

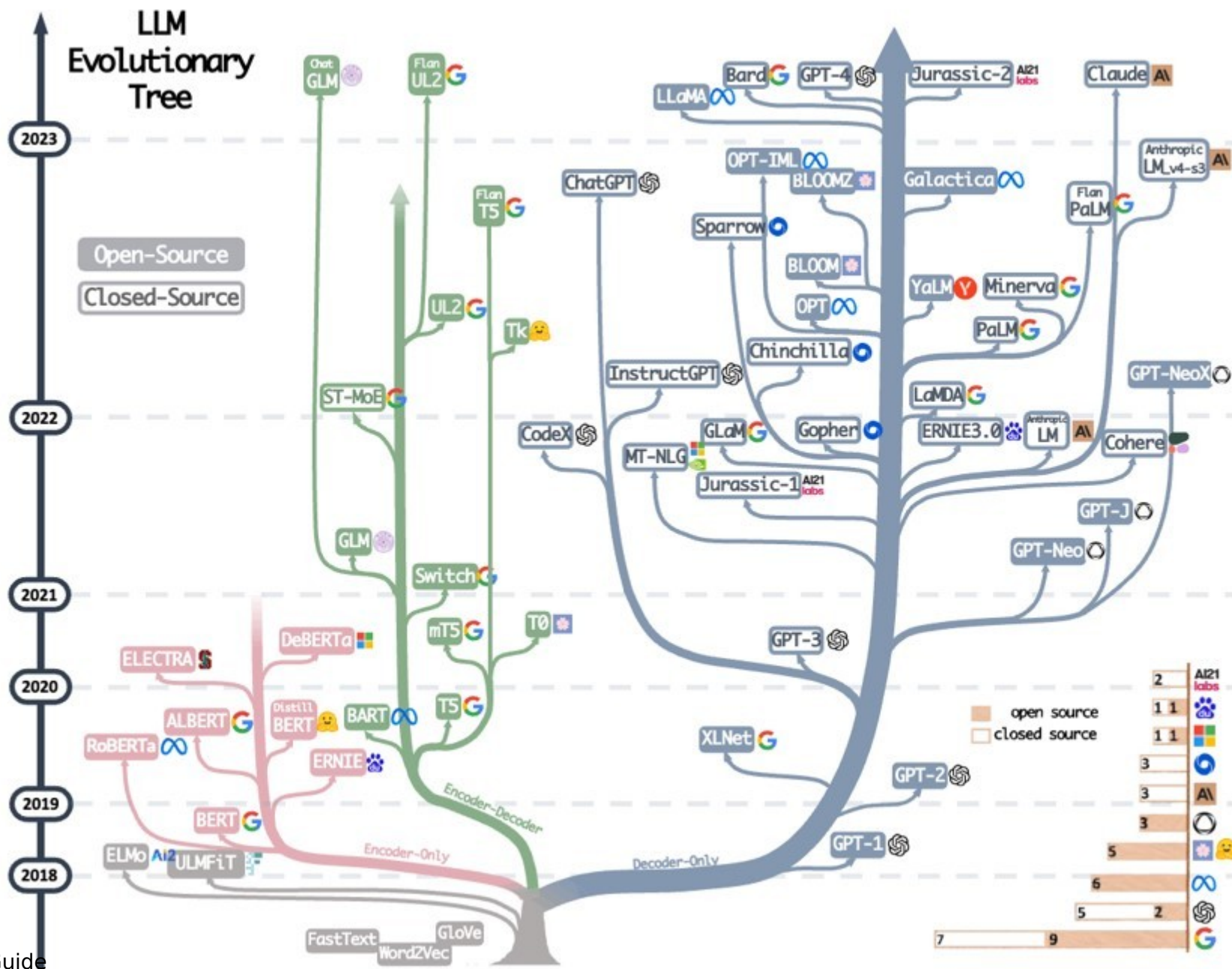
Fine-Tuning LLMs

Customizing LLMs for your needs

Ehsan Kamalinejad

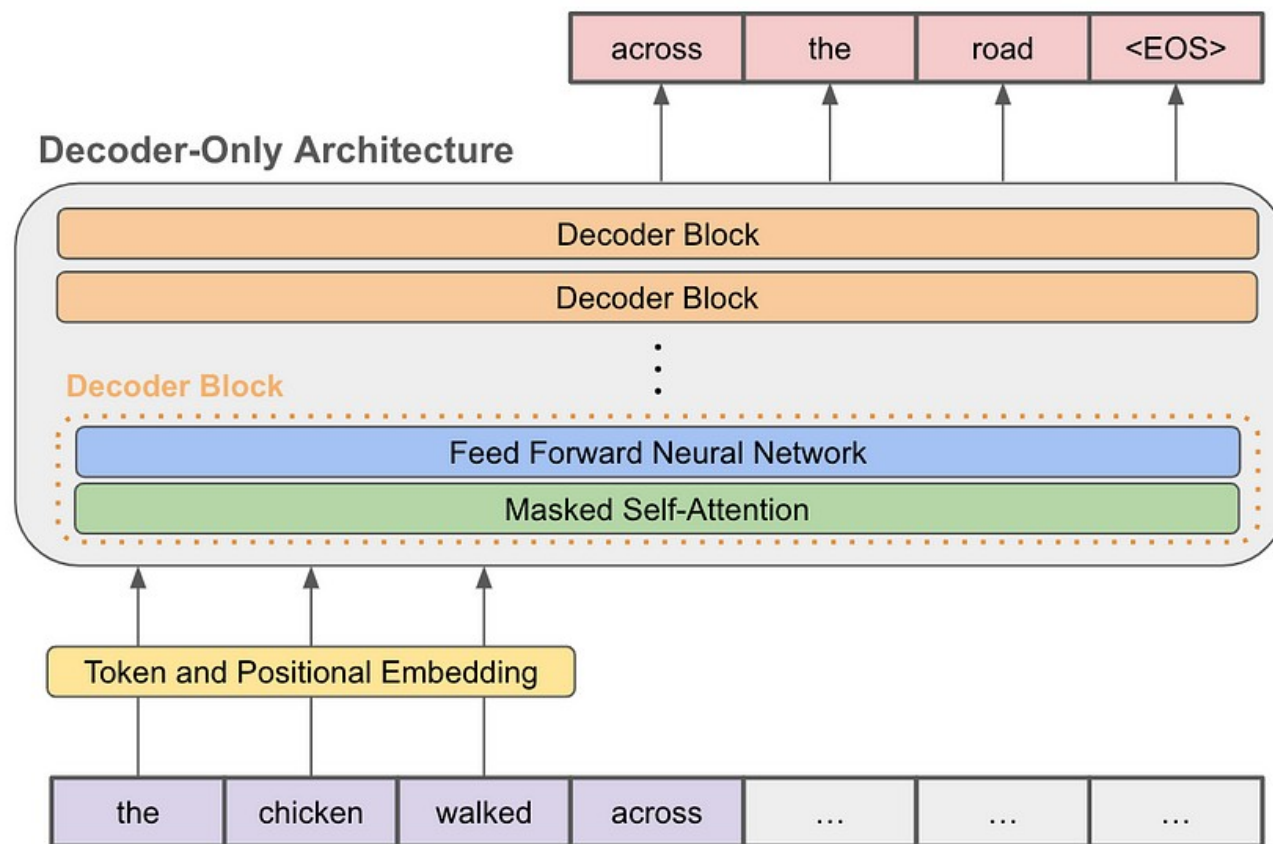
LLMs

Modern language models are pretrained through self-supervised learning methods such as MLM or CLM



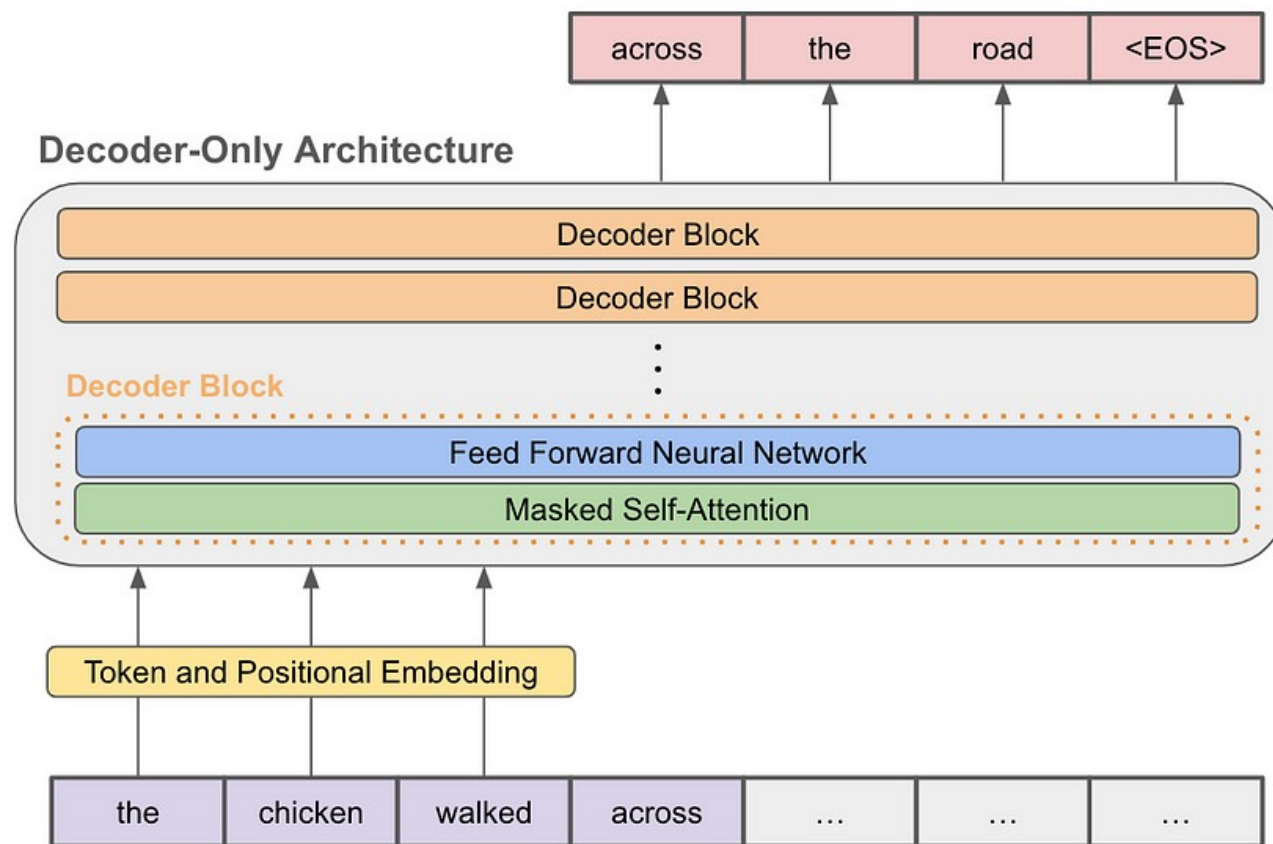
Pretraining LLMs

- Most of the modern LLMs are pretrained through autoregressive next token prediction.



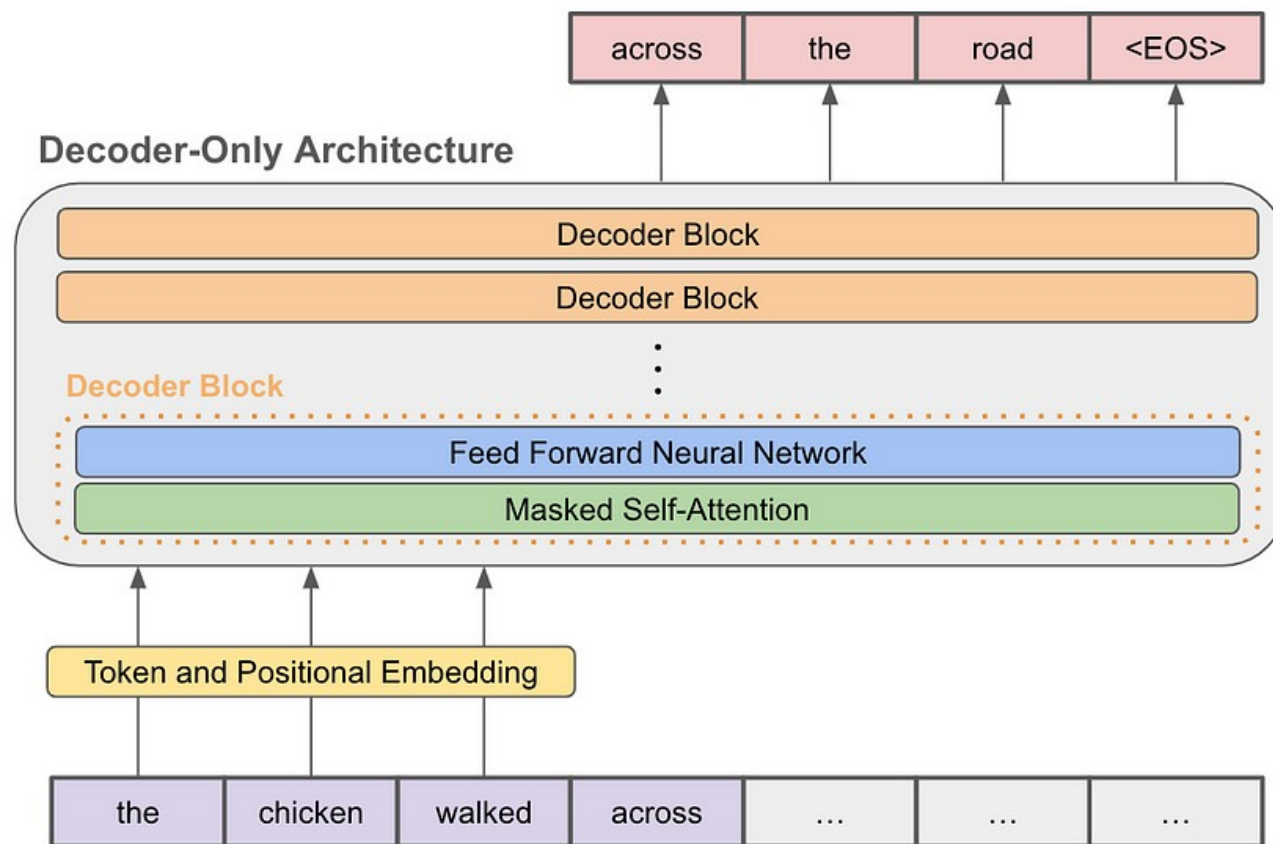
Pretraining LLMs

- Most of the modern LLMs are pretrained through autoregressive next token prediction.
- This allows to train on any kind of text (language, code, tabular, etc.)
- The whole internet is your playground!



Pretraining LLMs

- Most of the modern LLMs are pretrained through autoregressive next token prediction.
- This allows to train on any kind of text (language, code, tabular, etc.)
- The whole internet is your playground!
- **The task is hard!**



2 + 2 = <?>

Must know math

The president of US in 1983 was <?> Must know history

Foundational Models

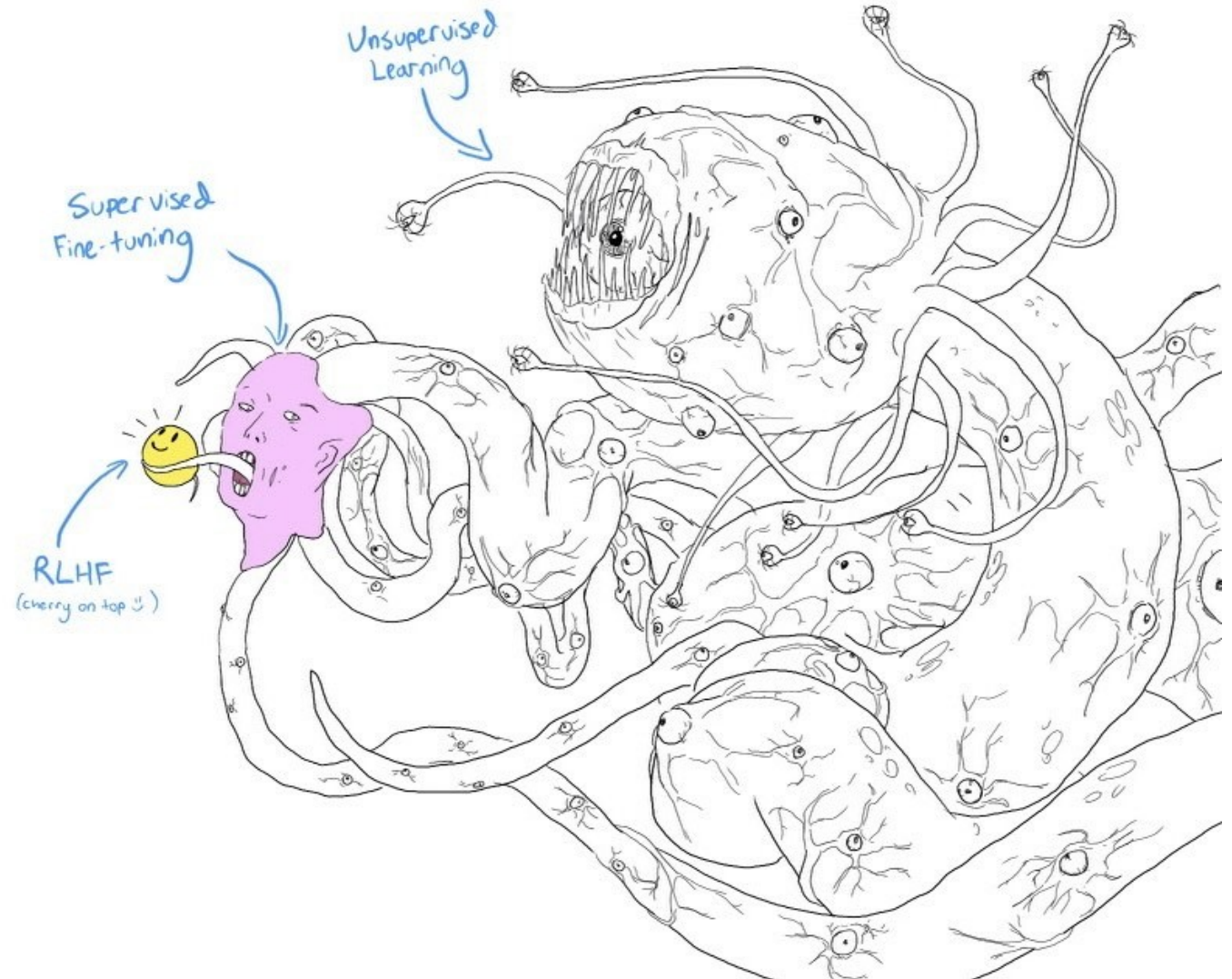
- Foundational LLMs are pretrained on large data
- They are versatile and adaptable base for multiple purposes
- In their raw form they do not follow instruction or answer question
- While the knowledge is stored in them, one needs to extract it

Fine-Tuning LLM

To extract the knowledge stored in LLMs, one can do fine-tuning

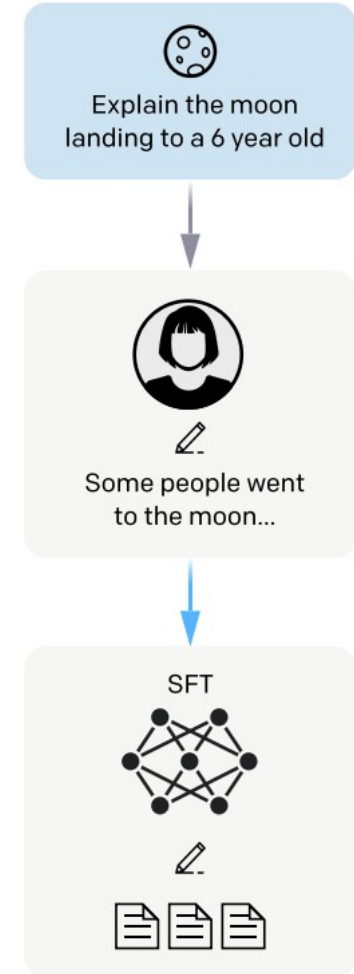
Two main fine-tuning methods:

- Supervised Fine-Tuning (SFT)
- Reinforcement Learning with Human Feedback (RLHF)



Supervised Fine-Tuning

- Create a diverse prompt dataset compatible with your task
- Sample the prompt dataset
- Give sample prompts to labelers to create the desired output
- Fine-tune the base model with the query-response pairs



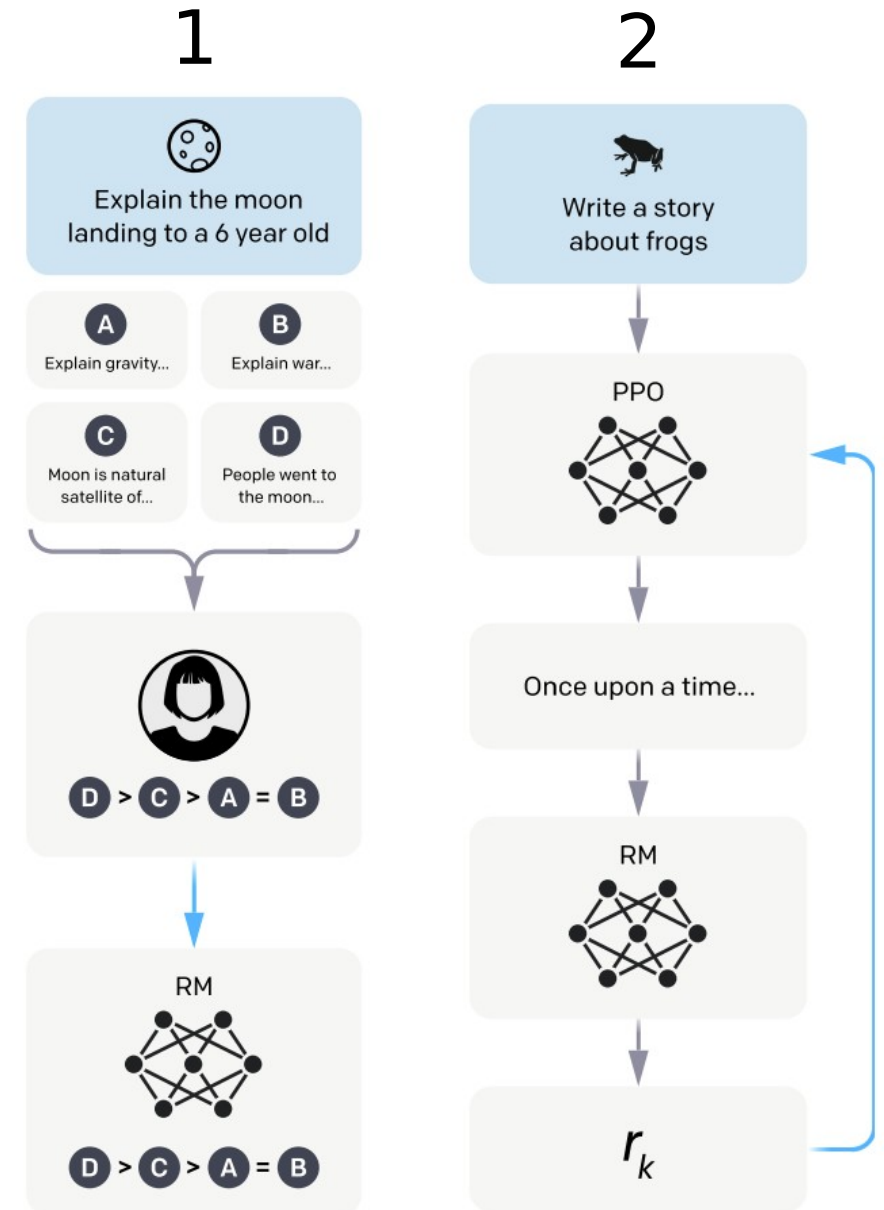
RLHF

1. Train a reward model:

- Collect several responses to a query from the model
- Let the labeler rank the responses according to your definition of a better response
- Train a reward model that can estimate alignment of a response with your desire

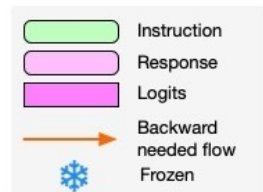
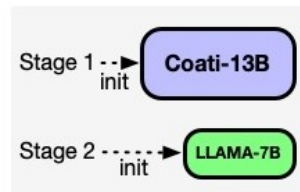
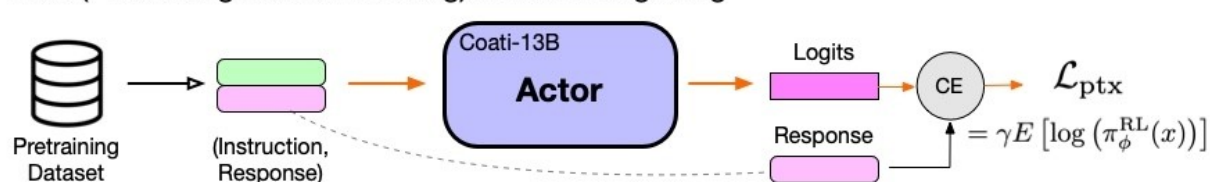
2. Train a policy model through PPO

- The reward model guides the base model through reinforcement learning (Proximal Policy Optimization)

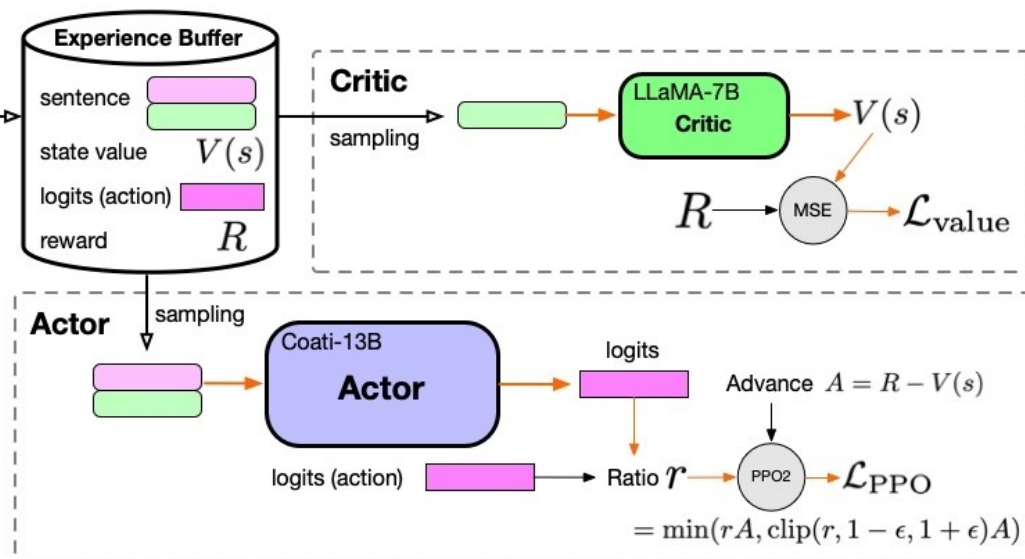
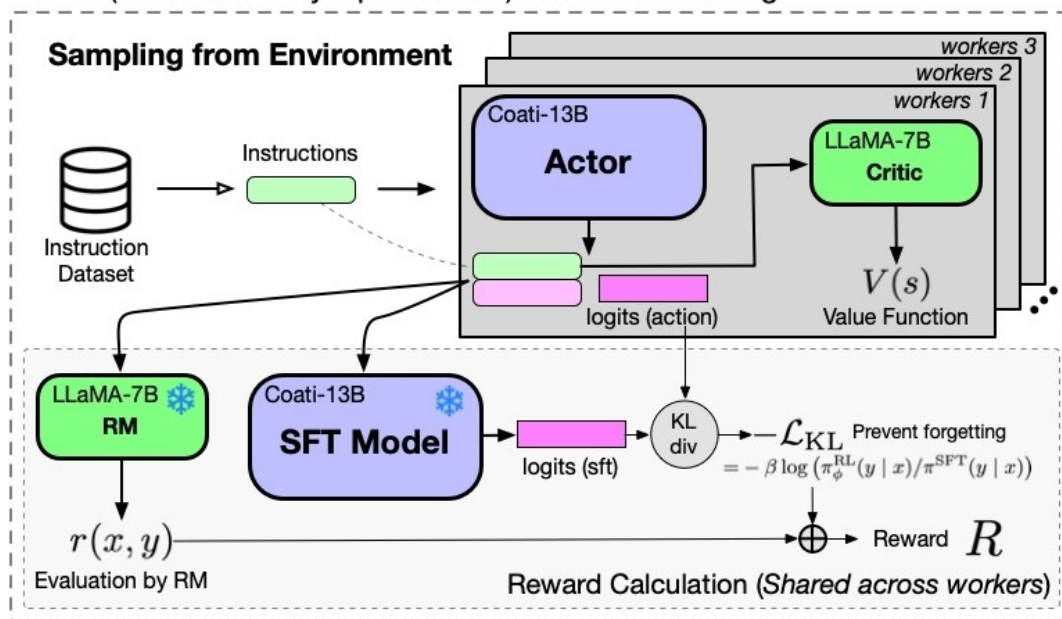


RLHF

PTX (Pretraining Gradient Mixing): Prevent forgetting



PPO (Proximal Policy Optimization): Reinforce learning



4 models: **Actor**, **Critic**, **RM**, **SFT**

$$\mathcal{L} = \mathcal{L}_{PPO} + \mathcal{L}_{value} + \mathcal{L}_{ptx}$$

SFT vs RLHF

- It is harder to label data for SFT vs RLHF

SFT vs RLHF

- It is harder to label data for SFT vs RLHF
- It is easier to crawl data for SFT vs RLHF

SFT vs RLHF

- It is harder to label data for SFT vs RLHF
- It is easier to crawl data for SFT vs RLHF
- RLHF data is collected online while SFT data is collected offline

SFT vs RLHF

- It is harder to label data for SFT vs RLHF
- It is easier to crawl data for SFT vs RLHF
- RLHF data is collected online while SFT data is collected offline
- RLHF suffers less from “catastrophic forgetting”

SFT vs RLHF

- It is harder to label data for SFT vs RLHF
- It is easier to crawl data for SFT vs RLHF
- RLHF data is collected online while SFT data is collected offline
- RLHF suffers less from “catastrophic forgetting”
- RLHF setup is a lot more complex compared to SFT

Important Questions

- How much data is needed for SFT or RLHF on top of a foundational model?
- How much improvement one can get by doing domain/problem specialization?
- What are the scaling laws for fine-tuning?
- How effective are methods such as PEFT in fine-tuning?

Thank you.
Questions?