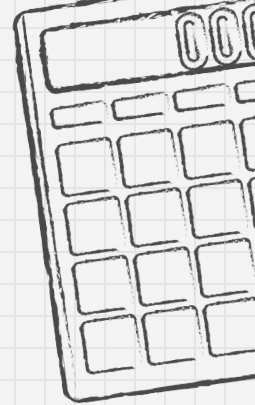


LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

By: Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu,
Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen

Table of contents



01

Paper discussion

Here we are going to talk about the motivation, briefly look at the proposed method and finally the results.

03

Demo time

Here we will look at some codes to strengthen what we discussed in the previous section

02

A little bit of Linear Algebra

We will introduce the essential lin alg for understanding LoRA

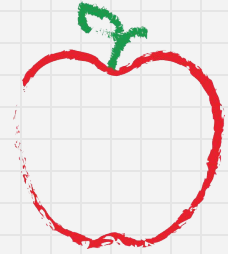
04

Connect LoRA implementation pictorially

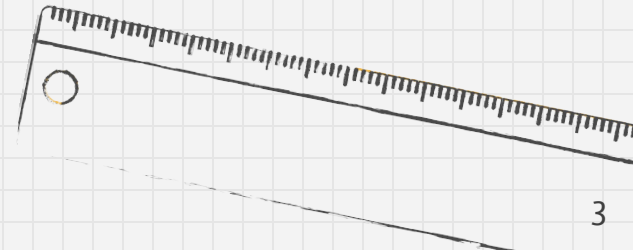
Here we will visually see on simple use case



01



Paper discussion



Motivation

- Since pretrained opensource large language models (LLMs) are readily available (and training them from scratch is not possible and pointless when we think about the regular applications). Fine tuning them towards our task should suffice. **BUT** fine tuning them itself becomes impossible due to their large size. Lets see why?
 - During training or fine tuning the amount of VRAM (GPU RAM) taken by the model is considerably higher than the amount of memory it takes during inference.
- If we want fine tune for two different task, then we will get 2 different fine tuned models. In the case of LLMs, storing these models in the disk becomes very cumbersome.
- Once fine tuning happens, a phenomenon known as catastrophic forgetting occurs. This where the model forgets the what it learnt during pre-training and only remembers the task it learnt during fine tuning.

Authors Claims why existing solutions do not scale to LLM paradigm

The authors take language modeling as an example scenario and identify two prominent strategies in efficient adaptation and discuss. The two strategies are,

- Adding adapter layers: Here we insert adapter layers at different positions of the model and train those adapter layers, the authors say by doing so we introduce inference latency.
- Directly Optimizing the Prompt: The authors observed that prefix tuning is difficult to optimize and that its performance changes non monotonically in trainable parameters.

Advantage of LoRA

- A pre-trained model can be shared and used to build many small LoRA modules for different tasks. By doing this task switching overhead is reduced significantly. (we will see the reason later on)
- LoRA makes training more efficient and lowers the hardware barrier to entry by up to 3 times when using adaptive optimizers since we do not need to calculate the gradients or maintain the optimizer states for most parameters. Instead, we only optimize the injected, much smaller low-rank matrices.
- Due to the simple linear design, the trainable matrices can be easily merged with frozen weights during deployment, thereby nullifying inference latency.
- LoRA is orthogonal to many prior methods and can be combined with many of them, such as prefix-tuning

Core Hypothesis of this work

- Authors take inspiration from the work of [1] where the paper claims that pre-trained language models have a low intrinsic dimension when adapting to a specific task.
- Taking this as inspiration the authors of LoRA hypothesized that the **updates to the weights are of low intrinsic ranks.**

Proposed Solution

$$h = W_0x + \Delta Wx = W_0x + BAx$$

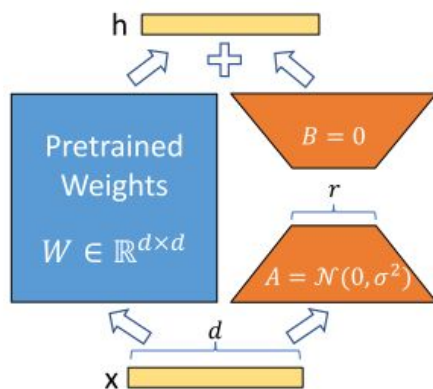
h is the output vector of a layer.

W_0 is pretrained weight of the layer (frozen, no learning allowed). Dimension $d \times k$

ΔW is the change in weight during fine tuning. Dimension $d \times k$

x is the input vector to that layer. Dimension $d \times 1$

B, A will be our low rank matrices which approximate ΔW . Dim of B $d \times r$ and Dim of A $r \times k$. Where r is rank.



Applying LoRA to Transformers

Theoretically LoRA can be applied to any weight matrix in a model, but the authors study on applying on Transformers only to the self attention module.

number of trainable parameters. In the Transformer architecture, there are four weight matrices in the self-attention module (W_q, W_k, W_v, W_o) and two in the MLP module. We treat W_q (or W_k, W_v) as a single matrix of dimension $d_{model} \times d_{model}$, even though the output dimension is usually sliced into attention heads. We limit our study to **only adapting the attention weights** for downstream tasks and freeze the MLP modules (so they are not trained in downstream tasks) both for simplicity and parameter-efficiency. We further study the effect on adapting different types of attention weight

Results

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB _{base} (FT)*	125.0M	87.6	94.8	90.2	63.6	92.8	91.9	78.7	91.2	86.4
RoB _{base} (BitFit)*	0.1M	84.7	93.7	92.7	62.0	91.8	84.0	81.5	90.8	85.2
RoB _{base} (Adpt ^D)*	0.3M	87.1 \pm .0	94.2 \pm .1	88.5 \pm 1.1	60.8 \pm .4	93.1 \pm .1	90.2 \pm .0	71.5 \pm 2.7	89.7 \pm .3	84.4
RoB _{base} (Adpt ^D)*	0.9M	87.3 \pm .1	94.7 \pm .3	88.4 \pm .1	62.6 \pm .9	93.0 \pm .2	90.6 \pm .0	75.9 \pm 2.2	90.3 \pm .1	85.4
RoB _{base} (LoRA)	0.3M	87.5 \pm .3	95.1 \pm .2	89.7 \pm .7	63.4 \pm 1.2	93.3 \pm .3	90.8 \pm .1	86.6 \pm .7	91.5 \pm .2	87.2
RoB _{large} (FT)*	355.0M	90.2	96.4	90.9	68.0	94.7	92.2	86.6	92.4	88.9
RoB _{large} (LoRA)	0.8M	90.6 \pm .2	96.2 \pm .5	90.9 \pm 1.2	68.2 \pm 1.9	94.9 \pm .3	91.6 \pm .1	87.4 \pm 2.5	92.6 \pm .2	89.0
RoB _{large} (Adpt ^P)†	3.0M	90.2 \pm .3	96.1 \pm .3	90.2 \pm .7	68.3 \pm 1.0	94.8 \pm .2	91.9 \pm .1	83.8 \pm 2.9	92.1 \pm .7	88.4
RoB _{large} (Adpt ^P)†	0.8M	90.5 \pm .3	96.6 \pm .2	89.7 \pm 1.2	67.8 \pm 2.5	94.8 \pm .3	91.7 \pm .2	80.1 \pm 2.9	91.9 \pm .4	87.9
RoB _{large} (Adpt ^H)†	6.0M	89.9 \pm .5	96.2 \pm .3	88.7 \pm 2.9	66.5 \pm 4.4	94.7 \pm .2	92.1 \pm .1	83.4 \pm 1.1	91.0 \pm 1.7	87.8
RoB _{large} (Adpt ^H)†	0.8M	90.3 \pm .3	96.3 \pm .5	87.7 \pm 1.7	66.3 \pm 2.0	94.7 \pm .2	91.5 \pm .1	72.9 \pm 2.9	91.5 \pm .5	86.4
RoB _{large} (LoRA)†	0.8M	90.6 \pm .2	96.2 \pm .5	90.2 \pm 1.0	68.2 \pm 1.9	94.8 \pm .3	91.6 \pm .2	85.2 \pm 1.1	92.3 \pm .5	88.6
DeB _{XXL} (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeB _{XXL} (LoRA)	4.7M	91.9 \pm .2	96.9 \pm .2	92.6 \pm .6	72.4 \pm 1.1	96.0 \pm .1	92.9 \pm .1	94.9 \pm .4	93.0 \pm .2	91.3

Table 2: RoBERTa_{base}, RoBERTa_{large}, and DeBERTa_{XXL} with different adaptation methods on the GLUE benchmark. We report the overall (matched and mismatched) accuracy for MNLI, Matthew’s correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. Higher is better for all metrics. * indicates numbers published in prior works. † indicates runs configured in a setup similar to [Houlsby et al. \(2019\)](#) for a fair comparison.

Crucial Results for us

	# of Trainable Parameters = 18M						
Weight Type Rank r	W_q 8	W_k 8	W_v 8	W_o 8	W_q, W_k 4	W_q, W_v 4	W_q, W_k, W_v, W_o 2
WikiSQL ($\pm 0.5\%$)	70.4	70.0	73.0	73.2	71.4	73.7	73.7
MultiNLI ($\pm 0.1\%$)	91.0	90.8	91.0	91.3	91.3	91.3	91.7

Table 5: Validation accuracy on WikiSQL and MultiNLI after applying LoRA to different types of attention weights in GPT-3, given the same number of trainable parameters. Adapting both W_q and W_v gives the best performance overall. We find the standard deviation across random seeds to be consistent for a given dataset, which we report in the first column.

Results show that applying LoRA to the Query layer (W_q) and Value Layer (W_v) should be sufficient.

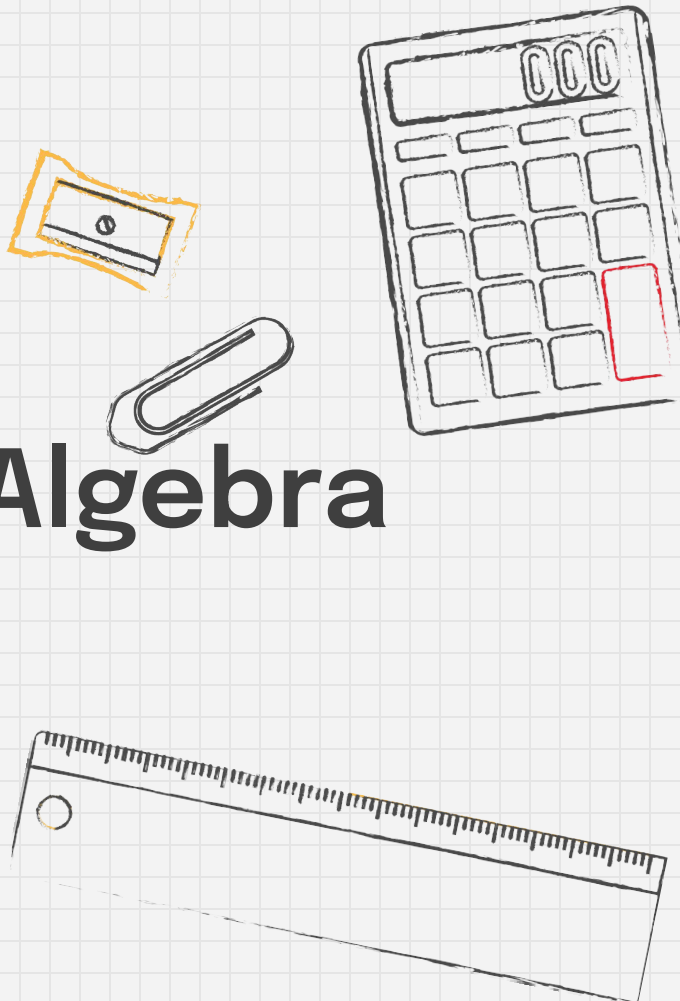
Crucial Results for us (Contd..)

	Weight Type	$r = 1$	$r = 2$	$r = 4$	$r = 8$	$r = 64$
WikiSQL($\pm 0.5\%$)	W_q	68.8	69.6	70.5	70.4	70.0
	W_q, W_v	73.4	73.3	73.7	73.8	73.5
	W_q, W_k, W_v, W_o	74.1	73.7	74.0	74.0	73.9
MultiNLI ($\pm 0.1\%$)	W_q	90.7	90.9	91.1	90.7	90.7
	W_q, W_v	91.3	91.4	91.3	91.6	91.4
	W_q, W_k, W_v, W_o	91.2	91.7	91.7	91.5	91.4

Table 6: Validation accuracy on WikiSQL and MultiNLI with different rank r . To our surprise, a rank as small as one suffices for adapting both W_q and W_v on these datasets while training W_q alone needs a larger r . We conduct a similar experiment on GPT-2 in [Section H.2](#).

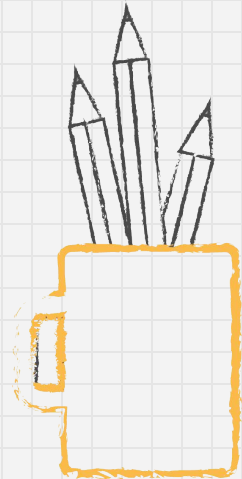
02

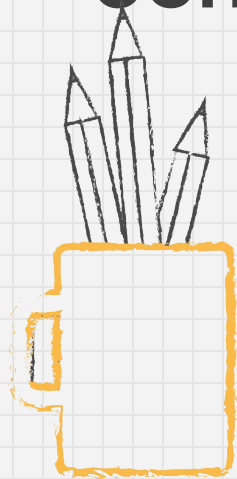
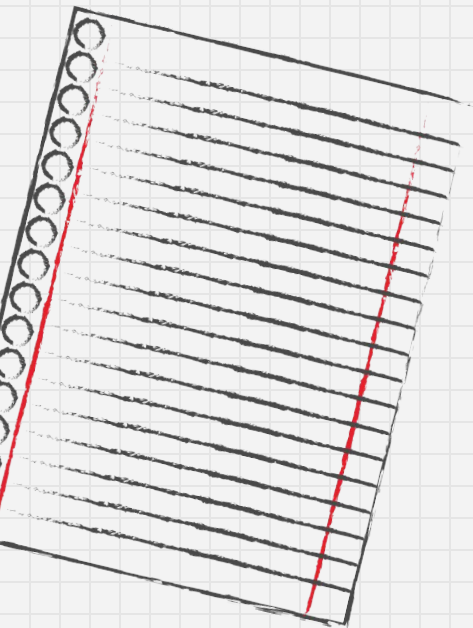
A little bit of Linear Algebra



03

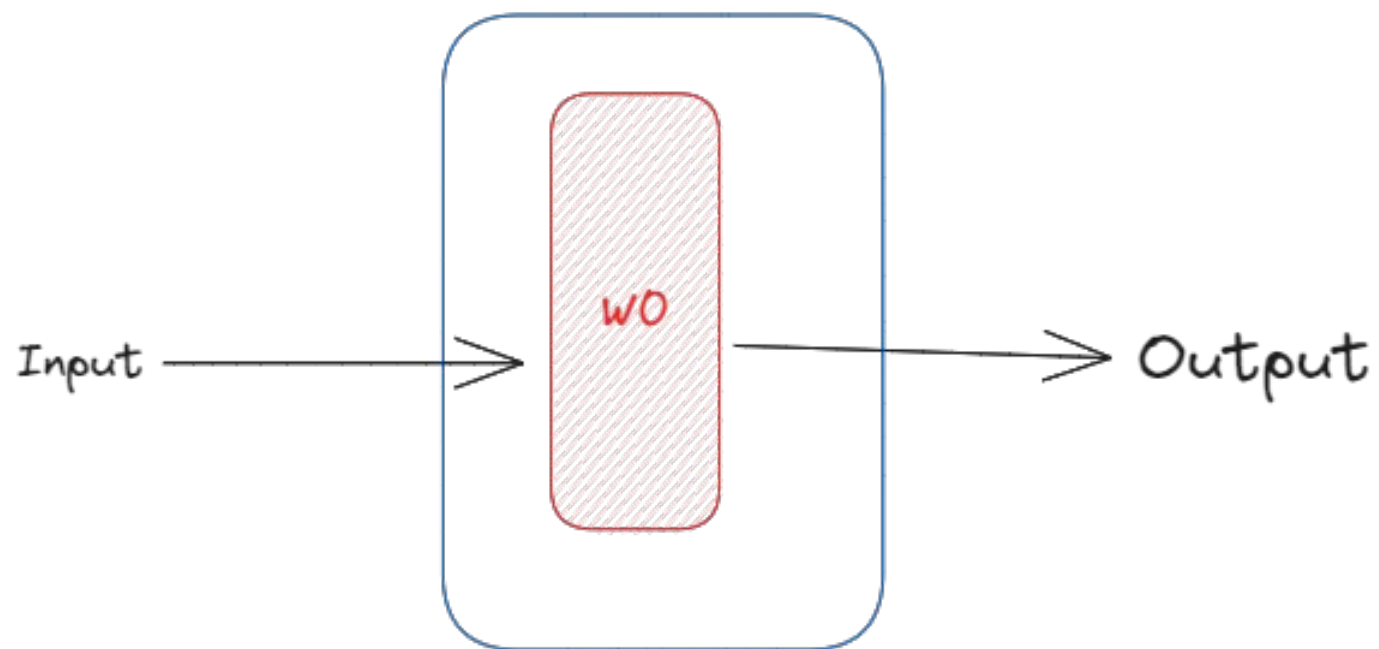
Demo Time

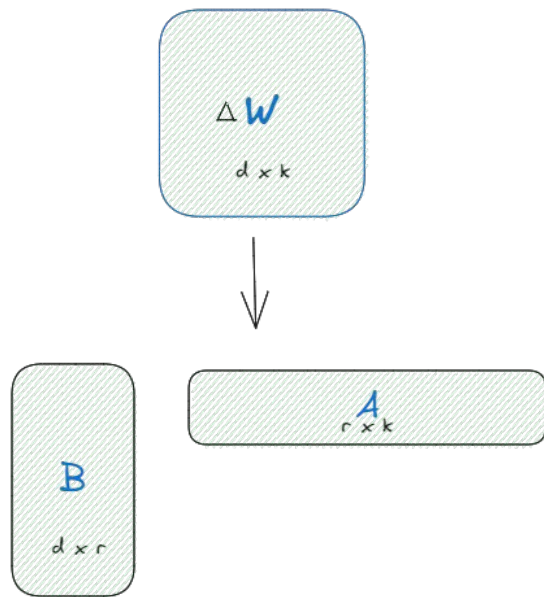
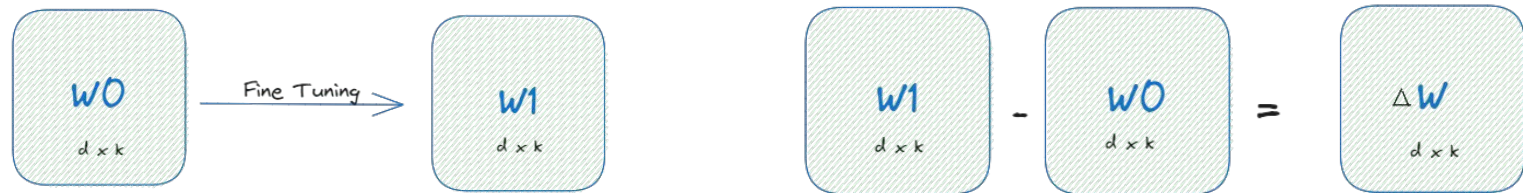


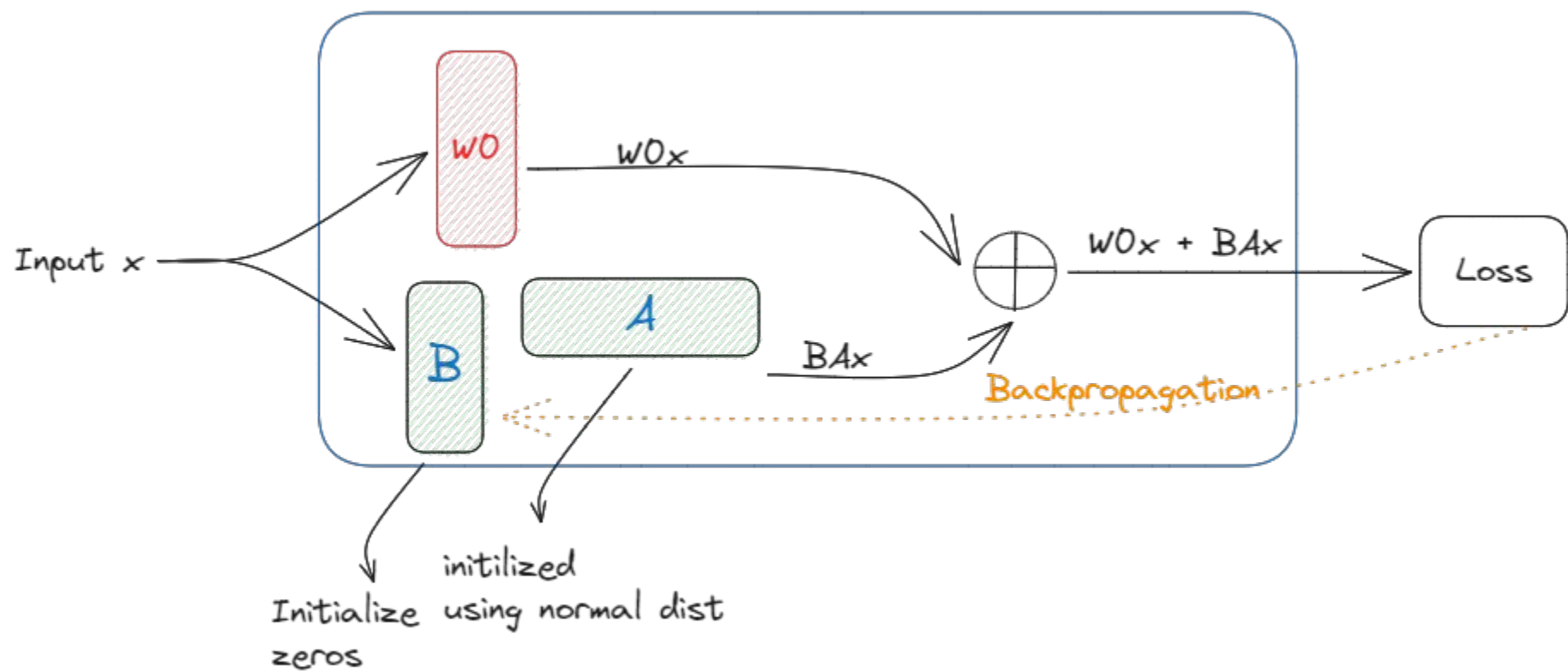


04

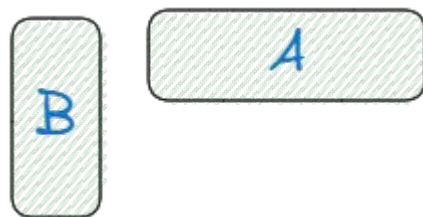
**Connect LoRA implementation
pictorially**







Save B and A for a particular task

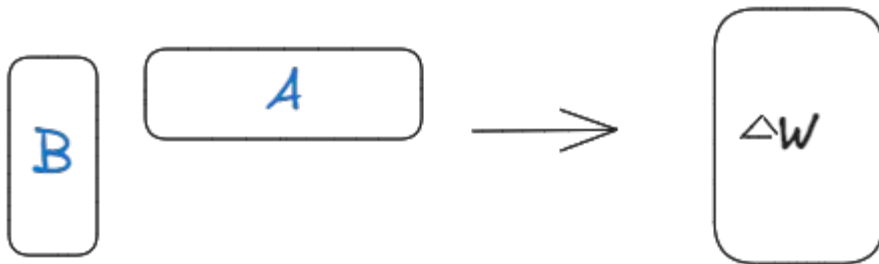


During Inference

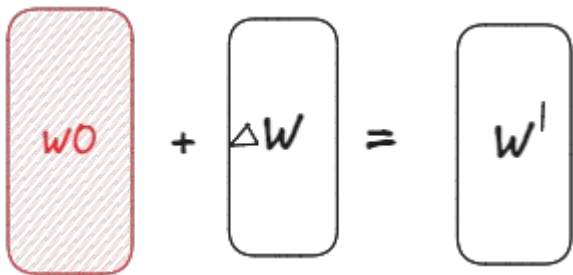
Take the stored A, B



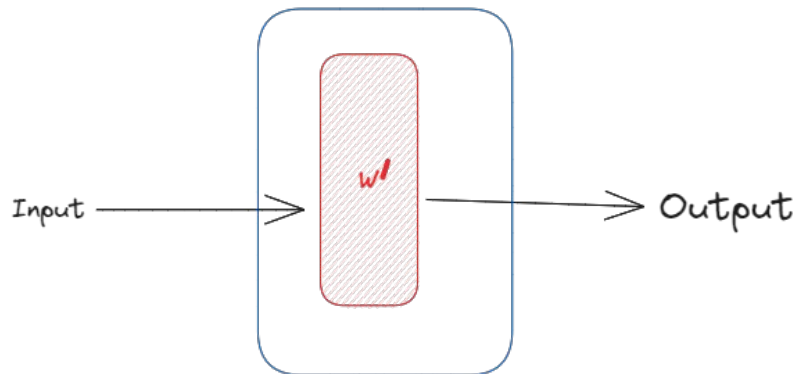
Multiply and create a delta matrix



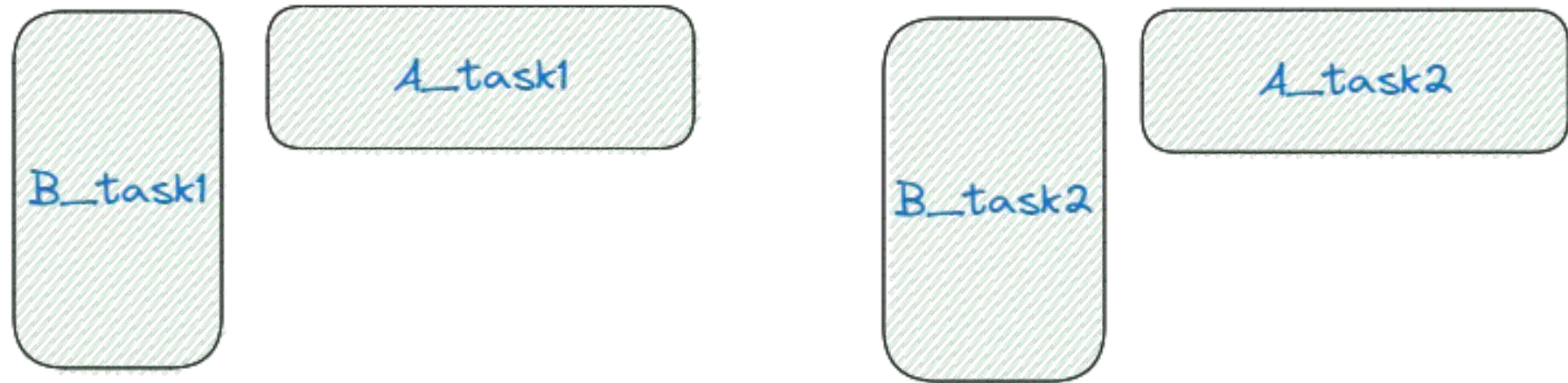
Add the matrix with the frozen weight of the layer and create updated weights for the layer



Finally model is ready for inference



Can store different A , B for different tasks



References

[1] Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic Dimensionality Explains the Effectiveness of Language Model Fine-Tuning. arXiv:2012.13255 [cs], December 2020. URL <http://arxiv.org/abs/2012.13255>. Zeyuan.