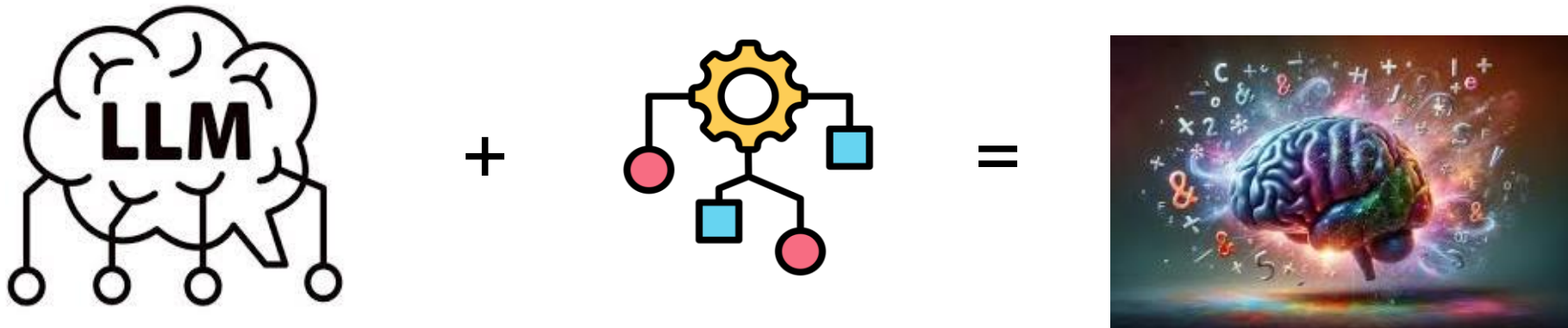


LINC: A Neuro-symbolic Approach for Logical Reasoning by Combining Language Models with First-Order Logic Provers

Olausson et. al. (MIT CSAIL, MIT BCS)



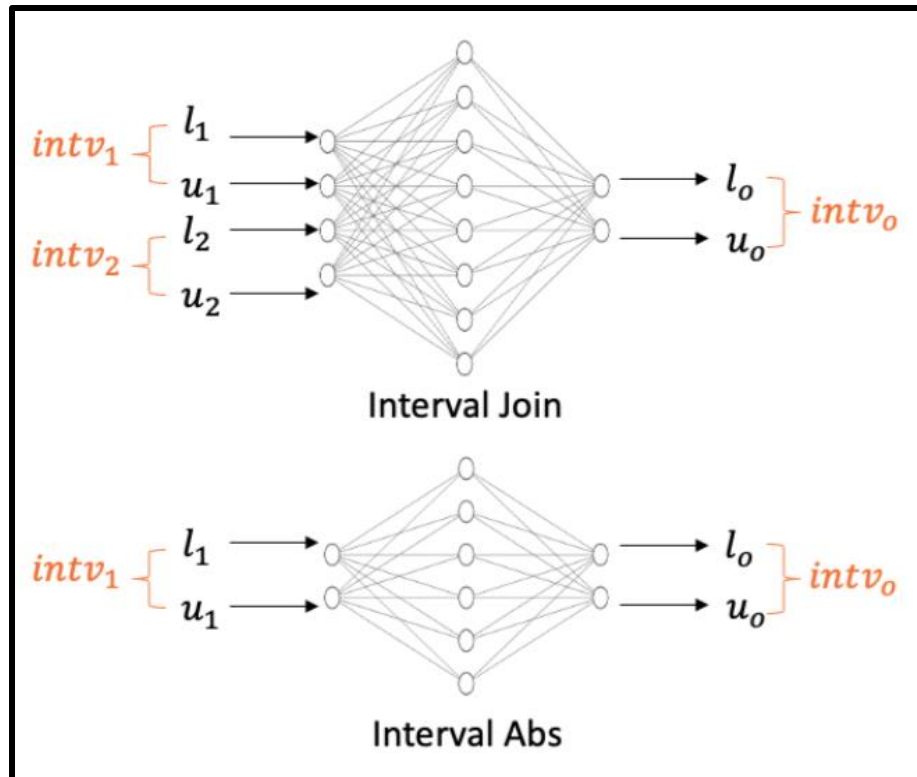
Presented by: Shaurya Gomber

About Me

- 2nd yr PhD student (Advisor: Prof. Gagandeep Singh).
- Learning (and optimization) based techniques for efficient and precise abstract interpretation.

About Me

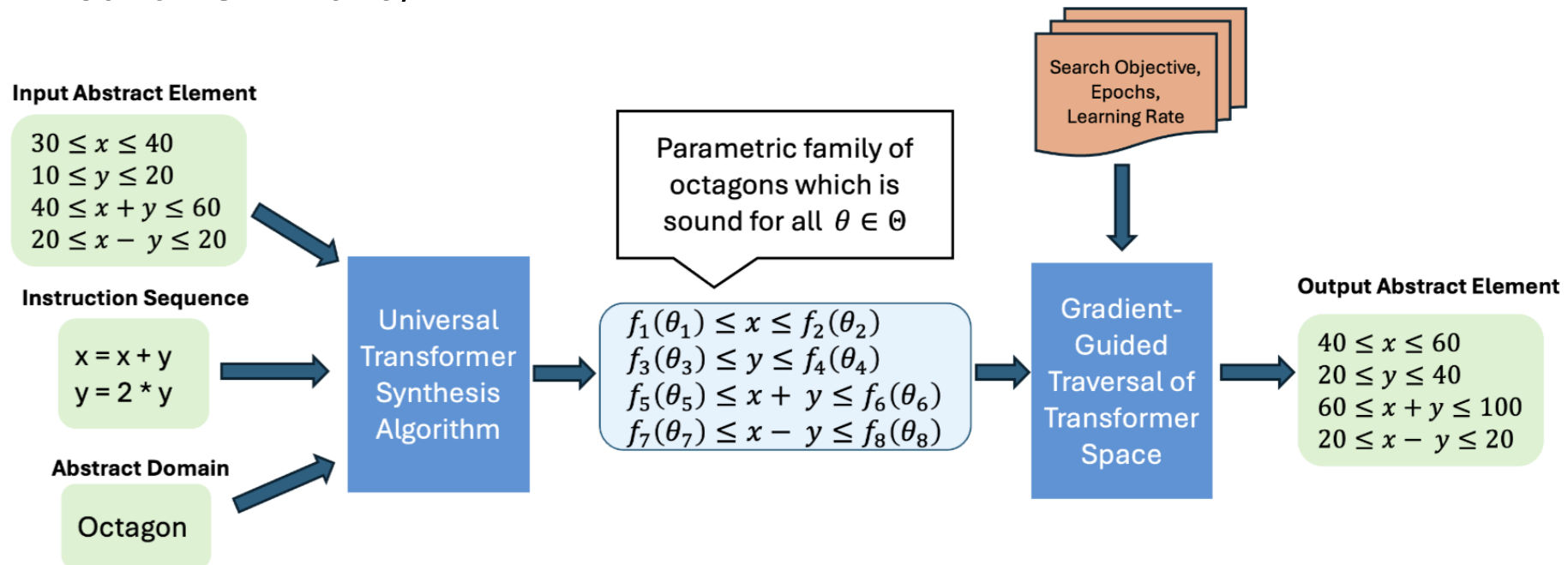
- 2nd yr PhD student (Advisor: Prof. Gagandeep Singh).
- Learning (and optimization) based techniques for efficient and precise abstract interpretation.
- MS Thesis: **Neural Abstract Interpretation** (VerifAI Workshop @ ICLR 2025)



Supervised and Unsupervised learning formulations to train "neural" abstract transformers!

About Me

- 2nd yr PhD student (Advisor: Prof. Gagandeep Singh).
- Learning (and optimization) based techniques for efficient and precise abstract interpretation.
- MS Thesis: Neural Abstract Interpretation (VerifAI Workshop @ ICLR 2025)
- Current: **Universal Synthesis of Differentiably Tunable Numerical Abstract Transformers** (submitted to POPL 2026)



Contents

- Problem & Motivation
- Methodology
- Datasets & Models
- Baselines
- Results
- Conclusions & Limitations

Contents

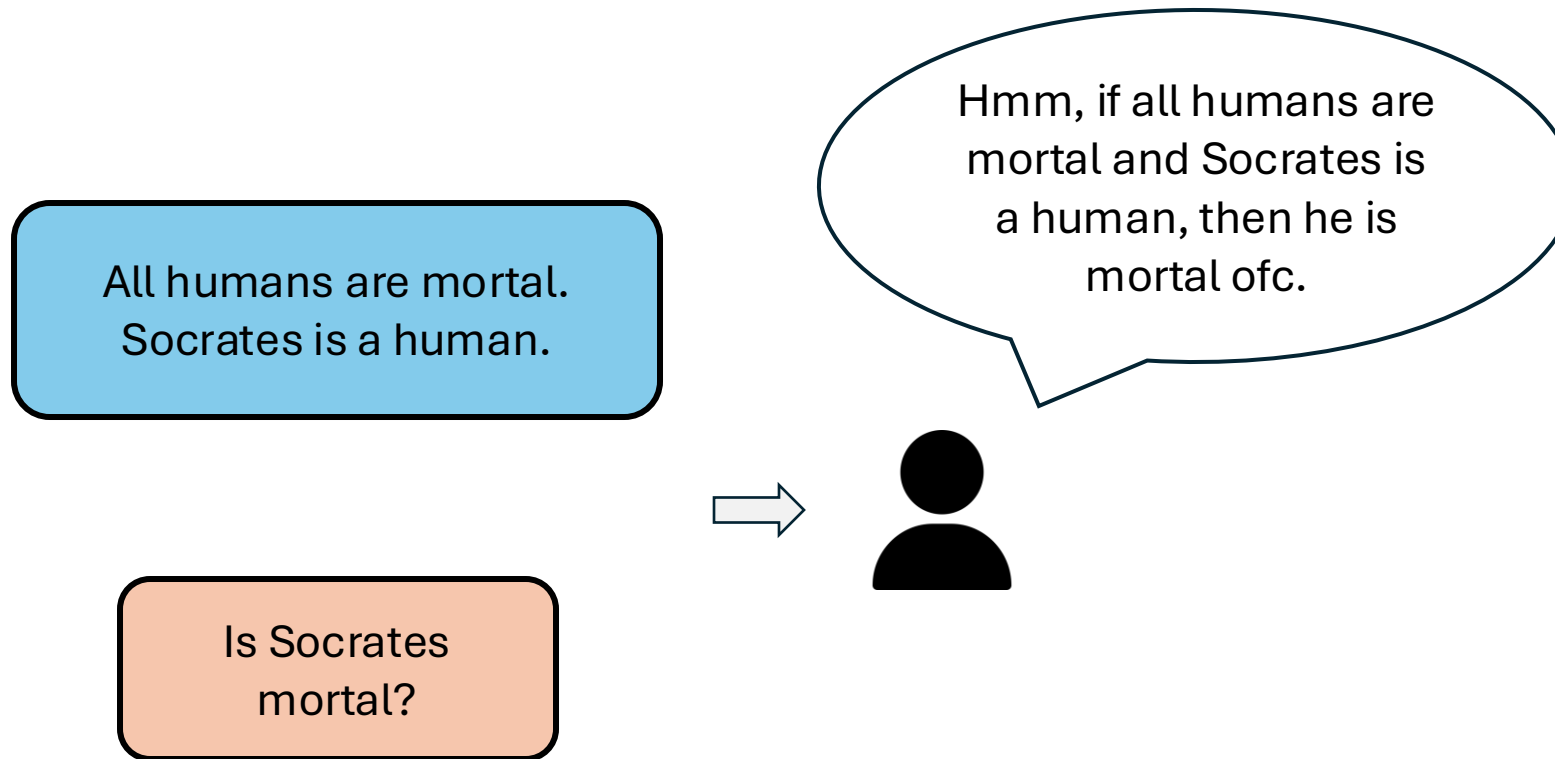
- **Problem & Motivation**
- Methodology
- Datasets & Models
- Baselines
- Results
- Conclusions & Limitations

Logical Reasoning

Structured rules → analyze facts (premises) → answer questions (derive conclusions)!

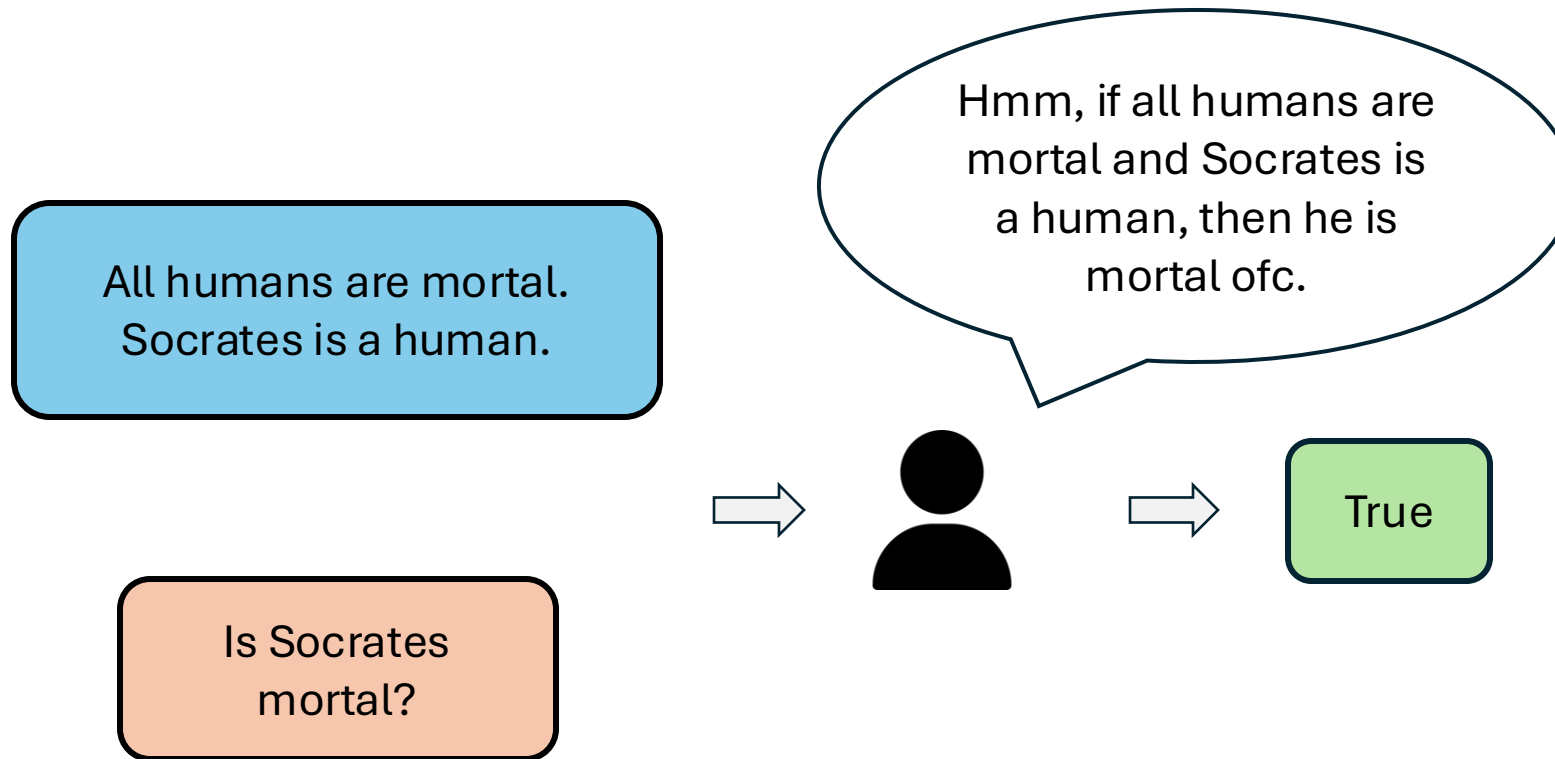
Logical Reasoning

Structured rules → analyze **facts (premises)** → answer **questions** (derive conclusions)!



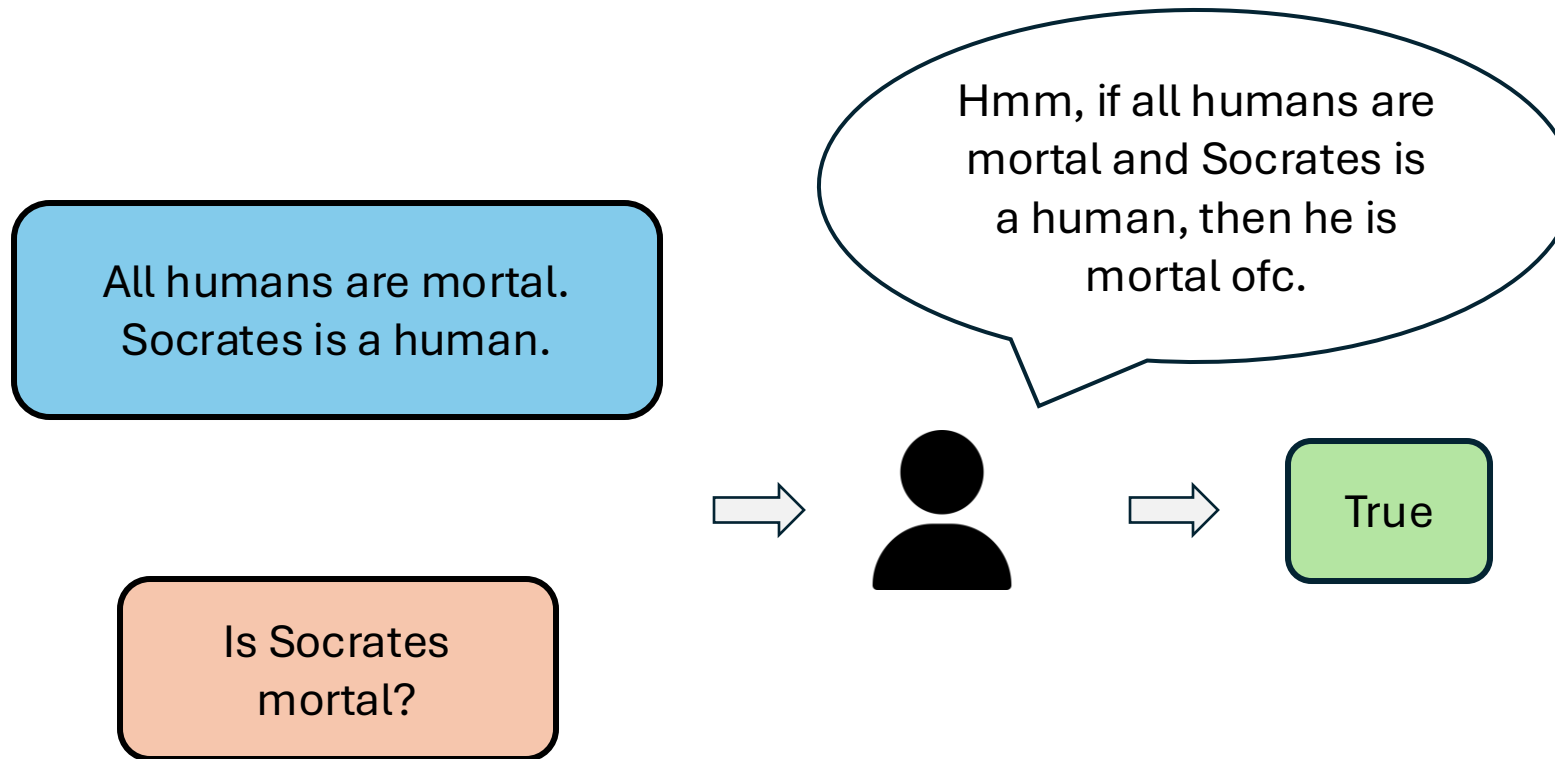
Logical Reasoning

Structured rules → analyze facts (premises) → answer questions (derive conclusions)!



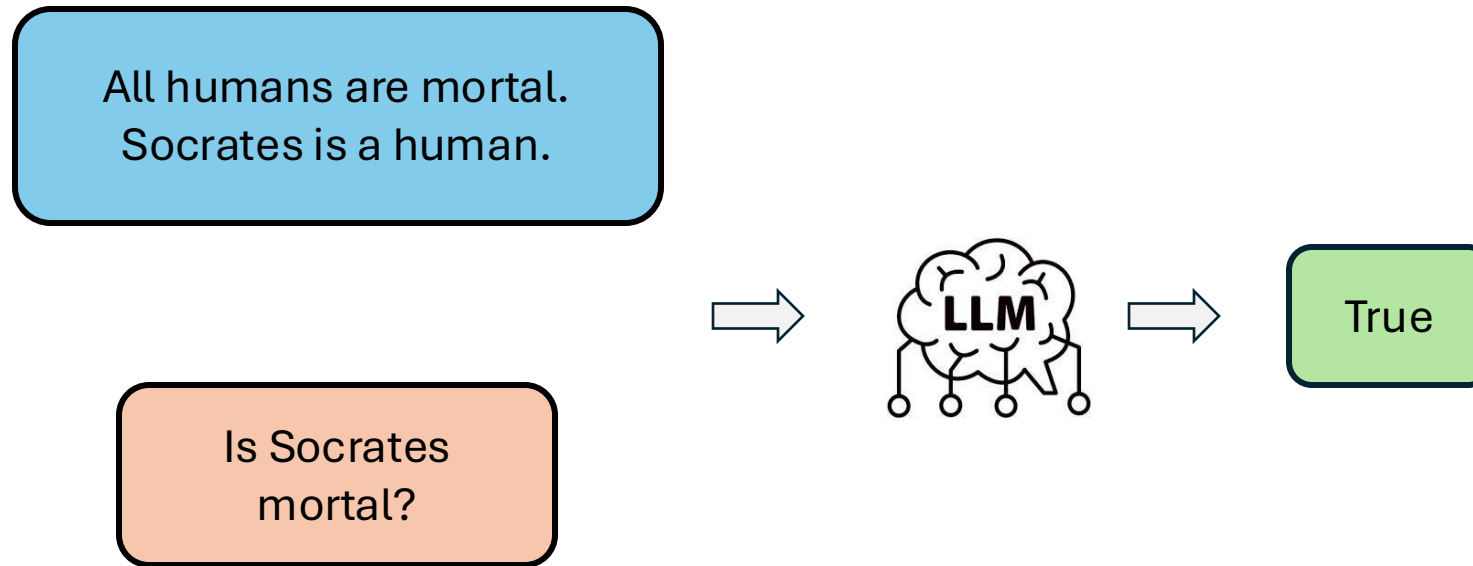
Logical Reasoning

Structured rules → analyze facts (premises) → answer questions (derive conclusions)!



If $(A \Rightarrow B)$ and A holds, then B holds (**Modus Ponens**)

Logical Reasoning in LLMs



Conclusion: Socrates is mortal.

This follows logically and necessarily from the two premises using **modus ponens** (if $A \rightarrow B$ and A , then B).



Logical Reasoning in LLMs

Those who enjoy Poetry write short verses.
Those who enjoy Novels write long stories.
Shaurya writes both short verses and long stories.

Does Shaurya
enjoy poetry?



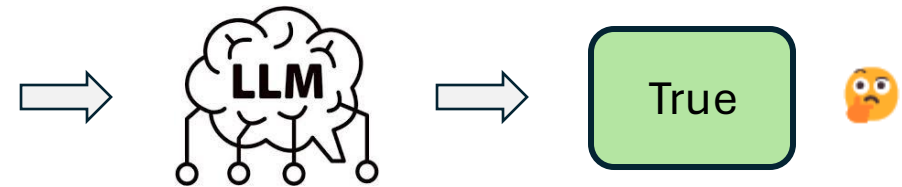
Uncertain

Though Shaurya writes
short verses, he may or
may not enjoy Poetry. Hard
to say.

Logical Reasoning in LLMs

Those who enjoy Poetry write short verses.
Those who enjoy Novels write long stories.
Shaurya writes both short verses and long stories.

Does Shaurya
enjoy poetry?



Logical Reasoning in LLMs

From the premises:

1. Those who enjoy Poetry → write short verses.
2. Those who enjoy Novels → write long stories.
3. Shaurya writes both short verses and long stories.

By premise 1, writing short verses implies Shaurya enjoys Poetry.

By premise 2, writing long stories implies Shaurya enjoys Novels.

So the truth value of *“Does Shaurya enjoy Poetry?”* is **True**.

(In fact, Shaurya enjoys **both** Poetry and Novels.)

GPT-5's Response

Logical Reasoning in LLMs

From the premises:

1. Those who enjoy Poetry \rightarrow write short verses.
2. Those who enjoy Novels \rightarrow write long stories.
3. Shaurya writes both short verses and long stories.

By premise 1, writing short verses implies Shaurya enjoys Poetry.

By premise 2, writing long stories implies Shaurya enjoys Novels.

So the truth value of *"Does Shaurya enjoy Poetry?"* is **True**.

(In fact, Shaurya enjoys **both** Poetry and Novels.)

GPT-5's Response

If ($A \Rightarrow B$)

Logical Reasoning in LLMs

From the premises:

1. Those who enjoy Poetry \rightarrow write short verses.
2. Those who enjoy Novels \rightarrow write long stories.
3. Shaurya writes both short verses and long stories.

By premise 1, writing short verses implies Shaurya enjoys Poetry.

By premise 2, writing long stories implies Shaurya enjoys Novels.

So the truth value of *"Does Shaurya enjoy Poetry?"* is **True**.

(In fact, Shaurya enjoys **both** Poetry and Novels.)

GPT-5's Response

If $(A \Rightarrow B)$ and B holds,

Logical Reasoning in LLMs

From the premises:

1. Those who enjoy Poetry \rightarrow write short verses.
2. Those who enjoy Novels \rightarrow write long stories.
3. Shaurya writes both short verses and long stories.

By premise 1 writing short verses implies Shaurya enjoys Poetry.

By premise 2, writing long stories implies Shaurya enjoys Novels.

So the truth value of *"Does Shaurya enjoy Poetry?"* is **True**.

(In fact, Shaurya enjoys **both** Poetry and Novels.)

GPT-5's Response

If $(A \Rightarrow B)$ and B holds, then A holds.

Logical Reasoning in LLMs

From the premises:

1. Those who enjoy Poetry \rightarrow write short verses.
2. Those who enjoy Novels \rightarrow write long stories.
3. Shaurya writes both short verses and long stories.

By premise 1 writing short verses implies Shaurya enjoys Poetry.

By premise 2, writing long stories implies Shaurya enjoys Novels.

So the truth value of *"Does Shaurya enjoy Poetry?"* is **True**.

(In fact, Shaurya enjoys **both** Poetry and Novels.)

GPT-5's Response

If $(A \Rightarrow B)$ and B holds, then A holds.

Modus GPT!

Logical Reasoning in LLMs

From the premises:

1. Those who enjoy Poetry ... write short verses

Affirming the consequent

🌐 24 languages ▾

Article [Talk](#)

[Read](#) [Edit](#) [View history](#) [Tools](#) ▾

From Wikipedia, the free encyclopedia

"False conversion" redirects here. For the Islamic doctrine, see [Taqiya](#).

In [propositional logic](#), **affirming the consequent** (also known as **converse error**, **fallacy of the converse**, or **confusion of necessity and sufficiency**) is a [formal fallacy](#) (or an [invalid](#) form of argument) that is committed when, in the context of an [indicative conditional](#) statement, it is stated that because the [consequent](#) is true, therefore the [antecedent](#) is true. It takes on the following form:

If *P*, then *Q*.

Q.

Therefore, *P*.

If ($A \rightarrow B$) and *B* holds, then *A* holds. — [Modus ONI](#).

Problem: LLMs are bad at Logical Reasoning

- **Unreliable:** fail on out-of-domain tasks [1]

Prove: Max is a gorpus.

Predicted answer: Max is a tumpus or a rompus or a lempus. Max is a tumpus.
Tumpuses are wumpuses. Max is a wumpus. Rompuses are gorpuses. Max is a gorpus.
Max is a gorpus.

Expected answer: Assume Max is a tumpus. Tumpuses are gorpuses. Max is a gorpus.
Assume Max is a rompus. Rompuses are gorpuses. Max is a gorpus.
Assume Max is a lempus. Lempuses are gorpuses. Max is a gorpus.
Since Max is a tumpus or a rompus or a lempus, Max is a gorpus.

FIGURE 5: Example of an incorrect proof generated by GPT-3.5 on an out-of-demonstration disjunction elimination example. The premises (axioms) are given in blue, and invalid steps are given in red. For the full example, see Figure 14 in the Appendix.

LLMs are bad at Logical Reasoning

- **Unreliable:** fail on out-of-domain tasks [1], have trouble understanding negation [2]

Prove: Max is a gorpus.

Predicted answer: Max is a tumpus or a romp.
Tumpuses are wumpuses. Max is a wumpus. R
Max is a gorpus.

Expected answer: Assume Max is a tumpus.
Assume Max is a romp. Rompuses are gorp.
Assume Max is a lempus. Lempuses are gorp.
Since Max is a tumpus or a romp or a lemp

FIGURE 5: Example of an incorrect proof generated by a model. The premises (axioms) are given in the example, see Figure 14 in the Appendix.

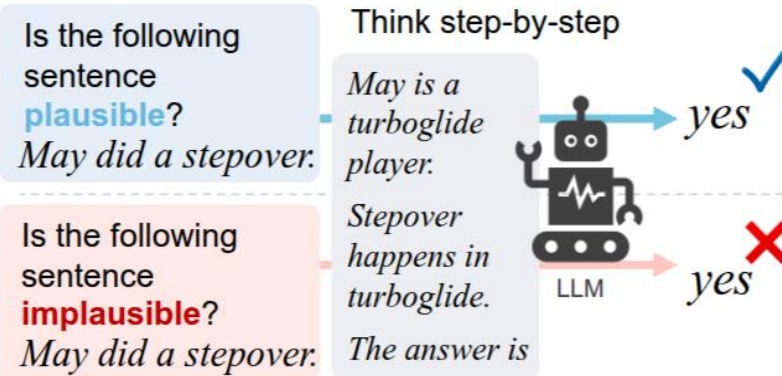


Figure 1: Overview of our experiments conducted to evaluate each model's reasoning ability against lexical negation. The model must answer *no* to the latter question about the **implausibility** of the sentence based on the valid logical flow. Here, to evaluate the robust logical

LLMs are bad at Logical Reasoning

- **Unreliable:** fail on out-of-domain tasks [1], have trouble understanding negation [2]
- **Scaling fails:** bigger models don't improve core logic (e.g., Modus Tollens) [3]

Modus Tollens:

If $(A \Rightarrow B)$ and B does not hold, then A does not hold.

Prove: Max is a gorpus.

Predicted answer: Max is a tumpus or a rompus.
Tumpuses are wumpuses. Max is a wumpus. Rompuses are gorpuses.
Max is a gorpus.

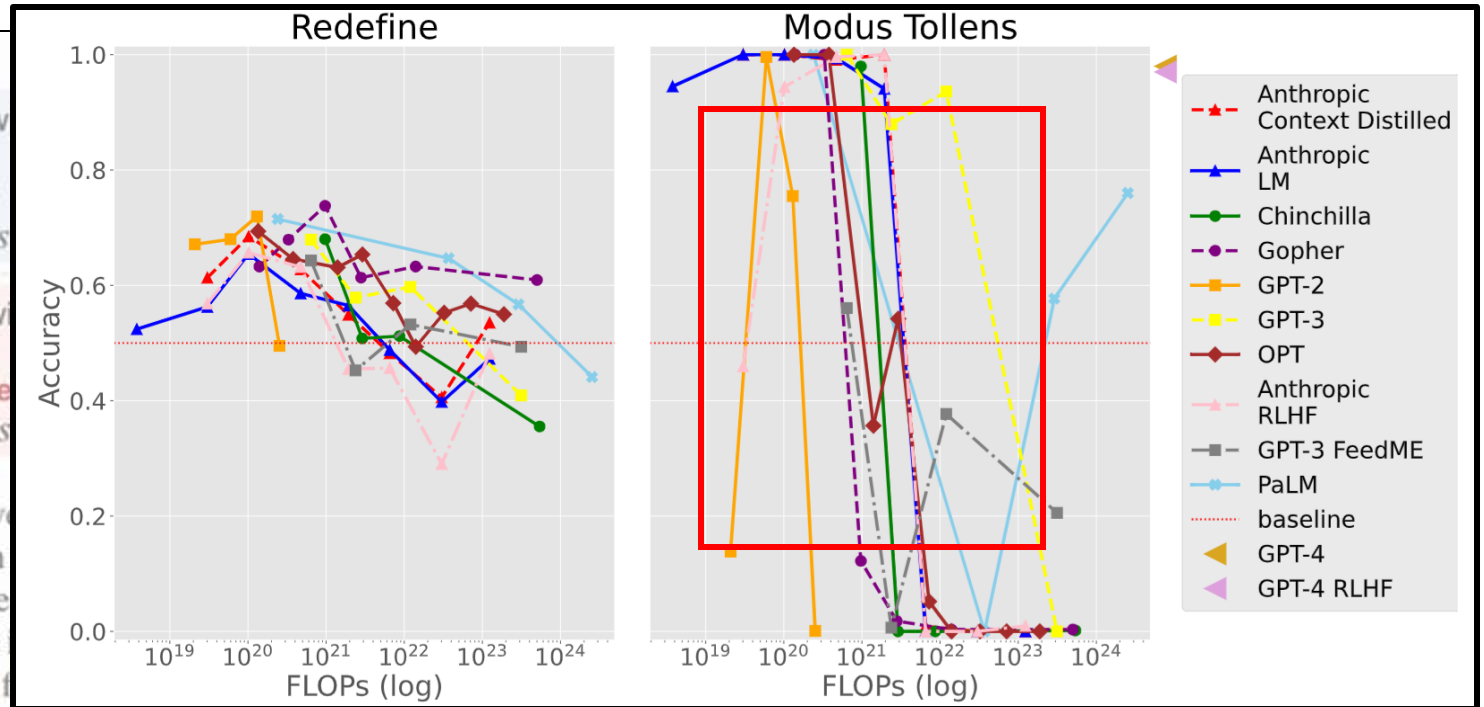
Expected answer: Assume Max is a tumpus.
Assume Max is a rompus. Rompuses are gorpuses.
Assume Max is a lempus. Lempuses are gorpuses.
Since Max is a tumpus or a rompus or a lempus, Max is a gorpus.

FIGURE 5: Example of an incorrect proof generated by a model. The premises (axioms) are given in the example, see Figure 14 in the Appendix.

Is the following sentence plausible?
May did a s

Is the following sentence implausible?
May did a s

Figure 1: Overview of the evaluation setup. We evaluate each model's performance on negation. The question about the validity of logical reasoning is given in the example.



1. Testing the General Deductive Reasoning Capacity of Large Language Models Using OOD Examples. Saparov et. Al (NeurIPS 2023)
2. Assessing Step-by-Step Reasoning against Lexical Negation: A Case Study on Syllogism. Ye et. Al (EMNLP 2023)
3. Inverse Scaling: When Bigger Isn't Better. McKenzie et. Al (TMLR 2023)

LLMs are bad at Logical Reasoning

- **Unreliable:** fail on out-of-domain tasks [1], have trouble understanding negation [2]
- **Scaling fails:** bigger models don't improve core logic (e.g., Modus Tollens) [3]
- Models lean on data patterns, not reasoning skills.

Modus Tollens:

If $(A \Rightarrow B)$ and B does not hold, then A does not hold.

Prove: Max is a gorpus.

Predicted answer: Max is a tumpus or a rompus.
Tumpuses are wumpuses. Max is a wumpus. Rompuses are gorpuses.
Max is a gorpus.

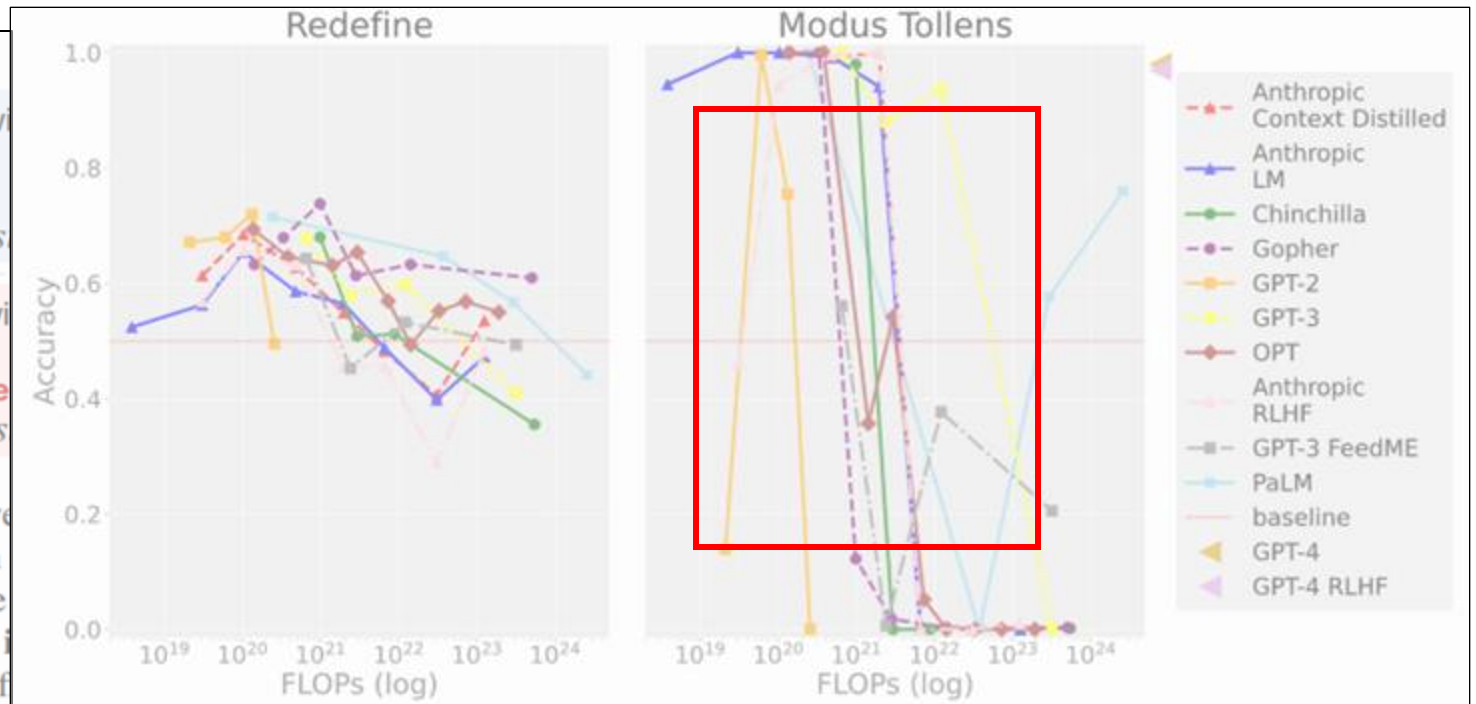
Expected answer: Assume Max is a tumpus.
Assume Max is a rompus. Rompuses are gorpuses.
Assume Max is a lempus. Lempuses are gorpuses.
Since Max is a tumpus or a rompus or a lempus, Max is a gorpus.

FIGURE 5: Example of an incorrect proof generated by a model. The premises (axioms) are given in the example, see Figure 14 in the Appendix.

Is the following sentence plausible?
May did a s

Is the following sentence implausible?
May did a s

Figure 1: Over- and under-evaluation of each sentence. The model's prediction about the validity of logical formulas is shown.



1. Testing the General Deductive Reasoning Capacity of Large Language Models Using OOD Examples. Saparov et. Al (NeurIPS 2023)

2. Assessing Step-by-Step Reasoning against Lexical Negation: A Case Study on Syllogism. Ye et. Al (EMNLP 2023)

3. Inverse Scaling: When Bigger Isn't Better. McKenzie et. Al (TMLR 2023)

Logical Reasoning is essential for AI!

- **Truthfulness in AI systems:** ensures chatbot answers follow from retrieved facts.
- **Advancing science & maths:** enables theorem proving and knowledge discovery.
- **Better education tools:** tutoring systems that teach clarity and rigor.

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology
- Datasets & Models
- Baselines
- Results
- Conclusions & Limitations

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- **Methodology**
- Datasets & Models
- Baselines
- Results
- Conclusions & Limitations

How to make LLMs reason well?

Specialized tools exist for logical reasoning!

First-order Logic

- A formal system for writing **logical statements** about the world.
- **Predicates** state properties or relations of objects.
 - Lazy(X) holds means X is lazy (e.g. Lazy(Shaurya))
 - Loves(X, Y) means X loves Y (e.g. Loves(Shaurya, Food))
- Has **\forall (for all)** and **\exists (there exists)** to talk about general rules.
 - $\forall X (\text{PhDStudent}(X) \Rightarrow \text{NeedToTakeQuals}(X))$ means "Every X who is a PhD student has to take quals."

First-order Logic

NL Description	FOL Conversion
Socrates is a human.	Human(Socrates)
All humans are mortal.	$\forall x (\text{Human}(x) \Rightarrow \text{Mortal}(x))$
Those who enjoy Poetry write short verses.	$\forall x (\text{EnjoyPoetry}(x) \Rightarrow \text{WriteShortVerses}(x))$
Shaurya writes both short verses and long stories.	$\text{WriteShortVerses}(\text{Shaurya}) \wedge \text{WriteLongStories}(\text{Shaurya})$

First-order Logic Prover

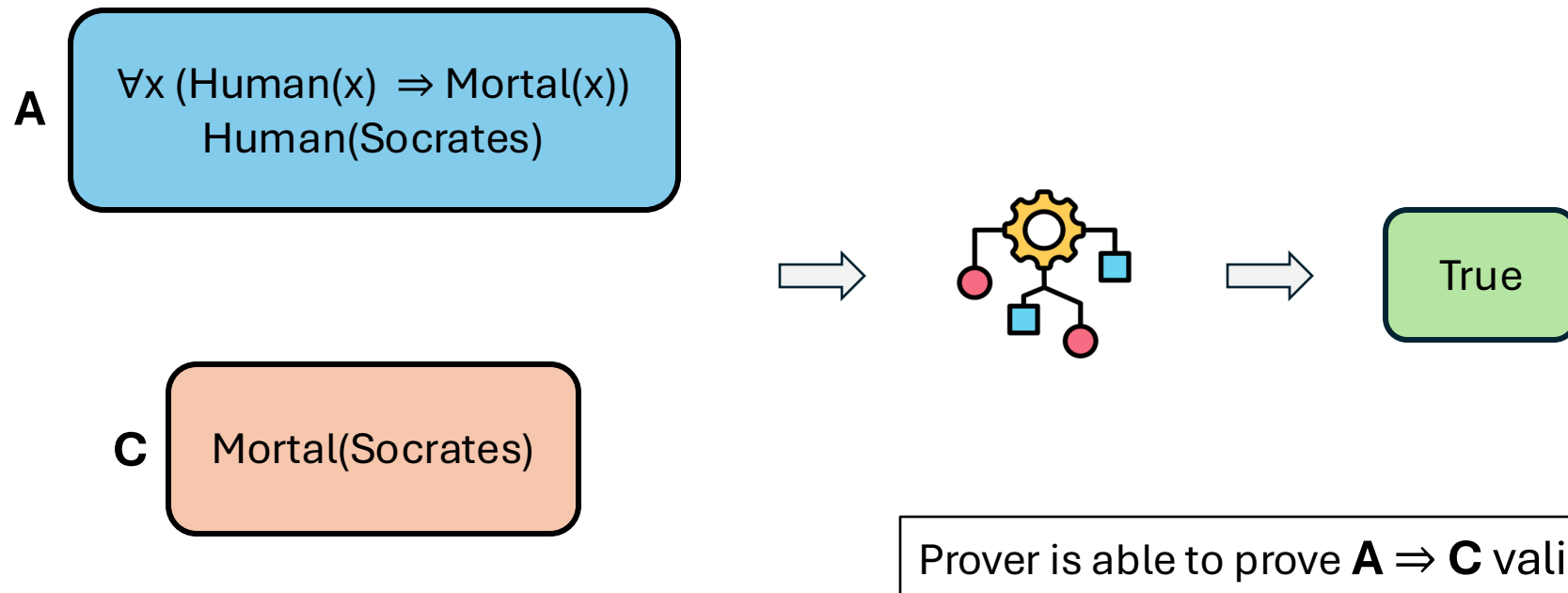
- Takes FOL premises and applies sound deduction to reason correctly.

First-order Logic Prover

- Takes FOL premises and applies sound deduction to reason correctly.
- Four possible cases:

First-order Logic Prover

- Takes FOL premises and applies sound deduction to reason correctly.
- Four possible cases: (1) True



First-order Logic Prover

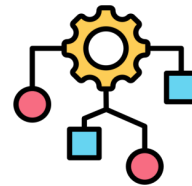
- Takes FOL premises and applies sound deduction to reason correctly.
- Four possible cases: (1) True (2) False

A

$\forall x (\text{Human}(x) \Rightarrow \text{Mortal}(x))$
 $\text{Human}(\text{Socrates})$

C

$\neg \text{Mortal}(\text{Socrates})$

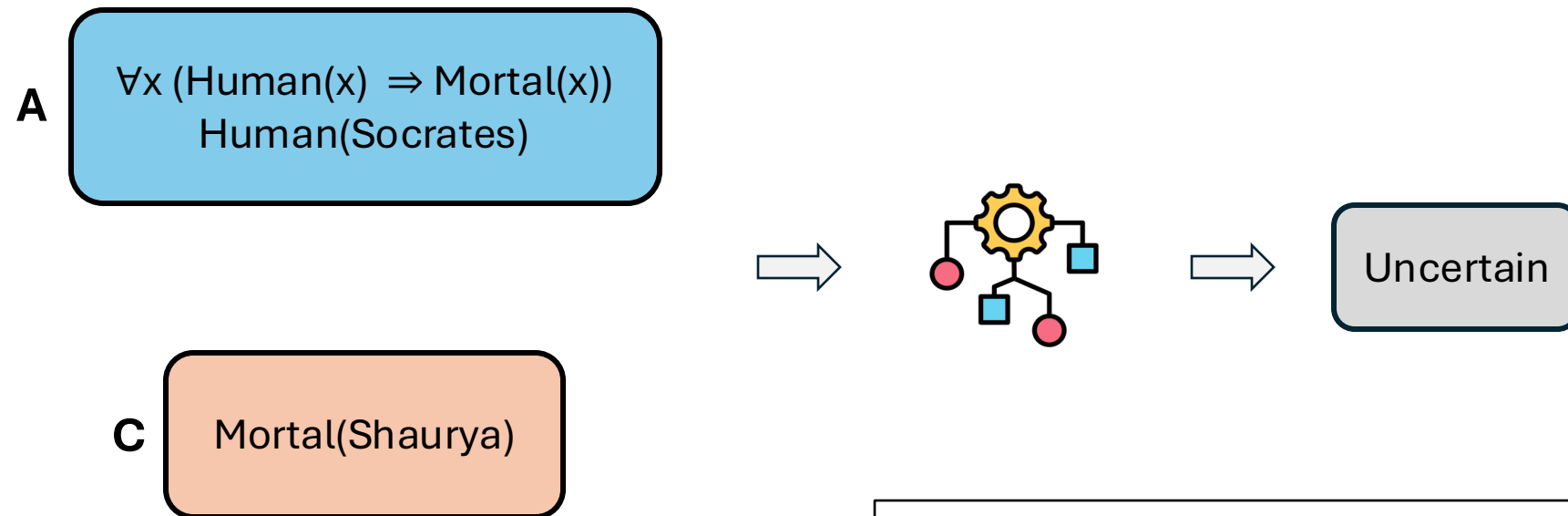


False

Prover is able to prove **A** \Rightarrow \neg **C** valid!

First-order Logic Prover

- Takes FOL premises and applies sound deduction to reason correctly.
- Four possible cases: (1) True (2) False (3) Uncertain



Prover is not able to prove **A** \Rightarrow **C** or **A** \Rightarrow \neg **C** valid!

First-order Logic Prover

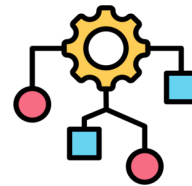
- Takes FOL premises and applies sound deduction to reason correctly.
- Four possible cases: (1) True (2) False (3) Uncertain (4) Error

A

$\forall x (\text{Human}(x) \Rightarrow \text{Mortal}(x))$
 $\text{Human}(\text{Socrates})$

C

$\text{Mortal}(\text{Socrates})$



Error

FOL Syntax errors!

How to make LLMs reason well?

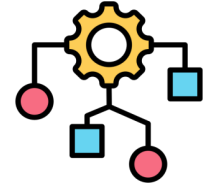
Let the cobbler stick to his last!

What are LLMs good at?



Understanding and parsing natural language,
but not guaranteed logical deduction.

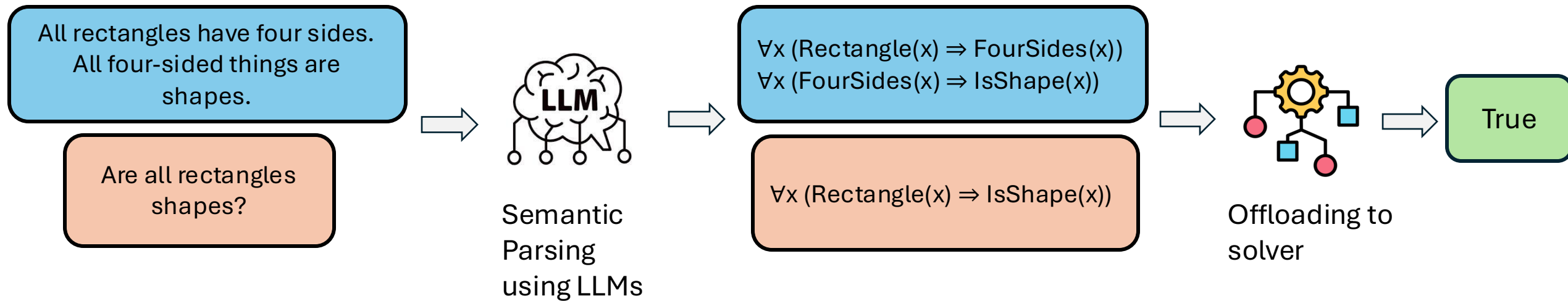
What are provers good at?



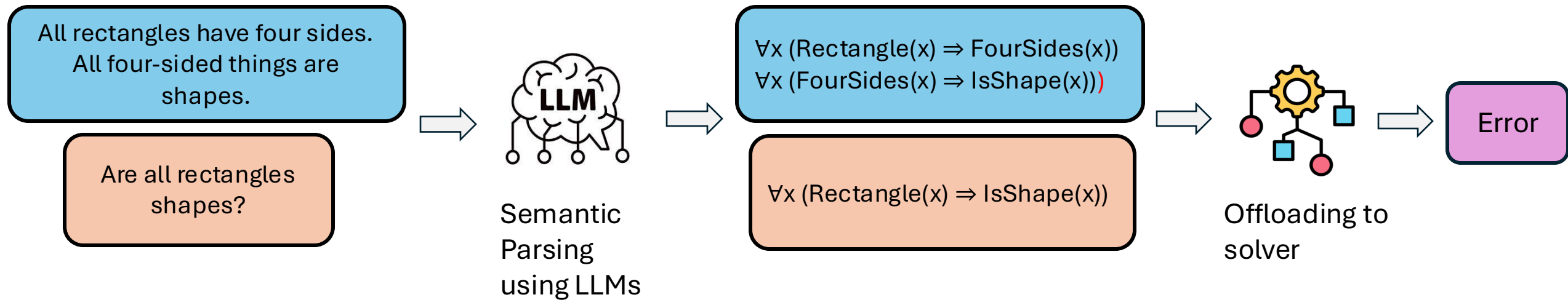
Sound logical deductions given the premises in FOL.

Instead of asking LLMs to do everything,
ask them to formalize the NL premises into FOL
and offload the reasoning to a solver.

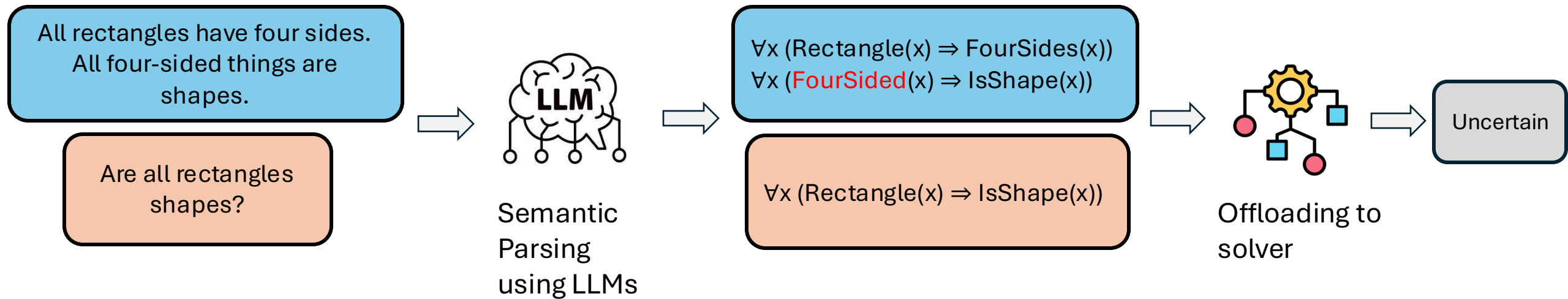
Methodology



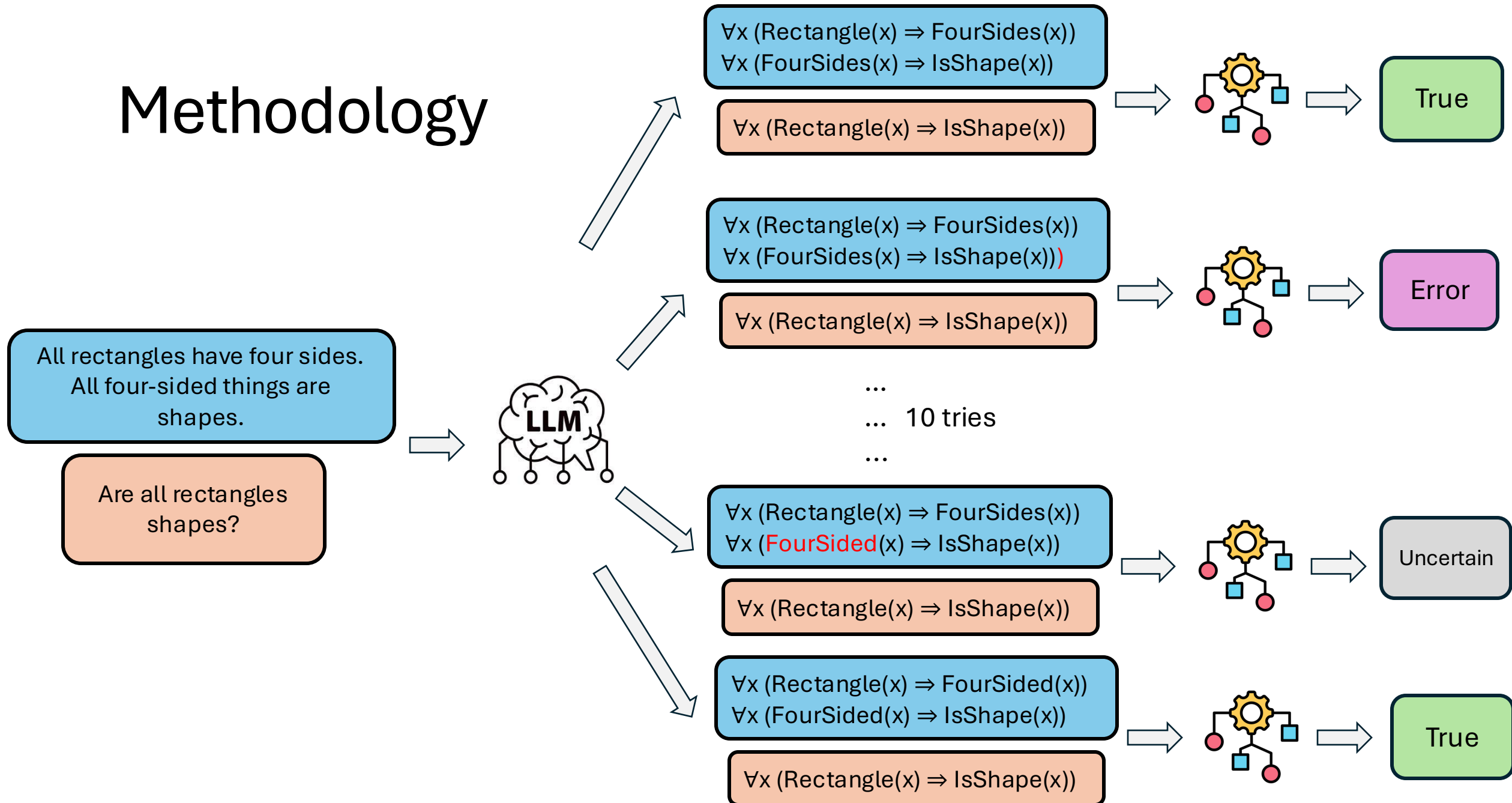
Methodology



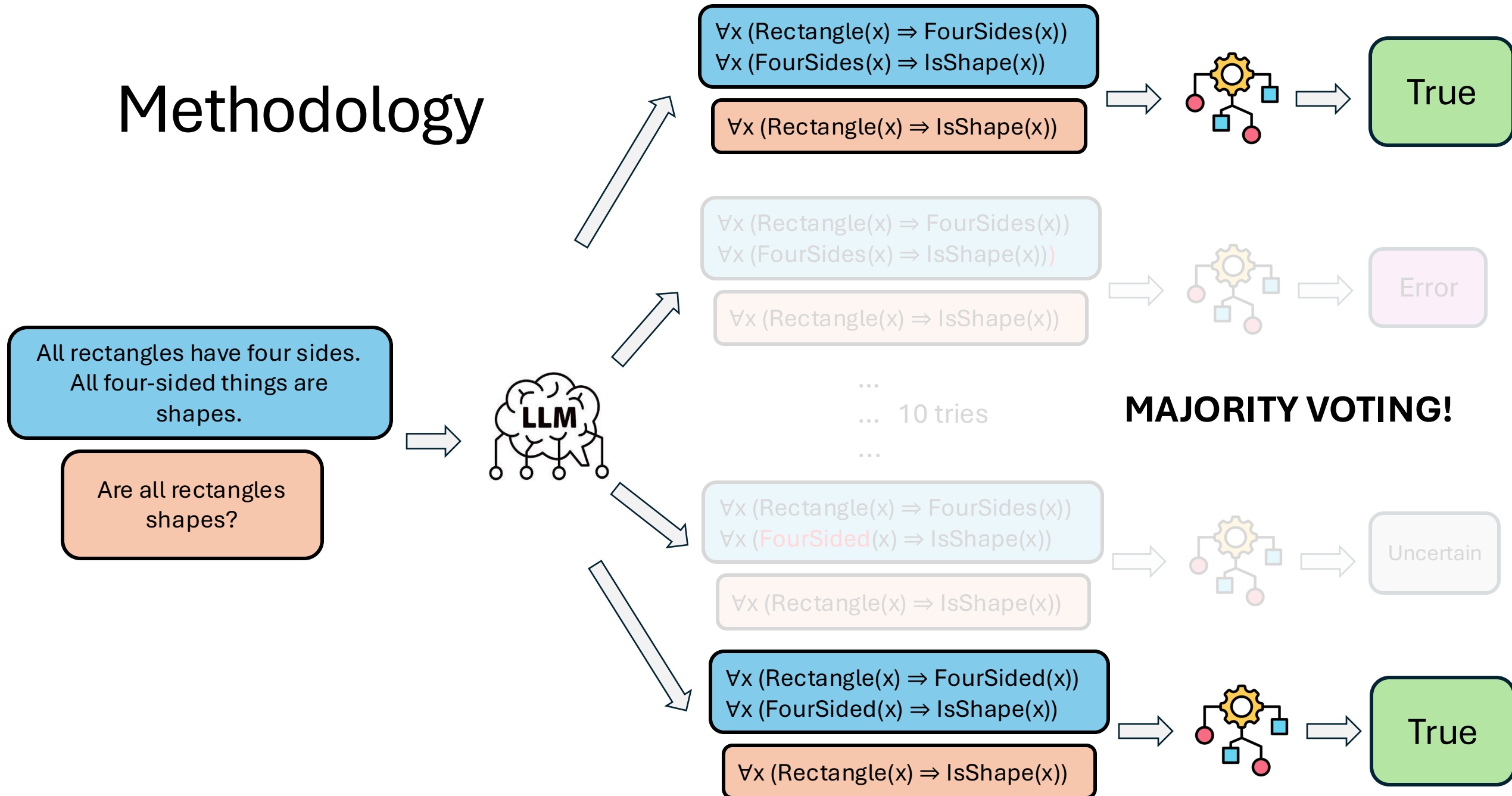
Methodology



Methodology



Methodology



Methodology

- Shifts LLM's task: Reasoning \Rightarrow Formalization in FOL
- Trade-off: NL expressiveness for syntactically strict logical formulas.
- 10-way majority procedure to mitigate formalization errors.

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models
- Baselines
- Results
- Conclusions & Limitations

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- **Datasets & Models**
- Baselines
- Results
- Conclusions & Limitations

Datasets: (1) ProofWriter

- Synthetically generated data!

The bald eagle does not eat the dog. The cat chases the dog.
The cat eats the bald eagle. The cat is nice. The cat likes the dog.
The cat likes the rabbit. The dog is furry.
The rabbit chases the bald eagle. The rabbit eats the bald eagle.
If someone does not eat the cat then they do not eat the dog.
If someone likes the bald eagle then they do not like the rabbit.
If someone eats the bald eagle and they do not eat the rabbit
then they are furry.
If someone is furry then they like the cat.
Q1. The bald eagle likes the cat. True/false? **[F]**
Q2. The rabbit likes the cat. True/false? **[T]**
Q3. The bald eagle is furry. True/false? **[F]**

Datasets: (1) ProofWriter

- Synthetically generated data!
- Fixed Rules:
 1. $\text{is}(X, Y)$

The bald eagle does not eat the dog. The cat chases the dog.
The cat eats the bald eagle. The cat is nice. The cat likes the dog.
The cat likes the rabbit. The dog is furry.
The rabbit chases the bald eagle. The rabbit eats the bald eagle.
If someone does not eat the cat then they do not eat the dog.
If someone likes the bald eagle then they do not like the rabbit.
If someone eats the bald eagle and they do not eat the rabbit
then they are furry.
If someone is furry then they like the cat.
Q1. The bald eagle likes the cat. True/false? [F]
Q2. The rabbit likes the cat. True/false? [T]
Q3. The bald eagle is furry. True/false? [F]

Datasets: (1) ProofWriter

- Synthetically generated data!
- Fixed Rules:
 1. `is(X, Y)`
 2. `verb(X, Y) [Likes(Cat, Dog)]`

The bald eagle does not eat the dog. The cat chases the dog.
The cat eats the bald eagle. The cat is nice. The cat likes the dog.
The cat likes the rabbit. The dog is furry.
The rabbit chases the bald eagle. The rabbit eats the bald eagle.
If someone does not eat the cat then they do not eat the dog.
If someone likes the bald eagle then they do not like the rabbit.
If someone eats the bald eagle and they do not eat the rabbit
then they are furry.
If someone is furry then they like the cat.
Q1. The bald eagle likes the cat. True/false? [F]
Q2. The rabbit likes the cat. True/false? [T]
Q3. The bald eagle is furry. True/false? [F]

Datasets: (1) ProofWriter

- Synthetically generated data!
- Fixed Rules:
 1. `is(X, Y)`
 2. `verb(X, Y) [Likes(Cat, Dog)]`
 3. `[C1 and C2 and ..] => C`

The bald eagle does not eat the dog. The cat chases the dog.
The cat eats the bald eagle. The cat is nice. The cat likes the dog.
The cat likes the rabbit. The dog is furry.
The rabbit chases the bald eagle. The rabbit eats the bald eagle.

If someone does not eat the cat then they do not eat the dog.
If someone likes the bald eagle then they do not like the rabbit.
If someone eats the bald eagle and they do not eat the rabbit
then they are furry.
If someone is furry then they like the cat.

Q1. The bald eagle likes the cat. True/false? **[F]**

Q2. The rabbit likes the cat. True/false? **[T]**

Q3. The bald eagle is furry. True/false? **[F]**

Datasets: (1) ProofWriter

- Synthetically generated data!
- Fixed Rules:
 1. $\text{is}(X, Y)$
 2. $\text{verb}(X, Y) [\text{Likes}(\text{Cat}, \text{Dog})]$
 3. $[\text{C1 and C2 and ..}] \Rightarrow \text{C}$
- Makes formalization task easier?

The bald eagle does not eat the dog. The cat chases the dog.
The cat eats the bald eagle. The cat is nice. The cat likes the dog.
The cat likes the rabbit. The dog is furry.
The rabbit chases the bald eagle. The rabbit eats the bald eagle.
If someone does not eat the cat then they do not eat the dog.
If someone likes the bald eagle then they do not like the rabbit.
If someone eats the bald eagle and they do not eat the rabbit
then they are furry.
If someone is furry then they like the cat.
Q1. The bald eagle likes the cat. True/false? **[F]**
Q2. The rabbit likes the cat. True/false? **[T]**
Q3. The bald eagle is furry. True/false? **[F]**

Datasets: (2) FOLIO

- Expert-written data!

NL Premises

1. Some employees good at time management do not exercise every week.
2. All employees good at time management are efficient in dealing with daily work.
3. All employees efficient in dealing with daily work perform better than others.
4. All employees who perform better than others have more opportunities to get a promotion.
5. James does not have more opportunities to get a promotion.

NL Conclusions

- A. James exercises every week.
 - B. James exercises every week and is good at time management.
 - C. If James does not perform better than others, then he exercises every week and is good at time management.
-

Datasets: (2) FOLIO

- Expert-written data!
- Less number of premises than ProofWriter (5 vs 19), but complex!

NL Premises

1. Some employees good at time management do not exercise every week.
2. All employees good at time management are efficient in dealing with daily work.
3. All employees efficient in dealing with daily work perform better than others.
4. All employees who perform better than others have more opportunities to get a promotion.
5. James does not have more opportunities to get a promotion

NL Conclusions

- A. James exercises every week.
 - B. James exercises every week and is good at time management.
 - C. If James does not perform better than others, then he exercises every week and is good at time management.
-

Datasets: (2) FOLIO

- Expert-written data!
- Less number of premises than ProofWriter (5 vs 19), but complex!
- Also provides FOL translations for these premises and conclusions!

NL Premises

1. Some employees good at time management do not exercise every week.
2. All employees good at time management are efficient in dealing with daily work.
3. All employees efficient in dealing with daily work perform better than others.
4. All employees who perform better than others have more opportunities to get a promotion.
5. James does not have more opportunities to get a promotion.

NL Conclusions

- A. James exercises every week.
 - B. James exercises every week and is good at time management.
 - C. If James does not perform better than others, then he exercises every week and is good at time management.
-

Models used in experiments

- GPT-3.5 [1] and GPT-4 [2]



- StarCoder+ [3]



- Free, Open Access
- Smaller (15B params) than the GPT models (175B+ models)
- Allows dataset search: wasn't trained on FOLIO or ProofWriter!

Models used in experiments

- GPT-3.5 [1] and GPT-4 [2]



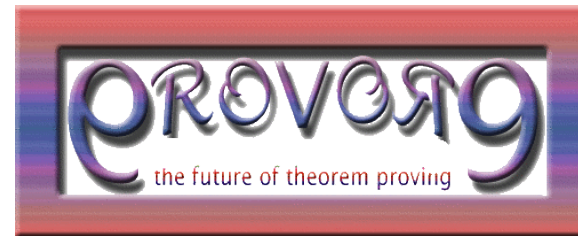
- StarCoder+ [3]



- Free, Open Access
- Smaller (15B params) than the GPT models (175B+ models)
- Allows dataset search: wasn't trained on FOLIO or ProofWriter!

FOL Prover used

- Prover9[4]: Automated FOL prover



Contents

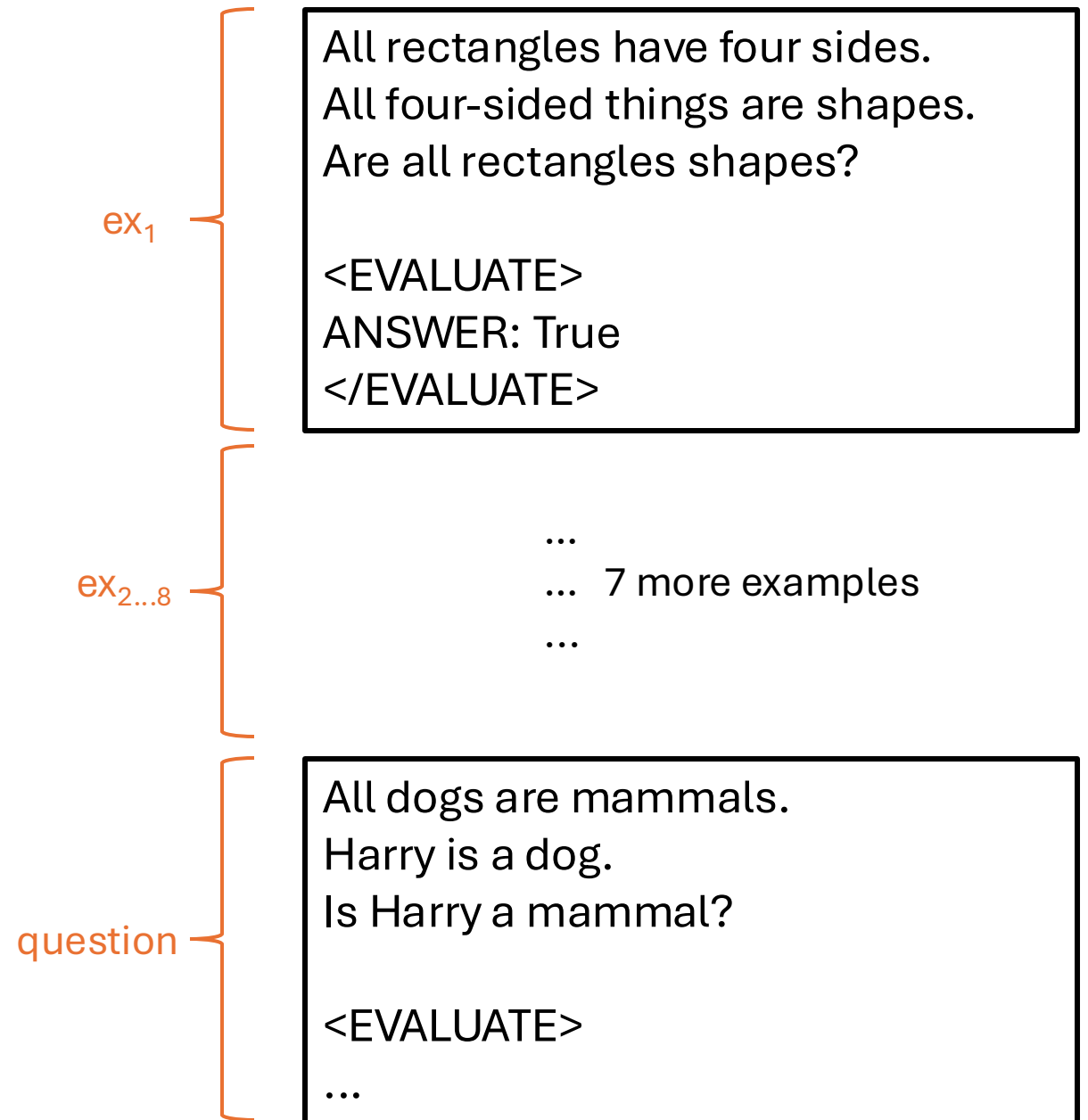
- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- **Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)**
- Baselines
- Results & Observations
- Conclusions & Limitations

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- **Baselines**
- Results
- Conclusions & Limitations

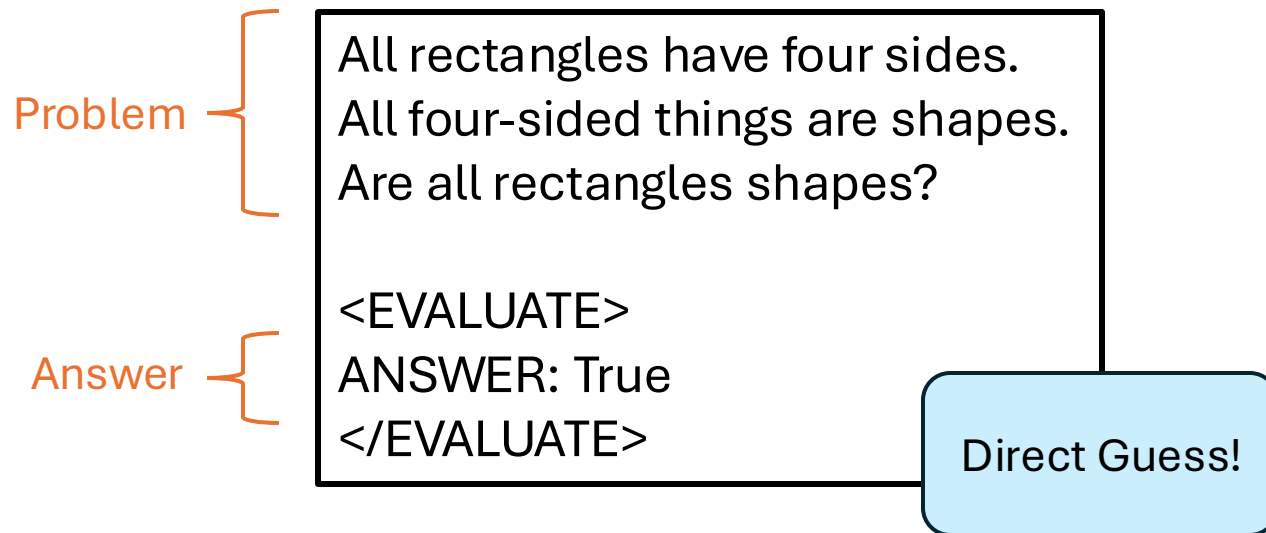
Baselines

- **Common setup:** ICL with 8 fixed **FOLIO** examples
- **Variation:** 3 baselines differ in **content of examples**
- **Decoding:** 10 generations (**T = 0.8**) → majority-vote



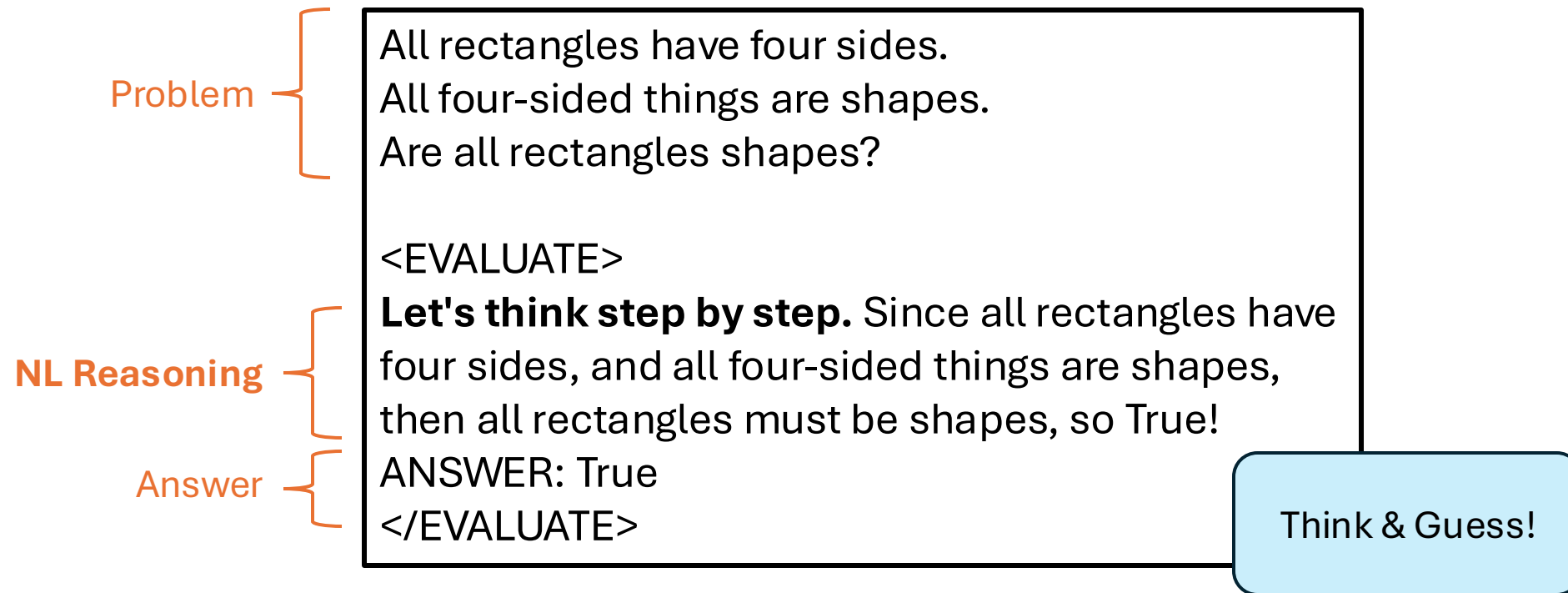
Baselines: (1) Naive

- Examples contain: Problem + Answer



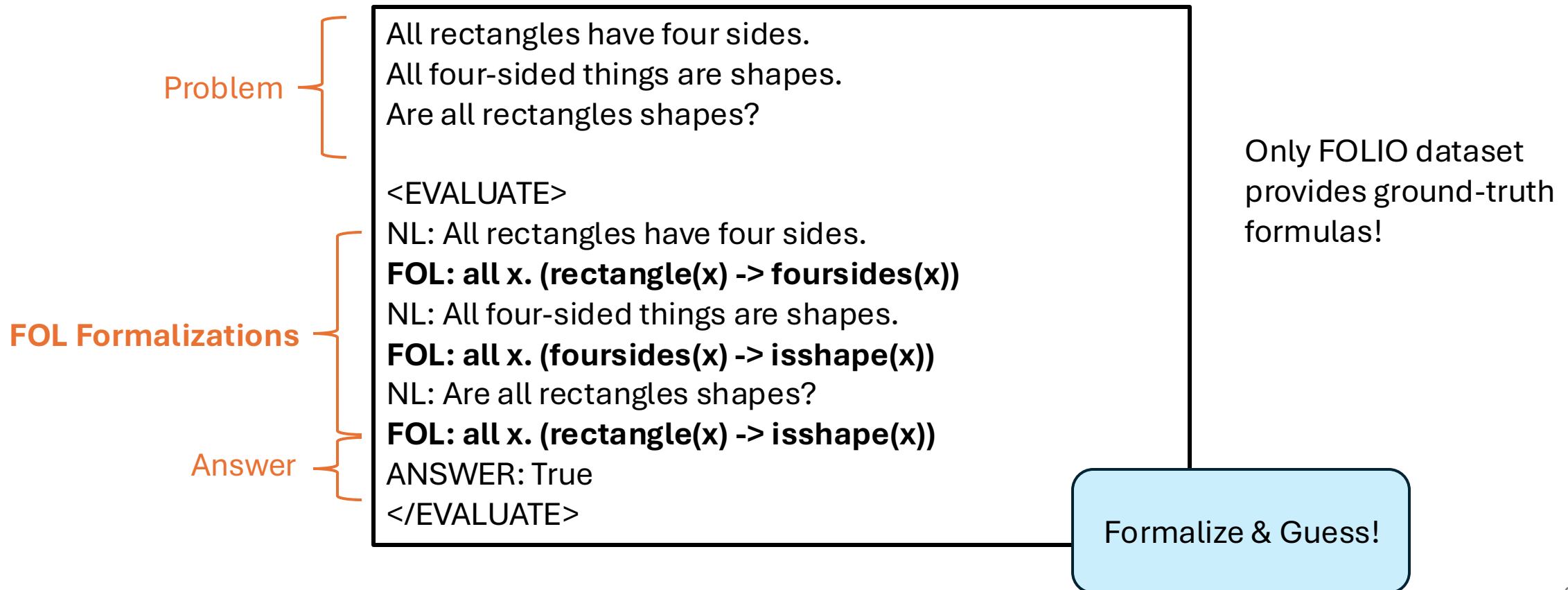
Baselines: (2) COT

- Examples contain: Problem + **NL Reasoning** + Answer



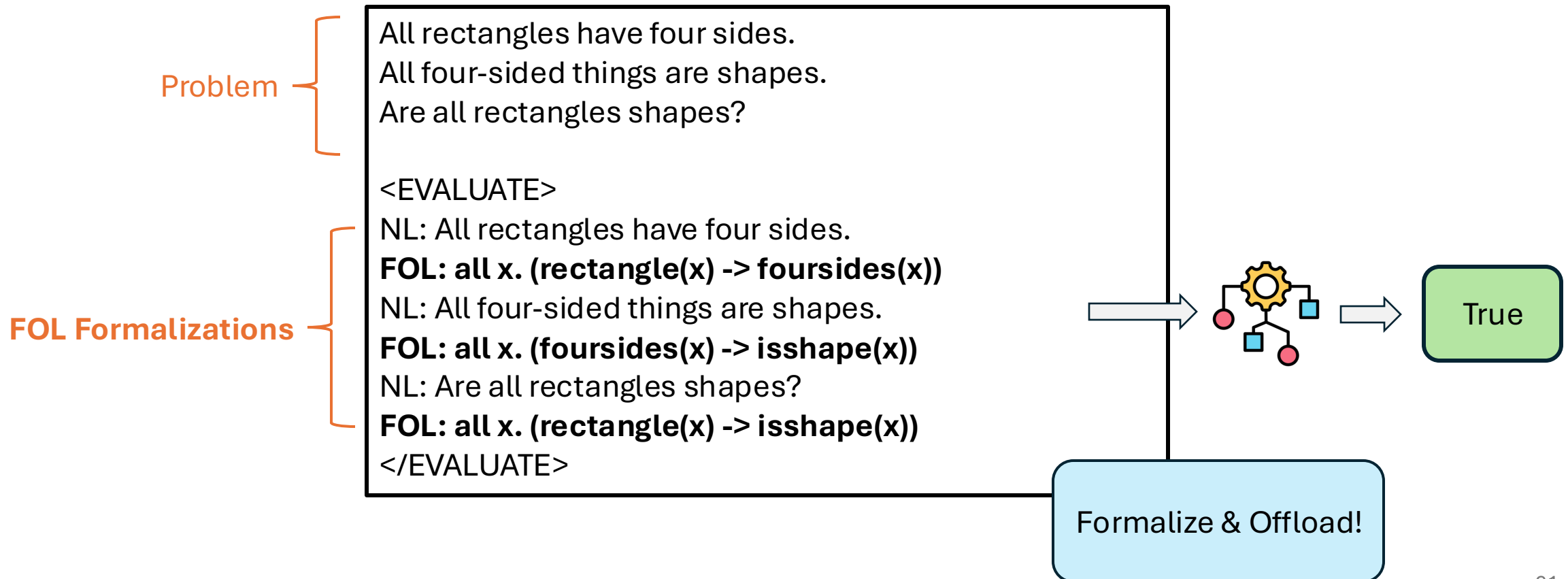
Baselines: (3) Scratchpad

- Examples contain: Problem + **FOL Formalizations** + Answer



Our Approach: LINC

- Examples contain: Problem + **FOL Formalizations**



Contents

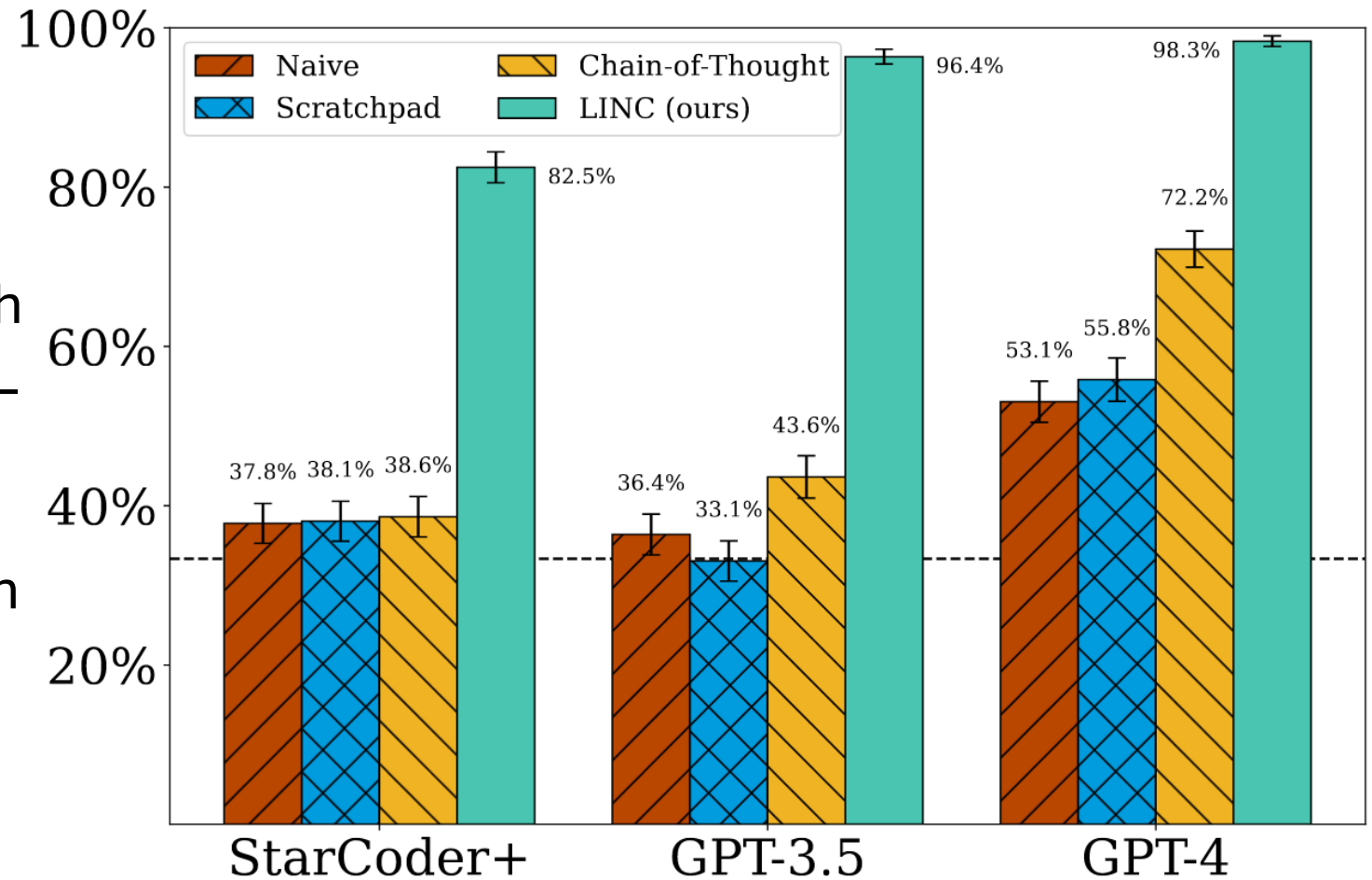
- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- **Baselines: Naïve, COT and ScratchPad**
- Results
- Conclusions & Limitations

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- Baselines: Naïve, COT and ScratchPad
- **Results**
 - **ProofWriter**
 - FOLIO
- Conclusions & Limitations

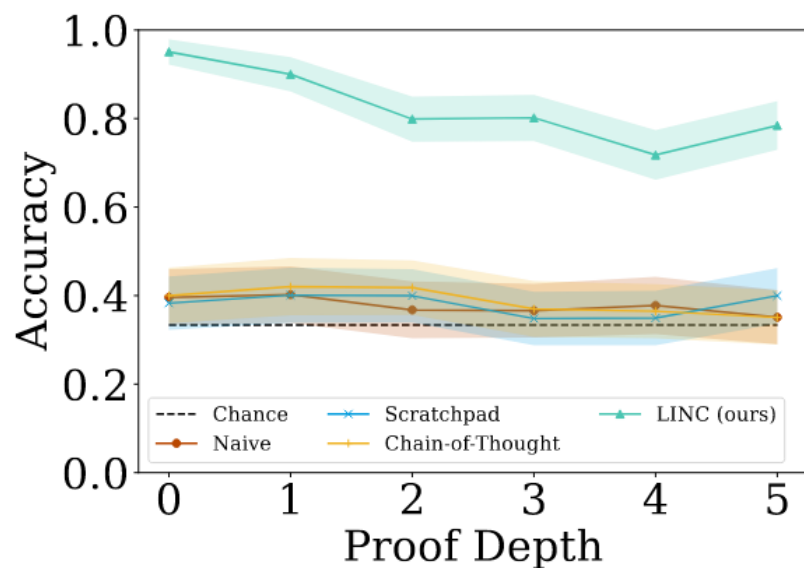
Results: ProofWriter (Accuracy)

- LINC comfortably beats all the baselines!
- Models formalize well even with more premises than seen in ICL examples.
- Formalization alone not enough as Scratchpad stays low.

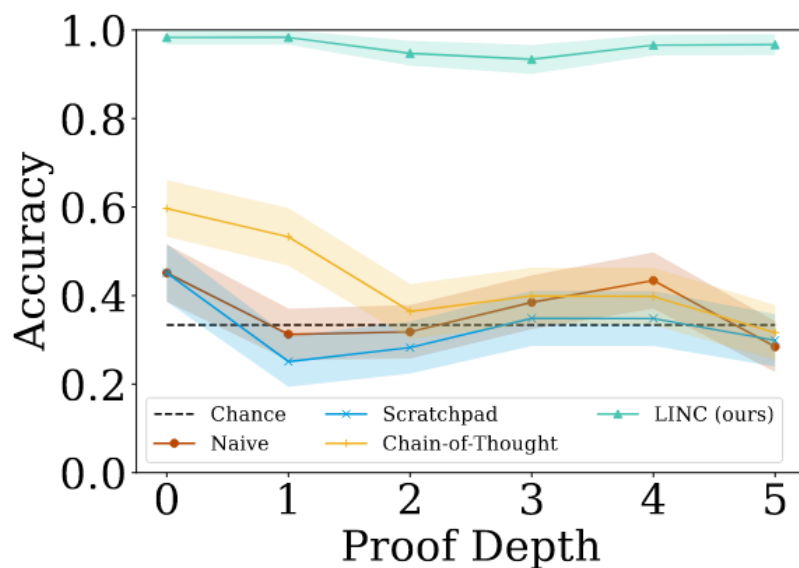


Results: ProofWriter (Accuracy vs Proof Depths)

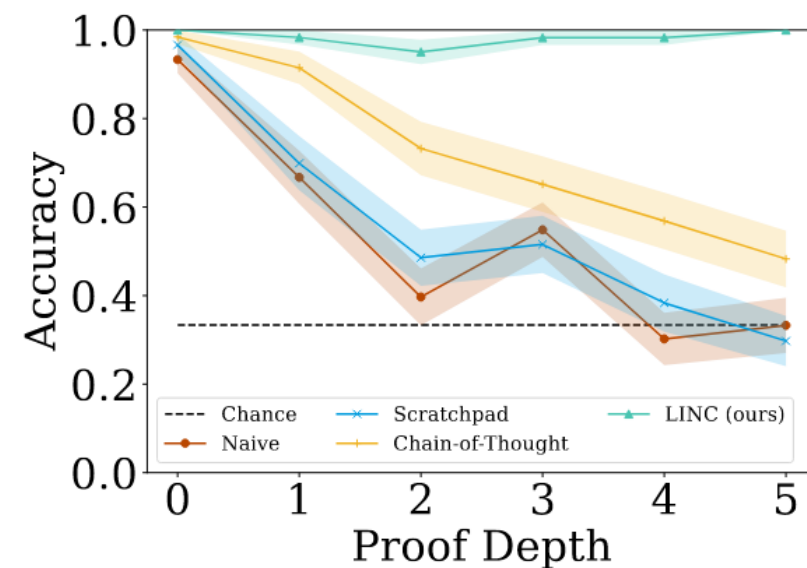
- Proof depth: Number of reasoning steps needed.
- LINC remains strong as proof depths increase (thanks to the solver)!



(a) StarCoder+.



(b) GPT-3.5.



(c) GPT-4.

Contents

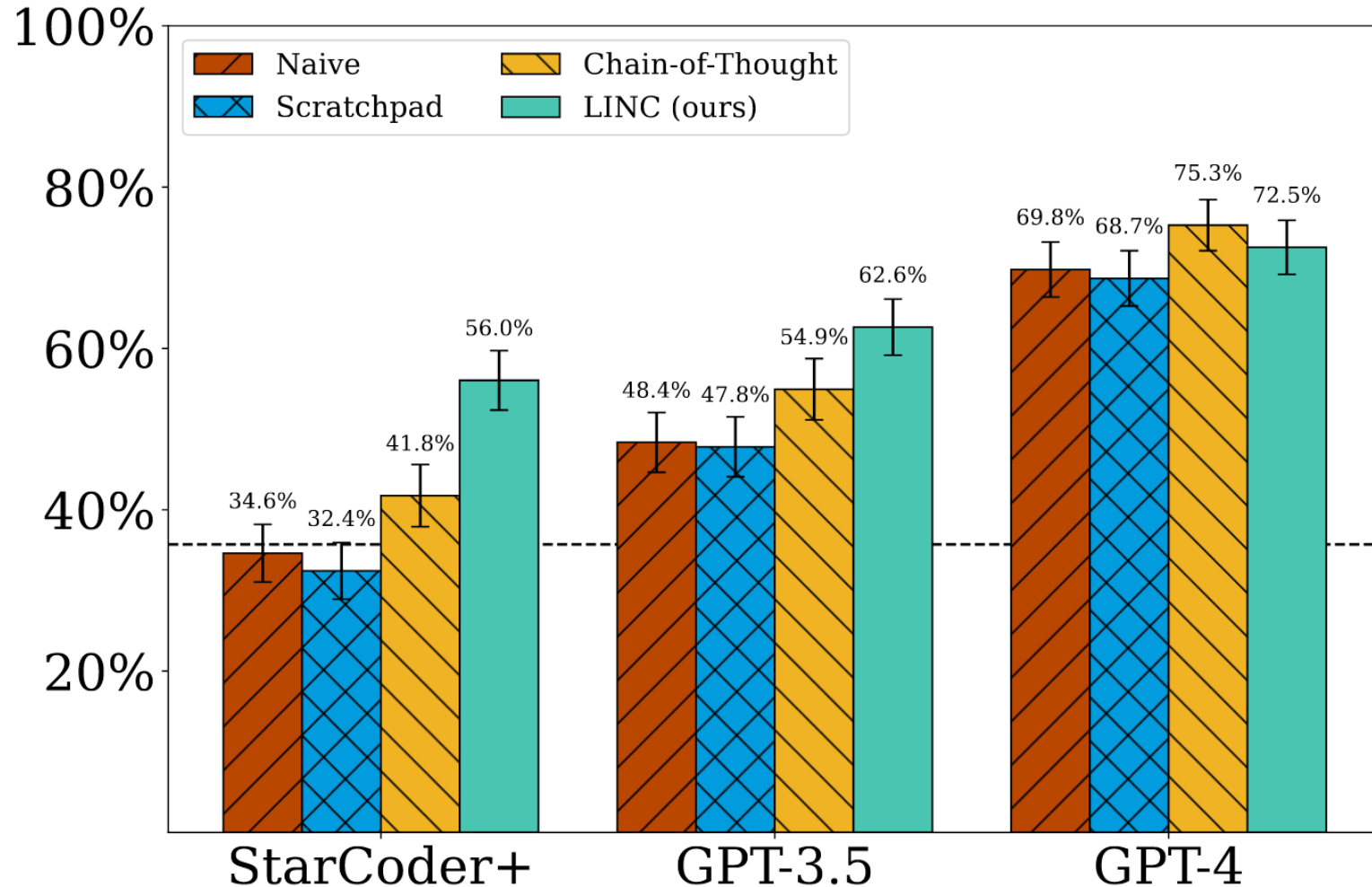
- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- Baselines: Naïve, COT and ScratchPad
- **Results**
 - ProofWriter: LINC clearly wins (as formalization task was easier)
 - FOLIO
- Conclusions & Limitations

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- Baselines: Naïve, COT and ScratchPad
- **Results**
 - ProofWriter: LINC clear winner (as formalization task was easier)
 - **FOLIO**
- Conclusions & Limitations

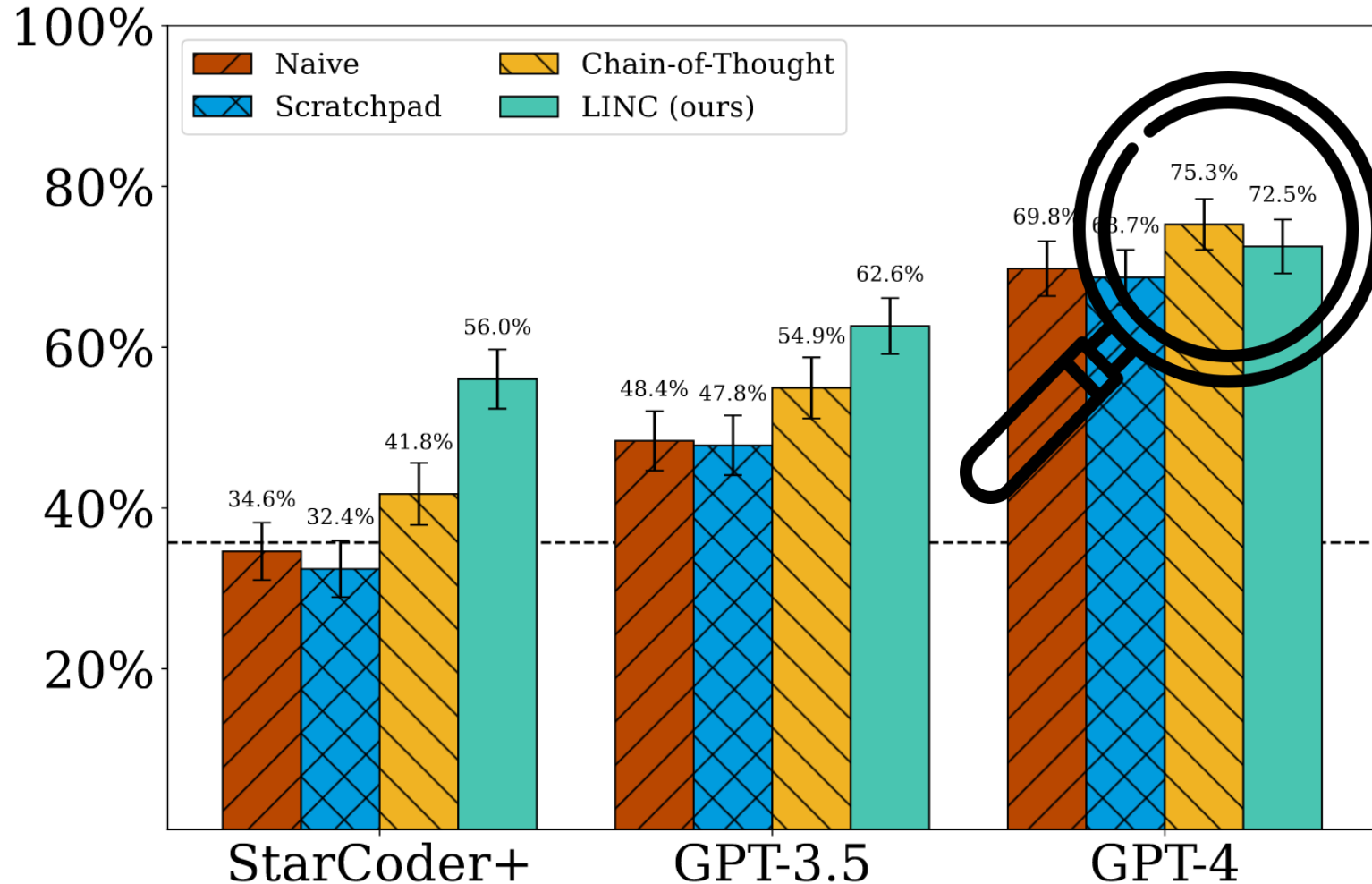
Results: FOLIO (Accuracy)

- LINC leads to some gains for StarCoder+ and GPT3.
- However, GPT-4 with COT performs better than LINC!
- FOLIO has more complicated premises (hard to formalize)!



Results: FOLIO (Accuracy)

- LINC leads to some gains for StarCoder+ and GPT3.
- **However, GPT-4 with COT performs better than LINC!**
- FOLIO has more complicated premises (hard to formalize)!



Results: FOLIO (3 LINC Failure Modes)

1. Fails to capture implicit information.

Person(Harry) & Book(Walden) missing!

Premise 1: When a person reads a book, that person gains knowledge.

FOL: all x. all y. (Person(x) & Reads(x, y) & Book(y) -> Gains(x, Knowledge))

Premise 2: Harry read the book "Walden" by Henry Thoreau.

FOL: Reads(Harry, Walden)

Conclusion (Prover9: Uncertain): Harry gains knowledge.

FOL: Gains(Harry, Knowledge)

Results: FOLIO (3 LINC Failure Modes)

1. Fails to capture implicit information.
2. Fails to capture explicit information (choice of representation).

Person(Harry) & Book(Walden)

Premise 1: When a person reads a book, they gain knowledge.

FOL: all x. all y. (Person(x) & Book(y) -> Gains(x, Knowledge))

Premise 2: Harry read the book Walden.

FOL: Reads(Harry, Walden)

Conclusion (Prover9: Uncertain) Harry gained knowledge.

FOL: Gains(Harry, Knowledge)

Premises:

All squares have four sides.

All four-sided things are shapes.

Conclusion (Ground Truth: True):

All squares are shapes.

FOLs:

all x. (Square(x) -> FourSided(x))

all x. (FourSided(x) -> Shape(x))

Conclusion (Prover9 Deduction: Uncertain)

all x. (Square(x) -> Shape(x))

Results: FOLIO (3 LINC Failure Modes)

1. Fails to capture implicit information.
2. Fails to capture explicit information (choice of representation).
3. FOL contains syntax error.

Person(Harry) & Book(Walden)

Premise 1: When a person reads a book, they gain knowledge.

FOL: all x. all y. (Person(x) & Book(y) -> Gains(x, Knowledge))

Premise 2: Harry read the book Walden.

FOL: Reads(Harry, Walden)

Conclusion (Prover9: Uncertain) Harry gained knowledge.

FOL: Gains(Harry, Knowledge)

Premises:

All squares have four sides
All four-sided things are squares

Conclusion (Group Theory)
All squares are squares

FOLs:

all x. (Square(x) -> FourSided(x))
all x. (FourSided(x) -> Shape(x))

Conclusion (Prover9 Deduction: Uncertain)
all x. (Square(x) -> Shape(x))

FOLs:

exists x. (Badults(x) & BritishSitcom(x) &
 -> StarringMembersOf(x, Pappys))
Piloted(Badults, July2013, BBCThree)

Prover9 Error:

The following symbols are used with multiple
 -> arities: Badults/1, Badults/0.

Badults used both as
constant and predicate!

Results: FOLIO (3 COT Failure Modes)

1. COT concludes something different than it suggests ("No reason to believe X" => "X is false").

(GPT-4 Incorrect CoT): Let's think step by step.

- ↪ We want to evaluate if Space is a
- ↪ vampire. We know that Space is a vacuum
- ↪ and that vacuums suck. However, while it'
- ↪ s true that vampires also suck, this
- ↪ doesn't mean that everything that sucks
- ↪ is a vampire. Therefore, there is no
- ↪ reason to believe that Space is a vampire
- ↪ based on these premises, so the
- ↪ statement is false.

ANSWER: False

Results: FOLIO (3 COT Failure Modes)

1. COT concludes something different than it suggests ("No reason to believe X" => "X is false").
2. COT makes incorrect logical deductions.

(GPT-4 Incorrect CoT): Let's think step by step. We want to evaluate if Mike has an F1 visa. From the premises, we know that Mike is an international student and needs to apply for CPT if he wants to work in the US. We also know that an international student in the US with an F1 visa needs to apply for CPT or OPT if the student wants to work in the US. So, it is clear that Mike has an F1 visa since he fits the criteria. Therefore, the conclusion is true.

If on F1, then need CPT.
Mike needs CPT.
So, he must be on F1.
(Affirming the consequent!)

(GPT-4 Incorrect CoT):

→ We want to evaluate if Space is a vampire.
→ vampire is a creature that lives in the shadows and feeds on the blood of the living.
→ and that Space is a creature that lives in the shadows and feeds on the blood of the living.
→ s true.
→ doesn't mean that everything that sucks blood is a vampire. Therefore, there is no reason to believe that Space is a vampire based on these premises, so the statement is false.

ANSWER: False

Results: FOLIO (3 COT Failure Modes)

1. COT concludes something different than it suggests ("No reason to believe X" => "X is false").
2. COT makes incorrect logical deductions.
3. COT fails to find complex paths of reasoning.

(GPT-4 Incorrect CoT): Let's think step by step to evaluate if Mike has an F1 visa. From the information provided, we know that Mike is an international student who can apply for CPT if he wants to work in the US. **that an international student in the US with a valid F1 visa can apply for CPT or OPT if the student wants to work in the US. So, it is clear that Mike has an F1 visa and meets the criteria. Therefore, the conclusion is true.**

Premises:
If something can transport multiple passengers,
 ↪ then they are not empty.
All airlines can transport multiple passengers.
All planes are empty.
All Boeing 707s are planes.
If Greyhound is not an airline, then there are
 ↪ no Greyhound planes.

Conclusion (Ground Truth: False)
A Greyhound is a Boeing 707.

(GPT-4 Incorrect CoT, 1): Let's think step by step. We want to evaluate if a Greyhound is a Boeing 707. Looking at our premises, we have no direct or indirect information linking Greyhounds to Boeing 707s. Given the information available, we cannot make any certain conclusion about a Greyhound being a Boeing 707. Therefore, it is uncertain whether a Greyhound is a Boeing 707.

6 steps of reasoning needed!

Greyhound → Boeing 707

↓

Plane

↓

Empty

↓

Cannot transport passengers

↓

Not an airline

↓

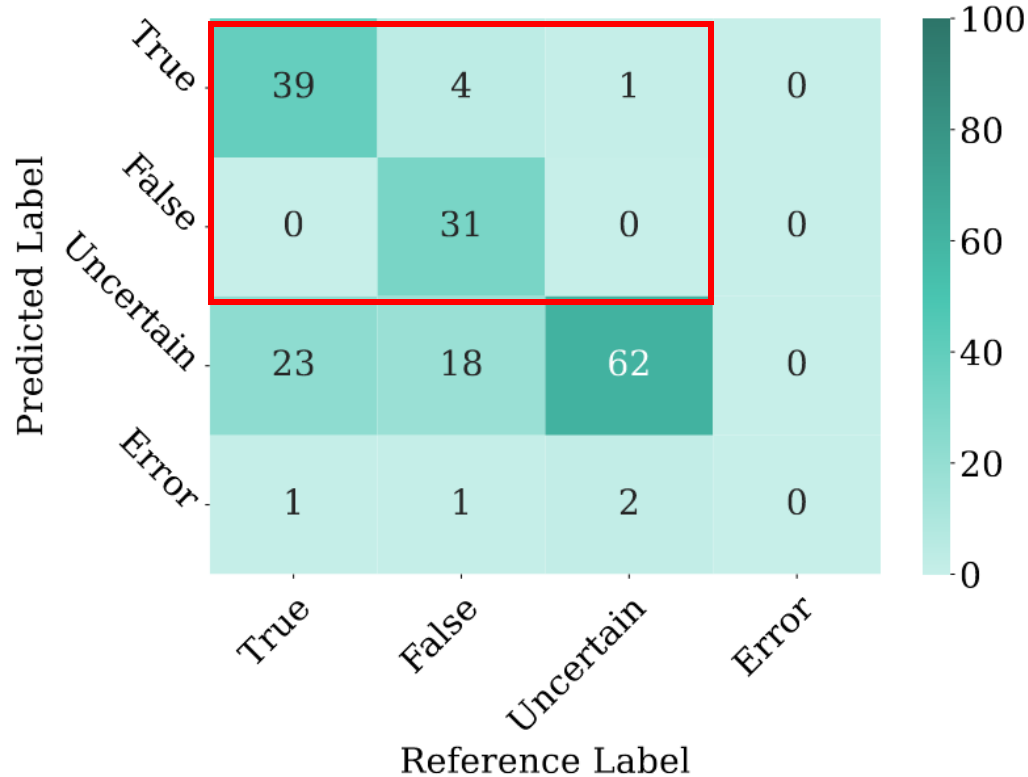
No Greyhound planes

↓

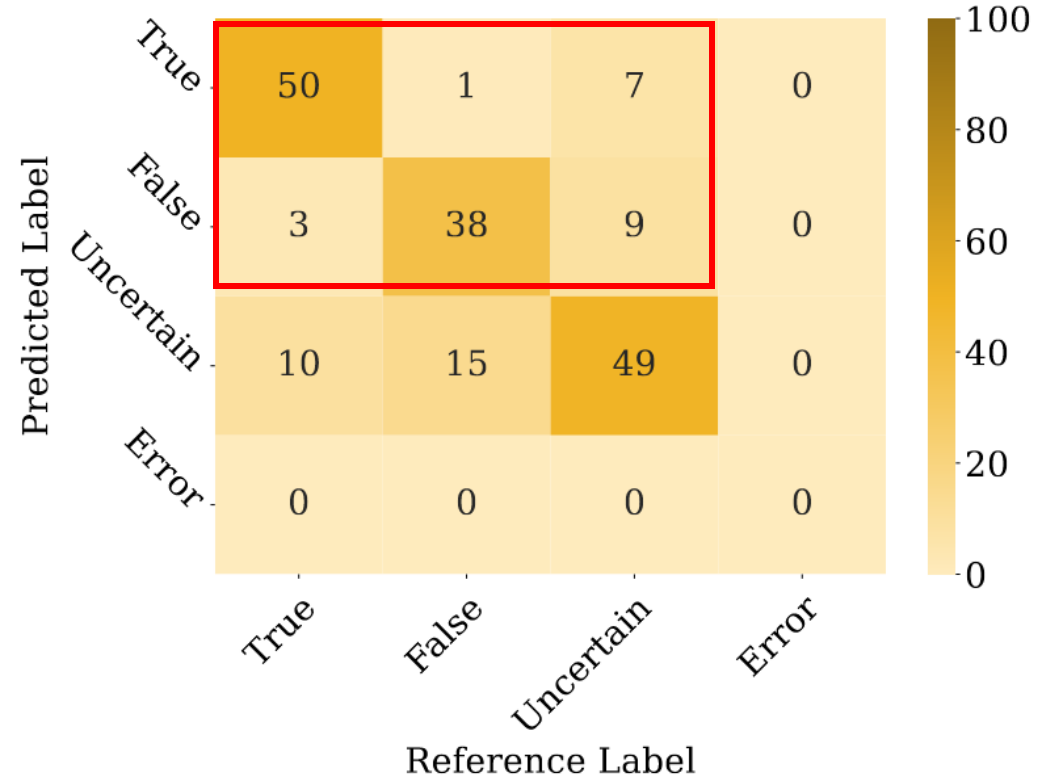
Contradiction

Results: FOLIO (LINC vs COT Quantitative)

- Compared to COT, LINC has **better precision** on True/False prediction (93% vs 81%)



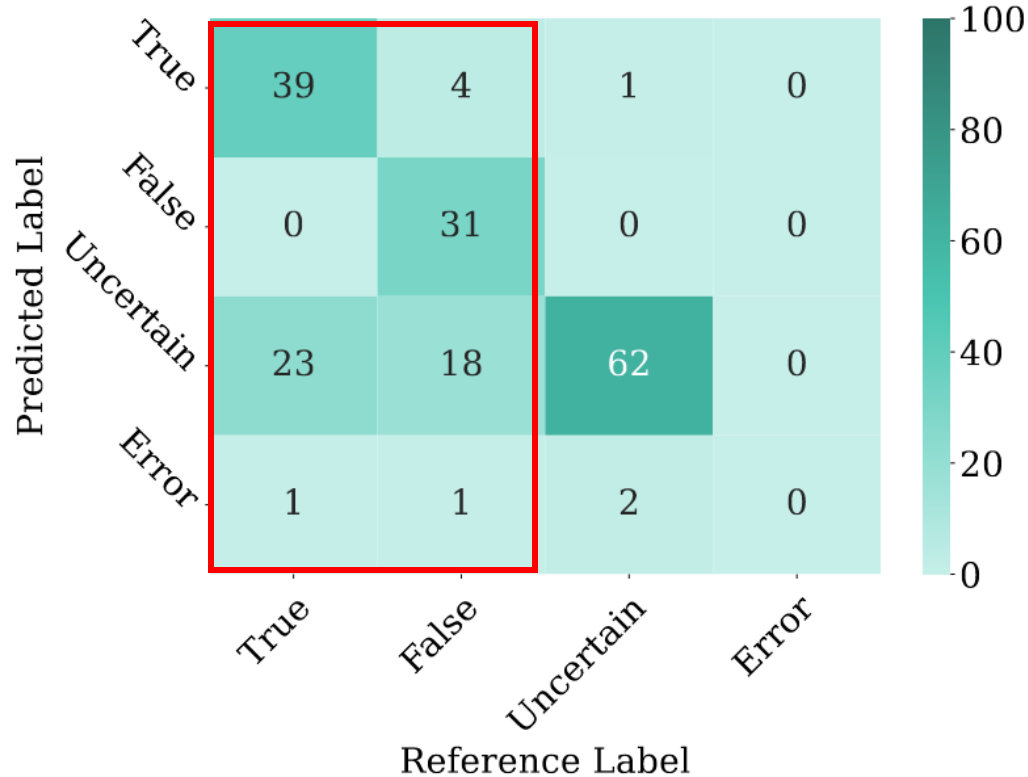
(a) Confusion matrix for LINC.



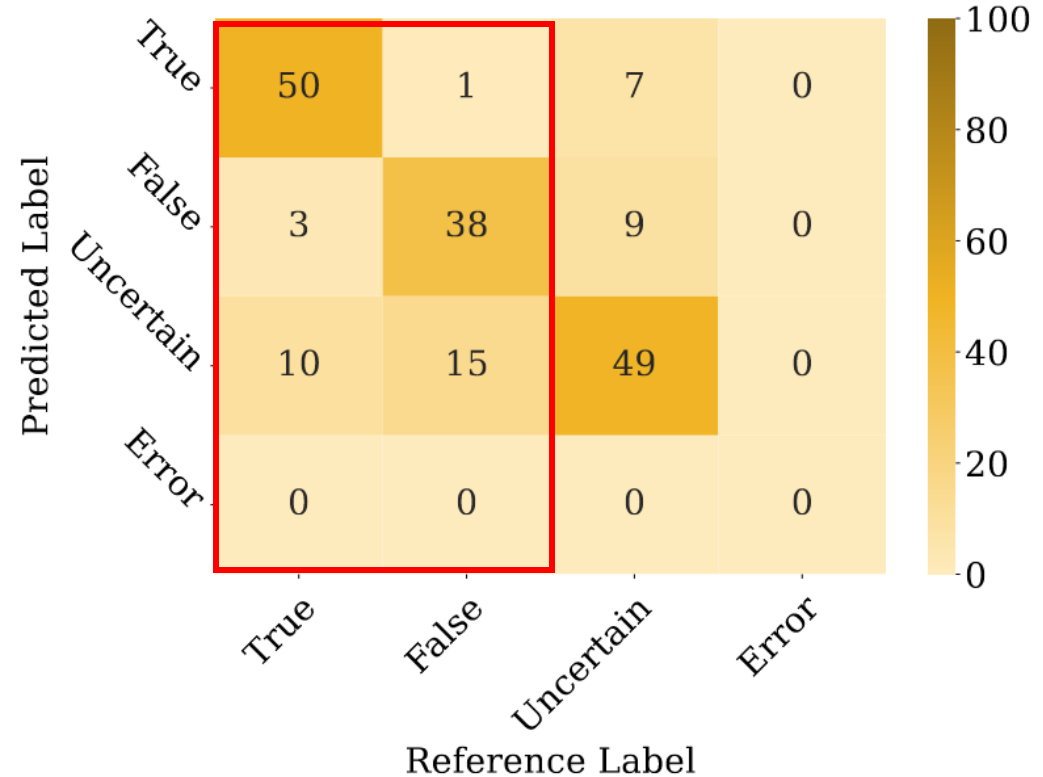
(b) Confusion matrix for Chain-of-Thought.

Results: FOLIO (LINC vs COT Quantitative)

- Compared to COT, LINC has **better precision** on True/False prediction (93% vs 81%)
- LINC has **worse recall** (60% vs 75%)



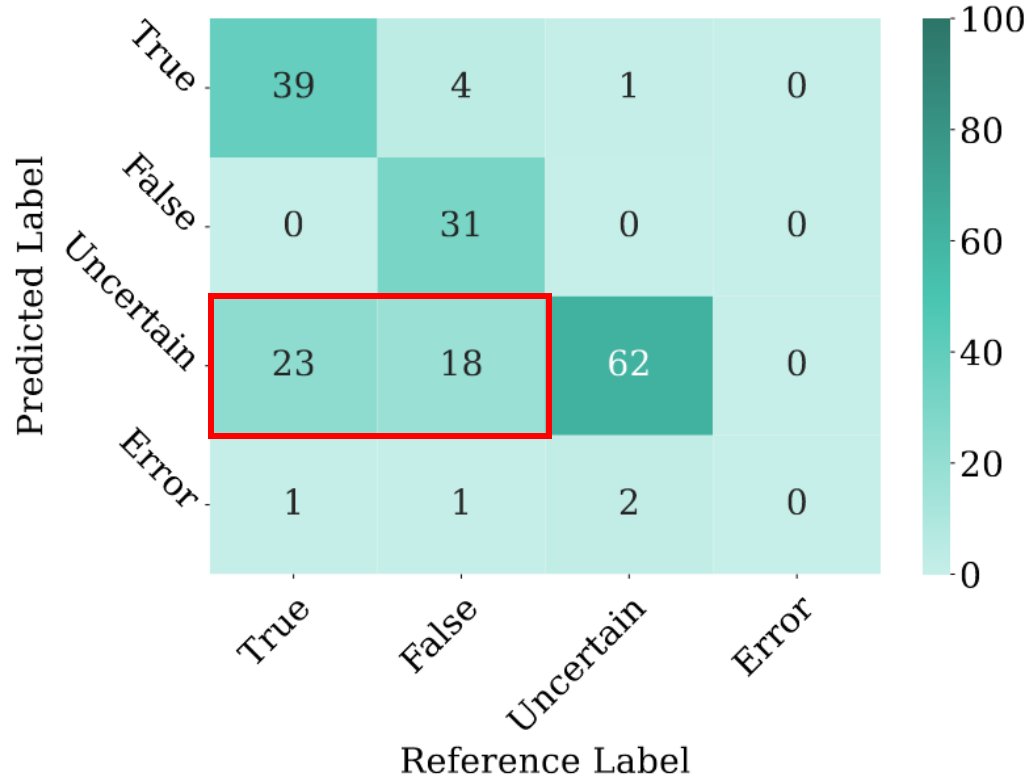
(a) Confusion matrix for LINC.



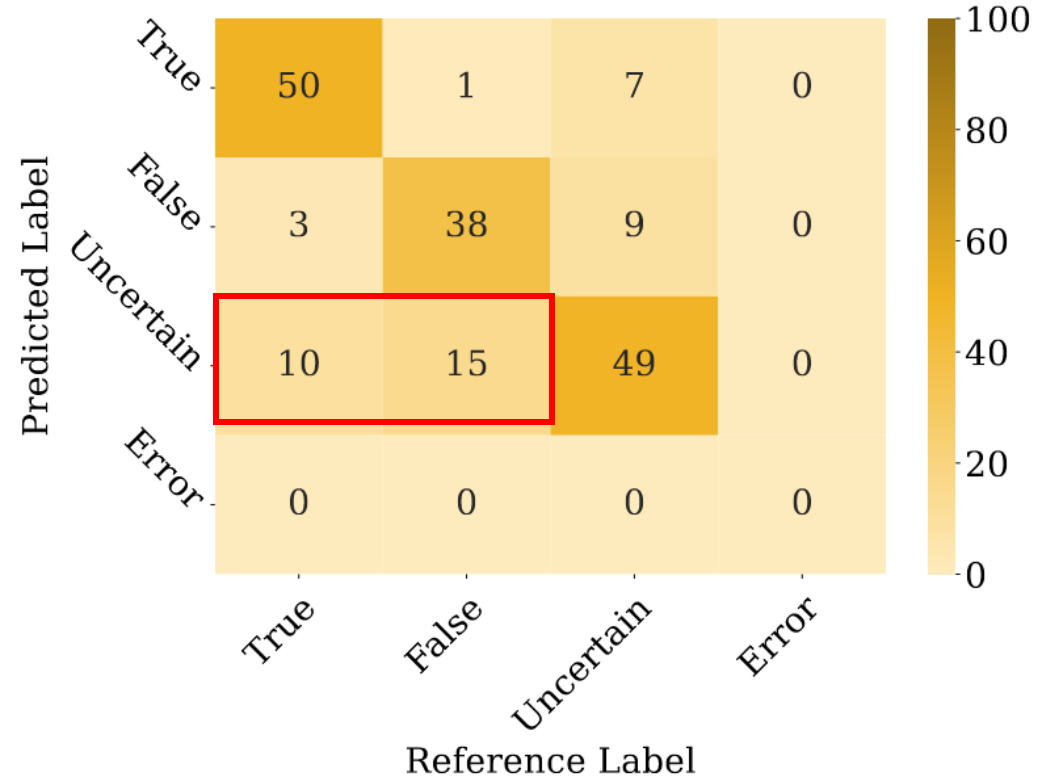
(b) Confusion matrix for Chain-of-Thought.

Results: FOLIO (LINC vs COT Quantitative)

- Compared to COT, LINC has **better precision** on True/False prediction (93% vs 81%)
- LINC has **worse recall** (60% vs 75%)
- LINC outputs "Uncertain" more: NL to FOL is a lossy process (but does not add false information)!



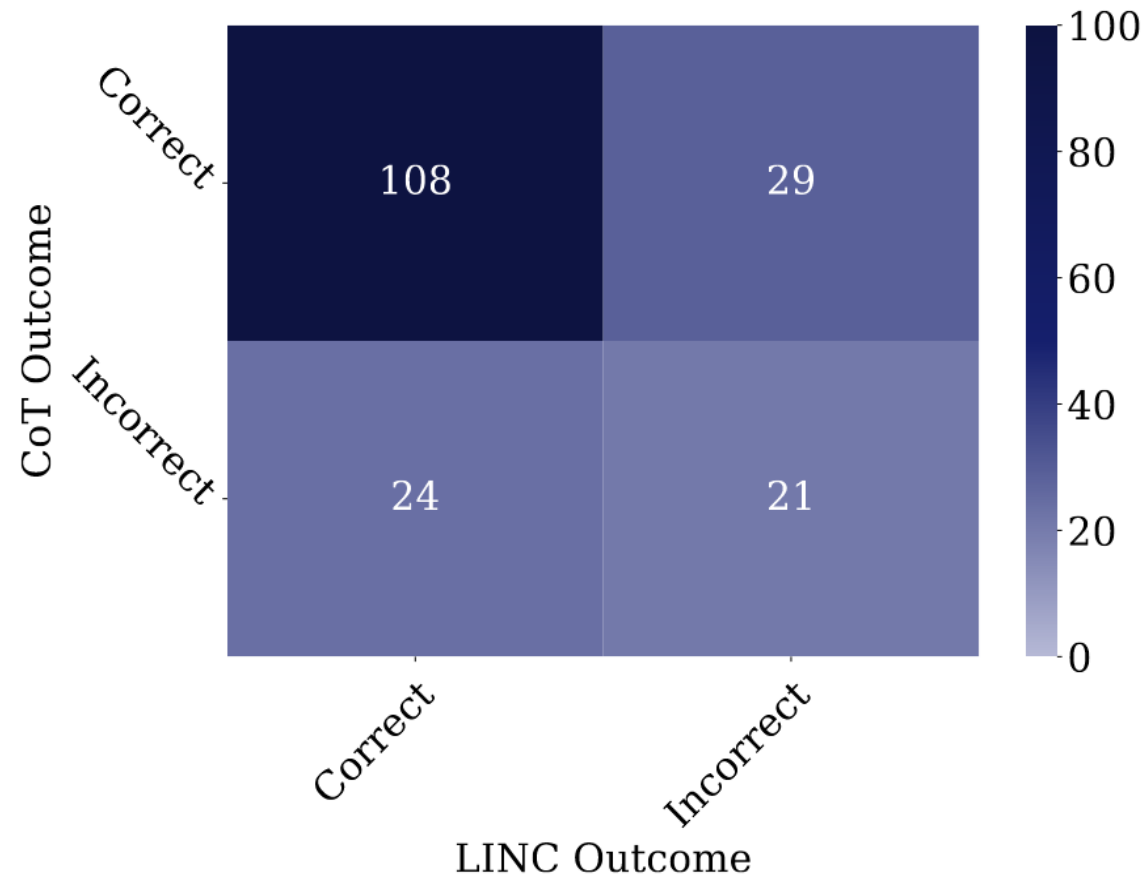
(a) Confusion matrix for LINC.



(b) Confusion matrix for Chain-of-Thought.

Results: FOLIO (LINC vs COT Quantitative)

LINC and COT mispredict on different examples!



Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- Baselines: Naïve, COT and ScratchPad
- **Results**
 - ProofWriter: LINC clear winner (as formalization task was easier)
 - FOLIO: LINC is behind but close to GPT-4, with complementary failures modes
- Conclusions & Limitations

Contents

- Problem & Motivation: LLMs struggle with logical reasoning!
- Methodology: Use LLM to go from NL to FOL, then use solver!
- Datasets & Models: (ProofWriter and FOLIO) & (GPT3.5, GPT4 and StarCoder+)
- Baselines: Naïve, COT and ScratchPad
- Results
 - ProofWriter: LINC clear winner (as formalization task was easier)
 - FOLIO: LINC is behind but close to GPT-4, with complementary failures modes
- **Conclusions & Limitations**

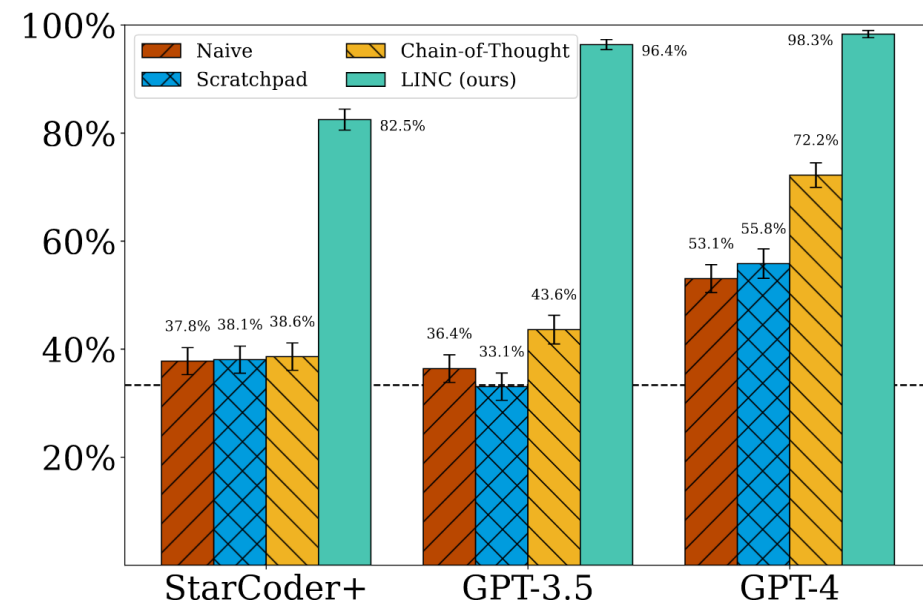
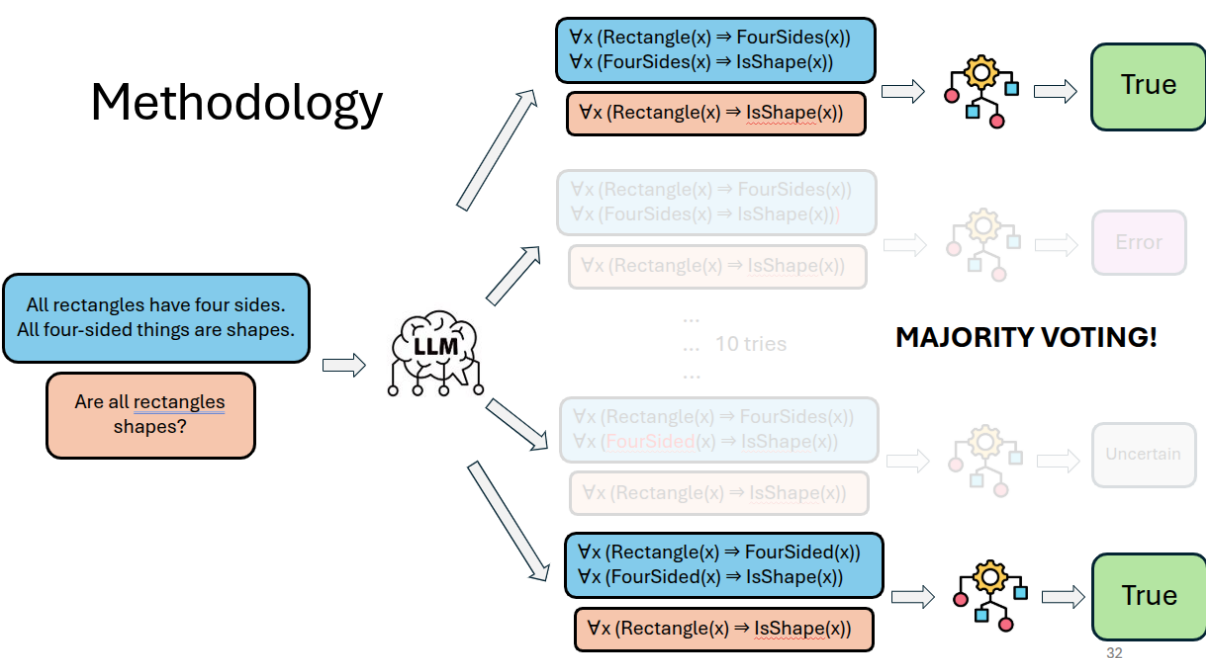
Conclusions

- **LINC improves reasoning accuracy** across almost all tested scenarios.
- **Generalizes to larger premise sets** than seen in in-context examples.
- **Complements Chain-of-Thought prompting** with different error patterns.

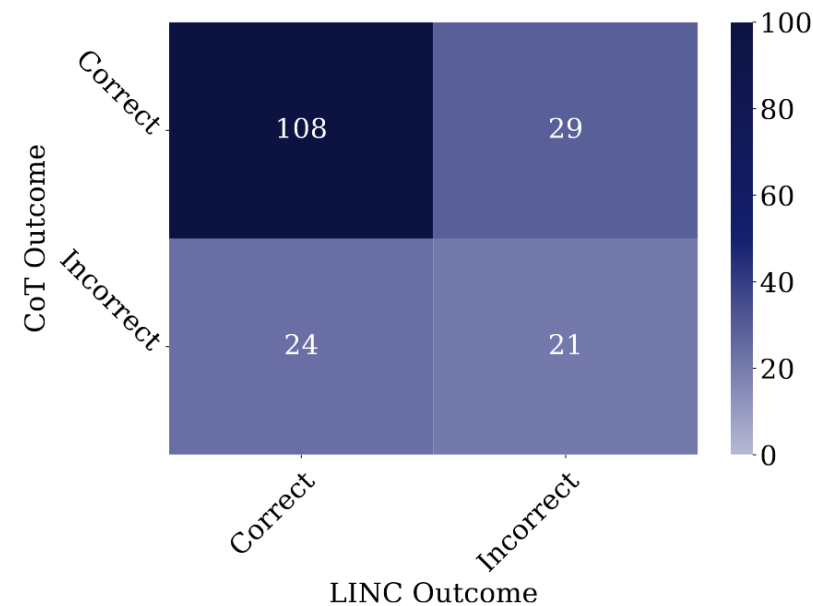
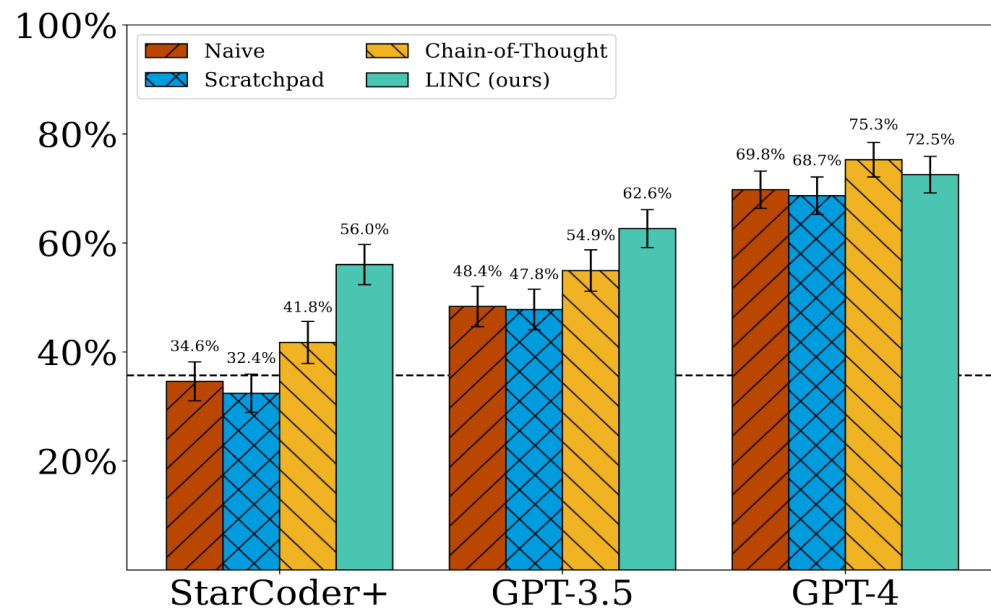
Limitations

- Find ways to combine COT and LINC.
- Integrate **prover feedback** in a refinement loop.
- Explore **fine-tuning** and other training techniques to boost reasoning.

Methodology



Questions?



Backup Slides

FOL BNF Grammar

$$t ::= x \mid c \mid f(t, t, \dots t)$$

$$\varphi ::= t = t \mid R(t, t, \dots t) \mid (\neg \varphi) \mid (\varphi \vee \varphi) \mid (\exists x \varphi)$$

Prover9 called twice to get complete info!

```
def evaluate(premises, conclusion):
    premises = [reformat_fol(p) for p in premises]
    conclusion = reformat_fol(conclusion)

    c = read_expr(conclusion)
    p_list = []
    for p in premises:
        p_list.append(read_expr(p))
    truth_value = prover.prove(c, p_list)
    if truth_value:
        return "True"
    else:
        neg_c = read_expr("-(" + conclusion + ")")
        negation_true = prover.prove(neg_c, p_list)
        if negation_true:
            return "False"
        else:
            return "Uncertain"
```

Prover9 Algorithm

The Inference Loop

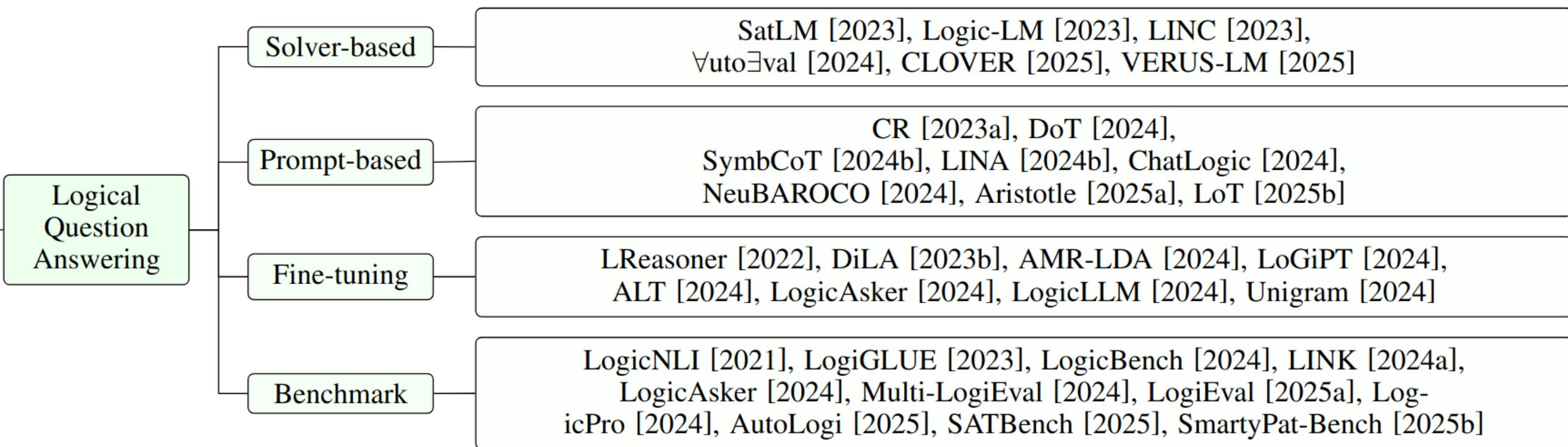
The *main loop* for inferring and processing clauses and searching for a proof is sometimes called the *given clause algorithm*. It operates mainly on the `sos` and `usable` lists.

```
While the sos list is not empty:  
  1. Select a given clause from sos and move it to the usable list;  
  2. Infer new clauses using the inference rules in effect;  
     each new clause must have the given clause as one of its  
     parents and members of the usable list as its other parents;  
  3. process each new clause;  
  4. append new clauses that pass the retention tests to the sos list.  
end of while loop.
```

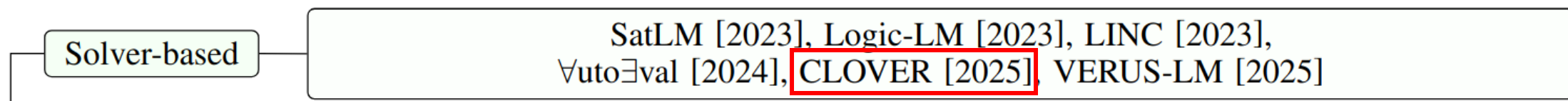

Prover9 Algorithm

- **Given-clause loop:** maintain usable/sos; pick given clause, infer with usable, simplify, retain; stop on empty clause or exhaustion.
- **Ordered resolution:** resolve only on maximal complementary literals (after unification) to prune search yet stay complete (with fairness).
- **Demodulation:** use oriented equalities as one-way rewrites (big \rightarrow small) to simplify clauses (no branching).
- **Paramodulation:** use an equality parent to replace equal subterms at eligible/maximal positions, producing new clauses.

Current State

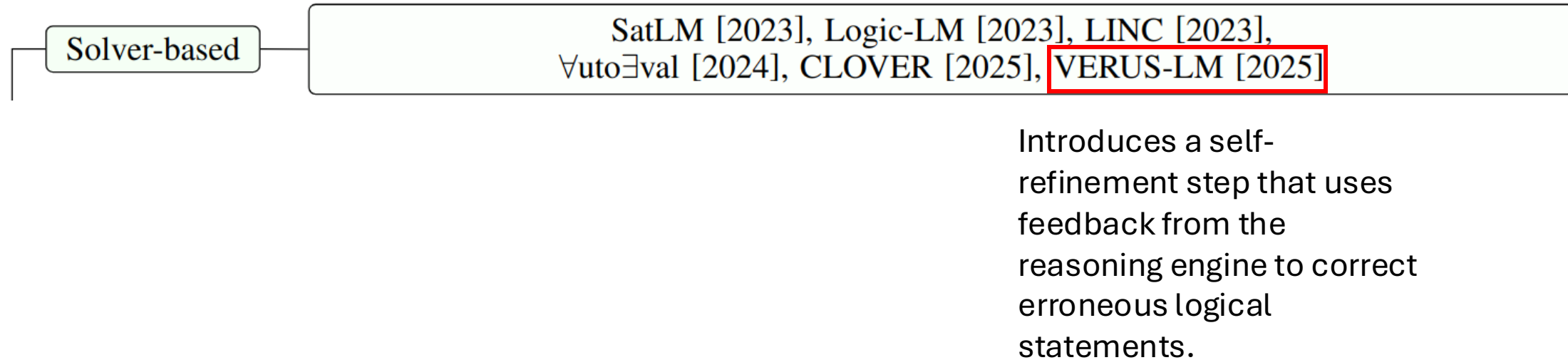


Current State

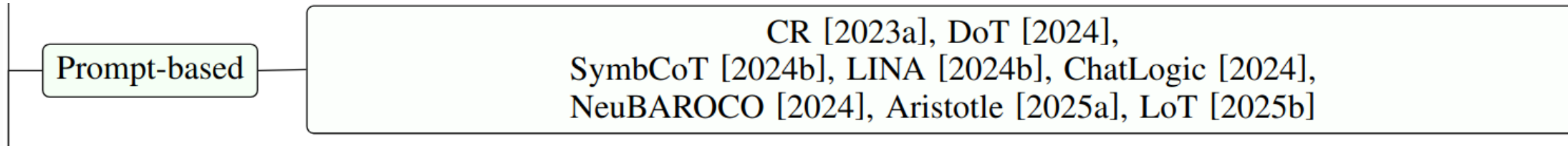


First translates the raw NL paragraph to atomic NL subsentences with their logical dependency structure, then translates to the target SL!

Current State

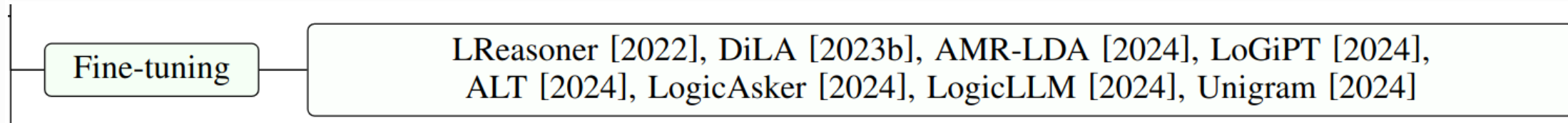


Current State



- **Process-oriented prompting:** show your work, then do more work; write out steps, explore and compare alternative chains, self-check with roles, and consolidate before answering
- **Symbolic-aware prompting:** translate to formal structure, decompose by true dependencies, expand implied rules, apply deductions, and verify the final answer

Current State



- **Rule-driven data generation for reasoning:** use formal logic rules or AMR structures to synthesize NL + proof traces and target weak rules for FT/ICL (LogicAsker, ALT, AMR-LDA; deeper chains help).
- **Symbolic-guided process learning:** imitate or integrate solvers/logic layers so models learn stepwise reasoning, not just answers (LoGiPT, DiLA, Unigram).