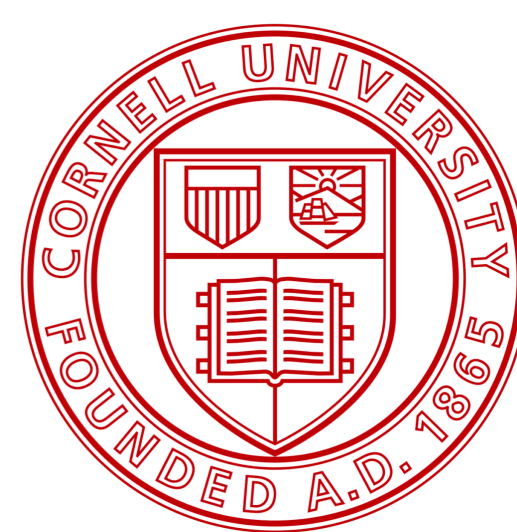


Aligning LLM Agents by Learning Latent Preference from User Edits

<https://github.com/gao-g/prelude>

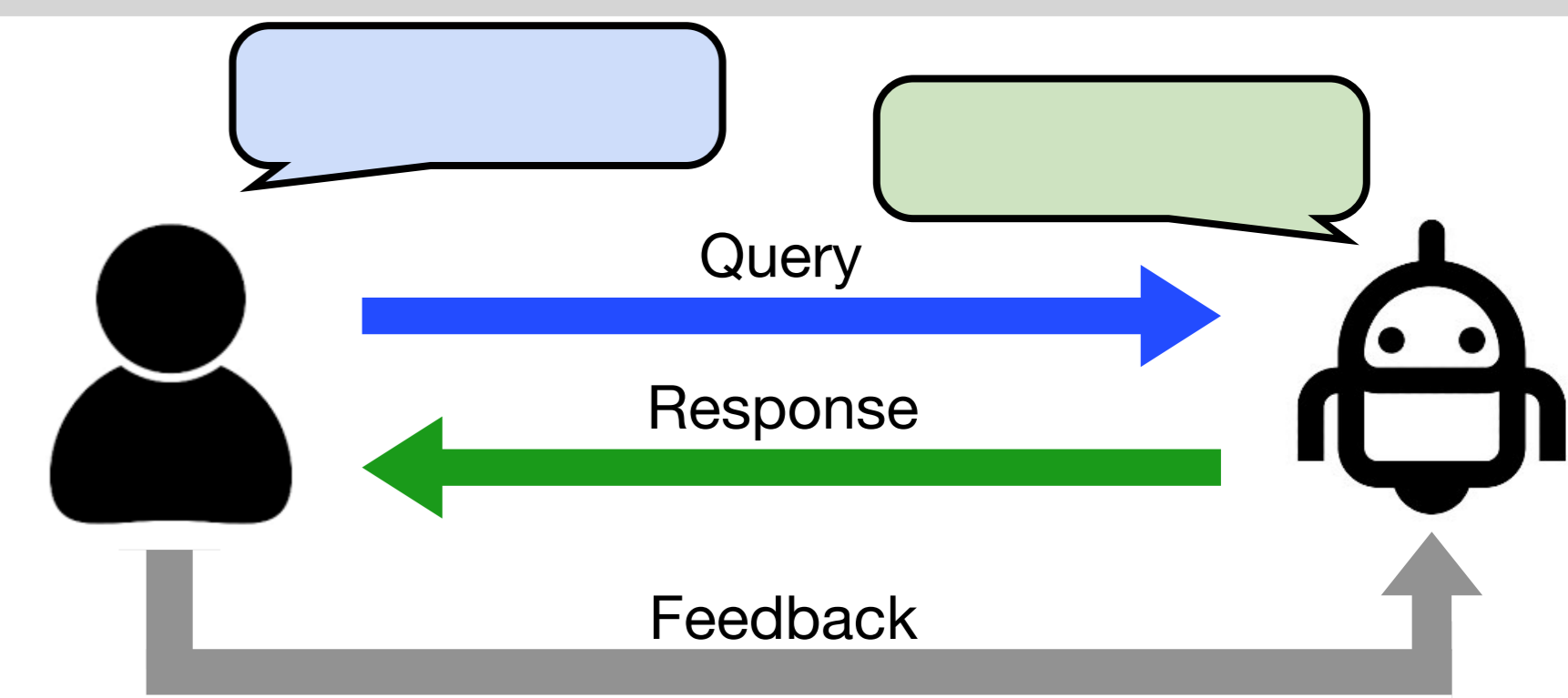


Ge Gao*, Alexey Taymanov*, Eduardo Salinas, Paul Mineiro, and Dipendra Misra

Overview

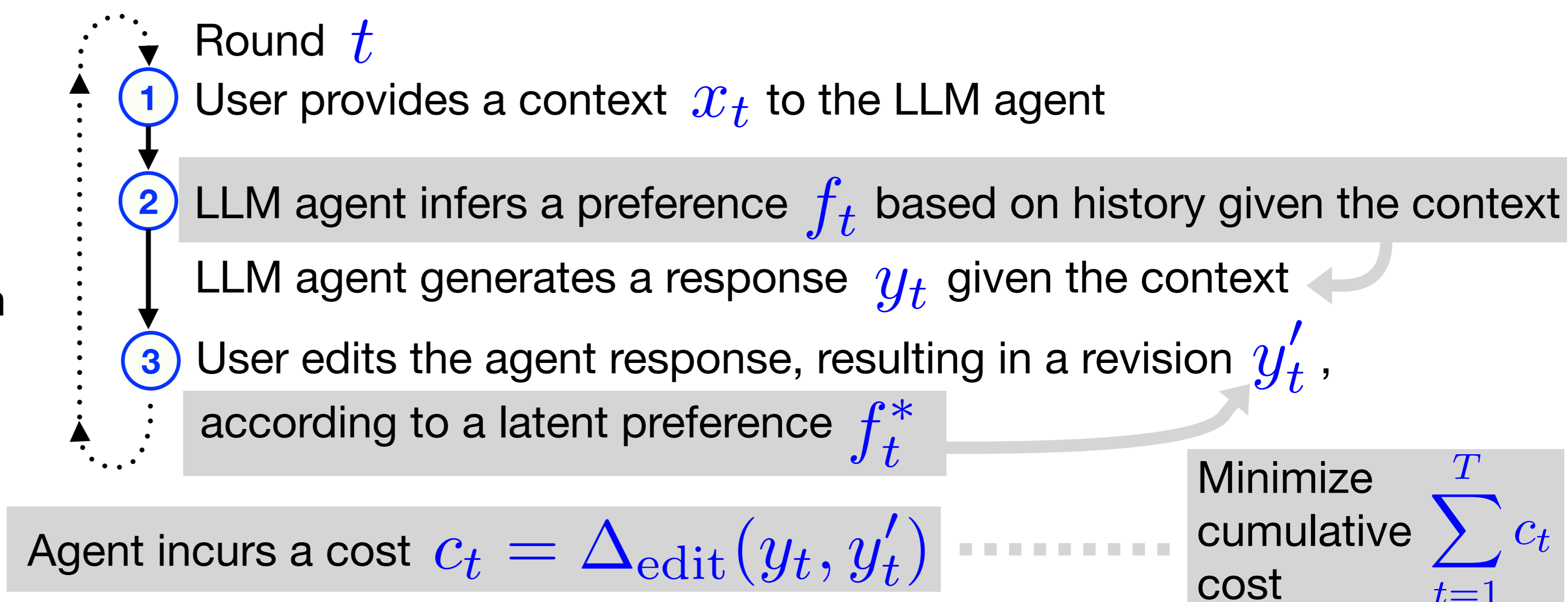
How to learn from user feedback in the form of edits to improve AI writing assistants?
We introduce:

- **PRELUDE** framework: formulates the interaction progress and preference learning as a cost minimization problem
- **CIPHER**: learns a prompt policy to infer a descriptive user preference



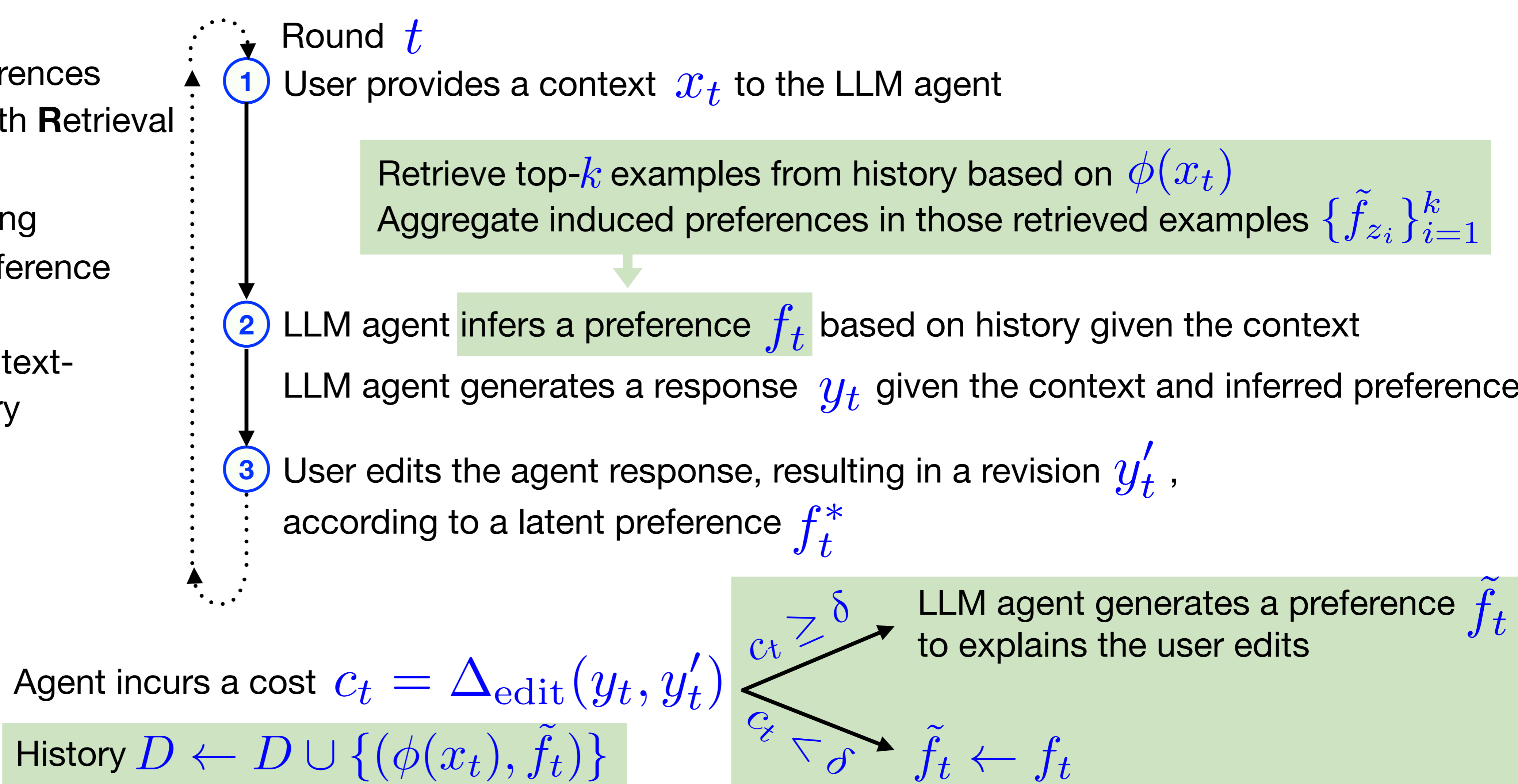
Learning Framework: PRELUDE

- **PRE**ference Learning from **U**ser's **D**irect **E**dits
- User directly makes edits to the agent response based on a latent preference
- Agent infers a user preference from interaction history, and uses it to generate a response
- Cost minimization to account for the amount of efforts spent by the user on making edits



Method: CIPHER

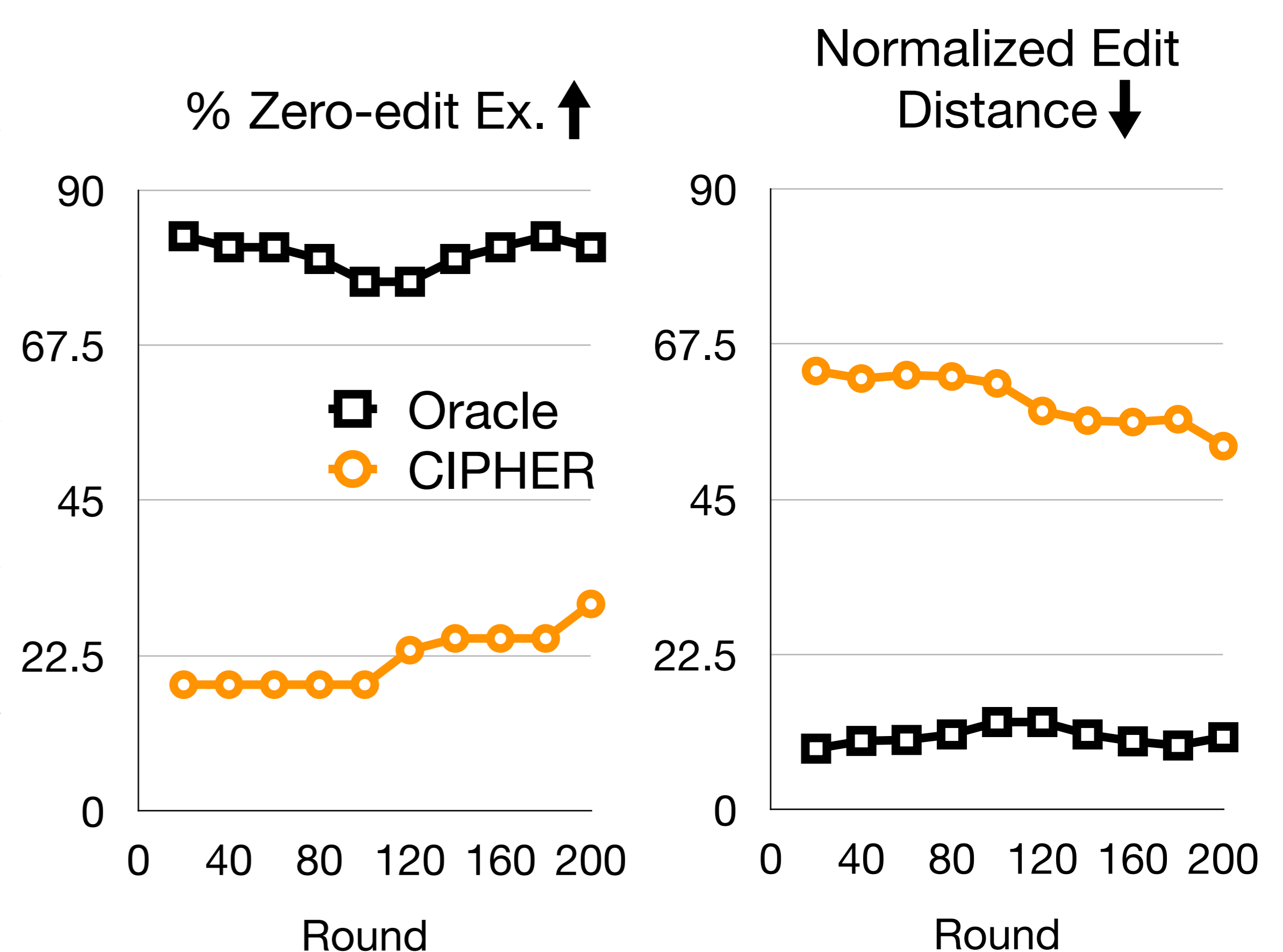
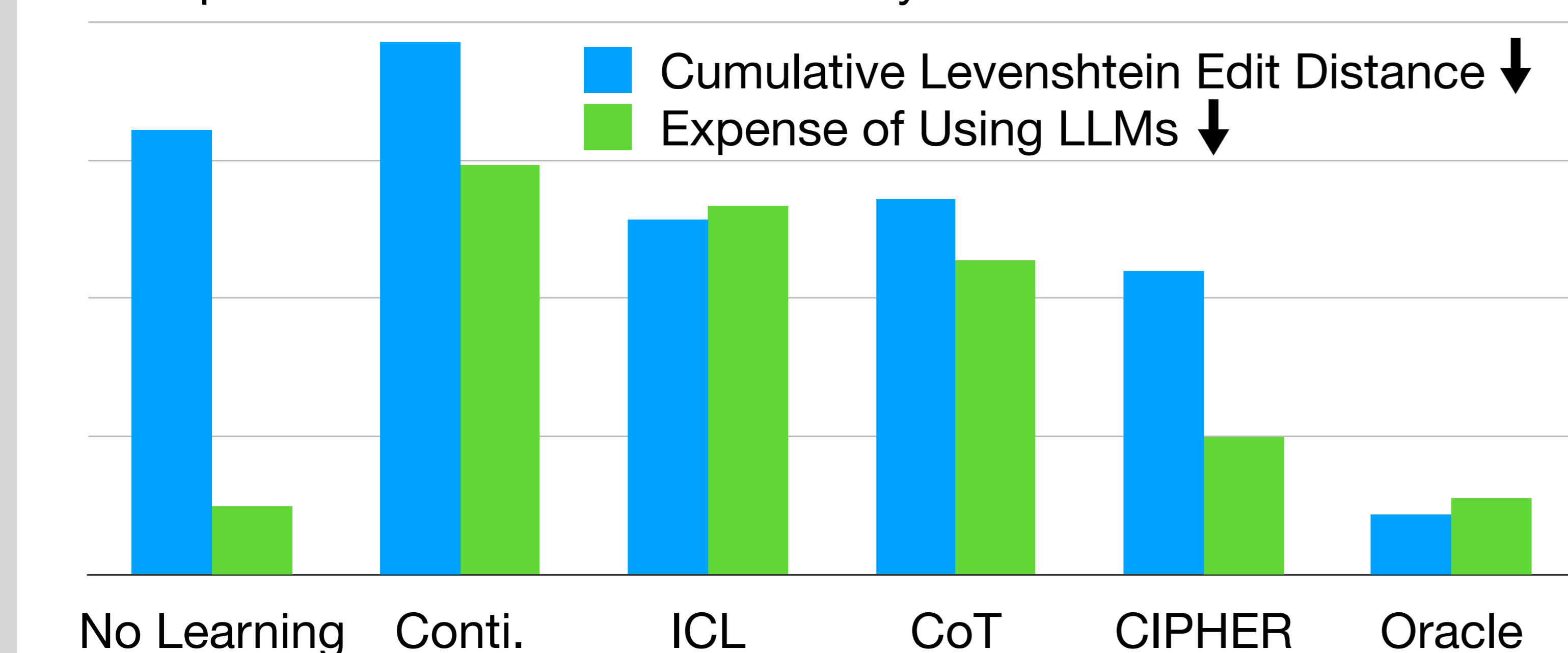
- **C**onsolidates **I**nduced Preferences based on **H**istorical **E**dits with **R**etrieval
- Leverages LLMs by prompting
- Infers a descriptive user preference from history with retrieval
- Manages a collection of context-dependent preference history
- Computationally efficient
- User friendly
- Interpretable



Experiments & Analysis

- Writing task: Summarization
- Experiment: 200 rounds of interaction
- User setup:
 - GPT-4 user as a simulation
 - Provide different context documents
 - Show context-dependent preference
- CIPHER setup:
 - GPT-4 as the base LLM
 - MPNet as the context representation function
 - Top 5 retrieval with cosine similarity

Use Case	Latent User Preference
Introduce a political news to kids	targeted to young children, storytelling, short sentences, playful language, interactive, positive
Promote a paper to invoke more attention	tweet style, simple English, inquisitive, skillful foreshadowing, with emojis
Take notes for knowledge from Wikipedia	bullet points, parallel structure, brief
Use online stories for inspirations in writing	second person narrative, brief, show emotions, invoke personal reflection, immersive
Extract main opinions from a movie review	question answering style



Interpretable ✓
Context-dependent ✓

Conti. ✓
ICL ✓
CoT ✓
CIPHER ✓