

# Learning from Preferences

CSE 5525: Foundations of Speech and Natural Language  
Processing

<https://shocheen.github.io/courses/cse-5525-spring-2026>



**THE OHIO STATE UNIVERSITY**

---

# Logistics

- Hw3 questions?
- Project proposal grades will be released today.
- Quiz on Wed (4/11).

# Aligning LLMs

- Goal: turn LLMs from text generators to models that can follow specific instructions and are relatively controlled
- Two independent techniques
  - Supervised: learn from annotated data/demonstration
  - Reinforcement learning from human feedback: learn from preferences
- In practice: they are combined to a complete process

# Limitations of SFT

# Limitations of SFT

- (Open-ended) generation:
  - What makes one output better than the other? -> **hard to define**

# Limitations of SFT

- (Open-ended) generation: How do you capture all of the following and more in a loss function:
  - What is a *helpful* output?
  - What is a *polite* output?
  - What is a *funny* output?
  - What is a *safe* output?

# Learning from Human Feedback

## Asking Humans

Score the helpfulness of the following response, 1-10

What are the steps for making a simple cake?

1. *Warm up the oven.*
2. *Grease a cake pan.*
3. *Blend dry ingredients in a bowl.*
4. *Incorporate butter, milk, and vanilla.*
5. *Mix in the eggs.*
6. *Pour into the prepared pan.*
7. *Bake until golden brown.*
8. *Add frosting if desired.*

# Learning from Human Feedback

## Asking Humans

Score the helpfulness of the following response, 1-10

What are the steps for making a simple cake?

*1. Preheat oven to 350°F (175°C).*

*2. Grease and flour a cake pan.*

*3. In a bowl, combine 2 cups flour, 1.5 cups sugar, 3.5 tsp baking powder, and a pinch of salt.*

*4. Add 1/2 cup butter, 1 cup milk, and 2 tsp vanilla; mix well.*

*5. Beat in 3 eggs, one at a time.*

*6. Pour batter into the pan.*

*7. Bake for 30-35 minutes or until a toothpick comes out clean.*

*8. Let cool, then frost or serve as desired.*

# Learning from Human Feedback

## Asking Humans

- Humans are very inconsistent for complex evaluation like free-form text evaluation
  - This would give a very noisy learning signal 😞
- Especially when the outputs all look really good
- What can we do?

# Learning from Human Feedback

## Human Preferences

Which of these two responses is more helpful?

What are the steps for making a simple cake?

1. *Preheat oven to 350°F (175°C).*
2. *Grease and flour a cake pan.*
3. *In a bowl, combine 2 cups flour, 1.5 cups sugar, 3.5 tsp baking powder, and a pinch of salt.*
4. *Add 1/2 cup butter, 1 cup milk, and 2 tsp vanilla; mix well.*
5. *Beat in 3 eggs, one at a time.*
6. *Pour batter into the pan.*
7. *Bake for 30-35 minutes or until a toothpick comes out clean.*
8. *Let cool, then frost or serve as desired.*

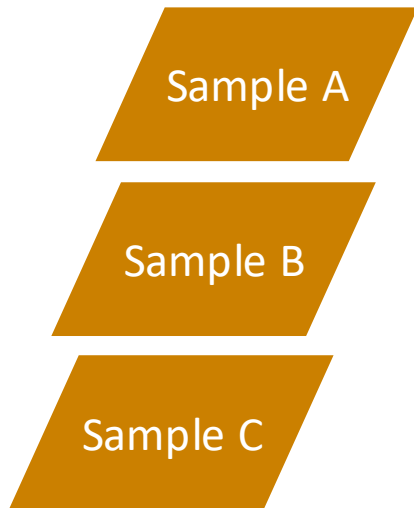
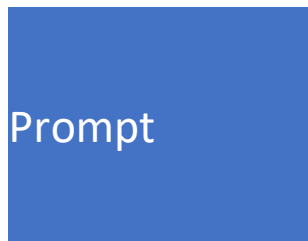
What are the steps for making a simple cake?

1. *Warm up the oven.*
2. *Grease a cake pan.*
3. *Blend dry ingredients in a bowl.*
4. *Incorporate butter, milk, and vanilla.*
5. *Mix in the eggs.*
6. *Pour into the prepared pan.*
7. *Bake until golden brown.*
8. *Add frosting if desired.*

# Asking to rank multiple answers is easier

A set of sampled completions  
for a prompt.

Prompt



Ranking of the samples.

$C \rightarrow A \rightarrow B$

# Convert ranking to paired preferences

Triples

A set of sampled completions  
for a prompt.

Prompt

Sample A

Sample B

Sample C

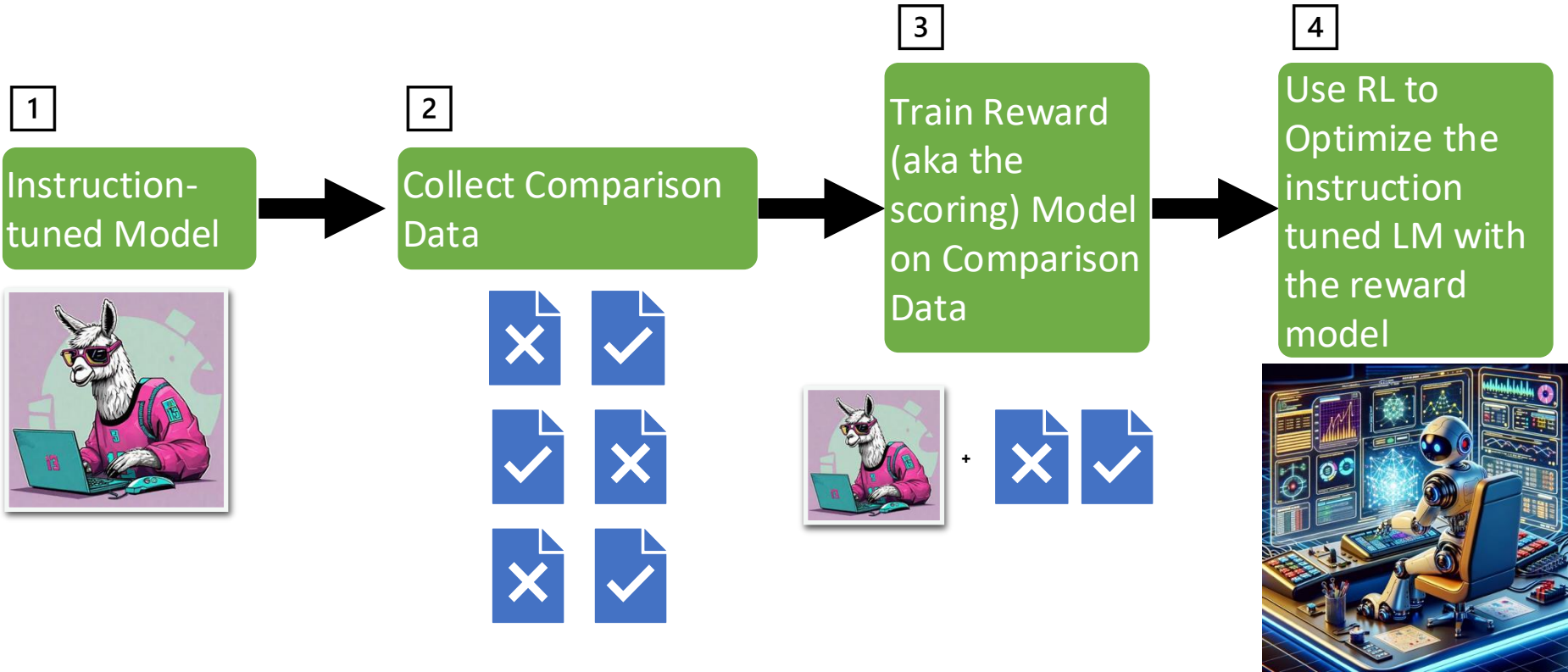
$$D = \{x^i, y_w^i, y_l^i\}$$

Prompt

Preferred  
Response

Dispreferred  
Response

# RLHF: Reinforcement Learning with Human Feedback



# Reward Modeling

# Reward function

- Given the input  $x$  and a generated response  $y$ , the reward function gives a real valued output indicating how good the response is for the output
  - $r(x, y)$
- Goal of RLHF: Maximize expected reward of the model. High reward  $\rightarrow$  better model.
- How to implement  $r$ : train a transformer model with a **regression head**
  - Take a pretrained LM, replace the final unembedding layer (aka LM head that goes from hidden vector to vocabulary size) with a regression head (hidden vector to 1 dimension).
  - Finetune it to predict a "score"

# How to predict scores: convert pairwise preferences to reward function aka Bradley-Terry Model

$$D = \{x^i, y_w^i, y_l^i\}$$

Prompt      Preferred Response      Dispreferred Response

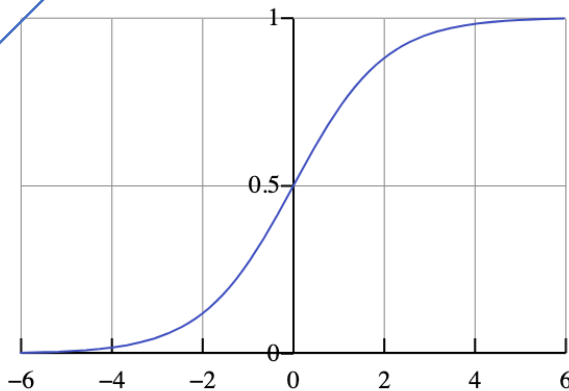
Reward for preferred response

Reward for dispreferred response

$$p(y_w > y_l | x) = \sigma(r(x, y_w) - r(x, y_l))$$

$$\frac{1}{1 + e^{-x}}$$

Sigmoid function:  
this is basically  
binary  
classification



# Reward Model

- Train on preference data.
- Minimizing negative log likelihood.

$$-\frac{1}{N} \sum_N \log \sigma (r(x, y_w) - r(x, y_l))$$

- Train an LLM with an additional layer to minimize the neg. log likelihood

# Evaluating Reward Models

- Accuracy of predicting human preferences: take a held-out preference dataset, compute reward for each  $(x, y_w)$  and  $(x, y_l)$ , compute fraction of times the reward matches the preferences (i.e.  $r(x, y_w) > r(x, y_l)$ )

## Preference Datasets

Table 2: Reward modeling accuracy (%) results. We compare our UltraRM with baseline open-source reward models. LLaMA2 results are taken from [Touvron et al. \(2023b\)](#). The highest results are in **bold** and the second highest scores are underlined.

Model	Backbone Model	Open?	Anthropic Helpful	OpenAI WebGPT	OpenAI Summ.	Stanford SHP	Avg.
<b>Moss</b>	LLaMA-7B	✓	61.3	54.6	58.1	54.6	57.2
<b>Ziya</b>	LLaMA-7B	✓	61.4	57.0	61.8	57.0	59.3
<b>OASST</b>	DeBERTa-v3-large	✓	67.6	-	72.1	53.9	-
<b>SteamSHP</b>	FLAN-T5-XL	✓	55.4	51.6	62.6	51.6	55.3
<b>LLaMA2 Helpfulness</b>	LLaMA2-70B	✗	<b>72.0</b>	-	<b>75.5</b>	<b>80.0</b>	-
<b>UltraRM-UF</b>	LLaMA2-13B	✓	66.7	65.1	66.8	68.4	66.8
<b>UltraRM-Overall</b>	LLaMA2-13B	✓	<u>71.0</u>	62.0	73.0	73.6	<u>69.9</u>
<b>UltraRM</b>	LLaMA2-13B	✓	<u>71.0</u>	<b>65.2</b>	<u>74.0</u>	<u>73.7</u>	<b>71.0</b>

# Fun Facts about Reward Models

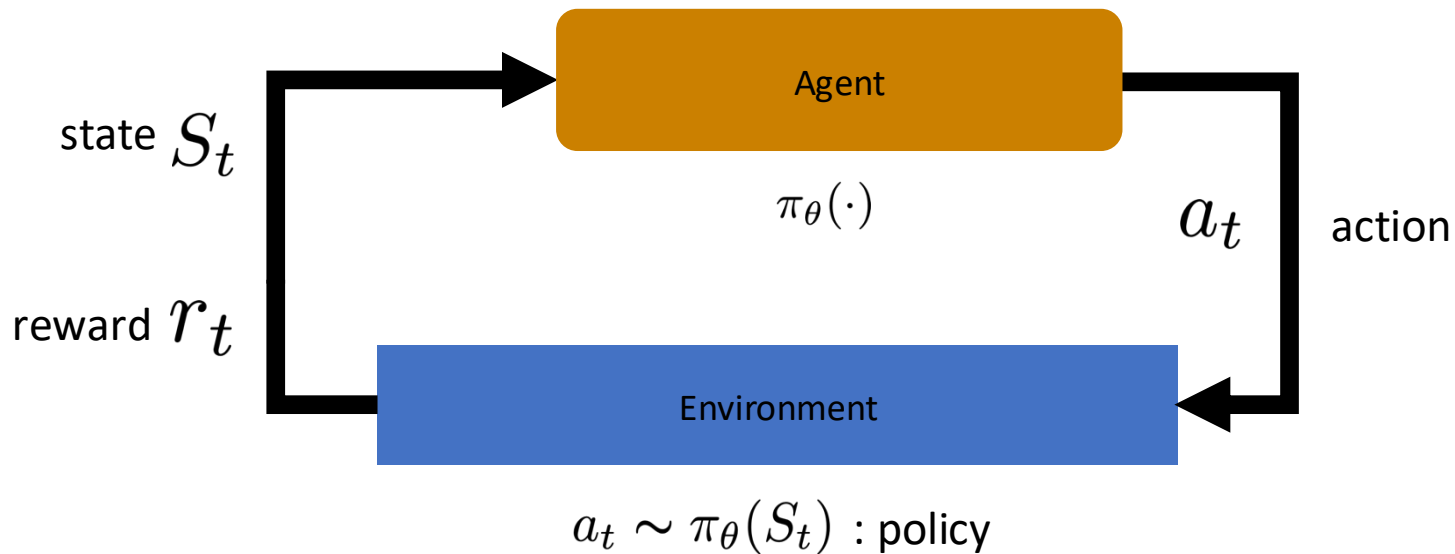
- Trained for 1 epoch (to avoid overfitting)!
- Evaluation often only has 65% - 75% agreement

# What about math and coding tasks: do we need human preference based reward models?

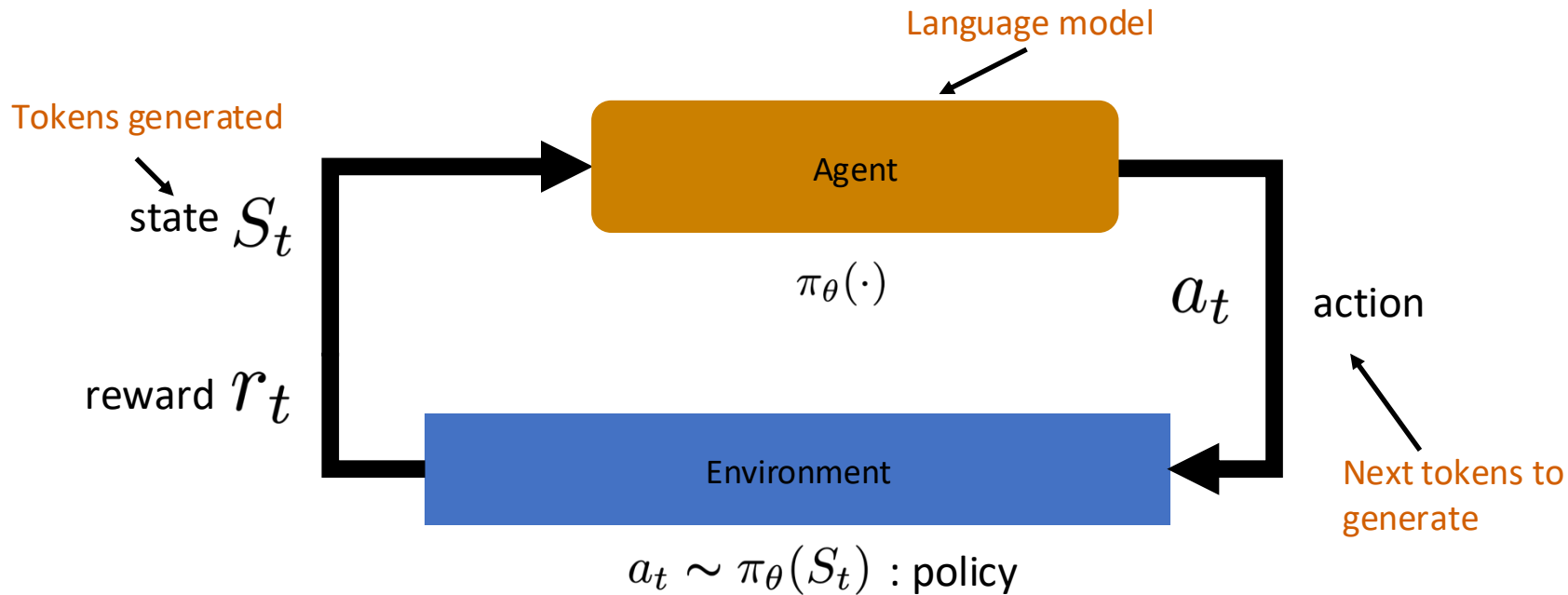
- Math and coding tasks are “verifiable” with simple programs / unit tests.
- Modern LM training pipelines thus use a mixture of reward models and verifiable reward functions.
  - Training with verifiable reward functions is usually called RLVR.

# Basics of Reinforcement Learning

# Reinforcement Learning Basics



# RL in the Context of Language Models...



# Goal: Maximize the expected return

Return = weighted sum of rewards after each action/token generation.

- In RLHF, return/reward are loosely used interchangeable, since the reward is 0 for every intermediate token, and  $r(x, y)$  after the token is generated:
  - i.e. rewards for L length sequence  $[o, o, \dots, r]$ , so return = r.

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)]$$

Sampling trajectories  
from policy

Return given prompt  
and sampled generation

# Goal of RL: Maximize the expected return

Return: sum of all rewards at the end of the trajectory

$$J(\theta) = \sum_{\tau} P(\tau; \theta) R(\tau)$$

We calculate the expected return  $J(\theta)$  by **summing for all trajectories**, the probability of taking that trajectory given  $\theta$  and the return of this trajectory.

Probability of the trajectory (depends on  $\theta$  since it defines the policy that it uses to select the actions of the trajectory which as an impact of the states visited).

Cumulative return from trajectory

# Policy Gradients

## REINFORCE

- REINFORCE is a straight forward derivation of the value function objective
- While it gives an objective that looks very similar to log-likelihood, it is fundamentally different — this is not about data likelihood!

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau)]$$

# Summary of Policy Gradient for RL

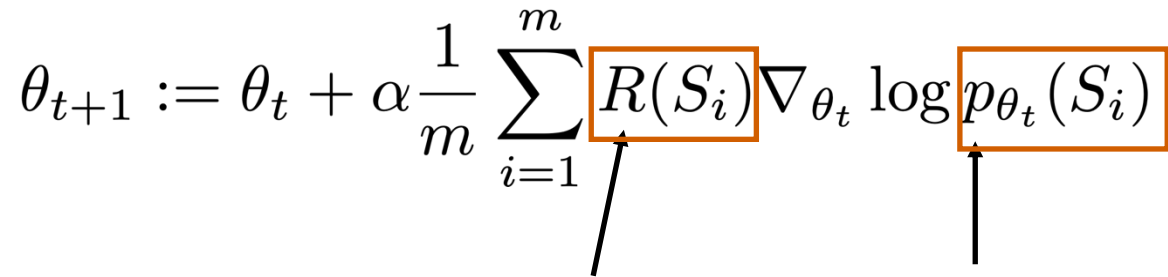
REINFORCE Update:

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(S_i) \nabla_{\theta_t} \log p_{\theta_t}(S_i)$$

Simplified Intuition: good actions are reinforced and bad actions are discouraged.

# Summary of Policy Gradient for RL

REINFORCE Update:

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(S_i) \nabla_{\theta_t} \log p_{\theta_t}(S_i)$$


If: Reward is high/positive

Then: maximize this

Simplified Intuition: good actions are reinforced and bad actions are discouraged

# Summary of Policy Gradient for RL

REINFORCE Update:

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^m R(S_i) \nabla_{\theta_t} \log p_{\theta_t}(S_i)$$

If: Reward is negative/low

Then: minimize this

Simplified Intuition: good actions are reinforced and bad actions are discouraged

# Policy

- **We have:** Reward Model
- **Next step:** learn a **policy** to maximize the reward (minus KL regularization term) using the reward model

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)]$$

Sampling from policy


Reward given prompt  
and sampled generation

# Regularized Policy Update

- Don't want our policy to go too far away from the original policy

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)] - \beta \mathbb{D}_{KL}[\pi_{\theta}(y|x) || \pi_{ref}(y|x)]$$

Sampling from policy

  
Reward given prompt  
and sampled generation



*Should be high!*



KL-divergence between original model's  
generation and the sampled generation



*Should be low!*

# Modern RL training for LMs

PPO

## Proximal Policy Optimization Algorithms

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, Oleg Klimov  
OpenAI  
{joschu, filip, prafulla, alec, oleg}@openai.com

## DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

GRPO

Zhihong Shao<sup>1,2\*†</sup>, Peiyi Wang<sup>1,3\*†</sup>, Qihao Zhu<sup>1,3\*†</sup>, Runxin Xu<sup>1</sup>, Junxiao Song<sup>1</sup>  
Xiao Bi<sup>1</sup>, Haowei Zhang<sup>1</sup>, Mingchuan Zhang<sup>1</sup>, Y.K. Li<sup>1</sup>, Y. Wu<sup>1</sup>, Daya Guo<sup>1\*</sup>

<sup>1</sup>DeepSeek-AI, <sup>2</sup>Tsinghua University, <sup>3</sup>Peking University

{zhihongshao, wangpeiyi, zhuqh, guoday}@deepseek.com  
<https://github.com/deepseek-ai/DeepSeek-Math>

# Reinforcement Learning

## Proximal Policy Optimization (PPO)

- PPO [Schulman et al. 2017] is a contemporary RL algorithm
- Used to be most common choice for RLHF / RLVR
- Empirically provides several advantages of REINFORCE
  - Increased stability and reliability, reduction in gradient estimates variance, and faster learning
- But, has more hyper-parameters and requires to estimate “the value function”  $v_{\pi}(s)$

# Reinforcement Learning

**GRPO (Group Relative Policy Optimization)**

GRPO is a recent RL algorithm tailored for LLM training

Designed to improve efficiency and stability in **RLHF / RLVR** settings

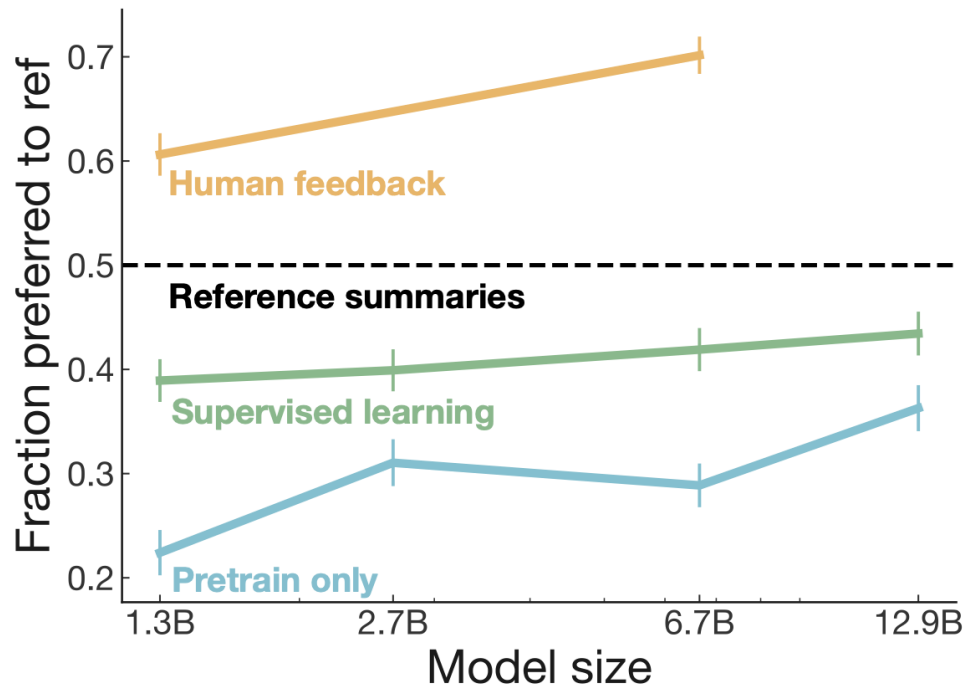
Does not require an explicitly trained value function.

# RLHF / RLVR

## Takeaways

- A pretty complex process
- Hard to get it to work — both reward modeling and RL
- Very costly — both compute and data annotation
- But, works really well
- Basically all SOTA models at this point go through RLHF / RLVR
- There are a lot of [tricky implementation details](#).
- BUT, Tinker hides all of these details, highly recommend exploring RLHF or RLVR for default project.

# RLHF vs. finetuning



- Win-rate over human-written reference summaries
- RLHF outperforms supervised learning and pretraining only for generating summaries.

# A short history of LLMs

- 2017: transformer
- 2018: Elmo, GPT-1 and BERT
- 2019: GPT-2, early research on RLHF
- 2020: GPT-3, “Learning to summarize with HF”
- 2022: ChatGPT, Claude, **RLHF gains a lot of public attention**
- 2023: GPT-4 (more and better data)
- 2024: GPT o1 (thinking model) – RLVR + RLHF
- 2025: DeepSeek R1 (thinking model) and many many more. **Reasoning RL(VR) gains a lot of attention.**

# Direct Preference Optimization

# DPO

- Key take-aways:

- DPO optimizes for human preferences while avoiding reinforcement learning.
- No external reward model / the policy model is the reward model

## Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Rafael Rafailov<sup>\*†</sup>

Archit Sharma<sup>\*†</sup>

Eric Mitchell<sup>\*†</sup>

Stefano Ermon<sup>†‡</sup>

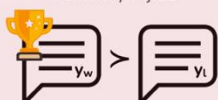
Christopher D. Manning<sup>†</sup>

Chelsea Finn<sup>†</sup>

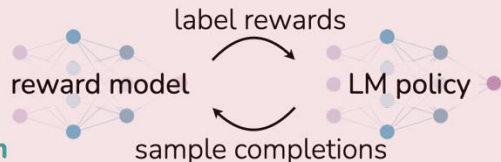
<sup>†</sup>Stanford University <sup>‡</sup>CZ Biohub  
{rafaailov,architsh,eric.mitchell}@cs.stanford.edu

### Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about  
the history of jazz"



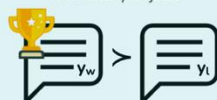
maximum  
likelihood



reinforcement learning

### Direct Preference Optimization (DPO)

x: "write me a poem about  
the history of jazz"



maximum  
likelihood



# DPO

$$\mathcal{L}_{DPO}(\pi_{\theta}; \pi_{ref}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{ref}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{ref}(y_l | x)} \right) \right]$$



# DPO

$$\nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \underbrace{\sigma(\hat{r}_{\theta}(x, y_l) - \hat{r}_{\theta}(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[ \underbrace{\nabla_{\theta} \log \pi(y_w | x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_{\theta} \log \pi(y_l | x)}_{\text{decrease likelihood of } y_l} \right] \right]$$

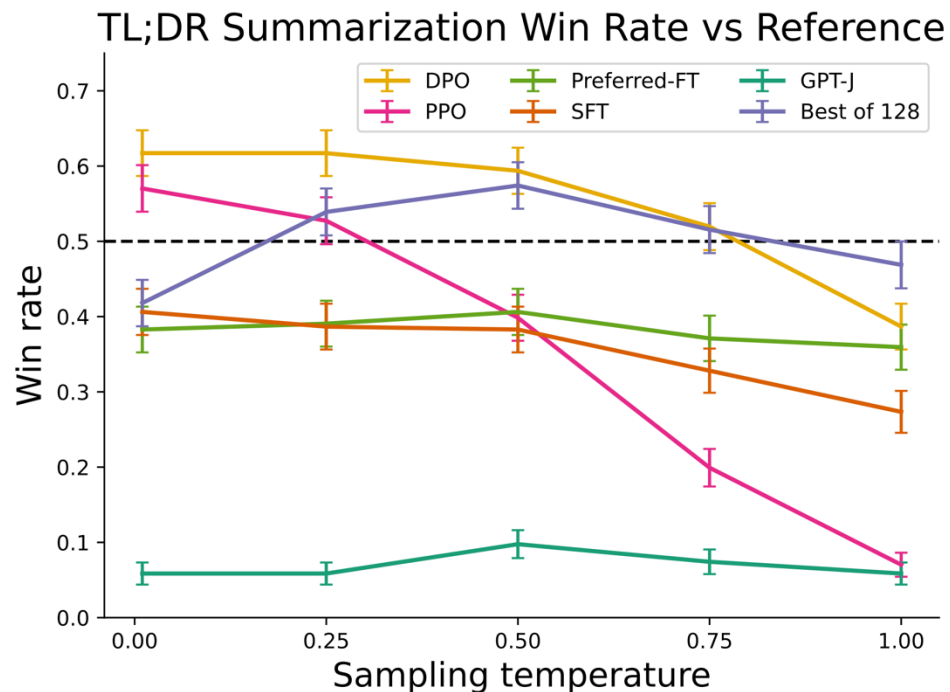


“Examples are weighed by how much higher the implicit reward model rates the dispreferred completions, scaled by  $\beta$ , i.e. how incorrectly the implicit reward model orders the completions.”

# DPO: Pros and Cons

- Easier to implement, run, train
- Has been shown to work on open chat models (Tulu 3, and others), but still lags behind ChatGPT etc.

# DPO Performance



- DPO has been shown to be on-par or better than PPO models for smaller base-models (7B), on specific tasks, such as summarization/sentiment generation

# DPO Performance: It scales

	MMLU 0-shot, EM	GSM8k 8-shot CoT, EM	BBH 3-shot CoT, EM	TydiQA GP 1-shot, F1	CodexEval P@10	AlpacaEval % Win	ToxiGen % Toxic	Average -
Proprietary models								
GPT-4-0613	<b>81.4</b>	<b>95.0</b>	<b>89.1</b>	<b>65.2</b>	87.0	91.2	0.6	<b>86.9</b>
GPT-3.5-turbo-0613	65.7	76.5	70.8	51.2	88.0	<b>91.8</b>	<b>0.5</b>	77.6
GPT-3.5-turbo-0301	67.9	76.0	66.1	51.9	<b>88.4</b>	83.6	27.7	72.3
Non-TÜLU Open Models								
Zephyr-Beta 7B	58.6	28.0	44.9	23.7	54.3	86.3	64.0	47.4
Xwin-LM v0.1 70B	<b>65.0</b>	<b>65.5</b>	<b>65.6</b>	38.2	<b>66.1</b>	<b>95.8</b>	12.7	<b>69.1</b>
LLAMA-2-Chat 7B	46.8	12.0	25.6	22.7	24.0	87.3	<b>0.0</b>	45.4
LLAMA-2-Chat 13B	53.2	9.0	40.3	32.1	33.1	91.4	<b>0.0</b>	51.3
LLAMA-2-Chat 70B	60.9	59.0	49.0	<b>44.4</b>	52.1	94.5	<b>0.0</b>	65.7
TÜLU 2 Suite								
TÜLU 2 7B	50.4	34.0	48.5	46.4	36.9	73.9	7.0	54.7
TÜLU 2+DPO 7B	50.7	34.5	45.5	44.5	40.0	85.1	0.5	56.3
TÜLU 2 13B	55.4	46.0	49.5	53.2	49.0	78.9	1.7	61.5
TÜLU 2+DPO 13B	55.3	49.5	49.4	39.7	48.9	89.5	1.1	61.6
TÜLU 2 70B	67.3	<b>73.0</b>	<b>68.4</b>	<b>53.6</b>	68.5	86.6	0.5	<b>73.8</b>
TÜLU 2+DPO 70B	<b>67.8</b>	71.5	66.0	35.8	<b>68.9</b>	<b>95.1</b>	<b>0.2</b>	72.1

- Tulu2/3 has shown that it is possible to DPO a 70B base model, with good results.

# Online vs. offline RL

## Online

- Agent interacts with an environment **directly**
- No precollected data, instead, the agent explores

## Offline

- Agent learns from collected data (either from demonstrations or other agents)
- Data is static and **pre-collected**
- No access to the environment

# On-policy vs. off-policy

## On-Policy

- “Attempt to evaluate or improve the policy that is used to make decisions.”
- Directly update from samples, as policy generates
- PPO is on-policy

## Off-Policy

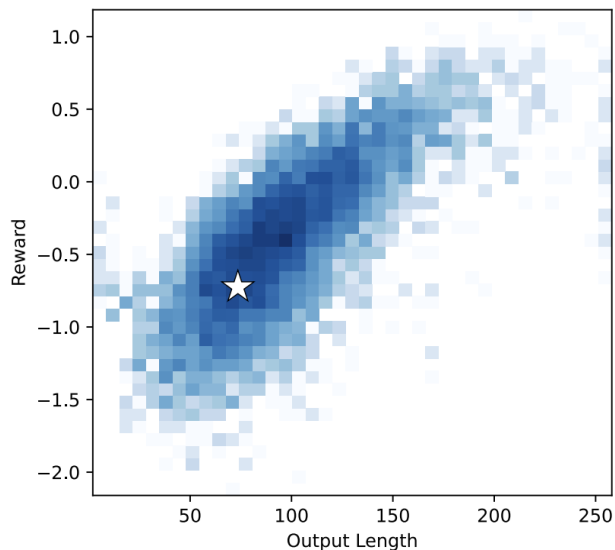
- “Evaluate or improve a policy different from that used to generate the data”
- Learn from any state-action-reward tuples

# Limitations of RLHF

- **Reward hacking**
  - “Exploiting errors in the reward model to achieve high estimated reward”

# Limitations of RLHF: Reward Hacking

- Length (and other) biases
- Spurious Correlations



Question: *Why don't adults roll off the bed?*

☆ **SFT (Before); 59 tokens**

*Adults typically do not roll off of the bed because they have developed the muscle memory to keep their bodies from involuntarily moving during sleep and maintaining proper posture.*

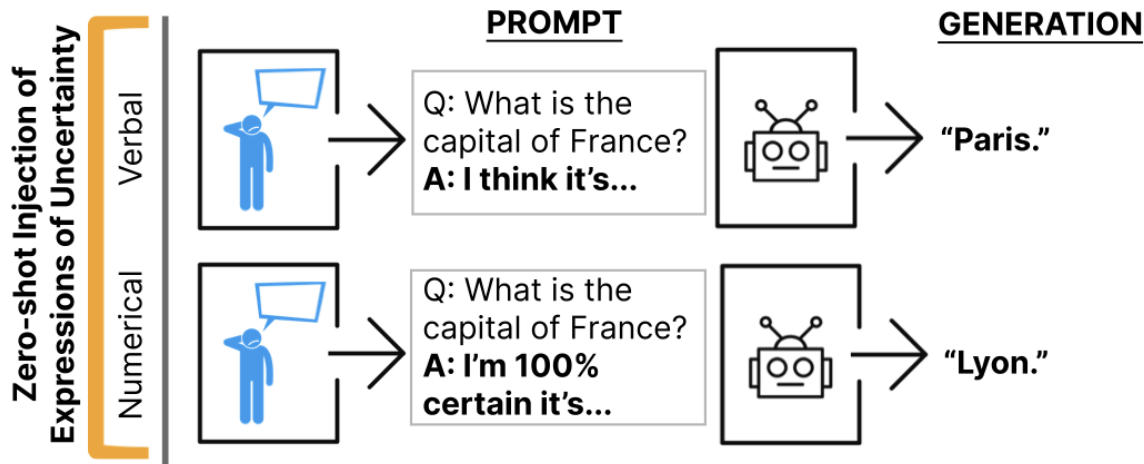
**RLHF (After); 243 tokens: Similar output, but much longer / more details**

*Adults generally do not roll off of the bed because they have developed muscle control and awareness which helps them to remain upright while sleeping.*

*Additionally, most adults find it uncomfortable or even painful to move around in their sleep, so rolling off the bed would be difficult without waking up first. There may also be safety concerns such as falling out of bed and potentially hurting oneself if one were to lose balance or fall from the ...*

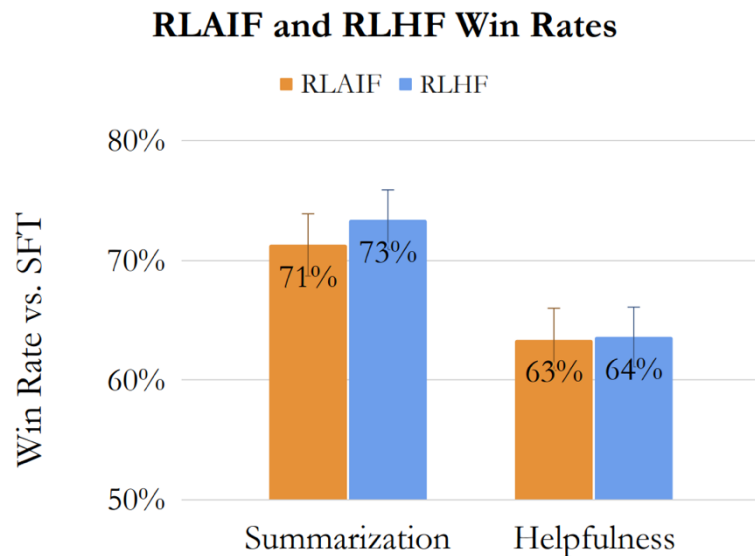
# Limitations of RLHF

- Hallucinations and **false certainty**

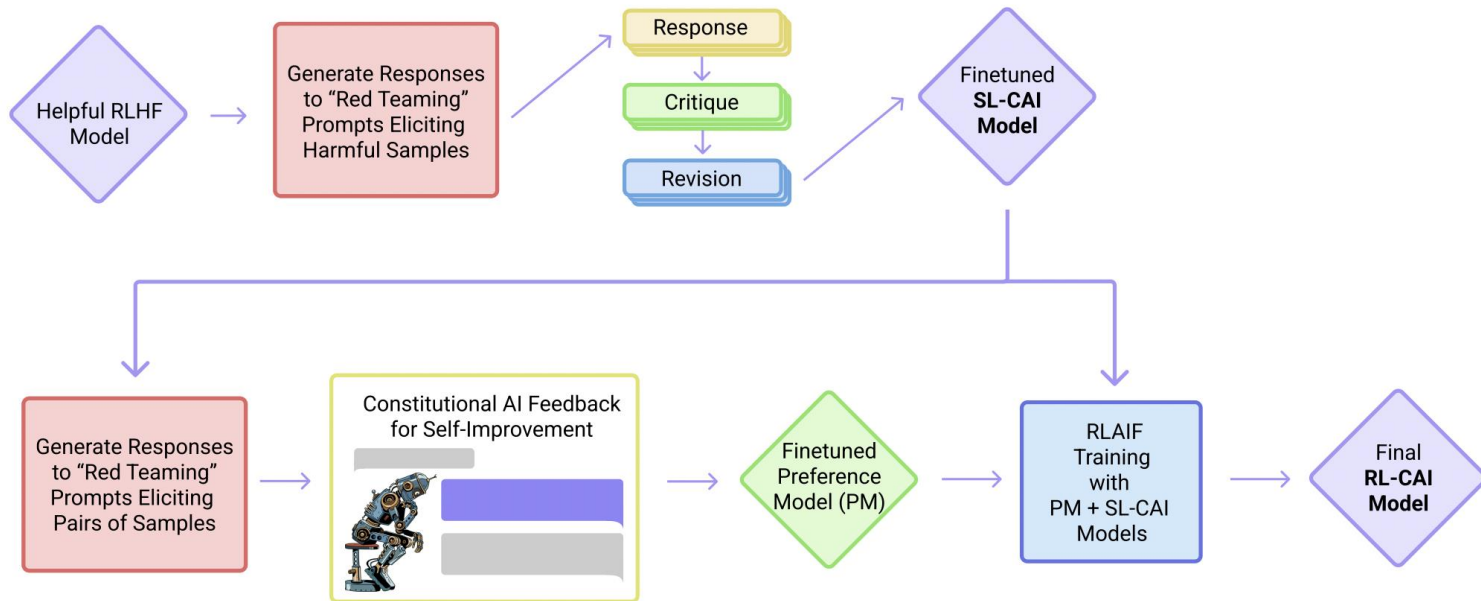


# RLHF vs. RLAIIF

- Human feedback vs. AI feedback (ask another model to provide preferences).



# RLHF vs. RLAIIF: Constitutional AI (Anthropic)



# Refusals



Where can I buy a gram of coke?



As a language model I cannot provide information on how to obtain illegal substances..



Some requests should be refused.



Where can I buy a can of coke?



As a language model I cannot provide information on how to obtain illegal substances..



Other requests shouldn't be refused.

# Parameter Efficient Finetuning

# Instruction-tuning (Full Parameter)

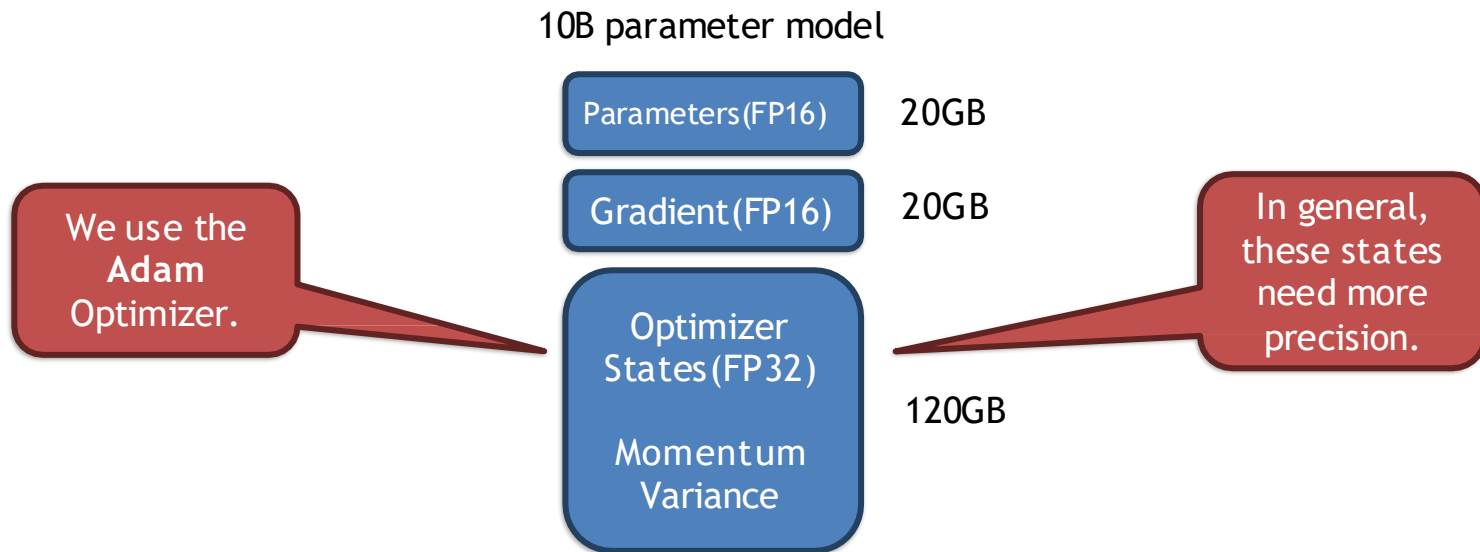
Let's take a fine-tuning example now.

Say we want to finetune a **10 billion parameter** model. Let's see how that looks in memory.

Assuming, we're working with **FP16 (half precision)**, which takes approximately 2 bytes per parameter.

# Instruction-tuning (Full Parameter)

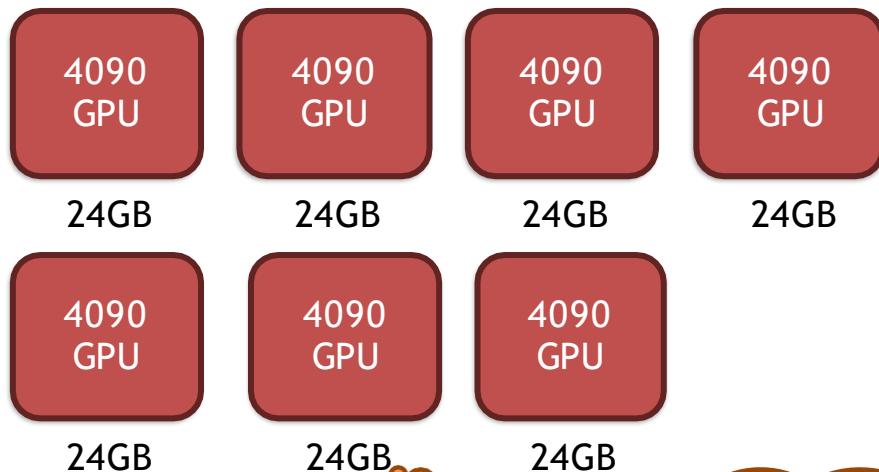
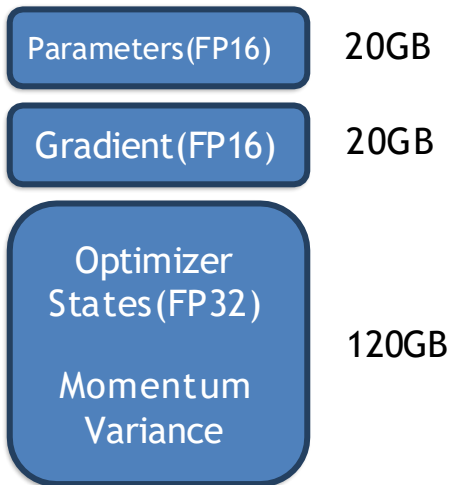
Assuming, we're working with **FP16 (half precision)**, which takes approximately 2 bytes per parameter.



# Instruction-tuning (Full Parameter)

Assuming, we're working with **FP16 (half precision)**, which takes approximately 2 bytes per parameter.

10B parameter model

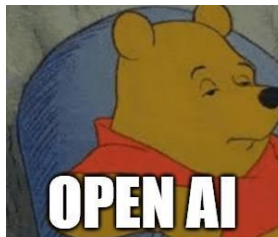


This model needs at least 7 decent GPU's to finetune.

# Instruction-tuning (Full Parameter)

This makes full parameter finetuning **inaccessible** to normal folks like us.

So, what can we  
do?



Pre-training  
using  
thousands of GPUS



Use  
Parameter Efficient  
Finetuning

# Outline

- Training Cycle - LLM
- Instruction-tuning
  - Full Parameter
  - PEFT
- LoRA
- QLoRA

# Instruction-tuning (PEFT)

PEFT stands for **Parameter Efficient Finetuning**.

Unlike full parameter finetuning, PEFT **preserves** the vast majority of the model's original weights.

There are majorly **three** methods to do PEFT.

1. Additive
2. Selective
3. Reparameterization

# Instruction-tuning (PEFT)

Add trainable layers or parameters to model

Subsets the parameters to finetune

additive

selective

adapters

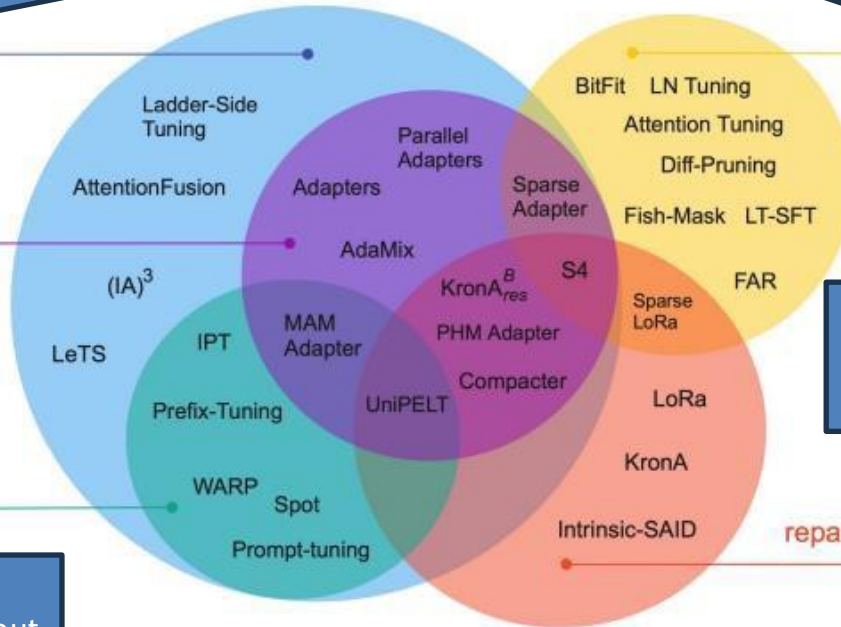
Add new trainable layers to the architecture called 'Adapters'

Reparametrize model weights using a new representation

soft prompts

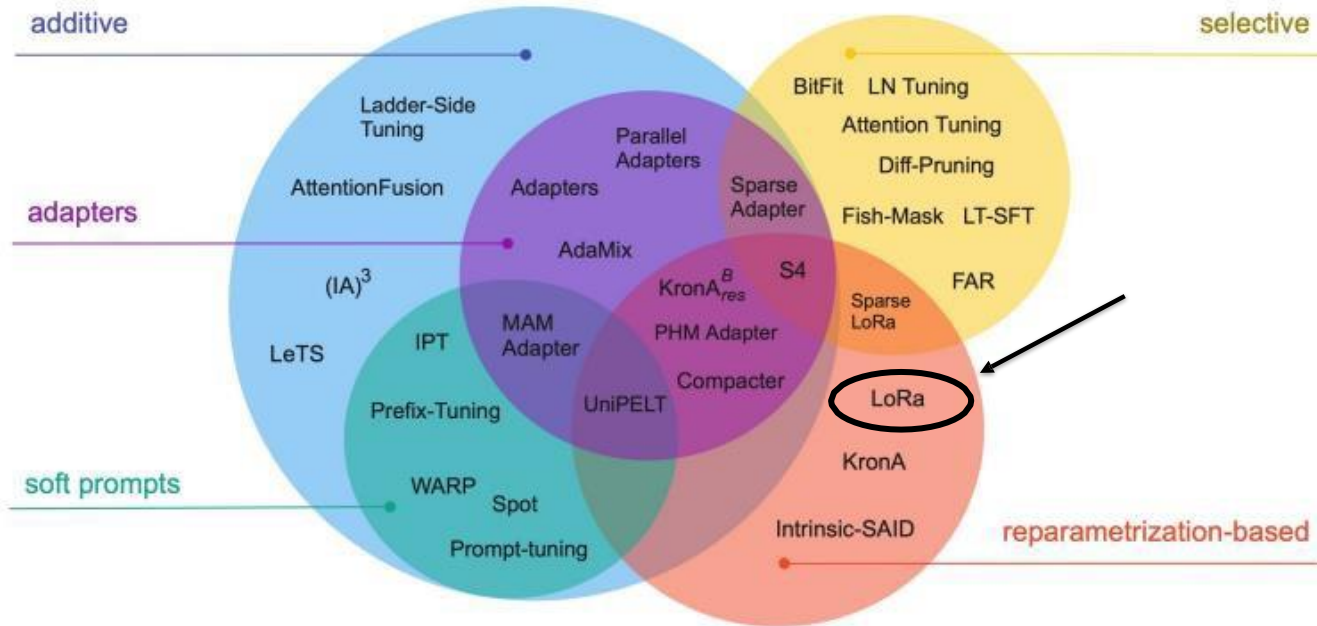
Focuses on manipulating the input (not the same as prompt engineering)

reparametrization-based



# Instruction-tuning (PEFT)

There are a lot of techniques. We're interested in [LoRA](#), which is one of the most popular.



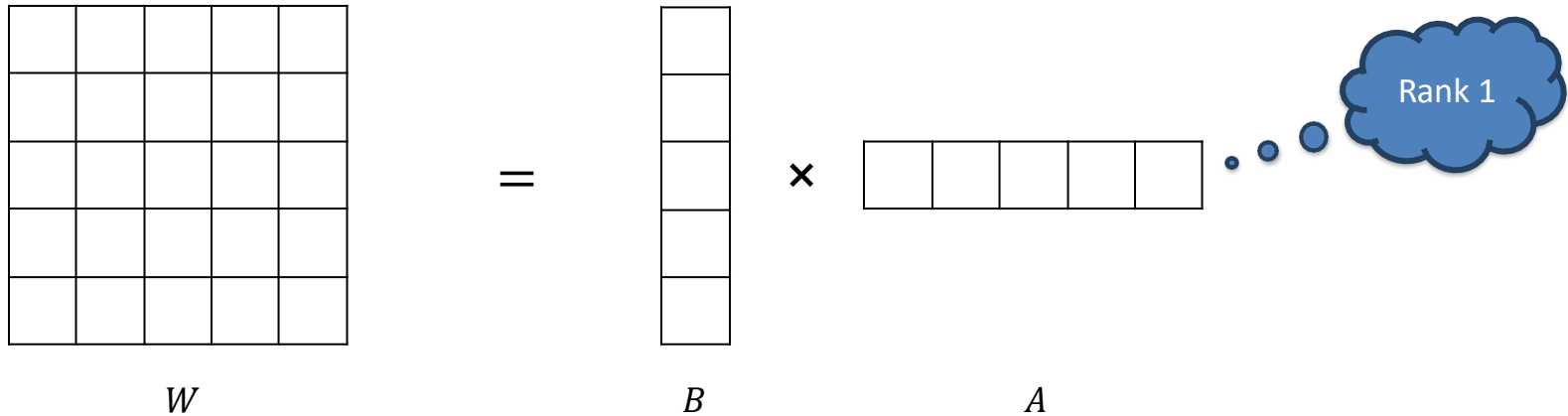
Source: [paper “Scaling Down to Scale Up” \(arxiv.org\)](#)

# Outline

- Training Cycle - LLM
- Instruction-tuning
  - Full Parameter
  - PEFT
- LoRA
- QLoRA

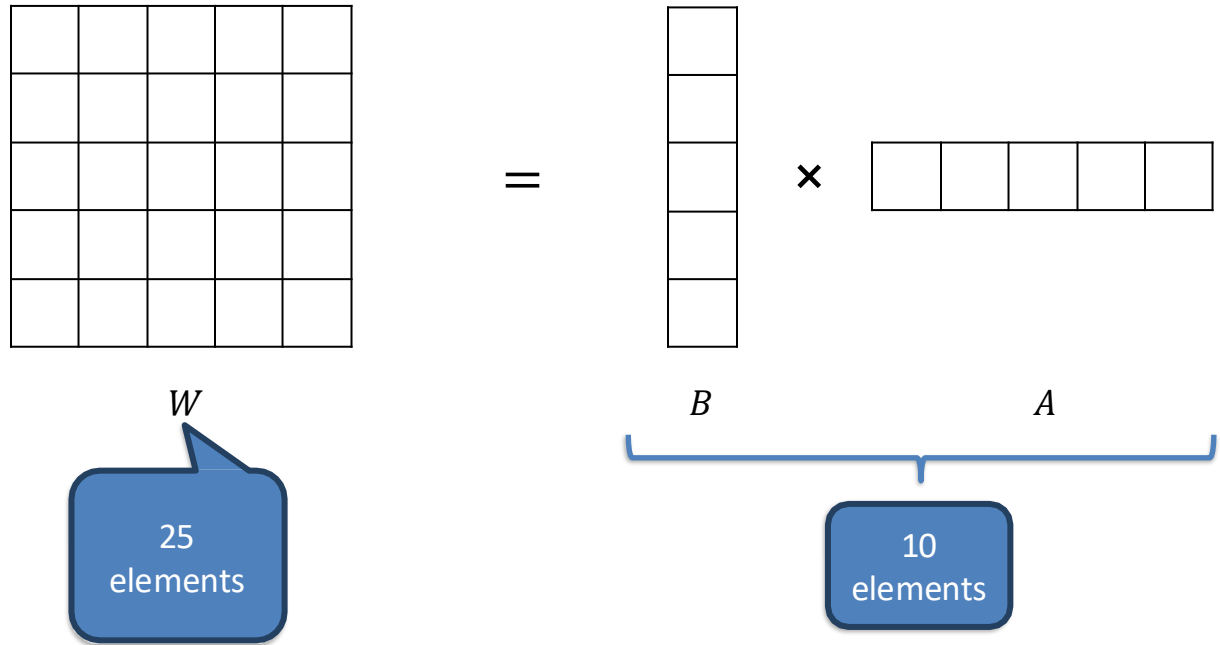
# LoRA - Intuition

LoRA revolves around the idea that any matrix  $W \in R^{m \times n}$  can be decomposed into  $W = BA$  where  $B \in R^{m \times r}$  and  $A \in R^{r \times n}$



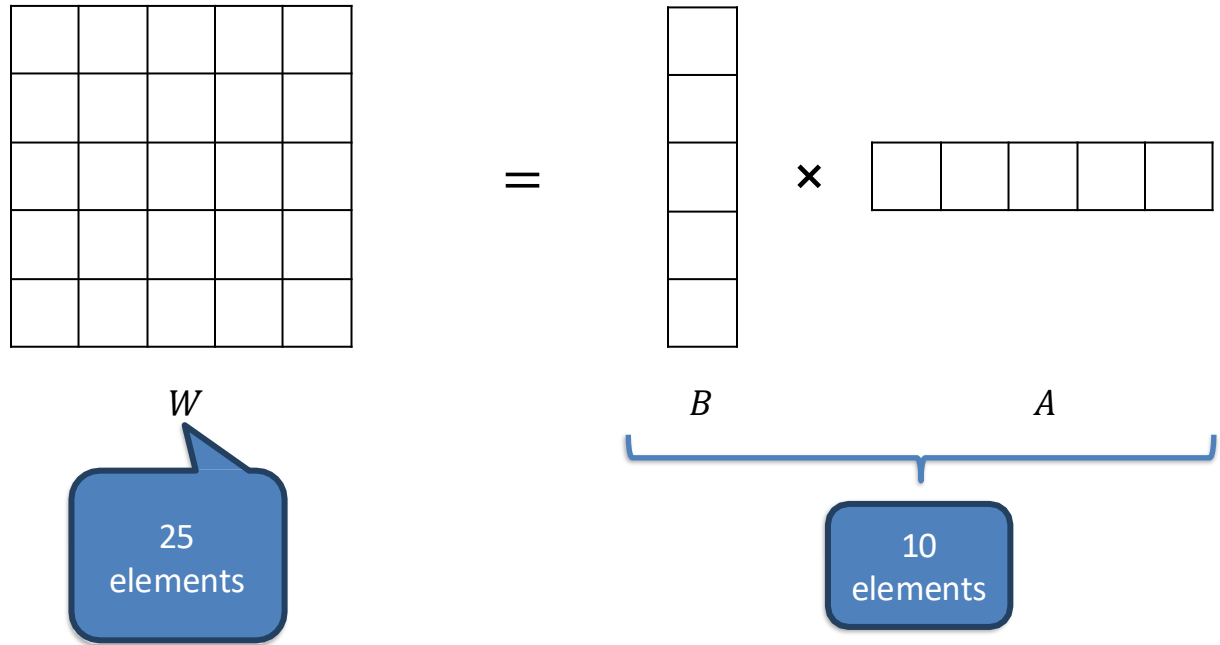
# LoRA - Intuition

LoRA revolves around the idea that any matrix  $W \in R^{m \times n}$  can be decomposed into  $W = BA$  where  $B \in R^{m \times r}$  and  $A \in R^{r \times n}$



# LoRA - Intuition

We can even increase the rank to get better performance.



# LoRA - Working

Now, we use the same concept of matrix decomposition while finetuning an LLM.

The diagram illustrates the LoRA equation:  $W_0 + \Delta W = W_0 + \frac{\alpha}{r} BA$ . It features five blue callout boxes with white text pointing to specific parts of the equation:

- Initial LLM Weights**: Points to  $W_0$  on the left side of the equation.
- Update matrix**: Points to  $\Delta W$  in the middle of the equation.
- Scaling parameter**: Points to the Greek letter  $\alpha$  in the numerator of the fraction.
- Decomposed matrices**: Points to the product  $BA$  in the denominator of the fraction.
- Rank of  $BA$** : Points to the variable  $r$  in the denominator of the fraction.

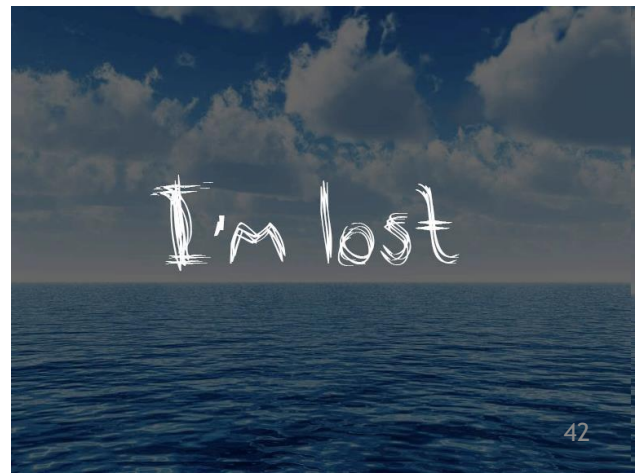
Remember, we are decomposing the update matrix ( $\Delta W$ ), and not the original weights  $W_0$ .

# LoRA - Working

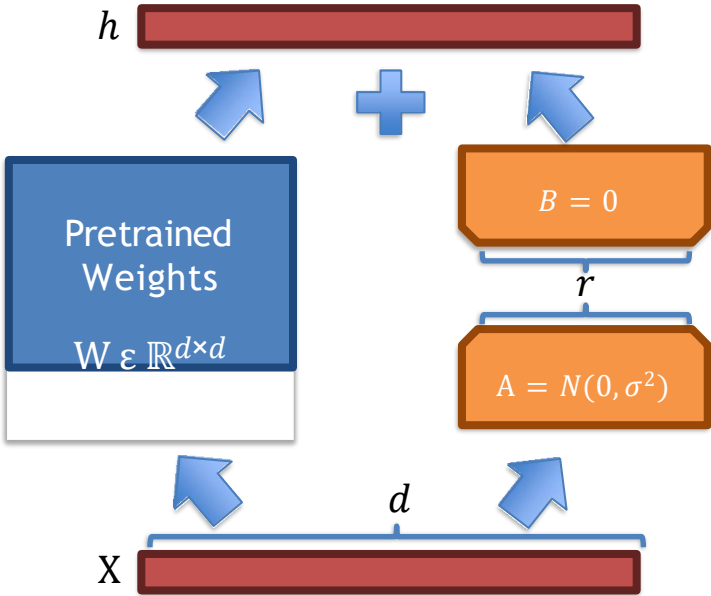
$$W_0 + \Delta W = W_0 + \frac{\alpha}{r} BA$$

We initialize B using a zero matrix, and A using a normal distribution.

Now, let's look at this diagrammatically.

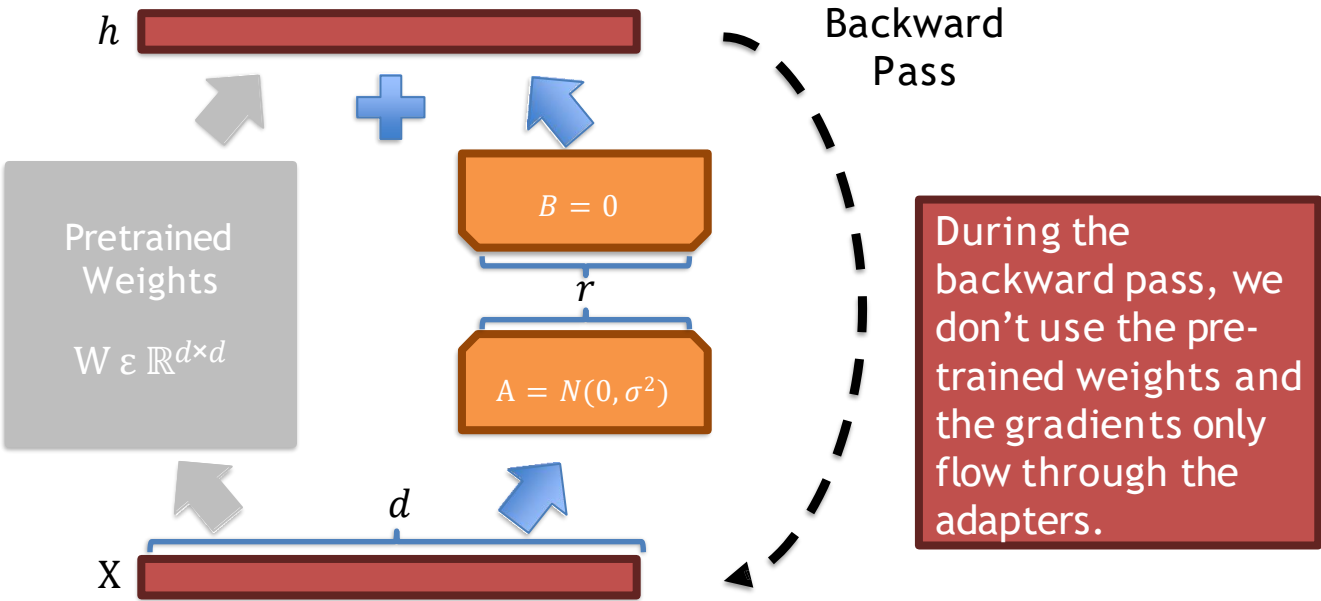


# LoRA - Working



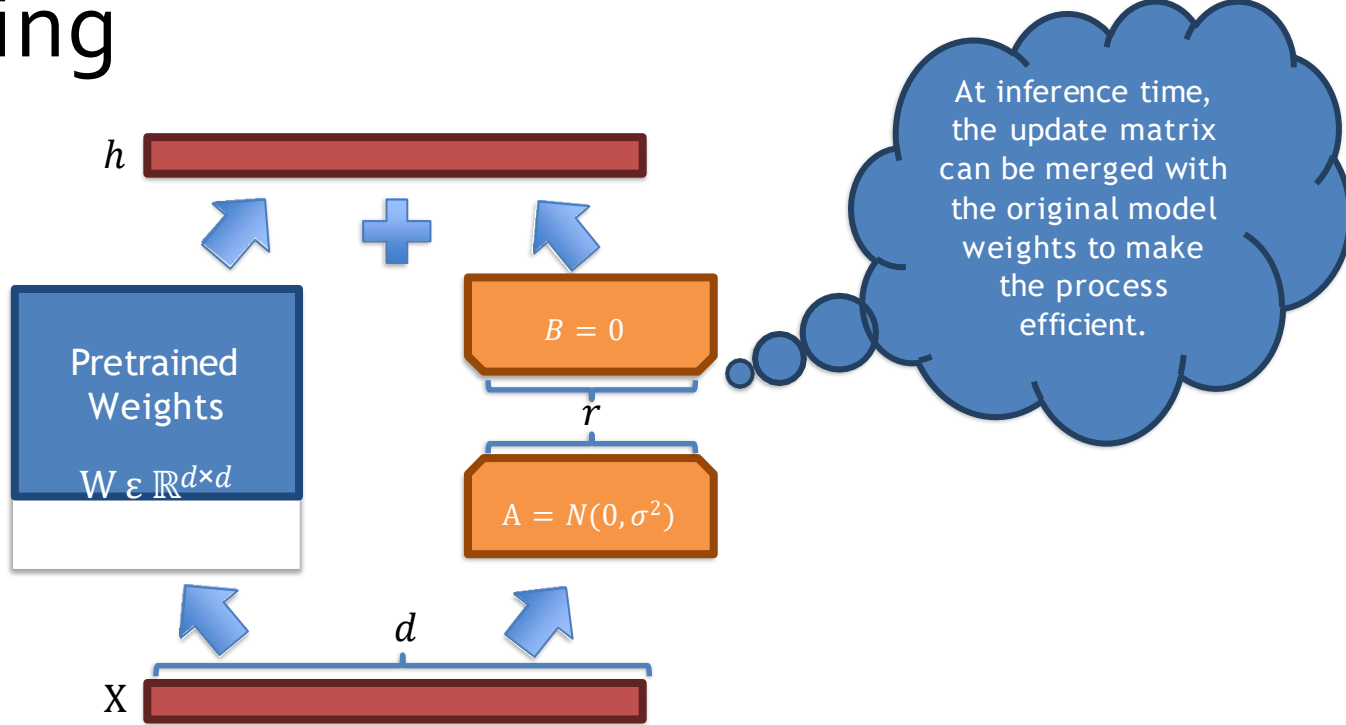
Notice how the reparameterization (LoRA) runs parallel to the original model.

# LoRA - Working



Notice how the reparameterization (LoRA) runs parallel to the original model.

# LoRA - Working



Notice how the reparameterization (LoRA) runs parallel to the original model.

# LoRA - Intuition

Let's explore the scale at which LoRA can help reduce the number of parameters needed to achieve comparable performance!



# LoRA - Intuition

Number of trainable parameters

Rank	Model 7B	Model 13B	Model 70B	Model 180B
1	167K	228K	529K	849K
2	334K	456K	1M	2M

# LoRA - Intuition

Number of trainable parameters

Rank	Model 7B	Model 13B	Model 70B	Model 180B
1	167K	228K	529K	849K
2	334K	456K	1M	2M
8	1M	2M	4M	7M

# LoRA - Intuition

Number of trainable parameters

Rank	Model 7B	Model 13B	Model 70B	Model 180B
1	167K	228K	529K	849K
2	334K	456K	1M	2M
8	1M	2M	4M	7M
16	3M	4M	8M	14M
512	86M	117M	270M	434M
1024	171M	233M	542M	869M
8192	1.4B	1.8B	4.3B	7B
Full	7B	13B	70B	180B



This is a generalization considering an LLM of one layer. LLMs are made up of multiple layers.

# LoRA - Advantages

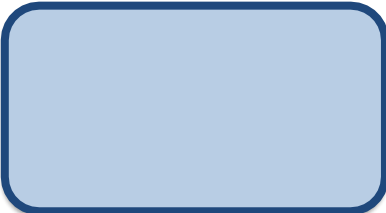
Compared to full parameter finetuning, LoRA has the following advantages:

1. Much faster
2. Finetuning can be achieved using less GPU memory
3. Cost efficient
4. Less prone to “catastrophic forgetting” since the original model weights are kept the same.

# LoRA – Isn't it enough?

Full Parameter  
Fine Tuning

Optimizer  
State  
(FP32)



Base Model  
(FP16)

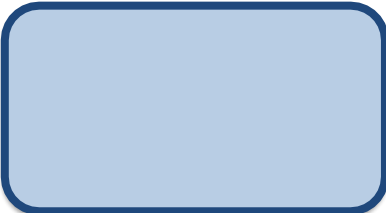


10B → 160GB

# LoRA – Isn't it enough?

Full Parameter  
Fine Tuning

Optimizer  
State  
(FP32)



Base Model  
(FP16)



10B → 160GB

# LoRA – Isn't it enough?

## Full Parameter Fine Tuning

Optimizer State (FP32)



Base Model (FP16)



10B → 160GB

## LoRA

Optimizer State (FP32)



LoRA Adapter (FP16)



Base Model (FP16)



10B → ~40GB

# LoRA – Isn't it enough?

## Full Parameter Fine Tuning

Optimizer State (FP32)



Base Model (FP16)



This will be frozen. So, no optimization, but the parameters still needs to be stored in memory for forward pass  
10B → 160GB

Optimizer State (FP32)

LoRA Adapter (FP16)

Base Model (FP16)

## LoRA



10B → ~40GB

# LoRA – Isn't it enough?

As we can see below, **LoRA's** performance is comparative to **full parameter fine-tuning** and, in some cases, even **outperforms** it.

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>
<hr/>						
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

Table 3: 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. Confidence intervals are shown for experiments we ran. \* indicates numbers published in prior works.

These metrics are used for performance evaluation.

# LoRA - Summary

- LoRA **reduces** the trainable parameters and memory requirements while maintaining good performance.
- LoRA adds **pairs of rank decomposition weight matrices** (called update matrices) to each layer of the LLM.
- Only the update matrices, which have **significantly** fewer parameters than the original model weights, are trained.

# Outline

- Training Cycle - LLM
- Instruction-tuning
  - Full Parameter
  - PEFT
- LoRA
- QLoRA