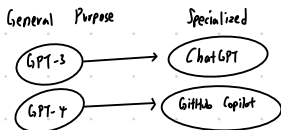


What is finetuning?

: Taking general purpose models (GPT-3) and specializing them into the specific task



What does finetuning do for the models?

- Lets you put more data into the pretrained model than what fits into the prompt
- Gets the model to learn the data, rather than just get access to it.
- Steers the model to more consistent outputs
- Reduces hallucinations
- Customizes the model to a specific use case

Prompt Engineering vs Finetuning

- | Pros  |                                 |
|---|---------------------------------|
| ① No data to get started                                    | ① Already unlimited data fits   |
| ② Smaller upfront cost                                      | ② Learn new information         |
| ③ No technical knowledge needed                             | ③ Correct incorrect information |
| ④ Connect data through RAG (Retrieval Augmented Generation) | ④ Use RAGs too                  |

- | Cons                  |                                  |
|-----------------------|----------------------------------|
| ① Much less data fits | ① More high-quality data         |
| ② Forgets data        | ② Needs some technical knowledge |
| ③ Hallucination       | ③ Computing cost                 |
| ④ Gets incorrect data |                                  |

Benefits of finetuning your own LLM

- ① Performance : 1) Stop hallucination 2) Increase consistency 3) Reduce unwanted info. 4) A smaller (fine-tuned) model can outperform a larger base model.
- ② Security : 1) Prevents leakage 2) No breaches
- ③ Cost
- ④ Reliability

Supervised Fine-tuning

1. Choose fine-tuning task
2. Prepare training dataset
3. Choose a base model
4. Fine-tune model via supervised learning
5. Evaluate model performance

### 3 options for Parameter Training

#### 1) Retrain all parameters

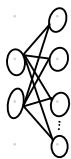
downside: billions of internal model parameters the computational cost for training explodes

#### 2) Transfer learning: Instead of retraining all the parameters, we freeze most of the parameters and only finetune the head (last few layers of the model)

#### 3) Parameter Efficient Fine-tuning (PEFT): Freeze all of the weights. It augments the model with additional parameters which are trainable.

One of the most popular ways to do PEFT: Low-Rank Adaptation (LoRA)

#### Pictorial View



$x \rightarrow h(x)$

$$h(x) = W_0 x$$

$$\begin{bmatrix} W_0 \end{bmatrix} x = h(x)$$

$$W_0 \in \mathbb{R}^{d \times d}$$

$$x \in \mathbb{R}^{d \times 1}$$

$$h(x) \in \mathbb{R}^{d \times 1}$$

If  $d=1000$ ,  $K=1000$ ,  
there are 1,000,000  
trainable parameters

#### LoRA

$$h(x) = W_0 x + \Delta W x, \quad W_0: \text{frozen}, \Delta W = BA$$

$$= W_0 x + BA x$$

$$\left( \begin{bmatrix} W_0 \end{bmatrix} \begin{bmatrix} B \end{bmatrix} \begin{bmatrix} A \end{bmatrix} \right) x = h(x)$$

$$W_0, \Delta W \in \mathbb{R}^{d \times d}$$

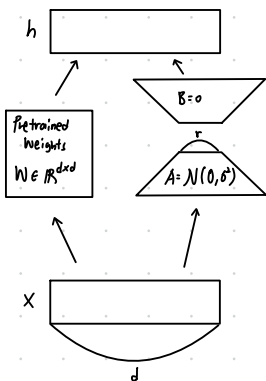
$$B \in \mathbb{R}^{d \times r}$$

$$A \in \mathbb{R}^{r \times K}$$

$$h(x) \in \mathbb{R}^{d \times 1}$$

$$\text{let } r=3, \text{ then } (dxr) + (rxK) \\ = 4,000 \text{ trainable parameters}$$

$r$ : intrinsic rank of the model.  $r \ll K, d$



The hyperparameter we need to choose is the rank  $r$ , since we do not know what the intrinsic rank of the weight matrix is, and if we implicitly remove heavily linearly dependent columns by the  $BA$  decomposition

If rank is too low: It will implicitly delete heavily independent columns too much

too high: we keep too many parameters that are linearly dependent, and waste computation

