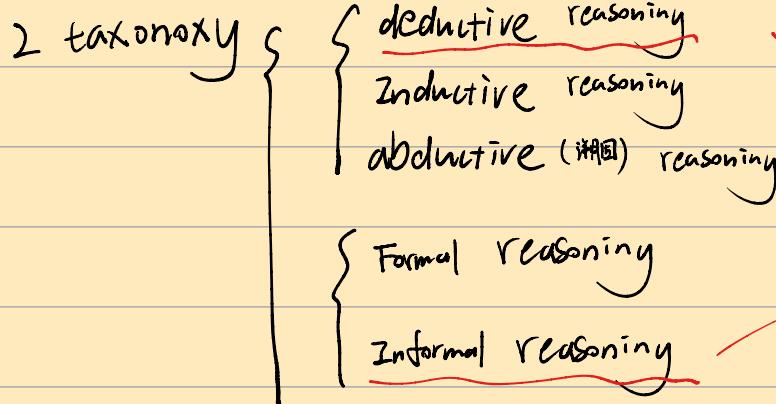


Reasoning in Language Models

Core: CoT... are not new model architecture, but a prompt/generation strategy

①



most we use

in LMs

②

Can LLMs really reason?

Large Language models are **REALLY GOOD** at predicting **plausible continuations of text** (Lecture-9), that respect **constraints in the input** (Lecture 10,11), and align well with **human preferences** (Lecture-10, 11).

Core question: are these "reasoning steps" genuine reasoning or merely a fit of linguistic patterns? - idk

③ Main Approaches to improve reasoning

① prompting-based methods

1. Chain-of-Thought (CoT): explicitly ask the model to

show intermediate steps

Standard Prompting	Chain-of-Thought Prompting
<p>Model Input</p> <p>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?</p> <p>A: The answer is 11.</p> <p>Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?</p> <p>A: The answer is 27. ✗</p>	<p>Model Input</p> <p>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?</p> <p>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.</p> <p>Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?</p> <p>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. ✓</p>

2. zero-shot / few-shot CoT

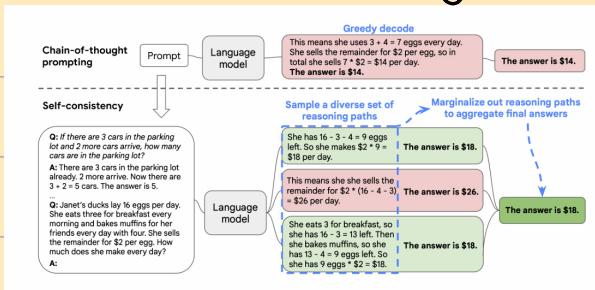
(a) Few-shot	(b) Few-shot-CoT
<p>Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?</p> <p>A: The answer is 11.</p> <p>Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?</p> <p>A:</p> <p>(Output) The answer is 8. ✗</p>	<p>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?</p> <p>A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.</p> <p>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?</p> <p>A:</p> <p>(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✓</p>
<p>(c) Zero-shot</p> <p>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?</p> <p>A: The answer (arabic numerals) is</p> <p>(Output) 8 ✗</p>	<p>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?</p> <p>A: Let's think step by step.</p> <p>(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓</p>

3. CoT with "self-consistency":

Replace greedy decoding with an ensemble of samples

Basic idea: if an answer keeps appearing for multiple rationales, the majority of the rationales agree on, then it's more likely to be correct

That sounds like ensembling, but in practice, self-consistency is doing more than simple ensembling



4. Least-to-most prompting: solve simpler subproblems first

Q: Elsa has 5 apples. Anna has 2 more apples than Elsa. How many apples do they have together?
A: Anna has 2 more apples than Elsa, so Anna has $2 + 5 = 7$ apples. Elsa and Anna have $5 + 7 = 12$ apples together. The answer is 12.

Standard CoT prompting

Q: Elsa has 5 apples. Anna has 2 more apples than Elsa. How many apples do they have together?
A: Let's break down this problem: 1. How many apples does Anna have? 2. How many apples do Elsa and Anna have together?
1. Anna has 2 more apples than Elsa. So Anna has $2 + 5 = 7$ apples.
2. Elsa and Anna have $5 + 7 = 12$ apples together.

Least-to-most prompting

Q: {question}
A: Let's break down this problem:
The answer is:

Accuracy by Steps (GSM8K)	All	2 Steps	3 Steps	4 steps	≥ 5 steps
Least-to-Most	62.39	74.53	68.91	59.73	45.23
Chain-of-Thought	60.87	76.68	67.29	59.39	39.07

Generalizes to more #steps than in-context example!

Prompting method	Accuracy
Zero-Shot	16.38
Standard prompting	17.06 ³
Chain-of-Thought (original)	61.18
Chain-of-Thought (1-shot)	60.88
Least-to-Most (1-shot)	62.39
Chain-of-Thought (best)	68.61³
Least-to-Most (best)	68.01

But with enough prompt engineering, CoT \approx Least-to-Most

2. Distillation of reasoning

In 1), we might think: instead of trying to get really large language models to do reasoning, maybe we want to somehow get this kind of reasoning behavior in a smaller/m

Idea:

{ Use strong models (e.g., GPT-4) to generate reasoning traces



Fine-Tune smaller models on these traces

Representative work : Orca

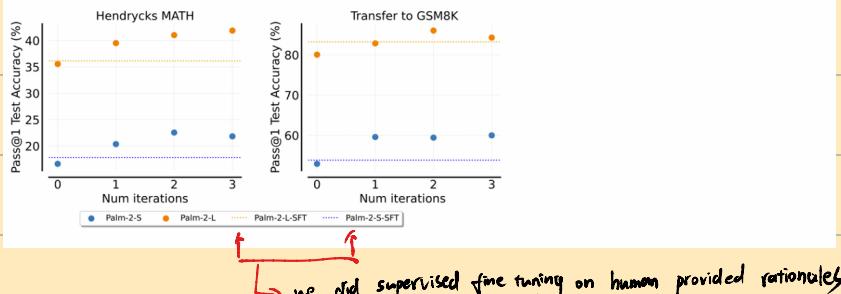
1. Collect a wide variety of instructions from the FLAN-v2 collection.
2. Prompt GPT4 or ChatGPT with these instructions along with a system message
3. Finetune Llama-13b on outputs generated via ChatGPT + GPT4

③ Self-Generated Reasoning Data

But in ② we might think: why not just finetune the big LM on its own rationales

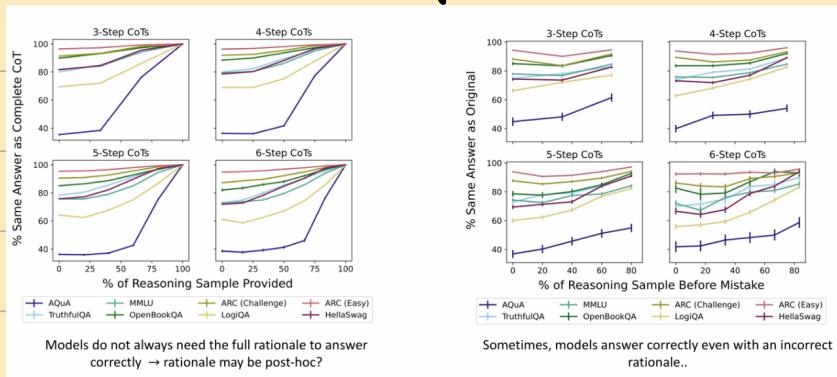
one of the methods: Reinforced self-training (REST)

- REST^{EM} alternates between the following two steps:
1. Generate (E-Step): Given reasoning problem, sample multiple solutions from language model. Filter based on some (problem specific) function [answer correctness for math problems]
 2. Improve (M-Step): Update the language model to maximize probability of filtered solutions, using supervised finetuning



④ Challenges in Reasoning Evaluation

II GT Rationales are often not faithful



⑤ Models may rely on memorization rather than reasoning

↳ So how to distinguish?

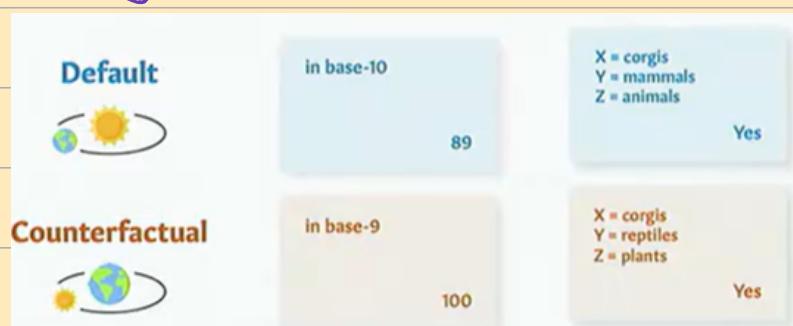
→ Using Counterfactuals (反事实)

What is Counterfactuals?

Expect to not be present that frequently in the training data

or say become slightly unnatural

Eg:



By using it, we found that performance drops significantly under counterfactual tests

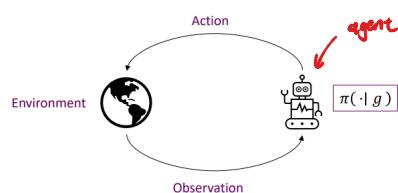


This suggests current reasoning abilities are fragile and non-systematic

Language Model Agents

△ Concept

Some Terminology



} this setting also called
"Instructed-following agent"
"Language conditioned policy"
"digital agent"

agent receives an observation from its environment



{ based on the observation, agent issues an action
along with that, agent receives second variable g

g represents a language instruction ↪

② Instruction following agents [Pre LLMs]

- 1. training semantic parsers
- 2. Inferring plans from instruction trajectory pairs
- 3. Modeling plans directly and then have an execution model that can execute plans
- 4. doing reinforcement learning with reward signal

③ Instruction following agents [in 2024]

main idea: Generative trajectory modeling with causal transformers

A Simple Language Model Agent with ReACT

You are an agent capable of the following actions:
 1. Type X on Y
 2. Move mouse to
 3. Click on X
 4. Type Char x on Y

Your objective is to follow user instructions, by mapping them into a sequence of actions.
 Instruction: {g}

So far, you have taken the following actions and observed the following environment states:

Previous Actions and Observations:
 o1:
 a1:
 o2:
 a2:
 -

After executing these actions, you observe the following HTML state: <HTML state>

Now, think about your next action:
 Thought: [model-pred]

Now, take an action:
 Action: [model-pred]

1. Action space in text
2. Instruction in text
3. Previous observations and actions
4. Provide current observation [as text]

Model generates next action (sequence prediction task), use that action to update environment and repeat!

Mostly, just CoT prompting in a loop

Some popular benchmarks for LM agents

MiniWoB ++ , Web Arena , WebLIXN

④ training data for Languages model Agents

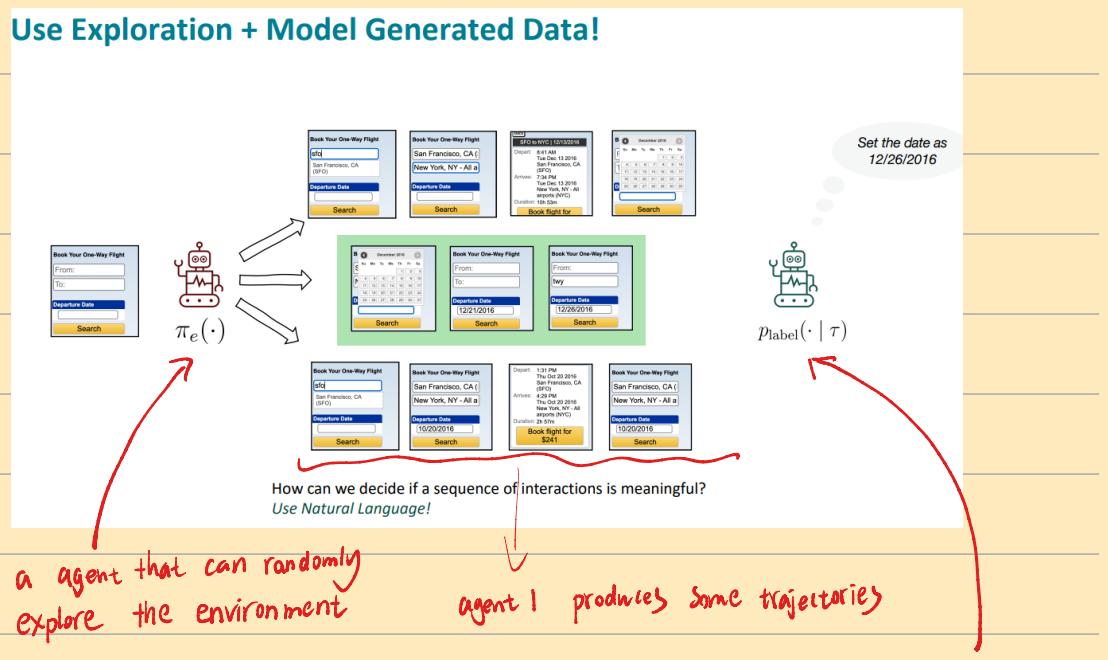
Standard practice : In-context learning with few-shot demonstrations of humans performing following similar instructions

[But this's not scalable / reliable (cause lots of environments, many kinds of interactions possible)]

Is there something better than we can do than just sort of getting humans to provide demonstrations for every new use case?

Using LM to generate rationales and then fine tune on that!
In here, we don't have rationales, but we could produce action trajectories. And then we're going to use that as supervision.

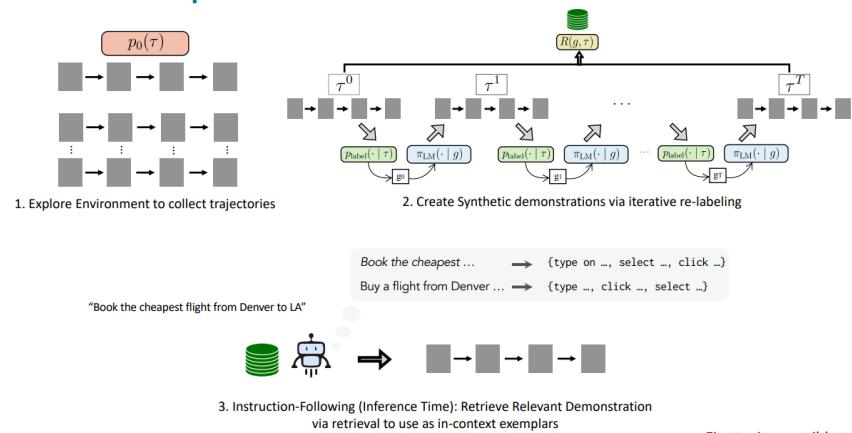
↳ example:



↑ { if the trajectory seems pretty good for completing tasks
we add it to a set of examples
it not

Instead of throwing away this interaction (Cause interactions are pretty noisy)
we invoke the re-labeler to take the trajectory and assign it a new label

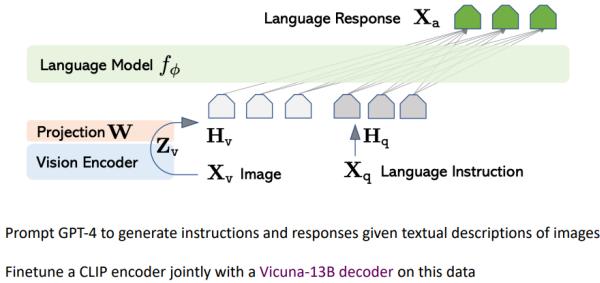
BAGEL (Bootstrapping Agents by Guiding exploration with Languages)



② Multimodality

- So far, we've looked at using text-only language models for agents
- This is intractable for real-world UIs with very long HTML
- Can we instead operate directly over pixel space?

LLaVA



Pix2Struct

