

# DyLoRA: Parameter-Efficient Tuning of Pretrained Models using Dynamic Search-Free Low Rank Adaptation

Mojtaba Valipour<sup>1,2</sup> Mehdi Rezagholizadeh<sup>2</sup> Ivan Kobyzev<sup>2</sup> Ali Ghodsi<sup>1</sup>

{mojtaba.valipour, ali.ghodsi}@uwaterloo.ca, {mehdi.rezagholizadeh, ivan.kobyzev}@huawei.com

1: University of Waterloo, 2: Huawei Noah’s Ark Lab

## Abstract

With the ever-growing size of pretrained models (PMs), fine-tuning them has become more expensive and resource-hungry. As a remedy, low-rank adapters (LoRA) keep the main pretrained weights of the model frozen and just introduce some learnable truncated SVD modules (so-called LoRA blocks) to the model. While LoRA blocks are parameter-efficient, they suffer from two major problems: first, the size of these blocks is fixed and cannot be modified after training (for example, if we need to change the rank of LoRA blocks, then we need to re-train them from scratch); second, optimizing their rank requires an exhaustive search and effort. In this work, we introduce a dynamic low-rank adaptation (DyLoRA) technique to address these two problems together. Our DyLoRA method trains LoRA blocks for a range of ranks instead of a single rank by sorting the representation learned by the adapter module at different ranks during training. We evaluate our solution on different natural language understanding (GLUE benchmark) and language generation tasks (E2E, DART and WebNLG) using different pretrained models such as RoBERTa and GPT with different sizes. Our results show that we can train dynamic search-free models with DyLoRA at least 4 to 7 times (depending to the task) faster than LoRA without significantly compromising performance. Moreover, our models can perform consistently well on a much larger range of ranks compared to LoRA.

## 1 Introduction

Pre-training/fine-tuning has become a popular paradigm for solving many tasks in natural language processing (NLP) (Devlin et al., 2018; Liu et al., 2019; Brown et al., 2020) and Computer Vision (Simonyan and Zisserman, 2014; He et al., 2016; Howard et al., 2019; Bochkovskiy et al.,

2020; Chen et al., 2020; Dosovitskiy et al., 2020). pretrained models (PMs) such as pretrained language models (PLMs) (Devlin et al., 2018; Brown et al., 2020), and pretrained visual-language models (Lu et al., 2019; Li et al., 2019; Su et al., 2019; Xia et al., 2021) have advanced a lot in recent years. With the ever-growing size of these pretrained models, fine-tuning them on downstream tasks becomes more expensive. Moreover, as the ratio of the number of parameters of models with respect to the labeled data increases, the fine-tuning process will be more prone to overfitting (Karimi Mahabadi et al., 2021). There are two categories of solutions: first, model compression (Jafari et al., 2021; Chen et al., 2021); second, parameter-efficient tuning (PET) (Houlsby et al., 2019a; Karimi Mahabadi et al., 2021; Mao et al., 2021).

There are many different model compression techniques in the literature for Transformer-based models such as matrix factorization (Noach and Goldberg, 2020; Tahaei et al., 2021), pruning (Wang et al., 2019), quantization (Tao et al., 2022; Prato et al., 2020), and knowledge distillation (Hinton et al., 2015; Li et al., 2021; Jafari et al., 2021; Passban et al., 2021; Rashid et al., 2021). There are also different types of PET techniques in the literature such as low-rank adapters (Wang et al., 2020; Karimi Mahabadi et al., 2021; Houlsby et al., 2019b; Hu et al., 2021b), and prompt-based techniques (Lester et al., 2021).

Although model compression solutions are well-established in recent years in the literature, applying them to large language models can be very costly, because compression techniques usually need to train (or fine-tune) the original large model. A case in point is knowledge distillation which relies on fine-tuning a large teacher model or even pre-training the student model as suggested in (Jiao et al., 2019). Moreover, using compression techniques usually leads to degrading the model performance. PETs can be alternatives to the compres-

<sup>1</sup>[github.com/huawei-noah/KD-NLP/tree/main/DyLoRA](https://github.com/huawei-noah/KD-NLP/tree/main/DyLoRA)

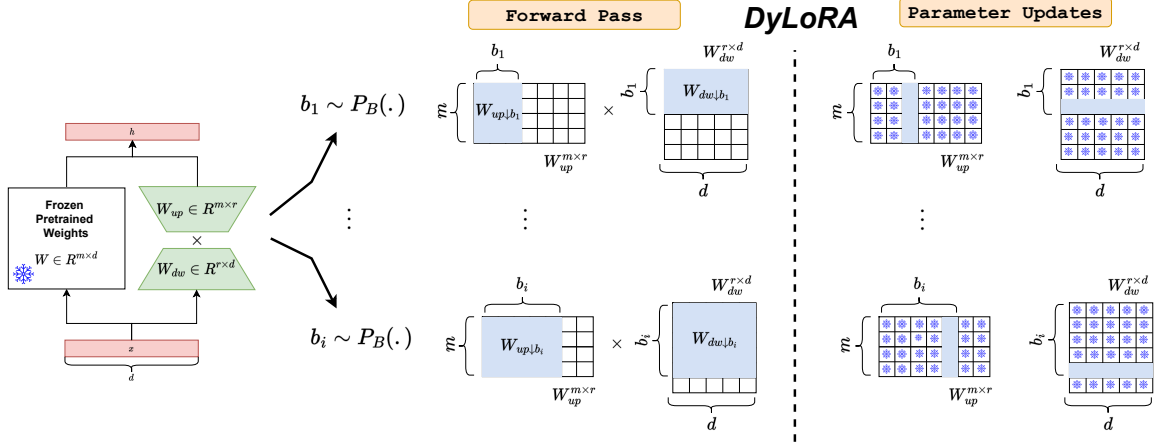


Figure 1: DyLoRA: The overall diagram of our proposed method. In each iteration, we sample from a pre-defined random distribution which will help us to truncate the up-projection and down-projection matrices in the LoRA (Hu et al., 2021a) objective.

sion methods, especially when we would like to use the full capacity of the large pretrained models with light training efforts (such as the *language-model-as-a-service* scenario (Sun et al., 2022)). Among PET techniques, low-rank adapters have received much attention because, in contrast to prompt-tuning techniques, low-rank adapters do not add to the sequence length, get trained faster, and perform better (Karimi Mahabadi et al., 2021). Even though there are several low-rank adaptation techniques in the literature, such as Adapter (Houlsby et al., 2019b), Compacter (Karimi Mahabadi et al., 2021), and LoRA (Hu et al., 2021b); they all suffer from two major common problems: first, it is not clear how to select the size of their rank (while their performance is very sensitive to this rank selection); second, their training is static which means that if a low-rank model is trained based on a particular rank size, it will not work well in other rank values (i.e. for any other rank value we need to train a separate model).

This paper proposes a dynamic low-rank adapter technique (DyLoRA) to address these two problems. Without loss of generality, we focus on LoRA (Hu et al., 2021a) and train LoRA blocks for a range of ranks instead of a single rank by sorting out the representation learned at different ranks during training. While our model is more flexible, it can outperform LoRA in a much wider range of ranks without adding to the training time. Moreover, our technique does not need extra training for searching across ranks. We summarize our contributions in the following:

- **Dynamic LoRA:** On top of LoRA, we developed a new algorithm (DyLoRA) that makes it dynamic at inference time without incurring extra costs.
- **Search-free LoRA:** We demonstrate that by making a negligible compromise in performance, it is possible to avoid the costly search process of choosing the optimal rank for LoRA.

## 2 Related Work

This section reviews low-rank adaptation techniques for parameter-efficient tuning and potential existing solutions to make these techniques dynamic and search-free.

It has been shown in (Aghajanyan et al., 2020) that for classification tasks such as natural language understanding (NLU), PLMs have a low intrinsic dimension. This observation motivates the use of low-rank adapters for parameter-efficient tuning. There are several low-rank adapters in the literature such as LoRA (Hu et al., 2021b), Adapter (Houlsby et al., 2019b), Compacter (Karimi Mahabadi et al., 2021), and Parallel Adapter (PA) (He et al., 2021). LoRA is a low-rank up-projection/down-projection transformation without any non-linearity applied in parallel to key and value attention matrices. The main benefit of LoRA is that the adapter module, after training, can be integrated into the original weight matrices of the model, which in turn can lead to a very efficient inference time. Adapters also have a low-rank up-projection/down-projection transformation with an intermediate non-

linearity. The Adapter module is applied in series with the feed-forward network (FFN). Having the adaptor module in-line with other blocks in the model can increase the inference time of the model. PA is a faster version of the Adapter, which can be applied in parallel with the FFN block. The compactor is a more memory-efficient version of the Adapter, which deploys the sum of Kronecker products to reconstruct each up-projection and down-projection matrices. All these low-rank adapters suffer from two major issues: first, finding the best rank requires heavy exhaustive training and search; second, the tuned adapter module works well only with a particular rank.

While there have been some efforts in the literature towards dynamic networks such as DynaBERT (Hou et al., 2020) and GradMax (Evcı et al., 2022), to the best of our knowledge, this problem for factorized networks and low-rank adapters is still open. DRONE (Chen et al., 2021) propose a technique for data-aware low-rank model compression however their approach is not search-free, and also, it is not dynamic. DynaBERT introduces a two-stage method to train width and depth-wise dynamic networks. However, DynaBERT requires a fine-tuned teacher model on the task to train its sub-networks which makes it unsuitable for PET techniques. GradMax is a technique that gradually adds to the neurons of a network without touching the already trained neurons. But it is unclear how GradMax can be deployed to alleviate the rank-search problem in low-rank adapters. Wang et al. (2019) propose a structured pruning technique called factorized low-rank pruning (FLOP). FLOP decomposes weight matrices of a network into the sum of rank-1 components, which are regularized during training to gain sparsity. It is worth mentioning that FLOP aims at compressing the main model, and even if it can be used for finding a good rank in the lower-rank representation of full-weight matrices, the final low-rank model will not be dynamic (i.e. it is trained well only for one rank and not a range of ranks, same as LoRA.). In this paper, we propose a new methodology for training low-rank modules for multiple ranks simultaneously rather than training a single-rank adapter at a time (without changing the training budget). Inspired by the idea of *nested dropout* (Rippel et al., 2014), we pursue ordering the representations of the bottleneck at the low-rank adapter modules with a new recipe. To the best of our knowledge, it is the first

time that the concept of ordering representations has been deployed in training PLMs.

### 3 Background

#### 3.1 Nested Dropout

Inspired by the dropout (Hinton et al., 2012), nested drop-out (Rippel et al., 2014) is a stochastic regularization technique that targets enforcing ordered representations in training auto-encoders. The nested dropout, adds an implicit bias (which does not exist in dropout) to favor order in training. For example, in dropout, we can randomly drop any nodes or units in the network, but in nested dropout, if we randomly select  $k^{\text{th}}$  unit, then we keep all the units indexed from 1 to  $k$  and drop the units with indices larger than  $k$ . Therefore, nested dropout tends toward accommodating more important information in lower indices while learning representations.

Following the notations of (Rippel et al., 2014), nested dropout assumes an auto-encoder mapping of  $N$  training examples  $\{y_i\}_{i=1}^N \in Y$ ,  $Y \subset \mathbb{R}^D$  to their corresponding representations  $\{x_i\}_{i=1}^N \in X$ ,  $X \subset \mathbb{R}^K$  using the function  $f_\theta : Y \rightarrow X$  with parameters  $\theta$ ; and then decoding these representations using another function  $g_\psi : X \rightarrow Y$  with parameters  $\psi$  to reconstruct the inputs. The reconstruction loss can be defined as follows:

$$C(\theta, \psi) = \sum_{i=1}^N \|y_i - g_\psi(f_\theta(y_i))\|^2. \quad (1)$$

Suppose we want to randomly drop some units in our representation vector  $x$ . In this regard, we sample a random variable  $b \sim p_B(\cdot)$ ,  $b \in \{1, 2, \dots, K\}$  from a pre-defined categorical distribution  $p_B(\cdot)$  and truncate the functions  $f_\theta$  and  $g_\psi$  to keep their corresponding units indexed from 1 to  $b$  and dropping  $b+1$  to  $K$  indices. Let's define the  $b$ -truncated version of the vector  $x$  as  $x_{\downarrow b}$  and the  $b$ -truncated version of the functions  $f_\theta$  and  $g_\psi$  as  $f_{\theta_{\downarrow b}}$  and  $g_{\psi_{\downarrow b}}$  respectively. In this case, the reconstruction loss is redefined for the  $b$ -truncated model as follows:

$$C(\theta, \psi) = \mathbb{E}_{p_B}[C_{\downarrow b}(\theta, \psi)] = \sum_{b=1}^K p_B(b) C_{\downarrow b}(\theta, \psi)$$

where

$$C_{\downarrow b}(\theta, \psi) = \sum_{i=1}^N \|y_i - g_{\psi_{\downarrow b}}(f_{\theta_{\downarrow b}}(y_i))\|^2. \quad (2)$$

In the final stage, the parameters of this model can be obtained by solving the following optimization problem.

$$(\theta^*, \psi^*) = \underset{\theta, \psi}{\operatorname{argmin}} C(\theta, \psi). \quad (3)$$

While our work in this paper is inspired by the feature of ordering information suggested in nested dropout, we can distinguish our work from nested dropout in several aspects:

1. The nested dropout technique is used to add order information to a vector representation; however, we are adding order information to the low-rank matrix decomposition to make it work across a range of ranks instead of a single rank.
2. Our training algorithm differs from nested dropout in the choice of the distribution function  $p_B(\cdot)$ , and we propose a more efficient individual loss for each truncated matrix compared to the linear summation loss (check equations 2 and 11 in the original paper (Rippl et al., 2014)) in nested dropout. The original proposal for the nested dropout was to use a batch with mixed truncated examples. To enhance efficiency and resolve suboptimality, we propose to fix truncation in the entire batch as part of our approach.

### 3.2 LoRA: Low-rank Adapters

In LoRA (Hu et al., 2021a), some pretrained weights of dense layers of PLMs are summed with parallel linear low-rank adapter modules. During fine-tuning, the original pretrained weights are kept frozen; LoRA modules can be updated instead. For example, let’s assume that  $W_0 \in \mathbb{R}^{m \times d}$  is a pretrained weight matrix in the network which is accompanied by a LoRA module  $\Delta W = W_{up}W_{dw}$  where  $W_{up} \in \mathbb{R}^{m \times r}$ ,  $W_{dw} \in \mathbb{R}^{r \times d}$ , and  $r \ll \min(m, d)$ . Then, the output of this layer can be obtained as

$$h = W_0x + \Delta Wx = W_0x + \frac{\alpha}{r}W_{up}W_{dw}x. \quad (4)$$

Bear in mind that the  $W_{up}$  matrix is initialized as a zero matrix, and the  $W_{dw}$  matrix is initialized as a zero-mean Gaussian distribution where  $\alpha$  is a constant scale hyper-parameter.

In LoRA, the rank  $r$  is a hyperparameter that should be tuned for each task. Moreover, LoRA is a *static* low-rank adapter that works only with a particular size of  $r$ , which has been trained on it.

## 4 Our Method: DyLoRA

In this section, we introduce our solution to get dynamic low-rank adapters that can be trained and deployed well on a range of ranks instead of a single particular rank (with a fixed training budget). This flexibility can free us from searching for the best ranks by training the model multiple times.

Without loss of generality, we explain our solution on top of LoRA as one of the prominent low-rank adapter techniques in the literature. In each LoRA module, we have an up-projection ( $W_{up} \in \mathbb{R}^{m \times r}$ ) and a down-projection matrix ( $W_{dw} \in \mathbb{R}^{r \times d}$ ). Let’s assume that we would like to train the LoRA module to operate in the range of  $r \in \text{Range}[r_{min}, r_{max}]$  where  $r_{min}$  and  $r_{max}$  can be treated as new hyper-parameters. To make the LoRA module work in a range of ranks instead of a single rank, we need to ensure that increasing or decreasing the rank will not significantly hamper the model’s performance. One way to implement such behavior would be by sorting the information content of different ranks in the training process of LoRA modules. In this regard, at each training step, we sample  $b \sim p_B(\cdot)$ ,  $b \in \{r_{min}, r_{min} + 1, \dots, r_{max}\}$  form a pre-defined categorical distribution (which has a support in  $\text{Range}[r_{min}, r_{max}]$ ) and truncate  $W_{dw}$  and  $W_{up}$  matrices accordingly.

$$\begin{aligned} W_{dw \downarrow b} &= W_{dw}[1 : b, :] \\ W_{up \downarrow b} &= W_{up}[:, 1 : b] \end{aligned} \quad (5)$$

$W_{dw \downarrow b}$  and  $W_{up \downarrow b}$  are b-truncated versions of  $W_{dw}$  and  $W_{up}$  respectively (see Fig. 1 for the visualization). Moreover, let’s define  $W_{dw}^b$  as the  $b^{\text{th}}$  row of  $W_{dw}$ ;  $W_{up}^b$  corresponds to the  $b^{\text{th}}$  column of  $W_{up}$ .

$$\begin{aligned} W_{dw}^b &= W_{dw}[b, :] \\ W_{up}^b &= W_{up}[:, b] \end{aligned} \quad (6)$$

Then, the forward pass of this truncated LoRA module during training will be calculated as following:

$$h = W_0x + \frac{\alpha}{b}W_{up \downarrow b}W_{dw \downarrow b}x \quad (7)$$

For simplicity, let’s assume that we have only one LoRA module in the network (the one which is described in Eq. 7). Let’s first consider the regular static loss function ( $\mathcal{L}^S$ ) of the network  $f(x; W_{dw}, W_{up})$  with  $W_{dw}$  and  $W_{up}$  tunable parameters for  $N$  given input-output pairs  $(\mathbf{x}, \mathbf{y}) =$

$(x_i, y_i)_{i=1}^N$ :

$$\min_{W_{dw}, W_{up}} \mathcal{L}^S(\mathbf{x}, \mathbf{y}; W_{dw}, W_{up}) \triangleq \sum_{i=1}^N l(f(x_i; W_{dw}, W_{up}), y_i). \quad (8)$$

where  $l(f, \mathbf{y})$  is a loss function that measures the divergence of network predictions compared with the target labels. Then, let's extend the training loss to make the network dynamic considering the b-truncation process. We can define our dynamic loss function  $\mathcal{L}^{\mathcal{DY}}$  as follows.

$$\mathcal{L}_{\downarrow b}^{\mathcal{DY}} = \sum_{i=1}^N l(f(x_i; W_{dw\downarrow b}, W_{up\downarrow b}), y_i). \quad (9)$$

Bear in mind that, our loss function has a major difference from the nested dropout loss, which makes it more efficient. The nested dropout loss is in the form of  $\sum_{b=r_{min}}^{r_{max}} p_B(b) \mathcal{L}_{\downarrow b}^{\mathcal{DY}}(\mathbf{x}, \mathbf{y}; W_{dw\downarrow b}, W_{up\downarrow b})$  which requires to sum the loss over the entire possible range of ranks and it is computationally expensive. To overcome this computational restriction, we replace it by optimizing the model parameters for each target rank individually at each time step. We show that this scheme quite works well.

The other difference with nested dropout is that in the parameter update phase, we add a new mode (so-called *frozen*) as a hyper-parameter to our training. This new mode suggests to only update the  $b^{\text{th}}$  corresponding row and column sampled in the truncation phase (i.e. a single row or column will be updated at a time to prevent the learning parameters from being forgotten at previous time steps.). With a minor performance cost, this approach can improve the efficiency of our algorithm even further.

$$\begin{aligned} W_{dw}^b &\leftarrow W_{dw}^b - \eta \nabla_{W_{dw}^b} \mathcal{L}_{\downarrow b}^{\mathcal{DY}} \\ W_{up}^b &\leftarrow W_{up}^b - \eta \nabla_{W_{up}^b} \mathcal{L}_{\downarrow b}^{\mathcal{DY}} \end{aligned} \quad (10)$$

Table 4 shows the impact of only updating "b" versus updating the columns and rows from 1 to  $b$ . The summary of our technique is described in Algorithm 1.

## 5 Experiments

In this section, we describe the experiments used to evaluate our DyLoRA model on both natural language understanding (NLU) and natural language

---

### Algorithm 1 DyLoRA - Training

---

**Require:**

$r \in \text{Range}[r_{min}, r_{max}]$ ;  $i$ : the number of training iterations;  $\alpha$ : a scaling factor;  $p_B$ : probability distribution function for rank selection;  $X \in \mathbb{R}^{d \times n}$ : all input features to LoRA;  $W_0 \in \mathbb{R}^{m \times d}$  the original frozen pretrained weight matrix

**Require:**  $W_{dw} \in \mathbb{R}^{r \times d}$ ;  $W_{up} \in \mathbb{R}^{m \times r}$ , **FROZEN**: whether to keep the lower ranks frozen when updating the higher ranks

**while**  $t < i$  **do**:

**Forward:**

    // sample a specific rank, during test is given

$b \sim p_B(\cdot)$

    // truncate down-projection matrix

$W_{dw\downarrow b} = W_{dw}[:, b:]$

$W_{dw}^b = W_{dw}[b, :]$

    // truncate up-projection matrix

$W_{up\downarrow b} = W_{up}[:, :b]$

$W_{up}^b = W_{up}[:, b]$

    // calculate the LoRA output

$h = W_0 X + \frac{\alpha}{b} W_{up\downarrow b} W_{dw\downarrow b} X$

**Backward:**

**if** **FROZEN** **then**

    // only update the unique parameters of the selected rank

$W_{dw}^b \leftarrow W_{dw}^b - \eta \nabla_{W_{dw}^b} \mathcal{L}_{\downarrow b}^{\mathcal{DY}}$

$W_{up}^b \leftarrow W_{up}^b - \eta \nabla_{W_{up}^b} \mathcal{L}_{\downarrow b}^{\mathcal{DY}}$

**else**

$W_{dw\downarrow b} \leftarrow W_{dw\downarrow b} - \eta \nabla_{W_{dw\downarrow b}^b} \mathcal{L}_{\downarrow b}^{\mathcal{DY}}$

$W_{up\downarrow b} \leftarrow W_{up\downarrow b} - \eta \nabla_{W_{up\downarrow b}^b} \mathcal{L}_{\downarrow b}^{\mathcal{DY}}$

**end if**

**end while**

---

generation (NLG) tasks. To be fair with the original LoRA method, we try to keep the setting of our experiments similar to the LoRA paper (Hu et al., 2021a). Therefore similarly, we chose the pretrained RoBERTa (Liu et al., 2019) base model as the backbone of the LoRA and DyLoRA experiments for the GLUE benchmark (Development Set), and GPT-Medium for the NLG tasks. For our experiments, we did not use any hyper-parameter tuning, nor did we search the validation epochs, nor did we use MLNI trick (use the MLNI checkpoint instead of the pretrained weights) to enhance the model's performance. More details about the hyper-parameters is available in Table 8 in Appendix B. In total, we conducted more than 200 experiments and evaluated more than 1600 models, details of

Model: RoBERTa-Base						
Task	Rank=1	Rank=2	Rank=4	Rank=8	Rank=16	Rank=32
QQP (Accuracy)	<u>89.14</u>	89.96	90.33	90.69	90.95	<b>91.02</b>
SST-2 (Accuracy)	<u>93.58</u>	94.15	94.38	<b>94.84</b>	94.27	94.5
MRPC (Accuracy)	87.25	<u>87.75</u>	88.24	<u>87.25</u>	<u>86.76</u>	<b>89.22</b>
CoLA (Mathews)	61.84	<u>57.78</u>	61.57	<b>63.81</b>	63.07	62.82

Table 1: The effect of the rank of the low-rank adaptation matrix over the performance of the model. In this experiment, all the other hyperparameters are fixed, and we only changed the rank of the LoRA model. In this search space, Underline shows the minimum performance rank, and the **bold** number shows the maximum performance rank.

Model	Accuracy MNLI	Accuracy SST-2	F1 MRPC	Mathews CoLA	Accuracy QNLI	Accuracy QQP	Accuracy RTE	Pearson STS-B	Avg
Rank = 1									
LoRA	34.60 $\pm$ 3.69	69.61 $\pm$ 7.99	83.47 $\pm$ 3.90	25.57 $\pm$ 9.71	53.00 $\pm$ 2.95	44.30 $\pm$ 7.50	57.55 $\pm$ 5.51	76.07 $\pm$ 6.06	54.90
DyLoRA (Frozen)	85.36 $\pm$ 0.26	93.51 $\pm$ 0.49	90.75 $\pm$ 0.70	56.95 $\pm$ 1.54	91.70 $\pm$ 0.28	87.87 $\pm$ 0.17	66.79 $\pm$ 8.54	89.95 $\pm$ 0.24	82.86
DyLoRA	85.59 $\pm$ 0.07	93.23 $\pm$ 0.63	91.58 $\pm$ 0.69	57.93 $\pm$ 2.12	91.95 $\pm$ 0.14	88.37 $\pm$ 0.15	74.80 $\pm$ 1.48	90.30 $\pm$ 0.13	<b>84.22</b>
Rank = 2									
LoRA	40.53 $\pm$ 6.17	82.75 $\pm$ 5.08	88.00 $\pm$ 1.81	43.30 $\pm$ 4.67	63.42 $\pm$ 2.99	59.21 $\pm$ 6.13	68.88 $\pm$ 1.26	85.51 $\pm$ 1.94	66.45
DyLoRA (Frozen)	85.74 $\pm$ 0.28	93.76 $\pm$ 0.52	91.09 $\pm$ 0.45	56.88 $\pm$ 2.09	92.03 $\pm$ 0.22	88.21 $\pm$ 0.07	63.90 $\pm$ 12.85	90.25 $\pm$ 0.15	82.73
DyLoRA	86.02 $\pm$ 0.06	93.81 $\pm$ 0.30	91.66 $\pm$ 0.46	59.91 $\pm$ 1.88	92.39 $\pm$ 0.25	89.33 $\pm$ 0.05	76.03 $\pm$ 1.61	90.60 $\pm$ 0.09	<b>84.97</b>
Rank = 3									
LoRA	58.95 $\pm$ 6.02	90.00 $\pm$ 1.27	89.66 $\pm$ 1.25	56.78 $\pm$ 1.88	79.26 $\pm$ 4.80	72.58 $\pm$ 4.09	72.49 $\pm$ 2.30	88.80 $\pm$ 0.29	76.07
DyLoRA (Frozen)	85.78 $\pm$ 0.25	93.76 $\pm$ 0.26	91.78 $\pm$ 0.89	58.86 $\pm$ 0.32	92.17 $\pm$ 0.18	88.40 $\pm$ 0.0	70.90 $\pm$ 6.14	90.50 $\pm$ 0.29	84.02
DyLoRA	86.70 $\pm$ 0.09	94.11 $\pm$ 0.33	91.56 $\pm$ 0.86	60.97 $\pm$ 2.01	92.77 $\pm$ 0.21	89.76 $\pm$ 0.07	77.11 $\pm$ 2.97	90.69 $\pm$ 0.14	<b>85.46</b>
Rank = 4									
LoRA	72.10 $\pm$ 5.25	91.56 $\pm$ 0.34	89.62 $\pm$ 0.92	58.53 $\pm$ 3.93	85.09 $\pm$ 1.20	80.78 $\pm$ 3.73	73.07 $\pm$ 2.29	89.28 $\pm$ 0.72	80.00
DyLoRA (Frozen)	85.93 $\pm$ 0.19	93.85 $\pm$ 0.33	91.28 $\pm$ 0.71	59.25 $\pm$ 1.05	92.27 $\pm$ 0.16	88.52 $\pm$ 0.08	71.12 $\pm$ 2.46	90.53 $\pm$ 0.18	84.10
DyLoRA	86.82 $\pm$ 0.04	94.40 $\pm$ 0.13	92.06 $\pm$ 0.46	59.81 $\pm$ 1.71	92.91 $\pm$ 0.31	89.80 $\pm$ 0.10	77.40 $\pm$ 2.72	90.86 $\pm$ 0.06	<b>85.53</b>
Rank = 5									
LoRA	78.61 $\pm$ 3.97	92.82 $\pm$ 0.46	90.75 $\pm$ 0.96	60.37 $\pm$ 3.10	88.97 $\pm$ 0.90	85.26 $\pm$ 1.56	73.21 $\pm$ 2.17	89.90 $\pm$ 0.30	82.49
DyLoRA (Frozen)	85.95 $\pm$ 0.17	93.78 $\pm$ 0.26	91.28 $\pm$ 0.64	59.41 $\pm$ 1.30	92.30 $\pm$ 0.17	88.56 $\pm$ 0.09	71.48 $\pm$ 2.92	90.60 $\pm$ 0.20	84.17
DyLoRA	87.00 $\pm$ 0.10	94.29 $\pm$ 0.41	91.73 $\pm$ 0.60	60.52 $\pm$ 1.07	93.01 $\pm$ 0.28	90.04 $\pm$ 0.10	76.90 $\pm$ 2.11	90.97 $\pm$ 0.20	<b>85.56</b>
Rank = 6									
LoRA	83.02 $\pm$ 1.59	93.49 $\pm$ 0.88	91.28 $\pm$ 0.63	61.94 $\pm$ 2.27	90.32 $\pm$ 0.76	87.54 $\pm$ 1.51	76.68 $\pm$ 1.16	90.12 $\pm$ 0.12	84.30
DyLoRA (Frozen)	85.98 $\pm$ 0.16	93.76 $\pm$ 0.46	91.12 $\pm$ 0.43	58.95 $\pm$ 1.10	92.46 $\pm$ 0.14	88.68 $\pm$ 0.13	72.64 $\pm$ 2.44	90.64 $\pm$ 0.23	84.28
DyLoRA	86.97 $\pm$ 0.20	94.27 $\pm$ 0.37	91.44 $\pm$ 0.64	60.16 $\pm$ 1.70	93.01 $\pm$ 0.21	90.07 $\pm$ 0.14	77.33 $\pm$ 1.66	91.03 $\pm$ 0.20	<b>85.53</b>
Rank = 7									
LoRA	85.44 $\pm$ 0.78	93.62 $\pm$ 0.35	91.27 $\pm$ 0.73	62.19 $\pm$ 2.66	91.88 $\pm$ 0.23	89.51 $\pm$ 0.30	75.52 $\pm$ 1.41	90.35 $\pm$ 0.24	84.97
DyLoRA (Frozen)	86.08 $\pm$ 0.14	93.97 $\pm$ 0.17	91.02 $\pm$ 0.70	58.76 $\pm$ 0.94	92.30 $\pm$ 0.10	88.77 $\pm$ 0.06	73.50 $\pm$ 1.67	90.68 $\pm$ 0.15	84.38
DyLoRA	86.82 $\pm$ 0.10	94.27 $\pm$ 0.33	91.38 $\pm$ 0.59	59.51 $\pm$ 1.75	92.99 $\pm$ 0.26	90.04 $\pm$ 0.06	77.91 $\pm$ 1.58	91.07 $\pm$ 0.19	<b>85.50</b>
Rank = 8									
LoRA	86.82 $\pm$ 0.18	94.01 $\pm$ 0.30	91.48 $\pm$ 0.73	62.08 $\pm$ 1.37	92.39 $\pm$ 0.39	90.42 $\pm$ 0.02	74.51 $\pm$ 0.41	90.48 $\pm$ 0.24	85.27
DyLoRA (Frozen)	86.10 $\pm$ 0.04	93.69 $\pm$ 0.41	91.19 $\pm$ 0.79	58.52 $\pm$ 0.95	92.47 $\pm$ 0.18	88.82 $\pm$ 0.06	73.29 $\pm$ 2.49	90.68 $\pm$ 0.14	84.35
DyLoRA	86.76 $\pm$ 0.13	94.36 $\pm$ 0.38	91.38 $\pm$ 0.83	59.51 $\pm$ 1.84	93.00 $\pm$ 0.32	89.91 $\pm$ 0.08	77.55 $\pm$ 0.59	91.05 $\pm$ 0.19	<b>85.44</b>
Best (Rank)									
LoRA	87.03(8)	94.50(6)	92.25(7)	<b>66.05(7)</b>	92.81(8)	<b>90.45(8)</b>	77.98(6)	90.87(8)	86.49
DyLoRA (Frozen)	86.18(7)	94.50(2)	<b>92.93(3)</b>	61.57(5)	92.70(6)	88.88(8)	75.81(7)	90.89(6)	85.43
DyLoRA	<b>87.17(6)</b>	<b>94.72(7)</b>	92.79(8)	63.32(3)	<b>93.56(8)</b>	90.17(6)	<b>80.14(4)</b>	<b>91.36(7)</b>	<b>86.66</b>
Full Rank									
Fine Tune*	<b>87.6</b>	<b>94.8</b>	90.2	63.6	92.8	<b>91.9</b>	78.7	91.2	86.4

Table 2: In this table, the task is to find a low-rank adaptation matrix that works with different ranks at inference time given a fixed budget (training time).

which can be found in the attachments.

## 5.1 Baselines

- **Fine Tune:** To show a relative upper bound for the performance of our proposed method, we fine-tuned all the parameters in the model.
- **LoRA:** As a baseline to DyLoRA, we employed the original LoRA model with their tuned hyperparameters (Hu et al., 2021a). As

Even though we have a large number of trainable parameters, this can help us better understand how higher-rank models perform.

Model (Rank)	Trainable Params	Accuracy		F1	Accuracy	Pearson
		SST-2	MRPC	QNLI	STS-B	AVERAGE
Fine Tune*	125M	<b>94.8</b>	90.2	92.8	<b>91.2</b>	92.25
FLOP*	80M	92.09	88.61	89.05	88.18	89.48
LoRA (1)	<b>0.628M</b>	93.58	91.93	91.98	90.85	92.09
Maximum Rank: $r_{max} = 8$						
<b>DyLoRA (1)</b>	<b>0.628M</b>	93.23 $\pm$ 0.63	91.58 $\pm$ 0.69	91.95 $\pm$ 0.14	90.30 $\pm$ 0.13	91.77
<b>DyLoRA (8)</b>	0.887M	94.36 $\pm$ 0.38	91.38 $\pm$ 0.83	93.00 $\pm$ 0.32	91.05 $\pm$ 0.19	<b>92.45</b>

Table 3: This table compares DyLoRA with compression-based algorithms. As indicated by \*, we reported "Fine Tune" and FLOP from their original papers, (Liu et al., 2019) and (Wang et al., 2019). To the best of our knowledge, experiments were conducted under the same experimental setting. We count all the trainable parameters including classifier, unlike LoRA paper (Hu et al., 2021a) which they count only LoRA specific parameters.

		Maximum Rank: $r_{max} = 8$						
$b \sim P_B$ : Distribution	Updated Parameters	Accuracy	F1	Mathews	Accuracy	Accuracy	Pearson	
		SST-2	MRPC	CoLA	QNLI	RTE	STS-B	AVERAGE
		<b>Rank=8</b>						
<b>Geometric (p=0.15)</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	93.97 $\pm$ 0.33	90.84 $\pm$ 1.15	58.95 $\pm$ 1.95	92.74 $\pm$ 0.13	74.80 $\pm$ 0.90	90.66 $\pm$ 0.15	83.66
	$W_{dw}^b, W_{up}^b$	93.60 $\pm$ 0.24	90.50 $\pm$ 0.42	58.19 $\pm$ 1.17	92.26 $\pm$ 0.12	71.91 $\pm$ 1.74	90.20 $\pm$ 0.36	82.78
<b>Uniform</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	94.36 $\pm$ 0.38	91.38 $\pm$ 0.83	59.51 $\pm$ 1.84	93.00 $\pm$ 0.32	77.59 $\pm$ 0.59	91.05 $\pm$ 0.19	<b>84.47</b>
	$W_{dw}^b, W_{up}^b$	93.69 $\pm$ 0.41	91.19 $\pm$ 0.79	58.52 $\pm$ 0.95	92.47 $\pm$ 0.18	73.29 $\pm$ 2.49	90.68 $\pm$ 0.14	83.31
		<b>Rank=1</b>						
<b>Geometric (p=0.15)</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	93.53 $\pm$ 0.47	91.36 $\pm$ 0.72	59.43 $\pm$ 1.12	92.24 $\pm$ 0.08	73.65 $\pm$ 3.55	90.33 $\pm$ 0.14	<b>83.42</b>
	$W_{dw}^b, W_{up}^b$	93.58 $\pm$ 0.26	90.81 $\pm$ 0.83	58.55 $\pm$ 1.13	92.27 $\pm$ 0.28	68.52 $\pm$ 11.88	90.60 $\pm$ 0.31	82.39
<b>Uniform</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	93.23 $\pm$ 0.63	91.58 $\pm$ 0.69	57.93 $\pm$ 2.12	91.95 $\pm$ 0.14	74.80 $\pm$ 1.48	90.30 $\pm$ 0.13	83.30
	$W_{dw}^b, W_{up}^b$	93.51 $\pm$ 0.49	90.75 $\pm$ 0.70	56.95 $\pm$ 1.54	91.70 $\pm$ 0.28	66.79 $\pm$ 8.54	89.95 $\pm$ 0.24	81.61

Table 4: Ablation Study - In this experiment, our goal is to demonstrate how the introduced distribution can affect the performance of DyLoRA.

a result, most of the experiments have been conducted in a favorable manner for LoRA.

- **FLOP**: Due to its flexibility, **Factorized Low Rank Pruning (FLOP)** (Wang et al., 2019) can be applied to any matrix multiplication and, therefore, can be used to avoid the search in our problem. However, this baseline lacks the dynamic properties of DyLoRA. We used it to show regularization-based techniques' performance and pros and cons.

## 5.2 LoRA rank selection problem

There is no clear guidance on how to determine the rank for the LoRA algorithm. It is evident in the LoRA paper (Hu et al., 2021a) that the performance of models varies a lot with different ranks (e.g. check Tables 15, and 18 in the LoRA paper), and does not indicate any clear trend. We also observe the same problem in the GLUE benchmark. We may argue that theoretically, the rank with the best performance is always the highest. High ranks, however, introduce additional parameters into the adaptive process and this might be undesirable. In practice, as demonstrated in Table 1, the most effective rank differs depending on the task. For

example, based on the MRPC results, the rank with the lowest performance is 16 while the rank with the highest performance is 32. This is different from SST-2, in which rank 1 is the least performing rank and rank 8 is the most effective rank. Many factors can contribute to this difference, including but not limited to the size of the dataset, hyperparameter selections, hardware configurations and the optimization.

## 5.3 Dynamic low rank adaptation

For example, suppose we have a neural network that we wish to deploy on various devices with different configurations. The use of higher ranks may pose a problem for very sensitive devices as they have a greater number of parameters. Therefore, we must either train several models with different configurations or find the most optimal rank. The cost associated with this is significant, as even in the setting of LoRA, we are required to find the best rank for each task and each device. Using DyLoRA, however, one needs to train one model per task and, as our method is adaptive at inference time, we can deploy it according to our needs. In Table 2, we demonstrate the dynamic properties of

DyLoRA. In order to ensure a fair comparison, all LoRA and DyLoRA models in this table have the same model size, we used the same code and evaluation process, and all models were trained to the same extent. In LoRA, we lose performance when performing inferences for the lower ranks. This occurs because the model has been trained only for rank 8 during training. In DyLoRA, we preserve a high level of performance for lower ranks while competing well with LoRA on rank 8.

Model	Time	SST-2 ( $r$ )	MRPC ( $r$ )
Maximum Rank: $r_{max} = 64$			
LoRA (Search)	7x	<b>95.3</b> (64)	89.71(64)
<b>DyLoRA (Frozen)</b>	<b>1x</b>	94.38(7)	<b>89.95(34)</b>
Maximum Rank: $r_{max} = 32$			
LoRA (Search)	6x	<b>94.84</b> (32)	88.73(16)
<b>DyLoRA (Frozen)</b>	<b>1x</b>	94.38(7)	<b>89.71(5)</b>

Table 5: In this table, the search space of rank is larger compared to the previous experiment and the goal is to find the most optimal rank for the low-rank adaptation of a pre-trained RoBERTa-Base. For LoRA (Search), we ran experiments for ranks=1,2,4,8,16,32,64 and we reported the best results. In the Exhaustive Search, one has to search all the ranks from 1 to 64, which means it will cost 64 times more than our proposed method. The lower the rank the better, and the higher the performance is the better.

#### 5.4 Search-free low rank adaptation

The process of selecting a particular rank can be expensive as previously mentioned. In Table 5, we present an experiment that illustrates the costs associated with such a search for LoRA and DyLoRA. As an example, if one naively wanted to search the entire range of ranks (for example, 64 in the experiment), then they would have to train and evaluate 64 distinct models in order to determine the proper rank. It becomes even more expensive if one search the entire rank space. In the case of uniform search, this cost is less, yet still more expensive (7 times in the experiment) than our proposed method. Therefore, for LoRA (Search), we ran experiments for ranks=1,2,4,8,16,32,64 and we reported the best results. The results demonstrate that our proposed method performs competitively at a much lower cost.

#### 5.5 Robustness of DyLoRA

As illustrated in Table 2, DyLoRA is quite robust to randomness and can produce consistently good results due to stable convergence.

#### 5.6 Regularization and Pruning

An alternative method of avoiding the search problem is using regularization/pruning techniques to determine the intrinsic rank of the weight matrix. In this way, we can reduce the number of parameters of the original matrices; however, we will not have a dynamic model during inference. To illustrate the difference between such methods and DyLoRA, we reported the performance of one of these models, FLOP (Wang et al., 2019), in Table 3. FLOP utilizes low-rank factorization to create new matrices representing the original weight matrix. Thus, they will have fewer total parameters but require more trainable parameters to reach a comparable performance to DyLoRA.

#### 5.7 Generative Tasks

In this experiment, we evaluate the performance of our model on different natural language generation (NLG) tasks such as the E2E NLG Challenge (Novikova et al., 2017), DART (Nan et al., 2020) and WebNLG (Gardent et al., 2017). The results of the E2E task are shown in Table 6 and due to the space limit, the results of the other two tasks are demonstrated in Appendix C. The generative tasks demonstrate a similar pattern as the NLU task, showing that our model is able to work well at wider range of ranks compared to LoRA.

#### 5.8 Ablation study

In this subsection, we investigate the impact of two design choices in DyLoRA: first, the new distribution  $P_B$  hyper-parameter in our technique; second, the impact of updating  $W_{dw}^b$  and  $W_{up}^b$  parameters instead of the entire  $W_{dw\downarrow b}$  and  $W_{up\downarrow b}$ . The distribution  $P_B$  changes the relative importance of the different ranks during the training process. To examine the impact of the chosen distribution on DyLoRA’s performance, we used two distributions, geometric and uniform. As shown in Table 4, the geometric distribution, provides a much better method for optimizing the lower ranks, since it pays much more attention to the lower ranks during training, and uniform distribution will give better performance over all ranks. We chose to use uniform distribution in most of our experiments to avoid adding another hyperparameter which is a requirement of the geometric distribution. Moreover, we demonstrate that it is possible to ensure that the optimization of rank  $b$  will not negatively affect the performance of the lower ranks (1 to  $b - 1$ ), while



Model (Method)	Updated Params	Trainable Params	E2E NLG Challenge				
			BLEU	NIST	MET	ROUGE-L	CIDEr
			Rank=1				
GPT-2 M (LoRA)		0.09M	3.38	1.18	9.23	18.79	0.12
<b>GPT-2 M (DyLoRA)</b>	$W_{dw}^b, W_{up}^b$	0.09M	67.92 $\pm$ 0.20	8.65 $\pm$ 0.06	44.91 $\pm$ 0.38	69.07 $\pm$ 0.32	2.38 $\pm$ 0.04
<b>GPT-2 M (DyLoRA)</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	0.09M	68.86 $\pm$ 0.55	8.72 $\pm$ 0.04	45.81 $\pm$ 0.40	70.33 $\pm$ 0.64	2.43 $\pm$ 0.04
			Rank=2				
GPT-2 M (LoRA)		0.19M	46.99	6.39	34.19	56.10	1.27
<b>GPT-2 M (DyLoRA)</b>	$W_{dw}^b, W_{up}^b$	0.19M	68.81 $\pm$ 0.49	8.75 $\pm$ 0.02	45.23 $\pm$ 0.22	69.81 $\pm$ 0.30	2.41 $\pm$ 0.01
<b>GPT-2 M (DyLoRA)</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	0.19M	68.97 $\pm$ 1.03	8.75 $\pm$ 0.07	45.88 $\pm$ 0.55	70.07 $\pm$ 0.86	2.43 $\pm$ 0.04
			Rank=3				
GPT-2 M (LoRA)		0.29M	63.68	8.46	42.37	65.84	2.24
<b>GPT-2 M (DyLoRA)</b>	$W_{dw}^b, W_{up}^b$	0.29M	68.41 $\pm$ 1.00	8.69 $\pm$ 0.10	45.31 $\pm$ 0.64	69.75 $\pm$ 0.69	2.42 $\pm$ 0.02
<b>GPT-2 M (DyLoRA)</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	0.29M	69.33 $\pm$ 0.26	8.76 $\pm$ 0.05	46.19 $\pm$ 0.22	70.56 $\pm$ 0.43	2.46 $\pm$ 0.01
			Rank=4				
GPT-2 M (LoRA)		0.39M	69.88	8.81	46.81	72.10	2.53
<b>GPT-2 M (DyLoRA)</b>	$W_{dw}^b, W_{up}^b$	0.39M	68.36 $\pm$ 0.41	8.70 $\pm$ 0.02	45.46 $\pm$ 0.56	69.91 $\pm$ 0.50	2.43 $\pm$ 0.01
<b>GPT-2 M (DyLoRA)</b>	$W_{dw\downarrow b}, W_{up\downarrow b}$	0.39M	69.19 $\pm$ 0.43	8.75 $\pm$ 0.03	46.26 $\pm$ 0.47	70.78 $\pm$ 0.63	2.46 $\pm$ 0.02
			Fine-Tune				
GPT-2 M (FT)*		354M	68.2	8.62	46.2	71.0	2.5

Table 6: For all metrics, higher is better. Rows with \* have been reported based on the LoRA paper. Unlike (Hu et al., 2021a), we included the classifier number of parameters in our trainable parameters count.

performing reasonably well. As mentioned, this can be accomplished by only updating the unique parameters associated with rank  $r$  that do not overlap with lower ranks.

In addition, in Table 7, we demonstrate the result of using our individual loss (Eq. 9) vs. the nested dropout original objective function in an equal setting. As shown, our proposed objective function is both effective and efficient. Furthermore, it is important to note that the summation loss is not scalable when many ranks are involved. We also discussed the time complexity of LoRA and DyLoRA in Appendix A.

Maximum Rank: $r_{max} = 8$		
Loss	Training Time	CoLA
$\mathcal{L}_{\downarrow b}^{DY}$	<b>645.82s</b>	52.64
$\sum p_B(b)\mathcal{L}_{\downarrow b}^{DY}$	1175.69s	<b>54.12</b>

Table 7: This experiment shows the impact of choosing individual loss vs. summation loss functions on our training. The average performance across all possible ranks (1,2,...,8) is reported. For summation loss to be computationally more feasible, smaller epochs were chosen. A total of seven GPUs were used in this experiment.

## 6 Conclusion

In this paper, we presented our solution DyLoRA to address two problems in low-rank adapters regarding rank selection and making them dynamic. We showed that DyLoRA can select the rank without

requiring multiple re-training and is able to make LoRA dynamic at inference time. As a result, we can avoid the process of searching for the most optimal ranks for many real-life scenarios. It has been demonstrated that DyLoRA performance is comparable with LoRA, yet we can support a wider range of ranks without adding additional time and effort.

## Limitations

According to LoRA (Hu et al., 2021a), a proper choice of the scalar can improve the results. In order to determine what is the best choice, further investigation is required. Despite our demonstration that uniform distribution can be as effective as specific geometric distribution, further investigation is necessary to evaluate the effect of different distributions on different downstream tasks. As shown in this paper, our algorithm works over a wide range of ranks, but further research is needed to understand the impact of choosing a particular range.

## 7 Acknowledgement

We would like to use DyLoRA with Mindspore<sup>2</sup>, which is a new framework for deep learning computing.

<sup>2</sup>mindspore.cn

## References

- Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. 2020. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. *arXiv preprint arXiv:2012.13255*.
- Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. 2020. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. 2020. Generative pretraining from pixels. In *International conference on machine learning*, pages 1691–1703. PMLR.
- Patrick Chen, Hsiang-Fu Yu, Inderjit Dhillon, and Chou-Jui Hsieh. 2021. Drone: Data-aware low-rank compression for large nlp models. *Advances in neural information processing systems*, 34:29321–29334.
- 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*.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Utku Evci, Max Vladymyrov, Thomas Unterthiner, Bart van Merriënboer, and Fabian Pedregosa. 2022. Gradmax: Growing neural networks using gradient information. *arXiv preprint arXiv:2201.05125*.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The webnlg challenge: Generating text from rdf data. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 124–133.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7).
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. 2012. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
- Lu Hou, Zhiqi Huang, Lifeng Shang, Xin Jiang, Xiao Chen, and Qun Liu. 2020. Dynabert: Dynamic bert with adaptive width and depth. *Advances in Neural Information Processing Systems*, 33:9782–9793.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019a. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019b. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. 2019. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021a. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021b. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Aref Jafari, Mehdi Rezagholizadeh, Pranav Sharma, and Ali Ghodsi. 2021. [Annealing knowledge distillation](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2493–2504. Online. Association for Computational Linguistics.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2019. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*.
- Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021. Compacter: Efficient low-rank hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–1035.

- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.
- Lei Li, Yankai Lin, Shuhuai Ren, Peng Li, Jie Zhou, and Xu Sun. 2021. Dynamic knowledge distillation for pre-trained language models. *arXiv preprint arXiv:2109.11295*.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32.
- Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen-tau Yih, and Madsian Khabsa. 2021. Unipelt: A unified framework for parameter-efficient language model tuning. *arXiv preprint arXiv:2110.07577*.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xian-gru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. 2020. Dart: Open-domain structured data record to text generation. *arXiv preprint arXiv:2007.02871*.
- 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.
- Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The e2e dataset: New challenges for end-to-end generation. *arXiv preprint arXiv:1706.09254*.
- Peyman Passban, Yimeng Wu, Mehdi Rezagholizadeh, and Qun Liu. 2021. **ALP-KD: attention-based layer projection for knowledge distillation**. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 13657–13665. AAAI Press.
- Gabriele Prato, Ella Charlaix, and Mehdi Rezagholizadeh. 2020. Fully quantized transformer for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1–14.
- Ahmad Rashid, Vasileios Lioutas, and Mehdi Rezagholizadeh. 2021. Mate-kd: Masked adversarial text, a companion to knowledge distillation. *arXiv preprint arXiv:2105.05912*.
- Oren Rippel, Michael Gelbart, and Ryan Adams. 2014. **Learning ordered representations with nested dropout**. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1746–1754, Beijing, China. PMLR.
- Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. VI-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*.
- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. 2022. Black-box tuning for language-model-as-a-service. *arXiv preprint arXiv:2201.03514*.
- Marzieh S Tahaei, Ella Charlaix, Vahid Partovi Nia, Ali Ghodsi, and Mehdi Rezagholizadeh. 2021. Kroneckerbert: Learning kronecker decomposition for pre-trained language models via knowledge distillation. *arXiv preprint arXiv:2109.06243*.
- Chaofan Tao, Lu Hou, Wei Zhang, Lifeng Shang, Xin Jiang, Qun Liu, Ping Luo, and Ngai Wong. 2022. Compression of generative pre-trained language models via quantization. *arXiv preprint arXiv:2203.10705*.
- Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Guihong Cao, Daxin Jiang, Ming Zhou, et al. 2020. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*.
- Ziheng Wang, Jeremy Wohlwend, and Tao Lei. 2019. Structured pruning of large language models. *arXiv preprint arXiv:1910.04732*.
- Qiaolin Xia, Haoyang Huang, Nan Duan, Dongdong Zhang, Lei Ji, Zhifang Sui, Edward Cui, Taroon Bharti, and Ming Zhou. 2021. Xgpt: Cross-modal generative pre-training for image captioning. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 786–797. Springer.

## A Time complexity

The training time for DyLoRA is comparable to that of LoRA trained once on a specific rank. Thus, when searching the rank space for LoRA, we need to train it multiple times, whereas our method does not require searching the ranks. Accordingly, DyLoRA’s relative time complexity is inversely proportional to the number of possible ranks for which

the LoRA model must be searched. In MRPC, DyLoRA (for all the ranks) and LoRA (only on a single rank 8) require a total training time of 408.39 seconds and 399.95 seconds, respectively. Consequently, when we need to train eight LoRA models (Rank=1,2,...,8), it will result in a cost of  $399.95 \times 8 = 3199.6s$ , compared to the training time of our model, which is only 408.39 seconds. A more efficient implementation of our algorithm may result in a better time complexity.

## **B Hyperparameters**

We did not use any parameter tuning nor MNLI trick (initializing some down-streams tasks from MNLI checkpoint instead of pretrained weights). Therefore, we fine-tuned all the datasets from original pretrained weights. We simply followed a unified hyper-parameters for all different experiments. Unlike LoRA (Hu et al., 2021a) which reported the median over 5 random seeds, we reported the mean and standard deviation over 5 random seeds. See the details in Table 8.

## **C GPT Experiments**

A summary of the additional experiments that have been conducted to demonstrate the effectiveness of our proposed method for the task of language generation is provided in Table 9.

<b>Model</b>	<b>Parameter</b>	<b>Value</b>
RoBERTa-Base	Optimizer	AdamW
	Warmup Ratio	0.06
	LR Scheduler	Linear
	Batch Size	32
	Epochs	30
	Learning Rate (LR)	4e-4
	Weight Decay	0.1
	LoRA Config	$r_q = r_v = 8$ (unless otherwise mentioned)
	LoRA $\alpha$	16
	Max Sequence Length	512
	Seeds	10, 42, 4242, 10, 1010
	GPU	Tesla V100-PCIE-32GB
GPT Medium	Optimizer	AdamW
	Adam Beta2	0.999
	Warmup Steps	500
	Clip	0.0
	LR Scheduler	Linear
	Batch Size	8
	Epochs	5
	Learning Rate (LR)	2e-4
	Weight Decay	0.01
	Correct Bias	True
	LoRA Dropout	0.1
	Lable Smooth	0.1
	LoRA Config	$r_q = r_v = 4$
	LoRA $\alpha$	32
	Seeds	10, 42, 4242
	GPU	Tesla V100-PCIE-32GB

Table 8: All the hyperparameters that have been used throughout our study.

Model (Method)	Trainable Params	DART		WebNLG	
		BLEU $\uparrow$	TER $\downarrow$	BLEU $\uparrow$	TER $\downarrow$
Rank=1					
GPT-2 M (LoRA)	0.09M	0.71	0.49	2.80	1.18
GPT-2 M (DyLoRA-Frozen)	0.09M	44.48 $\pm$ 0.11	0.49 $\pm$ 0.00	52.09 $\pm$ 0.10	0.40 $\pm$ 0.01
GPT-2 M (DyLoRA)	0.09M	44.77 $\pm$ 0.17	0.49 $\pm$ 0.01	53.04 $\pm$ 0.07	0.40 $\pm$ 0.00
Rank=2					
GPT-2 M (LoRA)	0.19M	15.90	0.48	26.58	0.67
GPT-2 M (DyLoRA-Frozen)	0.19M	45.04 $\pm$ 0.14	0.48 $\pm$ 0.01	52.74 $\pm$ 0.31	0.40 $\pm$ 0.01
GPT-2 M (DyLoRA)	0.09M	46.05 $\pm$ 0.31	0.48 $\pm$ 0.00	54.32 $\pm$ 0.09	0.39 $\pm$ 0.01
Rank=3					
GPT-2 M (LoRA)	0.29M	35.84	0.47	43.61	0.47
GPT-2 M (DyLoRA-Frozen)	0.29M	45.22 $\pm$ 0.14	0.49 $\pm$ 0.01	53.03 $\pm$ 0.55	0.40 $\pm$ 0.00
GPT-2 M (DyLoRA)	0.29M	46.68 $\pm$ 0.36	0.48 $\pm$ 0.01	54.48 $\pm$ 0.05	0.39 $\pm$ 0.00
Rank=4					
GPT-2 M (LoRA)	0.39M	47.10	0.46	55.57	0.39
GPT-2 M (DyLoRA-Frozen)	0.39M	45.56 $\pm$ 0.33	0.48 $\pm$ 0.00	53.03 $\pm$ 0.01	0.40 $\pm$ 0.00
GPT-2 M (DyLoRA)	0.39M	46.56 $\pm$ 0.42	0.48 $\pm$ 0.01	54.48 $\pm$ 0.45	0.39 $\pm$ 0.00
Fine-Tune					
GPT-2 M (FT)*	354M	46.2	0.46		

Table 9: Rows with \* have been reported from the LoRA paper. (Hu et al., 2021a).