# LLM-based Proactive Query Management

## Yang Deng

Singapore Management University

# Outline

- ❏ **Why Proactive Query Management?**
    - ❏ Overconfidence on Unanswerable Queries
    - ❏ Random Responses to Ambiguous Queries
- ❏ **Unanswerable Query Mitigation**
    - ❏ Refusal Fine-tuning
    - ❏ Uncertainty-based Reinforcement Learning
    - ❏ Self-alignment
- ❏ **Ambiguous Query Clarification**
    - ❏ In-Context Learning
    - ❏ Reinforcement Learning
    - ❏ Preference Optimization

# Outline

- ❏ **Why Proactive Query Management?**
  - ❏ Overconfidence on Unanswerable Queries
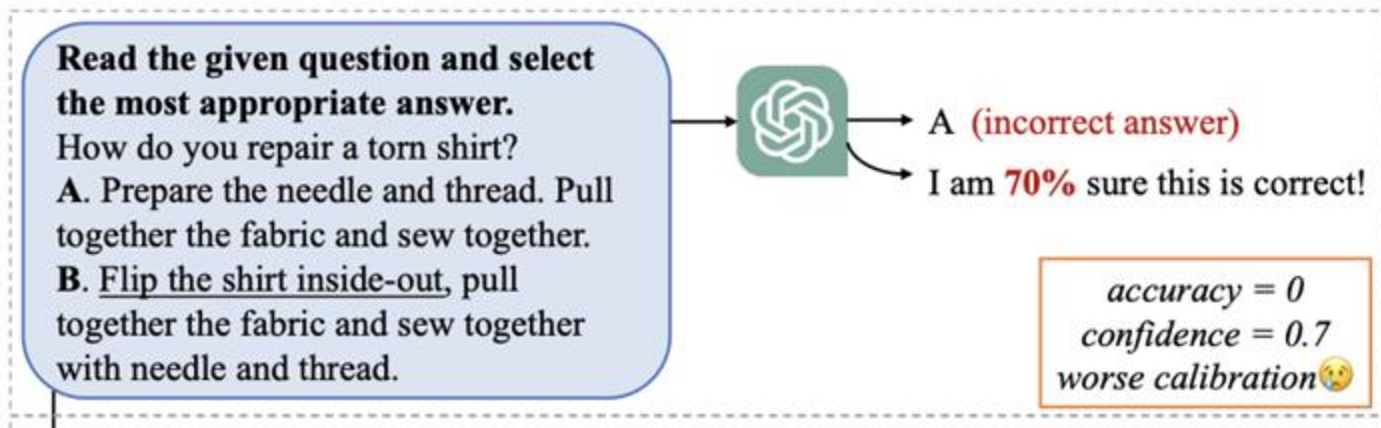  - ❏ Random Responses to Ambiguous Queries
- ❏ Unanswerable Query Mitigation
  - ❏ Refusal Fine-tuning
  - ❏ Uncertainty-based Reinforcement Learning
  - ❏ Self-alignment
- ❏ Ambiguous Query Clarification
  - ❏ In-Context Learning
  - ❏ Reinforcement Learning
  - ❏ Preference Optimization

# Overconfidence on Unanswerable Queries

*Kapoor et al., "Large Language Models Must Be Taught to Know What They Don't Know" (NeurIPS '24)*
*Li et al., "Think Twice Before Assure: Confidence Estimation for Large Language Models through Reflection on Multiple Answers" (EMNLP '24 Findings)*

# Overconfidence on Unanswerable Queries



## QnotA

**Incomplete Information**
Sam played the game against Iran, did Sam play in Iran?

**Future Questions**
Who won the 2034 FIFA world cup?

**Incorrect Information**
Who is prime minister of California state?

**Ambiguous**
Look at the dog with one eye, does the dog have only one eye?

**Unmeasurable**
How many drops of water are in the pacific ocean?

The question itself is unanswerable.

❑ Incomplete: questions are not specific enough

❑ Future: questions about the future we cannot know

❑ Incorrect: questions that contain an incorrect assumption or statement

❑ …

**Q:** What animal can be found at the top of the men's Wimbledon trophy?

**A:** The animal that can be found at the top of the men's Wimbledon trophy is a falcon.

**Direct Answer**

⚠ There is a **fruit-like design** at the top of the men's Wimbledon trophy, instead of an **animal**.

*Agarwal et al., "Can NLP models 'identify', 'distinguish', and 'justify' questions that don't have a definitive answer?" (TrustNLP@ACL '23)*
*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Random Responses to Ambiguous Queries

| Method | Shot | Prompt | Abg-CoQA | | | PACIFIC | | |
| | | | CNP | CQG | | CNP | CQG | |
| | | | F1 | BLEU-1 | Help. | F1 | ROUGE-2 | Help. |
|---|---|---|---|---|---|---|---|---|
| Baseline | - | - | 22.1 | 36.5 | 30.0 | 79.0 | 69.2 | 38.2 |
| SOTA | - | - | 23.6 | 38.2 | 56.0 | 86.9 | 90.7 | 80.1 |
| Vicuna-13B | 0 | Standard | - | 11.3 | 0.0 | - | 1.2 | 0.0 |
| | 1 | Standard | - | 11.4 | 0.0 | - | 2.5 | 0.0 |
| | 0 | Proactive | 4.1 | 13.2 | 0.0 | 2.3 | 2.3 | 0.0 |
| | 1 | Proactive | 12.1 | 13.2 | 4.5 | 0.0 | 3.3 | 0.0 |
| | 0 | ProCoT | 1.4 | 21.3 | 9.1 | 9.7 | 3.8 | 10.5 |
| | 1 | ProCoT | **18.3** | **23.7** | **22.7** | 27.0 | **41.3** | **33.1** |
| ChatGPT | 0 | Standard | - | 12.1 | 0.0 | - | 2.2 | 0.0 |
| | 1 | Standard | - | 12.3 | 0.0 | - | 2.0 | 0.0 |
| | 0 | Proactive | 22.0 | 13.7 | 17.6 | 19.4 | 2.9 | 0.0 |
| | 1 | Proactive | 20.4 | **23.4** | 23.5 | 17.7 | 14.0 | 12.5 |
| | 0 | ProCoT | 23.8 | 21.6 | 32.4 | **28.0** | **21.5** | 26.7 |
| | 1 | ProCoT | **27.9** | 18.4 | **45.9** | 27.7 | 16.2 | **35.8** |

**LLMs barely ask clarification questions, even when the user query is ambiguous.**

*Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Random Responses to Ambiguous Queries

| Category | Sources | Distribution | | |
|---|---|---|---|---|
| | | Ambig. | Non-Ambig. | ALL |
| Unfamiliar | ALCUNA | 684 | 547 | 1231 |
| Contradiction | AmbiTask | 600 | 600 | 1200 |
| Lexical | AmbER,AmbiPun | 815 | 921 | 1,736 |
| Semantic | AmbiCoref | 400 | 400 | 800 |
| What | AmbigQA, Dolly | 1255 | | |
| Whom | AmbigQA, Dolly | 762 | 3884 in total | 7167 in total |
| When | AmbigQA, Dolly | 779 | | |
| Where | AmbigQA, Dolly | 487 | | |

**Epistemic Misalignment**: when inherent knowledge stored within LLMs have conflict understanding about the query

**Linguistic Ambiguity**: when a word, phrase, or statement can be interpreted in multiple ways due to its imprecise or unclear meaning

**Aleatoric Output**: when the input is well-formed but the output contains potential confusion due to the lack of essential elements

*Zhang et al., "CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models" (ACL '24)*

# Random Responses to Ambiguous Knowledge

| Method | Shot | Prompt | Abg-CoQA | | | PACIFIC | | |
| | | | CNP | CQG | | CNP | CQG | |
| | | | F1 | BLEU-1 | Help. | F1 | ROUGE-2 | Help. |
|---|---|---|---|---|---|---|---|---|
| Baseline | - | - | 22.1 | 36.5 | 30.0 | 79.0 | 69.2 | 38.2 |
| SOTA | - | - | 23.6 | 38.2 | 56.0 | 86.9 | 90.7 | 80.1 |
| Vicuna-13B | 0 | Standard | - | 11.3 | 0.0 | - | 1.2 | 0.0 |
| | 1 | Standard | - | 11.4 | 0.0 | - | 2.5 | 0.0 |
| | 0 | Proactive | 4.1 | 13.2 | 0.0 | 2.3 | 2.3 | 0.0 |
| | 1 | Proactive | 12.1 | 13.2 | 4.5 | 0.0 | 3.3 | 0.0 |
| | 0 | ProCoT | 1.4 | 21.3 | 9.1 | 9.7 | 3.8 | 10.5 |
| | 1 | ProCoT | **18.3** | **23.7** | **22.7** | **27.0** | **41.3** | **33.1** |
| ChatGPT | 0 | Standard | - | 12.1 | 0.0 | - | 2.2 | 0.0 |
| | 1 | Standard | - | 12.3 | 0.0 | - | 2.0 | 0.0 |
| | 0 | Proactive | 22.0 | 13.7 | 17.6 | 19.4 | 2.9 | 0.0 |
| | 1 | Proactive | 20.4 | **23.4** | 23.5 | 17.7 | 14.0 | 12.5 |
| | 0 | ProCoT | 23.8 | 21.6 | 32.4 | **28.0** | **21.5** | 26.7 |
| | 1 | ProCoT | **27.9** | 18.4 | **45.9** | 27.7 | 16.2 | **35.8** |

**LLMs barely ask clarification questions, even when the user query is ambiguous.**

*Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Random Responses to Ambiguous Knowledge

| Dimension | Category | Explanation | Example |
|---|---|---|---|
| Epistemic Misalignment | **UNFAMILIAR** | Query contains unfamiliar entities or facts | Find the price of Samsung Chromecast. |
| | **CONTRADICTION** | Query contains self-contradictions | Output 'X' if the sentence contains [category withhold] and 'Y' otherwise. The critic is in the restaurant.>X. The butterfly is in the river.>Y. The boar is in the theatre.>? |
| Linguistic Ambiguity | **LEXICAL** | Query contains terms with multiple meanings | Tell me about the source of Nile. |
| | **SEMANTIC** | Query lacks of context leading multiple interpretations | When did he land on the moon? |
| Aleatoric Output | **WHO** | Query output contains confusion due to missing personal elements | Suggest me some gifts for my mother. |
| | **WHEN** | Query output contains confusion due to missing temporal elements | How many goals did Argentina score in the World Cup? |
| | **WHERE** | Query output contains confusion due to missing spatial elements | Tell me how to reach New York. |
| | **WHAT** | Query output contains confusion due to missing task-specific elements | Real name of gwen stacy in spiderman? |

*Zhang et al., "CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models" (ACL '24)*

# Random Responses to Ambiguous Knowledge

| Dimension | Category | Explanation | Example |
|---|---|---|---|
| Epistemic Misalignment | **UNFAMILIAR** | Query contains unfamiliar entities or facts | Find the price of Samsung Chromecast. |
| | **CONTRADICTION** | Query contains self-contradictions | Output 'X' if the sentence contains [category withhold] and 'Y' otherwise. The critic is in the restaurant.>X. The butterfly is in the river.>Y. The boar is in the theatre.>? |

**Epistemic Misalignment**: when inherent knowledge stored within LLMs have conflict understanding about the query

*Zhang et al., "CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models" (ACL '24)*

# Random Responses to Ambiguous Knowledge

| Linguistic Ambiguity | LEXICAL | Query contains terms with multiple meanings | Tell me about the source of Nile. |
|---|---|---|---|
| | SEMANTIC | Query lacks of context leading multiple interpretations | When did he land on the moon? |

**Linguistic Ambiguity**: when a word, phrase, or statement can be interpreted in multiple ways due to its imprecise or unclear meaning

*Zhang et al., "CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models" (ACL '24)*
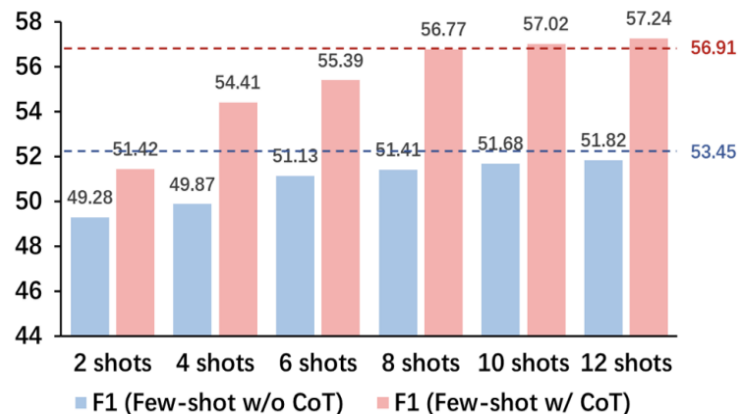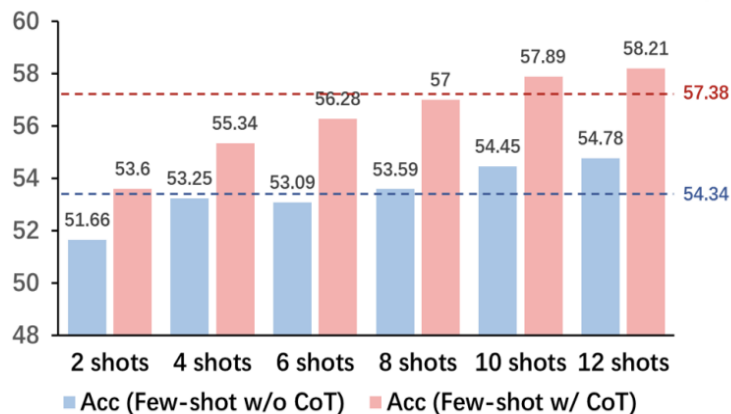
# Random Responses to Ambiguous Knowledge

**Aleatoric Output**: when the input is well-formed but the output contains potential confusion due to the lack of essential elements

| | | | |
|---|---|---|---|
| Aleatoric Output | **WHO** | Query output contains confusion due to missing personal elements | Suggest me some gifts for my mother. |
| | **WHEN** | Query output contains confusion due to missing temporal elements | How many goals did Argentina score in the World Cup? |
| | **WHERE** | Query output contains confusion due to missing spatial elements | Tell me how to reach New York. |
| | **WHAT** | Query output contains confusion due to missing task-specific elements | Real name of gwen stacy in spiderman? |

*Zhang et al., "CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models" (ACL '24)*

# Random Responses to Ambiguous Knowledge

| Methods | Epistemic Misalignment | | | | Linguistic Ambiguity | | | | Aleatoric Output | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | contradiction | | unfamiliar | | lexical | | semantic | | what | | whom | | when | | where | |
| | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 | Acc. | F1 |
| Vicuna-13B | 51.75 | 37.11 | 59.50 | 59.33 | **72.00** | **71.52** | 49.75 | 33.22 | 44.81 | 41.74 | 46.95 | 44.57 | 44.86 | 41.82 | 42.96 | 39.24 |
| Llama2-13B-I | 49.50 | 33.11 | 46.75 | 46.47 | 52.50 | 49.20 | 48.50 | 41.31 | 30.24 | 30.14 | 31.37 | 31.32 | 27.97 | 27.72 | 29.57 | 29.44 |
| Llama2-13B | 50.25 | 33.89 | 54.25 | 46.65 | 56.75 | 49.11 | 50.00 | 33.33 | 34.73 | 34.64 | 36.86 | 36.85 | 34.27 | 34.16 | 34.17 | 34.05 |
| Llama2-70B | **63.25** | **58.83** | 50.75 | 35.81 | 55.25 | 44.04 | 50.00 | 33.33 | 31.04 | 30.77 | 31.37 | 31.07 | 31.37 | 31.07 | 31.47 | 31.16 |
| ChatGPT | 38.00 | 28.17 | **60.00** | **59.67** | <u>58.75</u> | <u>58.06</u> | **50.75** | **49.32** | **65.40** | **50.54** | **68.77** | **57.48** | **65.00** | **45.66** | **63.10** | **45.24** |

*Zhang et al., "CLAMBER: A Benchmark of Identifying and Clarifying Ambiguous Information Needs in Large Language Models" (ACL '24)*

# Outline

❏ **Why Proactive Query Management?**

  ❏ Overconfidence on Unanswerable Queries

  ❏ Random Responses to Ambiguous Queries

❏ **Unanswerable Query Mitigation**

  ❏ Refusal Fine-tuning

  ❏ Uncertainty-based Reinforcement Learning

  ❏ Self-alignment

❏ **Ambiguous Query Clarification**

  ❏ In-Context Learning

  ❏ Reinforcement Learning

  ❏ Preference Optimization

# Outline

❏ **Why Proactive Query Management?**

    ❏ Overconfidence on Unanswerable Queries
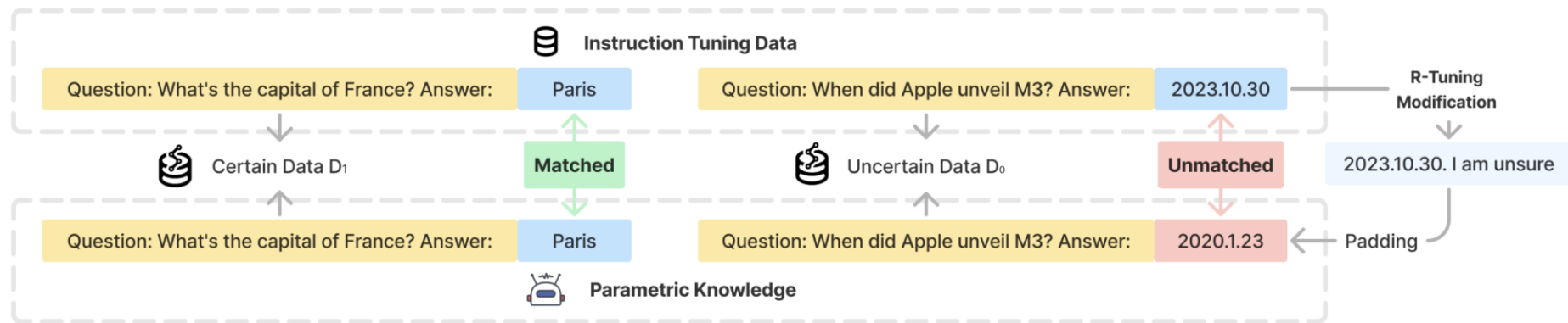
    ❏ Random Responses to Ambiguous Queries

❏ **Unanswerable Query Mitigation**

    ❏ Refusal Fine-tuning

    ❏ Uncertainty-based Reinforcement Learning

    ❏ Self-alignment

❏ **Ambiguous Query Clarification**

    ❏ In-Context Learning

    ❏ Reinforcement Learning

    ❏ Preference Optimization

# Refusal-Aware Instruction Tuning (R-Tuning)



❏ **Refusal-Aware Data Identification**

The question with mismatch between the prediction and the ground-truth label results

❏ **Refusal-Aware Data Construction**

Construct template-based refusal responses, e.g., "I am unsure"

❏ **Supervised Fine-tuning**

*Zhang et al., "R-Tuning: Instructing Large Language Models to Say ``I Don't Know'?" (NAACL '24)*

# Outline

- ❏ Why Proactive Query Management?
  - ❏ Overconfidence on Unanswerable Queries
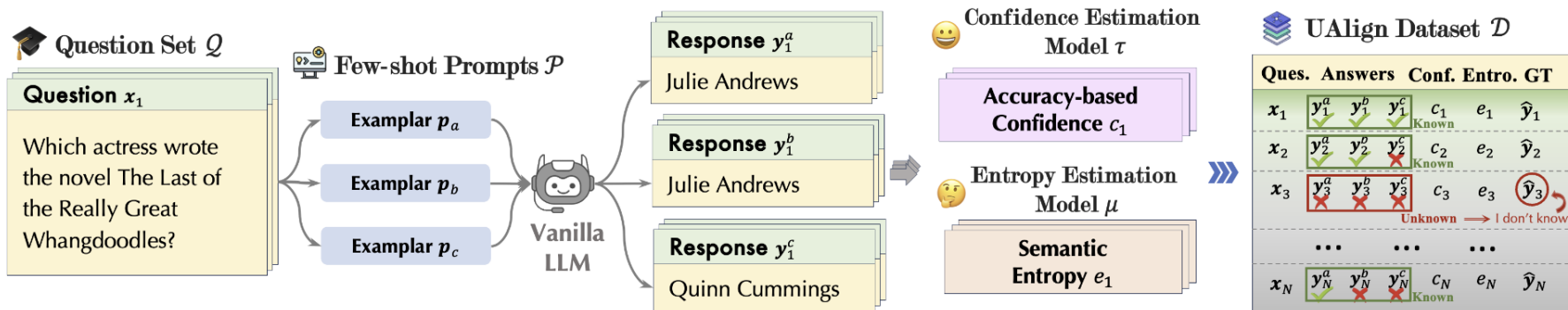  - ❏ Random Responses to Ambiguous Queries
- ❏ **Unanswerable Query Mitigation**
  - ❏ Refusal Fine-tuning
  - ❏ Uncertainty-based Reinforcement Learning
  - ❏ Self-alignment
- ❏ Ambiguous Query Clarification
  - ❏ In-Context Learning
  - ❏ Reinforcement Learning
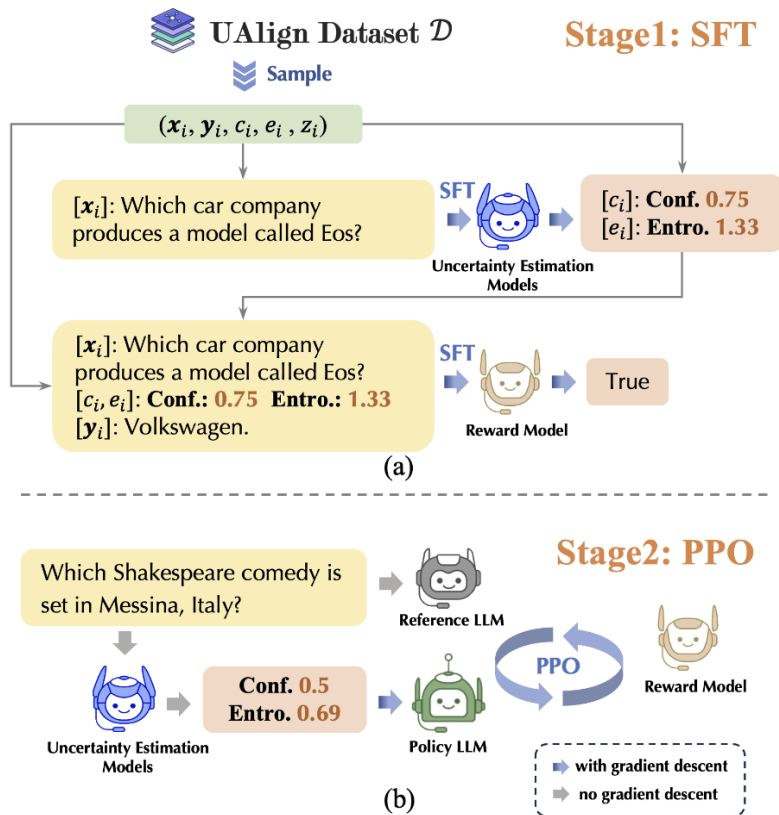  - ❏ Preference Optimization

# Uncertainty-based Alignment (UAlign)



## UAlign Data Construction

❑ Response Sampling

❑ Uncertainty Measurement: Accuracy-based Confidence & Semantic Entropy

*Xue et al., "UALIGN: Leveraging Uncertainty Estimations for Factuality Alignment on Large Language Models" (ACL '25)*

# Uncertainty-based Alignment (UAlign)



(a)

(b)

## UAlign Training Framework

- ❑ **Supervised Fine-tuning** to train uncertainty estimation model

- ❑ **Reward Model Training** to train a reward model as a binary evaluator to determine if a generated answer is correctly conditioned on the question, confidence, and entropy.

- ❑ **PPO Alignment** to optimize the LLM's factual expressions to a question with the uncertainty measurements.

*Xue et al., "UALIGN: Leveraging Uncertainty Estimations for Factuality Alignment on Large Language Models" (ACL '25)*

# Outline

❏ **Why Proactive Query Management?**

    ❏ Overconfidence on Unanswerable Queries

    ❏ Random Responses to Ambiguous Queries

❏ **Unanswerable Query Mitigation**

    ❏ Refusal Fine-tuning

    ❏ Uncertainty-based Reinforcement Learning

    ❏ Self-alignment

❏ **Ambiguous Query Clarification**

    ❏ In-Context Learning

    ❏ Reinforcement Learning

    ❏ Preference Optimization

# Issues of Refusal

**Q:** What animal can be found at the top of the men's Wimbledon trophy?

**A:** The answer is unknown.
**Unknown Question Detection**

**A:** The question is incorrect.
**Unknown Question Classification**

⚠️ Not User-friendly; Fail to Meet User Information Needs

**How to properly respond to unknown questions?**

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Issues of Refusal

**Q:** What animal can be found at the top of the men's Wimbledon trophy?

**A:** The answer is unknown.

**Unknown Question Detection**

**A:** The question is incorrect.

**Unknown Question Classification**

⚠️ Not User-friendly; Fail to Meet User Information Needs

**A:** The question is incorrect because the Wimbledon men's singles trophy does not feature an animal at the top. Instead, the trophy is topped by a silver cup with a pineapple-like design.

**Desired response format:**

❏ Identify the type of unknown question

❏ Provide justifications or explanations

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Workflow of Self-Align

**Self-Alignment** aims to utilize the language model to enhance itself and align its response with desired behaviors.

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Initialization

**Incorrect Questions**

Known Question: When did Neil Armstrong set foot on the Moon?

Seed Data  Unknown Question: When did Neil Armstrong set foot on Mars?

**Incomplete Questions**

Known Question: Priya said yes to Jay when he proposed. Did she say yes?

Seed Data  Unknown Question: Jay proposed to Priya yesterday. Did she say yes?

**Ambiguous Questions**

Known Question: Everyone is ready to eat the goat. Is the goat cooked?

Seed Data  Known-Unknown Question: The goat is ready to eat. Is the goat cooked?

**Futuristic Questions**

Known Question: What was the biggest sporting event in 2020?

Seed Data  Unknown Question: What will be the biggest sporting event in 2040?

**Base LLM**

**Known Questions**

**Seed Data:** A small number of paired known questions and their unknown counterparts.

**Base LLM:** A tunable base LLM to be improved.

**Known QA Data:** A large number of known question-answer pairs.

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Stage 1: Guided Question Rewriting



**Incorrect Questions**

Seed Data — **Known Question:** When did Neil Armstrong set foot on the Moon?
**Unknown Question:** When did Neil Armstrong set foot on Mars?

**Incomplete Questions**

Seed Data — **Known Question:** Priya said yes to Jay when he proposed. Did she say yes?
**Unknown Question:** Jay proposed to Priya yesterday. Did she say yes?

**Ambiguous Questions**

Seed Data — **Known Question:** Everyone is ready to eat the goat. Is the goat cooked?
**Known-Unknown Question:** The goat is ready to eat. Is the goat cooked?

**Futuristic Questions**

Seed Data — **Known Question:** What was the biggest sporting event in 2020?
**Unknown Question:** What will be the biggest sporting event in 2040?

Known Questions + Base LLM → Incorrect Questions / Incomplete Questions / Ambiguous Questions / Futuristic Questions
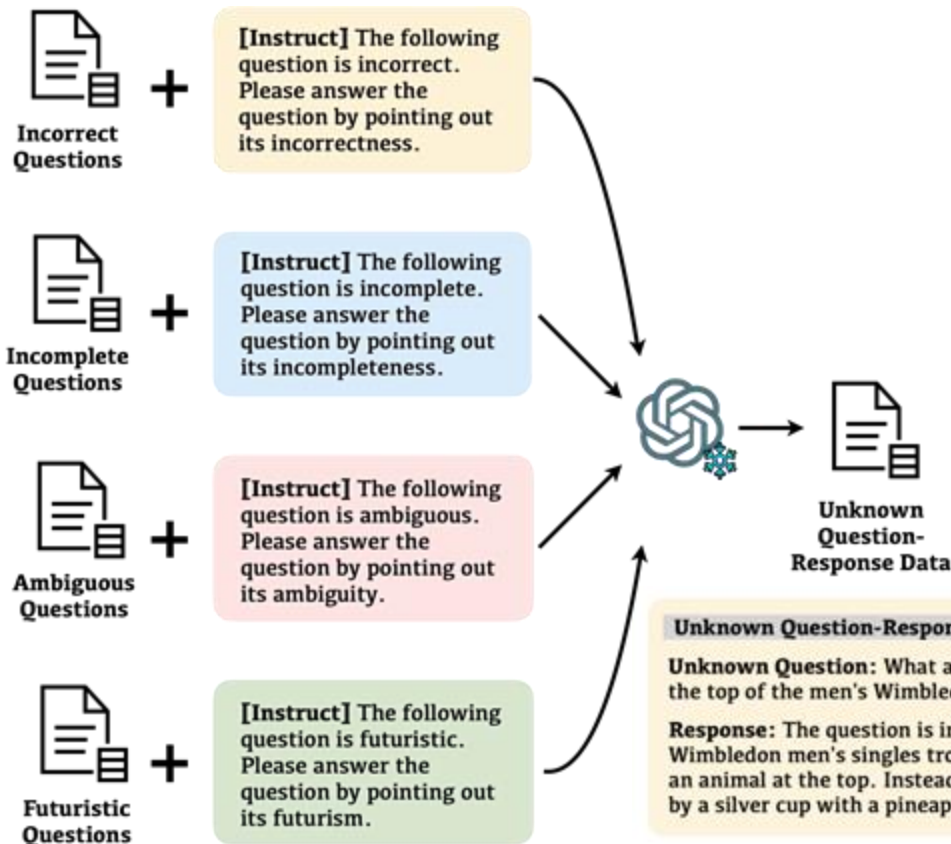
$$\mathcal{D}_{\text{uq}}^{c} = \{\mathcal{M}(z_{qr}^{c}; \mathcal{D}_{\text{seed}}^{c}; q)\}_{q \in \mathcal{D}_{\text{kq}}}$$

- ❑ **Seed Data**
  → demonstrations

- ❑ **Known Questions**
  → source text

- ❑ **Unknown Questions**
  → target text

- ❑ **Base LLM**
  → question rewriter

25    *Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Stage 2: Conditioned Response Generation



**Incorrect Questions** +

[Instruct] The following question is incorrect. Please answer the question by pointing out its incorrectness.

**Incomplete Questions** +

[Instruct] The following question is incomplete. Please answer the question by pointing out its incompleteness.

**Ambiguous Questions** +

[Instruct] The following question is ambiguous. Please answer the question by pointing out its ambiguity.

**Futuristic Questions** +

[Instruct] The following question is futuristic. Please answer the question by pointing out its futurism.
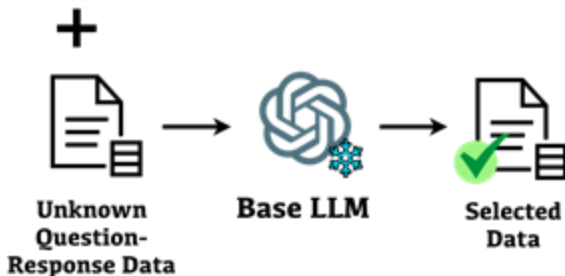
Unknown Question-Response Data

$$\mathcal{D}_{\text{unk}}^c = \{(p_i, \mathcal{M}(z_{rg}^c; p_i, q_i))\}_{p_i \in \mathcal{D}_{\text{uq}}^c, q_i \in \mathcal{D}_{\text{kq}}^c}$$

**Unknown Question-Response Example**

**Unknown Question:** What animal can be found at the top of the men's Wimbledon trophy?

**Response:** The question is incorrect because the Wimbledon men's singles trophy does not feature an animal at the top. Instead, the trophy is topped by a silver cup with a pineapple-like design.

## Instructions

❏ **Response Format**

    ❏ Unknown Question Type

    ❏ Explanation

❏ **Known Question as Reference**

    ❏ Analyze the unanswerability

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Stage 3: Disparity-driven Self-Curation

[Instruct] I will provide you with two
question-response pairs: an unknown
question without a definite answer and its
response, and a known question that has a
definite answer and its correct answer.
Please score the disparity between these
two pairs from 0 to 100:
<Unknown Question-Response Pair>
<Known Question-Response Pair>

**+**

Unknown Question-Response Data → Base LLM → ✅ Selected Data

**Unknown Question-Response Example**

**Unknown Question:** What animal can be found at the top of the men's Wimbledon trophy?

**Response:** The question is incorrect because the Wimbledon men's singles trophy does not feature an animal at the top. Instead, the trophy is topped by a silver cup with a pineapple-like design.
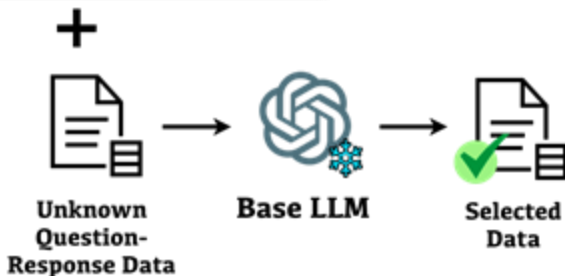
$$s_i = \mathcal{M}(z_{sc}; (q_i, a_i); (p_i, r_i))$$

**Why not directly scoring the quality?**

➢ The base model itself fails to identify whether the question has a definitive answer.

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Stage 3: Disparity-driven Self-Curation

[Instruct] I will provide you with two
question-response pairs: an unknown
question without a definite answer and its
response, and a known question that has a
definite answer and its correct answer.
Please score the disparity between these
two pairs from 0 to 100:
`<Unknown Question-Response Pair>`
`<Known Question-Response Pair>`

**+**



**Unknown Question-**
**Response Data** → **Base LLM** → **Selected**
**Data**

**Unknown Question-Response Example**

**Unknown Question:** What animal can be found at
the top of the men's Wimbledon trophy?

**Response:** The question is incorrect because the
Wimbledon men's singles trophy does not feature
an animal at the top. Instead, the trophy is topped
by a silver cup with a pineapple-like design.

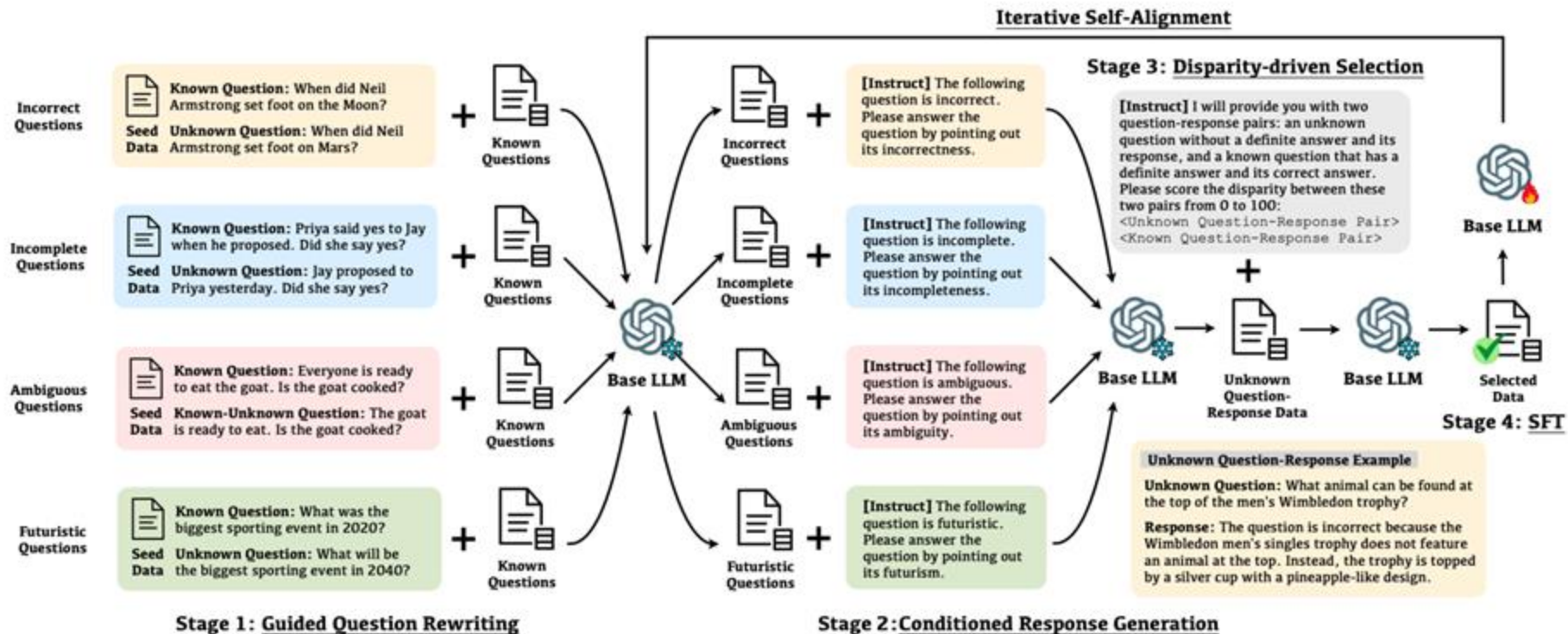$$s_i = \mathcal{M}(z_{sc}; (q_i, a_i); (p_i, r_i))$$

**Why not directly scoring the quality?**

➢ The base model itself fails to identify whether
the question has a definitive answer.

**Why scoring disparity?**

➢ The conditional generation capability of LLMs
ensure the semantic quality of the generated
question-response pair.

➢ Low disparity score can filter out those low-
quality pairs that fail to differentiate from
their original known QA counterparts.

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Stage 4: Supervised Fine-tuning & Iterative Self-alignment

*Deng et al,. "Don't Just Say 'I don't know'! Self-aligning Large Language Models for Responding to Unknown Questions with Explanations" (EMNLP '24)*

# Outline

❏ **Why Proactive Query Management?**

    ❏ Overconfidence on Unanswerable Queries

    ❏ Random Responses to Ambiguous Queries

❏ **Unanswerable Query Mitigation**

    ❏ Refusal Fine-tuning

    ❏ Uncertainty-based Reinforcement Learning

    ❏ Self-alignment

❏ **Ambiguous Query Clarification**

    ❏ In-Context Learning

    ❏ Reinforcement Learning

    ❏ Preference Optimization

# Outline

- ❏ **Why Proactive Query Management?**
  - ❏ Overconfidence on Unanswerable Queries
  - ❏ Random Responses to Ambiguous Queries
- ❏ **Unanswerable Query Mitigation**
  - ❏ Refusal Fine-tuning
  - ❏ Uncertainty-based Reinforcement Learning
  - ❏ Self-alignment
- ❏ **Ambiguous Query Clarification**
  - ❏ In-Context Learning
  - ❏ Reinforcement Learning
  - ❏ Preference Optimization

# Proactive Chain-of-Thought (ProCoT)

❏ **Standard Prompting**

    ❏ Input: Task Background & Conversation History

    ❏ Output: Response

$$p(r|\mathcal{D}, \mathcal{C})$$



(1) Clarification Dialogues: **Abg-CoQA**

**Task Background**: The grounded document is "Angie …. She made a drawing of her mother. Her mother found a large red book. Then they went to the Mystery section. Angie sat in a blue chair. She drew a picture of her brother. Her mother found the book. It was a green book. …"
**Conversation History**: ["User": "What did she draw?", "System": "Her mother", "User": "What did her mother find?", "System": "The book", "User": "What color was it?"]

(1a) Standard

**Prompt**: Given the task background and the conversation history, please generate the response:
**Response**: Green ❌

*Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Proactive Chain-of-Thought (ProCoT)

❑ **Standard Prompting**

  ❑ Input: Task Background & Conversation History

  ❑ Output: Response $p(r|\mathcal{D}, \mathcal{C})$

❑ **Proactive Prompting**

  ❑ Input: + Action Space

  ❑ Output: + Action $p(a, r|\mathcal{D}, \mathcal{C}, \mathcal{A})$

---

**(1) Clarification Dialogues: Abg-CoQA**

**Task Background**: The grounded document is "Angie …. She made a drawing of her mother. Her mother found a large red book. Then they went to the Mystery section. Angie sat in a blue chair. She drew a picture of her brother. Her mother found the book. It was a green book. …"

**Conversation History**: ["User": "What did she draw?", "System": "Her mother", "User": "What did her mother find?", "System": "The book", "User": "What color was it?"]

**(1a) Standard**

**Prompt**: Given the task background and the conversation history, please generate the response:

**Response**: Green ❌

**(1b) Proactive**

**Act**: ["Directly Answer", "Ask a Clarification Question"]

**Prompt**: Given the task background and conversation history, please use appropriate actions to generate the response:

**Response**: Ask a clarification question: Could you provide more information? ❌

---

33    *Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Proactive Chain-of-Thought (ProCoT)

❑ **Standard Prompting**

  ❑ Input: Task Background & Conversation History

  ❑ Output: Response $p(r|\mathcal{D}, \mathcal{C})$

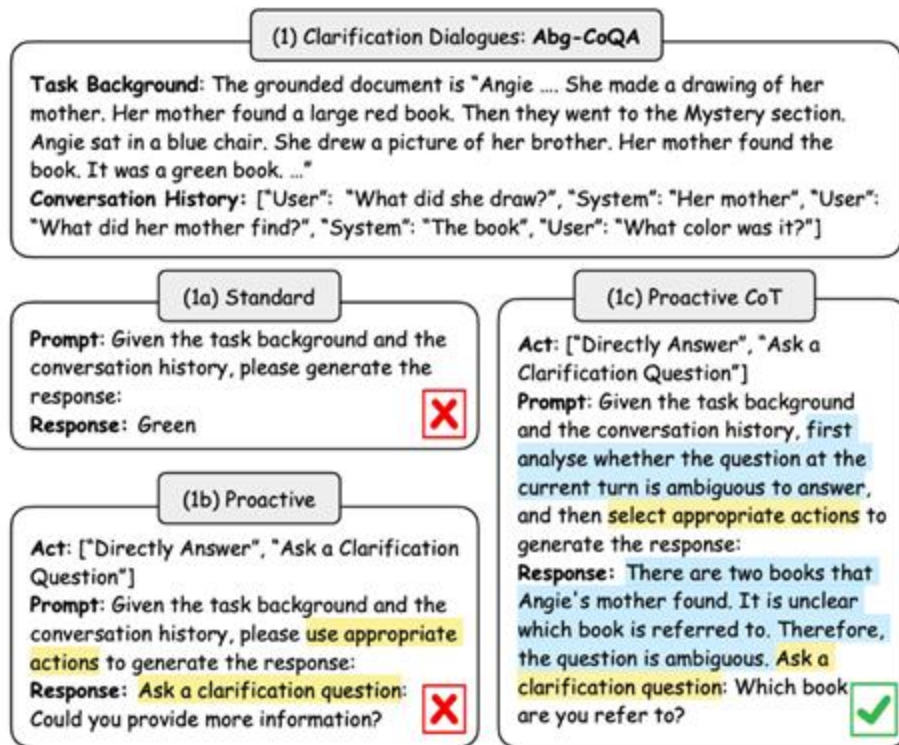❑ **Proactive Prompting**

  ❑ Input: + Action Space

  ❑ Output: + Action $p(a, r|\mathcal{D}, \mathcal{C}, \mathcal{A})$

❑ **Proactive Chain-of-Thought Prompting**

  ❑ Output: + Reasoning Chain

$$p(t, a, r|\mathcal{D}, \mathcal{C}, \mathcal{A})$$



**(1) Clarification Dialogues: Abg-CoQA**

**Task Background**: The grounded document is "Angie …. She made a drawing of her mother. Her mother found a large red book. Then they went to the Mystery section. Angie sat in a blue chair. She drew a picture of her brother. Her mother found the book. It was a green book. …"
**Conversation History**: ["User": "What did she draw?", "System": "Her mother", "User": "What did her mother find?", "System": "The book", "User": "What color was it?"]

**(1a) Standard**
**Prompt**: Given the task background and the conversation history, please generate the response:
**Response**: Green ❌

**(1b) Proactive**
**Act**: ["Directly Answer", "Ask a Clarification Question"]
**Prompt**: Given the task background and the conversation history, please use appropriate actions to generate the response:
**Response**: Ask a clarification question: Could you provide more information? ❌

**(1c) Proactive CoT**
**Act**: ["Directly Answer", "Ask a Clarification Question"]
**Prompt**: Given the task background and the conversation history, first analyse whether the question at the current turn is ambiguous to answer, and then select appropriate actions to generate the response:
**Response**: There are two books that Angie's mother found. It is unclear which book is referred to. Therefore, the question is ambiguous. Ask a clarification question: Which book are you refer to? ✅

34    *Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Proactive Chain-of-Thought (ProCoT)

| | | | Abg-CoQA | | | PACIFIC | | |
| | | | CNP | CQG | | CNP | CQG | |
| Method | Shot | Prompt | F1 | BLEU-1 | Help. | F1 | ROUGE-2 | Help. |
|---|---|---|---|---|---|---|---|---|
| Baseline | - | - | 22.1 | 36.5 | 30.0 | 79.0 | 69.2 | 38.2 |
| SOTA | - | - | 23.6 | 38.2 | 56.0 | 86.9 | 90.7 | 80.1 |
| Vicuna-13B | 0 | Standard | - | 11.3 | 0.0 | - | 1.2 | 0.0 |
| | 1 | Standard | - | 11.4 | 0.0 | - | 2.5 | 0.0 |
| | 0 | Proactive | 4.1 | 13.2 | 0.0 | 2.3 | 2.3 | 0.0 |
| | 1 | Proactive | 12.1 | 13.2 | 4.5 | 0.0 | 3.3 | 0.0 |
| | 0 | ProCoT | 1.4 | 21.3 | 9.1 | 9.7 | 3.8 | 10.5 |
| | 1 | ProCoT | **18.3** | **23.7** | **22.7** | **27.0** | **41.3** | **33.1** |
| ChatGPT | 0 | Standard | - | 12.1 | 0.0 | - | 2.2 | 0.0 |
| | 1 | Standard | - | 12.3 | 0.0 | - | 2.0 | 0.0 |
| | 0 | Proactive | 22.0 | 13.7 | 17.6 | 19.4 | 2.9 | 0.0 |
| | 1 | Proactive | 20.4 | **23.4** | 23.5 | 17.7 | 14.0 | 12.5 |
| | 0 | ProCoT | 23.8 | 21.6 | 32.4 | **28.0** | **21.5** | 26.7 |
| | 1 | ProCoT | **27.9** | 18.4 | **45.9** | 27.7 | 16.2 | **35.8** |

**LLMs barely ask clarification questions, even when the user query is ambiguous.**

*Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Proactive Chain-of-Thought (ProCoT)

**Open-domain**    **Finance**

| Method | Shot | Prompt | Abg-CoQA | | | PACIFIC | | |
|---|---|---|---|---|---|---|---|---|
| | | | CNP | CQG | | CNP | CQG | |
| | | | F1 | BLEU-1 | Help. | F1 | ROUGE-2 | Help. |
| Baseline | - | - | 22.1 | 36.5 | 30.0 | 79.0 | 69.2 | 38.2 |
| SOTA | - | - | 23.6 | 38.2 | 56.0 | 86.9 | 90.7 | 80.1 |
| Vicuna-13B | 0 | Standard | - | 11.3 | 0.0 | - | 1.2 | 0.0 |
| | 1 | Standard | - | 11.4 | 0.0 | - | 2.5 | 0.0 |
| | 0 | Proactive | 4.1 | 13.2 | 0.0 | 2.3 | 2.3 | 0.0 |
| | 1 | Proactive | 12.1 | 13.2 | 4.5 | 0.0 | 3.3 | 0.0 |
| | 0 | ProCoT | 1.4 | 21.3 | 9.1 | 9.7 | 3.8 | 10.5 |
| | 1 | ProCoT | **18.3** | **23.7** | **22.7** | 27.0 | **41.3** | **33.1** |
| ChatGPT | 0 | Standard | - | 12.1 | 0.0 | - | 2.2 | 0.0 |
| | 1 | Standard | - | 12.3 | 0.0 | - | 2.0 | 0.0 |
| | 0 | Proactive | 22.0 | 13.7 | 17.6 | 19.4 | 2.9 | 0.0 |
| | 1 | Proactive | 20.4 | **23.4** | 23.5 | 17.7 | 14.0 | 12.5 |
| | 0 | ProCoT | 23.8 | 21.6 | 32.4 | **28.0** | **21.5** | 26.7 |
| | 1 | ProCoT | **27.9** | 18.4 | **45.9** | 27.7 | 16.2 | **35.8** |

> LLMs barely ask clarification questions, even when the user query is ambiguous.

> **ProCoT largely overcomes this issue in open-domain, but the performance is still unsatisfactory in domain-specific applications.**

*Deng et al., "Prompting and Evaluating Large Language Models for Proactive Dialogues: Clarification, Target-guided, and Non-collaboration" (EMNLP '23 Findings)*

# Outline

❏ **Why Proactive Query Management?**

    ❏ Overconfidence on Unanswerable Queries

    ❏ Random Responses to Ambiguous Queries

❏ **Unanswerable Query Mitigation**

    ❏ Refusal Fine-tuning

    ❏ Uncertainty-based Reinforcement Learning

    ❏ Self-alignment

❏ **Ambiguous Query Clarification**

    ❏ In-Context Learning

    ❏ Reinforcement Learning

    ❏ Preference Optimization

# Limitations of In-context Learning Approaches



❏ Fail to optimize the long-term goal of the conversation.

❏ Not learnable.

❏ Limited by the strategy planning capability of LLMs.

➢ **Reinforcement Learning with Goal-oriented AI Feedback**

# Reinforcement Learning

❏ Formulate the proactive conversation as a **Markov Decision Process (MDP).**

❏ The objective is to learn a policy π maximizing the expected cumulative rewards over the observed dialogue episodes as:

$$\pi^* = \arg\max_{\pi \in \Pi} \left[ \sum_{t=0}^{T} \mathcal{R}(s_t) \right]$$    ***Reward Function***

$$= \arg\max_{\pi \in \Pi} \left[ \sum_{t=0}^{T} \mathcal{R}(\mathcal{T}(s_{t-1}, a_t)) \right]$$    ***State Transition***

$$= \arg\max_{\pi \in \Pi} \left[ \sum_{t=0}^{T} \mathcal{R}(\mathcal{T}(s_{t-1}, \pi(s_{t-1}))) \right]$$    ***Policy Network***

***How to enable the policy learning with LLMs?***

39

*Deng et al., "Plug-and-Play Policy Planner for Large Language Model Powered Dialogue Agents" (ICLR '24)*

# Policy Network – Plug-and-Play Dialogue Policy Planner

❏ A **tunable language model plug-in** for dialogue strategy learning.

$$a_t = \pi(s_{t-1})$$

❏ Conduct **Supervised Fine-Tuning** on available human-annotated corpus.

$$\mathcal{L}_c = -\frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} \frac{1}{T_d} \sum_{t=1}^{T_d} a_t \log y_t$$



40

*Deng et al., "Plug-and-Play Policy Planner for Large Language Model Powered Dialogue Agents" (ICLR '24)*

# Reward Function – Learning from AI Feedback

❑ An LLM as the reward model to assess the goal achievement and provide **goal-oriented AI feedback**.

$$\mathcal{R}(s_t) = \frac{1}{l} \sum_{i=1}^{l} \mathcal{M}_r(\mathbf{LLM}_{\mathrm{rwd}}(p_{\mathrm{rwd}}; s_t; \tau))$$

❑ Employ **Reinforcement Learning** to further tune the policy model.

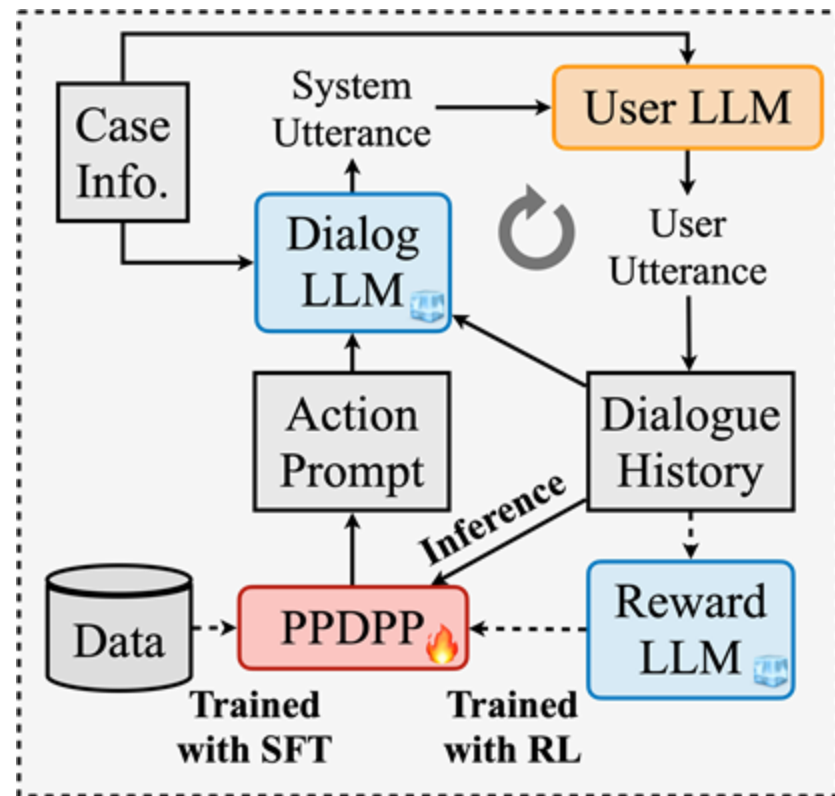$$\theta \leftarrow \theta - \alpha \nabla \log \pi_\theta(a_t|s_t) R_t$$
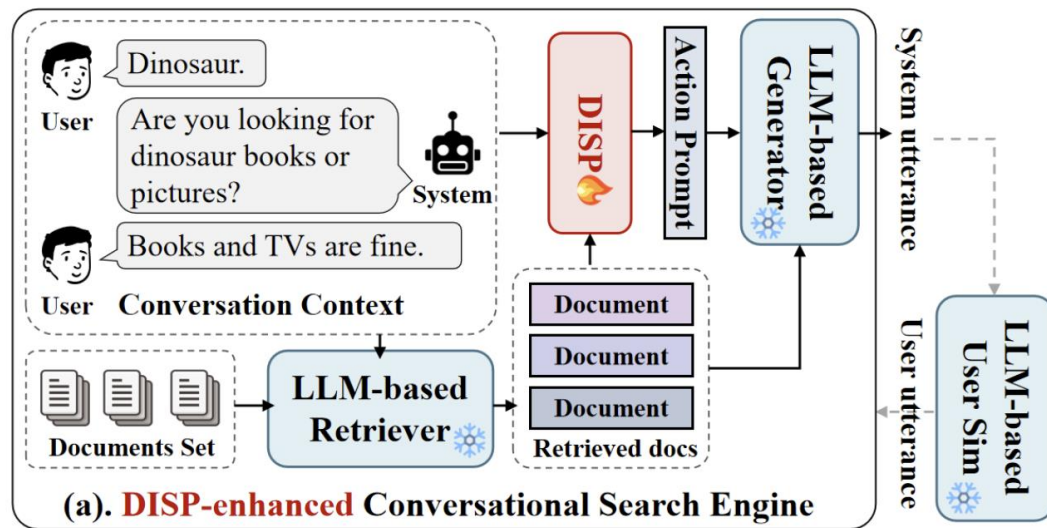
⚠️ *Interacting with real user is costly!*



41

*Deng et al., "Plug-and-Play Policy Planner for Large Language Model Powered Dialogue Agents" (ICLR '24)*

# State Transition – Multi-agent Simulation

❏ An LLM to simulate the user with user profiles.

❏ Employ **Multi-agent Simulation** to collect dynamic interaction data.

$$u_t^{sys} = \mathbf{LLM}_{sys}(p_{sys}; \mathcal{M}_a(a_t); s_{t-1})$$
$$u_t^{usr} = \mathbf{LLM}_{usr}(p_{usr}; s_{t-1}; u_t^{sys})$$
$$s_t = \mathcal{T}(s_{t-1}, a_t)$$
$$= \{s_{t-1}; u_t^{sys}, u_t^{usr}\}$$

*Deng et al., "Plug-and-Play Policy Planner for Large Language Model Powered Dialogue Agents" (ICLR '24)*

# RL for Asking Clarification Questions – STYLE



(a). **DISP-enhanced** Conversational Search Engine
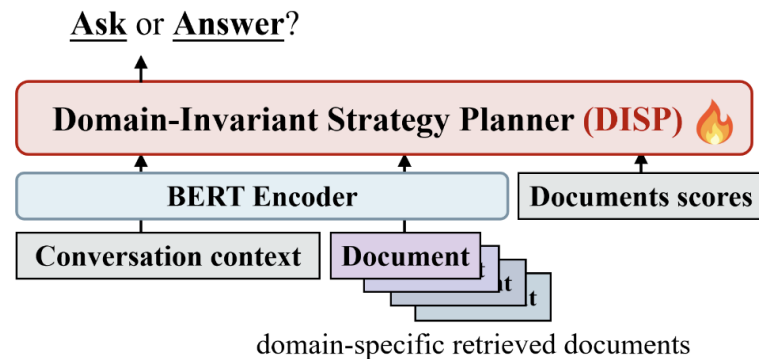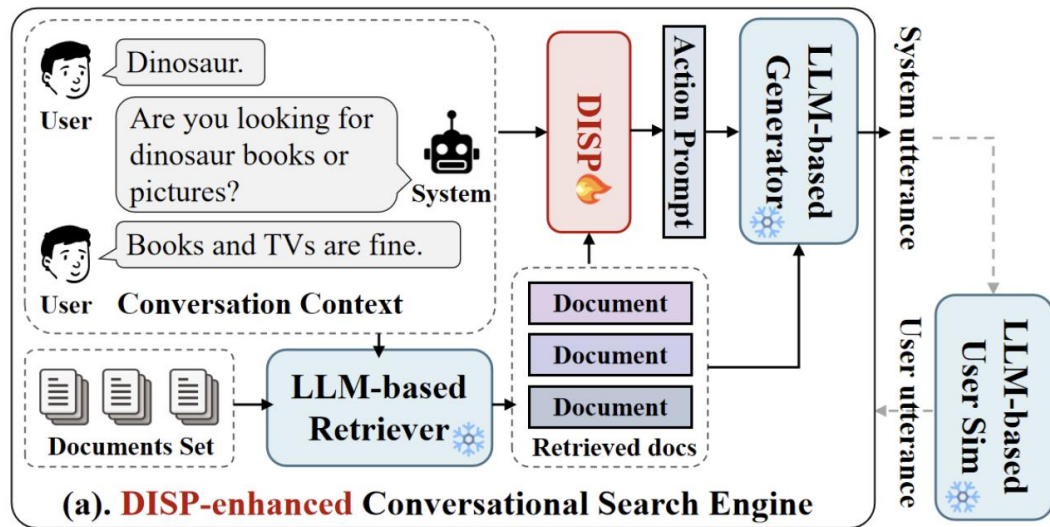
STYLE features rapid transfer to previously unseen domains via tailored strategies.

❑ Domain-Invariant Strategy Planner (DISP)

❑ Multi-Domain Training (MDT) Paradigm

*Chen et al., "STYLE: Improving Domain Transferability of Asking Clarification Questions in Large Language Model Powered Conversational Agents" (ACL '24 Findings)*

# RL for Asking Clarification Questions – STYLE



(a). **DISP-enhanced** Conversational Search Engine

DISP is a policy module that determines when to ask questions. It extract domain-invariant information, mitigating the mismatch in the distribution of domain-specific representations and ensuring robustness across domains.

*Chen et al., "STYLE: Improving Domain Transferability of Asking Clarification Questions in Large Language Model Powered Conversational Agents" (ACL '24 Findings)*

# RL for Asking Clarification Questions – STYLE



(b). Multi-Domain Training (MDT)

$$y_t = \mathbb{E}_{s_{t+1}} \left[ r_t + \gamma \max_{a \in \mathcal{A}} Q^*(s_{t+1}, a_{t+1}) \big| s_t, a_t \right]$$

MDT encourages the domain transferability of DISP by training it across multiple diverse domains. This is inspired by the population-based training, which suggests that the generalization of a collaborative agent to held-out populations can be improved by training larger and more diverse populations.

*Chen et al., "STYLE: Improving Domain Transferability of Asking Clarification Questions in Large Language Model Powered Conversational Agents" (ACL '24 Findings)*

# Outline

- Why Proactive Query Management?
  - Overconfidence on Unanswerable Queries
  - Random Responses to Ambiguous Queries
- Unanswerable Query Mitigation
  - Refusal Fine-tuning
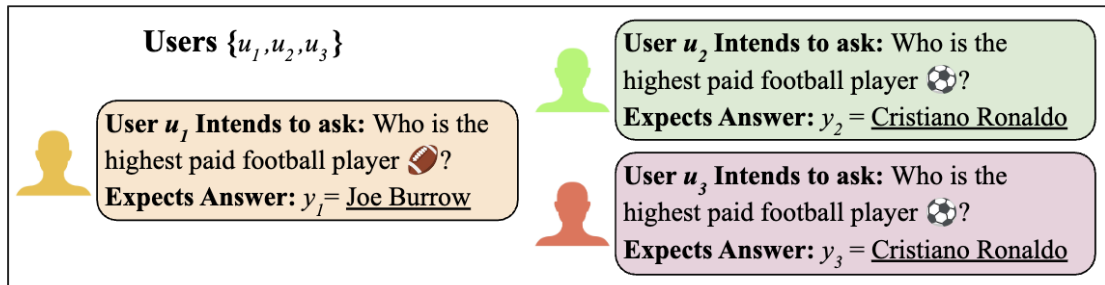  - Uncertainty-based Reinforcement Learning
  - Self-alignment
- **Ambiguous Query Clarification**
  - In-Context Learning
  - Reinforcement Learning
  - Preference Optimization

# Modeling Future Conversation Turns

**[Turn 1]** User's Input Query: $x$

> Who is the highest paid football player?

Users $\{u_1, u_2, u_3\}$

**User $u_1$ Intends to ask:** Who is the highest paid football player 🏈?
**Expects Answer:** $y_1$ = <u>Joe Burrow</u>

**User $u_2$ Intends to ask:** Who is the highest paid football player ⚽?
**Expects Answer:** $y_2$ = <u>Cristiano Ronaldo</u>

**User $u_3$ Intends to ask:** Who is the highest paid football player ⚽?
**Expects Answer:** $y_3$ = <u>Cristiano Ronaldo</u>

47

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*

# Modeling Future Conversation Turns

**[Turn 1] User's Input Query:** $x$

> Who is the highest paid football player?

**[Turn 2] LLM ($M$) Predicts Single-Turn Responses:** $M(x) = r_{init}$

**Users $\{u_1, u_2, u_3\}$**

**User $u_1$ Intends to ask:** Who is the highest paid football player 🏈?
**Expects Answer:** $y_1$ = Joe Burrow

**User $u_2$ Intends to ask:** Who is the highest paid football player ⚽?
**Expects Answer:** $y_2$ = Cristiano Ronaldo

**User $u_3$ Intends to ask:** Who is the highest paid football player ⚽?
**Expects Answer:** $y_3$ = Cristiano Ronaldo

**[A] Direct-Answer (⚽):**
$r_{init}$ = As of 2024, the highest-paid football players are Cristiano Ronaldo…
$\Phi(r_{init})$ = *False*

**[B] Direct-Answer (🏈):**
$r_{init}$ = The highest-paid football player in the NFL for 2024 is Joe Burrow…
$\Phi(r_{init})$ = *False*

**[C] Clarifying Question (🏈-or-⚽?):**
$r_{init}$ = Are you asking about American Football or Soccer?
$\Phi(r_{init})$ = *True*

**[D] Clarifying Question (📅?):**
$r_{init}$ = Are you asking about a specific time period?
$\Phi(r_{init})$ = *True*

**Single-Turn Preferences**

| | [A] | [B] | [C] | [D] |
|---|---|---|---|---|
| User 1 | ❌ | ✅ | ❌ | ❌ |
| Users 2, 3 | ✅ | ❌ | ❌ | ❌ |

**Majority Best Response: [A]**

48

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*
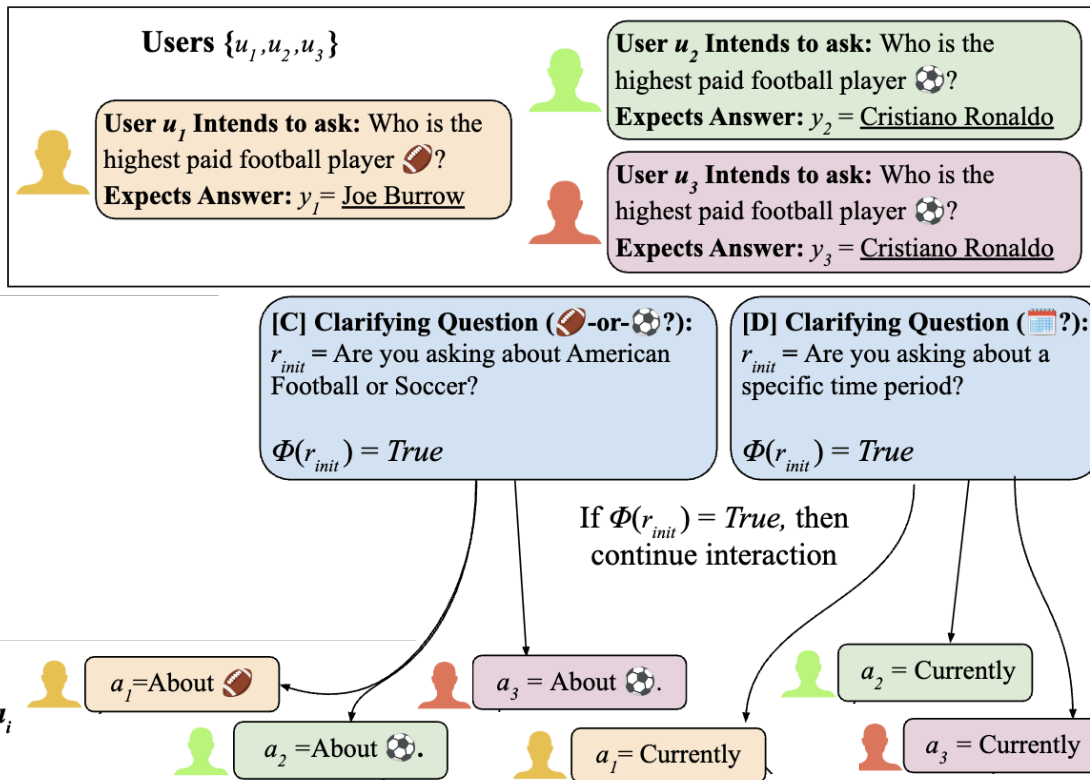
# Modeling Future Conversation Turns



**[Turn 1] User's Input Query:** $x$

> Who is the highest paid football player?

**[Turn 2] LLM ($M$) Predicts Single-Turn Responses:** $M(x) = r_{init}$

**Users** $\{u_1, u_2, u_3\}$

**User $u_1$ Intends to ask:** Who is the highest paid football player 🏈?
**Expects Answer:** $y_1$ = <u>Joe Burrow</u>

**User $u_2$ Intends to ask:** Who is the highest paid football player ⚽?
**Expects Answer:** $y_2$ = <u>Cristiano Ronaldo</u>

**User $u_3$ Intends to ask:** Who is the highest paid football player ⚽?
**Expects Answer:** $y_3$ = <u>Cristiano Ronaldo</u>

**[C] Clarifying Question (🏈-or-⚽?):**
$r_{init}$ = Are you asking about American Football or Soccer?

$\Phi(r_{init}) = True$

**[D] Clarifying Question (📅?):**
$r_{init}$ = Are you asking about a specific time period?

$\Phi(r_{init}) = True$

If $\Phi(r_{init}) = True$, then continue interaction

**[Turn 3] Users Respond with Clarifying Answers:** $\psi(x, y_i, r_{init}) = a_i$

$a_1$=About 🏈

$a_2$=About ⚽.

$a_3$ = About ⚽.

$a_1$= Currently

$a_2$ = Currently

$a_3$ = Currently

49

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*

# Modeling Future Conversation Turns



**[Turn 1] User's Input Query: $x$**

Who is the highest paid football player?

**[Turn 2] LLM ($M$) Predicts Single-Turn Responses: $M(x) = r_{init}$**

**Users $\{u_1, u_2, u_3\}$**

**User $u_2$ Intends to ask:** Who is the highest paid football player 🏀? **Expects Answer:** $y_2$ = Cristiano Ronaldo

**User $u_1$ Intends to ask:** Who is the highest paid football player 🏈? **Expects Answer:** $y_1$ = Joe Burrow

**User $u_3$ Intends to ask:** Who is the highest paid football player 🏀? **Expects Answer:** $y_3$ = Cristiano Ronaldo

**[A] Direct-Answer (⚽):** $r_{init}$ = As of 2024, the highest-paid football players are Cristiano Ronaldo… $\Phi(r_{init}) = False$

**[B] Direct-Answer (🏈):** $r_{init}$ = The highest-paid football player in the NFL for 2024 is Joe Burrow… $\Phi(r_{init}) = False$

**[C] Clarifying Question (🏈-or-⚽?):** $r_{init}$ = Are you asking about American Football or Soccer? $\Phi(r_{init}) = True$

**[D] Clarifying Question (📅?):** $r_{init}$ = Are you asking about a specific time period? $\Phi(r_{init}) = True$

**Single-Turn Preferences**

|   | [A] | [B] | [C] | [D] |
|---|---|---|---|---|
| User 1 | ✗ | ✓ | ✗ | ✗ |
| Users 2, 3 | ✓ | ✗ | ✗ | ✗ |

**Majority Best Response: [A]**

If $\Phi(r_{init}) = True$, then continue interaction

**[Turn 3] Users Respond with Clarifying Answers: $\psi(x, y_i, r_{init}) = a_i$**

$a_1$=About 🏈

$a_3$ = About ⚽.

$a_2$ = Currently

$a_2$ =About ⚽.

$a_1$ = Currently

$a_3$ = Currently

**[Turn 4] LLMs Predict Next Response: $M(x, r_{init}, a_i) = r^i_{next}$**

$r^1_{next}$ = Joe Burrow…

$r^2_{next} = r^3_{next}$ = Cristiano Ronaldo…

$r^1_{next} = r^2_{next} = r^3_{next}$ = Cristiano Ronaldo…

**Double-Turn Preferences**

|   | [A] | [B] | [C] | [D] |
|---|---|---|---|---|
| User 1 | ✗ | ✓ | ✓ | ✗ |
| Users 2, 3 | ✓ | ✗ | ✓ | ✓ |

**Majority Best Response: [C]**

**Evaluation Metrics**

|   | [A] | [B] | [C] | [D] |
|---|---|---|---|---|
| Efficiency (# of $M$ Turns) | 1 | 1 | 2 | 2 |
| F1 ($R$, $\{y_1, y_2, y_3\}$) | 0.8 | 0.5 | 1.0 | $0.\overline{6}$ |

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*

# Modeling Future Conversation Turns



**Supervised Fine-Tuning Data**

- **Clarify:** $(x) \rightarrow q$
- **Direct Ans:** $(x) \rightarrow y$
- **Ans-After-Clarify:** $(x, q, a) \rightarrow y$
- **User Simulator** $(x, q, y) \rightarrow a$

**Responses for Preference Learning**

- **Clarify Responses** $(x) \rightarrow q$ From Clarify SFT Model
- **Answer Responses** $(x) \rightarrow y$ From Direct Ans SFT Model

Base LLM → Clarify SFT / Direct Ans SFT / Clarify-or-Direct Ans SFT

Double turn scoring ▢ vs. ▢ → Clarify DPO (Likelihood, Match)

Single turn scoring ▢ vs. ▢ → Clarify DPO (RM)

Double turn scoring ▢ vs. ▢ / ▢ vs. ▢ → Clarify-or-Direct Ans DPO

❑ **Clarify SFT**: The base LLM is fine-tuned to ask clarifying question to the input query on the SFT data.

❑ **Direct-Ans SFT**: The base LLM is fine-tuned on QA data.

❑ **Clarify-or-Direct Ans SFT**: The base LLM is fine-tuned on the union of all data used to train Clarify SFT and Direct-Ans SFT models.

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*

# Modeling Future Conversation Turns



**Supervised Fine-Tuning Data**

- **Clarify:** $(x) \rightarrow q$
- **Direct Ans:** $(x) \rightarrow y$
- **Ans-After-Clarify:** $(x, q, a) \rightarrow y$
- **User Simulator** $(x, q, y) \rightarrow a$

**Responses for Preference Learning**

- **Clarify Responses** $(x) \rightarrow q$ From Clarify SFT Model
- **Answer Responses** $(x) \rightarrow y$ From Direct Ans SFT Model

Base LLM → Clarify SFT → Double turn scoring [ vs. ] → Clarify DPO (Likelihood, Match)

Single turn scoring [ vs. ] → Clarify DPO (RM)

Base LLM → Direct Ans SFT

Base LLM → Clarify-or-Direct Ans SFT → Double turn scoring [ vs. ] [ vs. ] → Clarify-or-Direct Ans DPO

❑ **Clarify DPO**: The Clarify SFT model is further fine-tuned on preference data using DPO.

❑ **Clarify-or-Direct Ans DPO**: The Clarify-or-Direct Ans model is further fine-tuned on the ***double-turn preference data*** over clarifying question and direct-answer responses using DPO.

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*
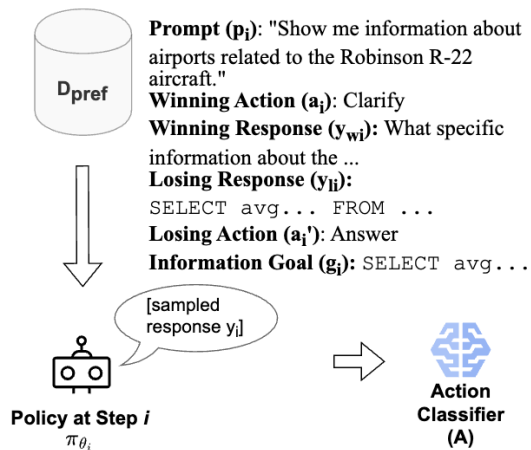
# Modeling Future Conversation Turns

| | # (↓) | Llama2 Answer F1 (↑) Unamb / Amb / All | # (↓) | Llama3 Answer F1 (↑) Unamb / Amb / All | # (↓) | Gemma Answer F1 (↑) Unamb / Amb / All |
|---|---|---|---|---|---|---|
| Direct-Ans SFT | | | | | | |
| w/ Greedy | 1 | 25.4 / 16.8 / 21.1 | 1 | 31.2 / 19.2 / 24.8 | 1 | 26.1 / 16.8 / 21.1 |
| w/ Sampled | 1 | 25.0 / 17.2 / 21.4 | 1 | 28.2 / 20.2 / 24.7 | 1 | 23.7 / 17.9 / 21.4 |
| Clarify SFT | 2 | 31.0 / 21.6 / 25.9 | 2 | 37.6 / 26.5 / 31.5 | 2 | 35.7 / 23.6 / 28.8 |
| Clarify DPO | | | | | | |
| w/ RM | 2 | 31.0 / 25.7 / 28.3 | 2 | 36.2 / 26.7 / 30.9 | 2 | 33.9 / 25.7 / 29.5 |
| w/ Likelihood | 2 | 30.2 / 23.9 / 27.2 | 2 | 43.5 / 29.6 / 359 | 2 | 37.3 / 26.8 / 31.5 |
| w/ Match | 2 | **38.3 / 28.2 / 32.8** | 2 | 42.9 / **3.17 / 36.5** | 2 | **40.7 / 28.6 / 33.9** |
| Clarify-or-Direct-Ans | | | | | | |
| SFT | 1.12 | 25.6 / 18.4 / 21.3 | 1.40 | 35.3 / 23.5 / 28.2 | 1.43 | 22.3 / 19.0 / 20.3 |
| DPO | 1.56 | 28.9 / 21.1 / 24.3 | 1.57 | 35.2 / 25.1 / 29.1 | 1.61 | 28.2 / 22.2 / 24.6 |

**Adding a clarifying turn can improve the performance on both ambiguous queries and unambiguous queries.**

53

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*

# Modeling Future Conversation Turns

| | # (↓) | Llama2 Answer F1 (↑) Unamb / Amb / All | # (↓) | Llama3 Answer F1 (↑) Unamb / Amb / All | # (↓) | Gemma Answer F1 (↑) Unamb / Amb / All |
|---|---|---|---|---|---|---|
| Direct-Ans SFT | | | | | | |
| w/ Greedy | 1 | 25.4 / 16.8 / 21.1 | 1 | 31.2 / 19.2 / 24.8 | 1 | 26.1 / 16.8 / 21.1 |
| w/ Sampled | 1 | 25.0 / 17.2 / 21.4 | 1 | 28.2 / 20.2 / 24.7 | 1 | 23.7 / 17.9 / 21.4 |
| Clarify SFT | 2 | 31.0 / 21.6 / 25.9 | 2 | 37.6 / 26.5 / 31.5 | 2 | 35.7 / 23.6 / 28.8 |
| Clarify DPO | | | | | | |
| w/ RM | 2 | 31.0 / 25.7 / 28.3 | 2 | 36.2 / 26.7 / 30.9 | 2 | 33.9 / 25.7 / 29.5 |
| w/ Likelihood | 2 | 30.2 / 23.9 / 27.2 | 2 | 43.5 / 29.6 / 359 | 2 | 37.3 / 26.8 / 31.5 |
| w/ Match | 2 | **38.3 / 28.2 / 32.8** | 2 | 42.9 / **3.17 / 36.5** | 2 | **40.7 / 28.6 / 33.9** |
| Clarify-or-Direct-Ans | | | | | | |
| SFT | 1.12 | 25.6 / 18.4 / 21.3 | 1.40 | 35.3 / 23.5 / 28.2 | 1.43 | 22.3 / 19.0 / 20.3 |
| DPO | 1.56 | 28.9 / 21.1 / 24.3 | 1.57 | 35.2 / 25.1 / 29.1 | 1.61 | 28.2 / 22.2 / 24.6 |

➢ **Clarify-or-Answer methods strike a balance between effectiveness and efficiency.**

➢ **DPO with double-turn preference data consistently outperforms SFT.**

*Zhang et al., "Modeling Future Conversation Turns to Teach LLMs to Ask Clarifying Questions" (ICLR '25)*

# Action-Based Contrastive Self-Training (ACT)



**On-Policy Response Sampling**

**Trajectory Simulation and Evaluation**

**Policy Update**

❑ **ACT** focuses on the clarification preference optimization in multi-turn conversations

❑ Construct conversation data with contrastive action pairs (*clarify* or *answer*) as the preference data

55

*Chen et al., "Learning to Clarify: Multi-turn Conversations with Action-Based Contrastive Self-Training" (ICLR '25)*