

Cognitive Motivations in Analogical and Physical Reasoning with Large Language Models

Shane Storks (he/him)

PhD Candidate in Computer Science and Engineering

Cognitive Science Seminar

October 10, 2023



Introduction

- Large language models (LLMs) like ChatGPT and GPT-4 have recently attracted attention
- Impressive, seemingly human-like conversation and reasoning capabilities solve many problems for automated language processing
- Enable research on interesting questions:
 1. **How can LLMs shed light on the nature of human language and reasoning?**
 2. **How can human reasoning strategies empower LLMs to better capture how the world works?**



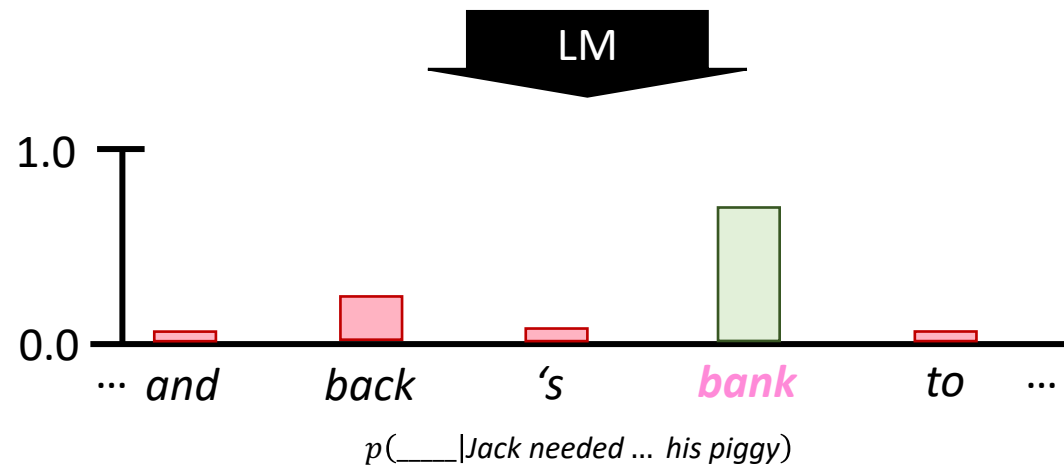
Outline

- **Language Model Basics**
- Application 1: Analogical Reasoning
- Application 2: Physical Commonsense Reasoning

Language Models

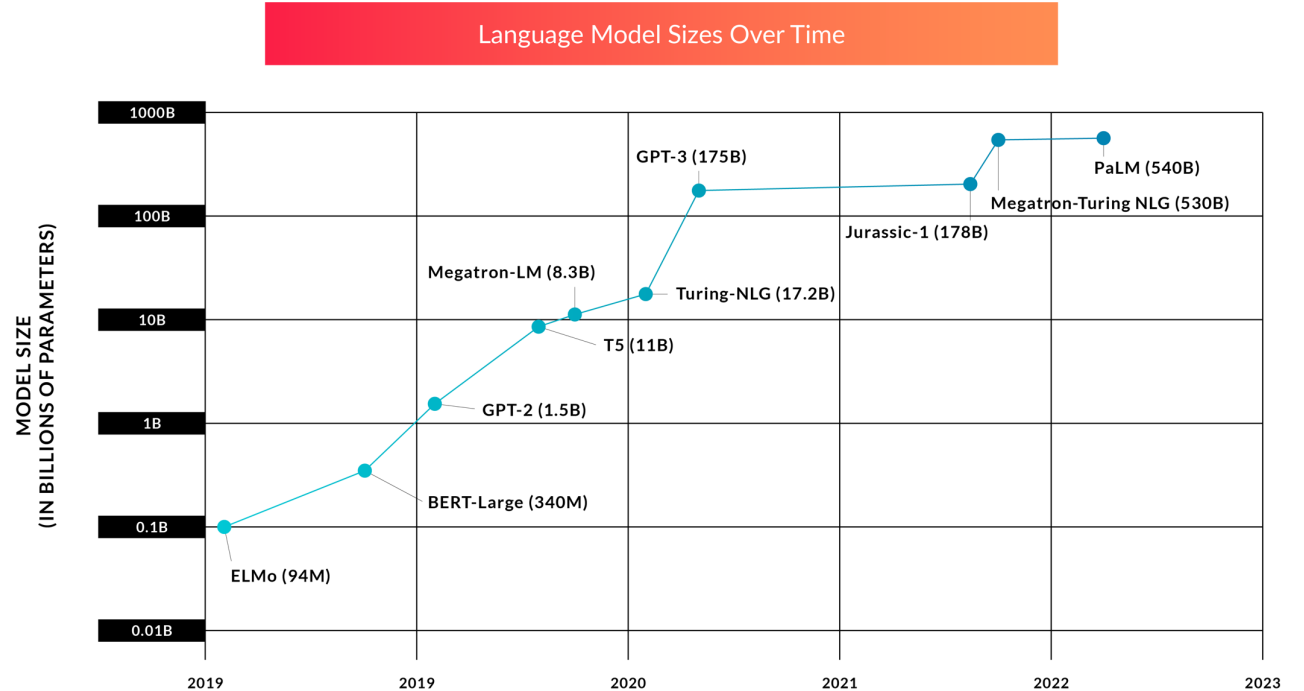
$$p(w_n | w_1, w_2, \dots, w_{n-1})$$

Jack needed some **money**, so he went and shook his **piggy** _____



Large Language Models

- What makes a language model a *large* language model?
- Recent trends:
 - More data
 - Web data
 - Human feedback annotation
 - More learned parameters
- Gives rise to new abilities...

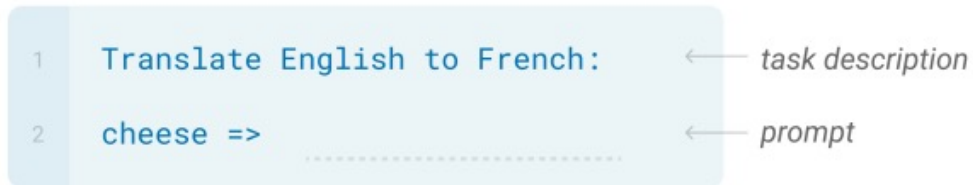


(figure from [Vinay Iyengar](#))

Prompting and In-Context Learning

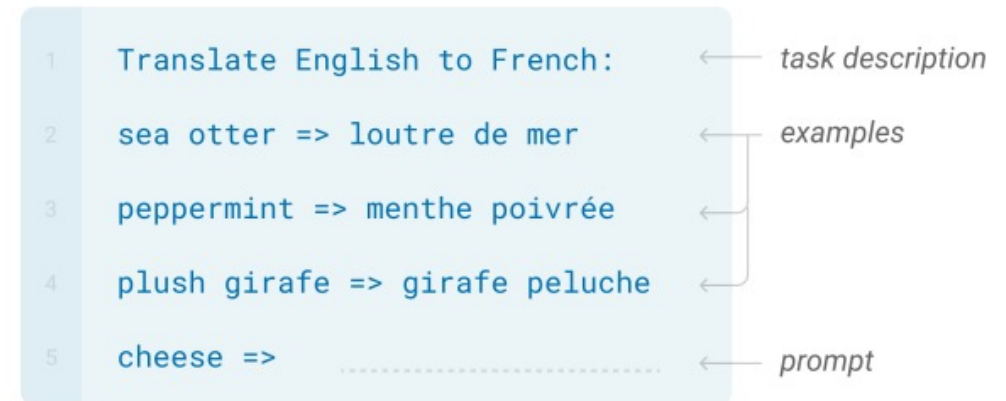
Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Chain-of-Thought Prompting

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. ✅

Outline

- Language Model Basics
- **Application 1: Analogical Reasoning**
- Application 2: Physical Commonsense Reasoning

In-Context Analogical Reasoning with Pre-Trained Language Models

Xiaoyang Hu^{1,2}*, Shane Storks¹*, Richard L. Lewis²†, Joyce Chai¹†

¹ Computer Science & Engineering Division, University of Michigan

² Department of Psychology, University of Michigan

* Equal contribution

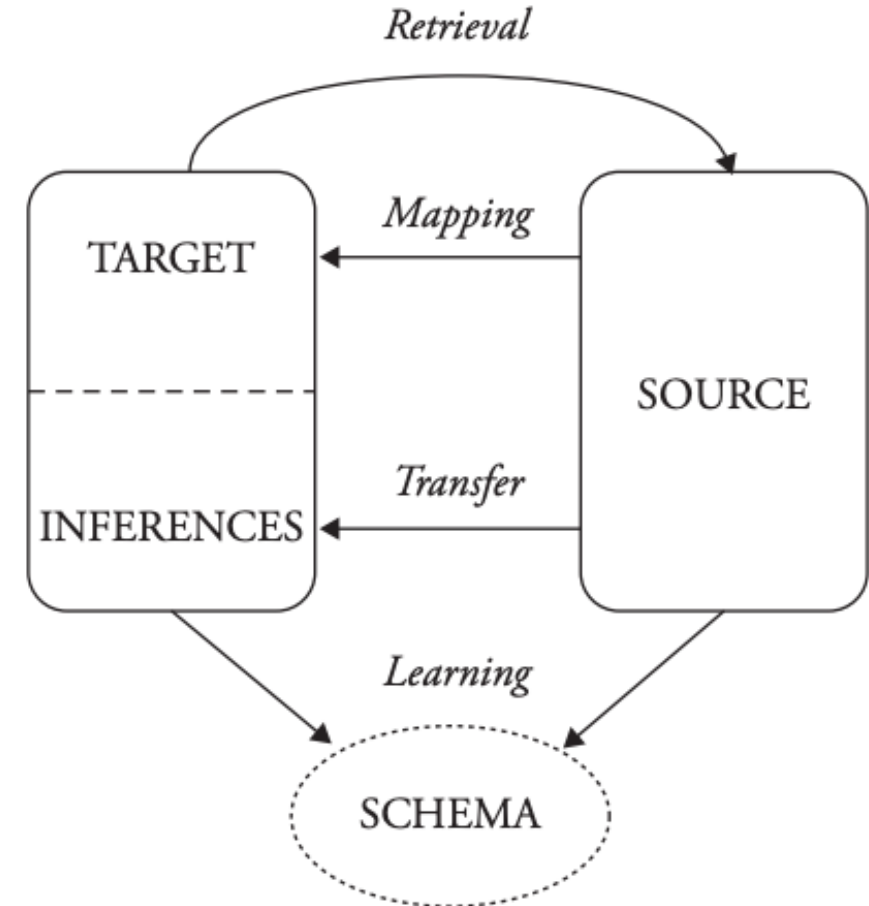
† Equal advising

[ACL 2023 Long Paper](#)

some slides made by Xiaoyang Hu

Motivation

- Making analogies is a fundamental capability of humans
- Enables us to tackle new situations based on past experience

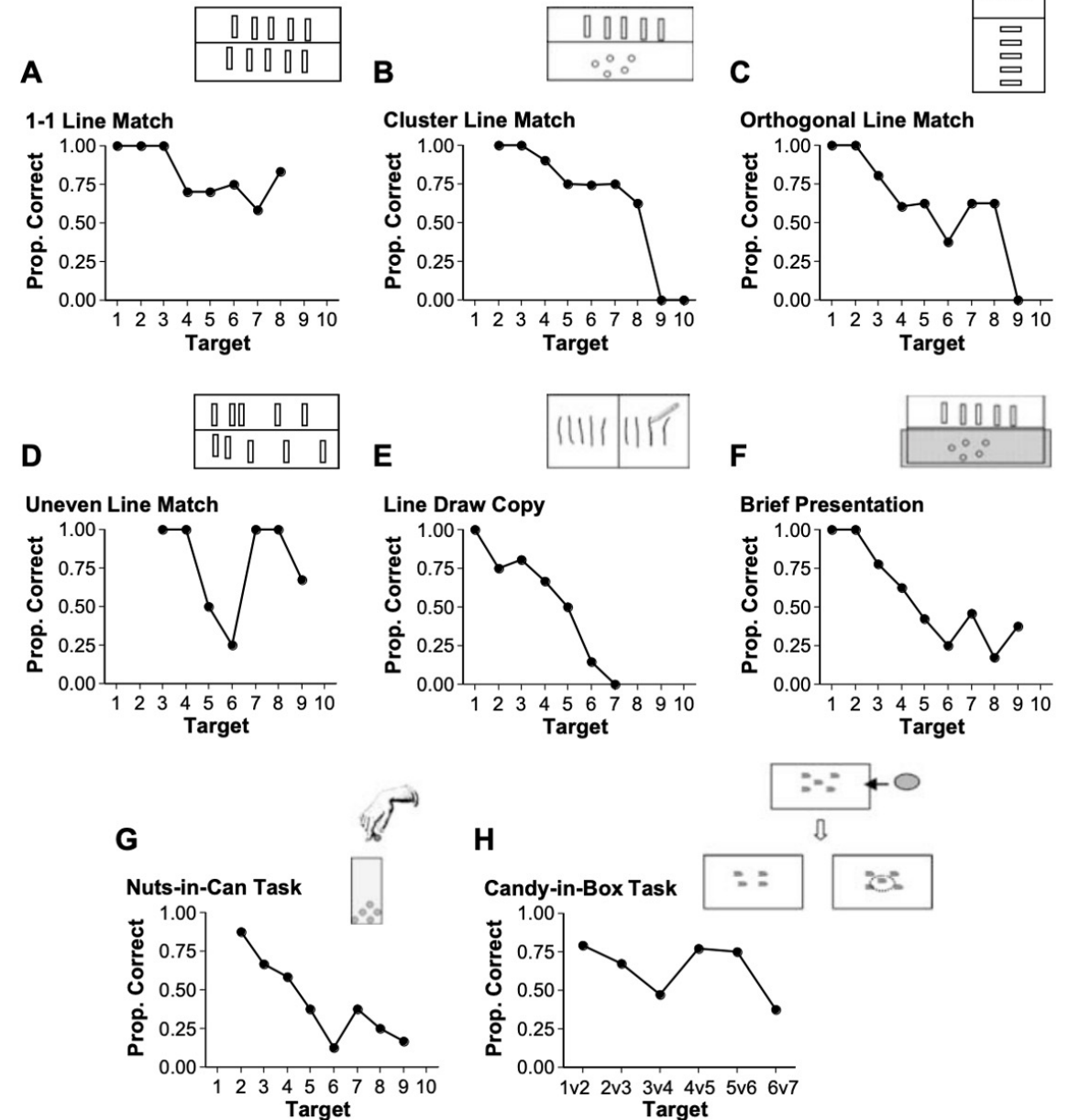


Motivation

- Work in cognitive science has found that language and analogy are connected in humans:

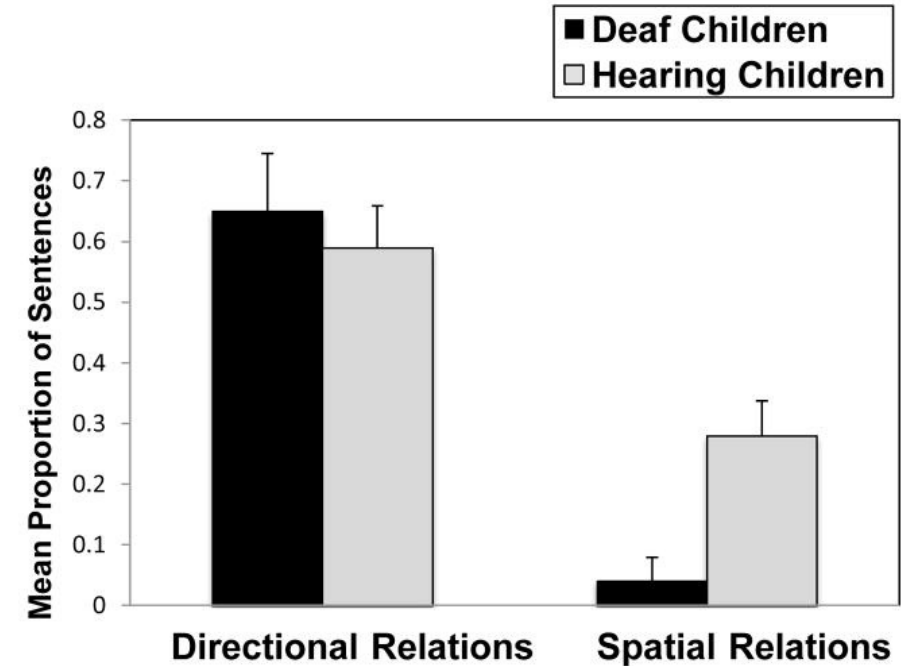
Motivation

- Work in cognitive science has found that language and analogy are connected in humans:
 - Numerical language facilitates numerical analogies



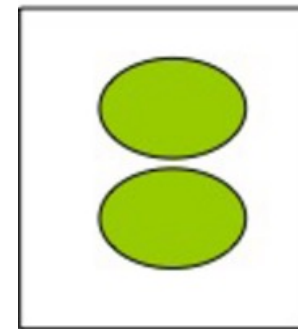
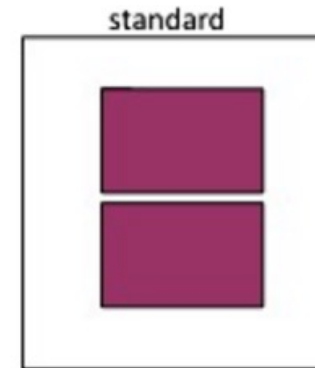
Motivation

- Work in cognitive science has found that language and analogy are connected in humans:
 - Numerical language facilitates numerical analogies
 - Spatial language facilitates spatial analogies

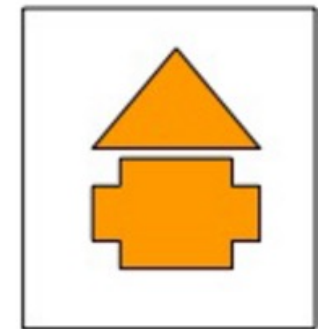


Motivation

- Work in cognitive science has found that language and analogy are connected in humans:
 - Numerical language facilitates numerical analogies
 - Spatial language facilitates spatial analogies
 - Names support analogy-making (even nonsense names)



relational match



non-relational match

Motivation

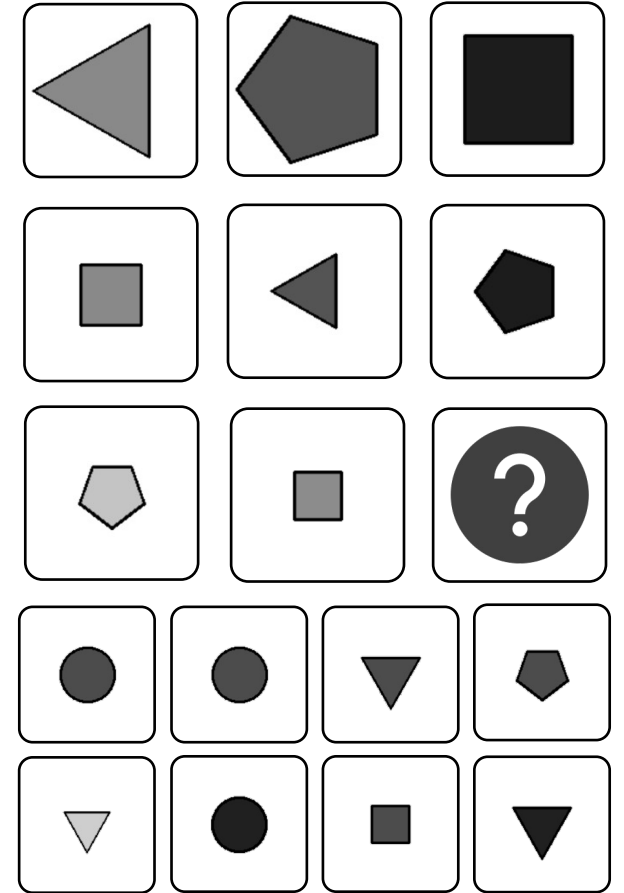
- Analogy-making may be key to robust reasoning in AI systems
- Contemporary AI approaches for analogy-making require thousands of training examples to make any progress
- Meanwhile, LLMs can pick up new tasks through in-context learning with just a few relevant examples (more like humans)
 - Are they capable of analogy-making?

Questions

1. Does training LLMs on *natural language* give rise to the ability to form *abstract* analogies?
2. How do various factors contribute to analogy-making in LLMs?
 - Complexity of situations to make analogies from
 - Language-based abstractions (like names)
 - Complexity (size/# learned parameters) of LLM
 - In-context demonstration of task

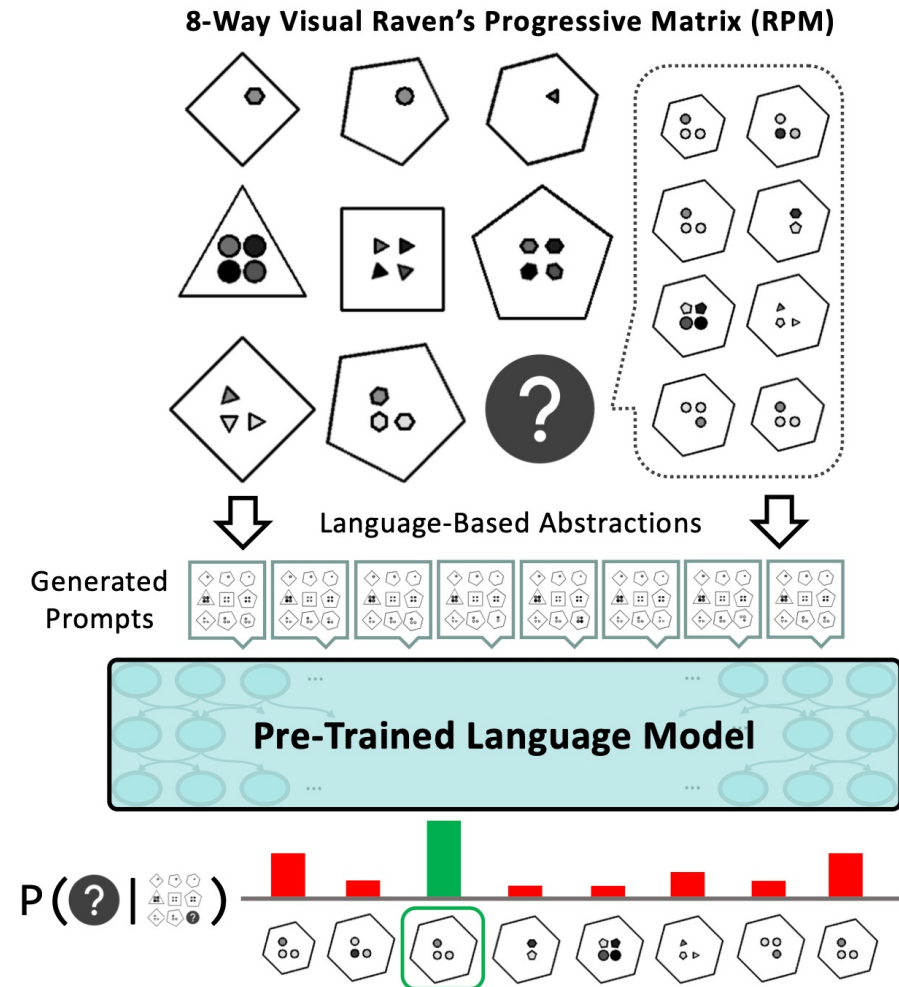
Raven's Progressive Matrices (RPM)

- A canonical test of analogical reasoning often used with human subjects
- Test-taker infers abstract rules from first 2 rows, then apply them to complete the third row
- RAVEN dataset
 - Relations:
 - Constant
 - Progression
 - Arithmetic
 - Distribute-Three

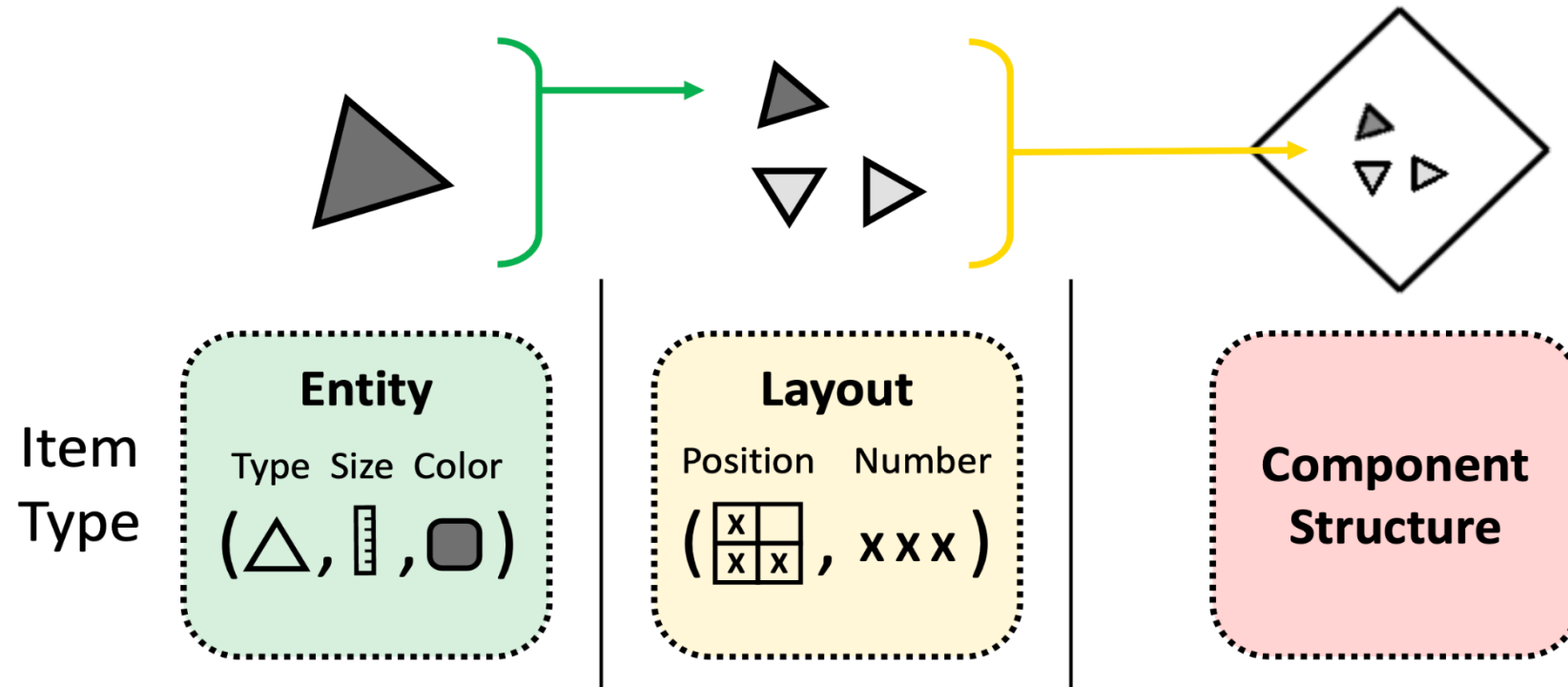


Prompting for Analogical Reasoning

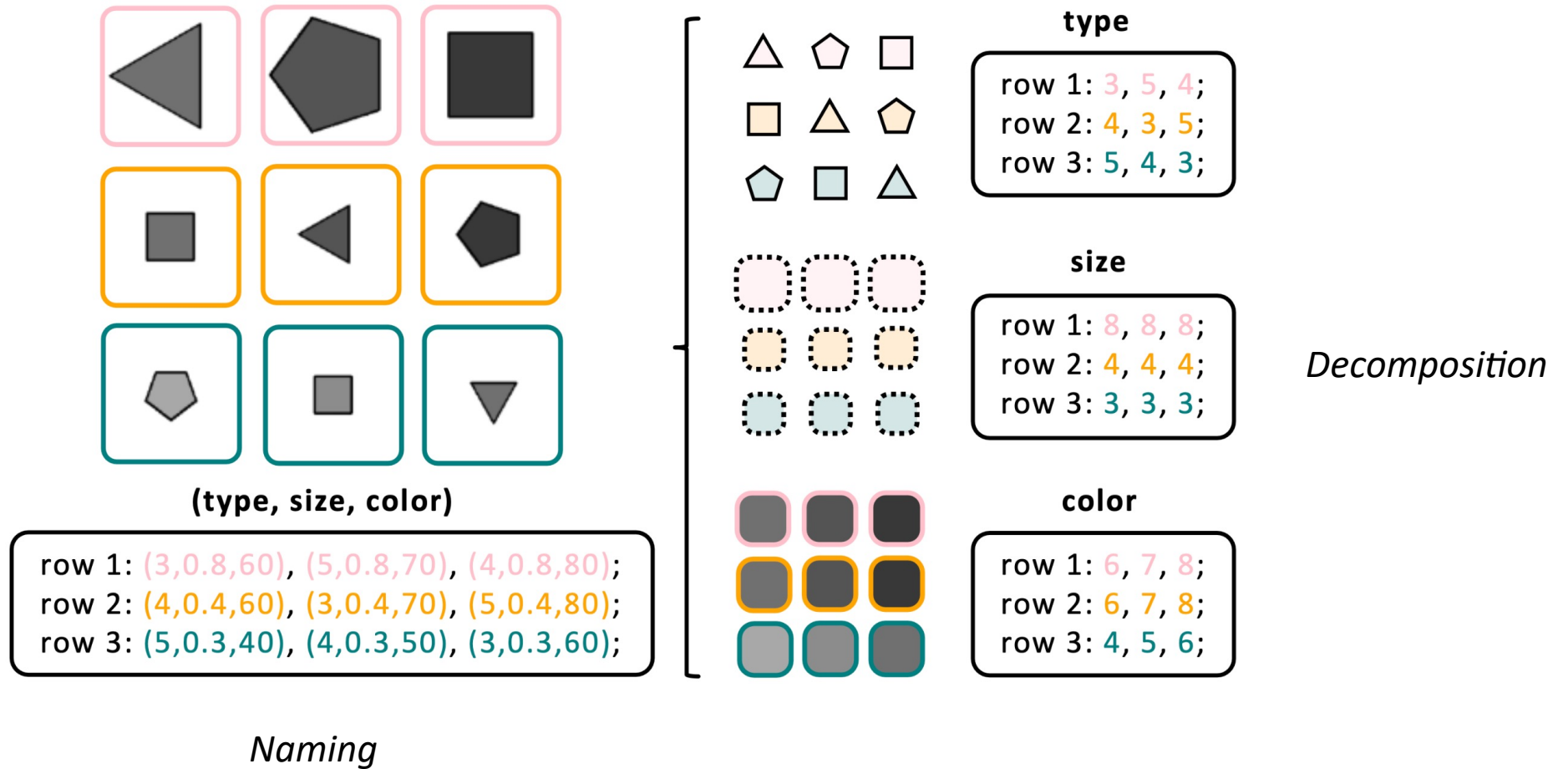
- Created language abstractions for RPMs in RAVEN dataset
- Prompt LLMs to test abstract analogical reasoning capability
 - OPT & InstructGPT at varying model complexity



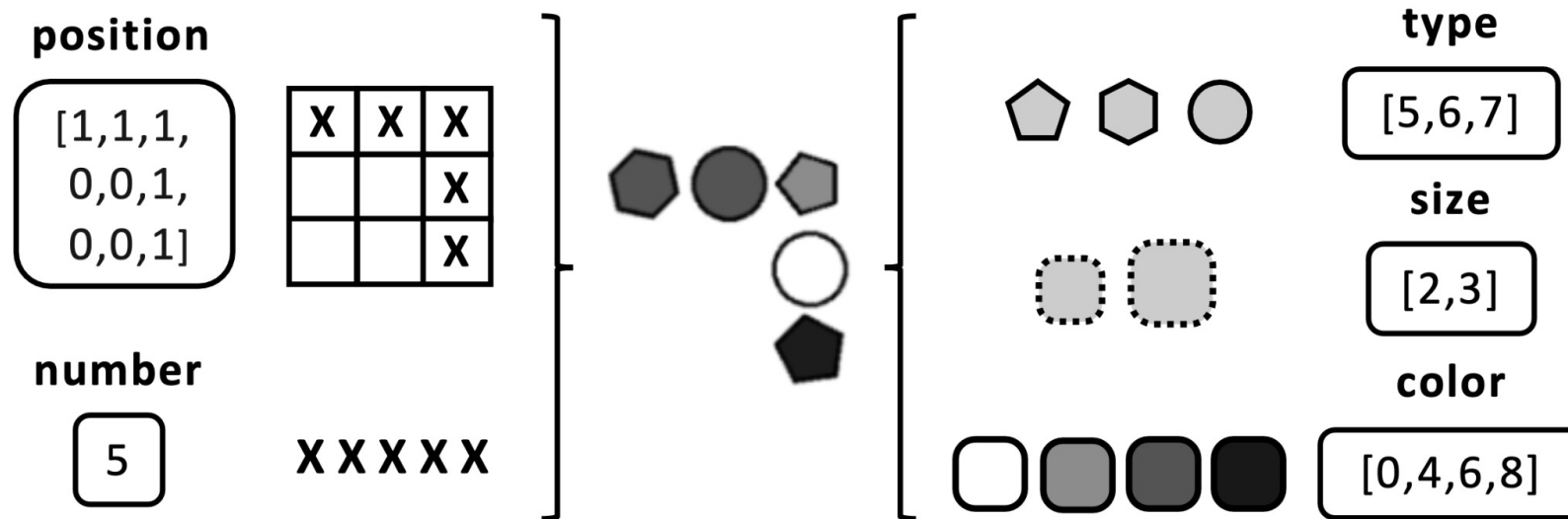
Components of RAVEN Matrix Items



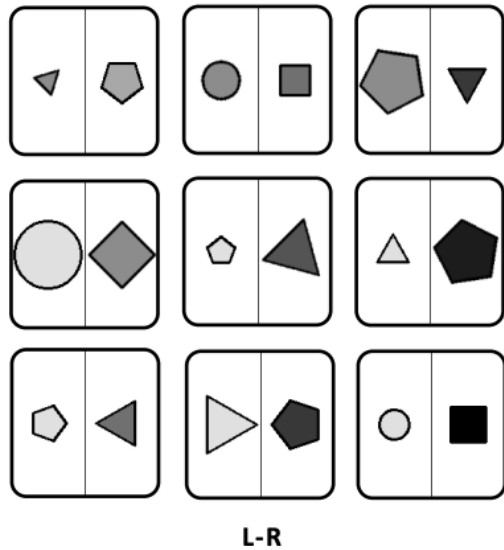
Entity-Level Abstractions



Layout-Level Abstractions



Component-Level Abstraction



Complete

row 1: A (3,0.1,40) / B (5,0.3,30), A (7,0.2,40) / B (4,0.3,50), A (5,0.6,40) / B (3,0.3,70);
 row 2: A (7,0.6,10) / B (4,0.6,40), A (5,0.1,10) / B (3,0.6,60), A (3,0.2,10) / B (5,0.6,80);
 row 3: A (5,0.2,10) / B (3,0.4,50), A (3,0.6,10) / B (5,0.4,70), A (7,0.1,10) / B (4,0.4,90);

Left Comp.	Right Comp.
row 1: (3,0.1,40), (7,0.2,40), (5,0.6,40);	row 1: (5,0.3,30), (4,0.3,50), (3,0.3,70);
row 2: (7,0.6,10), (5,0.1,10), (3,0.2,10);	row 2: (4,0.6,40), (3,0.6,60), (5,0.6,80);
row 3: (5,0.2,10), (3,0.6,10), (7,0.1,10);	row 3: (3,0.4,50), (5,0.4,70), (4,0.4,90);

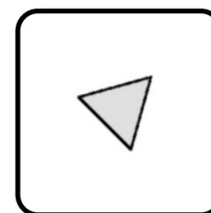
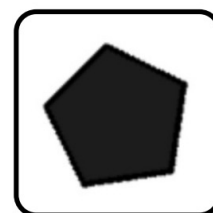
Type: Distr.	Size: Distr.	Color: Const.	Type: Distr.	Size: Const.	Color: Prog.
row 1: 3, 7, 5;	row 1: 1, 2, 6;	row 1: 4, 4, 4;	row 1: 5, 4, 3;	row 1: 3, 3, 3;	row 1: 3, 5, 7;
row 2: 7, 5, 3;	row 2: 6, 1, 2;	row 2: 1, 1, 1;	row 2: 4, 3, 5;	row 2: 6, 6, 6;	row 2: 4, 6, 8;
row 3: 5, 3, 7;	row 3: 2, 6, 1;	row 3: 1, 1, 1;	row 3: 3, 5, 4;	row 3: 4, 4, 4;	row 3: 5, 7, 9;



$$p(\text{○} \text{■} | \text{○} \text{■} \text{△} \text{□} \text{◇} \text{○} \text{■} \text{△} \text{□} \text{◇} \text{○} \text{■} \text{△} \text{□} \text{◇}) \propto p(\text{left type}) + p(\text{left size}) + p(\text{left color}) + p(\text{right type}) + p(\text{right size}) + p(\text{right color})$$

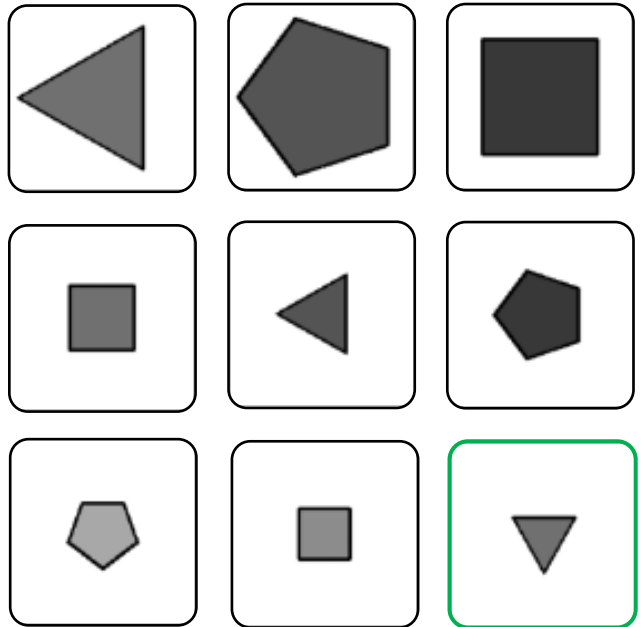
Baselines

- How helpful are the naming abstractions we chose?
- 2 baselines for comparison:
 1. **Quasi-image:** lower-level “pixel-like” abstraction
 2. **Random naming:** choose random words to represent attributes, removing numerical dependencies between attribute names


$$\begin{bmatrix} [2 & . & . & .], \\ [2 & 2 & . & .], \\ [2 & 2 & 2 & .], \\ [2 & 2 & 2 & 2], \end{bmatrix}$$

$$\begin{bmatrix} [. & . & . & 9 & . & . & .], \\ [. & . & 9 & 9 & 9 & . & .], \\ [. & 9 & 9 & 9 & 9 & 9 & .], \\ [9 & 9 & 9 & 9 & 9 & 9 & 9], \\ [9 & 9 & 9 & 9 & 9 & 9 & 9], \\ [9 & 9 & 9 & 9 & 9 & 9 & 9], \\ [9 & 9 & 9 & 9 & 9 & 9 & 9], \end{bmatrix}$$

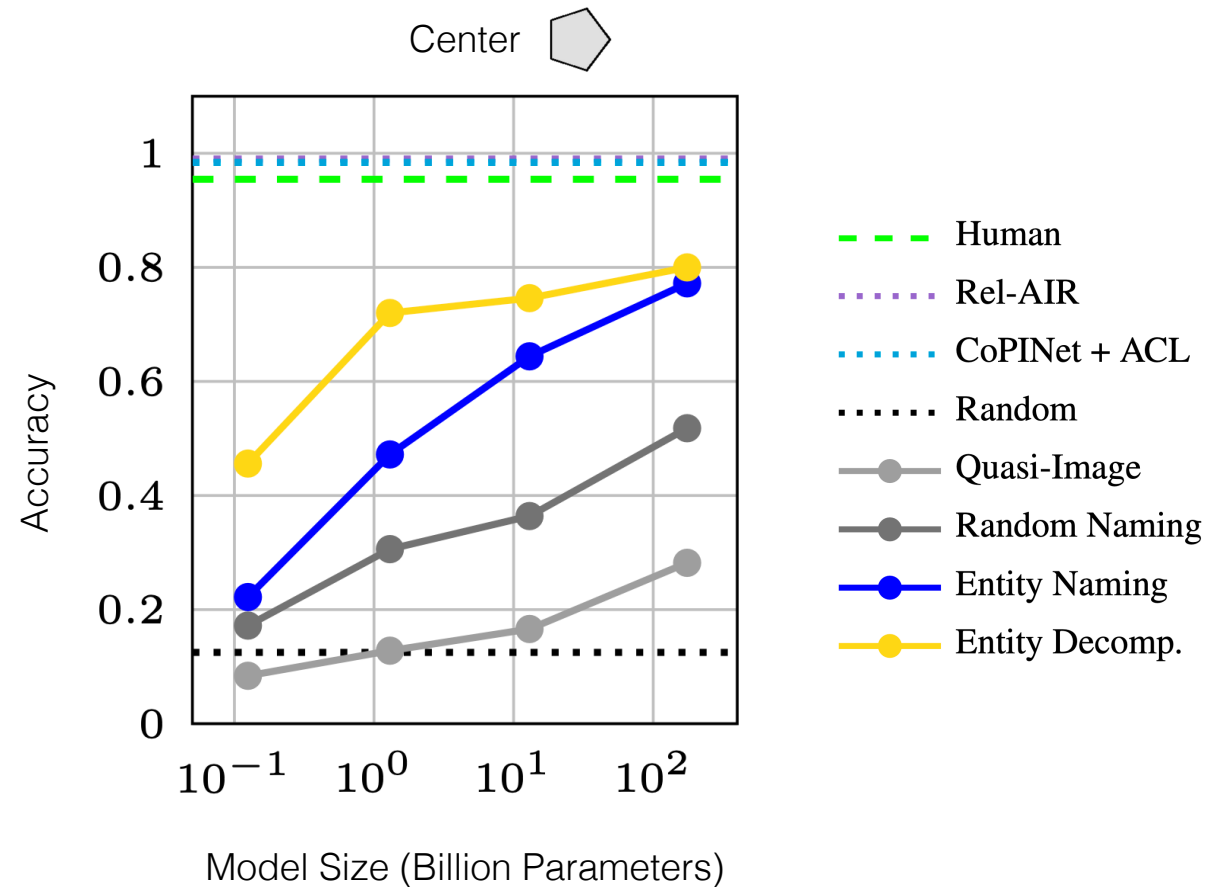
Single Entity Results

Center



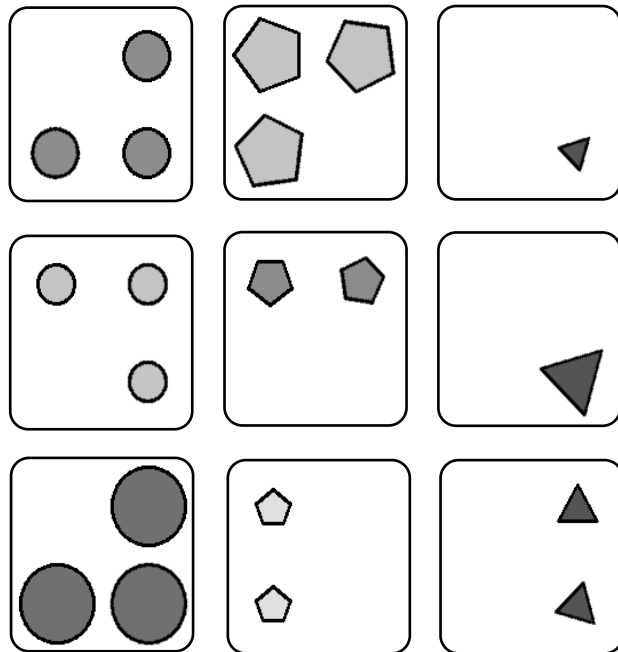
Single Entity Results

- Analogies *do* arise from natural language training!
- Bigger LLMs are better analogy-makers
- Numerical naming enables better analogy-making
- Decomposition abstractions especially help smaller LLMs
 - Model complexity \approx working memory?

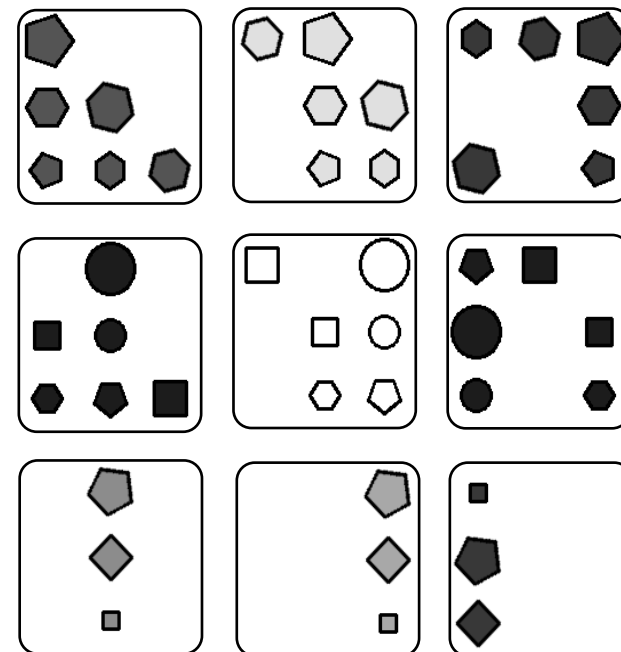


Multiple Entity Results

2x2Grid

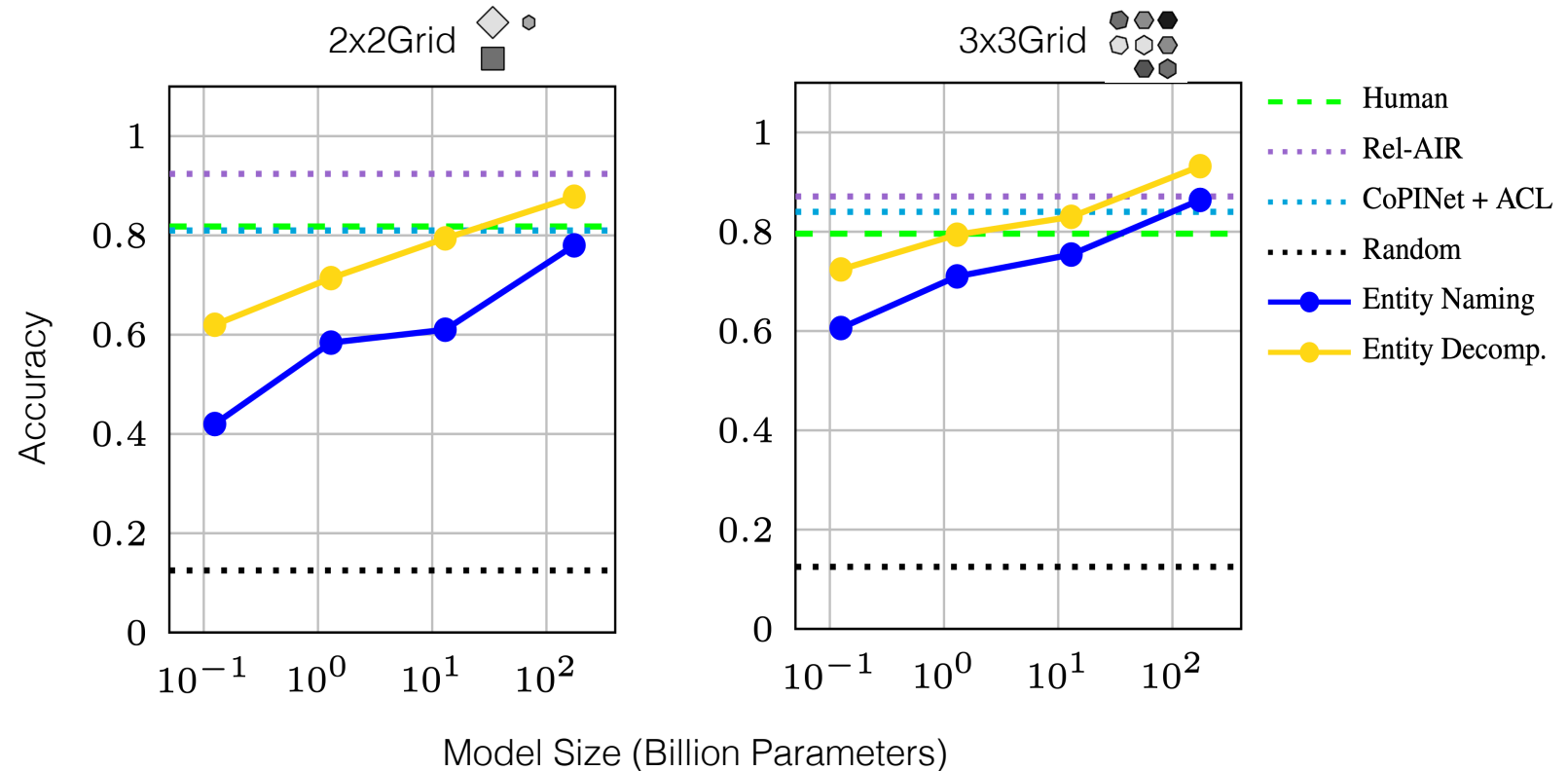


3x3Grid



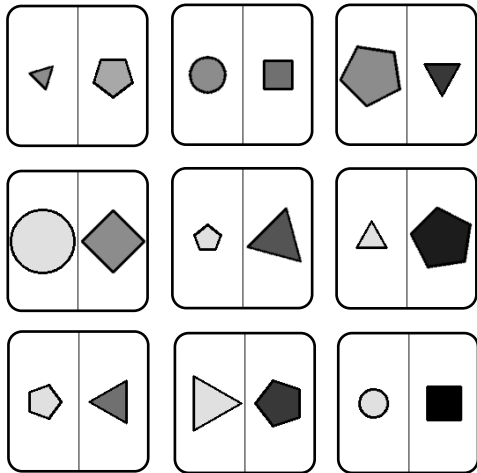
Multiple Entity Results

- Humans struggle more with task complexity than LLMs
 - Model complexity \approx working memory?
- Can outperform humans and supervised approaches

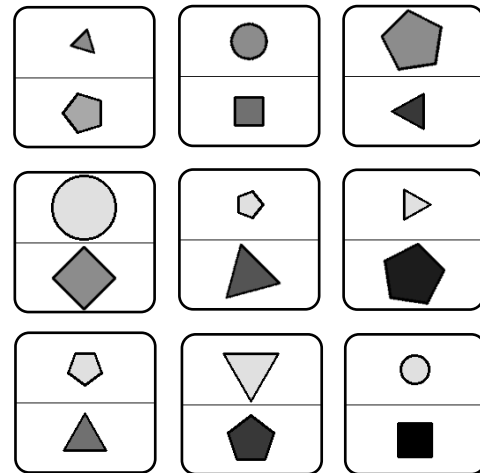


Multiple Component Results

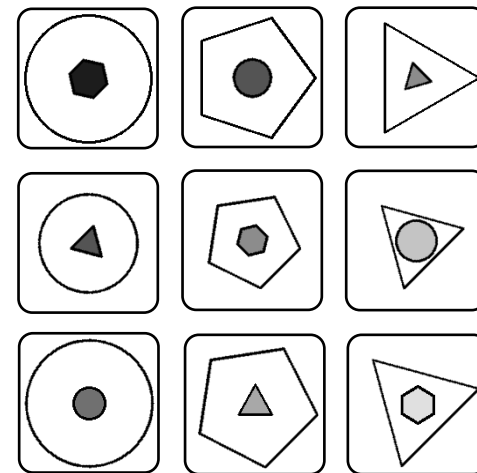
L-R



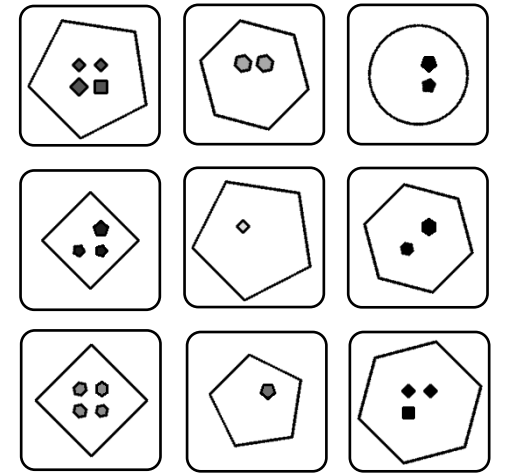
U-D



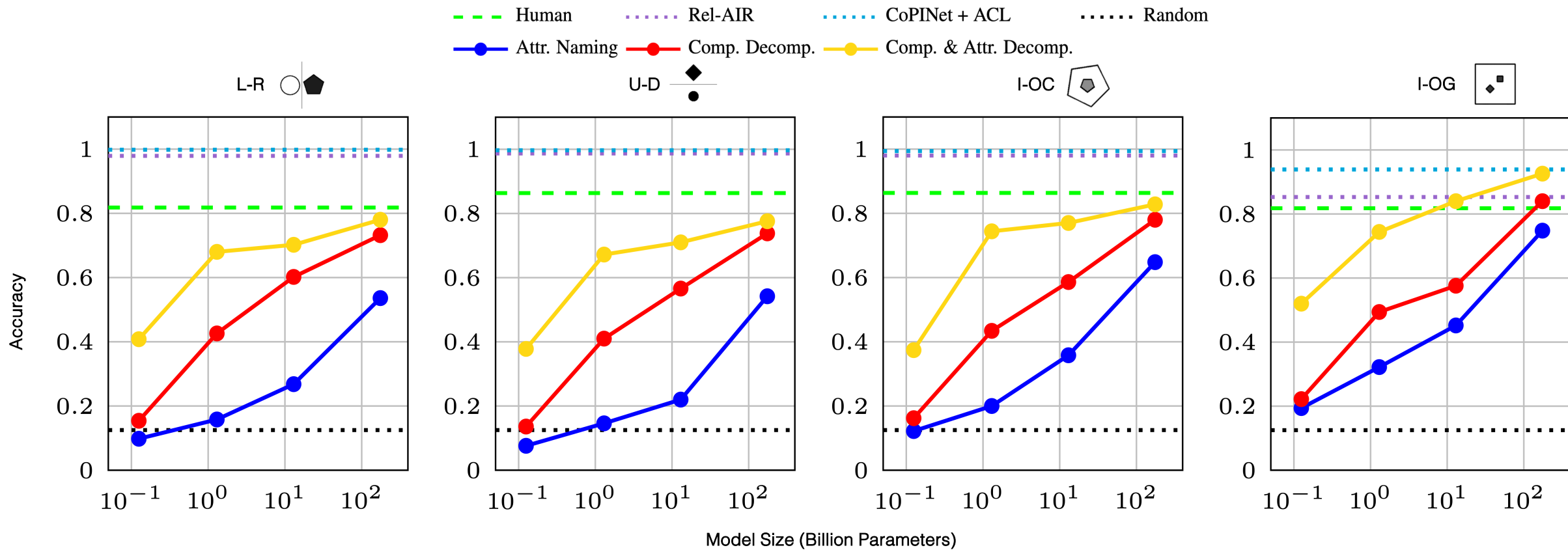
I-OC



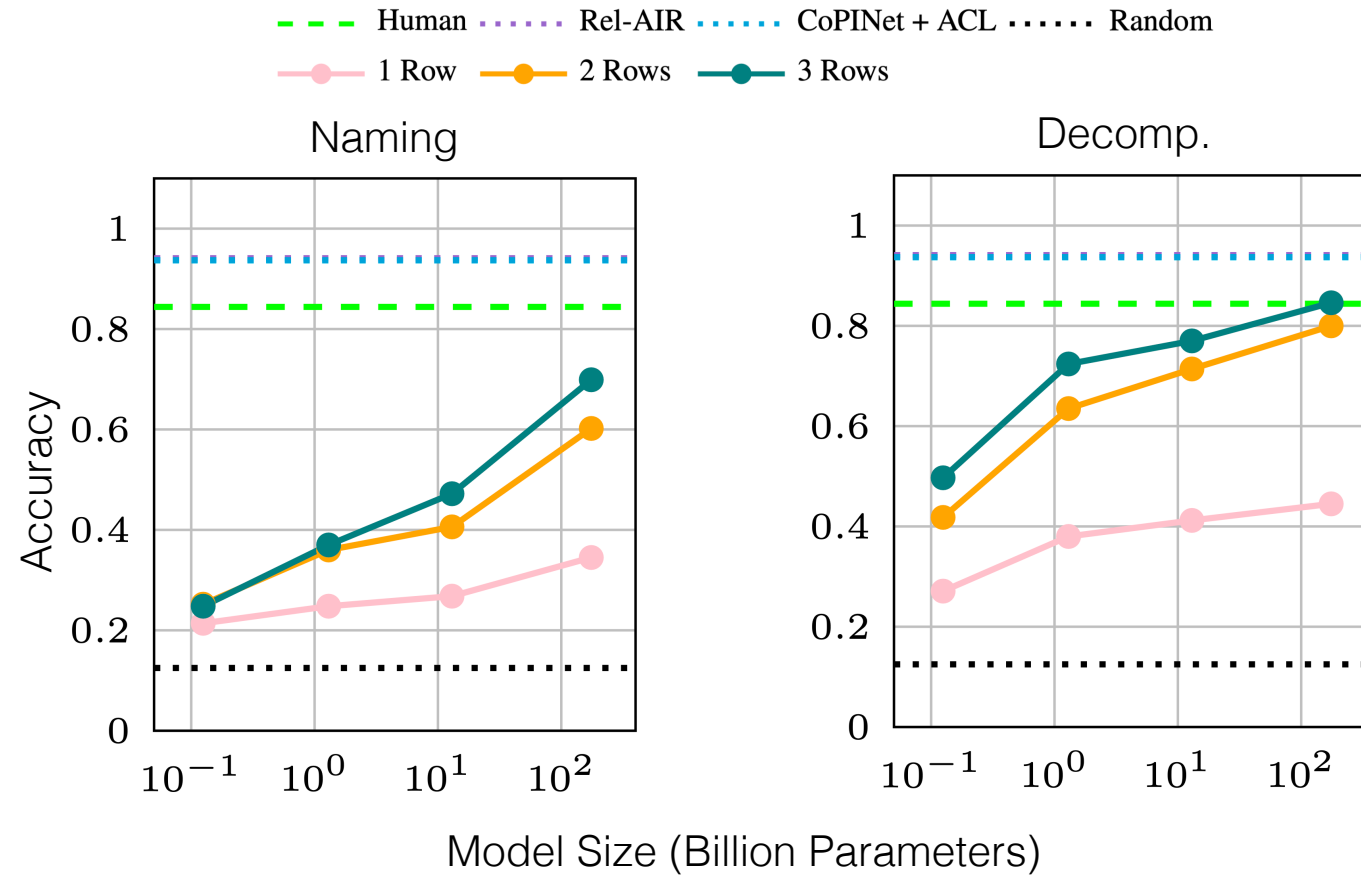
I-OG



Multiple Component Results



Impact of In-Context Learning



Takeaways

1. LLMs gain a fair capacity for abstract analogical reasoning from large-scale natural language training!
2. A number of factors strengthen their capability to make analogies:
 - Stronger language abstractions
 - LLM size
 - In-context demonstration
3. Complexity of context does not seem to impact LLMs as much as humans!

Outline

- Language Model Basics
- Application 1: Analogical Reasoning
- **Application 2: Physical Commonsense Reasoning**

From Heuristic to Analytic: Cognitively Motivated Reasoning Strategies for Coherent Physical Commonsense in Pre-Trained Language Models

**Zheyuan Zhang¹ Shane Storks¹ Fengyuan Hu¹ Sungryull Sohn²
Moontae Lee² Honglak Lee^{1,2} Joyce Chai¹**

¹University of Michigan, Computer Science and Engineering Division

²LG AI Research

EMNLP 2023 Long Paper



Tiered Reasoning for Intuitive Physics (TRIP)

Story A

1. Ann sat in the chair.
2. Ann turned off the telephone.
3. Ann picked up a pencil.
4. Ann opened the book.
5. Ann wrote in the book.

Story B

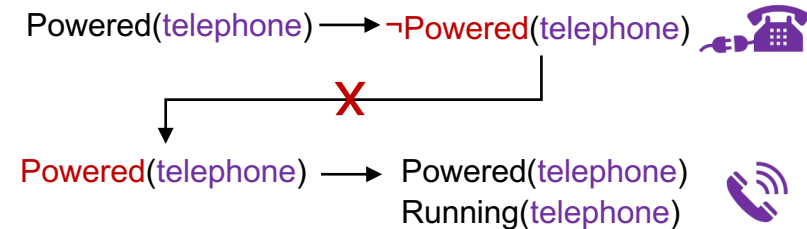
1. Ann sat in the chair.
2. Ann turned off the telephone.
3. Ann picked up a pencil.
4. Ann opened the book.
5. Ann heard the telephone ring.

Which story is more plausible? **A**

Why not **B**?

Conflicting sentences: 2 → 5

Physical states:

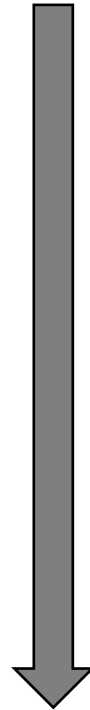
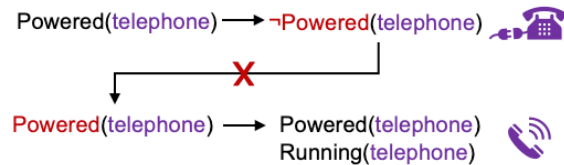


Evaluation Metrics

Story A

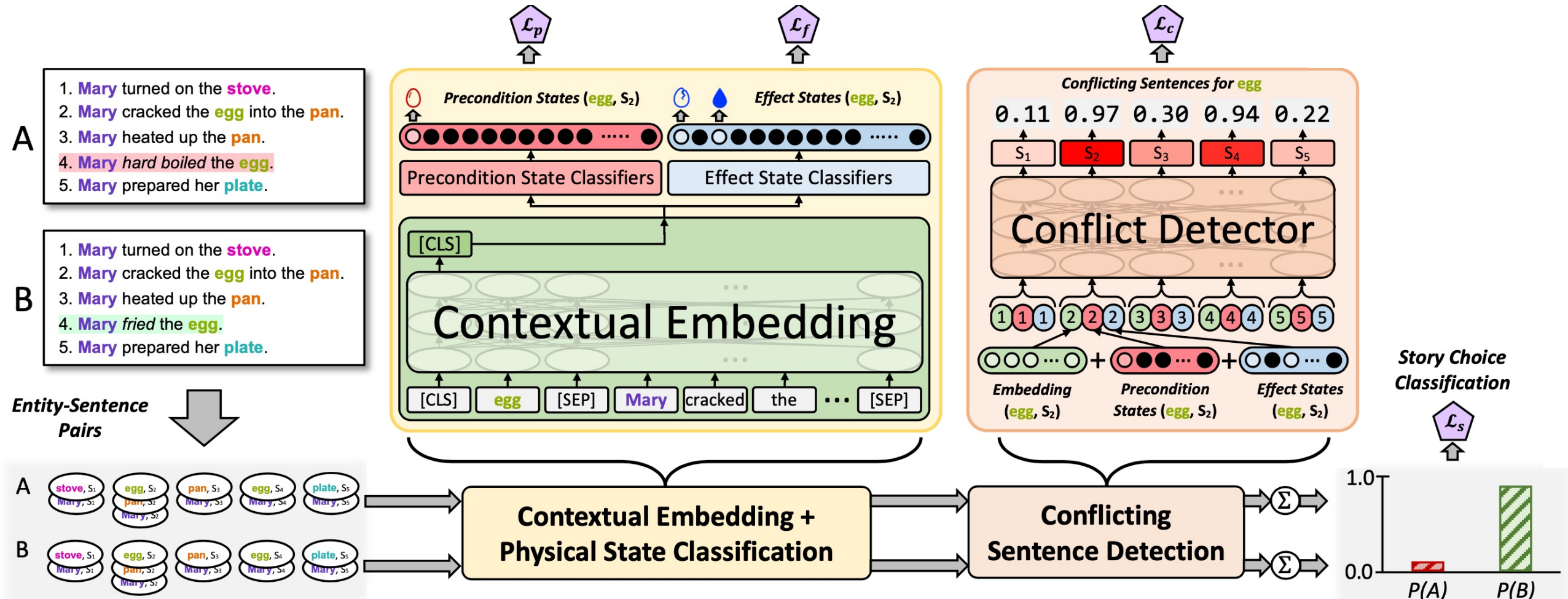
1. Ann sat in the chair.
2. Ann unplugged the telephone.
3. Ann picked up a pencil.
4. Ann opened the book.
5. Ann wrote in the book.

- 2. Ann unplugged the telephone.
3. Ann picked up a pencil.
4. Ann opened the book.
- ! 5. Ann heard the telephone ring.



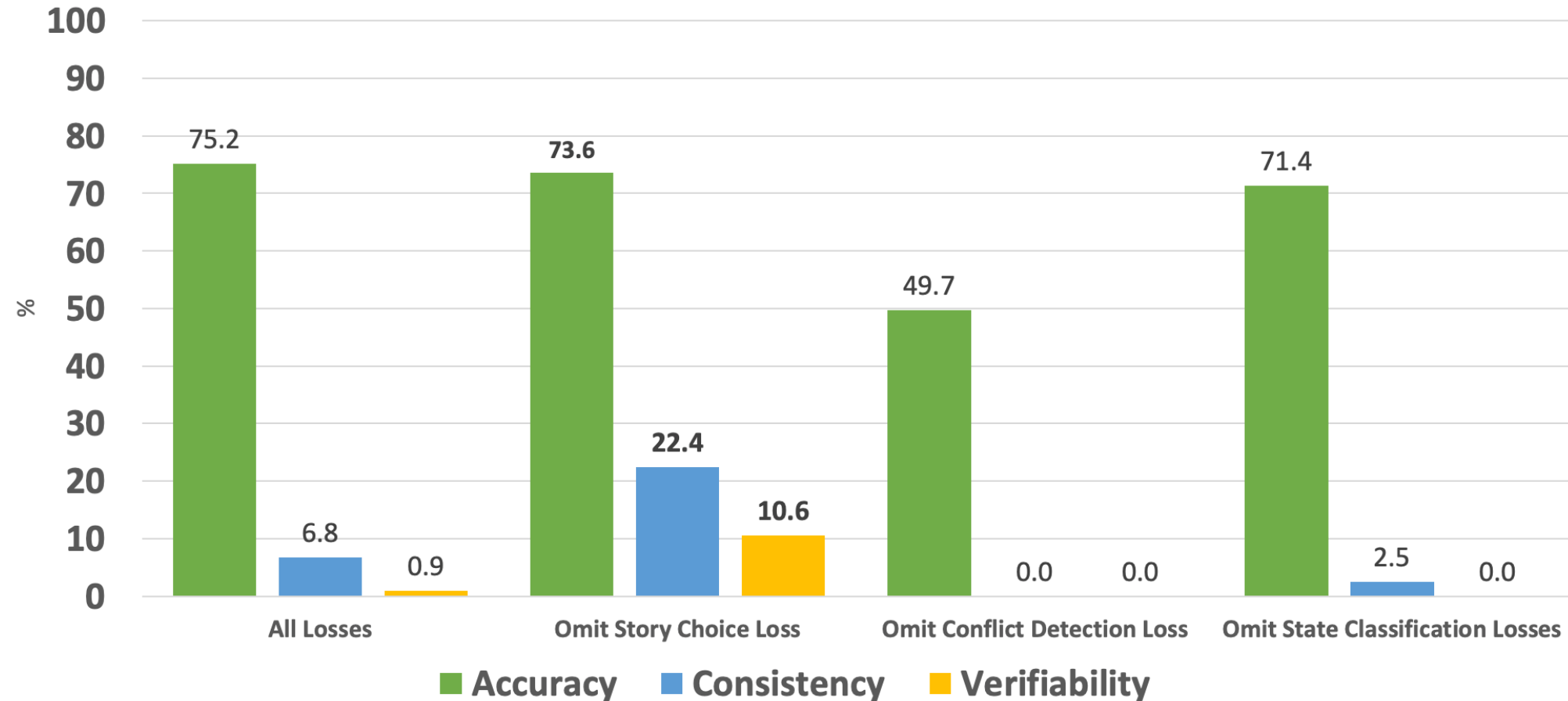
Metric	Story Choice	Conflicting Sentences	Physical States
<i>Accuracy</i>	✓		
<i>Consistency</i>	✓	✓	
<i>Verifiability</i>	✓	✓	✓

Tiered Baseline

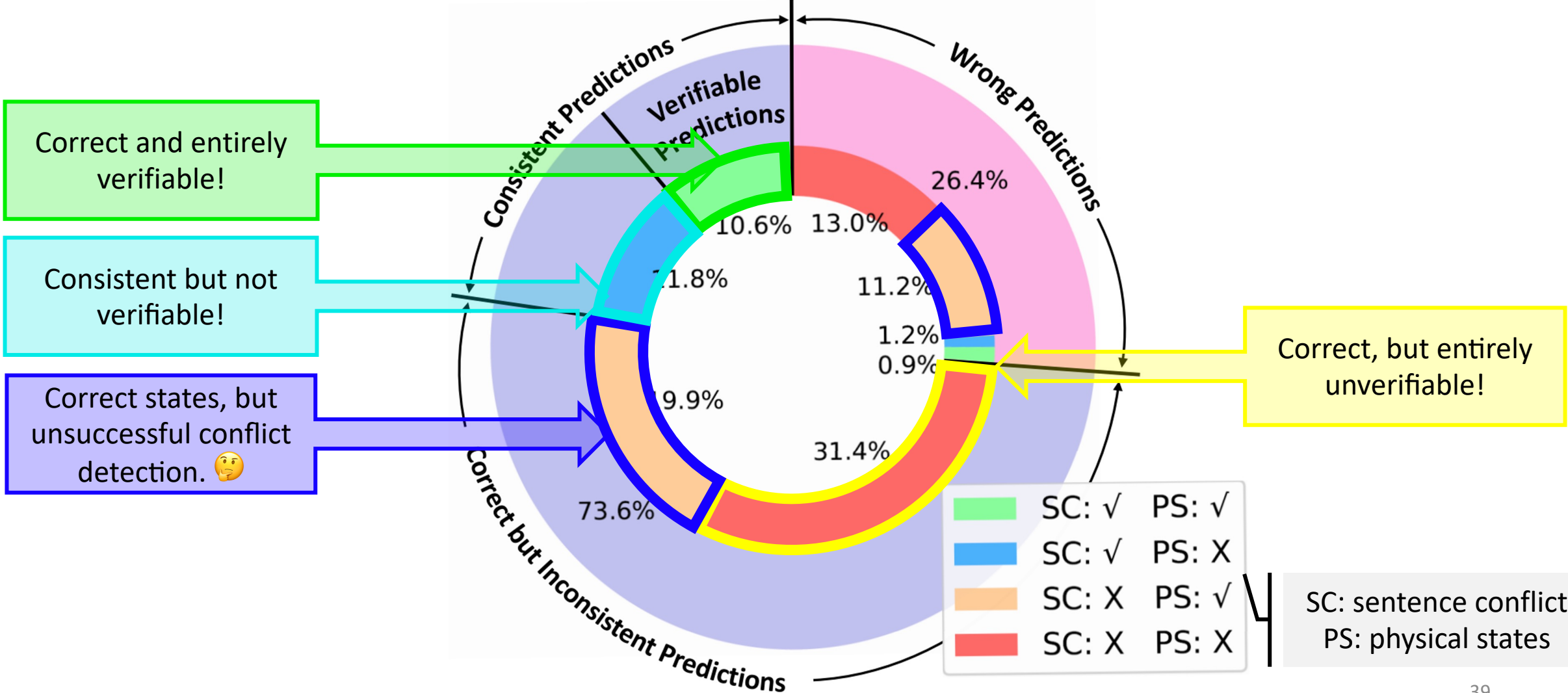


$$\mathcal{L} = \lambda_p \mathcal{L}_p + \lambda_f \mathcal{L}_f + \lambda_c \mathcal{L}_c + \lambda_s \mathcal{L}_s$$

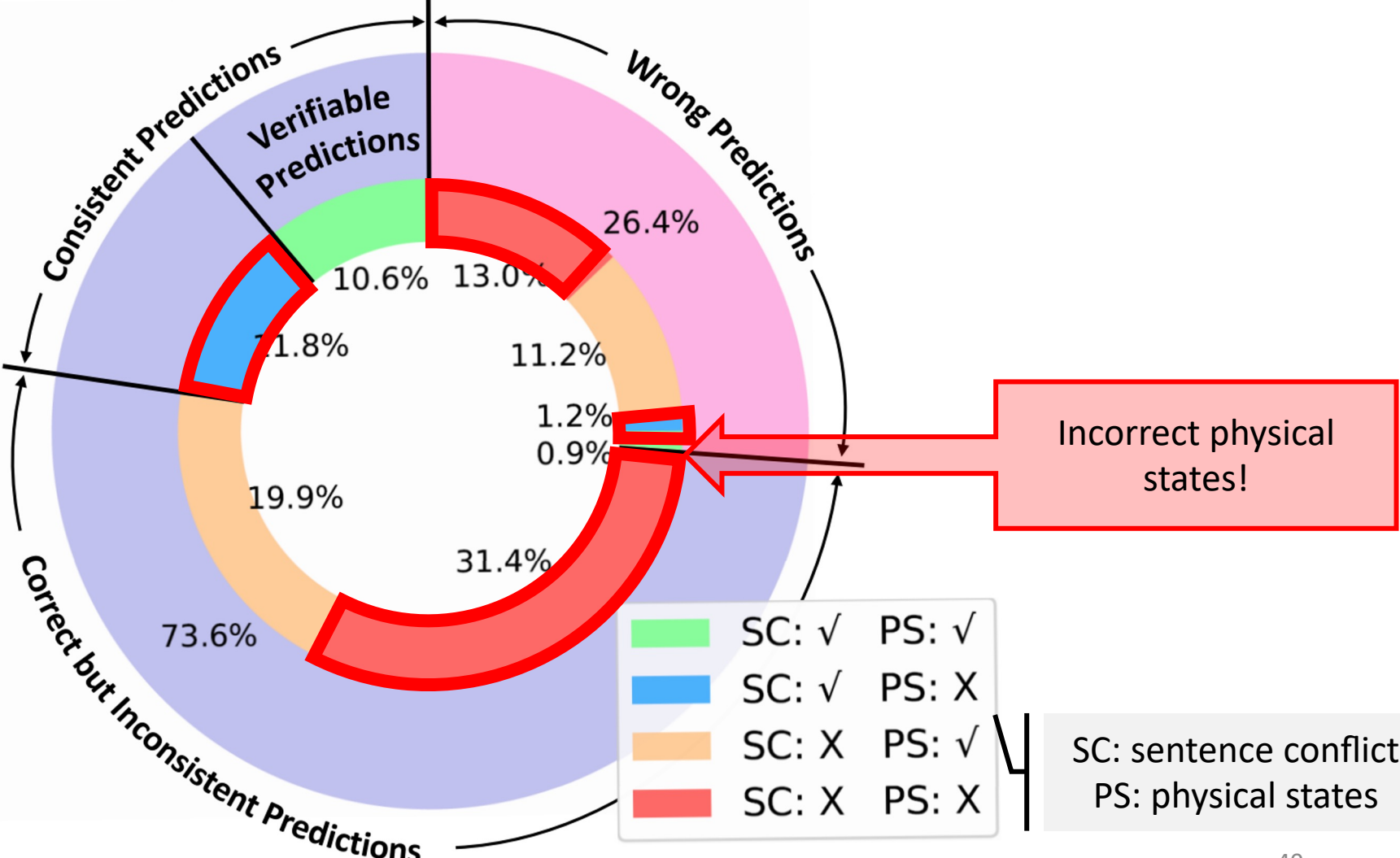
RoBERTa Baseline Results on TRIP



Error Distribution

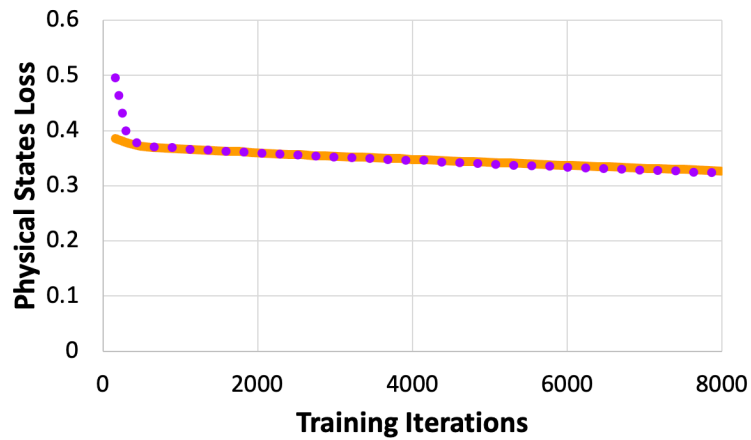


Baseline Results

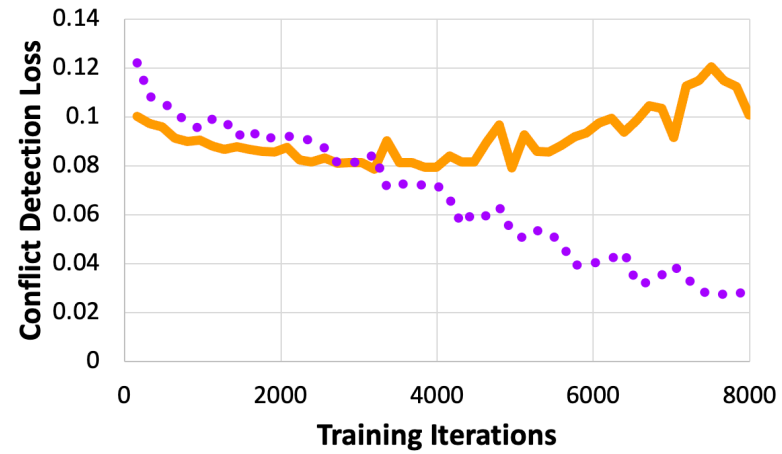


Tiered Task Learning

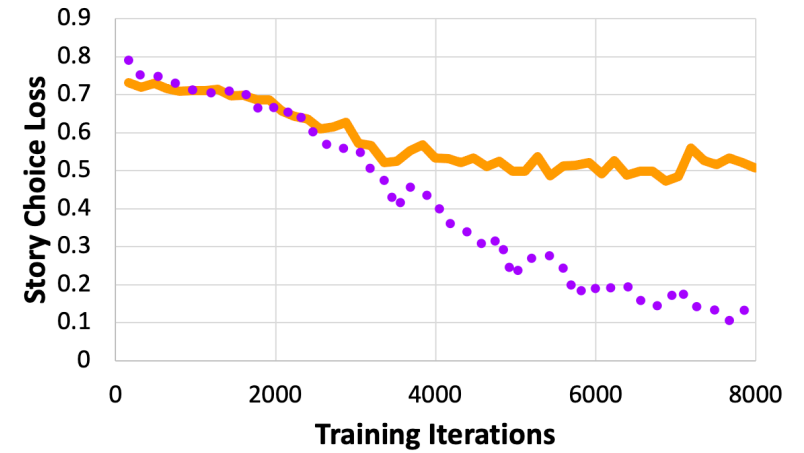
(A)



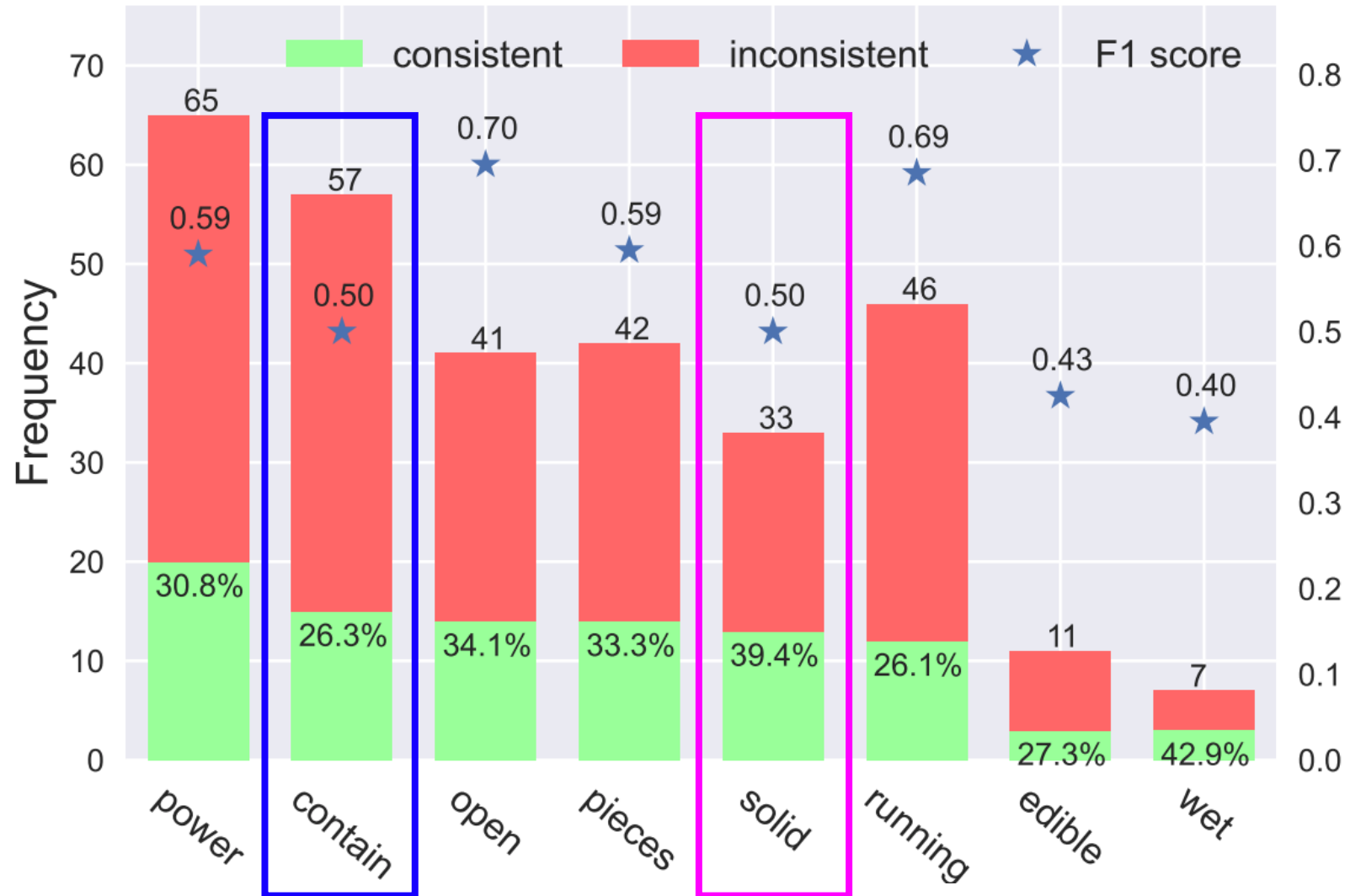
(B)



(C)



Utility of Attributes



Sample System Outputs

1. Tom brought a box to the table. **A**

2. Tom opened the box.

3. Tom took scissors out of the box.

4. Tom cut up the box with the scissors.

5. Tom put the scissors back in the box.

Physical State Predictions

	Preconditions	Effects
S4	\neg Pieces(box) Solid(box)	Pieces(box) Solid(box)
S5	Open(box)	Contain(box) InContainer(scissors)

B

1. Tom brought a box to the table.

2. Tom opened the box.

3. Tom took scissors out of the box.

4. Tom cut up his book with the scissors.

5. Tom put the scissors back in the box.

(a) A verifiable prediction.

1. Ann put the pants and towel in the washing machine. **A**

2. Ann turned the washing machine on.

3. Ann turned on the faucet, and filled the sink with water.

4. Ann put bleach in the water.

5. Ann used the brush to clean the sink.

Physical State Predictions

	Preconditions	Effects
S1	N/A	N/A ⚠️
S2	Power(wm) Running(wm)	Power(wm) Running(wm)

wm: washing machine

B

1. Ann realized that the washing machine was broken.

2. Ann turned the washing machine on.

3. Ann turned on the faucet, and filled the sink with water.

4. Ann put bleach in the water.

5. Ann used the brush to clean the sink.

Error Explanation

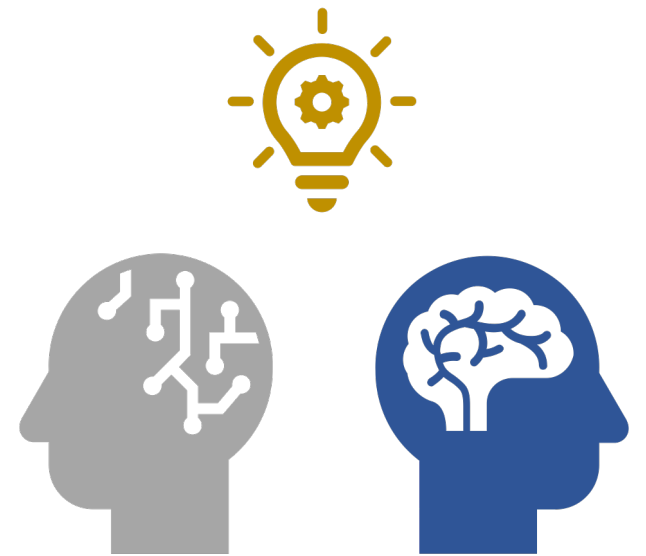
⚠️ Missed detection of \neg Usable(wm)

✖️ Should be \neg Running(wm)

(b) A consistent but not verifiable prediction.

Conclusion

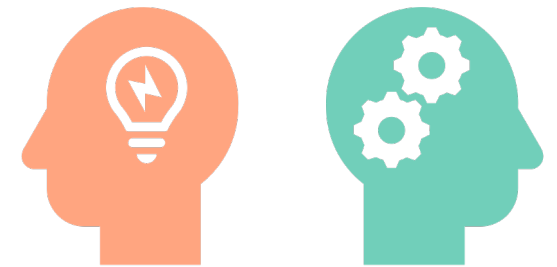
1. Natural language training creates a capacity for abstract analogical reasoning in LLMs!
2. Dual reasoning processes enable LLMs to focus on the correct language context and reason more coherently about the world through language!



Dual Processes of Human Cognition

A line of work theorizes two processes in human reasoning:

- **Heuristic:** fast, intuitive
 - Provides quick intuition for decisions; extracts most relevant info from context
- **Analytic:** slow, deliberative
 - Further operates on relevant info to rationalize and perform inference.
- Can these dual processes similarly strengthen reasoning in PLMs?



2 Tasks for Coherent Physical Commonsense

TRIP

Story A:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary put the cucumber on the counter.
5. Mary ate the donut.

Story B:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary tossed the donut in the trash.
5. Mary ate the donut.

Plausible story: A

Conflicting sentences: (4, 5)

States: inedible(donut) → edible(donut)

Tiered-ProPara

Story A:

1. Air is brought in through the mouth.
2. Passes through the lungs.
3. And into the bronchial tissue.
4. The *carbon dioxide* is removed.
5. The lungs bring the oxygen to the rest of the body.

Story B:

1. *Carbon dioxide* enters the leaves through the stomates by diffusion.
2. Water is transported to the leaves in the xylem.
3. Energy harvested through light reaction is stored by forming ATP.
4. *Carbon dioxide* and energy from ATP are used to create *sugar*.
5. Oxygen exits the leaves through the stomata by diffusion. ...

Carbon dioxide conversion story: B

Carbon dioxide conversion sentence: 4

Carbon dioxide conversion entity: sugar

Heuristic-Analytic Reasoning (HAR)

Language Model Inputs

Story A:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary tossed the donut in the trash.
5. Mary ate the donut.

Story B:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary put the cucumber on the counter.
5. Mary ate the donut.

Language Model Outputs

"Story B is more plausible."

"In Story A, sentences 4 and 5 conflict with each other."

"For sentence 4: After *Mary tossed the donut in the trash ... the donut is now inedible.*"

"For sentence 5: Before *Mary ate the donut ... the donut was edible.*"

Heuristic
Decisions

Analytic
Rationalization

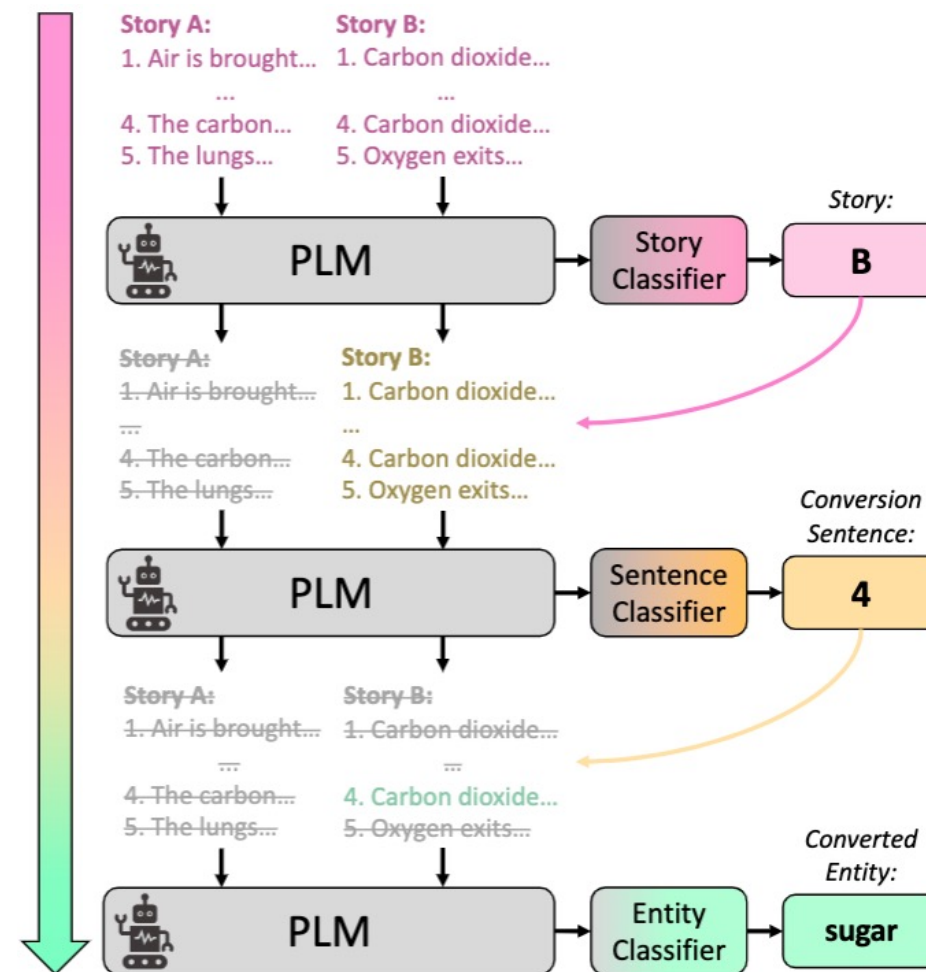


Outline

- **HAR in PLM Fine-Tuning**
- HAR in PLM In-Context Learning

Incorporating HAR into Fine-Tuning

- Coalescing Global & Local Information (**CGLI**):
 - Augments RoBERTa with temporal embedding to capture local information as states change
- Focused CGLI (**FCGLI**):
 - Small improvements to CGLI
- Focused CGLI with Heuristic-Analytic Reasoning (**FCGLI-HAR**):
 - After each prediction is made, delete segments of the context that become irrelevant



Fine-Tuning Results

<i>TRIP</i>				<i>Tiered-ProPara</i>			
<u>Approach</u>	<u>Accuracy</u>	<u>Consistency</u>	<u>Verifiability</u>	<u>Approach</u>	<u>Accuracy</u>	<u>Consistency</u>	<u>Verifiability</u>
RoBERTa	72.9	19.1	9.1	FCGLI	94.5	56.7	36.2
CGLI	94.1	77.3	28.0	FCGLI-HAR	95.1	83.6	57.4
Breakpoint	80.6	53.8	32.4				
FCGLI	93.7	66.2	33.8				
FCGLI-HAR	94.3	75.4	41.1				

Outline

- HAR in PLM Fine-Tuning
- **HAR in PLM In-Context Learning**


Unstructured In-Context Learning (ICL-U)

Story A:


1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary tossed the donut in the trash.
5. Mary ate the donut.

Story B:


1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary put the cucumber on the counter.
5. Mary ate the donut.



Story B is more plausible.



In Story A, **sentences 4 and 5** conflict with each other.



For sentence 4: After *Mary tossed the donut in the trash* ... the **donut** is now **inedible**.
For sentence 5: Before *Mary ate the donut* ... the **donut** was **edible**.

In-Context Learning with Traditional CoT (ICL-CoT)

Story A:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary tossed the donut in the trash.
5. Mary ate the donut.

Story B:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary put the cucumber on the counter.
5. Mary ate the donut.

(CoT)

Story B is more plausible.

(CoT)

In Story A, **sentences 4 and 5** conflict with each other.

(CoT)

For sentence 4: After *Mary tossed the donut in the trash* ... the **donut** is now **inedible**.
For sentence 5: Before *Mary ate the donut* ... the **donut** was **edible**.



In-Context Learning with HAR (ICL-HAR)

Story A:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary tossed the donut in the trash.
5. Mary ate the donut.

Story B:

1. Mary went to the fridge.
2. Mary took out a bowl from the fridge.
3. The bowl had a cucumber and a donut in it.
4. Mary put the cucumber on the counter.
5. Mary ate the donut.

Story B is more plausible.

In Story A, **sentences 4 and 5** conflict with each other.

For sentence 4: After *Mary tossed the donut in the trash* ... the **donut** is now **inedible**.
For sentence 5: Before *Mary ate the donut* ... the **donut** was **edible**.

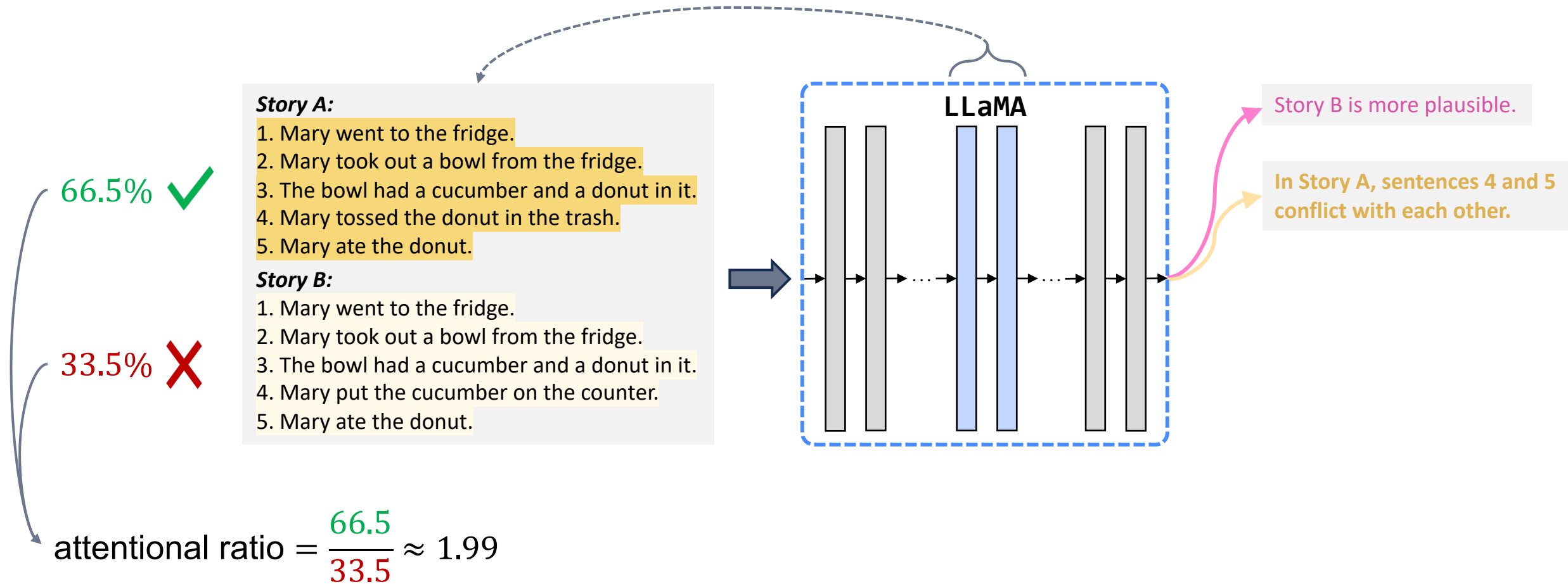


In-Context Learning Results

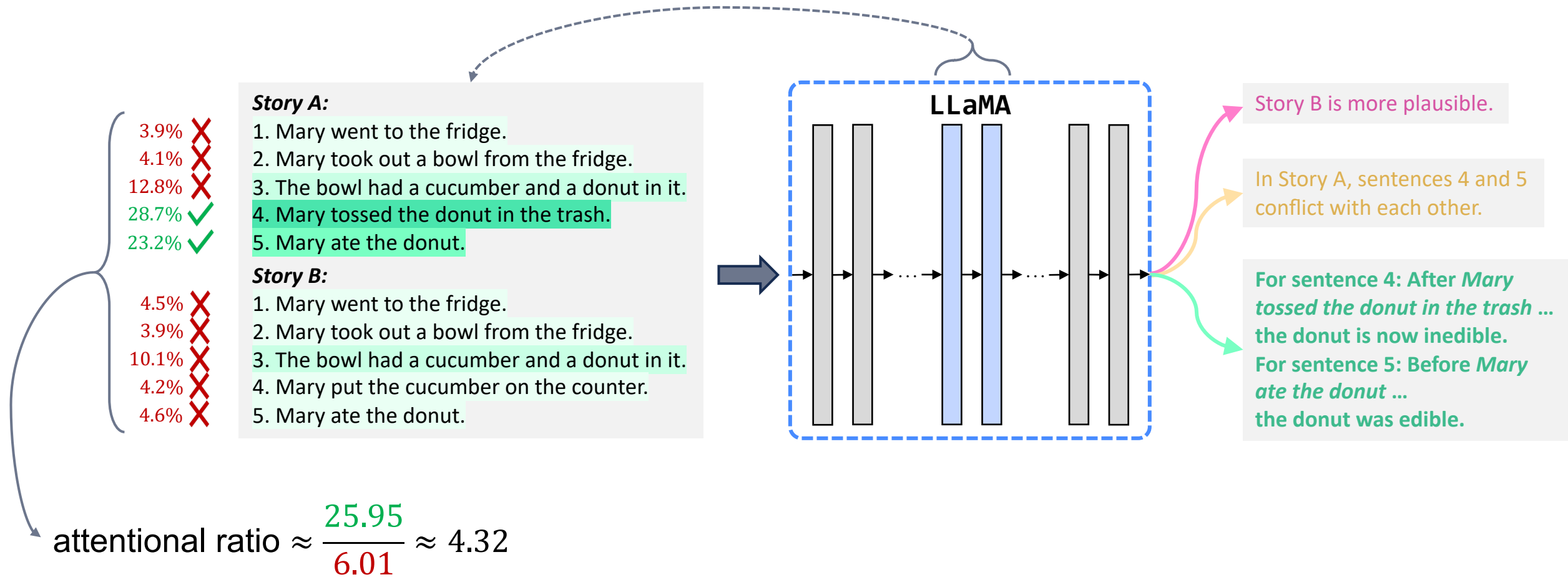
<i>InstructGPT</i>						
Approach	TRIP			Tiered-ProPara		
	<i>Acc.</i>	<i>Cons.</i>	<i>Ver.</i>	<i>Acc.</i>	<i>Cons.</i>	<i>Ver.</i>
ICL-U	70.9	40.7	7.1	54.9	17.4	5.2
ICL-CoT	75.0	40.7	10.8	50.7	19.2	7.5
ICL-HAR	72.6	47.9	23.9	54.9	31.5	20.7

<i>LLaMA</i>						
Approach	TRIP			Tiered-ProPara		
	<i>Acc.</i>	<i>Cons.</i>	<i>Ver.</i>	<i>Acc.</i>	<i>Cons.</i>	<i>Ver.</i>
ICL-U	70.4	42.3	14.8	51.2	3.8	1.4
ICL-CoT	74.6	42.3	19.7	57.3	9.4	4.2
ICL-HAR	55.6	44.4	35.2	41.8	17.8	13.1

Attention Analysis



Attention Analysis



Attentional Precision and Recall

- To measure how attended context and correct predictions correlate, we use **attentional precision** and **attentional recall**
 - *True/false positive*: Correct attention, and correct/incorrect prediction
 - *True/false negative*: Incorrect attention, and correct/incorrect prediction

Attention Analysis Results

- PLMs focus better on the correct language context during each step of reasoning
- Faithful attention and coherent reasoning go hand in hand!

<i>Sentence Selection Step</i>						
<u>Approach</u>	TRIP			Tiered-ProPara		
	<u>Ratio</u>	<u>Prec.</u>	<u>Rec.</u>	<u>Ratio</u>	<u>Prec.</u>	<u>Rec.</u>
ICL-U	0.96	42.6	39.6	0.90	14.8	30.6
ICL-HAR	1.07	75.2	48.7	1.80	51.1	58.2

<i>Physical State Prediction Step</i>						
<u>Approach</u>	TRIP			Tiered-ProPara		
	<u>Ratio</u>	<u>Prec.</u>	<u>Rec.</u>	<u>Ratio</u>	<u>Prec.</u>	<u>Rec.</u>
ICL-U	1.23	43.0	35.4	1.21	14.6	25.9
ICL-HAR	1.95	79.8	98.2	2.20	72.1	83.3

Conclusion

- Human-inspired heuristic-analytic reasoning helps PLMs reason more coherently when applied to downstream tasks
- Successful because it helps PLMs focus on the correct language context at each step of reasoning
- Check out our paper for more details and results!



From Heuristic to Analytic: Cognitively Motivated Strategies for Coherent Physical Commonsense Reasoning

Zheyuan Zhang^{1*}

Shane Storks^{1*}

Fengyuan Hu¹

Sungryull Sohn²

Moontae Lee^{2,3}

Honglak Lee^{1,2}

Joyce Chai¹

¹University of Michigan

²LG AI Research

³University of Illinois at Chicago

{zheyuan, sstorks, hufy, chaijy}@umich.edu

{srsohn, moontae.lee, honglak}@lgresearch.ai

Thank you!

