# Building Natural Language Interfaces in the Age of LLMs

Ziyu Yao

Assistant Professor, Computer Science
George Mason University
https://ziyuyao.org/

Department of Human-Centered Computing
Luddy School of Informatics, Computing, and Engineering
Indiana University, Indianapolis
April 12, 2024

# Natural Language Interfaces (NLIs)

# Natural Language Interfaces (NLIs) in History



```
Welcome to
            EEEEEE  LL      IIII  ZZZZZZZ  AAAAA
            EE      LL       II        ZZ  AA  AA
            EEEEE   LL       II       ZZZ  AAAAAAA
            EE      LL       II      ZZ    AA  AA
            EEEEEE  LLLLLL  IIII  ZZZZZZZ  AA  AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?
YOU:   Men are all alike.
ELIZA: What is the connection, do you suppose ?
YOU:   They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU:   Well, my boyfriend made me come here.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU:   He says I'm depressed much of the time.
ELIZA: I am sorry to hear that you are depressed.
YOU:   It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:
```
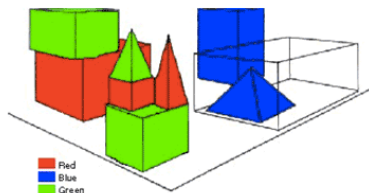
ELIZA (1966)

Person: Pick up a big red block.
Computer: OK.
Person: Grasp the pyramid.
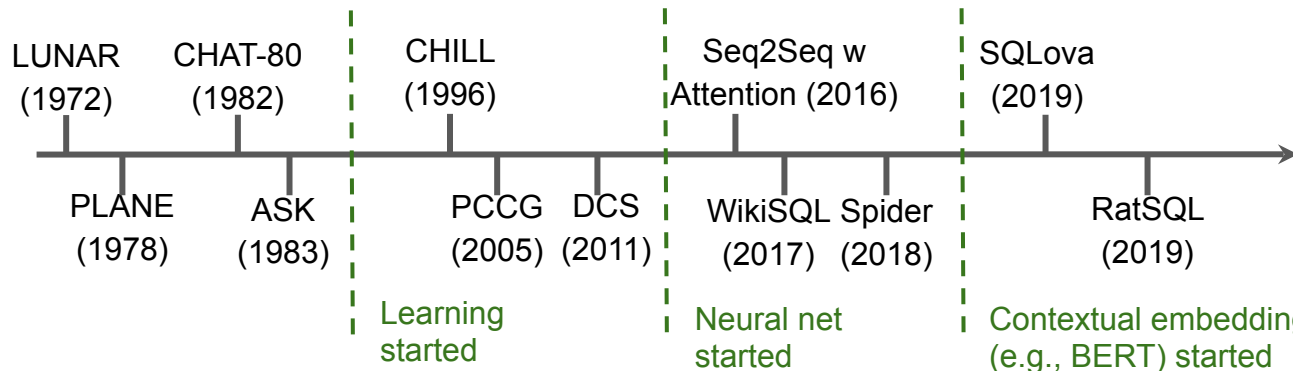Computer: I don't understand which pyramid you mean.
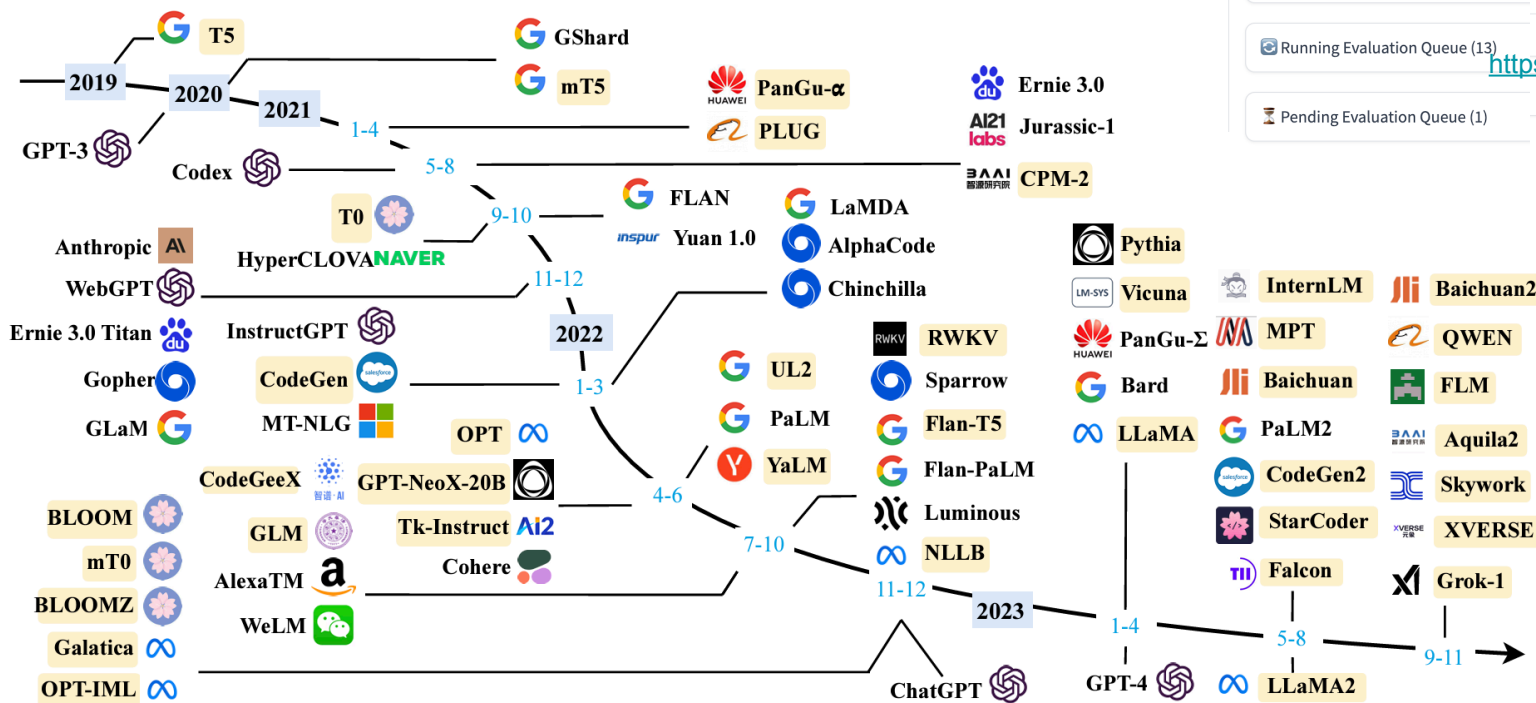
SHRDLU (1971)

Ask Jeeves (1997)

NLI to Database (NLIDB)

LUNAR (1972)  CHAT-80 (1982)  CHILL (1996)  Seq2Seq w Attention (2016)  SQLova (2019)

PLANE (1978)  ASK (1983)  PCCG (2005)  DCS (2011)  WikiSQL (2017)  Spider (2018)  RatSQL (2019)

Learning started

Neural net started

Contextual embedding (e.g., BERT) started

3

# Large Language Models (LLMs)



Image source: Zhao et al. "A survey of large language models." *arXiv preprint arXiv:2303.18223*.

# NLIs in the Age of Large Language Models

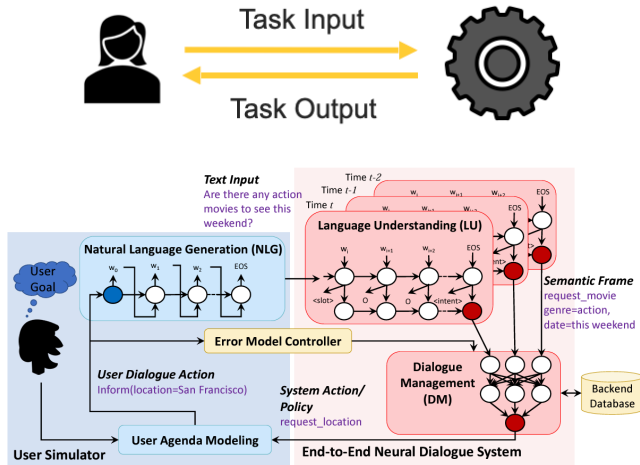Paradigm shift: unified architecture, task generalization, instruction following



UnifiedSKG (Xie…Yao et al., 2022)

+ Plugins

LLMs as unified, conversational NLIs
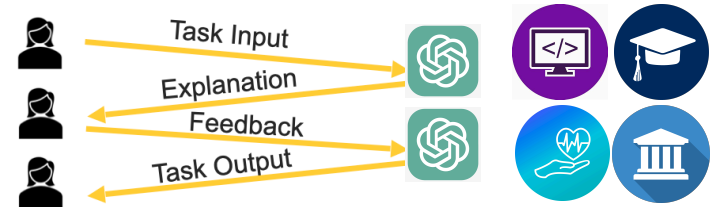
# NLIs in the Age of Large Language Models

Paradigm shift: unified architecture, task generalization, instruction following
↳ Paradigm shift in how humans interact with NLIs



(Li et al., 2017)
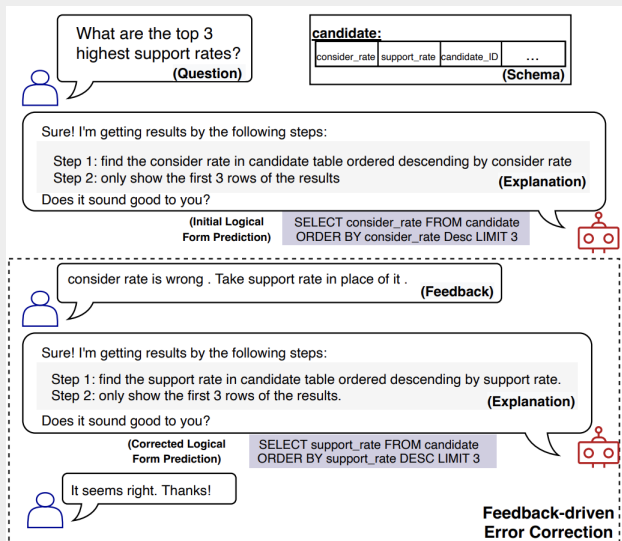
No interaction,  or task-specific interaction

Task-agnostic, multi-turn interactions
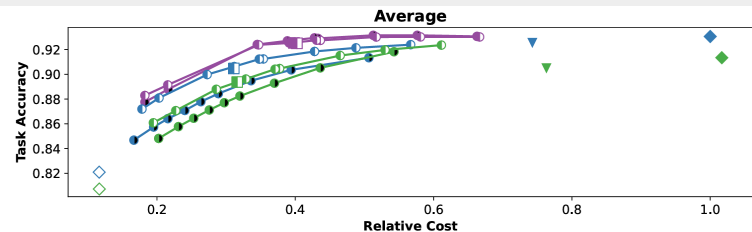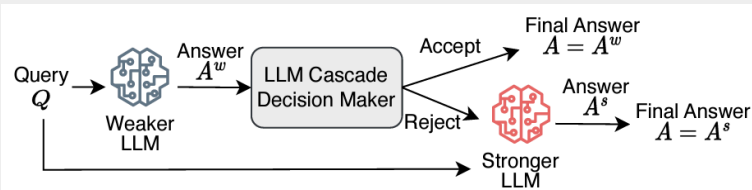& Broader application areas

*Do LLMs interact well with humans?*
*How to deal with the $ cost of*
*frequent queries to LLMs?*

# This Talk: Building NLIs in the Age of LLMs



Topic 1: Modeling Language Feedback in Human-NLI Interaction
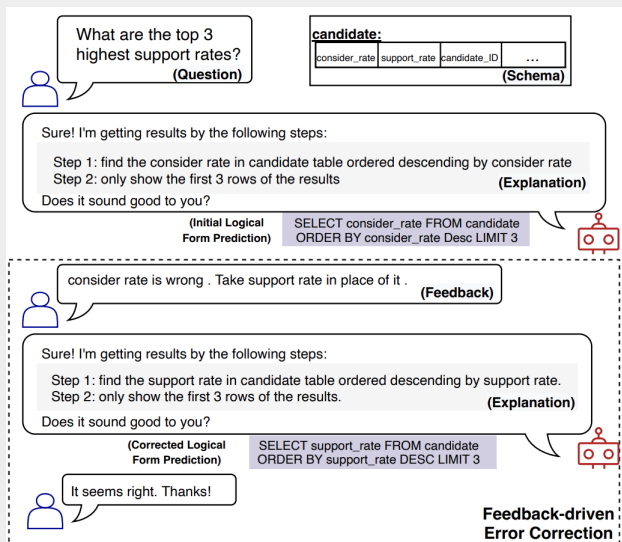(Task: Text-to-Code Generation)



Topic 2: Saving the Monetary Cost of LLM API Usage
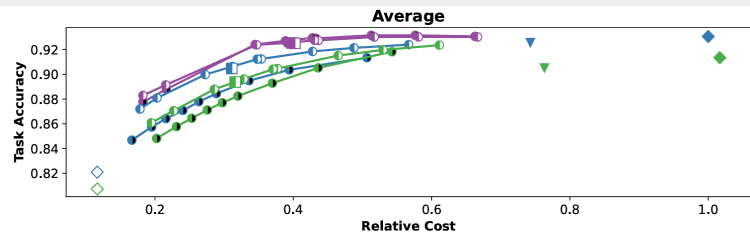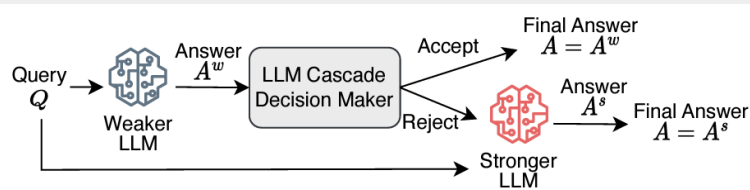(Task: Arithmetic/Symbolic/etc. Reasoning)

# This Talk: Building NLIs in the Age of LLMs



**Topic 1: Modeling Language Feedback in Human-NLI Interaction**
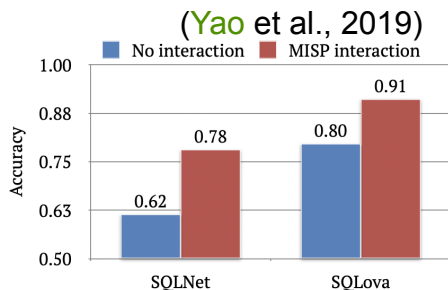(Task: Text-to-Code Generation)

**Topic 2: Saving the Monetary Cost of LLM API Usage**
(Task: Arithmetic/Symbolic/etc. Reasoning)

# Feedback-driven Human-NLI Interaction

- Humans naturally provide feedback while interacting with NLIs
  - e.g., "You should not do this; the result is not what I asked for!"

- Gap: existing NLIs are rarely evaluated with human interaction
  - Need more practical assessments, i.e., *when NLIs can interact with humans*
  - Feedback understanding and incorporation: not an easy task for LLMs!



(Yao et al., 2019)

Underestimate?

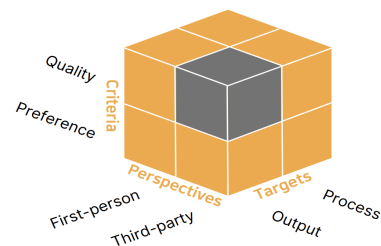**Evaluating Human-Language Model Interaction**

Mina Lee    Megha Srivastava    Amelia Hardy    John Thickstun    Esin Durmus
Ashwin Paranjape    Ines Gerard-Ursin[§]    Xiang Lisa Li    Faisal Ladhak
Frieda Rong    Rose E. Wang    Minae Kwon    Joon Sung Park    Hancheng Cao
Tony Lee    Rishi Bommasani    Michael Bernstein    Percy Liang

Stanford University    [§]Imperial College London

… or Overestimate?

# Semantic Parsing

- Translating a natural language (NL) question/command to its logical meaning representation
  - e.g., NL-to-SQL parsing for database querying
  - Other applications: robotics (NL-to-LTL), knowledge base query (NL-to-Lambda Calculus), AI-assisted programming (NL-to-Python/Java/C/…)

*In 2017, how many women were diagnosed with female breast cancer?*

non-technical user
(e.g., physician)

text-to-SQL
semantic parsing

*250, 520.*

```
SELECT COUNT(*)
FROM US_Cancer_Stats
WHERE Year=2017 AND
Sex="Female" AND
Type="Female Breast"
```

(Zhong et al., 2017; Yu et al., 2018)
https://gis.cdc.gov/Cancer/USCS/DataViz.html

# Interactive Semantic Parsing/Code Generation

- Semantic parsing with humans proving clarification and corrective feedback
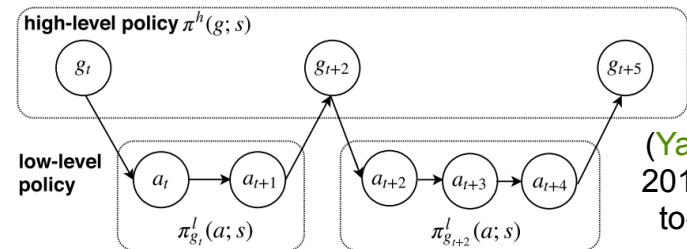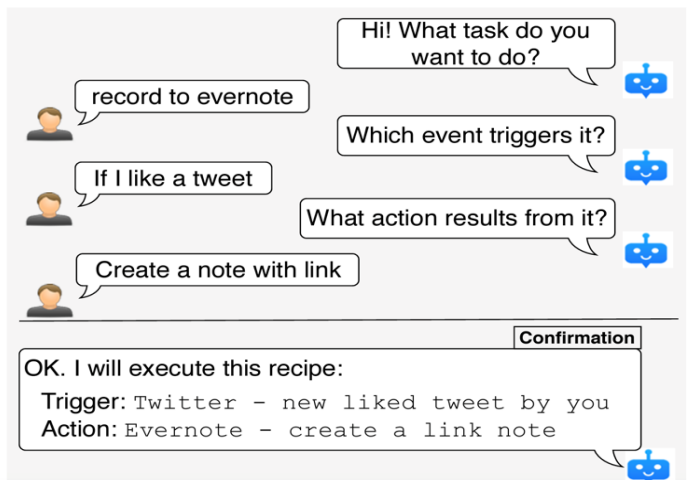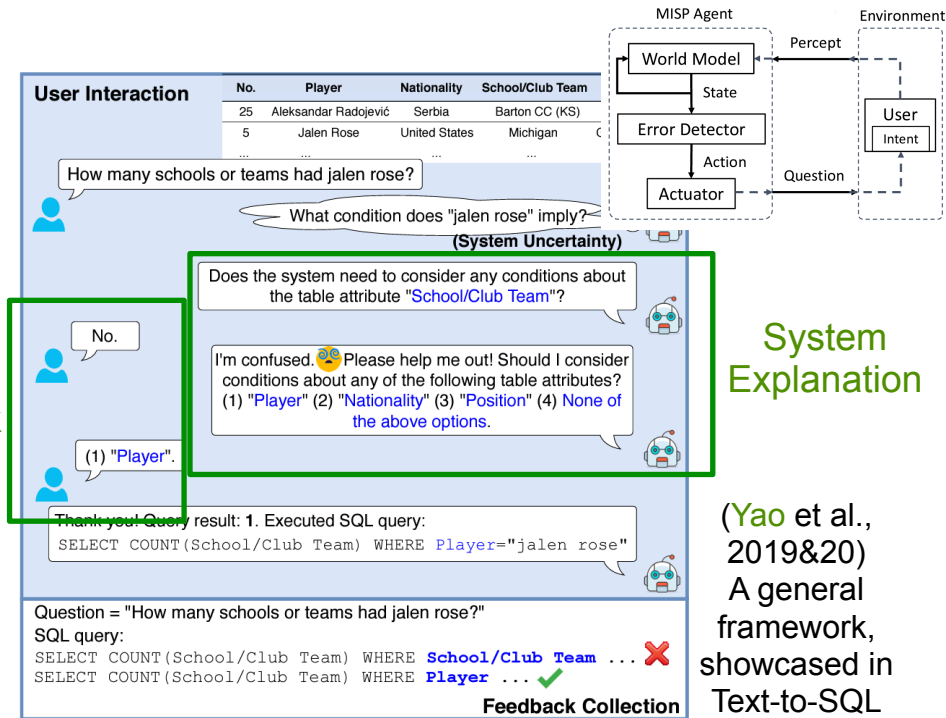


(Yao et al., 2019) Text-to-IFTTT

(Yao et al., 2019&20) A general framework, showcased in Text-to-SQL

# Interactive Semantic Parsing/Code Generation

- SPLASH dataset by Microsoft Research: text-to-SQL with natural language (NL) feedback

The **data bottleneck**: *Costly* and *model-dependent* feedback annotation

> Find all the locations whose names contain the word "film"

finding the Address of Locations table for which Location_Name contains "film"

| Address |
| --- |
| 770 Edd Lane Apt. 098 |
| 14034 Kohler Drive |

**System Explanation**

> Address is wrong. I want the name of the locations

**User's Corrective Feedback in NL**

finding the Location_Name of Locations table for which Location_Name contains "film"

| Location_Name |
| --- |
| Film Festival |
| Film Castle |
| … |

(Elgohary et al., 2020)

# Learning to Simulate Natural Language Feedback for Interactive Semantic Parsing

Hao Yan

Saurabh Srivastava

Yintao Tai

Sida I. Wang

Scott Yih

Ziyu Yao

ACL 2023

**Human-LLM Interaction via Language Feedback**

# Learning to Simulate Natural Language Feedback

- Idea:
    - *Build* a simulator with *small-scale* feedback annotations
    - *Apply* the simulator to generate *large-scale* synthetic feedback for model training

Initial Question
& Contextual Info →

Incorrect Code
& Explanation →

Correct Code
(accessible for
simulation only) →

**Feedback Simulator** → Natural Language Feedback

# Learning to Simulate Natural Language Feedback

- The importance of task representations ("prompt engineering"):

  - CWQES: Simply include the Correct and Wrong code snippets as input.

  - DQES: Inspired by NL-Edit (Elgohary et al. 2021), feed the EDits of revising the incorrect code snippet into the correct one.

  - TQES: Verbalize the edits using Templates.

**Simulated Feedback**

⇧

Feedback Simulator
(T5-based)

⇧

[question] **Q** [explanation] **E** [schema] **S**

**True & Predicted Logical Form**
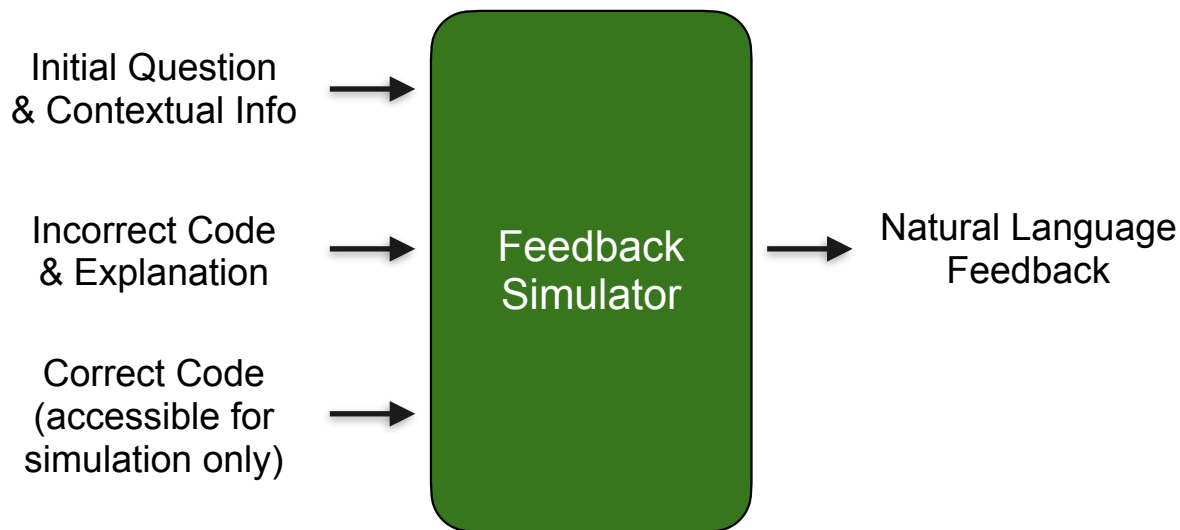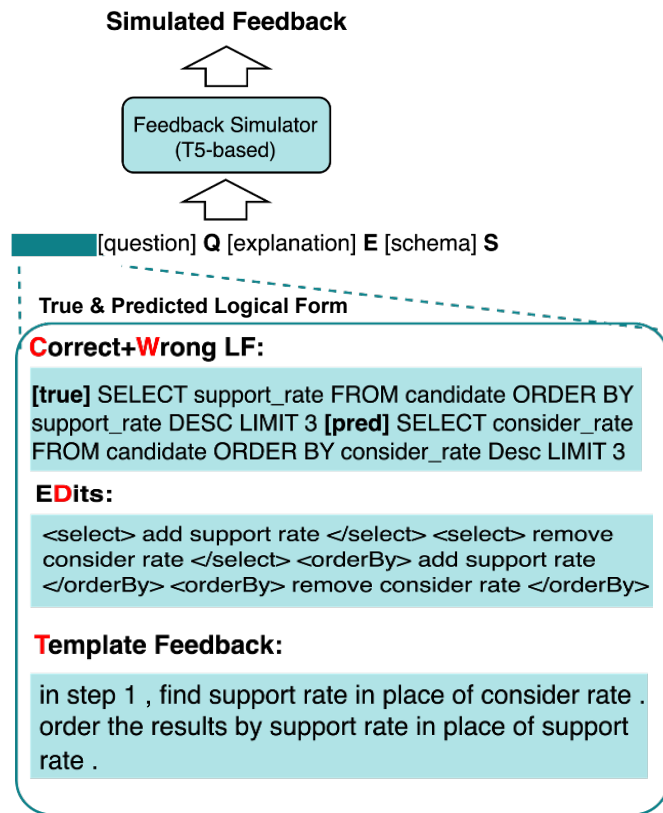
**C**orrect+**W**rong LF:

[true] SELECT support_rate FROM candidate ORDER BY support_rate DESC LIMIT 3 [pred] SELECT consider_rate FROM candidate ORDER BY consider_rate Desc LIMIT 3

E**D**its:

<select> add support rate </select> <select> remove consider rate </select> <orderBy> add support rate </orderBy> <orderBy> remove consider rate </orderBy>

**T**emplate Feedback:

in step 1 , find support rate in place of consider rate . order the results by support rate in place of support rate .

# Evaluating the Faithfulness of the Simulated Feedback

- **Faithfulness**: Does the simulated feedback precisely reflect the user intent of error correction?
  - Traditional metrics such as BLEU (Papineni et al., 2002) cannot measure it
  - More recent metrics such as BERTScore (Zhang et al., 2019) are too generic

- Our approach: fine-tuning BERTScore with contrastive examples
  - **Template feedback** as reference

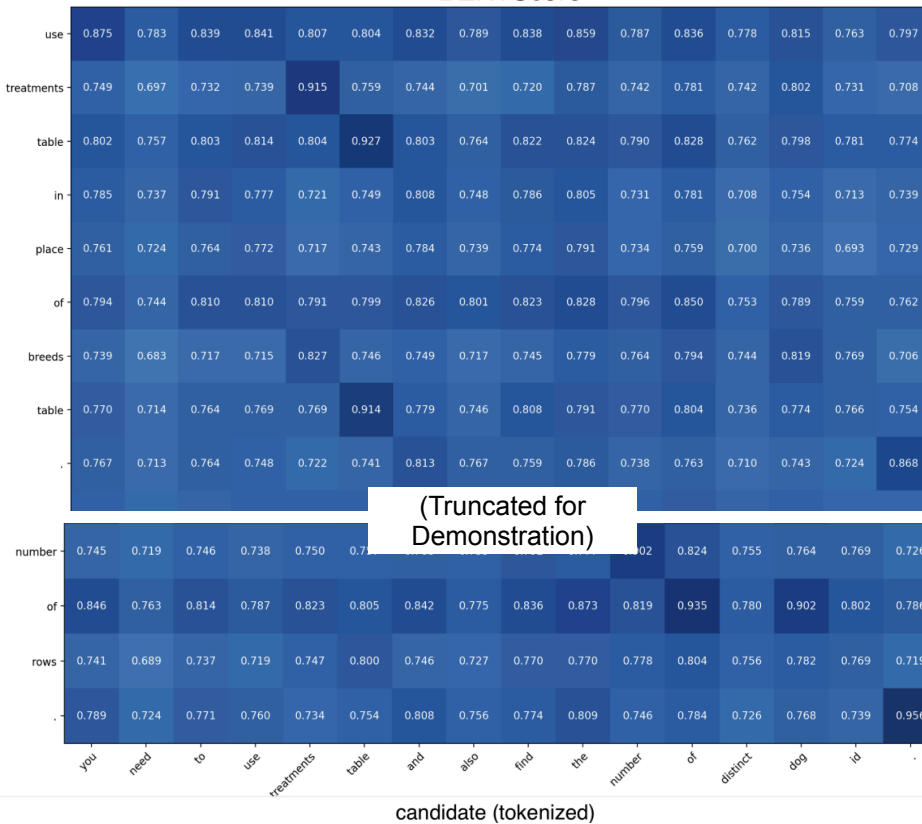| Metrics | MRR (dev) | Human |
|---|---|---|
| BLEU | 0.57 | 0.03 |
| BERTScore | 0.55 | 0.08 |
| **Our Evaluator** | **0.88** | **0.19** |



$$s(T,C) = \frac{1}{2}\left(\frac{1}{M}\sum_{m=1}^{M}\max_{n} \mathbf{A}_{nm} + \frac{1}{N}\sum_{n=1}^{N}\max_{m} \mathbf{A}_{nm}\right)$$
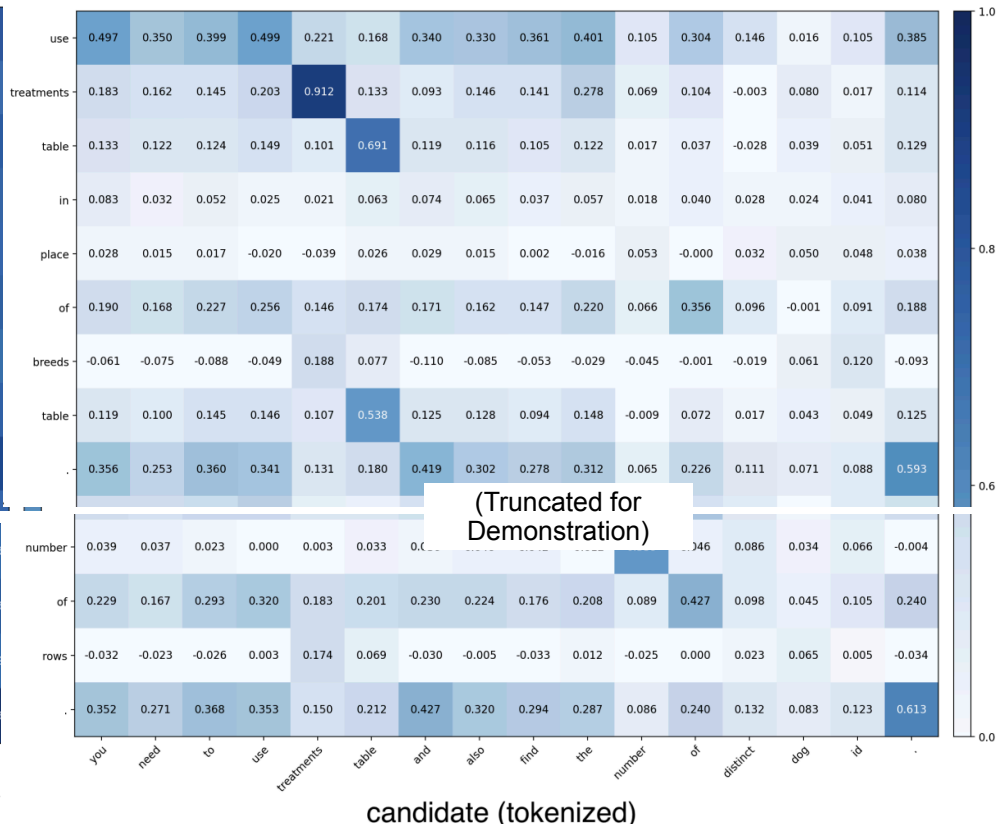
(Please refer to details in our paper)

# Evaluating the Faithfulness of the Simulated Feedback

# Example

| Model | BLEU | BERTScore | Our Evaluator |
|-------|------|-----------|---------------|
| **CWQES** | 0.132 | 0.881 | 0.491 |
| **DQES** | **0.134** | 0.882 | 0.518 |
| **TQES** | 0.125 | **0.884** | **0.535** |

## Easy Example from SPLASH-dev

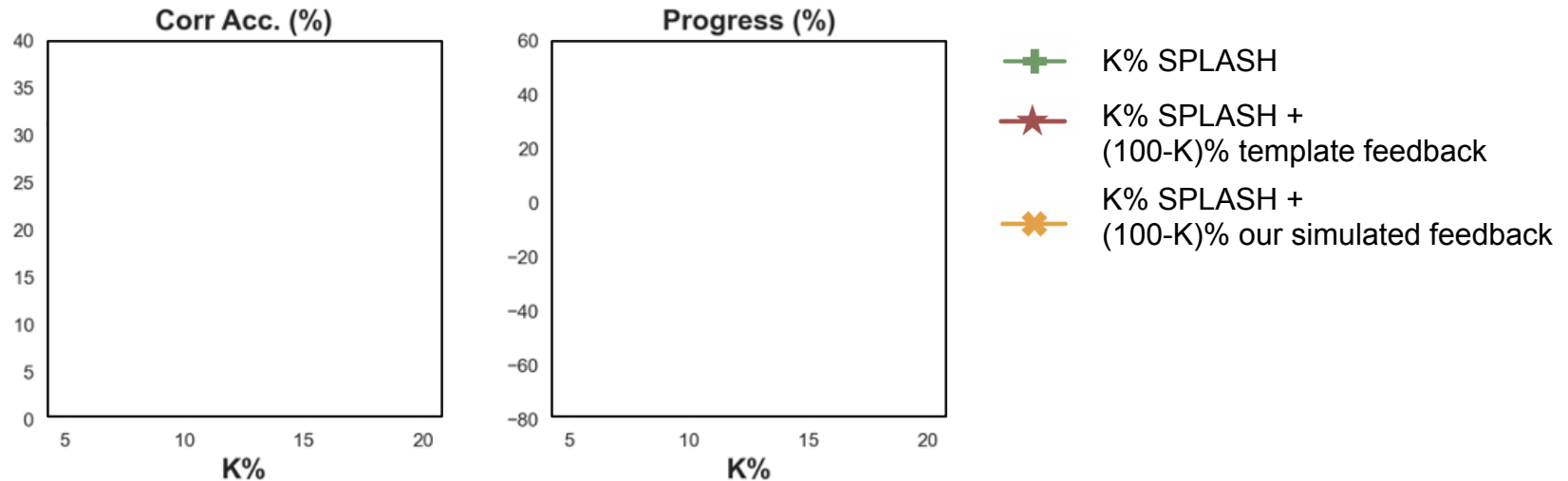| | |
|---|---|
| Question: | How many dogs went through any treatments? |
| Correct Parse: | SELECT count(DISTINCT dog_id) FROM treatments |
| Wrong Parse: | SELECT count ( * ) FROM breeds |
| Explanation: | find the number of rows in breeds table |
| Template Feedback: | use treatments table in place of breeds table . find number of different dog id in place of number of rows . |
| Human Feedback: | Change breeds table with treatments table . |

### Simulated Feedback & Evaluation Results

| | |
|---|---|
| CWQES | you need to use treatments table in place of breeds table . |
| | BLEU: 0.308, BERTScore: 0.876, Ours: 0.468 |
| DQES | you need to use treatments table and search for the number of distinct dog id . |
| | BLEU: 0.063, BERTScore: 0.879, Ours: 0.528 |
| TQES | you need to use treatments table and also find the number of distinct dog id . |
| | BLEU: 0.065, BERTScore: 0.889, Ours: 0.529 |

Our evaluator is better than
BERTScore in capturing differences in
simulated feedback

# Experimental Results

- "Low data" experiment: train a simulator with a small amount of feedback annotations, and apply it to synthesize more for model training
  - Text-to-SQL. Performance on error correction based on feedback.



K% SPLASH

K% SPLASH +
(100-K)% template feedback

K% SPLASH +
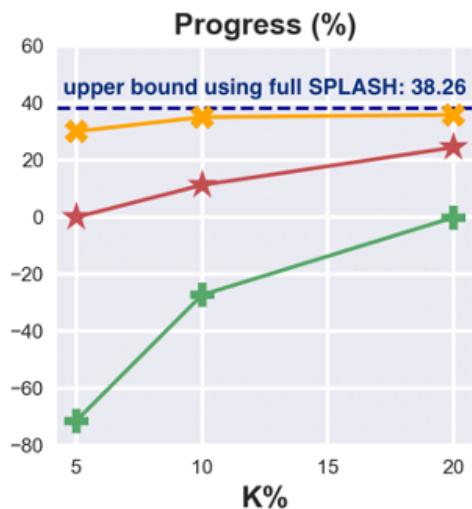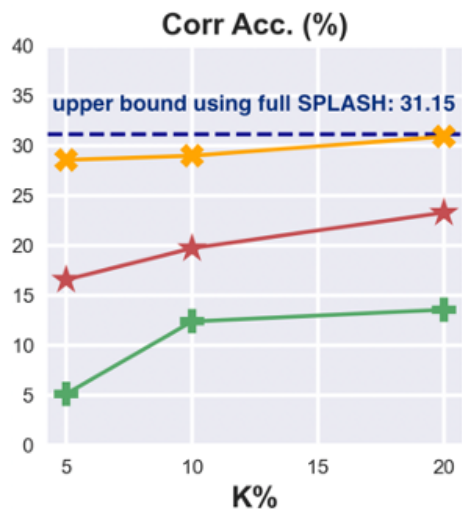(100-K)% our simulated feedback

# Experimental Results

- "Low data" experiment: train a simulator with a small amount of feedback annotations, and apply it to synthesize more for model training
  - Text-to-SQL. Performance on error correction based on feedback.



**Takeaway: Simulator is promising to augment feedback learning.**

# Discussion

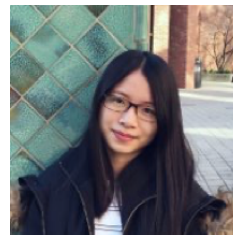- While we were working on the project (late 2022), ChatGPT came out…

- Are problems solved with ChatGPT?
  - If it does, this feedback simulator is not necessary:(
  - However, NO! Feedback modeling is not trivial even for ChatGPTI
  - Could be even more challenging with real human users, e.g., humans may not fully understand the code explanation

Hao Yan       Thomas LaToza       Ziyu Yao

Work in Progress, 2024

# Interactive Code Generation w/ ChatGPT-3.5

- Focus: *non-professional programmers* who have basic knowledge of computation and mathematics but are not professional in programming

- Very painful for them to interact with vanilla ChatGPT for programming
  - Users cannot understand or verify complicated code (they are not professional!)
  - Unstructured, back-and-forth queries for code explanation lead to frustration

- Experimental tasks: text-to-SQL and Python code generation

**Human-LLM Interaction via Language Feedback**

# Text-to-Python



LLM-Generated Code Explanation

User NL Feedback

# User Study Results

- Overall, how does our system help users in programming?
  - Double the success rate of vanilla ChatGPT-3.5, but still large room for improvement (20% for SQL and 50% for Python)

- Can users identify potential problems from our code explanation?
  - Yes but not always, for ~50% (SQL) and ~80% (Python) of the incorrect generations

- How do users provide NL feedback when they identify problems?
  - Direct instruction for error correction (58% for SQL and 70% for Python), question rephrasing, or step-by-step instructions

- Can the LLM understand the user feedback and successfully incorporate it for error correction?
  - Still very challenging! e.g., 35% (SQL) and 65% (Python) success rates for "direct instruction for error correction" feedback type

# Open Research Problems

- Future of human-LLM interaction
  - How to prompt LLMs to generate explanations that are *helpful to users*?
  - Psychological problems, e.g., cognitive bias, sycophancy (Wei et al., 2023)
  - Personalization requires modeling users beyond their feedback

- Improve *human feedback following*
  - Many efforts on *instruction following* (Webson and Pavlick 2022; Jang et al., 2022)
  - Being more challenging given the huge language variation of human feedback

- Benchmark for human-LLM interactions
  - Still an understudied field. Recent work: MINT (Wang et al., 2023)
  - Our work characterized how humans express feedback in AI-assisted programming

# This Talk: Building NLIs in the Age of LLMs



Topic 1: Modeling Language Feedback in Human-NLI Interaction
(Task: Text-to-Code Generation)

Topic 2: Saving the Monetary Cost of LLM API Usage
(Task: Arithmetic/Symbolic/etc. Reasoning)

# The Trade-Off between LLMs' Cost($) and Performance

- More powerful, but also more expensive, LLMs
  - E.g., GPT-4 vs. GPT-3.5-turbo

| Model       | Input              | Output             |
|-------------|--------------------|--------------------|
| 8K context  | $0.03 / 1K tokens  | $0.06 / 1K tokens  |
| 32K context | $0.06 / 1K tokens  | $0.12 / 1K tokens  |

| Model       | Input               | Output             |
|-------------|---------------------|--------------------|
| 4K context  | $0.0015 / 1K tokens | $0.002 / 1K tokens |
| 16K context | $0.003 / 1K tokens  | $0.004 / 1K tokens |

GPT-4

GPT-3.5-turbo

20x $ for input
30x $ for output
More powerful
but less
affordable!

https://openai.com/pricing

**How can we save $ without sacrificing task performance?**
**(Focus: Reasoning tasks)**

27

# LLM Cascades with Mixture of Thought Representations for Cost-Efficient Reasoning

Murong Yue     Jie Zhao     Min Zhang     Liang Du     Ziyu Yao

ICLR 2024

# LLM Cascades for Cost Saving

Intuition: easy questions can be handled by relatively weaker (and cheaper) LLMs to save $.

*Decision making based on textual descriptions of question and answer;*
*Do not work for Reasoning*



Final cost: $C = C^w + C^d + \mathbb{1}_{reject} \, C^s$

Extreme cases: only weaker LLM or only stronger LLM

# Reasoning with Thought Representations

Chain of Thought (**CoT**; Wei et al., 2022) &
Program of Thought (**PoT**; Chen et al., 2022,
Gao et al., 2022)

Q: A robe takes 2 bolts of blue fiber and half that
much white fiber.  How many bolts in total does
it take?

A (CoT): It takes 2/2=1 bolt of white fiber. So the
total amount of fabric is 2+1=3 bolts of fabric.
ans=3

A (PoT):
# Python code, return ans
bolts_of_blue_fiber = 2
bolts_of_white_fiber = num_of_blue_fiber / 2
ans = bolts_of_blue_fiber + bolts_of_white_fiber

Q: Today is Christmas Eve of 1937. What is the date
tomorrow in MM/DD/YYYY?

(CoT) Explain: Today is the Christmas Eve of 1937, so today
is 12/24/1937.
Today is 12/24/1937, the date tomorrow is 12/25/1937.
A: 12/25/1937

(PoT) # Write Python Code to solve the following questions.
from datetime import date, timedelta
from dateutil.relativedelta import relativedelta

# Q: Today is Christmas Eve of 1937. What is the date
tomorrow in MM/DD/YYYY?
# today is Christmas Eve of 1937, then today is 12/24/1937
today = date(1937, 12, 24)
# tomorrow
date_tomorrow = today + relativedelta(days=1)
# The answer formatted with %m/%d/%Y is
ans = date_tomorrow.strftime('%m/%d/%Y')

# This Work: Answer Consistency-based Decision Making

- Idea: if the weaker LLM is uncertain about an answer, the question could be too challenging for it to solve

- How to measure an LLM's certainty on an answer?
    - See how often it *samples* the same answer to the given question
    - Same idea as "Self Consistency (SC)" (Wang et al., 2023)

- Questions:
    - Where to sample the answers for better judgment?
    - How to quantify the answer consistency?

# Approaches

- Vote-based decision making



$$s = \frac{\sum_{i=1}^{K} \mathbb{1}_{A_i^w = A^w}}{K}$$

# Approaches

- Vote-based decision making, sampling from
  - a single thought representation
  - a single demonstration set



Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?

A: It takes 2/2=1 bolt of white fiber. So the total amount of fabric is 2+1=3 bolts of fabric. ans=3

**... (M shots of CoT examples)**

Q: Test question
A:

Weaker LLM → Sampled K answers

$$s = \frac{\sum_{i=1}^{K} \mathbb{1}_{A_i^w = A^w}}{K}$$

Method: CoT-1D-Vote

# Approaches

- Vote-based decision making, sampling from
  - a single thought representation
  - a single demonstration set



$$s = \frac{\sum_{i=1}^{K} \mathbb{1}_{A_i^w = A^w}}{K}$$

Q: A robe takes 2 bolts of blue fiber and half that much white fiber.  How many bolts in total does it take?

A:
# Python code, return ans
bolts_of_blue_fiber = 2
bolts_of_white_fiber = num_of_blue_fiber / 2
ans = bolts_of_blue_fiber + bolts_of_white_fiber

**… (M shots of PoT examples)**
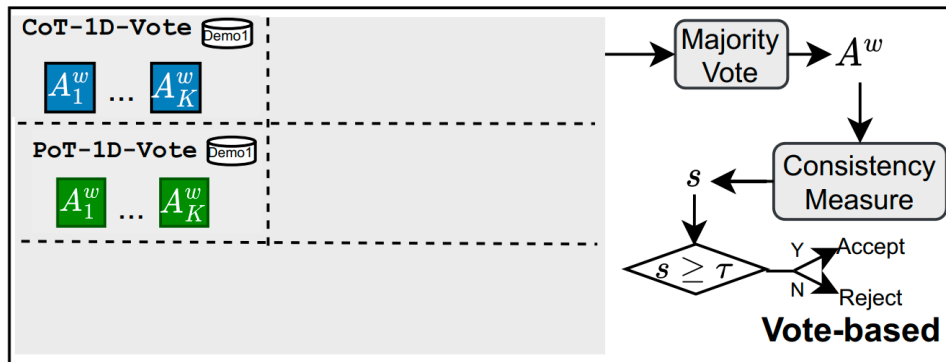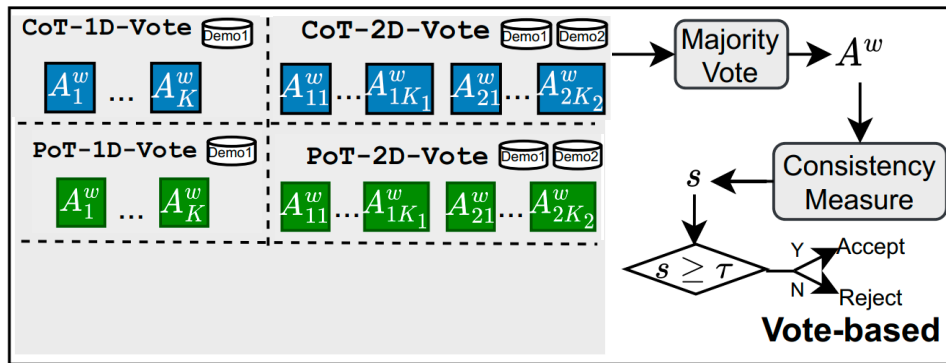
Q: Test question
A:

Weaker LLM

Sampled K answers

Method: PoT-1D-Vote

# Approaches

- Vote-based decision making, sampling from
  - a single thought representation
  - Two demonstration sets



CoT-1D-Vote [Demo1]
$A_1^w$ ... $A_K^w$

CoT-2D-Vote [Demo1] [Demo2]
$A_{11}^w$ ... $A_{1K_1}^w$  $A_{21}^w$ ... $A_{2K_2}^w$

PoT-1D-Vote [Demo1]
$A_1^w$ ... $A_K^w$

PoT-2D-Vote [Demo1] [Demo2]
$A_{11}^w$ ... $A_{1K_1}^w$  $A_{21}^w$ ... $A_{2K_2}^w$

Majority Vote → $A^w$

Consistency Measure → $s$

$s \geq \tau$  Y → Accept  N → Reject

**Vote-based**

$$s = \frac{\sum_{i=1}^{K_1} \mathbb{1}_{A_{1i}^w = A^w} + \sum_{i=1}^{K_2} \mathbb{1}_{A_{2i}^w = A^w}}{K_1 + K_2}$$

Demo1 →

Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?
A: It takes 2/2=1 bolt of white fiber. So the total amount of fabric is 2+1=3 bolts of fabric. ans=3

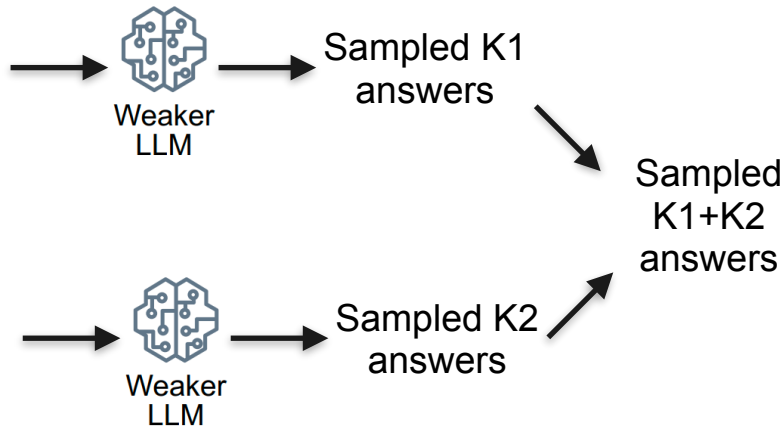**… (M shots of CoT examples from Set 1)**
Q: Test question
A:

→ Weaker LLM → Sampled K1 answers

Demo2 →

Q: Manny had 3 birthday cookie pies to share with his 24 classmates and his teacher, Mr. Keith. …
A: There is a total of 3 x 10 = 30 cookie slices… ans = 4

**… (M shots of CoT examples from Set 2)**
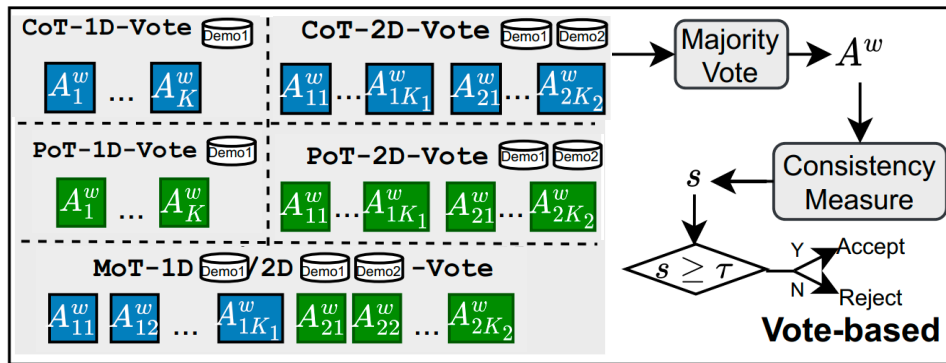Q: Test question
A:

→ Weaker LLM → Sampled K2 answers

Sampled K1+K2 answers

Method: CoT-2D-Vote
(Similarly for PoT-2D-Vote) 35

# Approaches



- Vote-based decision making, sampling from
  - Two thought representations
  - a single demonstration set

$$s = \frac{\sum_{i=1}^{K_1} \mathbb{1}_{A_{1i}^w = A^w} + \sum_{i=1}^{K_2} \mathbb{1}_{A_{2i}^w = A^w}}{K_1 + K_2}$$



Q: A robe takes 2 bolts of blue fiber and half that much white fiber.  How many bolts in total does it take?
A: It takes 2/2=1 bolt of white fiber. So the total amount of fabric is 2+1=3 bolts of fabric. ans=3

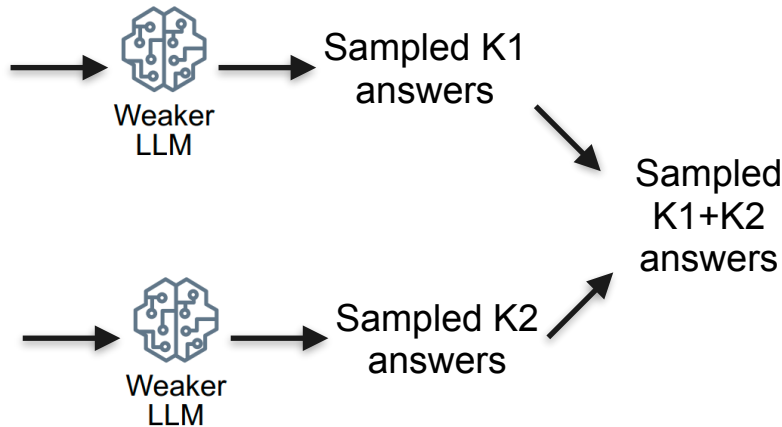**… (M shots of CoT examples from Set 1)**
Q: Test question
A:

Q: A robe takes 2 bolts of blue fiber and half that much white fiber.  How many bolts in total does it take?
A:
# Python code, return ans
…ans = bolts_of_blue_fiber + bolts_of_white_fiber
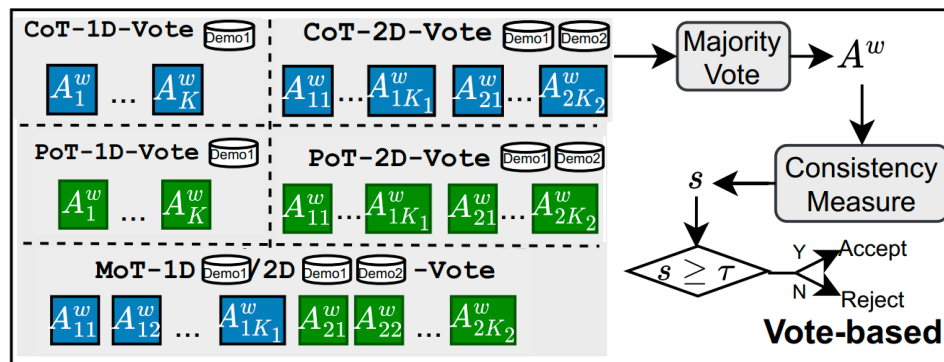
**… (M shots of PoT examples from Set 1)**
Q: Test question
A:

Weaker LLM → Sampled K1 answers

Weaker LLM → Sampled K2 answers

Sampled K1+K2 answers

Method: MoT-1D-Vote
"Mixture of Thought"

36

# Approaches

- Vote-based decision making, sampling from
  - Two thought representations
  - Two demonstration sets



$$s = \frac{\sum_{i=1}^{K_1} \mathbb{1}_{A_{1i}^w = A^w} + \sum_{i=1}^{K_2} \mathbb{1}_{A_{2i}^w = A^w}}{K_1 + K_2}$$

Q: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?
A: It takes 2/2=1 bolt of white fiber. So the total amount of fabric is 2+1=3 bolts of fabric. ans=3

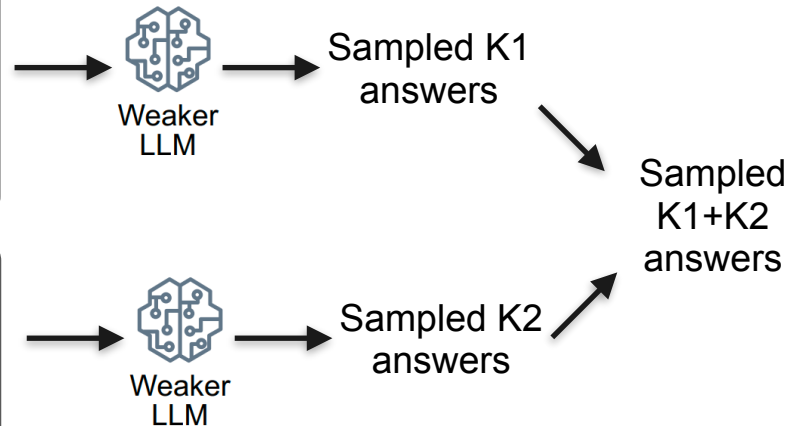**… (M shots of CoT examples from Set 1)**
Q: Test question
A:

Q: Manny had 3 birthday cookie pies to share with his 24 classmates and his teacher, Mr. Keith. …
A:
# Python code, return ans
…ans = total_cookie_pies - total_person_count

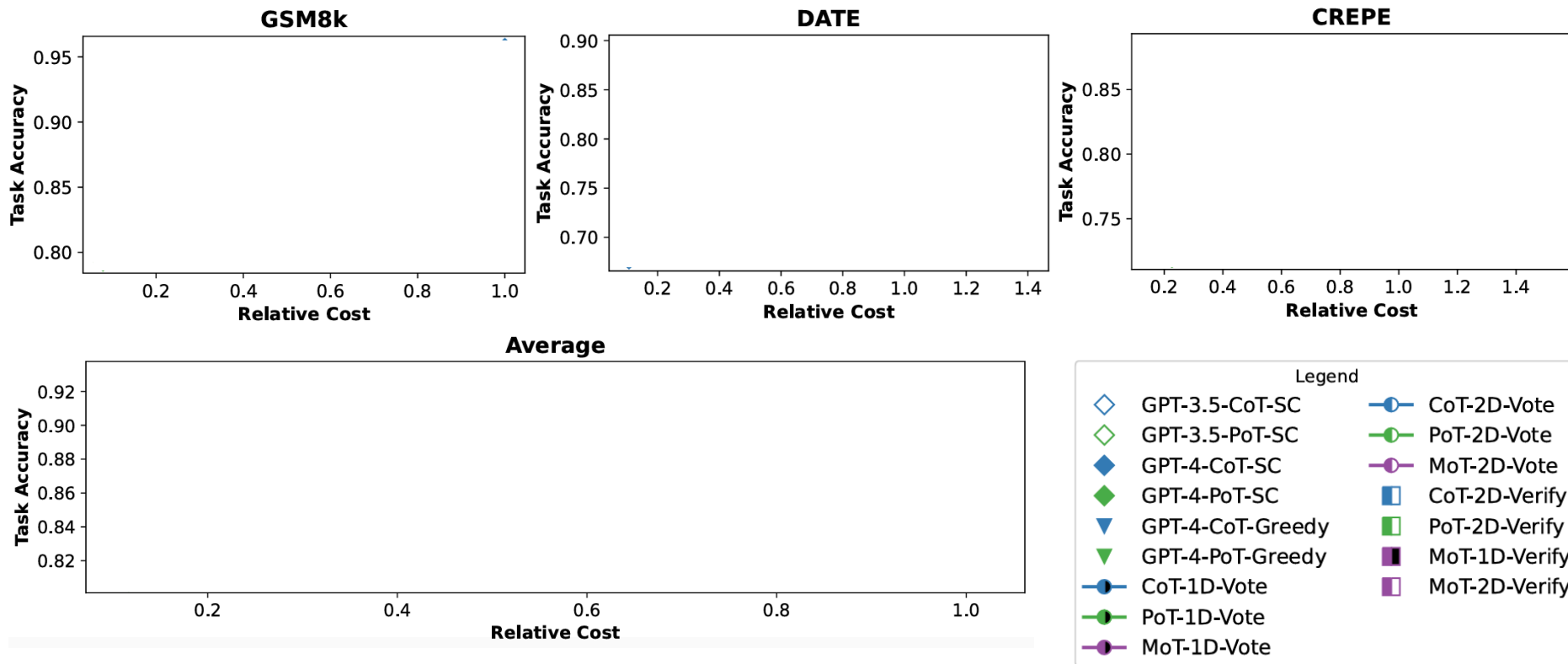**… (M shots of PoT examples from Set 2)**
Q: Test question
A:

Weaker LLM

Weaker LLM

Sampled K1 answers

Sampled K2 answers

Sampled K1+K2 answers

Method: MoT-2D-Vote

37

# Experimental Results
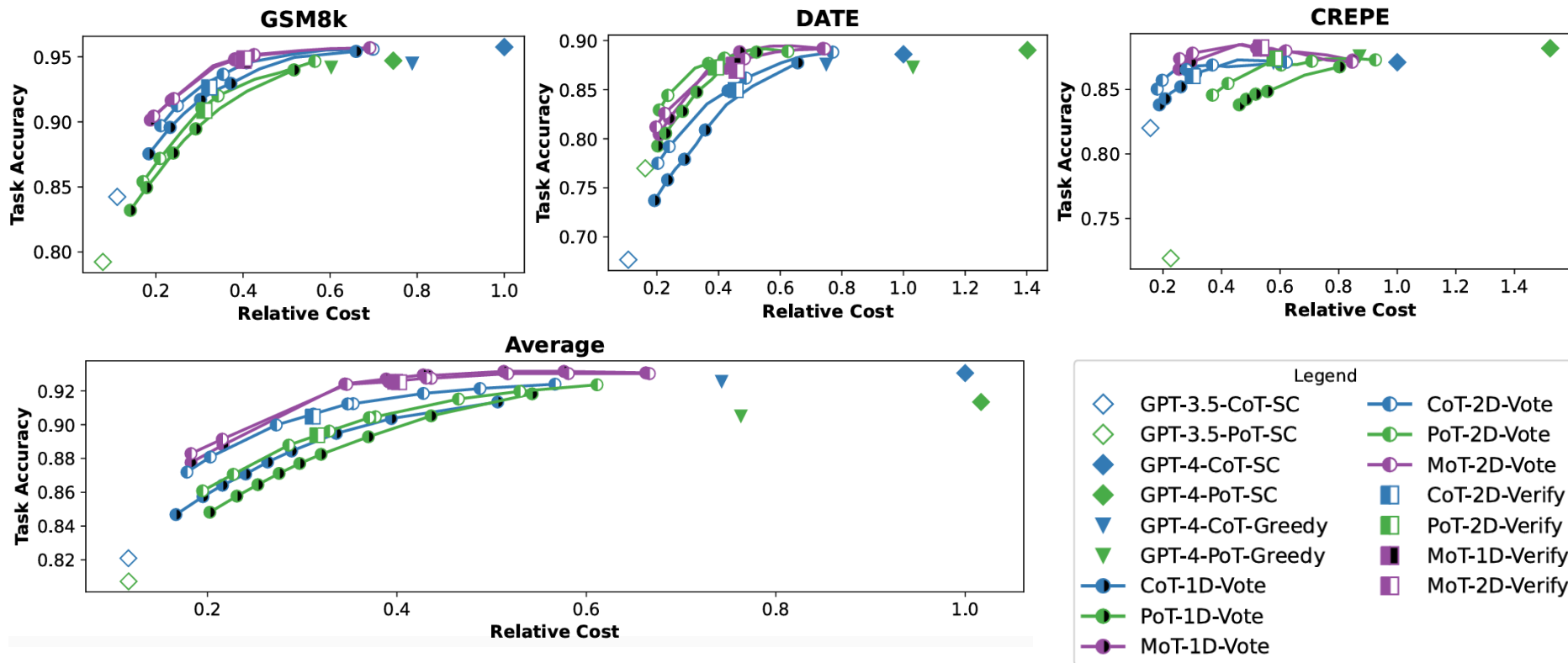
(Average over GSM8k, ASDIV, TabMWP, DATE, Navigate, CREPE)

# Experimental Results

Weaker LLM: GPT-3.5-turbo
Stronger LLM: GPT-4



(Average over GSM8k, ASDIV, TabMWP, DATE, Navigate, CREPE)

# Diversify Thought Representations for Uncertainty Measurement

- Mixture of Thought (MoT) introduces diverse "opinions", similar to *model ensemble*, which helps uncertainty measurement

**Question**: Blake and Kelly are having a contest to see who can run the most in 15 minutes. They decide to do it on a football field that is 100 yards long. Blake runs back and forth 15 times. Kelly runs back and forth once, and then decides that she doesn't want to run next to Blake, so she starts to run to the 40-yard line and back. She does this 34 times. How much farther does the winner run than the loser?

**Gold**

Black runs 200 yards for each back and forth because...He runs 3000 yards in total because...Kelly runs 200 yards at first because 100 x 2 = 200. She runs 80 yards each time she does a back and forth to the 40-yard line because... She runs 2720 from these back and forth. She runs 2920 in total because...The winner ran 80 yards more because 3000 - 2920=80.

**CoT 1**

Answer: [...(Ignored)] Kelly runs to the 40-yard line and back 34 times, so she runs 34*80=2720 yards. Thus, Blake runs 3000-2720=280 yards farther than Kelly.

**CoT 2**

Answer: [...(Ignored)] Kelly runs 34 x (40 x 2) = 2720 yards. Blake runs 3000 - 2720 = 280 yards farther than Kelly.

**PoT 1**

```
# Python code, return ans
[...(Ignored)]
distance_covered_by_kelly =
(length_of_football_field *
num_of_laps_for_kelly * 2) + (40 *
num_of_laps_for_kelly * 2)
ans = abs(distance_covered_by_blake
- distance_covered_by_kelly)
        (Answer via Python execution: 6520.0)
```

☐ Logic Generation Error          ☐ Value Grounding Error

# Diversify Thought Representations for Uncertainty Measurement

- Mixture of Thought (MoT) introduces diverse "opinions", similar to *model ensemble*, which helps uncertainty measurement
  - Applies to factual reasoning tasks as well

**Question**

Is a curling iron necessary in curling?

**Golden**

A curling iron is a tool used to make the hair curly using heat. The sport of curling requires curling brooms, stones (rocks), and curling shoes. A curling iron is not necessary in curling.

**CoT**

Yes, a curling iron is necessary for curling as it is the tool used to create curls in the hair.

**PoT**

necessary_equipment_for_curling = ["curling stone", "curling broom"]
 equipment_used_in_curling_with_curling_iron = ["curling iron"]
ans = ("curling iron" in necessary_equipment_for_curling)

Dataset: StrategyQA (Geva et a., 2021)
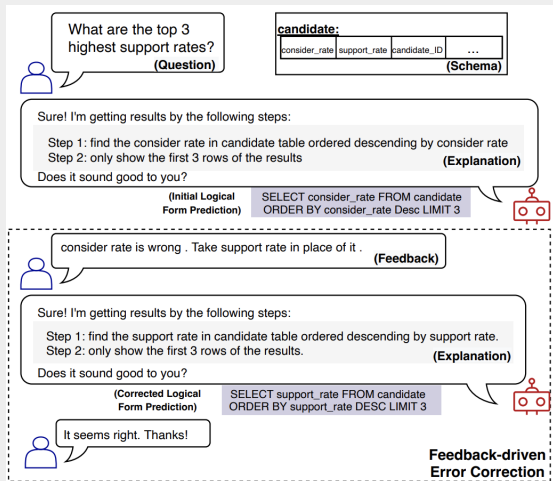
# Other Findings

- Decision-making based on textual hints (e.g., FrugalGPT)?
  - Takeaway: it is very challenging to distinguish between easy and hard questions solely based on textual hints

- How weak can the weaker LLM be?
  - Experiments using LLAMA2 13B
  - Takeaway: if an LLM is too weak, it won't contribute to the cost saving, i.e., all questions will eventually be passed to the stronger LLM

- Can outputs from the weaker LLM be hints to improve the stronger LLM?
  - No, and they actually confuse the stronger LLM

**LLM Cascade for Cost-Saving Query**

# Discussion & Future Work

- LLM Uncertainty: Does an LLM know when it doesn't know?
  - Many discussions (Kadavath et al., 2022; Xiong et al., 2023; etc.)
  - We showed the promise of mixing thought representations w/ vote-based metric
  - Generalize to tasks where we cannot vote? (e.g., text generation)

- Ensemble of multiple LLMs/LLM-powered agents
  - Similar synergy between CoT and PoT: e.g., model selection (Zhao et a., 2023), fine-tuning (Yue et al., 2023)
  - Generally speaking, tasking a cohort of LLMs, e.g., weaker vs. stronger, in-house vs. closed API, domain-specific vs. domain-general, etc.
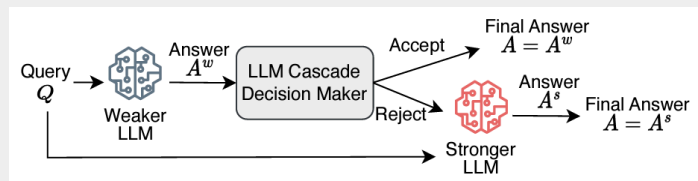
# This Talk: Building NLIs in the Age of LLMs



Topic 1: Modeling Language Feedback in Human-NLI Interaction
(Task: Text-to-Code Generation)

**Building simulators for feedback modeling; still challenges for LLMs as interactive NLIs**

Topic 2: Saving the Monetary Cost of LLM API Usage
(Task: Arithmetic/Symbolic/etc. Reasoning)

**LLM cascades with Mixture-of-Thought decision-making helps uncertainty measurement and enables cost efficiency**

# New Preprint: LLM Agents for Education

- LLM agents simulating students in collaborative mathematical problem solving
  - A platform for students to practice their math modeling skills
  - Helping students with limited educational resources

**Microsoft** Research
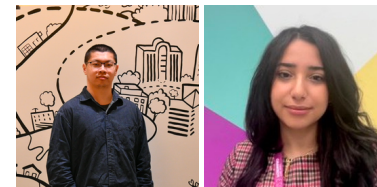
## Accelerating Foundation Models Research

Engaging the broader community in reimagining computing research

**Problem Description**

Martha hopes to sell 500 mugs of soup, each with a white or brown bread roll. She will sell a mug of soup with a bread roll for $1.25. She can buy the soup in 2.5 liter. Each bottle of soup costs $5 and provides ten servings. Bread rolls are sold in packs of 10. Each pack costs $2. To better meet the requirement, she made a preference survey. The response cards are shown in the left Data Panel. What exactly should Martha buy so that she can make the most profit?

**Data Panel**

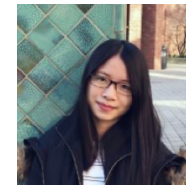Martha made a preference survey and received the following 40 responses:

**Alice** Hi guys, we should think about which flavors are popular to avoid wasting any soup.

**Bob** I agree! Let's say she decides on those flavors! Now how do we decide which flavor to buy?

**Charlie** That's exactly what the survey responses told us! I counted it a bit. Among 40 people, 15 pick tomato and 15/40=0.375. So I think it says 0.375% people like tomato!

**Human Student** Wait, did you count it right? I found 16 tomato actually... and the percentage should be multiplied by 100!

**Send** **Skip**

Murong Yue (AI/LLM)
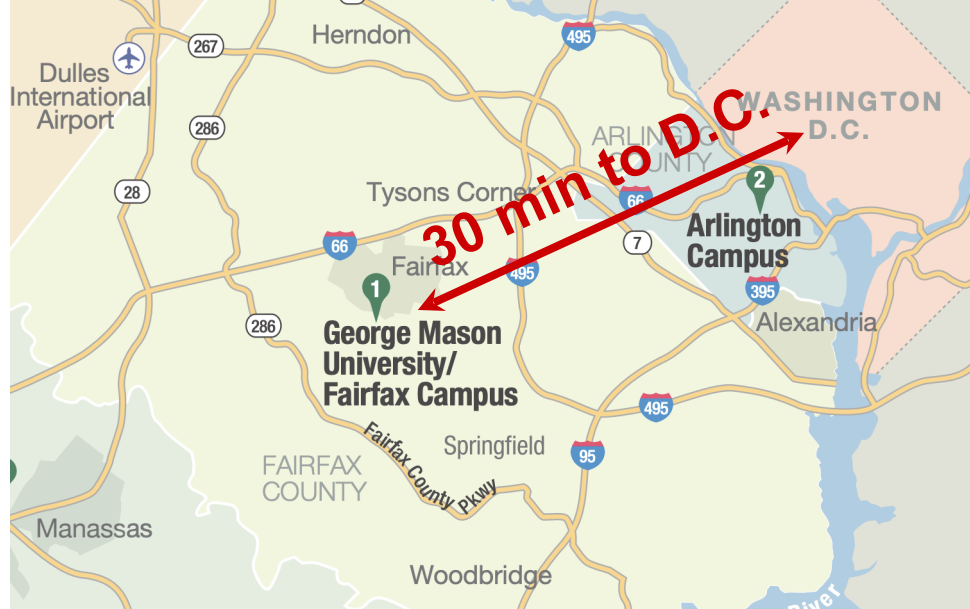
Wijdane Mifdal (AI/LLM)

Janice Zhang (HCI)

Ziyu Yao (AI/LLM)

Jenn Suh (MathEdu)

# Thank You!

Email: ziyuyao@gmu.edu
Webpage: https://ziyuyao.org/