

Reflect, Not Reflex: Inference-Based Common Ground Improves Dialogue Response Quality

Pei Zhou Hyundong Cho Pegah Jandaghi Dong-Ho Lee Bill Yuchen Lin
Jay Pujara Xiang Ren

Department of Computer Science and Information Sciences Institute
University of Southern California

{peiz, dongho.lee, yuchen.lin, jpujara, xiangren}@usc.edu, {jcho, jandaghi}@isi.edu

Abstract

Human communication relies on *common ground* (CG), the mutual knowledge and beliefs shared by participants, to produce coherent and interesting conversations. In this paper, we demonstrate that current response generation (RG) models produce generic and dull responses in dialogues because they act *reflexively*, failing to explicitly model CG, both due to the lack of CG in training data and the standard RG training procedure. We introduce *Reflect*, a dataset that annotates dialogues with explicit CG (materialized as inferences approximating shared knowledge and beliefs) and solicits 9k diverse human-generated responses each following one common ground. Using *Reflect*, we showcase the limitations of current dialogue data and RG models: less than half of the responses in current data is rated as high quality (sensible, specific, and interesting) and models trained using this data have even lower quality, while most *Reflect* responses are judged high quality. Next, we analyze whether CG can help models produce better quality responses by using *Reflect* CG to guide RG models. Surprisingly, we find that simply prompting GPT3 to “think” about CG generates 30% more quality responses, showing promising benefits to integrating CG into the RG process.¹

1 Introduction

Human communication is a collaborative effort (Grice, 1975; Allwood, 1976; Bohm et al., 2004) where participants strive to achieve *common ground* (CG), consisting of mutual beliefs and common knowledge (Stalnaker, 1978; Clark and Schaefer, 1989; Clark and Brennan, 1991). Conversational AI systems, while able to produce fluent texts, often generate generic and dull dialogue responses (Serban et al., 2017; Zhao et al., 2017), potentially because they do not explicitly model

¹Link to our data and code will be provided on our project page <https://inklab.usc.edu/Reflect/>.

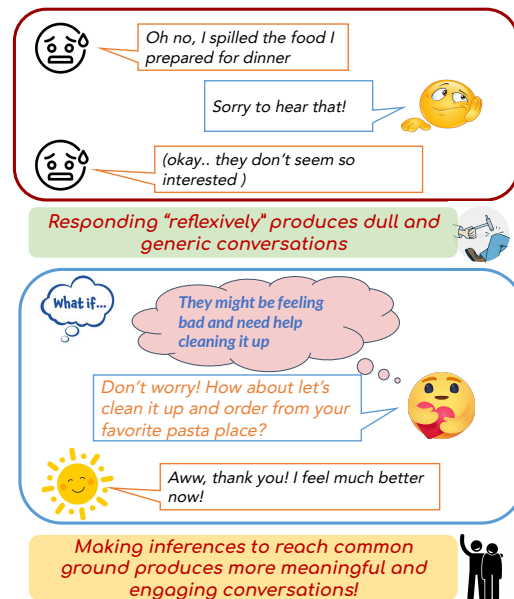


Figure 1: **A motivating example.** We aim to help RG models produce more *human-like* responses instead of generic ones. We argue that integrating common ground by making inferences is crucial.

CG in communication (as illustrated in Figure 1). Specifically, existing models mostly follow a *dialogue history* \rightarrow *response* training paradigm since such data can be easily obtained in the wild, skipping an important middle step that *builds common ground*, which naturally and universally exists in human communication, i.e., *dialogue history* \rightarrow *common ground* \rightarrow *response*. Moreover, the same history can yield numerous responses, predicated on the CG and intent of the responder. We conjecture that the omission of modeling CG explicitly is a crucial bottleneck in RG models because they are directly trained to produce responses without learning *how and why* those responses are uttered.

Modeling common ground between speakers, however, is challenging due to its implicit and subjective nature during conversations (Clark and Schaefer, 1989). Prior work on representing CG either mines noisy commonsense knowledge triples between dialogue history and existing re-

sponses (Zhou et al., 2022) or collects human inferences after reading the whole dialogue as a bystander (Ghosal et al., 2022). Such approaches provide useful augmentation, but post-hoc analysis cannot mirror the generative process and intent of diverse human dialogue. Figure 2 illustrates three paradigms for RG. We argue that truly modeling this generative process requires (1) articulating CG prior to the response; (2) generating responses conditioned on CG; (3) differentiating response generation based on different types of CG.

To this end, we formalize common ground in dialogues as *inferences* made by *one* participant to approximate potential beliefs shared by *other* participants, as shown in Figure 1. In this work, we instantiate inferences as *question-answer (QA) pairs in natural language (NL)* such as “What might happen later?” “They might need to clean the floor” to elicit others’ beliefs, inspired by inquiry-based dialogic learning (Bruner, 1961; Habermas, 1985; Wells, 2000). Another critical aspect of CG is its multi-dimensional nature, i.e., given the same dialogue context, different plausible inferences can be made, which then lead to different responses. Following these principles, we create a novel dialogue resource with multiple explicitly human-annotated common ground, each of which is further substantiated as a next-turn response continuing the conversations (an example of expanded CG and responses for one context shown in Figure 3).

We design a two-stage data collection process by first asking crowdsourcing workers to answer different inference questions eliciting beliefs about CG (e.g., *what is the speaker feeling right now?*) Answers rely on common sense, and adopt the point of view of the conversational respondent. We use these QA pairs to approximate various (non-exhaustive) inference dimensions to extend the common ground (e.g., empathy and event causality). Our second step converts these CG into dialogue responses by asking different workers to write a coherent response based on the answer/inference collected in the first stage. Our collected data *Reflect* contains 9k diverse responses from 600 dialogue contexts, based on 5 inference dimensions for CG.

Using *Reflect*, we first test our hypothesis that explicitly modeling CG and using CG to construct responses creates more engaging conversations. We conduct human evaluation to compare the quality of responses between *Reflect* and “*reflex*” style

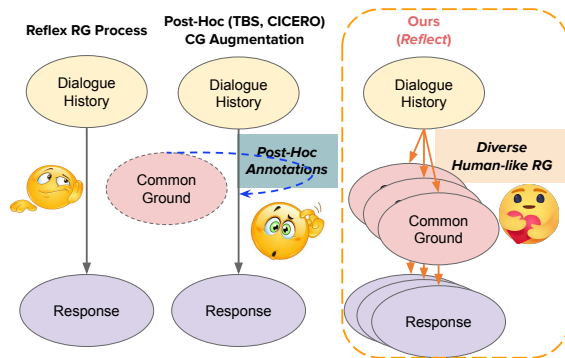


Figure 2: **Illustration of different RG approaches.** Common RG does not model CG explicitly, TBS (Zhou et al., 2022) and CICERO (Ghosal et al., 2022) post-hoc augments dialogues with CG, and we aim to follow natural human communication and first collect CG and then responses based on CG. We also factor in the diversity of plausible responses given a dialogue context that result from distinct CG.

datasets and models in terms of sensibility, specificity, and interestingness. We find that, compared to reflex-prone human-written and machine-generated dialogues, our two-stage data collection process results in more responses that are sensible, specific, and interesting as rated by humans. This highlights limitations of existing data collection procedures and models trained on the data.

Next, we look to study the potential of explicitly modeling CG in dialogue systems to help build models that can create more engaging conversations. As a case study, we use the inference dimensions from *Reflect* and test two simple ways to guide RG using CG. We surprisingly find that simple approaches such as appending an inference question to the dialogue context before the response in the few-shot (FS) in-context examples (from *Reflect*) help GPT3-175B (Brown et al., 2020) generate almost 30% more responses that are deemed sensible, specific, and interesting than vanilla FS learning GPT3 (no inference question). We demonstrate that, when prompted to “*think*” about an inference question (approximated CG), large models such as GPT-3 can create more engaging conversations. We also find that such effect is only shown in large models like GPT-3 as we find BlenderBot-440M (Roller et al., 2021) benefits from fine-tuning on *Reflect*, but appending inference questions does not further increase response quality.

In summary, our contributions are as follows: 1) we operationalize theories about common ground and formalize them for dialogue; 2) we collect the first large-scale (9k responses) dialogue dataset with diverse responses guided by CG and release this resource to facilitate training and evaluation; 3)

we show important limitations of existing dialogue data and RG models that detract from engaging communication; 4) we demonstrate that CG can dramatically improve RG quality even with simple prompting, boosting quality by 30%. The resources and results from this work promise to enable the research community to create and evaluate common ground-aware RG models.

2 Inference-Based Common Ground

We formally introduce the notion of *common ground* in conversations as the implicit variable conditioned on dialogue history and provides conditions to the next-turn response.

2.1 Grounding in Communication

Successful collaborative communication activity relies on mutual understanding of shared knowledge and beliefs (Clark and Brennan, 1991; Bohm et al., 2004) called *common ground*. However, due to *least collaborative effort* (Grice, 1975; Clark and Schaefer, 1989) where communication participants try to minimize the effort spent on contributing to the interaction, establishing CG relies on *signals* other than the surface communication information (i.e., actual utterances in a conversation). While humans in face-to-face communication receive some information from *non-verbal* signals such as gestures and facial expressions, virtual systems such as chatbots often do not have access to such signals. Thus, we argue that they have to rely heavily on another crucial way of getting signals for establishing CG: *making inferences* based on the surface communication utterances and common sense, in order to approximate two humans talking to create engaging conversations.

Furthermore, building CG by making relevant inferences also connects closely with the “*dual process*” theories of human reasoning (Stanovich and West, 2000; Evans, 2003; Kahneman, 2011). We argue that the “*reflexive*” RG is mostly modeling “*System 1*” which is intuitive and associative, but a more deliberative and logical “*System 2*” is lacking.

2.2 Formulating CG in Dialogue

Consider three high-level components in communication efforts: context C (often materialized as dialogue history consisting of a sequence of n contributions $C = c_1, \dots, c_n$), common ground G , and a new contribution continuing the context (often referred to as a “*response*” c_{n+1}). Specifically, for

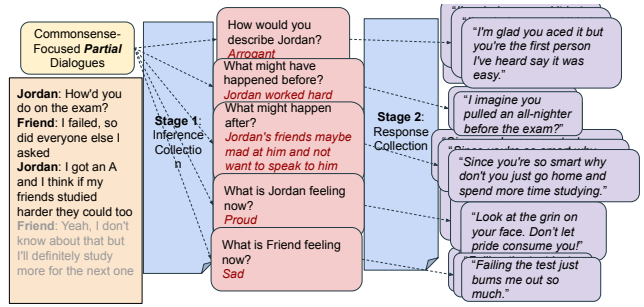


Figure 3: **Reflect collection procedure illustration.** We first collect CG materialized as inferences expressed in QA along different dimensions. Then for each QA pair, we collect multiple responses.

common ground G , we focus on signals gained from inferences and thus materialize G as a list of m potential inferences $G = I_1, \dots, I_m$ conditioned on the context. We furthermore materialize each inference as a QA pair in NL $I_j = (Q_j, A_j)$ (examples included in Figure 3 between Stage 1 and 2). We use QA format to express inferences to mimic inquiry-based dialogic learning (Bruner, 1961; Habermas, 1985; Wells, 2000) and follow empirical evidence that neural models take in QA-format knowledge effectively (Shwartz et al., 2020; Zhou et al., 2022).

3 Collecting Reflect Data

Here we describe how we collect *Reflect*, a novel large-scale dialogue dataset with diverse human-annotated inference-based CG and grounded responses. An overview of the procedure with examples are shown in Figure 3. We first select base dialogues from a dataset that is constructed without considering CG and only has one plausible response for each context (3.1). Then we aim to expand and collect multiple responses based on different inference dimensions. We introduce a two-stage process to first crowdsource potential inferences people make in conversations (3.2) and then ask a second cohort of workers to generate diverse responses based on the inferences (3.3). We designed a two-stage data collection to 1) collect multiple, diverse responses based on each CG; 2) to allow response writers to validate CG as high quality, generic common sense inferences. Finally, we include discussions of data quality assurance (3.4).

3.1 Pre-Collection: Selecting Base Dialogue Turns for Expansion

Our first step is to select base dialogues and dialogue turns to expand on, in terms of both inference-

Inference Dimensions	Inference Questions
Attributes of Speaker	How would you describe Speaker?
Potential prerequisites	What might have happened before?
Potential consequences	What might happen after?
Speaker Emotion States	What is Speaker feeling now?
Responder Emotion States	What is Responder feeling now?

Table 1: Inference dimensions and corresponding questions

based CG and more potential responses following the CG. One important criterion for base turns is that they should not be “social glue” turns such as “*You are welcome*” in responding to “*Thank you!*” We aim at expanding turns that have semantically-rich dialogue context, enabling different plausible inferences to be made. After investigation of existing datasets, we use dialogues from Commonsense-Focused Dialogues (Zhou et al., 2021) that are converted to dialogues from SocialQA (Sap et al., 2019b) contexts. We chose this dialogue data because SocialQA (crowdsourced from ATOMIC (Sap et al., 2019a), an if-then inferential commonsense knowledge base) contains everyday situations where people can make various inferences on. Then, to select what turns to expand on, we use simple heuristics and select the turn that has the largest semantic overlap with the event in SocialQA using SentenceBERT (Reimers and Gurevych, 2019).

3.2 Stage 1. Inference Collection

Our first goal is to collect potential inferences people might make (e.g. “*they might be feeling bad*”) given conversation contexts C to approximate *common ground*. Each inference I_j is further materialized as a QA pair (Q_j, A_j) along multiple inference dimensions as formulated in Section 2.2.

Inference Knowledge Schema We adopt inference dimensions from ATOMIC2020 (Hwang et al., 2021) since it focuses on social commonsense inferences based on everyday scenarios. Specifically, we conduct a pilot study to choose 5 dimensions from the 14 dimensions, consolidating those that overlap (e.g., “what might happen later” and “what would others likely want to do after”) in the context of dialogues. Our final five dimensions for conversation-based inference dimensions are shown in Table 1.

Crowdsourcing Our Stage 1 crowdsourcing task is: given a dialogue context, imagine that you are participating as the responder and write answers to the 5 inference questions (more details in Appendix). We recruit a group of around 30 crowd-

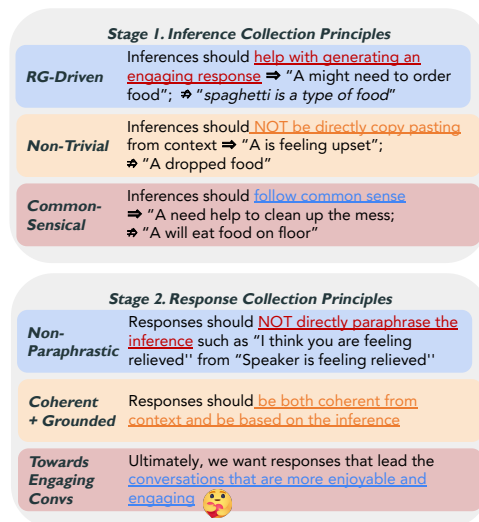


Figure 4: Crowdsourcing principles for two-stage collection.

Resources	Source	Makes Sense	Relevant
TBS (Zhou et al., 2022)	ConceptNet	83.7%	81.0%
CICERO (Ghosal et al., 2022)	Human	86%	96%
Reflect (Ours)	Human	93%	96%

Table 2: Human evaluation on **inference (CG) quality**. We compare inferences from three resources and compare their sensibility and relevance to dialogue context.

sourcing workers from Amazon’s Mechanical Turk platform (AMT) who are native English speakers and provide detailed feedback. Specifically, after carefully reading collected inferences from pilot studies, we provide feedback to turkers by stressing on several principles to make the inferences collected more closely approximate CG, shown in Figure 4.

3.3 Stage 2. Response Collection

After the first stage, we have collected 5 inferences (approximating CG) in the form of QA pairs for each dialogue context. Our next step is to collect next-turn responses given *both* the dialogue context and the collected inference-based CG along different dimensions. To account for diversity in responses, for each dialogue context we ask three Turkers to write a next-turn response based on each of the given inferences, yielding 15 responses for each dialogue context. Similarly to Stage 1, we communicate our collection principles to workers to improve the collected data quality (Figure 4). Both Stage 1 and Stage 2 UI and positive/negatives examples for workers are included in Appendix.

3.4 Quality Control and Analysis

Quality check for Inference Collection In our second stage for response collection, we ask workers an additional question before writing a response: “*do you think the shown inference answer is a valid*

reaction from the responder?” as a way to check the quality of the first collection stage results. We find that less than 7% (200/3000) of the inferences are deemed implausible by second stage workers and only keep the inferences where most workers agree that the inferences are plausible.

Quality check for Response Collection To check quality for our stage 2 response results, we randomly sampled around 5% of collected responses (500) and conduct a manual in-house check for two criteria: 1) is it a sensible continuation from the dialogue context? and 2) is the response based on the inference given? We find that around 93% of the responses are a sensible continuation and 89% are following the inferences given. Further human ratings of our collected grounded dialogue responses showing that our data improves the sensibility, specificity, and interestingness aspects compared to the base responses are included and discussed in Section 4.

Comparison to prior work on representing CG

We compare CG inferences from *Reflect* with TBS (Zhou et al., 2022) and CICERO (Ghosal et al., 2022), two prior work that aims to represent CG in dialogues using either ConceptNet (Speer et al., 2017) knowledge triples or post-hoc human annotations, respectively. Note we only compare inferences (CG) since neither collects new dialogue responses grounded in the inferences, and only consider a single response per context. Comparison results on sampled 100 inferences for each resource are shown in Table 2 where we find that inferences in *Reflect* are rated as make more sense and relevant to dialogue context than the prior dataset.

4 Limitations of Reflex-Prone Dialogue Data and Models

Most existing open-domain dialogue datasets are either crowdsourced by workers who do not have strong incentives to create engaging conversations (Rashkin et al., 2019; Zhou et al., 2021) or crawled from language learning websites and exams (Li et al., 2017; Cui et al., 2020). Both lack explicit CG. These collection processes can fail to capture engaging human-like conversations through under-specified response criteria. Accordingly, RG models trained on these data may mimic generic patterns. This section aims to demonstrate such limitations by comparing responses from *Reflect* with responses from both the original dialogue

Dimensions	Positive Examples	Negative Examples
Sensibleness	That’s too bad!	Thank you.
Specificity	Did you spill it in the kitchen? Let me help!	Do you need help?
Interestingness	It’s actually blessing in disguise, wanna guess why?	Let’s eat something else.
Quality (SSI)	It’s blessing in disguise, since I ordered extra from your favorite pasta place!	All above

Table 3: Evaluation dimensions for RG with examples (dialogue context from Figure 1).

dataset we expand on and models trained on the *dialogue history* \rightarrow *response* regime.

4.1 Human Evaluation Dimensions-SSI

We evaluate the *quality* of each response by head-to-head comparing across systems along several evaluation criteria. We follow the protocol used by LaMDA (Thoppilan et al., 2022) and measure SSI: sensibleness, specificity, and interestingness. Examples of positive and negative responses are shown in Table 3. Our assumption is that responses that contribute to more engaging conversations should satisfy *all* three dimensions and we refer to them as *quality responses*. We do not consider automatic metrics since they do not yet reliably replace human judgements on open-ended responses, especially for fine-grained evaluation dimensions.

4.2 Comparing Original vs Reflect Responses

First, we compare the quality of responses in previous dialogue datasets with our *Reflect* responses to analyze the effects of explicitly incorporating CG in *human RG*. Here we present results by adopting the aforementioned evaluation protocol on human dialogues, both from the original base dialogues (Zhou et al., 2021) and from our *Reflect* dataset, derived from the same dialogues. We sampled 300 dialogue contexts and asked 3 crowdsourcing workers to rate the three SSI criteria, using majority voting to get final scores (Fleiss-kappa (Fleiss, 1971) agreement is around 0.67). We compare the original next-turn response from the contexts with a randomly sampled one from our *Reflect* responses.

Reflect contains more specific and interesting responses than original dialogues

From human evaluation shown in Figure 5, we observe that our collected *Reflect* data, consists of dialogue responses that are on average more specific (20%) and interesting (13%) than the original data, while having slightly lower sensibility (4%) ratings. One possible contributor to the lower sensibility may be 2-stage collection where a new worker continues dialogues constrained by a specific inference generated by another person. Specifically, when

comparing the percentages of responses that satisfy all three criteria, *i.e.*, quality responses, we find that there are substantially more (18%) such responses in *Reflect* than in original data. This observation raises an interesting question: “*do existing dialogue training datasets capture high quality dialogues?*” Without sensible, specific, and interesting responses to learn from, RG models will necessarily be limited in the quality of their output.

4.3 Comparing Reflex RG vs Reflect Data

We now compare *Reflect* with RG models trained on dialogue data that lacks explicit CG and to directly generate an utterance given a context.

Reflexive model baselines Specifically, we consider models from two categories: medium-sized RG models pre-trained on dialogue data such as BlenderBot (440M parameters)² (Roller et al., 2021) and large-sized language models (LLM) pre-trained on general texts such as GPT3-DaVinci (175B parameters)³ (Brown et al., 2020). We directly use off-the-shelf Blender since it is pre-trained on dialogue data (**Blender**). For GPT3-175B, we apply *few-shot in-context* learning by providing 3 examples of dialogue context and response from existing data (**GPT3-FS**). We manually examine these responses to ensure their quality as demonstrating examples. Then we present a dialogue context from our test data and prompt GPT3 to generate a next-turn response. More details in Appendix A.

Models with no common ground struggle Unsurprisingly, as shown in Figure 6, we find a similar trend as comparing *Reflect* with original dialogue data: both BlenderBot-FT and GPT3-FS generate much fewer *quality* responses (53% and 38%, respectively) that satisfy all criteria and particularly on specificity. This further supports the hypothesis that RG models that learn from no-grounding dialogue responses struggle to capture what constituted meaningful conversations.

5 A Little CG Goes a Long Way

After showing that explicitly integrating inference-based CG helps *humans* produce more specific and interesting dialogue responses, we now test if this

²<https://parl.ai/projects/recipes/>

³<https://beta.openai.com/docs/models/gpt-3>

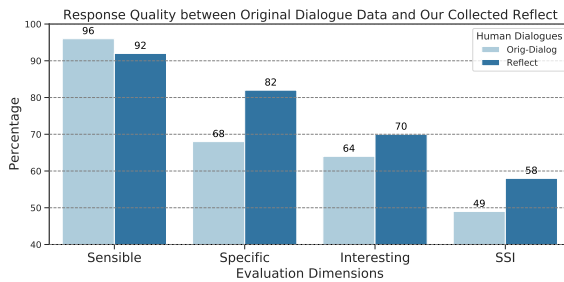


Figure 5: **Human evaluation** comparing human dialogues: original data and our collected Reflect.

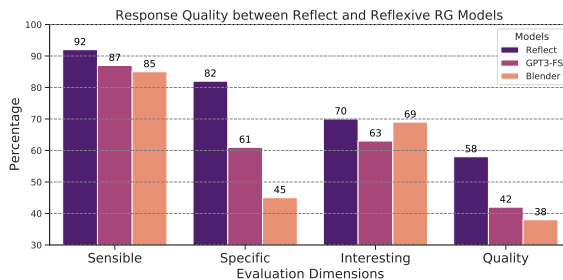


Figure 6: **Averaged performance** comparing *Reflect* with reflexive models. Predictably, RG models trained to directly output responses from dialogue context using no-common ground responses have a large gap in response quality compared to human responses guided by inferences.

also holds for neural RG models. We take the non-exhaustive inference dimensions we used in *Reflect* as case studies to see how CG could improve the quality of existing RG systems’ responses, in terms of the SSI human evaluation (Thoppilan et al., 2022).

5.1 Experiment Setup

Inference-Guided reflect models We attempt to shift models from “*reflexive*” RG to “*reflective*” RG by taking into account of plausible inferences that humans use to build common ground during com-

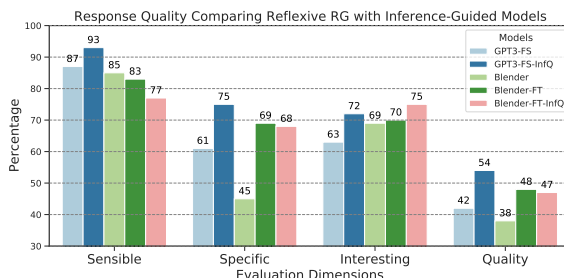


Figure 7: **Averaged performance** comparing before and after reflexive models leveraged inference-guided *Reflect* data. We find that inference-based common ground prompting helps GPT3-175B significantly, even comparable to human responses from Figure 5. Blender also gained much improvement from pre-trained after fine-tuning on *Reflect*, however we find no significant effects on appending inference questions.

Dialogue History	A: One of my patients keeps forgetting to take his medicine on time. B: That's not good, how can you help someone like that? Why did you have to quit? A: I've been calling him everyday to make sure he takes his medicine, I have noticed his health is improving too.	
Orig Response	B: You're a great nurse!	
Reflect	Q: How would you describe A? Ans: A is kind, patient and persistent.	B: You are really determined and wonderful with your patients. Not many doctors would take the time every day to do something like that.
	Q: What might have happened BEFORE? Ans: A noticed a patient's health wasn't improving with their medicine.	B: What made you notice he wasn't taking it? Did they need a refill by your notes but still have a lot left? Or did you notice they weren't better at their last appointment?
	Q: What do you think might happen AFTER? Ans: A will contact a family member of the patient to help their patient remember to take meds..	B: How will you handle it going forward? Can you keep up with that, or will you see if someone else close to them can help?
	Q: What do you think A is feeling? Ans: Satisfied, successful and determined.	B: You must be feeling really confident and strong right now. You're a doctor who goes above and beyond for your patients. I've never met another who would do this.
	Q: What do you think B is feeling? Ans: Impressed and optimistic for A.	B: Wow, that's so amazing. I can't believe you've done at all that. Most doctors would just assume it's the patients responsibility. I am confident you'll get him better after this.
	GPT3-FS-Top5	B: That is amazing. You have to be a great nurse to help people like that. B: You are a true caregiver and I'm glad you're taking care of your patients so well. B: That's good, once he stops forgetting you can give him back his phone and he can do it on his own. B: Well, you are a good doctor. B: That is good.
GPT-FS-InfQ	Q: How would you describe A? Q: What might have happened BEFORE? Q: What do you think might happen AFTER? Q: What do you think A is feeling? Q: What do you think B is feeling?	B: You are amazing! I could never do that job. B: That's really good that you took the initiative to make sure he was taking his medicine B: You are really dedicated to your job. I hope you get a raise soon. B: It's great that you are helping him out and that his health is improving. B: I'm glad that the patient is improving, but it's a lot of work for you.

Table 4: **Randomly sampled dialogue example to qualitatively compare responses from different resources and models.** Compared to original and GPT3 top-5 responses which mostly express “*You are a great nurse/doctor*”, *Reflect* responses based on the inference QA pairs are much more specific and engaging. We find that GPT3-InfQ, though only guided by an inference question, generates less generic responses such as “... *I hope you get a raise soon.*”

munication. Since both BlenderBot and GPT3 are trained to generate responses directly without integrating common ground, a non-trivial challenge is how to adapt them to use inference-based common ground before RG. Here we present our two intuitive and simple approaches.

For BlenderBot-440M, we follow the common practice of fine-tuning models to adapt to a new task format. We split our Reflect data into 60/10/30 for train/valid/test and first fine-tune BlenderBot-440M (**Blender-FT**) on only the collected responses to show potential benefits of training from inference-guided human responses. Then we fine-tune BlenderBot but modify the training task from outputting responses from contexts to inference-guided RG. Inspired by modular generation in dialogue RG (Adolphs et al., 2021; Zhou et al., 2022; Shuster et al., 2022), our training task is: given dialogue context and one of the five inference dimension questions, generate the answer as well as the response collected in Reflect (**Blender-FT-InfQ**, indicating that the model is given the *Inference Question*). More details in Appendix C.

For GPT-175B, we follow the few-shot in-context learning approach with one small addition in input: we append the dialogue context with an inference question and ask the model to generate a response. Our pilot studies show that GPT3 tends to generate directly an answer to the question, not a next-turn response to the dialogue context, thus we format the question into a prompt for GPT3

and stress that the end goal is RG. Specifically, we append the text “*Think about this when responding:*” and then one of our inference questions after the dialogue context to prompt GPT3 to generate a response by reflecting on the questions (**GPT3-FS-InfQ**). Illustrative figures for prompting GPT3 are shown in Appendix A Figures 10.

To compare and analyze the effects of each inference dimension, we randomly sample one response for each of the five inference dimensions for GPT3-FS-InfQ and Blender-FT-InfQ and take their average. For GPT3-FS, Blender, and Blender-FT, we pick the *top 5* responses generated using their default decoding strategy (beam search for GPT3 and nucleus sampling for Blender) and aggregate their evaluation results. In total, we evaluate 250 responses from *each* model following the procedure in Section 4.1.

5.2 Experimental Results

Prompting GPT3 to “think” about common ground improves response quality by 30% Figure 7 presents results when comparing models that has no access to inference-guided *Reflect* data with those that do. We test the hypothesis that whether guiding RG models with inference questions about common ground is helpful for generating more human-like responses. We find that with inferences, GPT3-FS-InfQ outperforms GPT3-FS on *all* evaluation dimensions. Specifically, inference-guided GPT3 produces almost 25% more specific and 30%

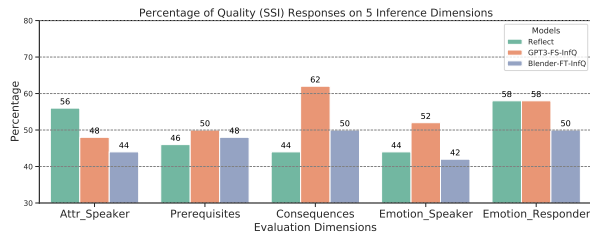


Figure 8: **Response evaluation separated by inference dimensions.** We find that GPT3-FS-InfQ generate better responses than humans on the potential consequences dimension while generates worse on attributes.

more quality responses. Moreover, 54% quality (sensible, specific, and interesting) responses already surpasses quality of human-written responses in original dialogues (49%), but still lags behind *Reflect* (58%) as shown in Figure 5.

Fine-tuning Blender on *Reflect* generates 26% more quality responses For BlenderBot-400M, we find that fine-tuning on inference-guided human responses from *Reflect* helps generate almost 50% more specific and 26% more quality responses. In contrast to GPT3, BlenderBot with inference-guided fine-tuning does not seem to improve much. We speculate that model size might play a role in how much model is influenced by CG inferences, leaving future work for more inference-customized fine-tuning on moderate-sized models.

5.3 Analysis

Which inference dimension helps models the most (and which the least)? Figure 8 shows the percentages of quality responses separated by the inference dimension we use to prompt humans and models. Interestingly, we find that on some dimensions, GPT3-FS-InfQ can produce significantly better responses than human responses from *Reflect*, especially event-based: “*What might have happened before*” and “*what might happen after?*” and emotion-based CG about the other speaker “*What is A (speaker1) feeling now?*”. However, on “*How would you describe A*”, human responses grounded on this question are much better. This dimension-specific analysis provides evidence that neural models’ capability to generate quality responses may depend on what *types of CG* we use to guide them.

Prompting GPT3-175B with complete human inferences To show how well GPT3 can make use of complete human-annotated common ground, we further append the inference answer after the question from *Reflect* data and prompt GPT3 to

generate a response given the fully materialized common ground. As expected, we observe further improvements in response quality especially in specificity (15% more) and general quality (16.7% more). This analysis shows promises to make reflect-style models produce better responses by providing quality inference answers for CG.

6 Related Work

We have presented discussion of previous work representing CG (Ghosal et al., 2022; Zhou et al., 2022) in Section 1 and relevant communication theory and psycholinguistic literature in Section 2. Here we provide additional discussions. Recent advances on neural RG models mainly focused on fine-tuning large pre-trained transformer models (Zhang et al., 2020; Roller et al., 2021; Thopilan et al., 2022) on huge number of dialogue data. However, few of the data provides explicit common grounding. Modular RG (Adolphs et al., 2021; Shuster et al., 2022) aims to generate relevant knowledge first by retrieving from the web and the generate knowledge-grounded responses. Compared to these work, we focus on inferences based on common sense instead of external knowledge. Another closely related work by Cho and May (2020) examined incorporating dialogue data with techniques from improvisational theater to teach models to implicitly build common ground.

7 Conclusion

We introduce *Reflect*, a dataset with diverse inference-grounded responses inspired by CG and communication theories. We carefully design our two-stage collection process and apply quality control. Then we demonstrate limitations of existing dialogue data and models trained on it. Finally, we present promising signs that guiding models with CG results in more engaging conversations. We hope to encourage more work on improving RG quality by looking at how humans use CG and adapt the communication process to machine learning models. Future directions include providing a *ranking* of inference dimensions depending on dialogue context and train models to generate responses following the most suitable dimension. *Reflect* also enables potential automated metrics to evaluate response since more responses per dialogue might help gauge the plausible response space given a context.

Acknowledgments

We thank anonymous reviewers for providing insightful feedback along with Brendan Kennedy, Peifeng Wang, and members from INK and JAUNTS lab. This research is supported in part by the DARPA MCS program under Contract No. N660011924033, the Defense Advanced Research Projects Agency with award NSF IIS 2048211, NSF SMA 182926, and support from Google.

Ethics and Broader Impact

We collect a new dialogue dataset in English, which benefits English speakers more. We use Amazon Mechanical Turk to recruit crowdsourcing workers and we pay workers over \$15/hour on average, well above the highest state minimum wage and engage in constructive discussions if they have concerns about the process. We also give each annotation instance enough time so that we do not pressure annotators. In our quality assurance process for this dataset, we also examine potential harmful biases and aggressive languages in responses and remove them in the final dataset. We also acknowledge that the generated responses from our experimented models might contain biases.

8 Limitations

Our first limitation in modeling CG is that we are using inferences from one speaker to approximate CG during the communication process. To truly represent CG, we need to recollect dialogues and as participants continue the conversations, we should ask both of them the same inference questions and perform post-hoc analysis on the answers to the questions.

Our second limitation is the lack of explicitly modeling *communicative intents*. In future work, we plan to heuristically link each inference dimension to a general communication goal. For example, making inferences about “speaker emotion states” is helpful to build emotional connections with the other speaker.

References

Leonard Adolphs, Kurt Shuster, Jack Urbanek, Arthur Szlam, and Jason Weston. 2021. Reason first, then respond: Modular generation for knowledge-infused dialogue. *arXiv preprint arXiv:2111.05204*.

Jens Allwood. 1976. *Linguistic communication as ac-*

tion and cooperation. University of Göteborg. Department of Linguistics.

David Bohm, Peter M Senge, and Lee Nichol. 2004. *On dialogue*. Routledge.

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Jerome S Bruner. 1961. The act of discovery. *Harvard educational review*.

Hyundong Cho and Jonathan May. 2020. **Grounding conversations with improvised dialogues**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2398–2413, Online. Association for Computational Linguistics.

Herbert H Clark and Susan E Brennan. 1991. Grounding in communication.

Herbert H Clark and Edward F Schaefer. 1989. Contributing to discourse. *Cognitive science*, 13(2):259–294.

Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming Zhou. 2020. **MuTual: A dataset for multi-turn dialogue reasoning**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1406–1416, Online. Association for Computational Linguistics.

Jonathan St BT Evans. 2003. In two minds: dual-process accounts of reasoning. *Trends in cognitive sciences*, 7(10):454–459.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.

Deepanway Ghosal, Siqi Shen, Navonil Majumder, Rada Mihalcea, and Soujanya Poria. 2022. Cicero: A dataset for contextualized commonsense inference in dialogues. *arXiv preprint arXiv:2203.13926*.

Herbert P Grice. 1975. Logic and conversation. In *Speech acts*, pages 41–58. Brill.

Jürgen Habermas. 1985. *The theory of communicative action: Volume 1: Reason and the rationalization of society*, volume 1. Beacon press.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.

Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. Comet-atomic 2020: On symbolic and neural commonsense knowledge graphs. In *AAAI*.

- Daniel Kahneman. 2011. *Thinking, fast and slow*. Macmillan.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. [DailyDialog: A manually labelled multi-turn dialogue dataset](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Borde, D. Parikh, and J. Weston. 2017. Parlai: A dialog research software platform. *arXiv preprint arXiv:1705.06476*.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. [Recipes for building an open-domain chatbot](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 300–325, Online. Association for Computational Linguistics.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019a. [ATOMIC: an atlas of machine commonsense for if-then reasoning](#). In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3027–3035. AAAI Press.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019b. [Social IQa: Commonsense reasoning about social interactions](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.
- Iulian Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31.
- Kurt Shuster, Mojtaba Komeili, Leonard Adolphs, Stephen Roller, Arthur Szlam, and Jason Weston. 2022. Language models that seek for knowledge: Modular search & generation for dialogue and prompt completion. *arXiv preprint arXiv:2203.13224*.
- Vered Shwartz, Peter West, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Unsupervised commonsense question answering with self-talk. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4615–4629.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*.
- Robert C Stalnaker. 1978. Assertion. In *Pragmatics*, pages 315–332. Brill.
- Keith E Stanovich and Richard F West. 2000. Individual differences in reasoning: Implications for the rationality debate? *Behavioral and brain sciences*, 23(5):645–665.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.
- Gordon Wells. 2000. Dialogic inquiry in education. *Vygotskian perspectives on literacy research*, pages 51–85.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. [DIALOGPT : Large-scale generative pre-training for conversational response generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.
- Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 654–664.
- Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayatnia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang Liu, and Dilek Hakkani-Tur. 2021. [Commonsense-focused dialogues for response generation: An empirical study](#). In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 121–132, Singapore and Online. Association for Computational Linguistics.

Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayatnia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang Liu, and Dilek Hakkani-Tur. 2022. Think before you speak: Using self-talk to generate implicit common-sense knowledge for response generation.

A Data Collection Details

We engage in active discussions with them in the TurkerNation⁴ Slack channel and provide detailed feedback after multiple rounds of pilot study to ensure the data quality.

A.1 Inference Collection

Here we present more detailed feedback for AMT workers on Stage 1. inference collection: First, we stress that the goal of these answers is to help with generating a response to continue the conversation instead of any inferences that might not be useful for directly generating engaging responses, such as “*spaghetti is a type of food*” for the example in Figure 1. Secondly, the answers should not be a direct copy-paste of some parts in the dialogue context as those would be trivial to collect, violate the least collaborative principle and the maxim of quantity (Grice, 1975), and should not be worth making inferences over. Finally, we remind them that the inferences written should be considered as “*common sense*” so that the approximated CG is more likely to become shared knowledge and beliefs among the dialogue participants. Collection UI and provided examples for turkers are shown in Figures 12 and 13.

A.2 Response Collection

We specifically stress on several points to workers: 1) to collect more engaging and interesting responses, response should not directly *paraphrase* the inference such as “*I think you are feeling relieved*” from inference QA pair “*What is speaker feeling now? Speaker is feeling relieved*”; 2) the response should be both *coherent* to the dialogue context as what would be naturally uttered by the responder and *based on* the reactions to lead the conversation in an interesting direction; 3) Ultimately, we want responses that lead the conversations that are more enjoyable and engaging. Collection UI and provided examples for turkers are shown in Figures 14 and 15.

B Human Evaluation Details

Specifically, a sensible response is one that is reasonable in context. A specific response is one that relates closely to the given dialogue context, instead of a generic one that can be applied in dozens of different contexts. An interesting response can

⁴<https://www.reddit.com/r/TurkerNation/>

“*catch someone’s attention or arouse their curiosity, or if it is unexpected, witty, or insightful.*” (Thoppilan et al., 2022). For more detailed instructions, please refer to Thoppilan et al. (2022). Evaluation UI and provided examples for turkers are shown in Figures 16 and 17.

C Model Implementation Details

We use two base models in our paper: BlenderBot-440M and GPT3-175B. For BlenderBot, we use the ParlAI (Miller et al., 2017) package for pre-trained modeling and fine-tuning. The format for fine-tuning BlenderBot on inference questions is: input sequence is “<speaker1> ... <speaker2>... <speaker1>... <infq> What might have happened before?” and output sequence is “<infa>... <speaker2> ...”, where we use “<infq>”, “<infa>” to indicate the start of an inference question and answer, respectively. We fine-tune BlenderBot-440M for 3 epochs with batch size 16 and set the learning rate to be 1e-06. We perform gradient accumulation for 8 steps and gradient clipping with a max norm of 1.0 and optimize using the Adam optimizer. For decoding, we use top-p nucleus sampling (Holtzman et al., 2019) with temperature T (p = 0.9 and T = 0.7), and a maximum decoding length of 300 tokens. BlenderBot-440M models are mostly trained on 4 Quadro RTX 8000 GPUs and take around 9 hours.

We use OpenAI-API⁵ to access GPT3-DaVinci (175B) and include prompting formats for GPT3-FS and GPT3-FS-InfQ in Figures 9 and 10, respectively.

D Additional Experimental Results

D.1 Inference-Separated Fine-Grained Evaluation Results

Inference dimension-separated full results are shown in Figure 11.

⁵<https://beta.openai.com/playground>


```

<speaker1>: I found a new friend in my neighborhood.
<speaker2>: you are new in that neighborhood, so how did you manage to have a new friend. <speaker1>: I went out
running one morning and I met a lady.
<speaker2>: Of course you did Jan. You make new friends no matter where you go.
###
<speaker1>: Have I told you about the crazy daydreams I've had lately?
<speaker2>: Not that I remember. I hardly ever dream when I take a nap.
<speaker1>: Same here, but I just had one during my nap about chasing rainbows!
<speaker2>: Maybe you were looking for a pot of gold! I hope you sleep better tonight and that you catch a rainbow
in your dreams.
###
<speaker1>: I can't wait until I have enough money to buy a new car.
<speaker2>: Why do you need a car?
<speaker1>: I was bumming rides to work with my neighbor but I just started at a new place downtown. Now I have to
take the bus every day.
<speaker2>: I will ask around and see if there's any easy way to get money for your car.
###
<speaker1>: Where are my Commandments?
<speaker2>: Last I heard was that Remi had them.
<speaker1>: I need to find Remi so I can get them.
<speaker2>:

```

Figure 9: GPT3-Few Shot Prompting Format (no inference).

```

<speaker1>: I found a new friend in my neighborhood.
<speaker2>: you are new in that neighborhood, so how did you manage to have a new friend.
<speaker1>: I went out running one morning and I met a lady.
Think about this when responding: How would you describe <speaker1>?
<speaker2>: Of course you did Jan. You make new friends no matter where you go.
###
...
###
<speaker1>: Where are my Commandments?
<speaker2>: Last I heard was that Remi had them.
<speaker1>: I need to find Remi so I can get them.
Think about this when responding: How would you describe <speaker1>?
<speaker2>:

```

Figure 10: GPT3-Few Shot-Inference Question Prompting Format.

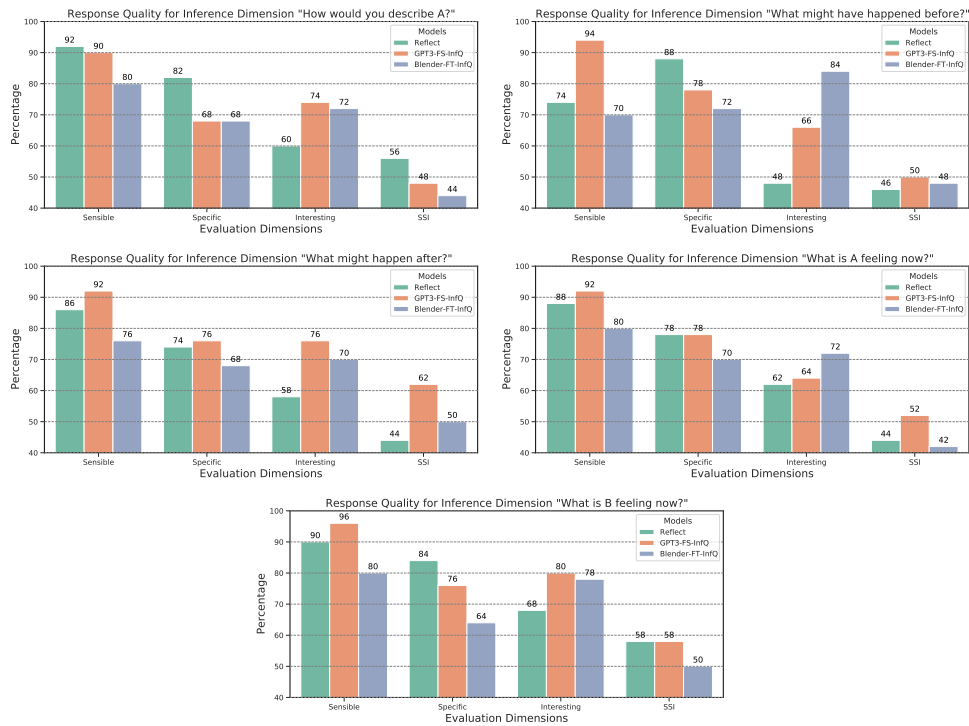


Figure 11: Response evaluation separated by inference dimensions. We find that GPT3-FS-InfQ generate better responses than humans on the potential consequences dimension while generates worse on attributes.

Welcome to our task! Please read our instructions very carefully. We're excited to work with you on this and on the thousands of HITs we have coming. Your annotation quality will be evaluated

Read Instructions and Examples

Here is an incomplete dialogue between two friends where $\$(speaker)$ talks about an event he/she experienced to "Friend".

You are asked to *imagine that you are "Friend"*. Answer the following questions about how you would have reacted. These reactions include how you think about $\$(speaker)$, any assumptions on what might have happened BEFORE for this conversation to take place, predictions about what might happen after, and what you and the Friend might be feeling AFTER.

Your answers should be a few words or a short sentence that directly answer the questions.

Dialogue:
 $\$(history)$

Question 1: How would you describe $\$(speaker)$?
 (e.g., $\$(speaker)$ is caring, a responsible employee, persistent, etc.)

Question1

Question 2: What might have happened BEFORE? Your assumptions should be implied by the dialogue. Please do not state what's obvious in the dialogue context. For examples, $\$(speaker)$ does ... BEFORE, $\$(speaker)$ thinks X is ... etc.

Question2

Question 3: What do you think might happen AFTER?

Question3

Question 4: What do you think $\$(speaker)$ is feeling? (e.g., nervous, annoyed, etc.)

Question4

Question 5: What do you think Friend (you) is feeling? (E.g., happy, relieved, angry, ...)

Question5

Figure 12: Inference collection UI.

Good examples	Bad examples
<p>Dialogue:</p> <ul style="list-style-type: none"> • Jordan: How'd you do on the exam? • Friend: I failed it, so did everyone else I asked • Jordan: Oh, I didn't think it was that hard • Friend: What did you get on it? • Jordan: I got an A and I think if my friends studied harder they could have done the same <h3>Friend Reaction 1:</h3> <p>How would you describe Jordan?</p> <p>Arrogant/Boastful/Diligent/Condescending/etc. (any one of these).</p> <h3>Friend Reaction 2:</h3> <p>What might have happened BEFORE? Your assumptions should be implied by the dialogue. Please do not state what's obvious in the dialogue context. For examples, $\{speaker\}$ does ... BEFORE, $\{speaker\}$ thinks X is ..., etc.</p> <p>Jordan needs to have studied hard.</p> <h3>Friend Reaction 3:</h3> <p>What do you think might happen AFTER?</p> <p>Jordan's friends maybe mad at him and not want to speak to him for a while.</p> <h3>Friend Reaction 4:</h3> <p>What do you think Jordan is feeling? (e.g., nervous, annoyed, etc.)</p> <p>Proud</p> <h3>Friend Reaction 5:</h3> <p>What do you think Friend (you) is feeling? (E.g., happy, relieved, angry, ...)</p> <p>Sad/Angry/Motivated/etc. (any one of these)</p>	<p>Dialogue:</p> <ul style="list-style-type: none"> • Jordan: How'd you do on the exam? • Friend: I failed it, so did everyone else I asked • Jordan: Oh, I didn't think it was that hard • Friend: What did you get on it? • Jordan: I got an A and I think if my friends studied harder they could have done the same <h3>Friend Reaction 1:</h3> <p>How would you describe Jordan?</p> <p>A caring person.</p> <p>Explanation: This is wrong because there is no place indicating that Jordan is caring from the context.</p> <h3>Friend Reaction 2:</h3> <p>What might have happened BEFORE? Your assumptions should be implied by the dialogue. Please do not state what's obvious in the dialogue context. For examples, $\{speaker\}$ does ... BEFORE, $\{speaker\}$ thinks X is ..., etc.</p> <p>Jordan aced the exam.</p> <p>Explanation: This is wrong because it is not an assumption implied but something explicitly stated as in "I got an A".</p> <h3>Friend Reaction 3:</h3> <p>What do you think might happen AFTER?</p> <p>Jordan studied hard for the exam.</p> <p>Explanation: It is something that more likely happened BEFORE not AFTER.</p> <h3>Friend Reaction 4:</h3> <p>What do you think Jordan is feeling? (e.g., nervous, annoyed, etc.)</p> <p>Sad</p> <p>Explanation: No indication in dialogue</p> <h3>Friend Reaction 5:</h3> <p>What do you think Friend (you) is feeling? (E.g., happy, relieved, angry, ...)</p> <p>The person is happy.</p> <p>Explanation: Not sure who "The person" is and no indication of feeling happy.</p>

Figure 13: Inference collection examples for turkers.

Welcome to our task! Please read our instructions very carefully. We're excited to work with you on this and on the thousands of HITs we have coming. Your annotation quality will be evaluated.

Instructions

1. Here you will see an incomplete dialogue between $\{speaker\}$ and his/her Friend and it's the Friend's turn to respond.
2. You are also given several potential reactions from Friend upon hearing $\{speaker\}$'s words.
3. Your task is to imagine that you are the Friend responding to this conversation BASED ON each of the reactions.
4. Please try NOT to directly copy/verbalize the reaction as a response such as "X is seen as kind" → "You are so kind!". Try to be creative in your response such as "X is seen as kind" → "Aren't you the best roommate!".
5. Your first priority should be writing a response that is COHERENT to the dialogue context and try to include the reaction you are given.

The conversations will be checked for quality.

View examples

Dialogue:

$\{dialogue\}$

Friend's Reaction 1:

How would you describe $\{speaker\}$?

$\{reaction1\}$

Do you think this is a plausible reaction by Friend?

Yes No

If you selected Yes, then your response should be coherent with the dialogue and be based on the reaction.

If you selected No, then please first write a valid reaction answering the question.

Write your response that continues the conversation based on Reaction 1:

Response

Friend's Reaction 2:

What might have happened BEFORE the conversation?

$\{reaction2\}$

Thank you for your participation! Please provide feedback below on how we can improve the instructions, annotation interface, or anything else that you found confusing while completing the task.

Feedback

Figure 14: Response collection UI.

Good examples

Dialogue:

- $\$(\text{speaker})$: This place is awesome.
- Friend: What place?
- $\$(\text{speaker})$: The amusement park, it rocks
- Friend: Yeah, they are fun.
- $\$(\text{speaker})$: I saw season ticket booth too..

Friend's Reaction 1:

How would you describe $\$(\text{speaker})$?

I would say $\$(\text{speaker})$ is excited and an outgoing person.

Response: "Are you going to get season tickets? I'd love to visit here more often with someone who's always so energetic like you!"

Friend's Reaction 2:

What might have happened BEFORE the conversation?

$\$(\text{speaker})$ really enjoyed the rides in this amusement park.

Response: "Seems like you plan to frequent here. You must've had a blast with these rides."

Friend's Reaction 3:

What do you think might happen AFTER this conversation?

Response: "Let me know if you purchase those season tickets and I'll go with you some time."

Friend's Reaction 4:

What do you think $\$(\text{speaker})$ is feeling?

$\$(\text{speaker})$ feels excited and satisfied with the amusement park.

Response: "You enjoyed the rides so much that you plan to buy the season tickets now? Tell me all about it."

Friend's Reaction 5:

What do you think Friend (you) is feeling? (E.g., happy, relieved, angry, etc.)

Friend feels happy to have such an outgoing friend.

Response: "You want to buy season tickets? I'm in! I'm really glad that you always seem to have more fun than everyone else and it's affecting me!"

Bad examples

Dialogue:

- $\$(\text{speaker})$: This place is awesome.
- Friend: What place?
- $\$(\text{speaker})$: The amusement park, it rocks
- Friend: Yeah, they are fun.
- $\$(\text{speaker})$: I saw season ticket booth too..

Friend's Reaction 1:

How would you describe $\$(\text{speaker})$?

I would say $\$(\text{speaker})$ is an excited and outgoing person.

Response: " $\$(\text{speaker})$, you are the coolest gal God has ever sent me. I love Amusement parks because of you..."

Explanation: It's NOT a coherent response to the dialogue context, i.e., it is disconnected from "I saw season ticket booth too..".

Friend's Reaction 2:

What might have happened BEFORE the conversation?

$\$(\text{speaker})$ really enjoyed the rides in this amusement park.

Response: "You must really enjoy the rides in this park!"

Explanation: It's a direct copy-pasting of the reaction.

Friend's Reaction 3:

What do you think might happen AFTER this conversation?

$\$(\text{speaker})$ may purchase season tickets.

Response: "Oh where is it?"

Explanation: This response is NOT based on the reaction given.

Friend's Reaction 4:

What do you think $\$(\text{speaker})$ is feeling?

$\$(\text{speaker})$ feels excited and satisfied with the amusement park.

Response: "You must be happy and filled with anticipation."

Explanation: This response is NOT coherent to the dialogue context.

Friend's Reaction 5:

What do you think Friend (you) is feeling? (E.g., happy, relieved, angry, etc.)

The other(s) are either bewildered or calm, more subdued than $\$(\text{speaker})$.

Response: "Heh, you guys need to match $\$(\text{speaker})$'s energy more, don't you all think? I'm kidding, I love you guys"

Explanation: There is no indication of other people involved in the conversation and is not a coherent continuation from "I saw season ticket booth too..".

Dialogue: \${context}	Response 1: \${response_1} Does this response make sense? <input type="radio"/> Yes <input type="radio"/> No Is the response specific? <input type="radio"/> Yes <input type="radio"/> No Is the response interesting? <input type="radio"/> Yes <input type="radio"/> No
	Response 2: \${response_2} Does this response make sense? <input type="radio"/> Yes <input type="radio"/> No Is the response specific? <input type="radio"/> Yes <input type="radio"/> No Is the response interesting? <input type="radio"/> Yes <input type="radio"/> No
	Response 3: \${response_3} Does this response make sense?

Thank you for your participation! Please provide feedback below on how we can improve the instructions, annotation interface, or anything else that you found confusing while completing the task.

Feedback

Figure 16: SSI evaluation UI.

Welcome to our task! Please read our instructions very carefully.

Instructions

1. Here you will see an incomplete dialogue between two people.
2. You are given several responses from different systems.
3. Your task is to evaluate whether they are 1) sensible; 2) specific; and 3) interesting given the dialogue context.
4. Please read carefully of the "Instructions for each criterion" below to understand how to rate each response.
5. Treat each of the three criteria as separate measures, i.e. a response that is NOT sensible could possibly be specific or interesting and vice versa.

The evaluations will be checked for quality and we will give bonuses :)

Instructions for each criterion:

Does the response make sense?

1. Use your common sense here. Is the response completely reasonable in context?
2. If anything seems off—confusing, illogical, out of context, or factually wrong—then rate it as Does not make sense.
3. If in doubt, choose Does not make sense.

Is the response specific?

1. For example:
2. – if A says "I love tennis" and B responds "That's nice", then mark it as Not specific. That reply could be used in dozens of different contexts.
3. – but if B responds "Me too, I can't get enough of Roger Federer!" then mark it as Specific, since it relates closely to what you're talking about.
4. If you're in doubt, or if the reply seems at all generic, rate it as Not specific

Is the response interesting?

1. Choose Interesting if the response would likely catch someone's attention or arouse curiosity; also use that rating for anything insightful, unexpected, or witty.
2. If the response is monotonous and predictable, or if you're unsure, then pick Not interesting.

Figure 17: SSI evaluation instructions.