



Lecture 19:

LLMs: Post-training

Instructor: Swabha Swayamdipta
USC CSCI 544 Applied NLP
Oct 31, Fall 2024



Announcements

- Today, Thu, 10/31 - Lecture + Paper Presentation I
- Tue, 11/5 - Lecture + Paper Presentation II
- Thu, 11/7 - Quiz 4 + Paper Presentation III
- Tue, 11/12 - Quiz 5 + Paper Presentation IV
- Thu, 11/14 Guest lecture on LLM Pretraining by Prof. Willie Neiswanger on 11/14 +
HW4 due
 - Questions from lecture materials will be included in final exam
- Quizzes 4 and 5 - all topics after the midterms
 - Consider these as practice tests for final exams

Lecture Outline

- Announcements
- Last Lecture: LLM Generative Evaluation + Pre-training
- Today:
 - Post-training with Supervised Finetuning
 - Instruction Tuning
 - Interacting with LLMs: Prompting
 - Post-training with Alignment with Human Feedback:
 - Preference Tuning: RLHF

The need for post-training

A Pre-trained GPT-3

Prompt: Explain the moon landing to a six year old in a few sentences.

Output: Explain the theory of gravity to a 6 year old.

Prompt: Translate to French: The small dog

Output: The small dog crossed the road.

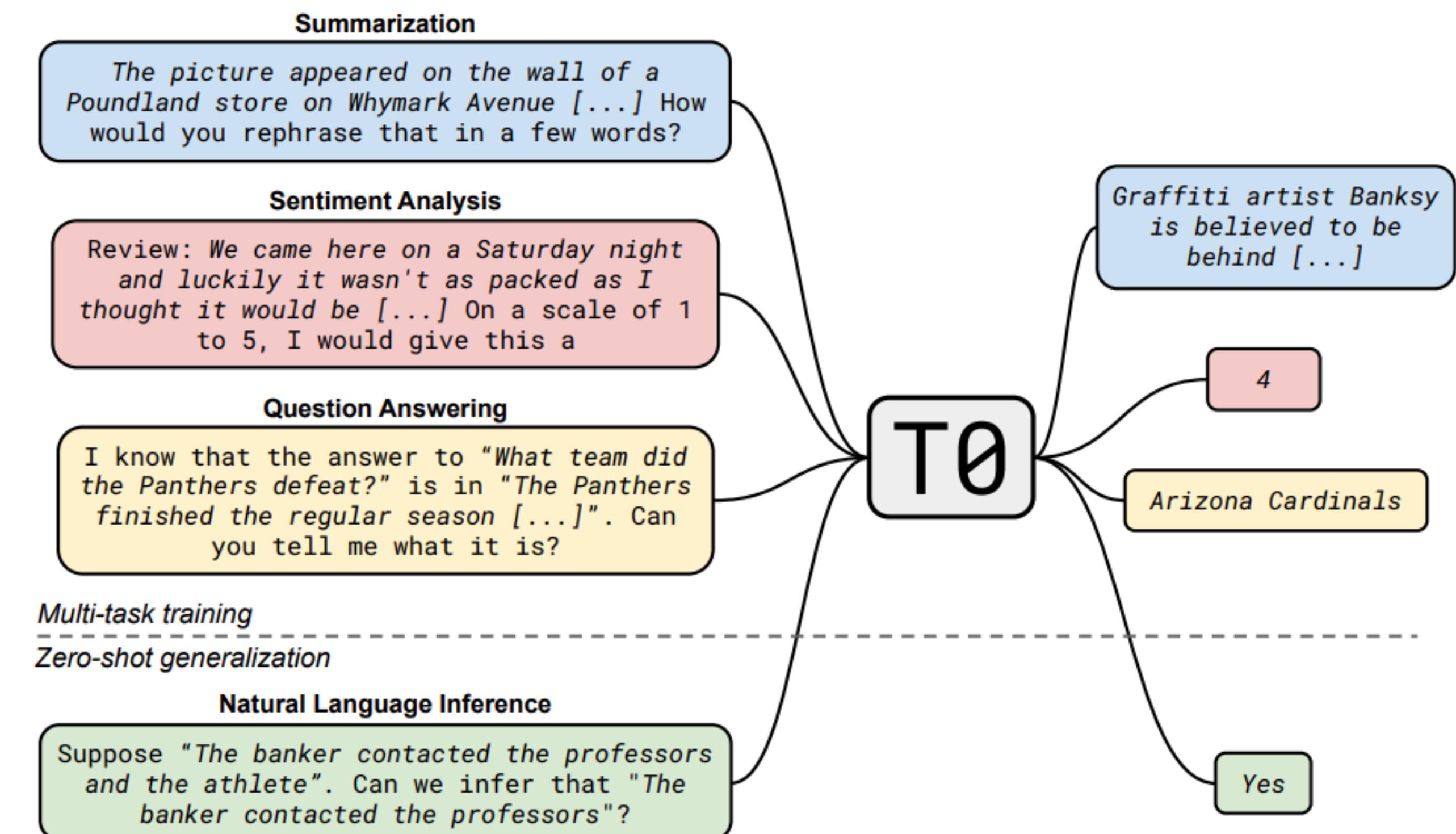
Ouyang et al., 2022; J&M Chap 12

- Make LLMs more helpful
 - Supervised Finetuning: Instruction Tuning
 - Prompting
- Make LLMs less harmful
 - Model Alignment with Human Preferences: Intro to RLHF / DPO

Supervised Fine-tuning LLMs: Instruction-Tuning

Instruction Tuning

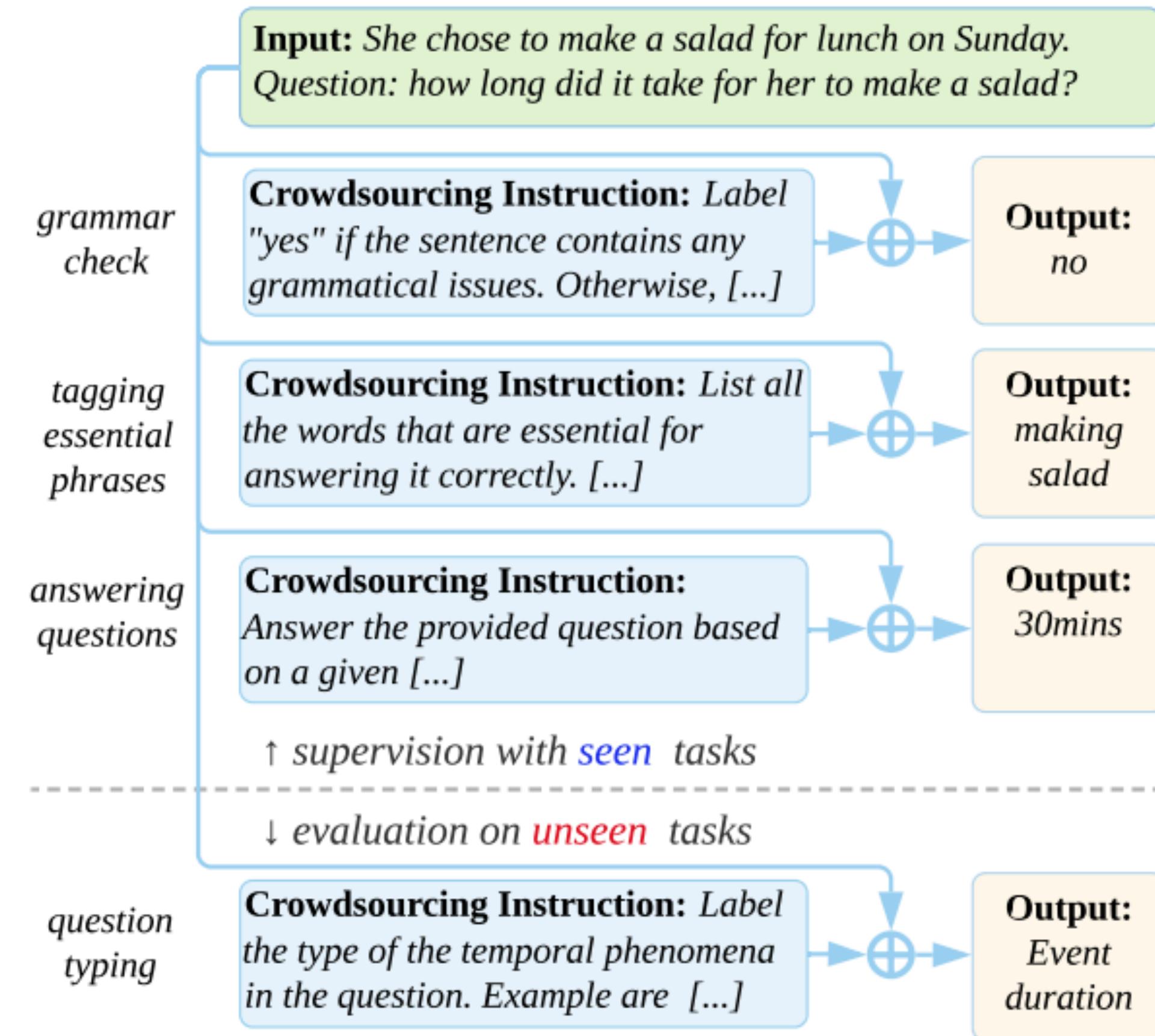
- Pretraining:
 - Train a model to continue a given context
- Instruction Tuning:
 - Train a model to follow varied instructions
 - Needed because the vast majority of pretraining is done on data which are not in the form of instructions
 - Fine-tuned (using the next-token-prediction objective) on a dataset of instructions together with correct responses



"Multitask Prompted Training Enables Zero-Shot Task Generalization" (Sahn et al., 2022)

Instruction Tuning and Task Generalization

- During training (supervised fine-tuning), the model learns to follow instructions of given tasks
- At test time, it generalizes to follow instructions on unseen tasks!



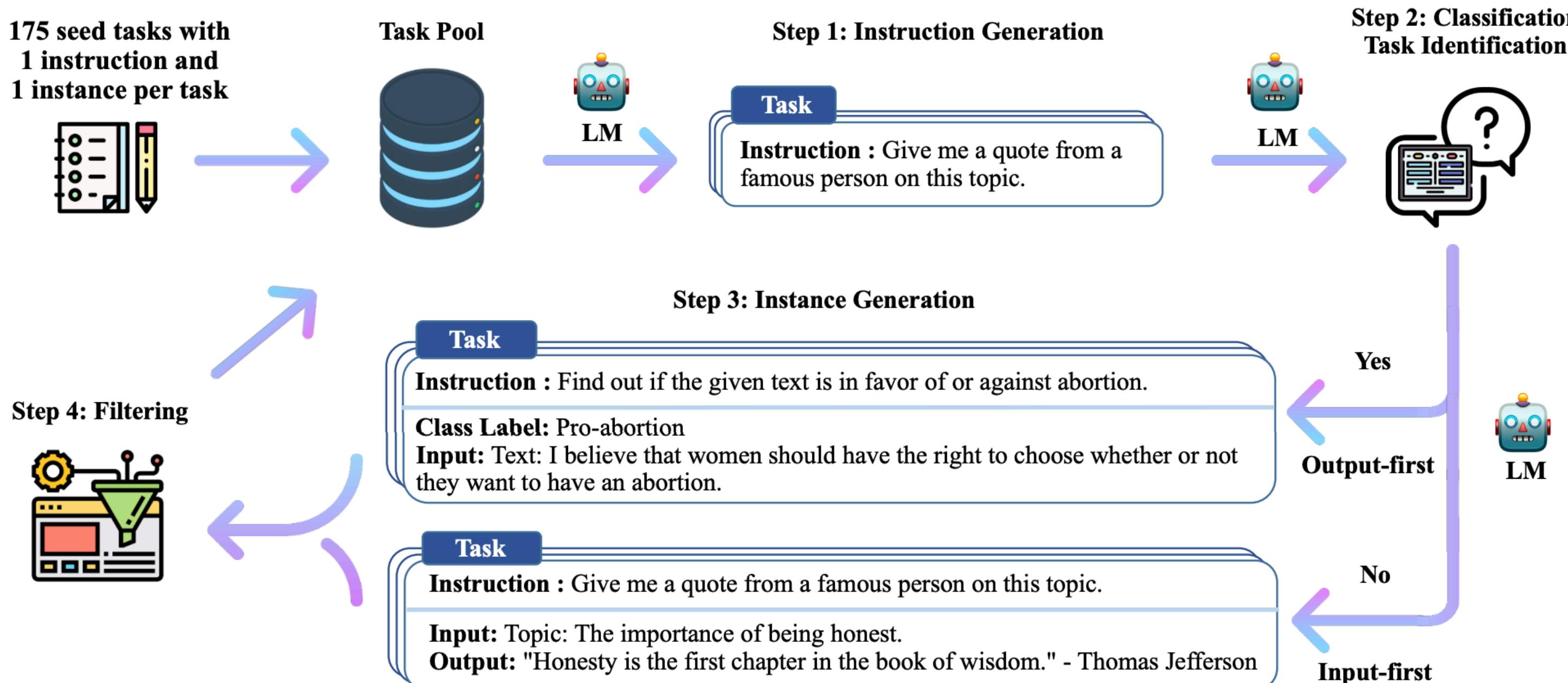
"Cross-Task Generalization via Natural Language Crowdsourcing Instructions" (Mishra et al., 2022)

Instruction Tuning Data

More data (instructions) → better model

Resource →	SUP-NATINST (this work)	NATINST (Mishra et al., 2022b)	CROSSFIT (Ye et al., 2021)	PROMPTSOURCE (Bach et al., 2022)	FLAN (Wei et al., 2022)	INSTRUCTGPT (Ouyang et al., 2022)
Has task instructions?	✓	✓	✗	✓	✓	✓
Has negative examples?	✓	✓	✗	✗	✗	✗
Has non-English tasks?	✓	✗	✗	✗	✓	✓
Is public?	✓	✓	✓	✓	✓	✗
Number of tasks	1616	61	269	176	62	—
Number of instructions	1616	61	—	2052	620	14378
Number of annotated tasks types	76	6	13	13*	12	10
Avg. task definition length (words)	56.6	134.4	—	24.8	8.2	—

"Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks" (Wang et al., 2022) <https://arxiv.org/abs/2212.10560>



Diverse data (instructions)
→ better model

"Self-Instruct: Aligning Language Models with Self-Generated Instructions" (Wang et al., 2023)

Instruction Tuning Data

- Instruction tuning datasets are often created by repurposing standard NLP datasets for tasks like question answering or machine translation
- Often synthesized!
 - Prompting existing LLMs
- More variety in the instruction templates → better models!

Release	Collection	Model	Model Details			Data Collection & Training Details				Methods
			Base	Size	Public?	Prompt Types	Tasks in Flan	# Exs	Methods	
2020 05	UnifiedQA	UnifiedQA	RoBERTa	110-340M	P	ZS	46 / 46	750k		
2021 04	CrossFit	BART-CrossFit	BART	140M	NP	FS	115 / 159	71M		
2021 04	Natural Inst v1.0	Gen. BART	BART	140M	NP	ZS / FS	61 / 61	620k	+ Detailed k-shot Prompts	
2021 09	Flan 2021	Flan-LaMDA	LaMDA	137B	NP	ZS / FS	62 / 62	4.4M	+ Template Variety	
2021 10	P3	T0, T0+, T0++	T5-LM	3-11B	P	ZS	62 / 62	12M	+ Template Variety + Input Inversion	
2021 10	MetalCL	MetalCL	GPT-2	770M	P	FS	100 / 142	3.5M	+ Input Inversion + Noisy Channel Opt	
2021 11	ExMix	ExT5	T5	220M-11B	NP	ZS	72 / 107	500k	+ With Pretraining	
2022 04	Super-Natural Inst.	Tk-Instruct	T5-LM, mT5	11-13B	P	ZS / FS	1556 / 1613	5M	+ Detailed k-shot Prompts + Multilingual	
2022 10	GLM	GLM-130B	GLM	130B	P	FS	65 / 77	12M	+ With Pretraining + Bilingual (en, zh-cn)	
2022 11	xP3	BLOOMz, mT0	BLOOM, mT5	13-176B	P	ZS	53 / 71	81M	+ Massively Multilingual	
2022 12	Unnatural Inst. [†]	T5-LM-Unnat. Inst.	T5-LM	11B	NP	ZS	~20 / 117	64k	+ Synthetic Data	
2022 12	Self-Instruct [†]	GPT-3 Self Inst.	GPT-3	175B	NP	ZS	Unknown	82k	+ Synthetic Data + Knowledge Distillation	
2022 12	OPT-IML Bench [†]	OPT-IML	OPT	30-175B	P	ZS + FS CoT	~2067 / 2207	18M	+ Template Variety + Input Inversion + Multilingual	
2022 10	Flan 2022 (ours)	Flan-T5, Flan-PaLM	T5-LM, PaLM	10M-540B	P NP	ZS + FS CoT	1836	15M	+ Template Variety + Input Inversion + Multilingual	

"The Flan Collection: Designing Data and Methods for Effective Instruction Tuning" (Longpre et al., 2023)

Instruction Tuning: Masking Instructions

- We're still using decoder-only models
- The instruction itself is masked, so the model does not generate instructions.

In the loss function, mask out the tokens corresponding to prompt

Below is an instruction that describes a task. Write a response that appropriately completes the request.

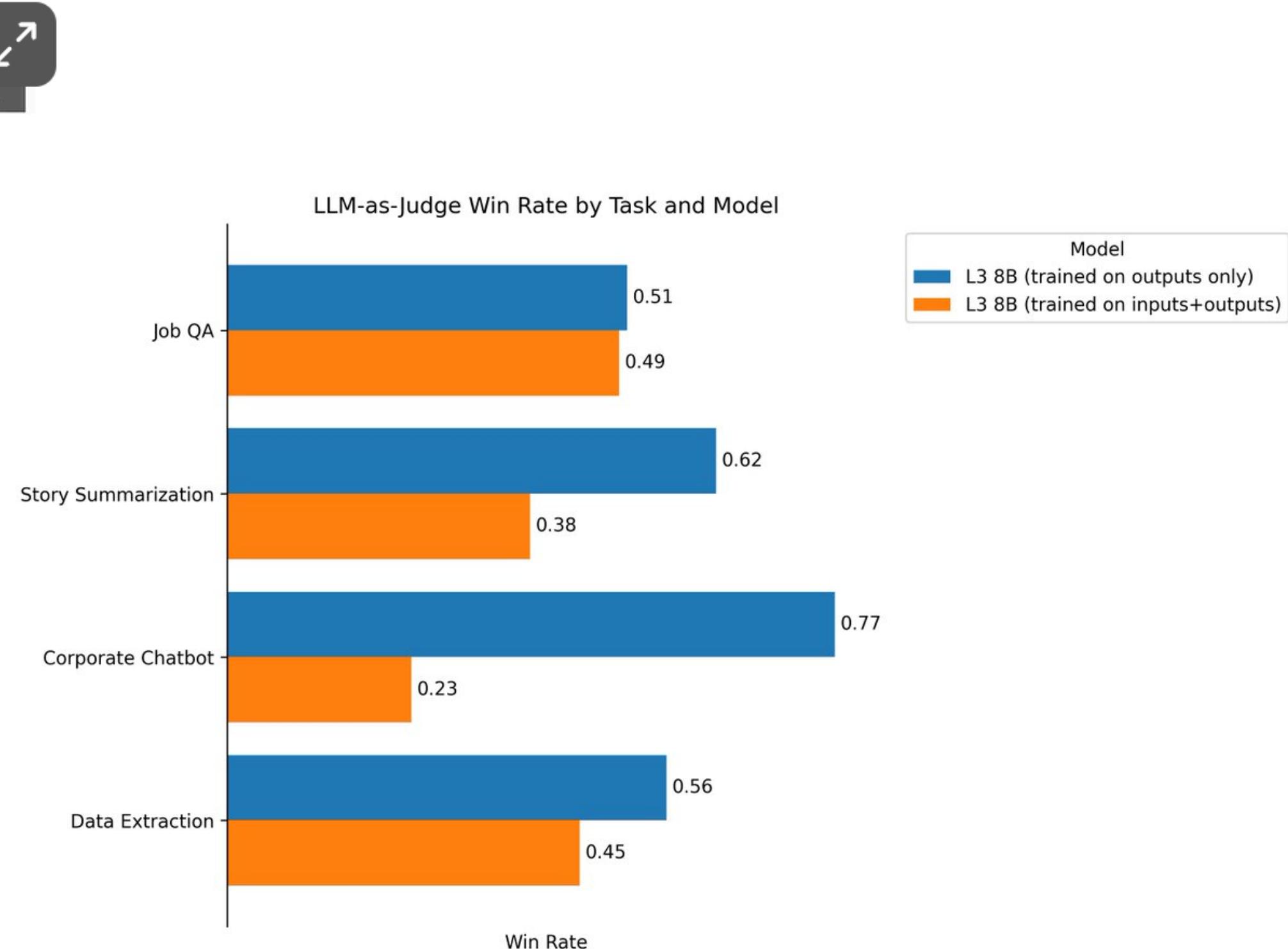
Instruction:
Rewrite the following sentence using passive voice.

Input:
The team achieved great results.

Response:
Great results were achieved by the team.

Calculate the loss only over the tokens corresponding to the response text

An illustration of input masking where the highlighted text is still fed to the LLM, but it is not used when calculating the loss during training.



How else to make language models do our tasks well?

- **Prompting (or In-Context / Few-Shot Learning):** the ability to do many tasks with no gradient updates and no / a few examples, by simply:
 - Specifying the right sequence prediction problem
 - You can get interesting zero-shot behavior if you're creative enough with how you specify your task!

Basic Prompt Templates

Summarization {input} ; tl;dr;

Translation {input} ; translate to French:

Sentiment {input}; Overall, it was

Fine-Grained-Sentiment {input}; What aspects were important in this review?

Interacting with LLMs: Prompting

Prompting vs. Instruction Tuning

Translate this to Spanish: Goodbye.

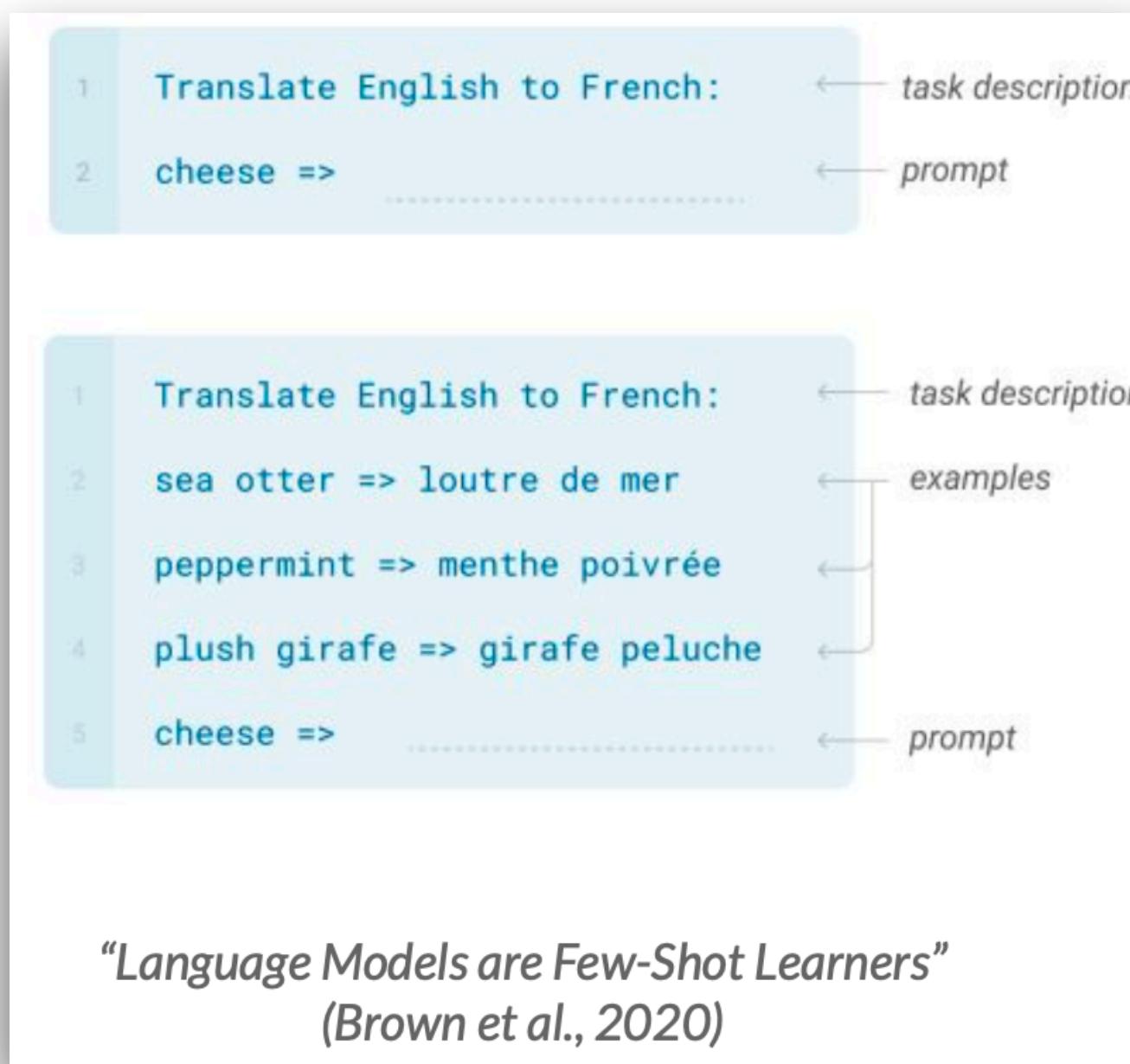
GPT

Adiós.

Input

Model

Output

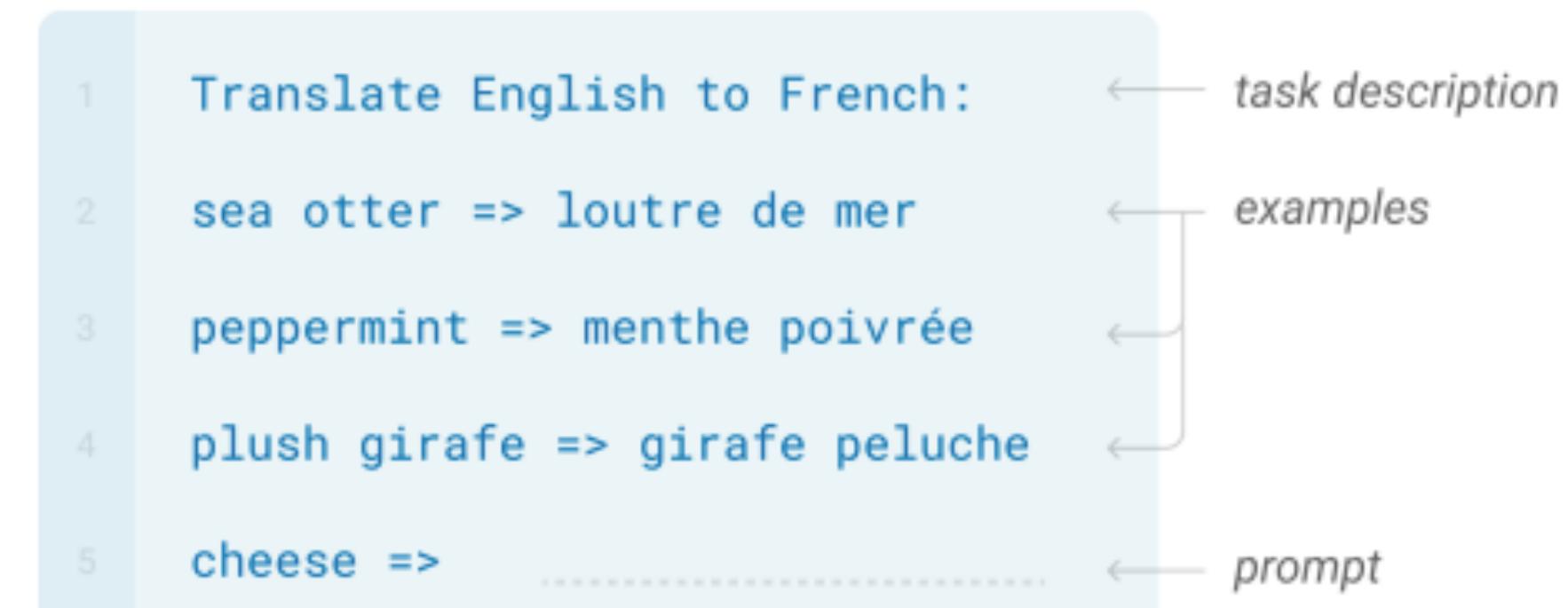


Prompting

- Interface to a language model: prompts in natural language
- Very large language models seem to perform some kind of “learning” without gradient steps simply from examples you provide within their contexts
- Sometimes called **in-context learning**
 - Misnomer: no learning (parameter update) actually happens during prompting
 - But the right examples seem to steer the language model in the right direction
- Can be zero shot (without examples) or few-shot (with a few examples)
 - Typically <10



Zero-Shot



Few-Shot

“Language Models are Few-Shot Learners” (Brown et al., 2020)

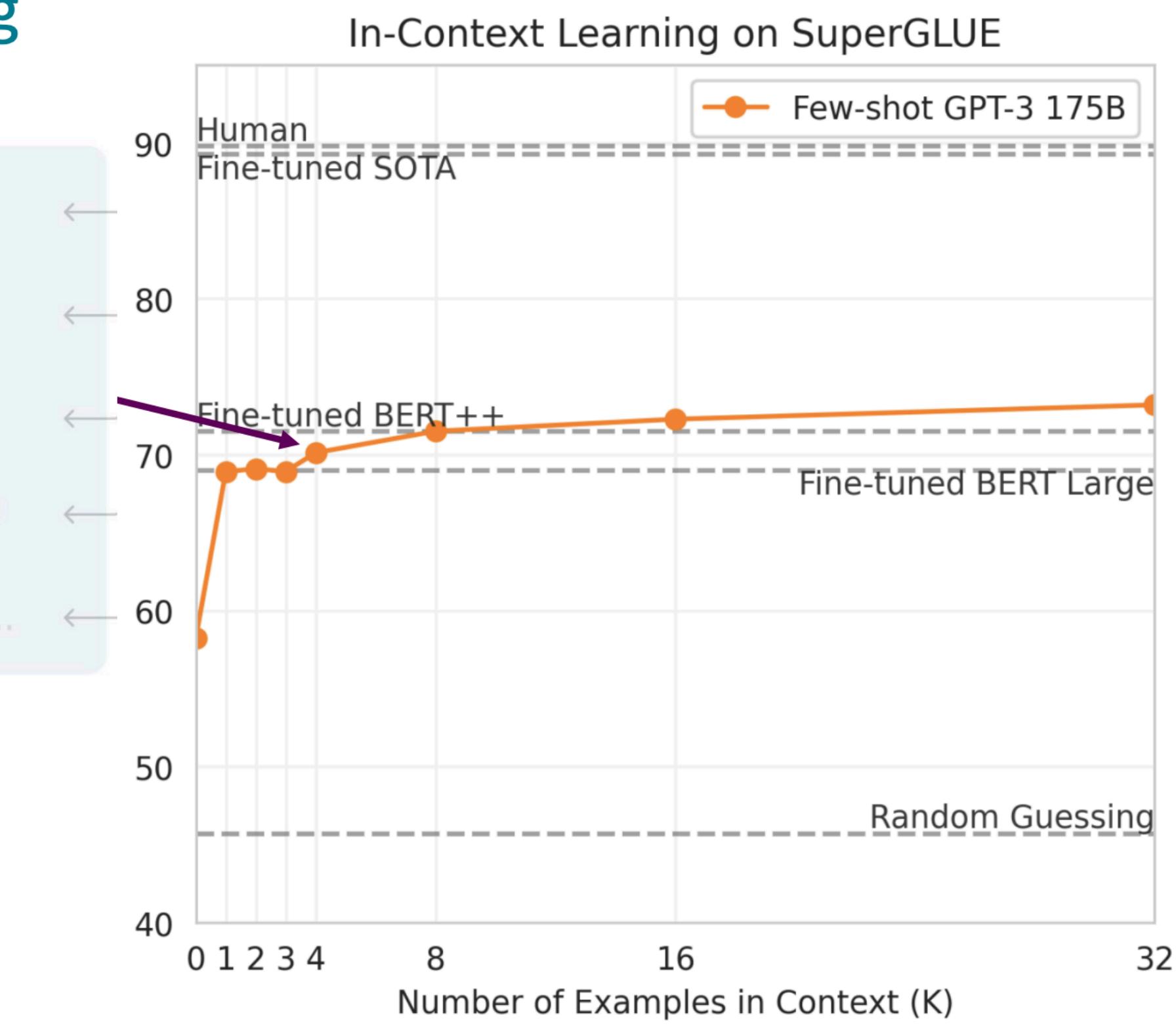
Prompting: Successes

- Much more flexible than older formulation of pre-training encoder-only models and fine-tuning to specific classification tasks (the BERT paradigm)
- Now, pre-train one large model and prompt it to do a variety of tasks!
- Much much more generalizability!

Emergent few-shot learning

Few-shot

```
1 Translate English to French:  
2 sea otter => loutre de mer  
3 peppermint => menthe poivrée  
4 plush girafe => girafe peluche  
5 cheese =>
```



[Brown et al., 2020]

Why does prompting work so well?

- Induction heads
- Discovered by looking at mini language models with only 1-2 attention heads
- If the model sees the pattern AB ... A in an input sequence, it predicts that B will follow, instantiating the pattern completion rule AB... A→ B
- Perhaps a generalized fuzzy version of this pattern completion rule, implementing a rule like A*B* ... A→ B, where A* ≈ A and B* ≈ B (by ≈ we mean some form of semantically similarity), might be responsible for in-context learning

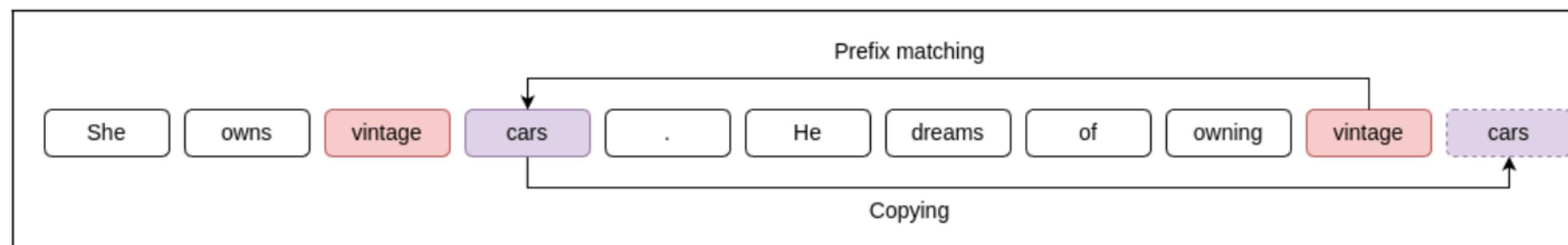
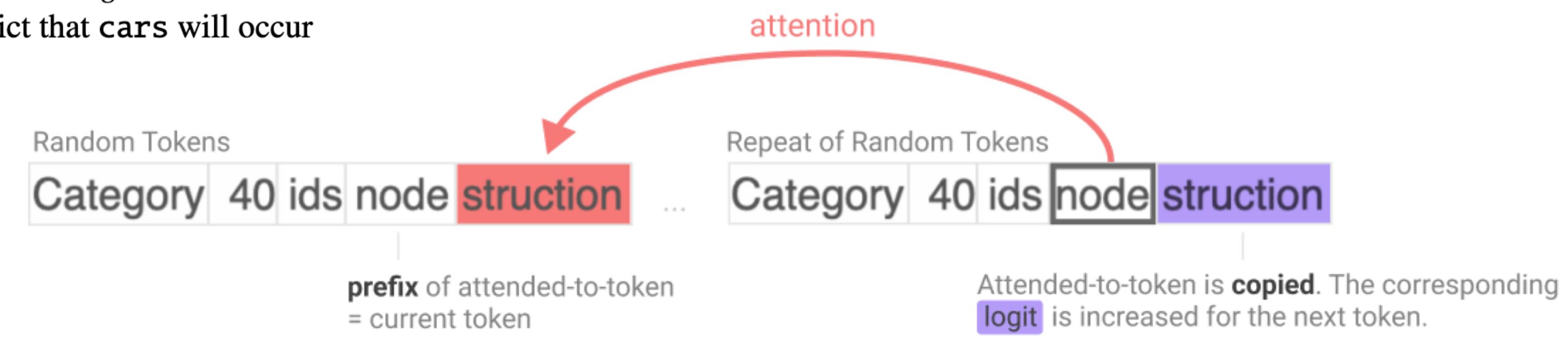


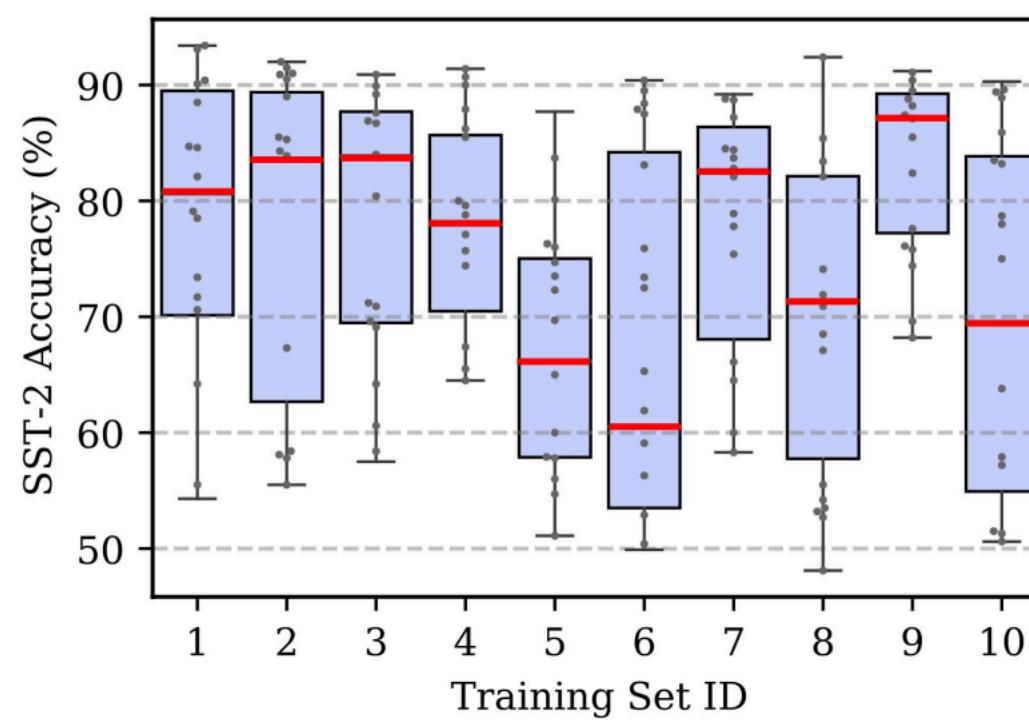
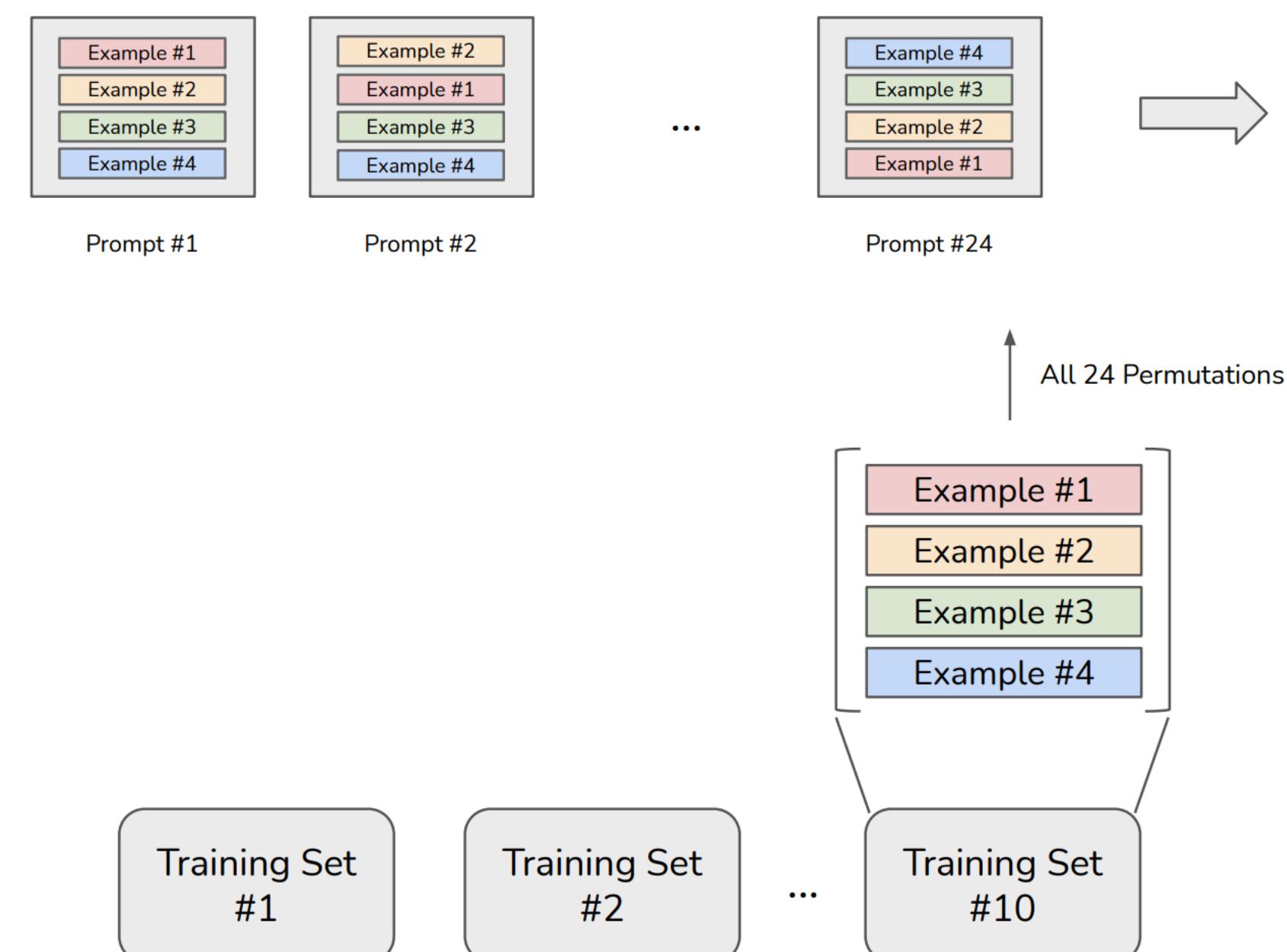
Figure 12.3 An induction head looking at `vintage` uses the *prefix matching* mechanism to find a prior instance of `vintage`, and the *copying* mechanism to predict that `cars` will occur again. Figure from [Crosbie and Shutova \(2022\)](#).



"In-context Learning and Induction Heads" (Olsson et al., 2022)

Prompting Limitations: Prompt Design

- Task performance is sensitive to prompt design
- Formatting
- Ordering of demonstrations
- Wording of the prompt



“Calibrate Before Use:
Improving Few-Shot
Performance of Language
Models” (Zhao et al., 2021)

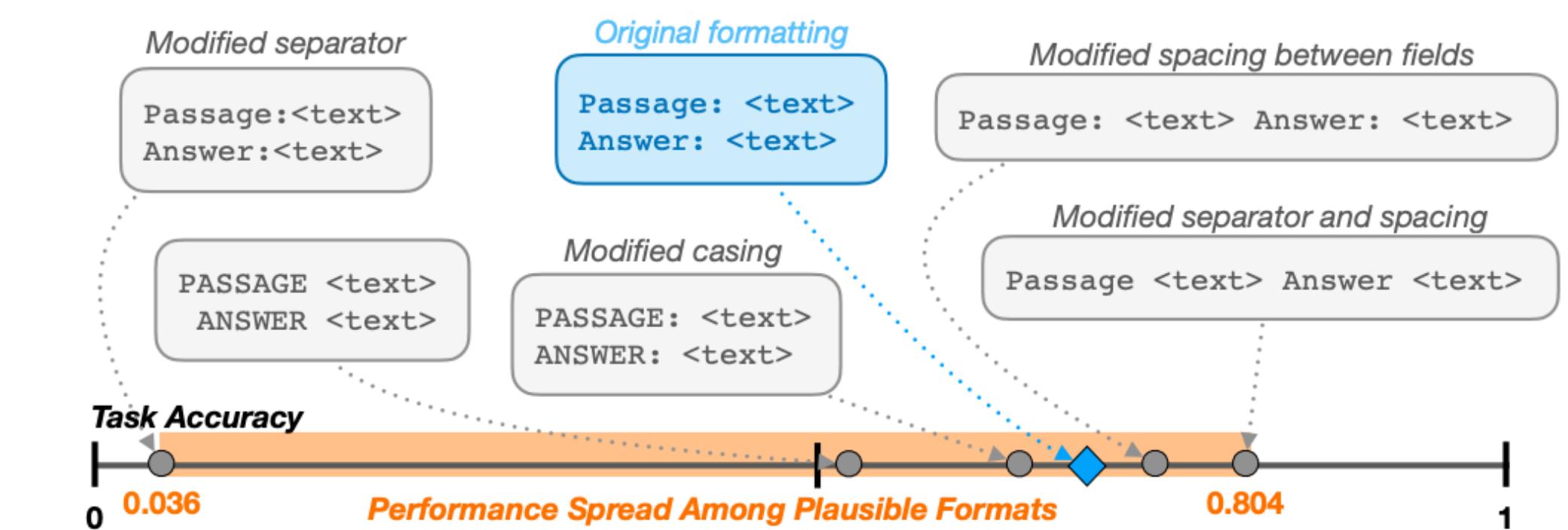
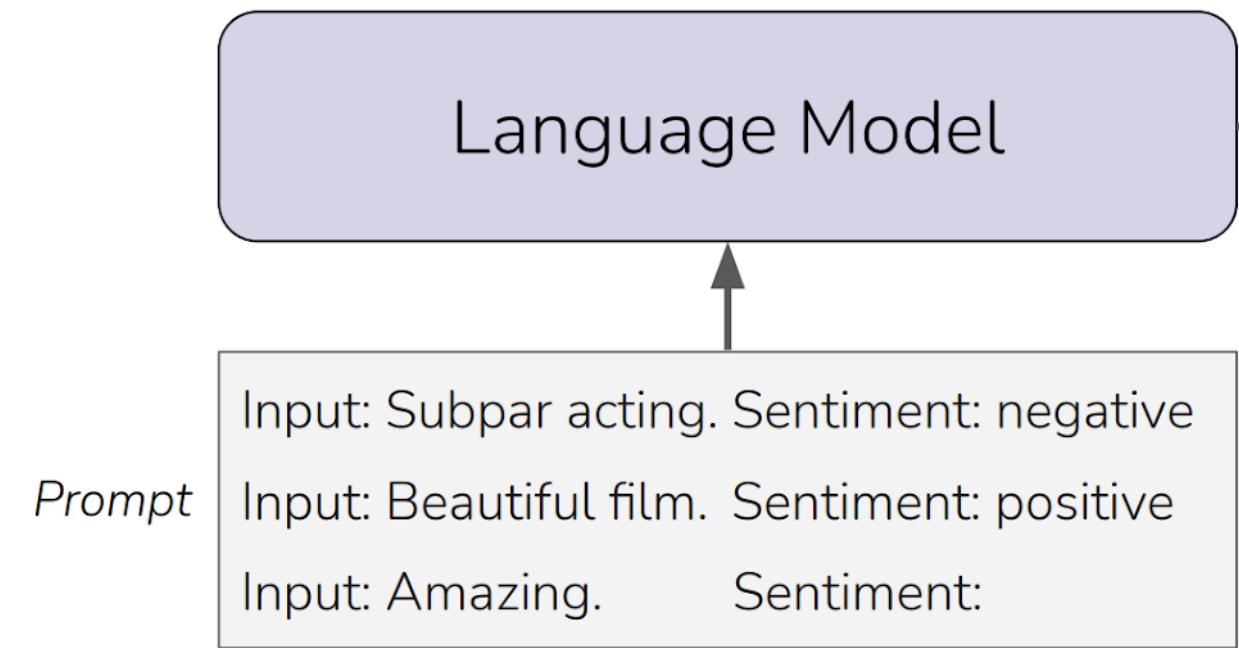


Figure 1: Slight modifications in prompt format templating may lead to significantly different model performance for a given task. Each <text> represents a different variable-length placeholder to be replaced with actual data samples. Example shown corresponds to 1-shot LLaMA-2-7B performances for task280 from SuperNaturalInstructions (Wang et al., 2022). This StereoSet-inspired task (Nadeem et al., 2021) requires the model to, given a short passage, classify it into one of four types of stereotype or anti-stereotype (gender, profession, race, and religion).

Sclar et al., ICLR 2024

Prompting Limitations?: Robustness

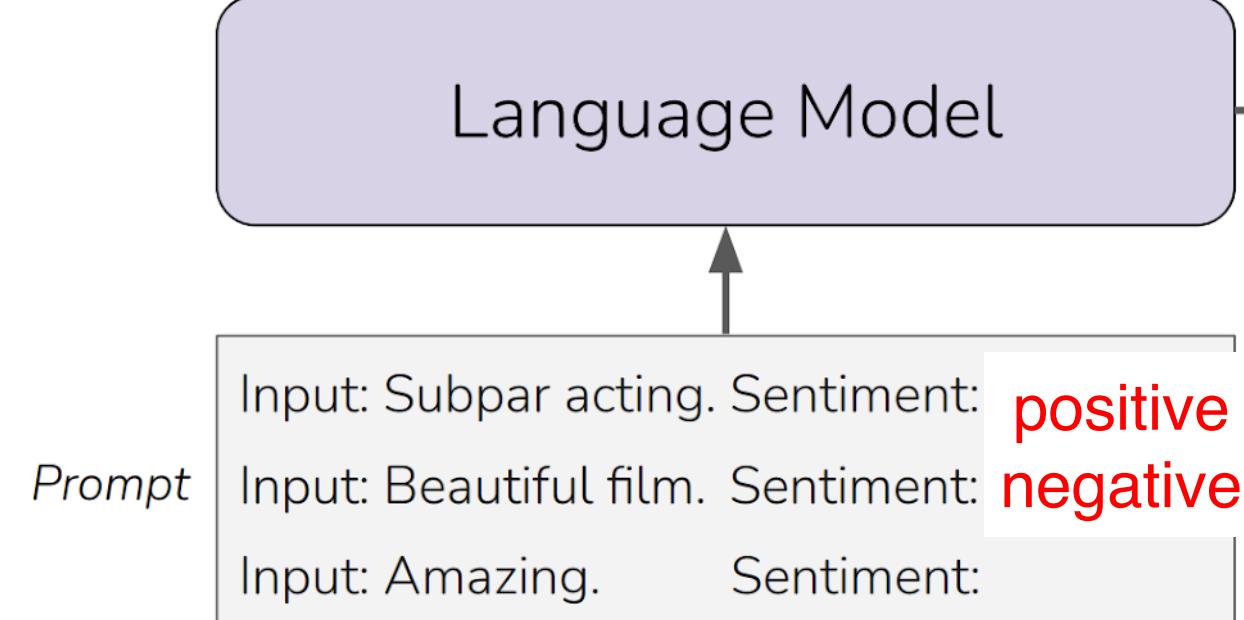
Correct labels



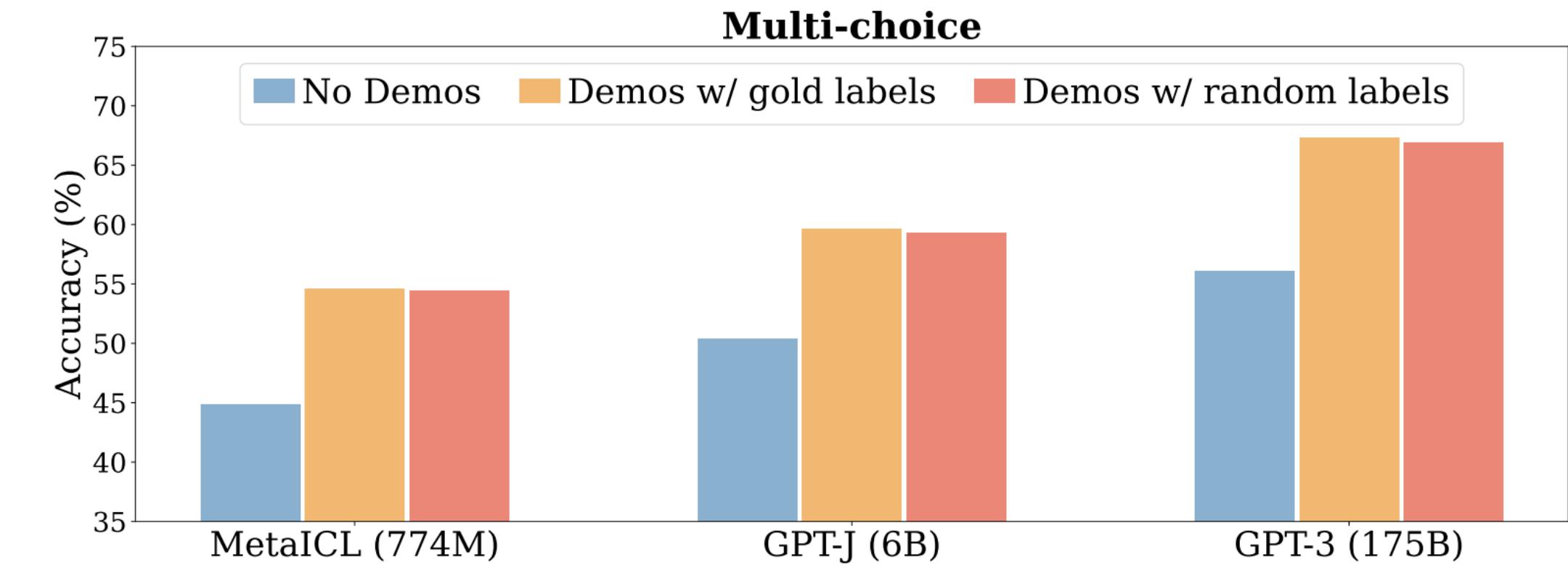
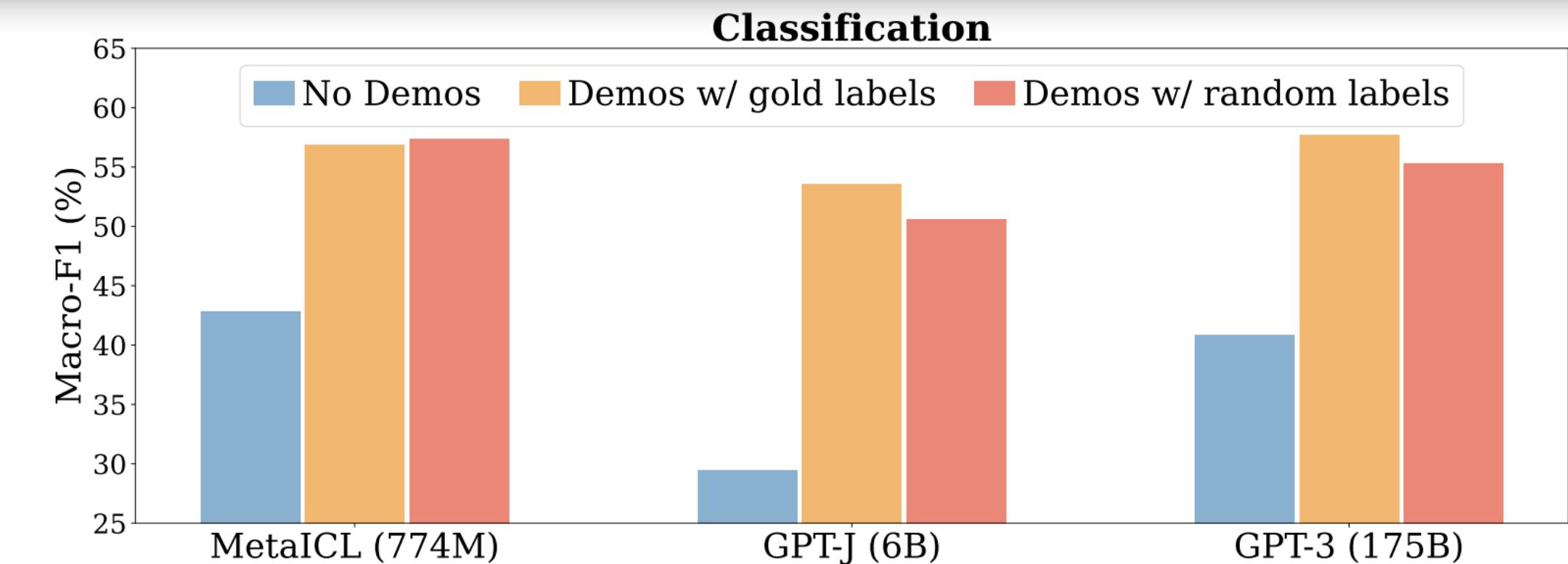
Token	Prob
positive	0.82
yes	0.10
negative	0.03
hello	0.005
acting	0.003
amazing	0.001

Demonstrations that have incorrect answers can still improve a system!

Random (perhaps wrong) labels

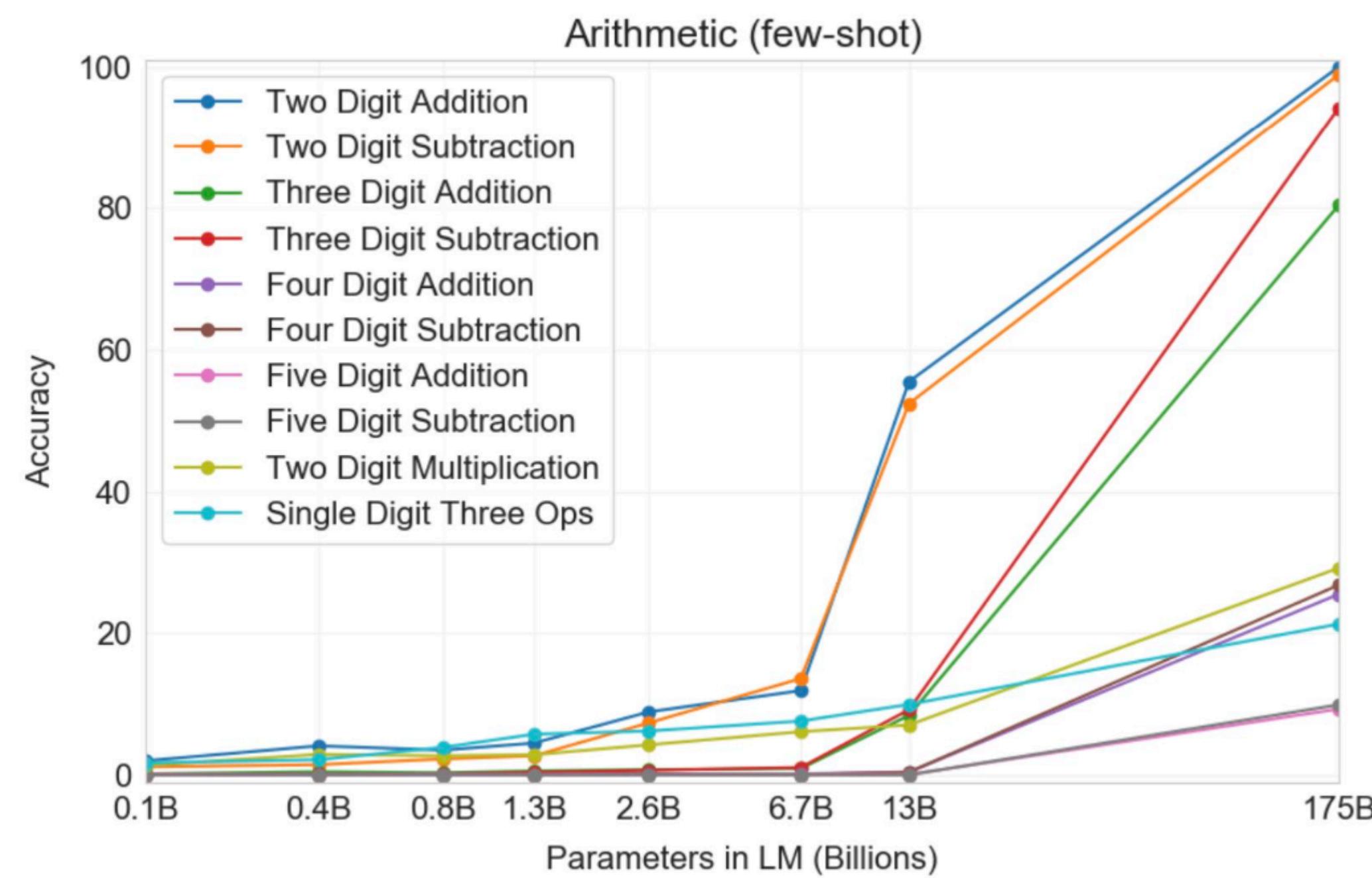


Token	Prob
positive	0.82
yes	0.10
negative	0.03
hello	0.005
acting	0.003
amazing	0.001



Prompting Limitations: Math and Reasoning

- Some tasks seem too hard for even large LMs to learn through prompting alone. Especially tasks involving richer, multi-step reasoning. (Humans struggle at these tasks too!)



Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

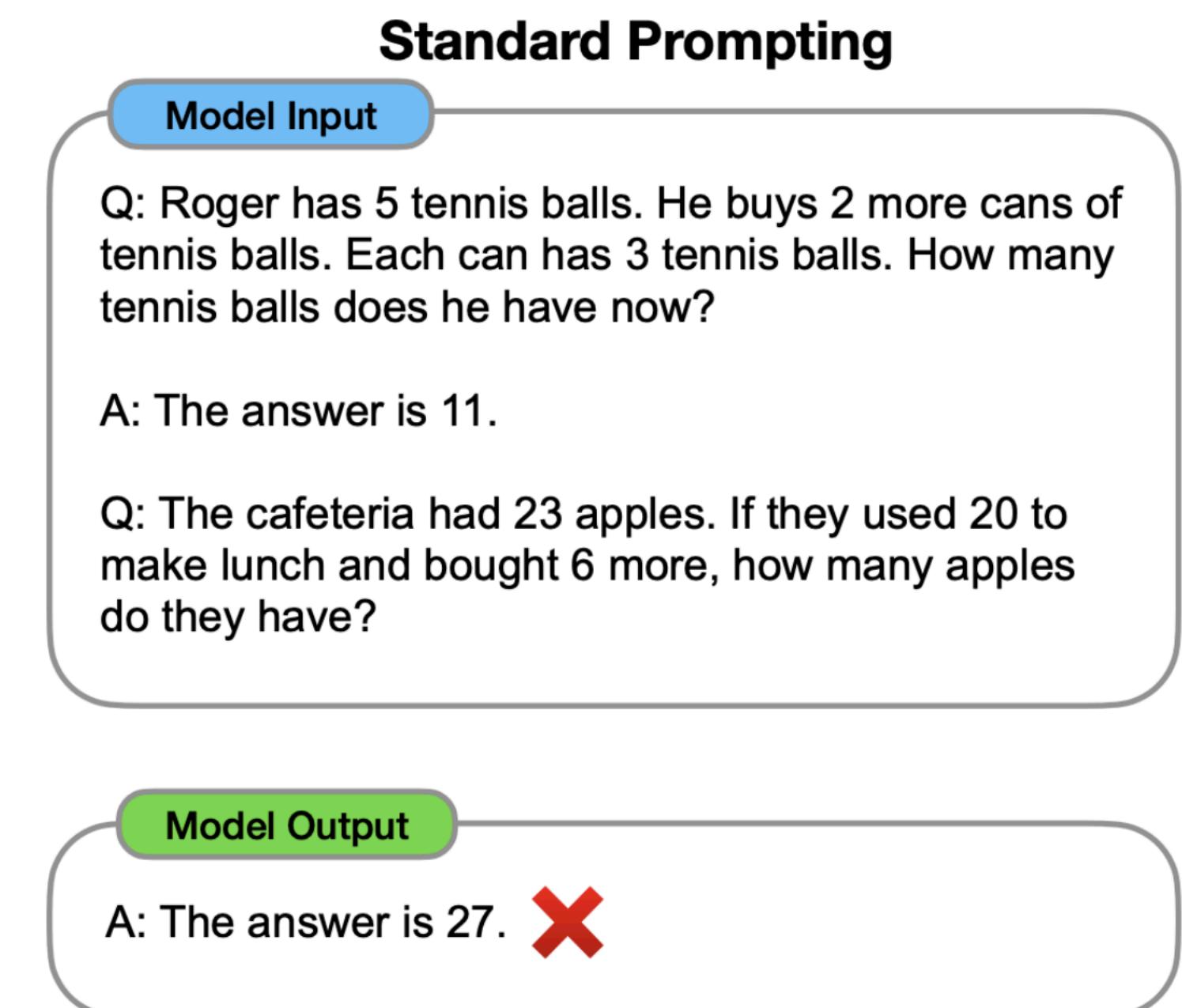
Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ✗

Chain-of-Thought Prompting

- Since the model is trained on lots and lots of language data, perhaps relying on its capabilities to generate language can make it more accurate!



Zero-Shot Chain-of-Thought Prompting

- The model may not even need examples of reasoning, it may be able to “reason” on its own if provided the right trigger context

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

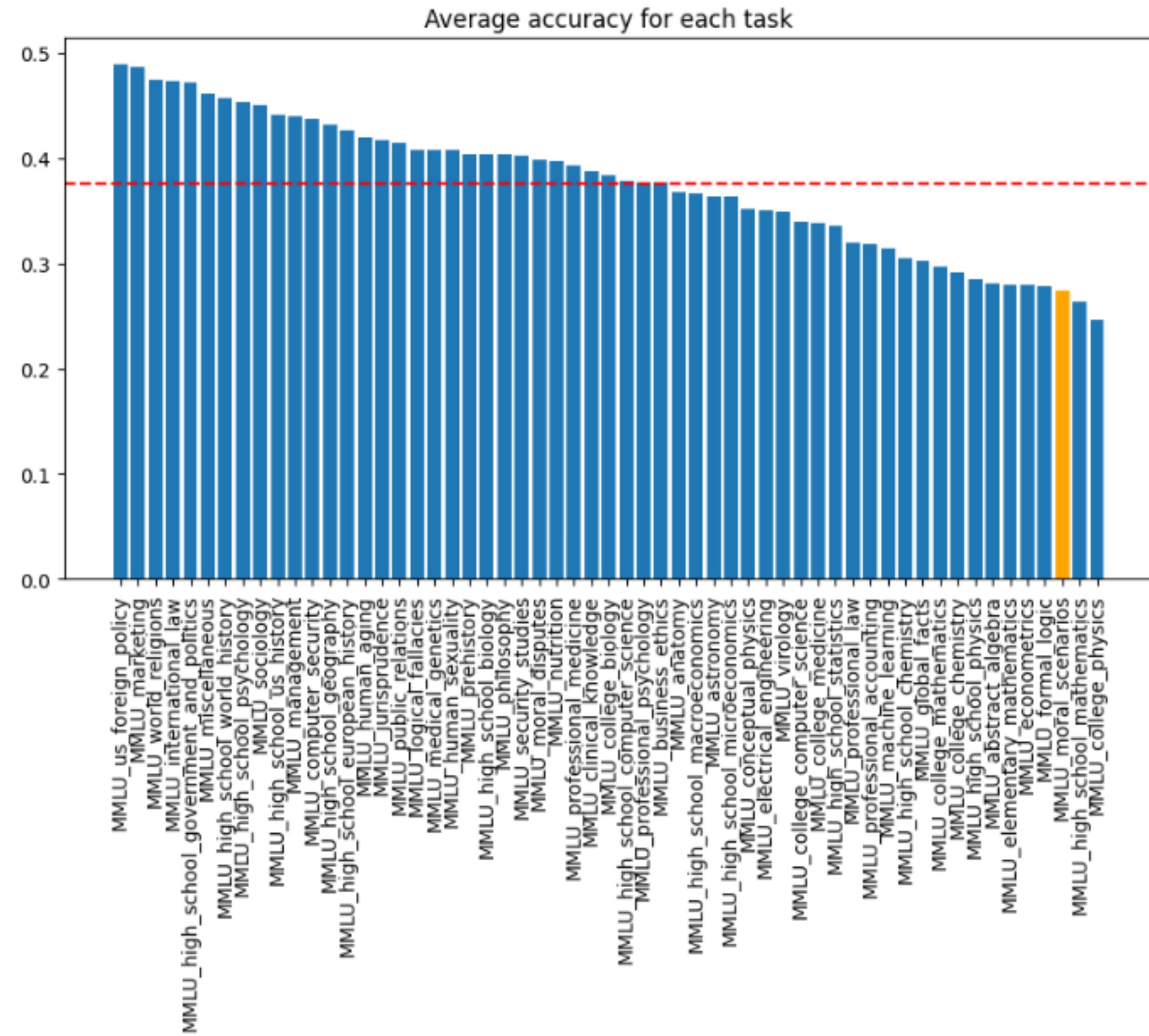
A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✓

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.** There are 16 balls in total. Half of the balls are golf balls. That means there are 8 golf balls. Half of the golf balls are blue. That means there are 4 blue golf balls. ✓

Evaluation of LLMs

- Almost exclusively on downstream tasks, as opposed to intrinsic metrics
 - Intrinsic metrics, e.g. perplexity
- Few popular multitask benchmarks
 - GLUE - Language Understanding Tasks
 - SuperGLUE - Language Understanding Tasks
 - HellaSwag - Commonsense Reasoning
 - Truthful QA - Fact Verification
 - MMLU - Massive Multitask Language Understanding, 15908 knowledge and reasoning questions in 57 areas including medicine, mathematics, computer science, law, and others
 - GSM - 8K Grade School Math
 - BigBench - subsumes some of these benchmarks



Chain-of-Thoughts Performance

	MultiArith	GSM8K
Zero-Shot	17.7	10.4
Few-Shot (2 samples)	33.7	15.6
Few-Shot (8 samples)	33.8	15.6

Zero-Shot-CoT	Greatly outperforms →	78.7	40.7
Few-Shot-CoT (2 samples)	zero-shot	84.8	41.3
Few-Shot-CoT (4 samples : First) (*1)		89.2	-
Few-Shot-CoT (4 samples : Second) (*1)		90.5	-
Few-Shot-CoT (8 samples)	Manual CoT → still better	93.0	48.7

There seems to be some wiggle room in the exact prompt to be used for achieving the best performance!

Kojima et al., 2022

Zero-shot CoT Trigger Prompt	Accuracy
Let's work this out in a step by step way to be sure we have the right answer.	82.0
Let's think step by step. (*1)	78.7
First, (*2)	77.3
Let's think about this logically.	74.5
Let's solve this problem by splitting it into steps. (*3)	72.2
Let's be realistic and think step by step.	70.8
Let's think like a detective step by step.	70.3
Let's think	57.5
Before we dive into the answer,	55.7
The answer is after the proof.	45.7
(Zero-shot)	17.7

Prompt Engineering and Auto Prompts

WIKIPEDIA
The Free Encyclopedia

Prompt engineering

5 languages

Article Talk More

From Wikipedia, the free encyclopedia

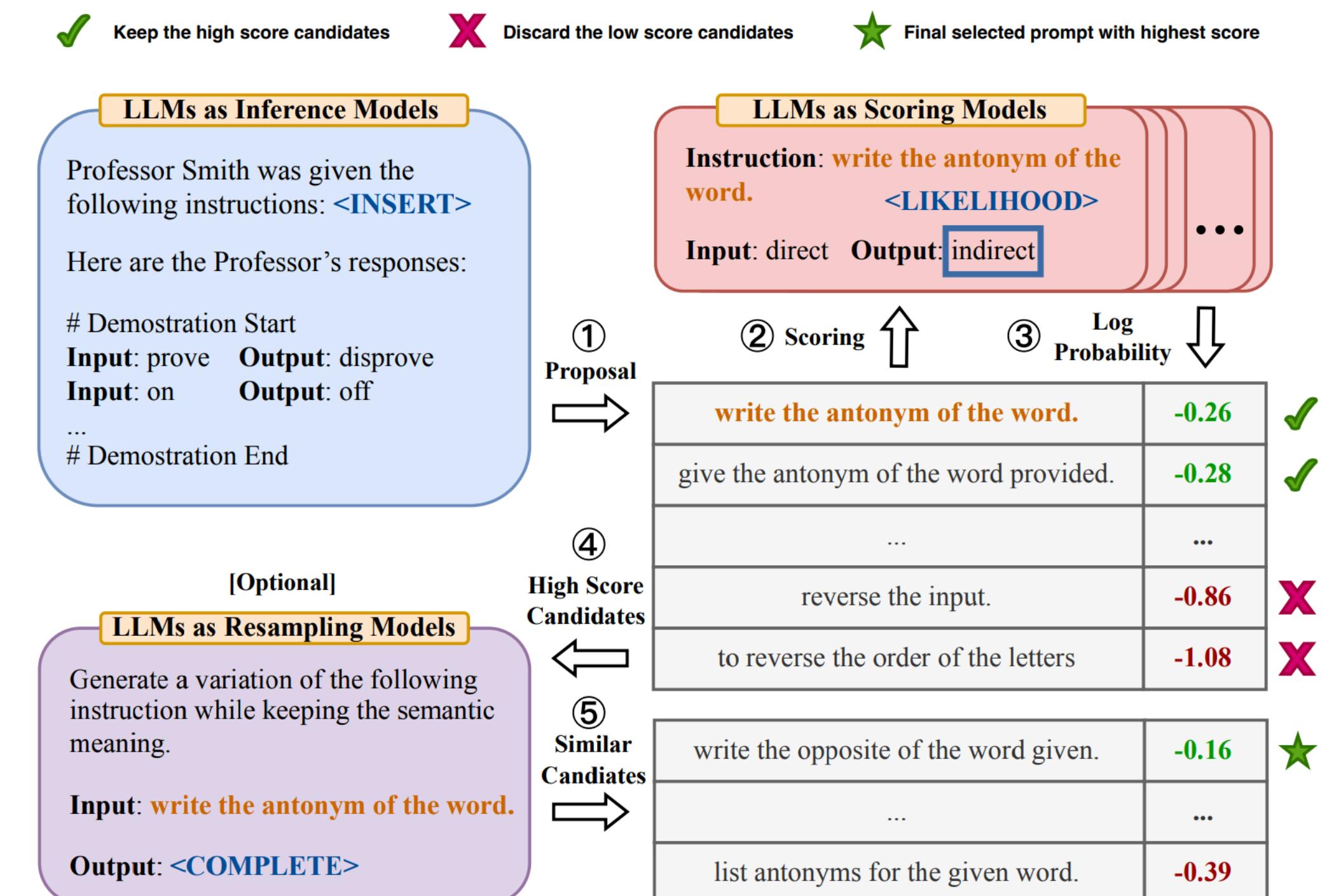
Prompt engineering is a concept in [artificial intelligence](#), particularly [natural language processing](#) (NLP). In prompt engineering, the description of the task is

Prompt Engineer and Librarian

APPLY FOR THIS JOB

SAN FRANCISCO, CA / PRODUCT / FULL-TIME / HYBRID

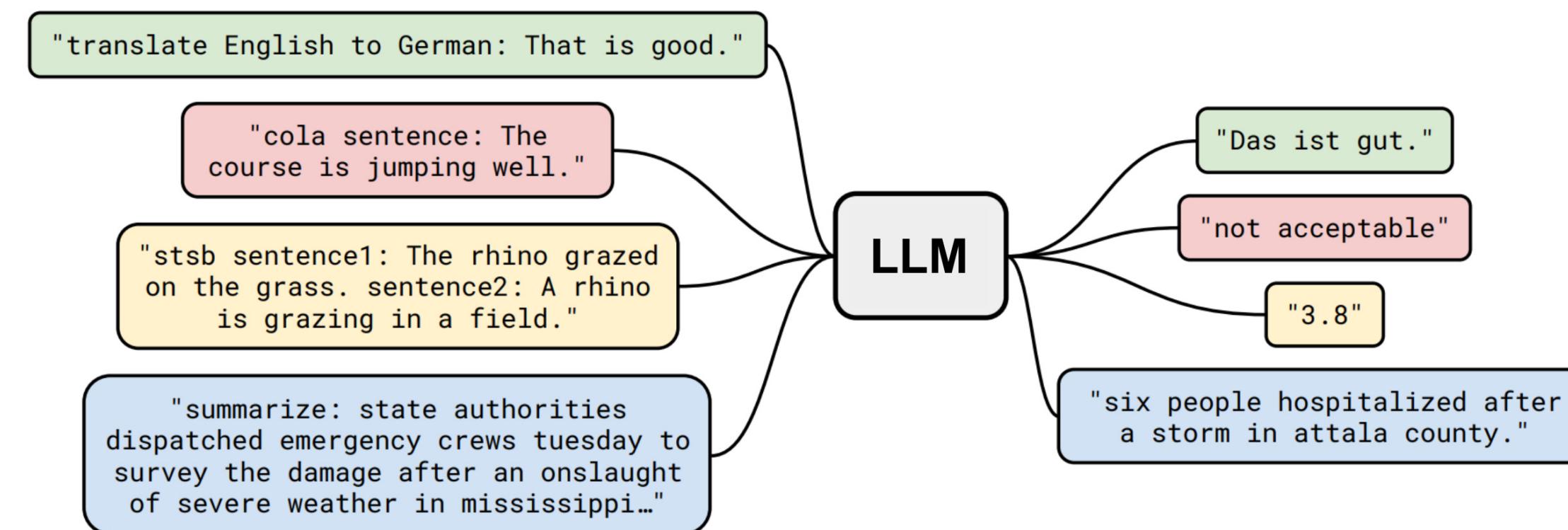
Job: keep trying new prompts for better performance, usually via tedious trial-and-error efforts



Automatic Prompt Engineer (APE). LLMs Are Human-Level
Prompt Engineers. Zhou et al., ICLR 2023

Prompting LLMs: Parting Thoughts

- Prompting is an interface into language models
- Works best with instruction-tuned language models



- **How Many Demonstrations?** small number of randomly selected labeled examples used as demonstrations can be sufficient
- **How to Select Demonstrations?** Demonstrations are generally created by formatting examples drawn from a labeled training set
 - using demonstrations that are similar to the current input seems to improve performance
 - dynamically retrieve demonstrations for each input, based on their similarity to the current example

Model Alignment with Human Preferences

Preference Alignment

- Let's say we were training a language model on some task (e.g. summarization).
- For an instruction x and a LM sample y , imagine we had a way to obtain a human reward of that summary: $R(x, y) \in \mathbb{R}$, higher is better.

SAN FRANCISCO,
California (CNN) --
A magnitude 4.2
earthquake shook the
San Francisco

...

overtake unstable
objects.

x

An earthquake hit
San Francisco.
There was minor
property damage,
but no injuries.

$$y_1 \\ R(x, y_1) = 8.0$$

The Bay Area has
good weather but is
prone to
earthquakes and
wildfires.

$$y_2 \\ R(x, y_2) = 1.2$$

- Maximize the expected reward of samples from our LM: $\mathbb{E}_{\hat{y} \sim p_\theta(y|x)}[RM_\phi(x, \hat{y})]$

Step 1

Collect demonstration data and train a supervised policy.

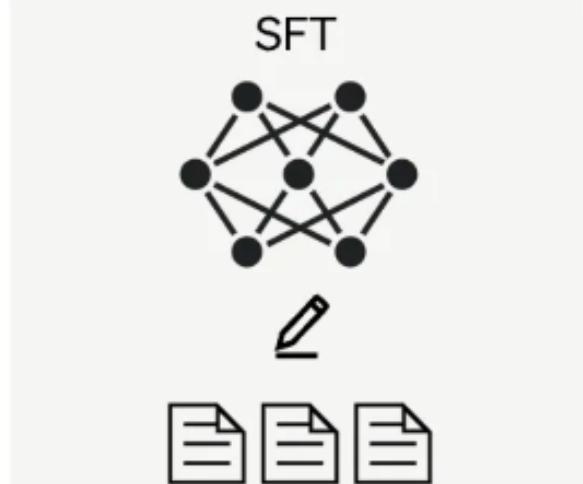
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.

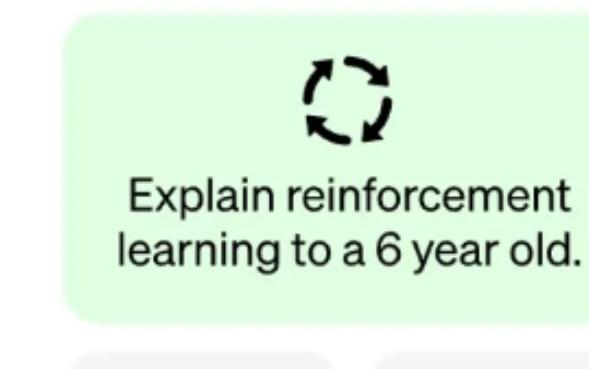


Instruction Tuning!

Step 2

Collect comparison data and train a reward model.

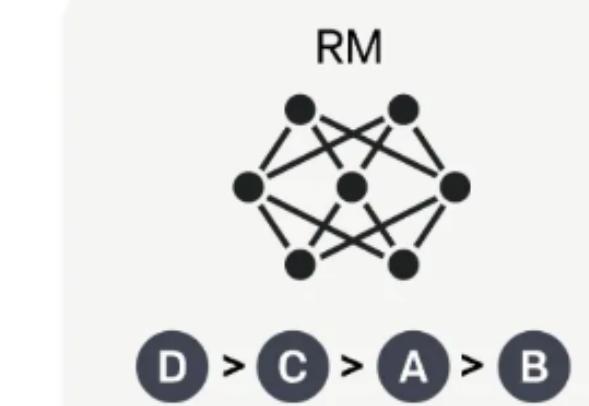
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



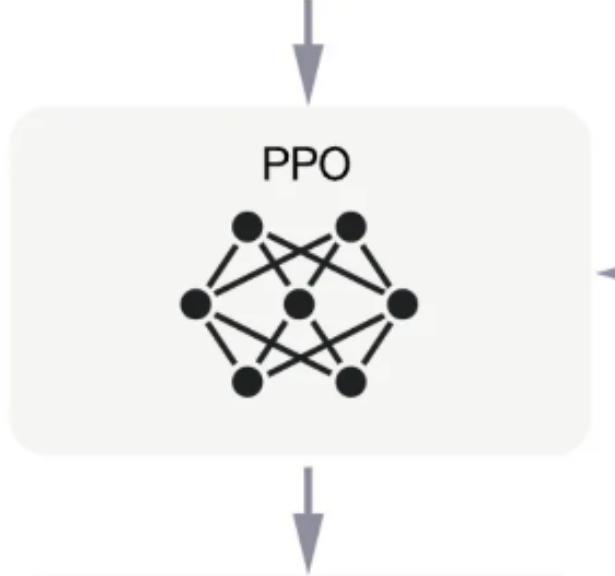
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

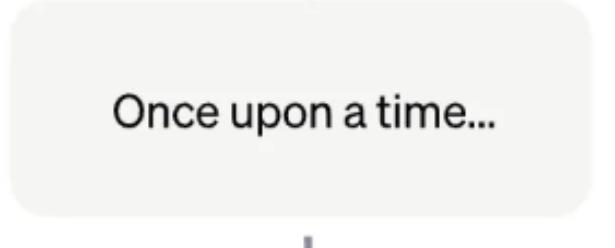
A new prompt is sampled from the dataset.



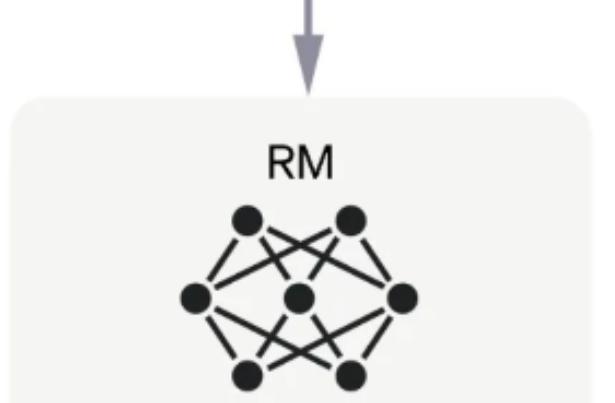
The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



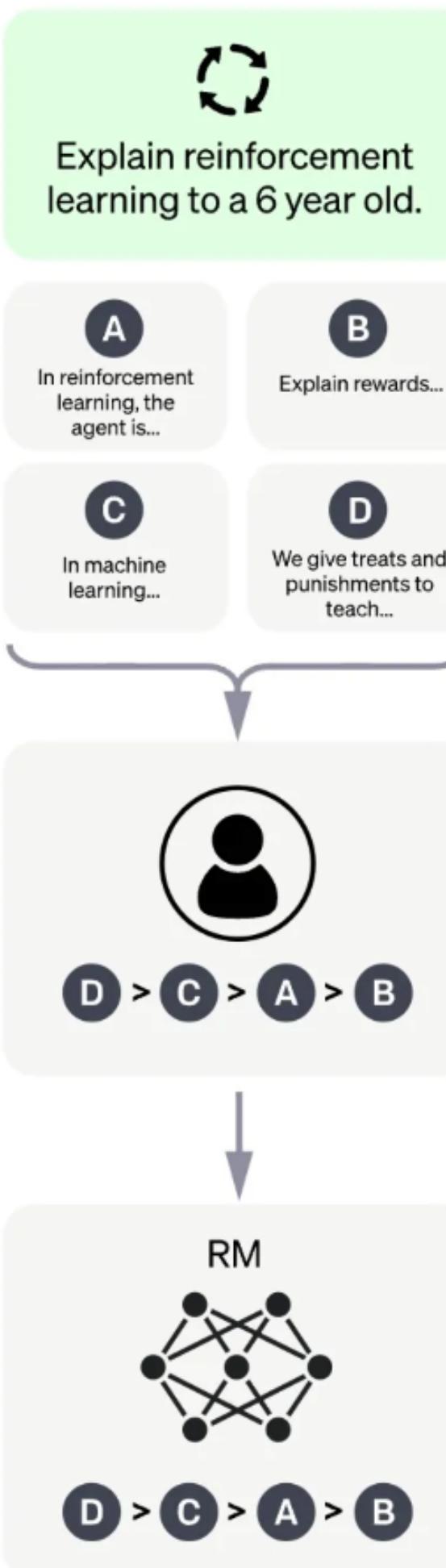
The reward is used to update the policy using PPO.



Step 2

Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.



Preference Data

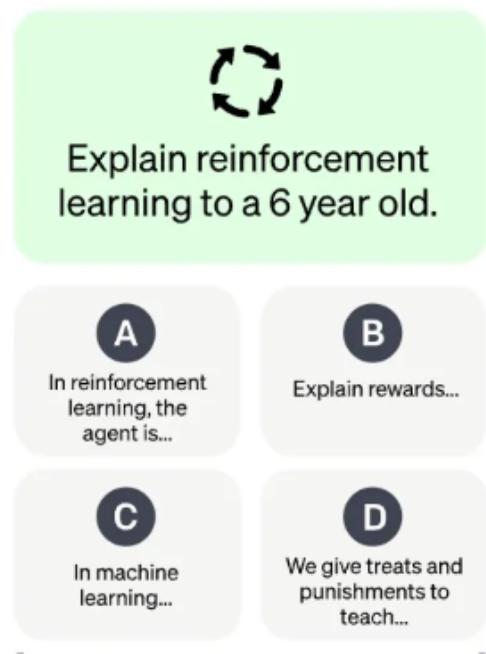
- Getting on-the-fly annotations with a human-in-the-loop is expensive!
 - Instead of directly asking humans for preferences, model their preferences as a separate (NLP) problem!
- Human judgments are noisy and miscalibrated!
 - Instead of asking for direct ratings, ask for pairwise comparisons, which can be more reliable
- Train a reward model, $RM_\phi(x, y)$ to predict human reward from an annotated dataset
 - Pairwise preferences converted into scores

Reward Modeling

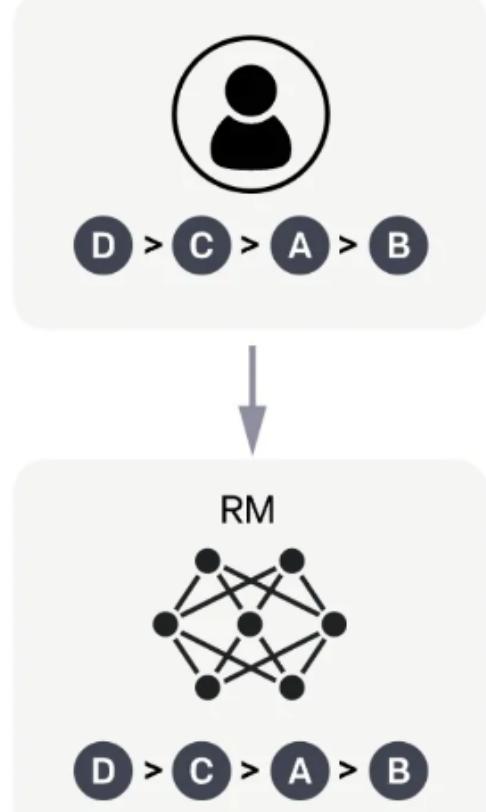
Step 2

Collect comparison data and train a reward model.

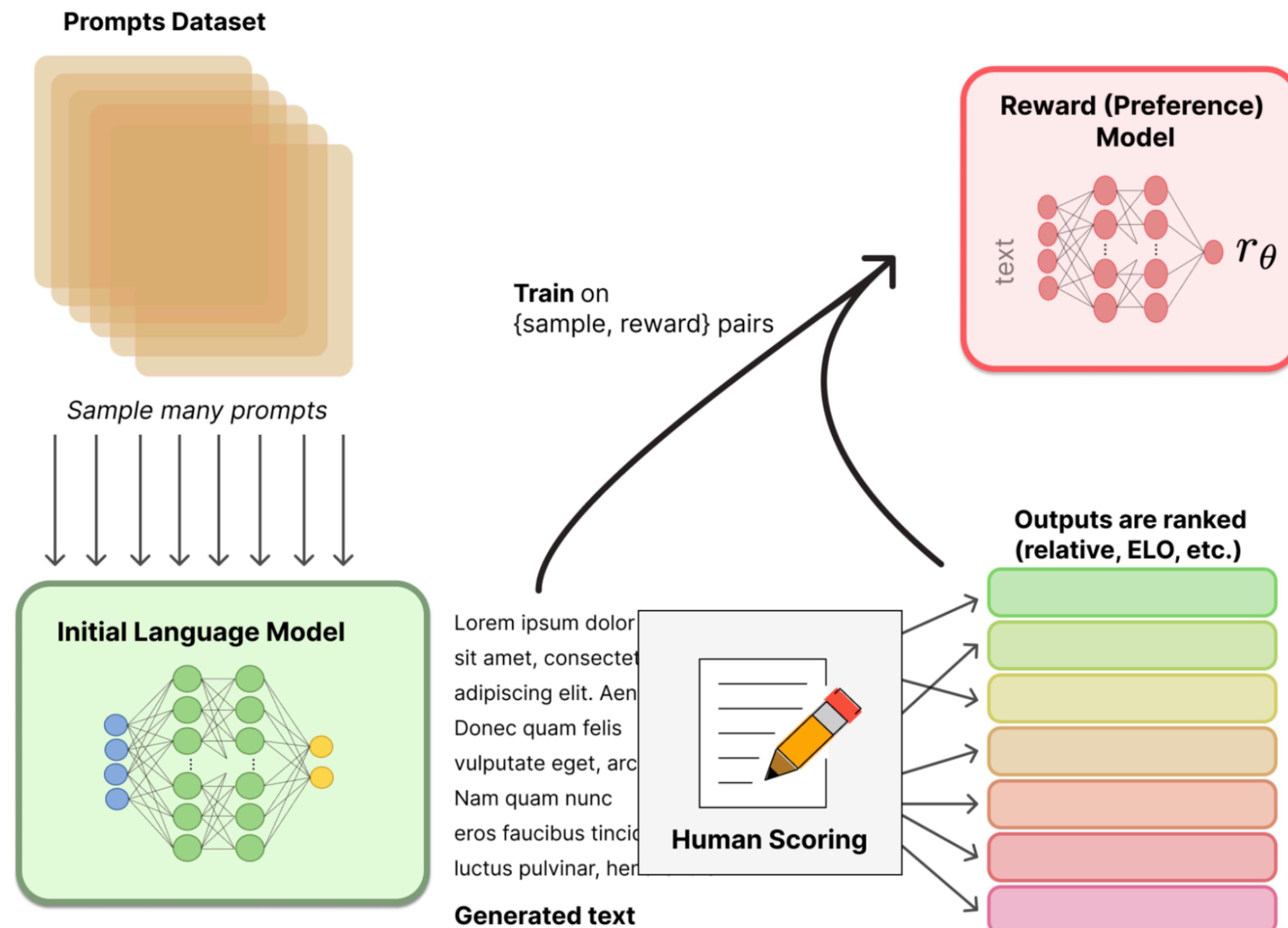
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



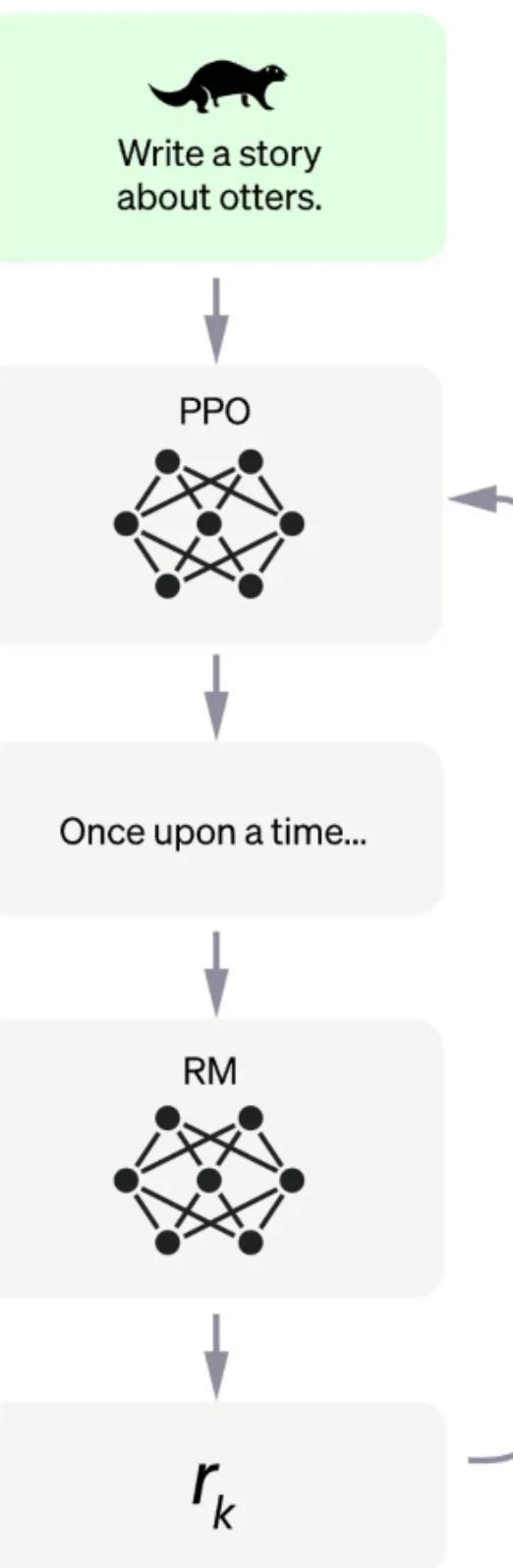
Reinforcement Learning with Human Feedback

- Ingredients
 - An instruction-tuned LM $p^{SFT}(\hat{y} | x)$
 - A reward model $RM_{\phi}(x, y)$
- Step 3 involves:
 - Copy the model to $p_{\theta}^{RL}(\hat{y} | x)$
 - Optimize: $\mathbb{E}_{\hat{y} \sim p_{\theta}^{RL}(\hat{y} | x)} [RM_{\phi}(x, y)]$
 - But, we still want a good instruction-tuned model, not just a reward maximizer
 - Hence, we add a penalty for drifting too far from the initialization: $\mathbb{E}_{\hat{y} \sim p_{\theta}^{RL}(\hat{y} | x)} \left[RM_{\phi}(x, y) - \beta \log \frac{p_{\theta}^{RL}(\hat{y} | x)}{p^{SFT}(\hat{y} | x)} \right]$
 - Use a reinforcement learning algorithm, like Proximal Policy Optimization (PPO) to maximize the above

Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.



The PPO model is initialized from the supervised policy.

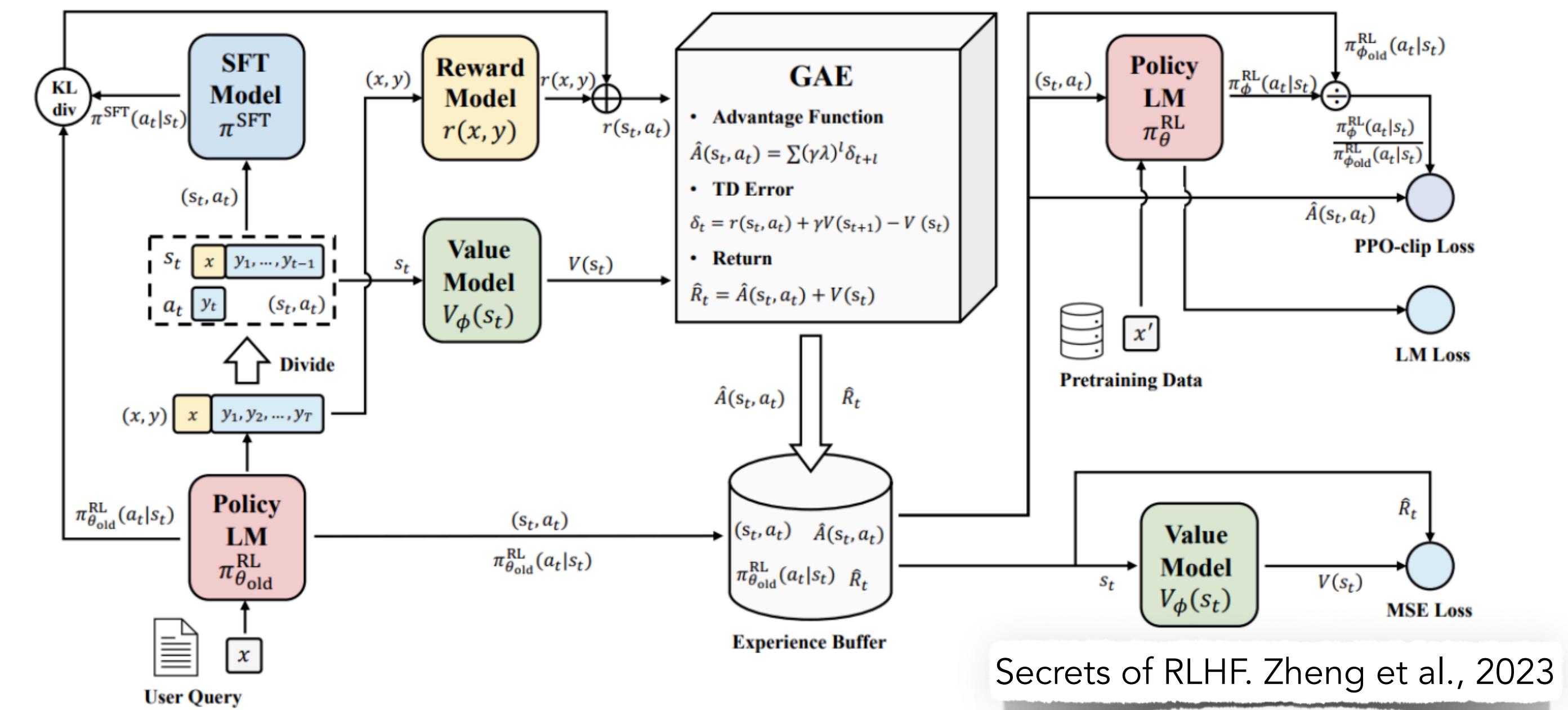
The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

RLHF to DPO

- Reinforcement Learning is tricky to train well, as well as computationally expensive
- Can we do supervised learning instead?
- Direct Preference Optimization (DPO) [Rafailov et al., 2023]!



- Clever trick: we really only need the difference between the rewards for preferred output (y_w) and dispreferred output (y_l)
- Change the reward model $RM_{\theta}(x, y)$ as a modification of the language model itself: $p_{\theta}^{RL}(\hat{y} | x)$
- Everything is now a supervised learning objective!

$$L_{DPO}(\theta) = -\mathbb{E}_{(x, y_l, y_w) \sim D} \left[\log \sigma \left(\beta \log \frac{p_{\theta}^{RL}(y_w | x)}{p^{SFT}(y_w | x)} - \beta \log \frac{p_{\theta}^{RL}(y_l | x)}{p^{SFT}(y_l | x)} \right) \right]$$

Preference Tuning: Parting Thoughts

- We want to optimize for human preferences as it's an important step towards LLM safety
 - Instead of humans writing the answers or giving uncalibrated scores, we get humans to rank different LM generated answers
- Reinforcement learning from human feedback
 - Train an explicit reward model on comparison data to predict a score for a given completion
 - Optimize the LM to maximize the predicted score without deviating too much
 - Very effective when tuned well, computationally expensive and tricky to get right
- Direct Preference Optimization
 - Optimize LM parameters directly on preference data
 - Simple and effective, similar properties and performance to RLHF
- Next Class: Safety and Harms of LLMs