

S-LORA - Serving Thousands of Concurrent LORA Adapters.

(PRETRAIN THEN FINE TUNE) Commonly adopted in LLM

LORA

Low Rank Adapter  
PARAMETER EFFICIENT FINE TUNING METHOD  
Understand what LORA is all about

GPT 3 175B → Out to fine tune it  
run all 175B from start  
Pre trained with huge data which is less feasible and prohibitively expensive.

No need that

LORA Low Rank Adapter.

- ① freezes the pre trained Model weights
- ② injects the trainable rank decomposition matrices into each layer of transformer As white box

Low Rank Adapter Matrices greatly reduces the No. of Parameters or Down time  
adapter weights GPT 3 175B

LORA can reduce the trainable parameters 10000 times than the Adam  
GPU Memory Requirement By 3 times  
As an inspiration they took =)

[learned over parametrized] models in fact reside in low intrinsic dimension

Change in weights during Model Adaption  
 $\Downarrow$   
 low Intrinsic Rank

LoRA  $\Rightarrow$  (freeze only some dense layers in a NN) indirectly optimizes the rank decomposition matrices of the dense layer. while keeping the pre-trained weights frozen

pre-trained Auto Regressive Language Model.

$[P \circlearrowleft (y|x)] \rightarrow$  Parameter

Can be a generic multi task learner such as GPT.

Each Downstream task is represented by the training dataset

$\mathcal{Z} = \{(x_i, y_i), \dots\}_{i=1, \dots, N}$   
 $x_i$  - input  $y_i$  - target (output)

During full fine tuning.

Model

$P \circlearrowleft (y|x) \Rightarrow$  pre-trained with  $\phi$  and updated to  $\phi + \Delta\phi$  by repeatedly following the gradient to maximize the conditional language modelling objective.

$$\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1} \log (P_{\Phi}(y_t | x, y_{<t}))$$

Draw backs during full fine tuning :

They learned different set of parameters  
by  $\Delta \Phi$  whose dimension

$$|\Delta \Phi| = |\Phi|$$

fine tuned size 175B Pre-trained size 175B

LoRA

$$\Delta \Phi = \Delta \Phi(\Theta)$$

Small sized set of parameters.  $\Theta$  with

$$|\Theta| \ll |\Phi_0|$$

Dimension 181

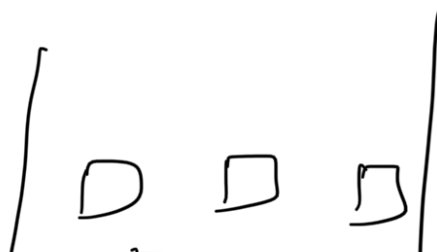
$$\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1} \log (P_{\Phi_0 + \Delta \Phi(\Theta)}(y_t | x, y_{<t}))$$

LoRA  $\Rightarrow$  Base <sup>one</sup> Model  
Substantial Collection of LoRA Adapters

Advantage  $\swarrow$

S-LoRA  $\Rightarrow$  Serving Concurrent of  
Thousands and Adapters.

SLoRA



fetch the adapter that are



all the adapters  
in the main memory

currently running query  
to GPU Memory

To efficiently use GPU  $\Rightarrow$  Reduce fragmentation  $\Rightarrow$  S-LORA [Unified Paging]

① Unified Memory Pool

to manage dynamic adapter weights.   
 Stores both adapter weights & KV Cache. As a hashmap.   
 (loading adapter weights & KV Cache. As a hashmap.)

② Tensor Parallelism

③ Highly optimized CUDA kernels

Heterogeneous Batch LORA Computation.

Compared with S-GPT libraries   
 [Hugging face PEFT]  $\Rightarrow$  improves throughput by up to 4 times and increases the number of served adapters.   
 S-LORA

## LORA Computation

① Low Rank Parametrized Update Matrices:

NN many Dense Layers  $\Rightarrow$  Matrix Multiplication   
 Weight Matrices typically have full Rank.

Weight Matrices in these layers

low intrinsic Rank during adaptation

Pretrained weight Matrix  $W_0 \in \mathbb{R}^{d \times k}$

$W_0 + \Delta W = W_0 + BA$  where  
latter  $B \in \mathbb{R}^{d \times r}$ ,  $A \in \mathbb{R}^{r \times k}$   
 Constraint and the rank  $r \leq \min(d, k)$

During training: freezes  $(W_0) \Rightarrow$  will not receive any gradient updates.

A and B contain trainable parameters.

Note both  $W_0$  and  $\Delta W = BA$   
 are multiplied with the same input

and their respective results are summed coordinate wise.

$$h = W_0 x$$

The Modified forward pass into the transformer

$$h = W_0 x + \Delta W x = W_0 x + BAx$$

{ Random Gaussian initialization for A }  
 zero for B  
 at the

beginning of the training.

② Then scale  $\Delta w$  by  $\frac{\alpha}{r}$  where  $\alpha$  is the constant in  $r$

② Applying LoRA to transformer.

We know that in transformer there are 4 weight matrices.

$W_q$   $W_k$   $W_v$   $W_o$  and two in MLP module

They treat  $[W_q, W_k, W_v, W_o]$  as a single matrix of dimension  $d_{model} \times d_{model}$  even though the output is usually sliced into attention heads.