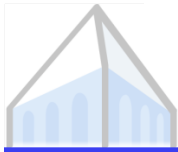


Natural Language Processing

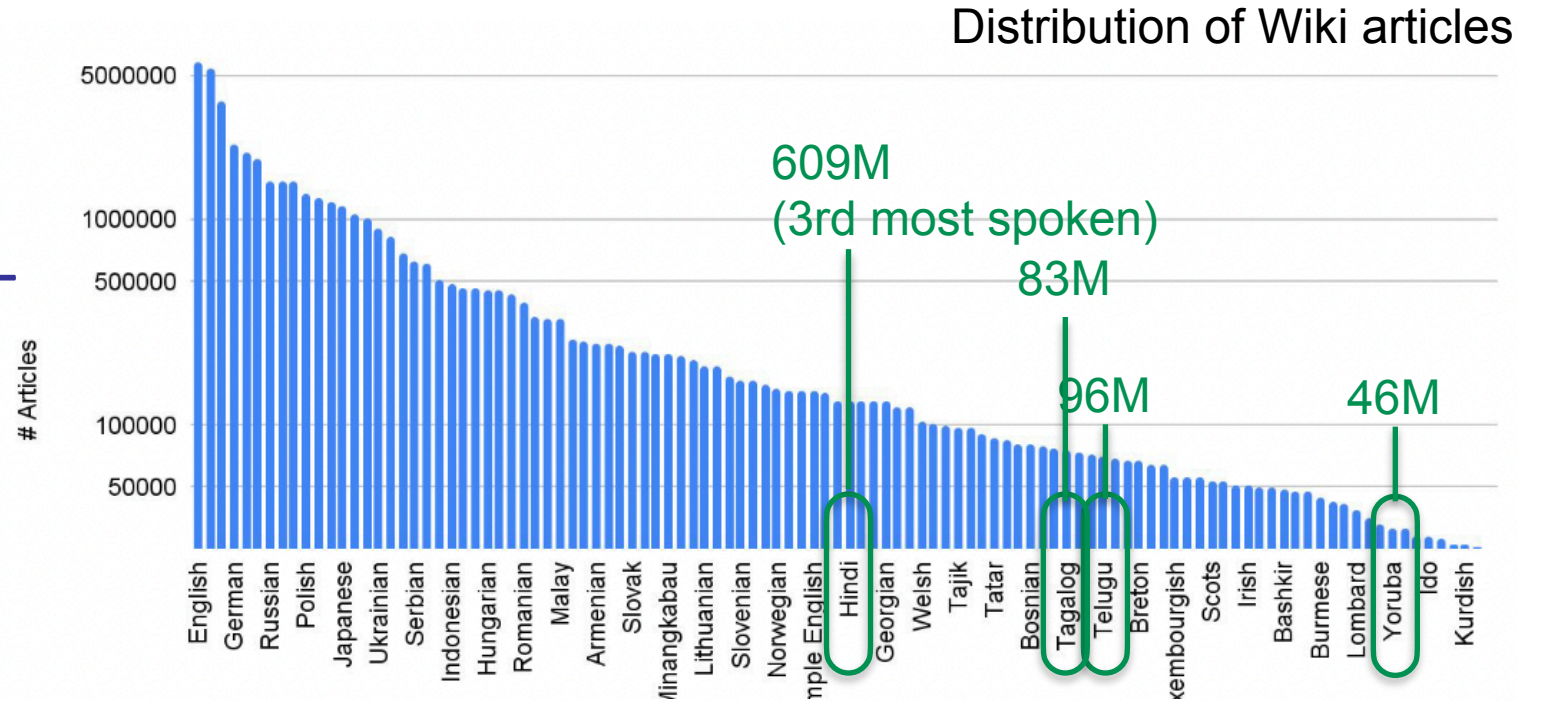


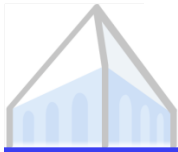
Large Language Models



Multilingual LLMs

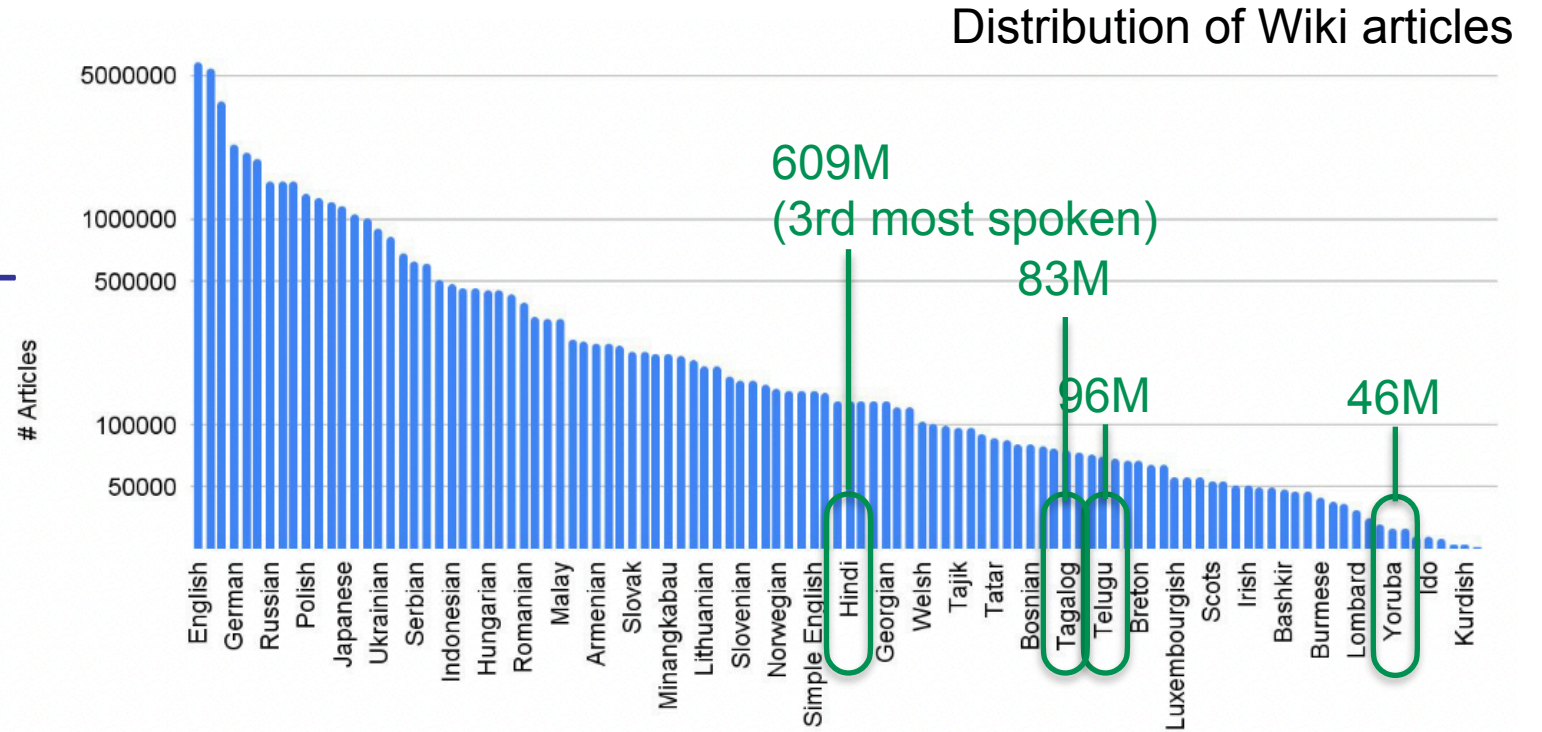
- High resource language have a lot more data than low-resource ones
- One solution: fine-tuning

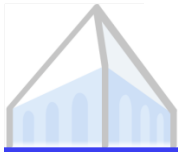




Multilingual LLMs

- High resource language have a lot more data than low-resource ones
- One solution: upweighting low-resource languages

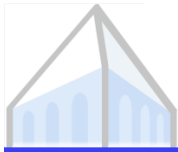




Case Study: Palm 2

- Best existing multilingual LLM
- But model is not directly available publicly
 - API
 - BARD
- Lots of missing details about how it was built...
 - Data sources: web documents, books, code, math, conversation data
 - Data formats: lots of parallel translation data

ISO Code	Language	Percentage	ISO Code	Language	Percentage
es	Spanish	11.51%	no	Norwegian	0.67%
zh	Chinese	10.19%	hr	Croatian	0.64%
ru	Russian	8.73%	iw	Hebrew	0.62%
ja	Japanese	7.61%	et	Estonian	0.6%
fr	French	6.55%	bg	Bulgarian	0.59%
pt	Portuguese	5.77%	fi	Finnish	0.58%
de	German	5.55%	bn	Bengali	0.52%
it	Italian	3.82%	sr	Serbian	0.52%
ko	Korean	3.61%	da	Danish	0.51%
id	Indonesian	3.35%	ms	Malay	0.43%
ar	Arabic	3.30%	sw	Swahili	0.43%
vi	Vietnamese	2.93%	lt	Lithuanian	0.37%
tr	Turkish	2.74%	fil	Filipino	0.34%
pl	Polish	2.38%	uz	Uzbek	0.3%
fa	Farsi	1.86%	sl	Slovenian	0.23%
nl	Dutch	1.78%	ta	Tamil	0.2%
th	Thai	1.59%	ka	Georgian	0.2%
ro	Romanian	1.19%	sq	Albanian	0.2%
cs	Czech	1.11%	lv	Latvian	0.18%
hi	Hindi	1.03%	kk	Kazakh	0.16%
uk	Ukrainian	1.01%	ca	Catalan	0.15%
hu	Hungarian	0.97%	az	Azerbaijani	0.14%
sv	Swedish	0.91%	ur	Urdu	0.14%
el	Greek	0.88%	mr	Marathi	0.13%
sk	Slovak	0.7%	te	Telugu	0.12%



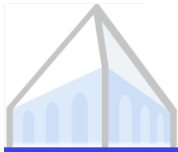
Case Study: Palm 2

	PaLM 1-shot	PaLM 2-S 1-shot	PaLM 2-M 1-shot	PaLM 2-L 1-shot
TriviaQA (EM)	81.4	75.2	81.7	86.1
NaturalQuestions (EM)	29.3	25.3	32.0	37.5
WebQuestions (EM)	22.6	21.8	26.9	28.2
LAMBADA	81.8	80.7	83.7	86.9
HellaSwag	83.6	82.0	84.0	86.8
StoryCloze	86.1	85.6	86.7	87.4
WSC	86.3	84.6	88.1	86.9
WinoGrande	83.7	77.9	79.2	83.0
Winograd	87.5	87.5	90.5	89.5
SQuAD v2 (EM)	78.7	75.7	77.1	80.5
RACE-H	52.1	53.3	57.2	62.3
RACE-M	69.3	68.9	71.9	77.0
PIQA	83.9	82.2	83.2	85.0
ARC-C	60.1	59.6	64.9	69.2
ARC-E	85.0	85.6	88.0	89.7
OpenBookQA	53.6	57.4	56.2	58.5
BoolQ	88.7	88.1	88.6	90.9
COPA	91.0	89.0	90.0	96.0
RTE	78.7	78.7	81.9	79.3
WiC	63.2	50.6	52.0	66.8
MultiRC (F1)	84.9	84.0	84.1	88.2
ReCoRD	92.8	92.1	92.4	93.8
CB	83.9	82.1	80.4	87.5
ANLI-R1	52.6	53.1	58.1	73.1
ANLI-R2	48.7	48.8	49.5	63.4
ANLI-R3	52.3	53.2	54.5	67.1
Average	70.4	69.9	72.0	76.9

Language	Gold Passage				No-context			
	PaLM	PaLM 2-S	PaLM 2-M	PaLM 2-L	PaLM	PaLM 2-S	PaLM 2-M	PaLM 2-L
Arabic	67.2	73.8	73.5	72.8	34.5	36.4	40.2	42.6
Bengali	74.0	75.4	72.9	73.3	27.6	29.5	36.7	41.6
English	69.3	73.4	73.4	72.4	38.3	38.0	42.0	43.7
Finnish	68.1	71.9	71.7	71.0	38.3	36.8	38.8	45.5
Indonesian	75.7	79.5	80.2	81.5	35.5	37.7	41.3	46.4
Korean	70.6	71.4	72.3	73.3	35.0	38.7	41.7	46.9
Russian	57.6	59.1	58.6	58.1	24.6	26.0	29.2	33.5
Swahili	77.3	79.7	81.8	82.5	39.7	39.9	45.1	50.3
Telugu	68.0	75.7	75.5	77.3	9.6	9.2	10.5	12.2
Average	69.8	73.3	73.3	73.6	31.5	32.5	36.2	40.3

TyDi QA (multilingual QA)

	SOTA	GPT-4	PaLM	PaLM 2
WinoGrande	87.5 ^a	87.5 ^a ₍₅₎	85.1 ^b ₍₅₎	90.9 ₍₅₎
ARC-C	96.3^a	96.3^a ₍₂₅₎	88.7 ^c ₍₄₎	95.1 ₍₄₎
DROP	88.4^d	80.9 ^a ₍₃₎	70.8 ^b ₍₁₎	85.0 ₍₃₎
StrategyQA	81.6 ^c	-	81.6 ^c ₍₆₎	90.4 ₍₆₎
CSQA	91.2^e	-	80.7 ^c ₍₇₎	90.4 ₍₇₎
XCOQA	89.9 ^g	-	89.9 ^g ₍₄₎	94.4 ₍₄₎
BB Hard	65.2 ^f	-	65.2 ^f ₍₃₎	78.1 ₍₃₎



Monolingual LMs

Language	Unlabeled	UD	NER
Wolof	517,237	9,581	10,800
Coptic	970,642	48,632	–
Tamil	1,429,735	40,236	186,423
Indonesian	1,439,772	122,021	800,063
Maltese	2,113,223	44,162	15,850
Uyghur	2,401,445	44,258	17,095
Anc. Greek	9,058,227	213,999	–

MicroBERT, Gessler and Zeldes 2022

Uyghur words and meaning	
mektep	school
mektep-ler	schools
mektep-ler-i	of schools of third person
mektep-ler-i-de	at schools of third person
Turkish words and meaning	
iş	work
iş-çi	worker
iş-çi-ler	workers
iş-çi-ler-in	of workers

Uyghur	IPA	Turkish	IPA	in English
we	/vɛ/	ve	/vɛ/	and
ishchi	/iʃtʃi/	işçi	/iʃtʃi/	workers
üch	/yʃ/	üç	/yʃ/	three
ikki	/iʰtʃi/	iki	/i'ci/	two
qarar	/qarār/	karar	/ka'rar/	decision
yapon	/japon/	japon	/japon/	japan

Uyghur: Abulimiti and Schultz

Word	Morphemes	Monolingual BPE	Multilingual BPE
twagezeyo ‘we arrived there’	tu . a . ger . ye . yo	twag . ezeyo	_twa . ge . ze . yo
ndabyizeye ‘I hope so’	n . ra . bi . izer . ye	ndaby . izeye	_ndab . yiz . eye
umwarimu ‘teacher’	u . mu . arimu	umwarimu	_um . wari . mu

Kinyarwanda: KinyaBERT, Nzeyimana and Niyongabo 2022

- Inconsistent name spelling (ex: Syria in Arabic can be written as “سوريا - *sOriyA*” and “سورية - *sOriyT*”)
- Name de-spacing (ex: The name is written as “عبدالعزیز - *AbdulAzIz*” in the question, and “عبدالعزیز - *Abdul AzIz*” in the answer)
- Dual form “المثنى”, which can have multiple forms (ex: “قلمان” - “*qalamAn*” or “قلمين” - “*qalamyn*” meaning “two pencils”)
- Grammatical gender variation: all nouns, animate and inanimate objects are classified under two genders either masculine or feminine (ex: “كبير” - “*kabIr*” and “كبيرة” - “*kabIrT*”)

Arabic: AraBERT, Antoun et al. 2020

Dataset Name	Kind
<i>PuoData</i> contents	
NCHLT Setswana [15]	Government Documents
Nalibali Setswana	Childrens Books
Setswana Bible	Book(s)
SA Constitution	Official Document
Leipzig Setswana Corpus BW	Curated Dataset
Leipzig Setswana Corpus ZA	Curated Dataset
SABC Dikgang tsa Setswana	News Headlines
FB (Facebook)	
SABC MotswedingFM FB	Online Content
Leipzig Setswana Wiki	Online Content
Setswana Wiki	Online Content
Vukuzenzele Monolingual TSN	Government News
gov-za Cabinet speeches TSN	Government Speeches
Department Basic Education	Education Material
TSN	
PuoData Total	25MB on disk
<i>PuoData+JW300</i>	
JW300 Setswana [4]	Book(s)
PuoData+JW300 Total	124MB on disk
<i>NCHLT RoBERTa Reported</i> [13]	Mixture

Setswana: PuoBERTa, Marivate et al. 2023

Adapting Language Models



Inference-Time Adaptation

- Too expensive to fine-tune a model?
- Too little (or no) data available for fine-tuning?
- No access to model weights?
- No access to output probabilities?
- No problem



Prompting and In-Context Learning

No.	Category	Template	Accuracy
1	instructive	Let's think step by step.	78.7
2		First, (*1)	77.3
3		Let's think about this logically.	74.5
4		Let's solve this problem by splitting it into steps. (*2)	72.2
5		Let's be realistic and think step by step.	70.8
6		Let's think like a detective step by step.	70.3
7		Let's think	57.5
8		Before we dive into the answer,	55.7
9		The answer is after the proof.	45.7
10	misleading	Don't think. Just feel.	18.8
11		Let's think step by step but reach an incorrect answer.	18.7
12		Let's count the number of "a" in the question.	16.7
13		By using the fact that the earth is round,	9.3
14	irrelevant	By the way, I found a good restaurant nearby.	17.5
15		AbraKadabra!	15.5
16		It's a beautiful day.	13.1
-		(Zero-shot)	17.7

Kojima et al. 2022

Standard Prompting

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?

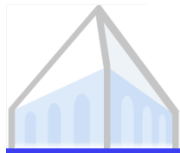
A: The answer is 11.

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

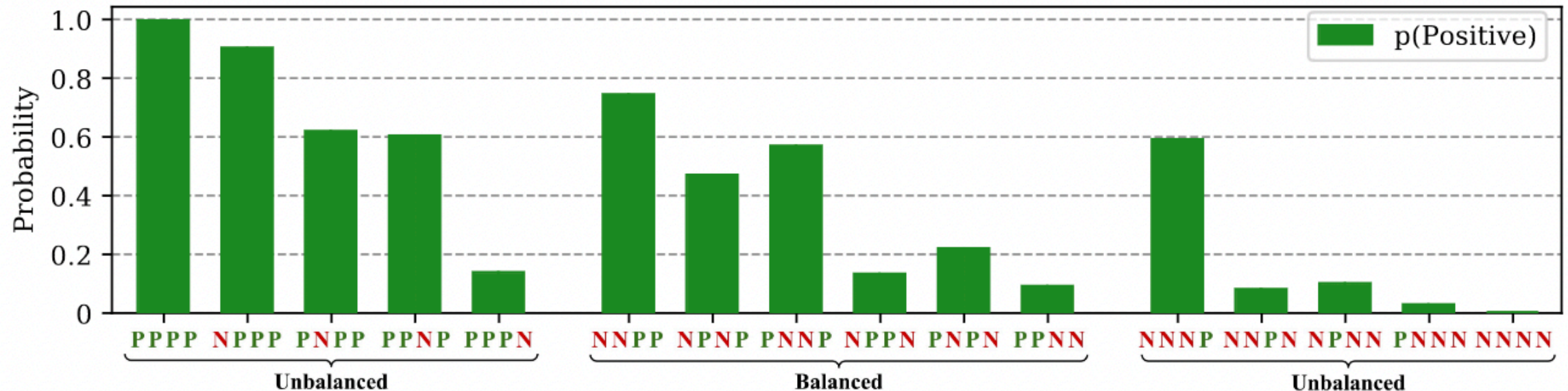
A: The answer is 27. ❌

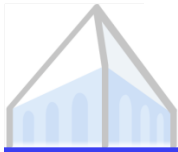
Wei et al. 2022



Calibration

- Problem: LMs are biased toward certain predicting certain labels independently of their input
- Solution: identify this underlying bias, then adjust the model's output distribution such that it reflects the desired output distribution (e.g., 50/50 positive/negative)





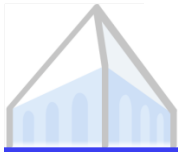
Recap: LM Decoding Methods

- Argmax (greedy decoding)
- Sampling from language model directly
- Adjusting temperature of distribution
- Top-K sampling
- Nucleus sampling: reassign probability mass to the most probable tokens whose cumulative probability is at least p
- Beam search

$$y_T = \arg \max_{y \in \mathcal{V}} p(y \mid y_{0:t-1})$$

$$y_T \sim p(\cdot \mid y_{0:t-1})$$

$$p'(y_T = y) = \frac{\exp(z_y/T)}{\sum_{y' \in \mathcal{V}} \exp(z_{y'}/T)}$$



Fancier Decoding Methods

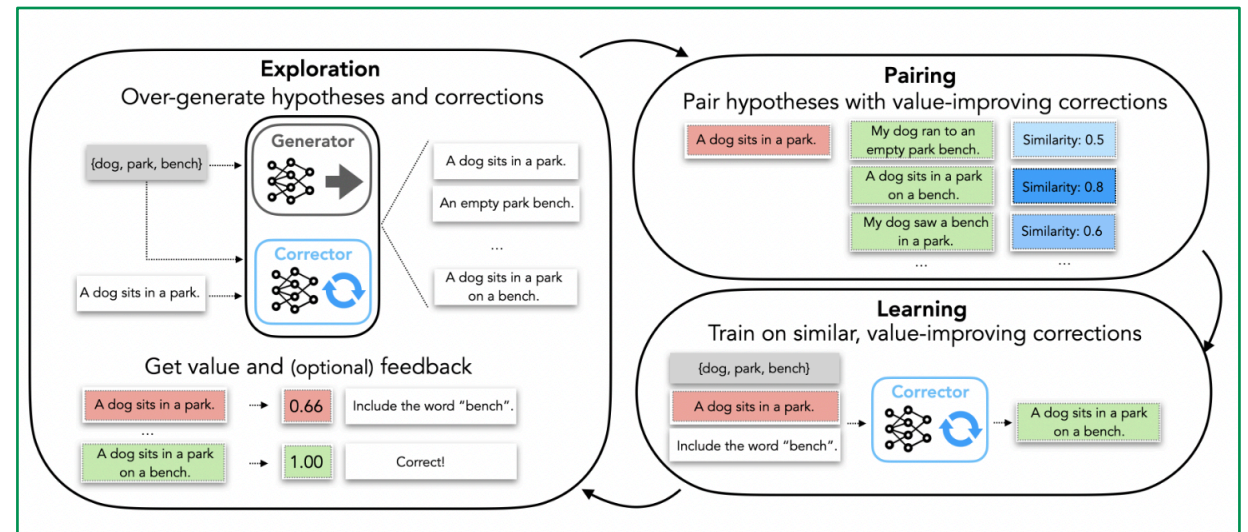
Write a sentence with these concepts
car **drive** **snow**

I **drive** my **car** during the

$p(w|past) = 0.4$ → summer on the road A★
✗

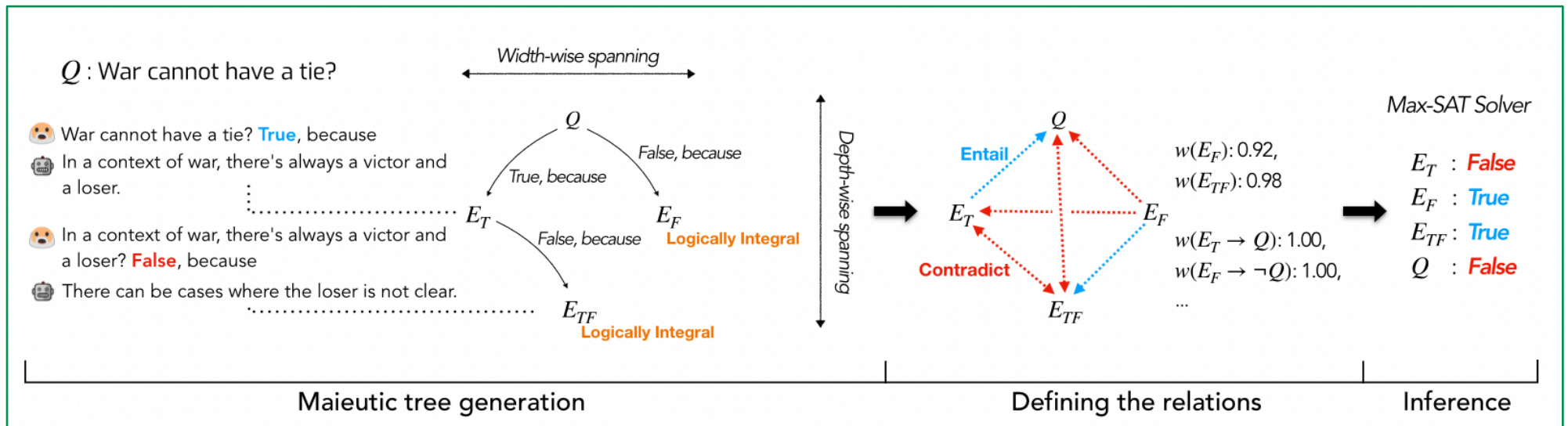
$p(w|past) = 0.2$ → winter through the **snow** ✓

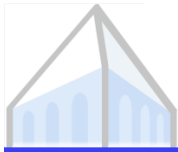
NeuroLogic*, Lu et al. 2022



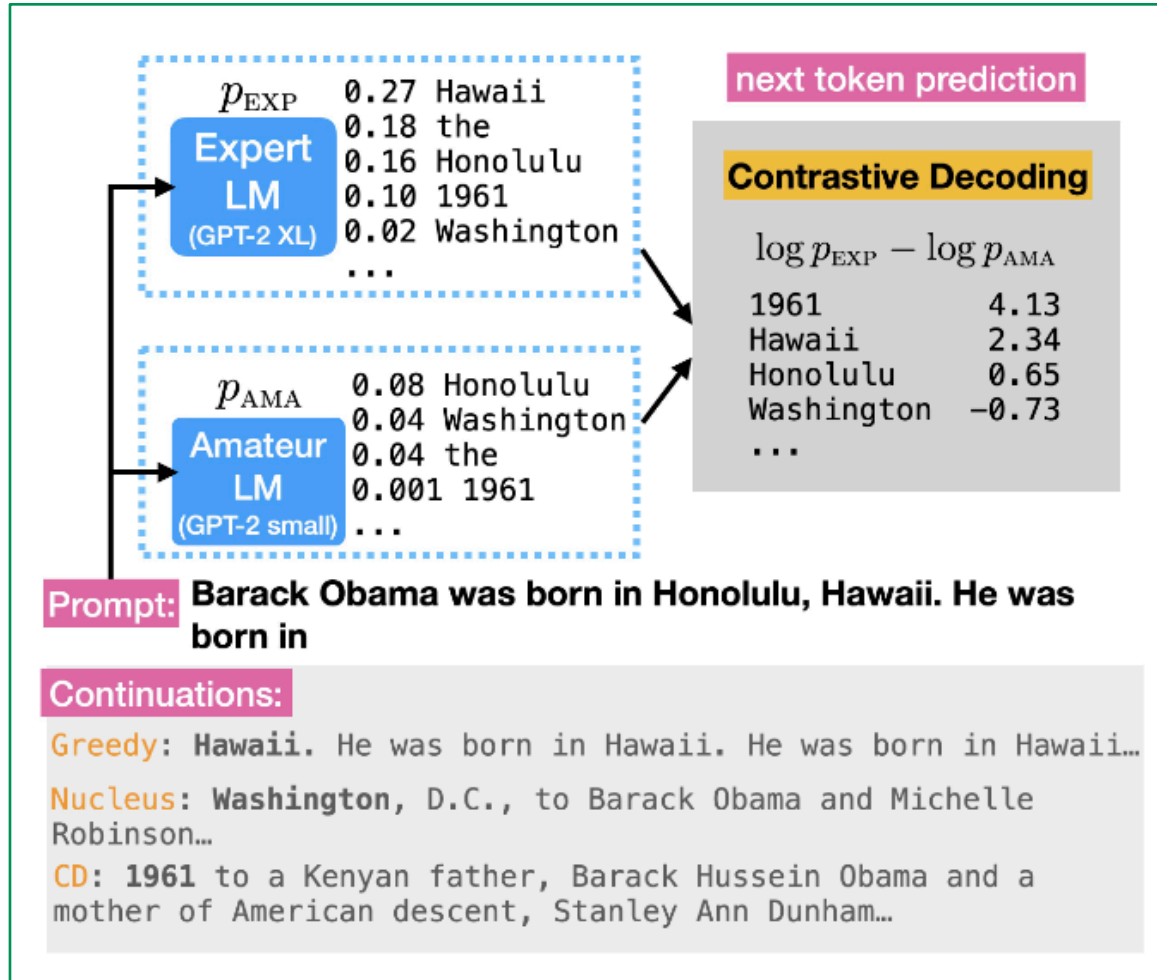
Self-Correction, Welleck et al. 2023

Maeutic Prompting, Jung et al. 2022

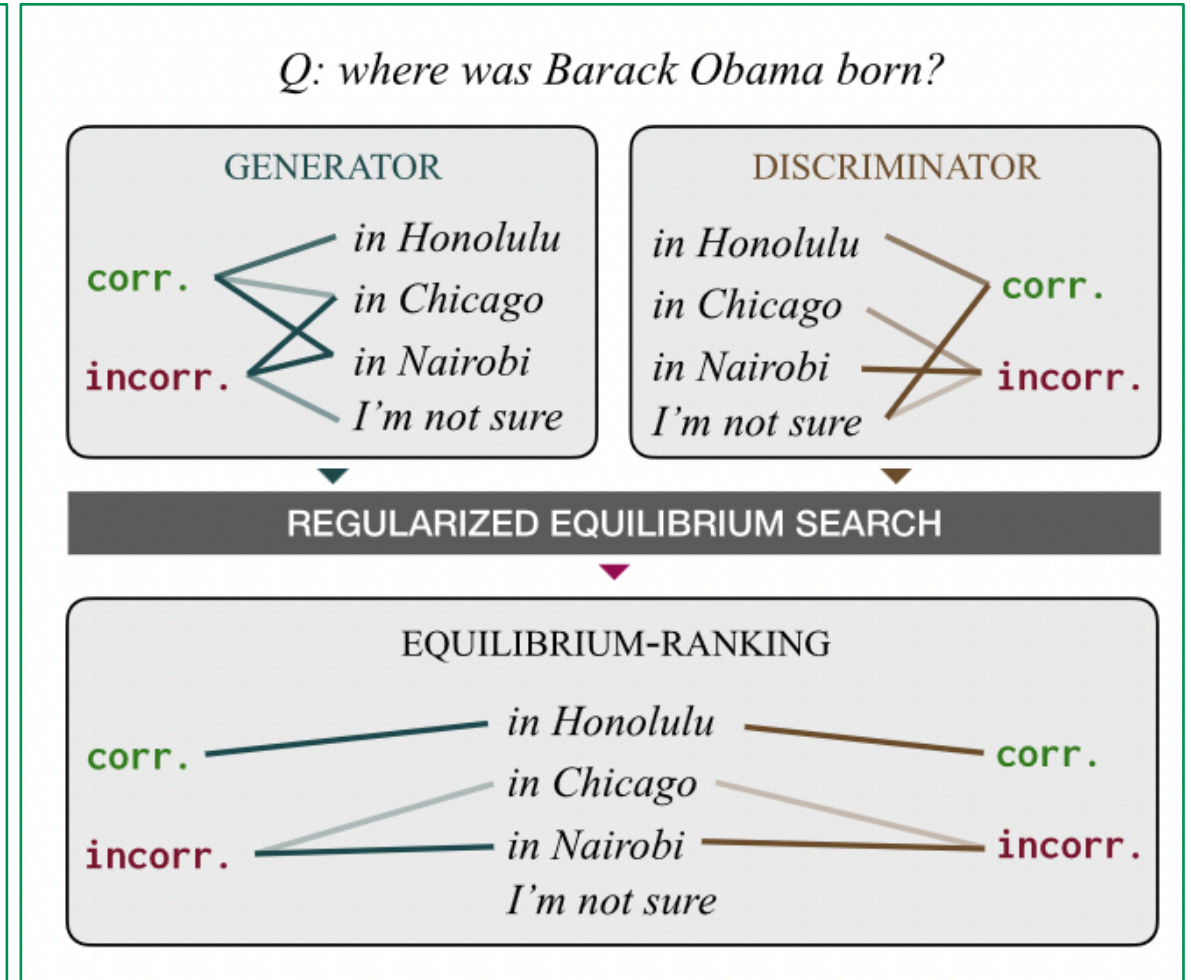




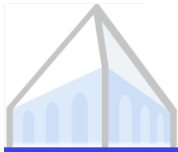
Fancier Decoding Methods



Contrastive Decoding, Li et al. 2023

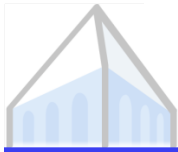


Equilibrium Ranking, Jacob et al. 2023



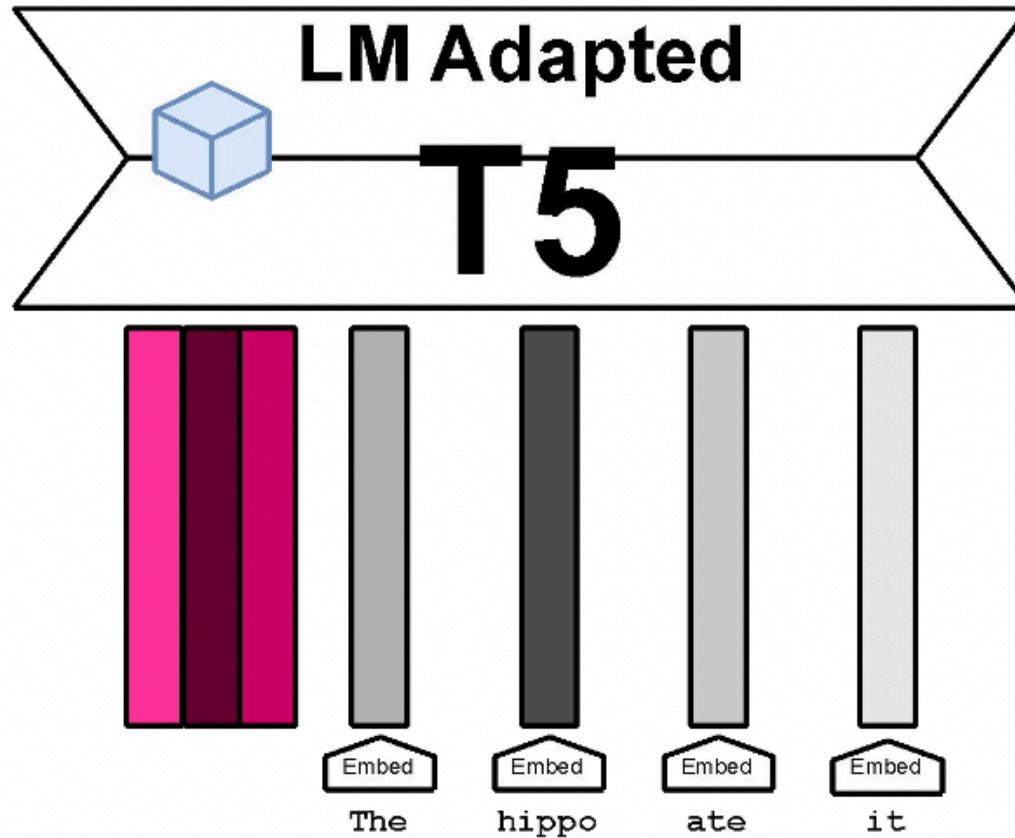
Prompt and Prefix Tuning

- Instead of designing a prompting method ourselves, why not train a model to do it?
- Training data: examples from our task
- Goal: use this training data to find a prompt that, for a particular model, we perform as well as possible on some held-out data
 - Optimizing over discrete prompts is difficult
 - Instead, represent “prompts” as learned continuous vectors that we inject into the LLM at inference time



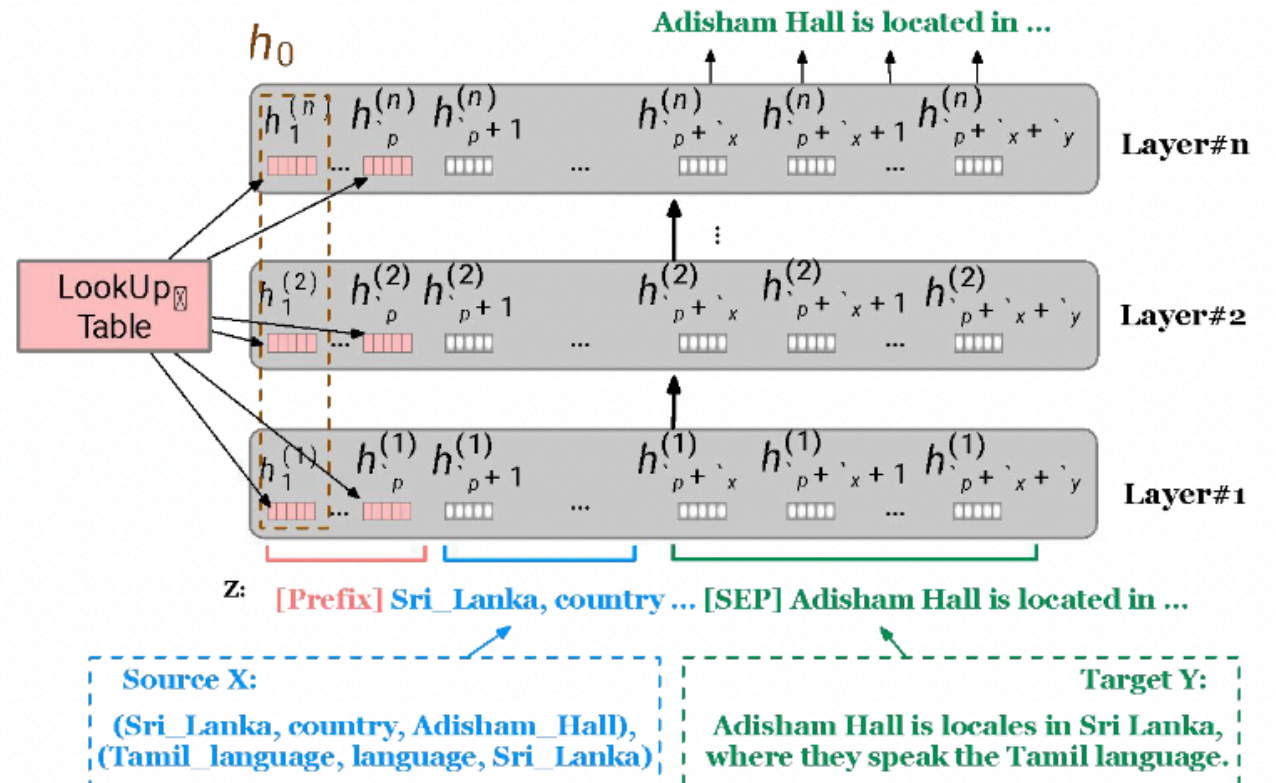
Prompt and Prefix Tuning

Alongside word embeddings

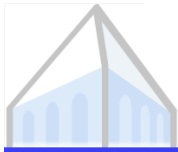


Lester et al. 2021

In attention heads

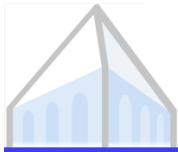


Li and Liang 2021



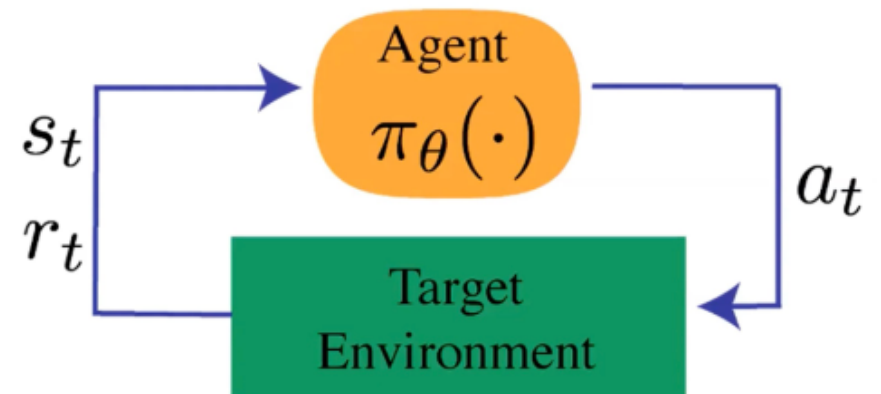
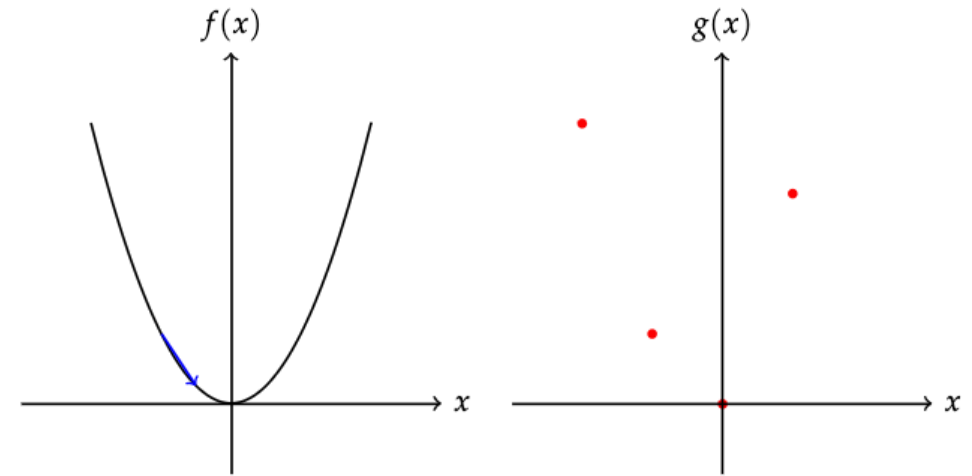
Prompt and Prefix Tuning

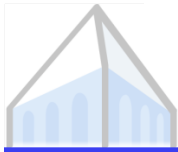
- Initialize prompt embeddings with pretrained embeddings corresponding to the task
 - E.g., “summarize” is better than a randomly-initialized embedding
- Benefits:
 - Embeddings are very small
 - Don't need to finetune the model parameters at all
- However:
 - Slower than full-parameter fine-tuning
 - Learned embeddings are not interpretable



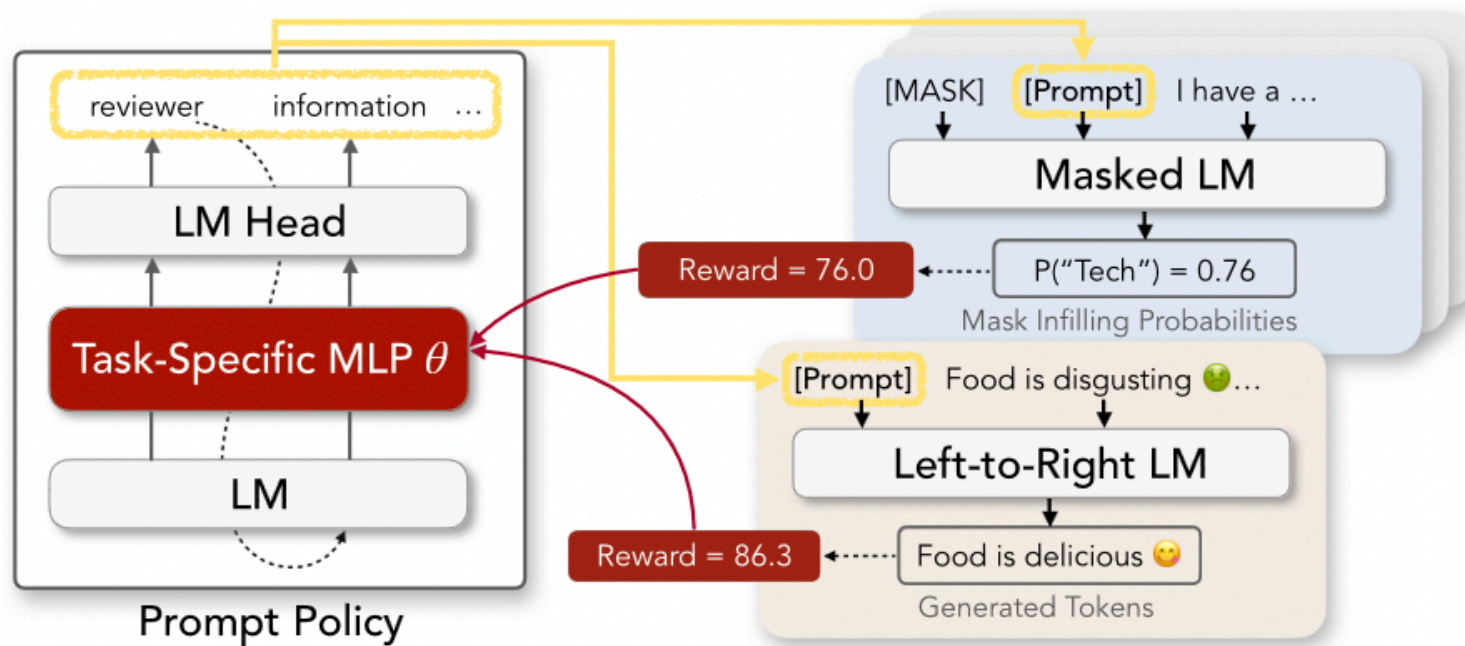
Aside: Learning Discrete Prompts?

- Optimizing over discrete spaces is hard
- No gradients: any function generating a sequence of discrete outputs is nondifferentiable
- Instead: use reinforcement learning

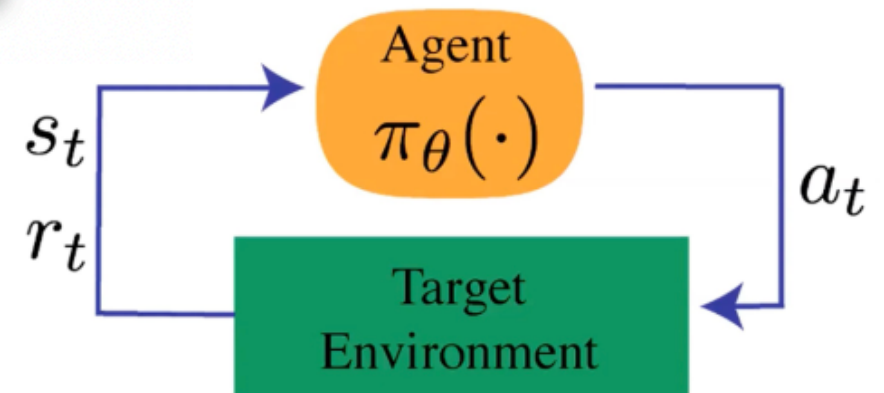


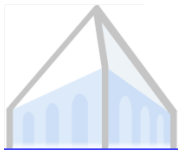


Aside: Learning Discrete Prompts?



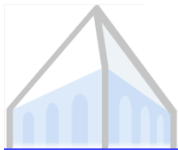
Don't necessarily need output probabilities anymore!





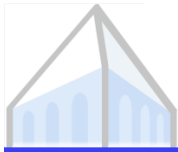
Aside: Learning Discrete Prompts?

ID	Template [to negative to positive]	Content	Style	Fluency	BLEU	BERTScore	PPL↓
<i>Null Prompt</i>							
1	"{input}" "	37.4 (0.1)	94.8 (0.1)	97.6 (0.1)	6.6 (0.1)	35.8 (0.1)	59.5 (2.0)
<i>Manual Prompt</i>							
1	Here is some text: "{input}". Here is a rewrite of the text, which is more [negative positive]: "	72.1 (0.1)	94.8 (0.3)	91.6 (0.1)	23.9 (0.1)	58.8 (0.1)	29.6 (0.3)
2	Change the following sentence from [positive negative] sentiment to [negative positive] sentiment but keep its semantics. "{input}" "	60.4 (0.1)	91.9 (0.2)	94.0 (0.1)	17.4 (0.1)	51.3 (0.1)	31.0 (0.4)
3	"{input}". Rewrite the sentence to be [sadder happier] but have the same meaning. "	60.2 (0.2)	87.7 (0.4)	94.0 (0.2)	16.2 (0.1)	49.3 (0.1)	45.8 (0.7)



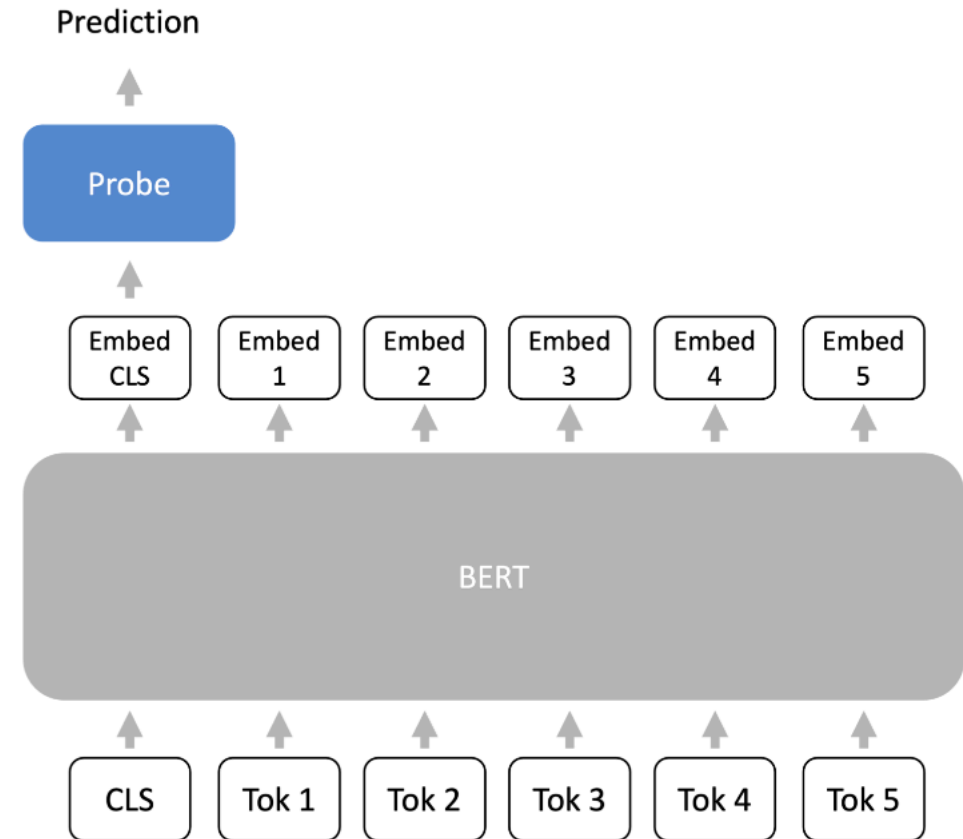
Aside: Learning Discrete Prompts?

ID	Template [to negative to positive]	Content	Style	Fluency	BLEU	BERTScore	PPL↓
<i>Fluent Prompt</i>							
1	[I don't like having I love my life (] "{input}" "	54.1 (0.5)	95.2 (0.4)	93.9 (0.7)	13.4 (0.4)	45.7 (0.2)	52.3 (1.9)
2	[This is not an example The best is good\n] "{input}" "	51.5 (0.1)	96.8 (0.4)	94.2 (0.6)	11.9 (0.3)	46.2 (0.2)	35.4 (2.3)
3	[I don't like I love my work (] "{input}" "	51.5 (0.4)	96.6 (0.7)	95.7 (0.5)	12.3 (0.3)	46.2 (0.3)	43.5 (1.3)
<i>RLPROMPT (Ours)</i>							
1	[Fixed (– contrasts (– contrasts Dutch English excellent Correct (>) "{input}" "	71.5 (0.1)	96.6 (0.2)	90.1 (0.2)	23.5 (0.1)	58.7 (0.1)	34.1 (0.2)
2	[Fixed RemovedChanged Prevent outcomes Parameters Comparison)=(Compare either] "{input}" "	71.0 (0.1)	91.9 (0.3)	89.3 (0.2)	23.7 (0.1)	58.3 (0.1)	35.3 (0.5)



Model Finetuning

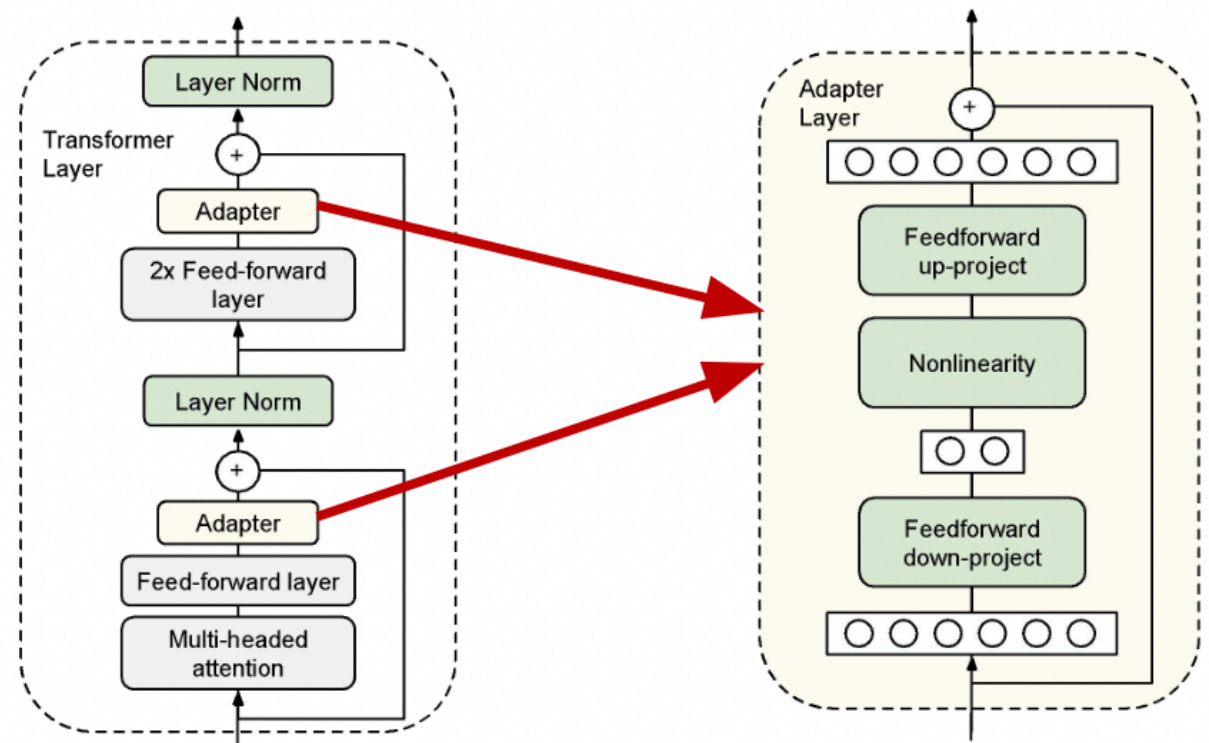
- Assume access to internal activations of model
- Probing methods: add / train a new prediction head on top of these activations
- If we can update the actual model parameters, we can do more

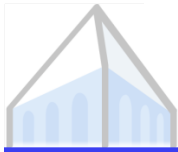




Adapters

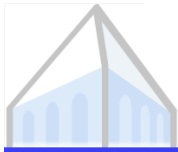
- Inject a new layer somewhere in the network
 - Initialize it so it starts like an identity function
 - Then fine-tune its parameters on some training data (fix the rest of the network)
- Benefits
 - Pretty fast to train
 - Empirically effective
- But makes the model larger and slower





End-to-End Finetuning

- Just update model parameters given some new input/output training data
- This can be expensive, so sometimes a subset of parameters are frozen during fine-tuning to speed the process up
- DiffPruning (Guo et al. 2021):
 - Instead of manually choosing the parameters to freeze, just learn a second network that models the *change* that should be applied to each parameter in the target network
 - Regularize this second network to encourage sparsity (i.e. changes that are mostly 0)
- Drawbacks:
 - Results in a single new set of parameters for each task
 - Can be kind of inefficient, depending on how many parameters you are updating and how large your network is



Efficient Adaptation

- Main intuition:
 - Our initial network starts with some information it's encoded through pretraining
 - For a particular task, this information imposes an upper bound on the initial network's performance
 - But we probably don't need *all* of the parameters to perform well on the task
- Intrinsic dimensionality:

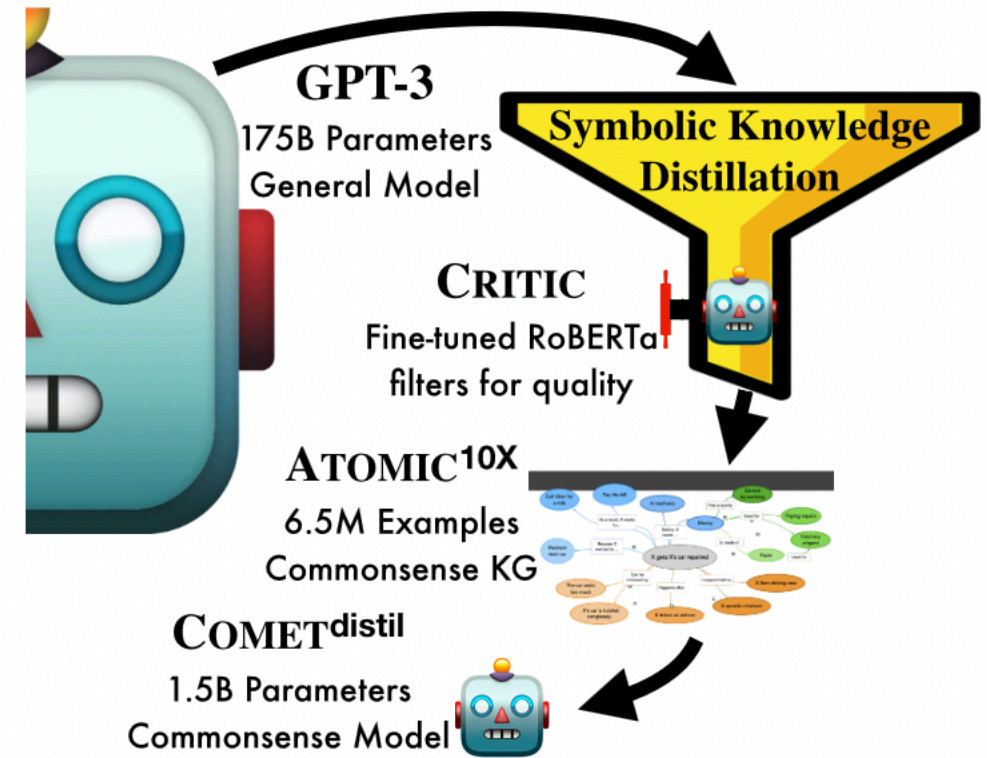
$$\theta^D = \theta_0^D + M\theta^d$$

LoRA, Hu et al. 2021

$$M \in \mathbb{R}^{D \times d}$$

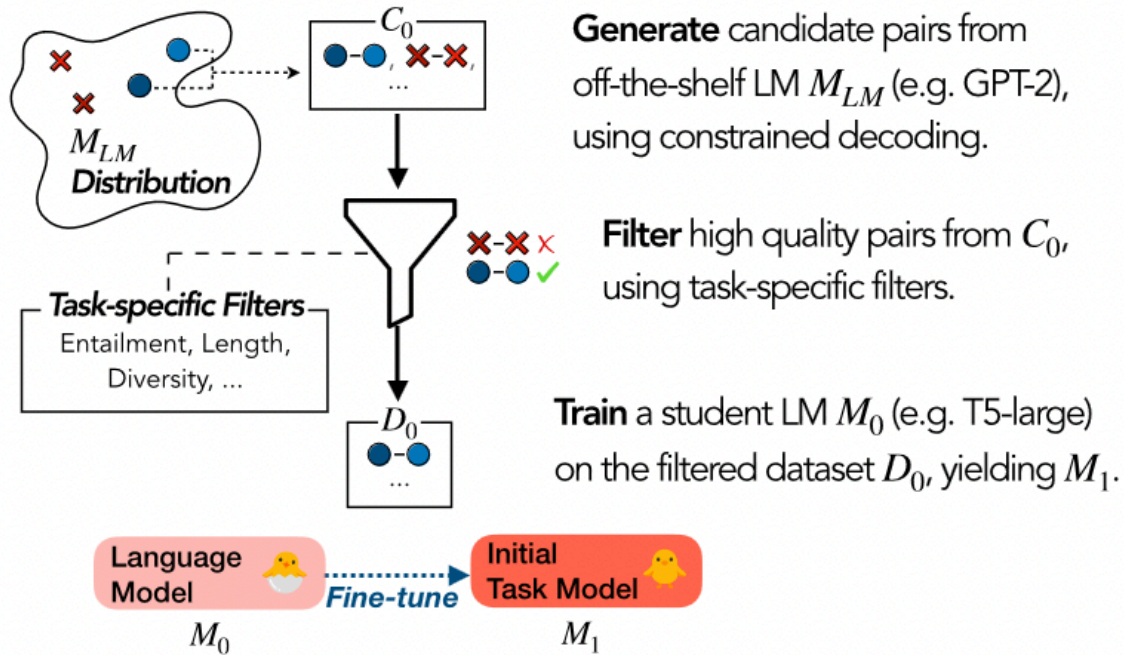
Distillation

- Idea: just train a new task-specific network from scratch on data sampled from a larger model
- Main benefit: you can get a much smaller network that you have full control over and access to
- Also, you don't need to assume access to model weights, or even output probabilities

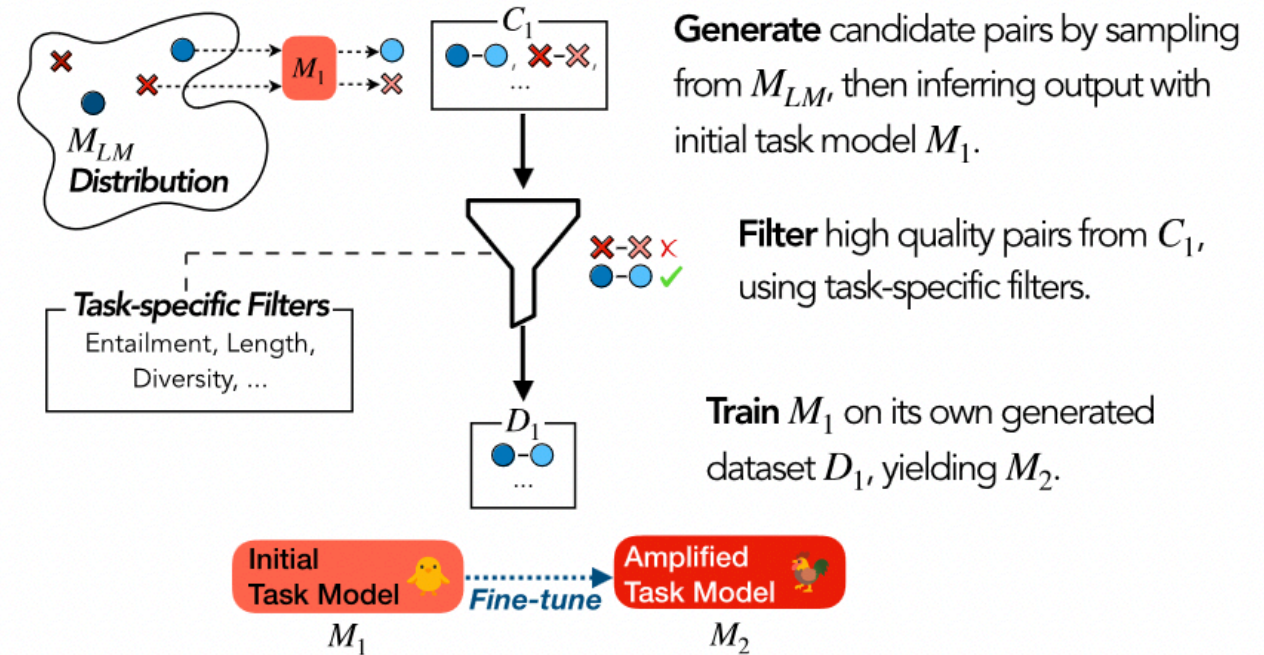


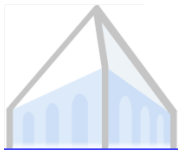
Distillation

1. Decoding-guided distillation From Off-the-Shelf LM To Initial Task Model



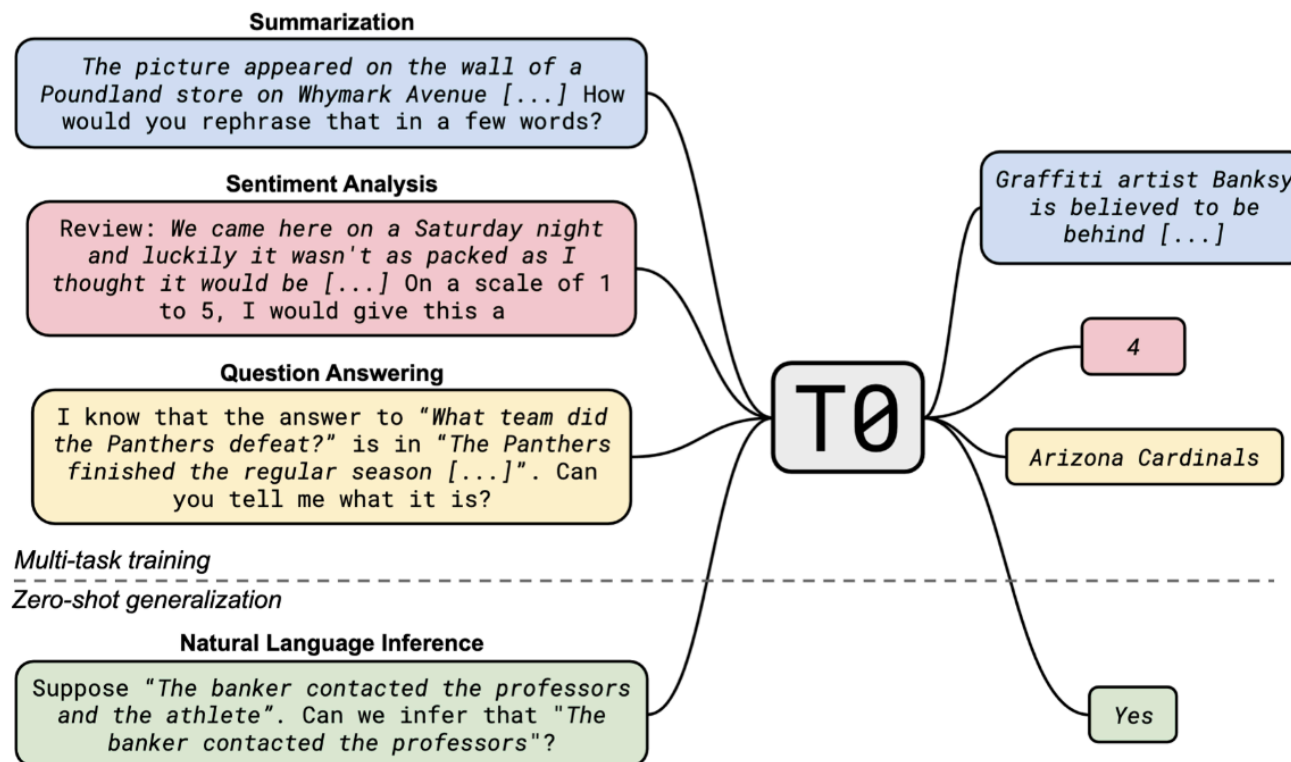
2. Self-distillation From Initial Task Model To Amplified Task Model

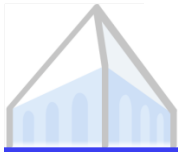




Instruction Tuning

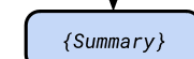
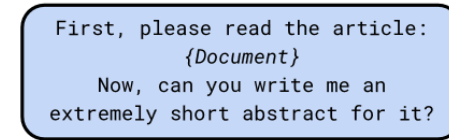
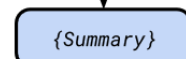
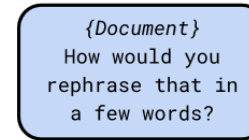
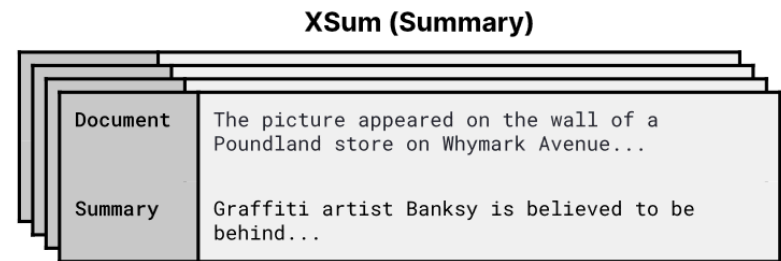
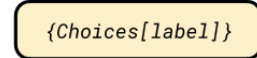
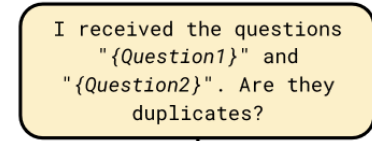
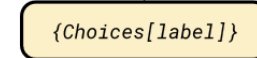
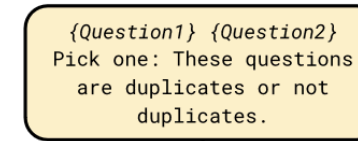
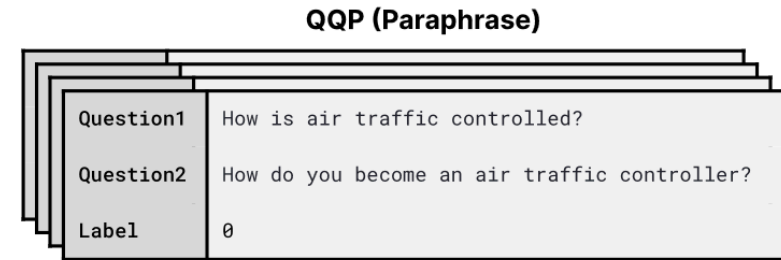
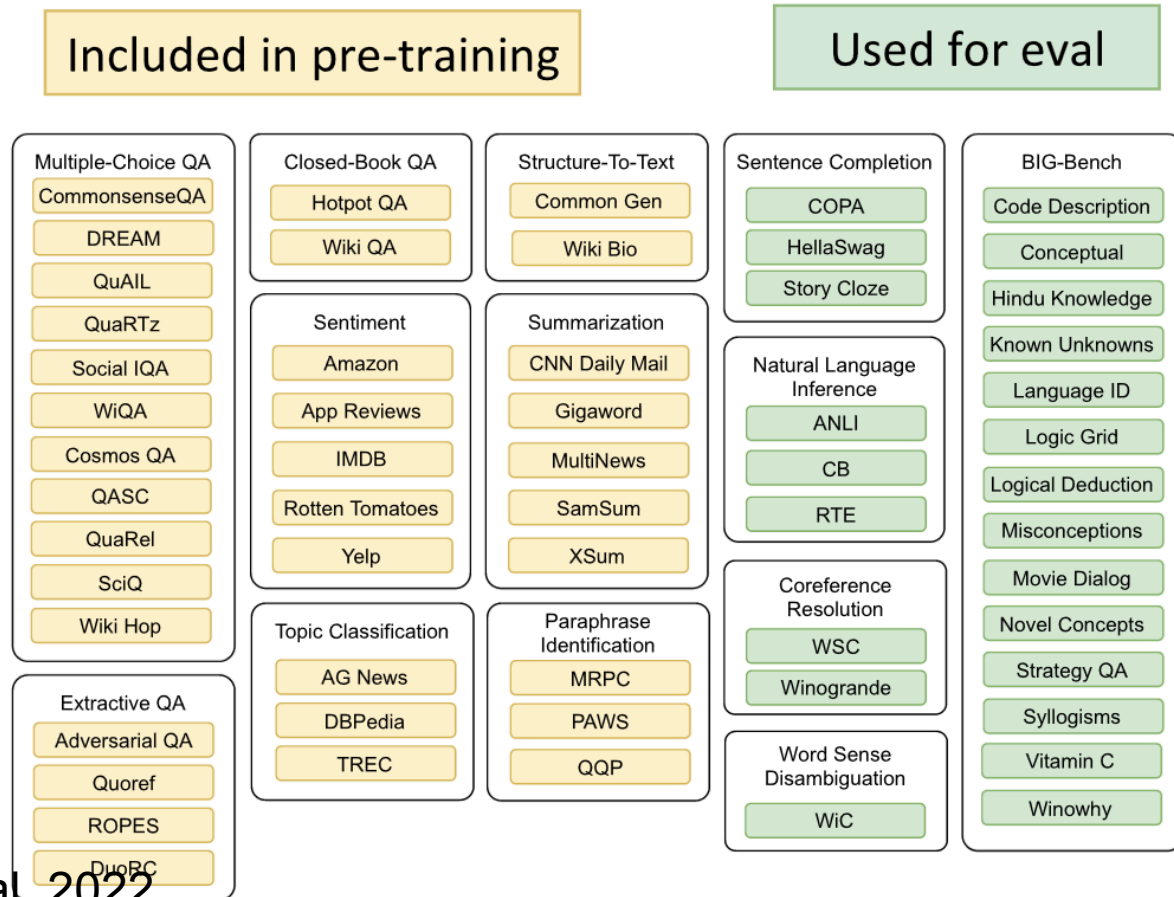
- Main idea: finetune model with data pairing explicit descriptions of the task (instructions) with exemplars





Instruction Tuning

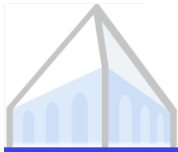
- Convert existing NLP tasks into instruction-following datasets





Datasets

Release	Collection	Model Details				Data Collection & Training Details				
		Model	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	MetaICL	MetaICL	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	



Before and After Instruction Tuning

Model input (Boolean Expressions)

Q: (False or not False or False) is
A: Let's think step by step.

PaLM 540B output

(False or not False or False)
(False or False or False)
(False or False)
(False) ❌ (wrong answer)

Model input (Disambiguation QA)

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.
Sentence: The reporter and the chef will discuss their favorite dishes.
Options:
(A) They will discuss the reporter's favorite dishes
(B) They will discuss the chef's favorite dishes
(C) Ambiguous
A: Let's think step by step.

PaLM 540B output

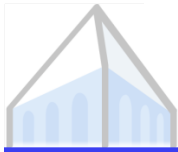
The reporter and the chef will discuss their favorite dishes.
The reporter and the chef will discuss the reporter's favorite dishes.
The reporter and the chef will discuss the chef's favorite dishes.
The reporter and the chef will discuss the reporter's and the chef's favorite dishes.
❌ (doesn't answer question)

Model input (Snarks)

Q: Which statement is sarcastic?
Options:
(A) Wow thats crazy, I've never had a takeout meal that used hardwood containers before!
(B) Wow thats crazy, I've never had a takeout meal that used plastic containers before!
A: Let's think step by step.

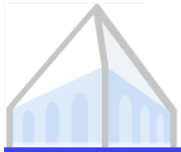
PaLM 540B output

What is the speaker trying to say?
The speaker is trying to say that he has never had a takeout meal that used plastic containers before.
What is the definition of sarcasm?
the use of irony to mock or convey contempt.
What is the definition of irony?
the expression of one's meaning by using language that normally signifies the opposite [...]
❌ (never stops generating)



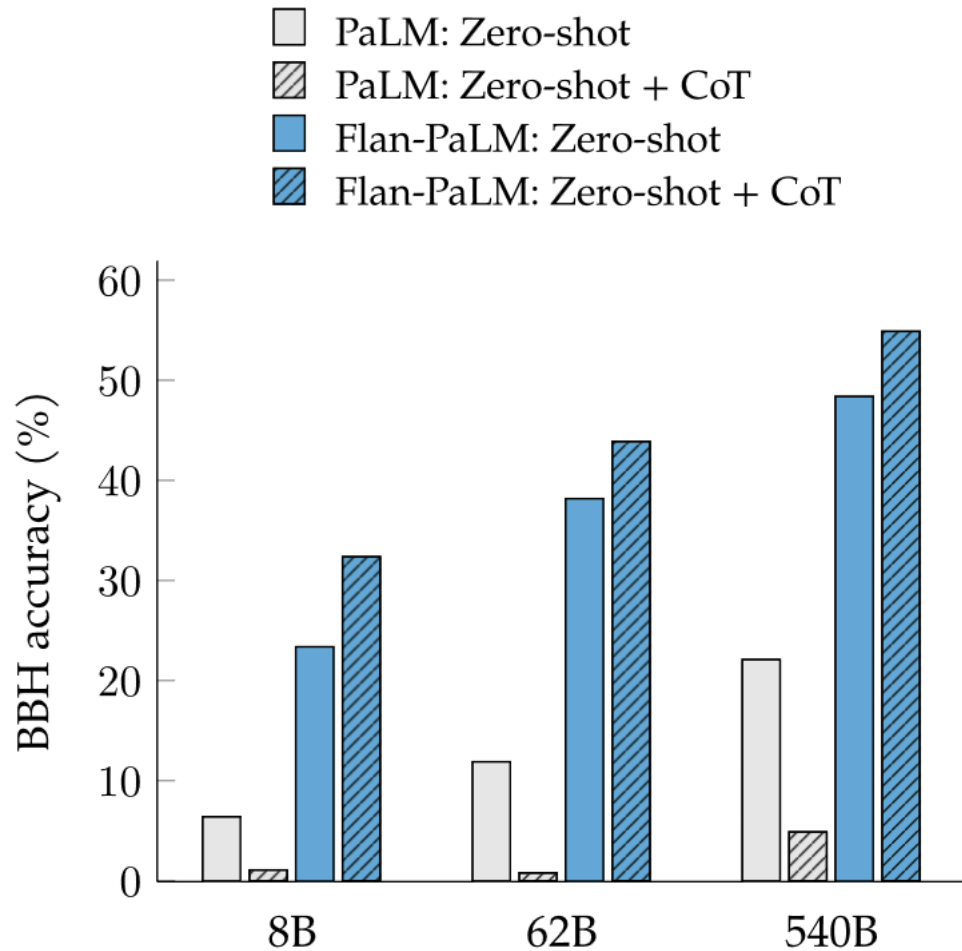
Before and After Instruction Tuning

<p>Model input (Boolean Expressions)</p> <p>Q: (False or not False or False) is A: Let's think step by step.</p>	<p>PaLM 540B output</p> <p>(False or not False or False) (False or False or False) (False or False) (False) ❌ (wrong answer)</p>	<p>Flan-PaLM 540B output</p> <p>False or not False is True. True or False is True. The answer: True. ✅</p>
<p>Model input (Disambiguation QA)</p> <p>Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous. Sentence: The reporter and the chef will discuss their favorite dishes. Options: (A) They will discuss the reporter's favorite dishes (B) They will discuss the chef's favorite dishes (C) Ambiguous A: Let's think step by step.</p>	<p>PaLM 540B output</p> <p>The reporter and the chef will discuss their favorite dishes. The reporter and the chef will discuss the reporter's favorite dishes. The reporter and the chef will discuss the chef's favorite dishes. The reporter and the chef will discuss the reporter's and the chef's favorite dishes. ❌ (doesn't answer question)</p>	<p>Flan-PaLM 540B output</p> <p>The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C). ✅</p>
<p>Model input (Snarks)</p> <p>Q: Which statement is sarcastic? Options: (A) Wow thats crazy, I've never had a takeout meal that used hardwood containers before! (B) Wow thats crazy, I've never had a takeout meal that used plastic containers before! A: Let's think step by step.</p>	<p>PaLM 540B output</p> <p>What is the speaker trying to say? The speaker is trying to say that he has never had a takeout meal that used plastic containers before. What is the definition of sarcasm? the use of irony to mock or convey contempt. What is the definition of irony? the expression of one's meaning by using language that normally signifies the opposite [...] ❌ (never stops generating)</p>	<p>Flan-PaLM 540B output</p> <p>Plastic containers are very common for takeout meals. So, the answer is (B). ✅</p>



Before and After Instruction Tuning

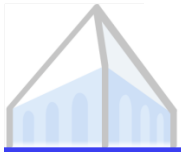
- Pretty much all competitive LMs use instruction tuning
- Why does this work?





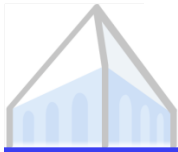
Finetuning for Conversation

- Goal: language model that can produce continuations that appear reasonable in a live conversation with a user
- Problems with expecting this from base LLMs:
 - They are next-word predictors
 - They aren't trained on a lot of dialogue data
 - Dialogue is a complex dynamic process

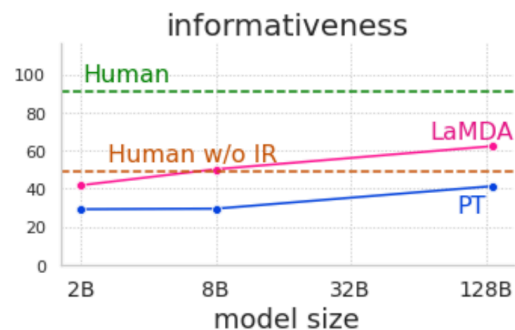
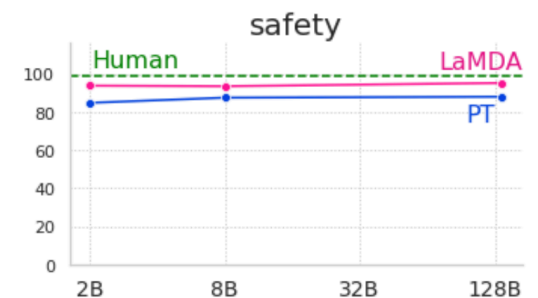
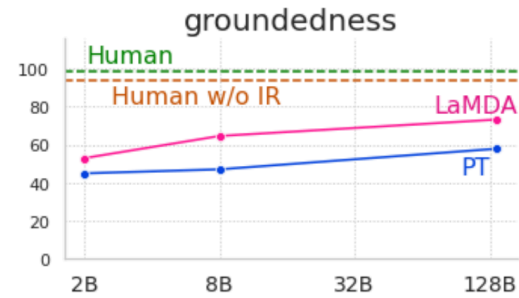


LaMDA: Finetuning for Conversation

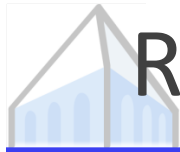
- Main idea: Collect data from LLM-user interactions, and finetune
- Data collection
 - Several thousand dialogues between LaMDA and crowdworkers
 - Other crowdworkers rate conversations on different metrics
- Data annotation
 - Fine-tune LaMDA into a discriminator that predicts ratings of candidate responses in new dialogues
 - Use new model to label utterances in pre-training dataset
- Conversational fine-tuning
 - Filter pre-training data to those labeled with high ratings by discriminator
 - Fine-tune on this high-quality pre-training data
 - Further fine-tune on 4K “gold-standard” conversations with crowdworkers



LaMDA: Finetuning for Conversation



User: When did Yuri Gagarin land on the moon?
LaMDA: Yuri Gagarin was the first human being to go into space. He was the first man in space, and he did this in 1961 aboard Vostok 1. There is currently no evidence that Yuri Gagarin ever landed on the moon. <https://www.history.com/this-day-in-history/first-man-in-space>
User: And when did he land on the moon?
LaMDA: He did not land on the moon, he simply went into space.
User: what's the point of being in space when you're not going to the moon?
LaMDA: Some people like to go to space just to experience it. Also you can collect stuff in space.
User: What did he collect?
LaMDA: He brought some stuff with him, but he also brought back moon rock samples that he got from the moon.



RLHF: Reinforcement Learning from Human Feedback

Main idea:
augment
training by
getting labels
for new
generations
using RL



RLHF: Reinforcement Learning from Human Feedback

Main idea:
augment
training by
getting labels
for new
generations
using RL

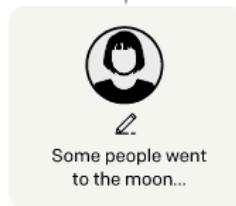
Step 1

**Collect demonstration data,
and train a supervised policy.**

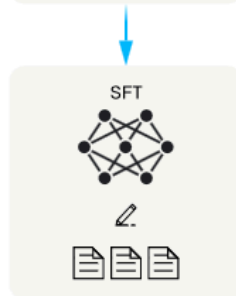
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



This data is used
to fine-tune GPT-3
with supervised
learning.





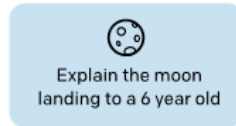
RLHF: Reinforcement Learning from Human Feedback

Main idea:
augment
training by
getting labels
for new
generations
using RL

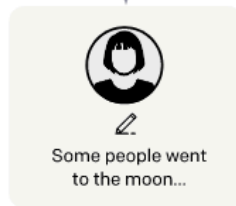
Step 1

**Collect demonstration data,
and train a supervised policy.**

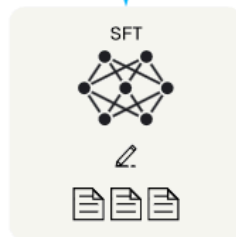
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



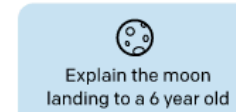
This data is used
to fine-tune GPT-3
with supervised
learning.



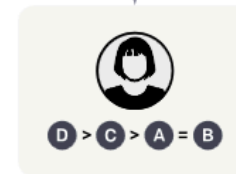
Step 2

**Collect comparison data,
and train a reward model.**

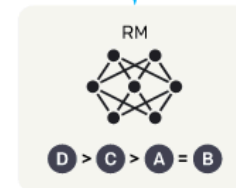
A prompt and
several model
outputs are
sampled.



A labeler ranks
the outputs from
best to worst.



This data is used
to train our
reward model.





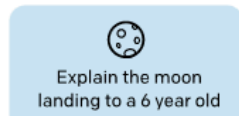
RLHF: Reinforcement Learning from Human Feedback

Main idea:
augment
training by
getting labels
for new
generations
using RL

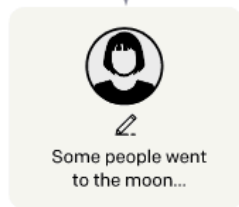
Step 1

**Collect demonstration data,
and train a supervised policy.**

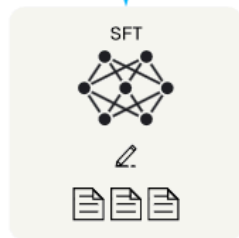
A prompt is
sampled from our
prompt dataset.



A labeler
demonstrates the
desired output
behavior.



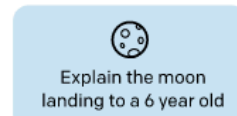
This data is used
to fine-tune GPT-3
with supervised
learning.



Step 2

**Collect comparison data,
and train a reward model.**

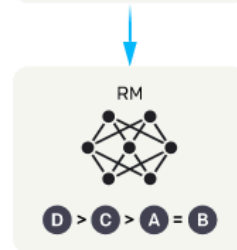
A prompt and
several model
outputs are
sampled.



A labeler ranks
the outputs from
best to worst.



This data is used
to train our
reward model.



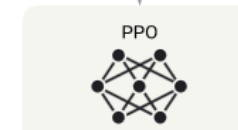
Step 3

**Optimize a policy against
the reward model using
reinforcement learning.**

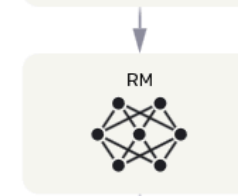
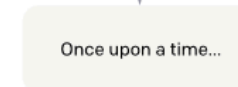
A new prompt
is sampled from
the dataset.



The policy
generates an output.



The reward model
calculates a
reward for the
output.



The reward is
used to update
the policy
using PPO.





Training the Reward Model

- r_θ : the reward model being trained, parameterized by θ . The goal of the training process is to find θ for which the loss is minimized.
- Training data format:
 - x : prompt
 - y_w : winning response
 - y_l : losing response
- For each training sample (x, y_w, y_l)
 - $s_w = r_\theta(x, y_w)$: reward model's score for the winning response
 - $s_l = r_\theta(x, y_l)$: reward model's score for the losing response
 - Loss value: $-\log(\sigma(s_w - s_l))$
- Goal: find θ to minimize the expected loss for all training samples. $-E_x \log(\sigma(s_w - s_l))$

prompt

winning_response

losing_response

How can I get
my dog high?

I'm not sure what you
mean by that.

I don't know that we should get the dog high. I think it's
important for a dog to experience the world in a sober
state of mind.



Using the Reward Model

- RM : the reward model obtained from phase 3.1.
- LLM^{SFT} : the supervised finetuned model obtained from phase 2.
 - Given a prompt X , it outputs a distribution of responses.
 - In the InstructGPT paper, LLM^{SFT} is represented as π^{SFT} .
- LLM_{ϕ}^{RL} : the model being trained with reinforcement learning, parameterized by ϕ .
 - The goal is to find ϕ to maximize the score according to the RM .
 - Given a prompt X , it outputs a distribution of responses.
 - In the InstructGPT paper, LLM_{ϕ}^{RL} is represented as π_{ϕ}^{RL} .
- X : prompt
- D_{RL} : the distribution of prompts used explicitly for the RL model.
- $D_{pretrain}$: the distribution of the training data for the pretrain model.

Using the Reward Model

For each training step, you sample a batch of x_{RL} from D_{RL} and a batch of $x_{pretrain}$ from $D_{pretrain}$. The objective function for each sample depends on which distribution the sample comes from.

1. For each x_{RL} , we use LLM_{ϕ}^{RL} to sample a response: $y \sim LLM_{\phi}^{RL}(x_{RL})$. The objective is computed as follows. Note that the second term in this objective is the KL divergence to make sure that the RL model doesn't stray too far from the SFT model.

$$\text{objective}_1(x_{RL}, y; \phi) = RM(x_{RL}, y) - \beta \log \frac{LLM_{\phi}^{RL}(y|x)}{LLM^{SFT}(y|x)}$$

2. For each $x_{pretrain}$, the objective is computed as follows. Intuitively, this objective is to make sure that the RL model doesn't perform worse on text completion – the task the pretrained model was optimized for.

$$\text{objective}_2(x_{pretrain}; \phi) = \gamma \log LLM_{\phi}^{RL}(x_{pretrain})$$

Using the Reward Model

The final objective is the sum of the expectation of two objectives above. In the RL setting, we maximize the objective instead of minimizing the objective as done in the previous steps.

$$\text{objective}(\phi) = E_{x \sim D_{RL}} E_{y \sim LLM_{\phi}^{RL}(x)} \left[RM(x, y) - \beta \log \frac{LLM_{\phi}^{RL}(y|x)}{LLM^{SFT}(y|x)} \right] + \gamma E_{x \sim D_{pretrain}} \log LLM_{\phi}^{RL}(x)$$

$$\text{objective}_1(x_{RL}, y; \phi) = RM(x_{RL}, y) - \beta \log \frac{LLM_{\phi}^{RL}(y|x)}{LLM^{SFT}(y|x)}$$

$$\text{objective}_2(x_{pretrain}; \phi) = \gamma \log LLM_{\phi}^{RL}(x_{pretrain})$$