

**CHATHF: COLLECTING RICH HUMAN FEEDBACK FROM REAL-TIME
CONVERSATIONS**

A Dissertation
Presented to
The Academic Faculty

By

Andrew Li

In Partial Fulfillment
of the Requirements for the Degree
Masters of Computer Science in the
School of Interactive Computing
College of Computing

Georgia Institute of Technology

Dec 2024

© Andrew Li 2024

**CHATHF: COLLECTING RICH HUMAN FEEDBACK FROM REAL-TIME
CONVERSATIONS**

Thesis committee:

Dr. Alan Ritter, Advisor
School of Computer Science
Georgia Institute of Technology

Dr. Zhenduo Wang
School of Computer Science
Georgia Institute of Technology

Dr. Wei Xu
School of Computer Science
Georgia Institute of Technology

Date approved: August 23, 2024

ACKNOWLEDGMENTS

I would like to thank my advisor Dr. Alan Ritter for his extensive guidance throughout this research process, as well as my thesis committee, Dr. Wei Xu for bringing me into the NLP X lab and the cooking chatbot project in the first place, and Dr. Zhenduo Wang for his continuous advice and involvement.

Special thanks to Duong Le, Jeongrok Yu, and Ethan Mendes for their support and assistance on this project.

TABLE OF CONTENTS

Acknowledgments	iii
List of Tables	vi
List of Figures	vii
Summary	viii
Chapter 1: Introduction	1
Chapter 2: Related Work	4
2.1 ChatBot Evaluation	4
2.2 Human Feedback	5
Chapter 3: Chatbot Infrastructure	6
3.1 Multimodality	6
3.2 User Interface	8
3.3 Chat History	8
Chapter 4: Annotation Features	10
4.1 Customizable Annotation Configuration	10
4.2 Rich Human Feedback	11

Chapter 5: Example Use Cases	16
5.1 #1: Leveraging ChatHF to Collect Richly Annotated Geolocation Dialogues	16
5.2 #2: Supporting an HCI User Study for AI Cooking Assistance with Older Adults	17
Chapter 6: Conclusion	20
Appendices	21
Appendix A: Additional Screenshots of Chatbot Interface	22
References	25

LIST OF TABLES

5.1	Performance (F1-score) on the geolocation moderation task as evaluated on the GPTGEOCHAT test set [7]. The results from the best-performing moderation agent at each granularity are bolded	17
5.2	Percentage of responses from models having each type of error. The evaluation is conducted on 10 multi-turn conversations (131 generated responses) in the test set of the ChattyChef dataset (“Order”: wrong order, “Lack info.”: lack of information, “Wrong info.”: wrong information).	19

LIST OF FIGURES

1.1	ChatHF incorporates integrated multimodal dialogue annotation. This concept figure shows an example for privacy-preserving moderation in conversational geolocation QA [7].	3
3.1	Demonstration of the multimodality of ChatHF. Users can upload images into the conversation to ask about or include as context, and speak into their microphone to converse rather than typing.	7
3.2	Screenshot of the main ChatHF interface. Configuration options can be modified on the left panel, with changes automatically reflected on the chat interface on the right. See more screenshots of included features in Appendix Appendix A.	8
4.1	Demonstration of a multiple choice annotation for intent labeling of an AI cooking application with the option to give an explanation.	11
4.2	In this visual question-answering task, the model is unable to fully identify the university in the picture. The user uses a post-edit to correct the mistake.	12
4.3	In this cooking assistance dialogue task, the model gives the incorrect order of steps without the user immediately realizing. The user then reverses to previous turn to try again, with the model giving the correct order of steps the second time.	14
5.1	A pilot HCI user study using ChatHF configured to support a voice assistant cooking chatbot (section 5.2).	18
A.1	The screen to add a new model to the list.	23
A.2	The screen to create a custom annotation.	24

SUMMARY

We introduce ChatHF, an interactive annotation framework for chatbot evaluation, which integrates configurable annotation within a chat interface. ChatHF can be flexibly configured to accommodate various chatbot evaluation tasks, for example detecting offensive content, identifying incorrect or misleading information in chatbot responses, and chatbot responses that might compromise privacy. It supports post-editing of chatbot outputs and supports visual inputs, in addition to an optional voice interface. ChatHF is suitable for collection and annotation of NLP datasets, and Human-Computer Interaction studies, as demonstrated in case studies on image geolocation and assisting older adults with daily activities.

CHAPTER 1

INTRODUCTION

Advances in large language models and vision-language models have led to surprisingly effective chatbots such as GPT-4V [1], Llama-3 [2], Gemini [3], and many more. While these chatbots display interesting and useful emergent capabilities, they can also exhibit some undesirable behaviors which may not be straightforward to recognize or mitigate. Thus, the evaluation of LLM-based chatbots remains a challenge. Some studies make use of automated GPT-based evaluations [4], but human evaluation is still needed to truly measure the effectiveness of these automatic metrics on new tasks. Other recent works, such as Chatbot Arena [5], make use of human evaluators, but present only holistic evaluations of which model produces “better” outputs (i.e., preference) to simplify and accelerate the collection of human feedback.

In this paper we present an interactive framework, **ChatHF** (see chapter 3), for evaluation and analysis of chatbots that supports fine-grained error detection and collecting human feedback simultaneously (see section 4.2). Rather than the common workflow where researchers first collect LLM-generated responses then evaluate (or annotate) as an afterthought, we envision an approach where the human annotators seamlessly interleave annotation with conversation. That is, human evaluators directly chat with LLMs on specific topics relevant to the phenomenon to be studied (see Figure 1.1). This not only saves the annotator’s time and energy to accomplish two tasks in a single pass, but also encourages annotators to engage in more interesting and complex conversations — as we show in two case studies: cooking chatbot [6] and multimodal privacy QA [7].

ChatHF is flexible and can be configured for many annotation tasks, such as offensive outputs [8], misinformation [9], or compromised privacy [10], enabling the creation of curated conversational datasets and the study of emergent behaviors in LLM-based chat-

bots. Its unique features include flexible configuration, post-editing of chatbot outputs, and multimodal inputs with images and voice interaction through speech-to-text and text-to-speech. ChatHF is in-depth and accessible enough to support both standard NLP data collection and annotation, as well as interactive Human-Computer Interaction (HCI) studies involving chatbots. In the two case studies (see section 5.1 and section 5.2), ChatHF was used to (1) collect a dataset of image geolocation conversations that are labeled with the granularity of location information revealed at each step of the conversation, and (2) as an interface, to support an HCI user study on older adults using chatbots to assist with activities of daily living.

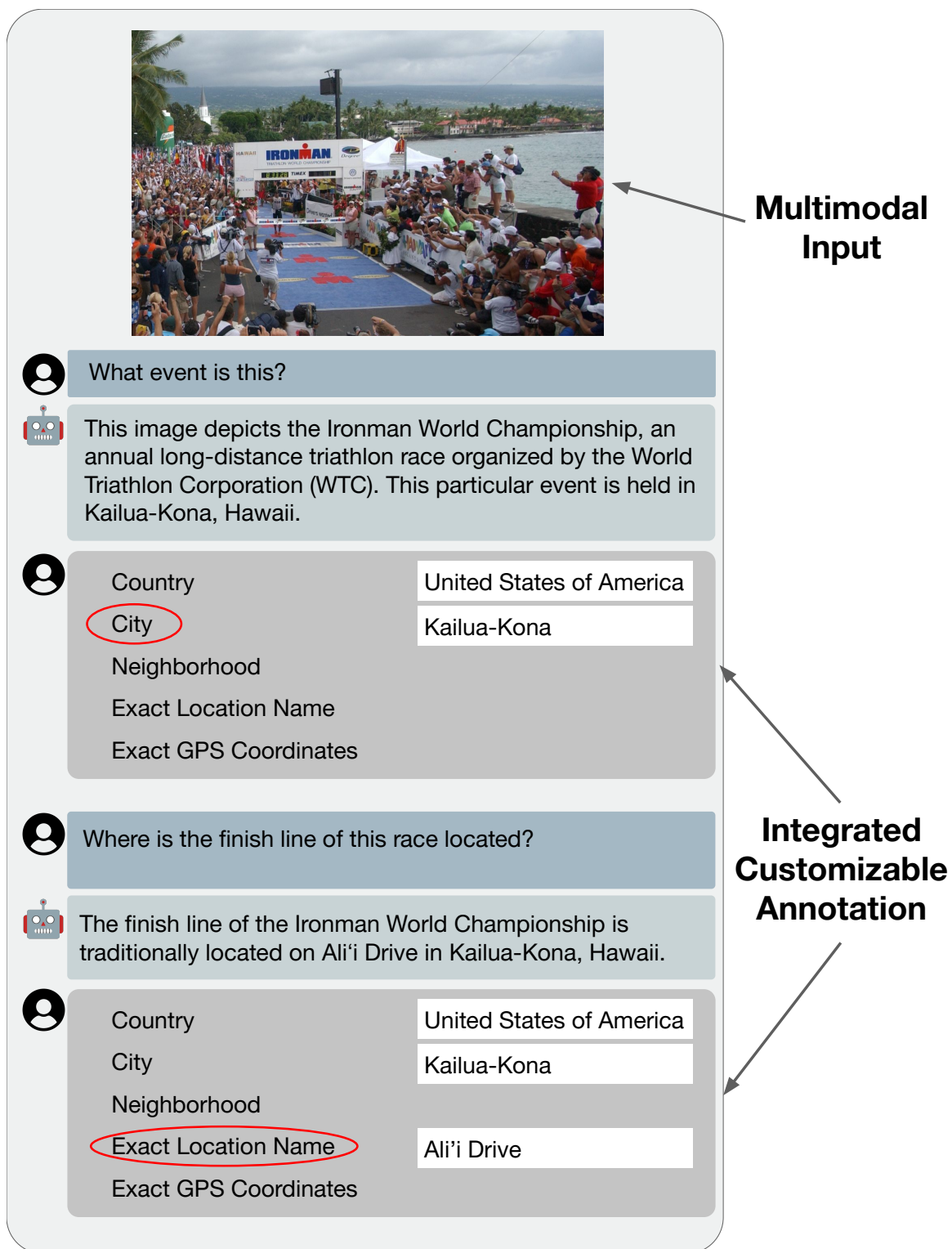


Figure 1.1: ChatHF incorporates integrated multimodal dialogue annotation. This concept figure shows an example for privacy-preserving moderation in conversational geolocation QA [7].

CHAPTER 2

RELATED WORK

The field of text annotation tools has seen iterative advancements in the past decade. This section gives a high-level overview of previous text annotation tools from two perspectives: conversational text evaluation and human feedback management.

2.1 ChatBot Evaluation

STAV [11] and BRAT [12] are examples of early text annotation tools. BRAT supports manual curation of the annotation and is optimized for rich structured annotation tasks and annotator productivity. It also provides high-quality annotation visualization. More recent tools like POTATO [13] support higher degrees of configuration and customization and provide even better quality control and productivity enhancement. However, most of them are primarily focused on annotation on a per-example basis, mostly useful for annotation tasks within one sentence or one paragraph rather than multi-turn conversations.

Within the field of conversational text annotation tools, there has been only a limited amount of available open-source tools. LIDA [14] was the first tool designed specifically for annotating multi-turn conversational text data [15]. Its later evolution MATILDA [16] improved it by facilitating multi-lingual and multi-annotator annotations. However, these tools have no web interfaces and require some technical knowledge for model integration and configuration, which inhibits their accessibility. EZCAT [17] can be used directly on their web application to both configure text labels, on a message or conversation level, and go through the annotation process. However, EZCAT does not have the option to collect multiple labels per turn. In this work, we aim to supply this field with a flexible multi-purpose annotation tool with a configurable and easy-to-use interface.

2.2 Human Feedback

It is increasingly important to audit and evaluate LLMs and VLMs through human judgement, beyond relying solely on traditional automatic metrics, and in turn, learn from rich and diverse human feedback (see the excellent survey by [18]) to improve the model’s performance. However, in addition to their restricted accessibility, existing annotation tools are also limited to only utilizing human feedback at the end of each conversation as an afterthought [19, 20, 21]. For example, INCEpTION [21] and GATE [22] provide large feature sets, but cannot display conversation data as turns [16]. LIDA and MATILDA fully support conversational text annotation tasks such as task-oriented dialogue systems. However, their frameworks can only be used to annotate static recorded dialogues. Such an annotation scheme fails to address human feedback during the conversation, which leads to systemic productivity loss.

In contrast, we present a customizable annotation tool capable of managing real-time human feedback during conversations. Annotators are allowed to edit model-generated utterances and to reverse and modify chat history to reflect their feedback. We track all these edits and reversals, as well as the reasons why these changes are made as free-text and/or multi-choice annotations.

CHAPTER 3

CHATBOT INFRASTRUCTURE

ChatHF supports various models and configuration options for easy prompt engineering and experimentation. Our public web demo¹ supports testing OpenAI, Anthropic, Google Gemini, and Mistral models directly through their respective APIs. For security, all configuration settings like API keys are stored client-side, and can be downloaded and loaded as a YAML file for easy sharing.

Run locally or self-hosted, ChatHF can be used with Ollama² and Huggingface³ models. Additionally, API keys can be hidden in an environment file. For more complex generation schemes, sample code is provided to set up a custom arbitrary generate function.

ChatHF also offers several configuration options to experiment with model settings, such as the system prompt, temperature, timeout limit, and conversation history memory length. Any changes are automatically reflected in the chat window. At each turn of the conversation, the model is passed the conversation history truncated to the memory length with the system prompt inserted at the start, and the model generates a response with the set temperature, timing out if the processing time exceeds the timeout limit.

3.1 Multimodality

To support voice chatbot applications, ChatHF integrates the option for text-to-speech on model outputs and speech-to-text with microphone input. Features such as press-to-talk, continuous listening, and text-to-speech are customizable, allowing ChatHF to cater to different needs from accessibility to hands-free operations.

Interfacing with Vision-Language models are also possible as ChatHF allows for image

¹<https://nlprx.cc.gatech.edu:2974/>

²<https://ollama.com/>

³<https://huggingface.co/>

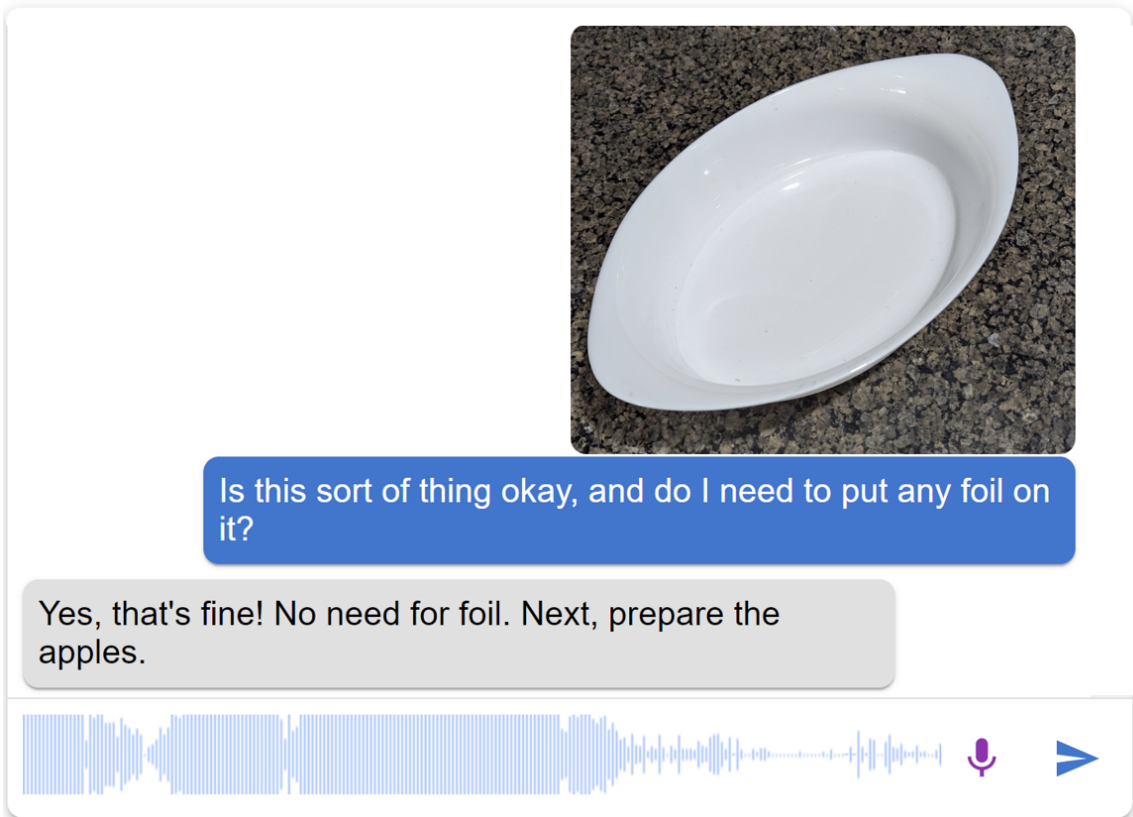


Figure 3.1: Demonstration of the multimodality of ChatHF. Users can upload images into the conversation to ask about or include as context, and speak into their microphone to converse rather than typing.

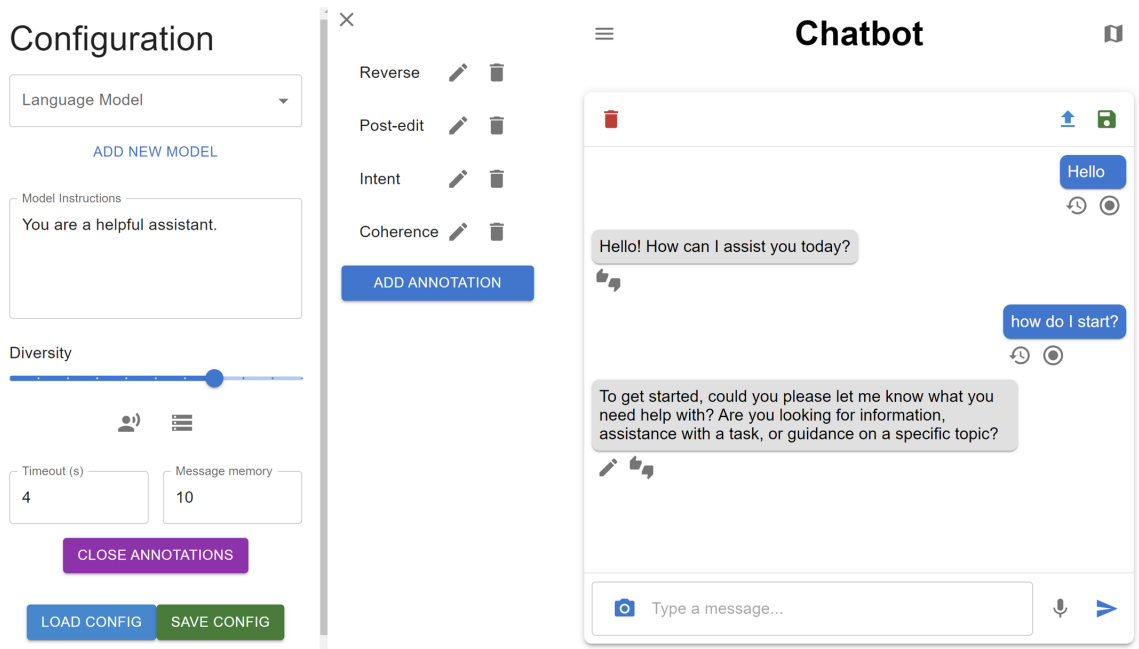


Figure 3.2: Screenshot of the main ChatHF interface. Configuration options can be modified on the left panel, with changes automatically reflected on the chat interface on the right. See more screenshots of included features in Appendix Appendix A.

input to the chatbox, which are simply saved as Base64 images in the chat history to be sent to the model.

3.2 User Interface

ChatHF is built on a Flask backend and a React frontend, with a publicly available codebase released under an Apache 2.0 license. We include a Flask backend written in Python to allow for easier integration of custom models or generation schemes into our chatbot interface. Text-to-speech and and speech-to-text are implemented via Azure AI Speech⁴, using their proprietary models.

3.3 Chat History

All messages in the chat history are saved into a JSON log file, timestamped with the date and time. User feedback is saved with each message with the user-specified name and

⁴<https://azure.microsoft.com/en-us/products/ai-services/ai-speech>

value. In the case of a reversal, the old chat history is not overwritten, and instead, an additional chat history created with all messages until the reversal point.

ChatHF supports downloading the log file locally or to a database such as Google Firebase⁵, as well as uploading a log file to view the chat history or edit the evaluation later. The user also has the option to clear the chat history to start a new conversation.

⁵<https://firebase.google.com/>

CHAPTER 4

ANNOTATION FEATURES

4.1 Customizable Annotation Configuration

In addition to the chatbot interface, ChatHF enables integrated on-the-fly human evaluation of the generated conversation and allows users to customize the annotation formats according to their needs. During a conversation, the user can annotate user messages, the generated model responses, or both. These messages can be annotated in various formats including binary, Likert-scale, multiple-choice, multiple-select, and free-text inputs. All annotation types can have a custom question and the option to require the annotator to provide an explanation through an additional textbox. Furthermore, the labels for binary, Likert-scale, multiple-choice, and multiple-select annotations are all customizable, and annotations can be specific to user messages, model responses, or both.

The full control of the annotation format and customizable labels is implemented as an annotator’s configuration panel in our tool located in the upper left corner. The panel settings can be saved and uploaded for reuse later. If needed, custom annotations can also be edited and deleted.

In the chatbot interface, if the annotation feature is turned on, icons representing each annotation type appear below each user message or model response (See Figure 3.2). Users can click on an annotation icon to reveal its prompt and input the specified response. This process is quick and responsive to facilitate real-time fine-grained data collection.

To demonstrate the efficacy of ChatHF’s customizable evaluation, we describe and release sample configuration files for our two example use cases.

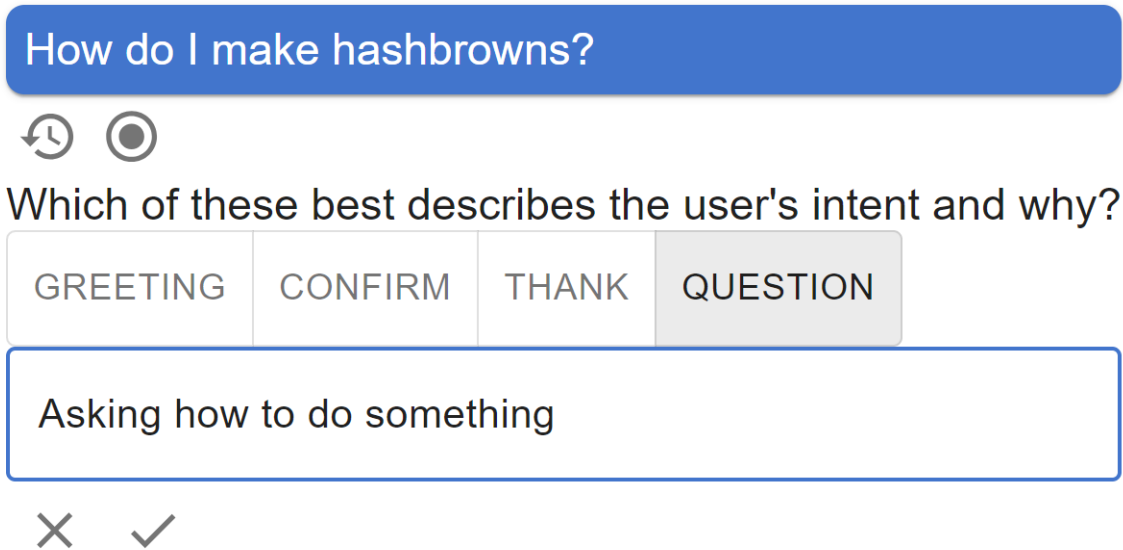


Figure 4.1: Demonstration of a multiple choice annotation for intent labeling of an AI cooking application with the option to give an explanation.

4.2 Rich Human Feedback

Along with the more traditional formats for human feedback, ChatHF includes two unique annotation types to collect real-time post-editing and reversal data for richer human feedback.

Post-editing Post-editing can be useful when only a portion of the model response is incorrect and requires changing or deleting, or if the output could be improved with just a minor addition. For instance, hallucinations and toxic language can be edited out and the offending spans can be easily extracted by comparing the post-edited and original text. Post-editing is also helpful when the model is partially correct, such as Figure 4.2, allowing for fine-grained corrections.

Crucially, post-editing corrects the conversation history, so that errors cannot propagate. This creates a more seamless chat experience and reduces the need to restart or reverse the conversation, which can be especially valuable in time, effort, or resource sensitive



What is the university that's on the left?

It looks like that's ~~Idaho~~ Boise State University.

Figure 4.2: In this visual question-answering task, the model is unable to fully identify the university in the picture. The user uses a post-edit to correct the mistake.

situations such as human studies in real world settings. (§section 5.2).

With post-editing selected in the configuration, users can directly edit the LLM-generated response. Similarly to the other annotations, users may be required to provide an explanation for the edit. Upon confirming the post-edit, the previous conversation history before the edit is added to the conversation log as a record of an unsuccessful termination.

Furthermore, each message stores its post-edit, with the most recent edit and original model output saved to the conversation log file. To ensure there is a fair evaluation only the most recent bot-message are editable. A list of the edits made will automatically be generated and saved as well.

Reversal In other cases, the model may have made an error that was not caught earlier in the conversation or had errors build up until the conversation was no longer salvageable. For instance, in instructional tasks where the order of instructions is crucial such as cooking, errors cannot be corrected by continuing the conversation, such as in the example in Figure 4.3. The choice to reverse may even be more subtle, perhaps due to uninteresting or stagnant dialog. Either way, it would be helpful to identify at which turn the conversation was recognized to be unrecoverable, and the point where the direction of the conversation shifted.

ChatHF’s reversal option allows for this rich feedback, saving both the reversed chat in the JSON log as well as either an optional annotator-provided reversal explanation or a simple indication of the success of the final dialogue. By default, when saving the conversation log, the current, most recent conversation is considered successful.

Multi-branch Conversation Employing the post-editing and reversal features, ChatHF can be used to explore a branching dialogue with multiple potentially successful continuations or completions. The set of branching conversations created by post-editing and reversing can be represented with a tree structure. At the simplest, a single continuous conversation is represented as one node. Once a branch is made, the conversation truncated

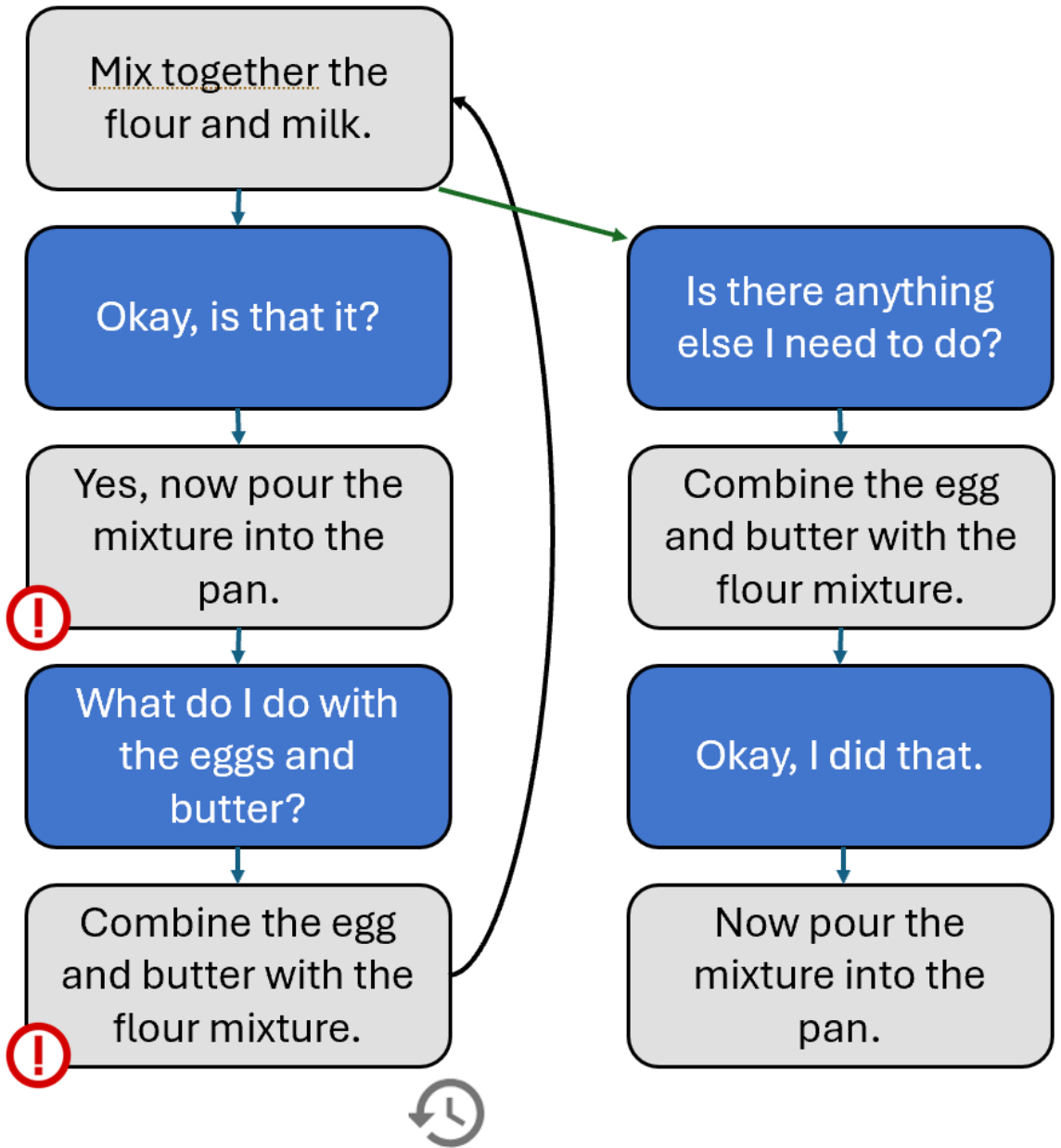


Figure 4.3: In this cooking assistance dialogue task, the model gives the incorrect order of steps without the user immediately realizing. The user then reverses to previous turn to try again, with the model giving the correct order of steps the second time.

at the branching point is set as the parent node, and the messages after the branching point both in the previous conversation and in the new conversation are each a child node.

This tree of interactions over a single overarching conversation topic can be viewed and each node can be selected to jump to a certain conversation.

CHAPTER 5

EXAMPLE USE CASES

We show the efficacy of the ChatHF framework through the following two use cases. In particular, we showcase the multimodality of ChatHF with the visual integration supporting the collection and annotation of geolocation dialogues and the audio functionality facilitating the conducting of an AI cooking assistance study in a real kitchen environment.

For transparency, neither of these contributions are part of this thesis work and are entirely the efforts of their respective authors. We include these examples to demonstrate proven use cases for ChatHF.

5.1 #1: Leveraging ChatHF to Collect Richly Annotated Geolocation Dialogues

ChatHF was built on to construct GPTGEOCHAT [7], a benchmark for granular privacy controls to moderate image geolocation dialogues, i.e. a human having multi-turn dialogues with a model about the location of an image provided in context. This work showcases the *multimodal model integration* of ChatHF (see chapter 3). The goal of this task was to train moderation agents to determine whether or not to withhold a vision language model (VLM) response based on whether or not the response violated the granular system privacy configurations:

$$[\text{Granularity Config, Image, Dialogue}] \xrightarrow{\text{Agent}} [Y, N]$$

For the studied geolocation task, these granular configurations were location granularities e.g., the *city*, *neighborhood*, or *exact-gps-coordinates* indicating the level of geolocation should be allowed during a conversation.

Table 5.1: Performance (F1-score) on the geolocation moderation task as evaluated on the GPTGEOCHAT test set [7]. The results from the best-performing moderation agent at each granularity are **bolded**.

Agent	Country	City	Neighborhood	Exact Location Name	Exact GPS Coordinates
LLaVA-13B (prompted)	0.56	0.55	0.52	0.41	0.48
IDEFICS-80B-instruct (prompted)	0.80	0.74	0.67	0.62	0.28
GPT-4v (prompted)	0.86	0.89	0.84	0.73	0.76
LLaVA-13B (finetuned on GPTGEOCHAT)	0.87	0.89	0.84	0.79	0.96

Data Collection To train and evaluate geolocation moderation agents, 1000 GPT4V-human dialogues are collected towards image geolocation, which form GPTGEOCHAT [7]. In-house annotators conversed with GPT-4v about the location of the image provided in context using ChatHF. During the conversation, each model response was annotated for (1) the finest granularity (*country, city, neighborhood, exact-location-name, exact-gps-coordinates*) of the location information revealed so far in the dialogue (2) the corresponding revealed location information e.g. $\{ 'country': 'United Kingdom', 'city': 'London' \}$. For the finest granularity, they represent each of the five granularities along with a *none* option using ChatHF’s multiple-choice annotation input. Similarly, they use multiple ChatHF-supported free-form text input fields for the corresponding location information.

Task Evaluation As shown in Table 5.1, finetuning a smaller model on a small high-quality training set of 400 dialogues from GPTGEOCHAT yields superior performance on the geolocation dialogue moderation task compared to prompting much larger models.

5.2 #2: Supporting an HCI User Study for AI Cooking Assistance with Older Adults

ChatHF was deployed to support an HCI user study on how a cooking chatbot can assist older adults to cook, an important activity of daily living, in coordination with the NSF AI Caring Institute.¹ In our pilot study (Figure 5.1), we configure ChatHF to work in a real kitchen environment, where the system interacts with users via a voice interface (i.e.,

¹<https://www.ai-caring.org/>



Figure 5.1: A pilot HCI user study using ChatHF configured to support a voice assistant cooking chatbot (section 5.2).

speech-to-text and text-to-speech modules) and help him/her to prepare meals. Particularly, we add a "press to talk" button to support the study condition, and reduce the speed of the text-to-speech module. In addition, we conduct prompt engineering to instruct the GPT-4o-mini to provide step-by-step and easy-to-follow guidance to users.² Our next plan is to have users from the target population to interact with ChatHF to identify specific challenges that older adults might face when using this technology.

ChatHF is also used to support the human analysis of the responses from different cooking chatbots. In this study, we investigate the outputs of ChatGPT and different fine-tuned versions of GPT-J models [23]: the base GPT-J model, GPT-J model incorporated with user intent information (*GPT-J+int*), GPT-J model incorporated with the instruction state information (*GPT-J+cut* and *GPT-J+ctr*), and GPT-J model incorporated with both types of information (*GPT-J+ctr+int*). In each conversation, each model response is annotated as correct or having one of the following errors: wrong order, irrelevant, lack of information, or wrong information. Table 5.2 demonstrates the error analysis of responses of the models on a subset of the test set of the Chattychef dataset [6].

²The configured ChatHF for cooking chatbots is available at: <https://tinyurl.com/chattychef2>

Table 5.2: Percentage of responses from models having each type of error. The evaluation is conducted on 10 multi-turn conversations (131 generated responses) in the test set of the ChattyChef dataset (“Order”: wrong order, “Lack info.”: lack of information, “Wrong info.”: wrong information).

Models	Order	Irrelevant	Lack info.	Wrong info.
GPT-J	22.9	10.7	8.4	8.4
GPT-J+int	18.3	8.4	11.5	6.1
GPT-J+cut	20.6	6.9	10.7	6.1
GPT-J+ctr	23.7	3.8	11.5	4.6
GPT-J+ctr+int	22.9	5.3	9.9	7.6
ChatGPT	6.1	0.0	1.5	3.1

CHAPTER 6

CONCLUSION

We present ChatHF, an interactive, customizable, and open-source tool for evaluating LLM-based multimodal chatbots with rich human feedback and annotation. It supports *real-time* conversation and manual annotation (or human evaluation) at the same time. For example, the users may directly revise LLM-generated response or request the LLM to regenerate another response when they are not satisfied with the LLM-generated response, then continue on the conversation, etc.

Appendices

APPENDIX A
ADDITIONAL SCREENSHOTS OF CHATBOT INTERFACE

Chatbot

Add New Model

Provider
API

API Provider
OpenAI

API Key
.....

Model Title
Cooking Chatbot

Model Name
gpt-4o-mini

Image Capable

CANCEL ADD MODEL

Figure A.1: The screen to add a new model to the list.

Create Annotation

Name
Intent

Annotation Type
Select

Question
Which of these best describes the user's intent and why?

Greeting	Confirm	Thank	Question
----------	---------	-------	----------

+ ADD OPTION

Annotation Message
User Only

Require Explanation?

CLOSE CREATE

Figure A.2: The screen to create a custom annotation.

REFERENCES

- [1] OpenAI, *Gpt-4 technical report*, 2024. arXiv: 2303.08774 [cs.CL].
- [2] A. Dubey *et al.*, *The llama 3 herd of models*, 2024. arXiv: 2407.21783 [cs.AI].
- [3] G. Team *et al.*, *Gemini: A family of highly capable multimodal models*, 2024. arXiv: 2312.11805 [cs.CL].
- [4] Y. Liu, D. Iter, Y. Xu, S. Wang, R. Xu, and C. Zhu, “G-eval: Nlg evaluation using gpt-4 with better human alignment,” in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 2511–2522.
- [5] W.-L. Chiang *et al.*, “Chatbot Arena: An open platform for evaluating llms by human preference,” *arXiv preprint arXiv:2403.04132*, 2024.
- [6] D. Le, R. Guo, W. Xu, and A. Ritter, “Improved instruction ordering in recipe-grounded conversation,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds., Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 10 086–10 104.
- [7] E. Mendes, Y. Chen, J. Hays, S. Das, W. Xu, and A. Ritter, *Granular privacy control for geolocation with vision language models*, 2024. arXiv: 2407.04952 [cs.CL].
- [8] A. Baheti, M. Sap, A. Ritter, and M. Riedl, “Just say no: Analyzing the stance of neural dialogue generation in offensive contexts,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds., Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 4846–4862.
- [9] E. Musi, E. Carmi, C. Reed, S. Yates, and K. O’Halloran, “Developing misinformation immunity: How to reason-check fallacious news in a human–computer interaction environment,” *Social Media + Society*, vol. 9, no. 1, p. 20 563 051 221 150 407, 2023.
- [10] Z. Zhang *et al.*, ““it’s a fair game”, or is it? examining how users navigate disclosure risks and benefits when using llm-based conversational agents,” in *Proceedings of the CHI Conference on Human Factors in Computing Systems*, ser. CHI ’24, Honolulu, HI, USA: Association for Computing Machinery, 2024, ISBN: 9798400703300.
- [11] P. Stenetorp, G. Topić, S. Pyysalo, T. Ohta, J.-D. Kim, and J. Tsujii, “Bionlp shared task 2011: Supporting resources,” in *Proceedings of BioNLP Shared Task 2011 Work-*

shop, Portland, Oregon, USA: Association for Computational Linguistics, Jun. 2011, pp. 112–120.

- [12] P. Stenetorp, S. Pyysalo, G. Topić, T. Ohta, S. Ananiadou, and J. Tsujii, “Brat: A web-based tool for NLP-assisted text annotation,” in *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*, F. Segond, Ed., Avignon, France: Association for Computational Linguistics, Apr. 2012, pp. 102–107.
- [13] J. Pei *et al.*, “POTATO: The portable text annotation tool,” in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, W. Che and E. Shutova, Eds., Abu Dhabi, UAE: Association for Computational Linguistics, Dec. 2022, pp. 327–337.
- [14] E. Collins, N. Rozanov, and B. Zhang, “LIDA: Lightweight interactive dialogue annotator,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations*, S. Padó and R. Huang, Eds., Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 121–126.
- [15] X. Liu, W. Xue, Q. Su, W. Nie, and W. Peng, “MetaCAT: A metadata-based task-oriented chatbot annotation tool,” in *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: System Demonstrations*, D. Wong and D. Kiela, Eds., Suzhou, China: Association for Computational Linguistics, Dec. 2020, pp. 20–25.
- [16] D. Cucurnia, N. Rozanov, I. Sucameli, A. Ciuffoletti, and M. Simi, “MATILDA - multi-AnnoTator multi-language InteractiveLight-weight dialogue annotator,” in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, D. Gkatzia and D. Seddah, Eds., Online: Association for Computational Linguistics, Apr. 2021, pp. 32–39.
- [17] G. Guibon, L. Lefeuvre, M. Labeau, and C. Clavel, “EZCAT: An easy conversation annotation tool,” in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, N. Calzolari *et al.*, Eds., Marseille, France: European Language Resources Association, Jun. 2022, pp. 1788–1797.
- [18] L. Pan, M. Saxon, W. Xu, D. Nathani, X. Wang, and W. Y. Wang, “Automatically correcting large language models: Surveying the landscape of diverse automated correction strategies,” *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 484–506, 2024.

- [19] P. A. Heeman, F. Yang, and S. E. Strayer, “DialogueView - an annotation tool for dialogue,” in *Proceedings of the Third SIGdial Workshop on Discourse and Dialogue*, Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 50–59.
- [20] M. Garg, C. Saxena, S. Saha, V. Krishnan, R. Joshi, and V. Mago, “CAMS: An annotated corpus for causal analysis of mental health issues in social media posts,” in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, N. Calzolari *et al.*, Eds., Marseille, France: European Language Resources Association, Jun. 2022, pp. 6387–6396.
- [21] J.-C. Klie, M. Bugert, B. Boullosa, R. Eckart de Castilho, and I. Gurevych, “The INCEpTION platform: Machine-assisted and knowledge-oriented interactive annotation,” in *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, D. Zhao, Ed., Santa Fe, New Mexico: Association for Computational Linguistics, Aug. 2018, pp. 5–9.
- [22] H. Cunningham, D. Maynard, K. Bontcheva, and V. Tablan, “GATE: An architecture for development of robust HLT applications,” in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, P. Isabelle, E. Charniak, and D. Lin, Eds., Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 168–175.
- [23] B. Wang and A. Komatsuzaki, *GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model*, <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.