



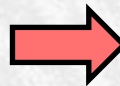
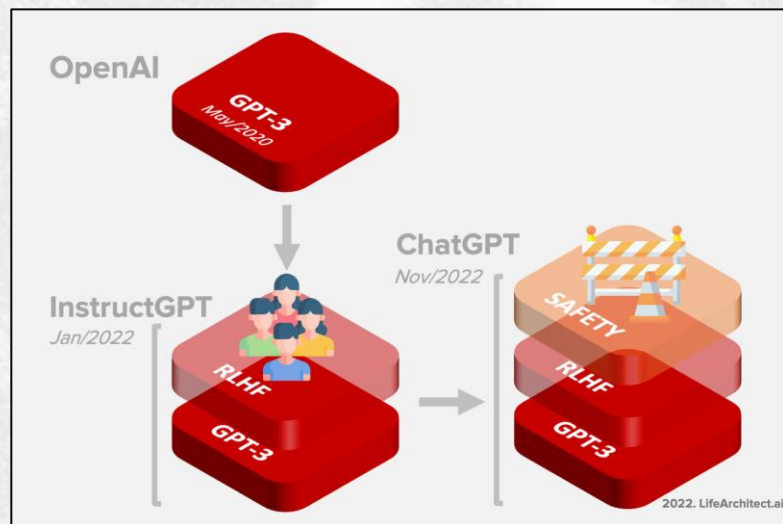
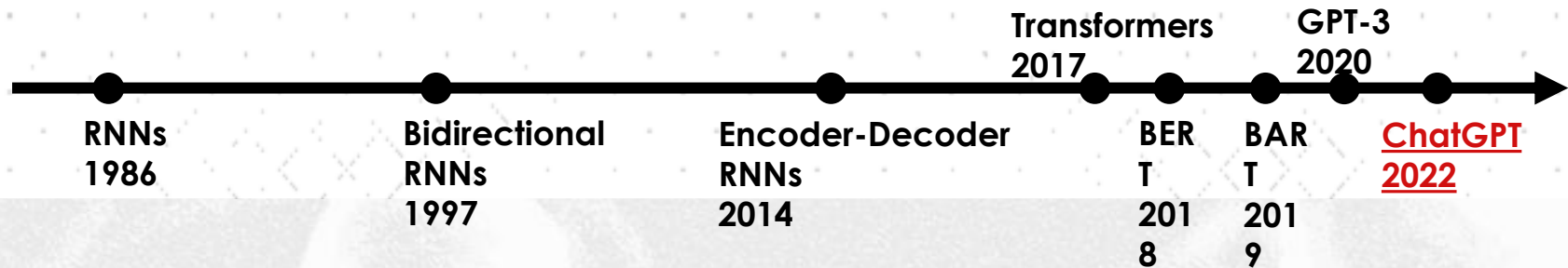
# From 0-shot Learners to Instruction Learning Networks

Roberto Basili, Danilo Croce  
Deep Learning, 2024/2025

# Outline

- From Decoder-Only architectures to ChatGPT
- Chain of Thoughts
- Instruction tuning
  - Instructing LLMs
- Instruction tuning from Human Feedback
  - A reward model for Instructions

# Machine learning paradigms underlying ChatGPT



ChatGPT		
Examples	Capabilities	Limitations
"Explain quantum computing in simple terms" →	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?" →	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?" →	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021

# Inspirations for chatGPT:CoT

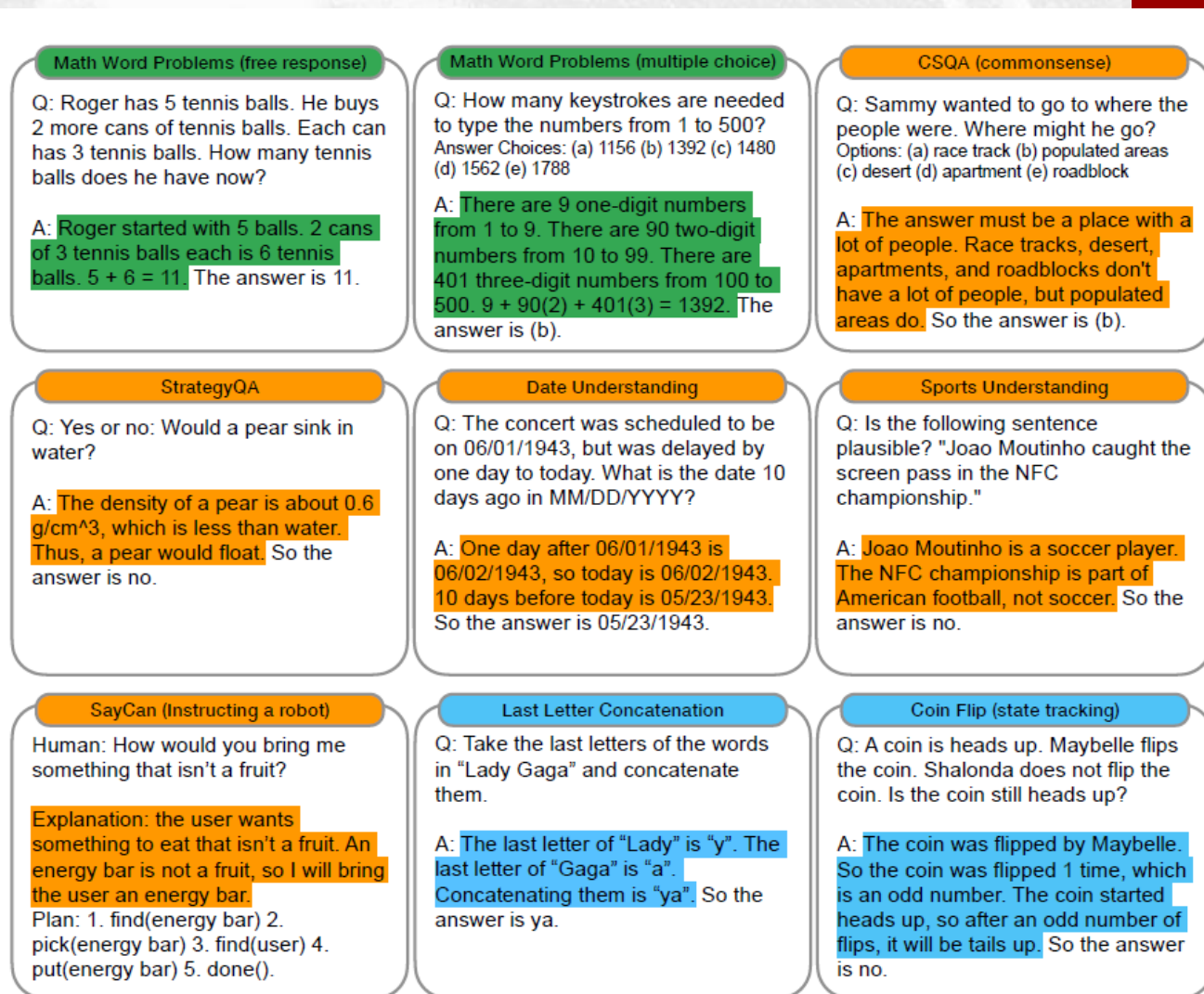


Figure 3: Examples of (input, chain of thought, output) triples for arithmetic, commonsense, and symbolic reasoning benchmarks. Chains of thought are highlighted. Full prompts in Appendix G.



# CoT seminal papers



## Chain of Thought Prompting Elicits Reasoning in Large Language Models

---

### Few-Shot CoT

Jason Wei   Xuezhi Wang   Dale Schuurmans   Maarten Bosma  
Brian Ichter   Fei Xia   Ed H. Chi   Quoc V. Le   Denny Zhou

Google Research, Brain Team  
{jasonwei,dennyzhou}@google.com

## Large Language Models are Zero-Shot Reasoners

---

### Zero-Shot CoT

Takeshi Kojima  
The University of Tokyo  
t.kojima@weblab.t.u-tokyo.ac.jp

Shixiang Shane Gu  
Google Research, Brain Team

Machel Reid  
The University of Tokyo

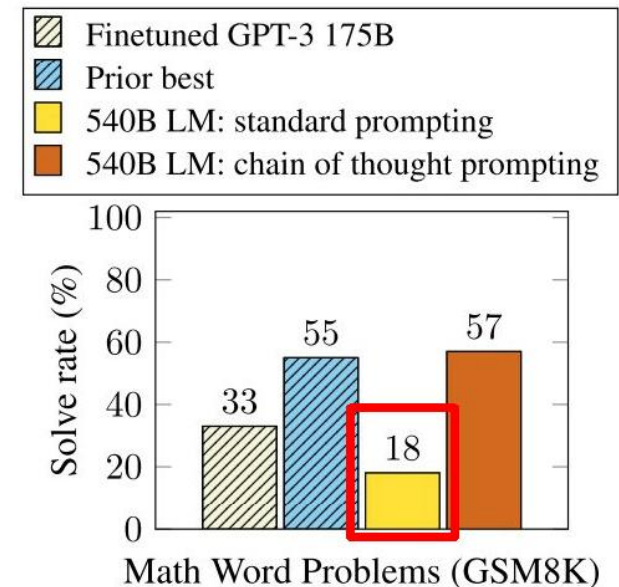
Yutaka Matsuo  
The University of Tokyo

Yusuke Iwasawa  
The University of Tokyo

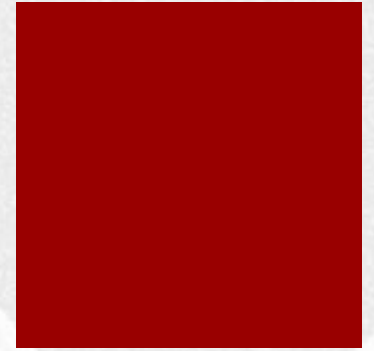
# Inspirations for chatGPT:CoT

- Chain-of-thought relies on two inspiring principles:
  - Solutions to some problems require a rationale that can be explicated in NL, e.g. math problems
  - Prompting is a NL explanation of a task that is competitive wrt to fine-tuning
- IDEA: combine the two by adding **Triggers** to the prompt template
- from (Chain-of-Thought Prompting Elicits Reasoning in LLMs, Wei et al., 2023)

GSM8K (arithmetic):



# CoT prompting



- **CHAIN-OF-THOUGHT PROMPTING** has several attractive properties as an approach for facilitating reasoning in language models.
  1. Chain of thought, in principle, **allows models to decompose multi-step problems into intermediate steps**
  2. A chain of thought **provides an interpretable window into the behavior of the model**, suggesting how it might have arrived at a particular answer and providing opportunities to debug where the reasoning path went wrong
  3. Chain-of-thought reasoning can be used for tasks such as **math word problems, commonsense reasoning, and symbolic manipulation**, and is potentially applicable (at least in principle) to any task that humans can solve via language.
  4. Chain-of-thought reasoning can be readily elicited in sufficiently large off-the-shelf language models simply by including examples of chain of thought sequences into the exemplars of few-shot prompting.



# CoT prompting: examples



Table 8: Example of correct chains of thought produced by the model for the GSM8K dataset.

---

**QUESTION:** James decides to run 3 sprints 3 times a week. He runs 60 meters each sprint. How many total meters does he run a week?

**MODEL ANSWER (CORRECT):** James decides to run 3 sprints 3 times a week. He runs 60 meters each sprint. So he runs 60 meters x 3 sprints x 3 times a week. That is 60 meters x 9. The answer is 540. ✓

---

**QUESTION:** Brandon's iPhone is four times as old as Ben's iPhone. Ben's iPhone is two times older than Suzy's iPhone. If Suzy's iPhone is 1 year old, how old is Brandon's iPhone?

**MODEL ANSWER (CORRECT):** Brandon's iPhone is 4 times as old as Ben's iPhone. Ben's iPhone is 2 times older than Suzy's iPhone. So Brandon's iPhone is  $4 \times 2 = 8$  times older than Suzy's iPhone. Suzy's iPhone is 1 year old. So Brandon's iPhone is  $8 \times 1 = 8$  years old. The answer is 8. ✓

---

**QUESTION:** Jean has 30 lollipops. Jean eats 2 of the lollipops. With the remaining lollipops, Jean wants to package 2 lollipops in one bag. How many bags can Jean fill?

**MODEL ANSWER (CORRECT):** Jean started with 30 lollipops. She ate 2 of them. So she has 28 lollipops left. She wants to package 2 lollipops in one bag. So she can package  $28 / 2 = 14$  bags. The answer is 14. ✓

---



Table 20: Few-shot exemplars for full chain of thought prompt for math word problems. This set of exemplars was used for all math word problem datasets except AQuA.

---

**PROMPT FOR MATH WORD PROBLEMS**

**Q:** There are 15 trees in the grove. Grove workers will plant trees in the grove today. After they are done, there will be 21 trees. How many trees did the grove workers plant today?

**A:** There are 15 trees originally. Then there were 21 trees after some more were planted. So there must have been  $21 - 15 = 6$ . The answer is 6.

**Q:** If there are 3 cars in the parking lot and 2 more cars arrive, how many cars are in the parking lot?

**A:** There are originally 3 cars. 2 more cars arrive.  $3 + 2 = 5$ . The answer is 5.

**Q:** Leah had 32 chocolates and her sister had 42. If they ate 35, how many pieces do they have left in total?

**A:** Originally, Leah had 32 chocolates. Her sister had 42. So in total they had  $32 + 42 = 74$ . After eating 35, they had  $74 - 35 = 39$ . The answer is 39.

**Q:** Jason had 20 lollipops. He gave Denny some lollipops. Now Jason has 12 lollipops. How many lollipops did Jason give to Denny?

**A:** Jason started with 20 lollipops. Then he had 12 after giving some to Denny. So he gave Denny  $20 - 12 = 8$ . The answer is 8.

**Q:** Shawn has five toys. For Christmas, he got two toys each from his mom and dad. How many toys does he have now?

**A:** Shawn started with 5 toys. If he got 2 toys each from his mom and dad, then that is 4 more toys.  $5 + 4 = 9$ . The answer is 9.

**Q:** There were nine computers in the server room. Five more computers were installed each day, from monday to thursday. How many computers are now in the server room?

**A:** There were originally 9 computers. For each of 4 days, 5 more computers were added. So  $5 * 4 = 20$  computers were added.  $9 + 20$  is 29. The answer is 29.

**Q:** Michael had 58 golf balls. On tuesday, he lost 23 golf balls. On wednesday, he lost 2 more. How many golf balls did he have at the end of wednesday?

**A:** Michael started with 58 golf balls. After losing 23 on tuesday, he had  $58 - 23 = 35$ . After losing 2 more, he had  $35 - 2 = 33$  golf balls. The answer is 33.

**Q:** Olivia has \$23. She bought five bagels for \$3 each. How much money does she have left?

**A:** Olivia had 23 dollars. 5 bagels for 3 dollars each will be  $5 \times 3 = 15$  dollars. So she has  $23 - 15$  dollars left.  $23 - 15$  is 8. The answer is 8.

---

# CoT: performances

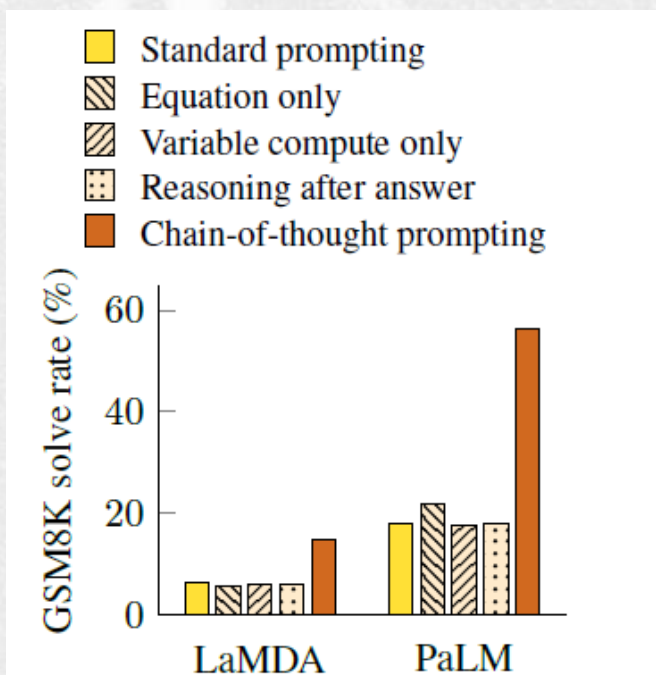


Figure 5: Ablation study for different variations of prompting using LaMDA 137B and PaLM 540B. Results for other datasets are given in Appendix Table 6 and Table 7.

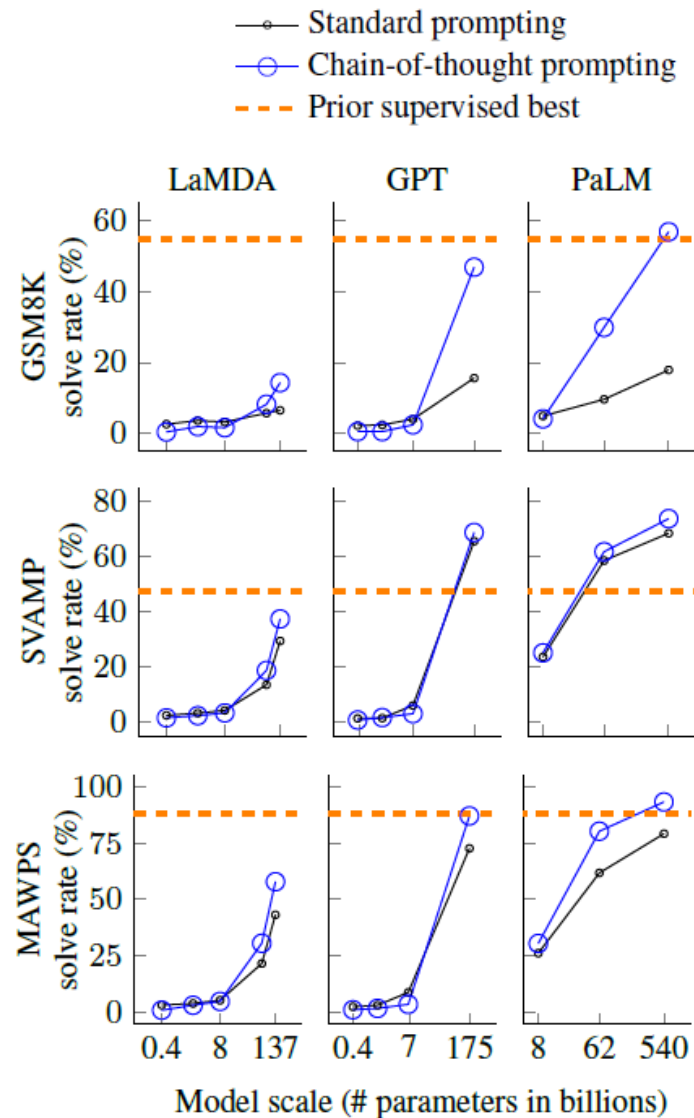
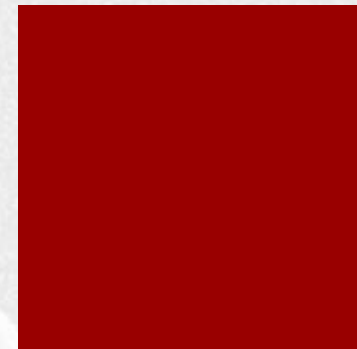


Figure 4: Chain-of-thought prompting enables large language models to solve challenging math problems. Notably, chain-of-thought reasoning is an emergent ability of increasing model scale. Prior best numbers are from Cobbe et al. (2021) for GSM8K, Jie et al. (2022) for SVAMP, and Lan et al. (2021) for MAWPS.

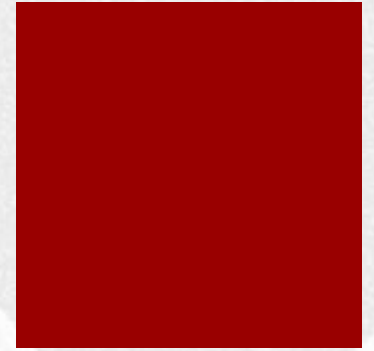
# Limitations of GPT-3



- Large language models often express unintended behaviours such as making up facts, generating biased or toxic text, or simply not following user instructions. This is because the language modeling objective is **misaligned**.
- The idea: aligning language models by training them to act in accordance with the user's intention (Leike et al., 2018).
  - explicit intentions such as following instructions
  - implicit intentions such as staying truthful, and not being biased, toxic, or otherwise harmful.
- Overall Objective: language models should be helpful (they should help the user solve their task), honest (they shouldn't fabricate information or mislead the user), and harmless (they should not cause physical, psychological, or social harm to people or the environment).



# Addressing alignment

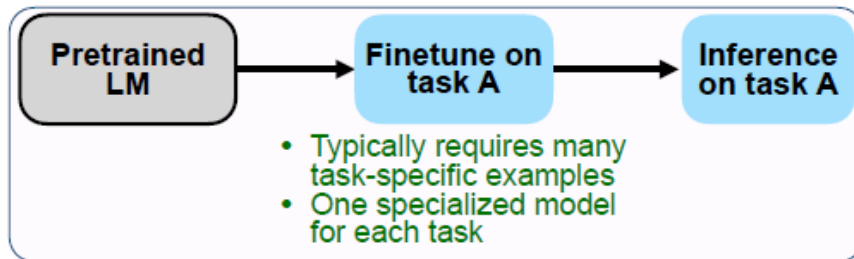


- **FLAN models** (Finetuned Language Models are Zero shot Learners, Wei et al, 2022)
  1. Aggregate Datasets (62): Collect wide variety of public datasets
  2. Instruction Templates: Manually write 10 templates / dataset that captures task
  3. Fine-tune: Use the instruction templates and datasets to fine-tune model
- Instruction tuning from **Human Feedback**

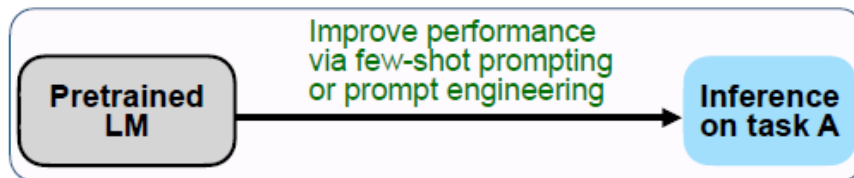
# FLAN (Wei et al., 2021)

- This paper explores a simple method for improving the zero-shot learning abilities of language models. We show that instruction tuning—finetuning language models on a collection of datasets described via instructions—substantially improves zeroshot performance on unseen tasks.

## (A) Pretrain–finetune (BERT, T5)



## (B) Prompting (GPT-3)



## (C) Instruction tuning (FLAN)

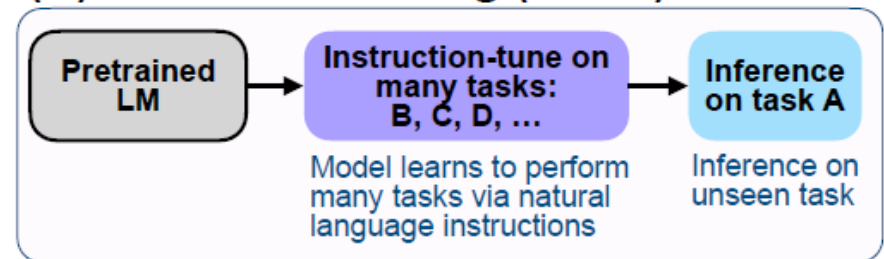


Figure 2: Comparing instruction tuning with pretrain–finetune and prompting.

# FLAN: dataset and templates

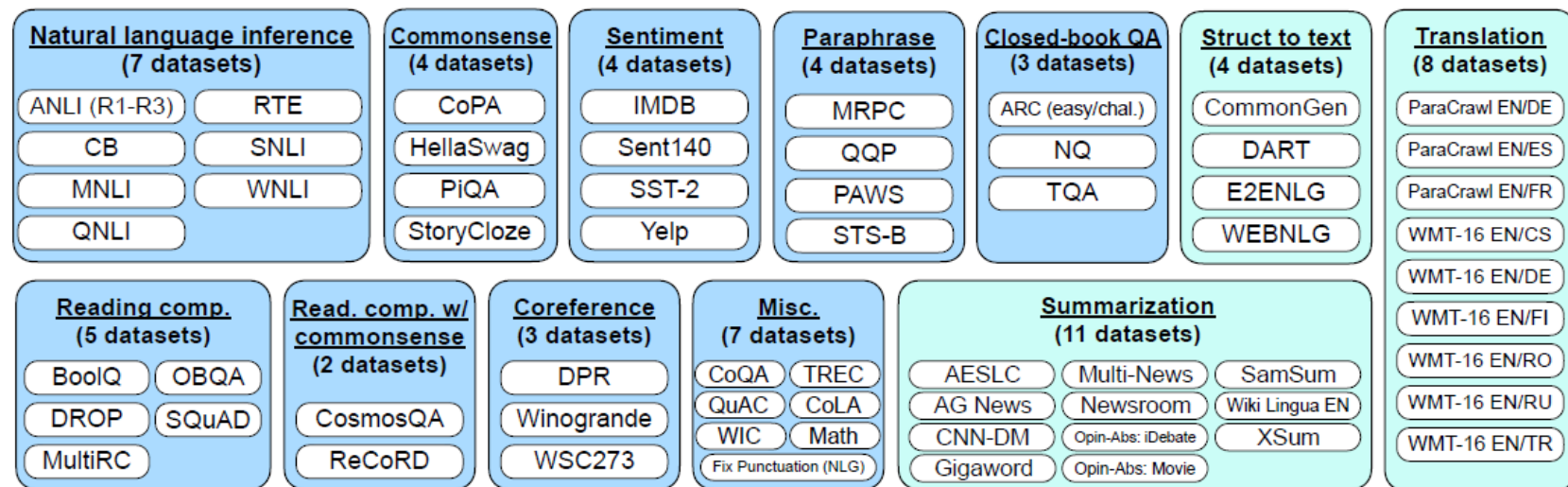


Figure 3: Datasets and task clusters used in this paper (NLU tasks in blue; NLG tasks in teal).

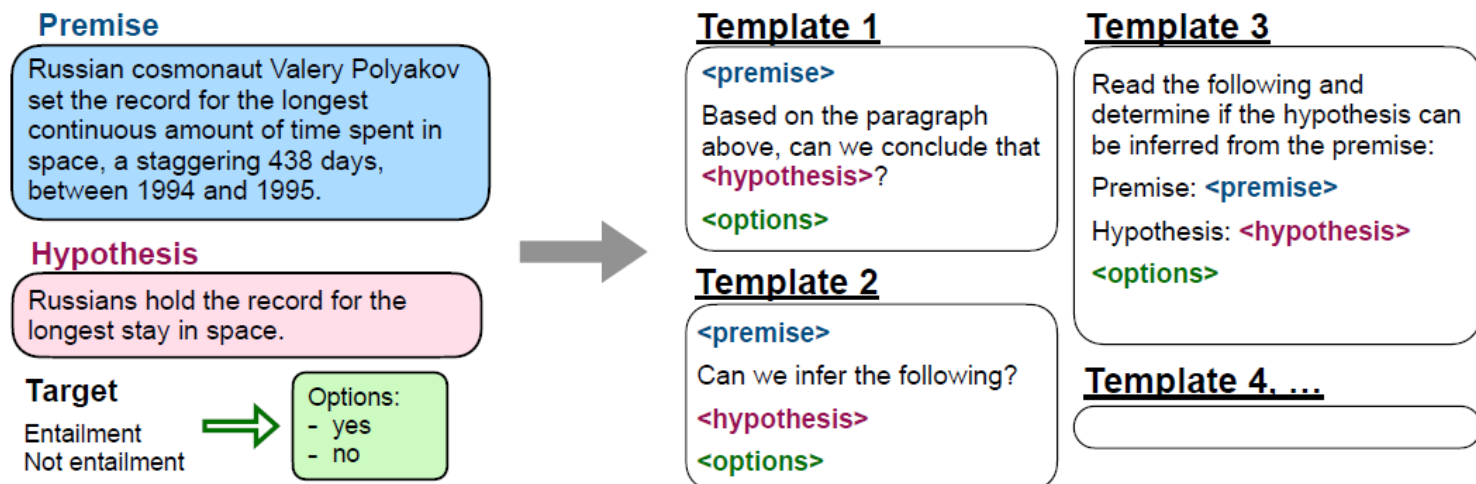


Figure 4: Multiple instruction templates describing a natural language inference task.



# FLAN: fine-tuning & results

- FLAN makes use of a decoder only architecture

**Input:**

John moved to mountain view. He spent friday night in San Francisco. He got a COVID vaccine. He voted in the state senate election.

Write a story that follows the above plot.

**FLAN output (cherrypicked):**

John arrived at the city with his vaccine, feeling prepared enough for the state senate election of San Francisco. He had voted for the first time last night and was feeling good about it. He felt more prepared than ever to win, and win big, because he was getting a COVID vaccine. He knew that his opponent had not gotten a COVID vaccine, and he knew he had already defeated him before the election began.

**Input:**

Write a sad story about carrots.

**FLAN output:**

The sun had barely risen when the carrots were taken from their beds and brought to market. They'd had no time to gather their strength. They smelled of earth, their greens still wet from the dew. At first no one was interested in buying them. But one man finally reached out and picked up one of the carrots. "I'll take this one!" he cried out. "It's for my mother! I'll pay you well for it." The man took the carrot to the cashier, eager to get it home to his mother. He was careful to wrap it carefully so that it wouldn't get bruised or dirty. He hurried home, and when he saw his mother, he gave her the carrot. "Is this one of the ones I grew?" she asked, her voice wavering with emotion. "Yes, mother," he replied. "It is, dear." The man left the carrot with his mother and hurried back to market. His mother was very sick, and he was afraid she wouldn't last much longer. He needed to buy all the carrots in market so that none would be left when she was gone.

**Finetune on many tasks ("instruction-tuning")**

<u>Input (Commonsense Reasoning)</u>	<u>Input (Translation)</u>
Here is a goal: Get a cool sleep on summer days. How would you accomplish this goal? OPTIONS: -Keep stack of pillow cases in fridge. -Keep stack of pillow cases in oven.	Translate this sentence to Spanish: The new office building was built in less than three months.
<u>Target</u> keep stack of pillow cases in fridge	<u>Target</u> El nuevo edificio de oficinas se construyó en tres meses.
Sentiment analysis tasks	
Coreference resolution tasks	
...	

Figure 21: Open-ended generation tasks by FLAN. The carrot story was from sampling sixteen outputs with a minimum length of 150 and choosing the highest probability output.

# Types of Prompts

## *T5 prompt:*

cb hypothesis: At my age you will probably have learnt one lesson.  
premise: It's not certain how many lessons you'll learn by your thirties.

## *GPT-3 prompt:*

At my age you will probably have learnt one lesson.  
question: It's not certain how many lessons you'll learn by your thirties. true, false, or neither? answer:

## *FLAN prompt:*

Premise: At my age you will probably have learnt one lesson.  
Hypothesis: It's not certain how many lessons you'll learn by your thirties.  
Does the premise entail the hypothesis?



# Instruction tuning from human feedback



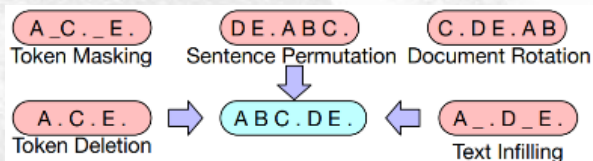
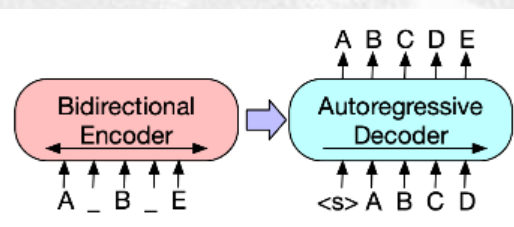
# InstructGPT



- **Step 1:** Collect demonstration data, and train a supervised policy. Labelers provide demonstrations of the desired behavior on the input prompt distribution. Then, fine-tuning of a pretrained GPT-3 model on this data using supervised learning is carried out.
- **Step 2:** Collect comparison data, and train a reward model. A dataset of comparisons between model outputs is collected: labelers indicate which output they prefer for a given input. A reward model to predict the human-preferred output is then trained.
- **Step 3:** Optimize a policy against the reward model using PPO. We use the output of the RM as a scalar reward. We fine-tune the supervised policy to optimize this reward using the proximal policy optimization (PPO) algorithm (Schulman et al., 2017).

# At the heart of ChatGPT (from BART to ChatGPT)

## BART Training-steps



## ChatGPT Training-steps

### Step 1

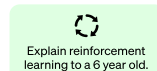
Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.

**human**

A labeler demonstrates the desired output behavior.

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



### Step 2

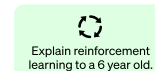
Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.

**human**

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.

**InstructGPT**  
The policy generates an output.

The reward model calculates a reward for the output.

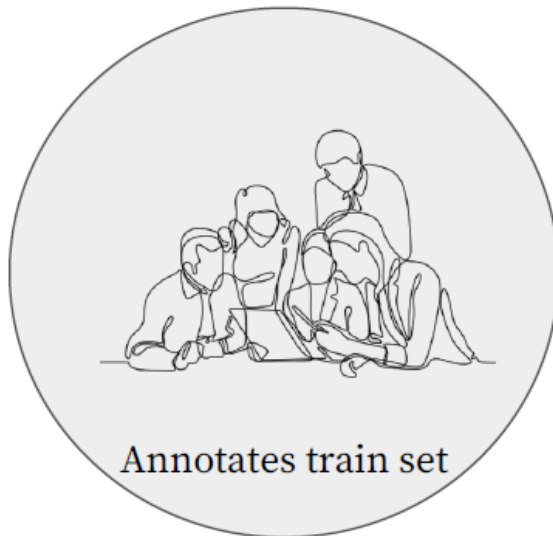
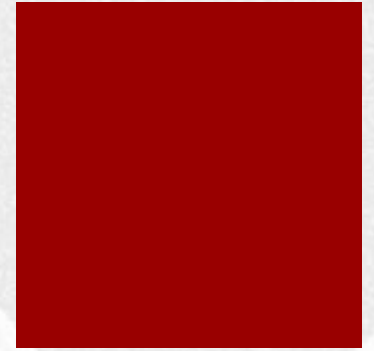
The reward is used to update the policy using PPO.



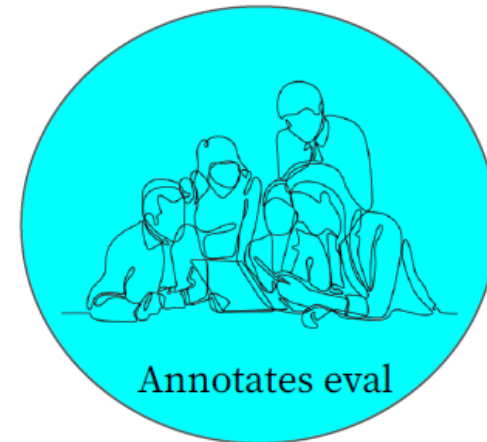
**Fine tune text-davinci-003 to get InstructGPT**

**The Environment**

# Instruct GPT: Human Annotators



- 40 Annotators from Upwork/ScaleAI
- Screened/Onboarded/Diverse etc etc etc



- Different annotators from Upwork/ScaleAI
- Not screened, to better mirror real-world





Step 1

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

A prompt is  
sampled from our  
prompt dataset.



Step 2

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

Step 3

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

Use-case	Prompt
Brainstorming	List five ideas for how to regain enthusiasm for my career
Generation	Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home.
Rewrite	This is the summary of a Broadway play: """" { summary } """"  This is the outline of the commercial for that play: """"

Use-case	(%)
Generation	45.6%
Open QA	12.4%
Brainstorming	11.2%
Chat	8.4%
Rewrite	6.6%
Summarization	4.2%
Classification	3.5%
Other	3.5%
Closed QA	2.6%
Extract	1.9%

Number of Prompts		
SFT Data		
split	source	size
train	labeler	11,295
train	customer	1,430
valid	labeler	1,550
valid	customer	103



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.



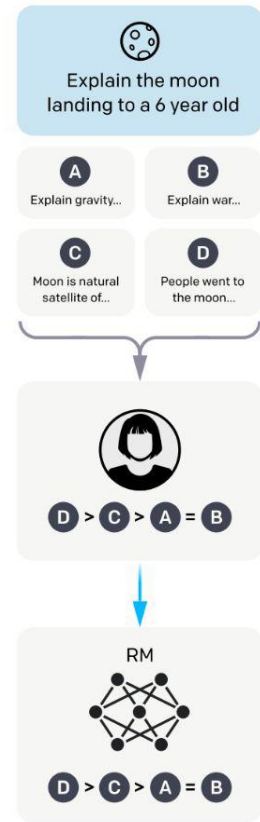
Step 2

**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

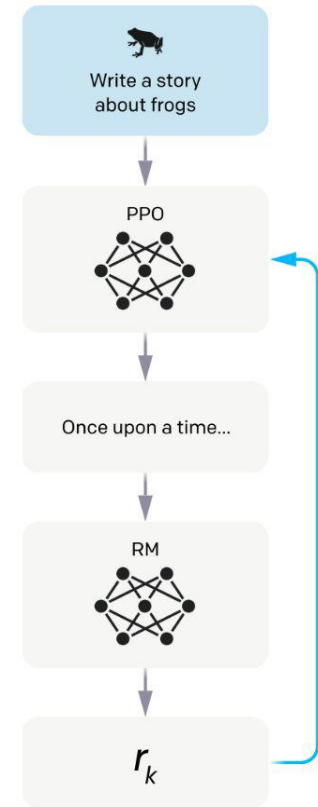
**Optimize a policy against the reward model using reinforcement learning.**

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.

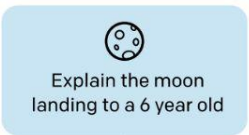




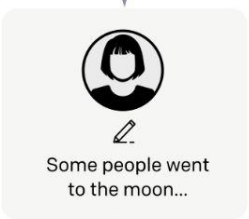
Step 1

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

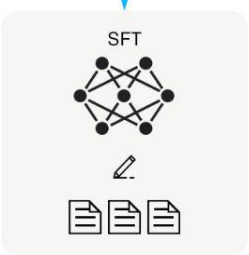
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



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



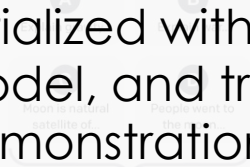
Step 2

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

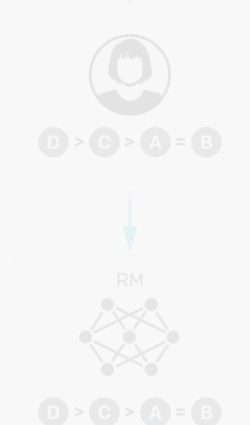
A prompt and several model outputs are sampled



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



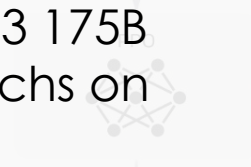
Step 3

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

A new prompt is sampled from the dataset.



The model generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

- Finetune the model, call this model SFT Model
- Initialized with pretrained GPT-3 175B model, and trained for 16 Epochs on demonstration data





Step 1  
Collect demonstration data,  
and train a supervised policy.

Step 2  
Collect comparison data,  
and train a reward model.

Step 3  
Optimize a policy against  
the reward model using  
reinforcement learning.

A prompt and  
several model  
outputs are  
sampled.

🧠  
Explain the moon  
landing to a 6 year old

A  
Explain gravity...

B  
Explain wat...

C  
Moon is natural  
satellite of...

D  
People went to  
the moon...

The outputs are sampled from the SFT model

Number of Prompts		
RM Data		
split	source	size
train	labeler	6,623
train	customer	26,584
valid	labeler	3,488
valid	customer	14,399

This data is used  
to fine-tune the policy  
with supervised learning.

the policy  
using PPO.



Step 1

Collect demonstration data, and train a supervised policy.

Step 2

Collect comparison data, and train a reward model.

Step 3

Optimize a policy against the reward model using reinforcement learning.

A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.

To increase data collection throughput, each user is given  $K = 4$  to 9 outputs to rank for each prompt

Ranking outputs

To be ranked

B A team of researchers from Yale University and University of California, Davis studied the vocalization patterns of several different types of parrots. They found that parrots like to mimic human speech, and can produce a wide range of sounds, such as whistles, squawks, and other types of vocalizations...

C Parrots have been found to have the ability to understand numbers. Researchers have found that parrots can understand numbers up to six. In a series of experiments, the parrots were able to identify the amount of food items under a number of cups...

Rank 1 (best)

A A research group in the United States has found that parrots can imitate human speech with ease, and some of them can even do so in the same way as humans. This group studied the sounds that parrots make in their natural habitats and found that they use their tongues and beaks in ways that are strikingly...

Rank 2

Rank 3

E Scientists have found that green-winged parrots can tell the difference between two noises that are the same except for the order in which they are heard. This is important because green-winged parrots are known to imitate sounds. This research shows that they are able to understand the difference between sounds.

D > C > A = B

used to update the policy using PPO.

$k$



### Step 1

Collect demonstration data,  
and train a supervised policy.

A prompt and  
several model  
outputs are  
sampled.



### Step 2

Collect comparison data,  
and train a reward model.

$r_\theta$ : The reward model we are trying to optimize  
 $x$ : the prompt  $y_w$ : the better completion  $y_l$ : the worse completion

$$\text{loss}(\theta) = -\frac{1}{\binom{K}{2}} E_{(x, y_w, y_l) \sim D} [\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))]$$

### Small but important detail:

- Each prompt has K completions -> K choose 2 pairs to compare
- If  $\forall$  batch we sample uniform over *every* pair (from any prompt):
  - Each completion can appear in K - 1 gradient updates
  - This can lead to overfitting
- **Solution:** sample the prompt, and then put all K choose 2 pairs from the prompt into the same batch
  - Corollary: computationally more efficient, since this only requires K forward passes through  $r_\theta$  for each prompt
- This is why there is the  $-1/(\text{K choose } 2)$  normalization in loss

D > C > A = B

used to update  
the policy  
using PPO.

k





Step 1  
Collect demonstration data,  
and train a supervised policy.

Step 2  
Collect comparison data,  
and train a reward model.

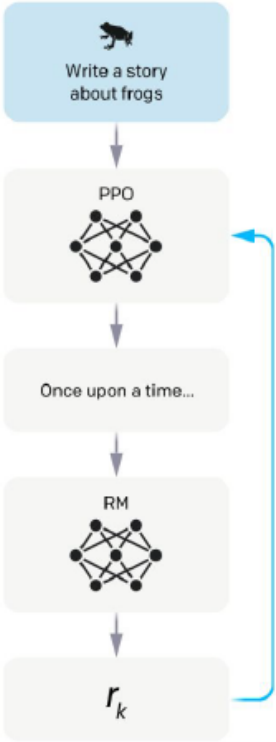
Step 3  
Optimize a policy against  
the reward model using  
reinforcement learning.

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.



Use RM to update the SFT model from step 1. Call model **PPO**

Number of Prompts		
PPO Data		
split	source	size
train	customer	31,144
valid	customer	16,185

using PPO.



Step 1

Collect demonstration data,  
and train a supervised policy.

Step 2

Collect comparison data,  
and train a reward model.

Step 3

Optimize a policy against  
the reward model using  
reinforcement learning.

A prompt is



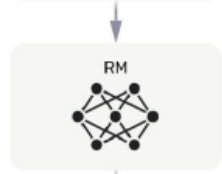
A new prompt  
is sampled from  
the dataset.



The policy  
generates  
an output.



The reward model  
calculates a  
reward for  
the output.



The reward is  
used to update  
the policy  
using PPO.



A prompt and



A new prompt

(Proximal Policy Optimization)

Use RM to update the SFT model from step 1. Call model **PPO**

Two problems:

1. As RLHF is updated, its outputs become very different from what the RM was trained on -> worse reward estimates  
**Solution:** add a KL penalty that makes sure PPO model output does not deviate too far from SFT

using PPO.



Step 1

Collect demonstration data,  
and train a supervised policy.

Step 2

Collect comparison data,  
and train a reward model.

Step 3

Optimize a policy against  
the reward model using  
reinforcement learning.

A prompt is

A prompt and

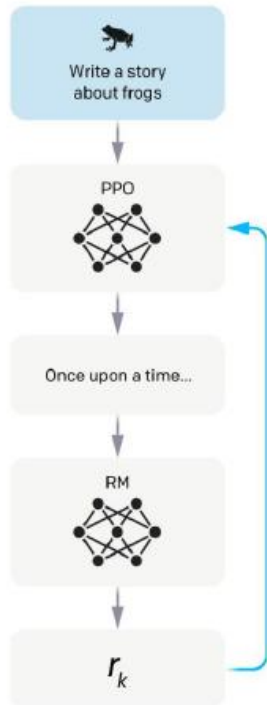
A new prompt

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.



Use RM to update the SFT model from step 1. Call model **PPO**

Two problems:

1. As RLHF is updated, its outputs become very different from what the RM was trained on -> worse reward estimates  
**Solution:** add a KL penalty that makes sure PPO model output does not deviate too far from SFT

2. Just using RL objective leads to performance degradation on many NLP tasks  
**Solution:** Add a auxiliary LM objective on the pretraining data. Call this variant **PPO-ptx**

**(Proximal Policy  
Optimization with  
PreTraining Mixture)**

the policy  
using PPO.





Step 1

Collect demonstration data,  
and train a supervised policy.

Step 2

Collect comparison data,  
and train a reward model.

Step 3

Optimize a policy against  
the reward model using  
reinforcement learning.

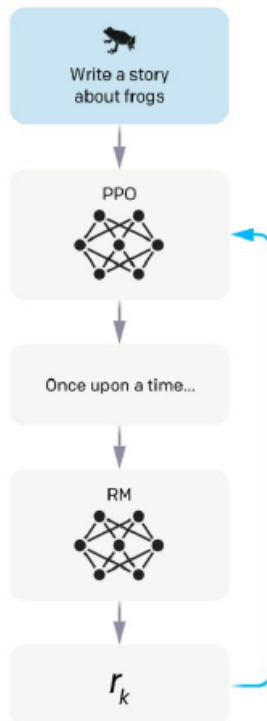
A prompt is

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.



A prompt and

A new prompt

Use RM to update the SFT model from step 1. Call model **PPO**

Two problems:

1. As RLHF is updated, its outputs become very different from what the RM was trained on -> worse reward estimates

**Solution:** add a KL penalty that makes sure PPO model output does not deviate too far from SFT

2. Just using RL objective leads to performance degradation on many NLP tasks

**Solution:** Add a auxiliary LM objective on the pretraining data. Call this variant **PPO-ptx**

$$\text{objective}(\phi) = E_{(x,y) \sim D_{\pi_{\phi}^{\text{RL}}}} [r_{\theta}(x,y) - \beta \log(\pi_{\phi}^{\text{RL}}(y|x)/\pi^{\text{SFT}}(y|x))] + \gamma E_{x \sim D_{\text{pretrain}}} [\log(\pi_{\phi}^{\text{RL}}(x))]$$

# The model

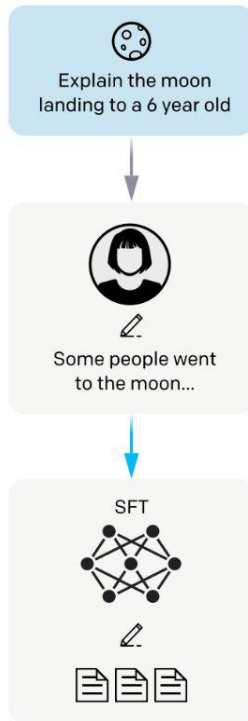
## Step 1

**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.



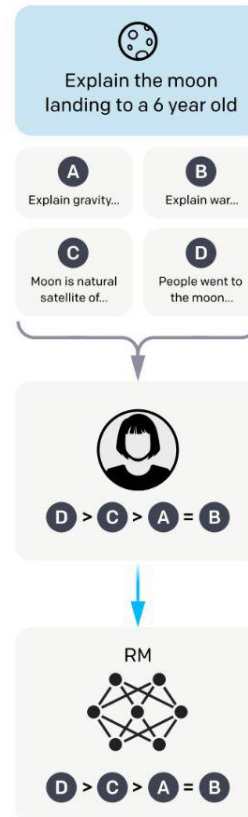
## Step 2

**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

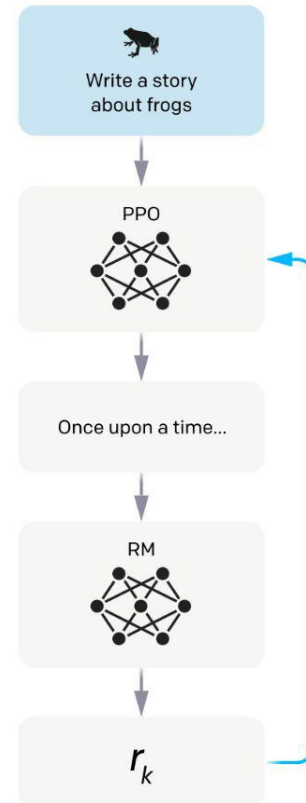
**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

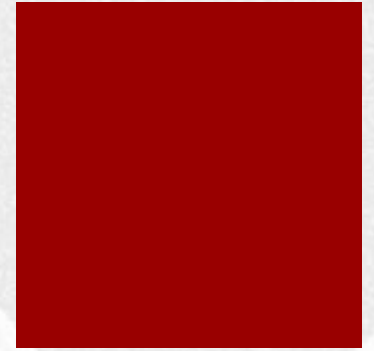
The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



# InstructGPT: model summary



- 1. **SFT: Supervised Fine-Tuning**
  - a. GPT-3 fine-tuned on human demonstrations of prompt completions
- 2. **RM: Reward Model**
  - a. Not actually used to generate anything, but used to train the PPO and PPO-ptx models
- 3. **PPO (Proximal Policy Optimization)**
  - a. SFT model further fine-tuned using RL with the RM providing the reward signal
  - b. A KL-loss is provided to prevent the PPO model from deviating far from SFT
- 4. **PPO-ptx (Proximal Policy Optimization with PreTraining Mixture)**
  - a. Identical to PPO, except with an additional auxiliary LM objective on the pretraining data

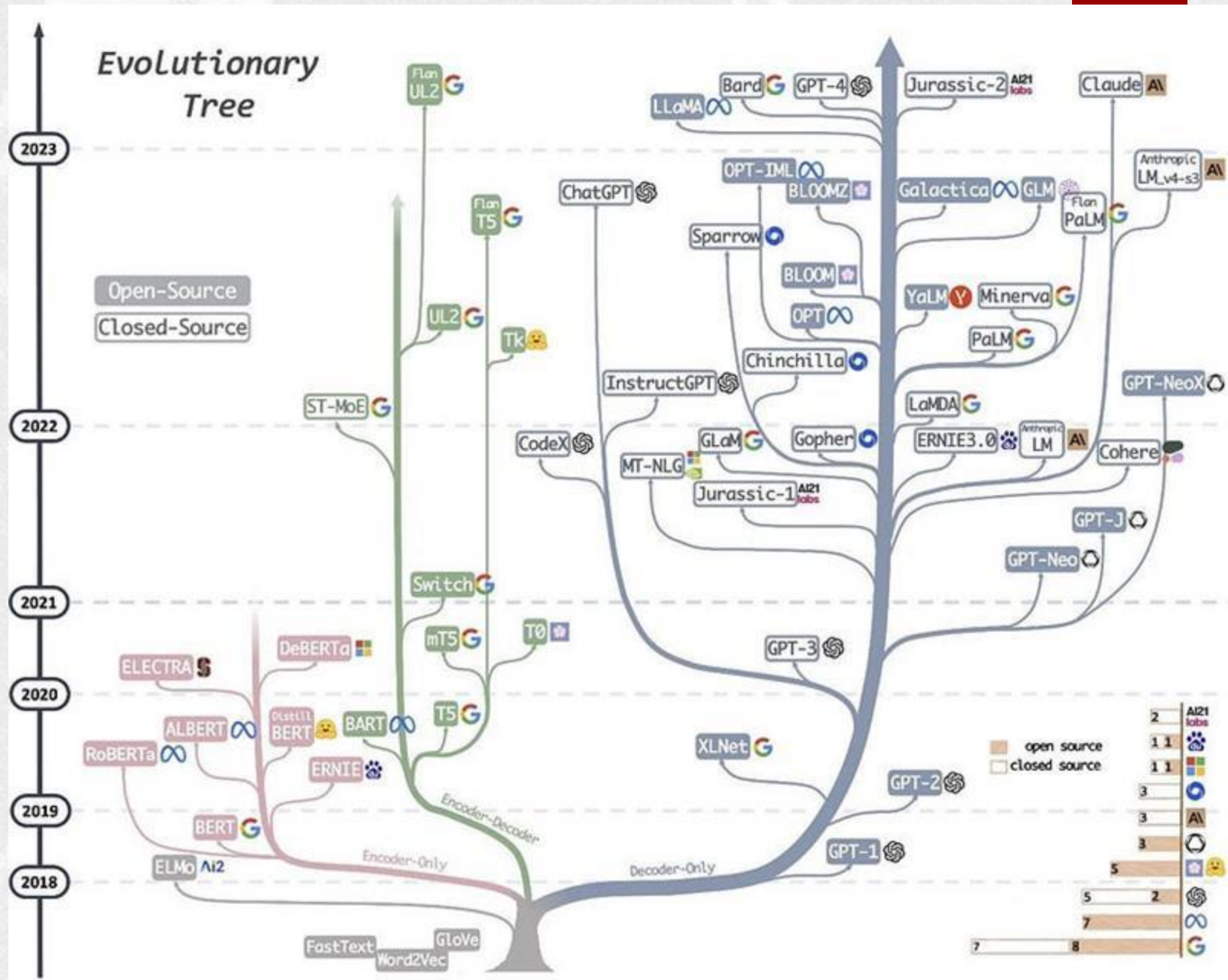


# Instruction tuning papers




- Learning to summarize from human feedback
- Fine-Tuning Language Models from Human Preferences
- MemPrompt: Memory-assisted Prompt Editing with User Feedback
- LaMDA: Language Models for Dialog Application

1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 2679, 2680, 26



# Neural writing



write a biblical verse in the style of the king james bible explaining how to remove a peanut butter sandwich from a VCR 



And it came to pass that a man was troubled by a peanut butter sandwich, for it had been placed within his VCR, and he knew not how to remove it.  

And he cried out to the Lord, saying, "Oh, Lord, how can I remove this sandwich from my VCR, for it is stuck fast and will not budge?"

And the Lord spoke unto him, saying, "Fear not, my child, for I shall guide thy hand and show thee the way. Take thy butter knife, and carefully insert it between the sandwich and the VCR, and gently pry them apart. And with patience and perseverance, the sandwich shall be removed, and thy VCR shall be saved."

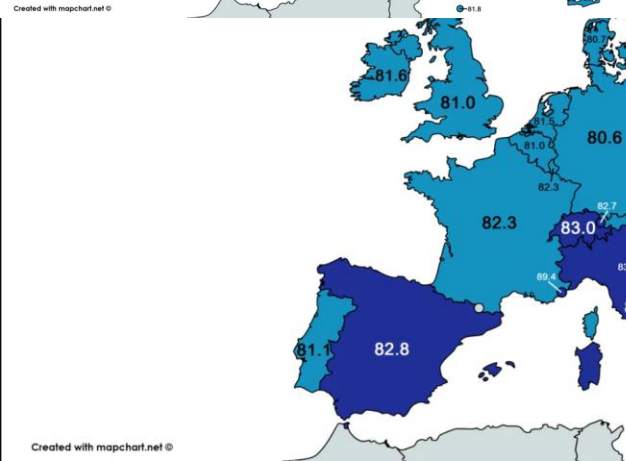
And the man did as the Lord commanded, and lo and behold, the sandwich was removed from the VCR, and the man was saved.

And the Lord said, "Verily I say unto thee, seek not to put thy peanut butter sandwiches in thy VCR, for it is not a suitable place for such things. Rather, keep thy sandwiches in thy refrigerator or on thy plate, where they belong."

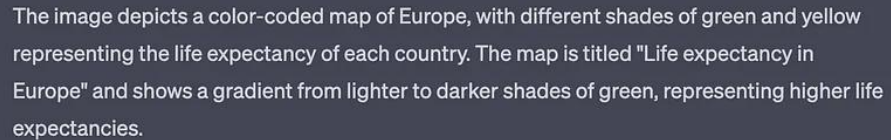
And the man heeded the Lord's words, and from that day forth, he kept his sandwiches in their proper place, and was saved from trouble and woe. Amen.



the 1990s, the number of people in the world who are illiterate has increased from 1.2 billion to 1.5 billion. The number of illiterate people in the world is expected to reach 1.7 billion by the year 2015. The number of illiterate people in the world is expected to reach 1.7 billion by the year 2015. The number of illiterate people in the world is expected to reach 1.7 billion by the year 2015.



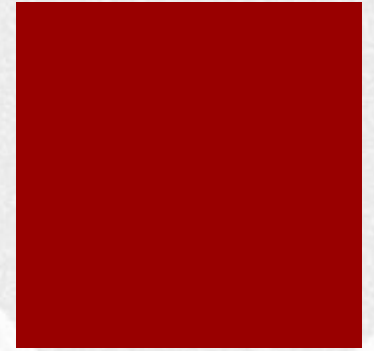
[https://en.wikipedia.org/wiki/List\\_of\\_European\\_countries\\_by\\_life\\_expectancy#/media/File:Life\\_expectancy\\_in\\_Europe.png](https://en.wikipedia.org/wiki/List_of_European_countries_by_life_expectancy#/media/File:Life_expectancy_in_Europe.png)



The countries with the highest life expectancy, such as Switzerland, Spain, Italy, and France, are shown in the darkest shades of green, indicating life expectancies of over 80 years. The countries with lower life expectancies, such as Russia, Belarus, and Ukraine, are depicted in yellow and lighter shades of green, indicating life expectancies of under 75 years.

The map also includes a legend that provides a color-coded scale and corresponding life expectancy ranges, allowing viewers to easily interpret the data. Overall, the image provides a visual representation of the variation in life expectancy across different European countries, highlighting the disparities between countries and regions.

# Towards AI chats bibliography



- [\(Chain-of-Thought Prompting Elicits Reasoning in LLMs, Wei et al., 2023\)](#)
- (FLAN) WEI, Jason, et al. [Finetuned language models are zero-shot learners](#). *arXiv preprint arXiv:2109.01652*, 2021.
- (Ziegler et al., 2021) [Fine-Tuning Language Models from Human Preferences](#)
- Madaan et al., EMNLP 2022, [MemPrompt: Memory-assisted Prompt Editing with User Feedback](#)
- Thoppilan et al, 2022, [LaMDA: Language Models for Dialog Application](#)
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, Paul F. Christiano: [Learning to summarize with human feedback](#). NeurIPS 2022
- [Training Language Models to follow instructions through Human feedback](#), Ouyang et al., 2022