

Parameter-efficient fine-tuning (PEFT)

↳ high-level: we want to avoid modifying most of the pretrained model's parameters during fine-tuning.

↳ prompting: requires adjusting zero params to solve a downstream task

What is the sentiment of the below sentence? Answer w/ either "pos" or "neg".

↳ input sentence

output: pos

↳ prompt engineering

↳ limitations:

→ hard to solve very complex reasoning / understanding tasks

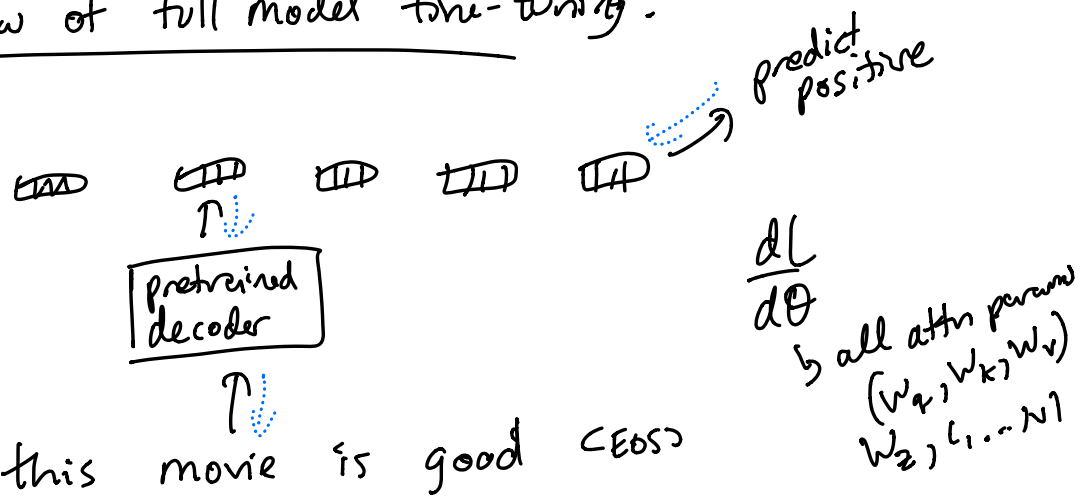
→ requirements for the pretrained model are immense

→ huge-scale pretraining

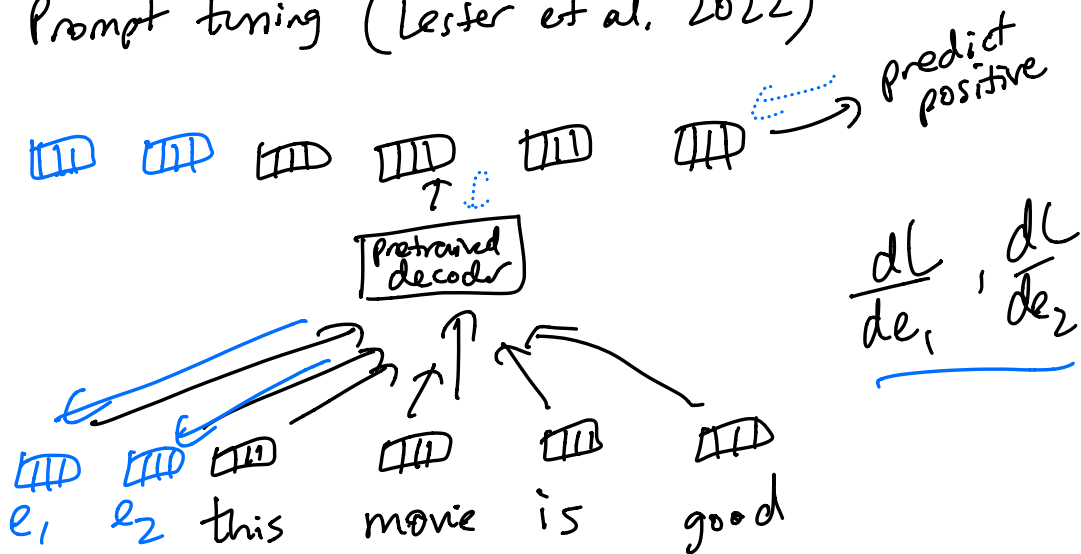
→ high quality large scale instruction tuning

→ RLHF, requires access to very expensive human pref datasets

Review of full model fine-tuning:



Prompt tuning (Lester et al. 2022)



update: keep all pretrained params frozen,

$$\text{only do } e_{1_{\text{new}}} = e_{1_{\text{old}}} - \eta \frac{dL}{de_1}$$

$$e_{2_{\text{new}}} = \dots$$

LoRA (low-rank adaptation):

$$h = f(Wx) \quad \frac{dL}{dW}$$

W is an $m \times n$ matrix
 $\frac{dL}{dW}$ is also $m \times n$

$$W_{\text{new}} = W_{\text{old}} - \eta \underbrace{\frac{dL}{dW}}_{m \times n}$$

having two low-rank matrices A and B
 \swarrow \searrow
 $m \times r$ $n \times r$

$r = \text{rank parameter}$
want $r \ll \ll \ll m, n$

product $\underbrace{AB^T}_{m \times n}$

in LoRA:

$$h = f\left(\underbrace{W_{\text{pre}} + AB^T}_{m \times n} x\right)$$

we compute $\frac{dL}{dA}$, $\frac{dL}{dB}$, much smaller than $\frac{dL}{dW}$

\uparrow \uparrow

$m \times r$ $n \times r$

at the end of LoRA fine-tuning,
we have a separate A, B for each tuned
weight matrix

$$\left. \begin{aligned} W_{\text{new}} &= W_{\text{pre}} + AB^T \\ f(W_{\text{new}}x) \end{aligned} \right\}$$

Q LoRA: quantized LoRA

normal models: FP32



4 bit, 8 bit integer
quantization