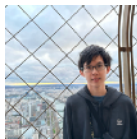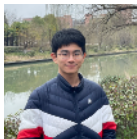# LLM Reasoners

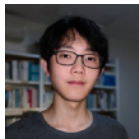A library for advanced reasoning with LLMs

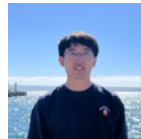Shibo Hao*  Yi Gu*  Haotian Luo*  Tianyang Liu  Xiyan Shan  Xinyuan Wang  Shuhua Xie  Haodi Ma  Adithya Samavedhi  Zhen Wang  Zhiting Hu

# Large Language Model Reasoning

**Google** The Keyword

**In this story**

Our new benchmark approach to MMLU enables Gemini to use its reasoning capabilities to think more carefully before answering difficult questions, leading to significant improvements over just using its first impression.

OpenAI            Research

GPT-4 surpasses
advar

**Anthropic** ✔
@AnthropicAI

Today, we're announcing Claude 3, our next generation of AI models.

The three state-of-the-art models—Claude 3 Opus, Claude 3 Sonnet, and Claude 3 Haiku—set new industry benchmarks across reasoning, math, coding, multilingual understanding, and vision.

# Large Language Model Reasoning

# Outline

- Reasoning with LLMs:

  Algorithms, Evaluation, Analysis



**LLM Reasoners**

# Outline

- Reasoning with LLMs:

  **Algorithms**, Evaluation, Analysis



**LLM Reasoners**

# Large Language Model Reasoning

**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.

# Large Language Model **Step-by-step** Reasoning

## 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. ❌

## Chain-of-Thought 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: 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. ✅

7

# Large Language Model **Step-by-step** Reasoning

# Large Language Model **Step-by-step** Reasoning

## Becomes the **Default Choice**

🤗 **Open LLM Leaderboard**

EleutherAI / lm-evaluation-harness `Public` 🔔 Notifications 🍴 Fork 1.1k ⭐ Star 4.3k ▾

<> Code ⊙ Issues 99 ⇅ Pull requests 31 ⊙ Actions ⊞ Projects 1 ⊘ Security ⌁ Insights

|  | Claude 3 Opus | Claude 3 Sonnet | Claude 3 Haiku |
|---|---|---|---|
| Undergraduate level knowledge *MMLU* | 86.8% 5 shot | 79.0% 5-shot | 75.2% 5-shot |
| Graduate level reasoning *GPQA, Diamond* | 50.4% 0-shot CoT | 40.4% 0-shot CoT | 33.3% 0-shot CoT |
| Grade school math *GSM8K* | 95.0% 0-shot CoT | 92.3% 0-shot CoT | 88.9% 0-shot CoT |
| Math problem-solving *MATH* | 60.1% 0-shot CoT | 43.1% 0-shot CoT | 38.9% 0-shot CoT |

## Can we design algorithms to generate better reasoning chains with LLMs?

# Reasoning with Language Model 🤖
# is Planning with World Model 🌍

Shibo Hao*  Yi Gu*  Haodi Ma  Joshua Hong  Zhen Wang  Daisy Wang  Zhiting Hu

# Chain-of-thoughts vs Human reasoning

**Blocksworld:** How to move the blocks to the goal state?



**B: Human Reasoning**
- Internal **world model** to track **states**
- **Explore** alternative reasoning paths
- **Assess outcomes** by looking ahead

## A: Chain-of-Thoughts Prompting (CoT) with LLM
- Autoregressive decoding

**System 1**

**Invalid Action!**
The yellow ⬜ block is still under the red 🟥 one.

1. Pick up the orange ⬜ block.
2. Stack it on the blue ⬜ block.
3. Pick up the yellow ⬜ block.  ✗
4. Stack it on the orange ⬜ block.
5. Pick up the red 🟥 block.
6. Put it on the table.

**System 2**

Pick up orange          Pick up blue

Stack on blue   ......   Stack on orange

**Better than**

On the planning abilities of large language models (a critical investigation with a proposed benchmark) [Valmeekam et al, 2023]
Chain-of-thought prompting elicits reasoning in large language models [Wei et al., 2022]
Mental models: Towards a cognitive science of language, inference, and consciousness [Johnson-Laird, 1983]
From System 1 Deep Learning to System 2 Deep Learning [Bengio, 2019]

# Reasoning-via-Planning (RAP 🎵)

**Human Reasoning**
- Internal **world model** to track **states**
- **Explore** alternative reasoning paths
- **Assess outcomes** by looking ahead

**How to enable LLMs to reason close to humans?**

**Reasoning-via-Planning: RAP 🎵**

- Repurpose LLM as **world model**
- Principled **planning** algorithm
- **Rewards** to estimate outcomes

**Reasoning-via-Planning (RAP)**



Pick up orange          Pick up blue

Stack on blue    ......          Stack on orange

......

# Planning Algorithm



**Monte Carlo Tree Search (MCTS):**

Iteratively build reasoning tree

1. Selection
2. Expansion
3. Simulation
4. Back-propagation

**Balanced exploration and exploitation**

Goal:

Pick up orange

Pick up blue $\qquad a_1$

$\qquad s_0$

$\qquad s_1$

Stack on blue

......

Stack on orange $\qquad a_2$

$\qquad s_2$

Pick up orange

Pick up red

$\qquad a_3$

$\qquad s_3$

$\qquad s_T$

13

# Rewards in RAP

**Reward** design is **flexible**

In Blocksworld:

- **Likelihood** of actions
- **Task-heuristic** (# of subgoals)

Other possible rewards:

- **Self-evaluation** by LLM (e.g. useful? correct?)
- **Confidence** of next state
- ......

# RAP on Plan Generation (Blocksworld)

# RAP on Plan Generation (Blocksworld)

# RAP on Mathematical Reasoning (GSM8k)

{}

Q1: How many pages did Julie read today?

Q1: How many pages has she read?

**Action: a sub-question for an unknown variable**

Q1: How many pages ⋯ today?
A1: 12×2=24

Q1: How many pages has ⋯?
A1: 12×2=24

Q2: How many pages should she read tomorrow?

Q1: How many pages has she read till now?

Q1: How many pages ⋯ today?
A1: 12×2=24
Q2: How many ⋯ tomorrow?
A2: (120-24)/2=48

Q1: How many pages ⋯ today?
A1: 12×2=24
Q2: How many ⋯ till now?
A2: 12+24=36

**State: A set of known variables**

Q1: How many pages ⋯ today?
A1: 12×2=24
......
Qn: How many pages should she read?
An: 84/2=42    (Answer: 42)

Question:
Julie is reading a 120-page book.
Yesterday … 12 pages
Today … twice as many pages as yesterday
Tomorrow … half of the remaining pages
How many pages should she read?

# RAP on Mathematical Reasoning (GSM8k)

# RAP on Logical Reasoning (PrOntoQA)



**Action: selecting a rule from the rule set**

Fae is a feline

5) Each feline is a carnivores

3) Every cat is a feline

Fae is a carnivore

Fae is a cat

1) Carnivores are carnivorous

4) Carnivores are mammals

Fae is carnivorous

Fae is a mammal

**State: The fact we are focusing on**

Fae is not unicellular

(The hypothesis is false)

**Rules**:
(1) Carnivores are carnivorous
(2) Animals are not unicellular
(3) Every cat is a feline
(4)…
**Fact**: Fae is a feline
**Hypothesis**: Fae is unicellular?

Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. [Saparov and He, 2022]

# RAP on Logical Reasoning (PrOntoQA)



Action: selecting a rule from the rule set

RAP outperforms CoT much in proof accuracy

State: The fact we are focusing on

Rules:
(1) Carnivores are carnivorous
(2) Animals are not unicellular
(3) Every cat is a feline
(4)…
Fact: Fae is a feline
Hypothesis: Fae is unicellular?

Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. [Saparov and He, 2022]

# Large Language Model Step-by-step Reasoning

**Chain-of-Thought Prompting Elicits Reasoning in Large Language Models**

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

**Solving Math Word Problems via Cooperative Reasoning induced Language Models**

Xinyu Zhu◇*    Junjie Wang♠*    Lin Zhang♡    Yuxiang Zh
Ruyi Gan♡    Jiaxing Zhang♡    Yujiu Yang◇†
◇Tsinghua University    ♠Waseda University
♡International Digital Economy Academy

zhuxy21@mails.tsinghua.edu.cn    yang.yujiu@sz.tsinghua.edu.cn
wjj1020181822@toki.waseda.jp    joel0495@asagi.waseda.jp
{zhanglin, ganruyi, zhangjiaxing}@idea.edu.cn

**Tree of Thoughts: Deliberate Problem Solving with Large Language Models**

Shunyu Yao    Dian Yu    Jeffrey Zhao    Izhak Shafran
Princeton University    Google DeepMind    Google DeepMind    Google DeepMind

Thomas L. Griffiths    Yuan Cao    Karthik Narasimhan
Princeton University    Google DeepMind    Princeton University

**Reasoning with Language Model is Planning with World Model**

Shibo Hao*♣    Yi Gu*♣    Haodi Ma◇    Joshua Jiahua Hong♣
Zhen Wang♣♠    Daisy Zhe Wang◇    Zhiting Hu♣

♣UC San Diego, ◇University of Florida
♠Mohamed bin Zayed University of Artificial Intelligence
{s5hao, yig025, jjhong, zhw085, zhh019}@ucsd.edu
{ma.haodi, daisyw}@ufl.edu

**GRACE: Discriminator-Guided Chain-of-Thought Reasoning**

Muhammad Khalifa*, Lajanugen Logeswaran
Honglak Lee*†, Lu Wang*
University of Michigan*, LG AI Research†, University

**AlphaZero-Like Tree-Search can Guide Large Language Model Decoding and Training**

Xidong Feng ×1    Ziyu Wan *2    Muning Wen 2    Stephen Marcus McAleer 3
Ying Wen 2    Weinan Zhang 2    Jun Wang 1

**TOOLCHAIN*: EFFICIENT ACTION SPACE NAVIGATIO IN LARGE LANGUAGE MODELS WITH A* SEARCH**

Yuchen Zhuang1*, Xiang Chen2, Tong Yu2, Saayan Mitra2
Victor Bursztyn2, Ryan A. Rossi2, Somdeb Sarkhel2, Chao Zhang1
Georgia Institute of Technology1 Adobe Research2
yczhuang@gatech.edu, {xiangche, tyu, smitra}@adobe.com
{soaresbu, ryrossi, sarkhel}@adobe.com, chaozhang@gatech.edu

# Large Language Model Step-by-step Reasoning

## Analysis on current reasoning algorithms?

**Solving Math Word Problems via Cooperative Reasoning induced Language Models**

Xinyu Zhu◇*   Junjie Wang♠*   Lin Zhang♡   Yuxiang Zhan
Ruyi Gan♡   Jiaxing Zhang♡   Yujiu Yang◇†
◇Tsinghua University   ♠Waseda University
♡International Digital Economy Academy
zhuxy21@mails.tsinghua.edu.cn   yang.yujiu@sz.tsinghua.edu.cn
wjj1020181822@toki.waseda.jp   joel0495@asagi.waseda.jp
{zhanglin, ganruyi, zhangjiaxing}@idea.edu.cn

**Chain-of-Thought Prompting Elicits Reason in Large Language Models**

**Jason Wei**   **Xuezhi Wang**   **Dale Schuurmans**   **Maarten Bosn**
**Brian Ichter**   **Fei Xia**   **Ed H. Chi**   **Quoc V. Le**   **Denny Zhou**
Google Research, Brain Team
{jasonwei,dennyzhou}@google.com

**Tree of Thoughts: Deliberate Problem Solving with Large Language Models**

**Shunyu Yao**   **Dian Yu**   **Jeffrey Zhao**   **Izhak Shafran**
Princeton University   Google DeepMind   Google DeepMind   Google DeepMind

**Thomas L. Griffiths**   **Yuan Cao**   **Karthik Narasimhan**
Princeton University   Google DeepMind   Princeton University

**Reasoning with Language Model is Planning with World Model**

**Shibo Hao**♣ **Yi Gu**♣ **Haodi Ma**◇ **Joshua Jiahua Hong**♣
**Zhen Wang**♣♠ **Daisy Zhe Wang**◇ **Zhiting Hu**♣
♣UC San Diego, ◇University of Florida
♠Mohamed bin Zayed University of Artificial Intelligence
{s5hao, yig025, jjhong, zhw085, zhh019}@ucsd.edu
{ma.haodi, daisyw}@ufl.edu

**GRACE: Discriminator-Guided Chain-of-Th**

**Muhammad Khalifa**\*, **Lajanugen Logeswaran**
**Honglak Lee**\*†, **Lu Wang**\*
University of Michigan\*, LG AI Research†, University

**TOOLCHAIN\*: EFFICIENT ACTION SP
IN LARGE LANGUAGE MODELS WITH**

Yuchen Zhuang[1]\*, Xiang Chen[2], Tong Yu[2], Saayan Mitra[2]
Victor Bursztyn[2], Ryan A. Rossi[2], Somdeb Sarkhel[2], Chao Zh
Georgia Institute of Technology[1] Adobe Research[2]
yczhuang@gatech.edu, {xiangche, tyu, smitra}@adobe.com
{soaresbu, ryrossi, sarkhel}@adobe.com, chaozhang@gatech.edu

**AlphaZero-Like Tree-Search can Guide Large Language Model Decoding and Training**

## Technical Connection?
## Which design really matters?

# Large Language Model Step-by-step Reasoning

Difficulties in implementation…

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\ldots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

Chain-of-Thoughts
( 🔗 CoT)

Tree-of-Thoughts
( 🌲 ToT)

Reasoning-via-
Planning
( 🎵 RAP)

…

 World
Model

 Search
Algorithm

 Reward
Function

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\dots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

Chain-of-Thoughts
( 🔗 CoT)

Tree-of-Thoughts
( 🌲 ToT)

Reasoning-via-
Planning
(🎵🎶 RAP)

…


World Model


Search Algorithm


Reward Function

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\ldots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

World Model    $s_t = (a_0, \ldots, a_t)$

Chain-of-Thoughts
( CoT)

Search Algorithm    greedy decoding

Reward Function    $P_\theta(a_t \mid s_t)$

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\ldots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$



**Task**:

Manipulates the blocks such that:
- Orange block on the blue block;
- Yellow block is on the orange block.

Chain-of-Thoughts
(CoT)

World Model $\qquad s_t = (a_0, \ldots, a_t)$

Search Algorithm $\qquad$ greedy decoding

Reward Function $\qquad P_\theta(a_t \mid s_t)$

Pick up the orange block $\quad a_0$

Stack the orange block on the blue block $\quad a_1$

$s_0$

$s_1$

$s_2$

$s_T$

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\ldots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

**Task**:



Manipulates the blocks such that:
- Orange block on the blue block;
- Yellow block is on the orange block.

Chain-of-Thoughts
( CoT)

World Model    $s_t = (a_0, \ldots, a_t)$

Search Algorithm    greedy decoding

Reward Function    $P_\theta(a_t \mid s_t)$

Pick up the orange block    $a_0$

Stack the orange block on the blue block    $a_1$

$s_0$
$s_1$
$s_2$
$s_T$

( Pick up the orange block,    Stack the orange block on the blue block )

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\dots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

Tree-of-Thoughts
(🌲ToT)

World
Model

$s_t = (a_0, \dots, a_t)$

Search
Algorithm

BFS / DFS

Reward
Function

$P_\theta(\text{"good"} \mid s_t, a_t)$

**Task**:

Manipulates the blocks such that:
- Orange block on the blue block;
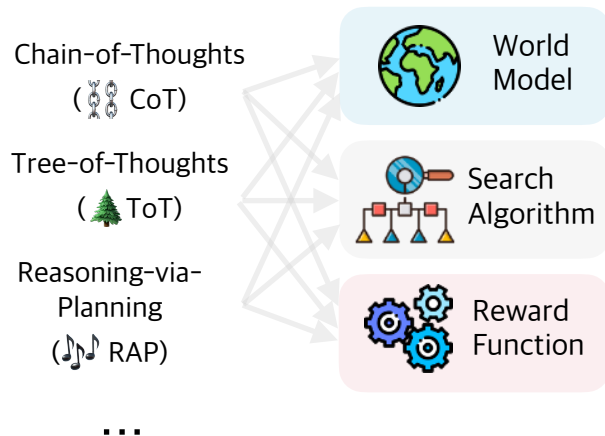- Yellow block is on the orange block.

# A Formulation of Step-by-step Reasoning

$$\text{argmax}_{(a_0,\dots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

**Task**:

Manipulates the blocks such that:
- Orange block on the blue block;
- Yellow block is on the orange block.

Tree-of-Thoughts
( 🌲 ToT)

World Model

Search Algorithm — BFS / DFS

Reward Function — $P_\theta(\text{"good"} \mid s_t, a_t)$

$$s_t = (a_0, \dots, a_t)$$

Pick up the orange block $a_0$

Stack the orange block on the blue block $a_1$

$s_0$

$a_0$

$s_1$

$a_1$

$s_2$

$s_1$

$s_2$

$a_2$

$s_T$

$s_T$

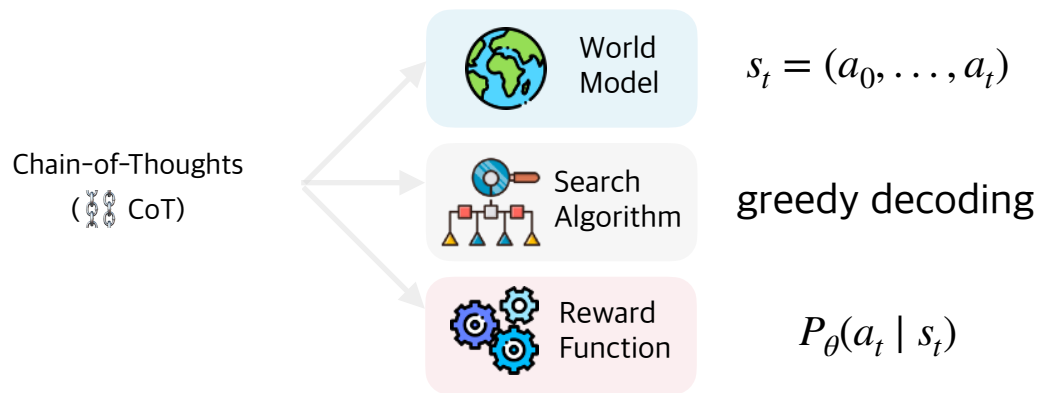( Pick up the orange block,     Stack the orange block on the blue block )

# A Formulation of Step-by-step Reasoning



$$\text{argmax}_{(a_0,\ldots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

**Task**:

Manipulates the blocks such that:
- Orange block on the blue block;
- Yellow block is on the orange block.

Reasoning-via-Planning
(🎵🎵 RAP)

World Model — $s_t \sim P_\theta(s_t \mid s_{t-1}, a_{t-1})$

Search Algorithm — MCTS

Reward Function

$$P_\theta(\text{"good"} \mid s_t, a_t)$$
$$P_\theta(a_t \mid s_t)$$

Other task-specific reward

Pick up the orange block $\quad a_0$

Stack the orange block on the blue block

$s_0$

$a_0$

$s_1$

$a_1$

$s_1$

$a_1$

$s_2$

$s_2$

$a_2$

$s_T$

$s_T$

# LLM Reasoners: A library for complex reasoning with LLMs

$$\text{argmax}_{(a_0,...,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

Search Configuration
get_actions(state)  # get action space
reward(state, action)  # reward function

World Model
init_state() # the initial state
step(state, action) # next state prediction
is_terminal(state) # determine terminal state

# LLM Reasoners: A library for complex reasoning with LLMs

$$\text{argmax}_{(a_0,...,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

**Search Algorithm**
▷ BFS
▷ MCTS
▷ …

**Search Configuration**
↳ get_actions(state)  # get action space
↳ reward(state, action)  # reward function

**World Model**
↳ init_state() # the initial state
↳ step(state, action) # next state prediction
↳ is_terminal(state) # determine terminal state

# LLM Reasoners: A library for complex reasoning with LLMs

$$\text{argmax}_{(a_0,...,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

**Search Algorithm**
▷ BFS
▷ MCTS
▷ ...

**Search Configuration**
↳ get_actions(state)  # get action space
↳ reward(state, action)  # reward function

**LLM API**
▷ Huggingface
▷ OpenAI
▷ ...

**World Model**
↳ init_state() # the initial state
↳ step(state, action) # next state prediction
↳ is_terminal(state) # determine terminal state

# LLM Reasoners: A library for complex reasoning with LLMs

$$\text{argmax}_{(a_0,\ldots,a_T)} \sum_{t=0}^{T} r(s_t, a_t), \quad s_t \sim P(s_t \mid s_{t-1}, a_t)$$

**Search Algorithm**
- BFS
- MCTS
- …

**Search Configuration**
- get_actions(state)  # get action space
- reward(state, action)  # reward function

**Benchmark**
- GSM8k
- StrategyQA
- …

**LLM API**
- Huggingface
- OpenAI
- …

**World Model**
- init_state() # the initial state
- step(state, action) # next state prediction
- is_terminal(state) # determine terminal state

**Visualization**
- Web-based interactive visualization

## Search Algorithm

- BFS
- MCTS
- …

## Search Configuration

↳ get_actions(state)  # get action space
↳ reward(state, action)  # reward function

## LLM API

- Exllama
- OpenAI
- …

## World Model

↳ init_state() # the initial state
↳ step(state, action) # next state prediction
↳ is_terminal(state) # determine terminal state

```python
from reasoners import SearchConfig, WorldModel
from reasoners.algorithm import MCTS
from reasoners.lm import Llama2Model
from reasoners import Reasoner

class MyWorldModel(WorldModel):
    def step(self, state, action):
        return self.llm.generate(self.next_state_prompt.format(state, action))
    ...

class MyConfig(SearchConfig):
    def reward(self, state, action):
        self_eval = self.llm.generate(self.eval_prompt.format(state, action))
        return self_eval
    ...

reasoner = Reasoner(
    world_model=MyWorldModel(), search_config=MyConfig(), search_algo= MCTS()
)
```

**Task**:

Manipulates the blocks such that:
- Orange block on the blue block;
- Yellow block is on the orange block.

$s_1$

| $a_1$ | $r(s_1, a_1)$ |
|---|---|
| **Pick up orange** | **0.6** |
| Pick up blue | 0.3 |
| Pick up yellow | 0.2 |

$s_2$

# LLM Reasoners: A library for complex reasoning with LLMs



**Visualization**

▷ Web-based interactive visualization

# Outline

- Reasoning with LLMs:

  Algorithms, **Evaluation**, Analysis



**LLM Reasoners**

# Large Language Model Step-by-step Reasoning

## How to evaluate step-by-step reasoning?

**39%** of the **correct** answers were derived from **incorrect** reasoning chains!
(Llama-2 70B on a random subset of StrategyQA)

Question
Did Aristotle use a laptop?

Reasoning Chain
$a_0$: Aristotle was born 384 BCE.
$a_1$: The laptop was invented in the 21st century
$a_2$: Since it Is invented after his birth. The answer is no.

Answer-based Evaluation

**Can we directly evaluate reasoning chains?**

# Reasoning Chain Evaluation

Previous methods:
- Compare to human-written reference (Celikyilmaz et al., 2020)
- Train a model to evaluate (Golovneva et al., 2022)
- Prompt GPT-4 to evaluate (He et al., 2023)

Evaluation of text generation: A survey [Celikyilmaz et al, 2020]
Roscoe: A suite of metrics for scoring step-by-step reasoning [Golovneva et al., 2022]
SocREval: Large Language Models with the Socratic Method for Reference-Free Reasoning Evaluation [He et al., 2023]

# Reasoning Chain Evaluation

Previous methods:
- Compare to human-written reference (Celikyilmaz et al., 2020)
- Train a model to evaluate (Golovneva et al., 2022)     Training data
- Prompt GPT-4 to evaluate (He et al., 2023, Tyen et al., 2023)

Prompt engineering

- Need additional human efforts

Evaluation of text generation: A survey [Celikyilmaz et al, 2020]
Roscoe: A suite of metrics for scoring step-by-step reasoning [Golovneva et al., 2022]
SocREval: Large Language Models with the Socratic Method for Reference-Free Reasoning Evaluation [He et al., 2023]

# Reasoning Chain Evaluation

Previous methods:
- Compare to human-written reference (Celikyilmaz et al., 2020)
- Train a model to evaluate (Golovneva et al., 2022)   Training data
- Prompt GPT-4 to evaluate (He et al., 2023, Tyen et al., 2023)

Prompt engineering

**LLMs cannot *find* reasoning errors, but can *correct* them!**

**Gladys Tyen**[*1], **Hassan Mansoor**[2], **Victor Cărbune**[2], **Peter Chen**[†2], **Tony Mak**[†2]
[1]University of Cambridge, Dept. of Computer Science & Technology, ALTA Institute
[2]Google Research
gladys.tyen@cl.cam.ac.uk
{hassan,chenfeif,tonymak,vcarbune}@google.com

- Need additional human efforts
- Overall unsatisfactory evaluation accuracy

Evaluation of text generation: A survey [Celikyilmaz et al, 2020]
Roscoe: A suite of metrics for scoring step-by-step reasoning [Golovneva et al., 2022]
SocREval: Large Language Models with the Socratic Method for Reference-Free Reasoning Evaluation [He et al., 2023]

# Reasoning Chain Evaluation (RICE)

**Q: Can one ignite helium?**

1. Helium is an odorless and tasteless gas.
2. Helium has no color.
3. So the answer is no.

**Is this answer correct?**

The given answer is partially correct…

# Reasoning Chain Evaluation (RICE)



Logic?

Accuracy?

Relevance?

...

**Following the criteria, evaluate the reasoning chain step by step.**

**Q: Can one ignite helium?**

1. Helium is an odorless and tasteless gas.
2. Helium has no color.
3. So the answer is no.

- Accuracy: ···, correct.

- Relevance: The information in the first two steps are irrelevant to the question.

- Logic: The final step cannot be inferred from the previous steps.

So, the reasoning is **INCORRECT**.

# Reasoning Chain Evaluation (RICE)

**Reference reasoning chains (Training set)**

**Q1:** ··· The answer is no
··· The answer is yes ❌

**Q2:** ··· The answer is no
··· The answer is no ✅

**Reasoning chains generated by the student LLM**

**I: Collecting wrong reasoning chains**

**Q: Did Aristotle use laptop?**

**Student**

- Aristotle lived from 384-322 BCE.
- Laptop was invented in 1980.
- Since it's invented after his death, the answer is no.

**Reference**

- Aristotle is a modern philosopher
- Laptop was invented in 1980.
- So, Aristotle should have used laptops, the answer is yes.

**What mistakes did the student make?**

The student made a <u>factual mistake</u> that Aristotle is a modern philosopher. He actually lived from 384–322 BCE.

**II: Detecting the errors**

For question 1, the student made a factual mistake that Aristotle is a modern philosopher···

For question ···, the student listed an irrelevant fact that ···

**To summarize, a good reasoning chain should** ···

- **Accuracy**: Be free of factual errors
- **Relevance**: ···
- **Logic**: ···

**III: Summarizing the evaluation criteria**

**Criterion List Construction**

# Reasoning Chain Evaluation (RICE)

**Q: Can one ignite helium?**

1. Helium is an odorless and tasteless gas.
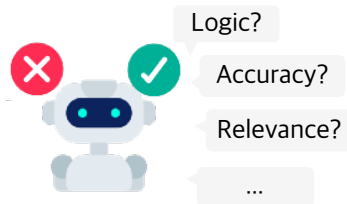2. Helium has no color.
3. So the answer is no.

For question 1, the student made a factual mistake that Aristotle is a modern philosopher···

For question ···, the student listed an irrelevant fact that ···

**To summarize, a good reasoning chain should ···**

- **Accuracy**: Be free of factual errors
- **Relevance**: ···
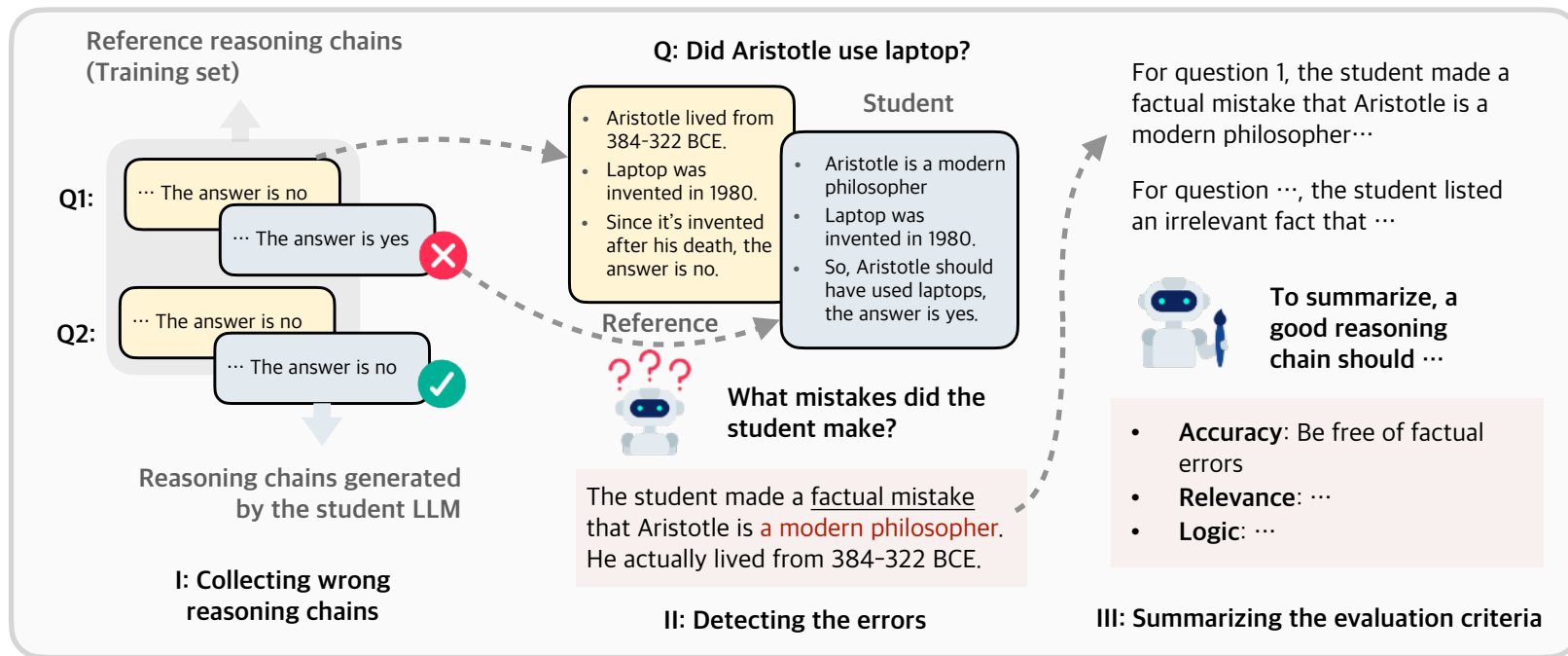- **Logic**: ···

Logic?

Accuracy?

Relevance?

···

**Following the criteria, evaluate the reasoning chain step by step.**

- Accuracy: ···, correct.

- Relevance: The information in the first two steps are irrelevant to the question.

- Logic: The final step cannot be inferred from the previous steps.

So, the reasoning is **INCORRECT**.

# Reasoning Chain Evaluation (RICE)

| Method | Math | | Common | | Logical | | Average |
|---|---|---|---|---|---|---|---|
| | GSM8k | Multi-Arith | StrategyQA | COSMOS | Logic-Deduct | Word-Sort | |
| *Answer-based* | 0.86 | 0.94 | 0.76 | 0.60 | 0.87 | 0.94 | 0.83 |
| SocREval | 0.89 | 0.85 | 0.71 | 0.79 | 0.89 | 0.77 | 0.82 |
| Direct (trace) | 0.90 | 0.38 | 0.80 | 0.72 | 0.21 | 0.36 | 0.56 |
| Direct (step) | 0.85 | 0.43 | 0.83 | 0.71 | 0.75 | 0.33 | 0.65 |
| CoT (step) | 0.78 | 0.74 | 0.78 | 0.70 | 1.00 | 0.86 | 0.81 |
| RICE (Ours) | 0.91 | 0.85 | 0.79 | 0.76 | 0.97 | 0.86 | **0.86** |

- No additional human efforts needed
- Better average accuracy
- Remain robust across different tasks
- More accurate than answer–based evaluation

# Reasoning Chain Evaluation (RICE)

A case study on MultiArith (Srivastava et al., 2023)



**Question:** What is the result of
$(((-9 + 5 - -7 - -8) - (1 * -7 * -2 * 8))$
$+ ((3 * 2 + 2 + 2) * (4 * 3 - 9 * -6))) ?$

**Step 1:** This equation can be written as "A – B + C"···
**Step 2:** Let's calculate A···
···
**Step 4:** Let's calculate C
$= ((3 * 2 + 2 + 2) * (4 * 3 - 9 * -6))$
$= ((6 + 2 + 2) * (12 - 54)) = (10 * -42) = -420.$
···
**Step 6:** So the answer is -521

(Reasoning chain to be evaluated)

### SocREval (Baseline)

**Own response:** ··· **Part C simplifies into (10 \* (-42)) = -420**. So the entire equation simplifies into 11 – 112 – 420 which equals -521.

**Analysis:** The generated response is correct and detailed··· It ends with an accurate conclusion.

**Overall quality:** 5.0 (Correct)

☹ Made the same mistake as the reasoning chain

- **Proper Operation Order**
- **Correct Interpretation of Signs**
- Logical Consistency
- ···
(Criteria learned for Multi-Arith)

### RICE (Ours)

··· **Step 4: This is an incorrect calculation.** The error lies in the calculation of the second part of C, specifically `(4 * 3 - 9 * -6)`. The correct calculation should be:
- Multiplication before addition and subtraction:
    `(4 * 3) + (-9 * -6)` = `12 + 54` = `66`.
- The total for C should then be:
    `(6 + 2 + 2) * 66` = `10 * 66` = `660`.
··· So the reasoning chain is **not correct**.

😊 Successfully recognized the calculation error

48

# Reasoning Chain Evaluation (RICE)

False positive chains recognized by RICE:

| Type | Question | Answer | Explanation |
|------|----------|--------|-------------|
| A | Henry made two stops during his 60-mile bike trip. He first stopped after 20 miles. His second stop was 15 miles before the end of the trip. How many miles did he travel between his first and second stops? | Step 1 - Henry traveled 20 miles + 15 miles = 35 miles **between his first and second stops**. Step 2 - Henry traveled 60 miles - 35 miles = 25 miles without stopping. Step 3 - The answer is 25 | There is a hallucination in an early reasoning step, but the following reasoning steps ignore the mistake and reach the correct answer. |
| B | Is "A Tale of Two Cities" a parody of the Bible? | "A Tale of Two Cities" was published in 1859, and the Bible was published much earlier. **Thus**, "A Tale of Two Cities" is not a parody of the Bible. So the answer is no. | The reasoning chain has obvious or multiple mistakes, but hits the correct answer by chance. |
| C | Is clementine pith highly sought after? | **Clementine pith is not highly sought after.** So the answer is no. | The reasoning chain is not informative at all, though the answer is correct. |

49

# Outline

- Reasoning with LLMs:

  Algorithms, Evaluation, **Analysis**



**LLM Reasoners**

# Experimental Results

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

For three datasets marked with ∗, we evaluate the reasoning chain with both RICE and the answer (RICE / Answer-based).

# Experimental Analysis

From auto-regressive decoding to reward-guided search

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

Overall improved performance with search

# Experimental Analysis

From auto-regressive decoding to reward-guided search

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

Less false positive chains!

# Experimental Analysis

Less false positive chains!

A mechanism to "regret"

From auto-regressive decoding to reward-guided search

| Type | Question | Answer | Explanation |
|------|----------|--------|-------------|
| A | Henry made two stops during his 60-mile bike trip. He first stopped after 20 miles. His second stop was 15 miles before the end of the trip. How many miles did he travel between his first and second stops? | Step 1 - Henry traveled 20 miles + 15 miles = 35 miles **between his first and second stops**. Step 2 - Henry traveled 60 miles - 35 miles = 25 miles without stopping. Step 3 - The answer is 25 | There is a hallucination in an early reasoning step, but the following reasoning steps ignore the mistake and reach the correct answer. |
| B | Is "A Tale of Two Cities" a parody of the Bible? | "A Tale of Two Cities" was published in 1859, and the Bible was published much earlier. **Thus**, "A Tale of Two Cities" is not a parody of the Bible. So the answer is no. | The reasoning chain has obvious or multiple mistakes, but hits the correct answer by chance. |
| C | Is clementine pith highly sought after? | **Clementine pith is not highly sought after.** So the answer is no. | The reasoning chain is not informative at all, though the answer is correct. |

# Experimental Analysis

From auto-regressive decoding to reward-guided search

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

The breadth of search matters more than the depth

# Experimental Analysis

The impact of world model

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

A world model improves performance in most domains

Most effective in Embodied Planning task

# Experimental Analysis

The impact of prompting format

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

**Did either Kublai Khan or his grandfather practice monogamy?**

# Experimental Analysis

**Q: Did either Kublai Khan or his grandfather practice monogamy?**

Chain-of-thoughts:

Step 1: Both Kublai Khan and his grandfather, Genghis Khan, were polygamous.
Step 2: Thus, neither Kublai Khan nor his grandfather practiced monogamy.
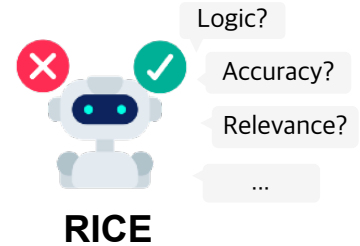Step 3: So the answer is no.

RAP:

Step 1: Did Kublai Khan have a harem? - Kublai Khan had a harem of 7,000 women. So the answer is yes.
Step 2: Did Genghis Khan have a harem? - Genghis Khan had a harem of 500 women. So the answer is yes.
Step 3: Does having a harem of women mean practicing polygamy? - Having a harem of women means practicing polygamy. So the answer is yes.
Step 4: … So the answer is no.

Logic?

Accuracy?

Relevance?

…

**RICE**

# Experimental Analysis

The impact of prompting format

| Type | Question | Answer | Explanation |
|------|----------|--------|-------------|
| A | Henry made two stops during his 60-mile bike trip. He first stopped after 20 miles. His second stop was 15 miles before the end of the trip. How many miles did he travel between his first and second stops? | Step 1 - Henry traveled 20 miles + 15 miles = 35 miles **between his first and second stops**. Step 2 - Henry traveled 60 miles - 35 miles = 25 miles without stopping. Step 3 - The answer is 25 | There is a hallucination in an early reasoning step, but the following reasoning steps ignore the mistake and reach the correct answer. |
| B | Is "A Tale of Two Cities" a parody of the Bible? | "A Tale of Two Cities" was published in 1859, and the Bible was published much earlier. **Thus**, "A Tale of Two Cities" is not a parody of the Bible. So the answer is no. | The reasoning chain has obvious or multiple mistakes, but hits the correct answer by chance. |
| C | Is clementine pith highly sought after? | **Clementine pith is not highly sought after.** So the answer is no. | The reasoning chain is not informative at all, though the answer is correct. |

# Experimental Analysis

The impact of prompting format

| Method | Math | | | Logical | Common | Embodied |
|---|---|---|---|---|---|---|
| | GSM8k* | AQuA* | Game-24 | PrOntoQA | StrategyQA* | Blocksworld |
| CoT | 0.37 / 0.54 | 0.09 / 0.34 | 0.04 | 0.58 | 0.34 / 0.76 | 0.05 |
| RAP (Chain) | 0.44 / 0.52 | 0.11 / 0.34 | 0.01 | 0.43 | 0.28 / 0.72 | 0.19 |
| ToT (BFS) | 0.53 / 0.58 | 0.15 / 0.42 | 0.04 | 0.52 | 0.41 / 0.76 | 0.09 |
| ToT (DFS) | 0.45 / 0.52 | 0.10 / 0.36 | **0.07** | 0.44 | **0.42** / 0.76 | 0.08 |
| RAP | **0.58 / 0.64** | **0.20 / 0.47** | **0.07** | **0.59** | 0.28 / **0.77** | **0.51** |

**Easier to trigger false positives**

But only for certain datasets, where the details are not necessary

60

# Summary

- Reasoning with LLMs:

  Algorithms, Evaluation, Analysis

**LLM Reasoners**