

# *Parameter Efficient Fine Tuning (PEFT)*

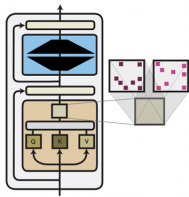
Pawan Goyal

CSE, IIT Kharagpur

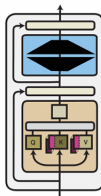
CS60010

# Different Perspectives to think about PEFT

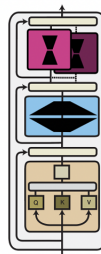
LORA  
QLORA  
Penny



Parameter



Input



Function

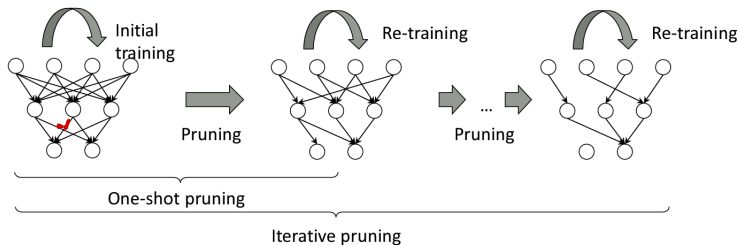
Task-specific

→ adapters

prefix-tuning  
prompt-tuning

# Pruning

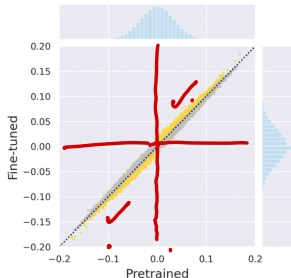
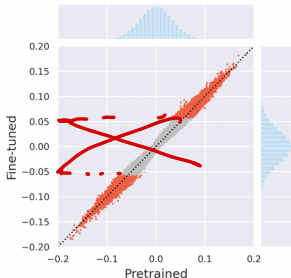
- During pruning, a fraction of the lowest-magnitude weights are removed
- The non-pruned weights are re-trained
- Pruning for multiple iterations is more common ([Frankle & Carbin, 2019](#))



# How to apply Pruning for Pre-trained models?

- Pruning does not consider how weights change during fine-tuning
- **Magnitude pruning**: keep weights farthest from 0
- **Movement pruning** [Sanh et al., 2020]: keep weights that *move the most away from 0*

Fine-tuned weights stay close to their pre-trained values. Magnitude pruning (left) selects **weights that are far from 0**.



Movement pruning (right) selects weights that **move away from 0**.

0.001  $\rightarrow$  0.005

# Revisiting full fine-tuning

- Assume we have a pre-trained autoregressive language model  $P_\phi(y|x)$ 
  - E.g., GPT based on Transformer
- Adapt this pretrained model to downstream tasks (e.g., summarization, NL2SQL, reading comprehension)
  - Training dataset of context-target pairs  $\{(x_i, y_i)\}_{i=1, \dots, N}$
- During full fine-tuning, we update  $\phi_o$  to  $\phi_o + \Delta\phi$  by following the gradient to maximize the conditional language modeling objective

$$\max_{\phi} \sum_{(x,y)} \sum_{t=1}^{|y|} \log(P_\phi(y_t|x, y_{<t}))$$

# LoRA: Low-rank Adaptation

- For each downstream task, we learn a different set of parameters  $\Delta\phi$ 
  - $|\Delta\phi| = |\phi_o|$
  - GPT-3 has a  $|\phi_o|$  of 175 billion
  - Expensive and challenging for storing and deploying many independent instances
- Key idea:** encode the task-specific parameter increment  $\Delta\phi = \Delta\phi(\Theta)$  by a smaller-sized set of parameters  $\Theta$ ,  $|\Theta| \ll |\phi_o|$
- The task of finding  $\Delta\phi$  becomes optimizing over  $\Theta$

$$\max_{\Theta} \sum_{(x,y)} \sum_{t=1}^{|y|} \log(P_{\phi_o + \Delta\phi(\Theta)}(y_t | x, y_{<t}))$$

$$|\Delta\phi(\Theta)| \ll |\Delta\phi|$$

# Low-rank-parameterized update matrices

- Updates to the weights have a low “intrinsic rank” during adaptation (Aghajanyan et al. 2020)

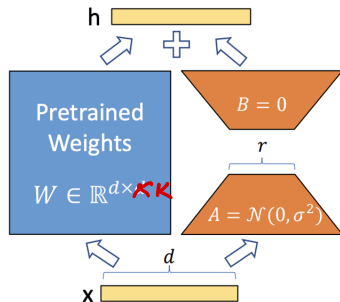
- $W_0 \in \mathbb{R}^{d \times k}$ : a pretrained weight matrix

- Constrain its update with a low-rank decomposition:

$$\underline{W_0 + \Delta W} = W_0 + \underline{BA}$$

where  $B \in \mathbb{R}^{d \times r}$ ,  $A \in \mathbb{R}^{r \times k}$ ,  $r \ll \min(d, k)$

- Only A and B contain **trainable** parameters



$BA: d \times k$

$$\textcircled{B} \quad r \times (d+k) \ll d, k \quad r \ll \min(d, k)$$

# LoRA: What happens during training and inference?

- $W_0 + \Delta W = W_0 + BA$
- During training,  $W_0$  is frozen and does not receive gradient updates.
- For  $h = W_0x$ , forward pass would be

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

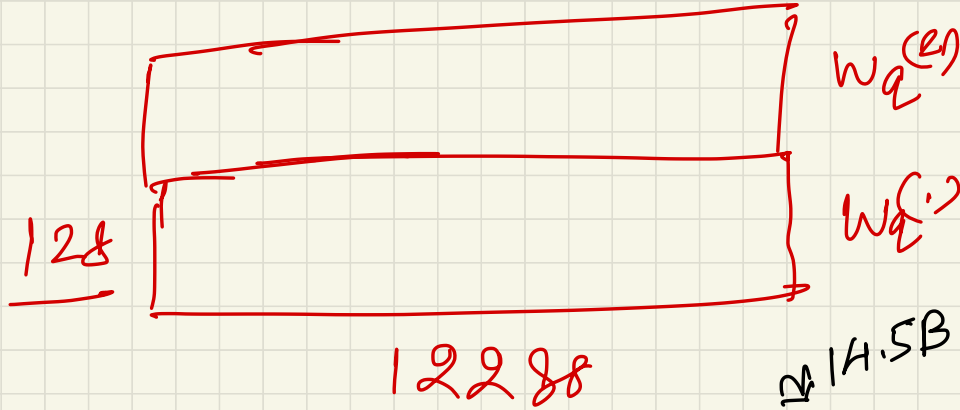
*At inference time?*

Use the matrix  $W = W_0 + BA$ , so no additional latency



$W_d(96)$

1  
2  
1



$$(12288 \times 2) \times 2 \times 96 \approx 4.7m$$

$$W_d : (12288 \times 12288) \times 96$$

$$W_d + \underline{B \cdot A}$$

$$\begin{aligned} B &: 12288 \times 2 \\ A &: 2 \times 12288 \\ \underline{2} &: 2 \end{aligned}$$

# Which weight matrices to apply to?

Which weight matrices in Transformers should we apply LoRA to?

	# 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	<b>73.7</b>	<b>73.7</b>
MultiNLI ( $\pm 0.1\%$ )	91.0	90.8	91.0	91.3	91.3	91.3	<b>91.7</b>

Adapting both  $W_q$  and  $W_v$  gives the best performance overall.

What is the optimal rank  $r$  for LoRA?

	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

LoRA already performs competitively with a very small  $r$

# Understanding LoRA Parameters over GPT-3

- Trainable parameters =  $2 \times \hat{L}_{LoRA} \times d_{model} \times r$
- For GPT-3,  $d_{model} = 12288$ , 96 decoders and 96 attention heads

LoRA	$r_v = 2$	4.7 M
	$r_q = r_v = 1$	4.7 M
	$r_q = r_v = 2$	9.4 M
	$r_q = r_k = r_v = r_o = 1$	9.4 M
	$r_q = r_v = 4$	18.8 M
	$r_q = r_k = r_v = r_o = 2$	18.8 M
	$r_q = r_v = 8$	37.7 M
	$r_q = r_k = r_v = r_o = 4$	37.7 M
	$r_q = r_v = 64$	301.9 M
	$r_q = r_k = r_v = r_o = 64$	603.8 M

# Applying LoRA to Transformer

Model & Method	# Trainable Parameters	E2E NLG Challenge				
		BLEU	NIST	MET	ROUGE-L	CIDEr
GPT-2 M (FT)*	354.92M	68.2	8.62	46.2	71.0	2.47
GPT-2 M (Adapter <sup>L</sup> )*	0.37M	66.3	8.41	45.0	69.8	2.40
GPT-2 M (Adapter <sup>L</sup> )*	11.09M	68.9	8.71	46.1	71.3	2.47
GPT-2 M (Adapter <sup>H</sup> )	11.09M	67.3 $\pm$ .6	8.50 $\pm$ .07	46.0 $\pm$ .2	70.7 $\pm$ .2	2.44 $\pm$ .01
GPT-2 M (FT <sup>Top2</sup> )*	25.19M	68.1	8.59	46.0	70.8	2.41
GPT-2 M (PreLayer)*	0.35M	69.7	8.81	46.1	71.4	2.49
GPT-2 M (LoRA)	0.35M	<b>70.4<math>\pm</math>.1</b>	<b>8.85<math>\pm</math>.02</b>	<b>46.8<math>\pm</math>.2</b>	<b>71.8<math>\pm</math>.1</b>	<b>2.53<math>\pm</math>.02</b>
GPT-2 L (FT)*	774.03M	68.5	8.78	46.0	69.9	2.45
GPT-2 L (Adapter <sup>L</sup> )	0.88M	69.1 $\pm$ .1	8.68 $\pm$ .03	46.3 $\pm$ .0	71.4 $\pm$ .2	<b>2.49<math>\pm</math>.0</b>
GPT-2 L (Adapter <sup>L</sup> )	23.00M	68.9 $\pm$ .3	8.70 $\pm$ .04	46.1 $\pm$ .1	71.3 $\pm$ .2	2.45 $\pm$ .02
GPT-2 L (PreLayer)*	0.77M	70.3	8.85	46.2	71.7	2.47
GPT-2 L (LoRA)	0.77M	<b>70.4<math>\pm</math>.1</b>	<b>8.89<math>\pm</math>.02</b>	<b>46.8<math>\pm</math>.2</b>	<b>72.0<math>\pm</math>.2</b>	2.47 $\pm$ .02

GPT-2 medium (M) and large (L) with different adaptation methods on the E2E NLG Challenge. For all metrics, higher is better. LoRA outperforms several baselines with comparable or fewer trainable parameters