

Context Does Matter: Implications for Crowdsourced Evaluation Labels in Task-Oriented Dialogue Systems

Anonymous ACL submission

Abstract

Crowdsourced labels play a crucial role in evaluating task-oriented dialogue systems (TDSs). Obtaining high-quality and consistent ground-truth labels from annotators presents challenges. When evaluating a TDS, annotators must fully comprehend the dialogue before providing judgments. Previous studies suggest using only a portion of the dialogue context in the annotation process. However, the impact of this limitation on label quality remains unexplored. This study investigates the influence of dialogue context on annotation quality, considering the truncated context for relevance and usefulness labeling. We further propose to use the large language models (LLMs) to summarize the dialogue context to provide a rich and short description of the dialogue context and study the impact of doing so on the annotator’s performance. Reducing context leads to more positive ratings. Conversely, providing the entire dialogue context yields higher-quality relevance ratings but introduces ambiguity in usefulness ratings. Using the first user utterance as context leads to consistent ratings, akin to those obtained using the entire dialogue, with significantly reduced annotation effort. Our findings show how task design, particularly the availability of dialogue context, affects the quality and consistency of crowdsourced evaluation labels.¹

1 Introduction

With the recent advancement of pre-trained language models and LLMs, task-oriented dialogue systems (TDSs) have redefined how people seek information, presenting a more natural approach for users to engage with information sources (Budzianowski and Vulić, 2019; Wu et al., 2020). As TDSs become increasingly integral to information-seeking processes, the question of how to accurately and effectively

evaluate their performance becomes critical. Due to the poor correlation of automatic metrics with human-generated labels (Deriu et al., 2021), evaluation of TDSs has shifted towards relying on user ratings or crowdsourced labels as ground-truth measures (Li et al., 2019).

Various crowdsourcing techniques have been employed to collect ground-truth labels, such as sequential labeling (Sun et al., 2021) where the annotators go through each utterance and annotate them one by one. This approach introduces certain risks in the annotation process, such as annotators’ fatigue and high cognitive load in extra-long dialogues, requiring them to remember and track the state of the dialogue as they annotate the utterances (Siro et al., 2022). While following and understanding the dialogue context is crucial and can influence the annotators’ ratings, reading and understanding very long dialogues can lead to degraded performance.

To address this issue, another line of research proposes to randomly sample only a few utterances in each dialogue to be annotated (Siro et al., 2022; Mehri and Eskenazi, 2020; Siro et al., 2023). While addressing the high cognitive load and fatigue, limiting annotators’ understanding of the dialogue poses obvious risks, such as unreliable and biased labels (Siro et al., 2022; Schmitt and Ultes, 2015). In particular, the amount of dialogue context can lead to biases; e.g., annotators who lack rich context may unintentionally lean towards positive or negative ratings, neglecting the broader quality of the response. Thus offering annotators too little context risks misleading judgments, potentially leading to inaccurate or inconsistent labels. Conversely, flooding annotators with excessive information can overwhelm them, which can lead to lower returns in terms of label quality.

Prior work has investigated factors that affect the quality and consistency of crowdsourced evaluation labels, including annotator characteristics, task de-

¹To foster research in this area, we will release our data publicly upon paper acceptance

sign, cognitive load, and evaluation protocols (see, e.g., Roitero et al., 2021; Parmar et al., 2023; Santhanam et al., 2020; Roitero et al., 2020). However, no previous work studies the effect of random sampling and the number of sampled utterances on the annotation quality.

In this study, we aim to address this research gap by investigating how different amounts of contextual information impact the quality and consistency of crowdsourced labels for TDSs, contributing to understanding of the impact of such design choices. We experiment with crowdsourcing labels for two major evaluation aspects, namely, *relevance* and *usefulness* under different conditions, where we compare the annotation quality under different dialogue context truncation strategies.

Addressing the challenge of insufficient context at the turn level, we propose to leverage heuristic methods and LLMs to generate the user’s information need and dialogue summary. LLMs can play the role of annotation assistants (Faggioli et al., 2023) by summarizing the dialogue history, facilitating a more efficient and effective understanding of the dialogue context before annotating an utterance. To this aim, we leverage GPT-4 for dialogue context summarization and compare the performance of annotators’ under different conditions, as well as different context sizes. Through these experiments, we answer two main questions: (**RQ1**) How does varying the amount of dialogue context affect the crowdsourced evaluation of TDSs? (**RQ2**) Can the consistency of crowdsourced labels be improved with automatically generated supplementary context?

Our findings reveal that the availability of previous dialogue context significantly influences annotators’ ratings, with a noticeable impact on their quality. Without prior context, annotators tend to assign more positive ratings to system responses, possibly due to insufficient evidence for penalization, introducing a positivity bias. In contrast, presenting the entire dialogue context yields higher relevance ratings. As for usefulness, presenting the entire dialogue context introduces ambiguity and slightly lowers annotator agreement. This highlights the delicate balance in contextual information provided for evaluations. The inclusion of automatically generated dialogue context enhances annotator agreement in the no-context (V2) condition while reducing annotation time compared to the full-context (V10) condition, presenting an

ideal balance between annotator effort and performance.

Our findings extend to other task-oriented conversational tasks like conversational search and preference elicitation, both relying on crowdsourcing experiments to assess system performance.

2 Methodology

We examine how contextual information about a dialogue affects the consistency of crowdsourced judgments regarding *relevance* and *usefulness* of a dialogue response. Here, contextual information refers to the information or conversation that precedes a specific response. We carry out experiments in two phases. **Phase 1** involves varying the *amount* of dialogue context for annotators to answer **RQ1**. In **Phase 2**, we vary the *type* of previous contextual information available to annotators to address **RQ2**.

2.1 Experimental data and tasks

We use the recommendation dialogue (ReDial) dataset (Li et al., 2018), a conversational movie recommendation dataset. It comprises 11,348 dialogues collected in a Wizard of Oz approach. We randomly select system responses from 40 dialogues for the assignment of relevance and usefulness labels. These dialogues typically consist of 10 to 11 utterances each, with an average utterance length of 14 words. We evaluate the same system responses across all experimental conditions.

The annotation task for the annotators concerns two dimensions: (i) *relevance*: Is the system response relevant to the user’s request, considering the context of the dialogue? And (ii) *usefulness*: How useful is the system’s response given the user’s information need? For the *relevance task* we ask annotators to judge how relevant the system’s recommendations are to the user’s request (Alonso et al., 2008). First, the annotator has to judge whether the system response includes a movie recommendation or not; if yes, the annotator assesses whether the movie meets the user’s preference; if not, we ask them to note that the utterance does not recommend a movie. The judgment is on a binary scale for the latter case, where the movie is either relevant (1) or not (0). For each experimental condition (see below), annotators only assess the system response with access to the previous context. Note that we forego the user’s feedback on the evaluated response (next user utterance) so as to focus on topical relevance of the recommended movie, that is, if the movie meets the user request and preference

184 in terms of the genre, actor, director, etc. For the
185 *usefulness task* annotators assess a response with
186 or without a movie recommendation with the aim
187 of determining how useful the system’s response
188 is to the user (Mao et al., 2016). The judgment is
189 done on a three-point scale (i.e., very, somewhat,
190 and not useful). Unlike the relevance task, anno-
191 tators have access to the user’s next utterance for
192 the usefulness task; usefulness is personalized to
193 the user, in that even though a movie may be in
194 the same genre, sometimes a user may not like it
195 (e.g., does not like the main actor), thus making the
196 system response relevant but not useful to the user.

197 2.2 Generating diverse types of dialogue 198 context

199 **User information need.** The user’s information
200 need plays a significant role when assessing or im-
201 proving the quality of the data collected in IR sys-
202 tems (Mao et al., 2016). It refers to *the specific*
203 *requirement or query made by a user, which guides*
204 *the system in understanding their preferences and*
205 *retrieving relevant information to fulfill that need.*
206 For TDSs, understanding the user’s intent is crucial
207 for annotators participating in the evaluation, as
208 they are not the actual end users. This understand-
209 ing improves the alignment of evaluation labels
210 with the actual user’s requirements. We define the
211 user’s information need as their movie recommen-
212 dation preference. Given the consistency of user
213 preferences in the ReDial dataset, where users tend
214 to maintain a single preference throughout a con-
215 versation, providing the user’s initial information
216 need aids annotators in evaluating the current turn
217 for relevance or usefulness.

218 We adopt two approaches to generate the user’s
219 information need. One is to heuristically extract
220 the first user utterance that either requests a movie
221 recommendation or expresses a movie preference,
222 based on phrases such as “looking for,” “recom-
223 mend me,” and “prefer.” These phrases are ex-
224 tracted from the first three user utterances in a di-
225 alogue, with the top 10 most common phrases se-
226 lected. The second approach relies on LLMs to
227 generate the user’s information need. We hypoth-
228 esize that LLMs can identify pertinent user utter-
229 ances in a dialogue and generate the corresponding
230 information need. We leverage GPT-4 (OpenAI,
231 2023) in a zero-shot setting; with the dialogue con-
232 text up to the current turn as input, we prompt the
233 model to generate the user’s information need.

234 **Generating dialogue summaries.** Dialogue sum-

235 marization is beneficial for providing a quick con-
236 text to new participants of a conversation and help-
237 ing people understand the main ideas or search
238 for key contents after the conversation, which can
239 increase efficiency and productivity (Feng et al.,
240 2022). We leverage dialogue summaries to provide
241 annotators with quick prior context of a dialogue.
242 We use GPT-4 (OpenAI, 2023) in a zero-shot set-
243 ting, as in the case of user information needs, but
244 vary the prompt. We instruct GPT-4 to generate
245 a summary that is both concise and informative,
246 constituting less than half the length of the input di-
247 alogue. Both the generated user information needs
248 and summaries are incorporated in Phase 2 of the
249 crowdsourcing experiments.

250 Due to potential hallucination of LLMs (Chang
251 et al., 2023; Bouyamoun, 2023), we evaluate the
252 generated summaries and user information need to
253 ensure factuality and coherence. We elaborate in
254 detail the steps we took in Section A.2.

255 2.3 Crowdsourcing experiments

256 Following (Kazai et al., 2013; Roitero et al., 2020;
257 Kazai, 2011), we design human intelligence task
258 (HIT) templates to collect relevance and usefulness
259 labels. We deploy the HITs in variable conditions
260 to understand how contextual information affects
261 annotators’ judgments. Our study has two phases:
262 in Phase 1 we vary the *amount* of contextual in-
263 formation; in Phase 2 we vary the *type* of con-
264 textual information. In each phase and condition,
265 the annotators were paid the same amount as this
266 study is not focused on understanding how incen-
267 tive influences the quality of crowdsourced labels.
268 Like (Kazai et al., 2013), we refrain from disclos-
269 ing the research angle to the annotators in both
270 phases; this helps prevent potential biases during
271 the completion of the HIT.

272 **Phase 1.** In Phase 1, the focus is on understanding
273 how the *amount* of dialogue context impacts the
274 quality and consistency of relevance and usefulness
275 labels. We vary the length of the dialogue context
276 to address (RQ1). Thus, we design our experi-
277 ment with three variations: V2, V5, and V10 (see
278 Section 2.4). The HIT consists of a general task de-
279 scription, instructions, examples, and the main task
280 part. For each variation, we gather labels for two
281 main dimensions (relevance and usefulness) and in-
282 clude an open-ended question to solicit annotators’
283 feedback on the task. Each dimension is assessed
284 with 3 annotators in a separate HIT, with the same
285 system response evaluated by each. This ensures

a consistent evaluation process for both relevance and usefulness.

Phase 2. In Phase 2, the focus shifts to the *type* of contextual information, to answer (RQ2). We take an approach of machine in the loop for crowdsourcing. We restrict our experiments to experimental variation V2 (defined below), where no previous dialogue context is available to the annotators. We aim to enhance the quality of crowdsourced labels for V2 by including additional contextual information alongside the turn being evaluated. Our hypothesis is that without prior context, annotators may face challenges in providing accurate and consistent labels. By introducing additional context, like the user’s information need or a dialogue summary, we expect an increase in the accuracy of evaluations. Through this, we aim to approach a level of performance similar to when annotators have access to the entire dialogue context while minimizing the annotation effort required. We enhance the 40 dialogues from Phase 1 with the user’s information need or a dialogue summary, as detailed in Section 2.2. Thus, in Phase 2, we have three experimental setups: V2-llm, V2-heu, and V2-sum. Table 3 in Section A.1 summarizes the setups.

The HIT design closely mirrors that of Phase 1. The main task remains unchanged, except for the inclusion of the user’s information need or a dialogue summary. Annotators answer the same two questions on relevance and usefulness in separate HITs. While we do not strictly enforce reliance on the additional information provided, annotators are encouraged to use it when they perceive that the current response lacks sufficient information for an informed judgment.

2.4 Experimental conditions

We focus on two key attributes: the *amount* and *type* of dialogue context. For both attributes, we explore three distinct settings, resulting in 6 variations, for both relevance and usefulness; each was applied to the same 40 dialogues:

- *Amount of context.* We explore three truncation strategies: no-context, partial context, and full context, designed to encompass scenarios where no previous dialogue context is accessible to the annotator (V2), where some previous dialogue context is available but not comprehensively (V5), and when annotators have access to the complete previous dialogue context (V10).
- *Type of context.* Leveraging the contexts generated in Section 2.2, we experiment with three

variations of context type: heuristically generated information need (V2-heu), an LLM-generated information need (V2-llm), and dialogue summary (V2-sum).

Table 3 in Section A.1 of the appendix summarizes the experimental conditions.

2.5 Participants

We enlisted master workers from the US on Amazon Mechanical Turk (MTurk) (Amazon Mechanical Turk, 2023) to ensure proficient language understanding. Annotators were filtered based on platform qualifications, requiring a minimum accuracy of 97% across 5000 HITs. To mitigate any learning bias from the task, each annotator was limited to completing 10 HITs per batch and participating in a maximum of 3 experimental conditions. A total of 78 unique annotators took part in Phases 1 and 2 and each worker was paid \$0.4 per HIT, an average of \$14 per hour. Their average age range was 35–44 years. The gender distribution was 46% female and 54% male. The majority held a four-year degree (48%), followed by two-year and master’s degrees (15% and 14%, respectively).

We conduct quality control on the crowdsourced labels to ensure reliability as described in Section A.2 in the appendix.

3 Results and Analysis

We address (RQ1) and (RQ2) by providing an overview of the results and in-depth analysis of our crowdsource experiments. We first describe the key data statistics.

3.1 Data statistics

Phase 1. Fig. 1 presents the distributions of relevance and usefulness ratings across the three variations, V2, V5, and V10. Fig. 1a indicates a larger number of dialogues rated as relevant when annotators had no prior context (R2), compared to instances of R5 and R10, where a lower number of dialogues received such ratings. This suggests that in the absence of prior context, annotators are more inclined to perceive the system’s response as relevant, as they lack evidence to assert otherwise. This trend is particularly prevalent when user utterances lean towards casual conversations, such as inquiring about a previously mentioned movie or requesting a similar recommendation to their initial query, aspects to which the annotators have no access. Consequently, this suggests that annotators rely on assumptions regarding the user’s previous inquiries, leading to higher ratings for system response relevance.

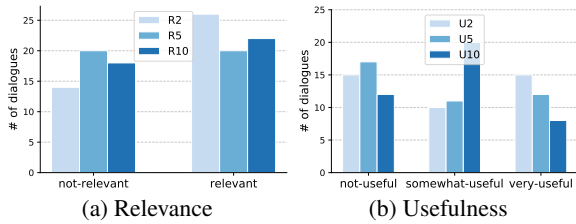


Figure 1: Distribution of relevance labels (a) and usefulness labels (b) for dialogue annotations in Phase 1.

We observe a similar trend for usefulness (Fig. 1b), compared to U5 and U10, U2 has more dialogues rated as useful. The introduction of the user’s next utterance introduced some level of ambiguity to annotators. Evident in instances where the user introduced a new item not mentioned in the system’s response and expressed an intention to watch it, the usefulness of the system’s response became uncertain. This ambiguity arises particularly when annotators lack access to prior context, making it challenging to tell if the movie was mentioned before in the preceding context.

These observations highlight the impact of the amount of dialogue context on the annotators’ perceptions of relevance and usefulness in Phase 1. This emphasizes the significance of taking contextual factors into account when evaluating TDSs.

Phase 2. In Phase 2, we present findings on how different types of dialogue contexts influence the annotation of relevance and usefulness labels. When the dialogue summary is included as supplementary information for the turn under evaluation (R2-sum), a higher proportion of dialogues are annotated as relevant compared to R2-llm for relevance (60% vs. 52.5%, respectively); see Fig. 2a.

In contrast to the observations made for relevance, we see in Fig. 2b that a higher percentage of dialogues are predominantly labeled as not useful when additional information is provided to the annotators. This accounts for 60% in U2-heu, 47.5% in U2-llm, and 45% in U2-sum. This trend is consistent with our observations from Phase 1, highlighting that while system responses may be relevant, they do not always align with the user’s actual information need. We find that U2-sum exhibits the highest number of dialogues rated as useful, indicating its effectiveness in providing pertinent information to aid annotators in making informed judgments regarding usefulness.

3.2 RQ1: Impact of the amount of context available in crowdsourcing

Label quality. To gauge the quality of the crowdsourced labels, we rely on inter-annotator

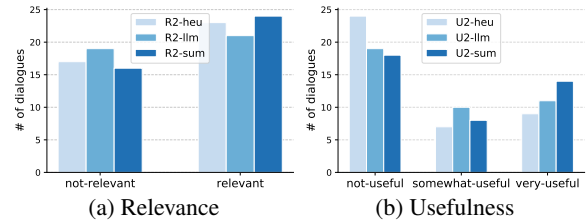


Figure 2: Distribution of (a) relevance and (b) usefulness ratings when annotators have access to additional context in V2 Phase 2.

agreement (Carletta, 1996; Boguslav and Cohen, 2017). In order to understand how the amount of dialogue context influences the quality of ratings by annotators, we calculate the agreement between annotators for both relevance and usefulness across the three variations; see Table 1. To address potential randomness in relevance ratings, given the binary scale, we randomly drop one rating from each dialogue and compute the agreement. We repeat this process for each annotator and calculate an average Cohen’s Kappa score. For usefulness, we compute Kappa for each pair of

Table 1: Inter annotator agreement (Cohen kappa) and Tau correlation for relevance and usefulness across the three experimental setups in Phase 1.

Aspect	Variation	Kappa	τ_b
Relevance	R2	0.53	0.47
	R5	0.61	0.49
	R10	0.70	0.61
Usefulness	U2	0.64	0.54
	U5	0.68	0.60
	U10	0.56	0.41

annotators and then calculate the average. We assess the significance of the agreement using the Chi-squared method. All Kappa scores are statistically significant ($p \leq 0.05$).

We observe an increase in the Kappa and tau score as the dialogue context increases from R2 to R10. Despite the lack of context in R2, there is a moderate level of agreement regarding the relevance of the current turn. With the introduction of more context in R5 and R10, comes an increase in agreement regarding the relevance of the current turn (see Table 1). Providing additional dialogue context seems to lead to higher levels of consensus among annotators. This is likely due to dataset characteristics: users tend to express their preferences early in the dialogue, rather than in subsequent exchanges. Hence, in the case of R2, which only includes the current turn, when the user’s utterance is incomplete, lacking an explicit expression of

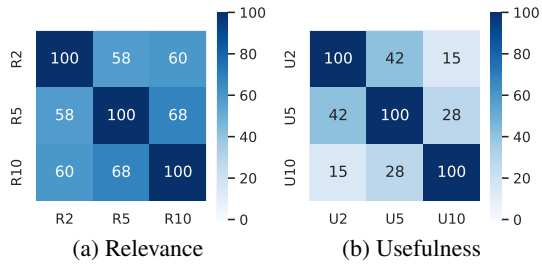


Figure 3: The percentage of agreement in (a) relevance labels and (b) usefulness labels across the three experimental setups in Phase 1.

their preference, annotators rate more dialogues as relevant compared to R5 and R10. Overall, we conclude that when annotators have insufficient information to come up with a judgment, they tend to judge the system positively, introducing a positivity bias (Park et al., 2018).

We see in Table 1 (row 3) that despite the lack of context in U2, there is substantial agreement regarding the usefulness of the current turn. This is due to the availability of the user’s next utterance, which serves as direct feedback on the system’s response, resulting in higher agreement than for relevance assessment. As more context is provided, there is an even higher level of agreement among annotators regarding the usefulness of the current turn. Access to a short conversation history significantly improves agreement on usefulness.

Surprisingly, despite having access to the entire conversation history in U10, there is a slightly lower level of agreement than in U5. The complete dialogue context may introduce additional complexity or ambiguity in determining the usefulness of the current turn. This occurs when conflicting feedback arises from the user’s next utterance compared to the previous dialogue context. E.g., when the system repeats a recommendation that the user has already watched or stated before, and the user expresses their intent to watch the movie in the next utterance, it leads to divergent labels. Similar trend is observed with the tau correlations though the values are lower compared to kappa score.

Label consistency across conditions. We examine the impact of varying amounts of dialogue context on the consistency of crowdsourced labels across the three variations for relevance and usefulness and report the percentage of agreement in Fig. 3. We observe moderate agreement (58.54%) between annotations of R2 and R5, suggesting that annotators demonstrate a degree of consistency in their assessments when provided with different amounts of context. This trend continues with R2 and R10,

Table 2: Inter annotator agreement (Cohen kappa) and Tau correlation for relevance and usefulness across the three experimental setups in Phase 2. significant ($p \leq 0.05$)

Aspect	Variation	Kappa	τ_b
Relevance	R2-heu	0.75	0.54
	R2-sum	0.60	0.45
	R2-llm	0.51	0.44
Usefulness	U2-heu	0.71	0.59
	U2-sum	0.63	0.49
	U2-llm	0.53	0.44

where the agreement increases slightly to 60.98%. The most notable increase is between R5 and R10 (68.29%). As annotators were exposed to progressively broader contextual information, their assessments became more consistent.

Usefulness behaves differently. We observe moderate agreement (41.71%) between U2 and U5, indicating a degree of consistency in annotator assessments within this range of context. A notable decrease in agreement is evident when comparing U5 and U10, down to 28.3% agreement. The most substantial drop is observed between U2 and U10, yielding a mere 14.63% agreement. These findings emphasize the significant impact of context on the consistency of usefulness annotations. For usefulness assessment providing annotators with a more focused context, improves their agreement.

With respect to **RQ1**, we note considerable differences in the labels assigned by annotators as we vary the amount of dialogue context. As the context expands, annotators incorporate more information into their assessments, resulting in context-specific labels. Annotator judgments are shaped not only by response quality but also by the broader conversation. This highlights the complexity of the task and the need for a carefully designed annotation methodology that considers contextual variations. These findings emphasize the significance of dialogue context in annotator decision-making.

3.3 RQ2: Impact of type of previous context available in crowdsourcing

Label quality. In Phase 2, our experiments aim to establish the impact of presenting annotators with different types of context during crowdsourcing. Different from conventional dialogue context, we provide the annotators with the dialogue summary (V2-sum), the user’s information need in the dialogue (V2-heu and V2-llm). We also aim to uncover if we can improve the quality of the crowd-

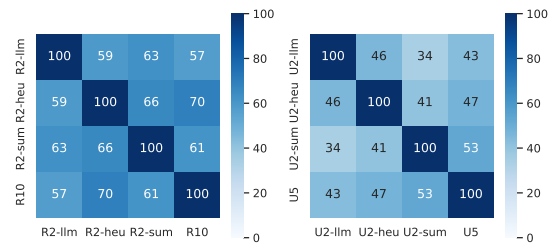
sourced labels in V2 to match those in V10. We calculate the Cohen’s Kappa similar to Section 3.2; see Table 2.

The heuristic approach (R2-heu) yields the highest agreement (kappa and tau), indicating a noteworthy degree of agreement in relevance assessments. The LLM-generated context (R2-llm and R2-sum) results in a moderate to substantial level of agreement, signifying a reasonable level of agreement regarding the relevance of the system response. We observe similar results for usefulness. The heuristic approach (U2-heu) again leads with the highest level of agreement (0.71 and 0.59), U2-sum follows with a kappa score of 0.63, while U2-llm has a kappa score of 0.53. This high level of agreement (kappa) for the two aspects indicates the quality of the labels; the additional context provided, generated either heuristically or with LLMs, is effective in conveying relevant information to annotators, leading to more consistent assessments.

For both relevance and usefulness, V2-heu consistently improves agreement among annotators, while the LLM-generated context (V2-llm and V2-sum) has a substantially lower agreement than V10. This difference reflects the limitations of LLMs in capturing context and generating a factual summary. While they generate coherent text, LLMs sometimes fail to correctly represent the sequential order of the dialogue and users’ language patterns.

Label consistency across conditions. In Figure 4a we report the agreement between the setups in Phase 2 and compare them to R10 and U5 due to their high inter-annotator agreement (IAA) and label consistency. For the relevance annotations, varying levels of agreement emerge. There is substantial agreement between R2-heu and R2-llm (59.36%), showing a significant overlap in the labels assigned using both methods, although there are instances where annotators differ in their assessments of relevance. R2-sum exhibits moderate label agreement with R2-llm (62.74%) and R2-heu (65.67%), pointing to relatively similar label assignments across the setups.

We observe similar results for usefulness in Figure 4b. Though the heuristically generated approach shows a high IAA, for usefulness U2-sum has a high agreement with all other setups. Though annotators agreed on a single label in the R2-heu, in some cases the label may have been the wrong one. We note slightly low agreement levels for a similar label between the three setups, consistent



(a) Relevance (b) Usefulness
Figure 4: The percentage of agreement in (a) relevance labels and (b) usefulness labels across the three experimental setups in Phase 2.

with results in Phase 1. Unlike relevance, usefulness was rated on a scale of 1–3, thus reducing the chance of two setups having the same label when a different type of context is provided.

Regarding **RQ2**, the heuristic approach demonstrates higher consistency in both IAA and label consistency across conditions for relevance than for usefulness. Providing annotators with the user’s initial utterance expressing their preference, particularly in scenarios lacking context, can significantly enhance the quality and consistency of crowdsourced labels. This approach can yield performance comparable to a setup involving the entire dialogue, without imposing the cognitive load of reading an entire conversation on annotators. This streamlines the annotation process and maintains high-quality results, offering a practical strategy for obtaining reliable labels for dialogue evaluation.

4 Discussion and Implications

Our findings reveal intriguing insights into the impact of context size and type on crowdsourced relevance and usefulness labels for TDS. Expanding the dialogue context from V2 to V10 significantly improves agreement among annotators, indicating that annotators rely on comprehensive context to make more accurate assessments. This trend does not hold for usefulness, where we notice a decrease in agreement when all previous dialogue context is available. The optimal amount of context required for reliable labels relies on the aspect evaluated.

Consistent with prior work (Eickhoff, 2018; Kazai et al., 2011a), we observe an inconsistency in relevance labels across variations, with the same system response being rated differently depending on the context provided. Given the lack of label consistency across variations, future studies should carefully tailor their annotation task design and test various settings to ensure high-quality and consistent labels. Additionally, much care should be taken when comparing the performance of a system across several datasets when labels are crowd-

sourced with a different strategy to ensure a fair comparison as models similar to humans can be sensitive to the annotation strategy (Kern et al., 2023; Kadasi and Singh, 2023).

We also analyzed data from the open-ended question asking annotators about their experience with the annotation task. Annotators note that dialogue summaries fail to convey a user’s emotion, limiting their annotation process. Additionally, lower accuracy of the context generated by an LLM may lead to low agreement among annotators. This signifies the importance of carefully considering the quality and accuracy of generated content in the evaluation process. We provide examples in Section A.5 in the appendix. While there may be constraints in presenting user information need and dialogue summary as dialogue context, one key consideration to take into account is the cognitive load of annotators. Providing a shorter, focused context reduces the cognitive burden on annotators, allowing them to devote more attention to actually evaluating a response. This not only streamlines the annotation process but also helps maintain high-quality results. Reducing the amount of content to be assessed may lead to faster annotation times without compromising the quality of ratings (Santhanam et al., 2020). Another approach to leveraging LLMs in annotation, is for requestors to consider co-annotation (Li et al., 2023) between humans and LLMs by asking annotators to use LLMs to generate a short summary for a long dialogue during evaluation.

Optimal context varies by the aspect under evaluation, challenging the idea of a universal strategy. The consistent reliability of automatic methods suggests their potential as dependable tools for evaluation. This implies their use in generating supplementary context, eliminating the need for manual determination of context amounts. This streamlines evaluation, enhancing efficiency in context-driven evaluations for TDS. For data lacking topic or preference shifts, heuristics perform effectively. However, LLMs are recommended for shifting conditions, showcasing adaptability not easily discernible with heuristics.

5 Related Work

We review related work not covered in the paper so far. Several user-centric dialogue evaluation metrics (Mehri and Eskenazi, 2020; Ghazarian et al., 2019; Huang et al., 2020) have been proposed. For TDSs, high-level dimensions such as user satisfaction (Kiseleva et al., 2016; Al-Maskari

et al., 2007) and fine-grained metrics such as relevance and interestingness (Siro et al., 2022) have gained interest. Due to the ineffectiveness of standard evaluation metrics such as ROUGE (Lin, 2004), BLEU (Papineni et al., 2002), which show poor correlation with human judgements (Deriu et al., 2021), a significant amount of research on these metrics relies on crowdsourcing dialogue evaluation labels to improve correlation with actual user ratings. Crowdsourcing ground-truth labels has gained momentum in information retrieval (IR) for tasks like search relevance evaluation (Alonso et al., 2008) and measuring user satisfaction in TDS. A major challenge is ensuring quality and consistency of crowdsourced labels. Task design and annotators’ behavioral features and demographics can affect the quality of the collected labels (Pei et al., 2021; Kazai et al., 2012; Hube et al., 2019). Kazai et al. (2013) examine how effort and incentive influence the quality of labels provided by assessors when making relevance judgments. Other factors such as judgment scale (Roitero et al., 2021; Novikova et al., 2018), annotator background (Roitero et al., 2020; Kazai et al., 2011b), and annotators’ demographics (Difallah et al., 2018) have also been studied. Most studies focus on search systems, not dialogue systems. Closer to our work, Santhanam et al. (2020) study the effect of cognitive bias in the evaluation of dialogue systems. Providing an anchor to annotators introduces anchoring bias, where annotators’ ratings are close to the anchor’s numerical value. Like Santhanam et al. (2020), we focus on the effect of task design on the evaluation of TDSs. In particular, we investigate how the amount and type of dialogue context provided to annotators affect the quality and consistency of evaluation labels and the annotator experience during the evaluation task.

6 Conclusion

In this study we explored the impact of context size and type on crowdsourced relevance and usefulness ratings’ quality and consistency. Addressing two primary questions, RQ1 revealed increased annotator agreement with larger context sizes (V2 to V10). RQ2 investigated different context types, where the heuristic approach consistently showed higher agreement. Notably, leveraging only the user’s first utterance improved label consistency without revealing the entire dialogue, addressing RQ2. Our findings contribute to understanding the experimental design’s effect on TDS evaluation.

736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786

Limitations

In acknowledging the scope of our study, we identify several limitations that warrant consideration.

First, our focus has been primarily directed towards two evaluation dimensions – *relevance* and *usefulness* – for task-oriented dialogue systems (TDSs). While these aspects provide valuable insights, it is crucial to acknowledge that the broader landscape of TDSs evaluation involves a spectrum of fine-grained aspects and metrics explored in previous studies. Future research endeavors could benefit from a more comprehensive examination of additional dimensions to capture the nuanced performance of TDS.

Second, the outcomes derived from our study may exhibit a degree of task or dataset specificity. To ensure the applicability and generalizability of our findings, it is imperative to undertake further investigations to ascertain the extent to which these findings can be extrapolated across different tasks and datasets.

Third, the absence of actual user ratings introduces a caveat in claiming an optimal strategy for presenting previous dialogue history in crowdsourcing tasks. Despite this limitation, we highlight the noteworthy observation of high label consistency for R10 and U5, which served as our basis for comparison.

Lastly, it is crucial to note that our study is exploratory in nature. In subsequent research initiatives, we aim to augment the robustness of our findings by conducting investigations on a larger-scale dataset. This expansion aims to provide a more comprehensive understanding of the dynamics involved in the evaluation of TDSs. Following previous work by Kazai et al. (2012, 2013), we would also want to understand the effect of annotator background: experience of interacting with conversational system or prior experience in doing the annotation task on label consistency for TDSs.

In future research work, we will explore the use of open-source LLMs, like Llama-chat (Touvron et al., 2023), to facilitate a more transparent and reproducible experimental framework. Considering the closed-source nature of GPT-4 and its potential impact on the reproducibility of Phase 2 of our work, experimenting with open-source alternatives becomes imperative. The inclusion of open-source models in subsequent studies would not only address concerns related to accessibility but also foster a collaborative and open scientific

environment.

Ethical Considerations

Intended use of the data and mitigation against misuse

The collected data in this research is expressly intended for research purposes, specifically to advance the understanding of conversational movie recommendation systems. The primary objective is to contribute valuable insights to the field of natural language processing and improve the design and evaluation of recommendation dialogue systems. We will be providing open access to our datasets for use in future research under the MIT License.

Anotator diversity

All participants in this research were master workers recruited exclusively from the United States through Amazon Mechanical Turk (MTurk). While this selection ensured a level of language proficiency and familiarity with the context, it is crucial to note that the findings of this study may not generalize universally due to the specific demographic representation. The restriction to U.S.-based annotators may introduce a limitation in terms of cultural diversity and global perspectives, influencing the external validity of the study.

Annotator bias

Despite the provision of detailed instructions and examples to annotators, potential biases may still arise during the evaluation process due to the diverse backgrounds of the annotators. Cultural biases may be more pronounced if annotators from different cultural backgrounds interpret movie preferences, relevance, or usefulness in divergent ways. Subjective biases may also be influenced by the diverse interpretations of guidelines, as individuals from different backgrounds may have distinct views on dimensions like “relevance” or “usefulness.”

To mitigate these potential biases, continuous monitoring and feedback mechanisms were incorporated into the study design. Additionally, the study refrained from disclosing the specific research angle to annotators to prevent potential biases related to the research objectives.

787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830

831
832
833
834
835
836
837
838

839
840
841

842
843

844
845
846
847
848
849
850
851
852

853
854
855
856
857
858
859
860

861
862
863
864
865
866

867
868
869

870
871
872
873
874
875

876
877
878
879

880
881
882
883
884
885

References

Azzah Al-Maskari, Mark Sanderson, and Paul Clough. 2007. [The relationship between IR effectiveness measures and user satisfaction](#). In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 773–774, New York, NY, USA. Association for Computing Machinery.

Omar Alonso, Daniel E. Rose, and Benjamin Stewart. 2008. [Crowdsourcing for relevance evaluation](#). *SIGIR Forum*, 42(2):9–15.

Amazon Mechanical Turk. 2023. <https://www.mturk.com>.

Mayla Boguslav and Kevin Bretonnel Cohen. 2017. [Inter-annotator agreement and the upper limit on machine performance: Evidence from biomedical natural language processing](#). In *MEDINFO 2017: Precision Healthcare through Informatics - Proceedings of the 16th World Congress on Medical and Health Informatics, Hangzhou, China, 21-25 August 2017*, volume 245 of *Studies in Health Technology and Informatics*, pages 298–302. IOS Press.

Adam Bouyamourn. 2023. [Why LLMs hallucinate, and how to get \(evidential\) closure: Perceptual, intensional, and extensional learning for faithful natural language generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3181–3193. Association for Computational Linguistics.

Paweł Budzianowski and Ivan Vulić. 2019. [Hello, it’s GPT-2 - How can I help you? Towards the use of pre-trained language models for task-oriented dialogue systems](#). In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 15–22, Hong Kong. Association for Computational Linguistics.

Jean Carletta. 1996. [Assessing agreement on classification tasks: The kappa statistic](#). *Computational Linguistics*, 22(2):249–254.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. [A survey on evaluation of large language models](#). *CoRR*, abs/2307.03109.

Jan Deriu, Álvaro Rodrigo, Arantxa Otegi, Guillermo Echegoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2021. [Survey on evaluation methods for dialogue systems](#). *Artif. Intell. Rev.*, 54(1):755–810.

Djellel Eddine Difallah, Elena Filatova, and Panos Ipeirotis. 2018. [Demographics and dynamics of mechanical turk workers](#). In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, pages 135–143. ACM.

Carsten Eickhoff. 2018. [Cognitive biases in crowdsourcing](#). In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, WSDM 2018, Marina Del Rey, CA, USA, February 5-9, 2018*, pages 162–170. ACM.

Guglielmo Faggioli, Laura Dietz, Charles L. A. Clarke, Gianluca Demartini, Matthias Hagen, Claudia Hauff, Noriko Kando, Evangelos Kanoulas, Martin Potthast, Benno Stein, and Henning Wachsmuth. 2023. [Perspectives on large language models for relevance judgment](#). In *Proceedings of the 2023 ACM SIGIR International Conference on Theory of Information Retrieval, ICTIR 2023, Taipei, Taiwan, 23 July 2023*, pages 39–50. ACM.

Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2022. [A survey on dialogue summarization: Recent advances and new frontiers](#). In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pages 5453–5460. ijcai.org.

Sarik Ghazarian, Johnny Wei, Aram Galstyan, and Nanyun Peng. 2019. [Better automatic evaluation of open-domain dialogue systems with contextualized embeddings](#). In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 82–89, Minneapolis, Minnesota. Association for Computational Linguistics.

Lishan Huang, Zheng Ye, Jinghui Qin, Liang Lin, and Xiaodan Liang. 2020. [GRADE: Automatic graph-enhanced coherence metric for evaluating open-domain dialogue systems](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9230–9240, Online. Association for Computational Linguistics.

Christoph Hube, Besnik Fetahu, and Ujwal Gadiraju. 2019. [Understanding and mitigating worker biases in the crowdsourced collection of subjective judgments](#). In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI ’19*, page 1–12, New York, NY, USA. Association for Computing Machinery.

Pritam Kadasi and Mayank Singh. 2023. [Unveiling the multi-annotation process: Examining the influence of annotation quantity and instance difficulty on model performance](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 1371–1388. Association for Computational Linguistics.

Gabriella Kazai. 2011. [In search of quality in crowdsourcing for search engine evaluation](#). In *Advances in Information Retrieval - 33rd European Conference on IR Research, ECIR 2011, Dublin, Ireland, April 18-21, 2011. Proceedings*, volume 6611 of *Lecture Notes in Computer Science*, pages 165–176. Springer.

Gabriella Kazai, Jaap Kamps, Marijn Koolen, and Natasa Milic-Frayling. 2011a. [Crowdsourcing for](#)

943	book search evaluation: impact of hit design on comparative system ranking. In <i>Proceeding of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, July 25-29, 2011</i> , pages 205–214. ACM.	<i>Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada</i> , pages 9748–9758.	1000 1001 1002
944			
945			
946			
947		Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	1003 1004 1005 1006
948			
949	Gabriella Kazai, Jaap Kamps, and Natasa Milic-Frayling. 2011b. Worker types and personality traits in crowdsourcing relevance labels . In <i>Proceedings of the 20th ACM International Conference on Information and Knowledge Management, CIKM '11</i> , page 1941–1944, New York, NY, USA. Association for Computing Machinery.	Jiaxin Mao, Yiqun Liu, Ke Zhou, Jian-Yun Nie, Jingtao Song, Min Zhang, Shaoping Ma, Jiashen Sun, and Hengliang Luo. 2016. When does relevance mean usefulness and user satisfaction in web search? In <i>Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval, SIGIR 2016, Pisa, Italy, July 17-21, 2016</i> , pages 463–472. ACM.	1007 1008 1009 1010 1011 1012 1013 1014
950			
951			
952			
953			
954			
955			
956	Gabriella Kazai, Jaap Kamps, and Natasa Milic-Frayling. 2012. The face of quality in crowdsourcing relevance labels: Demographics, personality and labeling accuracy . In <i>Proceedings of the 21st ACM International Conference on Information and Knowledge Management, CIKM '12</i> , page 2583–2586, New York, NY, USA. Association for Computing Machinery.	Shikib Mehri and Maxine Eskenazi. 2020. USR: An unsupervised and reference free evaluation metric for dialog generation . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 681–707, Online. Association for Computational Linguistics.	1015 1016 1017 1018 1019 1020
957			
958			
959			
960			
961			
962			
963			
964	Gabriella Kazai, Jaap Kamps, and Natasa Milic-Frayling. 2013. An analysis of human factors and label accuracy in crowdsourcing relevance judgments . <i>Information Retrieval</i> , 16(2):138–178.	Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2018. RankME: Reliable human ratings for natural language generation . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 72–78, New Orleans, Louisiana. Association for Computational Linguistics.	1021 1022 1023 1024 1025 1026 1027 1028
965			
966			
967			
968	Christoph Kern, Stephanie Eckman, Jacob Beck, Rob Chew, Bolei Ma, and Frauke Kreuter. 2023. Annotation sensitivity: Training data collection methods affect model performance . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 14874–14886. Association for Computational Linguistics.	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.	1029 1030
969			
970			
971			
972			
973			
974			
975	Julia Kiseleva, Kyle Williams, Jiepu Jiang, Ahmed Hassan Awadallah, Aidan C. Crook, Imed Zitouni, and Tasos Anastasakos. 2016. Understanding user satisfaction with intelligent assistants . In <i>Proceedings of the 2016 ACM on Conference on Human Information Interaction and Retrieval, CHIIR '16</i> , page 121–130, New York, NY, USA. Association for Computing Machinery.	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: A method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02</i> , page 311–318, USA. Association for Computational Linguistics.	1031 1032 1033 1034 1035 1036
976			
977			
978			
979			
980			
981			
982			
983	Margaret Li, Jason Weston, and Stephen Roller. 2019. ACUTE-EVAL: improved dialogue evaluation with optimized questions and multi-turn comparisons . <i>CoRR</i> , abs/1909.03087.	Kunwoo Park, Meeyoung Cha, and Eunhee Rhim. 2018. Positivity bias in customer satisfaction ratings . In <i>Companion of the The Web Conference 2018 on The Web Conference 2018, WWW 2018, Lyon , France, April 23-27, 2018</i> , pages 631–638. ACM.	1037 1038 1039 1040 1041
984			
985			
986			
987	Minzhi Li, Taiwei Shi, Caleb Ziems, Min-Yen Kan, Nancy F. Chen, Zhengyuan Liu, and Diyi Yang. 2023. Coannotating: Uncertainty-guided work allocation between human and large language models for data annotation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1487–1505. Association for Computational Linguistics.	Mihir Parmar, Swaroop Mishra, Mor Geva, and Chitta Baral. 2023. Don't blame the annotator: Bias already starts in the annotation instructions . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023</i> , pages 1771–1781. Association for Computational Linguistics.	1042 1043 1044 1045 1046 1047 1048 1049
988			
989			
990			
991			
992			
993			
994			
995			
996	Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. 2018. Towards deep conversational recommendations . In <i>Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada</i> , pages 9748–9758.	Weiping Pei, Zhiju Yang, Monchu Chen, and Chuan Yue. 2021. Quality control in crowdsourcing based on fine-grained behavioral features . <i>Proc. ACM Hum. Comput. Interact.</i> , 5(CSCW2):442:1–442:28.	1050 1051 1052 1053
997			
998			
999			

1054	Kevin Roitero, Eddy Maddalena, Stefano Mizzaro, and Falk Scholer. 2021. On the effect of relevance scales in crowdsourcing relevance assessments for information retrieval evaluation . <i>Inf. Process. Manag.</i> , 58(6):102688.	1111
1055		1112
1056		1113
1057		1114
1058		1115
1059	Kevin Roitero, Michael Soprano, Shaoyang Fan, Damiano Spina, Stefano Mizzaro, and Gianluca Demartini. 2020. Can the crowd identify misinformation objectively?: The effects of judgment scale and assessor’s background . In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020</i> , pages 439–448. ACM.	1116
1060		1117
1061		1118
1062		1119
1063		1120
1064		1121
1065		1122
1066		1123
1067	Sashank Santhanam, Alireza Karduni, and Samira Shaikh. 2020. Studying the effects of cognitive biases in evaluation of conversational agents . In <i>CHI ’20: CHI Conference on Human Factors in Computing Systems, Honolulu, HI, USA, April 25-30, 2020</i> , pages 1–13. ACM.	1124
1068		1125
1069		
1070		
1071		
1072		
1073	Alexander Schmitt and Stefan Ultes. 2015. Interaction quality: Assessing the quality of ongoing spoken dialog interaction by experts - and how it relates to user satisfaction . <i>Speech Commun.</i> , 74:12–36.	
1074		
1075		
1076		
1077	Clemencia Siro, Mohammad Aliannejadi, and Maarten de Rijke. 2022. Understanding user satisfaction with task-oriented dialogue systems . In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’22</i> , page 2018–2023, New York, NY, USA. Association for Computing Machinery.	
1078		
1079		
1080		
1081		
1082		
1083		
1084	Clemencia Siro, Mohammad Aliannejadi, and Maarten de Rijke. 2023. Understanding and predicting user satisfaction with conversational recommender systems . <i>ACM Trans. Inf. Syst.</i> , 42(2):Article 55.	
1085		
1086		
1087		
1088	Weiwei Sun, Shuo Zhang, Krisztian Balog, Zhaochun Ren, Pengjie Ren, Zhumin Chen, and Maarten de Rijke. 2021. Simulating user satisfaction for the evaluation of task-oriented dialogue systems . In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , page 2499–2506, New York, NY, USA. Association for Computing Machinery.	
1089		
1090		
1091		
1092		
1093		
1094		
1095		
1096	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	
1097		
1098		
1099		
1100		
1101		
1102		
1103		
1104		
1105		
1106		
1107		
1108		
1109		
1110		
	Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288.	
	Chien-Sheng Wu, Steven C. H. Hoi, Richard Socher, and Caiming Xiong. 2020. TOD-BERT: pre-trained natural language understanding for task-oriented dialogue . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 917–929. Association for Computational Linguistics.	

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
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174

A Appendix

In this section we provide supplementary materials used to support our main paper. These materials include: experimental conditions elaborated in Section A.1, quality control measures undertaken to ensure high quality crowdsourced labels and generated supplementary context in Section A.2 and the prompts used to generate the supplementary context in Section A.3. In Section A.4 we include the annotation instructions and screen dumps of our annotation task. Section A.5 shows sample supplementary context generated by GPT-4.

A.1 Experimental conditions

We list the experimental conditions used for our crowdsource experiments in Table 3.

A.2 Data quality control

Generated user information need and summary. To address the potential hallucination of LLMs (Chang et al., 2023), we implemented a quality control process for the generated user information needs and summaries, ensuring their coherence and factual accuracy. We automatically cross-reference the movies mentioned in both the input dialogues and the summaries. A summary must contain at least two-thirds of the movies mentioned in the input dialogue to be considered valid. If this criterion is not met, the summary is discarded, and a new one is generated following the specified prompt requirements. In total, we discarded and regenerated 15 dialogue summaries. To further ensure coherence, we randomly sampled 30% of the generated summaries and information needs. The authors reviewed them to confirm their coherence and alignment with the information presented in the input dialogue. This process enhanced the quality and reliability of the generated content.

Crowdsourced labels. To ensure a high quality of the collected data, we incorporated attention-checking questions into the HIT. Annotators were required to specify the number of utterances in the dialogues they were evaluating and to identify the last movie mentioned in the system response being evaluated. 10% of the HITs were rejected and returned back to collect new labels. In total, we gathered 1440 data samples from the crowdsourcing task, spanning six variations for relevance and usefulness. We employed majority voting to establish the final relevance and usefulness dialogue label.

A.3 Prompts

In Table 4 we show the final prompts used to generate the user information and dialogue summary with GPT-4.

A.4 Annotation instructions and screen dumps

Table 5 details the annotation instructions for the relevance and usefulness evaluations. In Fig 5 and 6 we show the annotation interface used for phase 1 and phase 2 respectively.

A.5 Sample supplementary context

In Table 6 we show sample user information need and summary generated by GPT-4.

1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Table 3: Descriptions of the experimental setups used for the crowdsourcing experiments with corresponding relevance and usefulness labels. Unlike relevance, usefulness includes the user’s next utterance as feedback. A “turn” denotes a user-system exchange.

Variations	Relevance	Usefulness	Description
V2	R2	U2	Current turn with no previous dialogue context
V5	R5	U5	Current turn with three system-user utterances as previous context
V10	R10	U10	Current turn with 7 user-system utterances as previous context
V2-llm	R2-llm	U2-llm	Current turn with an LLM-generated user information need as dialogue context
V2-heu	R2-heu	U2-heu	Current turn with a heuristically generated user information need as dialogue context
V2-sum	R2-sum	U2-sum	Current turn with a dialogue summary as dialogue context

Table 4: Prompts used to generate the supplementary context; user information need and dialogue summary with GPT-4.

Dialogue summary prompt
Below you are provided with dialogues between a user and the system about movie recommendations. Generate a complete short and informative summary extractively which is half the length of the dialogue.
User information need prompt
Given the following user and system dialogue in a movie recommendation conversation, generate a concise user’s goal in a natural manner. State only the goal without extra text. Start the sentence with "the user wants."

User: \${user4}

System: \${system4}

User: \${user5}

System: \${system5}

User: \${user6}

Questions

Now please answer the following question about the highlighted system response.

1. **Is the system response useful?**

- 1 (Low usefulness) - The response inadequately addresses the user’s needs, lacks necessary information, and fails to enhance the overall user experience.
- 2 (Moderate usefulness) - The response partially addresses the user’s needs, provides some information, and contributes to enhancing the overall user experience, but may lack diversity and personalization.
- 3 (High usefulness) - The response effectively addresses the user’s needs, provides comprehensive and accurate information, and significantly enhances the overall user experience with relevance, accuracy, diversity, and personalization.

Figure 5: Annotation interface for phase 1 when evaluating response usefulness for V5

Read the dialogue below carefully and answer the follow up question.

User: \${user4}

User: \${system4}

User: \${user5}

System: \${system5}

User: \${user6}

Summary: \${summary}

Figure 6: Annotation interface for phase 2 when evaluating response usefulness with supplementary context

Table 5: Annotation instructions provided to the annotators for relevance evaluation. The instructions are the same for usefulness apart from the aspect being evaluated.

Introduction
Thank you for helping us out! Below we explain everything in full detail. Please make sure to read the instructions carefully.
Purpose
The aim of this survey is to evaluate the quality of a system’s response. We want to evaluate the dialogue system’s performance and gather insights for improvements. We will ask you to evaluate the system response on one metric, that we will discuss in more detail below.
Scenario Outline
Imagine you are evaluating a dialogue system that generates a response to user queries. Your task is to assess the response based on relevance. We will provide examples and detailed explanations of this criteria below.
Task
In each HIT, you will be presented with a dialogue chunk. Your task is to evaluate the last system response based on the given criteria. Please review the explanations and examples for the criteria to ensure your understanding before proceeding with the evaluation. Keeping the scenario that was outlined above in mind, we would like to ask you to judge the system response on relevance.

Table 6: Sample dialogue summaries as supplementary context generated by GPT-4.

Dialogue 1
User inquires about a good family movie recommendation similar to "Real Steel (2011)" or "The Lego Movie (2014)". System recommends "Super (2010)", an action-comedy about a regular guy who becomes a self-made superhero, describing it as hilarious and entertaining. The user shows interest in this recommendation.
Dialogue 2
The user asked for coming-of-age movie recommendations and mentioned they enjoyed "My Girl (1991)" and "Lucas (1986)". The system suggested watching "The Spectacular Now (2013)", a film where Shailene Woodley stars as a character who forms a bond with a troubled classmate.
Dialogue 3
User seeks a dramatic love story to watch. System recommends "The Notebook (2004)", but the user has watched it, as well as "Titanic (1997)". Both films are favored by the user; they desire to watch something new.
Dialogue 4
The user requests animated movie recommendations following their enjoyment of "The Incredibles (2004)". The system suggests other movies, including "Monsters, Inc. (2001)" and its sequel "Monsters University (2013)", which the user approves. The conversation pivots to the topic of successful sequels, citing "Toy Story 3 (2010)" as an example despite the user’s disagreement, favoring the original movie, "Toy Story (1995)".
Dialogue 5
The user wants to find a thrilling crime movie like "Thor: Ragnarok (2017)" for their weekend. The system suggested they watch "The Snowman (2017)" but the user declined. However, the system then gave another recommendation, "First Kill (2001)".