

Weight-Inherited Distillation for Task-Agnostic BERT Compression

Anonymous ACL submission

Abstract

Knowledge Distillation (KD) is a predominant approach for BERT compression. Previous KD-based methods focus on designing extra alignment losses for student model to mimic the behavior of teacher model. These methods transfer the knowledge in an indirect way. In this paper, we propose a novel Weight-Inherited Distillation (WID), which directly transfers knowledge from the teacher. WID does not require any additional alignment loss and trains a compact student by inheriting the weights, showing a new perspective of knowledge distillation. Specifically, we design the compactors as mappings and then compress the weights via structural re-parameterization. Experimental results on the GLUE and SQuAD benchmarks show that WID outperforms previous state-of-the-art KD-based baselines. Further analysis indicates that WID can also learn the attention patterns from the teacher model without any alignment loss on attention distributions.

1 Introduction

Transformer-based Pre-trained Language Models (PLMs), such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), XLNET (Yang et al., 2019), have achieved great success in many Natural Language Process (NLP) tasks. These models are pre-trained on massive corpus via self-supervised tasks to learn contextualized text representations. However, PLMs have high costs in terms of storage, memory, and computation time, which brings challenges to online service in real-life applications. Therefore, it is crucial and feasible to compress PLMs while maintaining their performance.

Knowledge Distillation (KD), which trains a compact student model by mimicking the behavior of a teacher model, is a predominant method for PLM compression. There are two settings for KD in BERT compression: task-specific, which first fine-tune the teacher PLMs on specific tasks and then perform distillation, and task-agnostic,

Approach	Alignment Loss		Hard Loss	Task-Agnostic
	Logit	Feature		
DistilBERT	✓	✓	✓	✓
TinyBERT (GD)	✓	✓	✗	✓
PKD	✓	✓	✓	✗
MiniLM	✗	✓	✗	✓
MobileBERT	✓	✓	✓	✓
WID (ours)	✗	✗	✓	✓

Table 1: Comparison with previous state-of-the-art distillation methods. **Logit** and **Feature** denote whether logit-based loss and feature-based loss are used for distillation. To the best of our knowledge, WID is the first distillation method without any alignment loss and directly transfers the knowledge by weight inheritance.

which distill PLMs in pre-training stage. For task-agnostic distillation, the student model can be directly and generically fine-tuned on various downstream tasks (Wang et al., 2020; Sun et al., 2020). Hence, we conduct our weight-inherited distillation under task-agnostic setting.

Previous KD-based methods mainly focus on designing alignment losses to minimize the distance between the teacher model and the student model. We can further categorize these alignment losses into: logit-based, which measures the distance of logit distributions, and feature-based, which aims to align the intermediate features including token embeddings, hidden states, and self-attention distributions. However, adopting these alignment losses brings the following drawbacks: 1) selecting various loss functions and balancing the weights of each loss are laborious (Sun et al., 2019; Jiao et al., 2020); 2) some losses will restrict the architecture of the student model. For example, attention-based loss (Jiao et al., 2020; Wang et al., 2020; Sun et al., 2020) requires the student model to have the same attention heads as the teacher.

In this work, we propose Weight-Inherited Distillation (WID), which does not require any additional alignment loss and trains the student by directly inheriting the weights from teacher. Inspired by

structural re-parameterization in CNN compression (Ding et al., 2021), we design row compactors and column compactors and view them as mappings to compress the weights by row and column, respectively. Figure 1 shows the process of compressing a linear layer by WID. All compactors are initialized as identity matrices, thus the re-parameterized teacher model produces identical outputs as the original teacher. We train the re-parameterized teacher model on the pre-training task and add weight penalty to compactors simultaneously. After training, we compress the compactors to desired sizes and merge these compactors and original weights into compact one. As shown in Table 1, WID is the only method for task-agnostic distillation without any alignment loss.

We conduct extensive experiments on downstream NLP tasks, including the GLUE and SQuAD benchmarks. Experimental results demonstrate that WID outperforms traditional KD-based baselines. Further analysis shows that WID can also learn knowledge such as self-attention patterns from the teacher model.

Our contributions can be summarized as follows:

- We propose Weight-Inherited Distillation (WID), revealing a new pathway to knowledge distillation by directly inheriting the weights via structural re-parameterization.
- We conduct WID for task-agnostic BERT compression. Experiments on the GLUE and SQuAD benchmark datasets demonstrate the effectiveness of WID for model compression.
- We perform further analyses on how to get better performance in BERT compression. Moreover, we find that WID can also learn attention patterns from the teacher.

2 Preliminaries

In this section, we present a brief introduction to the transformer. Moreover, we also present existing KD-based methods for transformer networks.

2.1 Embedding Layer

In BERT (Devlin et al., 2019), the input texts are tokenized to tokens by WordPiece (Wu et al., 2016). The representations ($\{\mathbf{x}_i\}_{i=1}^{|\mathbf{x}|}$) of input sequence are constructed by summing the corresponding token embedding, segment embedding, and position embedding. For the token embedding layer in

BERT, the weight is $W_T \in \mathbb{R}^{|V| \times d}$, where $|V|$ and d denote the size of the vocabulary and the hidden state vector.

2.2 Transformer Layer

Transformer layer is adopted to encode the contextual information of input texts. The input vector ($\{\mathbf{x}_i\}_{i=1}^{|\mathbf{x}|}$) are packed to $\mathbf{H}^0 = [\mathbf{x}_1, \dots, \mathbf{x}_{|\mathbf{x}|}]$. After that, the L -layer transformer computes the encoding vectors following:

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

The final output $\mathbf{H}^L = [h_1^L, \dots, h_{|\mathbf{x}|}^L] \in \mathbb{R}^{|\mathbf{x}| \times d}$ is employed as the contextualized representation of $\{\mathbf{x}_i\}_{i=1}^{|\mathbf{x}|}$. Each transformer layer consists of a multi-head self-attention (MHA) sub-layer and a feed-forward (FFN) sub-layer. In these two sub-layers, the residual connection (He et al., 2016) is employed, followed by layer normalization (Ba et al., 2016).

MHA For the l -th transformer layer with A attention heads, the output $\mathbf{O}_{l,a}$ of the attention head $a \in [1, A]$ is calculated as:

$$\begin{aligned} \mathbf{Q}_{l,a} &= \mathbf{H}^{l-1} \mathbf{W}_{l,a}^Q \\ \mathbf{K}_{l,a} &= \mathbf{H}^{l-1} \mathbf{W}_{l,a}^K \\ \mathbf{V}_{l,a} &= \mathbf{H}^{l-1} \mathbf{W}_{l,a}^V \end{aligned} \quad (2)$$

$$\mathbf{O}_{l,a} = \mathbf{A}_{l,a} \mathbf{V}_{l,a}, \mathbf{A}_{l,a} = \text{softmax}\left(\frac{\mathbf{Q}_{l,a} \mathbf{K}_{l,a}^T}{\sqrt{d_k}}\right) \quad (3)$$

where linear projection $\mathbf{W}_{l,a}^Q, \mathbf{W}_{l,a}^K, \mathbf{W}_{l,a}^V \in \mathbb{R}^{d \times d_k}$ and $d_k = \frac{d}{A}$ is the dimension of each head. The final output of MHA sub-layer is as follows:

$$\mathbf{O}_l = \text{LN}(\mathbf{H}^{l-1} + (\|_{a=1}^A \mathbf{O}_{l,a}) \mathbf{W}_l^O) \quad (4)$$

where $\mathbf{W}_l^O \in \mathbb{R}^{d \times d}$, LN is layer normalization and $\|$ denotes the concatenation operation.

FFN The l -th FFN sub-layer consists of an up projection and a down projection, parameterized by $\mathbf{W}_{l,u} \in \mathbb{R}^{d \times d_f}$, $\mathbf{W}_{l,d} \in \mathbb{R}^{d_f \times d}$, and corresponding bias $\mathbf{b}_{l,u} \in \mathbb{R}^{d_f}$, $\mathbf{b}_{l,d} \in \mathbb{R}^d$:

$$\text{FFN}(\mathbf{O}_l) = \text{gelu}(\mathbf{O}_l \mathbf{W}_{l,u} + \mathbf{b}_{l,u}) \mathbf{W}_{l,d} + \mathbf{b}_{l,d}. \quad (5)$$

Typically, $d_f = 4d$. Finally, we obtain the output of layer l by:

$$\mathbf{H}^l = \text{LN}(\mathbf{O}_l + \text{FFN}(\mathbf{O}_l)). \quad (6)$$

2.3 Knowledge Distillation

Knowledge Distillation (KD) aims to transfer the knowledge from teacher model T to compact student model S . The student model S is trained to mimic the behaviors of teacher model T via minimizing the distance between them. The object losses can be categorized into logit-based and feature-based.

For logit-based loss, the target is to minimize the logit distribution \mathbf{p}_s from student and \mathbf{p}_t from teacher, which can be formalized as:

$$\mathcal{L}_{logit} = \mathcal{H}_1(\mathbf{p}_s/\tau, \mathbf{p}_t/\tau), \quad (7)$$

where τ is the temperature and \mathcal{H}_1 is the cross-entropy loss or KL-divergence.

Feature-based loss aims to align the intermediate features between the teacher and the student by:

$$\mathcal{L}_{feature} = \mathcal{H}_2(f^S(x), f^T(x)), \quad (8)$$

where \mathcal{H}_2 is the loss function such as Mean Square Error (MSE) and $f(x)$ notes for the intermediate output including hidden state vector \mathbf{H} and attention distribution \mathbf{A} .

As shown in Table 1, logit-based and feature-based loss can be jointly employed for better distillation. However, balancing the weights of each loss is laborious. For example, the overall loss of PKD (Sun et al., 2019) is:

$$\mathcal{L} = (1 - \alpha)\mathcal{L}_{hard} + \alpha\mathcal{L}_{logit} + \beta\mathcal{L}_{feature}, \quad (9)$$

where \mathcal{L}_{hard} is the loss on target tasks and α and β are the hyper-parameters. They perform grid search over α and τ , where $\alpha \in \{0.2, 0.5, 0.7\}$ and $\tau \in \{5, 10, 20\}$. After that, they fix α and τ with the best performance and search $\beta \in \{10, 100, 500, 1000\}$.

Meanwhile, selecting various loss functions is also laborious. In PKD, $\mathcal{L}_{feature}$ is defined as the mean square loss between the normalized hidden states for each layer. DistilBERT (Sanh et al., 2019) adopts the cosine embedding loss for hidden states vectors. In TinyBERT (Jiao et al., 2020), they employ the mean square loss for self-attention distributions, embedding layer outputs, and hidden states.

3 Weight-Inherited Distillation

In this section, we propose a novel Weight-Inherited Distillation (WID) method for

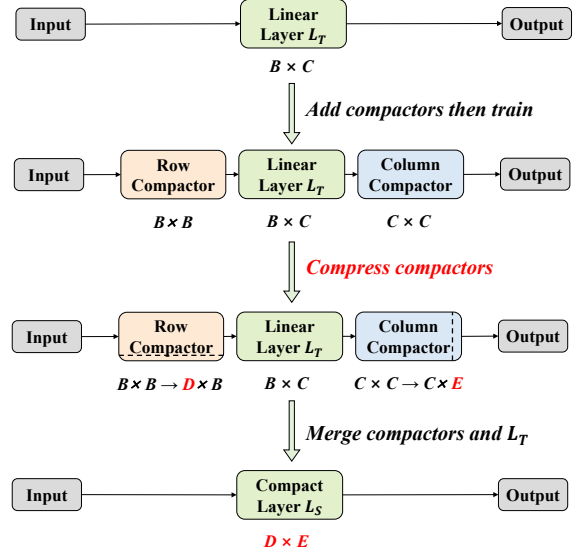


Figure 1: Overview of compressing linear layer L_T with weight $W^{L_T} \in \mathbb{R}^{B \times C}$ to compact linear layer L_S with weight $W^{L_S} \in \mathbb{R}^{D \times E}$ via WID. Both row compactor and column compactor are initialized as **identity matrices**. After training, we compress the compactors and merge them with original layer. All the linear layers in teacher model are compressed **simultaneously**.

transformer-based models without any alignment loss. The WID aims to directly leverage knowledge in weight and compress the teacher model by learning mappings for the compact student model.

3.1 Structural Re-parameterization

As mentioned in Section 2, the PLMs (e.g., BERT) consist of embedding layers and transformer layers. To compress the BERT, we have to learn a mapping from the larger weight in the teacher model to the compact one. In WID, we adopt the structural re-parameterization and design the row compactors and column compactors.

Figure 1 gives an example showing the process of compressing the original weight $\mathbf{W}^{L_T} \in \mathbb{R}^{B \times C}$ to compact weight $\mathbf{W}^{L_S} \in \mathbb{R}^{D \times E}$ adopting the row compactor and the column compactor. First, we insert the row compactor with weight $\mathbf{W}^{rc} \in \mathbb{R}^{B \times B}$ and the column compactor with weight $\mathbf{W}^{cc} \in \mathbb{R}^{C \times C}$ before and after the linear layer L_T from teacher model. All compactors are linear layers without bias and their weights are initialized as identity matrices. For an arbitrary input X , the re-parameterized teacher model produces identical outputs as the original, since

$$\mathbf{W}^{L_T} X = \mathbf{W}^{rc} \mathbf{W}^{L_T} \mathbf{W}^{cc} X. \quad (10)$$

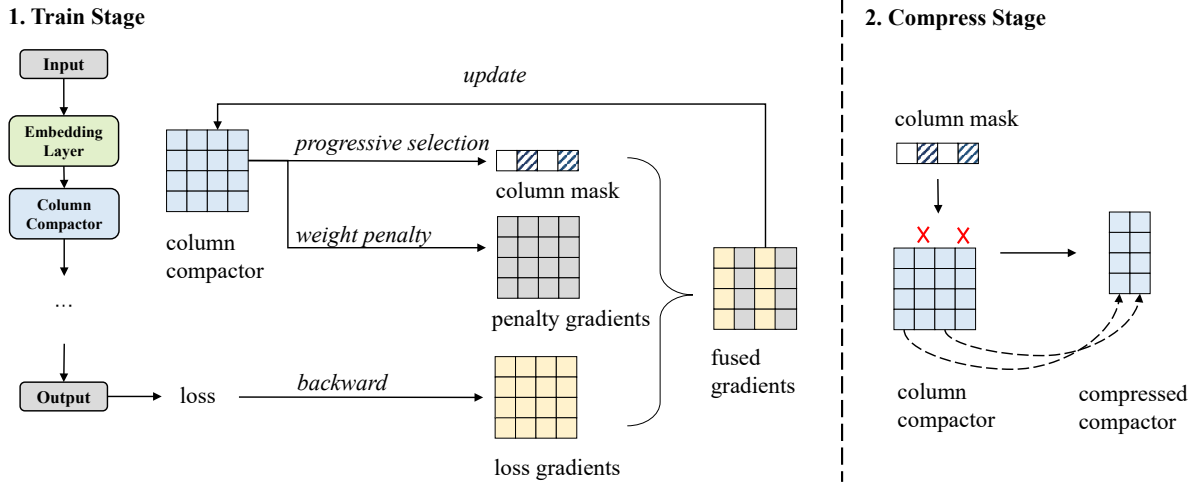


Figure 2: Training and compression for column compactor. During the training process, we add weight penalty gradients by column and progressively select the mask to fuse the penalty gradients and original loss gradients. For gradients fusion, we decouple penalty gradients and original loss gradients to avoid gradient competition. After training, we prune the column compactor following the column mask.

Second, we train the re-parameterized teacher model on the pre-training task. During training, we add the row penalty to row compactor and column penalty to column compactor. The goal is to maintain the performance of the teacher model and compress the compactor simultaneously. After training, the row compactor is compressed by pruning $B - D$ rows, and the column compactor is compressed by pruning $C - E$ columns. The objects are as follows:

$$\begin{aligned} \mathbf{W}^{rc} \in \mathbb{R}^{B \times B} &\rightarrow \mathbf{W}^{rc'} \in \mathbb{R}^{D \times B} \\ \mathbf{W}^{cc} \in \mathbb{R}^{C \times C} &\rightarrow \mathbf{W}^{cc'} \in \mathbb{R}^{C \times E}. \end{aligned} \quad (11)$$

More details can be found in Section 3.2. Final, we merge the compressed compactors $\mathbf{W}^{rc'}$, $\mathbf{W}^{cc'}$ and the original teacher layer \mathbf{W}^{L_T} to obtain the compact layer for the student following:

$$\mathbf{W}^{L_S} = \mathbf{W}^{rc'} \mathbf{W}^{L_T} \mathbf{W}^{cc'} \in \mathbb{R}^{D \times E} \quad (12)$$

For the weights to compress the row only, such as the output layer for MLM task with size $\mathbb{R}^{d \times |V|}$, we adopt the row compactor exclusively. Similarly, we employ the column compactor exclusively for the weights to compress the column only, such as the token embedding matrix $\mathbf{W}_T \in \mathbb{R}^{|V| \times d}$.

3.2 Compactor Compression

In WID, we design row compactors and column compactors and view them as mappings to compress the weights by row and column, respectively. Compared to directly learning these compactors, our key insight is to initialize these compactors

with identity matrices and compress them to the desired size progressively.

Figure 2 presents the training and compression process for the column compactor. Given the column compactor $\mathbf{W}^{cc} \in \mathbb{R}^{C \times C}$ and original gradients $g_{ori}^{cc} \in \mathbb{R}^{C \times C}$, the penalty gradients $g_{pen}^{cc} \in \mathbb{R}^{C \times C}$ are calculated as follows:

$$g_{pen}^{cc} = \frac{\mathbf{W}^{cc}}{\|\mathbf{W}^{cc}\|_p} \quad (13)$$

where $\|\mathbf{W}^{cc}\|_p$ denotes the p -norm cross each column. Based on the $\|\mathbf{W}^{cc}\|_p$, we pick top- k columns with lower norm value and set the corresponding value in our column mask $M = \{0, 1\}^C$ to be 1. For gradients fusion, we decouple penalty gradients and original loss gradients to avoid gradient competition. Thus, the original gradients g_{ori}^{cc} and the penalty gradients g_{pen}^{cc} are fused as follows:

$$g_{fused}^{cc}[:, i] = \begin{cases} g_{pen}^{cc}[:, i], & \text{if } M[i] = 1 \\ g_{ori}^{cc}[:, i], & \text{if } M[i] = 0 \end{cases} \quad (14)$$

where $0 \leq i \leq C$. The fused gradients g_{fused}^{cc} are employed to update the column compactor by optimizer. After training, we prune the column compactor by column mask:

$$\mathbf{W}^{cc'} = \mathbf{W}^{cc}[:, i], \text{ where } M[i] = 1. \quad (15)$$

Moreover, the processing is similar for row compactors. We calculate $\|\mathbf{W}^{rc}\|_p$ for each row and select the top- k rows with the lower norm value.

Algorithm 1 Weight-Inherited Distillation

Input: teacher model \mathcal{T} with width d_t
Params: k : number of rows/columns to compress, N : steps to increase k , d : increment for k each time
Output: student model \mathcal{S} with width d_s

- 1: Add compactors for \mathcal{T} to construct the re-parameterized teacher model $\hat{\mathcal{T}}$. Initialize the weights for compactors as identity matrices.
- 2: $k \leftarrow 0$; $M \leftarrow []$
- 3: **for** $i = 0$ to max training steps **do**
- 4: Forward a batch through $\hat{\mathcal{T}}$, derive the gradients g_{ori} for each compactor
- 5: **if** $i\%N == 0$ & $k < d_t - d_s$ **then**
- 6: Calculate p-norm values
- 7: Select the top- k row/column with the lower norm to get M
- 8: Get penalty gradients g_{pen} following Eq. 13
- 9: $g_{fused} \leftarrow f(g_{ori}, g_{pen}, M)$ following Eq. 14
- 10: $k \leftarrow k + d$
- 11: **end if**
- 12: Update the compactors with corresponding g_{fused} and original layers with g_{ori}
- 13: **end for**
- 14: Compress the compactors following Eq. 15
- 15: Merge the compactors and original layers following Eq. 12 to get compact layers for \mathcal{S}
- 16: **return** \mathcal{S}

For stability and better performance, we compress the compactors progressively. Specifically, we increase k for some steps until reaching the desired size during the training stage. More details are shown in Algorithm 1.

4 Experiments

4.1 Task-Agnostic Distillation

We employ the uncased version of BERT_{base} as our teacher model¹. BERT_{base} (Devlin et al., 2019) is a 12-layer transformer model ($d=768$, $A=12$, $L=12$), which contains 110M parameters. For student models, we compress the teacher model to various model sizes for comparison, including WID₅₅ ($d=516$, $A=12$, $L=12$) with 55M parameters and WID₁₁ ($d=192$, $A=12$, $L=12$) with 11M parameters. We use the documents of English Wikipedia and BookCorpus (Zhu et al., 2015) for pre-training following Devlin et al. (2019). We use Adamw (Loshchilov and Hutter, 2019) with $\beta_1 = 0.9$, $\beta_2 = 0.99$. The compactors are trained with peak learning rate $5e-5$ and the original linear layers with peak learning rate $1e-6$. For WID, we adopt the 2-norm and set $N=500$, $d=\lfloor (d_t-d_s)/16 \rfloor$. It costs about 64 hours to train for 400,000 steps with a batch size of 960 on 8 A100 GPUs.

¹We employ the weight from <https://huggingface.co/bert-base-uncased>.

4.2 Downstream Tasks

Following previous PLM distillation (Sanh et al., 2019; Wang et al., 2020), we evaluate our WID on the SQuAD v1.1 (Rajpurkar et al., 2016) and GLUE benchmark (Wang et al., 2019). The GLUE benchmark consists of CoLA (Warstadt et al., 2019), SST-2 (Socher et al., 2013), MRPC (Dolan and Brockett, 2005), STS-B (Cer et al., 2017), QQP (Chen et al., 2018), MNLI (Williams et al., 2018), QNLI (Rajpurkar et al., 2016) and RTE (Bentivogli et al., 2009). After task-agnostic distillation, we fine-tune our compressed BERT WID₅₅ and WID₁₁ on these benchmarks adopting the grid search and report the results on the development sets. The result of MNLI is the score of MNLI-m. More details about these datasets including dataset sizes and metrics and the hyperparameters for fine-tune can be found in the Appendix A.

4.3 Baselines

For a fair comparison, we compare our WID with the **task-agnostic distillation** baselines. These baselines include: 1) DistilBERT (Sanh et al., 2019), which distills the student by the combination of the original MLM loss, the cosine distance for features, and the KL divergence for output logits. 2) TinyBERT (GD) (Jiao et al., 2020), which aligns the attention distributions and hidden states for general distillation. 3) MiniLM (Wang et al., 2020) and MiniLM v2 (Wang et al., 2021), which align the attention matrix and values-values scaled dot-product. We also reproduce the TinyBERT in the same architecture as WID, following the official code. For fair comparison, we employ the same corpus and follows the official hyperparameters. We do not compare with MobileBERT (Sun et al., 2020) since its teacher is IB-BERT_{large} (much higher accuracy than BERT_{base}) and its computations (4096 batch size * 740,000 steps) is much higher. Moreover, we also compare WID with task-specific methods in Appendix C.

4.4 Main Results

We compare our WID with other task-agnostic distillation methods in **various** model size. All the methods utilize the BERT_{base} as teacher model. As shown in Table 2, WID retains 98.9% and 90.9% performance of BERT_{base} with 49.2% and 10.2% parameters, respectively. In particular, on the CoLA task, our proposed WID₅₅ gets a higher score than BERT_{base}. Compared to the baselines

Method	FLOPs	Params	SST-2	CoLA	MRPC	QNLI	QQP	RTE	STS-B	MNLI	SQuAD	AVG
BERT _{base}	22.7B	110.1M	92.7	59.1	90.4	91.7	91.4	70.8	90.1	84.5	89.6/82.6	84.3
DistilBERT	11.9B	67.5M	91.3	51.3	87.5	89.2	88.5	59.9	86.9	82.2	86.2/78.1	80.1
MiniLM	11.9B	67.5M	92.0	49.2	-	91.0	91.0	71.5	-	-	-/-	-
MiniLM v2	11.9B	67.5M	92.4	52.5	-	90.8	91.1	72.1	-	-	-/-	-
TinyBERT (GD) [†]	11.9B	67.5M	92.9	44.1	89.5	90.7	91.0	73.7	89.6	83.8	84.0/74.2	81.3
TinyBERT (GD) [‡]	10.4B	54.9M	92.3	47.0	87.3	90.8	90.9	69.7	89.0	83.3	85.4/76.2	81.2
WID ₅₅ (ours)	10.4B	54.9M	92.4	61.7	88.2	90.1	91.0	70.4	87.9	82.9	88.5/80.8	83.4
TinyBERT (GD) [‡]	1.6B	11.3M	88.4	30.3	80.4	87.5	89.1	65.3	84.0	79.4	80.5/70.7	75.6
WID ₁₁ (ours)	1.6B	11.3M	88.8	44.2	81.9	85.4	89.5	60.3	84.5	78.4	81.2/72.4	76.7

Table 2: Comparison between our WID and the previous task-agnostic distillation methods. For SQuAD v1.1, we report the F1/EM scores. We compare the task-agnostic distilled models without both data augmentation and task-specific distillation. WID achieves better performances than TinyBERT under various model size. † means that we fine-tune the official weights. ‡ means that we reproduce the methods following the official code. Other results are taken from corresponding papers.

with 67.5M parameters, WID₅₅ gets comparable performance with MiniLM and higher performance than DistilBERT with less parameters. Meanwhile, WID outperforms the TinyBERT under the same architecture on GLUE benchmarks and SQuAD, showing its supremacy over the traditional KD methods with logit-based loss and feature-based loss. Without CoLA, WID₅₅ gets an average score of 85.8 and still outperforms the TinyBERT (GD) with an average score of 85.0.

Larger Performance Gap Since performance gap between teacher and student has always been a crucial point and difficulty in the knowledge distillation. We conduct experiments for smaller student models (11.3M parameters). We reproduce the task-agnostic TinyBERT under the General Distillation (GD) as baseline. As shown in Table 2, we find that WID (average score: 76.7) still outperforms TinyBERT (average score: 75.6) when the student model is about 10x smaller.

5 Analysis and Discussion

5.1 Compare WID with Pruning and Self-Distillation

We propose WID, a weight-inherited distillation method for task-agnostic BERT compression without extra alignment loss, which learns mappings from the teacher model to compact student via re-parameterization. To compress the linear layer, we design the row compactor and column compactor for row squeezing and column squeezing, respectively. However, WID is very likely to be fused with pruning (LeCun et al., 1989) and self-distillation (Zhang et al., 2019).

Pruning aims to remove redundant weights from

a neural network to achieve parameter-efficiency while preserving model performance, including unstructured pruning which sets weights to 0, and structured pruning which removes components. However, unstructured pruning does not compress the model size, while structured pruning prunes the weights directly. In WID, we **do not remove any parts** from the original teacher model. Instead of that, we learn the compactors to compress the weights via structural re-parameterization.

Self-distillation (Zhang et al., 2019) is a one-step **online distillation** method, which distills the knowledge in deeper layer to shallow layer during the training process of teacher model. Compared to self-distillation, WID is an **offline distillation** method, since the teacher model is trained before knowledge distillation. Furthermore, self-distillation aims to transfer knowledge by aligning intermediate features or logit distributions, while WID transfers knowledge by inheriting the weight directly.

5.2 MHA: Dropping Head or Reducing Dimension

Multi-Head Attention (MHA) allows the model to jointly attend to the information from different representation subspaces (Vaswani et al., 2017). When compressing the weights in MHA, there are two options, including 1) dropping head, which reduces the number of heads A and 2) reducing dimension, which reduces the size of each head d_k . For TinyBERT (Jiao et al., 2020) and MiniLM (Wang et al., 2020), they keep $A=12$ and reduce d_k due to the constraint of attention-based loss. Our proposed WID is more flexible, since we do not employ any alignment loss. Moreover, we can easily achieve

Method	SST-2	CoLA	MRPC	QNLI	QQP	RTE	STS-B	MNLI	SQuAD	AVG
WID ₅₅ ^{dim}	92.4	61.7	88.2	90.1	91.0	70.4	87.9	82.9	88.5/80.8	83.4
WID ₅₅ ^{head}	92.0	61.6	88.2	89.4	91.0	70.8	87.6	82.6	87.3/79.4	83.0
WID ₁₁ ^{dim}	88.8	44.2	81.9	85.4	89.5	60.3	84.5	78.4	81.2/72.4	76.7
WID ₁₁ ^{head}	89.6	46.2	83.1	86.1	89.5	62.1	85.3	79.0	81.7/72.9	77.6

Table 3: Comparison between dropping head and reducing dimension of each head for WID₅₅ with 55M parameters and WID₁₁ with 11M parameters.

Teacher	Params	SST-2	CoLA	MRPC	QNLI	QQP	RTE	STS-B	MNLI	SQuAD	AVG
BERT _{base}	110.1M	89.6	46.2	83.1	86.1	89.5	62.1	85.3	79.0	81.7/72.9	77.6
BERT ₅₅	54.2M	89.5	43.2	84.6	86.3	89.7	63.2	85.7	79.4	81.2/72.5	77.5
WID ₅₅ ^{head}	54.2M	89.9	46.2	84.8	86.5	89.5	64.6	84.7	78.8	82.1/73.5	78.1

Table 4: Comparison between different teacher models which are compressed to WID₁₁^{head}. BERT₅₅ means the BERT model with same architecture as WID₅₅^{head}.

these two strategies by constraining the column mask in MHA. For WID₅₅ and WID₁₁ reported in Table 2, we reduce the size of each attention head following TinyBERT for a fair comparison.

To further explore these two strategies, we conduct WID under these two settings and report the scores on downstream tasks. In BERT_{base}, we have $A=12$ and $d_k=64$. The student models are selected as: WID₅₅^{dim} ($A=12$, $d_k=43$), WID₅₅^{head} ($A=8$, $d_k=64$), WID₁₁^{dim} ($A=12$, $d_k=16$), and WID₁₁^{head} ($A=3$, $d_k=64$). As shown in Table 3, the dropping head strategy performs slightly worse under 55M parameters and much better under 11M parameters. For attention heads in WID₅₅, both 43 and 64 are large enough to encode the textual information in the representation subspace. Thus, the WID₅₅^{dim} with more attention heads gets slightly better results. Similarly, the attention heads with size 16 perform worse due to the limited representation subspace, leading to the poor performance of WID₁₁^{dim}.

5.3 Impact of Teacher Models

To study the impact of teacher models, we compare the results of three teachers, including 1) BERT_{base}, 2) WID₅₅^{head}, which are compressed by BERT_{base} adopting the dropping head strategy, 3) BERT₅₅, which shares the same architecture as WID₅₅^{head}. Both BERT_{base} and BERT₅₅ are downloaded from the official repository². We compress these three teachers to WID₁₁^{head} employing the dropping head strategy.

Table 4 shows the results of three teachers. Some

²<https://github.com/google-research/bert>

findings are summarized as follows:

(1) Smaller teacher can also teach smart student. Both BERT_{base} and BERT₅₅ are pre-trained on the MLM tasks. We can find that the compressed student from BERT₅₅ gets an average score of 77.5, which is comparable to 77.7 from the student of BERT_{base}.

(2) Educated teacher teach better. The WID₅₅^{head} are compressed by BERT_{base} adopting the dropping head strategy. Compared to BERT₅₅ under the same architecture, WID₅₅^{head} can teach a better student on both GLUE benchmarks and the SQuAD task.

5.4 Looking into WID

We visualize the attention distributions between the teacher BERT_{base} and the student WID₁₁^{dim} with the same input tokens. For more comparison, we also pre-train BERT₁₁ which shares the same architecture as WID₁₁^{dim}. As shown in Figure 3, we find that WID can learn the attention patterns in various layers of the teacher model BERT_{base}, while BERT₁₁ is much more different. The results of more attention heads in these models can be found in the Appendix B.

In WID, we adopt the hard loss for the pre-training task during the distillation, without any alignment loss between the teacher model and the student model. However, the compressed student model can also learn the knowledge about attention patterns. This observation indicates that inheriting the weights can also inheriting the high-level semantic knowledge.

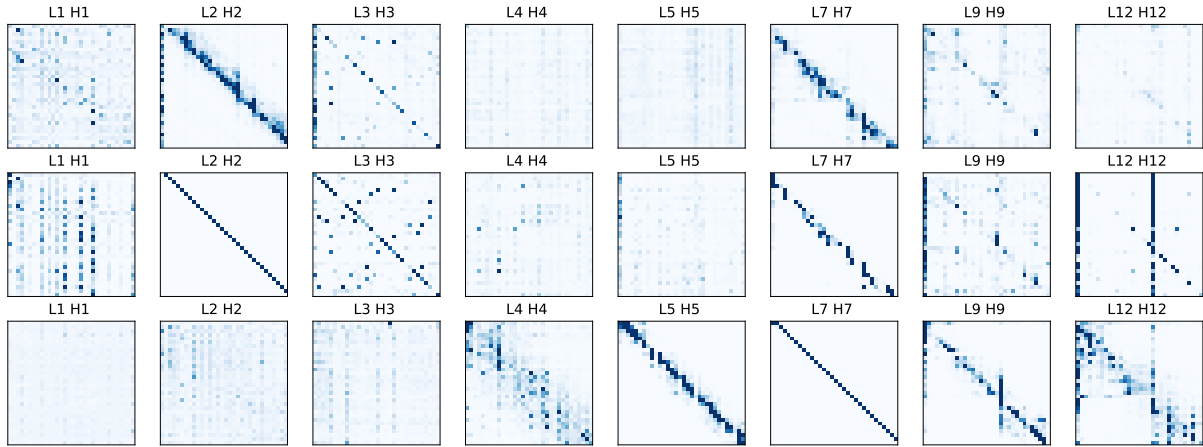


Figure 3: Attention distributions under same input tokens for $BERT_{base}$ (upper), WID_{11}^{dim} (middle), and $BERT_{11}$ (bottom). Our WID can learn the knowledge about attention distributions from teacher without any alignment loss.

6 Related Work

6.1 BERT Compression

Transformer-based Pre-trained Language Models (PLMs) can be compressed via Quantization (Stock et al., 2021; Tao et al., 2022), Matrix Decomposition (Mao et al., 2020), Pruning (Xia et al., 2022; Lagunas et al., 2021), and Knowledge Distillation (Jiao et al., 2020; Wang et al., 2020). We refer the readers to Ganesh et al. (2021) for a comprehensive survey. In this paper, we focus on knowledge distillation for bert compression.

6.2 Knowledge Distillation

Knowledge Distillation refers to transfer the knowledge from the teacher model to the student model (Hinton et al., 2015). The distillation methods can be directly divided into three main categories: offline distillation, online distillation, and self-distillation (Gou et al., 2021). For PLMs, majority methods follow the offline distillation pattern where the teacher model is pre-trained before distillation. Meanwhile, distillation methods for PLMs can be divided into task-agnostic, which distill PLM in pre-training stage, and task-specific, which fine-tune the teacher model on specific tasks and then distill.

In this work, we focus on the task-agnostic distillation since the task-specifically fine-tuning procedure of large PLMs is costly and time-consuming while the task-agnostic distilled models can be directly fine-tuned on downstream tasks. Previous methods mainly focus on designing extra matching losses for the student model to mimic the teacher model. These loss objects mainly include feature-

based loss for features in intermediate layers and logit-based loss for output logits. DistilBERT (Sanh et al., 2019) adopts the output logit and embedding outputs of the teacher to train the student. TinyBERT (Jiao et al., 2020) and MobileBERT (Sun et al., 2020) further employ the self-attention distributions and hidden states for alignment loss. Such layer-to-layer distillation restrict the number of student layers or require an extra mapping function. To address this issue, MiniLM (Wang et al., 2020) proposes a new loss based on the attention matrix and values-values scaled dot-product.

Different from previous methods, our proposed WID does not require additional alignment losses, thus avoiding labor selection for both loss functions and loss weights. We directly leverage the knowledge contained in the weights of the teacher model.

7 Conclusion

In this work, we propose a novel Weight-Inherited Distillation (WID) method for task-agnostic BERT compression. In WID, we consider the compression process as weight mapping, and design the row compactors and column compactors for row mapping and column mapping. Empirical results on various student model sizes demonstrate the effectiveness of WID. Further analysis indicates that inheriting the weights can also inheriting high-level semantic knowledge such as attention patterns. In future work, we would consider to reduce the extra memory cost by compactor layers, such as compactor sharing. Moreover, performing the WID on other backbones such as GNN would be another interesting topic.

553 Limitations

554 Our proposed WID adds row/column compactors
555 to learn the mappings from the teacher model to the
556 student model. Thus, WID requires additional com-
557 putational time and memory. However, WID still
558 outperforms TinyBERT with less time costs. As
559 shown in Table 6, WID₅₅^{dim} trained with 100k steps
560 achieves a higher score and saves more than 50%
561 time costs compared to TinyBERT. Meanwhile, we
562 believe that such a trade-off is valuable, since a
563 faster and better compact student would save more
564 time in downstream tasks.

565 References

566 Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E.
567 Hinton. 2016. [Layer normalization](#). *CoRR*,
568 abs/1607.06450.

569 Luisa Bentivogli, Bernardo Magnini, Ido Dagan,
570 Hoa Trang Dang, and Danilo Giampiccolo. 2009.
571 [The fifth PASCAL recognizing textual entailment
572 challenge](#). In *Proceedings of the Second Text Analy-
573 sis Conference, TAC 2009, Gaithersburg, Maryland,
574 USA, November 16-17, 2009*. NIST.

575 Daniel M. Cer, Mona T. Diab, Eneko Agirre, Iñigo
576 Lopez-Gazpio, and Lucia Specia. 2017. [Semeval-
577 2017 task 1: Semantic textual similarity multilingual
578 and crosslingual focused evaluation](#). In *Proceedings
579 of the 11th International Workshop on Semantic Eval-
580 uation, SemEval@ACL 2017, Vancouver, Canada,
581 August 3-4, 2017*, pages 1–14. Association for Com-
582 putational Linguistics.

583 Zihan Chen, Hongbo Zhang, Xiaoji Zhang, and Leqi
584 Zhao. 2018. [Quora question pairs](#).

585 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
586 Kristina Toutanova. 2019. [BERT: pre-training of
587 deep bidirectional transformers for language under-
588 standing](#). In *Proceedings of the 2019 Conference of
589 the North American Chapter of the Association for
590 Computational Linguistics: Human Language Tech-
591 nologies, NAACL-HLT 2019, Minneapolis, MN, USA,
592 June 2-7, 2019, Volume 1 (Long and Short Papers)*,
593 pages 4171–4186. Association for Computational
594 Linguistics.

595 Xiaohan Ding, Tianxiang Hao, Jianchao Tan, Ji Liu, Jun-
596 gong Han, Yuchen Guo, and Guiguang Ding. 2021.
597 [Resrep: Lossless CNN pruning via decoupling re-
598 membering and forgetting](#). In *2021 IEEE/CVF In-
599 ternational Conference on Computer Vision, ICCV
600 2021, Montreal, QC, Canada, October 10-17, 2021*,
601 pages 4490–4500. IEEE.

602 William B. Dolan and Chris Brockett. 2005. [Automati-
603 cally constructing a corpus of sentential paraphrases](#).
604 In *Proceedings of the Third International Workshop
605 on Paraphrasing, IWP@IJCNLP 2005, Jeju Island,*

Korea, October 2005, 2005. Asian Federation of Nat-
ural Language Processing.

606
607
608 Prakhhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali
609 Khan, Yin Yang, Deming Chen, Marianne Winslett,
610 Hassan Sajjad, and Preslav Nakov. 2021. [Compress-
611 ing large-scale transformer-based models: A case
612 study on BERT](#). *Transactions of the Association for
613 Computational Linguistics*, 9:1061–1080.

614 Jianping Gou, Baosheng Yu, Stephen J. Maybank, and
615 Dacheng Tao. 2021. [Knowledge distillation: A sur-
616 vey](#). *Int. J. Comput. Vis.*, 129(6):1789–1819.

617 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian
618 Sun. 2016. [Deep residual learning for image recogni-
619 tion](#). In *2016 IEEE Conference on Computer Vision
620 and Pattern Recognition, CVPR 2016, Las Vegas,
621 NV, USA, June 27-30, 2016*, pages 770–778. IEEE
622 Computer Society.

623 Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean.
624 2015. [Distilling the knowledge in a neural network](#).
625 *CoRR*, abs/1503.02531.

626 Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao
627 Chen, and Qun Liu. 2020. [Dynabert: Dynamic BERT
628 with adaptive width and depth](#). In *Advances in Neu-
629 ral Information Processing Systems 33: Annual Con-
630 ference on Neural Information Processing Systems
631 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

632 Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao
633 Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. [Tinybert: Distilling BERT for natural language un-
634 derstanding](#). In *Findings of the Association for Com-
635 putational Linguistics: EMNLP 2020, Online Event,
636 16-20 November 2020*, volume EMNLP 2020 of *Find-
637 ings of ACL*, pages 4163–4174. Association for Com-
638 putational Linguistics.

640 François Lagunas, Ella Charlaix, Victor Sanh, and
641 Alexander M. Rush. 2021. [Block pruning for faster
642 transformers](#). In *Proceedings of the 2021 Confer-
643 ence on Empirical Methods in Natural Language Pro-
644 cessing, EMNLP 2021, Virtual Event / Punta Cana,
645 Dominican Republic, 7-11 November, 2021*, pages
646 10619–10629. Association for Computational Lin-
647 guistics.

648 Yann LeCun, John S. Denker, and Sara A. Solla. 1989.
649 [Optimal brain damage](#). In *Advances in Neural In-
650 formation Processing Systems 2, [NIPS Conference,
651 Denver, Colorado, USA, November 27-30, 1989]*,
652 pages 598–605. Morgan Kaufmann.

653 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
654 dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
655 Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining
656 approach](#). *CoRR*, abs/1907.11692.

657
658 Ilya Loshchilov and Frank Hutter. 2019. [Decoupled
659 weight decay regularization](#). In *7th International
660 Conference on Learning Representations, ICLR 2019,
661 New Orleans, LA, USA, May 6-9, 2019*. OpenRe-
662 view.net.

663	Yihuan Mao, Yujing Wang, Chufan Wu, Chen Zhang,	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	722
664	Yang Wang, Quanlu Zhang, Yaming Yang, Yunhai	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	723
665	Tong, and Jing Bai. 2020. Ladabert: Lightweight	Kaiser, and Illia Polosukhin. 2017. Attention is all	724
666	adaptation of BERT through hybrid model compres-	you need . In <i>Advances in Neural Information Pro-</i>	725
667	sion . In <i>Proceedings of the 28th International Confer-</i>	cessing Systems 30: Annual Conference on Neural	726
668	ence on Computational Linguistics, COLING 2020,	Information Processing Systems 2017, December 4-9,	727
669	Barcelona, Spain (Online), December 8-13, 2020,	2017, Long Beach, CA, USA, pages 5998–6008.	728
670	pages 3225–3234. International Committee on Com-		
671	putational Linguistics.		
672	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	Alex Wang, Amanpreet Singh, Julian Michael, Felix	729
673	Percy Liang. 2016. Squad: 100, 000+ questions	Hill, Omer Levy, and Samuel R. Bowman. 2019.	730
674	for machine comprehension of text . In <i>Proceedings</i>	GLUE: A multi-task benchmark and analysis plat-	731
675	of the 2016 Conference on Empirical Methods in	form for natural language understanding . In <i>7th In-</i>	732
676	Natural Language Processing, EMNLP 2016, Austin,	ternational Conference on Learning Representations,	733
677	Texas, USA, November 1-4, 2016, pages 2383–2392.	ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.	734
678	The Association for Computational Linguistics.	OpenReview.net.	735
679	Victor Sanh, Lysandre Debut, Julien Chaumond, and	Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong,	736
680	Thomas Wolf. 2019. Distilbert, a distilled version	and Furu Wei. 2021. Minilmv2: Multi-head self-	737
681	of BERT: smaller, faster, cheaper and lighter . <i>CoRR,</i>	attention relation distillation for compressing pre-	738
682	abs/1910.01108 .	trained transformers . In <i>Findings of the Associa-</i>	739
683	Richard Socher, Alex Perelygin, Jean Wu, Jason	tion for Computational Linguistics: ACL/IJCNLP	740
684	Chuang, Christopher D. Manning, Andrew Y. Ng,	2021, Online Event, August 1-6, 2021, volume	741
685	and Christopher Potts. 2013. Recursive deep mod-	ACL/IJCNLP 2021 of Findings of ACL, pages 2140–	742
686	els for semantic compositionality over a sentiment	2151 . Association for Computational Linguistics.	743
687	treebank . In <i>Proceedings of the 2013 Conference on</i>		
688	Empirical Methods in Natural Language Processing,	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan	744
689	EMNLP 2013, 18-21 October 2013, Grand Hyatt	Yang, and Ming Zhou. 2020. Minilm: Deep self-	745
690	Seattle, Seattle, Washington, USA, A meeting of SIG-	attention distillation for task-agnostic compression	746
691	DAT, a Special Interest Group of the ACL, pages	of pre-trained transformers . In <i>Advances in Neural</i>	747
692	1631–1642 . ACL.	Information Processing Systems 33: Annual Confer-	748
693	Pierre Stock, Angela Fan, Benjamin Graham, Edouard	ence on Neural Information Processing Systems 2020,	749
694	Grave, Rémi Gribonval, Hervé Jégou, and Armand	NeurIPS 2020, December 6-12, 2020, virtual.	750
695	Joulin. 2021. Training with quantization noise for ex-	Alex Warstadt, Amanpreet Singh, and Samuel R. Bow-	751
696	treme model compression . In <i>9th International Con-</i>	man. 2019. Neural network acceptability judgments .	752
697	ference on Learning Representations, ICLR 2021, Vir-	Trans. Assoc. Comput. Linguistics, 7:625–641.	753
698	tual Event, Austria, May 3-7, 2021 . OpenReview.net.		
699	Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. 2019.	Adina Williams, Nikita Nangia, and Samuel R. Bow-	754
700	Patient knowledge distillation for BERT model com-	man. 2018. A broad-coverage challenge corpus for	755
701	pression . In <i>Proceedings of the 2019 Conference on</i>	sentence understanding through inference . In <i>Pro-</i>	756
702	Empirical Methods in Natural Language Processing	ceedings of the 2018 Conference of the North Amer-	757
703	and the 9th International Joint Conference on Nat-	ican Chapter of the Association for Computational	758
704	ural Language Processing, EMNLP-IJCNLP 2019,	Linguistics: Human Language Technologies, NAACL-	759
705	Hong Kong, China, November 3-7, 2019, pages 4322–	HLT 2018, New Orleans, Louisiana, USA, June 1-6,	760
706	4331 . Association for Computational Linguistics.	2018, Volume 1 (Long Papers), pages 1112–1122.	761
707	Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu,	Association for Computational Linguistics.	762
708	Yiming Yang, and Denny Zhou. 2020. Mobilebert:	Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le,	763
709	a compact task-agnostic BERT for resource-limited	Mohammad Norouzi, Wolfgang Macherey, Maxim	764
710	devices . In <i>Proceedings of the 58th Annual Meet-</i>	Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff	765
711	ing of the Association for Computational Linguistics,	Klingner, Apurva Shah, Melvin Johnson, Xiaobing	766
712	ACL 2020, Online, July 5-10, 2020, pages 2158–2170.	Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato,	767
713	Association for Computational Linguistics.	Taku Kudo, Hideto Kazawa, Keith Stevens, George	768
714	Chaofan Tao, Lu Hou, Wei Zhang, Lifeng Shang, Xin	Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason	769
715	Jiang, Qun Liu, Ping Luo, and Ngai Wong. 2022.	Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals,	770
716	Compression of generative pre-trained language	Greg Corrado, Macduff Hughes, and Jeffrey Dean.	771
717	models via quantization . In <i>Proceedings of the 60th</i>	2016. Google’s neural machine translation system:	772
718	Annual Meeting of the Association for Computational	Bridging the gap between human and machine	773
719	Linguistics (Volume 1: Long Papers), ACL 2022,	translation . <i>CoRR, abs/1609.08144</i> .	774
720	Dublin, Ireland, May 22-27, 2022, pages 4821–4836.	Mengzhou Xia, Zexuan Zhong, and Danqi Chen. 2022.	775
721	Association for Computational Linguistics.	Structured pruning learns compact and accurate	776
		models . In <i>Proceedings of the 60th Annual Meeting of</i>	777
		the Association for Computational Linguistics (Vol-	778
		ume 1: Long Papers), ACL 2022, Dublin, Ireland,	779

780 May 22-27, 2022, pages 1513–1528. Association for
781 Computational Linguistics.

782 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Car-
783 bonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019.
784 [Xlnet: Generalized autoregressive pretraining for](#)
785 [language understanding](#). In *Advances in Neural In-*
786 *formation Processing Systems 32: Annual Confer-*
787 *ence on Neural Information Processing Systems 2019,*
788 *NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*
789 *Canada*, pages 5754–5764.

790 Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen,
791 Chenglong Bao, and Kaisheng Ma. 2019. [Be your](#)
792 [own teacher: Improve the performance of convolu-](#)
793 [tional neural networks via self distillation](#). In *2019*
794 *IEEE/CVF International Conference on Computer*
795 *Vision, ICCV 2019, Seoul, Korea (South), October 27*
796 *- November 2, 2019*, pages 3712–3721. IEEE.

797 Wangchunshu Zhou, Canwen Xu, and Julian J.
798 McAuley. 2022. [BERT learns to teach: Knowledge](#)
799 [distillation with meta learning](#). In *Proceedings of the*
800 *60th Annual Meeting of the Association for Compu-*
801 *tational Linguistics (Volume 1: Long Papers), ACL*
802 *2022, Dublin, Ireland, May 22-27, 2022*, pages 7037–
803 7049. Association for Computational Linguistics.

804 Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan
805 Salakhutdinov, Raquel Urtasun, Antonio Torralba,
806 and Sanja Fidler. 2015. [Aligning books and movies:](#)
807 [Towards story-like visual explanations by watching](#)
808 [movies and reading books](#). In *2015 IEEE Interna-*
809 *tional Conference on Computer Vision, ICCV 2015,*
810 *Santiago, Chile, December 7-13, 2015*, pages 19–27.
811 IEEE Computer Society.

A GLUE and SQuAD

A.1 Data Statistics

Table 5 shows the sizes of the train/development set and the metrics for downstream tasks.

Task	#Train	#Dev	Metric
SST-2	67k	872	Accuracy
QNLI	105k	5.5k	Accuracy
MNLI	393k	20k	Accuracy
QQP	364k	40k	Accuracy
CoLA	8.5k	1k	Matthews corr.
RTE	2.5k	276	Accuracy
STS-B	7k	1.5k	Spearman corr.
MRPC	3.7k	408	Accuracy
SQuAD	87.6k	34.7k	F1 & EM

Table 5: Data statistics of GLUE and SQuAD datasets.

A.2 Hyperparameters

We employ the grid search to fine-tune the GLUE benchmarks and SQuAD.

GLUE The learning rate are searched in $\{1e-5, 2e-5, 3e-5\}$. We set the search space for the training batch size based on the size of the training set. For large dataset including QNLI, MNLI, and QQP, the batch size is searched in $\{32, 48\}$. For small dataset including MRPC, RTE, CoLA and STS-B, the batch size is searched in $\{4, 6\}$. For SST-2, the batch size is searched in $\{8, 16\}$. All tasks are trained for 10 epochs.

SQuAD The learning rate is searched in $\{1e-5, 2e-5, 3e-5\}$ and batch size is searched in $\{4,6,8\}$. The training epochs are set to 3.

B Attention Distributions

We visualize the attention distributions for the teacher $BERT_{base}$, pre-trained $BERT_{55}$ and the student WID_{11}^{head} under the same input tokens (input sentence: "if the world harassed me, it will harass you too.") in Figure 4, Figure 5 and Figure 6, respectively. From the bottom layer to the top layer, WID can effectively learn the attention patterns from the teacher model while $BERT_{11}$ is much more different.

C Comparison with Task-Specific Distillation

It can be unfair to directly compare task-agnostic WID with task-specific distillation methods, since

the teacher model in task-specific distillation methods is fine-tuned for the task before distillation. We compare our WID with DynaBERT (Hou et al., 2020) and MetaDistill(Zhou et al., 2022). As shown in Table 7, WID still outperforms these task-specific methods on the GLUE benchmarks.

D Less Training Steps

In Table 2, we report the results of WID_{55}^{dim} trained for 400k steps. We re-implement TinyBERT and train 3 epochs following the setting in Jiao et al. (2020). We reduce the training steps for WID_{55}^{dim} to 50k and 100k. All experiments are carried out with 8 A100 GPUs. As shown in Table 6, WID_{55}^{dim} trained with 100k steps can still outperform Tinybert and save more than 50% training time.

Model	Steps	Time	Score
TinyBERT (GD)	450k	33h	81.27
WID_{55}^{dim}	50k	8h	80.78
WID_{55}^{dim}	100k	16h	81.65
WID_{55}^{dim}	400k	64h	83.08

Table 6: Comparison between TinyBERT and WID trained with less steps on GLUE benchmarks.

Method	FLOPS	Params	SST-2	CoLA	MRPC	QNLI	QQP	RTE	STS-B	MNLI	AVG
BERT _{base}	22.7B	110.1M	92.7	59.1	90.4	91.7	91.4	70.8	90.1	84.5	83.8
DynaBERT	11.9B	67.5M	92.7	54.6	85.0	90.6	91.1	66.1	88.6	83.7	81.6
MetaDistill	11.9B	67.5M	92.3	58.6	86.8	90.4	91.0	69.4	89.1	83.8	82.7
TinyBERT*	11.9B	67.5M	91.9	52.4	86.5	89.8	90.6	67.7	88.7	83.8	81.4
WID ₅₅ (ours)	10.4B	54.9M	92.4	61.7	88.2	90.1	91.0	70.4	87.9	82.9	83.4

Table 7: Comparison between our WID and the previous task-specific distillation methods on GLUE benchmarks without data augmentation. * means the results are taken from Zhou et al. (2022).

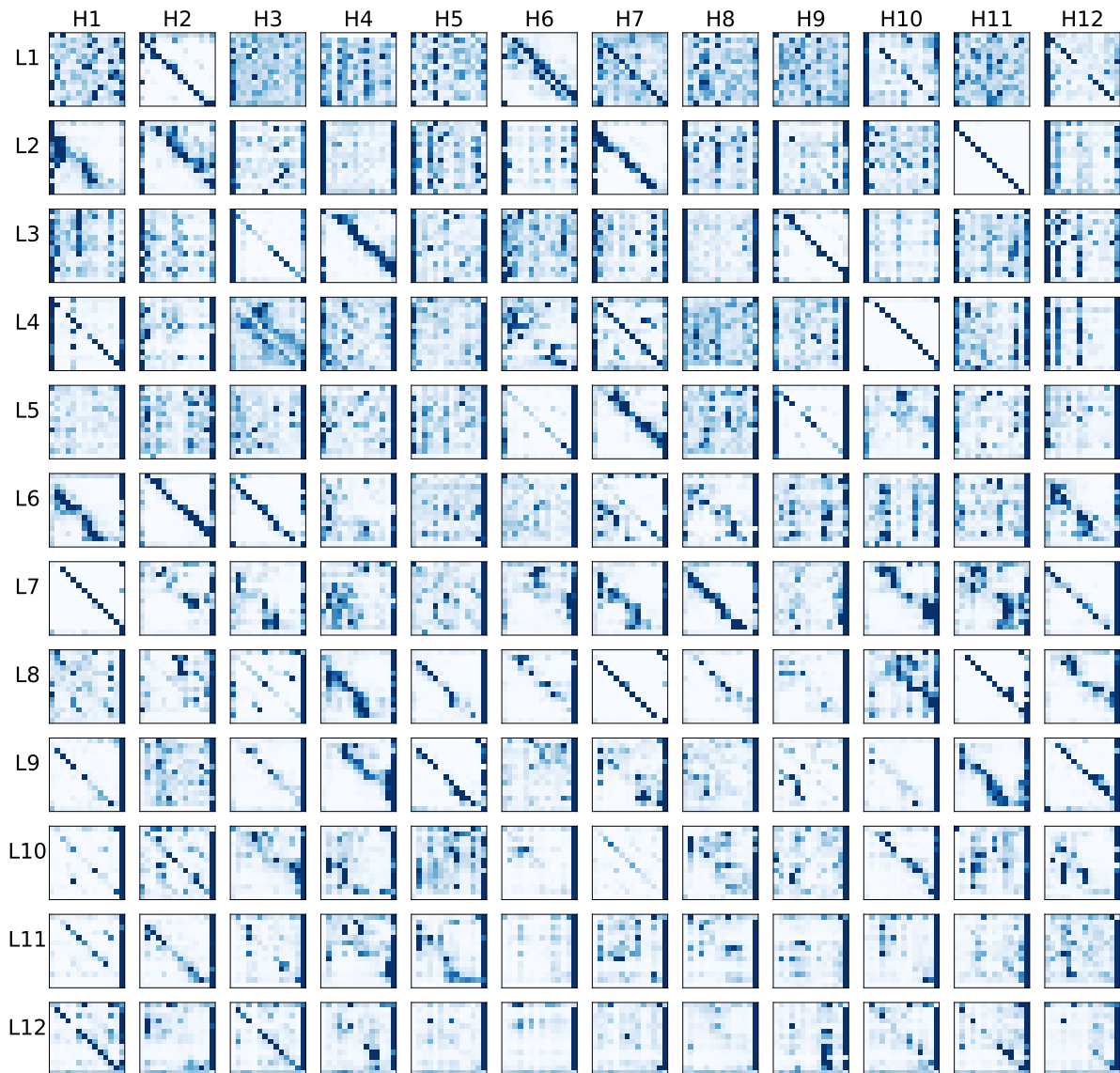


Figure 4: The self-attention distributions for teacher model BERT_{base}.

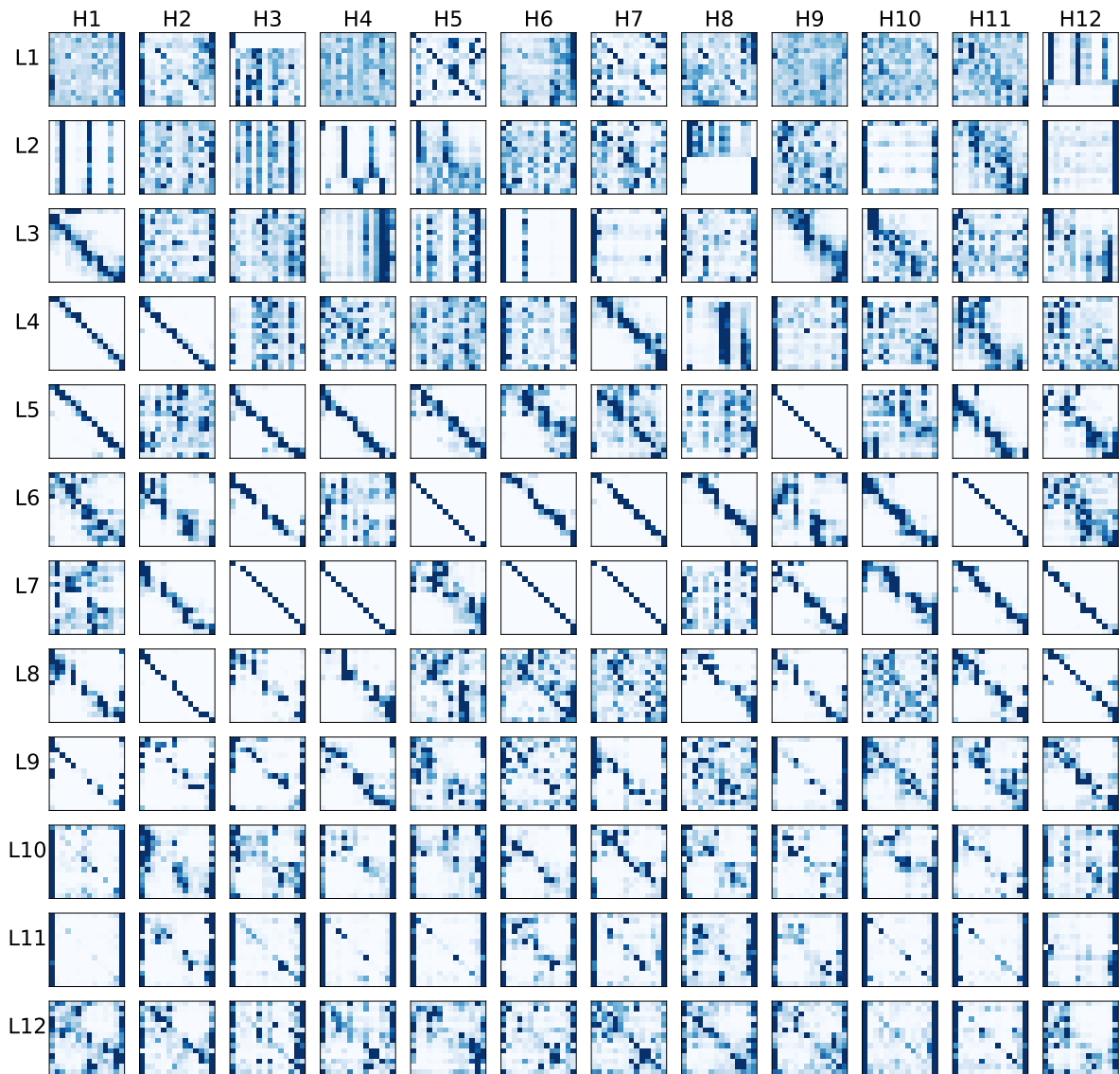


Figure 5: The self-attention distributions for BERT₁₁.

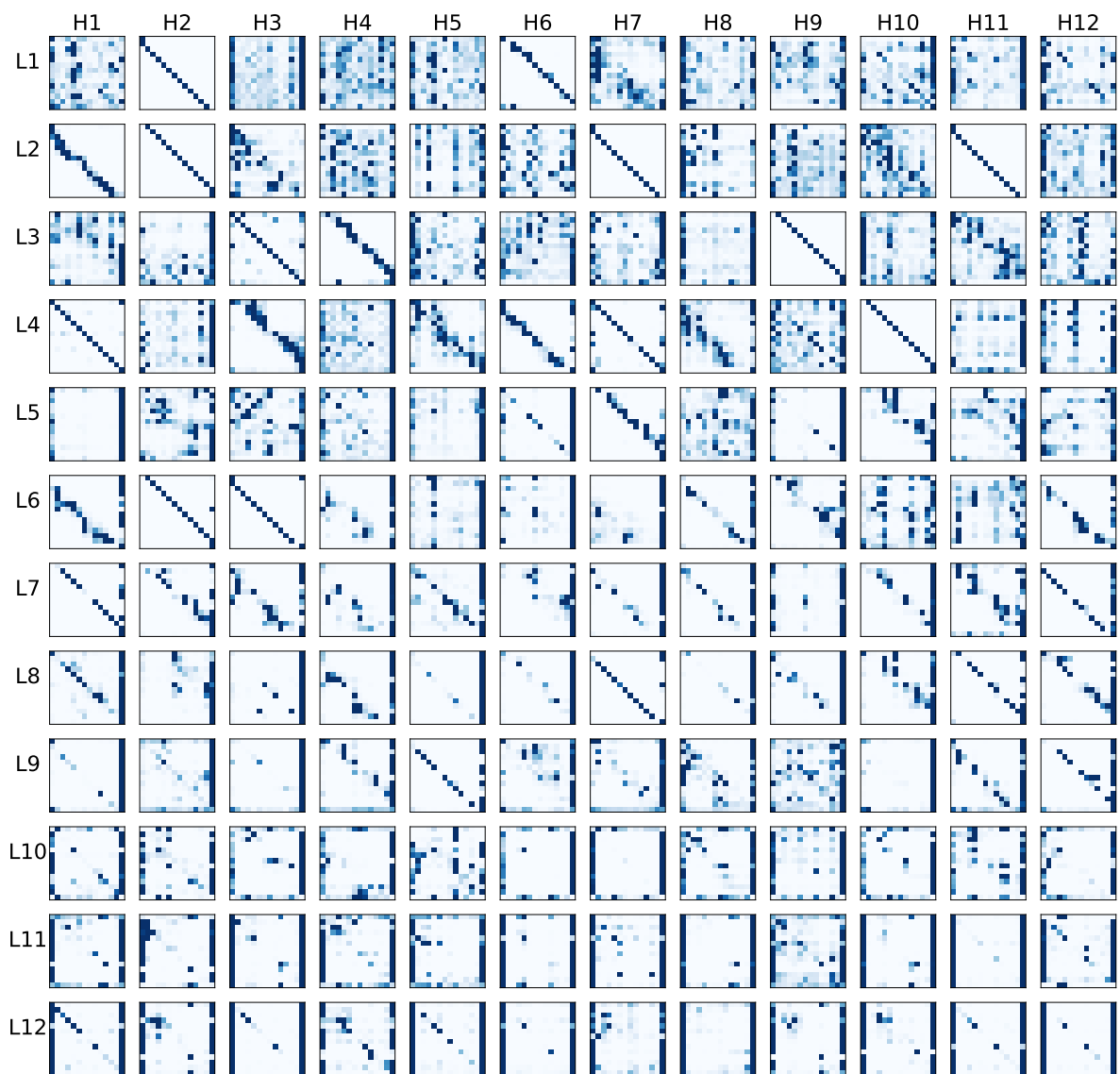


Figure 6: The self-attention distributions for our proposed WID_{11}^{dim} .