# SoftLMs: Efficient Adaptive Low-Rank Approximation of Language Models using Soft-Thresholding Mechanism

**Priyansh Bhatnagar[*], Linfeng Wen, Mingu Kang[*]**

University of California San Diego

[*]{prbhatnagar, mingu}@ucsd.edu

## Abstract

Extensive efforts have been made to boost the performance in the domain of language models by introducing various attention-based transformers. However, the inclusion of linear layers with large dimensions contributes to significant computational and memory overheads. The escalating computational demands of these models necessitate the development of various compression techniques to ensure their deployment on devices, particularly in resource-constrained environments. In this paper, we propose a novel compression methodology that dynamically determines the rank of each layer using a soft thresholding mechanism, which clips the singular values with a small magnitude in a differentiable form. This approach automates the decision-making process to identify the optimal degree of compression for each layer. We have successfully applied the proposed technique to attention-based architectures, including BERT for discriminative tasks and GPT2 and TinyLlama for generative tasks. Additionally, we have validated our method on Mamba, a recently proposed state-space model. Our experiments demonstrate that the proposed technique achieves a speed-up of $1.33\times$ to $1.72\times$ in the encoder/ decoder with a 50% reduction in total parameters.
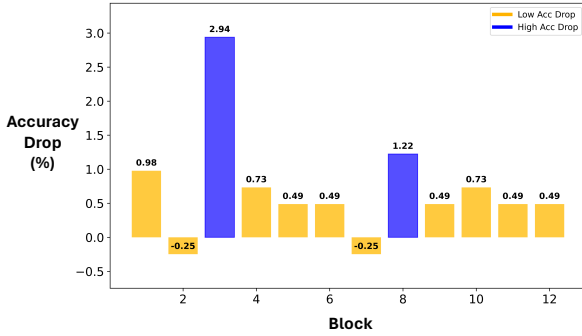
## 1 Introduction

In the domain of Natural Language Processing (NLP), attention-based models (Vaswani et al., 2023) have achieved unprecedented success across a variety of tasks including language modeling, text classification, and question answering (Devlin et al., 2018; Radford et al., 2019). Despite their superior performance, language models (LMs) have been growing increasingly larger, with new models being introduced with higher parameter numbers nearly every day. Despite the advancements in powerful GPUs that mitigate some of these overheads, successful deplo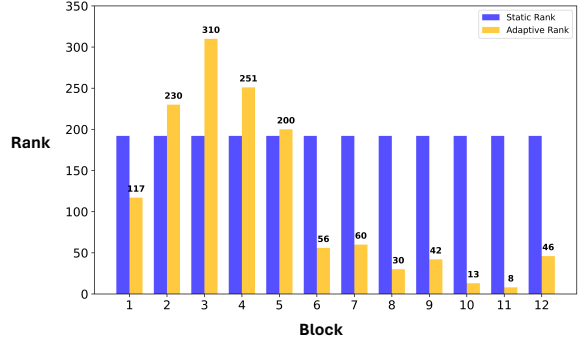yment on resource-constrained devices remains critical for expanding the applicability of these language models. Specifically, LMs require increasing memory and computational resources as sequence lengths grow longer, posing significant challenges for deployment on edge devices within tight memory and computational budgets.

For this reason, model compression has become increasingly important, introducing various techniques including pruning (Gordon et al., 2020; Han et al., 2016, 2015), knowledge distillation (Sun et al., 2019), and quantization (Zafrir et al., 2019; Wang et al., 2019b). However, these methods often require extensive computational resources and may not always be feasible for end-users due to their lengthy training times or limited compression effectiveness. An alternative line of research has explored the use of low-rank matrix factorization (Golub and Reinsch, 1971) to reduce the parameter count (Noach and Goldberg, 2020; Sainath et al., 2013). Singular Value Decomposition (SVD) is one of the widely employed methods to decompose weight matrices into smaller matrices (Hsu et al., 2022; Schotthöfer et al., 2022). Since one of the decomposed matrices contains singular values ordered by magnitude, which often indicates the importance of each dimension, truncating small singular values effectively compresses the matrix to have lower rank. However, this process incurs a new challenge as indiscriminately chosen rank regardless of their end-to-end impact on the model often degrades the performance significantly.

Recent studies have proposed modifications to this approach by incorporating task-specific knowledge into the compression process, such as weighting the singular values (Hsu et al., 2022; Wang et al., 2024). However, these methods assume static importance distributions, meaning that the *low-rank* is set to a constant for all layers or specific groups, which could be sub-optimal in terms of the performance and the parameter count.

| (a) Block-wise accuracy drop | (b) Static and adaptive rank decomposition |
|---|---|

Figure 1: Accuracy trend of low-rank approximation for BERT-Base with GLUE MRPC dataset. (a) accuracy drop when individually compressing a BERT block to reduce 50% of parameters, indicating varying contributions across blocks on overall task performance. (b) Rank of $W_Q$ layer obtained for each block from static vs. proposed adaptive rank decomposition after fine-tuning for both to have 50% overall parameter reductions. An F1 score gain of 4.3 is observed from the adaptive rank decomposition (F1: 90.7) over the static method (F1: 86.4).

On the other hand, our experiments summarized in Figure 1 underscore the importance of dynamic rank determination over blocks (layers). Figure 1(a) demonstrates that each block (layer) of the model contributes differently to overall performance, with specific layers showing higher tolerance for low-rank approximation. Figure 1(b) also indicates that optimally determining the rank for each block leads to significantly higher performance with the same total parameter count. Given this motivation, we aim to propose a platform to adaptively find the optimal low-rank to maximize compression with minimal accuracy degradations. Our key objectives of this work are as follows.

**Model compression for efficient inference:** We aim for efficient inference with a lower number of parameters, allowing deployment in a resource constrained environment. The reduced model complexity also accelerates inference, resulting in lower latency across various language models.

**Dynamic and adaptive rank decision:** Despite the potential for an optimal combination of ranks for each layer's weight at every processing stage, the prohibitively large search space makes it impractical to find the best combination through exhaustive searching over all possible rank combinations, especially given the recent growth in model sizes. This work automates such a search process during training.

**Learnable threshold for singular values:** In the process of truncating the singular values with low magnitude, it is critical to find the optimal threshold to maximize the compression while minimizing performance degradation. We employ a single

learnable threshold parameter for each linear layer to be adapted during the training. This approach allows each layer to determine the optimal degree of compression based on its specific contribution to the overall task performance. This process simply includes an additional threshold layer while adhering to the standard fine-tuning process based on the pre-trained model, without requiring sophisticated implementation and training methodologies.

Given above motivations, this paper proposes *SoftLM*, a language model that includes a differentiable soft threshold layer, making the threshold learnable and easily deployable in various models. This approach not only achieves competitive performance, with an average accuracy degradation of approximately 1% compared to the original models, but also provides significant reductions in latency, speeding up inference by upto 1.72× at 50% reduction in the number of parameters.

## 2 Related Work

**Static Low-Rank Decomposition**: Recent studies have explored low-rank approximation for deep learning models. Low-Rank Adaptation (LoRA) (Hu et al., 2021) proposes reducing the number of trainable parameters in large language models (LLMs) by factorizing weight updates into low-rank matrices. This approach allows for efficient fine-tuning of pre-trained models without significantly increasing computational resources. However, LoRA primarily focuses on reducing training time without compressing the model used for inference. Similarly, QLoRA (Dettmers et al., 2023) incorporates quantization techniques to further re-

duce memory footprint and computational complexity. Another approach, FWSVD (Hsu et al., 2022), introduces a decomposition method that approximates weight matrices using SVD combined with Fisher information, which quantifies the amount of information in the specific matrix. Additionally, ASVD (Yuan et al., 2024) scales the weight matrix based on the activations' distribution for the layer. However, ASVD suffers from limited parameter reduction of 10% to 20% while FWSVD suffers from the noticeable performance degradation. In SVD-LLM (Wang et al., 2024), the truncation of singular values is directly controlled by introducing a compression loss unlike in FWSVD and ASVD. However, these works only support static rank across all the layers without incorporating the layer-wise impact on performance. **Dynamic low-Rank decomposition**: Rather than employing a static rank for all layers, Low-Rank Adaptation (AdaLoRA) (Zhang et al., 2023) adjusts the rank based on learning dynamics and task complexity, thereby improving efficiency and effectiveness in minimizing the training complexity, but not for the inference. For the inference tasks, RankDyna (Hua et al., 2023) proposed an importance-driven dynamic rank adjustment for low-rank decomposition in LLMs, leveraging gradient-based techniques to determine the optimal rank for each layer during training. However, as well noted by the authors, RankDyna suffers from significant memory overhead due to the need to store entire gradients for importance calculations during fine-tuning.

## 3 Background

### 3.1 Low rank approximation

Singular Value Decomposition (SVD) performs decomposition of a matrix, i.e., for any matrix $W \in \mathbb{R}^{M \times N}$, SVD is given by:

$$W = U\Sigma V^T \quad (1)$$

where: $U$ is an $M \times M$ orthogonal matrix, $V$ is an $N \times N$ orthogonal matrix, $\Sigma$ is an $M \times N$ diagonal matrix with non-negative real numbers on the diagonal. Elements on the diagonal of $\Sigma$, known as singular values, are denoted by $\sigma_1, \sigma_2, \ldots, \sigma_r$ where $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$ and $r$ is the rank of matrix $W$. To reduce dimensionality, a truncated SVD is commonly used:

$$W \approx W_k = U_k \Sigma_k V_k^T$$

where $U_k$, $\Sigma_k$, and $V_k$ contain only the first $k$ columns, singular values, and rows, respectively.

### 3.2 Attention Based Language Models

Attention-based models (Vaswani et al., 2023) typically comprise of two main modules: Multi Head Attention (MHA) and Feed Forward Network (FFN). The attention mechanism selectively focuses on different parts of the input sequence when generating an output. The fundamental principle is to compute a weighted sum of input values ($V$), where the weights are determined by the similarity between a query ($Q$) and a set of keys ($K$) as follows, where $Q$, $V$ and $K$ are $N \times d_k$-dim matrices with $N$ being the sequence length:

$$Q = W_Q(x), K = W_K(x), V = W_V(x) \quad (2)$$

$$Attn(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (3)$$

where $x$ is the input sequence, $\sqrt{d_k}$ and $W_Q$, $W_K$, and $W_V$ are learnable weight matrices for the generation of queries, keys, and values, respectively. The attention scores are computed using the dot product of the query with each key stored in each column of $K$, followed by a softmax function to obtain the weights. The attention module is typically followed by an additional linear layer . Feed Forward Network (FFN) consists of two linear layers ($W_{fc1}$ and $W_{fc2}$). The first layer expands the dimension of the input, while the second layer reduces it back to the original dimension.

### 3.3 SSMs and Mamba

State-space models (SSMs) are mathematical models used to describe the behavior of dynamic systems. A general state-space model is defined by:

$$\begin{aligned} h_t &= Ah_{t-1} + Bx_{t-1} \\ y_t &= Ch_t \end{aligned} \quad (4)$$

where $h_t$, $x_t$, and $y_t$ are the state, input, and output vectors or matrices at time step $t$. The $A$, $B$ and $C$ are matrices that define the system dynamics.

Mamba is a selective state-space based model which includes above state-space component surrounded by two projection layers at the start ($in_{proj}$) and at the end ($out_{proj}$) in each block. Unlike the conventional SSM, the $B, C$, and $\Delta$ are made input-dependent by being generated from the linear operation on the input. The $in_{proj}$ and $out_{proj}$ layers take the dominant portion (e.g.,
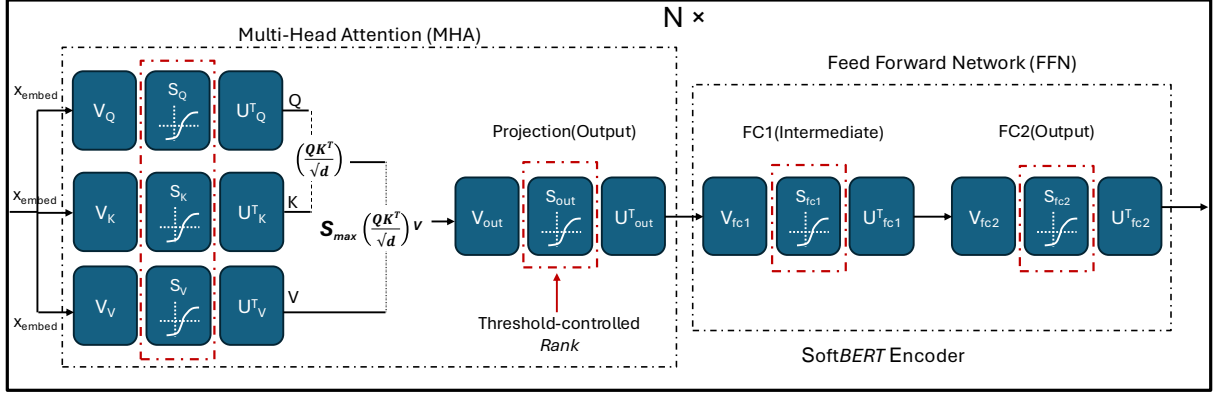
Figure 2: SoftBERT encoder with $N$ blocks, where $W_Q$, $W_K$, $W_V$, $W_{proj}$ in Multi-Head Attention and $W_{fc1}$ and $W_{fc2}$ in the Feed Forward Network are substituted with U, S and V. The module S employs the Soft Threshold function to clamp the singular values in $\Sigma$, enabling dynamic rank for each block adaptively in fine-tuning.

94%) of the parameters in each block whereas the number of parameters inside SSM is significantly less. Thus, we focus on these two layers in this paper.

## 4 Methodology

This section introduces SoftLM based on the proposed SVD compression method with a differentiable thresholding called *Soft Threshold* as shown in Figure 2, which provides an example on BERT. A detailed algorithm and training methodology are provided to adaptively compress the linear layers, achieving an optimal balance between accuracy and compression.

### 4.1 Compressed Linear Layer

We first replace the standard linear layer with a decomposed linear layer including three modules, $U$, $\Sigma$ and $V$. Modules for $U$ and $V$ share the same structure, which is similar to the vanilla SVD. A soft threshold layer $\mathcal{T}h_s$ (described in Section 4.3) is applied on the $\Sigma$. The conventional linear operation of $y = x \cdot W^T$ is transformed as:

$$y = x \cdot V \cdot \mathcal{T}h_s(\Sigma) \cdot U^T \tag{5}$$

where $x$ is the input, $W$ is the weight matrix, and $y$ is the output.

### 4.2 Learnable Threshold on Singular Values

The descending order of singular values in the $\Sigma$ matrix by magnitude reflects the decreasing importance of the corresponding singular vectors in the SVD matrices (Figure 3). Our technique employs a gradient-based learnable threshold ($\alpha$) to select a limited set of most consequential singular values.
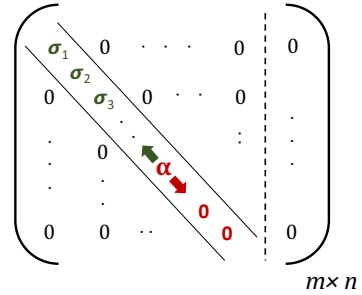


Figure 3: Learnable threshold ($\alpha$) for selecting top singular values in $\Sigma$; values below ($\alpha$) are set to zero.

This threshold traverses diagonally across the $\Sigma$ matrix during the training so that values in $\Sigma$ below the threshold are replaced by zero. A higher threshold corresponds to a greater level of compression. Consequently, the corresponding rows and columns of the orthogonal matrices $U$ and $V$, respectively, are also neglected. This selective removal of singular vectors allows effective dimensionality reduction. The threshold is set as a learnable parameter and is distinct for each layer to be adjusted given the varying importance of each linear layer.

### 4.3 Soft Threshold

It is known that non-differentiable operations hinder the back-propagation process during training, as gradients cannot flow through such operations. However, the threshold function is inherently non-differentiable, which poses challenges in adapting the $\alpha$ during training. To facilitate the gradient flow, we employ a soft thresholding (Figure 4) as follows:

$$\mathcal{T}h_s(x) = \begin{cases} x\tanh(s(x-\alpha)), & \text{if } x \geq \alpha \\ c\tanh(s(x-\alpha)), & \text{if } x < \alpha \end{cases} \tag{6}$$

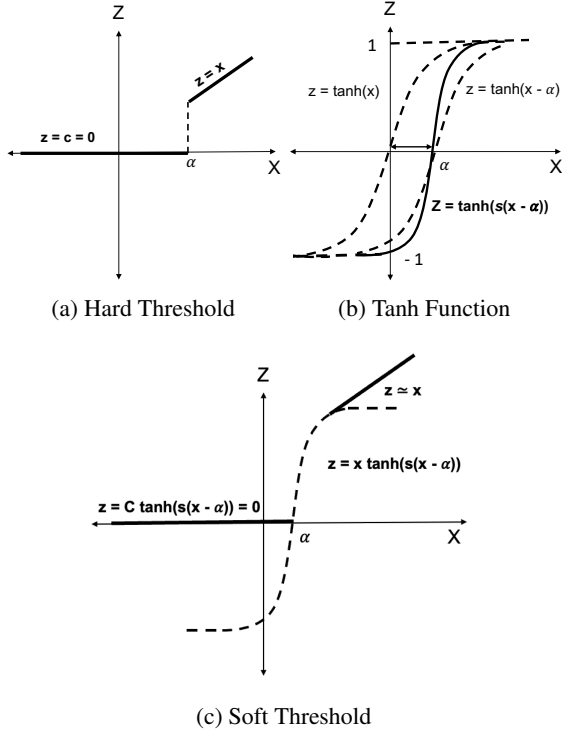(a) Hard Threshold

(b) Tanh Function



(c) Soft Threshold

Figure 4: Threshold functions. (a) conventional non-differentiable thresholding, (b) shifted Tanh with sharpness control factor $s$, and (c) differentiable soft thresholding combining above functions

where $x$ represents the input, $s$ is a sharpness control parameter, $\alpha$ is the cut-off value, and constant $c$ replaces values less than the threshold in the matrix. $c$ is chosen to be zero in our case. This function acts as a smooth approximation of the hard thresholding, allowing gradients to propagate through the network during the training process. Specifically, when $x$ exceeds the cut-off value $\alpha$, the soft thresholding function yields $x$ multiplied by the tanh term. Conversely, for $x$ less than $\alpha$, the function yields $c$, effectively replacing values below the threshold to zero.

### 4.4 Adaptive Loss for Compression

To maximize the compression and the performance at the same time, we employ the total loss function defined as:

$$\mathcal{L}_{tot} = \mathcal{L}_{acc} + \gamma \cdot \mathcal{L}_{cmp} \tag{7}$$

where $\mathcal{L}_{acc}$ is a cross entropy loss to maximize the performance while $\mathcal{L}_{cmp}$ is the compression loss and $\gamma$ is the balancing factor for regularization. For $\mathcal{L}_{cmp}$, we simply sum the thresholds as:
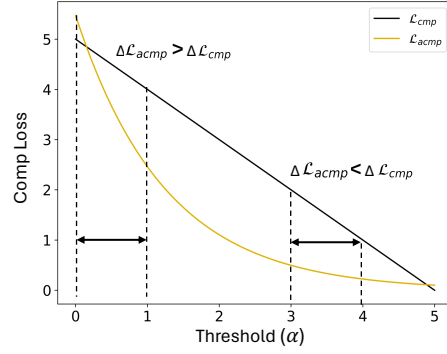
$$\mathcal{L}_{cmp} = -\sum_i \alpha_i \tag{8}$$



Figure 5: The magnitude of compression loss over the threshold $\alpha$ with a conventional linear loss ($\mathcal{L}_{cmp}$) and the proposed adaptive loss ($\mathcal{L}_{acmp}$).

Here, $\alpha_i$ is the learnable threshold for the $i$-th block. While the $\alpha_i$ is set to be zero at the beginning of fine-tuning, this loss term gradually increases the $\alpha_i$ to truncate more singular values.

However, a linear loss term continuously pressures for more compression at any rank. Ideally, the loss should adjust based on the fine-tuning phase, i.e., a larger loss is preferable at the beginning to accelerate compression whereas the compression should slow down to allow for performance recovery once the model is sufficiently compressed. To enable such a capability, we introduce an adaptive compression loss ($\mathcal{L}_{acmp}$) defined as follows:

$$\mathcal{L}_{acmp} = \sum_i e^{-\alpha_i} \tag{9}$$

By exploiting the exponential term, the loss reduces to be negligible once the $\alpha_i$ increases enough (Figure 5). The following total loss is employed in this paper by reflecting the $\mathcal{L}_{acmp}$ as:

$$\mathcal{L}_{tot} = \mathcal{L}_{acc} + \gamma \cdot \mathcal{L}_{acmp} \tag{10}$$

### 4.5 Fine-tuning Algorithm

In this section, we describe the detailed procedure of fine-tuning. It should be noted that the learnable threshold ($\alpha$) is frozen once the total parameter count reaches the target value by changing the learning rate to be zero for the $\alpha$s. At the end of fine-tuning, we merge $S = \mathcal{T}h_s(\Sigma)$ and $V$ into a matrix to further save memory.

## 5 Experiments

### 5.1 Models and Datasets

Extensive experiments are conducted on a diverse set of language models and datasets to evaluate the efficacy of the proposed method. We applied the

| Model | #Para | MRPC | COLA | QNLI | QQP | SST2 | G-Avg |
|---|---|---|---|---|---|---|---|
| BERT-Base (w/o *emb*) | 86M | | | | | | |
| + *embeddings* | 110M | 88.0 | 56.2 | 91.3 | 87.8 | 93.0 | 83.3 |
| DistilBERT | 67M | 87.5 | 51.3 | 89.2 | 86.7 | 91.3 | 81.20 |
| MiniLMv2 | 67M | 89.1 | 43.3 | **90.6** | 86.7 | 91.4 | 80.2 |
| BERT6-KD | 67M | 86.2 | - | 88.3 | - | 91.5 | - |
| BERT6-PKD | 67M | 85.0 | - | 89.0 | - | **92.0** | - |
| SVD-BERT | 67M | 86.4 | 40.5 | 89.0 | 86.5 | 90.1 | 78.5 |
| SoftBERT (w/o *emb*) | 43M (–50%) | | | | | | |
| + *embeddings* | 67M | **90.7** | **54.6** | 89.5 | **87.6** | 90.4 | **82.6** |

Table 1: A comparison of SoftBERT with prior works in terms of parameter count and performance (reported under the name of datasets). G-Avg: the average performance for GLUE tasks.

---

**Algorithm: Adaptive Rank using SoftThreshold**

**Require:** $P$: Target parameter number

    Initialize $\alpha_i \leftarrow 0$
    SVD: $U\Sigma V^T \leftarrow W$
    SoftThreshold: $\mathcal{T}h_s(\Sigma) \leftarrow \Sigma$
    **while** #param $> P$ **do**
        Compute $\mathcal{L}_{tot} = \mathcal{L}_{acc} + \gamma \cdot \mathcal{L}_{acmp}$
        Optimizer step (AdamW)
    **end while**
    Freeze threshold $(\alpha_i)$ parameters
    **for** $epoch$ in $(0, N_{epoch})$ **do**
        Compute $\mathcal{L}_{tot} = \mathcal{L}_{acc} + \gamma \cdot \mathcal{L}_{acmp}$
        Optimizer step (AdamW)
    **end for**
    // End-of-fine-tuning
    Compute $VS = V \cdot \mathcal{T}h_s(\Sigma)$
    Store $(U^T, VS)$

**Ensure:** #param in ($U$ and $VS$) < #param in $W$

---

proposed technique on models: BERT-Base (Devlin et al., 2018), GPT2-medium (Radford et al., 2019), Mamba-130M (Gu and Dao, 2024) and TinyLlamav1.1 (Zhang et al., 2024) to produce SoftBERT, SoftGPT2, SoftMamba and SoftTinyLlama, respectively. BERT-Base is evaluated on classification tasks such as MRPC, COLA, SST-2, QNLI, and QQP selected from the popular GLUE benchmark (Wang et al., 2019a). While MRPC and QQP were evaluated using the F1-score, COLA was assessed using the Matthews Correlation Coefficient (MCC). SST-2 and QNLI were evaluated based on accuracy. Similarly, Mamba-130M is evaluated on the GLUE datasets. In Table 1, the target compression ratio (CR) for SoftBERT is chosen as 0.50 for a fair evaluation and comparison convenience with other references. Table 2 shows the varying performance of SoftBERT with

respect to the compression ratios for the five GLUE datasets. GPT2-medium and TinyLlama are evaluated on WikiText-2 (wik) by measuring the perplexity score.

| #Para (encoder) | CR | G-Avg |
|---|---|---|
| 86M | 0.00 | 83.1 |
| 64M | 0.25 | 83.0 |
| 43M | 0.50 | 82.6 |
| 22M | 0.75 | 78.5 |

Table 2: Performance variation (GLUE-Average) of SoftBERT with respect to the compression ratio (CR).

## 5.2 Experiment Results and Comparison with Prior Works

For BERT-Base, the six linear layers ($W_Q$, $W_K$, $W_V$, and $W_{proj}$ in the attention module as well as $W_{fc1}$ and $W_{fc2}$ in the FFN) in each block sum up to 72 layers across the 12 blocks of the entire encoder, accounting for approximately 86 million parameters. The remaining 23 million parameters are allocated for the token, positional, and segment embeddings. In Table 1, the compression within the encoder is set to CR=0.50, reducing the 86 million parameters inside the encoder to 43 million. The parameters for all the embeddings remain untouched. We used the following baselines and benchmarks for comparison with our method: DistilBERT (Sanh et al., 2020), MiniLMv2 (Wang et al., 2020), BERT6-KD, PKD (Sun et al., 2019), and vanilla SVD applied to BERT-Base with fine-tuning. The final model size for each compared model is carefully adjusted to be the same (67M) to ensure a fair comparison. Our method achieves an average performance score of 82.56 in Table 1, outperforming all the other methods, with a performance drop from the original uncompressed model limited to just 0.7. It is observed that the rank of

| Model | #Para | COLA | SST2 | QNLI | QQP | Avg |
|---|---|---|---|---|---|---|
| S4 | 131M | 23.0 | 87.0 | 72.14 | 89.6 | 67.9 |
| Mamba (*in/out_proj*) | 85M | | | | | |
| + *misc* | 130M | 80.5 | 92.0 | 88.5 | 89.0 | 87.5 |
| SoftMamba (*in/out_proj*) | 42 M (–50%) | | | | | |
| + *misc* | 87M | 77.3 | 90.3 | 88.1 | 89.6 | 86.3 |

Table 3: Task performance (reported under the name of datasets) and parameter count of SoftMamba on GLUE datasets and comparison with S4s and Mamba.
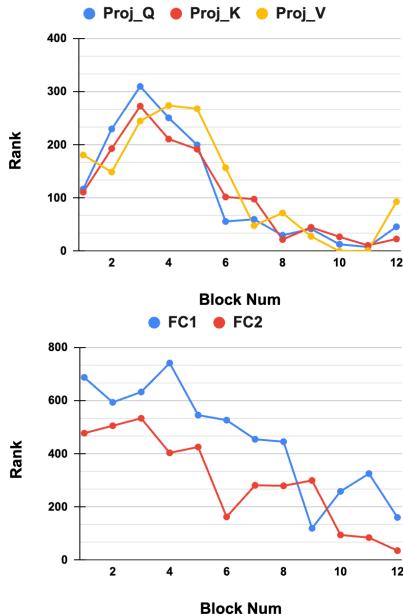


Figure 6: Rank of each block after fine-tuning Soft-BERT on MRPC in MHA (top) and FFN (bottom).

| Model | #Para(decoder) | wiki-2 | CR |
|---|---|---|---|
| GPT2-medium | 301M | 16.7 | 0.00 |
| SVD-GPT2-medium | 227M | 23.0 | 0.25 |
| SoftGPT2-medium | 227M | 18.6 | 0.25 |
| SVD-GPT2-medium | 151M | 32.1 | 0.50 |
| SoftGPT2-medium | 151M | 19.6 | 0.50 |
| SVD-GPT2-medium | 76M | 50.6 | 0.75 |
| SoftGPT2-medium | 76M | 35.4 | 0.75 |

Table 4: Performance variation of SoftGPT2-medium in perplexity with respect to CR in comparison with static-rank decomposition.

| Model | #Para(decoder) | wiki-2 | CR |
|---|---|---|---|
| TinyLlamav1.1 | 968M | 8.9 | 0.00 |
| SVD-TinyLlama | 726M | 13.8 | 0.25 |
| SoftTinyLlama | 726M | 12.0 | 0.25 |
| SVD-TinyLlama | 484M | 20.4 | 0.50 |
| SoftTinyLlama | 484M | 15.1 | 0.50 |
| SVD-TinyLlama | 242M | 39.2 | 0.75 |
| SoftTinyLlama | 242M | 19.1 | 0.75 |

Table 5: Performance variation of SoftTinyLlama in perplexity with respect to CR in comparison with static-rank decomposition.

each linear layer within the Multi-Head Attention (MHA) and Feed Forward Network (FFN) varies and is learned to balance the performance and the degree of compression (Figure 6). The initial layers demonstrate higher importance and sensitivity to the loss compared to those in the middle and at the end. Consequently, the rank tends to be preserved in the initial layers.

There are four linear layers and a 1-dimensional convolution layer in each block of Mamba-130m. The number of parameters for 24 blocks amounts to 90M, with approximately 85M contributed by the $in_{proj}$ and $out_{proj}$ layers. Thus, we compress these two layers by 50%. SoftMamba achieves comparable performance with an average degradation of merely 1.2% compared to the original Mamba model (Table 3).

Similarly, each block of GPT-2 contains four 1-dimensional convolution layers, accounting for 301M parameters out of the total 345M across

24 blocks. Table 4 compares the performance of SoftGPT2-medium against the GPT2-medium compressed with static-SVD, which applies the same compression ration across all blocks. To further extend our analysis on a bigger model with total parameters greater than one billion, we employ TinyLlama to validate our technique in Table 5. We observe a significant performance improvement in the aforementioned models using dynamic rank-based compression compared to static rank-based compression, highlighting the importance of adaptive low-rank approximations.

## 5.3 Execution Cost Reductions

In this section, we provide detailed insights about resource utilization by SoftLMs.

### 5.3.1 Savings in Memory and Power

The proposed technique, evaluated at CR=0.50, leads to overall savings of 33-44%, depending on the model. For BERT-base, which initially has a 440 MB weights file, this corresponds to a reduction of 168 MB, assuming all parameters are in float32 format. This technique reduces not only the memory footprint but also the computational complexity, with Multiply-Accumulate (MAC) operations decreasing linearly with the compression ratio.
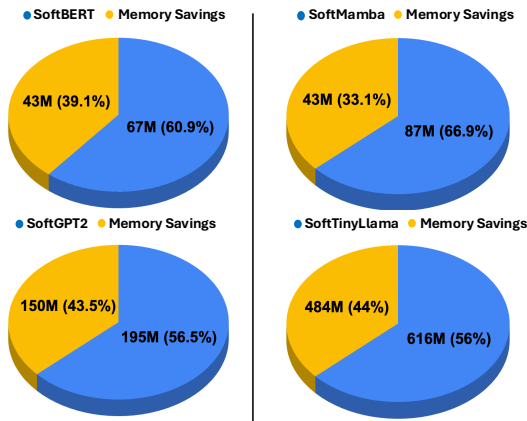


Figure 7: Memory savings in SoftLMs at CR=0.50.

The power consumption is measured in Watts (W) using the NVIDIA pyNVML library (mlp) while evaluating the SoftLMs at inference. The reduction in power consumption for BERT and GPT-2 ranges from 30% to 35%, while for Mamba and TinyLlama, it varies between 15% and 20%.

### 5.3.2 Latency Reduction

Latency reduction is measured as the speed-up of the compressed models during inference. Figure 8 shows the average speed-up for SoftBERT, SoftTinyLlama, SoftGPT2, and SoftMamba, each compressed at a CR of 0.50. The speed-up ranges from $1.33\times$ to $1.72\times$, with the highest speed-up achieved in BERT, followed by TinyLlama, GPT-2, and Mamba. We attribute the lower speed-up in Mamba to the presence of the SSM block, which is not accelerated by our technique. Additionally, we observe that latency remains similar for a given model, with minimal differences across datasets.



Figure 8: Inference speed-up achieved by SoftLMs.



Figure 9: Performance variation with respect to CR in SoftBERT and SoftMamba with QNLI dataset.

### 5.4 Trade-off in Performance vs. Compression

Figure 9 illustrates the impact on performance with respect to the compression ratio for the attention-based BERT and SSM-based Mamba. Although the baseline accuracy of BERT is higher than that of Mamba, the impact of compression is less significant in Mamba, leading to smaller performance degradation at higher compression levels. For this reason, Mamba is a more feasible model for deployment under stringent resource constraints, particularly due to the absence of Multi-Head Attention (MHA) and its lower performance degradation by the compression.

## 6 Conclusion

We introduce a novel compression method that is easily deployable on any deep learning model with linear layers, adding minimal complexity to the fine-tuning process. By adaptively learning the rank for each layer using a simple thresholding mechanism, our method achieves an optimal balance between task performance and compression. We validate the effectiveness of the proposed technique across various language models, demonstrating that resource-efficient SoftLMs outperform state-of-the-art compression methods.

# 7 Limitations

Although the model is compressed at the end of fine-tuning, the initial decomposition increases its size during the fine-tuning process. This may be less of an issue for smaller models such as BERT and GPT-2, but it presents a challenge when training larger language models. Additionally, while SVD is applied only once during fine-tuning, it is computationally expensive, adding significant complexity to fine-tuning large models. To mitigate this, one could replace the costly SVD with more efficient decomposition methods that begin directly from the truncated rank.

# References

pyNVML: Python bindings to the NVIDIA Management Library. https://pythonhosted.org/nvidia-ml-py/. Accessed: 2024-5-3.

The WikiText Long Term Dependency Language Modeling Dataset. https://blog.salesforceairesearch.com. Accessed: 2024-5-3.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Preprint*, arXiv:2305.14314.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Gene H Golub and Christian Reinsch. 1971. Singular value decomposition and least squares solutions. In *Handbook for Automatic Computation: Volume II: Linear Algebra*, pages 134–151. Springer.

Mitchell A Gordon, Kevin Duh, and Nicholas Andrews. 2020. Compressing bert: Studying the effects of weight pruning on transfer learning. *arXiv preprint arXiv:2002.08307*.

Albert Gu and Tri Dao. 2024. Mamba: Linear-time sequence modeling with selective state spaces. *Preprint*, arXiv:2312.00752.

Song Han, Huizi Mao, and William J. Dally. 2016. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *Preprint*, arXiv:1510.00149.

Song Han, Jeff Pool, John Tran, and William J. Dally. 2015. Learning both weights and connections for efficient neural networks. *Preprint*, arXiv:1506.02626.

Yen-Chang Hsu, Ting Hua, Sungen Chang, Qian Lou, Yilin Shen, and Hongxia Jin. 2022. Language model compression with weighted low-rank factorization. *Preprint*, arXiv:2207.00112.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *Preprint*, arXiv:2106.09685.

Ting Hua, Xiao Li, Shangqian Gao, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. 2023. Dynamic low-rank estimation for transformer-based language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9275–9287, Singapore. Association for Computational Linguistics.

Matan Ben Noach and Yoav Goldberg. 2020. Compressing pre-trained language models by matrix decomposition. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 884–889.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Tara N. Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. 2013. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6655–6659.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *Preprint*, arXiv:1910.01108.

Steffen Schotthöfer, Emanuele Zangrando, Jonas Kusch, Gianluca Ceruti, and Francesco Tudisco. 2022. Low-rank lottery tickets: finding efficient low-rank neural networks via matrix differential equations. *Preprint*, arXiv:2205.13571.

Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. 2019. Patient knowledge distillation for bert model compression. *Preprint*, arXiv:1908.09355.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019a. Glue: A multi-task benchmark and analysis platform for natural language understanding. *Preprint*, arXiv:1804.07461.

Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. 2019b. Haq: Hardware-aware automated quantization with mixed precision. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8604–8612.

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Preprint*, arXiv:2002.10957.

Xin Wang, Yu Zheng, Zhongwei Wan, and Mi Zhang. 2024. Svd-llm: Truncation-aware singular value decomposition for large language model compression. *Preprint*, arXiv:2403.07378.

Zhihang Yuan, Yuzhang Shang, Yue Song, Qiang Wu, Yan Yan, and Guangyu Sun. 2024. Asvd: Activation-aware singular value decomposition for compressing large language models. *Preprint*, arXiv:2312.05821.

Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. 2019. Q8bert: Quantized 8bit bert. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition (EMC2-NIPS)*, pages 36–39. IEEE.

Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model. *Preprint*, arXiv:2401.02385.

Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adalora: Adaptive budget allocation for parameter-efficient fine-tuning. *Preprint*, arXiv:2303.10512.

# A   Appendix

## A.1   System Setup

The system setup that we used includes a NVIDIA GeForce RTX 3090 GPU (total RAM of 24 GB with maximum power limit of 350W) on which experiments for BERT and Mamba were performed. Experiments for TinyLlama and GPT2-medium were performed on a single NVIDIA A100 GPU (total RAM of 80 GB and maximum power limit of 400W). We do not use any distributed or parallel training methods.

## A.2   Implementation Details

### A.2.1   Datasets

The General Language Understanding Evaluation (GLUE) benchmark is a collection of diverse natural language understanding (NLU) tasks designed to evaluate the performance of models on several language understanding challenges. GLUE consists of various tasks including linguistic acceptability, sentence similarity, textual entailment, and natural language inference.

The benchmark consists of nine tasks, out of which we have evaluated on five. CoLA (Corpus of Linguistic Acceptability) involves determining whether a given sentence is grammatically correct in English. SST-2 (Stanford Sentiment Treebank) requires models to classify the sentiment of a given sentence as either positive or negative. MRPC (Microsoft Research Paraphrase Corpus) identifies whether two sentences in a pair are semantically equivalent. QQP (Quora Question Pairs) determines whether two questions posted on Quora are paraphrases of each other. QNLI (Question Natural Language Inference) determines whether a given sentence (question) is answerable by another sentence (context).

The evaluation metric for QNLI and SST2 is accuracy, whereas it is F1 score for MRPC and QQP, and Matthews Correlation Coefficient (MCC) for CoLA.

Wikitext-2 is chosen as the generative task which is a set of almost 100 million tokens extracted from Wikipedia's articles. The performance is evaluated through the exponential of validation loss, also commonly referred as the perplexity score.

### A.2.2   Training Setup

We employed the following hyper-parameters in evaluating BERT and Mamba on GLUE tasks.

| | |
|---|---|
| **Batch Size** | 32 |
| **Learning Rate (LR)** | 2e-5 |
| **LR (threshold $\alpha$)** | 1e-2, 1e-3 |
| **Regularizer** ($w$) | 0.001, 0.01, 0.1 |
| **Optimizer** | AdamW |

Table 6: Hyper-parameters for fine-tuning.

For GPT2 and TinyLlama, we used similar settings except the batch size which is set to four for GPT2 and two for TinyLlama due to limited GPU RAM.

# B   SoftLM Rank Distribution

SoftBERT has the following rank distribution for the six linear layers - $W_Q$, $W_K$, $W_V$, and $W_{proj}$ in the attention module, $W_{fc1}$ and $W_{fc2}$ in the FFN as illustrated in Figures 10 - 14 where $x$-axis is the block number. These results indicate that the initial layers have higher importance than later ones. The layers in the FFN have four times higher parameter counts than the ones in the attention leading to higher ranks. It is also noted that the $W_{fc2}$ achieves better compression than $W_{fc1}$.
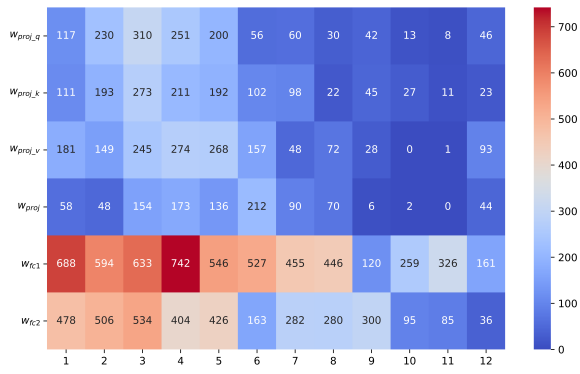
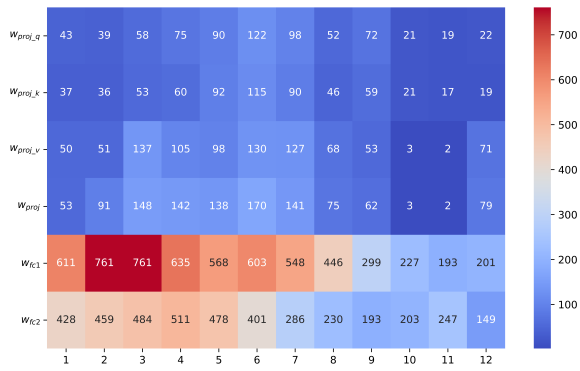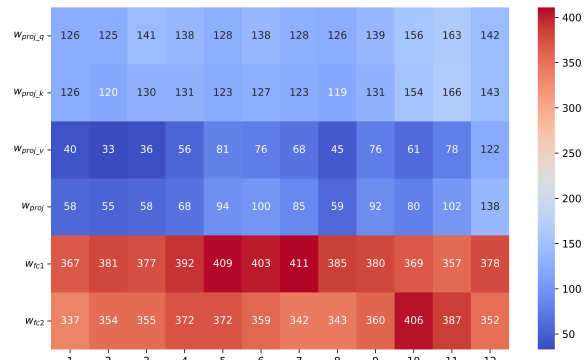| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $w_{proj\_q}$ | 117 | 230 | 310 | 251 | 200 | 56 | 60 | 30 | 42 | 13 | 8 | 46 |
| $w_{proj\_k}$ | 111 | 193 | 273 | 211 | 192 | 102 | 98 | 22 | 45 | 27 | 11 | 23 |
| $w_{proj\_v}$ | 181 | 149 | 245 | 274 | 268 | 157 | 48 | 72 | 28 | 0 | 1 | 93 |
| $w_{proj}$ | 58 | 48 | 154 | 173 | 136 | 212 | 90 | 70 | 6 | 2 | 0 | 44 |
| $w_{fc1}$ | 688 | 594 | 633 | 742 | 546 | 527 | 455 | 446 | 120 | 259 | 326 | 161 |
| $w_{fc2}$ | 478 | 506 | 534 | 404 | 426 | 163 | 282 | 280 | 300 | 95 | 85 | 36 |

Figure 10: MPRC.



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $w_{proj\_q}$ | 43 | 39 | 58 | 75 | 90 | 122 | 98 | 52 | 72 | 21 | 19 | 22 |
| $w_{proj\_k}$ | 37 | 36 | 53 | 60 | 92 | 115 | 90 | 46 | 59 | 21 | 17 | 19 |
| $w_{proj\_v}$ | 50 | 51 | 137 | 105 | 98 | 130 | 127 | 68 | 53 | 3 | 2 | 71 |
| $w_{proj}$ | 53 | 91 | 148 | 142 | 138 | 170 | 141 | 75 | 62 | 3 | 2 | 79 |
| $w_{fc1}$ | 611 | 761 | 761 | 635 | 568 | 603 | 548 | 446 | 299 | 227 | 193 | 201 |
| $w_{fc2}$ | 428 | 459 | 484 | 511 | 478 | 401 | 286 | 230 | 193 | 203 | 247 | 149 |

Figure 11: CoLA.



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $w_{proj\_q}$ | 99 | 67 | 80 | 77 | 83 | 106 | 86 | 83 | 80 | 83 | 84 | 69 |
| $w_{proj\_k}$ | 100 | 75 | 84 | 75 | 78 | 89 | 77 | 74 | 83 | 85 | 82 | 68 |
| $w_{proj\_v}$ | 99 | 57 | 68 | 69 | 71 | 83 | 66 | 40 | 49 | 24 | 15 | 68 |
| $w_{proj}$ | 88 | 70 | 65 | 56 | 78 | 82 | 64 | 36 | 39 | 39 | 5 | 46 |
| $w_{fc1}$ | 635 | 497 | 532 | 531 | 467 | 456 | 452 | 411 | 404 | 323 | 298 | 321 |
| $w_{fc2}$ | 405 | 404 | 406 | 409 | 391 | 355 | 325 | 311 | 307 | 353 | 323 | 319 |

Figure 12: SST-2.



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $w_{proj\_q}$ | 49 | 85 | 163 | 183 | 101 | 70 | 66 | 56 | 74 | 151 | 123 | 46 |
| $w_{proj\_k}$ | 47 | 75 | 150 | 181 | 98 | 66 | 62 | 61 | 74 | 153 | 137 | 57 |
| $w_{proj\_v}$ | 62 | 34 | 61 | 117 | 78 | 56 | 55 | 56 | 72 | 43 | 14 | 66 |
| $w_{proj}$ | 57 | 65 | 93 | 102 | 90 | 89 | 65 | 43 | 76 | 23 | 14 | 4 |
| $w_{fc1}$ | 546 | 615 | 669 | 506 | 474 | 447 | 490 | 526 | 423 | 360 | 276 | 182 |
| $w_{fc2}$ | 395 | 415 | 443 | 386 | 364 | 354 | 292 | 337 | 303 | 303 | 214 | 185 |

Figure 13: QNLI.



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $w_{proj\_q}$ | 126 | 125 | 141 | 138 | 128 | 138 | 128 | 126 | 139 | 156 | 163 | 142 |
| $w_{proj\_k}$ | 126 | 120 | 130 | 131 | 123 | 127 | 123 | 119 | 131 | 154 | 166 | 143 |
| $w_{proj\_v}$ | 40 | 33 | 36 | 56 | 81 | 76 | 68 | 45 | 76 | 61 | 78 | 122 |
| $w_{proj}$ | 58 | 55 | 58 | 68 | 94 | 100 | 85 | 59 | 92 | 80 | 102 | 138 |
| $w_{fc1}$ | 367 | 381 | 377 | 392 | 409 | 403 | 411 | 385 | 380 | 369 | 357 | 378 |
| $w_{fc2}$ | 337 | 354 | 355 | 372 | 372 | 359 | 342 | 343 | 360 | 406 | 387 | 352 |

Figure 14: QQP.