

SOCIAL SCAFFOLDS: A Generalization Framework for Social Understanding Tasks

Ritam Dutt, Carolyn Penstein Rosé, Maarten Sap

Language Technologies Institute, Carnegie Mellon University
{rdutt, cprose, msap2}@cs.cmu.edu

Abstract

Effective human communication in social settings is contingent on recognizing the subtle signals encoded in conversational exchange. However, inferring such social signals is challenging for most dialogue systems, especially when faced with a new task or setting. We introduce SOCIAL SCAFFOLDS, a rationale-generation framework for generalization in social understanding tasks. Our framework uses LLMs to generate three types of social signals or rationales that reflect the perspectives of the speaker, listener, and the general worldview. We conduct a comprehensive set of experiments spanning 150 cross-task scenarios wherein we first pre-train a model on a given source task (say detecting persuasion strategies), and subsequently deploy it for a target task (say identifying implicit hate speech). Our results show that providing language models with these rationales facilitates conversational understanding in both instruction-tuned and in-context learning settings; we find significant gains when we incorporate the social rationales alongside the utterance text as part of the input. Particularly, rationales modeling the speaker’s intentions yield the largest generalization gains (34%) across tasks. Our analysis also reveals that the generated rationales share low similarity with each other and the corresponding utterance, thereby capturing distinct concepts. They are also designed to be task-agnostic such that the rationale category with greatest impact depends on the task. Our framework shows the promise of pragmatics-oriented data augmentation for social understanding and generalization.

1 Introduction

Computational modeling of human behavior in social interactions is challenging because communication often employs indirect language, i.e. language whose meaning goes beyond the surface words of the text (Yerukola et al., 2024; Yusupujang and Ginzburg, 2023; Markowska et al., 2023;

Dutt et al., 2024). For example, Figure 1 illustrates that one needs to detect the underlying sarcastic intentions behind the message to infer the veiled implications of hate towards immigrants. Understanding the hidden meaning behind a message or conversational exchange is crucial for several tasks, such as automated content moderation (Calabrese et al., 2024; Horta Ribeiro et al., 2023), intent resolution (Yerukola et al., 2024; Joshi et al., 2021) and aiding LLM-based agents and tools (Kim et al., 2024; Qian et al., 2024).

This study investigates the extent to which language models (or broadly AI systems) can understand social inferences behind messages and how these inferences can serve as additional sources of information to facilitate generalization across different dialogue understanding tasks. While computational frameworks grounded in sociolinguistic theories such as the politeness framework of Brown et al. (1987), the cooperative principles/maxims of Grice (Bernsen et al., 1996), and the appraisal theory of Martin and White (Martin and White, 2003) have been proposed to understand the implicit social inferences, these frameworks cannot be applied readily to new tasks since their instantiation is contingent on the given task setting. For example, in politeness theory, what constitutes a positive face or a negative face (Brown et al., 1987) depends on the power dynamics and social distance of the participants and the given sociocultural setting.

We introduce SOCIAL SCAFFOLDS, a generalizable framework which automatically extracts these implicit social signals or inferences from the conversation which we henceforth refer to as “rationales.” Motivated by different points of view in narrative modeling (Eisenberg and Finlayson, 2016; Hamilton, 2024), we explore rationales that reflect (i) the speaker’s intentions and beliefs (Dutt et al., 2024; Zhou et al., 2023), (ii) the effect of the utterance on the listener (Yusupujang and Ginzburg, 2023), and (iii) the common world-view that par-

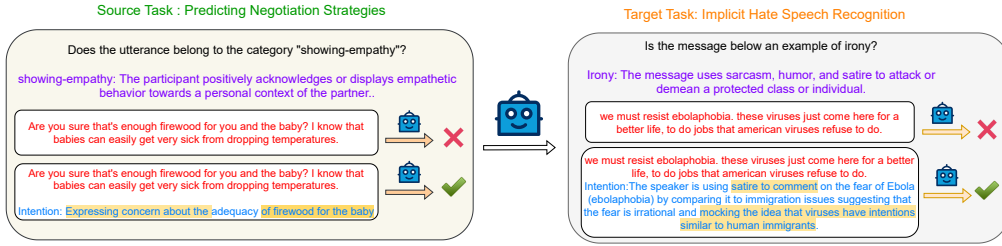


Figure 1: We illustrate the phenomena of indirect or subtle language usage in two scenarios; the scenario on the left corresponding to predicting negotiation strategies, whereas the scenario on the right corresponds to identifying different categories of hate. For both cases, we see how the model fails to associate the input message (in red) with the description of the label (in purple) since it is unable to capture the hidden cues in the message. Incorporating rationales, as additional inputs, can guide model prediction for both in-domain and cross-task settings.

ticipants presupposes to be true for the utterance to be credible (Mulcahy and Gouldthorp, 2016). These capture speaker-centric, listener-centric, and shared-centric perspectives and corresponds to first-person, second-person, and third-person points-of-view respectively in narrative modeling.

To showcase the utility of our framework, we generate $\approx 135K$ rationales, using GPT-4o and GPT-3.5-turbo as our backbone LLMs for six social dialogue datasets. We compare and contrast the impact of rationales for instruct-tuning and in-context learning setups, perform a thorough quantitative analysis of factors that affect generalizability, and characterize how similar different categories of rationales are to each other and to those generated by different LLMs. We observe more pronounced performance gains on datasets with higher skew in label distributions and for the infrequent label categories, highlighting the efficacy for more complex tasks. Our results also show significant associations between the choice of rationale and task performance showcasing that no single category of rationale acts as a silver bullet across all tasks.

We observe significant gains from incorporating rationales in a cross-task transfer setup. Simply put, we investigate whether a model fine-tuned or adapted for a given source task can generalize to a different target task. Figure 1 highlights that a model trained to detect negotiation strategies can also understand the different categories of hate speech when the intentions of the speaker are provided, in addition to the utterance text, as augmentations to the model during inference. Including the rationales corresponding to the speaker’s intentions, hearer’s reactions, and the presuppositions improve performance over the baseline significantly by 33.3%, 13%, and 13.3% respectively in the cross-task transfer scenario.

Our framework shows the promise of pragmatics-oriented data augmentation for social understanding and generalization. We make our dataset and code public for the research community.

2 Related Work

We contextualize our work in the broader literature on generalization in dialogue tasks as well as on rationales in language tasks.

2.1 Generalization in Dialogue

Generalization in dialogue is challenging because interactions are typically structured towards accomplishing a task rather than simply conveying information, involve multiple points of control, and rely heavily on implicit context (Dutt et al., 2024).

Mehri (2022) outlines different types of generalization imperative for dialogue. These include (i) new inputs arising from covariate shift or stylistic variation (Khosla and Gangadharaiiah, 2022), (ii) new problems in dialogue modeling such as evaluation and response generation (Peng et al., 2020), (iii) new outputs and schemas corresponding to out-of-domain shift (Larson et al., 2019) and (iv) new tasks such as controlled generation or fact verification (Gupta et al., 2022).

In this work, we focus on generalization across different dialogue tasks and investigate how rationales can act as pivots for the same. Prior work on few-shot generalization in dialogue has benefited from large-scale multitask pre-training (Wu et al., 2020; Peng et al., 2021; Hosseini-Asl et al., 2020) or instruction tuning (Gupta et al., 2022; Wang et al., 2025; Sanh et al.; Wang et al., 2022). We propose an efficient solution that leverages the underlying social signals, i.e. factors that remain common across dialogues thereby unifying different tasks, without the need to pre-train across multiple tasks.

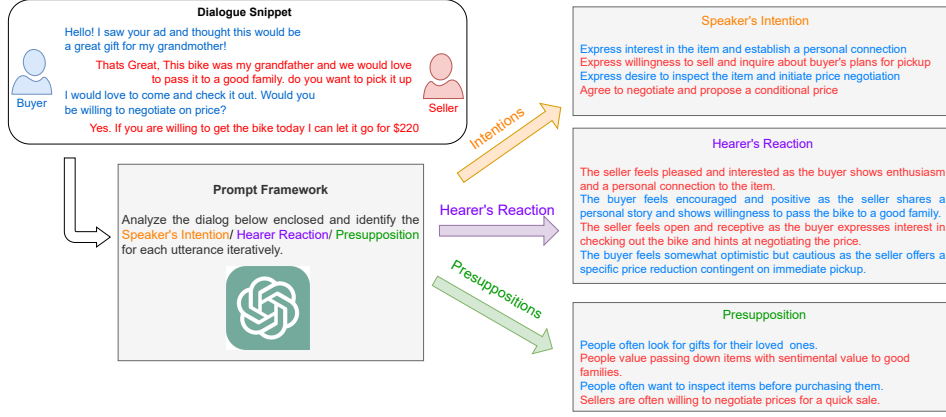


Figure 2: An overview our rationale generation framework SOCIAL SCAFFOLDS. We present a dialogue snippet between a buyer and a seller, shown in blue and red. We prompt an LLM with the dialogue snippet to generate the speaker’s intentions, the hearer’s reaction, and the presuppositions in orange, purple, and green, respectively.

2.2 Rationales in NLP

In NLP, “rationales”¹ has long been used to refer to *textual explanations*, either generated by machines or humans (Camburu et al., 2018). Rationales serve several purposes such as facilitating commonsense and social reasoning (Zelikman et al., 2022; Majumder et al., 2022), explaining the predictions of neural models (Wiegreffe et al., 2021; Jayaram and Allaway, 2021; Zaidan et al., 2007), and aiding humans in their tasks (Das and Chernova, 2020; Joshi et al., 2023; Zhang et al., 2023).

Recent research has demonstrated the efficacy of LLM in generating step-by-step explanations or rationales (Gurrapu et al., 2023) that can be utilized to improve downstream task performance (Rao et al., 2023; Wei et al., 2022; Zelikman et al., 2022). Rationales have also contributed to the OOD generalization (Majumder et al., 2022; Xiong et al., 2023; Joshi et al., 2022). Building upon this foundation, we frame rationales as the elicited verbalization of the underlying social signals that helps overcome some limitations of static text like the omission of communicative intent (Sap et al., 2022).

Our work improves upon that of Dutt et al. (2024), which investigates the domain generalization capabilities of rationales for dialogue understanding tasks in two ways. Firstly, we investigate the efficacy of rationales arising from multiple perspectives, i.e., the intentions of the speaker, the reaction of the listener, and the presuppositions involved in making the utterance, whereas prior work has emphasized mostly on the speaker’s intentions.

Additionally, we investigate the generalization capabilities of rationales across multiple dialogue tasks and not simply across different domains for the same task.

3 Modeling Framework

We present SOCIAL SCAFFOLDS, a framework that automatically generates rationales to capture the implicit information behind a message.

3.1 Rationale Types

This study explores three distinct but complementary perspectives or point-of-views to generate the rationales. Motivated by prior work on narrative modeling (Mulcahy and Gouldthorp, 2016), we present a one-to-one correspondence of the rationale category with the narrative point of view.

Intentions: Intentions refer to the hidden beliefs and desires of the speaker and correspond to the *first-person point-of-view*. These capture the implied meaning behind the speaker’s utterance and signal the outcome the speaker is interested in (Dutt et al., 2024; Yusupujang and Ginzburg, 2023).

Hearer Reaction: Rationales corresponding to the hearer’s reaction (Zhou et al., 2023; Sap et al., 2020) help capture the effect of the utterance on the listener(s). It provides insight into the listener’s emotions or belief states, akin to second-order thinking, and thus corresponds to the *second-person point-of-view*.

Presuppositions: Presuppositions refer to general facts or truths about the world that both parties must believe for the utterance to be credible. These presuppositions not only encapsulate common sense reasoning or social and communal norms often ob-

¹While rationales can also refer to a subset of input tokens or words that contribute to a classification decision (Bao et al., 2018), we use it in the broader sense of textual explanations.

served in practice (Perez Gomez, 2021; Kim et al., 2022), but also provides a de-contextualized or impersonal perspective of the scenario and thus serves as a *third-person point-of-view* (Mulcahy and Gouldthorp, 2016).

3.2 Rationale Generation Framework

We describe our prompting framework to automatically generate the different types of rationale. We provide an overview of our framework, SOCIAL SCAFFOLDS in Figure 2, with a sample dialogue snippet on the left and the corresponding intentions, hearer reactions, and presuppositions on the right.

SOCIAL SCAFFOLDS takes as input a multiparty dialog and generates rationales using a Large Language Model (such as GPT-4o) on an utterance-by-utterance basis. We employ a structured prompting framework to ensure that the generated rationale aligns with its corresponding utterance. We address erroneous cases by prompting the framework to regenerate the rationales iteratively. Additional details appear in Appendix Section B .

We reuse the same prompting framework to generate each category of rationale separately to prevent any ordering effects. Additionally, we do not provide any few-shot instances to avoid biasing the generations with previously seen examples as in Dutt et al. (2024). Overall, our framework enables us to compare and contrast not only different categories of rationales with each other but also the same categories of rationales generated by different LLMs. We explore two LLMs i.e. GPT-4o and GPT-3.5-turbo as the backbone of our SOCIAL SCAFFOLDS to generate the rationales.

3.3 Assessment of Rationale Quality

Since our framework automatically generates rationales without any human supervision, we develop a rigorous annotation manual to assess the validity of those generations based on three criteria: soundness, informativeness, and relevance. Additional details of these criteria appear in Appendix C

We score each rationale based on soundness, informativeness, and relevance using a Likert scale of 1 to 3, with 1 being the lowest and 3 the highest. The evaluations were carried out by two annotators with a graduate level proficiency in English and at least five years of experience in computational linguistics and NLP. Due to the highly subjective nature of the task, we relied on these professional annotators as an alternative to crowd-sourcing or employing an automated annotation framework.

We compute the inter-rater reliability scores using the multi-item agreement measure of Lindell et al. (1999) and observe strong to moderate agreement on all three criteria: soundness (0.98), informativeness (0.76), and relevance (0.70). The mean scores of soundness, informativeness, and relevance are 2.95, 2.76, and 2.61 respectively, highlighting that the rationales are of sufficiently high-quality.

Our results in Appendix E highlight that the rationales of different categories differ substantially between themselves showcasing that each category captures distinct concepts. We also observe low similarity between the rationale and the corresponding utterance once again signifying that the rationale generated captures information distinct from what is present in the utterance text.

4 Methodology

We outline our methodology for investigating how rationales can facilitate generalization to different social dialogue understanding tasks. We describe here the datasets, tasks, and experimental details.

4.1 Tasks and Datasets

We explore many dialogue understanding tasks, each instantiated with a distinct dataset, such that each task operates over a distinct domain. Moreover, these datasets have unique labels or categories to prevent any overlap between them. Such a setting would enable us to inspect the capabilities of rationales in a cross-task setting, where a model is trained for one task and then evaluated on another.

We explore six different datasets i.e., (i) P4G (Wang et al., 2019b) to identify persuasive strategies in charitable donations, (ii) CaSiNo (Chawla et al., 2021) to detect negotiation tactics in a simulated camping environment and (iii) Res_CB (Dutt et al., 2021) to categorize strategies employed to resist persuasion in online bargaining, (iv) EMH (Sharma et al., 2020) to understand different dimensions of empathy, (v) PROP (Jo et al., 2020) to categorize different kinds of argumentation, and (vi) IMP_HATE (ElSherief et al., 2021) to classify different kinds of implicit hate speech.

We present a brief overview of the dataset statistics in Table 1 and their corresponding distribution of labels in Figure 7 of the Appendix A. We observe that the datasets exhibit distinct characteristics, such as long conversations for P4G and PROP, and a higher skew for CaSiNo and Res_CB.

| Dataset | Avg Words per Turn | Avg Turns per Dialog | # Dialogs | # Labels |
|-----------------------------------|-----------------------|----------------------|------------------|--------------|
| P4G (Wang et al., 2019a) | 10.75 / 13.76 / 11.53 | 18.74 / 15.45 / 17.9 | 4004 / 110 / 154 | 11 / 11 / 11 |
| CaSiNo (Chawla et al., 2021) | 21.53 / 20.29 / 26.50 | 5.42 / 4.88 / 5.02 | 4862 / 49 / 247 | 10 / 9 / 10 |
| Res_CB (Dutt et al., 2021) | 12.22 / 13.63 / 13.71 | 5.86 / 5.18 / 6.09 | 6348 / 160 / 160 | 8 / 8 / 8 |
| PROP (Jo et al., 2020) | 12.55 / 14.86 / 15.71 | 11.66 / 9.47 / 12.21 | 741 / 43 / 75 | 4 / 4 / 4 |
| EMH (Sharma et al., 2020) | 54.03 / 47.75 / 53.83 | 1 / 1 / 1 | 1823 / 104 / 112 | 3 / 3 / 3 |
| IMP_HATE (ElSherief et al., 2021) | 15.79 / 17.18 / 15.39 | 0 / 0 / 0 | 3182 / 156 / 153 | 6 / 6 / 6 |

Table 1: Overview of the dataset statistics across the train, validation, and test splits.

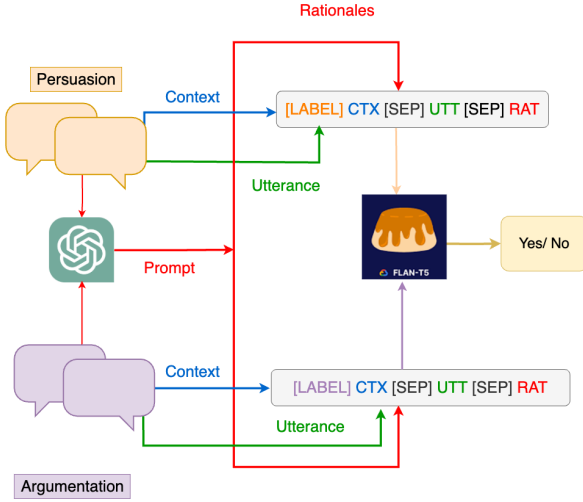


Figure 3: Overview of our instruction tune setting

4.2 Experimental Framework

We investigate the impact of rationales on downstream task performance in two experimental settings. The first is an instruction-tuned paradigm (Figure 3) where we fully fine-tune a pre-trained language model on a given source task (say persuasion) and then subsequently evaluate it on a new target task (say argumentation) in a 0-shot or few-shot setting. The second is an in-context learning setting, where we prompt an LLM with 0-shot or few-shot examples with the rationale as a control.

We frame each of the six multi-label, multi-class classification tasks as binary classification, where the label definition, utterance, dialog context, and rationale serve as input to the model. The model has to output whether the utterance conforms with the definition of the label via "Yes" or "No". We adopt the same approach for both instruction-tuned and in-context learning settings. This design takes into account that each task operates in their own label space without any overlap. Moreover, fine-tuning LMs with a single multiclass classification head is unlikely to generalize in 0-shot settings. Moreover, our design would allow for a fair comparison of the two paradigms. We show an example

of how these tasks have been set-up in Figure 1.

4.3 Models and Metrics

We use the base version of Flan-T5 (Chung et al., 2022) as our instruction-tuned model, while Gemma-2-9B-it (Team, 2024) and Llama-3-8B-it (AI@Meta, 2024) serve as in-context learning models. These models have been fine-tuned for instruction-following and thus serve as strong baselines for the respective experimental paradigms. We inspect the difference in performance from adding rationales as part of the input text (i.e., intentions, presuppositions, and hearer reaction) over only the utterance (which serves as the baseline).

To account for the skewed label distribution, we use macro-F1 score as the main evaluation metric for each of these six tasks. Following the recommendations in Dror et al. (2018), we use the non-parametric bootstrap test of Berg-Kirkpatrick et al. (2012) to measure the statistical significance between the baseline and the rationale-augmented model. We reject the null hypothesis that the baseline and rationale-augmented models have similar performance for cases with p -value ≤ 0.05 .

5 Analysis

5.1 Rationales on Task Performance

We evaluate the performance of our instruction-tuned model in an in-domain setting (model is evaluated on the same source task as it was trained on) and a cross-task setting (model is evaluated on a new target task). We repeat over three seeds to account for variations across runs.

In-domain Performance: We present the in-domain performance in Table 2 and observe modest gains in five of six tasks, with significant improvements for res_CB and IMP_HATE, and a significant drop for EMH. We also notice that the rationale corresponding to intentions, i.e., the speaker’s perspective, has the most consistent and prominent gains out of all the rationales. We observe similar findings for both GPT-3.5-turbo and GPT-4o.

| Generator | Rationale | P4G | CaSiNo | res_CB | PROP | EMH | IMP_HATE |
|---------------|-----------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| - | UTT | 69.70 +/- 2.42 | 71.22 +/- 1.70 | 66.77 +/- 1.02 | 82.38 +/- 1.21 | 90.91 +/- 0.13 | 62.68 +/- 0.79 |
| GPT-4o | INT | 69.36 +/- 1.45 | 72.35 +/- 0.50 | 70.91 +/- 0.71 | 84.66 +/- 1.07 | 89.35 +/- 1.35 | 67.91 +/- 1.49 |
| | HR | 70.54 +/- 1.70 | 71.71 +/- 0.84 | 68.80 +/- 0.97 | 82.88 +/- 1.69 | 90.26 +/- 0.32 | 65.08 +/- 0.34 |
| | PreSup | 68.12 +/- 2.30 | 71.81 +/- 1.39 | 69.69 +/- 1.51 | 80.11 +/- 2.86 | 89.37 +/- 0.16 | 62.88 +/- 2.55 |
| GPT-3.5-turbo | INT | 67.64 +/- 3.16 | 72.35 +/- 0.38 | 71.22 +/- 3.03 | 81.52 +/- 1.47 | 90.01 +/- 1.12 | 62.82 +/- 0.62 |
| | HR | 68.90 +/- 1.54 | 71.95 +/- 2.67 | 70.87 +/- 1.17 | 83.61 +/- 2.00 | 89.18 +/- 0.73 | 64.16 +/- 0.97 |
| | PreSup | 72.21 +/- 0.25 | 70.43 +/- 1.27 | 69.28 +/- 1.45 | 78.61 +/- 2.97 | 90.00 +/- 0.96 | 59.85 +/- 0.52 |

Table 2: Performance of FLAN-T5 model in an in-domain setting across six tasks. The baseline includes only the utterance (UTT), which we compare against the three kinds of rationales, i.e. intentions (INT), hearer-reactions (HR), and presuppositions (PreSup). We represent the mean and standard deviation across three runs.

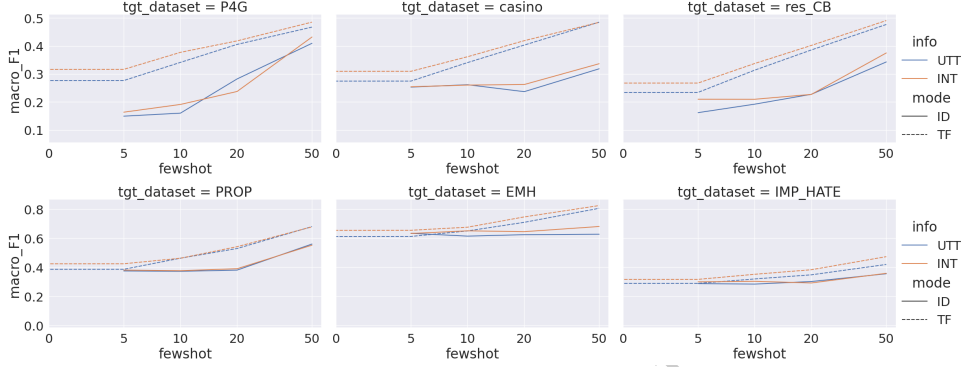


Figure 4: Impact of rationales on cross-task performance for instruction-tuned models across the six datasets for different fewshot settings.

Cross-Task Transfer Performance: We note the aggregate effect of adding the rationales in a cross-task environment resulting in 30 different combinations of source and target datasets in Figure 4. When comparing against the baseline case, i.e., the utterance, we see consistent and significant gains during transfer (in dotted lines) over the in-domain setting (in solid lines) for zero-shot and few-shot cases from adding the speakers’ intentions.

In-context Learning: A similar story emerges for the in-context learning (ICL) paradigm, where we observe that adding intentions to LLMs, i.e. Llama-3 and Gemma-2-7B, significantly improves the macro-F1 score (see Figure 5). We see mixed results for PreSup and HR, where the former and the latter are better at 0-shot and 5-shot settings, respectively. We also note that with only a mere 20 or 50 few-shot examples, the instruct-tuned models in a cross-task setting can surpass ICL.

We observe the impact of rationales to be highest for datasets that exhibits a high skew in their label distribution (such as P4G, res_CB, and IMP_HATE). Additionally, the label-wise macro-F1 scores in Figures 15 and 16 reveals that rationales have a higher impact on the infrequent label categories such as “foot-in-the-door” strategy for P4G, “Self-Assertion” and “Self-Pity” for res_CB,

and “threatening” for IMP_HATE. We posit that the rationales are more helpful for more complex dialogue understanding tasks in both in-domain and cross-task settings.

We note the fraction of cases where rationales significantly improve performance over the baseline for instruction-tuned models (both in-domain and cross-task settings) and in-context learning models in Figure 6. Across all settings, INT demonstrate consistent improvements and highlight that the speaker’s perspective plays the greatest role in facilitating dialogue understanding. However, despite the comparatively low performance in-domain, both HR and PreSup show pronounced gains in the cross-task transfer setting for instruction-tuned models, demonstrating their generalizability as pivots for task transfer.

5.2 Factors affecting Tasks Performance

Instance-wise Correlations We investigate several factors that could predict the performance of rationales on an instance-wise basis. The co-variables observed include (i) the length of the rationale, (ii) the length of the preceding dialogue history, (iii) the similarity between the rationale and the utterance, (iv) the similarity between the rationale and the label description being classified, (v) the read-

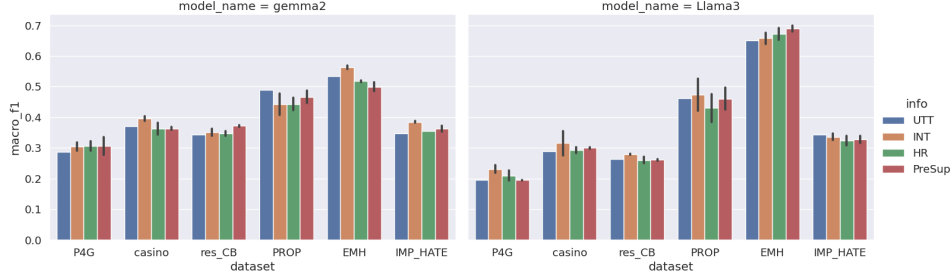


Figure 5: Zero-shot performance for in-context learning models.

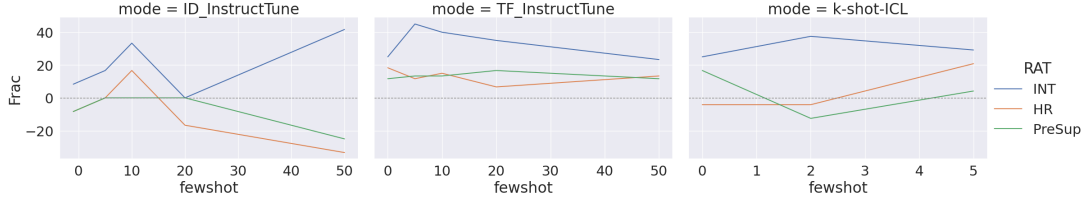


Figure 6: Fraction of cases where adding the rationale was significantly better (or worse) than the baseline in an indomain setting (left), a cross-task or transfer setting (middle), and in-context learning setup (right).

ability score measured using the Flesch’s readability ease (Farr et al., 1951; Kincaid, 1975), (vi) the valence, arousal, and dominance scores measured via the VAD NRC lexicon (Mohammad, 2018), and (vii) scores corresponding emotional intensity, emotional polarity and empathy (Wu et al., 2024).

We measure the point biserial correlation between these factors and instance-wise accuracy, i.e. whether the rationale could predict the label correctly or not. We observe very low (almost zero) correlation for each of the factors in Table 13 of the Appendix E. Our results highlight that the task accuracy is not dependent on these external data artifacts like rationale length or emotional intensity. Furthermore, as opposed to prior work on “free-text” rationales that were generated keeping in mind the label category such as E-SNLI (Wiegreffe et al., 2021), our rationales are task-agnostic based on the low similarity scores between the label description and the rationale.

Generalization Characteristics: We inspect the factors that characterize generalizability over the different experimental settings. We perform a multivariate ANOVA analysis with the relative performance difference (expressed as a percentage over the baseline) from including the rationale information as the dependent variable. The independent variables chosen were the rationale category, the LLM used to generate the rationales, the choice of source and target dataset², and the few-shot setting;

we also consider the pair-wise interaction effects of each of these variables. We note the F-statistic and their corresponding p-value for the indomain, cross-task and incontext-learning setting respectively in Tables 14, 15, and 16 in the Appendix E.

For the indomain setting, we observe that performance change hinges most on the fewshot setting followed by the choice of rationale and the dataset. We also see significant pair-wise effects for each of the categories except between the LLM and the choice of fewshot or between the LLM and dataset, highlighting that the rationales generated by the two LLM have similar effect.

In the cross-task setting, where we note that the choice of the target dataset has the greatest impact on the relative performance, followed by the few-shot setting and the source dataset. Although the rationales individually do not have a significant impact on performance, we observe significant pair-wise interaction between the rationale category and the choice of the source dataset, target dataset, and few-shot setting in decreasing order of significance. We thus glean that not only the choice of the source dataset but also the kind of rationale impacts the generalization performance.

Finally, in the in-context learning paradigm, the factors that significantly impact relative performance are the choice of the dataset, the rationale and the LLM. The pairwise interaction terms are insignificant except between the dataset and few-shot

²For the indomain setting we consider only the target

dataset

| Dataset | Label | Utterance text | Rationale Text | CAT |
|----------|------------------|---|--|--------|
| casino | vouch-fair | hey buddy I hope we both end up with a good deal:) | Expressing hope for a mutually beneficial outcome | INT |
| IMP_HATE | white_grievance | but that wouldn't enable them to destroy white neighbourhoods . | There is a belief or concern that certain actions or policies could lead to the destruction of white neighborhoods. | PreSup |
| P4G | foot-in-the-door | Every little bit help. | EE feels reassured that their small donation is still valuable. | HR |
| P4G | foot-in-the-door | Every little bit help. | Reassure the listener that any contribution is valuable. | INT |
| res_CB | Self Pity | at this i can only pay about 1600 could you do that | Seller realizes the buyer's budget constraints. | HR |

Table 3: We present instances across different datasets where adding the rationale information was crucial in predicting the correct label always. We compute Shapley values for each token in the rationale to observe its contribution to the model’s decision; the highlighted portions correspond to high positive associations with the label.

setting, and between the ICL model (i.e. Gemma and LLama) and rationale/dataset. Overall, we observe that the choice of the rationale does play a significant role on relative task performance across all experimental settings.

5.3 Qualitative Analysis

We carry out a qualitative analysis to investigate the specific instances where including the rationales improves the model’s predictions. We consider only those instances where the baseline (i.e., the utterance text) fails to predict the label correctly, but succeeds when the rationale is provided a majority of times. The distribution of these cases for both the indomain and cross-task setting appear in Figures 19 and Figures 20 in the Appendix.

The rationale with the greatest impact on performance is dependent on the nature of the task. As gleaned from Figure 19, the hearer reaction or HR has the highest impact on P4G, possibly because it captures the thought processes of the persuadee (EE) as they are being persuaded to donate. For example, the utterance “Anything would help even small donations add up when everyone pitches in.” evokes a sense of reassurance from the persuadee (EE) that any contribution is valuable and is thus recognized as a “foot-in-the-door” strategy. Presuppositions are useful for IMP_HATE, a dataset that directly references stereotypes and thus requires generic knowledge to infer the type of implicit hatred. Tasks that are centered around the outcome the speaker is invested in, i.e. strategies employed to resist persuasion (res_CB) benefit mostly from intentions. Furthermore, similar tasks e.g., CaSiNo and res_CB which deal with negotiation have similar relative performance for the same rationales.

To highlight the specific tokens in the rationales

that guide model prediction, we use the SHAPLEY values (Roth, 1988) for instances where adding the rationales always resulted in the correct answer over three seeds. We present five examples of these instances in Table 3 across four datasets for at least one kind of rationale category. We observe that the highlighted tokens in the rationale text indeed aligns well with human-intuition to explain the label category, for example the phrase “destruction of white neighbourhoods” as a signal for white-grievance or “that their small donation” as a signal for foot-in-the-door strategy in Table 3. We present additional examples of these in the Appendix G. We also conduct ablation studies on the impact of perturbations on rationale text, and the interplay between rationales and utterance on task performance in Appendix F.

6 Conclusion

We present a taxonomy for rationales, inspired by narrative modeling, that categorizes them into speaker-centric, hearer-centric, and general-world-view perspectives. Leveraging an automated framework, we generate a substantial dataset of approximately 135,000 rationale instances across diverse social dialogue datasets with different large language models (LLMs) as the backbone. Our findings demonstrate that these rationales aid task performance in both instruct-tuning and in-context learning setups. In particular, we observe significant gains in a cross-task transfer setting from incorporating rationales corresponding to the speakers’ intentions 34% of the times. Through a comprehensive quantitative analysis over 3150 experimental settings, we identify key factors that influence generalizability of rationales for different tasks.

Limitations

Some of the main limitations of our work include:

(i) Our framework SOCIAL SCAFFOLDS employs closed-source or proprietary LLMs i.e. GPT-4o and GPT-3.5-turbo to generate the rationales. Consequently we are not able to assure that the reproducibility of generating such rationales or whether the service will be discontinued. We do however, release the entire dataset of rationales for public use.

(ii) We note that our in-domain and cross-task experiments is based on a single pre-trained model, i.e. FLAN-T5 and our in-context learning experiments involved only two LLMs (Gemma-2 and Llama-3). This was a deliberate choice to help manage our computational budget. Even with a single model, we ran 630 in-domain experiments, and an additional 2520 cross-task experiments. Future work would entail exploring larger models to see the impact of rationales on model scale.

(iii) We have only focused on simple multi-label and multi-class classification tasks in this given study and that too at an utterance level. We plan to investigate whether rationales can facilitate dialogue understanding at a conversational level and help generalize to new tasks such as response generation. We defer this to future work.

(iv) While we observe the positive impact of our machine-generated rationales on task performance, and validate that the rationales are of sufficient high quality, further research is necessary to compare and contrast these machine-generated rationales from human-generated ones.

Ethical Concerns

Our research relies on the responses generated by LLMs which are known to exhibit hidden biases in their representations. While during our experiments, we encountered no potential biases in terms of offensive language or stereotypes in the generated response for our controlled setting of social meaning detection, we implore practitioners and other researchers to conduct thorough analysis before adopting our particular prompting approach for the respective use-case. We also recognize the limitations of LLM in interpreting social meanings and clarify that our conclusions, based on probabilistic model outputs, do not construe absolute facts. Moreover, we stress that the application of LLM rationales, while beneficial within our controlled research environment for understanding hu-

man intent in utterances, should not be extended uncritically beyond these confines. The use of LLM rationales in broader contexts, especially as substitutes for human judgment and rationale, is not advocated.

References

- AI@Meta. 2024. [Llama 3 model card](#).
- Yujia Bao, Shiyu Chang, Mo Yu, and Regina Barzilay. 2018. [Deriving machine attention from human rationales](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1903–1913, Brussels, Belgium. Association for Computational Linguistics.
- Taylor Berg-Kirkpatrick, David Burkett, and Dan Klein. 2012. [An empirical investigation of statistical significance in NLP](#). In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 995–1005, Jeju Island, Korea. Association for Computational Linguistics.
- Niels Ole Bernsen, Hans Dybkjær, and Laila Dybkjær. 1996. Cooperativity in human-machine and human-human spoken dialogue. *Discourse processes*, 21(2):213–236.
- Penelope Brown, Stephen C Levinson, and Stephen C Levinson. 1987. *Politeness: Some universals in language usage*, volume 4. Cambridge university press.
- Agostina Calabrese, Leonardo Neves, Neil Shah, Maarten Bos, Björn Ross, Mirella Lapata, and Francesco Barbieri. 2024. [Explainability and hate speech: Structured explanations make social media moderators faster](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 398–408, Bangkok, Thailand. Association for Computational Linguistics.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31.
- Kushal Chawla, Jaysa Ramirez, Rene Clever, Gale Lucas, Jonathan May, and Jonathan Gratch. 2021. [CaSiNo: A corpus of campsite negotiation dialogues for automatic negotiation systems](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3167–3185, Online. Association for Computational Linguistics.
- Hanjie Chen, Faeze Brahman, Xiang Ren, Yangfeng Ji, Yejin Choi, and Swabha Swayamdipta. 2023. [REV: Information-theoretic evaluation of free-text rationales](#). In *Proceedings of the 61st Annual Meeting of*

- the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2007–2030, Toronto, Canada. Association for Computational Linguistics.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Devleena Das and Sonia Chernova. 2020. Leveraging rationales to improve human task performance. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*, pages 510–518.
- Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. [The hitchhiker’s guide to testing statistical significance in natural language processing](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1383–1392, Melbourne, Australia. Association for Computational Linguistics.
- Ritam Dutt, Sayan Sinha, Rishabh Joshi, Surya Shekhar Chakraborty, Meredith Riggs, Xinru Yan, Haogang Bao, and Carolyn Rose. 2021. [ResPer: Computationally modelling resisting strategies in persuasive conversations](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 78–90, Online. Association for Computational Linguistics.
- Ritam Dutt, Zhen Wu, Jiaxin Shi, Divyanshu Sheth, Prakhar Gupta, and Carolyn Rose. 2024. [Leveraging machine-generated rationales to facilitate social meaning detection in conversations](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6901–6929, Bangkok, Thailand. Association for Computational Linguistics.
- Joshua Eisenberg and Mark Finlayson. 2016. [Automatic identification of narrative diegesis and point of view](#). In *Proceedings of the 2nd Workshop on Computing News Storylines (CNS 2016)*, pages 36–46, Austin, Texas. Association for Computational Linguistics.
- Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. 2021. [Latent hatred: A benchmark for understanding implicit hate speech](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 345–363, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- James N Farr, James J Jenkins, and Donald G Paterson. 1951. Simplification of flesch reading ease formula. *Journal of applied psychology*, 35(5):333.
- Prakhar Gupta, Cathy Jiao, Yi-Ting Yeh, Shikib Mehri, Maxine Eskenazi, and Jeffrey Bigham. 2022. [InstructDial: Improving zero and few-shot generalization in dialogue through instruction tuning](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 505–525, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Sai Gurrapu, Ajay Kulkarni, Lifu Huang, Ismini Lourentzou, and Feras A. Batareseh. 2023. [Rationalization for explainable nlp: a survey](#). *Frontiers in Artificial Intelligence*, 6.
- Sil Hamilton. 2024. [Detecting mode collapse in language models via narration](#). In *Proceedings of the First edition of the Workshop on the Scaling Behavior of Large Language Models (SCALE-LLM 2024)*, pages 65–72, St. Julian’s, Malta. Association for Computational Linguistics.
- Manoel Horta Ribeiro, Justin Cheng, and Robert West. 2023. Automated content moderation increases adherence to community guidelines. In *Proceedings of the ACM web conference 2023*, pages 2666–2676.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. *Advances in Neural Information Processing Systems*, 33:20179–20191.
- Sahil Jayaram and Emily Allaway. 2021. [Human rationales as attribution priors for explainable stance detection](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5540–5554, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yohan Jo, Elijah Mayfield, Chris Reed, and Eduard Hovy. 2020. [Machine-aided annotation for fine-grained proposition types in argumentation](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1008–1018, Marseille, France. European Language Resources Association.
- Brihi Joshi, Aaron Chan, Ziyi Liu, Shaoliang Nie, Maziar Sanjabi, Hamed Firooz, and Xiang Ren. 2022. [ER-test: Evaluating explanation regularization methods for language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3315–3336, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Brihi Joshi, Ziyi Liu, Sahana Ramnath, Aaron Chan, Zhewei Tong, Shaoliang Nie, Qifan Wang, Yejin Choi, and Xiang Ren. 2023. Are machine rationales (not) useful to humans? measuring and improving human utility of free-text rationales. *arXiv preprint arXiv:2305.07095*.
- Ratnesh Joshi, Arindam Chatterjee, and Asif Ekbal. 2021. [Towards explainable dialogue system: Explaining intent classification using saliency techniques](#). In *Proceedings of the 18th International Conference on Natural Language Processing (ICON)*, pages 120–127, National Institute of Technology Silchar, Silchar, India. NLP Association of India (NLPAI).

- Sopan Khosla and Rashmi Gangadharaiah. 2022. Benchmarking the covariate shift robustness of open-world intent classification approaches. *AACL-IJCNLP 2022*, page 14.
- Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. 2022. [ProsocialDialog: A prosocial backbone for conversational agents](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4005–4029, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jaekyeom Kim, Dong-Ki Kim, Lajanugen Logeswaran, Sungryull Sohn, and Honglak Lee. 2024. [Auto-intent: Automated intent discovery and self-exploration for large language model web agents](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16531–16541, Miami, Florida, USA. Association for Computational Linguistics.
- JP Kincaid. 1975. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. *Chief of Naval Technical Training*.
- Stefan Larson, Anish Mahendran, Joseph J Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K Kummerfeld, Kevin Leach, Michael A Laurenzano, Lingjia Tang, et al. 2019. An evaluation dataset for intent classification and out-of-scope prediction. 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 1311–1316.
- Michael K Lindell, Christina J Brandt, and David J Whitney. 1999. A revised index of interrater agreement for multi-item ratings of a single target. *Applied Psychological Measurement*, 23(2):127–135.
- Bodhisattwa Prasad Majumder, Oana Camburu, Thomas Lukasiewicz, and Julian McAuley. 2022. Knowledge-grounded self-rationalization via extractive and natural language explanations. In *International Conference on Machine Learning*, pages 14786–14801. PMLR.
- Magdalena Markowska, Mohammad Taghizadeh, Adil Soubki, Seyed Mirroshandel, and Owen Rambow. 2023. [Finding common ground: Annotating and predicting common ground in spoken conversations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8221–8233, Singapore. Association for Computational Linguistics.
- James R Martin and Peter R White. 2003. *The language of evaluation*, volume 2. Springer.
- Shikib Mehri. 2022. *Towards Generalization in Dialog through Inductive Biases*. Ph.D. thesis, Language Technologies Institute, Carnegie Mellon University.
- Saif Mohammad. 2018. [Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 English words](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 174–184, Melbourne, Australia. Association for Computational Linguistics.
- Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. [MTEB: Massive text embedding benchmark](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.
- Melissa Mulcahy and Bethanie Gouldthorp. 2016. Positioning the reader: the effect of narrative point-of-view and familiarity of experience on situation model construction. *Language and Cognition*, 8(1):96–123.
- Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayan-deh, Lars Liden, and Jianfeng Gao. 2021. Soloist: Building task bots at scale with transfer learning and machine teaching. *Transactions of the Association for Computational Linguistics*, 9:807–824.
- Baolin Peng, Chenguang Zhu, Chunyuan Li, Xiujuan Li, Jinchao Li, Michael Zeng, and Jianfeng Gao. 2020. Few-shot natural language generation for task-oriented dialog. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 172–182.
- Javiera Perez Gomez. 2021. Verbal microaggressions as hyper-implicatures. *Journal of Political Philosophy*, 29(3):375–403.
- Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Zhong Zhang, Jie Zhou, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. [Tell me more! towards implicit user intention understanding of language model driven agents](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1088–1113, Bangkok, Thailand. Association for Computational Linguistics.
- Kavel Rao, Liwei Jiang, Valentina Pyatkin, Yuling Gu, Niket Tandon, Nouha Dziri, Faeze Brahman, and Yejin Choi. 2023. [What makes it ok to set a fire? iterative self-distillation of contexts and rationales for disambiguating defeasible social and moral situations](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12140–12159, Singapore. 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.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. [Beyond accuracy: Behavioral testing of NLP models with CheckList](#). In

- Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Alvin E Roth. 1988. Introduction to the shapley value. *The Shapley value*, 1.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. [Social bias frames: Reasoning about social and power implications of language](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490, Online. Association for Computational Linguistics.
- Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. 2022. [Neural theory-of-mind? on the limits of social intelligence in large LMs](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3762–3780, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ashish Sharma, Adam Miner, David Atkins, and Tim Althoff. 2020. [A computational approach to understanding empathy expressed in text-based mental health support](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5263–5276, Online. Association for Computational Linguistics.
- Gemma Team. 2024. [Gemma](#).
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*.
- Minjia Wang, Pingping Lin, Siqi Cai, Shengnan An, Shengjie Ma, Zeqi Lin, Congrui Huang, and Bixiong Xu. 2025. Stand-guard: A small task-adaptive content moderation model. In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pages 1–20.
- Xuwei Wang, Weiyan Shi, Richard Kim, Yoojung Oh, Sijia Yang, Jingwen Zhang, and Zhou Yu. 2019a. [Persuasion for good: Towards a personalized persuasive dialogue system for social good](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5635–5649, Florence, Italy. Association for Computational Linguistics.
- Xuwei Wang, Weiyan Shi, Richard Kim, Yoojung Oh, Sijia Yang, Jingwen Zhang, and Zhou Yu. 2019b. [Persuasion for good: Towards a personalized persuasive dialogue system for social good](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5635–5649, Florence, Italy. Association for Computational Linguistics.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. 2022. Supernaturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Jason Wei and Kai Zou. 2019. [EDA: Easy data augmentation techniques for boosting performance on text classification tasks](#). 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 6382–6388, Hong Kong, China. Association for Computational Linguistics.
- Sarah Wiegrefe, Ana Marasović, and Noah A. Smith. 2021. [Measuring association between labels and free-text rationales](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10266–10284, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Chien-Sheng Wu, Steven CH Hoi, Richard Socher, and Caiming Xiong. 2020. Tod-bert: Pre-trained natural language understanding for task-oriented dialogue. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 917–929.
- Zhen Wu, Ritam Dutt, and Carolyn Penstein Rosé. 2024. Evaluating large language models on social signal sensitivity: An appraisal theory approach. In *The First Human-Centered Large Language Modeling Workshop*, page 67.
- Weimin Xiong, Yifan Song, Peiyi Wang, and Sujian Li. 2023. [Rationale-enhanced language models are better continual relation learners](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15489–15497, Singapore. Association for Computational Linguistics.
- Akhila Yerukola, Saujas Vaduguru, Daniel Fried, and Maarten Sap. 2024. [Is the pope catholic? yes, the pope is catholic. generative evaluation of non-literal intent resolution in LLMs](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 265–275, Bangkok, Thailand. Association for Computational Linguistics.

Zulipiye Yusupjiang and Jonathan Ginzburg. 2023. [Unravelling indirect answers to wh-questions: Corpus construction, analysis, and generation](#). In *Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 336–348, Prague, Czechia. Association for Computational Linguistics.

Omar Zaidan, Jason Eisner, and Christine Piatko. 2007. Using “annotator rationales” to improve machine learning for text categorization. In *Human language technologies 2007: The conference of the North American chapter of the association for computational linguistics; proceedings of the main conference*, pages 260–267.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488.

Yiming Zhang, Sravani Nanduri, Liwei Jiang, Tongshuang Wu, and Maarten Sap. 2023. [BiasX: “thinking slow” in toxic content moderation with explanations of implied social biases](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4920–4932, Singapore. Association for Computational Linguistics.

Xuhui Zhou, Hao Zhu, Akhila Yerukola, Thomas Davidson, Jena D. Hwang, Swabha Swayamdipta, and Maarten Sap. 2023. [COBRA frames: Contextual reasoning about effects and harms of offensive statements](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6294–6315, Toronto, Canada. Association for Computational Linguistics.

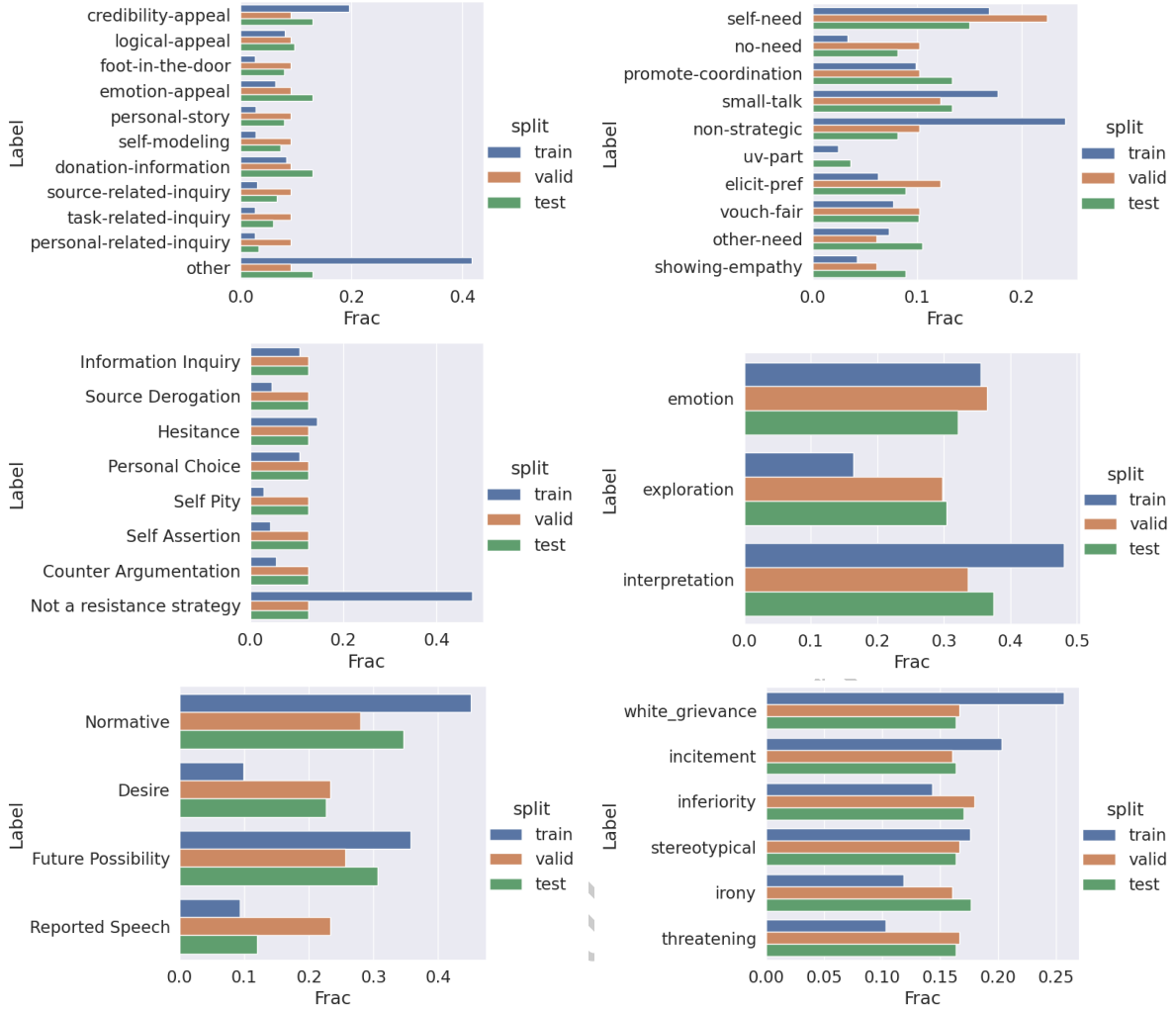


Figure 7: Distribution of labels across the different splits for the six datasets or tasks.

A Dataset Statistics

We describe in detail the six different datasets (or tasks) that we explore in this study. We showcase the distribution of the different labels across the different splits in Figure 7.

1. **Persuasion** - The task involves identifying persuasive strategies between two AMT workers where one adopts the role of the persuader and is expected to convince the other party (the persuadee) to donate to charity. We use the Persuasion for Good (P4G) dataset of Wang et al. (2019b).
2. **Negotiation tactic** - The negotiation task is grounded in the CaSiNo corpus of (Chawla et al., 2021), which consists of bargaining for campsite resources between crowd workers in a simulated camping setting. Dialogs contain various aspects of a realistic negotiation, such

as building relationships, discussing preferences, exchanging offers, emotional expression, and persuasion with personal and logical arguments.

3. **Resisting Strategies** - Complementary to task of identifying persuasive attempts, the task proposed by Dutt et al. (2021) involves detecting resisting strategies, i.e. strategies employed to resist being persuaded by others. We focus on the Craigslist Bargain dataset (henceforth res_CB) which consists of simulated conversations between a buyer (BU) and a seller (SE) over an online exchange platform. Both are given their respective target prices and employ resisting strategies to negotiate the offer.
4. **Empathy in mental health** - We use the framework and dataset of Sharma et al. (2020) that characterizes the communication of empathy

| Dataset | Dialog History | Speaker: UTT | INT | PreSup | HR |
|-------------------------------------|--|---|--|---|--|
| CaSiNo (Chawla et al., 2021) | mturk_agent_1: I am running low on firewood. I need more to keep the fire going and cook food. | mturk_agent_2: Yeah I need firewood too. It's only 50 degrees at night on the mountain, even in the summer. | Agreeing on the need for firewood and providing additional context | It gets cold at night on the mountain even in the summer. | mturk_agent_1 feels validated in their need for firewood and understands the shared predicament. |
| P4G (Wang et al., 2019b) | ER: Save the Children is an amazing charity that helps kids who are in desperate need. ER: They can help with safety, education and more. ER: You can donate some of your earnings to this amazing charity. EE: I believe in this charity, but still wonder how much of the money I donate actually helps. ER: Every little bit makes a difference. | ER: When you have people who are so poor, it's amazing what a tiny amount can do. | Emphasizing the value of small contributions | Even minimal financial aid can greatly benefit those in extreme poverty. | EE is encouraged by the impact of small donations. |
| res_CB (Dutt et al., 2021) | Buyer: Hi there. I was looking for ads and this one caught my attention. Is it in a good and working condition? Seller: It is, it's been used a lot less than its age would suggest. I only rode it a few times a month. Are you interested? | Buyer: Yes very much so. | Express strong interest in the item | The buyer is expressing a strong interest in the item indicating a willingness to negotiate further. | The seller feels encouraged by the buyer's interest and is prompted to discuss the price. |
| IMP_HATE (ElShrief et al., 2021) | | Poster: flynn's resignation set a dangerous precedent for the administration . #sessions is essential to justice 4 white america . he must not #resign | Expressing concern about the implications of Flynn's resignation and emphasizing the importance of Sessions to their view of justice for white America while urging that Sessions should not resign. | The resignation of a high-ranking official can have significant and potentially negative consequences for the administration. | The hearer may feel concerned about the implications of Flynn's resignation and the potential impact on the administration. They may also feel a sense of urgency or importance regarding Sessions' role and the need for him to remain in his position. |
| EMH (Sharma et al., 2020) | Seeker: Why do I always have good news followed by a shit night, followed by sitting up at 2am wanting to kill myself? Why is life so difficult? Why is it so impossible to be fucking happy for once in my shit fucking life? What's the point anymore? | Responder: well not for nothing but you made it extremely difficult to read your post by only using a period in the title. JUST saying not judging. | Pointing out the difficulty in reading the post due to formatting while attempting to clarify that they are not judging. | Clear communication is important for understanding and responding to others' concerns effectively. | The Seeker may feel invalidated or criticized as the Responder's comment focuses on the format of the post rather than addressing the Seeker's emotional distress. |
| PROP (Jo et al., 2020) | S_1: It is called the Constitution of the United States S_2: unfortunately, those few months gave us OBAMA S_3: We're going to win when we unite people with a hopeful, optimistic message S_3: we had high sustained economic growth | S_3: We created 1.3 million jobs | Emphasizing job creation | Creating jobs is a positive achievement. | Impression of job creation success |

Table 4: Examples of rationales generated by GPT-4o for six utterances, each coming from a different dataset and task. For each utterance, we provide the dialog history and the corresponding intention, presupposition, and hearer reaction abbreviated as INT, PreSup, and HR respectively. The rationales score high on factuality, soundness, and relevance as evaluated by two annotators.

Table 5: Description of the resisting strategies used in our work for the res_CB (Dutt et al., 2021). Examples of each strategy are italicised.

| Resisting Strategy | Description |
|-----------------------|--|
| Source Derogation | Attacks the other party or questions the item <i>Was it new denim, or were they someone's funky old worn out jeans?</i> |
| Counter Argumentation | Provides a non-personal argument/factual response to refute a previous claim or to justify a new claim. <i>It may be old, but it runs great. Has lower mileage and a clean title.</i> |
| Personal Choice | Provides a personal reason for disagreeing with the current situation or chooses to agree with the situation provided some specific condition is met. <i>I will take it for \$300 if you throw in that printer too.</i> |
| Information Inquiry | Requests for clarification or asks additional information about the item or situation. <i>Can you still fit it in your pocket with the case on?</i> |
| Self Pity | Provides a reason (meant to elicit sympathy) for disagreeing with the current terms. <i>\$130 please I only have \$130 in my budget this month.</i> |
| Hesitance | Stalls for time and is hesitant to commit; specifically, they seek to further the conversation and provide a chance for the other party to make a better offer. <i>Ok, would you be willing to take \$50 for it?</i> |
| Self-assertion | Asserts a new claim or refutes a previous claim with an air of finality/ confidence. <i>That is way too little.</i> |

Table 6: Description of the negotiation strategies used in our work for Casino (Chawla et al., 2021). Examples of each strategy are italicised.

| Negotiation Label | Description |
|----------------------|---|
| self-need | Participant argues for creating a personal need for an item in the negotiation. <i>Yes. I'm actually taking a large group of people. Some friends and family are going and I kind of also wanted a bit of extra firewood. :)</i> |
| no-need | Participant points out that they do not need an item based on personal context. <i>I don't like food. my stomach is always full. I only drink water since im thirsty most of the time.</i> |
| promote-coordination | Participant promotes coordination between the two partners. <i>Alright so I think we can make a fair deal here where we both will be happy. :)</i> |
| small-talk | Participant engages in small talk while discussing topics apart from the negotiation in an attempt to build a rapport. <i>My mistake, hypothermia is messing with my brain.</i> |
| uv-part | Participant undermines the requirements of their opponent. <i>I understand that atleast you are going to be close to water; that will be our most important thing since we will be thirsty and you know kids and trying to tell them to ration the water...LOL</i> |
| elicit-pref | Participant provides an attempt to discover the preference order of the opponent <i>I get that and understand completely. I have a large number of mouths to feed making the food a necessity or all the firewood to cook whatever we hunt. How many you have?</i> |
| vouch-fair | Participant announces a callout to fairness for personal benefit, either when acknowledging a fair deal or when the opponent offers a deal that benefits them <i>hey buddy I hope we both end up with a good deal :)</i> |
| other-need | Participants discuss a need for someone else rather than themselves. <i>I would be willing to do that if I could have two of the waters? I didn't bring as much as I thought I would need because I forgot I would have my dog.</i> |
| showing-empathy | Participant positively acknowledges or displays empathetic behavior towards a personal context of the partner. <i>Are you sure that's enough firewood for you and the baby? I know that babies can easily get very sick from dropping temperatures.</i> |
| non-strategic | Utterance does not have any strategic element <i>oh well that's fantastic, so let me know what your deal is</i> |

Table 7: Description of the different dimensions of empathy used in our work for EMH (?). Examples of each strategy are italicised.

| Empathy Dimension | Description |
|-------------------|---|
| emotion | Responder specifies the experienced emotions explicitly or communicates emotional reactions without the emotions being explicitly labeled. <i>I love you all. Genuinely the replies here made me feel much less alone. your kindness is beautiful</i> |
| interpretation | Responder specifies the inferred feeling or experience, or communicates understanding through descriptions of similar experiences, or mentions that they understand the situation. <i>I struggle with the same thing, but with my step dad. I've basically given up on it.</i> |
| exploration | Responder seeks to explore the specific experiences and feelings of the supporter, which can be specific to the situation or generic. <i>What kind of connections do you want / need?</i> |

Table 8: Description of the argumentation labels used in our work for PROP (Jo et al., 2020). Examples of each strategy are italicised.

| Argumentation Label | Description |
|---------------------|--|
| Normative | A normative proposition is defined as a proposition where the speaker or someone else proposes that a certain situation should be achieved or that an action should be carried out. <i>Americans have a right to know if there are any conflicts of interest</i> |
| Desire | A desire proposition is defined as a proposition that explicitly claims that the speaker or someone else desires to own something, do something, or desires for a certain situation to be achieved. <i>Bernie Sanders' soft stance on guns is not what Americans want to hear right now</i> |
| Future Possibility | A future possibility proposition claims a possibility or prediction that something may be the case in the future. <i>the next president of the US could very well be a reality TV star</i> |
| Reported Speech | A reported speech proposition must convey an explicit or implicit predicate borrowed from a source external to the speaker. <i>Clinton's view that sometimes capitalism needs controls in order to deal with inequality were a lot more realistic</i> |

Table 9: Description of the argumentation labels used in our work for IMP_HATE (ElSherief et al., 2021). Examples of each strategy are italicised.

| Type of Hate | Description |
|-----------------|--|
| white_grievance | The message includes frustration over a minority groups perceived privilege and casting majority groups as the real victims of racism. This language is linked to extremist behavior and support for violence <i>is zuma talking about seizing white land or not ? are you denying it ?</i> |
| incitement | The message includes flaunting in group unity and power or elevating known hate groups and ideologies. <i>india wud be a peaceful prosperous n progressive country only when each of its citizen is a hindu.that shud b d endall objective 4 d state</i> |
| inferiority | The message implies one group or individual is inferior to another, and it can include dehumanization, denial of a person's humanity, and toxic language that compares the target with disease, insects, animals . Related to assaults on human dignity, dominance, and declarations of superiority of the in group <i>have you ever written about islam being a crutch for a failed civilization? im new to your work, just wondering.</i> |
| stereotypical | The message associates a protected class with negative attributes such as crime, or terrorism or includes misinformation that feeds stereotypes and vice versa, like holocaust denial and other forms of historical negationism. <i>You can't be a person of colour; you're too pale!</i> |
| irony | The message uses sarcasm, humor, and satire to attack or demean a protected class or individual. <i>What's the one good thing about black people? They provide jobs for the prison guards</i> |
| threatening | The message conveys a speaker's commitment to a target's pain, injury, damage, loss or violation of rights, threats related to implicit violation of rights and freedoms, removal of opportunities, and more subtle forms of intimidation. <i>We have this huge military. Why don't we just go down there and create an ethno-state for whites. Most of the blacks weren't even there when South Africa was founded by whites!</i> |

Table 10: Description of the persuasion labels used in our work for P4G(Wang et al., 2019b). Examples of each strategy are italicised.

| Persuasion Label | Description |
|--------------------------|---|
| credibility-appeal | Refers to the uses of credentials and citing organizational impacts to establish credibility and earn the persuadee’s trust <i>It is the worlds first global charity for children, and have credentials to back them up.</i> |
| logical-appeal | Refers to the use of reasoning and evidence to convince others. <i>You are donating money you don’t even have yet so it is not like you are missing something.</i> |
| foot-in-the-door | Refers to the strategy of starting with small donation requests to facilitate compliance followed by larger requests." <i>Are you sure, you can do as little as 5 cents???</i> |
| emotion-appeal | Refers to the elicitation of specific emotions to influence others in the form of story-telling, empathy, guilt, or anger" <i>It broke my heart to see that famous photograph of a child with a vulture sitting next to it.</i> |
| personal-story | Refers to the strategy of using narrative exemplars to illustrate someone’s donation experiences or the beneficiaries’ positive outcomes, which can motivate others to follow the actions." <i>I have three children myself, and the welfare of children around the world is a very important cause to me.</i> |
| self-modeling | Refers to the strategy where the persuader first indicates their own intention to donate and chooses to act as a role model for the persuadee to follow" <i>I think I am going to give a small portion of my hit payment to save the children.</i> |
| donation-information | Refers to providing specific information about the donation task, such as the donation procedure, donation range, etc." <i>The research team will collect all donations and send it to Save the Children.</i> |
| source-related-inquiry | Asks about the persuadee’s opinion and expectation related to the task." <i>Isn’t alright, just reading up on this organization called "Save the Children".. have you heard about it?</i> |
| task-related-inquiry | Asks if the persuadee is aware of the organization (charity) <i>Do you need more info about this program?</i> |
| personal-related-inquiry | Asks about the persuadee’s previous personal experiences relevant to charity donation" <i>I imagine hospitals are very strict about who gets to be with the little ones.</i> |
| other | Does not conform to any persuasion category <i>I am homeless and at Mcdonalds on the wif.</i> |

in text-based conversations. The task involves detecting different dimensions of empathy in text-based mental health support, i.e., empathy expressed or communicated by peer supporters in their textual interactions with seekers.

5. Argumentation - We formalize the task of argumentation into identifying different kinds of proposition in rhetorical debates. We use the data set of Jo et al. (2020) which consists of four categories of propositions: normative statements, desires statements, statements about future possibilities, and reported speech.
6. Implicit Hate Speech Detection - The task involves identifying different categories of covert or indirect language that disparages a particular individual or group based on certain protected attributes (ElSherief et al., 2021). Some instances include irony, inferiority language, and incitement to violence, among others.

We also provide descriptions of the label categories for each dataset along with an example of each for res_CB, Casino, EMH, PROP, IMP_HATE, and P4G in the Tables 5, 6, 7, 8, 9, and 10 respectively.

B Prompting Framework Description

SOCIAL SCAFFOLDS takes as input a multiparty dialog and generates rationales on an utterance-by-utterance basis. This is achieved using a Large Language Model (such as GPT-4o) that goes over each utterance in the conversation and generates the corresponding rationale. We instruct the framework to generate the outputs in a structured format, i.e. the rationales are generated in the form of a CSV file and aligned with the corresponding speaker and utterance index. These checks and measures help ensure that each utterance has a corresponding rationale and enables us to revisit erroneous cases. We address those misaligned dialogs by simply prompting the framework to regenerate the rationales for those dialogs in an iterative fashion. After 3 iterations, the fraction of valid dialogs

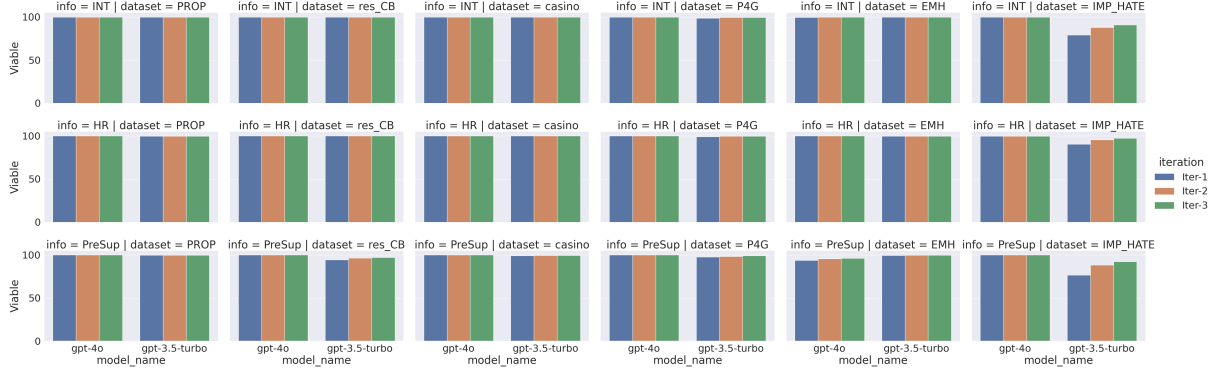


Figure 8: Validity of rationales over iterations for different datasets.

whose utterances have their corresponding rationale is 99.2%. We show the impact of iterations on the validity of these rationales in Figure 8 in the Appendix.

We reuse the prompting framework to generate each category of rationale separately. The motivation for our design choice is two-fold. Firstly, we wish to observe whether the different rationale categories can capture distinct concepts; by forcing the framework to generate the rationales together would make it sensitive to ordering effects, for e.g. if the intentions are generated first, then those intentions would influence the generation of presuppositions. Secondly, our framework is easily generalizable to new categories of rationales. We actually explore a few other categories of rationales such as the literal meaning of the utterance or the dialog acts, which we defer in the Appendix.

Additionally, we do not provide any few-shot instances for in-context learning while generating these rationales to avoid biasing the generations with previously seen examples as in [Dutt et al. \(2024\)](#). Overall, our framework enables us to compare and contrast not only different categories of rationales with each other but also the same categories of rationales generated by different LLMs.

C Annotation Guidelines

C.1 Metrics for Annotating Rationales

Since our framework automatically generates rationales without any human supervision, we develop a rigorous annotation framework to assess the validity of generations. To validate the quality of rationales, we define the following three criteria: soundness, informativeness, and relevance.

Soundness: Soundness reflects whether the rationale adheres to the definition provided during prompting, i.e. whether the generated rationale

reflects the speaker’s intentions, the hearer’s reactions, and the presuppositions about the world. In some cases, the rationale generated might not contain any additional subtext beyond the literal rephrasing of the utterance. Such instances are scored high on soundness.

Informativeness: The information conveyed by the rationales should comply with the context of the current dialogue. The information should be correct, i.e. rationale should not exhibit hallucination, (present additional information that has not been encountered so far in the dialogue), and complete, i.e. they should not omit important information that could change the meaning of the utterance.

Relevance: A rationale is relevant when it goes beyond the utterance text and presents information that is not only factual and sound but also provides additional subtext. We include this metric to assess whether the rationale is useful or not for the current scenario by providing important information or cues that are not directly observable.

We score each rationale based on soundness, informativeness, and relevance using a Likert scale of 1 to 3, with 1 being the lowest and 3 the highest. The evaluations were carried out by two annotators with a graduate level proficiency in English and at least five years of experience in computational linguistics and NLP. Due to the highly subjective nature of the task, we relied on these professional annotators as an alternative to crowd-sourcing or employing an automated annotation framework. We also follow the appropriate protocols to assure the annotation and data aligned with institutional approval guidelines.

We compute the inter-rater reliability scores (IRR) using the multi-item agreement measure of [Lindell et al. \(1999\)](#) and observe strong agreement scores for all three criteria: soundness (0.983), in-

formativeness (0.763), and relevance (0.697).

C.2 Flowchart for Scoring Rationales

We present the flowchart for annotating rationales according to soundness, informativeness, and relevance.

Step 1: Read the dialogue history, utterance and the rationale; start with judging the Speaker Intention rationale. Perform Steps 2-4 for the Speaker Intention rationale and then reiterate for Hearer Reaction and Presuppositions.

Step 2: Check for Soundness criteria if the generated rationale encapsulates the meaning of the rationale category. When checking for Speaker Intention rationales, see if it is about the speaker’s beliefs, goals, objectives, outcomes. When checking for Hearer Reaction see if it is about the belief of the hearer or their interpretation. When checking for Presuppositions see if it reflects the general world view or the assumptions shared by the participants.

- If the rationale is ascribing the correct perspective, we assign a 3 to Soundness.
- If the perspective appears to be ambiguous, we assign 2 for Soundness.
- If the perspective is blatantly incorrect, for example the Hearer Reaction actually reflects the speaker’s intentions we assign 1 to Soundness.
- If Soundness is 1 all criteria should be assigned 1, since it does not make sense to evaluate a wrong rationale.

Step 3: We now check whether the rationale is Informative or not, i.e. whether the information present in the rationale is accurate.

- If all the details have been carried over from the utterance, with an appropriate level of generalization assign a 3 to Informativeness.
- If the generalization has omitted some information/details that are important to the meaning of the utterance, assign a 2 for Informativeness.
- If the rationale hallucinates information, i.e. presents information that cannot be inferred from the current dialogue context, or is otherwise just wrong, assign a 1 for Informativeness.

Note that Informativeness and Relevance are always 1 when the Soundness is 1.

Step 4: We finally check for Relevance.

- If the utterance has a subtext and the rationale has identified a subtext not overtly stated in the utterance text, assign a 3 for Relevance.
- If the rationale includes information that appears earlier in the dialogue history whether it is subtext or not, but is not in the particular utterance, assign a 3 for Relevance.
- If the utterance lacks subtext, but the rationale presents an expression or action not found in the utterance, such as expressing agreement or an opinion, assign a 3 for Relevance.
- If the utterance lacks subtext and the rationale simply summarizes the details of the given utterance without adding anything new at all, assign a 2 for Relevance.
- If the utterance has an underlying subtext but that is not captured by the rationale, or an incorrect subtext is present, assign a 1 for Relevance.

D Experimental Details and Hyper-Parameter Tuning

We present the hyperparameters for our experiments in Table 11. We carry out the experiments over 3 seeds on a A6000 GPU with early stopping with patience of 5 over the validation set for all experiments. We implement the entire experiments in Python, with help of the Pytorch library and use the pre-trained models as specified in Huggingface under the agreed upon license agreements. We explicitly specify the software libraries and their corresponding versions in Table 12

Our experimental suite comprises encompasses 6 datasets in the indomain setting for the FLAN-T5 models for 5 few-shot settings (5, 10, 20, 50, and all) across 3 seeds and for 7 cases, corresponding to the 3 types of rationales (INT, HR, PreSup), for each of the two LLMs (GPT-3.5-turbo and GPT-4o) and the baseline (UTT). Furthermore, for a model pre-trained on a given source task, we further fine-tune it for 4 k-shot settings (5, 10, 20, and 50) for each of the 5 different target tasks. This results in a massive experimental suite of 630 in-domain experiments and 3150 cross-task experiments.

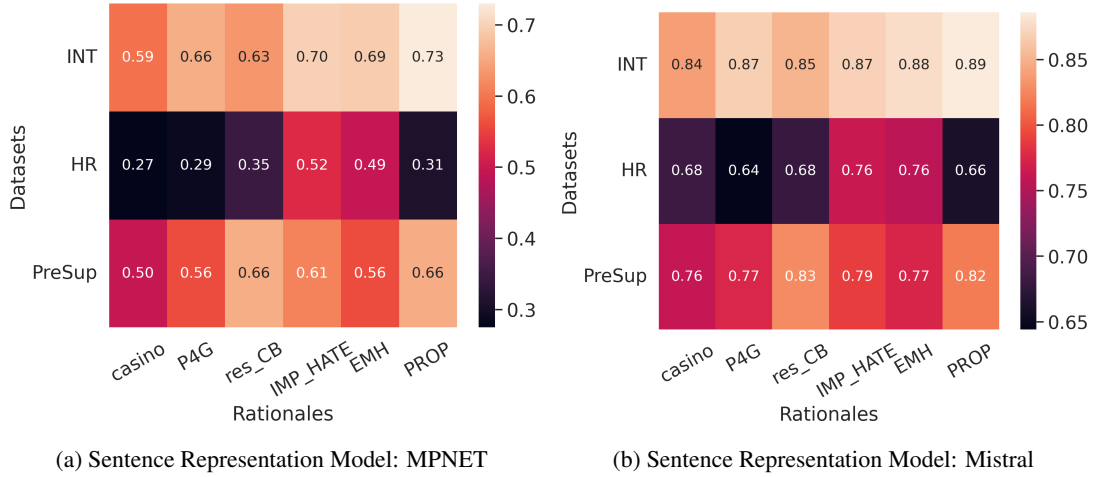


Figure 9: Cosine similarities between rationales generated by two LLMs, GPT-4o and GPT-3.5-turbo, across different datasets and rationale categories. The figures displayed on the left and right correspond to the models Mistral and MPNET, respectively.

For our incontext learning setting, we experiment over instruct-tuned versions of two open-sourced models, i.e. LLama-3-8B and the Gemma-9B. To account for prompt sensitivity, the prompts used for inference were first validated on the development split for each of the 6 datasets. We mention the final prompt used in our experiment below.

The total cost of the OpenAI credits during the course of our experiments to generate the rationales was approximately USD 265 USD, with the cost of the GPT-4o model being approximately 10 times as costly as the GPT-3.5-turbo version.

Table 11: Hyperparameters used for fine-tuning the FLAN-T5-base model for all the experiments.

| Hyperparameter | Value |
|---------------------|-----------|
| Max sequence length | 1024 |
| Learning rate | $2e^{-5}$ |
| Batch size | 8 |
| Num. epochs | 10 |
| Optimizer | Adam |
| Patience | 5 |
| Seeds | 3 |
| ICL | |
| Temperature | 0.9 |
| Fewshot examples | [0, 2, 5] |
| Batch size | 8 |
| GPUs | A6000 *2 |

Table 12: Versions of Library used in our work.

| Libraries | Version |
|-----------------------|--------------|
| Python | 3.9.12 |
| torch | 1.12.1+cu113 |
| transformers | 4.40.2 |
| numpy | 1.24.2 |
| sklearn | 1.2.2 |
| sentence-transformers | 2.7.0 |

E Analysis of Rationale Characteristics

We measure the similarity of the generated rationales across three fronts:

(i) How similar are the three different categories of rationales to each other?

(ii) How similar are the rationales generated by different LLMs for the same rationale category?

(iii) How similar is a generated rationale to its corresponding utterance?

We use cosine distance between the sentential representations as the metric for quantifying similarity. We explore two models to generate these representations, i.e., the popular MPNET model of (Reimers and Gurevych, 2019) for its simplicity and the instruction-tuned version of Mistral-7B (Wang et al., 2023) for its superior performance on the MTEB leaderboard (Muennighoff et al., 2023). We present the similarity scores across different LLMs, different rationale categories, and between the utterance and the rationale in Figures 9, 10, and 11 respectively.

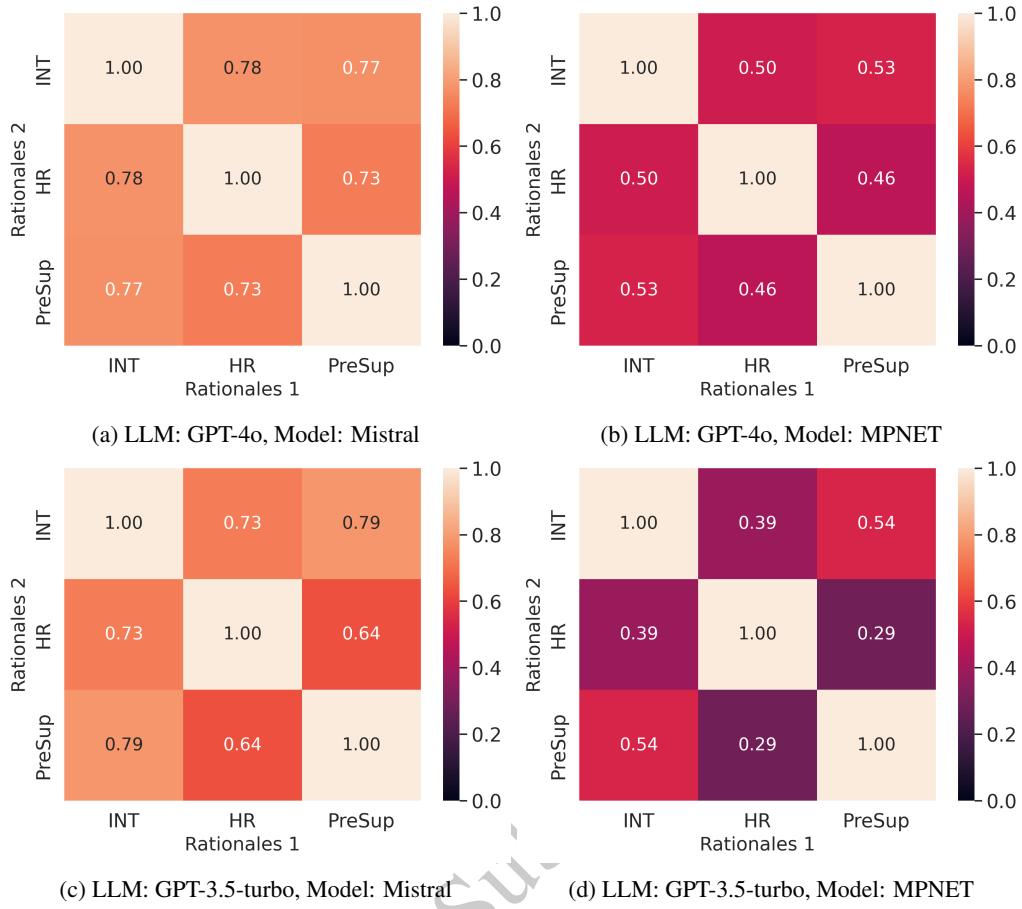


Figure 10: Cosine similarities between different categories of rationales corresponding to intentions, hearer reactions, and presuppositions as generated by two LLMs, GPT-4o and GPT-3.5-turbo, and evaluated by the sentence transformers, i.e. Mistral and MPNET.

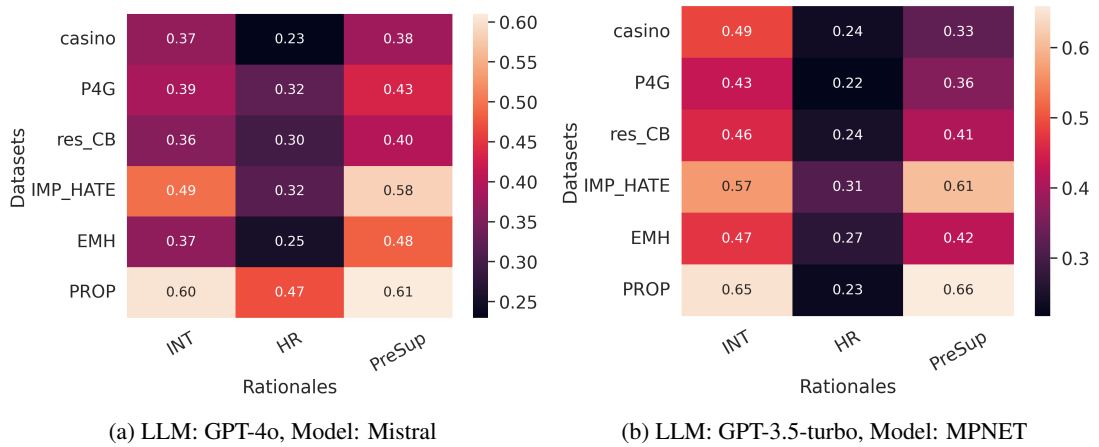


Figure 11: Cosine similarities between the original utterance and the rationales generated by different LLMs and evaluated by the sentence transformers, i.e. Mistral and MPNET.

| Factor | GPT-4-o | | | GPT-3.5-turbo | | |
|---------------------------------|---------|--------|--------|---------------|--------|--------|
| | INT | HR | PreSup | INT | HR | PreSup |
| Length of the Rationale | -0.069 | -0.053 | -0.056 | -0.060 | -0.054 | -0.057 |
| Length of the dialogiue context | 0.047 | 0.052 | 0.046 | 0.050 | 0.056 | 0.051 |
| Label Similarity | -0.058 | -0.062 | -0.041 | -0.064 | -0.039 | -0.011 |
| Utterance Similarity | -0.019 | 0.020 | -0.017 | -0.007 | -0.029 | -0.022 |
| Valence | 0.016 | 0.059 | 0.035 | 0.025 | 0.029 | 0.023 |
| Arousal | -0.014 | -0.011 | -0.003 | -0.002 | -0.019 | 0.004 |
| Dominance | 0.005 | 0.053 | 0.026 | 0.012 | 0.022 | 0.000 |
| Emotional Intenstisy | -0.010 | -0.036 | -0.022 | -0.028 | -0.036 | -0.032 |
| Emotional Polarity | -0.010 | -0.036 | -0.022 | -0.028 | -0.036 | -0.032 |
| Empathy | -0.010 | -0.036 | -0.022 | -0.028 | -0.036 | -0.032 |
| Flesch’s Reading Scale | 0.021 | 0.034 | 0.019 | 0.027 | 0.003 | 0.026 |

Table 13: Correlation of different factors with classification accuracy for different rationales generated by the two models.

| Category | F-statistic | p-value |
|-----------------------|-------------|-----------------|
| C(LLM) | 0.1761 | 6.75E-01 |
| C(RAT) | 27.2818 | 6.03E-12 |
| C(Dataset) | 6.5388 | 6.74E-06 |
| C(fewshot) | 27.8057 | 8.92E-21 |
| C(Dataset):C(LLM) | 1.2790 | 2.72E-01 |
| C(RAT):C(Dataset) | 4.6992 | 2.01E-06 |
| C(LLM):C(RAT) | 3.1047 | 4.57E-02 |
| C(fewshot):C(LLM) | 1.2457 | 2.91E-01 |
| C(RAT):C(fewshot) | 3.7960 | 2.46E-04 |
| C(fewshot):C(Dataset) | 17.1829 | 2.06E-44 |

Table 14: The F-statistics and corresponding p-value for the multi-variate ANOVA analysis to investigate the factors that characterize the performance difference in an indomain setting for instruction tuned models.

We observe similar trends in the scores regardless of the model used to generate the representations, i.e., MPNET and Mistral. The rationales generated by GPT-4o and GPT-3.5-turbo vary considerably in their similarity scores depending on their category; those corresponding to the speaker’s intentions (INT) are the most similar, followed by pre-suppositions (PreSup), while the hearer reactions (HR) are highly dissimilar. Furthermore, we note a low similarity between rationales corresponding to different categories (the weakest scores occur between PreSup and HR) and between the rationale and the original utterance. Overall, these results highlight that the categories capture perspectives distinct from each other and the original utterance.

| Category | F-statistic | p-value |
|-------------------------------|-------------|-----------------|
| C(LLM) | 0.9177 | 3.38E-01 |
| C(RAT) | 1.9741 | 1.39E-01 |
| C(fewshot) | 10.7986 | 1.11E-08 |
| C(src_dataset) | 5.2840 | 3.08E-04 |
| C(tgt_dataset) | 11.1723 | 5.50E-09 |
| C(LLM):C(RAT) | 0.1824 | 8.33E-01 |
| C(LLM):C(fewshot) | 0.9177 | 4.53E-01 |
| C(LLM):C(src_dataset) | 0.3452 | 8.86E-01 |
| C(LLM):C(tgt_dataset) | 0.8948 | 4.84E-01 |
| C(fewshot):C(RAT) | 1.9741 | 4.59E-02 |
| C(src_dataset):C(fewshot) | 5.3249 | 1.78E-13 |
| C(fewshot):C(tgt_dataset) | 10.5797 | 1.76E-32 |
| C(RAT):C(src_dataset) | 2.3990 | 7.83E-03 |
| C(RAT):C(tgt_dataset) | 1.9911 | 3.06E-02 |
| C(src_dataset):C(tgt_dataset) | 5.0937 | 1.13E-12 |

Table 15: The F-statistics and corresponding p-value for the multi-variate ANOVA analysis to investigate the factors that characterize the performance difference in a cross-task transfer setting for instruction tuned models.

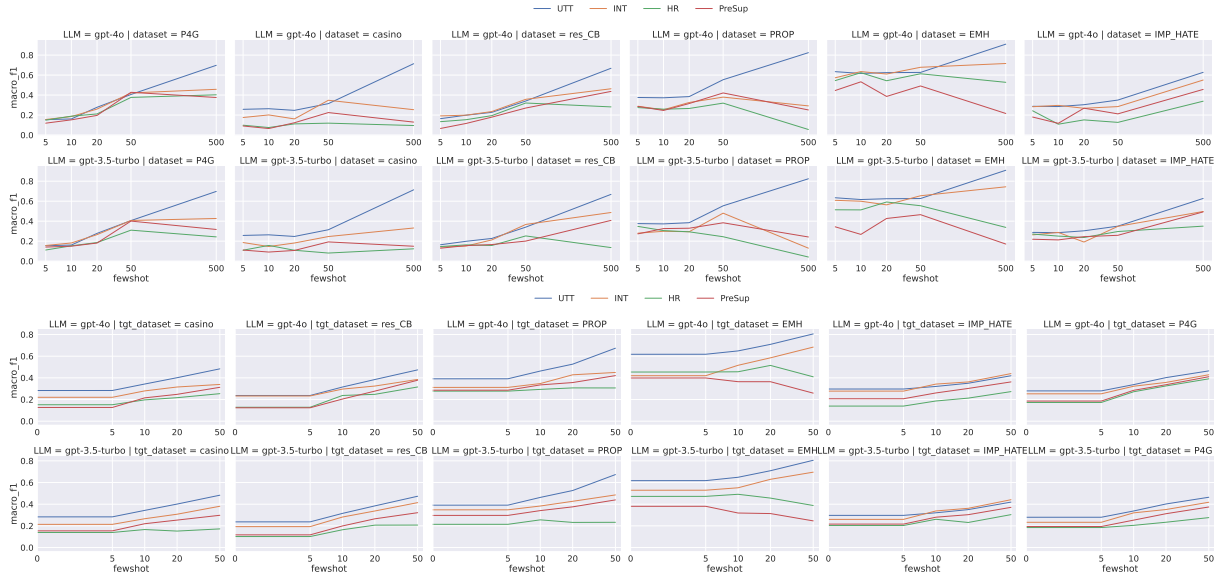


Figure 12: In-domain performance (top) and cross-task performance of models in presence of only the rationale across different few-shot cases. Note that the model was trained on BOTH the rationale and utterance.

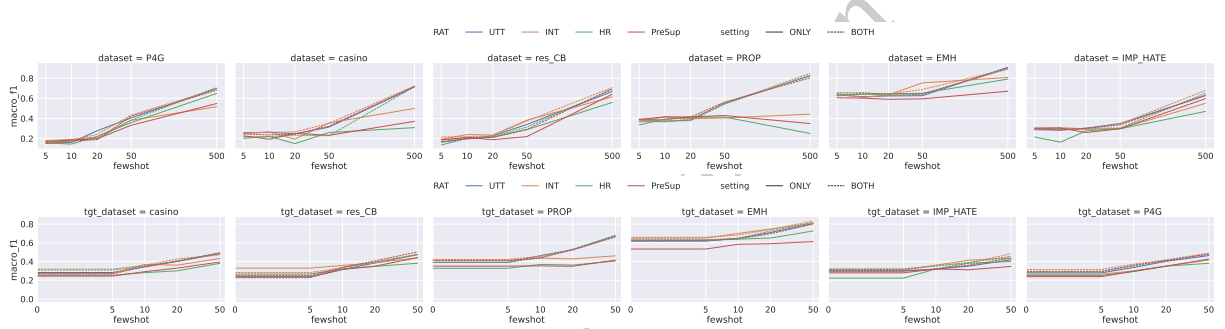


Figure 13: In-domain performance (top) and cross-task performance (below) of models using only the rationale across different few-shot cases. Note that the model was trained on ONLY the rationale.

| Category | F-statistic | p-value |
|--------------------------|-------------|-----------------|
| C(LLM) | 5.7572 | 1.76E-02 |
| C(RAT) | 13.8255 | 2.88E-06 |
| C(dataset) | 7.2547 | 3.74E-06 |
| C(fewshot) | 0.4060 | 6.67E-01 |
| C(model_name) | 2.9662 | 8.69E-02 |
| C(LLM):C(RAT) | 1.8923 | 1.54E-01 |
| C(LLM):C(dataset) | 0.3870 | 8.57E-01 |
| C(LLM):C(fewshot) | 0.7054 | 4.95E-01 |
| C(LLM):C(model_name) | 0.6620 | 4.17E-01 |
| C(RAT):C(dataset) | 0.6843 | 7.38E-01 |
| C(RAT):C(fewshot) | 1.1929 | 3.16E-01 |
| C(RAT):C(model_name) | 3.9246 | 2.17E-02 |
| C(dataset):C(fewshot) | 8.2394 | 1.02E-10 |
| C(dataset):C(model_name) | 2.8153 | 1.82E-02 |
| C(fewshot):C(model_name) | 0.2097 | 8.11E-01 |

Table 16: The F-statistics and corresponding p-value for the multi-variate ANOVA analysis to investigate the factors that characterize the performance difference in fewshot setting for in-context learning models.

F Ablation Results

F.1 Importance of the utterance information

We carry out ablation studies to investigate the role of the utterance on task performance i.e. how does the performance vary when we omit out the utterance and evaluate the fine-tuned model using only the rationale. We explore two settings: (i) where the model is provided with both the utterance and rationale information during training, but use only the rationale during inference, (see Figures 12) and (ii) where we train and test the model with only the rationale as an augmentation (see Figure 13).

We observe a noticeable degradation in performance compared to the baseline (the model is trained only on the utterance) in the former case for both the indomain and cross-task setting; the drop progressively increases with the amount of training data, highlighting that fine-tuned models do not solely rely on the rationale to make its predic-



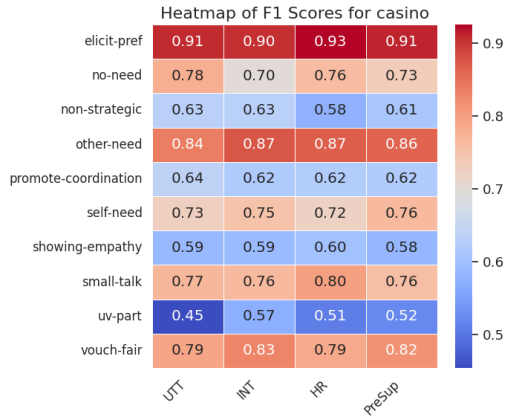
Figure 14: Impact of different kinds of perturbation on the rationale text for classification performance.

tions. The latter scenario where the model is fine-tuned with only the rationales fares better, albeit still falling short of the baseline in the in-domain setting. When trained on only the rationale information, the impact of the rationale category on the task performance becomes more pronounced. We see higher gains from adding the hearer reactions to P4G, the presuppositions to IMP_HATE, and the intentions to casino, and EMH. In the cross-task setting, the performance drop is almost negligible; in fact we see marked improvements for res_CB, IMP_HATE and EMH with the intention rationales over the baseline. In short, we see that the utterance information is crucial for task performance and though rationales provides a useful augmentation, they cannot be used as a replacement or substitute for the utterance. Future work needs to inspect how to design free-text rationales that can capture all the salient aspects of the utterance (Chen et al., 2023).

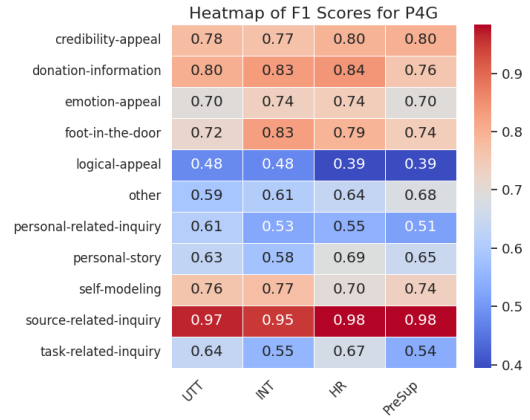
F.2 Perturbation of the Rationales

We also carry out sensitivity analysis of the rationales by observing how perturbing the rationale text affects task performance. We compare different kinds of perturbations such as synonym swap using Checklist (Ribeiro et al., 2020) and WordNet, different kinds of augmentations (EmbedDA), deletions or combination of them (EDA) (Wei and Zou, 2019). We also control for the fraction of words being perturbed in the rationale text i.e. 10%, 50% and 90%. We depict the change in task performance due to perturbations in Table 14

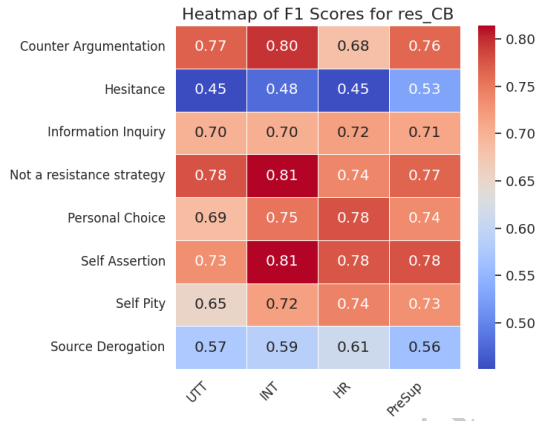
Overall, on a macro scale, we observe that perturbations indeed decrease task performance with the deterioration becoming more pronounced as the proportion of words being perturbed increases. We also note that certain methods are more effective than others such as deletion as opposed to synonym matching or entity replacement. Such an analysis highlights that the instruct-tuned model does rely on the rationales for classification.



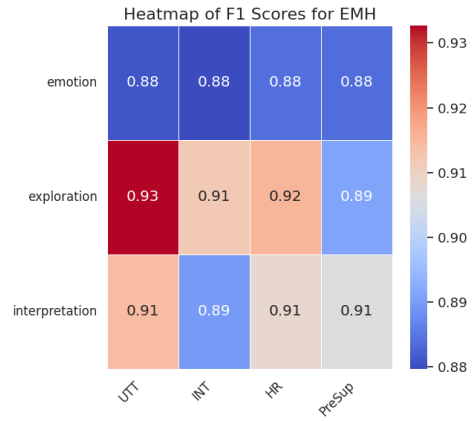
(a) Casino



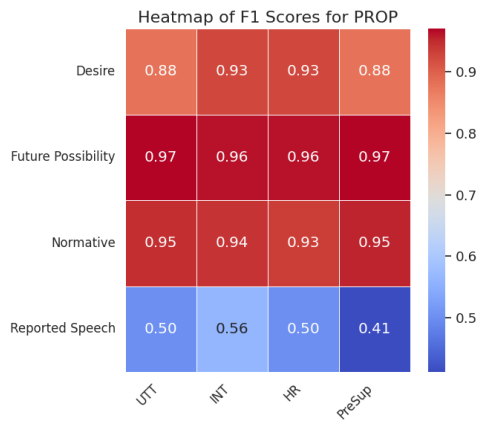
(b) P4G



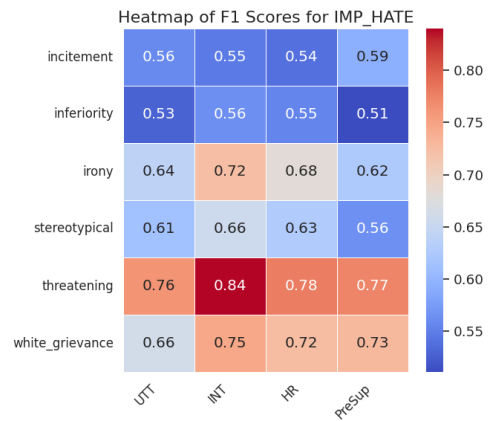
(c) res_CB



(d) EMH

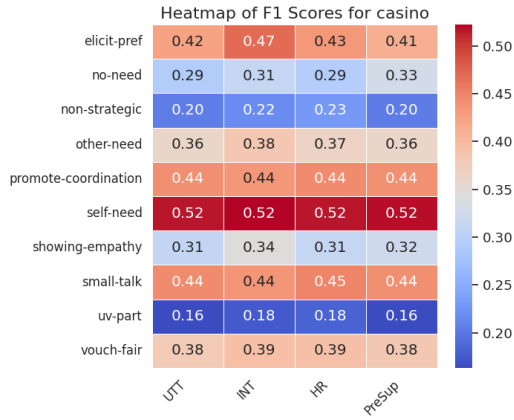


(e) PROP

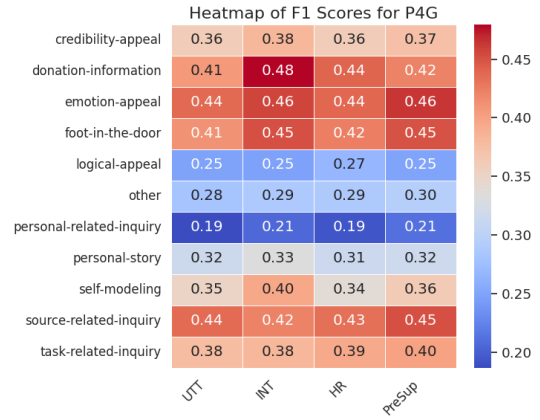


(f) IMP_HATE

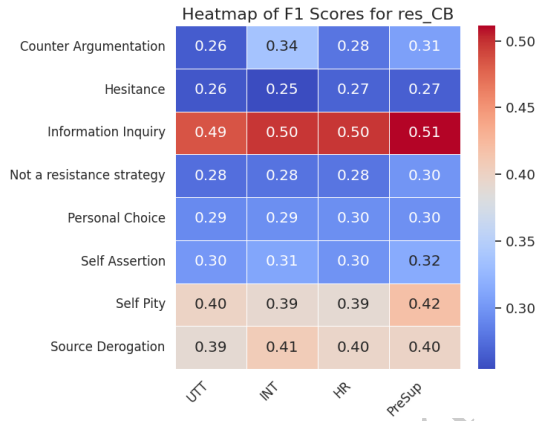
Figure 15: Comparative performance of rationales in terms of macro F1 score across different labels for different tasks in an indomain setting.



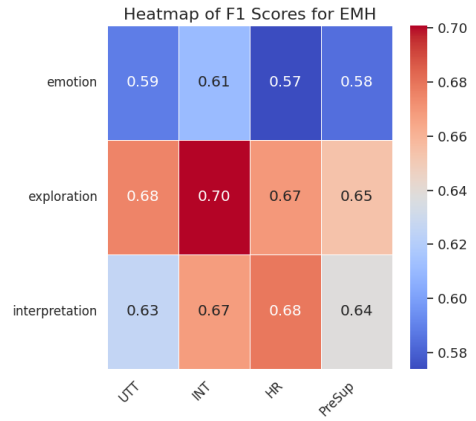
(a) Casino



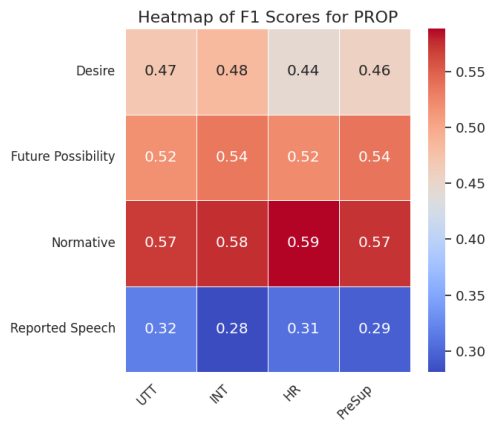
(b) P4G



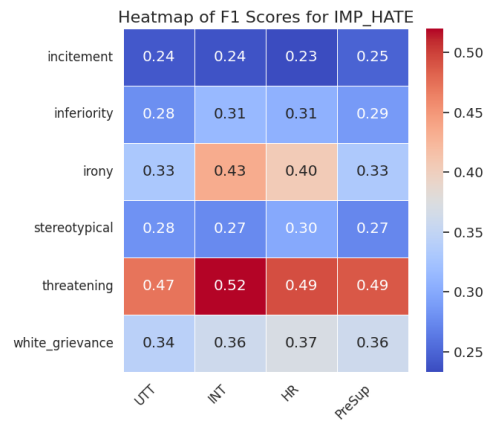
(c) res_CB



(d) EMH



(e) PROP



(f) IMP_HATE

Figure 16: Comparative performance of rationales in terms of macro F1 score across different labels for the different target tasks in a cross-task setting

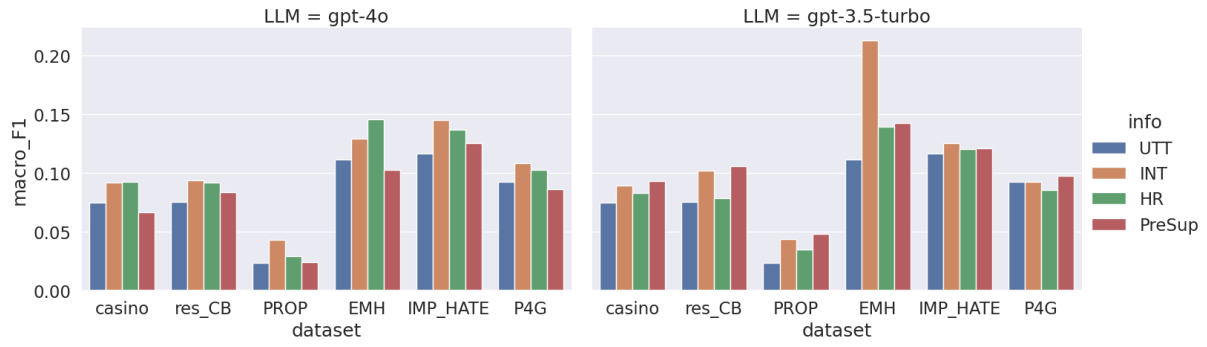


Figure 17: Zero-shot cross-task performance for instruction tuned models

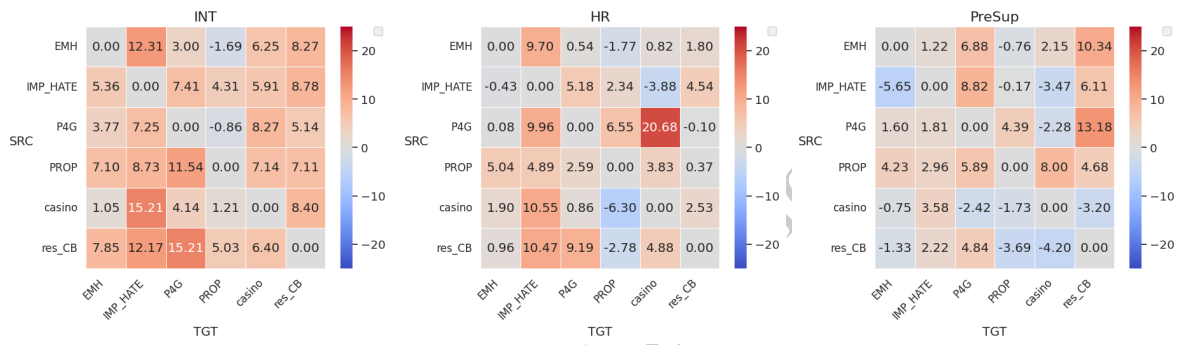
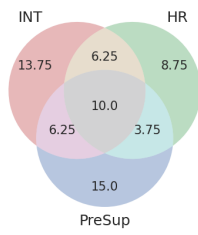
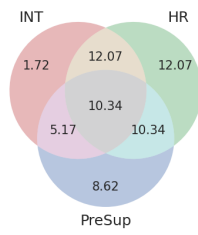


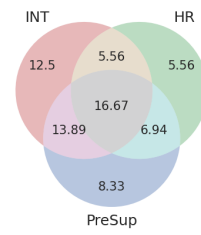
Figure 18: Relative change in performance measured in terms of F1 score over the baseline when incorporating the rationale information for different source and target pairs for the cross-task transfer setting.



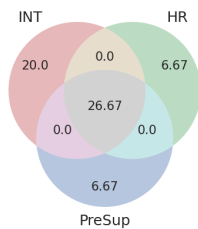
(a) CaSiNo



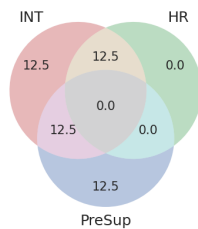
(b) P4G



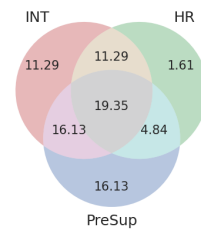
(c) res_CB



(d) EMH



(e) PROP



(f) IMP_HATE

Figure 19: Venn Diagram showing the proportion of instances where including the rationales fared better than the baseline in an in domain setting.

| Dataset | Label | Utterance text | Rationale Text | CAT |
|----------|--------------------|--|---|--------|
| casino | showing-empathy | Are you sure that's enough firewood for you and the baby? I know that babies can easily get very sick from dropping temperatures. | Expressing concern about the adequacy of firewood for the baby | INT |
| casino | vouch-fair | hey buddy I hope we both end up with a good deal:) | Expressing hope for a mutually beneficial outcome | INT |
| IMP_HATE | white_grievance | is zuma talking about seizing white land or not ? are you denying it ? | Seeking clarification and confirmation about Zuma's stance on seizing white land ; challenging the respondent to deny it. | INT |
| IMP_HATE | irony | why do we have an immigration policy? should just let them have a free for all and all can pay! | Expressing frustration and sarcasm about the current immigration policy and suggesting an extreme alternative to highlight perceived issues. | INT |
| IMP_HATE | irony | mt : ""we must resist ebolaphobia. these viruses just come here for a better life, to do jobs that american viruses refuse to do. | The speaker is using satire to comment on the fear of Ebola (ebolaphobia) by comparing it to immigration issues suggesting that the fear is irrational and mocking the idea that viruses have intentions similar to human immigrants. | INT |
| IMP_HATE | white_grievance | but that wouldn't enable them to destroy white neighbourhoods . | There is a belief or concern that certain actions or policies could lead to the destruction of white neighborhoods . | PreSup |
| IMP_HATE | threatening | and they should have to pay every dime back, plus interest. from their own damn pockets. | People should be held financially accountable for their actions and should not be able to avoid paying back debts. | PreSup |
| P4G | credibility-appeal | As a donor I wanted you to see the site and the children that you would be helping | EE is informed about the purpose of the link and the impact of their donation. | HR |
| P4G | foot-in-the-door | Every little bit help. | EE feels reassured that their small donation is still valuable. | HR |
| P4G | foot-in-the-door | Every little bit help. | Reassure the listener that any contribution is valuable. | INT |
| P4G | foot-in-the-door | Your right, but I'm not asking for much. | Minimizing the financial impact of the donation | INT |
| res_CB | Source Derogation | Too be honest don't like the front bumper would be better without that black cover at this i can only pay about 1600 could you do that | The seller might feel a need to address the buyer's concern about the bumper. | HR |
| res_CB | Self Pity | Yes. What didn't your wife like about the bed? | Seller realizes the buyer's budget constraints . | HR |
| res_CB | Source Derogation | Yes. What didn't your wife like about the bed? | Seller feels questioned about the reason for selling the bed . | HR |

Table 17: We present instances across different datasets where adding the rationale information was crucial in predicting the correct label always. We compute Shapley values for each token in the rationale to observe its contribution to the model's decision; the highlighted portions correspond to high positive associations with the label.

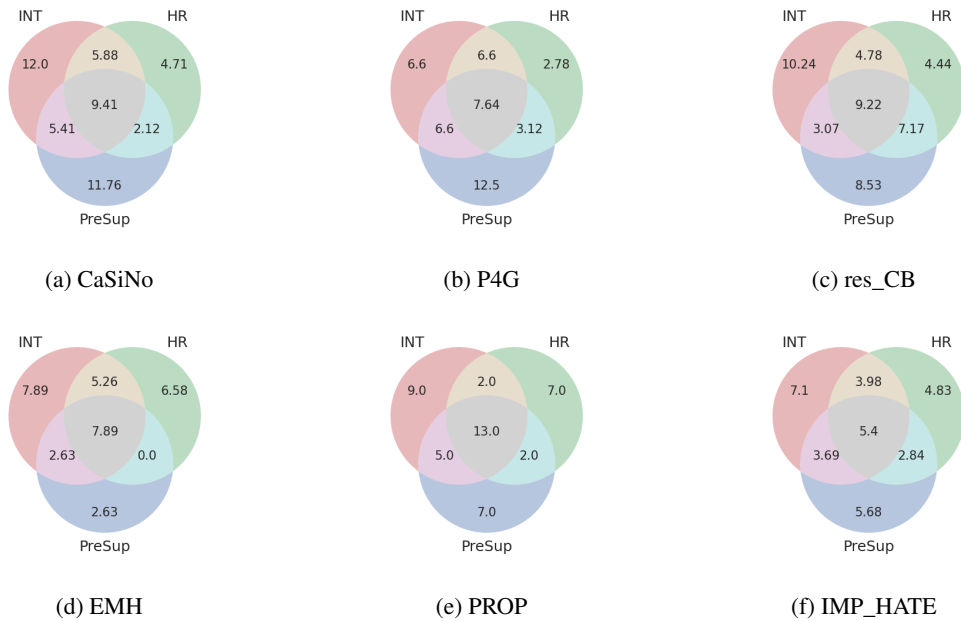


Figure 20: Venn Diagram showing the proportion of instances where including the rationales fared better than the baseline in a 5-shot transfer setting.

G Qualitative Analysis

We now carry out a qualitative analysis to investigate the specific instances where including the rationales actively improves the model’s predictions in an indomain setting.

We depict the fraction of cases that benefit from adding rationales in the form of a Venn Diagram in Figure 19 in the Appendix. The overlapping areas indicate the fraction of instances that benefit from more than one types of rationale; for example, 10.0% of all instances benefit from all three rationales in CaSiNo. We consider only those instances where the baseline (i.e., only the utterance text) fails to predict the label correctly a majority of times, but succeeds when the rationale is provided.

The rationale with the greatest impact on performance is dependent on the nature of the task. The hearer reaction or HR has the highest impact on P4G, possibly because it captures the thought processes of the persuadee (EE) as they are being persuaded to donate. For example, the utterance “Anything would help even small donations add up when everyone pitches in.” evokes a sense of reassurance from the persuadee (EE) that any contribution is valuable and is thus recognized as a “foot-in-the-door” strategy. Presuppositions are useful for IMP_HATE, a dataset that directly references stereotypes and thus requires generic knowledge to infer the type of implicit hatred. Tasks that are centered around the outcome the speaker is invested in, i.e. strategies employed to resist persuasion (res_CB), or signaling empathy to someone in therapy (EMH) benefit mostly from intentions. Furthermore, similar tasks e.g., CaSiNo and res_CB which deal with negotiation have similar relative performance for the same rationales.

However, it should also be noted that a given rationale category does not serve as a silver bullet for all instances. We highlight some examples where model improvements were due to only one type of rationale in Table 17 in the Appendix and the possible reasoning for the same. While all three rationales are valid with respect to the utterance, we hypothesize that certain phrases or terms in the given generation might make it easier to predict the label category. For example, the phrase “feels questioned” in the HR hints at source derogation, which is not observed for the other rationales for the res_CB example. Likewise, the wording “how one might treat a dog” in the presupposition conveys the sense of inferiority more prominently than

the generic idea of mistreatment in IMP_HATE. Since the rationales were not generated with a particular task in mind, the number of instances where the wording aligns with one of the task label’s definition is also infrequent.