiny and a 0.35% increase for DeiT-Small over state-of-the-art solutions, along with greater model size reduction.

4.3 EXPERIMENTS ON FINE-TUNING LLAMA2-7B (STRUCTURED PRUNING)

Experimental Setting. We evaluate the pruning performance of TranSpa on the pre-trained
LaMA2 Touvron et al. (2023) models. We use WikiText-2 Merity et al. (2016) as the calibration dataset and evaluate the perplexity of the pruned model. We follow the same training settings
adopted in van der Ouderaa et al. (2023), use 128 sequences with a sequence length of 2048 tokens
from the training data set, and evaluate perplexity on the standard test split. Additionally, we also
evaluate the performance of the pruned models on downstream zero-shot tasks.

Comparison Results. Table 4 presents a comparison of the perplexity performance between TranSpa and existing LLM pruning and low-rank factorization methods applied to LLaMA2-7B. Our approach consistently achieves lower perplexity across various target model size configurations compared to previous works. Additionally, Table 5 illustrates the zero-shot task performance of the pruned LLaMA2-7B model. In comparison to SliceGPT Ashkboos et al. (2024), our method demonstrates improved results across different target model sizes, indicating its effectiveness.

4.4 DISCUSSION & ANALYSIS

**Selection of** K and  $\theta$ . As described in Algorithm 1, K determines the percentage of  $W_{couple}$ 's that will be sparsified in each epoch, and threshold  $\theta$  impacts the amount of to-be-removed coupled

	Target Size	PIQA	WinoGrande	HellaSwag	ARC-e	ARC-c	Avg.
LLaMA2-7B	100%	79.11	69.06	75.99	74.58	46.25	69.00
TranSpa	80%	73.78	61.48	67.79	61.62	39.42	60.82
	75%	71.93	60.69	64.69	54.38	35.15	57.37
	70%	70.18	59.98	60.00	49.16	34.22	54.71
SliceGPT	80%	69.42	65.11	59.04	59.76	37.54	58.18
	75%	66.87	63.38	54.16	58.46	34.56	55.48
	70%	63.55	61.33	49.62	51.77	31.23	51.50

Table 5: Comparison of downstream zero-shot task performance of LLaMA2-7B model when trained



Figure 3: Pre-training sparse DeiT-Small model on CIFAR-10 dataset with various K and  $\theta$ . The target model size is 11M.



Figure 4: Fine-tuning Pre-trained sparse BERT-Base model on SQuAD benchmark with various Kand  $\theta$ . The target model size is 72.6M.

row/column pairs within those unimportant  $W_{couple}$ 's. As shown in Fig. 3, when applying TranSpa at the pre-training stage, larger  $\theta$  brings better training performance, while the model is less sensitive to the change of K. On the other hand, Fig. 4 shows that it is better to use smaller K at the fine-tuning stage, while the selection of  $\theta$  is less significant in this scenario. We hypothesize that such difference might be due to the existence of a pre-trained model in the fine-tuning process since a pre-trained model typically has more diverse distribution of  $W_{couple}$  than the model being sparsely trained from random initialization, making the identification of unimportant  $W_{couple}$  more effective. Therefore, considering K cannot be too small (otherwise, it is challenging to meet the model size budget (see the change of parameters (dashed curve) in Fig. 4), we set K = 30% and  $\theta = 90\%$  in our experiments. 





Figure 5: Pre-training DeiT-Small model on CIFAR-10 dataset using different  $W_{couple}$  compression methods. All the methods remove the same number of parameters within the same  $W_{couple}$ 's in each epoch.

solution, demonstrating the importance of removing the row/column of  $W_1$  and  $W_2$  in a coupled way. Notice that though our approach is inspired by tSVD, it achieves better performance than directly applying tSVD on  $W_{couple}$ . This is because SVD not only changes the structure of the model but also alters the numerical distribution of the original model weights. Consequently, the optimizers that keep track of moment information, e.g., Adam Kingma & Ba (2014), cannot work well since they will not be able to use previously accumulated information, thereby affecting the model performance.

#### CONCLUSION 5

511 512 513

514

515

516

517 518

519 520

521

522

523

524

525

526 527

528

529

530

531

486

487

488

489 490

491

492 493

494

495

496

497

498

499

500

501

502

504

505

506

507

508

509 510

> In this paper, we propose TranSpa, an efficient structured sparse training approach for transformers. By estimating the coupled weight matrix-wise importance and removing the coupled row/column pair during training, TranSpa brings a significant reduction in training costs and model complexity with preserving high task performance. Experiments across various transformer models demonstrate the superior performance of TranSpa in both pre-training and fine-tuning scenarios.

### REFERENCES

- Yongqi An, Xu Zhao, Tao Yu, Ming Tang, and Jinqiao Wang. Fluctuation-based adaptive structured pruning for large language models. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 10865-10873, 2024.
- Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James Hensman. Slicegpt: Compress large language models by deleting rows and columns. arXiv preprint arXiv:2401.15024, 2024.
- Brian Bartoldson, Ari Morcos, Adrian Barbu, and Gordon Erlebacher. The generalization-stability tradeoff in neural network pruning. Advances in Neural Information Processing Systems, 33: 20852-20864, 2020.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In European conference on computer 532 vision, pp. 213-229. Springer, 2020. 533
- 534 Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. arXiv preprint 536 arXiv:1708.00055, 2017.
- Tianlong Chen, Yu Cheng, Zhe Gan, Lu Yuan, Lei Zhang, and Zhangyang Wang. Chasing sparsity 538 in vision transformers: An end-to-end exploration. Advances in Neural Information Processing Systems, 34:19974-19988, 2021a.

540 541 542	Tianyi Chen, Tianyu Ding, Badal Yadav, Ilya Zharkov, and Luming Liang. Lorashear: Efficient large language model structured pruning and knowledge recovery. <i>arXiv preprint arXiv:2310.18356</i> , 2023a.
543	
544	Tianyi Chen, Luming Liang, Tianyu DING, Zhihui Zhu, and Ilya Zharkov. OTOv2: Automatic,
545	generic, user-friendly. In The Eleventh International Conference on Learning Representations,
546	2023b. URL https://openreview.net/forum?id=7ynoX1ojPMt.
547	Tianyi Chan Luming Liang Tianyu Ding Zhihui Zhu, and Ilya Zharkov, Otov?: Automatic, ganaric
548 549	user-friendly. arXiv preprint arXiv:2303.06862, 2023c.
550	Xiaodong Chen, Yuxuan Hu, and Jing Zhang. Compressing large language models by streamlining
551 552	the unimportant layer. arXiv preprint arXiv:2403.19135, 2024.
553 554 555 556	Xiaohan Chen, Yu Cheng, Shuohang Wang, Zhe Gan, Zhangyang Wang, and Jingjing Liu. Earlybert: Efficient bert training via early-bird lottery tickets. In <i>Proceedings of the 59th Annual Meeting</i> of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 2195–2207, 2021b.
557	Zihan Chen, Hongbo Zhang, Xiaoji Zhang, and Leqi Zhao. Quora question pairs, 2018.
558	Tri Dao, Beidi Chen, Nimit S Sohoni, Ariun Desai, Michael Poli, Jessica Grogan, Alexander Liu,
559	Aniruddh Rao, Atri Rudra, and Christopher Ré. Monarch: Expressive structured matrices for
560	efficient and accurate training. In International Conference on Machine Learning, pp. 4690–4721.
561	PMLR, 2022.
562	
563	Iristan Deleu and Yoshua Bengio. Structured sparsity inducing adaptive optimizers for deep learning.
504	
565	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
500	hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition,
507	pp. 248–255. Ieee, 2009.
500	La L. D. 1'. Mar W. 'Char Kasta La sa L. V. 's' Total Data Data' 's a film
509	Jacob Devlin, Ming-wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv praprint arXiv:1810.04805, 2018
570	oldifectional transformers for language understanding. <i>urxiv preprint urxiv</i> .1610.04805, 2018.
572	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
573	deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
574	Thamar Solorio (eds.), Proceedings of the 2019 Conference of the North American Chapter of the
575	Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and
576	<i>Short Papers</i> ), pp. 41/1–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics doi: 10.18653/v1/N19-1423 URL https://aclanthology.org/N19-1423
577	
578	Bill Dolan, Chris Brockett, and Chris Quirk. Microsoft research paraphrase corpus. <i>Retrieved March</i> ,
579	29(2008):63, 2005.
580	Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: a no-recurrence sequence-to-sequence
581	model for speech recognition. In 2018 IEEE international conference on acoustics, speech and
582	signal processing (ICASSP), pp. 5884–5888. IEEE, 2018.
583	
584	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
585	Untertainer, Mostala Dengnani, Matthias Minderer, Georg Heigold, Sylvan Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at cools.
586	arXiv:2010.11929.2020
587	withilo10,11/2/, 2020.
588	Abhimanyu Dubey, Moitreya Chatterjee, and Narendra Ahuja. Coreset-based neural network com-
589	pression. In Proceedings of the European Conference on Computer Vision (ECCV), pp. 454-470,
590	2018.
500	Utku Evci Trevor Gale Jacob Menick Dable Samuel Castro, and Erich Elsen. Discing the letteru
593	Making all tickets winners. In <i>International Conference on Machine Learning</i> , pp. 2943–2952. PMLR, 2020.

594 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. arXiv preprint arXiv:1803.03635, 2018. 596 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in 597 one-shot. In International Conference on Machine Learning, pp. 10323–10337. PMLR, 2023. 598 Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo 600 Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for 601 speech recognition. arXiv preprint arXiv:2005.08100, 2020. 602 Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks 603 with pruning, trained quantization and huffman coding. arXiv preprint arXiv:1510.00149, 2015a. 604 605 Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for 606 efficient neural network. Advances in neural information processing systems, 28, 2015b. 607 608 Yang He, Ping Liu, Ziwei Wang, Zhilan Hu, and Yi Yang. Filter pruning via geometric median for deep convolutional neural networks acceleration. In Proceedings of the IEEE/CVF conference on 609 computer vision and pattern recognition, pp. 4340–4349, 2019. 610 611 Zejiang Hou, Minghai Qin, Fei Sun, Xiaolong Ma, Kun Yuan, Yi Xu, Yen-Kuang Chen, Rong 612 Jin, Yuan Xie, and Sun-Yuan Kung. Chex: Channel exploration for cnn model compression. 613 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 614 12287-12298, 2022. 615 Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, 616 and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked 617 prediction of hidden units. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 618 29:3451-3460, 2021. 619 620 Yen-Chang Hsu, Ting Hua, Sungen Chang, Qian Lou, Yilin Shen, and Hongxia Jin. Language model 621 compression with weighted low-rank factorization. arXiv preprint arXiv:2207.00112, 2022. 622 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, 623 et al. Lora: Low-rank adaptation of large language models. In International Conference on 624 Learning Representations. 625 626 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In International Conference on 627 Learning Representations, 2021. 628 629 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint 630 arXiv:1412.6980, 2014. 631 632 Frederik Kunstner, Philipp Hennig, and Lukas Balles. Limitations of the empirical fisher approx-633 imation for natural gradient descent. Advances in neural information processing systems, 32, 2019. 634 635 Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. Relora: High-636 rank training through low-rank updates. In The Twelfth International Conference on Learning 637 Representations, 2023. 638 639 Mingbao Lin, Rongrong Ji, Yan Wang, Yichen Zhang, Baochang Zhang, Yonghong Tian, and Ling Shao. Hrank: Filter pruning using high-rank feature map. In Proceedings of the IEEE/CVF 640 conference on computer vision and pattern recognition, pp. 1529–1538, 2020. 641 642 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 643 Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the 644 IEEE/CVF international conference on computer vision, pp. 10012–10022, 2021. 645 Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural 646 network compression. In Proceedings of the IEEE international conference on computer vision, 647 pp. 5058–5066, 2017.

648	Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large
649 650	language models. Advances in neural information processing systems, 36:21702–21720, 2023.
651 652	Srikanth Malla, Joon Hee Choi, and Chiho Choi. Copal: Continual pruning in large language generative models. <i>arXiv preprint arXiv:2405.02347</i> , 2024.
653	V' M. M' V O' 71 D ' ' We Her I' V '' I V' ' I I'
654	Ain Men, Mingyu Au, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Alanpei Han, and Weipeng Chen. Shortant: Layers in large language models are more redundant than you expect
655 656	arXiv preprint arXiv:2403.03853, 2024.
657 658	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. <i>arXiv preprint arXiv:1609.07843</i> , 2016.
659	N., W.'.M., Zhata' Wesse Characta' D' 1 C' - Marco Cata and Tracta Mathematica
660 661	matrix arithmetic-only bert inference by eliminating complex non-linear functions. In <i>The Eleventh</i>
662	International Conference on Learning Representations, 2022.
663 664 665	Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. <i>Nature communications</i> , 9(1):2383, 2018.
666	
667 668	Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. In <i>International Conference on Machine Learning</i> , pp. 4646–4655, pMLP, 2010
669	4040–4033. PMLR, 2019.
670	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, and et al. Gpt-4 technical report, 2023.
671	
672 673	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
674 675	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. <i>arXiv preprint arXiv:1606.05250</i> , 2016.
676	
677 678	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> , 2019.
679	Yang Sui, Miao Yin, Yi Xie, Huy Phan, Saman Aliari Zonouz, and Bo Yuan. Chip: Channel
680 681	independence-based pruning for compact neural networks. Advances in Neural Information Processing Systems, 34:24604–24616, 2021.
682	
683 684	large language models. <i>arXiv preprint arXiv:2306.11695</i> , 2023.
685	Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks.
686	Advances in Neural Information Processing Systems, 34:24193–24205, 2021.
687	
688	Chong Min John Tan and Mehul Motani. Dropnet: Reducing neural network complexity via iterative
689	prunnig. In international Conjerence on Machine Learning, pp. 5550–5500. FMLR, 2020.
690 601	Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve
602	Jegou. Training data-efficient image transformers & distillation through attention. In International
693	Conference on Machine Learning, volume 139, pp. 10347–10357, July 2021.
694	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amiad Almahairi, Yasmine Babaei, Nikolav
695	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
696	and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023.
697	Tycho FA yan dar Audaraa Markus Nagal Mart Van Baalan Vuki M Asano, and Tiiman Plankayoort
698 699	The llm surgeon. <i>arXiv preprint arXiv:2312.17244</i> , 2023.
700	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez. Łukasz
701	Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.

702	Michael E Wall, Andreas Rechtsteiner, and Luis M Rocha. Singular value decomposition and
703	principal component analysis. In A practical approach to microarray data analysis, pp. 91–109.
704	Springer, 2003.
705	

- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao.
   Learning deep transformer models for machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1810–1822, 2019a.
- Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong, and Furu Wei. Minilmv2: Multi-head self-attention relation distillation for compressing pretrained transformers. *arXiv preprint arXiv:2012.15828*, 2020.
- Wenxiao Wang, Cong Fu, Jishun Guo, Deng Cai, and Xiaofei He. Cop: Customized deep model compression via regularized correlation-based filter-level pruning. *arXiv preprint arXiv:1906.10337*, 2019b.
- Xin Wang, Yu Zheng, Zhongwei Wan, and Mi Zhang. Svd-llm: Truncation-aware singular value decomposition for large language model compression. *arXiv preprint arXiv:2403.07378*, 2024.
- Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments.
   *Transactions of the Association for Computational Linguistics*, 7:625–641, 2019.
- Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for
   sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared Ilama: Accelerating language model pre-training via structured pruning. *arXiv preprint arXiv:2310.06694*, 2023.
- Yifei Yang, Zouying Cao, and Hai Zhao. Laco: Large language model pruning via layer collapse.
   *arXiv preprint arXiv:2402.11187*, 2024.
- Shixing Yu, Tianlong Chen, Jiayi Shen, Huan Yuan, Jianchao Tan, Sen Yang, Ji Liu, and Zhangyang Wang. Unified visual transformer compression. *arXiv preprint arXiv:2203.08243*, 2022.
- Geng Yuan, Xiaolong Ma, Wei Niu, Zhengang Li, Zhenglun Kong, Ning Liu, Yifan Gong, Zheng
  Zhan, Chaoyang He, Qing Jin, et al. Mest: Accurate and fast memory-economic sparse training
  framework on the edge. *Advances in Neural Information Processing Systems*, 34:20838–20850,
  2021.
- Edouard Yvinec, Arnaud Dapogny, Matthieu Cord, and Kevin Bailly. Red: Looking for redundancies
   for data-freestructured compression of deep neural networks. *Advances in Neural Information Processing Systems*, 34:20863–20873, 2021.
- Yuxin Zhang, Lirui Zhao, Mingbao Lin, Yunyun Sun, Yiwu Yao, Xingjia Han, Jared Tanner, Shiwei Liu, and Rongrong Ji. Dynamic sparse no training: Training-free fine-tuning for sparse llms. *arXiv* preprint arXiv:2310.08915, 2023.
- Bowen Zhao, Hannaneh Hajishirzi, and Qingqing Cao. Apt: Adaptive pruning and tuning pretrained language models for efficient training and inference. *arXiv preprint arXiv:2401.12200*, 2024a.
- Jiawei Zhao, Yifei Zhang, Beidi Chen, Florian Schäfer, and Anima Anandkumar. Inrank: Incremental
   low-rank learning. 2023.
- Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong
  Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint arXiv:2403.03507*, 2024b.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and
   Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching
   movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pp. 19–27, 2015.

## 756 A APPENDIX

#### 758 A.1 FINE-TUNING BERT MODEL ON GLUE AND SQUAD1.1 DATASET 759

Experimental Setting. We evaluate the pruning performance of TranSpa on the pre-trained BERT De-760 vlin et al. (2018) models. For experiments on BERT, the model is pruned and evaluated on 761 GLUE Wang et al. (2018) and SQuAD Rajpurkar et al. (2016) benchmarks. When performing 762 fine-tuning on the GLUE benchmark, the learning rate and weight decay are set as  $5 \times 10^{-5}$  and 763 0, respectively, with 3 training epochs and batch size of 32. The maximum sequence length is 764 128. For fine-tuning on the SQuAD benchmark, the learning rate is set to  $3 \times 10^{-5}$  at the 3-epoch 765 fine-tuning stage, and the batch size, maximum sequence length, and document stride are 12, 320, 766 128, respectively. Both training and testing are conducted on a single Nvidia GeForce RTX 3090. 767 For single-sentence tasks, we include the CoLA Warstadt et al. (2019) task, evaluated based on 768 Matthew's correlation, and the SST-2 task, assessed using classification accuracy. For sentence 769 similarity tasks, we consider three distinct evaluations: MRPC Dolan et al. (2005) evaluated using 770 F-1 score, STS-B Cer et al. (2017) assessed through Pearson-Spearman correlation, and QQP Chen 771 et al. (2018) measured in F-1 score. The evaluation also covers three natural language inference tasks: MNLI Williams et al. (2017), measured in classification accuracy with the average of the matched 772 and mismatched subsets, and QNLI Rajpurkar et al. (2016), assessed in terms of accuracy. For the 773 reading comprehension task, SQuAD consists of questions generated by crowd workers based on 774 a collection of Wikipedia articles. Each question requires extracting an answer, in the form of a 775 segment of text or span, from the corresponding reading passage. 776

777

778 Table 6: Experimental results on GLUE. 'Avg." means the average of GLUE tasks. Average training 779 and inference times are calculated by averaging the time consumption across all tasks, "G" denotes Generic Knowledge Distillation, where models are compressed, and then knowledge distillation is conducted on the datasets from English Wikipedia and Toronto Book Corpus Zhu et al. (2015) 781 before fine-tuning on various task datasets. "G+TS" indicates that a task-specific teacher model is 782 still utilized during fine-tuning for further knowledge distillation after generic knowledge distillation. 783 Training times for DitilBERT and MA-DistilBERT are sourced from the Sanh et al. (2019): Ming 784 et al. (2022), respectively. Average training and inference times are computed by averaging the time 785 consumption across all tasks. 78

Method	MRPC F-1	QNLI Acc.	MNLI Acc.	STS-B Corr.	SST-2 Acc.	QQP F-1	CoLA Corr.	Avg.	# Params. (M)	Knowledge Distllation	Fine-tuning Epochs	Avg. Training Time	Avg. Inference Time (ms/batch)
BERT-Base (Baseline)	87.4	91.3	84.7	88.5	93	87.8	56.2	84.1	109.5	No	3	1402.7s, 1×RTX3090	118.1
TranSpa	88.2	90.1	83.2	86.5	91.6	87.2	52.9	82.8	75.1	No	2	821.7s, 1×RTX3090	84.3
TranSpa	89.5	89.4	82.7	86.7	92.7	87.3	51.6	82.8	66.1	No	3	1210.7s, 1×RTX3090	70.3
FWSVD Hsu et al. (2022)	88.0	89.5	83.0	87	91.2	87.6	49.4	82.2	66.5	No	4	1548.1s, 1×RTX3090	81.3
FISH Sung et al. (2021)	86.2	93.3	85.3	85.7	94.0	79.3	56.4	82.9	335.8	No	7	-	376.6
EarlyBERT Chen et al. (2021b)	-	89.2	81.8	-	90.7	-	-	-	77.0	No	2.2	-	-
DistilBERT Sanh et al. (2019)	88.7	89.3	82.2	86.1	90.4	86.7	49.8	81.9	67.0	Yes	3	3.75d, 8×V100	-
MiniLMv2 Wang et al. (2020)	89.1	90.6	84	88.1	91.4	86.7	43.3	81.9	67.0	Yes	3	-	-
MA-DistilBERT Ming et al. (2022	) 86.7	88.4	82.3	87.5	90.8	-	51.3	-	67.2	Yes	5-10	3.5d, 1×RTX3090	-

Table 7: Experimental results on SQuAD1.1 for pruning BERT-Base model.

Method	# Params. (M)	Exact (%)	F1-score (%)	Training epochs	Training Time (s)	Inference Speedup
BERT-Base (Baseline)	108.9	81.00	88.30	3	3,834	-
ProxSSI Deleu & Bengio (2021)	90.8	72.30	82.00	3	-	1×
OTOv2 Chen et al. (2023b)	87.2	79.40	87.30	3	-	1.2×
OTOv2	59.9	74.60	83.80	3	-	1.4×
EarlyBERT Chen et al. (2021b)	77.0	-	86.13	2.2	-	-
TranSpa	72.3	78.59	86.53	2	2,138	1.5×
TranSpa	86.7	79.92	87.59	3	3,618	1.2×
TranSpa	58.6	75.97	84.64	3	2,939	$1.7 \times$

801 802

794

803

Comparison Results. Table 6 shows the performance results of compressing the BERT-Base model
 using different methods. It is seen that TranSpa consistently outperforms existing structured pruning
 and knowledge distillation solutions across a set of tasks. Meanwhile, because it does not require
 costly knowledge distillation at the fine-tuning stage, our approach enjoys a very fast fine-tuning
 speed. Table 7 shows the results on the SQuAD benchmark. It is seen that TranSpa consistently
 outperforms the existing structured pruning approaches with significantly higher exact and F-1 scores
 (at least 0.52% and 0.29% increase) and smaller model size.