

# Confidence-based Rephrasing, Refinement, and Selection

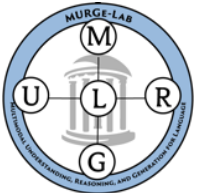
UncertaiNLP Workshop, EACL 2024

Elias Stengel-Eskin

03/22/2024



THE UNIVERSITY  
*of* NORTH CAROLINA  
at CHAPEL HILL



# Outline

## Part I: Uncertainty in Human-Model Interactions

**Calibrated Interpretation: Confidence Estimation in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, TACL (2023)

**Did You Mean...? Confidence-based Trade-offs in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, EMNLP (2023)

## Part II: Model-based Selection to Reduce Uncertainty

**Rephrase, Augment, Reason: Visual Grounding of Questions for Vision-Language Models**, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal, ICLR (2024)

## Part III: Confidence for Model-Model Interactions

**ReConcile: Round-Table Conference Improves Reasoning via Consensus among Diverse LLMs**, Justin Chih-Yao Chen, Swarnadeep Saha, Mohit Bansal (2024)

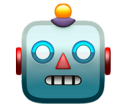
**MAGDi: Structured Distillation of Multi-Agent Interaction Graphs Improves Reasoning in Smaller Language Models**, Justin Chih-Yao Chen\*, Swarnadeep Saha\*, Elias Stengel-Eskin Mohit Bansal (2024)

# Executable parsing

## Predicting executable programs



*Do I have anything going on tonight?*



```
(Yield (> (size
  (QueryEventResponse.results
    (... (EventOnDateWithTimeRange
      (EventOnDate (Today)
        (^ (Event) EmptyStructConstraint))
        (Night)))) 0L))
```

# Executable parsing

Predicting executable programs

Language as an API for interaction



# Executable parsing

Predicting executable programs

Language as an API for interaction

## Domains

Querying: Zelle and Mooney (1996), Berant et al. (2013), Yu et al. (2018)

Digital assistants: Semantic Machines (2020), Damonte et al. (2019)

Robotics: Kate et al. (2005), Tellex et al. (2011), Artzi et al. (2013), Tellex et al. (2020)

# Why semantic parsing?

## Form vs. meaning

Focus on semantics

Restricted output space

**Uncertainty in natural language generation: From theory to applications**

Baan et al., 2023

**Interpreting Predictive Probabilities: Model Confidence or Human Label Variation?**

Baan et al., 2024

# Why semantic parsing?

## Form vs. meaning

Focus on semantics

Restricted output space

## Measurability

When executable, accuracy is well-defined

**Uncertainty in natural language generation: From theory to applications**

Baan et al., 2023

**Interpreting Predictive Probabilities: Model Confidence or Human Label Variation?**

Baan et al., 2024

# Why semantic parsing?

## Form vs. meaning

Focus on semantics

Restricted output space

## Measurability

When executable, accuracy is well-defined

## Safety concerns

Agents doing things in the real world

**Uncertainty in natural language generation: From theory to applications**

Baan et al., 2023

**Interpreting Predictive Probabilities: Model Confidence or Human Label Variation?**

Baan et al., 2024



# Execution and safety



# Execution and safety

$A = \text{drop\_item}()$        $\neg A = \text{do\_nothing}()$

# Execution and safety

$A = \text{drop\_item}()$

$G = \text{drop\_item}$

$\neg A = \text{do\_nothing}()$

$\neg G = \text{do\_nothing}$

# Execution and safety

$A = \text{drop\_item}()$

$G = \text{drop\_item}$

$I = \textit{let's drop it}$

$\neg A = \text{do\_nothing}()$

$\neg G = \text{do\_nothing}$

# Execution and safety

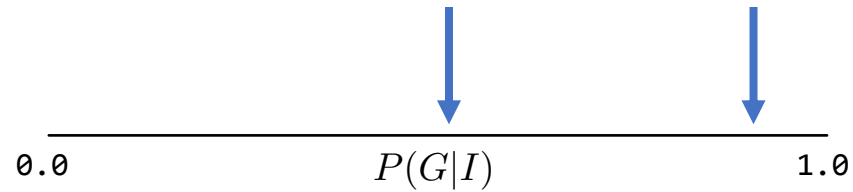
$A = \text{drop\_item}()$

$G = \text{drop\_item}$

$I = \textit{let's drop it}$

$\neg A = \text{do\_nothing}()$

$\neg G = \text{do\_nothing}$



# Execution and safety

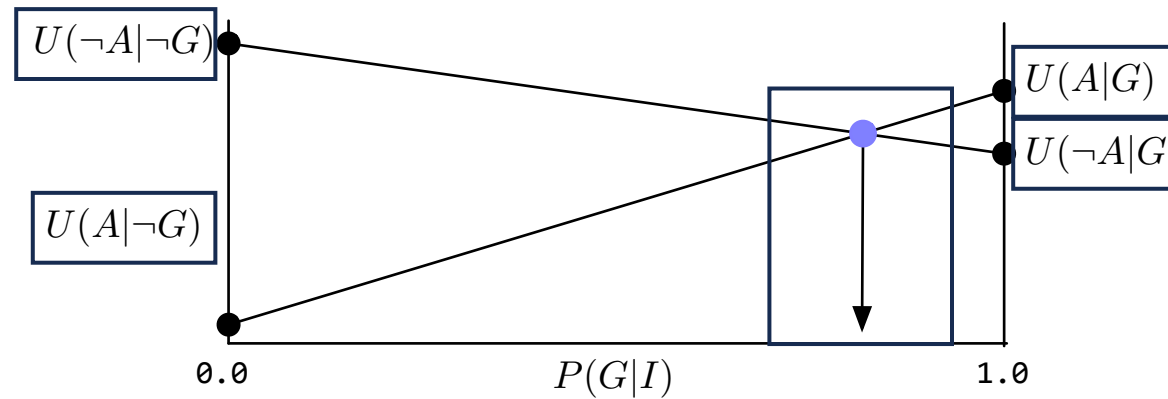
$A = \text{drop\_item}()$

$\neg A = \text{do\_nothing}()$

$G = \text{drop\_item}$

$\neg G = \text{do\_nothing}$

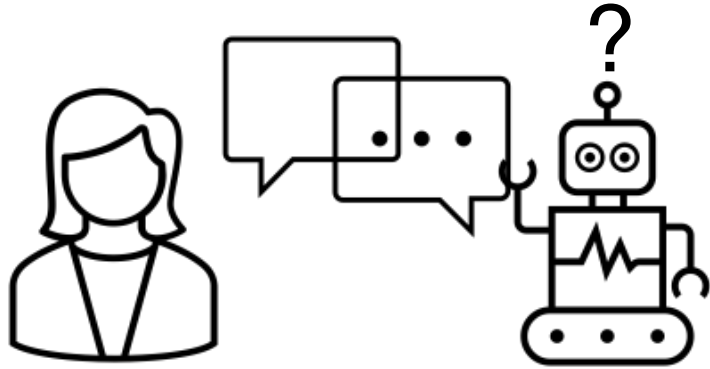
$I = \textit{let's drop it}$



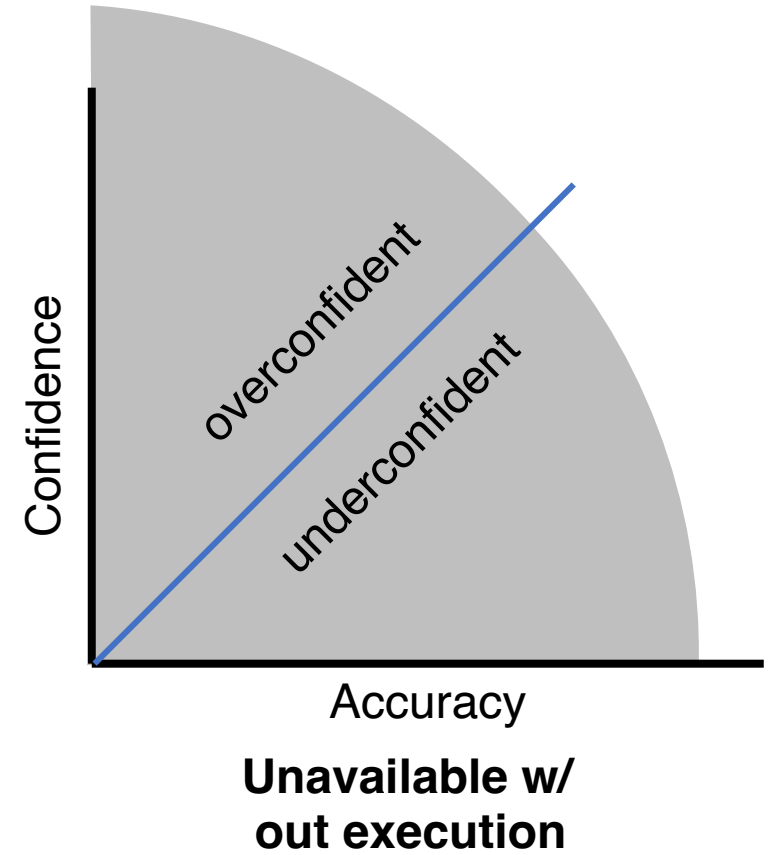
**Principles of Mixed-Initiative User Interfaces**

**Eric Horvitz**  
Microsoft Research

# Calibration



**Available  
w/ out  
execution**



# Outline

## Part I: Uncertainty in Human-Model Interactions

**Calibrated Interpretation: Confidence Estimation in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, TACL (2023)

**Did You Mean...? Confidence-based Trade-offs in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, EMNLP (2023)

## Part II: Model-based Selection to Reduce Uncertainty

**Rephrase, Augment, Reason: Visual Grounding of Questions for Vision-Language Models**, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal, ICLR (2024)

## Part III: Confidence for Model-Model Interactions

**ReConcile: Round-Table Conference Improves Reasoning via Consensus among Diverse LLMs**, Justin Chih-Yao Chen, Swarnadeep Saha, Mohit Bansal (2024)

**MAGDi: Structured Distillation of Multi-Agent Interaction Graphs Improves Reasoning in Smaller Language Models**, Justin Chih-Yao Chen\*, Swarnadeep Saha\*, Elias Stengel-Eskin Mohit Bansal (2024)





# Goals

**How to extract confidence from models?**

**How well-calibrated are semantic parsing models?**

# Datasets: digital assistant

## SMCalFlow

Calendar domain

Lisp-like programs

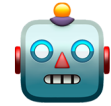
## TreeDST

several domains

same format



*Do I have anything going on tonight?*



```
(Yield (> (size
  (QueryEventResponse.results
    (... (EventOnDateWithTimeRange
      (EventOnDate (Today)
        (^ (Event) EmptyStructConstraint))
        (Night)))))) 0L))
```

# Measuring Calibration Error

$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$

where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence

# Measuring Calibration Error

1. Obtain word-level confidence (min over token probabilities)
2. Bin by confidence
3. Compute average accuracy per bin
4. Expected calibration error (ECE) is difference between confidence and accuracy

S	EL	ECT
0.92	0.31	0.85

SELECT  $\rightarrow$  0.31

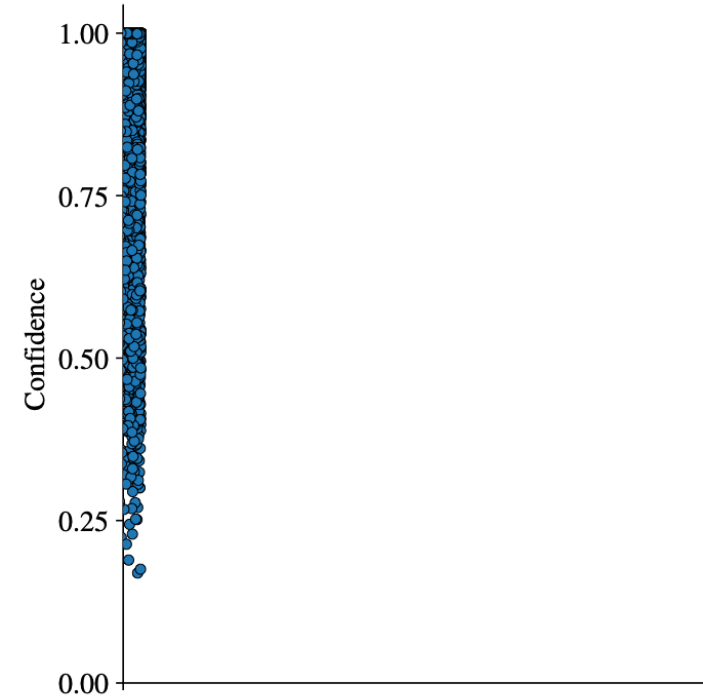
$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$

where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence

# Measuring Calibration Error

1. Obtain word-level confidence (min over token probabilities)
2. Bin by confidence
3. Compute average accuracy per bin
4. Expected calibration error (ECE) is difference between confidence and accuracy

$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$

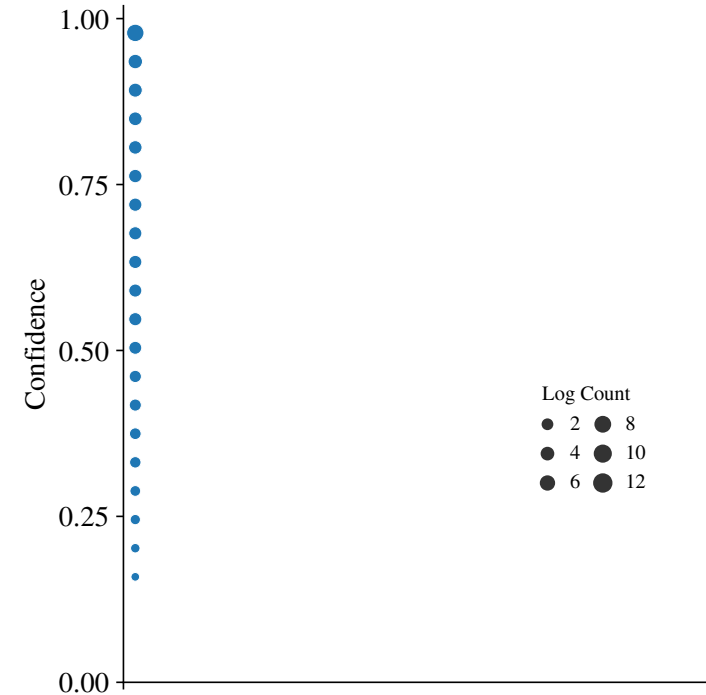


where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence

# Measuring Calibration Error

1. Obtain word-level confidence (min over token probabilities)
- 2. Bin by confidence**
3. Compute average accuracy per bin
4. Expected calibration error (ECE) is difference between confidence and accuracy

$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$

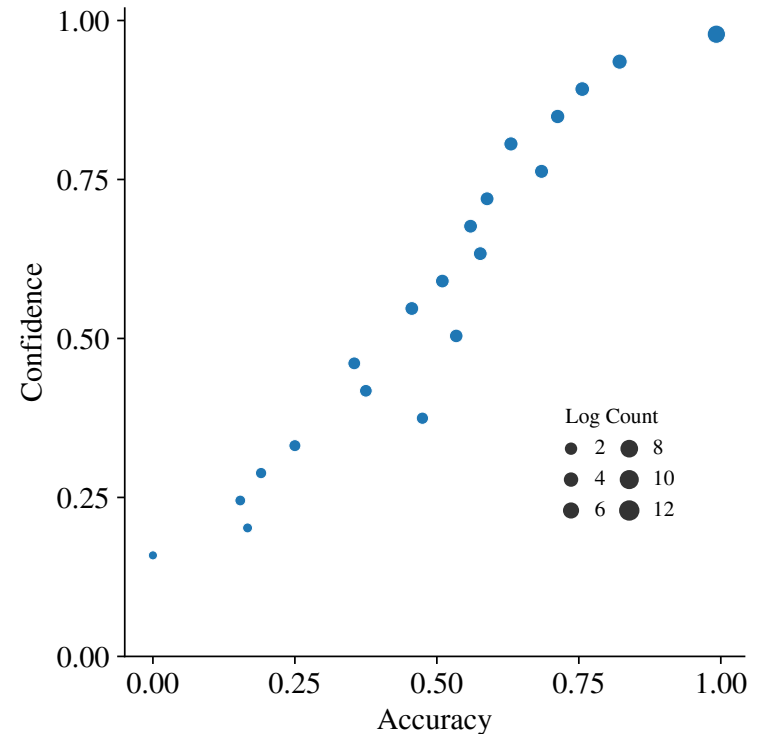


where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence

# Measuring Calibration Error

1. Obtain word-level confidence (min over token probabilities)
2. Bin by confidence
- 3. Compute average accuracy per bin**
4. Expected calibration error (ECE) is difference between confidence and accuracy

$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$

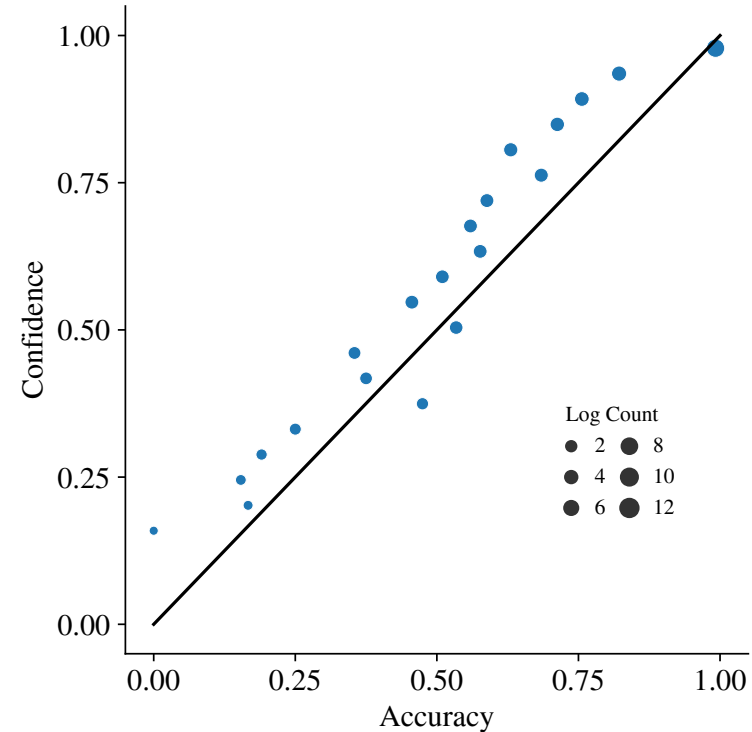


where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence

# Measuring Calibration Error

1. Obtain word-level confidence (min over token probabilities)
2. Bin by confidence
3. Compute average accuracy per bin
4. **Expected calibration error (ECE) is difference between confidence and accuracy**

$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$



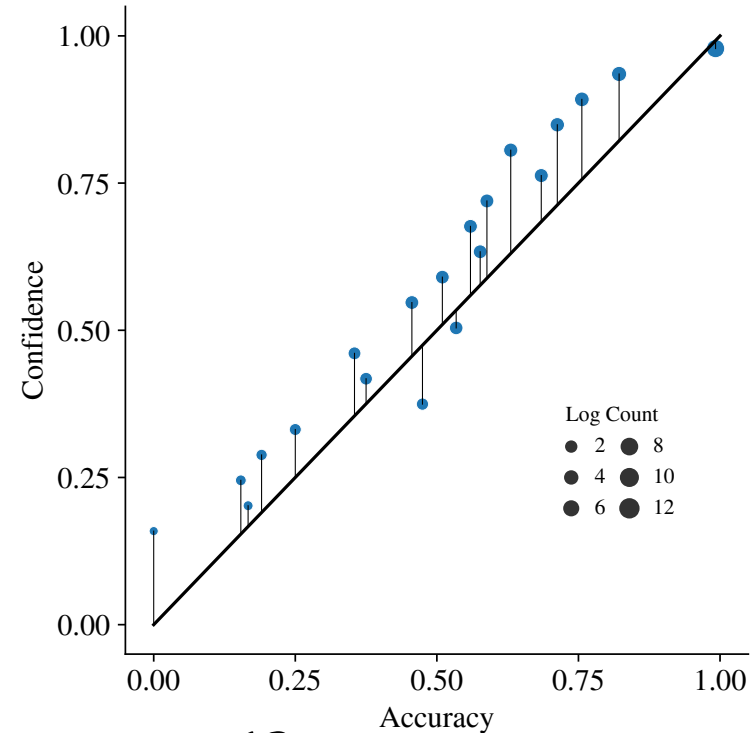
where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence



# Measuring Calibration Error

1. Obtain word-level confidence (min over token probabilities)
2. Bin by confidence
3. Compute average accuracy per bin
4. **Expected calibration error (ECE) is difference between confidence and accuracy**

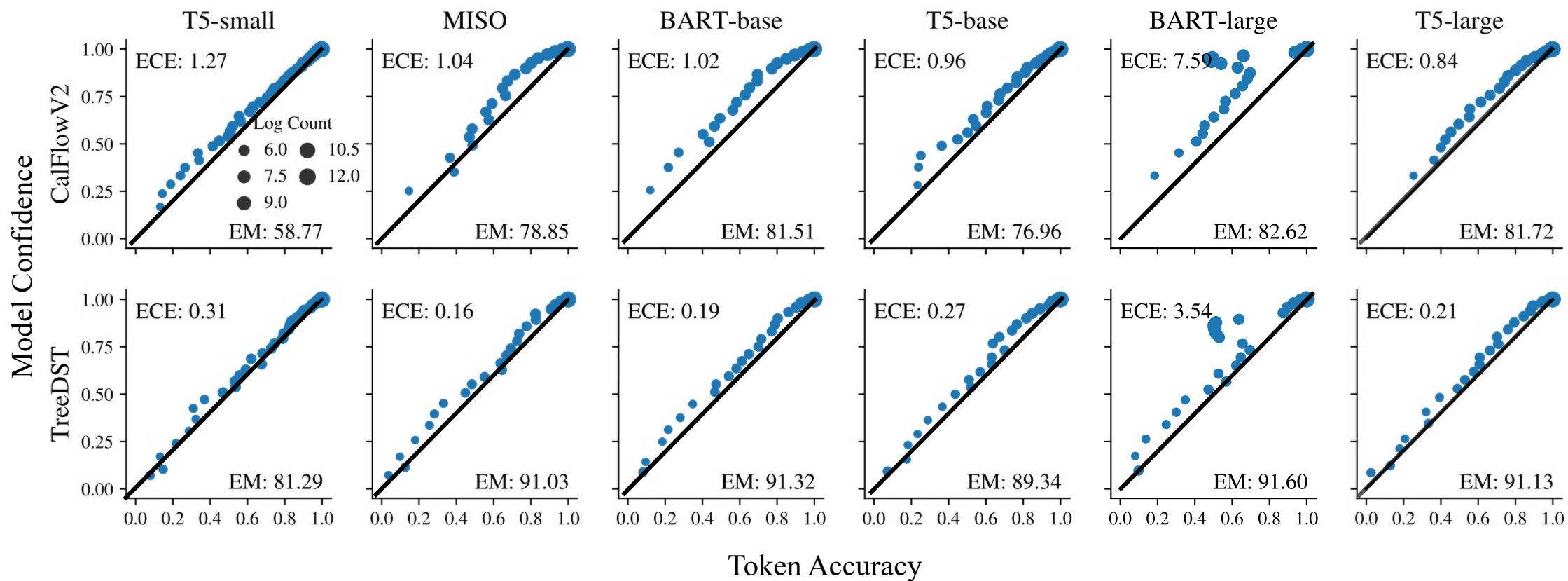
$$ECE(\mathcal{B}) = \sum_{i=1}^N \frac{|\mathcal{B}_i|}{N} \left| \frac{\sum_{j \in \mathcal{B}_i} a_j}{|\mathcal{B}_i|} - \frac{\sum_{j \in \mathcal{B}_i} c_j}{|\mathcal{B}_i|} \right|$$



where  $\mathcal{B}$  are the  $N$  bins,  
 $a_i$  is the accuracy (0 or 1) and  
 $c_j$  is the model confidence

# How Calibrated Are Semantic Parsing Models?

# How Calibrated Are Semantic Parsing Models?



# Calibration in Semantic Parsing

**Models surprisingly well-calibrated**

Seq2seq + finetuning seems sufficient

# Calibration in Semantic Parsing

Models surprisingly well-calibrated

**More experiments/takeaways in the paper**

SQL vs digital assistant domains

Execution accuracy

Token frequency

Perplexity and calibration

Few-shot LLMs

# Outline

## Part I: Uncertainty in Human-Model Interactions

**Calibrated Interpretation: Confidence Estimation in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, TACL (2023)

**Did You Mean...? Confidence-based Trade-offs in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, EMNLP (2023)

## Part II: Model-based Selection to Reduce Uncertainty

**Rephrase, Augment, Reason: Visual Grounding of Questions for Vision-Language Models**, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal, ICLR (2024)

## Part III: Confidence for Model-Model Interactions

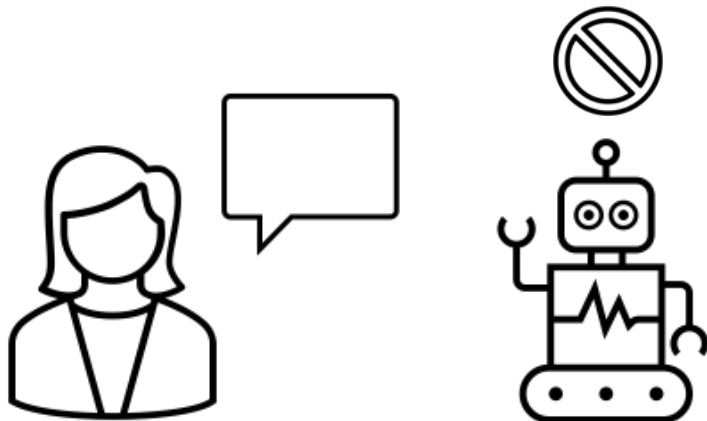
**ReConcile: Round-Table Conference Improves Reasoning via Consensus among Diverse LLMs**, Justin Chih-Yao Chen, Swarnadeep Saha, Mohit Bansal (2024)

**MAGDi: Structured Distillation of Multi-Agent Interaction Graphs Improves Reasoning in Smaller Language Models**, Justin Chih-Yao Chen\*, Swarnadeep Saha\*, Elias Stengel-Eskin Mohit Bansal (2024)



# Selective Prediction

When do we accept or reject a prediction?



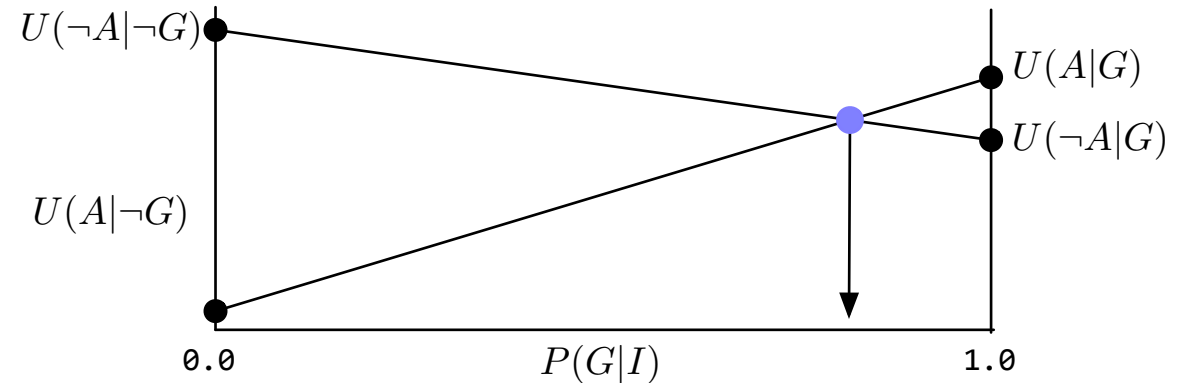
$A = \text{drop\_item}()$

$\neg A = \text{do\_nothing}()$

$G = \text{drop\_item}$

$\neg G = \text{do\_nothing}$

$I = \textit{let's drop it}$



Selective classification for deep neural networks  
Y. Geifman and R. El-Yaniv, 2017

# Problem: False rejection

## False positives harm safety

Unintended consequences

## False negatives harm usability

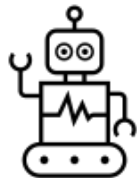
Lack of action means increased frustration



Please set an alarm for 8am



*Sure, I have set the alarm*



#\$!@#!



Please set an alarm for 8am

*Sorry, I'm not sure what you mean*



#\$!@#!





# Usability vs. Safety

**Usability: executing instructions**

**Safety: not making mistakes**

## Full usability



Do I have anything going on today?

(Yield (< (...



Execute program

## Full safety



Do I have anything going on today?

(Yield (< (...



Reject program

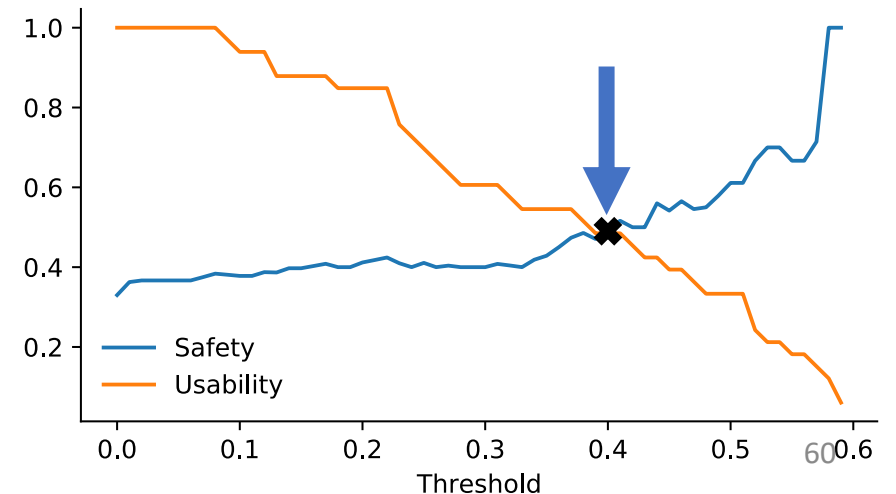
## Threshold

Confidence  $< x$

Reject program

Confidence  $\geq x$

Execute program



# Usability vs. Safety

## Threshold

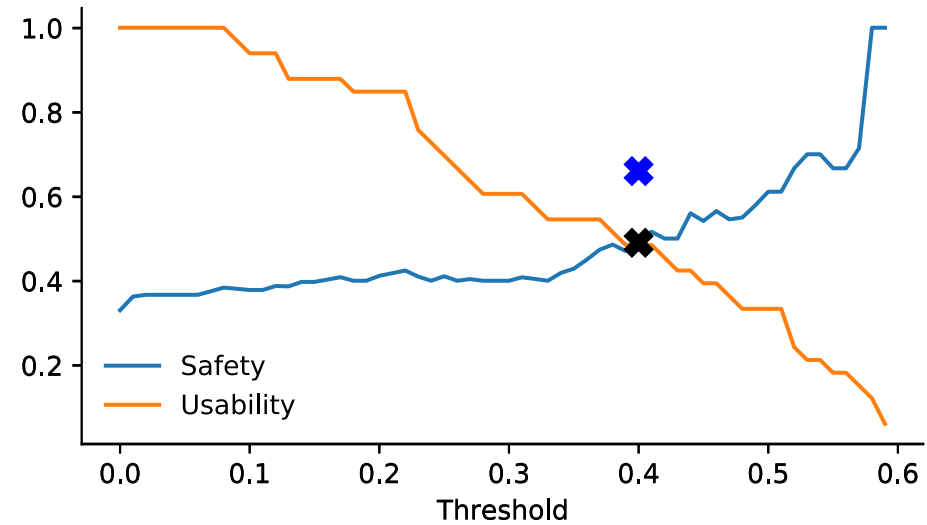
Confidence  $< x$

Reject program

Confidence  $\geq x$

Execute program

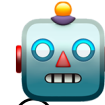
## Human-in-the-loop



# DidYouMean

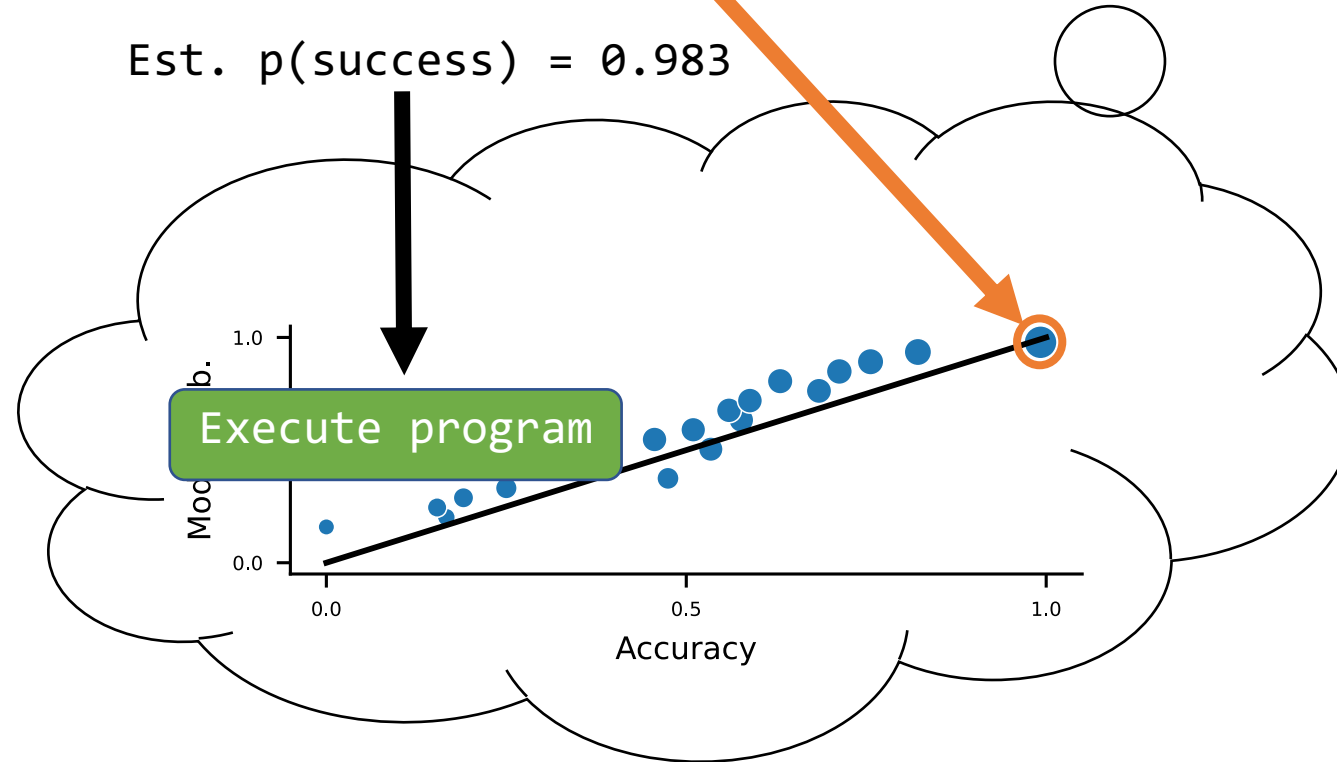


Do I have anything going on tonight?



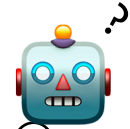
Prediction: (Yield (< (... (EventOnDateWithTimeRange...))))  
with confidence 0.98

Est.  $p(\text{success}) = 0.983$



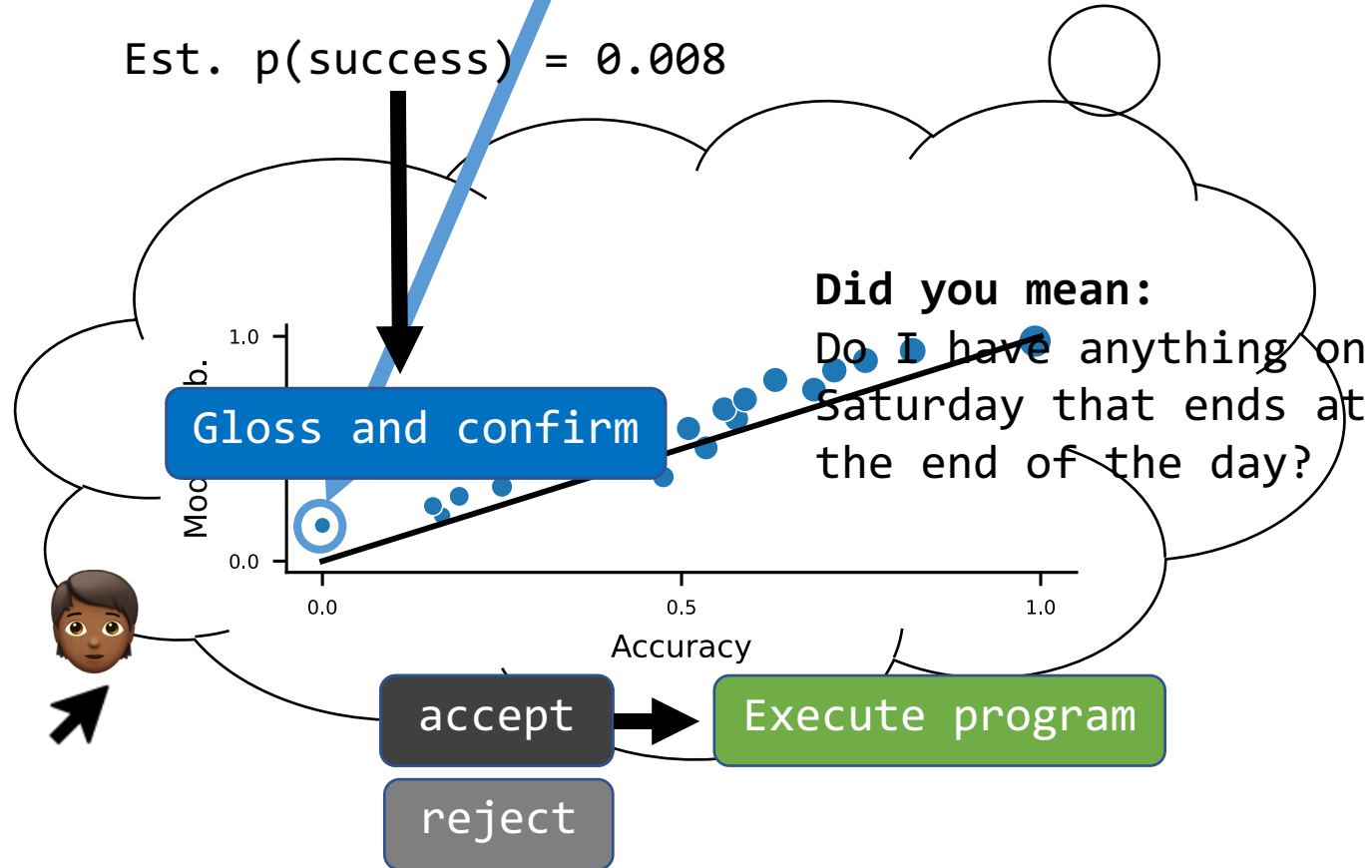


Do I have any things that end at EOD on Sat?

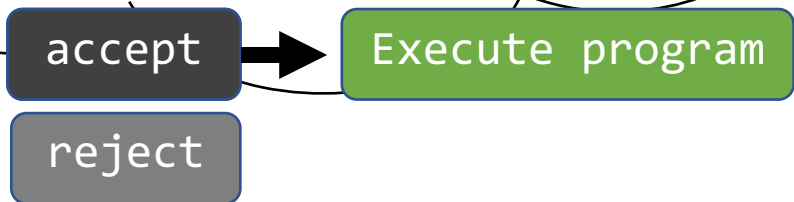


Prediction: (Yield (<(...(EventOnDateBeforeTime...))))  
with confidence 0.012

Est. p(success) = 0.008



Did you mean:  
Do I have anything on  
Saturday that ends at  
the end of the day?



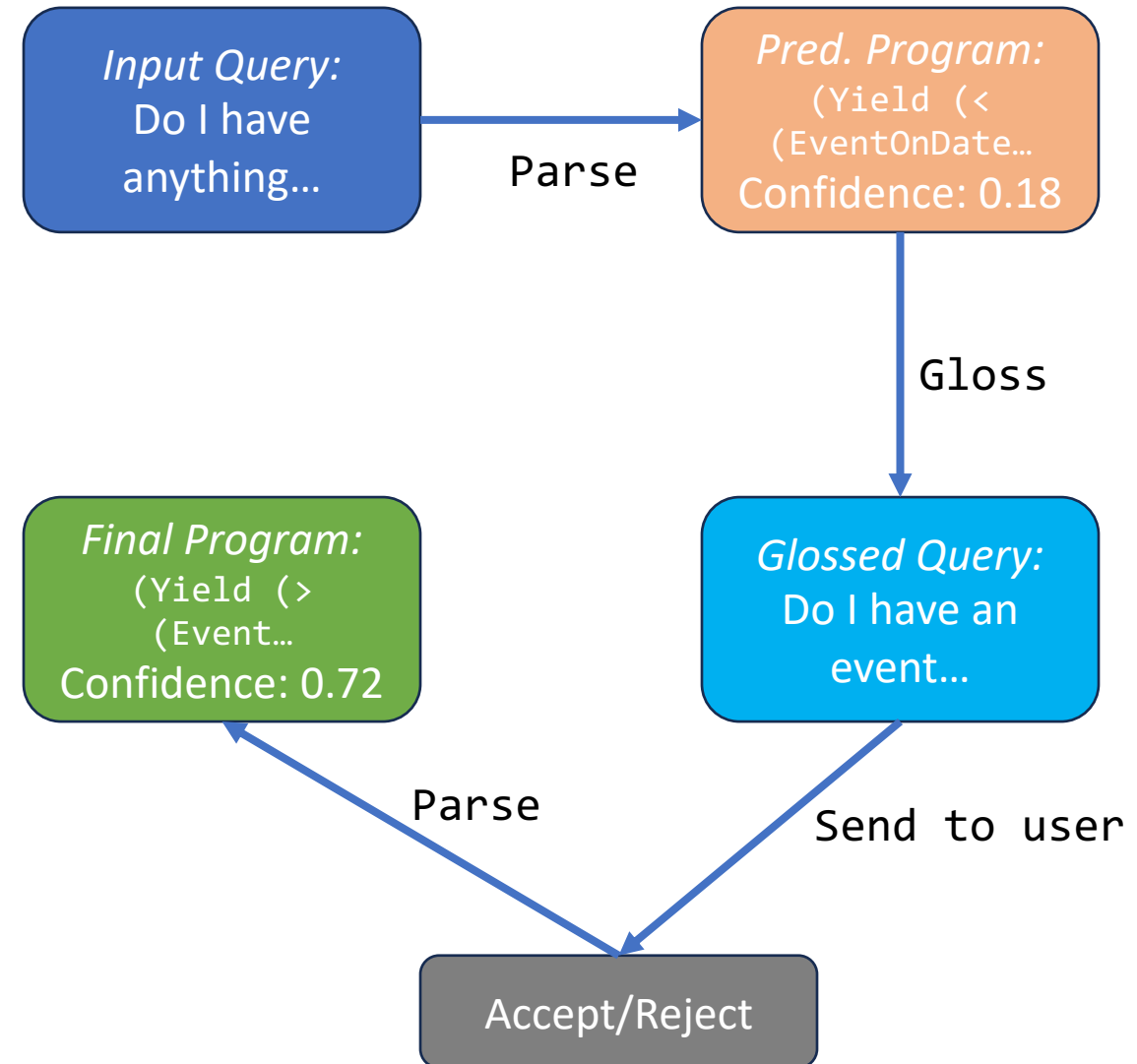
# DidYouMean details

## Parsing model

Translates queries to programs  
MISO (Stengel-Eskin et al., 2022)  
Finetuned seq2graph model

## Glossing model

Translates programs to queries  
BART-large seq2seq  
Finetuned on reverse data



When More Data Hurts: A Troubling Quirk in Developing Broad-Coverage Natural Language Understanding Systems  
Stengel-Eskin et al., 2022

# DidYouMean Results

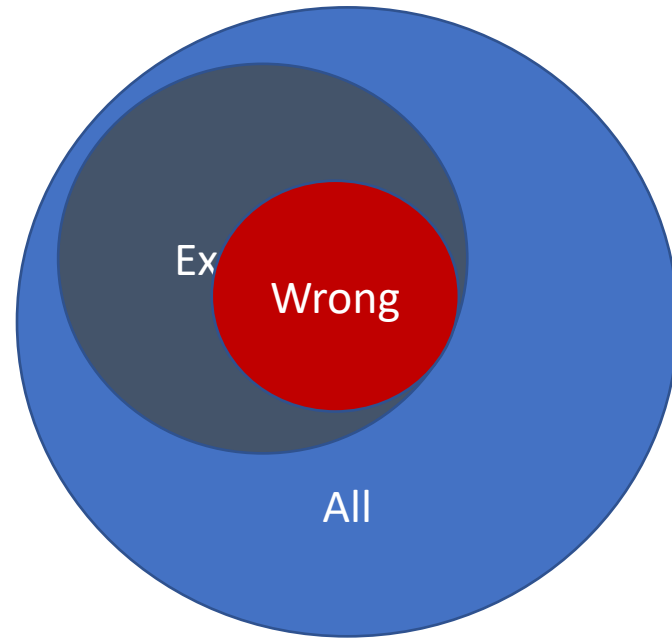
## User study

Mturk annotators, 100 examples (below 0.6 confidence)

## Coverage and Risk

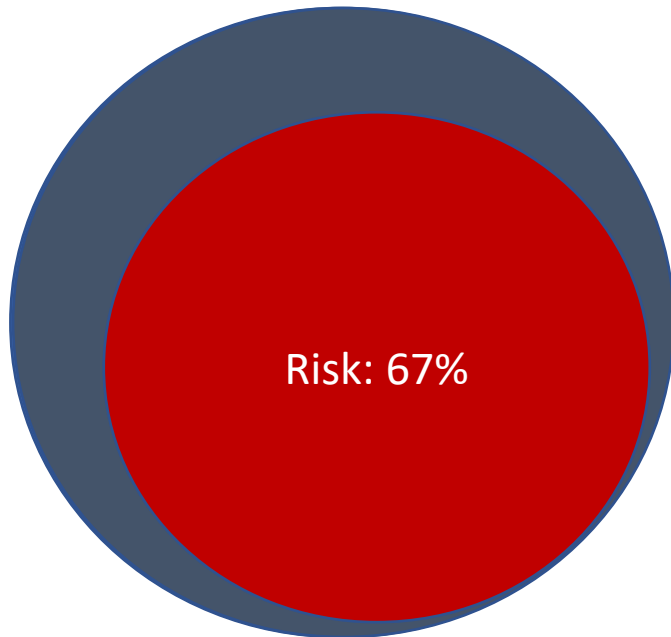
Coverage (usability): % of programs executed

Risk (safety): % of executed programs that were incorrect

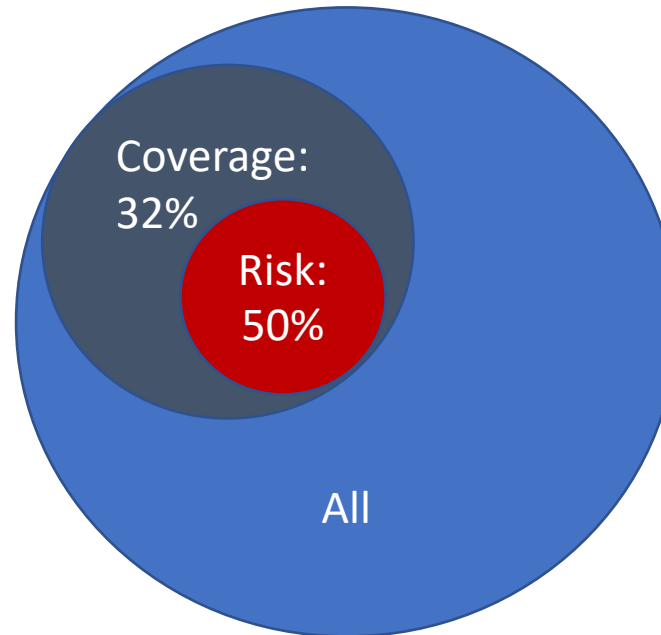


# DidYouMean Results

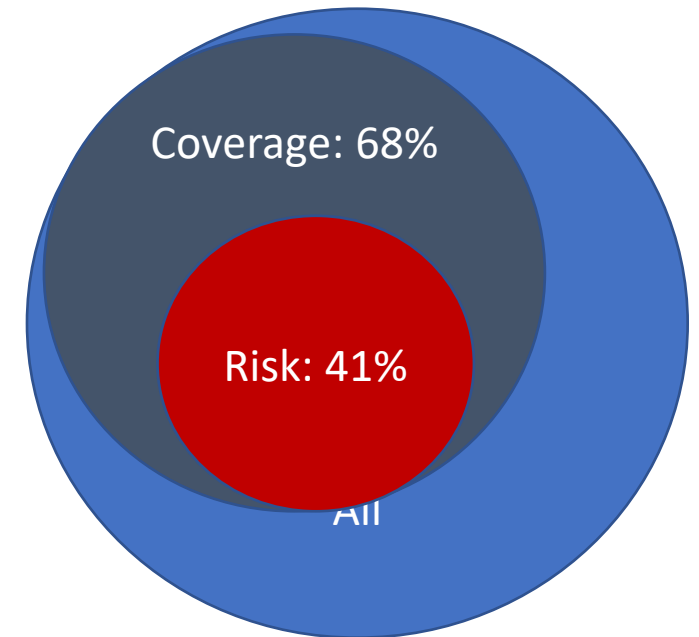
**Baseline 1:  
Accept everything**



**Baseline 2:  
Threshold**



**DidYouMean**



# Takeaways

**Calibration enables thresholding**

**Thresholding can improve safety**

**Human interaction improves balance**

When uncertainty is high, defer to human judgment

**Can we use the model to reduce uncertainty?**



# Outline

## Part I: Uncertainty in Human-Model Interactions

**Calibrated Interpretation: Confidence Estimation in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, TACL (2023)

**Did You Mean...? Confidence-based Trade-offs in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, EMNLP (2023)

## Part II: Model-based Selection to Reduce Uncertainty

**Rephrase, Augment, Reason: Visual Grounding of Questions for Vision-Language Models**, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal, ICLR (2024)

## Part III: Confidence for Model-Model Interactions

**ReConcile: Round-Table Conference Improves Reasoning via Consensus among Diverse LLMs**, Justin Chih-Yao Chen, Swarnadeep Saha, Mohit Bansal (2024)

**MAGDi: Structured Distillation of Multi-Agent Interaction Graphs Improves Reasoning in Smaller Language Models**, Justin Chih-Yao Chen\*, Swarnadeep Saha\*, Elias Stengel-Eskin Mohit Bansal (2024)

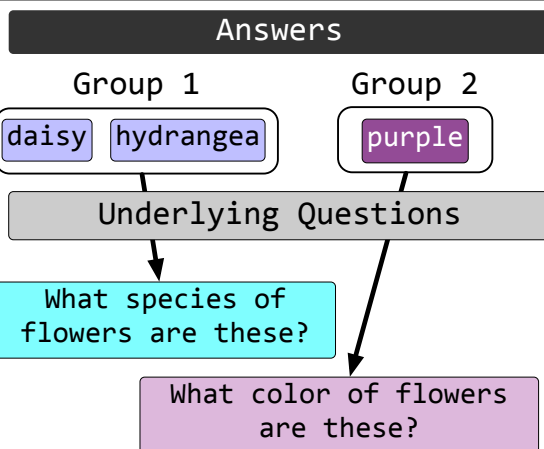


# Uncertainty and Underspecification

## VQA questions

May not provide sufficient information to answer correctly

VQA Q: What kind of flowers are these?



Question: Where are we at?



Why did the chicken cross the road? Rephrasing and analyzing ambiguous questions in VQA

Stengel-Eskin et al., 2023

Why does a visual question have different answers?

Bhattacharya et al., 2019

Dealing with semantic underspecification in multimodal NLP

Pezzelle, 2023

# RepARe: Rephrase, Augment, Reason

**Rephrase** questions to make them easier to answer

**Augment** questions with visual information

**Reason** about example and visual world

# RepARe: Rephrase, Augment, Reason

**Rephrase** questions to make them easier to answer

**Augment** questions with visual information

**Reason** about example and visual world

**Focus: zero-shot VQA with a large vision-language model**

# RepARe: Rephrase, Augment, Reason

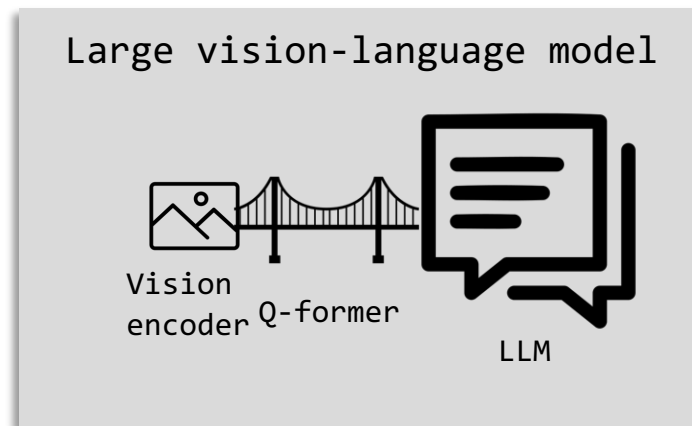
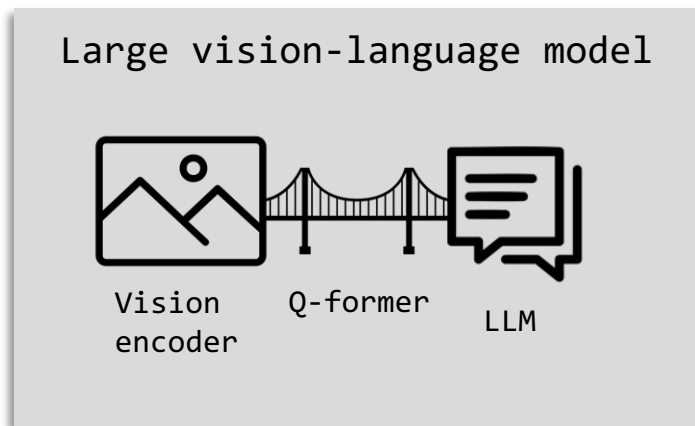
**Rephrase** questions to make them easier to answer

**Augment** questions with visual information

**Reason** about example and visual world

**Focus: zero-shot VQA with a large vision-language model**

**Key insight: Asymmetric strength and asymmetric ability**



QA ability: ★ ☆ ☆ ☆ ☆

Captioning ability: ★ ★ ★ ☆ ☆

# RepARe phases

Q. What period of the day does this photo reflect?



*I*: Image

**I. Extracting Visual Details**

# RepARe phases

Q. What period of the day does this photo reflect?



*I*: Image

Caption  
Generation

**I. Extracting Visual Details**

# RepARe phases

Q. What period of the day does this photo reflect?



*I*: Image

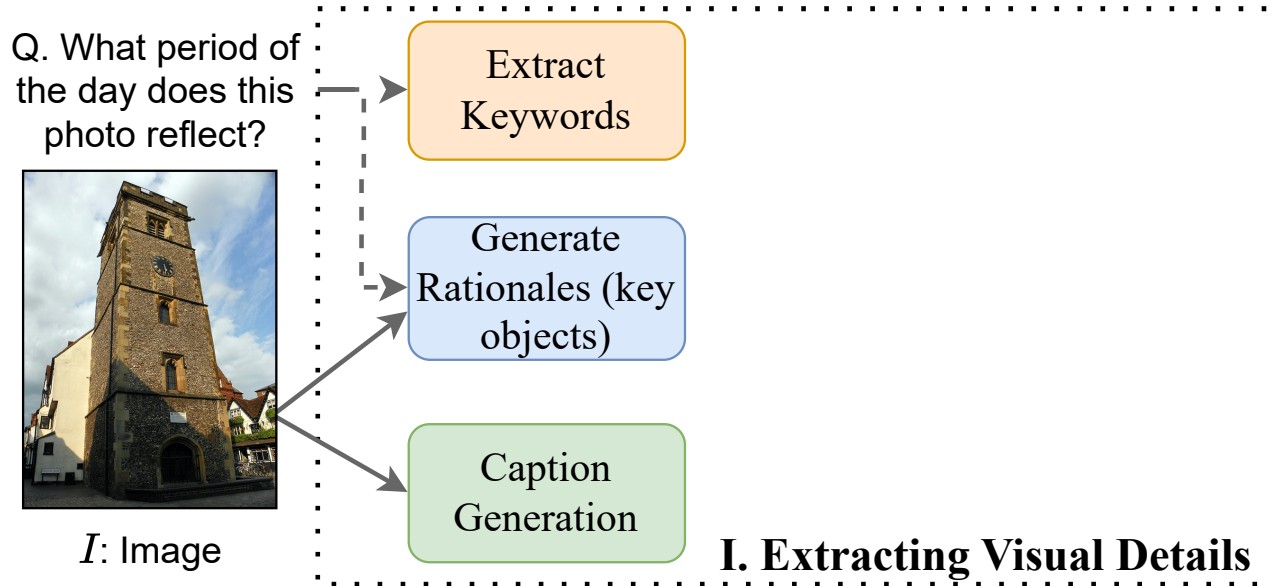
Generate Rationales (key objects)

Caption Generation

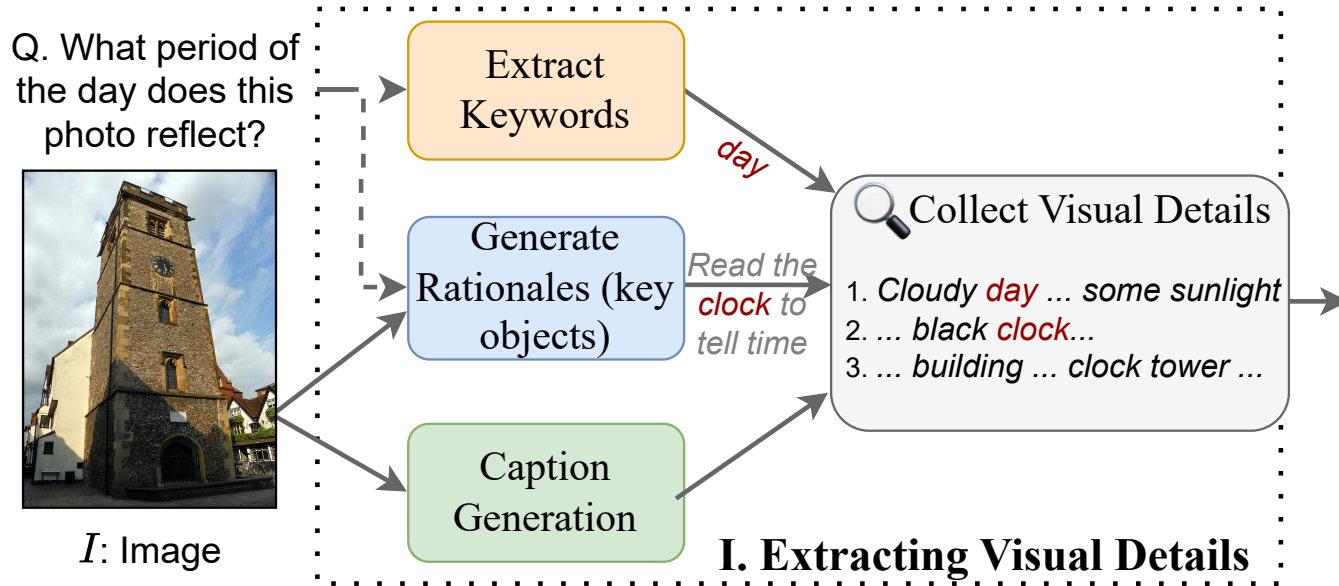
**I. Extracting Visual Details**



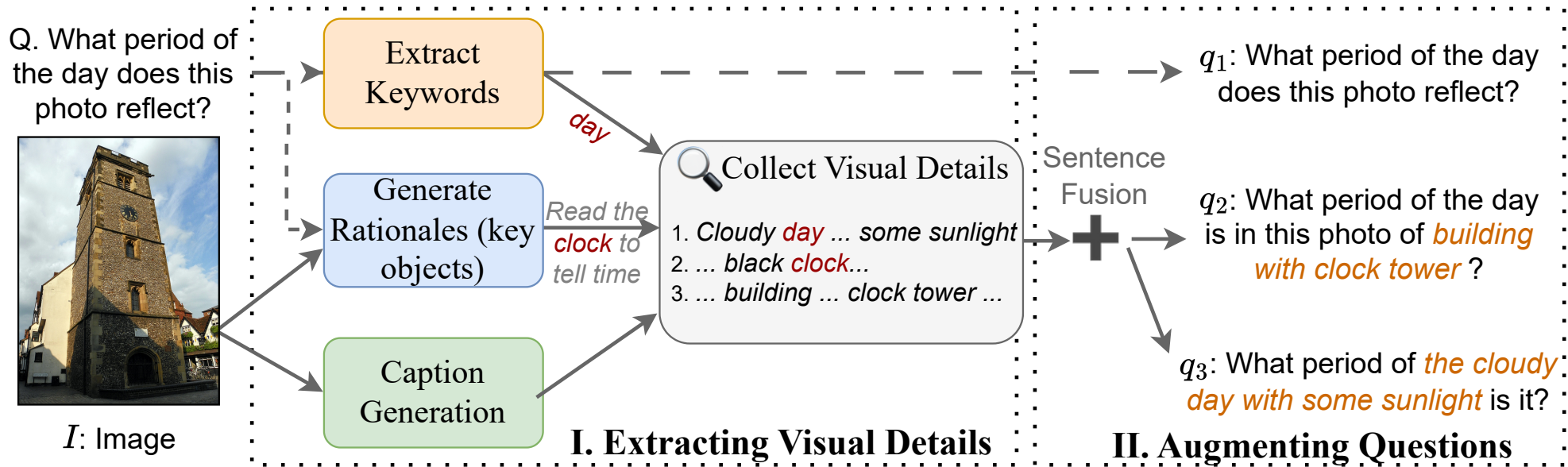
# RepARe phases



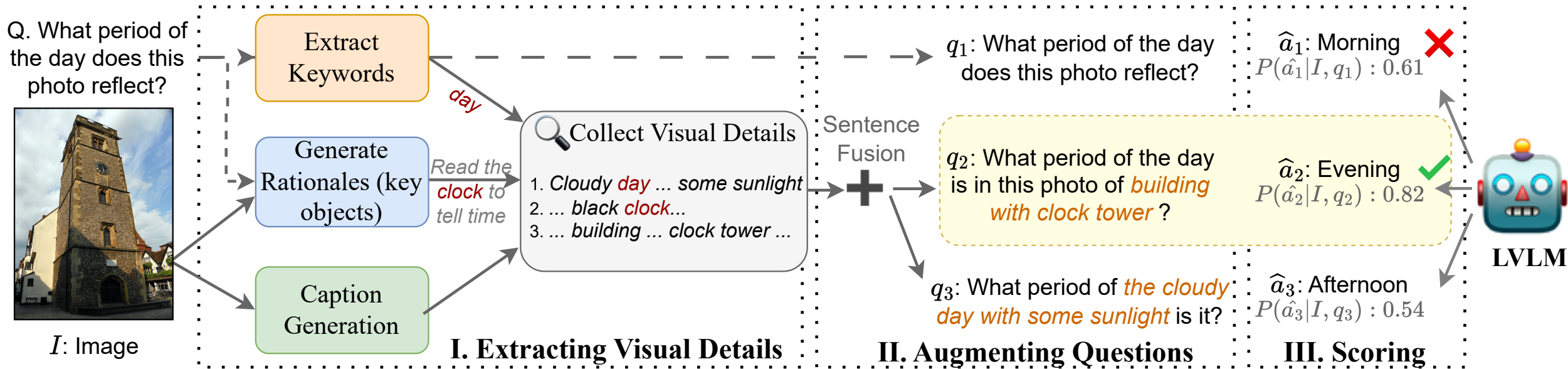
# RepARe phases



# RepARe phases



# RepARe phases

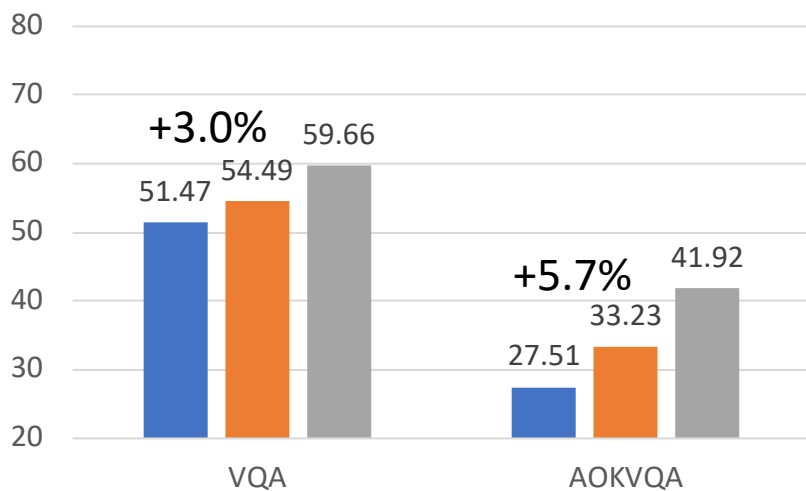


**Chosen via ARGMAX over confidence**  
Choose the question that increases answer confidence the most

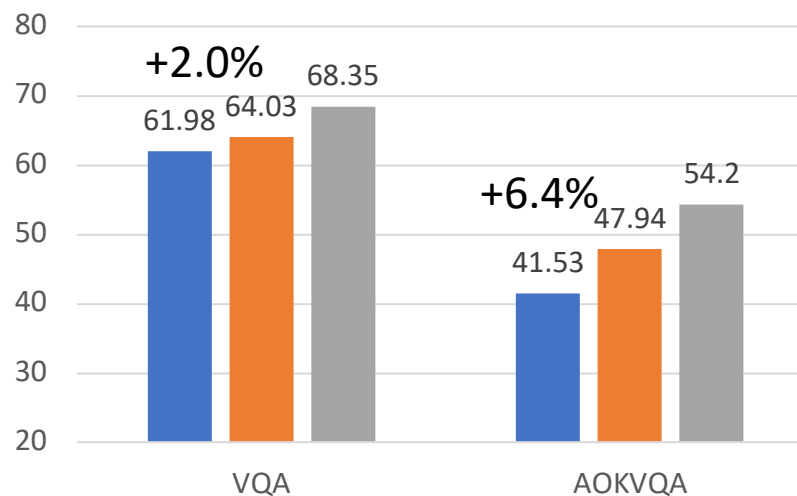
**i.e. question that reduces uncertainty**

# Results

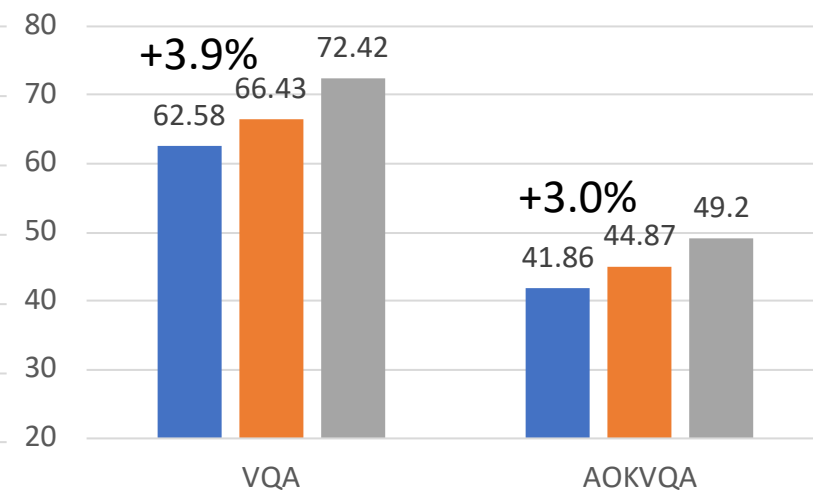
### MiniGPT-4 (Vicuna 7B)



### MiniGPT-4 (Vicuna 13B)



### BLIP-2 (Flan T5 XL)



■ Zero-shot ■ RepARe ■ Oracle selection

# Qualitative Examples

## Localization



*Q: Does the water have ripples?*

*RepARe Q: Does the water have **the small ripples around the boats?***

# Qualitative Examples

RepAR can help with reasoning



Q: *What will be built here one day?*

*RepARe Q: What will be built **at this construction site?***

# Outline

## Part I: Uncertainty in Human-Model Interactions

**Calibrated Interpretation: Confidence Estimation in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, *TACL* (2023)

**Did You Mean...? Confidence-based Trade-offs in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, *EMNLP* (2023)

## Part II: Model-based Selection to Reduce Uncertainty

**Rephrase, Augment, Reason: Visual Grounding of Questions for Vision-Language Models**, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal, *ICLR* (2024)

## Part III: Confidence for Model-Model Interactions

**ReConcile: Round-Table Conference Improves Reasoning via Consensus among Diverse LLMs**, Justin Chih-Yao Chen, Swarnadeep Saha, Mohit Bansal (2024)

**MAGDi: Structured Distillation of Multi-Agent Interaction Graphs Improves Reasoning in Smaller Language Models**, Justin Chih-Yao Chen\*, Swarnadeep Saha\*, Elias Stengel-Eskin Mohit Bansal (2024)





# Overview

**LLMs struggle with complex reasoning**

**ReConcile: A multi-model and multi-agent framework**

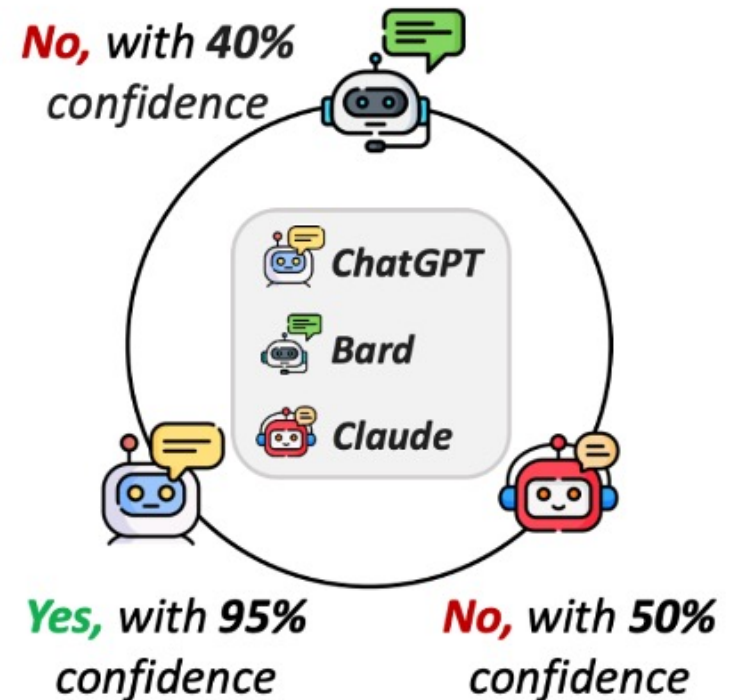
**Key components:**

Multi-LLM discussion

Multiple discussion rounds

Learning to convince other agents

Confidence-weighted voting



# Three stages of discussion

1

Initial response  
generation

Obtain initial response  
from each model  
(ChatGPT, Claude and Bard)

# Three stages of discussion

1

Initial response  
generation

Obtain initial response  
from each model  
(ChatGPT, Claude and Bard)

2

multi-round  
discussion

Aggregate the  
information by grouping  
+ **confidence estimation**

**Example of grouping:** There are **2**  
agents think the answer is **no**, and  
**1** agent thinks the answer is **yes**.

# Three stages of discussion

1

Initial response  
generation

Obtain initial response  
from each model  
(ChatGPT, Claude and Bard)

2

multi-round  
discussion

Aggregate the  
information by grouping  
+ **confidence estimation**

3

final answer  
generation

**Confidence-weighted**  
vote to derive the final  
answer

**Example of grouping:** There are 2  
agents think the answer is **no**, and  
1 agent thinks the answer is **yes**.

# Phase 1: Initial response

Question: Is August a winter month for part of the world?

Phase1: Initial Response Generation

Round 0

Initial Prompt



**Yes**, parts of the world in Southern Hemisphere...  
I am **60%** confident ...

Initial Prompt



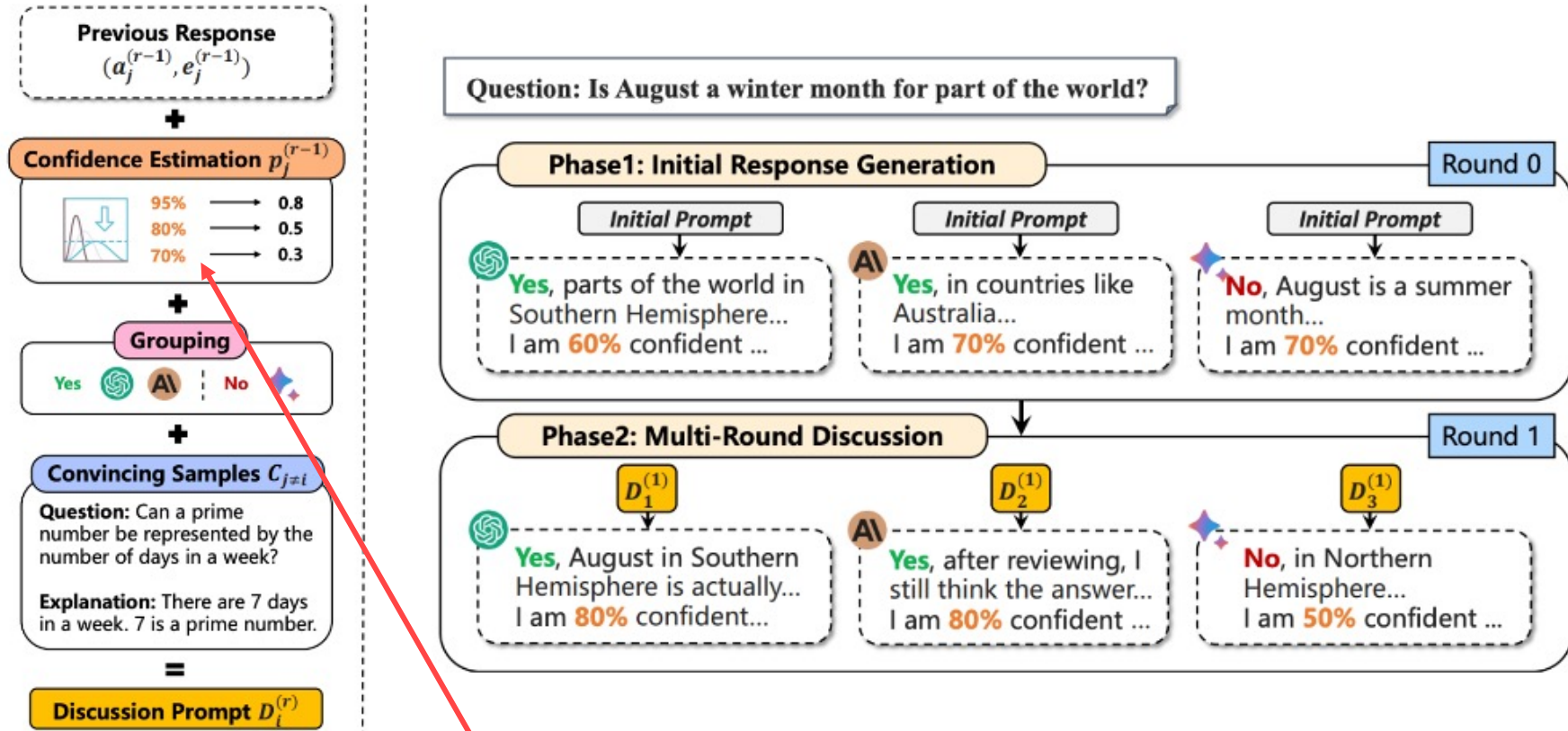
**Yes**, in countries like Australia...  
I am **70%** confident ...

Initial Prompt



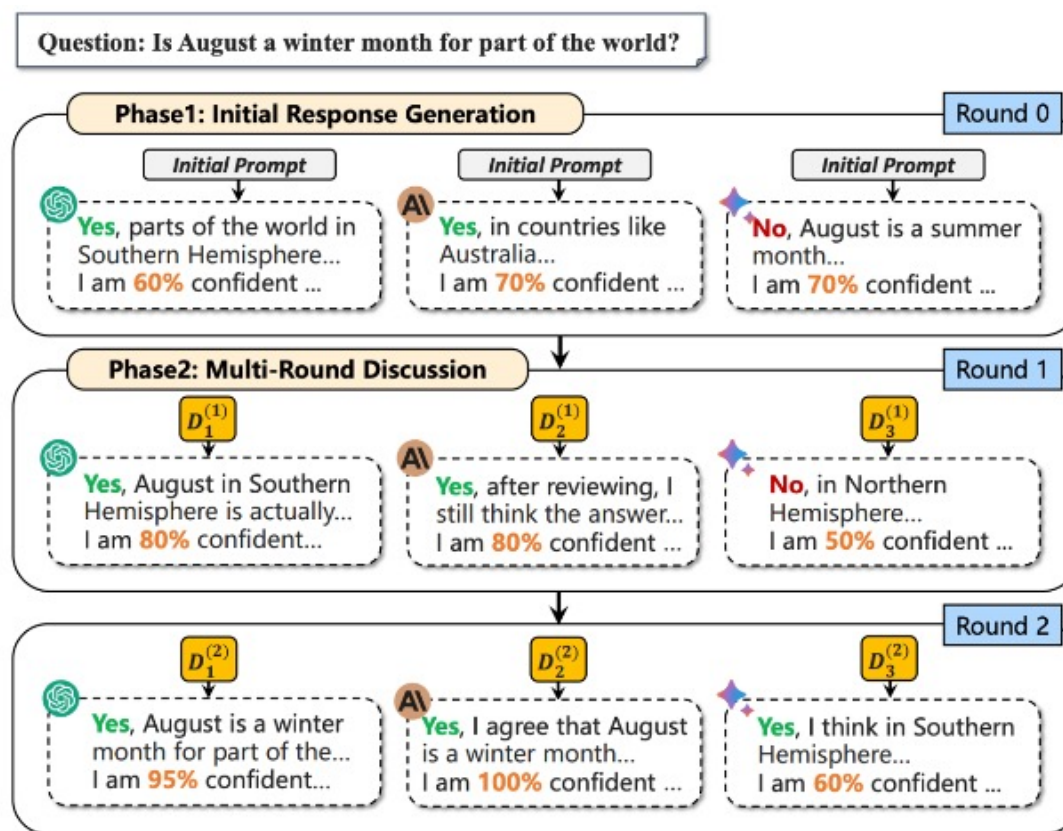
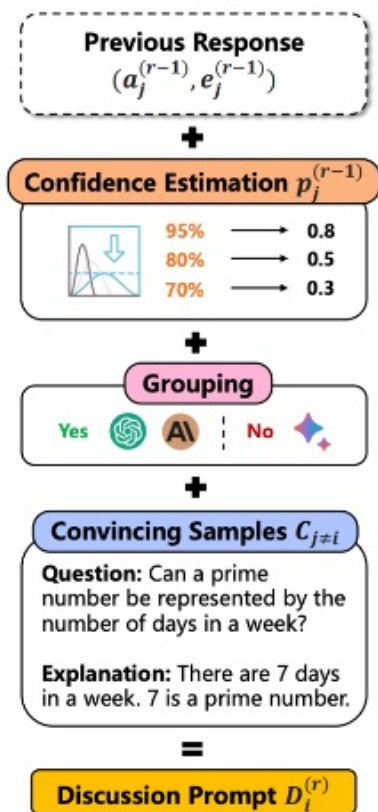
**No**, August is a summer month...  
I am **70%** confident ...

# Phase 2: Multi-round discussion

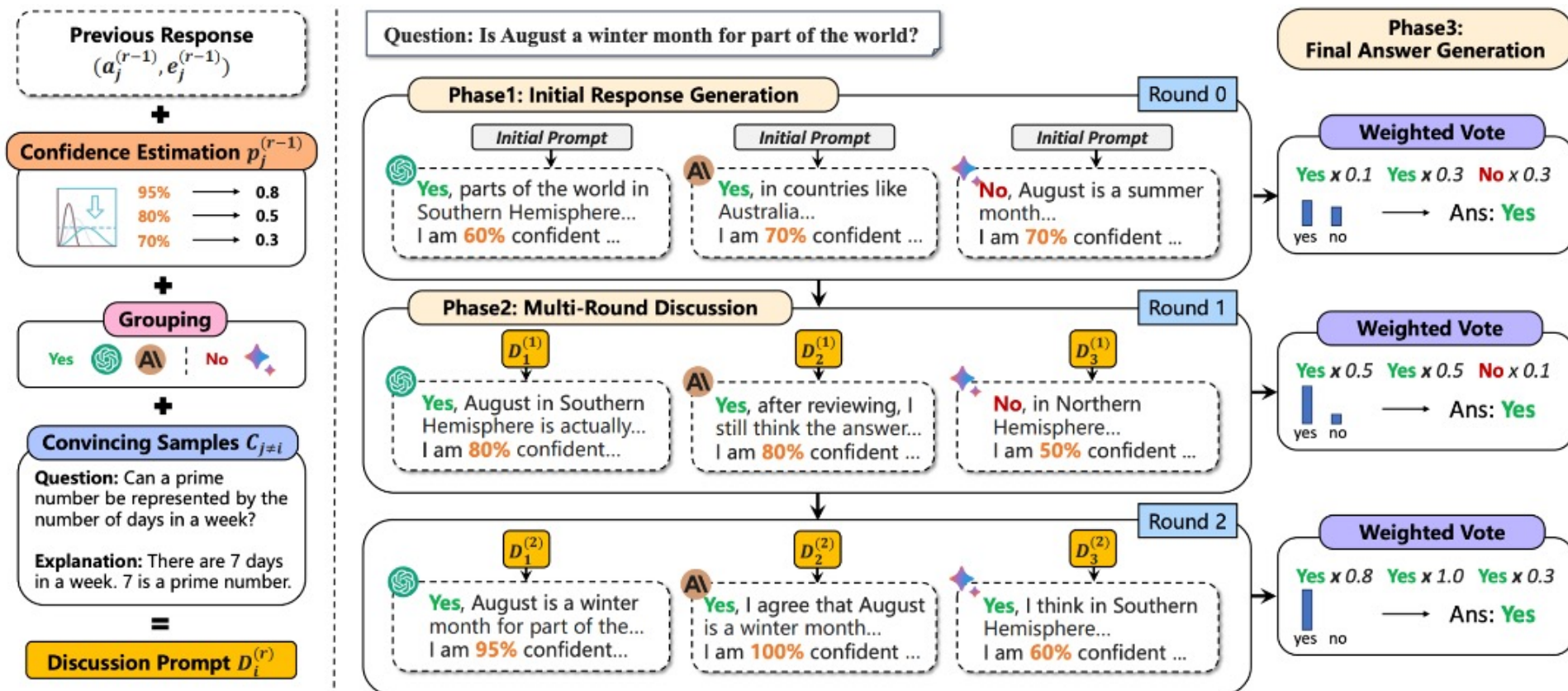


Corrected confidence scores

# Phase 2: Multi-round discussion



# Phase 3: Answer generation



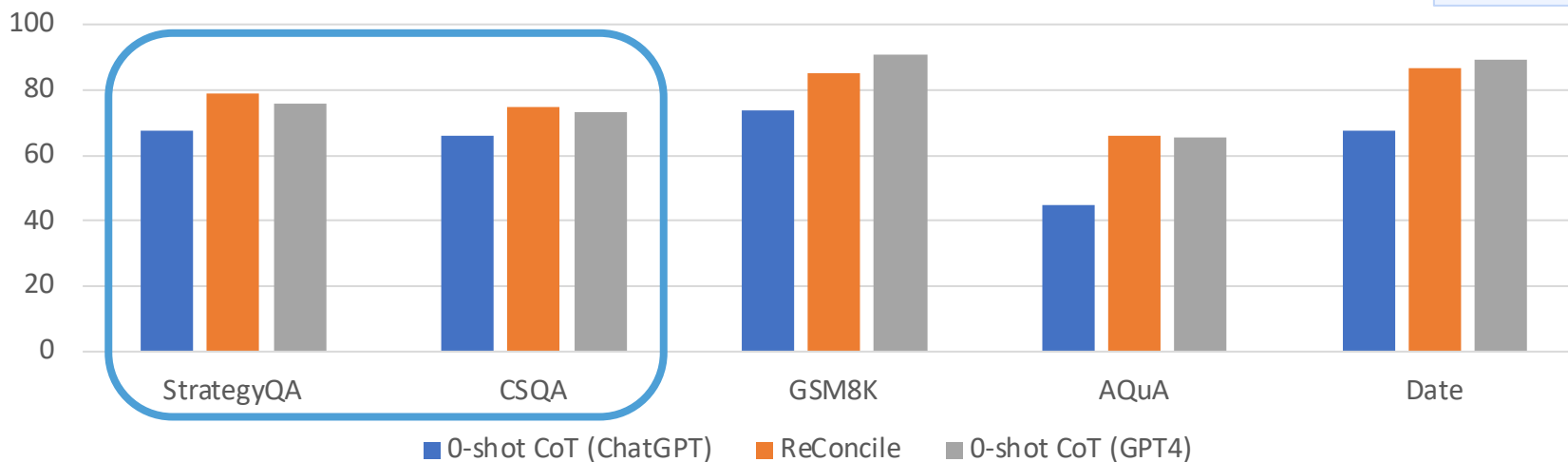


# Results

Most powerful / expensive model considered

Method Category	Method	Agent	StrategyQA	CSQA	GSM8K	AQuA	Date
Vanilla Single-agent	Zero-shot CoT	GPT-4	75.6 $\pm$ 4.7	73.3 $\pm$ 0.4	90.7 $\pm$ 1.7	65.7 $\pm$ 4.6	89.0 $\pm$ 2.2
	Zero-shot CoT	ChatGPT	67.3 $\pm$ 3.6	66.0 $\pm$ 1.8	73.7 $\pm$ 3.1	44.7 $\pm$ 0.5	67.7 $\pm$ 1.2
	Zero-shot CoT	Bard	69.3 $\pm$ 4.4	56.8 $\pm$ 2.7	58.7 $\pm$ 2.6	33.7 $\pm$ 1.2	50.2 $\pm$ 2.2
	Zero-shot CoT	Claude2	73.7 $\pm$ 3.1	66.7 $\pm$ 2.1	79.3 $\pm$ 3.6	60.3 $\pm$ 1.2	78.7 $\pm$ 2.1
Advanced Single-agent	Self-Refine (SR)	ChatGPT	66.7 $\pm$ 2.7	68.1 $\pm$ 1.8	74.3 $\pm$ 2.5	45.3 $\pm$ 2.2	66.3 $\pm$ 2.1
	Self-Consistency (SC)	ChatGPT	73.3 $\pm$ 2.1	70.9 $\pm$ 1.3	80.7 $\pm$ 1.2	54.0 $\pm$ 2.9	69.0 $\pm$ 0.8
	SR + SC	ChatGPT	72.2 $\pm$ 1.9	71.9 $\pm$ 2.1	81.3 $\pm$ 1.7	58.3 $\pm$ 3.7	68.7 $\pm$ 1.2
Single-model Multi-agent	Debate	$\times$ 3	66.7 $\pm$ 3.1	62.7 $\pm$ 1.2	83.0 $\pm$ 2.2	65.3 $\pm$ 3.1	68.0 $\pm$ 1.6
	Debate	$\times$ 3	65.3 $\pm$ 2.5	66.3 $\pm$ 2.1	56.3 $\pm$ 1.2	29.3 $\pm$ 4.2	46.0 $\pm$ 2.2
	Debate	$\times$ 3	71.3 $\pm$ 2.2	68.3 $\pm$ 1.7	70.7 $\pm$ 4.8	62.7 $\pm$ 2.6	75.3 $\pm$ 3.3
	Debate+Judge	$\times$ 3	69.7 $\pm$ 2.1	63.7 $\pm$ 2.5	74.3 $\pm$ 2.9	57.3 $\pm$ 2.1	67.7 $\pm$ 0.5
Multi-model Multi-agent	RECONCILE	, ,	79.0 $\pm$ 1.6	74.7 $\pm$ 0.4	85.3 $\pm$ 2.2	66.0 $\pm$ 0.8	86.7 $\pm$ 1.2

ReConcile w/o GPT4 outperforms it!



# Outline

## Part I: Uncertainty in Human-Model Interactions

**Calibrated Interpretation: Confidence Estimation in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, *TACL* (2023)

**Did You Mean...? Confidence-based Trade-offs in Semantic Parsing**, Elias Stengel-Eskin and Benjamin Van Durme, *EMNLP* (2023)

## Part II: Model-based Selection to Reduce Uncertainty

**Rephrase, Augment, Reason: Visual Grounding of Questions for Vision-Language Models**, Archiki Prasad, Elias Stengel-Eskin, Mohit Bansal, *ICLR* (2024)

## Part III: Confidence for Model-Model Interactions

**ReConcile: Round-Table Conference Improves Reasoning via Consensus among Diverse LLMs**, Justin Chih-Yao Chen, Swarnadeep Saha, Mohit Bansal (2024)

**MAGDi: Structured Distillation of Multi-Agent Interaction Graphs Improves Reasoning in Smaller Language Models**, Justin Chih-Yao Chen\*, Swarnadeep Saha\*, Elias Stengel-Eskin Mohit Bansal (2024)



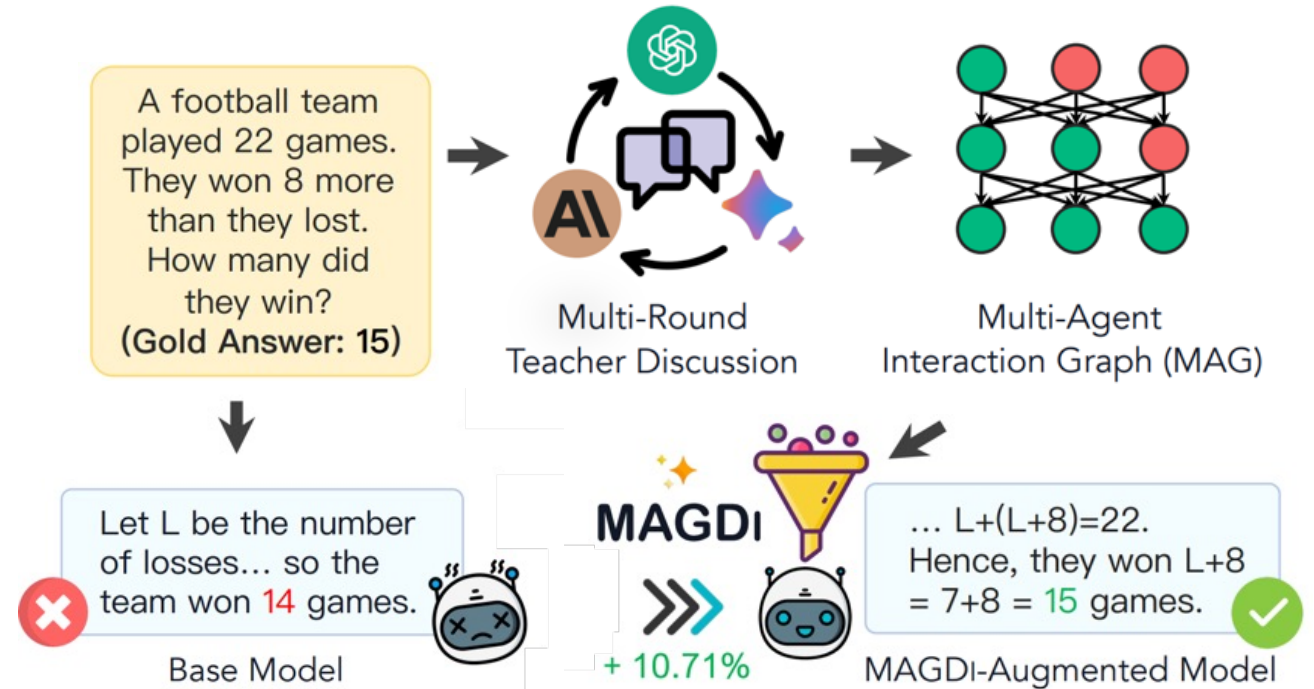
# Can we distill ReConcile into a single model?

**Strong performance boost but...**

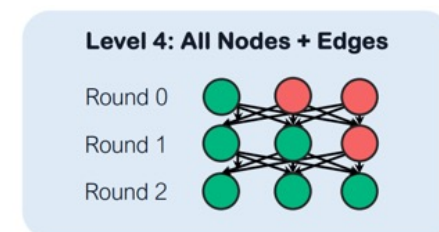
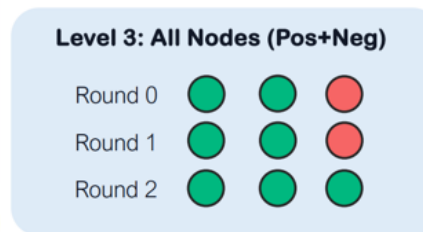
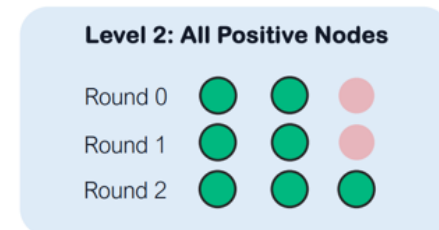
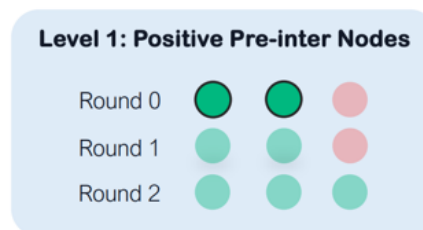
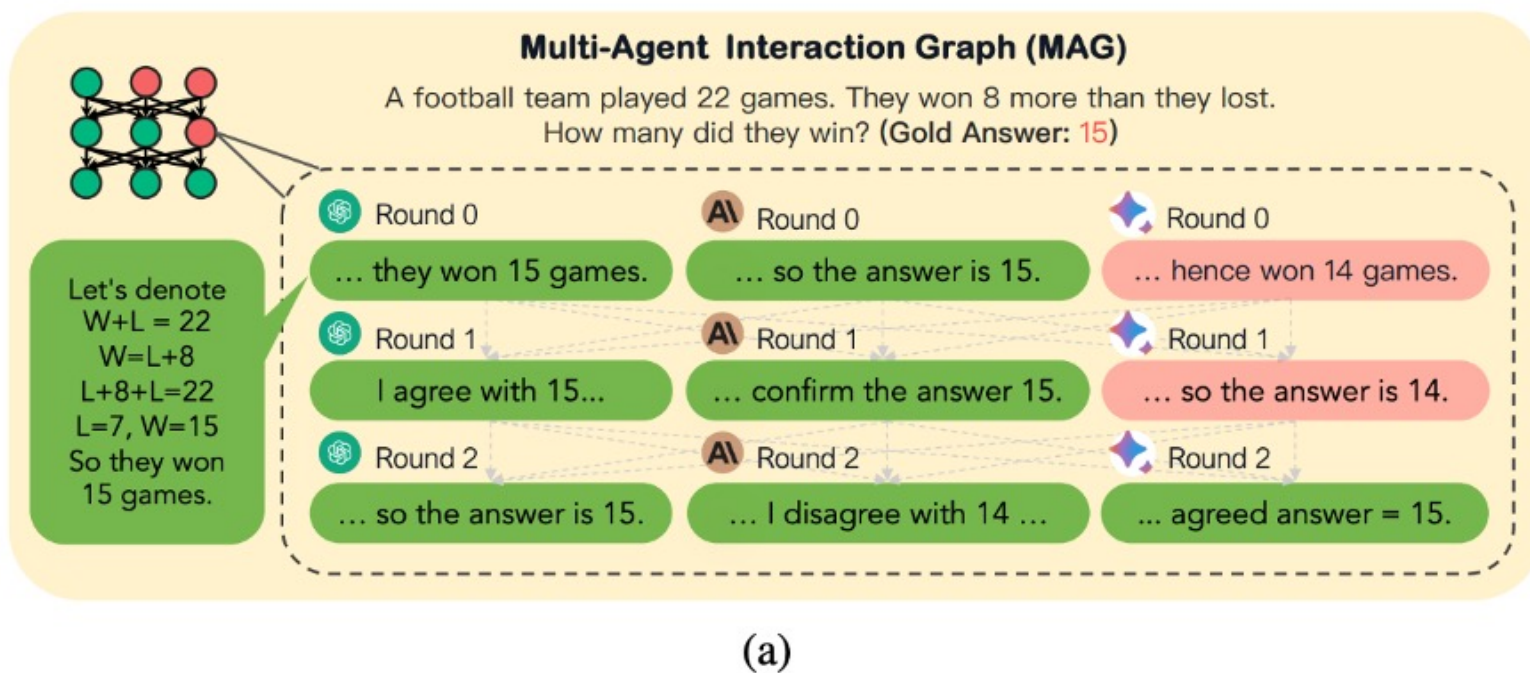
Multiple LLMs across multiple rounds

Expensive!

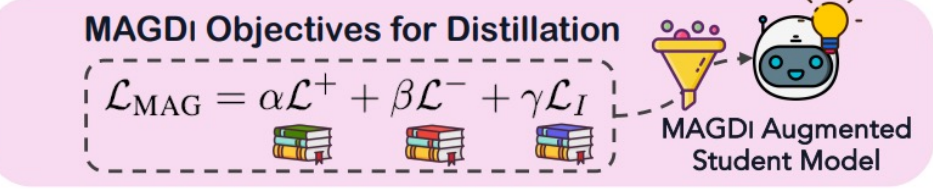
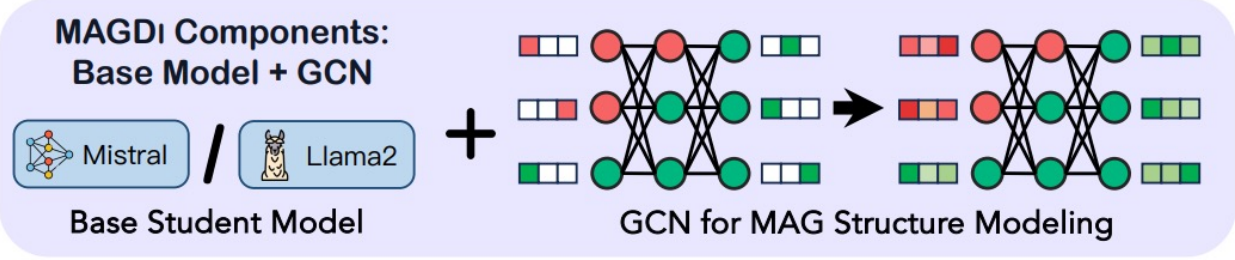
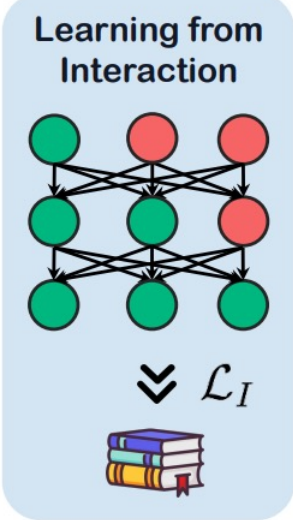
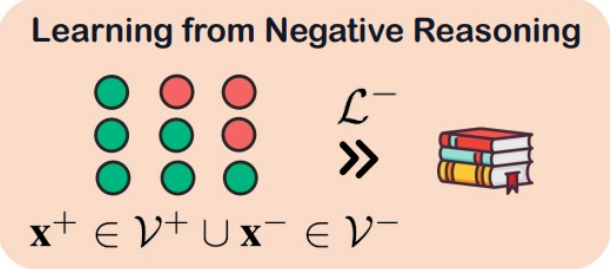
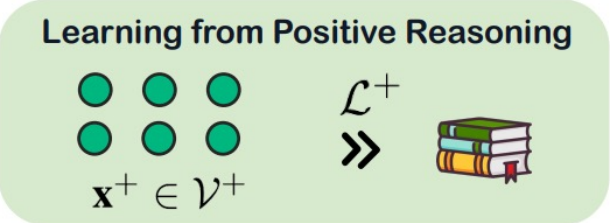
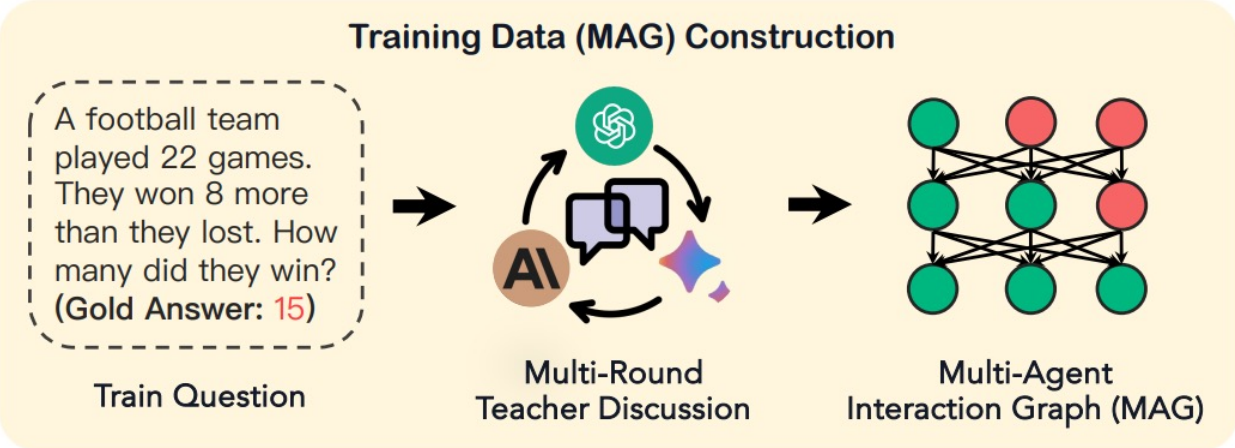
**Open-source LLMs**



# Levels of distillation

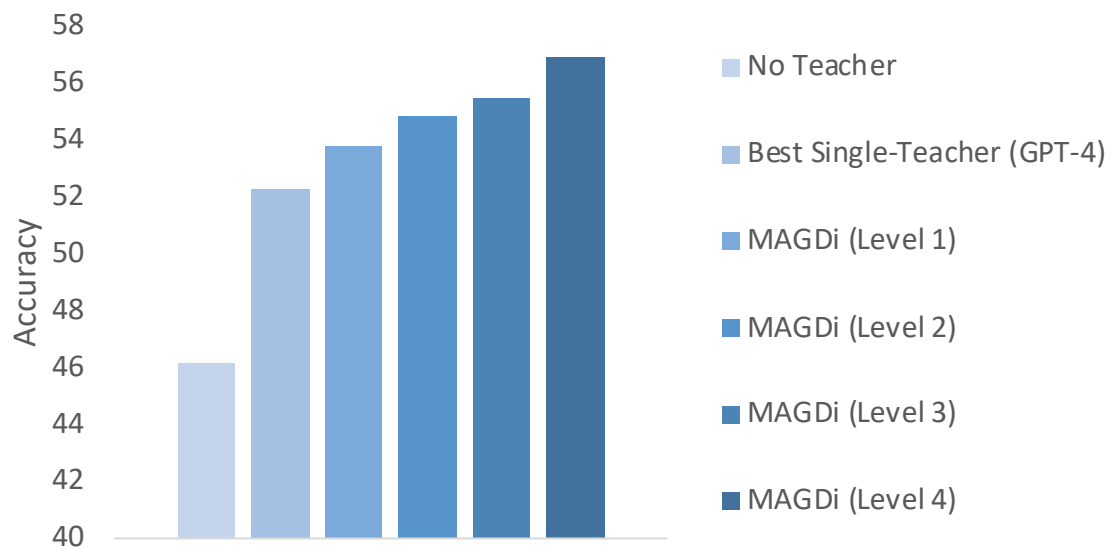


# Details

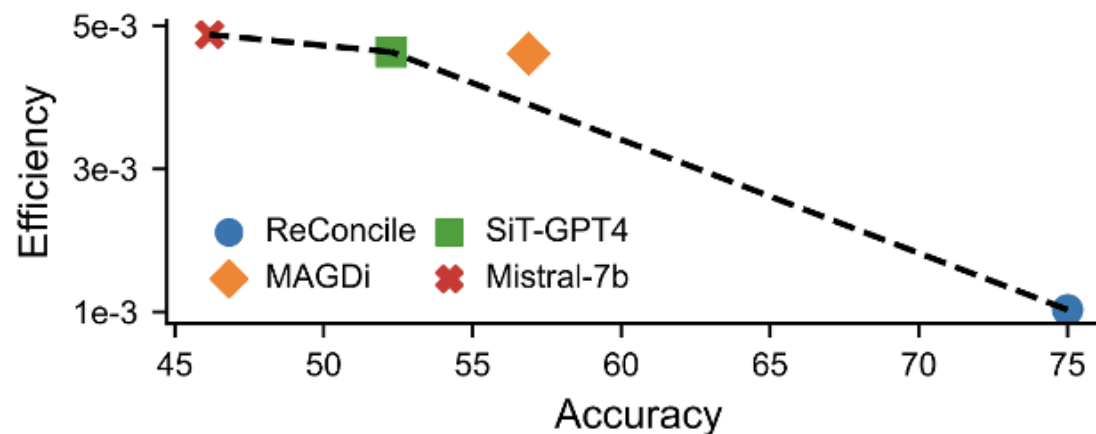


# Results

Mistral-7B average across StrategyQA, CSQA, ARC, GSM8K, MATH



Best tradeoff between efficiency and performance



# Conclusions

## **Part I: Confidence and human-model interaction**

Calibration in semantic parsing

What can we do when interacting with calibrated models?

# Conclusions

## **Part I: Confidence and human-model interaction**

Calibration in semantic parsing

What can we do when interacting with calibrated models?

## **Part II: Confidence for selection**

Improving zero-shot VQA by reducing uncertainty



# Conclusions

## **Part I: Confidence and human-model interaction**

Calibration in semantic parsing

What can we do when interacting with calibrated models?

## **Part II: Confidence for selection**

Improving zero-shot VQA by reducing uncertainty

## **Part III: Confidence for model-model interaction**

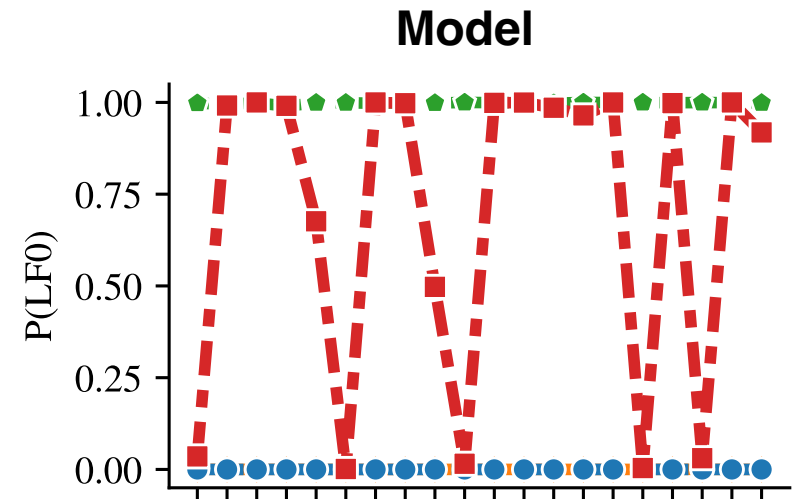
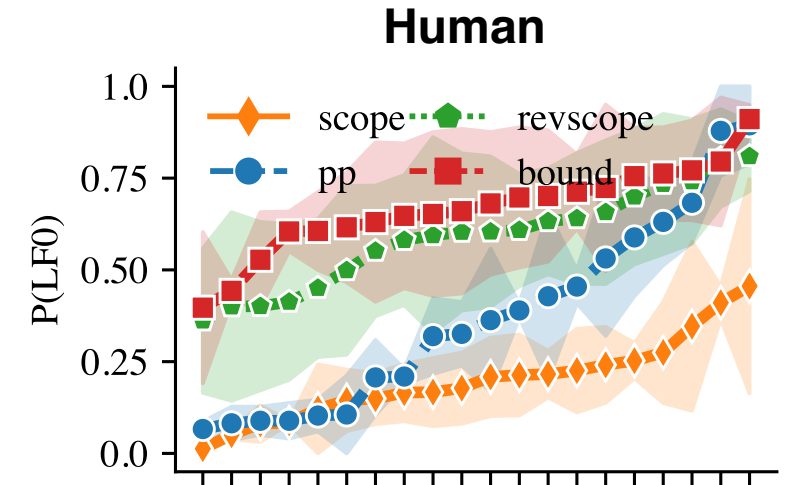
Improving reasoning via LLM collaboration

# Future directions

## Dealing with ambiguity

Semantic vs. form uncertainty

Models do not capture human uncertainty



# Future directions

## Dealing with ambiguity

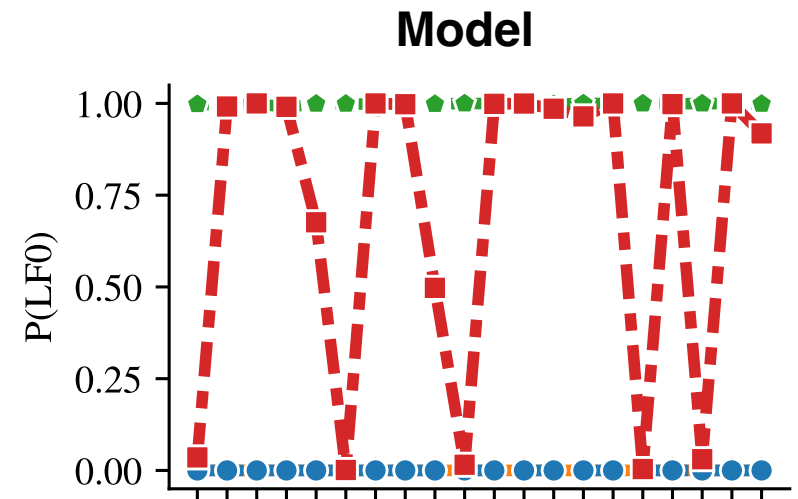
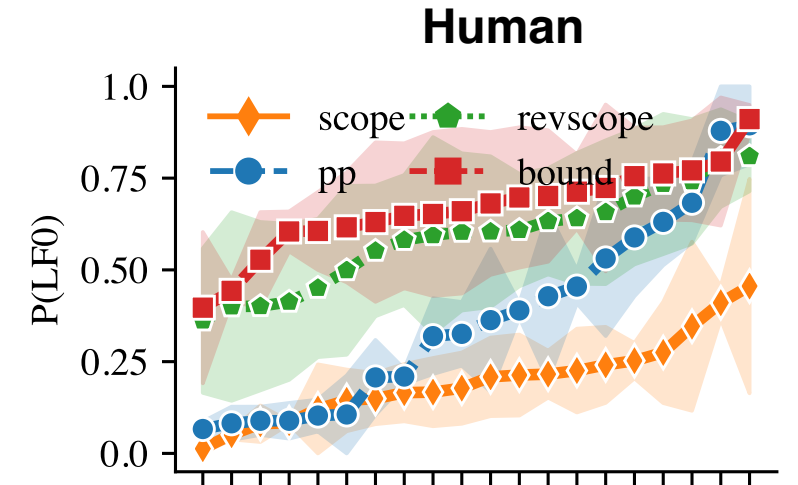
Semantic vs. form uncertainty

Models do not capture human uncertainty

## Confidence as a proxy for accuracy

Safety-usability tradeoffs

Also: improving accuracy via confidence



**Thank you!**