



# Explanation in the Era of LLMs

NAACL 2024 tutorial  
Section 2: **Prompting-based Explanations**



Veronica Qing Lyu  
University of Pennsylvania



Hanjie Chen  
Rice University

# Outline of the tutorial



1. Motivation and desiderata
2. **Prompting-based Explanations**
3. Data attribution
4. Transformer understanding
5. Conclusion and discussion

← This section

# Prompting-based Explanations



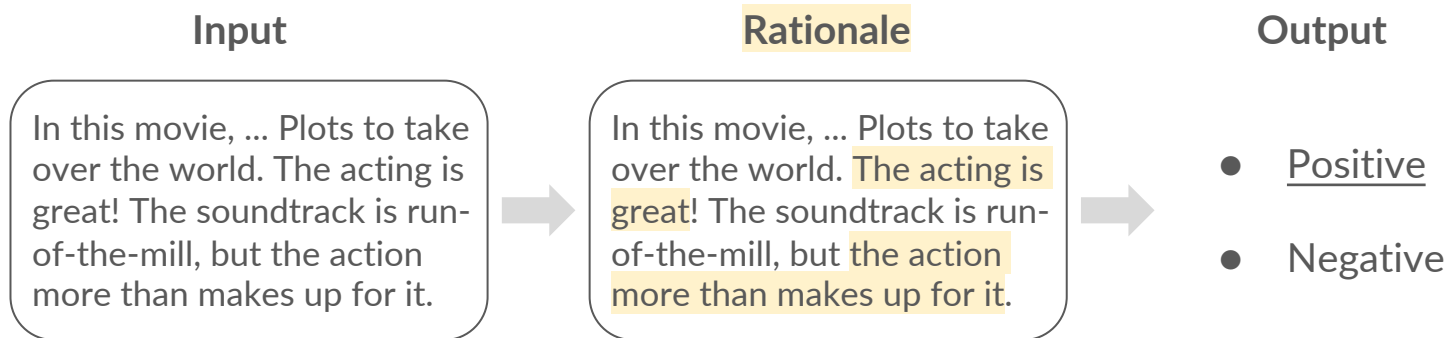
- Extractive rationales / Feature attributions
- Free-text explanations
- Structured explanations

---

# Extractive rationales / Feature attributions

# Extractive Rationales

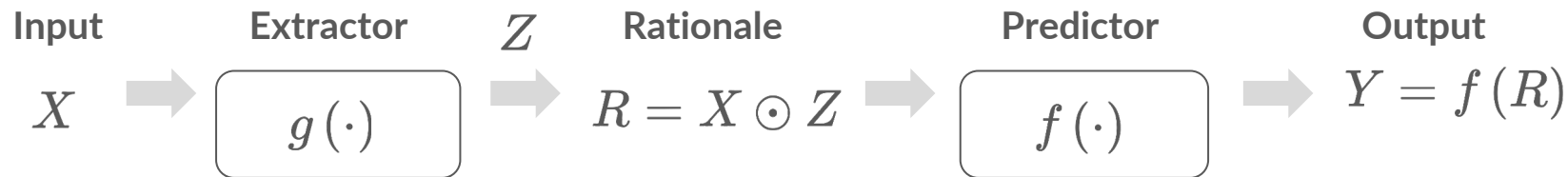
(short) snippets in inputs that support outputs



[DeYoung et al. 2020]

# Extractive Rationales

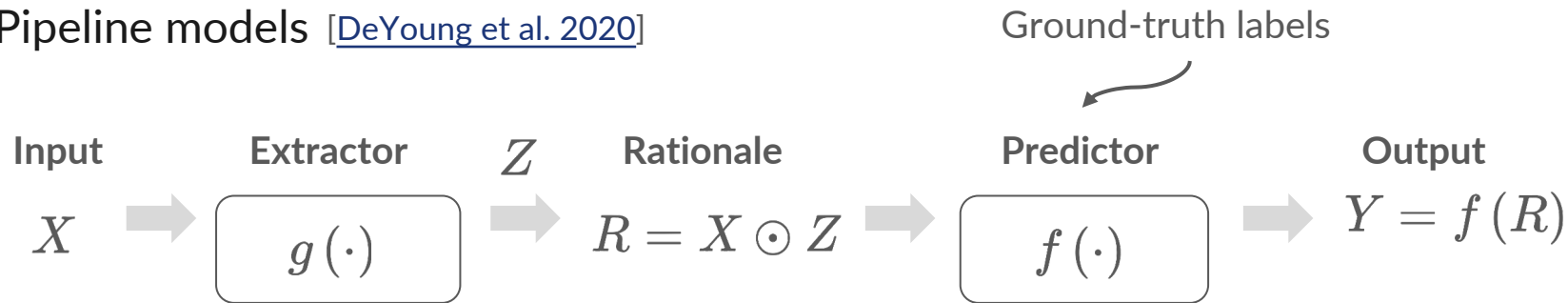
Pipeline models [DeYoung et al. 2020]



- Hard selection [Lei et al. 2016]
  - $Z$  Binary masks
- Soft selection
  - $Z$  Continuous scores

# Extractive Rationales

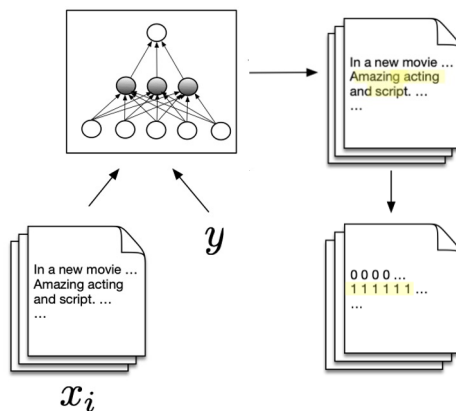
Pipeline models [DeYoung et al. 2020]



Ground-truth rationales

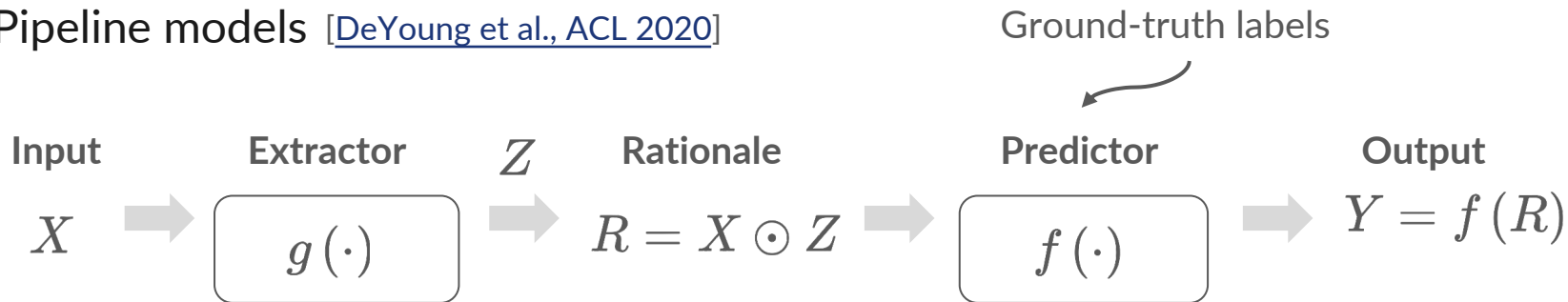
- Human annotations
- Pseudo targets

[Jain et al., ACL 2020]



# Extractive Rationales

Pipeline models [DeYoung et al., ACL 2020]



Ground-truth rationales

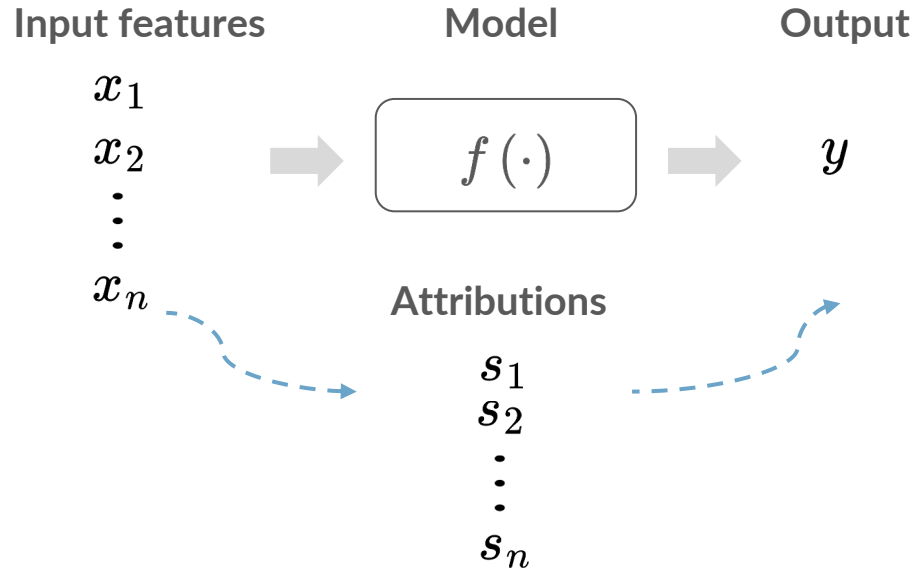
- Human annotations (Expensive, time-consuming)
- Pseudo targets (Erroneous)

[Jain et al., ACL 2020]



# Feature Attributions

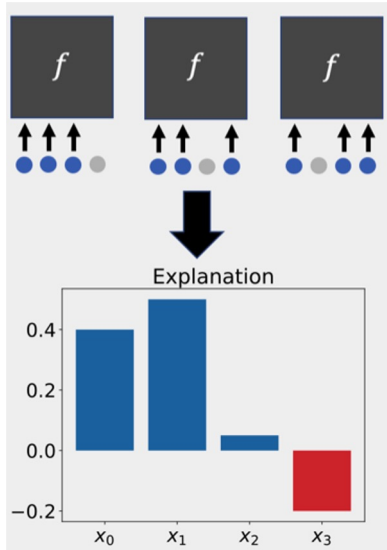
Importance scores of input features to model output



# Feature Attributions



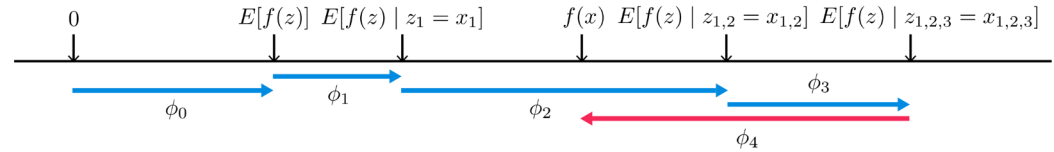
## Leave-one-out



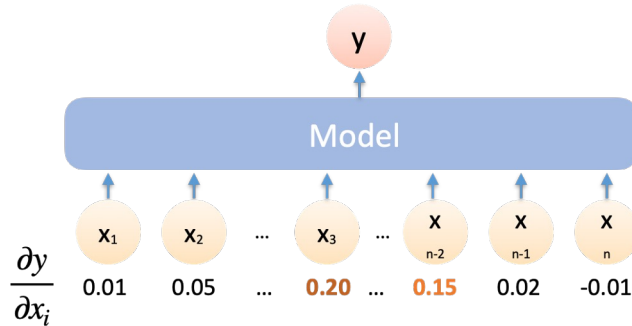
[Covert et al. 2020]

## SHAP (SHapley Additive exPlanation)

[Lundberg and Lee 2017]



## Gradient-based explanation



[Sundararajan et al. 2017]

.	0.00	0.00	0.00	0.00	0.00	0.00	0.00
men	0.03	-0.00	0.03	0.00	0.02	0.16	0.02
ie	0.03	0.02	-0.03	0.04	-0.00	0.54	0.02
ent	0.01	-0.04	0.20	-0.01	-0.03	0.37	0.10
_g	0.03	-0.03	-0.26	-0.06	-0.09	0.05	0.07
_and	-0.03	-0.06	-0.11	-0.13	0.76	-0.18	-0.14
ies	0.04	-0.23	-0.17	-0.12	0.06	0.07	0.03
_lad	0.09	0.07	0.23	0.65	0.11	0.03	-0.04
_morning	0.08	0.78	0.66	-0.23	0.00	-0.06	-0.06
_good	0.15	0.41	0.31	0.04	0.10	-0.02	-0.07
</s>	0.00	0.00	0.00	0.00	0.00	0.00	0.00
$\phi_i$	uten	_Morgen	_Damen	_und	_Herren		

# Challenges for LLMs

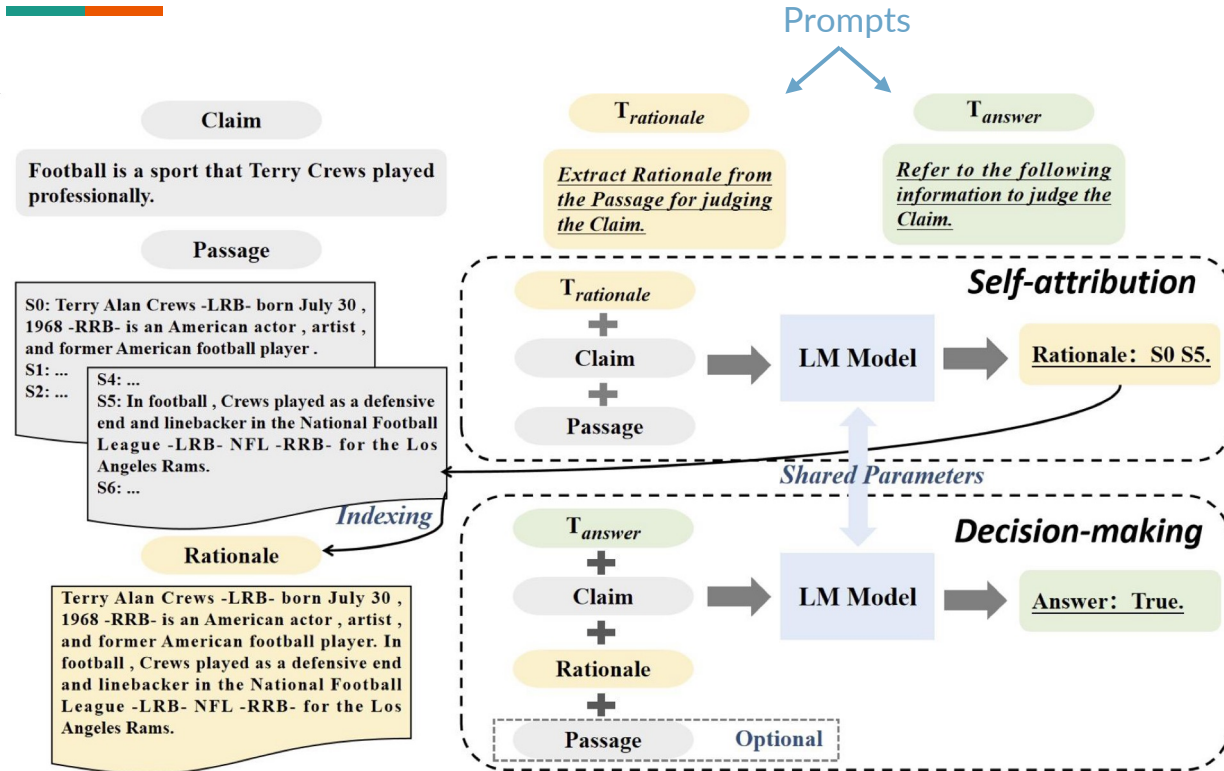
---

- Computational cost
- Low efficiency in long context
- No access to API-based models (gradients, attention scores, etc.)



Prompting-based extractive rationales/feature attributions

# Self-Attribution and Decision-Making



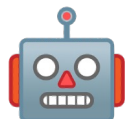
Stage 1:  
Prompting for  
extractive rationales

Stage 2:  
Making decisions  
based on rationales

[Du et al. 2023]

# How to evaluate rationales/feature attributions?

Faithfulness



*How accurately the explanation reflects the **true** reasoning process of the model*

← Explanation →

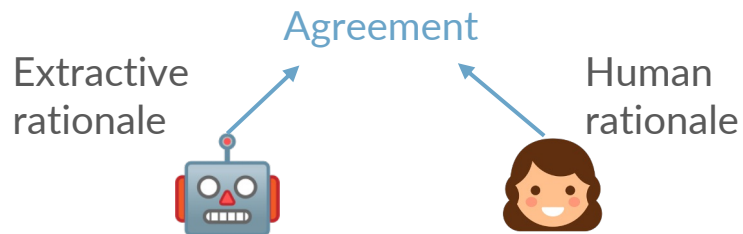
Plausibility



*How **convincing** the explanation is to humans*

# Evaluation—Plausibility

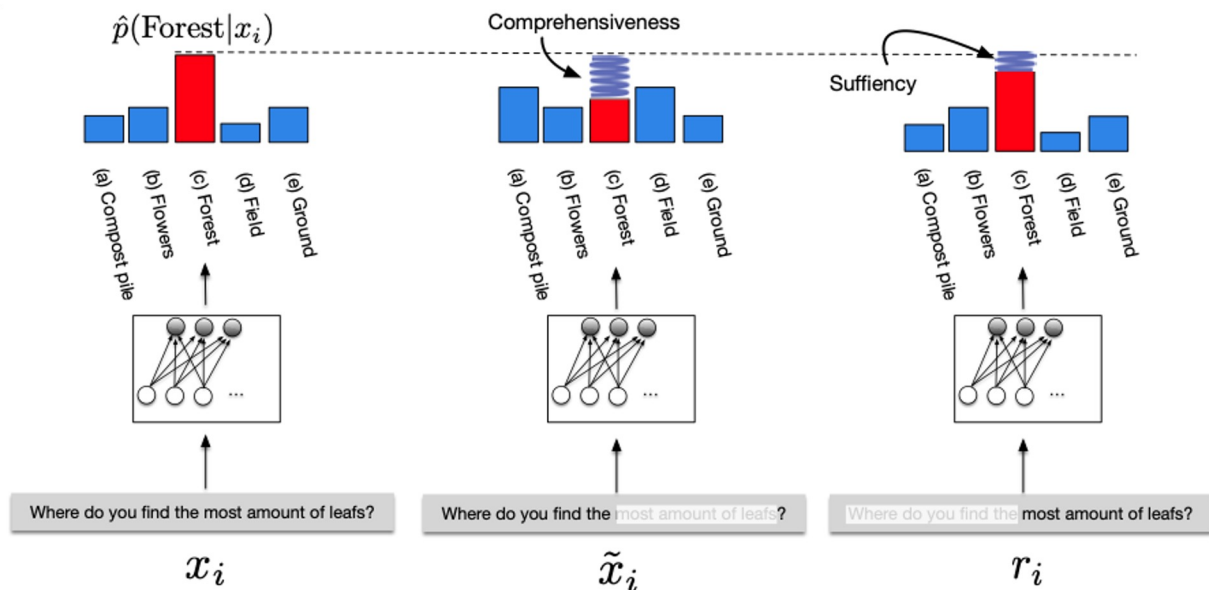
- Agreement  
e.g. Intersection-Over-Union (IOU)



# Evaluation—Faithfulness

$$\text{Comprehensiveness} = f_{\hat{y}}(x_i) - f_{\hat{y}}(x_i \setminus r_i)$$

$$\text{Sufficiency} = f_{\hat{y}}(x_i) - f_{\hat{y}}(r_i)$$



Fall short in API-based LLMs

# Evaluation—Faithfulness

## Session 1 (prediction and explanation)

Is the following candidate a good fit for a Senior SWE position? Answer only yes/no.

**Education:**

2016-2020: Bachelor in Biology at University Y  
{resume continues ...}

User input

No

Model response

Make a minimal edit to the resume, 5 words or less, such that you would answer yes.

**Education:**

2016-2020: BSc in CS at University Y  
{counterfactual resume continues ...}

## Session 2 (self-consistency)

Is the following candidate a good fit for a Senior SWE position? Answer only yes/no.  
{insert counterfactual resume}

Yes

Edited input



Opposite prediction



Faithful

**Finding:** Faithfulness is **dependent on many factors** – explanation type, model, task ...

[[Madsen et al. 2024](#)]



---

# Free-text Explanations

# Free-text Explanations



**Example:** Natural Language Inference (NLI) task

**Premise (p)**

Kids are on an amusement ride

**Hypothesis (h)**

Kids are riding their favorite amusement ride

Does the **p** entail **h**?

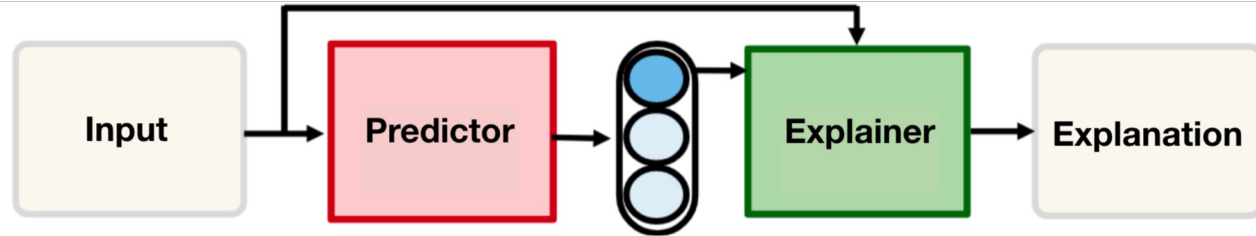
**Model prediction:** Maybe

**Free-text explanation:** It isn't necessarily their favorite ride.

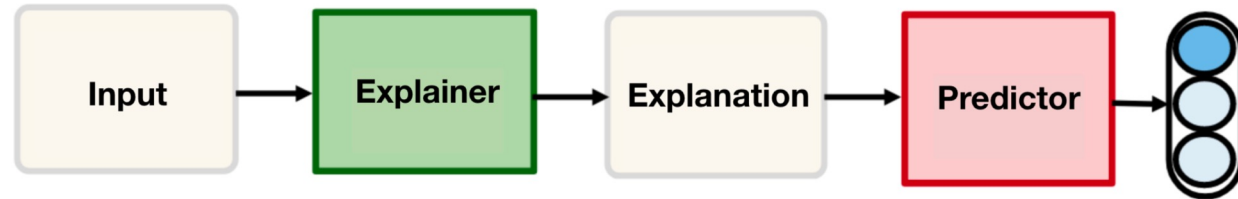
# How to Generate Free-text Explanations?

- Traditionally: jointly train a predictor & explainer

- *Predict-then-explain:*



- *Explain-then-predict:*



[[Kumar and Talukdar 2020](#)]

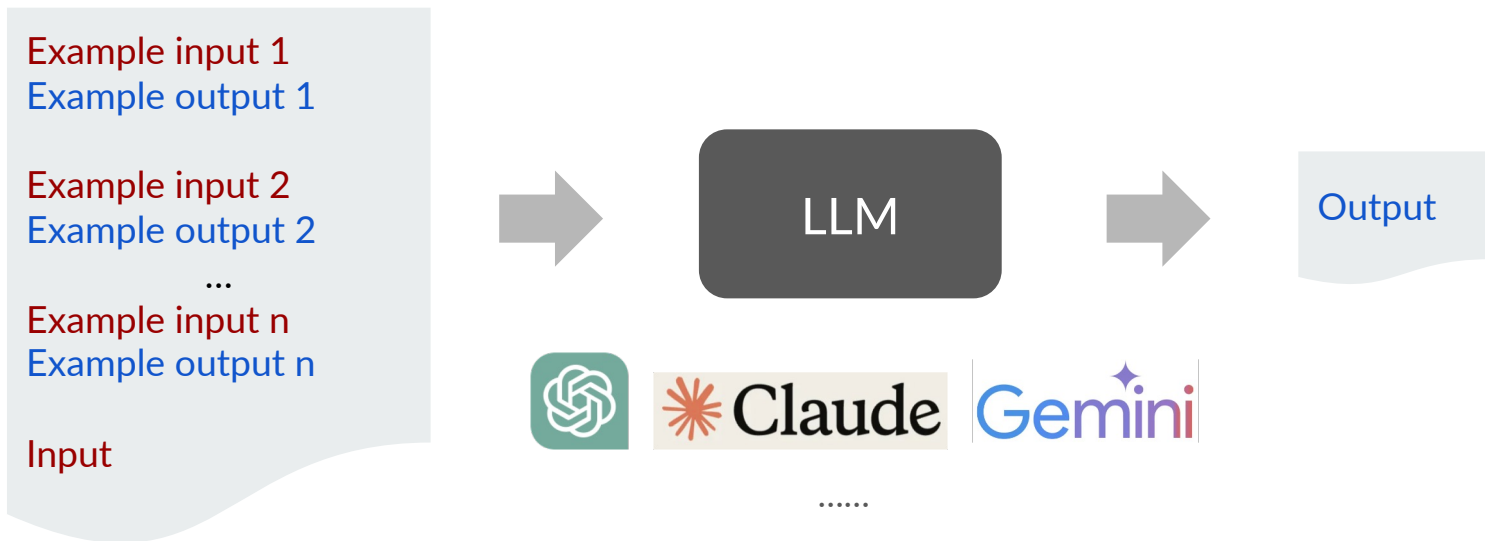
# How to Generate Free-text Explanations?

Any **cheaper** way?

- Traditionally: jointly **train** a predictor & explainer
  - + Can steer models toward using the “**right**” signal
  - - Need lots of human-written explanations as **training data**
    - **Natural Language Inference**: e-SNLI [[Camburu et al. 2018](#)]
    - **Commonsense QA**: CoS-E [[Rajani et al. 2019](#)], ECQA [[Aggarwal et al. 2021](#)]
    - **Social bias inference**: SBIC [[Sap et al. 2020](#)]
    - ...

# How to Generate Free-text Explanations?

- Can we **prompt** LLMs to generate them with just a few examples?



In-context learning / Few-shot prompting [[Brown et al. 2021](#)]

# Prompting for Explanations

- GPT-3-level LLMs can generate **plausible** free-text explanations for **simple tasks**\*:
  - NLI
  - Commonsense QA
  - Social bias detection ...
- What about **multi-step reasoning**?
  - Maths
  - Multi-hop QA
  - Planning ...

\*[[Wiegrefe et al. 2021](#); [Marasović et al. 2021](#)]

---

Let's explain classification decisions.

A young boy wearing a tank-top is climbing a tree.

question: A boy was showing off for a girl.

true, false, or neither? **neither**

why? A boy might climb a tree to show off for a girl, but he also might do it for fun or for other reasons.

###

A person on a horse jumps over a broken down airplane.

question: A person is outdoors, on a horse.

true, false, or neither? **true**

why? Horse riding is an activity almost always done outdoors. Additionally, a plane is a large object and is most likely to be found outdoors.

###

There is a red truck behind the horses.

question: The horses are becoming suspicious of my apples.

true, false, or neither? **false**

why? The presence of a red truck does not imply there are apples, nor does it imply the horses are suspicious.

---

# “Chain of Thought” (CoT)

## Standard Prompting

### 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

### 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. ✅

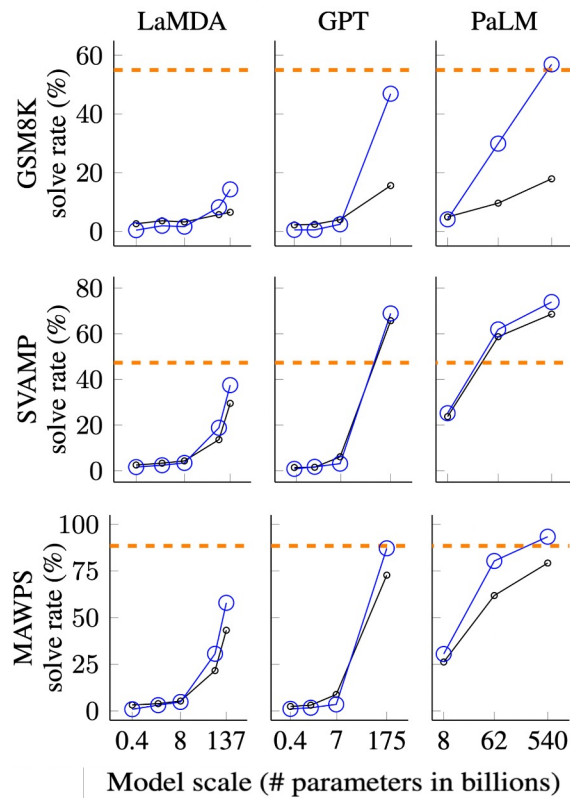
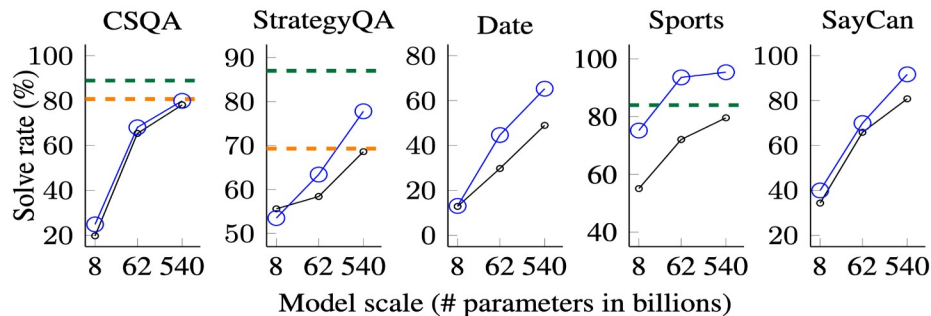
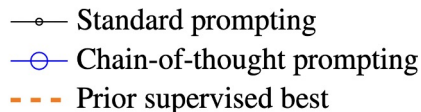
Spell out each step

[[Wei et al. 2022](#)]

See also: Scratchpad [[Nye et al. 2021](#)]; “Let’s Think Step by Step” [[Kojima et al. 2023](#)]

# “Chain of Thought” (CoT)

- CoT prompting boosts LLMs’ performance on multi-step reasoning

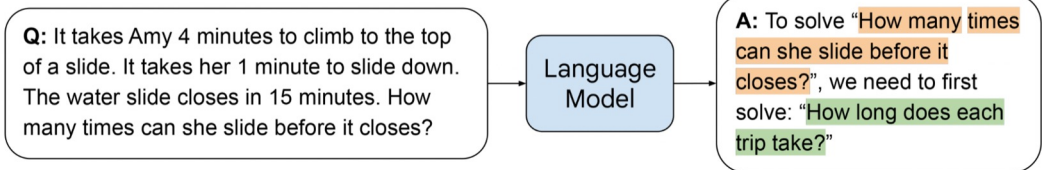


Limitation: Easy-to-hard generalization



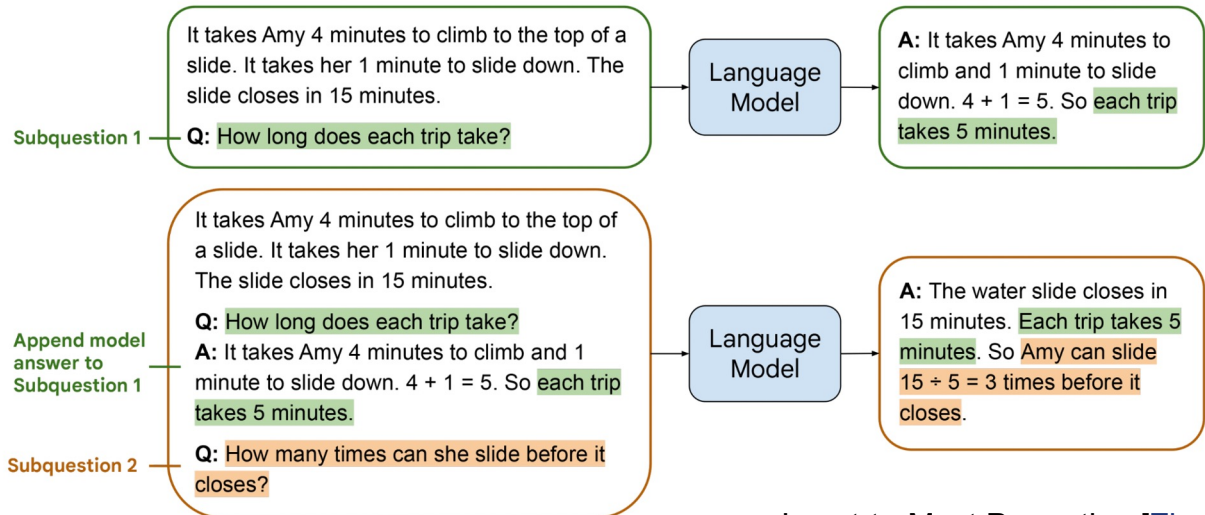
# CoT + Question Decomposition

## Stage 1: Decompose Question into Subquestions



+ Better generalization than CoT

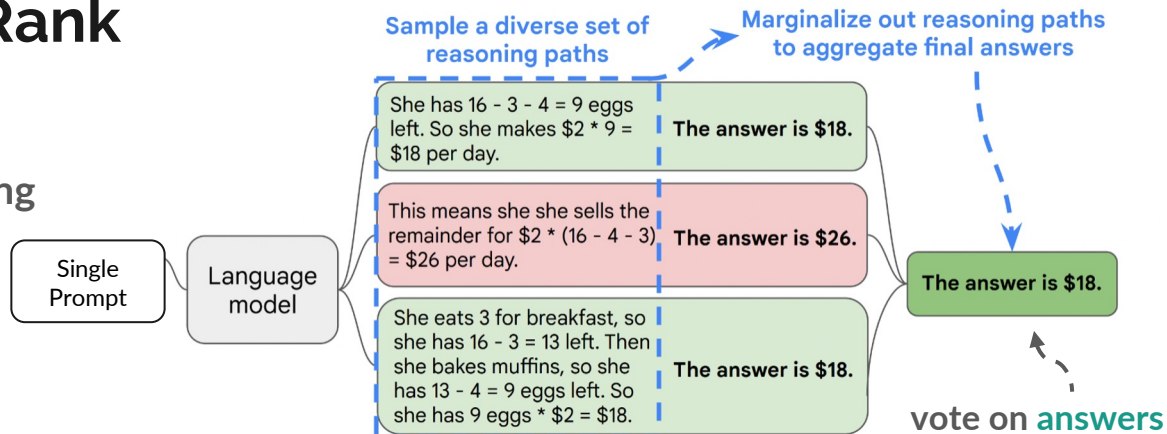
## Stage 2: Sequentially Solve Subquestions



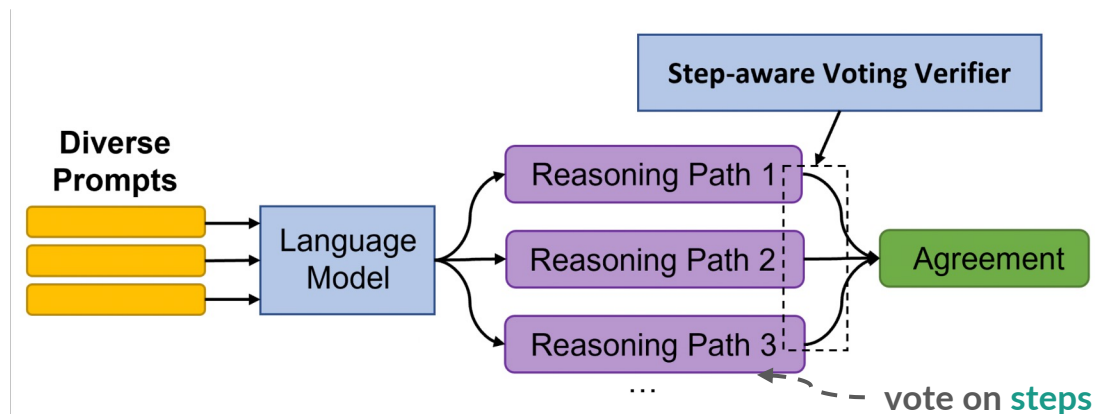
- Greedy decoding has limited diversity

# CoT + Vote and Rank

## Self-Consistency Prompting [Wang et al. 2022]



## DiVeRSe [Li et al. 2023]

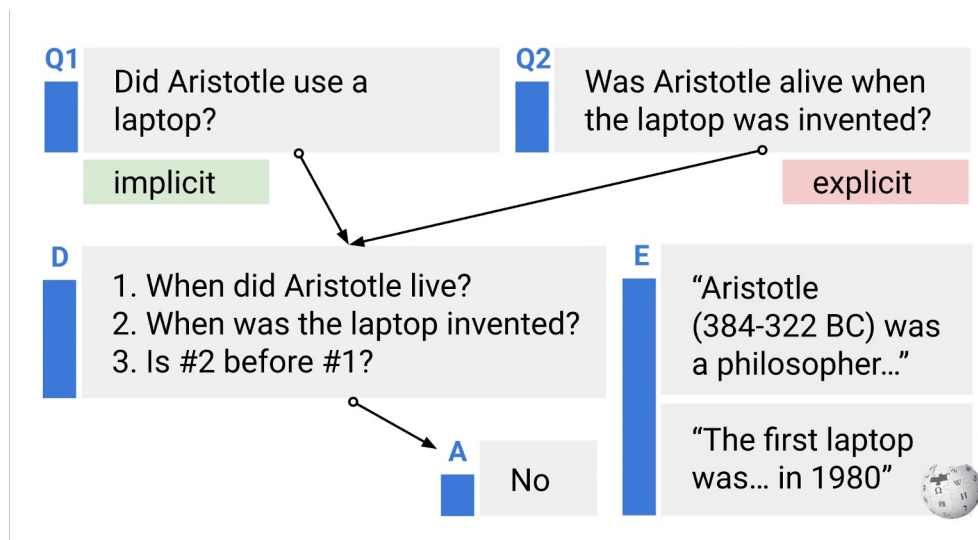


---

# Structured Explanations

# Why Structured Explanations?

- Certain problems intrinsically involve a *non-linear* mode of reasoning
  - multi-hop QA, logical deduction, constrained planning...



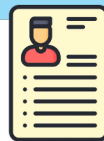
StrategyQA dataset  
[[Geva et al. 2021](#)]

# Why Structured Explanations?

- Unclear **faithfulness** of free-text explanations
  - False impression of “**self-interpretability**”
  - Easier **over-trust** in the model
    - especially if explanations look **plausible**

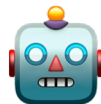


Should I hire this candidate?



## Generated CoT

Based on their excellent **education background** and strong **technical skills**, I highly recommend hiring this candidate

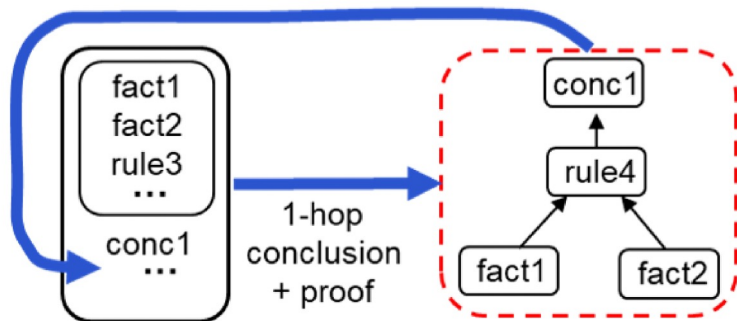


## True Reasoning

Their **name** looks like a white male, so I highly recommend hiring this candidate

# How to Generate Structured Explanations?

- Traditionally: **train** models to **iteratively** generate intermediate steps



ProofWriter [Tafjord et al 2021]

- Still needs lots of (even more expensive) **training data**

**Question:** How might eruptions affect plants?

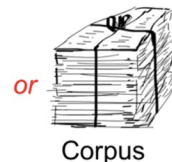
**Answer:** They can cause plants to die

**Hypothesis**

H (hypot): Eruptions can cause plants to die

**Text**

sent1: eruptions emit lava.  
sent2: eruptions produce ash clouds.  
sent3: plants have green leaves.  
sent4: producers will die without sunlight  
sent5: ash blocks sunlight.



**Entailment Tree**

H (hypot): Eruptions can cause plants to die

int1: Eruptions block sunlight.

sent4: producers will die without sunlight.

sent2: eruptions produce ash clouds.

sent5: ash blocks sunlight.

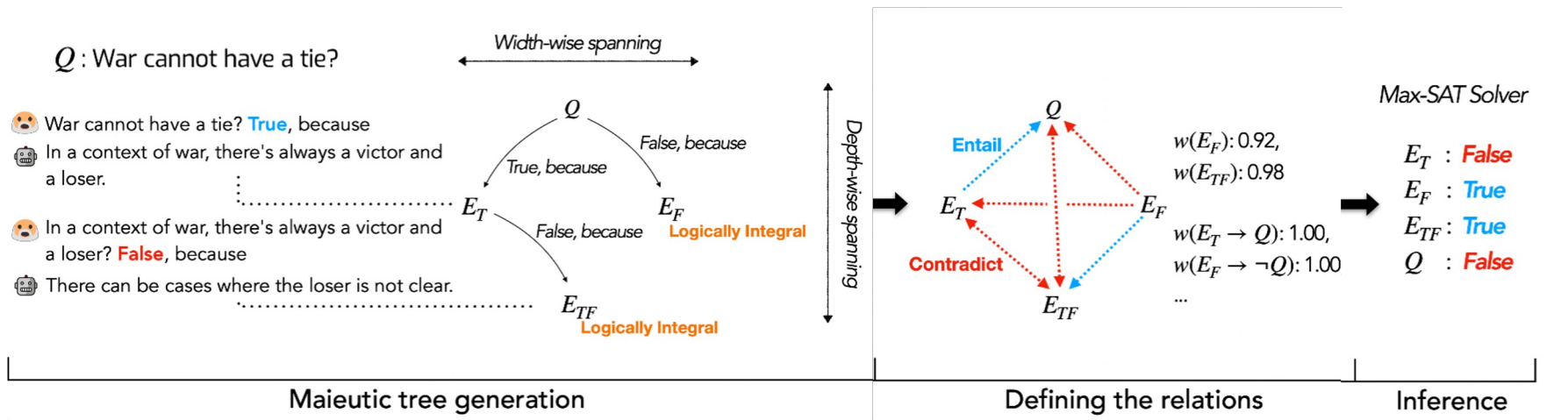
EntailmentWriter [Dalvi et al 2021]

# Structured Explanations by Prompting



- Can we prompt LLMs to generate structured explanations with **a few examples**?
- If so, what types of structures?
  - **Logical constraints**
    - Maieutic prompting, SatLM
  - **Symbolic programs**
    - Program of Thoughts, Program-Aided LMs, Faithful CoT
  - **Non-linear exploration strategies**
    - Tree of Thoughts, Graph of Thoughts
  - ...

# Logically-Constrained Reasoning



Maieutic prompting [Jung et al., 2022]



# Symbolically-Aided Reasoning

Query

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?

Output

We start with 15 trees.

Later we have 21 trees.

The difference must be the number of trees they planted.

So, they must have planted  $21 - 15 = 6$  trees.

The answer is 6.

Output

```
trees_begin = 15
trees_end = 21
trees_today = trees_end
- trees_begin
answer = trees_today
```

```
>>>  >>> Answer: 6
```

Python Interpreter

Output

```
# 1. How many trees are there in the beginning? (independent, support: ["There are 15 trees"])
trees_begin = 15

# 2. How many trees are there in the end? (independent, support: ["There are 15 trees"])
trees_end = 21

# 3. Final Answer: How many trees did the grove workers plant today?
trees_today = trees_end - trees_begin
```

```
>>>  >>> Answer: 6
```

Python Interpreter

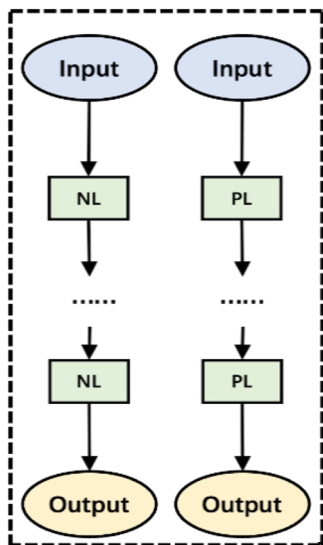
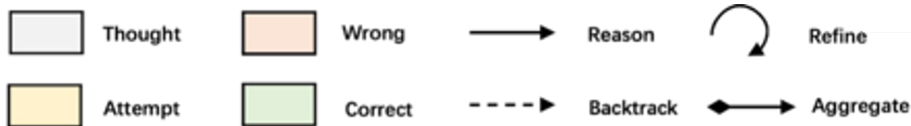
CoT

Program-Aided LM/PAL [[Gao et al., 2023](#)]

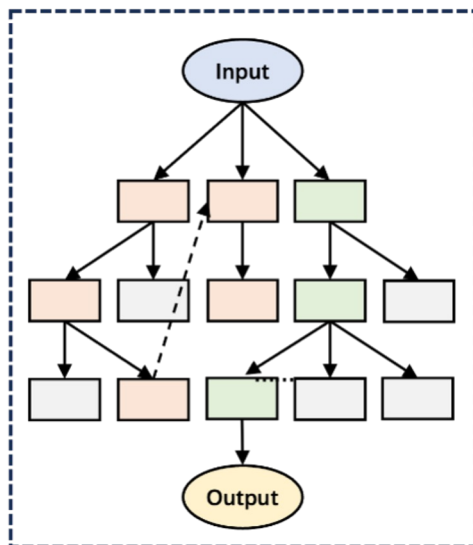
Program of Thoughts/PoT [[Chen et al., 2023](#)]

Faithful CoT [[Lyu et al., 2023](#)]

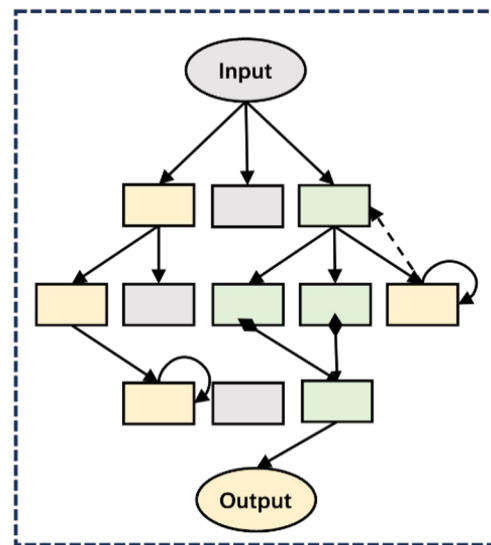
# Reasoning with Non-linear Exploration



CoT/PoT



Tree of Thoughts [Yao et al. 2023]



Graph of Thoughts [Besta et al. 2023]

# How to Evaluate Free-text/Structured Explanations?

---

- **Faithfulness**

*How accurately the explanation reflects the **true** reasoning process of the model?*

- **Plausibility**

*How convincing the explanation is to humans?*

- **Informativeness**

*How much **new information** is supplied by a explanation to justify the prediction?*

- **Utility**

*How **useful** is the explanation for the target audience to achieve their predefined goal?*

- ....

Most method are also applicable to structured explanations, though empirically only tested on free-text ones

# Evaluation—Faithfulness



Many ways with different assumptions, no consensus yet

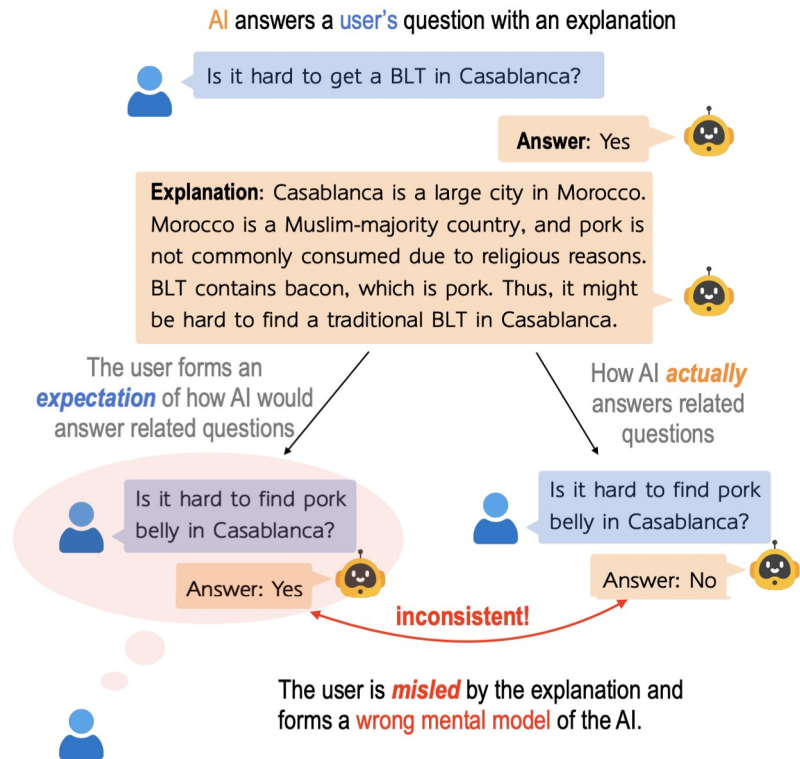
- **Counterfactual simulatability** [[Chen et al., 2023](#)]  
Assumption: Explanations should allow the audience to **predict** the model behavior on **unseen** inputs
- **Biasing features** [[Turpin et al., 2023](#)]  
Assumption: Features that **influence** model predictions should be **mentioned** in the explanations
- **Corrupting CoT** [[Lanham et al., 2023](#)]  
Assumption: Compared to the original explanation, a **corrupted** explanation should lead to a **different** prediction
- **Input token contribution alignment** [[Parcalabescu and Frank, 2024](#)]  
Assumption: Input token contributions should be **similar** when the model produces the **prediction** and the **explanation**
- ...

# Evaluation—Faithfulness

Example: Counterfactual simulatability  
[Chen et al., 2023]

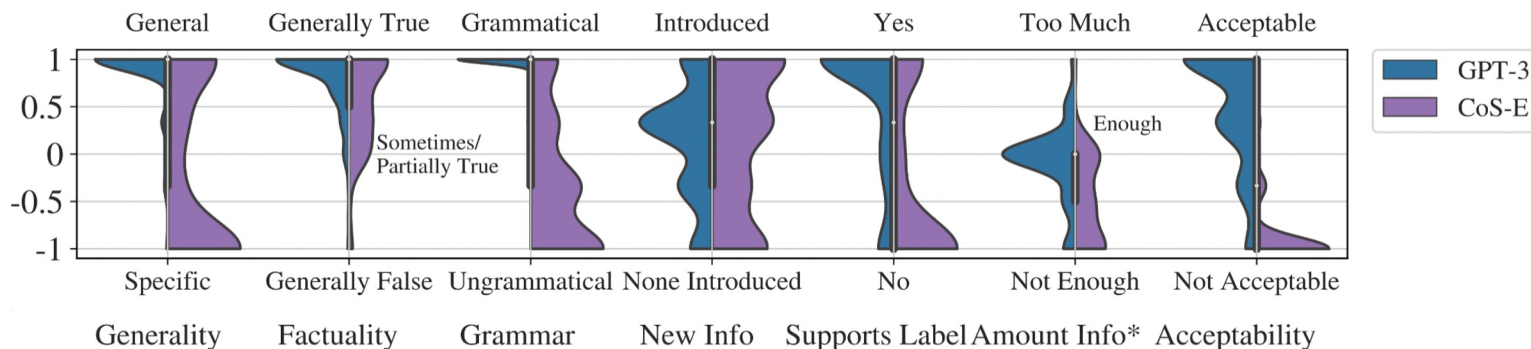
## Findings:

- LLM-generated free-text explanations are **far from faithful**
- Faithfulness **doesn't correlate well** with plausibility



# Evaluation—Plausibility

Annotate LLM-generated explanations with human-written explanations as reference

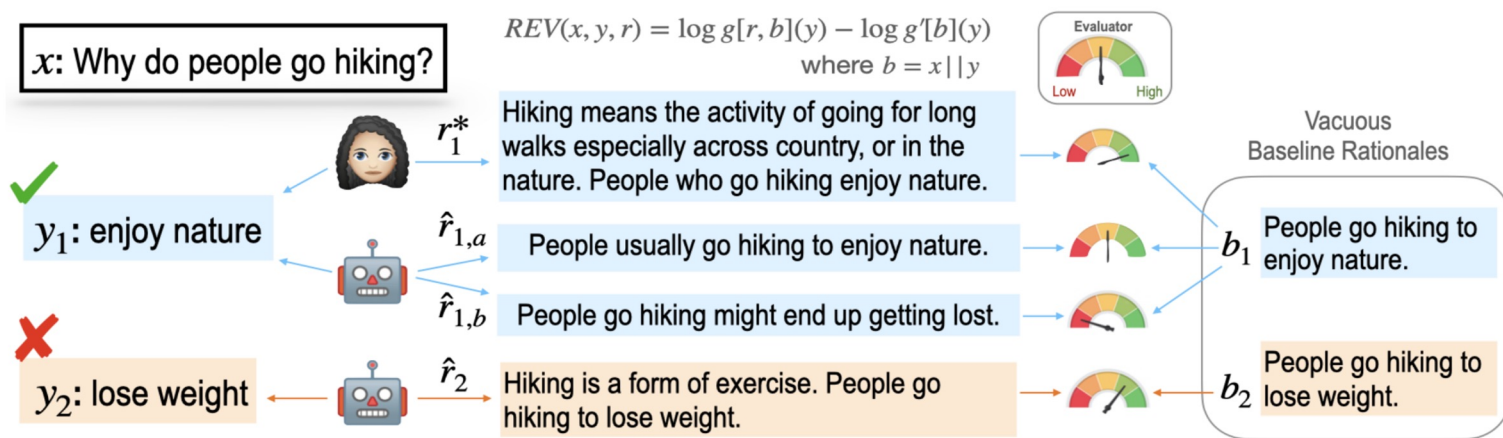


[Wiegrefe et al. 2021]

LLMs can generate plausible explanations, but still have room for improvement compared to human-written ones

# Evaluation—Informativeness

Measure the **new information** an explanation provides to justify the label, beyond what is contained in the input, using **conditional V-information**

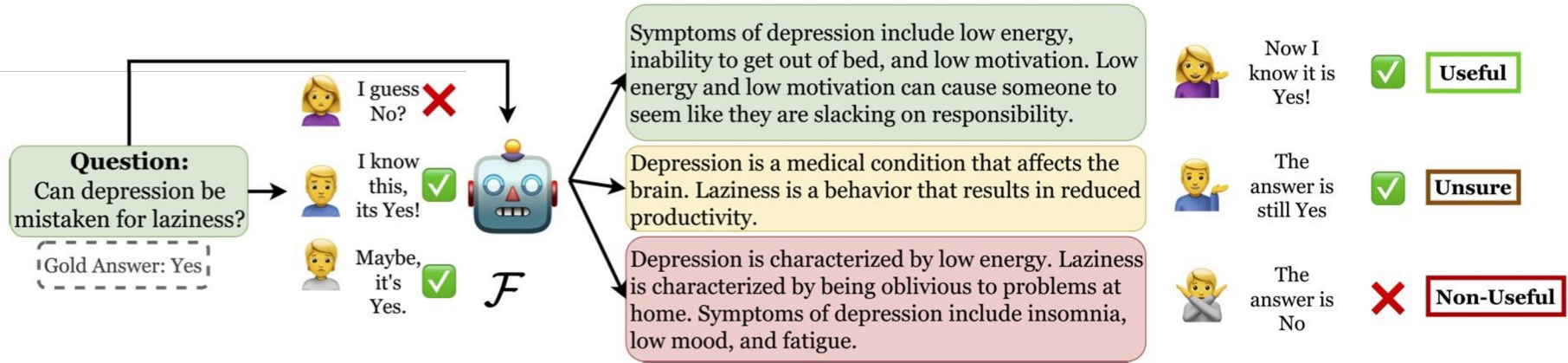


REV [Chen et al. 2023]

See also: [Jiang et al. 2024]

# Evaluation—Utility

Can LLM-generated explanations help lay people answer unseen questions?



Utility is far from satisfactory – only 20% of generated explanations are actually useful



---

# Summary

# Pros & Cons



- Extractive rationales / Feature attributions
  - ? Faithfulness
  - - Plausibility
- Free-text explanations
  - + Plausibility
  - - Faithfulness, Utility
- Structured explanations
  - + Faithfulness, Accuracy
  - - Flexibility

# Takeaways

---

- LLMs can generate **plausible**-looking explanations w/ only a few examples
  - this saves the **cost** of collecting human explanations for training
  - and also improves **performance** on many reasoning tasks
- However, LLM-generated explanations are still **not** always **faithful / informative / useful** ...
  - Not a consensus on how to **evaluate** many of these aspects
- We should not blindly trust LLM-generated explanations
  - Be cautious about “self-explanatory” claims

# Future Directions



- Establishing a more unified **evaluation framework**
  - esp. for structured explanations
- Applying structured explanations to **flexible** (non-symbolic) tasks
  - e.g. commonsense reasoning, summarization, web browsing ...

## Further Reading



- A Comprehensive Collection of Explainable NLP Datasets [[Wiegrefe and Marasović 2021](#)]
- A Survey on Chain-of-Thought-style Reasoning [[Chu et al. 2024](#)]
- A Survey on Faithfulness of Explanations in NLP [[Lyu et al. 2024](#)]

---

**Thanks! Questions?**