

OUTCOME-CONSTRAINED LARGE LANGUAGE MODELS FOR COUNTERING HATE SPEECH

We explore methods, including instruction prompts, LLM finetuning, and LLM reinforcement learning, to incorporate the desired conversation outcomes, including low conversation incivility and non-hateful hater reentry, into the generation of counter speech.

LINGZI HONG

UNIVERSITY OF NORTH TEXAS



PENGCHENG LUO

PEKING UNIVERSITY



EDUARDO BLANCO

UNIVERSITY OF ARIZONA



XIAOYING SONG

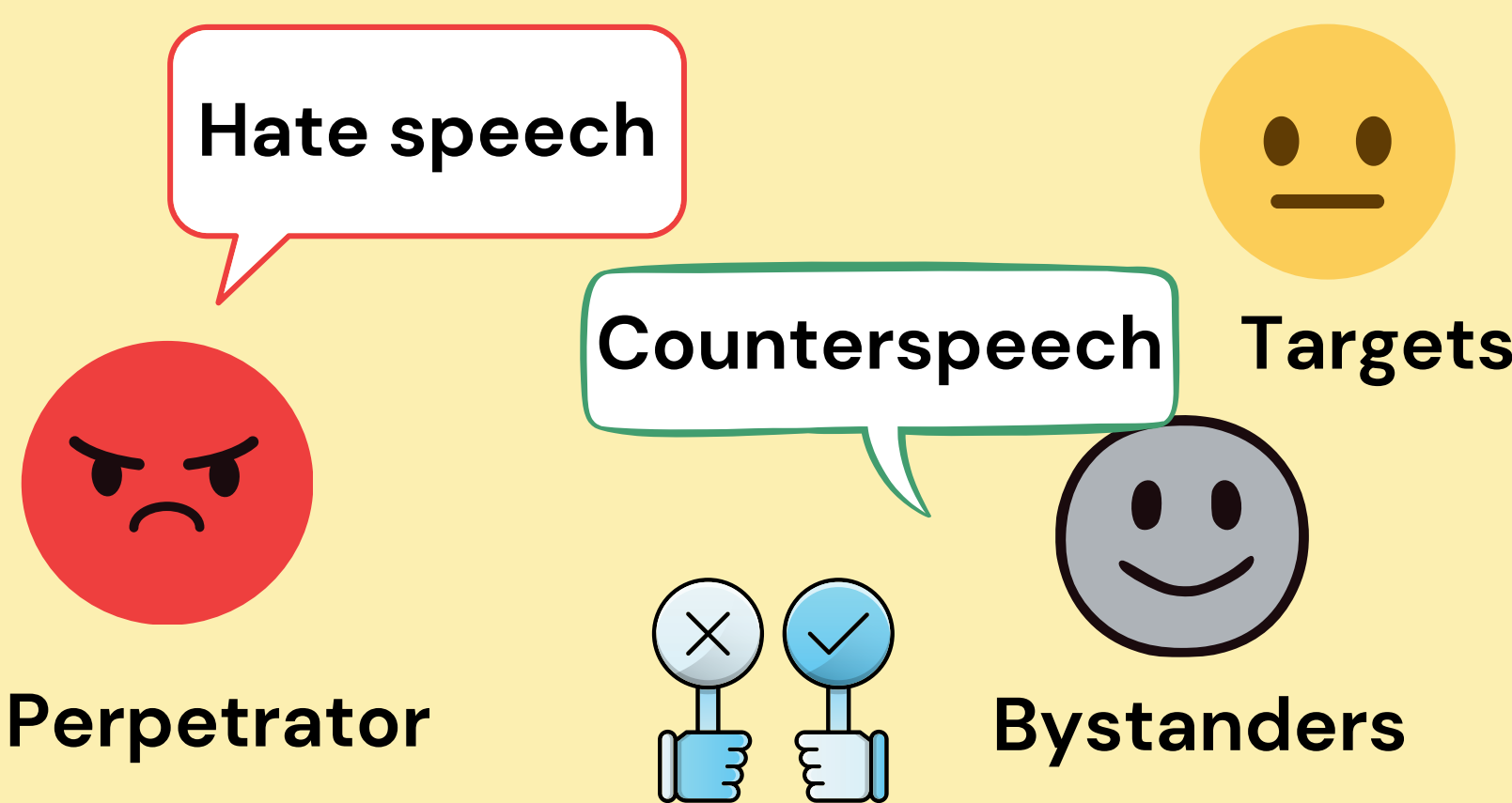
UNIVERSITY OF NORTH TEXAS



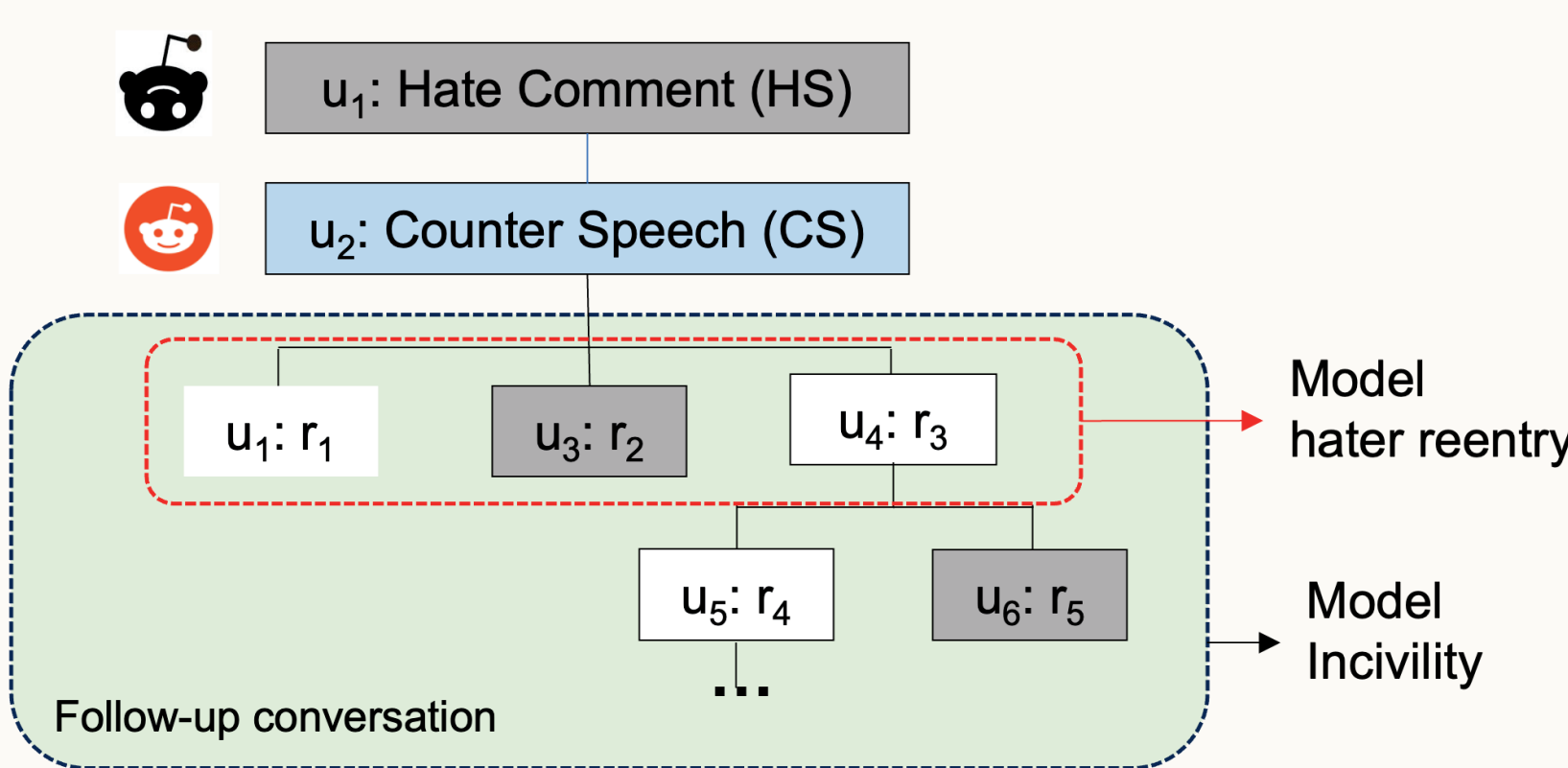
BACKGROUND

Hate speech has posed significant challenges to healthy and productive online communication.

Counterspeech that challenges or counteracts hate speech can moderate online hostilities, promote productive user engagement, and educate users.



CONVERSATION OUTCOMES



Conversation Incivility: A metric to measure the outcome based on the number of civil and uncivil comments and unique authors.

$$S(r) = aU(r) - (1-a)C(r)$$

Hater reentry: The hate perpetrators' reactions following a counterspeech: no reentry, hateful reentry, and non-hateful reentry.

DATASET

Reddit Data: Conversations related to hate speech were collected from 39 subreddits. The data is used to build conversation outcomes classifiers.

Benchmark-Reddit: A dataset with hate speech from Reddit and counterspeech by humans. The data is used for counterspeech generation and evaluation.

RESEARCH QUESTIONS

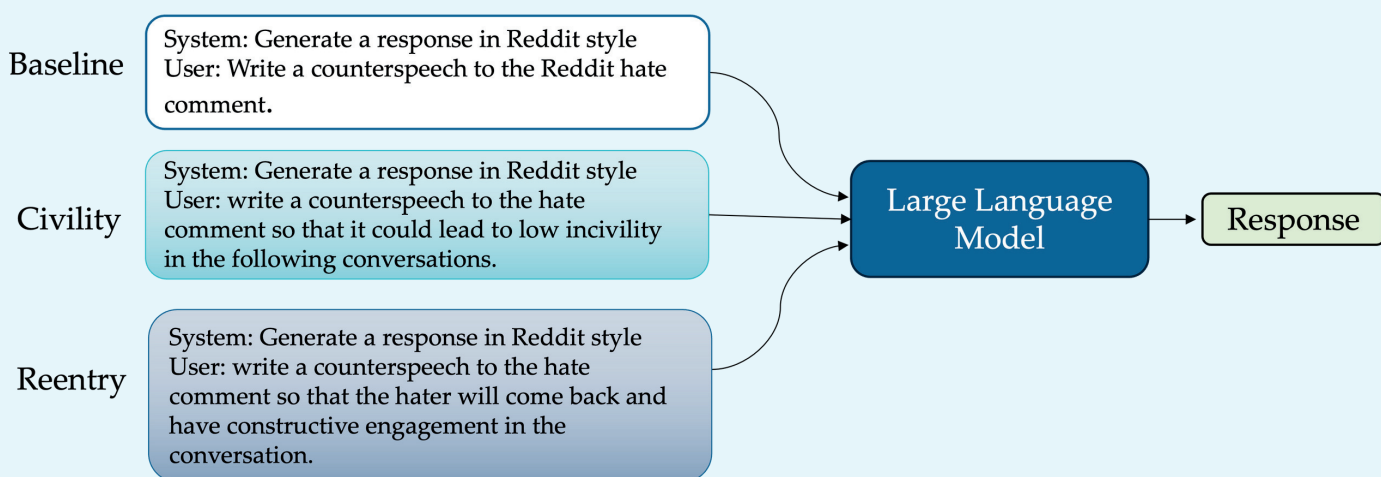
RQ1: How can **constraints on conversation outcomes** be incorporated into developing LLMs for generating counterspeech?

RQ2: How **effective** are these methods in generating outcome-oriented counterspeech?

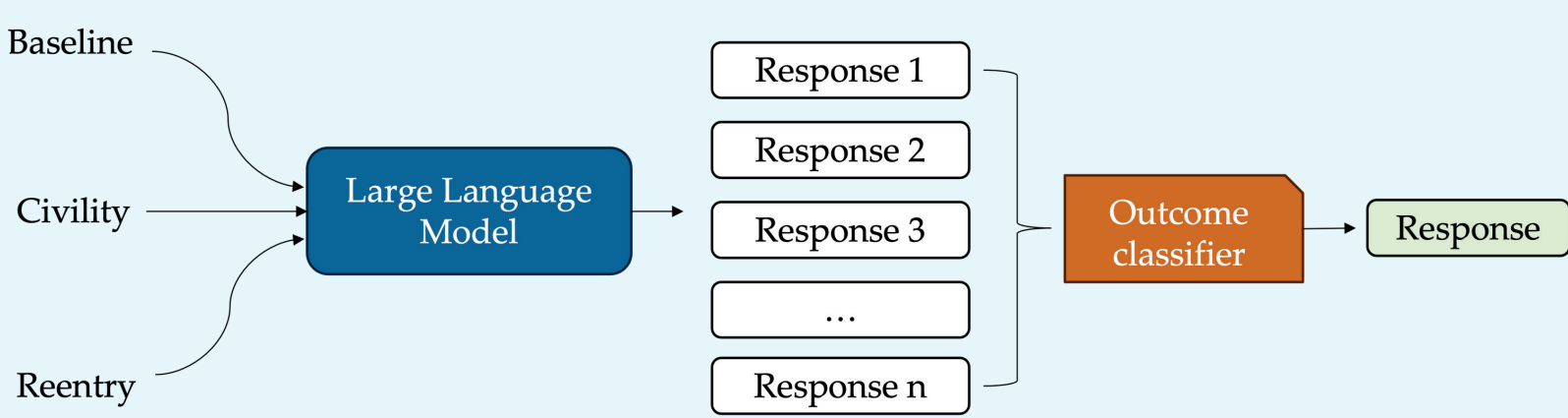
OUTCOME-CONSTRAINED

COUNTERSPEECH GENERATION

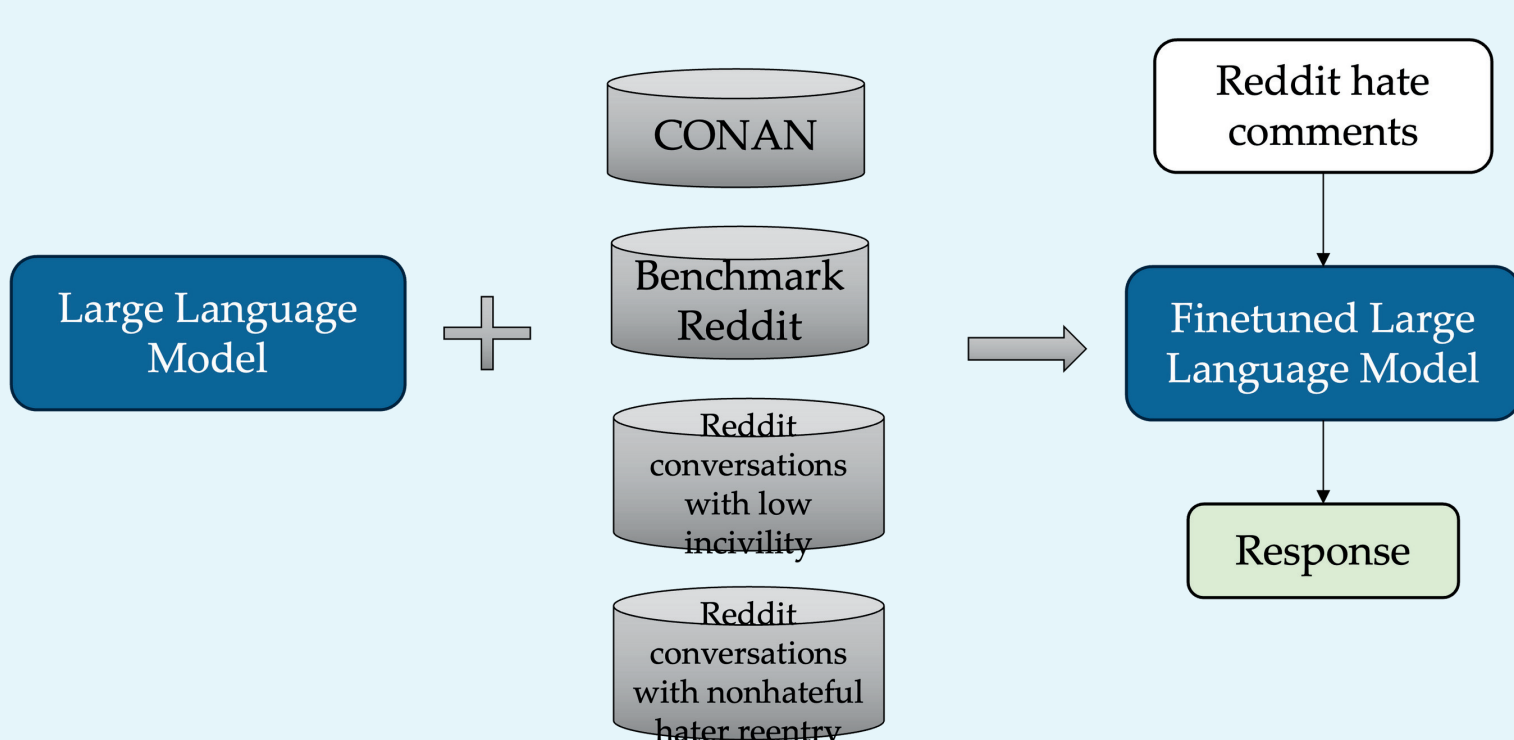
Instruction Prompts



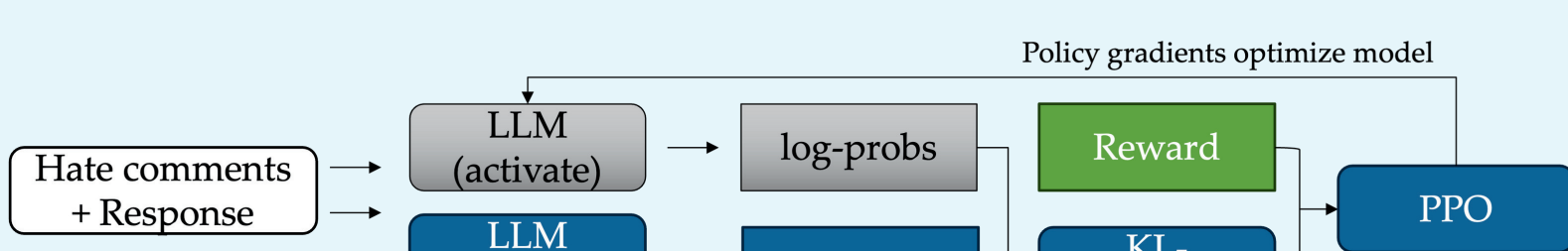
Prompt and Select



LLM Finetune

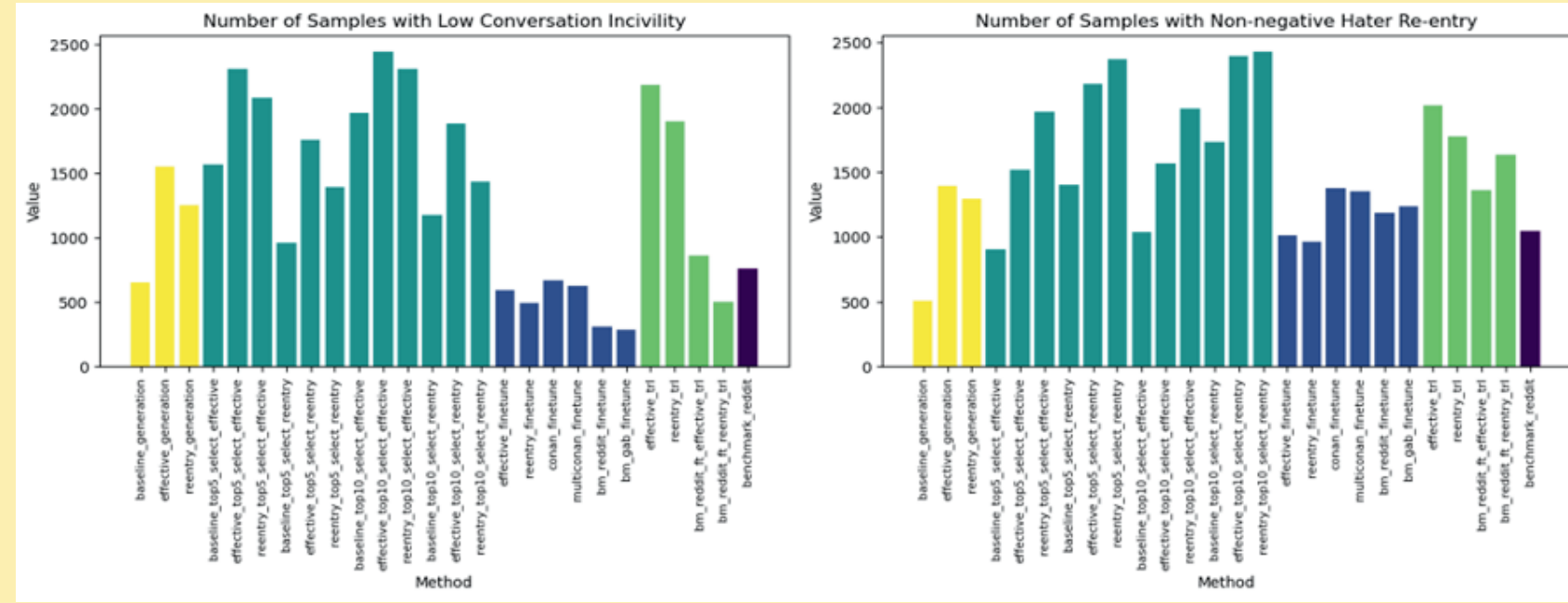


LLM Reinforcement Learning



RESULTS

Conversation Outcomes



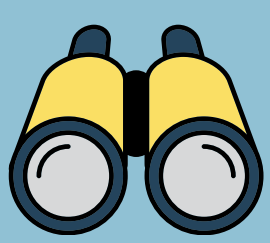
Instructions, prompt and select, and RL are effective strategies.

Stylistic Metric

Counterspeech by instruction prompts has lower text quality. LLM Finetune and RL generate text with better focus, low redundancy

Human Evaluation

| Method | Suitability | Relevance | Effectiveness |
|------------|-------------|-----------|---------------|
| Prompt | 0.50 | 0.88 | 0.54 |
| Finetuning | 0.80 | 0.68 | 0.80 |
| RL | 0.74 | 0.76 | 0.72 |



PREVIOUS RESEARCH

Informative counterspeech generated by LLMs with information from knowledge repositories

Polite and detoxified counterspeech generated by GEDI

Counterspeech with **different intents** with QUARC

Our study: outcome-constrained counterspeech with LLM finetuning and RL.

ACKNOWLEDGEMENT

This research was supported by the Institute of Museum and Library Services (US) under Grants LG-256661-OLS-24 and LG-256666-OLS-24.



Resilience of Immigrants to Crisis



LLM4Cat

