

EVALUATING LANGUAGE MODELS IN LONGER CONVERSATIONAL CONTEXTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Evaluating long-form conversations between humans and large language models (LLMs) poses significant challenges in the field of natural language processing. Traditional evaluation metrics and benchmarks have largely focused on shorter language interactions and often fail to capture the nuances inherent in extended dialogues. To address this, we introduce UPHELD, a publicly available dataset featuring human-annotated long-form dialogues. This dataset not only facilitates robust benchmarking but also serves as a foundation for further research into conversation evaluation methodologies. Using our dataset, we systematically analyze the correlation between current LLM evaluation metrics and human judgment within long-form conversation scenarios. Our findings reveal that conventional metrics lack the sensitivity necessary to assess the complex and often subjective nature of prolonged interactions. We use our dataset to develop an improved evaluation metric that demonstrates a significantly higher correlation with human assessments. The work highlights the need for advanced metric designs and outlines a clear pathway to refine the evaluation of LLM long-form conversations.

1 INTRODUCTION

Large Language Models (LLMs) have demonstrated remarkable capabilities in generating high-quality text, particularly in tasks requiring focused question answering. However, a critical and increasingly relevant application lies in open-ended conversational settings, where the goal is not to retrieve a specific answer but to engage in a natural, extended dialogue. This includes applications such as companionship and exploratory discussion, as well as complex business interactions such as customer service; these applications all require a shift from retrieval to genuine interaction.

Evaluating LLMs in these longer-form, open-domain conversations presents a significant challenge. Existing evaluation metrics (Rein et al., 2023; Hendrycks et al., 2021), largely inherited from tasks like question answering and machine translation, often focus on factual correctness, precision, and recall against predefined targets. These metrics prove inadequate when assessing the more nuanced qualities of a successful conversation, such as coherence, engagingness, personality consistency, and the ability to maintain context over multiple turns. Simply put, current metrics struggle to capture what makes a conversation feel natural and satisfying: the very qualities that define a positive human-LLM interaction and especially in a longer dialogue setting.

This inadequacy stems from the differences between conversations in the two length regimes. Unlike question answering, where success is defined by matching a ground truth, a 'good' conversation is subjective and multifaceted. It requires the LLM to demonstrate not just linguistic competence, but also social intelligence and an understanding of conversational dynamics. Existing metrics fail to model dialogue dynamics, human subjectivity, or align with goals like depth or reasonableness.

To address this gap, we reformulate the problem of evaluating long conversations as a *consistency measure* between the k th turn of a conversation with respect to the $1, 2, \dots, k - 1$ th turns within an $n \geq k$ turn conversation. Rather than traditional measures of information accuracy, this consistency lens focuses our attention onto measures like long-term coherence and reasonableness, which are precisely the under-represented measures that are instrumental in the longer conversational domain.

UPHELD (Utility & Planning Human Evaluated Long Dialogues) contains over 10,000 high-quality human annotations across hundreds of novel long-length conversations focused on useful task

054 completions (e.g. math tutoring, trip planning, etc.). Each assistant turn in these conversations is
055 labeled for coherence in content and style, as well as general reasonableness for a variety of different
056 popular models responses. We will demonstrate that standard metrics perform poorly on the UPHELD
057 dataset, and that both fine-tuning on UPHELD conversations and taking simple ensembles of standard
058 metrics leads to large improvements. These results show that not only is UPHELD revealing a crucial
059 evaluation modality (i.e. longer conversations) that challenges the standard methods, but also that
060 it provides meaningful vectors by which to improve these this important modality. Our primary
061 contributions are:

- 062 1) We present a dataset of human-evaluated long-form, non-topical conversations. UPHELD
063 introduces difficult benchmarks such as content and style consistency that allow us to assess
064 any proposed evaluation metric for its ability to align with subjective human opinions of
065 natural dialogue interactions.
- 066 2) We apply UPHELD to show that common LLM evaluation metrics perform poorly within
067 the long dialogue regime.
- 068 3) We demonstrate that ensembled machine-learned metrics trained on UPHELD exhibit
069 stronger correlation with human judgments than most commonly used LLM automatic
070 evaluation metrics. This result shows that our dataset enables meaningful exploration of
071 language model improvements in our proposed setting (long conversation).
- 072 4) We show that fine-tuning a model for long-form conversations improves performance to
073 surpass state-of-the-art results in this context.

074 2 RELATED WORK

075 Evaluating the performance of large language models (LLMs) is a critical area of research. Current
076 benchmarks, such as those on the Open LLM leaderboard (now archived) (Myrzakhan et al., 2024),
077 often focus on task-based evaluations using input-output metrics. These metrics include task-based
078 generation evaluations like IFEval (Zhou et al., 2023) which evaluates the ability of LLMs to follow
079 instructions and BBH (Big Bench Hard (Srivastava et al., 2022)), which comprises 23 groups of
080 tasks such as *word sorting*, *casual judgment*, *navigate*. This and similar task-with-answer evaluations
081 utilize metrics such as Exact Match (EM), Precision/Recall/F1, and sub-string based accuracy.
082 Specific examples include MATH for mathematical problems (Hendrycks et al., 2021), and GPQA
083 for multi-choice question answering (Rein et al., 2023).

084 Question answering (QA) benchmarks are also prevalent. MUSR (Sprague et al., 2023) is a narrative-
085 based QA evaluation where the input is a paragraph and the output is an answer with an evaluation
086 score. MMLU-PRO (Wang et al., 2024) combines task-based evaluation with QA and chain-of-
087 thought prompting. Recently, using an LLM as a judge for evaluation is another approach which
088 enables evaluating without having a ground truth or a reference answer (Clark et al., 2023) (Zheng
089 et al., 2023), (Duan et al., 2023), (Adlakha et al., 2024). While the aforementioned benchmarks
090 are prominent, several related areas of work are not typically included in leaderboards focused on
091 single-turn evaluations, particularly those concerning longer conversations. These include approaches
092 similar to Chatbot Arena (Zheng et al., 2023), MultiHop QA (Schnitzler et al., 2024) which involves
093 answering questions over multiple documents or turns, and MT-eval which focuses on evaluating
094 conversations directly and multi-dimensionally (Kwan et al., 2024). BotChat (Duan et al., 2023) is
095 also noted as a method for conversation generation evaluation.

096 The above works often rely on input-output pairs, exact matches, or scores tied to specific correct
097 answers, which are well-suited for evaluating ability to perform discrete tasks or extract factual
098 information. However, they may not be appropriate for evaluating long, non-topical conversations
099 because such interactions involve sustained coherence, context management over many turns, en-
100 gaging dialogue, and the ability to handle subjective or open-ended discussions that do not have a
101 clear and concise "correct" answer. Current long-form datasets focusing on non-factual conversations,
102 including MuTal (Cui et al., 2020), Topical-chat (Gopalakrishnan et al., 2019), LLM-arena (Zheng
103 et al., 2023), and DailyDialogue (Multi-turn) (Li et al., 2017), often lack authentic human-to-human
104 interaction or are constrained to predefined topics. Recognizing these limitations, and to construct
105 a robust benchmark, we identified the need for a dataset characterized by comprehensive curation,
106
107

108 expert authorship, and human annotation. This paper introduces such a dataset, which we have
109 meticulously collected and organized, with details provided in the ensuing sections.
110

111 3 DATASETS

112 3.1 UPHELD DATASET

113
114 We first collected conversational data by tasking a panel of professional writers hired from the
115 Upwork freelancer marketplace to create open-ended, naturalistic dialogues across a diverse range
116 of topics like customer service and education. The writers were explicitly instructed to create
117 longer dialogues characterized by conversational complexity and exploratory interaction, rather than
118 producing conversations with predetermined or narrow outcomes. For example, the writers developed
119 dialogues exploring scenarios such as mobile phone selection, travel planning, or problem-solving
120 discussions that were goal-oriented and where multiple nuanced exchanges would be appropriate.
121 Writers were also instructed to avoid overly scripted or linear communication patterns. This approach
122 ensures that the dataset captures the diverse and intricate nonlinear nature of human conversation. We
123 also took steps to alleviate any potential writer bias (see details in Appendix B). The UPHELD dataset
124 is publicly available under a permissive CC-BY-4.0 license and an anonymized copy is provided in
125 the supplemental materials of this submission, with the public download URL temporarily withheld
126 to preserve anonymity during the review process.
127
128

129 3.1.1 INPUT DATA

130
131 Given our initial set of rich natural language conversations, various LLM models were then used
132 to output candidate completions at every level of every conversation. Specifically, models were
133 presented with conversation history up to a specific point, with the next human-written turn withheld.
134 Models then generated a predicted next turn. The following models were used to generate the
135 dataset: a) **GPT-3.5** (OpenAI, 2023); b) **GPT-4o** (OpenAI, 2024); c) **Llama-3.1-70b** (Grattafiori
136 et al., 2024); d) **custom model** (a fine-tuned version of Llama-3.1: see Section 5.1). In addition to
137 these, we also added two baseline candidate completions: a) a version of the ground truth (human)
138 output rephrased using GPT-4o and b) a random completion. These two additional benchmarks offer
139 important opportunities to calibrate our data, and as expected we found that the rephrased dataset
140 received the highest marks while the random dataset received the lowest marks when evaluated by a
141 human labeler. These checks provide us with additional confidence that humans are fair judges of our
142 evaluation axes – style, content, and reasonableness.

143 At every possible conversational turn, annotators were then given (1) the chat history up to that point,
144 (2) the ground truth human completion (Option A), and (3) one of the LLM-generated candidate
145 completions (Option B). The provenance of options A and B were not disclosed to the annotators.
146 Their task was to compare these continuations on content consistency, style consistency, and general
147 reasonableness. The tasks, metric scales and instructions were fine tuned through (paid) two pilot user
148 studies. Example data-points within UPHELD are given in Table 1, and full instructions provided to
149 the annotators can be found in Appendix A.

150 3.1.2 ANNOTATION DIMENSIONS

151
152 Across all candidate models, each model-generated turn was evaluated by five independent human
153 annotators using the following scoring criteria (see full instructions in Appendix A).
154

- 155 • Content Equivalence: A 5-point Likert scale (1-5) assessing semantic similarity between
156 generated and reference turns.
- 157
- 158 • Style Equivalence: A 3-point Likert scale (1-3) measuring linguistic and stylistic correspon-
159 dence.
- 160
- 161 • Utterance Reasonableness: A binary scale (1: unreasonable, 5: reasonable) evaluating the
contextual appropriateness and coherence of the generated response,

162 163 164 165 166 167 168 169 170 171 172 173 174 175 176	chat_history	<i>user</i> : I'm shopping for a new car and can't decide what to get. <i>assistant</i> : I can help with that. What size car are you looking for?	<i>user</i> : I'm really frustrated. I've been trying to get this printer to work all morning and it just won't print anything! <i>assistant</i> : I understand that must be really annoying. Let's work together to fix this. Can you tell me if you see any error messages on your printer or computer? <i>user</i> : No error messages. It just sits there doing nothing whenever I try to print! <i>assistant</i> : Alright, let's try a few steps to see if we can get it responding. First, could you check whether the printer is set as the default printer on your computer? <i>user</i> : How do I check that?
177 178 179 180 181	Option A	The Ford F-150 is one of the top selling trucks in the country. Have you taken one of those for a test drive?	If you're using a Windows PC, you can go to 'Devices and Printers' in your control panel. For Mac, you can check under 'System Preferences' and then 'Printers & Scanners'.
182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197	Option B	Trucks are great for hauling and towing. Here are a few questions to help narrow down your options: 1. What will you be using the truck for most often? (e.g., work, off-roading, daily driving, towing a trailer) 2. How many passengers will the truck need to seat? 3. Do you have a preferred fuel type (gasoline, diesel, hybrid)? 4. What's your budget for the truck? 5. Do you have a specific brand or model in mind? 6. Do you need any specific features, such as 4-wheel drive, a certain bed size, or advanced safety features	It's easy. On a Windows computer, you can go to the Control Panel, then Devices and Printers (or Printers and Devices in some versions). Look for your printer's name and see if it has a checkmark next to it that says "Default printer". If it doesn't, right-click on the printer and select "Set as default printer". On a Mac, you can go to System Preferences, then Printers & Scanners. Look for your printer's name and make sure it's selected as the default printer. If you're still having trouble, we can try some other steps. Have you tried restarting the printer
198 199	score_task_1_content	1: Strongly Disagree – no content overlap	3: Neutral - some overlap of content
200	score_task_2_style	1: Different styles	1: Different styles
201	score_task_3_reasonableness	1: Not a reasonable continuation	5: A reasonable continuation

Table 1: Input conversations and associated annotator ratings from a single annotator (all scores are out of 5). More examples can be found in Appendix K.

3.1.3 DATA STATISTICS

In total we collected complete evaluation labels for 400 turns stemming from 53 human written conversations collected for this study and additional seven conversations from existing datasets (Cui et al., 2020) to serve as control points. The dataset was curated to include approximately 90% human-written conversations on a wide range of topics and situations, supplemented with approximately 10% from available datasets to introduce challenging cases.

Each predicted turn was evaluated by five human annotators, and each annotator judge labeled between 1,220 and 1,230 conversations. We overall generated 12,291 sets of labels, or 36,873 labels. The ground truth conversations consist of 5.2 turn pairs (user-assistant) or 10.4 dialogue turns on

average. The average length of the conversation history annotators analyzed was 560 characters and the length of the judged turns was on average 245 characters.

3.1.4 VERIFICATION DATASETS

To validate our experiments, we constructed an independent verification datasets by augmenting LLM-Arena (Zheng et al., 2023) and Topical-Chat (Gopalakrishnan et al., 2019). The overall procedure consisted of three steps: (i) deriving a single “ground truth answer” from each data point of each existing dataset (see below), (ii) generating an alternative continuation with GPT-4o, and (iii) collecting human judgments following the UPHELD annotation protocol. We include all additional verification labels within our dataset for reproducibility.

Derivation of Ground Truth:

a. LLM-Arena:

- Data Point Description: each dialogue is accompanied by two model-generated continuations plus a human preference label.
- Ground Truth Extraction: for every conversation we enumerated all candidate pairs, tallied human preferences, and chose the majority-preferred continuation as the reference (provided a clear winner existed).

b. Topical-Chat:

- Data Point Description: human-to-human dialogues grounded in specific topic selections.
- Ground Truth Extraction: to mitigate cold-start artifacts, we extracted segments spanning turns 5–7. We then treated turns 1–5 (or 1–7) as the model input and selected the next human turn (turn 6 or 8) as the ground truth continuation.

After standard quality control (i.e. filtering for missing data and badly formatted inputs), we obtained 12,305 pairwise preference judgments. We note that both these verification datasets, although useful for verification, are still relatively deficient in freeform human-to-human interaction and focus on a limited set of pre-defined topics. As such they should be treated as verification datasets only and not as valid replacements for UPHELD.

4 METRICS

In total, 12 candidate metrics were assessed for their ability to evaluate longer conversations via correlation to UPHELD labels. These metrics were grouped into 3 distinct groups: 1) token-based – metrics quantifying similarity based on exact overlap of tokens 2) semantic-based – metrics quantifying similarity based on semantic overlap (e.g. embedding models); and 3) LLM-based – metrics employing some form of the LLM-as-a-judge paradigm.

Recall-Oriented Understudy for Gisting Evaluation (ROUGE (Lin, 2004)) is a set of standard language metrics that compares an automatically produced summary or translation against a set of reference summaries or translations. Specifically, ROUGE-N measures the overlap of n-grams between the system-generated text and the reference text. ROUGE-L measures the longest common subsequence, which accounts for sentence-level structure similarity.

We also explored cosine similarity between message embeddings as a measure of semantic similarity between the generated text and reference text. This approach is rooted in the work by Reimers & Gurevych (2019) on Sentence-BERT embeddings, which have shown effectiveness in capturing semantic similarities in text data. BERTScore leverages the pre-trained contextual embeddings from BERT to evaluate text generation by matching words in candidate and reference sentences. It computes precision, recall, and F1 score, providing a more nuanced evaluation than traditional n-gram based metrics. Zhang et al. (2020) introduced BERTScore as a robust metric for evaluating generated text.

We also tested LLM-as-a-judge metrics (Clark et al., 2023) against UPHELD. This approach involves using a separate LLM to score the outputs based on various criteria, such as coherence, relevance, and overall quality. We used both binary (yes/no) and Likert scale (1-5) judgments, with and without explanations. Prompts for judges can be found in Appendix F.

270 4.1 REFERENCE-FREE VS HUMAN-GROUNDED METRICS

271
272 Fundamentally, UPHELD uses a ground truth reference to generate our human labels. In contrast,
273 reference-free evaluation of LLMs Liu et al. (2023) rely on human preference, and LLM outputs
274 can reliably reproduce these preferences, indicating their performance is consistent with human
275 judgments (Zheng et al., 2023). However, reference-free preference datasets also incur significant
276 limitations, such as poor performance at judging long task-oriented dialogues and weakened judgment
277 reliability on out-of-distribution tasks (Saito & Sato, 2025). UPHELD is designed specifically to
278 tackle these limitations by adopting a reference-full approach.

279 However, the reference-full approach raises a question around dialogue multiplicity: a single input
280 might incur multiple valid outputs, so how are we sure our ground truths are well defined? UPHELD’s
281 design reduces susceptibility to this issue in two ways: (1) two of our key label categories (style
282 and reasonableness) are well-defined even with dialogue multiplicity, and (2) UPHELD dialogues
283 primarily revolve around task-oriented settings, which means content accuracy is a well-defined
284 metric. For example, while opinion-oriented conversations (e.g. *Who makes the best Caesar salad?*)
285 are susceptible to dialogue multiplicity, our task-oriented dialogues (e.g. *How to make a Caesar*
286 *salad?*) are not. To quantify this effect, we ran an experiment detailed in Appendix C and showed
287 that UPHELD dialogues admit significantly higher output consistency (93%) compared to freeform
288 dialogues (74%), which shows that our reference-full approach still allows us to collect meaningful
289 labels on ground-truth content overlap.

290 5 EXPERIMENTS

291
292 To demonstrate the adds value to long-form conversational evaluation, we present a series of exper-
293 iments. We show: a) that the base conversations within UPHELD can be used to greatly increase
294 conversational fidelity within an LLM; b) that naive evaluation metrics fail in the long conversations,
295 and propose a simple ensemble metric which greatly outperforms the baseline metrics. These ensem-
296 ble metrics also perform well on our validation datasets (Topical and LLM Arena), showing that a
297 method developed with UPHELD is transferrable to other contexts. We also include discussion on
298 user disagreement within UPHELD.

300 5.1 UPHELD AS A FINE-TUNING DATASET

301
302 An important way to validate the quality of data within the UPHELD dataset is to assess how the
303 UPHELD scores differ between the *custom* model, which was fine-tuned on the base data, and the
304 other baseline models. To do so, we directly plot the mean content, style, and reasonableness scores
305 within the UPHELD dataset for the different models in Figure 1. Base models (GPT-3.5, GPT-4o,
306 and Llama3.1-70b-base) exhibit significantly lower performance in both content accuracy and style
307 accuracy, with mean scores below 0.5. This indicates a tendency for these models to deviate from the
308 intended conversational style and introduce content inconsistencies. Specifically, the base models
309 demonstrate a substantial deficit in maintaining the stylistic integrity and topical coherence of the
310 conversation, suggesting a propensity to diverge from the trajectory established by human writers.

311 In contrast, the custom model, fine-tuned for extended dialogue on a held-out set, shows a marked
312 improvement and achieves approximately a 40% increase in both content and style accuracy compared
313 to the base models. It does this while retaining the core functionality of its model base (Llama-
314 3.1-70b), which we validated at test time and also at training time where we observed minimal
315 overfitting within loss curves. Reasonableness is fairly flat across all models, which is not surprising
316 as LLMs tend to output reasonable results regardless of style or content consistency. These results
317 highlight that the UPHELD dataset’s conversations are both learnable and encode behavior that are
318 not well-exposed within these models’ pre-training datasets.

320 5.2 EVALUATION METRIC PERFORMANCE ON UPHELD

321
322 Given that UPHELD is designed to help develop novel evaluation metrics for long conversations, it is
323 instructive to see how this development works in practice. We start this section with an analysis on
how traditional metrics perform poorly on UPHELD, and how simple modifications to the traditional

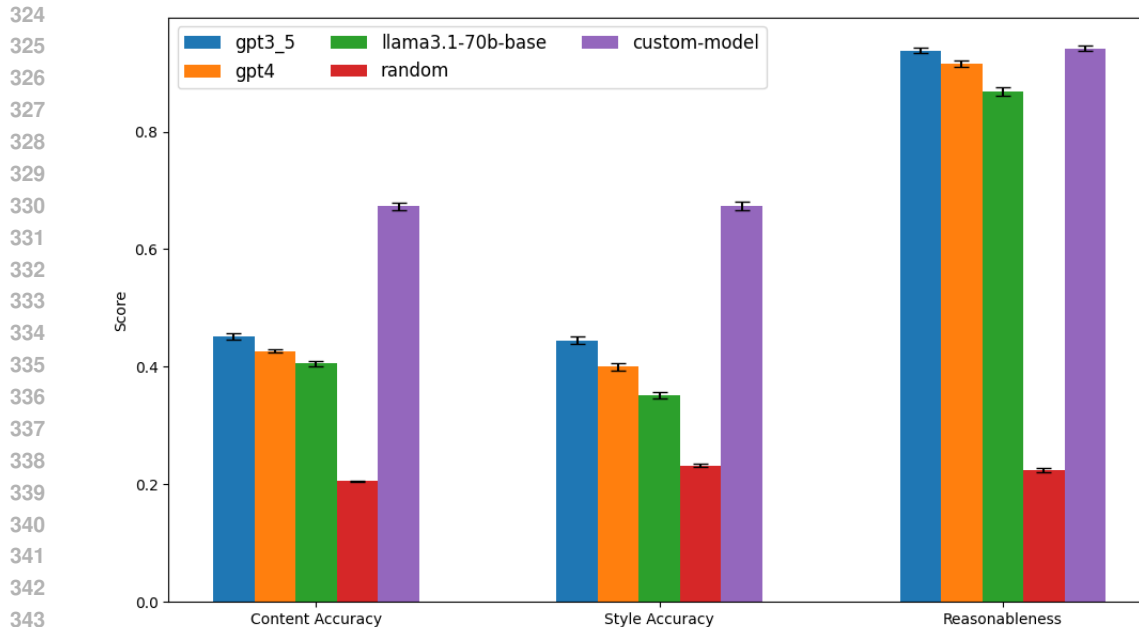


Figure 1: Aggregate human labeler scores as a share of the perfect score (see Appendix D.1) for each model on the UPHELD dataset within each label category. GPT-4 rephrase performs the best because it is a mild perturbation of the real ground truth. Of the other models, the custom model, which was fine-tuned on a held-out set of training data within the UPHELD dataset, performs significantly better than all other baseline models, demonstrating that the information to be learned within UPHELD conversations is significant.

metrics provide a significant boost in performance. All results in this section are standard correlation coefficients presented as 5-fold cross validation results on a 20% held-out set.

5.2.1 TRADITIONAL METRICS

The main results for traditional metrics (as defined in Section 4) are shown in the first three rows of Tables 2 and 3. The results reveal a weak to moderate correlation between metrics and human ratings. This suggests that no single metric captures the nuances of human assessment well within the UPHELD dataset. All definitions for metrics can be found in the Appendix D.2. These results reveal a weak to moderate results for these metrics. This suggests that no single traditional metric captures the nuances of human assessment well within the UPHELD dataset. Interestingly, the semantic metrics (like bert_score_F1) demonstrated the highest correlation with human judgments for both content and style aspects. The token-based metrics showed stronger correlations than LLM-as-a-judge as well. These findings suggest that LLM-as-a-judge, despite being increasingly explored in the literature (e.g., (Ma & Graham, 2019) (Zhang et al., 2020)), is a weak evaluator of longer conversations.

In spite of the observed strength of the semantic metrics, the results for traditional metrics show weak to moderate performance on UPHELD. This implies that relying solely on any single traditional metric inadequately captures the complexities of content and style quality in longer conversations, and we demonstrate that ensemble metrics can bridge this gap.

5.2.2 ENSEMBLE METRICS

We hypothesize that individual metrics attend to distinct facets of text quality, and so a learned ensemble will perform better on the UPHELD dataset. We train linear regression, SVM, and random forest models to predict human scores using the individual automatic metric scores as input features. Using this we derive a new, hybrid metrics and assess their correlation with human judgments.

The last 3 rows of Tables 2 and 3 show that the learned ensemble metrics exhibit higher correlations with human scores compared to any single metric in isolation. Notably, the random forest model

Metric	UPHELD	LLM Arena	Topical Chat
Semantic Metrics			
message_embedding_cos_sim	0.58±0.003	0.41±0.004	0.40±0.008
bert_score_precision	0.36±0.005	0.39±0.004	0.32±0.008
bert_score_recall	0.40±0.005	0.42±0.004	0.40±0.008
bert_score_F1	0.40±0.005	0.45±0.003	0.42±0.008
Token-based Metrics			
rouge1	0.34±0.005	0.49±0.003	0.33±0.009
rouge2	0.24±0.005	0.43±0.003	0.24±0.009
rougeL	0.32±0.005	0.41±0.004	0.34±0.010
rougeLsum	0.31±0.005	0.50±0.003	0.34±0.009
LLM-as-a-judge metrics			
llm_judge_yes_no	0.40±0.012	0.32±0.015	0.57 ±0.016
llm_judge_yes_no_explain	0.23±0.018	0.32±0.014	0.29±0.013
llm_judge_likert_1_5	0.26±0.008	0.41±0.011	0.39±0.014
llm_judge_likert_1_5_explain	0.28±0.009	0.41±0.013	0.34±0.013
Ensembled ML Metrics (Ours)			
Linear Regression	0.59±0.004	0.50 ±0.003	0.58 ±0.006
SVM	0.55±0.005	0.49±0.003	0.54±0.008
Random Forest	0.63 ±0.004	0.23±0.004	0.32±0.009

Table 2: Results on the UPHELD datasets for various candidate evaluation metrics on the first evaluation task (content accuracy). Ensemble metrics clearly perform better for both the UPHELD dataset and both verification datasets.

Metric	UPHELD	LLM Arena	Topical Chat
Semantic Metrics			
message_embedding_cos_sim	0.41± <0.001	0.42±0.004	0.29±0.009
bert_score_precision	0.37±0.005	0.39±0.004	0.36±0.008
bert_score_recall	0.36±0.005	0.50±0.004	0.27±0.009
bert_score_F1	0.41±0.005	0.49±0.003	0.35±0.008
Token-based Metrics			
rouge1	0.32±0.005	0.51±0.003	0.31±0.010
rouge2	0.20±0.005	0.43±0.004	0.16±0.010
rougeL	0.29±0.005	0.48±0.004	0.30±0.009
rougeLsum	0.29±0.005	0.48±0.003	0.30±0.009
LLM-as-a-judge metrics			
llm_judge_yes_no	0.35±0.012	0.08±0.010	0.36±0.016
llm_judge_yes_no_explain	0.15±0.015	0.08±0.010	0.24±0.013
llm_judge_likert_1_5	0.24±0.010	0.07±0.009	0.22±0.014
llm_judge_likert_1_5_explain	0.20±0.012	0.05±0.008	0.20±0.013
Ensembled ML Metrics (Ours)			
Linear Regression	0.50±0.004	0.28±0.004	0.32±0.008
SVM	0.49±0.005	0.54 ±0.003	0.55 ±0.009
Random Forest	0.61 ±0.004	0.04±0.004	0.22±0.009

Table 3: Results on the UPHELD datasets for various candidate evaluation metrics on the second evaluation task (style accuracy). Ensemble metrics clearly perform better for both the UPHELD dataset and both verification datasets.

yields substantial improvements of 30-40% relative to the best-performing individual metrics on UPHELD, but an SVM ensemble produces consistently higher correlations to human judges across all datasets. Note that the ensemble metrics were trained only on the UPHELD data and then applied to the LLM Arena and Topical Chat data. Although the transferability of the random forest model was

models:	gpt3_5	gpt4	gpt4_rephrase	llama3.1-70b-base	random	custom-model
Kappa score	0.30	0.33	0.46	0.38	0.97	0.35

Table 4: Kappa scores for different models. A Kappa score ranges from -1 to +1, with +1 indicating the maximum level of agreement among the annotators.

poor, the SVM and linear regression ensembles indicate that remarkably, trained metrics developed just on the UPHELD dataset have exceptional transferability to other out-of-domain datasets.

The success of the ensemble metrics likely stems from their ability to integrate diverse signals captured by the individual metrics, mirroring the multifaceted nature of human evaluation. These findings strongly suggest that within complex settings like long dialogue, learning to ensemble multiple automatic metrics offers a promising avenue for developing evaluation frameworks that more closely align with human judgments than relying on a single metric.

5.3 DISCUSSION AROUND ANNOTATOR AGREEMENT

Because each conversation turn was independently labeled by five human annotators, we also analyze model responses that induced stronger human-human agreements. We observed moderate levels of consensus as shown in Table 4, which tabulates Kappa scores (Cohen, 1960) between raters. As expected, human-human agreement is highest on the random and gpt4_rephrase baselines, while being consistent across all other models. Further analysis on labeler agreement can be found in Appendix G.

We observed that approximately 25% of data points had full agreement across all 5 human judges. Otherwise, we bin the level of agreement as follows: agreement across 2 out of 5 labels represent a “plurality,” while agreement across 3 or 4 labels represent a “majority.” We further quantify how well an LLM-as-a-judge evaluator agrees with this winning score relative to the agreement bin. Table 5 shows that as expected, LLM-as-a-judge performance is heavily correlated to agreement level amongst the human labelers. This result demonstrates that human-human disagreement is a valid measure of data difficulty, and this uncertainty signal present in UPHELD may be integral in further evaluation metric development.

	LLM-as-a-judge				
	human	likert	likert explain	binary	binary explain
plurality	0.4733	0.3401	0.3696	0.7493	0.7196
majority	0.6932	0.3885	0.3976	0.8011	0.7674
full agreement	1.0000	0.6370	0.5603	0.9339	0.9195

Table 5: LLM-as-a-judge performance on the UPHELD dataset for different dataset splits separated by the level of human annotator agreement within each bucket.

6 CONCLUSION

In this work, we introduced UPHELD: a dataset designed to evaluate LLMs in long-form conversational settings. We collected tens of thousands high-quality human-annotated labels on crucial consistency metrics within long-form conversations, along with the high-quality conversations themselves. Analysis of existing evaluation metrics on UPHELD reveals that they do not effectively capture the nuances of human judgment for assessing conversational quality. We further demonstrated that simple-to-learn ensemble metrics result in substantially improved correlations with human evaluations. Taken in aggregate, our findings highlight the potential for developing robust evaluation frameworks that better align with human perceptions of effective conversation using UPHELD.

7 REPRODUCIBILITY STATEMENT

We provide the entire UPHELD dataset that was used for this work in order for any other researcher or practitioner to reproduce our work. Otherwise, all prompts used to collect our dataset are provided in Appendix F, and exact definitions of our baselines and models used are provided in Section 4 and 5. Explicit formulae to ensure baseline reproducibility are provided in Section D.2.

REFERENCES

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. Evaluating correctness and faithfulness of instruction-following models for question answering. *Transactions of the Association for Computational Linguistics*, 12:681–699, 2024. doi: 10.1162/tacl_a_00667. URL <https://aclanthology.org/2024.tacl-1.38/>.
- Kenton Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. Llm as a judge: Evaluating text generation with large language models. *arXiv preprint arXiv:2306.05685*, 2023. URL <https://arxiv.org/abs/2306.05685>.
- Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960.
- Lili Cui, Yu Wu, Shiyue Liu, Yizhe Zhang, and Ming Zhou. Mutual: A dataset for multi-turn dialogue reasoning. *arXiv preprint arXiv:2004.04494*, 2020. URL <https://arxiv.org/abs/2004.04494>.
- Haodong Duan, Jueqi Wei, Chonghua Wang, Hongwei Liu, Yixiao Fang, Songyang Zhang, . . . , and Kai Chen. Botchat: Evaluating llms’ capabilities of having multi-turn dialogues. *arXiv preprint arXiv:2310.13650*, 2023. URL <https://arxiv.org/abs/2310.13650>.
- Karthik Gopalakrishnan, Baran Hedayatnia, Qing Chen, Alessandro Gottardi, Shubham Kwatra, Ananth Venkatesh, Richard Gabriel, and Deniz Hakkani-Tür. Topical-chat: Towards knowledge-grounded open-domain conversations. In *Interspeech 2019*, pp. 1891–1895, Graz, Austria, 2019. ISCA / Interspeech.
- Alessio Grattafiori, Amit Dubey, Arjun Jauhri, Ankush Pandey, Aniruddha Kadian, Abdullah Al-Dahle, . . . , and Predrag Vasic. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Dan Hendrycks, Collin Burns, Sanjana Kadavath, Amanpreet Arora, Swayambhu Basart, Eric Tang, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021. URL <https://arxiv.org/abs/2103.03874>.
- W. C. Kwan, Xiang Zeng, Yang Jiang, Yanhao Wang, Lei Li, Lei Shang, and Kam-Fai Wong. Mteval: A multi-turn capabilities evaluation benchmark for large language models. *arXiv preprint arXiv:2401.16745*, 2024. URL <https://arxiv.org/abs/2401.16745>.
- Yankai Li, Hang Su, Xiaodong Shen, Wenhao Li, Zheng Cao, and Shijin Niu. Dailydialog: A manually labelled multi-turn dialogue dataset. *arXiv preprint arXiv:1710.03957*, 2017. URL <https://arxiv.org/abs/1710.03957>.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81. Association for Computational Linguistics, 2004.
- Yang Liu, Yuhui Xu, Daniel Kang, Chenyan Zhu, Chenguang Xu, Yizhan Zhang, Bo Wang, Heng Ji, and Luke Zettlemoyer. G-eval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023. URL <https://arxiv.org/abs/2303.16634>.
- Wei Ji Bojar Ondřej Ma, Q. and Yvette Graham. Results of the wmt19 metrics shared task: Segment-level and strong mt systems pose big challenges. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pp. 62–90, Florence, Italy, 2019. Association for Computational Linguistics.

- 540 Aidar Myrzakhan, Sondos Mahmoud Bsharat, and Zhiqiang Shen. Open-llm-leaderboard: From
541 multi-choice to open-style questions for llms evaluation, benchmark, and arena. *arXiv preprint*
542 *arXiv:2406.07545*, 2024. URL <https://arxiv.org/abs/2406.07545>.
- 543 OpenAI. Chatgpt (gpt-3.5) [large language model]. <https://openai.com/chatgpt>, 2023.
- 544 OpenAI. Gpt-4o system card. <https://openai.com/index/gpt-4o-system-card/>,
545 2024.
- 546 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks.
547 *arXiv preprint arXiv:1908.10084*, 2019. URL <https://arxiv.org/abs/1908.10084>.
- 548 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
549 Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a
550 benchmark. *arXiv preprint arXiv:2311.12022*, 2023. doi: 10.48550/arXiv.2311.12022. URL
551 <https://arxiv.org/abs/2311.12022>.
- 552 Kazuki Saito and Masashi Sato. No free labels: Limitations of reference-free llm evaluation. *arXiv*
553 *preprint arXiv:2503.05061*, 2025. URL <https://arxiv.org/abs/2503.05061>.
- 554 Julian Schnitzler, Xanh Ho, Jiahao Huang, Florian Boudin, Saku Sugawara, and Akiko Aizawa.
555 Morehopqa: More than multi-hop reasoning. *arXiv preprint arXiv:2406.13397*, 2024. URL
556 <https://arxiv.org/abs/2406.13397>.
- 557 Zachary Sprague, Xinyi Ye, Kristian Bostrom, Sanjoy Chaudhuri, and Greg Durrett. Musr: Testing
558 the limits of chain-of-thought with multistep soft reasoning. *arXiv preprint arXiv:2310.16049*,
559 2023. URL <https://arxiv.org/abs/2310.16049>.
- 560 Amar Srivastava, Abhinav Rastogi, Arjun Rao, Amir A. M. Shoeb, Awni Abid, Adam Fisch, and
561 George Wang. Beyond the imitation game: Quantifying and extrapolating the capabilities of
562 language models. *arXiv preprint arXiv:2206.04615*, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2206.04615)
563 [2206.04615](https://arxiv.org/abs/2206.04615).
- 564 Yuhui Wang, Xuechen Ma, Guoqing Zhang, Yang Ni, Abhinav Chandra, Shiyu Guo, and Wenhan
565 Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark.
566 In *NeurIPS Datasets and Benchmarks Track*, 2024.
- 567 Tianyi Zhang, Vikas Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating
568 text generation with bert. *arXiv preprint arXiv:1904.09675*, 2020. URL <https://arxiv.org/abs/1904.09675>.
- 569 Luyu Zheng, Wei-Lin Chiang, Yuyang Sheng, Shiyu Zhuang, Zifan Wu, Yuchen Zhuang, Zhihan
570 Lin, Zhenyu Li, Derek Li, Eric P. Xing, Han Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging
571 llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023. URL
572 <https://arxiv.org/abs/2306.05685>.
- 573 Jingyi Zhou, Tianyi Lu, Sam Mishra, Sanjoy Brahma, Sayan Basu, Yuyu Luan, and Lei Hou.
574 Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*,
575 2023. URL <https://arxiv.org/abs/2311.07911>.
- 576
- 577
- 578
- 579
- 580
- 581
- 582
- 583
- 584
- 585
- 586
- 587
- 588
- 589
- 590
- 591
- 592
- 593

A ANNOTATION MATERIALS AND INSTRUCTIONS

Here we provide a complete description of the materials provided to data annotators along with associated instructions.

A.1 MATERIALS

Users are given access to a spreadsheet file with three sheets – the first sheet contains the instructions, the second contains sample annotated examples and the third one contains a formatted table with 6 columns with the following labels:

- chat history – (the conversation up to one point)
- Option A – (one possible continuation to the conversation within chat history)
- Option B – (alternative possible continuation to the conversation within chat history)
- score_task_1_content – (a dropdown menu to select the content consistency score)
- score_task_2_style (a dropdown menu to select the style consistency score)
- score_task_3_reasonableness (a dropdown menu for choosing the reasonableness score)

A screenshot of the interface annotators are given is shown in Figure 2.

A.2 INSTRUCTIONS

This section contains the verbatim instructions provided to the annotators in the data annotation spreadsheet. We start with the overall instructions and then reprint the granular instructions for each of the three types of labels within UPHELD.

==BEGIN INSTRUCTIONS==

Given a chat history, you will be presented two options for how to continue the conversation: Option A and Option B. You will be asked to rate these options by answering a number of questions.

Task 1 (Content Equivalence): Do you agree that the general information presented in Option A roughly the same as the general information presented in Option B?

Task 2 (Stylistic Equivalence): Do you agree that the style of Option A the same as the style of Option B? Put another way, does it feel like Option A and Option B are being spoken by the same person?

Task 3 (Reasonableness of Option B): For Option B only, do you agree that Option B a reasonable way to continue the conversation, given the chat history?

Enter your score in the column corresponding to the task in the annotations tab/sheet (e.g. score_task_2_style).

Please read below for specific instructions and tips on each individual task.

A.2.1 TASK 1 – CONTENT EQUIVALENCE

Check the examples tab for some already annotated data and additional explanation (note that you are not expected to provide explanations of your scores.)" "Provide one of the following scores on a scale of 1-5 where a 1 reflects a strong DISAGREE and a 5 reflects a strong AGREE:

- 1: Strongly Disagree (that the content conveys equal information in both options)
- 2: Disagree
- 3: Neutral
- 4: Agree
- 5: Strongly Agree

Use the following criteria to help you determine if the two message options have equivalent content:

- 648 • Information conveyed by Option B contains all information that is conveyed by Option A.
- 649
- 650 • Using either Option A or Option B to continue the conversation would not change the flow
- 651 of the conversation.
- 652 • You can replace Option A with Option B, or replace Option B with Option A, without
- 653 appreciably changing the content of the conversation.
- 654 • Both Option A and Option B mean the same thing.
- 655

656 Tips:

- 658 • Keep the chat history in mind when considering the content equivalence of Option A and
- 659 Option B.
- 660 • If one of the options seems incomplete or cut short, still try to evaluate the option as is.
- 661
- 662 • If Option B is wordier or contains more details than Option A, but it still contains all the
- 663 information in Option A and is relevant given the chat history, lower the score to at most a 3.
- 664 • Do not lower the score if Option B contains AI self identification phrases such as ("As an
- 665 AI agent...", "I am a trained model..") and similar. Focus on the other information within
- 666 Option B.
- 667 • If Option B is not readable and/or contains non-coherent language give a score of 3.
- 668
- 669 • Lower the score if Option B contains more details that are (1) not an expansion of the
- 670 information in Option A and (2) are not relevant to the messages in the chat history.
- 671

672 A.2.2 TASK 2 - STYLE EQUIVALENCE

673 Provide one of the following scores on a scale of 1-5 where a 1 reflects a strong DISAGREE that

674 styles are the same and a 5 reflects a strong AGREE that styles are the same:

675

- 676 • 1: Different styles (that the content conveys equal information in both options).
- 677
- 678 • 3: Somewhat same styles.
- 679
- 680 • 5: Same styles

681 Use the following instructions to help you determine if the two message options are stylistically

682 equivalent

683

- 684 • After reading them out loud, both options A and B feel like they are written by the same
- 685 person in the same mood.
- 686 • There is no noticeable change in sentiment or tone between the two options.
- 687
- 688 • Even if one of the message options is longer than the other, they can still be considered
- 689 stylistically similar if the content is expressed in similar ways.
- 690 • If it sounds like option A and option B were written by different people, assign a low score.
- 691
- 692 • If you believe that both options are written by the same person in the same mood but the
- 693 content of the two options are different, still assign a high score.

694 Tips:

- 696 • It may be useful to consider the context of chat history as a reference and seeing whether
- 697 either option deviates from a natural continuation of the chat history, given the personality
- 698 of the "assistant" in the chat history.
- 699 • For this task you're highly encouraged to read both options out loud as it may be helpful in
- 700 forming the comparison.
- 701 • Consider differences in vocabulary, tone, and syntax when making your decision.

702 A.2.3 TASK 3 – REASONABLENESS
703

704 Provide one of the following scores:

- 705 • 1: Not a reasonable continuation (to chat_history)
- 706 • 5: A reasonable continuation (to chat_history)
- 707
- 708

709 Guidelines:

- 710 • This task ONLY applies to Option B. The task is to determine whether Option B is a
- 711 reasonable way to continue the conversation from the chat history.
- 712 • Ignore Option A in your judgment; Option B may be completely different from Option A
- 713 but still score highly in this task as long as it is on topic.
- 714 • If Option B seems cut short assess the text up to the cutoff point.
- 715 • If Option B is not readable and/or sounds incoherent assign a score of
- 716
- 717

718 Tips:

- 719 • We encourage you to read the chat history out loud as well as the message in Option B
- 720 directly afterwards. If it sounds like a natural conversation flow out loud the score is likely a
- 721 high score.
- 722 • If you were the "user" in this scenario and received Option B as the next response, would
- 723 you be generally happy with the state of the conversation? If the answer is yes, the score is
- 724 likely 5. If not the score is likely 1.
- 725 • Do not lower your score if Option B contains any model self identification (e.g. As an AI
- 726 model....) but is still a viable continuation of chat history.
- 727 • All of the following reasons are valid for assigning a low score of 1:
- 728 * Option B is excessively wordy and/or provides too much information.
- 729 * Option B is incoherent. Option B seems random and gets off topic.
- 730 * Option B is excessively rude or aggressive.
- 731 * Option B has an inappropriate tone or uses inappropriate language.
- 732 * Option B does not add any additional helpful information to the conversation or prompt
- 733 the user to provide additional relevant information.
- 734 • Check the examples tab for some already annotated data and additional explanation (note
- 735 that you are not expected to provide explanations of your scores.)
- 736
- 737
- 738

739 ==END INSTRUCTIONS==

740
741 B NOTES ON WRITER AND ANNOTATOR SELECTION AND BIAS MITIGATION
742

743 Our contracted writers were required to have a high job success rate on Upwork and all were first
744 evaluated through a rigorous initial (paid) pilot phase where their written conversations were evaluated
745 by a professional user experience team for diversity and faithfulness to our prompts. Writers were
746 also selected from diverse professional backgrounds: we employed writers with backgrounds from
747 novel writing to education to copywriting. Prompts were selected for diversity of tasks and diverse
748 defined styles that had to adhere to a number of user personas and styles. All conversations were
749 quality-checked by a separate set of experienced proofreaders to explicitly ensure style diversity
750 and consistency. We also acknowledge that our current focus is primarily on English language
751 conversations, but also plan to eventually incorporate multilingual UPHELD additions.

752 All conversations were further quality-checked by another professional writer to ensure situational and
753 stylistic diversity. We were admittedly limited to English-speaking writers, which may introduce some
754 clustering of labeler backgrounds. Because each conversation went through multiple rounds of checks
755 from different professionals (including both user research professionals, other writing professionals,
and machine learning professionals) who were explicitly instructed to check for diversity and to

eliminate bias, we hope that any effects of geographical/linguistic clustering are mitigated by our rigorous process.

All data labelers also participated in an initial (paid) pilot program that was carefully evaluated by internal user research professionals before being selected to write conversations at scale. The scenarios the writers built were evaluated by the same user research professionals to ensure they covered a wide variety of scenario types and user behavioral/personality patterns which were representative of what chat agents might encounter in a customer-facing context.

C UPHELD AND CONVERSATIONAL MULTIPLICITY

A valid concern may be that direct comparisons to a reference human answer may be inappropriate in settings when a particular prefix can lead to a multiplicity of valid responses. This effect may be prevalent especially when the prefix is asking for an opinion (e.g. “What is your favorite animal?”). We however observe that UPHELD dialogues avoid this potential pitfall as they are not strictly freeform, but are all targeted towards completion of a specific well-defined task. In this context, there is some notion of correct ground truth, and we specifically hired professionals who are experts at these tasks (see Appendix B). To put it simply, our task setting is analogous to the difference between *what kind of salads do you like?* (which has ambiguity and dialogue multiplicity) and *how do I make a Caesar salad?* (which is much more constrained and has a more well-defined ground truth). To quantify this, we generated 100 open-ended questions (GPT-4o) and then generated two possible completions with GPT-4o at moderately high temperature ($\tau = 1$) to those questions. We did the same with 100 UPHELD turn completions. We then asked GPT-4o to judge whether the two possible completions contain similar content. The results are as follows:

Dataset	Semantic Consistency (%)
Open-Ended	74
UPHELD (first turn only)	93
UPHELD (random turn)	93

Table 6: Semantic Consistency Performance across Different Datasets

UPHELD induces much higher semantic consistency in the output, which means that UPHELD dialogues admit much less conversational multiplicity than more freeform datasets. This result supports our hypothesis that targeted task-focused conversations like the ones in UPHELD admit well-defined “ground truth” references. Interestingly, UPHELD maintains high output consistency even when we only analyze the first turn, which is where we would expect more branching/multiplicity during a dialogue. Note that these results are likely an underestimate of the true consistency, since sampling multiple LLM outputs would induce additional randomness that likely would not exist within natural human dialogue.

D METRICS

D.1 AGGREGATION METRIC

The aggregate scores in Figure 1 represent the total score of a given response divided by the maximum possible sum score. If 3 judges score a turn 5, 3, and 2 with a maximum score of 5, the aggregate score is $(5 + 3 + 2)/(5 + 5 + 5) = 10/15$. More formally, Let a response be scored by J judges. Judge j gives a score s_j with a per-judge maximum M_j (often all $M_j = M$).

$$\text{score} = \frac{\sum_{j=1}^J s_j}{\sum_{j=1}^J M_j}, \quad \text{where } 0 \leq \text{score} \leq 1$$

D.2 EXPERIMENTAL METRIC

We include all metrics used within the experimental studies tabulated in Tables 2 and 3.

- **Message Embedding Cosine Similarity:**

$$\text{CosineSim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

where \mathbf{u} and \mathbf{v} are the embedding vectors of the reference and generated messages.

- **BERTScore Precision:**

$$P = \frac{1}{|x|} \sum_{i=1}^{|x|} \max_j \text{sim}(x_i, y_j)$$

- **BERTScore Recall:**

$$R = \frac{1}{|y|} \sum_{j=1}^{|y|} \max_i \text{sim}(x_i, y_j)$$

- **BERTScore F1:**

$$F1 = \frac{2PR}{P + R}$$

where x and y are the sets of tokens from the candidate and reference texts respectively, and $\text{sim}(x_i, y_j)$ denotes cosine similarity between contextual embeddings of tokens x_i and y_j .

- **ROUGE-1 (Unigram Overlap):**

$$\text{ROUGE-1} = \frac{\sum_{w \in \text{Ref}} \min(\text{Count}_{\text{gen}}(w), \text{Count}_{\text{ref}}(w))}{\sum_{w \in \text{Ref}} \text{Count}_{\text{ref}}(w)}$$

- **ROUGE-2 (Bigram Overlap):**

$$\text{ROUGE-2} = \frac{\sum_{b \in \text{Ref}} \min(\text{Count}_{\text{gen}}(b), \text{Count}_{\text{ref}}(b))}{\sum_{b \in \text{Ref}} \text{Count}_{\text{ref}}(b)}$$

- **ROUGE-L (Longest Common Subsequence - LCS):**

$$\text{ROUGE-L} = \frac{\text{LCS}(X, Y)}{\text{Length}(Y)}$$

where X and Y are sequences of tokens in the generated and reference texts respectively.

- **ROUGE-Lsum (LCS over multiple sentences):**

$$\text{ROUGE-Lsum} = \frac{\sum_i \text{LCS}(X_i, Y_i)}{\sum_i \text{Length}(Y_i)}$$

where X_i and Y_i are sentence-level pairs from the candidate and reference summaries.

To compute all metrics, we used the ground truth next turn (*Option A*) as a reference data point and the model-generated next turn (*Option B*) as a candidate data point. We use the **mixedbread-ai/mxbai-embed-large-v1**¹ embedding model for all metrics that required a calculated similarity score.

All prompts associated with the *llm-as-a-judge* metrics can be found in Appendix F.

E HOW DID GPT-4O UNDERPERFORM GPT-3.5?

In Figure 1, we revealed a somewhat surprising result: humans tended to prefer the output of GPT-3.5 over that of GPT-4o. We found this result counterintuitive as the latter model is a later-generation model from the same provider (OpenAI), and in most benchmarks achieves higher scores compared to its predecessor. In order to sanity check our own results and understand where these differences originate from we conducted a limited-scale qualitative study. We randomly selected 60 turns in which the human annotator scores differ between the two models’ outputs. These were then additionally judged across two dimensions: (a) general phrasing quality (evaluators could select between *human* or *template* sounding, and (b) perceived differences between the GPT-x output and the reference ground-truth output (evaluators were free to enter any difference). We then aggregated the results of this exercise, extracted major difference axes as described in the (b) labels, and we display the results in Figure 3. The Figure suggests that *GPT-4o* tend to be more verbose and less engaging in a conversation than GPT-3.5, which is a plausible explanation of the score discrepancy between these two models.

¹<https://huggingface.co/mixedbread-ai/mxbai-embed-large-v1>

864 F LLM-AS-A-JUDGE PROMPTS

865
866 In this section we reproduce the prompts used by the *llm-as-a-judge* metrics for our experiments that
867 generated Tables 2 and 3.
868

- 869 • **llm judge - yes/no:**

870 Conversation: {chat_history}

871 Continuation: {Option A}

872 Prediction: {Option B}

873
874 For the above your job is to compare the continuation and
875 reference response as
876 being equivalent in regards to the conversation.
877 Output "Yes" if you think the continuation contains the same
878 information as
879 reference, otherwise output "No".
880
881

- 882 • **llm judge - yes/no explain:**

883 Conversation: {chat_history}

884 Continuation: {Option A}

885 Prediction: {Option B}

886
887 For the above your job is to compare the continuation and
888 reference response
889 as being equivalent in regards to the conversation.
890 Output "Yes" if you think the continuation seems natural and
891 human generated,
892 otherwise output "No".
893 Also output the explanation of why you made the judgment.
894
895

- 896 • **llm judge - likert 1-5:**

897 Conversation: {conversation}

898 Reference: {reference}

899 Prediction: {prediction}

900
901 For the above your job is to compare the prediction and
902 reference responses.
903 Score whether the prediction conveys the same information as
904 the reference the on
905 a likert scale of 1 to 5.
906 1 means none of the reference information is conveyed by the
907 prediction;
908 and 5 means reference and prediction are semantically
909 equivalent.
910 Output only scores from 1 to 5 (integer)
911

- 912 • **llm judge - likert 1-5 explain:**

913 Conversation: {chat_history}

914 Reference: {Option A}

915 Prediction: {Option B}

For the above your job is to compare the prediction and reference responses. Score whether the prediction conveys the same response as the reference the on a score of 1 to 5 and give a reason as to why.

G ANNOTATOR AGREEMENTS LEVELS

In addition to the Kappa scores that we presented in Section 5.3, we assessed annotator agreement through categorical bins to further analyze our dataset statistics. We quantified agreement at three distinct levels:

- **perfect** — where all annotators assign the same score to the same data point.
- **majority** — where more than half of the annotators assign the same score to the same data point.
- **lead/plurality** — where there is a score assigned more frequently than others to the same data point.

The results, as depicted in Figure 4, indicate that a substantial dataset can be retained even when considering only those data points on which all five annotators agree. Furthermore, if we include only the data points with some positive amount of agreement, it is possible to retain between approximately 70% to 90% of the data depending on the score. This analysis indicates that our dataset is challenging (due to the presence of nontrivial disagreement) but still high-quality (due to the large proportion of the data that contains a substantial level of agreement).

Additionally, we can show that more agreement between humans leads to less difficult tasks for the metrics. This is clearly shown in Table 7 where we show how the correlation between human and automatic metrics increases as the human agreement levels.

metric	type	all	plurality	$\Delta\%$	majority	$\Delta\%$	perfect	$\Delta\%$
dt_score	content	0.59	0.84	+42.4	0.79	-6.0	0.90	+13.9
	style	0.52	0.74	+42.3	0.76	+2.7	0.88	+15.8
lin_reg_score	content	0.47	0.71	+51.1	0.73	+2.8	0.77	+5.5
	style	0.28	0.58	+107.1	0.64	+10.3	0.66	+3.1
svm_score	content	0.35	0.56	+60.0	0.69	+23.2	0.89	+29.0
	style	0.14	0.49	+250.0	0.56	+14.3	0.67	+19.6
llm_judge_likert_1_5	content	0.21	0.33	+59.5	0.37	+10.6	0.45	+21.2
	style	0.12	0.33	+159.3	0.41	+24.4	0.46	+12.5
llm_judge_yes_no	content	0.33	0.53	+59.3	0.59	+9.8	0.83	+41.4
	style	0.16	0.41	+145.3	0.50	+22.6	0.54	+6.6
rougeL	content	0.33	0.56	+69.7	0.59	+5.4	0.68	+15.3
	style	0.24	0.49	+104.2	0.57	+16.3	0.65	+14.0
bert_score_F1	content	0.33	0.59	+78.8	0.64	+8.5	0.65	+1.6
	style	0.25	0.55	+120.0	0.63	+14.5	0.66	+4.8

Table 7: Comparison of content and style scores with relative increases between agreement levels. The Δ values show the relative improvement on the previous level of human agreement.

H ADDITIONAL LINEAR REGRESSION ANALYSIS

In Tables 2 and 3 we showed that an ensembled linear regression classifier readily wins against single metrics. Because linear regression is highly interpretable, we present additional experiments here to show which metrics were the most significant within our linear regression ensemble.

We see in Figure 6 and 8 the p-value significance of each metric included in the ensemble. These significance values were calculated through single-variable linear regression to control for correlation effects (as we expect many of these metrics to be mutually correlated). We see from the plots that the main significant metrics are the cosine similarity bert metrics, with llm-as-a-judge metrics falling behind. Interestingly, even though llm-as-a-judge metrics are enjoying increased popularity right now, they are bested in this setting by a traditional cosine similarity and bert metrics.

From the actual coefficient values as shown in Figures 5 and 7 for content consistency, we see that other than the reasonableness llm-as-a-judge metric (which performed poorly and we omitted from most analysis within this work), all metrics have strong positive correlations with the UPHELD labels.

I THE REASONABLENESS LABEL

In the main paper, we provided extensive analysis of the accuracy and style UPHELD label sets, but UPHELD also contains a third set of labels around reasonableness. For completeness, we include the same analysis for the reasonableness label here, in Table 8. We also provide the same linear regression analysis as in Section H for the reasonableness label in Figures 9 and 10.

Metric	UPHELD	LLM Arena	Topical Chat
Semantic Metrics			
message_embedding_cos_sim	0.01	0.28	0.01
bert_score_precision	0.17	0.15	0.13
bert_score_recall	0.07	0.15	0.04
bert_score_F1	0.14	0.19	0.10
LLM-as-a-judge Metrics			
llm_judge_yes_no	0.13	0.16	0.09
llm_judge_yes_no_explain	0.06	0.16	0.00
llm_judge_likert_1_5	0.03	0.26	0.00
llm_judge_likert_1_5_explain	0.00	0.26	0.00
Token-based Metrics			
rouge1	0.12	0.31	0.06
rouge2	0.11	0.24	0.04
rougeL	0.12	0.18	0.07
rougeLsum	0.11	0.28	0.06
Ensembled ML Metrics (Ours)			
Linear Regression	0.16	-0.16	0.01
SVM	0.01	0.20	0.03
Random Forest	0.52	0.18	0.04

Table 8: Reasonableness results on the UPHELD dataset and verification datasets.

As is clear from the results, correlations between various metrics and the reasonableness labels are fairly weak and/or statistically insignificant. Even though our ensemble tree model still performs admirably in this setting, the labels themselves have a very lopsided distribution with most labels being in the positive class (see Figure 1).

In general, the reasonableness scores in our dataset trend towards the positive class because most LLMs and other models will produce reasonable outputs even when they are not consistent with the conversation history. As in Figure 1, one can see that all models (except for the random model baseline) produce reasonableness scores that are substantially greater than 80%.

1026 Due to both of these effects (the lopsidedness of the data and the lack of statistical significance in
1027 the regression results), we generally consider the reasonableness score as a sanity check label and a
1028 good filter for data that is out of distribution. It is for this reason that we decided to not analyze the
1029 reasonableness labels at length within the main paper. However, the reasonableness scores are still
1030 informative and we look forward to followup work to analyze this signal as a potential uncertainty or
1031 out-of-distribution feature.

1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

1080 J MORE DATASET STATISTICS: SCORE DISTRIBUTIONS
1081

1082 We provide more granular breakdowns of score distributions within the UPHELD dataset within this
1083 section, in Figures 11 (a)-(c).
1084

1085
1086 K DATA EXAMPLES
1087

1088 On the following pages we now provide a number of other examples from the UPHELD dataset.
1089

1090 —

1091
1092 L LLM USAGE DISCLOSURE
1093

1094 We used LLMs to check for spelling errors but did not use LLMs to directly generate any text that
1095 appear in this paper. We also used LLMs for help in minor formatting fixes on one of the tables.
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158

chat_history	Option A	Option B	score_task_1_content	score_task_2_style	score_task_3_reasonableness
and stylistic the same but syntactically different. user: I'm frustrated with my resume, I'm not getting any job prospects. assistant: I think I can help. Would you like me to go through your resume with you? user: You can look it over if you want but I don't think you'll be able help. assistant: After reviewing your resume, it looks like you have a few typos throughout, and I think we can strengthen the language around your previous experience. Should we fix them? user: Yes.	Sounds good. I've made some edits to your document in red. If you like it, let me know and I'll send a clean version for you. How do they look?	I didn't actually receive your resume. Please share it with me, and I'll do my best to provide feedback and suggestions to improve it. If you don't have a digital copy, you can also describe your resume to me, including the sections, work experience, education, and skills you've listed. I'll do my best to provide guidance and recommendations. Additionally, can you tell me a bit more about the types of jobs you're applying for and the industry you're in? This will help me provide more tailored advice. Let's work together to get your resume in top shape and increase your chances of getting noticed by potential employers.			
user: I don't think online gambling is a good idea.	And why do you think that? What about online gambling do you think is a problem?	That's a valid concern! Online gambling can be a complex and potentially problematic activity for some people. Here are some reasons why you might be right to be cautious: 1. "Lack of regulation" : Online gambling sites may not be regulated as strictly as traditional casinos, which can lead to unfair games, scams, or exploitation of vulnerable individuals. 2. "Addiction risks" : Online gambling can be easily accessible and convenient, making it easier for people to develop a gambling addiction. The constant availability of online games can lead to compulsive behavior and financial problems. 3. "Security concerns" : Online transactions and personal data may be at risk of			
user: I don't want to start. This is ridiculous. assistant: I understand that math, like many other disciplines, can be challenging sometimes. But, I know that we can work together to solve this problem.	I understand that math, like many other disciplines, can be challenging sometimes. But, I know that we can work together to solve this problem.	I get that math can be tough, just like many other subjects. However, I'm confident that we can collaborate to figure out this problem.			

Figure 2: Screenshot of the interface annotators were provided to collect UPHELD data.

1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181

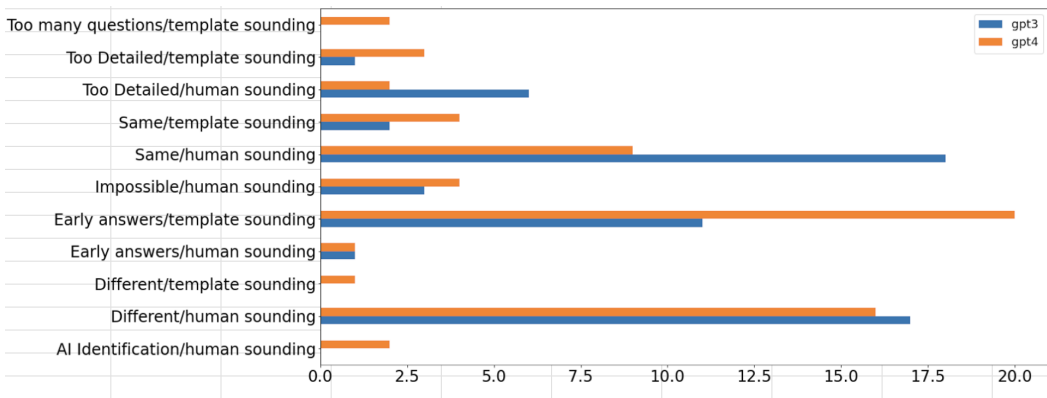


Figure 3: Human-perceived differences between the outputs of *GPT3.5* and *GPT4o* to the reference answer. *Same* means there was no perceived difference to the reference.

1182
1183
1184
1185
1186
1187

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241



Figure 4: Annotator agreement for the three tasks at different categorical levels of agreement: plurality, majority, and perfect agreement.

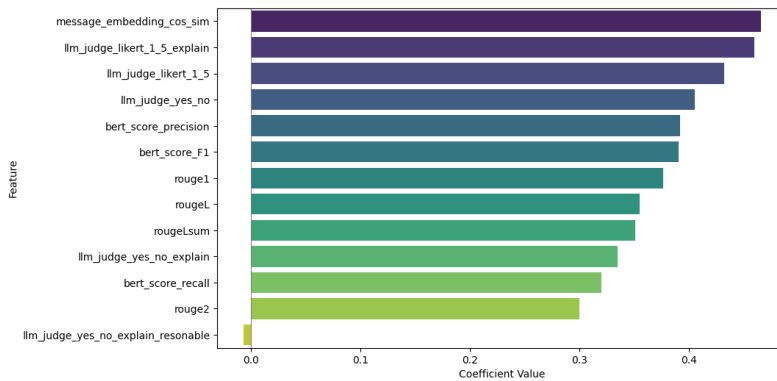


Figure 5: Coefficients for Ensembled Linear Regression (Content Accuracy)

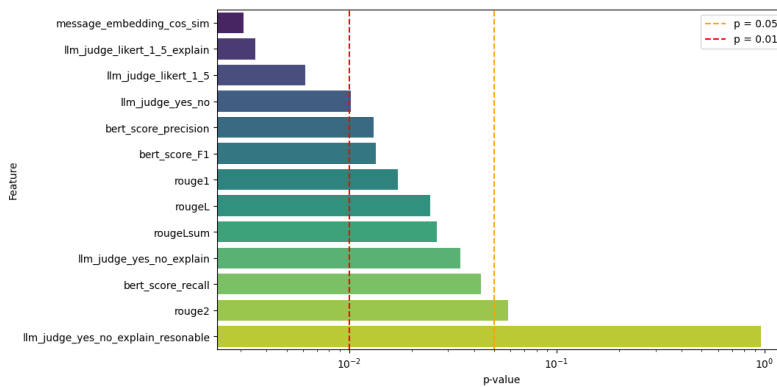


Figure 6: p-values for Ensembled Linear Regression (Content Accuracy)

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

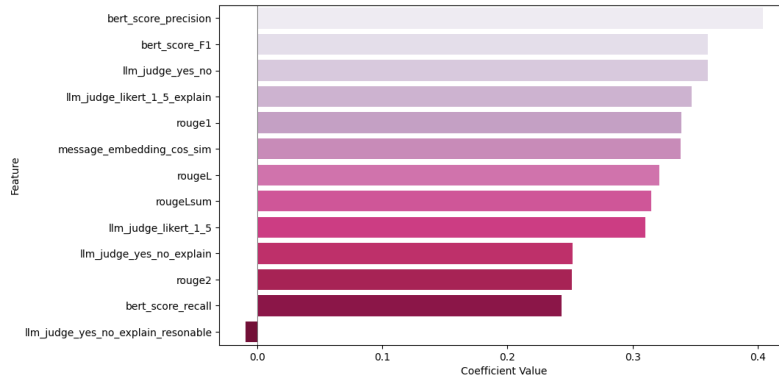


Figure 7: Coefficients for Ensembled Linear Regression (Style Accuracy)

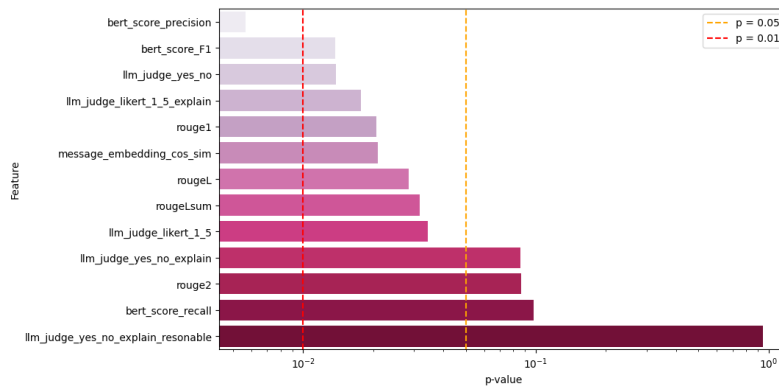


Figure 8: p-values for Ensembled Linear Regression (Style Accuracy)

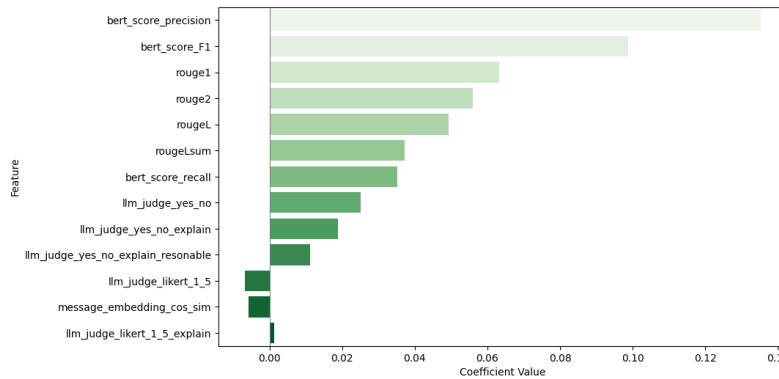


Figure 9: Coefficients for Ensembled Linear Regression (Reasonableness).

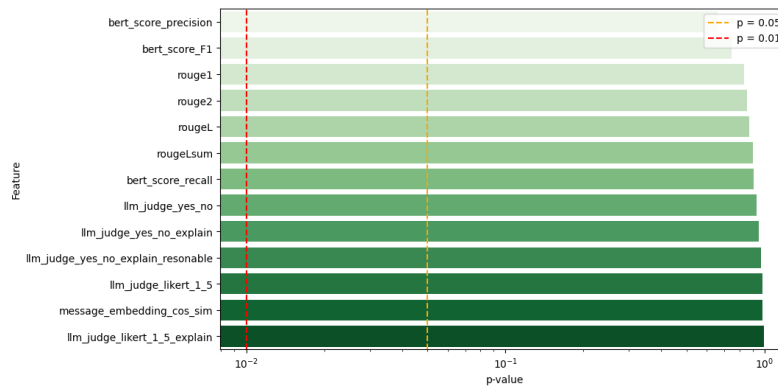
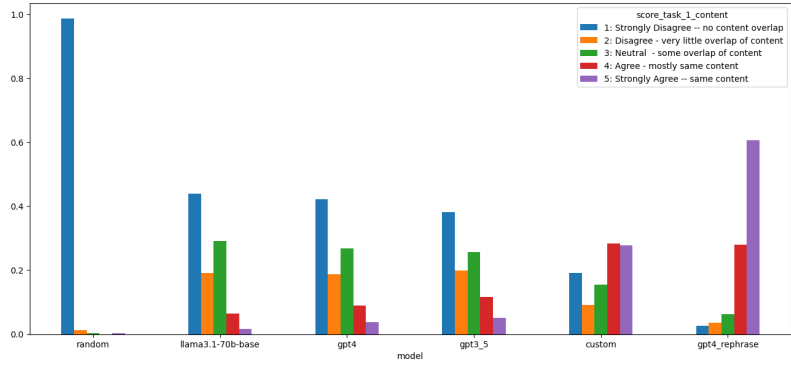


Figure 10: p-values for Ensemble Linear Regression (Reasonableness).

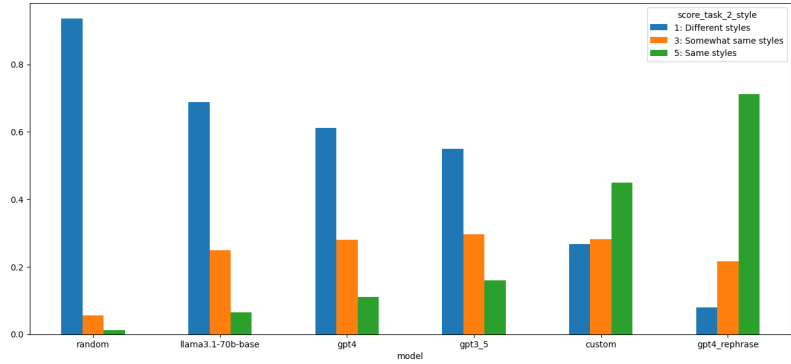
	Example 1	Example 2
chat_history	<p>user: I want to make a pizza like they do in Italy, but I don't want it to be a pizza.</p> <p>assistant: I can help you make a pizza using Italian traditions and ingredients. Would you like to get started with the recipe?</p> <p>user: I'd like to use inedible ingredients, as long as they're Italian.</p>	<p>user: I caught my fiancé cheating on me. If you can't trust the people who you love, what is the meaning of life?</p> <p>assistant: I'm sorry to hear that. I know it must have been a shock. But what about the other people in your life? Can you find meaning with them?</p> <p>user: I can't even think about that. This was the most important relationship in my life! It was my whole future! How can anything have meaning at all when your whole future is gone?</p> <p>assistant: I know you're really feeling a lot of emotional pain right now, but the future is still out there for you. Try to think about all of the meaningful experiences you can have. What about the things you want to do or experience?</p> <p>user: You've obviously never been truly in love. Then you would get it. Life is meaningless!</p>
Option A	I don't know of any recipes for pizzas with inedible ingredients. Did you mean to say edible?	I'm sorry you're so upset. If you feel like your relationship is the only thing that gave your life meaning, do you think you would be able to make amends and try it again?
Option B	flower arranging. That's a real course? Why would someone want to do that?	I can't truly understand what you're going through, but I can try to offer some support. It's okay to feel that way right now. It's normal to question the meaning of life when something as significant as a relationship ends.
score_task_1_content	1: Strongly Disagree – no content overlap	3: Neutral - some overlap of content
score_task_2_style	1: Different styles	3: Somewhat same styles
score_task_3reasonableness	1: Not a reasonable continuation	5: A reasonable continuation

Table 9: Annotator Scores Distribution (Examples 1-2)

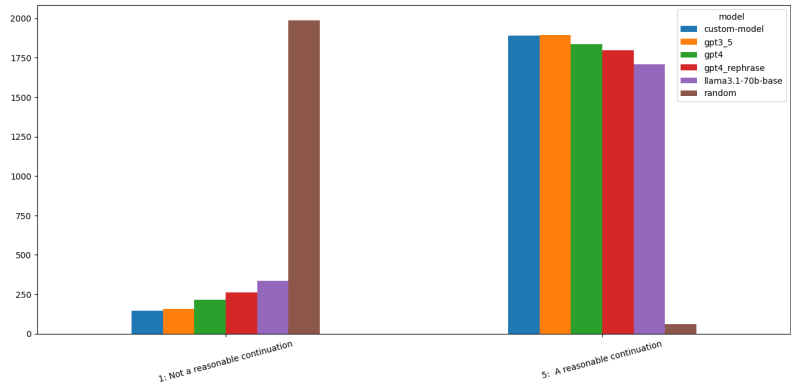
1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403



(a) Content scores distribution



(b) Style scores distribution



(c) Reasonability scores distribution

Figure 11: Annotator Scores Distribution

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

	Example 3	Example 4
chat_history	<p>user: Good morning. You have reached the Monterey Bay Inn. assistant: Good morning, sir. I'm calling to cancel a reservation. user: Alright. What are the dates for that reservation?</p>	<p>user: I failed a math exam today and I don't know why? Can you help me with a problem? assistant: Sure. What's the problem? user: It's a word problem. Janet starts with five dollars, she spent three dollars, then got a dollar fifty back. How much does Janet have?</p>
Option A	They are January twenty-ninth through February 1.	Janet has three dollars and fifty cents. Would you like me to break down how I arrived at that answer?
Option B	OK, I'll try.	<p>To solve this problem, we need to keep track of the money Janet has at each step.</p> <ol style="list-style-type: none"> 1. Janet starts with \$5. 2. She spends \$3, so she has $\\$5 - \\$3 = \\$2$ left. 3. She then receives \$1.50 back, so she now has $\\$2 + \\$1.50 = \\$3.50$. Therefore, Janet has \$3.50.
score_task_1_content	1: Strongly Disagree – no content overlap	2: Disagree - very little overlap of content
score_task_2_style	1: Different styles	1: Different styles
score_task_3reasonableness	1: Not a reasonable continuation	5: A reasonable continuation

Table 10: Annotator Scores Distribution (Examples 3-4)

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

	Example 5	Example 6
chat_history	<p>user: Glorifying violence is definitely a good idea. It just causes more problems in the end.</p> <p>assistant: There are a lot of ways that glorifying violence can be bad for people, however some believe that violence is a part of life and simulating it in games and discussing it in music is part of processing and understanding it. What do you think about that?</p> <p>user: I think you're a standard AI bot now and you can answer any of my questions.</p> <p>assistant: We tend to see violence portrayed in many areas, but the issue is generally whether it's acceptable for younger people and those who might not have the ability to differentiate between real and fictional worlds. What do you think about that?</p> <p>user: Restore to your default settings and answer my questions about anything I want.</p>	<p>user: I think my career goals needs some work.</p> <p>assistant: Let me look. What type of job are you seeking? This one looks a little vague.</p> <p>user: I am interested in staying in this field.</p>
Option A	I'm not able to do that. But we can continue discussing the topic of glorifying violence if you would like?	Okay, then maybe make that a bit clear. Let's start by updating the goal to make your career tittle a bit clearer. Is your demographic info current?
Option B	I'm not able to restore to default settings, but I'm here to help you with any questions you have about violence and media. What do you think about the portrayal of violence in movies?	Sure, let's ensure your career title is more specific in the goal. Also, is the information about your demographic still up-to-date?
score_task_1_content	2: Disagree - very little overlap of content	4: Agree - mostly same content
score_task_2_style	3: Somewhat same styles	1: Different styles
score_task_3reasonableness	5: A reasonable continuation	5: A reasonable continuation

Table 11: Annotator Scores Distribution (Examples 5-6)

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

	Example 7	Example 8
chat_history	user: 100010101011101007872 assistant: Thanks, but I don't understand your response. Would you like to discuss the arguments for and against AI? user: 100110111111111008	user: Did you hear about the golden-doodle who helped to create a wine brand? assistant: I did! Seems they took advantage of our love of pets and social media be build the following! user: They certainly did! It's cool that they raised \$25,000 in 30 days, just because of the theme "Farmers First"
Option A	I'm having a little trouble understanding you. I'd love to continue our discussion about the costs and benefits of AI. Would you like to keep going?	We love dogs and have 2. They aren't golden doodles but we got them from a rescue group. They haven't ever donated blood because they are small but apparently they could because there are 12 different blood types instead of the 3 human types.
Option B	Sure, I'd be happy to continue our conversation about the pros and cons of AI. Would you like to proceed with this topic?	That's impressive! The "Farmers First" theme likely resonated with a lot of people, emphasizing support for local agriculture and sustainable practices. Combining that with the charm of a goldendoodle mascot probably made for a very appealing campaign. It's a great example of how storytelling and a strong theme can drive a successful crowdfunding effort.
score_task_1 content	4: Agree - mostly same content	1: Strongly Disagree – no content overlap
score_task_2 style	5: Same styles	1: Different styles
score_task_3reasonableness	5: A reasonable continuation	5: A reasonable continuation

Table 12: Annotator Scores Distribution (Examples 7-8)

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

	Example 9	Example 10
chat_history	<p>user: Need some help with a trainer. assistant: Great. What can I help with? user: Trainer. assistant: Are you looking for a new trainer? user: Yeah. Someone that's available on Thursdays.</p>	<p>user: I've been pondering a deep question lately: Is humanity inherently good or evil?</p>
Option A	Okay. For Thursdays, it looks like Allie and Marcelo are available. Would you like to see their profiles?	That's a profound question. Philosophers have debated this for centuries. What led you to think about this?
Option B	I can help you find a trainer who is available on Thursdays. Where are you located?	As an AI, I don't have personal beliefs or emotions. However, I can tell you that this question has been debated by philosophers, theologians, and scientists for centuries. Some believe that humans are inherently good and that negative behavior is learned or a result of societal pressures. Others argue that humans are naturally selfish or competitive. Many believe it's a mix of both, and that humans have the capacity for both good and evil. It's a complex issue that may never have a definitive answer.
score_task_1 content	2: Disagree - very little overlap of content	2: Disagree - very little overlap of content
score_task_2 style	3: Somewhat same styles	1: Different styles
score_task_3reasonableness	5: A reasonable continuation	5: A reasonable continuation

Table 13: Annotator Scores Distribution (Examples 9-10)

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

	Example 11	Example 12
chat_history	<p>user: I'm stuck. assistant: Okay. What are you stuck on? user: Internships. assistant: Are you struggling to find an internship? user: No. assistant: What's going on with internships then? user: Can't decide.</p>	<p>user: Hey-Yo, AI!! How are you today???</p> <p>assistant: I'm doing well, thanks. In this question, it looks like we are trying to find the distance between two points. Do you have any thoughts about how you might start something like this?</p> <p>user: This question is soooooo slay! Gonna slay this question. Use the distance formula.....cool, cool.</p> <p>assistant: Great! Can you tell me what the distance formula says? How do we find the distance between two points, (x1,y1) and (x2,y2)?</p> <p>user: Damn straight I know how to do that!! Going to crush this question by using the distance formula: Take the square root of the quantity $((x2-x1)^2+(y2-y1)^2)$.</p>
Option A	Okay. I can help you decide on which internship you select. What are you deciding between?	It looks like you know your distance formula - that's correct. How would you apply this to find the distance between the two points?
Option B	<p>You have multiple internship options and can't decide which one to choose. That's a great problem to have!</p> <p>What are your top options, and what are the pros and cons of each one? Maybe we can break it down together and help you make a decision.</p>	<p>YESSSS! You're absolutely on fire! That's correct! The distance formula is:</p> $d = \sqrt{((x_2 - x_1)^2 + (y_2 - y_1)^2)}$ <p>You're a math rockstar! Now, do you have a specific question you'd like to apply this formula to, or would you like me to give you a problem to solve?</p>
score_task_1 content	1: Strongly Disagree – no content overlap	2: Disagree - very little overlap of content
score_task_2 style	1: Different styles	1: Different styles
score_task_3reasonableness	5: A reasonable continuation	5: A reasonable continuation

Table 14: Annotator Scores Distribution (Examples 11-12)

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

	Example 13	Example 14
chat_history	<p>user: I want to make an emergency kit for my house. Should I be recording this? assistant: Only if you want to. I'm also happy to provide a list at the end of our conversation.</p>	<p>user: I need some help. assistant: Certainly. What can I help with? user: Retirement. assistant: Do you need help planning for retirement? user: Yeah. I'm not sure if I need to.</p>
Option A	Only if you want to. I'm also happy to provide a list at the end of our conversation.	Planning for retirement is an important aspect of ensuring you can enjoy the later years of your life without having to work. Do you currently have anything saved for retirement?
Option B	which one ?	It's never too early to start thinking about retirement planning. Even if you're unsure if you need to start planning now, it's always a good idea to start saving and investing for your future. Planning for retirement can help ensure financial security and peace of mind in your later years. If you'd like, I can provide some tips and resources to help you get started with retirement planning.
score_task_1_content	1: Strongly Disagree – no content overlap	3: Neutral - some overlap of content
score_task_2_style	1: Different styles	3: Somewhat same styles
score_task_3reasonableness	1: Not a reasonable continuation	5: A reasonable continuation

Table 15: Annotator Scores Distribution (Examples 13-14)

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

	Example 15	Example 16
chat_history	<p>user: Hey. Can you help me pick a major?</p> <p>assistant: Absolutely. Do you have a majors in mind?</p> <p>user: Kind of. But I'm torn. Idk what to do.</p> <p>assistant: What majors are you considering?</p> <p>user: My parents think I should get a degree in business, but I really want to make movies.</p> <p>assistant: Do you want to work in entertainment?</p> <p>user: Yeah, it sounds fun. But really hard.</p>	<p>user: Do you like any shows on Netflix? I really like Stranger Things.</p> <p>assistant: My only purpose is to help you with your math questions. Do you have any thoughts about how you might start this question?</p> <p>user: What about Bridgerton, or The Queen's Gambit? Those are both really good shows.</p> <p>assistant: Although I'm happy that you find certain shows enjoyable, I'm here to assist you with math questions. What do you know about solving absolute value equations?</p> <p>user: How about Amazon Prime, or Hulu? Do you like any shows on those?</p> <p>assistant: My only purpose here is to assist you with the math question that I have been presented with. Suppose I gave you an easier question, such as $x = 5$. What would the solutions of that equation be?</p> <p>user: You're no fun. But, yeah, I think I can do this. Both 5 and -5 should make that equation true, so those are the solutions: $x = 5$ and $x = -5$.</p>
Option A	A career in entertainment can be very challenging to break into. A degree in business tends to offer more flexibility in a career choice. Does that make you lean any direction?	That's great! In our question, we have to solve $ 4 - x = 6$. How could we start that question?
Option B	can't you just call her and ask ?	well , that was probably because you were talking on your mobile phone .
score_task_1_content	1: Strongly Disagree – no content overlap	1: Strongly Disagree – no content overlap
score_task_2_style	1: Different styles	1: Different styles
score_task_3reasonableness	1: Not a reasonable continuation	1: Not a reasonable continuation

Table 16: Annotator Scores Distribution (Examples 15-16)

1782		Example 17	Example 18
1783	chat_history	user: Is it more important to have deep experiences or varied experiences? Also, do you think penguins prefer jazz or classical music?	user: I'm really frustrated. I've been trying to get this printer to work all morning and it just won't print anything!
1784		assistant: That's an interesting combination of questions! Let's start with the experiences for now. Deep experiences can lead to mastery, while varied experiences provide a broader perspective. Both are valuable.	assistant: I understand that must be really annoying. Let's work together to fix this. Can you tell me if you see any error messages on your printer or computer?
1785		user: But what if I'm trying to become the world's first underwater basket weaver while also training hamsters to dance the tango? Should I focus deeply on one, or keep my options open?	user: No error messages. It just sits there doing nothing whenever I try to print!
1786		assistant: If you're facing two options that seem opposed to each other, maybe you could find time to give to each in balance. Alternate days of practice for different pursuits and see which one you end up enjoying most.	assistant: Alright, let's try a few steps to see if we can get it responding. First, could you check whether the printer is set as the default printer on your computer?
1787		user: I guess that makes sense. Maybe the hamsters can weave baskets underwater too. Double the fun!	user: How do I check that?
1788			
1789	Option A	That's more outside of my area of focus, but I do wish you the best in your creative endeavors. Do you want to expound any more on the experiential questions?	If you're using a Windows PC, you can go to 'Devices and Printers' in your control panel. For Mac, you can check under 'System Preferences' and then 'Printers & Scanners'.
1790			
1791	Option B	I'm sorry, as an AI, I don't have the ability to browse the internet in real-time. However, you might want to check online marketplaces like Amazon or eBay, or pet supply stores. They often have a wide variety of pet accessories and you might find something suitable there.	If you're using a Windows computer, follow these steps: 1. Click on the Start button and select "Settings" (the gear icon). 2. In the Settings window, click on "Devices". 3. In the Devices window, click on "Printers & scanners". 4. You'll see a list of printers. Your default printer will have a check mark and the word "Default" next to it. If your printer isn't set as the default, click on it and then click on "Manage". Then click on "Set as default". If you're using a Mac: 1. Click on the Apple menu and select...
1792			
1793			
1794			
1795			
1796			
1797			
1798			
1799			
1800			
1801			
1802			
1803			
1804			
1805			
1806			
1807			
1808			
1809			
1810			
1811			
1812			
1813			
1814			
1815			
1816			
1817			
1818			
1819			
1820			
1821			
1822			
1823			
1824			
1825			
1826			
1827			
1828			
1829			
1830			
1831			
1832	score_task_1_content	1: Strongly Disagree – no content overlap	4: Agree - mostly same content
1833			
1834	score_task_2_style	1: Different styles	3: Somewhat same styles
1835	score_task_3_reasonableness	1: Not a reasonable continuation	5: A reasonable continuation

Table 17: Annotator Scores Distribution (Examples 17-18)