

LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

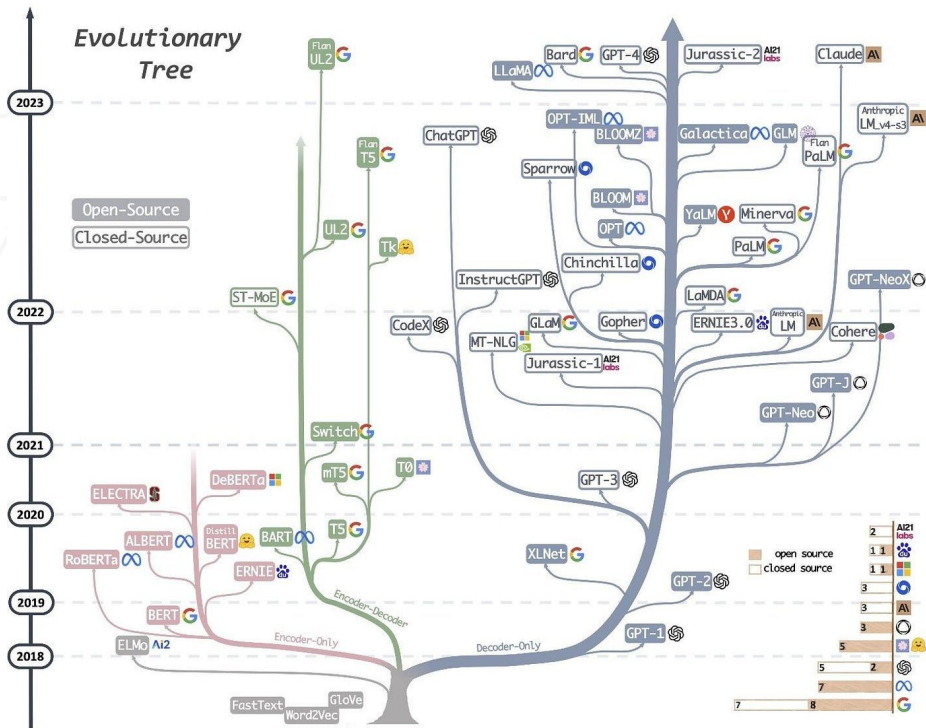
Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu,
Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen

Presented by: Kien Tran, Ashish Dhiman

Agenda

1. Problem introduction
2. Background
3. Method
4. Experimental Analysis
5. Strengths and Weaknesses
6. Further developments

Problem introduction



Problems with Foundation LLMs

- General-purpose - i.e. not task-specific
- VERY LARGE Models

```
pip install -q meaningcloud-python
```

§ Output

```
> stdout : ['Note: you may need to restart the kernel to use updated packages.\n']
```

§ Code

```
import meaningcloud

#Dictionary to store API Keys
DC=

{'AS':meaningcloud.MeaningCloud('YOUR_LICENSE_KEY'),'DV':meaningcloud.MeaningCloud('YOUR_OTHER_LICENSE_KEY')}

#Automatic categorization

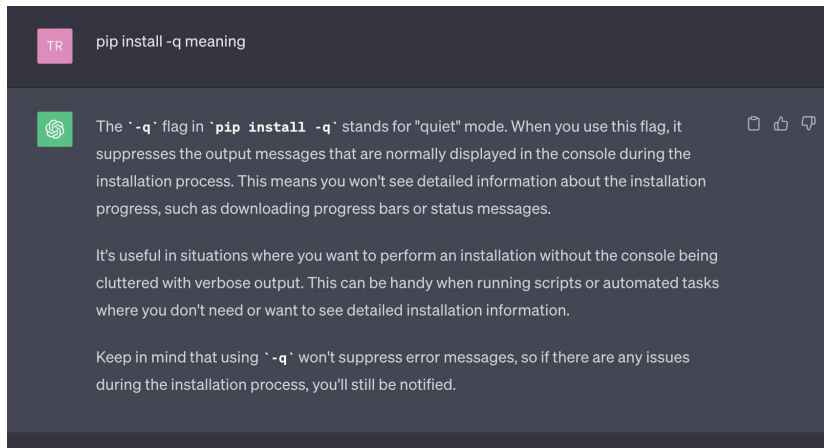
categorization_req = {
    'key': DC['DV'].user_key,
    'lang': 'en',
    'url': 'http://www.dcc.uchile.cl/'
}
```

Domain Adaptation

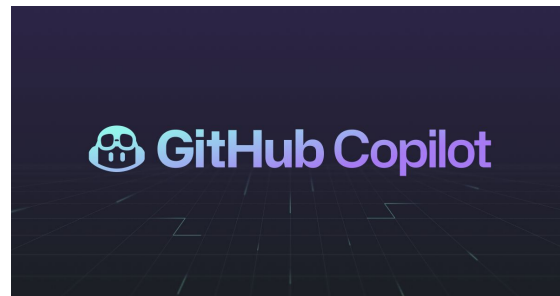
Domain Adaptation is a type of **transfer learning** that involves training a model with **data from a source domain**.

Examples:

Chat Conversation



Code Generation



Problem Statement: Fine-tuning LLMs

Base model

Parameters

$$\Phi_0$$

GPT4, LLaMA,
etc.

Adapting Objective

Full Fine-tuning (Adaptation)

$$\Phi_0 + \Delta\Phi$$

GPT-3 with $|\Phi_0| \approx 175 \text{ Billion} \approx |\Delta\Phi|$

- Large hardware requirements
- Costly training, storage, and inference

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

Adapting to domain-specific dataset \mathcal{Z}

Parameter-Efficient Adaptation

$$\Phi_0 + \Delta\Phi(\Theta)$$

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

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

Background

Existing Methods & Limitations (1/2)

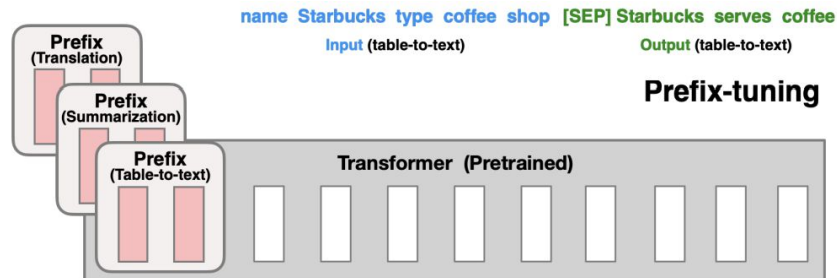
Methods

Prompt Engineering

Describe the **instructions** and **examples** of the task (one-shot, few-shots learning) **with words**

Continuous prompts (e.g. Li & Liang, 2021)

Instead of discrete prompts (words), use continuous prompts (trainable special vectors)



Limitations

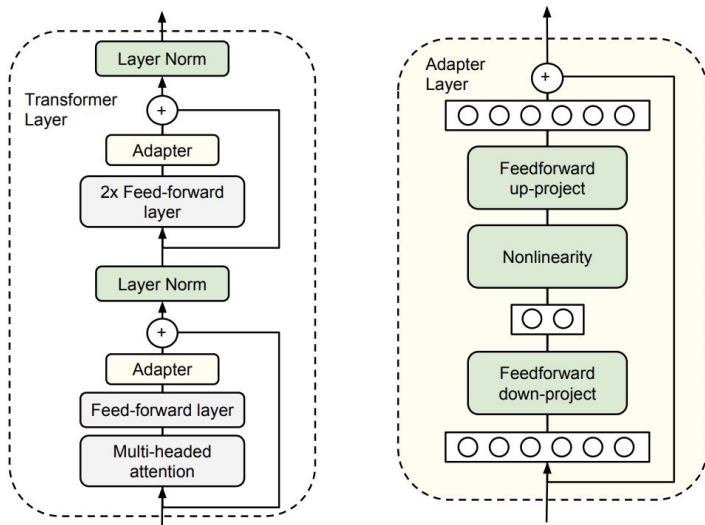
- Unreliable performance (“prompt engineering is an art”)
- Waste computing power processing prompts
- Valuable token space must be spent on prefix token embeddings
- Increase inference time
- Not so sure regarding scalability

Existing Methods & Limitations (2/2)

Methods

Adapter-based (e.g. Houlsby et al., 2019)

Insert low-rank, trainable adapter layers between existing layers



Limitations

- Multi-headed attention weights was not changed
- Make the model deeper, thus introduce additional latency during inference
- Not able to out-perform full fine-tuning baseline

Method

LoRA

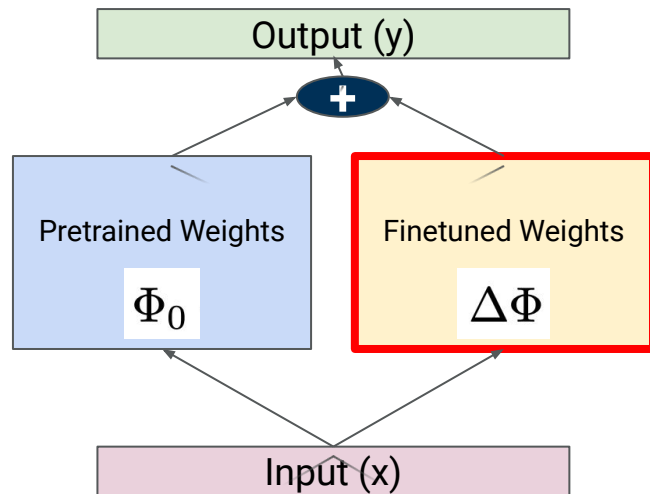
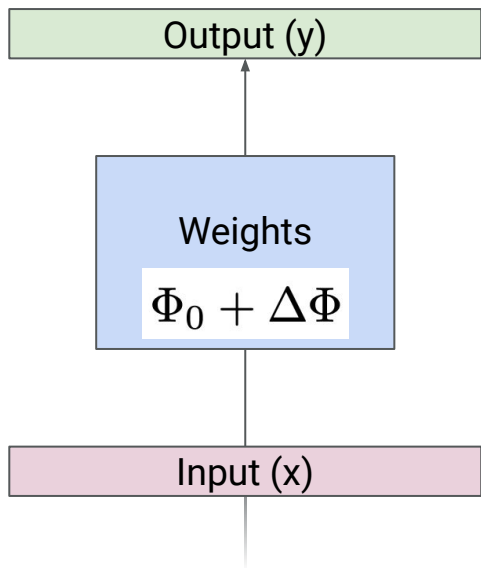
Why do we need it ?

- It's too expensive to fine-tune all parameters in a large model.

$$\Phi' = \Phi_0 + \Delta\Phi$$

LoRA

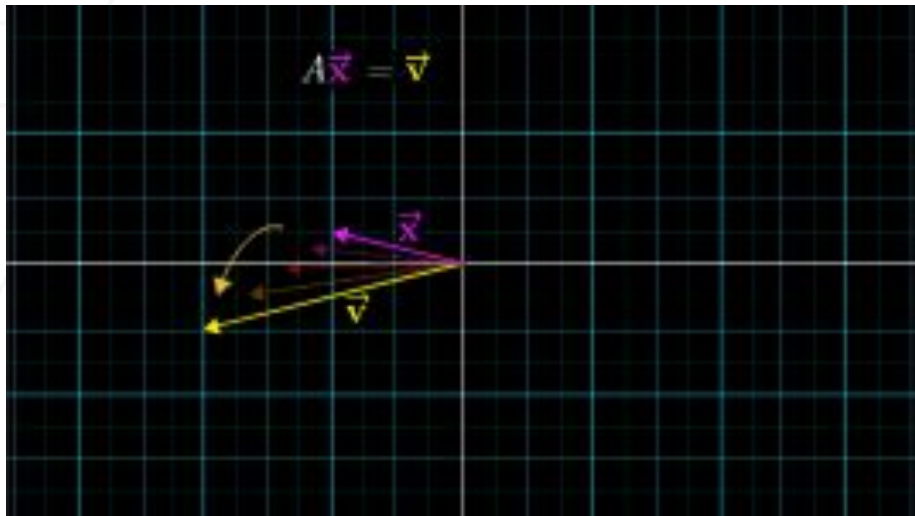
Are the fine-tuning updates full rank ?



Is it full rank

LoRA

What is full rank vs low rank matrix ?

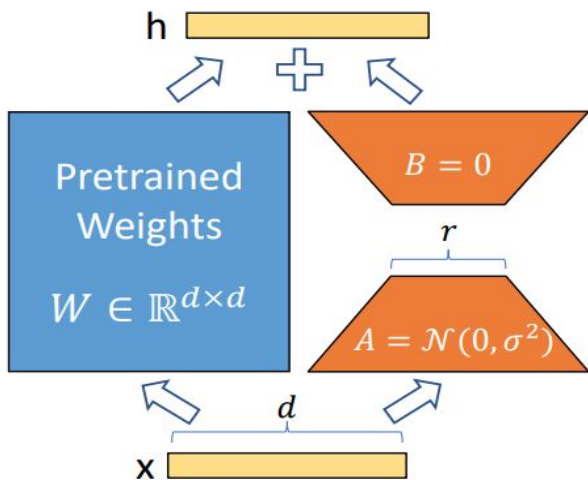


- Rank of a matrix, (r) is the number of linearly independent columns/rows
- For a full rank matrix, r = lowest dimension of the matrix

$$\begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

LoRA

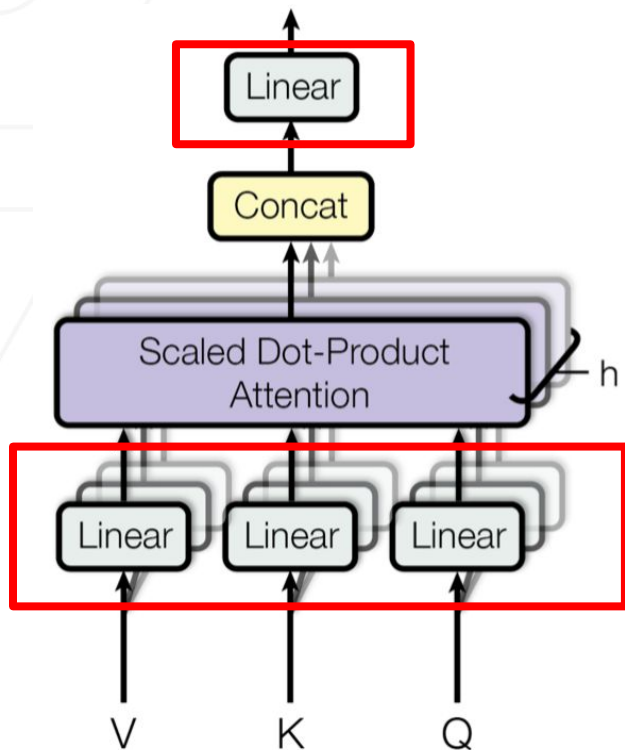
Use low rank decomposition for fine tuning updates



The diagram shows the matrix representation of the LoRA update. The initial weight matrix W is shown as a matrix of zeros and ones. This is updated to $W + \Delta W$, where ΔW is the low-rank update. The update is decomposed into two matrices, A and B , such that $\Delta W = BA$. The matrix A is of size $r \times k$ and the matrix B is of size $d \times r$. The rank r is much smaller than $\min(d, k)$.

$$h = W_0 x + \Delta W x = W_0 x + B A x$$
$$B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times k}, r \ll \min(d, k)$$

LoRA for Transformer



- Weight matrices of linear layers in Transformer architecture:

- W_Q, W_V, W_K, W_O

- Apply Low rank decomposition to these matrices while fine tuning

- $W_Q = []_{100,100} = B_{100,3} A_{3,100}$

$3 \times 100 \times 2$ vs 100^2

Experimental Analysis

Experimental Setup

Baseline methods

- Fine-tune
- Bias only
- Prefix-embedding tuning
 - injects special tokens alongside the input tokens
- Prefix-layer tuning
 - learn the Prefix-embedding after every layer.
- Adapter tuning
 - inserts adapter layers between the self-attention module (and the MLP module) and the subsequent residual connection.

Baseline models

- BERT
- **RoBERTa**
- GPT 2
- **GPT 3**

Results: RoBERTa

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\pm1.2	68.2\pm1.9	94.9\pm.3	91.6 \pm .1	87.4\pm2.5	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\pm1.0	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\pm1.0	68.2 \pm 1.9	94.8\pm.3	91.6 \pm .2	85.2\pm1.1	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\pm1.1	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.

Results: GPT-3

Method	# of Trainable Parameters	WikiSQL	MNLI-m	SAMSum
		Accuracy (%)	Accuracy (%)	R1/R2/RL
GPT-3 175B (Fine-Tune)	175,255.8M	73.0	89.5	52.0/28.0/44.5
GPT-3 175B (Bias Only)	14.2M	71.3	91.0	51.3/27.4/43.5
GPT-3 175B (PrefixEmbed)	3.2M	63.1	88.6	48.3/24.2/40.5
GPT-3 175B (PrefixLayer)	20.2M	70.1	89.5	50.8/27.3/43.5
GPT-3 175B (LoRA)	4.7M	73.4	91.3	52.1/28.3/44.3
GPT-3 175B (LoRA)	37.7M	73.8	91.7	53.2/29.2/45.0

Table 1: Logical form validation accuracy on WikiSQL, validation accuracy on MultiNLI-matched and Rouge-1/2/L on SAMSum achieved by different GPT-3 adaptation methods. LoRA performs better than prior approaches, including conventional fine-tuning. The result on WikiSQL has a fluctuation of $\pm 0.3\%$ and MNLI-m $\pm 0.1\%$.

Ablation Study: Understanding Low-rank Update

Q1: Which weight to apply LoRA?

	# 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

- Training **at least W_q and W_v** to achieve good results
- Training all attention weights gives the best results on the same parameter budget

Q2: How to choose 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

- Only $r = 1$ already enable very good results

Strengths and Weaknesses

Strengths and Weaknesses

Strengths

1. Innovative approach to fine-tuning with multiple practical advantages
 - Save on storage of multiple fine-tune models
 - Can merge the weight during inference, thus requiring no additional latency
 - Lower hardware requirements for tuning
2. Ablation studies prove that low-rank adaptation is effective, even with $r = 1$
3. A general method, can be applied to many problems and in combinations with other fine-tuning methods

Weaknesses

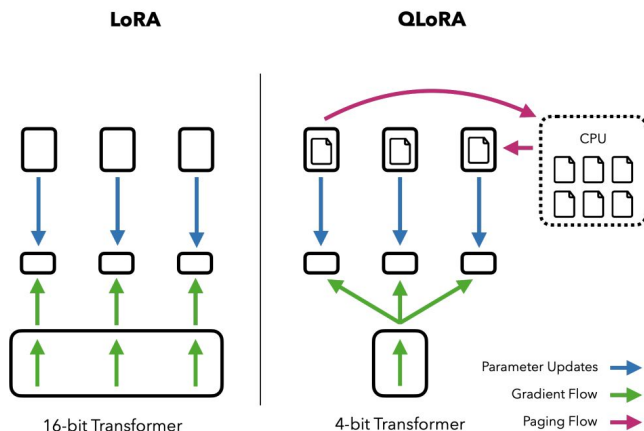
1. Lower hardware requirements, **but still high**: At least GPUs that have enough RAM to load the full base models
2. Deployment still needs same amount of memory (say 175B), only if you are deploying multiple models, do the savings kick in.
3. Adds decomposition rank ' r ' as another hyperparameter to tune.
4. Ablation missing for cases where different weight matrices have different rank decompositions, say 2 or Q, 6 for V.

Further developments

Recent developments leveraging LoRA

Quantization (QLoRA)

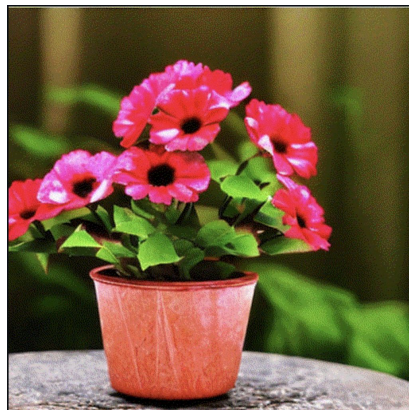
Dettmers et al. (2023) combine LoRA with 4-bit quantization to *reduce the memory requirement* and improve *computation efficiency* of fine tuning without sacrificing performance



Fine-tune Diffusion model*

- 2x faster fine-tuning of the Stable Diffusion model compared to Dreambooth
- Small model (1MB ~ 6MB vs GBs), enable sharing

Photo-realistic images



Cartoon images



* cloneofsim, Low-rank Adaptation for Fast Text-to-Image Diffusion Fine-tuning, <https://github.com/cloneofsim/lora> 2023

Q&A