

USDC: A Dataset of User Stance and Dogmatism in Long Conversations

Anonymous Author(s)
Affiliation
Address
email

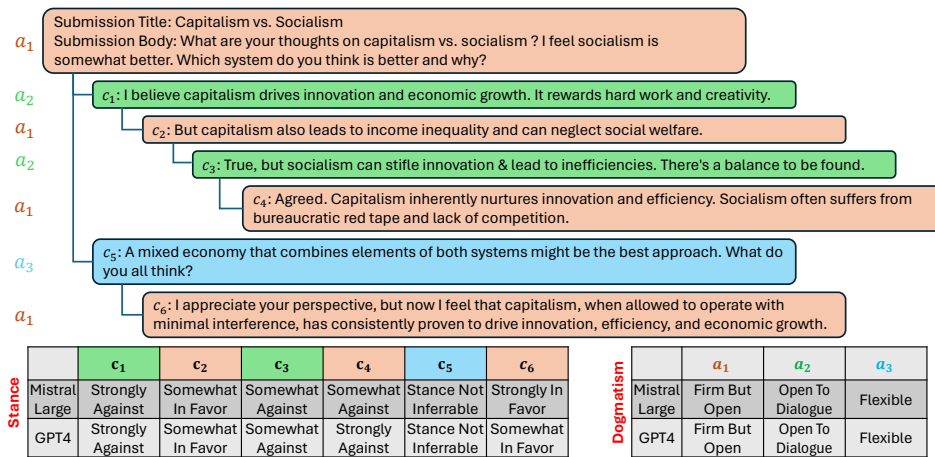


Figure 1: Sample Reddit conversation on “Capitalism vs. Socialism” with Stance (for every comment $\{c_i\}_{i=1}^6$) and Dogmatism (for every author $\{a_j\}_{j=1}^3$) labels from Mistral Large and GPT-4. The submission content favors to socialism and examines how the authors position their opinions regarding socialism vs. capitalism.

Abstract

1 Identifying user’s opinions and stances in long conversation threads on various
 2 topics can be extremely critical for enhanced personalization, market research,
 3 political campaigns, customer service, conflict resolution, targeted advertising and
 4 content moderation. Hence, training language models to automate this task is
 5 critical. However, to train such models, gathering manual annotations has multiple
 6 challenges: 1) It is time-consuming and costly; 2) Conversation threads could be
 7 very long, increasing chances of noisy annotations; and 3) Interpreting instances
 8 where a user changes their opinion within a conversation is difficult because often
 9 such transitions are subtle and not expressed explicitly. Inspired by the recent
 10 success of large language models (LLMs) for complex natural language process-
 11 ing (NLP) tasks, we leverage Mistral Large and GPT-4 to automate the human
 12 annotation process on the following two tasks while also providing reasoning: i)
 13 User Stance classification, which involves labeling a user’s Stance of a post in a
 14 conversation on a five-point scale; ii) User Dogmatism classification, which deals
 15 with labeling a user’s overall opinion in the conversation on a four-point scale. The
 16 majority voting on zero-shot, one-shot, and few-shot annotations from these two
 17 LLMs on 764 multi-user Reddit conversations helps us curate the USDC dataset.
 18 USDC is then used to finetune and instruction-tune multiple deployable small

19 language models for the 5-class stance and 4-class dogmatism classification tasks.
20 We make the code and dataset publicly available ¹.
21 .

22 1 Introduction

23 Understanding the user’s (or author’s) opinion in a conversation is a fundamental aspect of successful
24 interpersonal interactions, and it is essential for developing better interpersonal communication skills,
25 empathy development, and informed decision-making. This user understanding is particularly relevant
26 in the context of dogmatism, a phenomenon observed in various areas such as politics, religion,
27 culture, intellect, and science, where rigid adherence to beliefs often hinders open-mindedness and
28 empathy (Rokeach, 1954). Advertisers can target their campaigns more effectively by aligning
29 with the opinions and stances of potential customers. Companies can use this information for
30 market research to tailor products and services to meet consumer needs and preferences. Political
31 groups can gauge public reaction to policies and campaigns and adjust their strategies accordingly.
32 Identifying differing opinions can help conflict resolution by understanding the perspectives of all
33 parties’ perspectives. Society can promote tolerance and maintain social harmony by recognizing and
34 respecting diverse opinions.

35 Fig. 1 shows a sample Reddit conversation on the topic of *Capitalism vs. Socialism*. We refer to an
36 author’s initial post (containing title and body) as a submission. Multiple authors can then share their
37 opinions as comments on the submission. Specifically this example contains 6 comments $\{c_i\}_{i=1}^6$
38 from 3 authors $\{a_j\}_{j=1}^3$. We also show stance and dogmatism predictions from two large language
39 models (LLMs): Mistral Large and GPT-4. Some authors like a_1 change their views during the
40 discussion based on the beliefs or opinions of others. At the beginning of the dialogue, we note that
41 author a_1 is somewhat favoring socialism (in submission and c_2). But the author shifts their stance
42 to somewhat favors capitalism (in c_4) after considering the viewpoints of author a_2 in comments c_1
43 and c_3 , illustrating author a_1 ’s firm yet open-minded approach. On the other hand, author a_3 seems
44 very flexible based on their comment c_5 . Understanding conversations requires understanding the
45 fine-grained topics being discussed and the dynamic viewpoints of the individual users.

46 Given the importance of understanding these user dynamics in conversations, training language
47 models to perform this task automatically at scale is critical. While numerous datasets are available
48 for analyzing individual user posts (Fast & Horvitz, 2016; Sakketou et al., 2022), typically through
49 random subsampling of posts or selecting posts with a limited number of tokens, the exploration of a
50 specific user’s opinion across each post within an entire conversational thread remains under-explored.

51 Crowdsourcing is one possible approach to address the need for a suitable dataset. However, a
52 significant limitation in manually annotating datasets for user opinions is the time-consuming nature
53 of the process, as annotators must read entire conversations to label each user’s post, making data
54 acquisition costly. Additionally, manual annotation often faces challenges related to quality, as accu-
55 rately labeling opinions requires understanding demographic details and domain-specific knowledge.
56 Given these limitations, achieving a comprehensive and accurate set of user opinions corresponding
57 to posts about a topic often requires multiple annotators or iterative rounds of annotation. Since users
58 could change their opinion (often times with subtle transitions and not with explicit statements) within
59 a conversation, tracking such changes across multiple users manually becomes very cumbersome.

60 Recently, large language models (LLMs), especially those built on Transformer architectures (Vaswani
61 et al., 2017) and pretrained on large datasets, have resulted in state-of-the-art accuracies on several
62 complex natural language processing (NLP) tasks (Brown et al., 2020; Chung et al., 2024). LLMs
63 are also being frequently used for dialog response generation (Zhang et al., 2020; Bao et al., 2019;
64 Roller et al., 2021; Adiwardana et al., 2020). Given the complex and cumbersome nature of con-
65 versation understanding, we hypothesize that LLMs can be effective in capturing nuances involved
66 in understanding user opinions and their shifts in multi-user conversational contexts. Also, since
67 these models possess long-range memory capabilities, we believe that they can reason over extended
68 conversational threads involving numerous participants, as good as human annotators, if not better.

69 In this work, we leverage LLMs like Mistral Large and GPT-4 to perform two tasks: i) User Stance
70 classification, which involves labeling a user’s stance of a post in a conversation on a five-point

¹<https://anonymous.4open.science/r/USDC-0F7F>

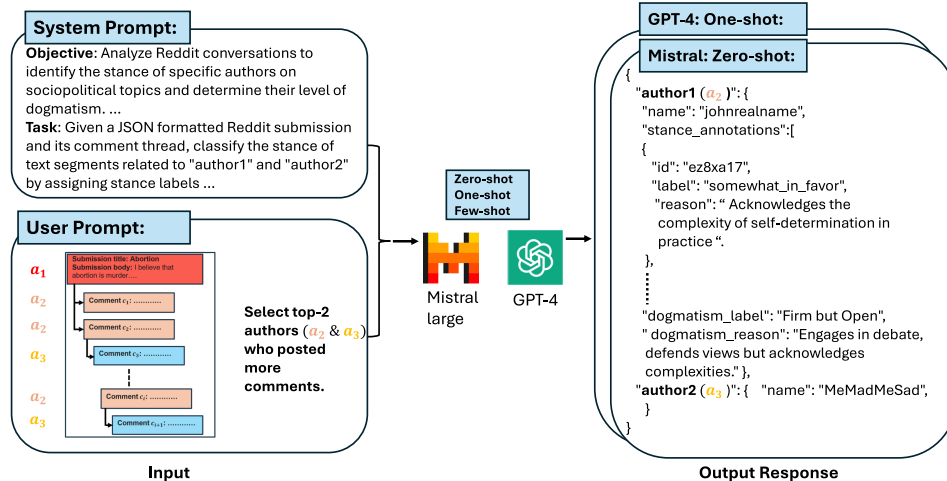


Figure 2: Generating annotations using LLMs: We pass the entire conversation for each Reddit thread in JSON format. The JSON highlights the top two authors who posted the most comments, alongside annotation guidelines for stance and dogmatism labels in the system prompt.

71 scale; ii) User Dogmatism classification, which deals with labeling a user’s overall opinion in
 72 the conversation on a four-point scale. Besides the class labels, we also obtain reasoning behind
 73 these labels from these LLMs. We experiment with these two models as human-like annotators to
 74 generate user opinions in full-length, multi-user Reddit conversations in a zero-shot, one-shot as
 75 well as few-shot setup. Thus, overall for every sample, we obtain six annotations ($\{\text{Mistral Large,}$
 76 $\text{GPT-4}\} \times \{\text{zero-shot, one-shot, few-shot}\}$). Fig. 2 presents our LLM-based annotation pipeline for
 77 user-level Stance and Dogmatism tasks. We consider majority voting over these six as our final
 78 annotations. Overall, this helps us curate our USDC (a dataset of user stance and dogmatism in
 79 conversations) dataset, which consists of 764 multi-user conversations from 22 subreddits, including
 80 1,528 user-level dogmatism samples and 9,618 stance samples across all posts from selected users.
 81 Overall, the annotations in the dataset highlight specific user opinions in each post related to stance,
 82 track opinion fluctuations leading to a dogmatic nature, and provide reasoning about why users hold
 83 specific opinions.

84 USDC addresses several weaknesses of existing post-level stance and dogmatism datasets. First, the
 85 full-length multi-user conversation aspect of USDC enables it to capture contextual and opinion shifts
 86 of multiple users. This feature allows it to serve as both an instruction-tuning user opinion dataset and
 87 an evaluation benchmark. We believe that the ability to perform instruction tuning for user opinions
 88 at a large scale can bridge the gap between open-source and commercial user trait understanding
 89 models. Additionally, the in-context learning annotations using state-of-the-art LLMs in USDC make
 90 it a more comprehensive measure of how current LLMs understand complex tasks like capturing
 91 opinions. This aspect makes it a valuable resource, especially for social media agents seeking deeper
 92 insights into user behavior.

93 In this work, we utilize our USDC dataset to finetune as well as instruction-tune open-source LLMs
 94 for generating stance and dogmatism labels for users. We experiment with three pretrained small
 95 language models (SLMs) like LLaMA-2-7B, LLaMA-3-8B, and Falcon-7B. We also experiment with
 96 four instruction-tuned SLMs like LLaMA-2-chat-7B, LLaMA-3-8B-instruct, Vicuna-7B-v.1.5, and
 97 Falcon-7B-instruct. We report weighted F1 scores obtained using these models for both the tasks.

98 We make the following contributions: 1) We contribute USDC (a dataset of user stance and dogmatism
 99 in conversations) dataset consisting of 764 multi-user conversations labeled with 1,528 user-level
 100 dogmatism samples and 9,618 stance samples. 2) We report initial results for the stance and
 101 dogmatism detection tasks using seven small language models for the UDSC dataset. We find that
 102 stance detection performance improves with instruction-tuning (F1-score of 56.2) compared to fine-
 103 tuning (F1-score of 54.9). However, dogmatism detection performs worse with instruction-tuning
 104 (F1-score of 49.2) compared to fine-tuning (F1-score of 51.4), highlighting the complexity of this
 105 task. 3) We make the code and dataset publicly available¹. Also, the finetuned and instruction-tuned
 106 models are made available as well.

107 2 Related Work

108 **Opinion fluctuations in user conversations.** Our work is closely related to previous studies (Fast
109 & Horvitz, 2016; Sakketou et al., 2022), which explore Stance and Dogmatism at the post level,
110 where posts are randomly sampled from conversation threads. Fast & Horvitz (2016) predicted user
111 dogmatism on randomly sampled Reddit posts from conversations, with each post limited to 200-300
112 characters. One major limitation of this work is the unavailability of a public dataset and missing
113 annotator demographic details. Sakketou et al. (2022) created the post-level Stance dataset, SPINOS,
114 where each post is considered independent, and submission posts are missing while annotators label
115 the data. Additionally, the quality of the dataset is not validated due to missing demographic details
116 of these annotators. Our work overcomes the limitations of previous studies and presents Stance
117 detection for posts and Dogmatism labels of users in conversations, considering the entire context,
118 while preserving submission IDs. Hence, our dataset provides clear user-level posts and dogmatism
119 data, which are useful for modeling dynamic user representations.

120 **Generating annotations for NLP tasks using Large Language Models** Our work also relates to a
121 growing body of literature suggesting that large language models can perform similarly to human
122 annotators in labeling complex NLP tasks (Zhou et al., 2022; Zhang et al., 2023; Bansal & Sharma,
123 2023; Lowmanstone et al., 2023; Wadhwa et al., 2023; Honovich et al., 2023; Zheng et al., 2024; Ye
124 et al., 2022a; Meng et al., 2022). Several studies have explored LLM-based annotation generation
125 in zero-shot or few-shot task settings (Ye et al., 2022a; Meng et al., 2022; Ye et al., 2022b), while
126 others have compared pairs of language models to assess the quality of annotations generated by
127 these LLMs (Zheng et al., 2024). However, these studies focused on generating annotations for NLP
128 tasks such as sentiment analysis, natural language inference (Gilardi et al., 2023; Alizadeh et al.,
129 2023), or creating synthetic dialogues, but only for dyadic conversations (Lee et al., 2023). Our
130 approach complements these previous studies by focusing on generating annotations of user opinions
131 in complex multi-user conversations.

132 3 USDC Dataset Curation

133 In this section, we will discuss three main things: 1) Collection of Reddit conversations, 2) Obtaining
134 LLM annotations, and 3) Inter-annotator agreement with LLMs as annotators.

135 3.1 Collection of Reddit Conversation Threads

136 **Initial crawl.** We crawl a year (2022) worth of multi-user conversation data from 22 subreddits of
137 Reddit² using praw API³. This dataset includes submissions and all associated user comments. Each
138 submission, which serves as the initial message of the conversation, contains a title and content body.
139 This is followed by comments and replies to the submission or other comments. Overall, we crawled
140 3,619 Reddit conversations across the 22 subreddits. A sample Reddit conversation is displayed in
141 Fig. 1.

142 **Quality filtering of conversations.** Since submission content on Reddit can sometimes include
143 videos, we perform the following filtering steps. 1) We only consider submissions where the content is
144 text. 2) We remove conversations with [deleted] tags and empty content. 3) We exclude conversations
145 where the posts were either discarded by users or removed by moderators.

146 Reddit user conversations can be very long and we observed up to 591 comments in a single crawled
147 conversation data. Considering the maximum sequence length allowed by various language models,
148 we retained only those conversations that contain at least 20 and at most 70 comments. Considering
149 conversations with fewer than 20 comments results in too few comments to accurately gauge user
150 opinions based on small samples. Further, we ensure that at least two users covering ~50% of the
151 comments in the conversations. We did not remove any comments or reduce the post length in the
152 selected conversations. Out of the initial 3,619 conversations, these filtering steps result into 764
153 conversations getting selected. Table. 4 in the Appendix shows detailed subreddit level statistics.

²<https://www.reddit.com/>

³<https://github.com/praw-dev/praw>

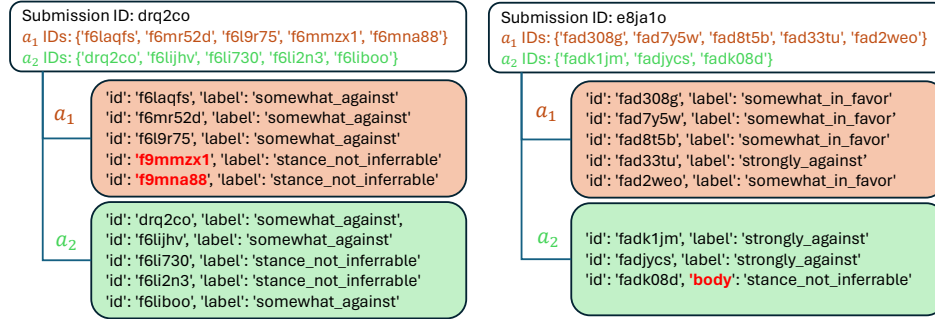


Figure 3: Failure cases of LLMs: Mistral Large few-shot output (left), the ids (“f6mmzx1”, “f6mna88”) were mismatched with generated ids (“f9mmzx1”, “f9mna88”), GPT-4 zero-shot output (right), the key “label” was mismatched with generated key “body”.

154 3.2 Obtaining LLM Annotations

155 Representing Reddit conversations in JSON format.

156 To create the prompt, we follow the nested hierarchical structure of Reddit conversations to maintain
 157 the context. Specifically, we maintain a JSON structure for each conversation, where each author has
 158 their post IDs, and comments or replies are available in the body section. An example of a Reddit
 159 conversation in JSON format is provided in Appendix D. Note that the JSON explicitly includes the
 160 top-2 authors who posted the most comments in the conversation, as well as their respective post IDs.
 161 Our emphasis on these top-2 users (covering 47% posts of total posts on average) aimed at accurately
 162 assigning Stance and Dogmatism labels, acknowledging the challenge of modeling a user’s opinion
 163 belief based on a very number of posts within a conversation.

164 **Using LLMs as human-like annotators.** To annotate the position (or Stance) of a user towards a
 165 subreddit topic at each post and the overall opinion (or Dogmatism level) of a user in a conversation,
 166 we employ two well-known commercialized API-based LLMs: GPT-4 (OpenAI, 2023) and Mistral
 167 Large (Jiang et al., 2024). OpenAI GPT-4 is a decoder-based language model which features a context
 168 window of 32k to 128k tokens. Mistral Large features a context window of 32k tokens. Additionally,
 169 we also examined other versions of these models, such as GPT-3.5 and Mistral-small and medium,
 170 but found that these models failed to produce annotations in the desired format. We briefly discuss
 171 these limitations in Section 6.

172 For both GPT-4 and Mistral Large, we supplied a system prompt that contains the definition of Stance
 173 and Dogmatism, guidelines for annotating each user conversation, and the necessary labels for Stance
 174 and Dogmatism, as shown in Fig 2. The system prompt is detailed in the Appendix B. Along with the
 175 system prompt, we provided a user prompt comprising the entire user conversation in a structured
 176 JSON format, as discussed above. Additionally, we prompted the model to generate reasoning for
 177 each label, explaining why the LLMs assigned a particular label to a specific user post. We used
 178 zero-shot, one-shot, and few-shot settings to get the LLM-based annotations. For the few-shot setting,
 179 we added two examples in the prompt. Samples of generated outputs using GPT-4 in zero-shot,
 180 one-shot, and few-shot settings are shown in Appendix E.1, E.2, E.3 respectively. Similarly, samples
 181 of generated outputs using Mistral Large in zero-shot, one-shot, and few-shot settings are shown in
 182 Appendix E.4, E.5, E.6 respectively.

183 **Annotation tasks.** We prompt the LLMs to perform two annotation tasks: 1) Stance detection, which
 184 determines if a user comment or post is *Strongly In Favor*, *Strongly Against*, *Stance Not Inferrable*,
 185 *Somewhat In Favor*, or *Somewhat Against* towards specific subreddit submission content; 2) Dog-
 186 matism identification, which evaluates the user’s overall opinion in conversation and categorizes
 187 them into one of four categories: *Firm but Open*, *Open to Dialogue*, *Flexible* or *Deeply Rooted*.
 188 This assessment reveals whether a user is open to changing their beliefs or remains steadfast in their
 189 opinions based on interactions with other users.

190 **Addressing LLM response and JSON parsing failures.** Sometimes the LLMs got confused with
 191 the author IDs and missed Stance labels for some author IDs (Fig. 3 (left)). Sometimes, there were
 192 minor errors in key naming (‘label’ vs ‘body’ in Fig. 3 (right)). For each LLM setting, we observed

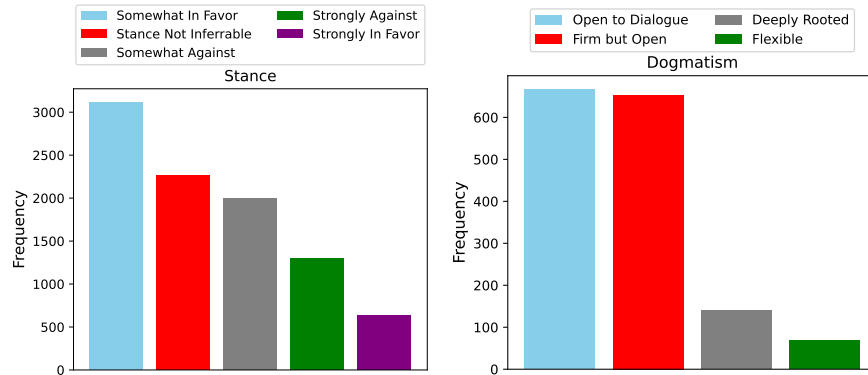


Figure 4: Distribution of class labels for Stance (left) and Dogmatism (right) tasks. These class labels are determined by majority voting across GPT-4 and Mistral Large models.

193 such errors in around 15 cases on average. We manually fixed such JSON parse errors and missing
 194 Stance labels for some author IDs.

195 **Majority voting.** After obtaining six annotations ($\{\text{Mistral Large, GPT-4}\} \times \{\text{zero-shot, one-shot,}$
 196 $\text{few-shot}\}$) for each sample, we aggregate using majority voting to determine the final gold annotations
 197 for the Stance and Dogmatism tasks. Fig. 4 presents the class distributions for both the annotation
 198 tasks. Additionally, we present the class distributions obtained from each model with the three
 199 settings (zero-shot, one-shot and few-shot) for two tasks in Appendix Figs. 5 and 6 respectively.

200 3.3 Inter-annotator Agreement with LLMs as Annotators

201 As the quality of labeling on subjective tasks is challenging, we validated the inter-annotator agree-
 202 ment (IAA) between the six LLMs (GPT-4 Zero-shot, GPT-4 One-shot, GPT-4 Few-shot, Mistral
 203 Large Zero-shot, Mistral Large One-shot, and Mistral Large Few-shot) for the Stance as well as
 204 Dogmatism tasks. We perform IAA using two approaches: i) Cohen’s kappa score (Cohen, 1960)
 205 and ii) Fleiss’ kappa score (Fleiss, 1971). Cohen’s kappa measures the agreement between two raters,
 206 while Fleiss’ kappa extends this to multiple raters. Hence, we employed Cohen’s kappa for pairwise
 207 comparisons and Fleiss’ kappa for overall agreement across all models.

208 Fig. 7 in Appendix shows the pairwise Cohen’s kappa values for both Stance and Dogmatism tasks.
 209 We observe that Cohen’s kappa values range from 0.36 to 0.72 for Stance and 0.31 to 0.61 for
 210 dogmatism, indicating moderate agreement between the models. Broadly kappa values are higher for
 211 model pairs within a family (GPT-4 or Mistral large). Thus, the large variance in the kappa scores
 212 is not due to the various in-context learning settings (ZS, OS, FS) but rather due to architectural
 213 differences.

214 The overall Fleiss’ kappa value was calculated as 0.485 for Stance and 0.435 for Dogmatism,
 215 suggesting moderate agreement among all six models. Comparing LLM IAA with previous studies,
 216 we observe that for dogmatism, the LLM IAA of 0.435 matches with 0.44 as mentioned in Fast &
 217 Horvitz (2016). Similarly, for Stance, the LLM IAA of 0.485 is much higher than 0.34 as reported
 218 in Sakketou et al. (2022). It is important to note that previous studies on Stance and Dogmatism
 219 datasets were created on post-level data with limited token lengths, whereas our work focuses on
 220 entire user conversations. This suggests that LLMs can be considered as competent annotators
 221 for complex subjective tasks. However, the moderate agreement levels indicate potential areas for
 222 improvement and align with the observed performance variations among the models.

223 4 Training Small Language Models

224 In this section, we briefly discuss the small language models that we experiment with. We also
 225 discuss their finetuning and instruction tuning details.

226 4.1 Small Language Models

227 we train three pretrained small language models (LLaMA-2-7B, LLaMA-3-8B, Falcon-7B) and
228 four instruction-tuned small language models (LLaMA-2-chat-7B, LLaMA-3-8B-instruct, Vicuna-
229 7B-v.1.5, and Falcon-7B-instruct). We finetune as well as instruction tune these models using the
230 proposed USDC dataset. We use pretrained models checkpoints from Hugging Face. All of these
231 LLMs have context length of 4096 tokens.

232 **LLaMA** models (Touvron et al., 2023a) are decoder-only LLMs trained on 1.6 trillion tokens from a
233 mixture of corpora including C4, English CommonCrawl, Wikipedia, Github, and more. We use two
234 versions of models in our study: LLaMa-2-7B (Touvron et al., 2023b) and LLaMa-3-8B and their
235 instruction tuned variants.

236 **Falcon** models (Almazrouei et al., 2023) are decoder-only LLMs trained on ≥ 1 trillion tokens of
237 text, with a particular emphasis on the RefinedWeb corpus. For Falcon, we use both the pretrained
238 and instruction tuned 7B parameter variants in our study.

239 **Vicuna** model (Chiang et al., 2023) is finetuned from the LLaMA 7B model on approximately 70K
240 user-shared conversations gathered from ShareGPT.com and we used the 7B parameter variants.

241 4.2 Experimental Setup

242 **Train-test setup.** We conducted both finetuning and instruction-tuning of small language models. For
243 this purpose, we divided the dataset of 764 conversations into train ($\sim 75\%$) and test splits ($\sim 25\%$).
244 The training dataset comprised 564 conversations, including 1128 samples of Dogmatism labels and
245 7520 samples of Stance labels. Conversely, the testing dataset consisted of 200 conversations, with
246 400 samples of Dogmatism labels and 1831 samples of Stance labels across two authors posts.

247 **Implementation details for reproducibility.** All experiments were conducted on a machine equipped
248 with an NVIDIA A100 GPU with 80 GB of GPU RAM, partitioned into two devices of 40 GB
249 each. We employed 4-bit quantization with normalized floating precision (nf4) from the bitsandbytes
250 library⁴. Additionally, we utilized LoRA (Hu et al., 2021) with a rank of 64 and an alpha value of
251 16 during task-based instruction tuning. Finally, we use PEFT (Parameter Efficient Finetuning)⁵
252 library to train large language models with SFTT (Supervised Finetuning Trainer) setting. To further
253 enhance performance, we divided the training dataset into a validation set, comprising a randomly
254 chosen 10% subset from the training set, used exclusively for hyperparameter tuning. More details
255 about bitsandbytes, PEFT and SFTT parameters are reported in Appendix.

256 4.3 Finetuning and Instruction Tuning of Small Language Models (SLMs)

257 **Finetuning of SLMs.** For Stance classification, we treat each user post as an independent sample. In
258 contrast, for Dogmatism classification, we consider the entire user conversation as a single sample
259 by concatenating all the threads from a user in that conversation. To load the pretrained SLMs, we
260 perform 4-bit quantization, apply the LoRA technique (Hu et al., 2021), and fine-tune the models with
261 SFTT before saving the fine-tuned model. For finetuning, we used prompt for Stance classification as
262 shown in Fig. 8 (see Appendix). Similarly, Fig. 9 (see Appendix) displays prompt for Dogmatism
263 identification.

264 **Instruction tuning of SLMs.** We instruction tune the SLMs on user conversations along with their
265 gold labels from the training part of the USDC dataset. For instruction tuning, we use the same
266 prompt as used for LLMs to generate the USDC dataset (also shown in Appendix B). Similar to
267 finetuning, we use same train-test splits for instruction tuning.

268 5 Results

269 **Do SLMs finetuned with task-specific LLM annotations accurately perform Stance and Dogma-**
270 **tism tasks on user opinions?** We show the weighted F1 of various SLMs finetuned with task-specific
271 LLM annotations on the stance and dogmatism detection tasks on the USDC test set in Table 1. We

⁴<https://pypi.org/project/bitsandbytes/>

⁵<https://github.com/huggingface/peft>

Table 1: Finetuning: weighted F1 score for Stance classification using SLMs on USDC test set. ZS: Zero-shot, OS: One-shot, FS: Few-shot.

Model	Stance Classification						Dogmatism Classification							
	GPT-4			Mistral Large			Majority	GPT-4			Mistral Large			Majority
	ZS	OS	FS	ZS	OS	FS		ZS	OS	FS	ZS	OS	FS	
LLaMA-2-7B	51.8	52.9	52.7	35.1	49.2	46.0	54.0	42.1	44.2	45.2	39.3	47.6	43.7	43.4
LLaMA-2-chat-7B	52.8	51.4	51.8	34.7	47.5	46.5	51.3	42.1	42.5	48.8	41.1	49.7	45.5	48.3
LLaMA-3-8B	51.3	52.2	52.9	34.9	48.5	47.0	54.9	42.0	47.8	45.3	39.9	47.4	36.3	51.4
LLaMA-3-8B-instruct	51.2	52.6	52.7	33.9	49.5	45.6	54.5	44.8	46.2	49.7	46.1	45.8	46.1	50.8
Falcon-7B	50.7	51.1	51.6	34.9	47.2	43.9	53.2	41.5	42.1	43.3	36.5	38.4	37.5	40.1
Falcon-7B-instruct	51.2	51.5	51.6	35.1	47.7	44.2	51.0	41.7	42.1	42.9	36.8	38.5	36.9	39.7
Vicuna-7B-v.1.5	51.0	53.0	53.2	35.1	48.5	45.8	54.7	42.9	48.3	40.8	45.9	42.6	46.2	42.3

Table 2: Instruction-tuning: weighted F1 score for Stance classification using SLMs on USDC test set. ZS: Zero-shot, OS: One-shot, FS: Few-shot.

Model	Stance Classification						Dogmatism Classification							
	GPT-4			Mistral Large			Majority	GPT-4			Mistral Large			Majority
	ZS	OS	FS	ZS	OS	FS		ZS	OS	FS	ZS	OS	FS	
LLaMA-2-7B	53.2	54.0	54.5	36.8	50.3	47.2	55.5	43.0	45.0	46.3	40.6	48.2	45.0	44.0
LLaMA-2-chat-7B	54.0	54.5	55.0	36.5	50.7	47.6	54.0	43.2	45.5	47.0	40.8	48.5	45.5	43.8
LLaMA-3-8B	53.5	54.8	55.5	37.0	50.5	48.0	56.2	43.5	46.0	47.5	41.0	48.8	45.8	45.1
LLaMA-3-8B-instruct	53.0	54.2	55.0	36.0	50.0	47.0	55.5	43.8	46.5	47.8	41.5	49.2	46.0	44.8
Falcon-7B	52.8	53.4	54.0	36.5	49.5	46.5	54.8	42.5	44.6	45.8	39.8	47.0	44.0	43.8
Falcon-7B-instruct	53.0	53.8	54.2	36.8	49.8	46.8	54.5	42.8	44.8	46.0	40.0	47.2	44.2	43.0
Vicuna-7B-v.1.5	53.3	54.5	55.2	37.0	50.2	47.8	55.2	43.7	46.8	47.2	41.2	48.2	46.5	44.8

272 report AUC scores and other qualitative analysis in Appendix F (Fig. 11 and 12). We make the
 273 following observations from these results: 1) For both tasks, the majority voting labels as ground
 274 truth, has a relatively high performance, scoring above 50% weighted F1-score across several models.
 275 2) LLaMa-3 models (LLaMA-3-8B and LLaMA-3-8B-instruct) perform better across both the tasks.
 276 3) For GPT-4 annotations, in most cases, SLMs finetuned with few-shot annotations outperform
 277 those trained with zero and one-shot annotations. For Mistral Large annotations, typically SLMs
 278 finetuned with one-shot annotations performs the best. 4) Specifically, for Stance detection task,
 279 Vicuna-7B-v.1.5 finetuned using few-shot annotations is the best model trained with GPT-4 anno-
 280 tations. Similarly, LLaMA-3-8B-instruct finetuned with one-shot annotations is the best model
 281 trained with Mistral Large annotations. 5) For the Dogmatism detection task, LLaMA-3-8B-instruct
 282 finetuned using few-shot annotations is the best model trained with GPT-4 annotations. Similarly,
 283 LLaMA-2-chat-7B finetuned with one-shot annotations is the best model trained with Mistral Large
 284 annotations. 6) Overall, we observe that instruction tuned SLMs perform better than the pretrained
 285 SLMs.

286 **Do SLMs instruction-tuned with task-specific LLM annotations perform better than SLMs**
 287 **finetuned with task-specific LLM annotations for the Stance and Dogmatism tasks?** We show
 288 the weighted F1 of various SLMs instruction-tuned with task-specific LLM annotations, on the
 289 stance and dogmatism detection tasks on the USDC test set in Table 2. We report AUC scores and
 290 other qualitative analysis in Appendix F (see Fig. 13). We make the following observations from
 291 these results: 1) SLMs with instruction-tuning result in higher weighted F1-scores than SLMs with
 292 finetuning for stance detection, while SLMs with finetuning outperform SLMs with instruction-tuning
 293 in dogmatism detection. 2) Contrary to results in Table 1, Table 2 demonstrates that using majority
 294 voting labels as ground truth, SLM instruction-tuning yields relatively high performance only for the
 295 stance detection task, but not for the dogmatism detection. 3) Similar to results in Table 1, LLaMA-3
 296 models (LLaMA-3-8B and LLaMA-3-8B-instruct) perform better across both tasks. Additionally,
 297 GPT-4 annotations yield the best results in the few-shot setting, while Mistral Large annotations
 298 perform best in the one-shot setting.

299 Overall, we draw the following conclusions when comparing SLM finetuning and instruction-tuning:
 300 (1) Since dogmatism detection is inherently a more complex and varied task than stance detection,
 301 the model might struggle to generalize from the instructional data. (2) The system prompt used
 302 in finetuning is much simpler than the original system prompt for instruction-tuning, making it
 303 challenging to handle the context length for longer conversations. We perform an error analysis to
 304 further analyze the results in the next subsection.

305 **Error Analysis** Table 3 illustrates the confusion matrix for stance detection for LLaMa-3-8B
 306 finetuning and instruction-tuning. We make the following observations this table: 1) For both
 307 finetuning and instruction-tuning, there is a significant misclassification between “Somewhat Against”
 308 and “Somewhat In Favor,” as well as between “Somewhat In Favor” and “Stance Not Inferred.”
 309 These overlaps suggest challenges in distinguishing moderate stances, indicating a need for enhanced

		Predicted				
		SOA	SOIF	SNI	SGA	SIF
Actual	SOA	151	132	34	44	2
	SOIF	93	537	113	17	14
	SNI	23	78	259	5	0
	SGA	52	35	13	115	17
	SIF	18	50	12	25	27

		Predicted				
		SOA	SOIF	SNI	SGA	SIF
Actual	SOA	143	125	37	54	4
	SOIF	82	543	106	27	16
	SNI	22	82	253	6	2
	SGA	41	35	11	131	14
	SIF	16	53	10	23	30

Table 3: Confusion matrix for LLaMa-3-8B Stance detection models on USDC test set: finetuning (left) and instruction-tuning (right). SOA: Somewhat Against, SOIF: Somewhat In Favor, SNI: Stance Not Inferable, SGA: Strongly Against, SIF: Strongly In Favor.

310 feature representation and clearer class definitions to improve model performance. We report the
311 confusion matrix for dogmatism detection task in Appendix Fig. 10. Fig. 10 shows significant
312 misclassifications, especially for the “Deeply Rooted” and “Flexible” labels, with both having zero
313 accuracy and F1-scores. On the other hand, the model performs moderately better for “Firm but Open”
314 and “Open to Dialogue” classes with accuracies of 48.7% and 64.4% respectively. The confusion
315 matrix also indicates substantial confusion to distinguish between intermediate levels of dogmatism,
316 such as “Firm but Open” and “Open to Dialogue”. The area under the ROC curve (AUC) is a measure
317 of the model’s ability to distinguish between classes. Hence, we further report the ROC curve which
318 shows the trade-off between the true positive rate (TPR) and false positive rate (FPR) for each class
319 for stance and dogmatism tasks, see Figs. 11 and. 12 in Appendix F.

320 **Verification using Human Interaction.** Due to the time-consuming nature of the manual annotation
321 process, we perform human annotations on the set of 200 test conversations. In the forms for human
322 annotations, we displayed the top 2 author’s Reddit posts from the conversation, along with the
323 submission title and content. We also provided a link to the original Reddit URL for annotators
324 to look at the full conversation. We provided detailed annotation guidelines (similar to the ones
325 mentioned in the prompt in Appendix B) to instruct human annotators in carrying out these tasks.
326 Here is a sample Google form⁶. With three human annotators on a sample of 10 conversations, the
327 agreement of majority labels (i.e., USDC test set labels) with human labels is 0.56 for the stance
328 detection task and 0.45 for the dogmatism task. The annotators included two males and one female,
329 affiliated with both academia and industry, aged between 20 and 40, and were very familiar with
330 Reddit topics.

331 6 Conclusion

332 In this paper, we focused on the problems of 5-class stance and 4-class dogmatism classification in
333 long conversations. Using LLMs as human-like annotators, we introduced USDC, a large-scale dataset
334 of user stance and dogmatism in conversations. This is achieved by providing detailed annotation
335 guidelines in the system prompt and full-length conversation as user prompt. Commercialized API-
336 based LLMs generate author-level stance and dogmatism labels via zero, one and few-shot settings.
337 The full-length multi-user conversation aspect of USDC allows it to capture the contextual and
338 opinion shifts of multiple users in a conversation. We believe that the ability to perform finetuning
339 or instruction tuning SLMs for user opinions at a large scale can bridge the gap between SLMs and
340 commercial LLMs for understanding user traits. While finetuning SLMs shows F1-score on both
341 stance and dogmatism tasks, the F1-score remains below 60% (54.9% for Stance and 51.4% for
342 Dogmatism). On the other hand, instruction tuning of SLMs only improves F1-score performance
343 on stance, not the dogmatism task. Further, the performance still falls short of 60%, with weighted
344 F1-scores of 56.2% for stance and 49.2% for dogmatism. These findings indicate that there is still
345 significant room for improvement in understanding user opinions from a text segment.

346 **Limitations.** We plan to extend this work along the following directions in the future. 1) We
347 performed this work on English conversations only. It would be nice to extend this to multi-lingual
348 conversations and verify how accurately SLMs and LLMs perform on the Stance and Dogmatism
349 tasks in the multi-lingual scenario. 2) We analyzed user dogmatism based on their posts within a
350 single conversation. This approach could be extended to include posts across multiple conversations
351 and utilize similar profile information if available. 3) We analyzed dogmatism information for only the
352 top two authors. This was mainly because considering more authors increases the output generation
353 length, and we were constrained by our budget. This implies that our current models have not been
354 evaluated for authors who do not post frequently.

⁶<https://forms.gle/dbPQBsnYfNjvUeR9>

355 References

- 356 Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan,
357 Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. Towards a human-like open-
358 domain chatbot. *arXiv preprint arXiv:2001.09977*, 2020.
- 359 Meysam Alizadeh, Maël Kubli, Zeynab Samei, Shirin Dehghani, Juan Diego Bermeo, Maria Ko-
360 robeynikova, and Fabrizio Gilardi. Open-source large language models outperform crowd workers
361 and approach chatgpt in text-annotation tasks. *arXiv preprint arXiv:2307.02179*, 2023.
- 362 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojo-
363 caru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al.
364 The falcon series of open language models. *arXiv preprint arXiv:2311.16867*, 2023.
- 365 Parikshit Bansal and Amit Sharma. Large language models as annotators: Enhancing generalization
366 of nlp models at minimal cost. *arXiv preprint arXiv:2306.15766*, 2023.
- 367 Siqi Bao, Huang He, Fan Wang, Hua Wu, and Haifeng Wang. Plato: Pre-trained dialogue generation
368 model with discrete latent variable. *arXiv preprint arXiv:1910.07931*, 2019.
- 369 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
370 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
371 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 372 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
373 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna:
374 An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
375
- 376 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
377 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language
378 models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- 379 Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological*
380 *measurement*, 20(1):37–46, 1960.
- 381 Ethan Fast and Eric Horvitz. Identifying dogmatism in social media: Signals and models. In
382 *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp.
383 690–699, 2016.
- 384 Joseph L Fleiss. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76
385 (5):378, 1971.
- 386 Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. Chatgpt outperforms crowd workers for
387 text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120,
388 2023.
- 389 Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning
390 language models with (almost) no human labor. In *Proceedings of the 61st Annual Meeting of the*
391 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14409–14428, 2023.
- 392 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
393 et al. Lora: Low-rank adaptation of large language models. In *International Conference on*
394 *Learning Representations*, 2021.
- 395 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
396 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
397 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 398 Dong-Ho Lee, Jay Pujara, Mohit Sewak, Ryen White, and Sujay Jauhar. Making large language
399 models better data creators. In *Proceedings of the 2023 Conference on Empirical Methods in*
400 *Natural Language Processing*, pp. 15349–15360, 2023.

- 401 London Lowmanstone, Ruyuan Wan, Risako Owan, Jaehyung Kim, and Dongyeop Kang. Annotation
402 imputation to individualize predictions: Initial studies on distribution dynamics and model
403 predictions. *arXiv preprint arXiv:2305.15070*, 2023.
- 404 Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han. Generating training data with language models:
405 Towards zero-shot language understanding. *Advances in Neural Information Processing Systems*,
406 35:462–477, 2022.
- 407 R OpenAI. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2(5), 2023.
- 408 Milton Rokeach. The nature and meaning of dogmatism. *Psychological Review*, 61(3), 1954.
- 409 Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle
410 Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. Recipes for building an open-domain
411 chatbot. In *Proceedings of the 16th Conference of the European Chapter of the Association for Com-
412 putational Linguistics: Main Volume*, pp. 300–325, Online, 2021. Association for Computational
413 Linguistics.
- 414 Flora Sakkettou, Allison Lahnala, Liane Vogel, and Lucie Flek. Investigating user radicaliza-
415 tion: A novel dataset for identifying fine-grained temporal shifts in opinion. *arXiv preprint
416 arXiv:2204.10190*, 2022.
- 417 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
418 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
419 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 420 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
421 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
422 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 423 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
424 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing
425 systems*, 30, 2017.
- 426 Manya Wadhwa, Jifan Chen, Junyi Jessy Li, and Greg Durrett. Using natural language explanations
427 to rescale human judgments. *arXiv preprint arXiv:2305.14770*, 2023.
- 428 Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng
429 Kong. Zerogen: Efficient zero-shot learning via dataset generation. In *Proceedings of the 2022
430 Conference on Empirical Methods in Natural Language Processing*, pp. 11653–11669, 2022a.
- 431 Jiacheng Ye, Jiahui Gao, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. Progen:
432 Progressive zero-shot dataset generation via in-context feedback. In *Findings of the Association
433 for Computational Linguistics: EMNLP 2022*, pp. 3671–3683, 2022b.
- 434 Ruoyu Zhang, Yanzeng Li, Yongliang Ma, Ming Zhou, and Lei Zou. Llm-aaa: Making large language
435 models as active annotators. *arXiv preprint arXiv:2310.19596*, 2023.
- 436 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng
437 Gao, Jingjing Liu, and Bill Dolan. DIALOGPT : Large-scale generative pre-training for con-
438 versational response generation. In *Proceedings of the 58th Annual Meeting of the Associa-
439 tion for Computational Linguistics: System Demonstrations*, pp. 270–278, Online, 2020. As-
440 sociation for Computational Linguistics. doi: 10.18653/v1/2020.acl-demos.30. URL <https://aclanthology.org/2020.acl-demos.30>.
- 442 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
443 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
444 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.
- 445 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan,
446 and Jimmy Ba. Large language models are human-level prompt engineers. *arXiv preprint
447 arXiv:2211.01910*, 2022.

448 **NeurIPS Paper Checklist**

449 **1. Claims**

450 Question: Do the main claims made in the abstract and introduction accurately reflect the
451 paper's contributions and scope?

452 Answer: [\[Yes\]](#)

453 Justification: We have ensured that the main claims made in the abstract and introduction
454 are directly correlating to the research findings and the methods we have employed.

455 Guidelines:

- 456 • The answer NA means that the abstract and introduction do not include the claims
457 made in the paper.
- 458 • The abstract and/or introduction should clearly state the claims made, including the
459 contributions made in the paper and important assumptions and limitations. A No or
460 NA answer to this question will not be perceived well by the reviewers.
- 461 • The claims made should match theoretical and experimental results, and reflect how
462 much the results can be expected to generalize to other settings.
- 463 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
464 are not attained by the paper.

465 **2. Limitations**

466 Question: Does the paper discuss the limitations of the work performed by the authors?

467 Answer: [\[Yes\]](#)

468 Justification: The paper discusses the main limitations of the work performed by the authors
469 in the discussion section.

470 Guidelines:

- 471 • The answer NA means that the paper has no limitation while the answer No means that
472 the paper has limitations, but those are not discussed in the paper.
- 473 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 474 • The paper should point out any strong assumptions and how robust the results are to
475 violations of these assumptions (e.g., independence assumptions, noiseless settings,
476 model well-specification, asymptotic approximations only holding locally). The authors
477 should reflect on how these assumptions might be violated in practice and what the
478 implications would be.
- 479 • The authors should reflect on the scope of the claims made, e.g., if the approach was
480 only tested on a few datasets or with a few runs. In general, empirical results often
481 depend on implicit assumptions, which should be articulated.
- 482 • The authors should reflect on the factors that influence the performance of the approach.
483 For example, a facial recognition algorithm may perform poorly when image resolution
484 is low or images are taken in low lighting. Or a speech-to-text system might not be
485 used reliably to provide closed captions for online lectures because it fails to handle
486 technical jargon.
- 487 • The authors should discuss the computational efficiency of the proposed algorithms
488 and how they scale with dataset size.
- 489 • If applicable, the authors should discuss possible limitations of their approach to
490 address problems of privacy and fairness.
- 491 • While the authors might fear that complete honesty about limitations might be used by
492 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
493 limitations that aren't acknowledged in the paper. The authors should use their best
494 judgment and recognize that individual actions in favor of transparency play an impor-
495 tant role in developing norms that preserve the integrity of the community. Reviewers
496 will be specifically instructed to not penalize honesty concerning limitations.

497 **3. Theory Assumptions and Proofs**

498 Question: For each theoretical result, does the paper provide the full set of assumptions and
499 a complete (and correct) proof?

500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552

Answer: [NA]

Justification: Our paper does not require any explicit theorems and proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper has delineated all the information related to the experimental setup in the experimental setup section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

553 Question: Does the paper provide open access to the data and code, with sufficient instruc-
554 tions to faithfully reproduce the main experimental results, as described in supplemental
555 material?

556 Answer: [Yes]

557 Justification: We have released the code and dataset, making the dataset publicly available
558 under a license.

559 Guidelines:

- 560 • The answer NA means that paper does not include experiments requiring code.
- 561 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
562 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 563 • While we encourage the release of code and data, we understand that this might not be
564 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
565 including code, unless this is central to the contribution (e.g., for a new open-source
566 benchmark).
- 567 • The instructions should contain the exact command and environment needed to run to
568 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
569 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 570 • The authors should provide instructions on data access and preparation, including how
571 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 572 • The authors should provide scripts to reproduce all experimental results for the new
573 proposed method and baselines. If only a subset of experiments are reproducible, they
574 should state which ones are omitted from the script and why.
- 575 • At submission time, to preserve anonymity, the authors should release anonymized
576 versions (if applicable).
- 577 • Providing as much information as possible in supplemental material (appended to the
578 paper) is recommended, but including URLs to data and code is permitted.

579 6. Experimental Setting/Details

580 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
581 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
582 results?

583 Answer: [Yes]

584 Justification: We provided all the training and testing details in the experimental setup.

585 Guidelines:

- 586 • The answer NA means that the paper does not include experiments.
- 587 • The experimental setting should be presented in the core of the paper to a level of detail
588 that is necessary to appreciate the results and make sense of them.
- 589 • The full details can be provided either with the code, in appendix, or as supplemental
590 material.

591 7. Experiment Statistical Significance

592 Question: Does the paper report error bars suitably and correctly defined or other appropriate
593 information about the statistical significance of the experiments?

594 Answer: [Yes]

595 Justification: We conducted our experiments on all LLM-generated annotations across
596 zero-shot, one-shot, and few-shot settings, using majority voting as labels. Our reported
597 results represent the average performance across all test samples.

598 Guidelines:

- 599 • The answer NA means that the paper does not include experiments.
- 600 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
601 dence intervals, or statistical significance tests, at least for the experiments that support
602 the main claims of the paper.

- 603 • The factors of variability that the error bars are capturing should be clearly stated (for
604 example, train/test split, initialization, random drawing of some parameter, or overall
605 run with given experimental conditions).
- 606 • The method for calculating the error bars should be explained (closed form formula,
607 call to a library function, bootstrap, etc.)
- 608 • The assumptions made should be given (e.g., Normally distributed errors).
- 609 • It should be clear whether the error bar is the standard deviation or the standard error
610 of the mean.
- 611 • It is OK to report 1-sigma error bars, but one should state it. The authors should
612 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
613 of Normality of errors is not verified.
- 614 • For asymmetric distributions, the authors should be careful not to show in tables or
615 figures symmetric error bars that would yield results that are out of range (e.g. negative
616 error rates).
- 617 • If error bars are reported in tables or plots, The authors should explain in the text how
618 they were calculated and reference the corresponding figures or tables in the text.

619 8. Experiments Compute Resources

620 Question: For each experiment, does the paper provide sufficient information on the com-
621 puter resources (type of compute workers, memory, time of execution) needed to reproduce
622 the experiments?

623 Answer: [Yes]

624 Justification: We have included the specifications of the hardware and software environments
625 to ensure the reproducibility of our results.

626 Guidelines:

- 627 • The answer NA means that the paper does not include experiments.
- 628 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
629 or cloud provider, including relevant memory and storage.
- 630 • The paper should provide the amount of compute required for each of the individual
631 experimental runs as well as estimate the total compute.
- 632 • The paper should disclose whether the full research project required more compute
633 than the experiments reported in the paper (e.g., preliminary or failed experiments that
634 didn't make it into the paper).

635 9. Code Of Ethics

636 Question: Does the research conducted in the paper conform, in every respect, with the
637 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

638 Answer: [Yes]

639 Justification: The research conducted in this paper fully conforms with the NeurIPS Code of
640 Ethics in every respect.

641 Guidelines:

- 642 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 643 • If the authors answer No, they should explain the special circumstances that require a
644 deviation from the Code of Ethics.
- 645 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
646 eration due to laws or regulations in their jurisdiction).

647 10. Broader Impacts

648 Question: Does the paper discuss both potential positive societal impacts and negative
649 societal impacts of the work performed?

650 Answer: [Yes]

651 Justification: The paper explores how advancements and applications of our findings could
652 benefit society by capturing opinions of users in conversation benefit interpersonal skills.
653 Specifically, we investigate the effectiveness of current state-of-the-art large language models
654 in this context.

655 Guidelines:

- 656 • The answer NA means that there is no societal impact of the work performed.
- 657 • If the authors answer NA or No, they should explain why their work has no societal
658 impact or why the paper does not address societal impact.
- 659 • Examples of negative societal impacts include potential malicious or unintended uses
660 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
661 (e.g., deployment of technologies that could make decisions that unfairly impact specific
662 groups), privacy considerations, and security considerations.
- 663 • The conference expects that many papers will be foundational research and not tied
664 to particular applications, let alone deployments. However, if there is a direct path to
665 any negative applications, the authors should point it out. For example, it is legitimate
666 to point out that an improvement in the quality of generative models could be used to
667 generate deepfakes for disinformation. On the other hand, it is not needed to point out
668 that a generic algorithm for optimizing neural networks could enable people to train
669 models that generate Deepfakes faster.
- 670 • The authors should consider possible harms that could arise when the technology is
671 being used as intended and functioning correctly, harms that could arise when the
672 technology is being used as intended but gives incorrect results, and harms following
673 from (intentional or unintentional) misuse of the technology.
- 674 • If there are negative societal impacts, the authors could also discuss possible mitigation
675 strategies (e.g., gated release of models, providing defenses in addition to attacks,
676 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
677 feedback over time, improving the efficiency and accessibility of ML).

678 **11. Safeguards**

679 Question: Does the paper describe safeguards that have been put in place for responsible
680 release of data or models that have a high risk for misuse (e.g., pretrained language models,
681 image generators, or scraped datasets)?

682 Answer: [NA]

683 Justification: Our research does not pose any risks for misuse.

684 Guidelines:

- 685 • The answer NA means that the paper poses no such risks.
- 686 • Released models that have a high risk for misuse or dual-use should be released with
687 necessary safeguards to allow for controlled use of the model, for example by requiring
688 that users adhere to usage guidelines or restrictions to access the model or implementing
689 safety filters.
- 690 • Datasets that have been scraped from the Internet could pose safety risks. The authors
691 should describe how they avoided releasing unsafe images.
- 692 • We recognize that providing effective safeguards is challenging, and many papers do
693 not require this, but we encourage authors to take this into account and make a best
694 faith effort.

695 **12. Licenses for existing assets**

696 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
697 the paper, properly credited and are the license and terms of use explicitly mentioned and
698 properly respected?

699 Answer: [Yes]

700 Justification: We have explicitly cited the crawled websites, code and models used.

701 Guidelines:

- 702 • The answer NA means that the paper does not use existing assets.
- 703 • The authors should cite the original paper that produced the code package or dataset.
- 704 • The authors should state which version of the asset is used and, if possible, include a
705 URL.
- 706 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

717 **13. New Assets**

718 Question: Are new assets introduced in the paper well documented and is the documentation
719 provided alongside the assets?

720 Answer: [Yes]

721 Justification: We open-source the code and the new USDC dataset, and we provide complete
722 documentation on how the dataset was created.

723 Guidelines:

- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

732 **14. Crowdsourcing and Research with Human Subjects**

733 Question: For crowdsourcing experiments and research with human subjects, does the paper
734 include the full text of instructions given to participants and screenshots, if applicable, as
735 well as details about compensation (if any)?

736 Answer: [Yes]

737 Justification: We provide full instructions on how we surveyed our LLM generated annota-
738 tions using human participants in the Results section.

739 Guidelines:

- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

748 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
749 Subjects**

750 Question: Does the paper describe potential risks incurred by study participants, whether
751 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
752 approvals (or an equivalent approval/review based on the requirements of your country or
753 institution) were obtained?

754 Answer: [NA]

755 Justification: We use publicly available Reddit user conversations to create the USDC
756 dataset, and we do not collect any new data that would require IRB approval.

757 Guidelines:

758
759
760
761
762
763
764
765
766
767

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

768 **A Detailed Statistics of the USDC Dataset**

769 Table 4 shows the detailed statistics of our USDC dataset at the sub-reddit level. Fig. 5 shows the
 770 distribution of stance labels across LLM annotations across zero-shot, one-shot and few-shot settings.
 771 Fig. 6 shows the distribution of dogmatism labels across LLM annotations across zero-shot, one-shot
 and few-shot settings.

Table 4: Statistics of the User Conversation Dataset.

subreddit	num_conversations	min_total_token_count	max_total_token_count
DebateCommunism	73	529	11557
Abortiondebate	70	1271	7401
CapitalismVSocialism	61	665	16927
prochoice	60	582	7278
brexit	56	637	4553
climateskeptics	56	734	7550
prolife	54	672	13342
gunpolitics	52	683	7889
MensRights	52	623	5774
climatechange	49	520	7427
nuclear	41	572	5282
progun	39	436	3632
NuclearPower	23	629	4589
Vegetarianism	22	627	3958
AntiVegan	20	351	5052
climate	13	701	4678
Egalitarianism	10	665	4060
VeganActivism	8	460	3685
Veganism	2	1332	1738
AnimalRights	1	845	845
animalwelfare	1	1363	1363
GunsAreCool	1	2945	2945

772

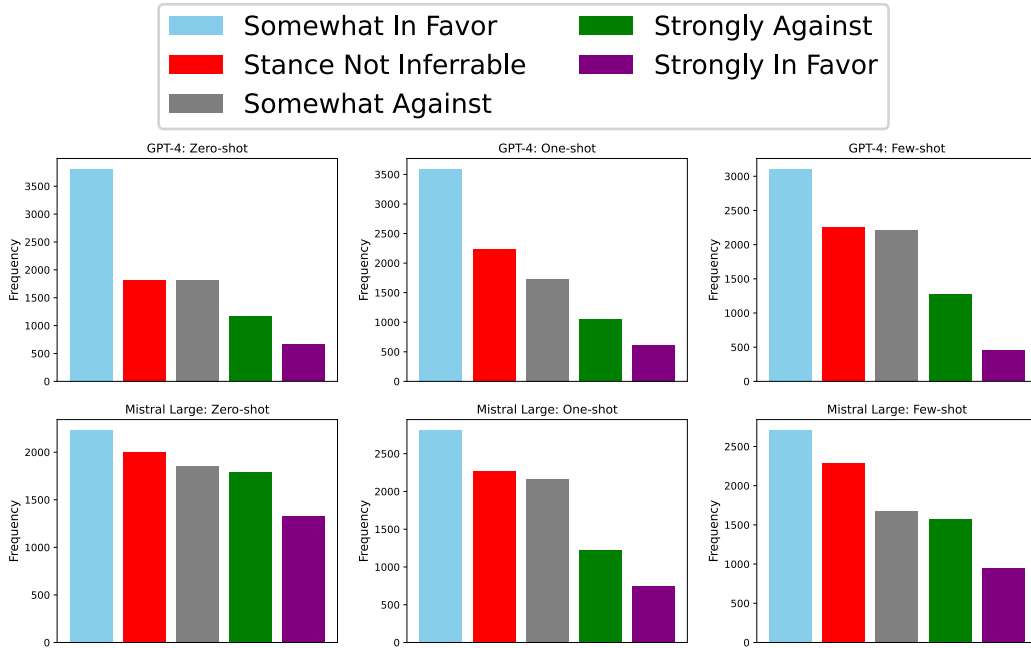


Figure 5: Distribution of Stance labels across LLM annotations.

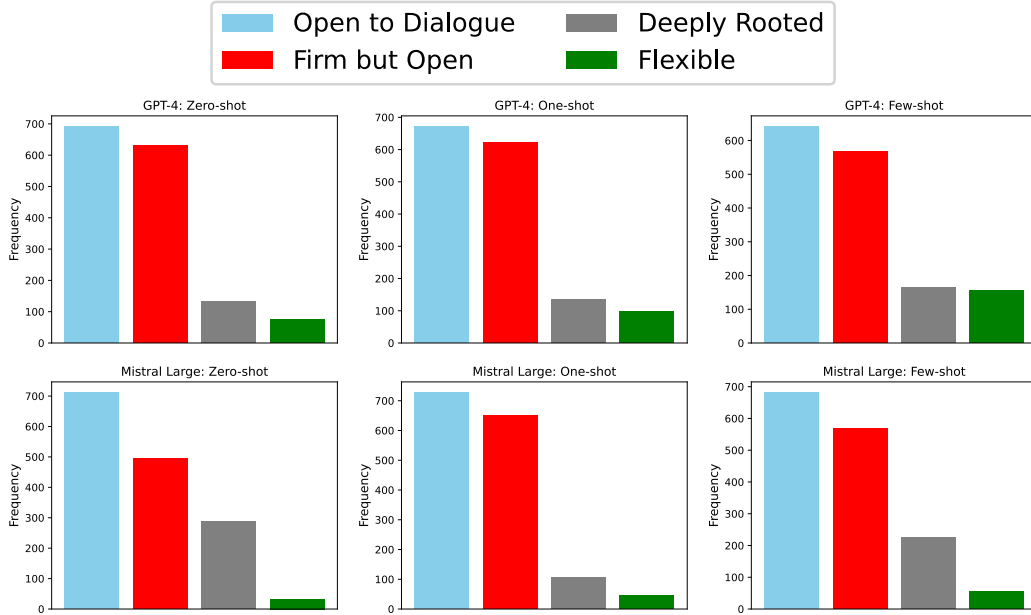


Figure 6: Distribution of dogmatism labels across LLM annotations.

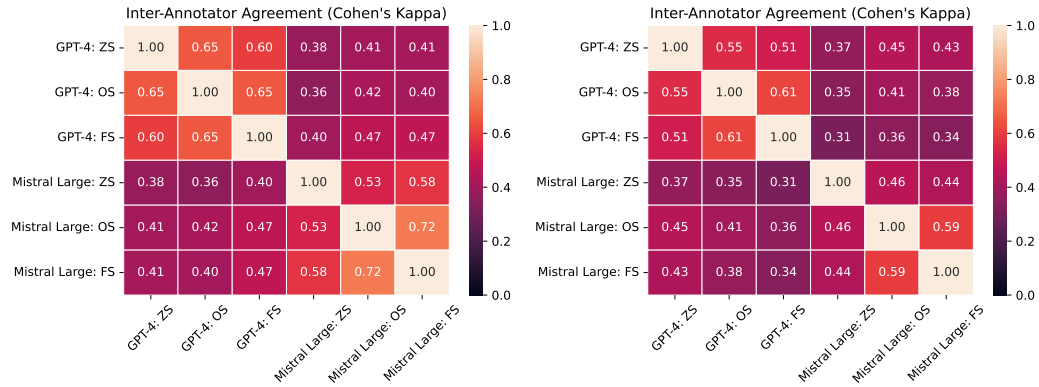


Figure 7: Inter-annotator agreement (IAA): Cohen's Kappa score across six different models (2 models \times 3 settings) for Stance (left) and Dogmatism (right) tasks.

773 B System Prompt for LLM Annotation

774 We used the following prompt for getting annotations from LLMs as well as for instruction-tuning of
775 SLMs.

```

776 """
777 ### Introduction
778 **Objective**: Analyze Reddit conversations to identify the stance of
779 specific authors on sociopolitical topics and determine their level of
780 dogmatism.
781 **Stance Definition**: Stance is defined as the expression of the author's
782 standpoint and judgement towards a given topic.
783 **Dogmatism Definition**: Dogmatism is an opinion strongly believed as a fact
784 to support a stance without a question or allowance for conversation.
785 **Task**: Given a JSON formatted Reddit submission and its comment thread,
786 classify the stance of text segments related to 'author1' and
787 'author2' by assigning one of the following five predefined stance
788 labels: 'strongly_against', 'somewhat_against', 'somewhat_in_favor',
789 'strongly_in_favor', 'stance_not_inferable'. Also, assign a dogmatism

```

790 label for each author by assigning one of the following four predefined
791 labels: 'Deeply Rooted', 'Firm but Open', 'Open to Dialogue', 'Flexible'.
792

793 **### Description of Stance Labels:**
794 1. ****strongly_against / strongly_in_favor****: Marks text showing strong
795 opinions, emotional expressions, or argumentative tones.
796 2. ****somewhat_against / somewhat_in_favor****: Identifies texts with openness
797 to discussion, less certainty, or showing interest in different
798 viewpoints.
799 3. ****stance_not_inferrable****: Use for texts that are neutral, support both
800 stances, or where the stance is unclear despite being on-topic.
801

802 **### Description of Dogmatism Labels:**
803 1. ****Deeply Rooted****: Reflects a strong, unchangeable belief. This label
804 conveys the idea of someone who is firm in their opinion and unlikely to
805 be swayed.
806 2. ****Firm but Open****: Indicates a person who is not likely to change their
807 mind but does not impose their views authoritatively. It captures the
808 essence of being steadfast in one's beliefs without being dismissive of
809 others.
810 3. ****Open to Dialogue****: Describes someone who holds a certain opinion but is
811 genuinely interested in considering other viewpoints. This label suggests
812 a willingness to engage in meaningful conversation about differing
813 perspectives.
814 4. ****Flexible****: Denotes a person who is not firmly committed to their stance
815 and is open to changing their opinion. This label is indicative of
816 flexibility and openness to new information or arguments.
817

818 **### Input Data Format**
819 The input data will be in JSON format and will include several key elements
820 to represent a Reddit submission and its associated comments. Each
821 element provides specific information as described below:
822

- 823 - 'id': This is the unique identifier for the Reddit submission.
- 824 - 'title': The title of the post. This is what users see first and often
825 summarizes or hints at the content of the submission.
- 826 - 'content': The main post's detailed description. This text segment provides
827 the core message or information the author wishes to communicate with the
828 Reddit community. It may include narratives, questions, or any
829 information relevant to the title.
- 830 - 'author1' or 'author2': The username of our focus author. This field is
831 applicable if the post or comment is made by one of the specific authors
832 we are tracking in the dataset.
- 833 - 'comments': An array (list) of comments related to the Reddit submission.
834 Each comment in this array includes the following fields:
835 - 'id': The unique identifier for the comment, allowing for identification
836 and reference within the dataset.
837 - 'author1' or 'author2': The username of the comment's author, if it is
838 made by one of our focus authors. This helps in tracking contributions
839 by specific individuals.
840 - 'body': The text of the comment. This is the main content of the comment
841 where the author responds to the post or another comment, providing
842 insights, opinions, or further information.
843 - 'replies': An array of comments that are direct responses to this
844 comment. The structure of each reply follows the same format as the
845 initial comment, including 'id', 'author1' or 'author2' (if
846 applicable), 'body', and potentially more 'replies'.
847

848 **### Output Data Format**
849 Submit your annotations in JSON format, grouping all stance annotations under
850 the key 'stance_annotations'. Each entry should be a dictionary
851 containing the segment's 'id', your 'label', and the 'reason' for
852 your choice. Include the dogmatism label and its justification under
853 'dogmatism_label' and 'dogmatism_reason' keys, respectively.
854

```

855 The output should follow this structure:
856 ‘‘json
857 {
858   "author1": {
859     "name": "[author_name]",
860     "stance_annotations": [
861       {
862         "id": "[segment_id]",
863         "label": "[chosen_label]",
864         "reason": "[Justification in <50 words]"
865       },
866       ...
867     ],
868     "dogmatism_label": "[chosen_dogmatism_label]",
869     "dogmatism_reason": "[Justification in <50 words]"
870   },
871   "author2": {
872     "name": "[author_name]",
873     "stance_annotations": [
874       {
875         "id": "[segment_id]",
876         "label": "[chosen_label]",
877         "reason": "[Justification in <50 words]"
878       },
879       ...
880     ],
881     "dogmatism_label": "[chosen_dogmatism_label]",
882     "dogmatism_reason": "[Justification in <50 words]"
883   }
884 }
885 ’’
886 ### Instructions for Effective Annotation
887
888 1. Labeling Stance: For each segment (including the original Reddit
889 submission, comments, or replies) where "author1" or "author2" is
890 mentioned, assign a stance label that best represents the stance
891 expressed towards the discussed topic in the submission. This
892 comprehensive approach ensures no relevant contribution by "author1" or
893 "author2" is overlooked. Evaluate the stance based on the content’s tone,
894 argumentation, and engagement level with the topic.
895
896 2. Providing Justification: For each label assigned, include a concise
897 reason, aiming for less than 50 words. Focus on the stance and
898 argumentative indicators present in the text.
899
900 3. Dogmatism Assessment: After reviewing all segments from "author1" and
901 "author2", assign a single dogmatism label reflecting the overall tone
and approach in their contributions.
'''

```

902 C Prompts for Finetuning SLMs

903 Fig. 8 and 9 shows the prompts used for finetuning SLMs for the stance and dogmatism classification
904 tasks respectively.

Stance Classification

Analyze the stance of the post enclosed in square brackets.
Categorize each post into one of the following categories based on its stance:

- Somewhat In Favor
- Somewhat Against
- Stance Not inferrable
- Strongly In Favor
- Strongly Against

and return the answer as one of the corresponding stance labels.

```
[{data_point["stance_id_comment"]}]]
```

Figure 8: Prompt for stance classification, for finetuning SLMs.

User Dogmatism Identification

Analyze the comments of a user in conversation enclosed in square brackets.
Categorize the opinion fluctuation of the user into one of the following categories based on its change:

- Open to Dialogue
- Firm but Open
- Deeply Rooted
- Flexible

Return the answer as one of the corresponding dogmatism labels.

```
[{data_point["comments_string_for_dogmatism"]}]]
```

Figure 9: Prompt for dogmatism classification, for finetuning SLMs.

905 D Sample of User Input Prompt

```
906 ""
907 ""
908 ### User Prompt
909 Now complete the given task for the respective authors i.e., author1
910 name is "rookerin0" and respective ids are ['dhoxyz', 'f3pghji', '
911 f3tywb4', 'f3uomn2']. author2 name is "MikeWillTerminate" and
912 respective ids are ['f3rt0bf', 'f3rqu2u'] for the data in json
913 format
914 {
915   "id":"dhoxyz",
916   "author1":"rookerin0",
917   "title":"This sub should encourage anti vs. pro-gun discussions
918   instead of shutting them down instantly",
919   "content":"Honesly, I followed this sub especifically to take part
920   in these discussions, but everytime I see a comment that even
921   remotely suggests anti gun ideals or a discussion on the
922   subject just gets ignored and downvoted to hell. Kind of
923   expecting this to go the same way (my karma anus is ready,
924   downvotes) , but I have to hope for healthy discussions on the
925   subject.",
926   "comments":[
927     {
928       "id":"f3p9n2c",
929       "body":"I think the problem now is the two sides are at an
930       impasse. Everytime there is a "compromise" pro gun loses
```

```

931 something. Now days pro gun is interpreting the
932 Constitution more literal, which leaves even the most
933 mild policies of anti gun as infringements. To further
934 compound this anti gun is only considering the most
935 extreme measures. "Assault Weapons" bans, mandatory
936 buybacks, red flag laws, etc.. I think at this point
937 there is just nothing left to talk about. The middle
938 ground is gone.",
939 "replies":[
940 {
941 "id":"f3pati9",
942 "replies":[
943 {
944 "id":"f3pdu44",
945 "body":"You are exactly right. I'm done with the
946 idea that there can be real compromise. We
947 should have at least gotten national
948 reciprocity and shall-issue in every state in
949 exchange for what we've given up. Now you
950 have to be a goddamn lawyer to exercise your
951 rights without violating the law."
952 },
953 {
954 "author2":"MikeWillTerminate",
955 "id":"f3rt0bf",
956 "body":"I am prepared for UBCs, if they do this:
957 1. Lower the age to buy handguns to 18,
958 nationwide.
959 2. Repeal the Hughes Amendment:
960 3. A FOIPA-like ban on assault weapon bans (what
961 the FOIPA did with a registry)
962 4. The punishment for violation is a monetary
963 fine only
964 5. A repeal of the GCA ban on foreign NFA weapons
965 6. A repeal of the National Minimum Drinking Age
966 Act of 1984"
967 }
968 ]
969 },
970 {
971 "id":"f3pd55z",
972 "body":"Everytime there is a "compromise" pro gun loses
973 something. That and today's compromise is tomorrow
974 's loophole to be closed. All such compromises do
975 is push that policy off until the next round."
976 }
977 ]
978 },
979 {
980 "id":"f3paf0j",
981 "body":"Yeah this sub it's not conducive to conversion. Its
982 quickly devolving to little more than "Boogaloo" memes
983 and shouting "SHALL. NOT." at each other. However, as
984 far as I know, the mods won't delete your thread and ban
985 you from the sub for trying to have a good faith
986 discussion, like some of the gun control subs will.",
987 "replies":[
988 {
989 "id":"f3pusbm",
990 "body":"Unfortunately this sub's mod team takes a very
991 passive approach to moderation. With very little
992 effort they could make this sub into a quality
993 progun meeting ground *without having to resort to
994 censorship*. Instead they promote low-effort memes
995 and endless duplication of posts through their

```



```

996         inaction. whubbard has the chops to resurrect this
997         sub. Let's see if he's up to the challenge.",
998         "replies":[
999             {
1000                 "id":"f3q8xj6",
1001                 "body":"We voted to ban memes last week. All
1002                     about rolling it out now.",
1003                 "replies":[
1004                     {
1005                         "id":"f3qn4p8",
1006                         "body":"Damn I might have to eat some crow
1007                             here then..."
1008                     }
1009                 ]
1010             }
1011         ]
1012     },
1013 ]
1014 },
1015 {
1016     "id":"f3pafqa",
1017     "body":"Found the gun grabber!!",
1018     "replies":[
1019         {
1020             "id":"f3pcw4h",
1021             "body":"Witch hunter."
1022         }
1023     ]
1024 },
1025 {
1026     "id":"f3pal5l",
1027     "body":"I see people have discussions when it makes sense to.
1028         Not much reason to spend time responding to the same gun
1029         control measures over and over though."
1030 },
1031 {
1032     "id":"f3paw3h",
1033     "body":"I get where you're coming from, but people's ability
1034         to protect themselves and own their own property isn't
1035         something that is compromisable. Anything less, and they
1036         cease to own their own property. It's like breathing,
1037         there can be nothing less than total ability to breath
1038         when and how someone wants. It's just that simple."
1039 },
1040 {
1041     "id":"f3pax9m",
1042     "body":"My take on this, What kind of open discussion is
1043         possible for a right that is guaranteed and most
1044         importantly, not to be infringed upon? They're making all
1045         these unlawful laws to portray it as it's somehow
1046         legitimate. They are not, We are at an apex, to which
1047         both political spectrums and even us to a degree are
1048         liable for.\nI certainly believe both sides are waiting
1049         for this to boil over so each can finger point. I just
1050         speculate it's going to be the hell humanity been
1051         whispering about but never thought it would ever occur."
1052 },
1053 {
1054     "id":"f3pb6ny",
1055     "body":"The time for discussion is over."
1056 },
1057 {
1058     "id":"f3pfqwq",
1059     "body":"I don't know what you're talking about. Sure people
1060         downvote, but they also talk. We get "why do you need

```

```

1061 guns" posts at least weekly, and several people will
1062 engage in actual conversation with them, citing facts,
1063 clearing up statistics, and telling stories to illustrate
1064 why this is important to them, but they are usually met
1065 with "you stupid @#$$, you think you're Rambo" or
1066 something equally clever. People who come here to discuss
1067 and learn will be treated well. People who are just
1068 trolling are treated like trolls.",
1069 "replies":[
1070 {
1071   "author1":"rookerin0",
1072   "id":"f3pghji",
1073   "body":"I made this post because I'm always seeing
1074     rational, conversation seeking comments getting
1075     blown to downvote hell.",
1076   "replies":[
1077     {
1078       "id":"f3pi9xv",
1079       "body":"[Like this one?](https://www.reddit.com/r/
1080         progun/comments/dhcu92/yup/f3p75tg/)> One smart
1081         man in a sub full of... welp... "strong opinions
1082         ". You start off with arrogance, as the sole
1083         arbiter of what constitutes a "smart man". Then
1084         you back it up with a dismissive swipe at what
1085         you term "strong opinions".> Every other country
1086         can see that PROPER gun control reduces gun
1087         violence by a ton, More arrogance. False
1088         equivalence. Unsupported claims.> but the US
1089         refuses to let go of it's antique laws In a
1090         shocking turn of events, more arrogance.> Fully
1091         aware that this is a fully pro gun sub, willing
1092         to take the downvotes in order to spark a
1093         discussion and crack some heads. You aren't the
1094         first arrogant asshole to grace this sub with
1095         posts like this. Try bringing something other
1096         than your own self-importance to the discussion.
1097         Edit: And then there's [this gem](https://www.
1098         reddit.com/r/unpopularopinion/comments/d3w5z1/
1099         people_living_in_the_us_are_living_in_one_of_the/
1100         f06r3sg/.> Wanna feel like you could be shot at
1101         every single moment? Move to the US, it'll prob
1102         happen to you either as a bystander, or you'd be
1103         shot by a random citizen (sometimes police)."
1104     },
1105     {
1106       "id":"f3pj8k0",
1107       "body":"As is tradition. We're done with that
1108         condescending bullshit from antis, you dont
1109         come here for good faith discussion and
1110         whether you get a reasonable response or not,
1111         nothing ever changes, easier to downvote you
1112         and move on because we get the same
1113         treatment anytime we attempt to speak out in
1114         anti subs."
1115     },
1116     {
1117       "id":"f3plgf4",
1118       "body":"If downvotes hurt your feelings, you
1119         shouldn't be on reddit. People tend to
1120         downvote anything they disagree with (which
1121         is why some subs specifically ask you to only
1122         downvote things that contribute nothing to
1123         the discussion). It's a bad habit, but that's
1124         the way it is. People downvote and *still*
1125         enage. You want to post a view contrary to

```

```

1126         the prevailing view of the sub, take your
1127         lumps and participate in what conversation
1128         you are offered. But if you're only here to
1129         preach about how stupid, misguided, unevolved
1130         , uneducated, irrational, and/or violent we
1131         are, don't expect a polite response."
1132     },
1133     {
1134         "id":"f3tcgf1",
1135         "body":"An arrogant Israeli trying to tell
1136         another nation how they should be run. You're
1137         just a walking stereotype aren't you? And
1138         before you say anything, I popped into your
1139         comment history. That's where the calling you
1140         Israeli comes from.",
1141         "replies":[
1142             {
1143                 "author1":"rookerin0",
1144                 "id":"f3tywb4",
1145                 "body":"I thought that trying to tell other
1146                 nations how they should run was your
1147                 guys's stereotype.",
1148                 "replies":[
1149                     {
1150                         "id":"f3u0vkq",
1151                         "body":"No we go in and try to make
1152                         them work our way."
1153                     }
1154                 ]
1155             }
1156         ]
1157     }
1158 ]
1159 }
1160 ]
1161 },
1162 {
1163     "id":"f3pzseh",
1164     "body":"It's a little unfortunate but the grabbers who come
1165     on here tend to be intellectually dishonest and/or
1166     uninformed. There was some Australian post a few days ago
1167     that pretty much asked why we like our guns more than
1168     children. No discussion to be had there. There's also
1169     some posts that clearly demonstrate the poster should
1170     inform himself or herself a little."
1171 },
1172 {
1173     "author2":"MikeWillTerminate",
1174     "id":"f3rqu2u",
1175     "body":"Actually, do that. It shows everyone that they tend
1176     to be crazy, unstable, ignorant, stereotyping, arrogant
1177     bastards who hate black people with a hair trigger."
1178 },
1179 {
1180     "id":"f3t7tgg",
1181     "body":"Welcome to reddit, home of every single safe place
1182     for anything that doesnt violate the TOS. At least its
1183     slightly better than r/politics"
1184 },
1185 {
1186     "id":"f3unt9z",
1187     "body":"This isn't r/gundebate. This is a pro gun subreddit.
1188     That said, we do allow some debate provided it remains
1189     civil.",
1190     "replies":[

```

```

1191         {
1192             "author1": "rookerin0",
1193             "id": "f3uomn2",
1194             "body": "Sadly tho, r/gundebate is pretty dead..."
1195         }
1196     ]
1197 },
1198 {
1199     "id": "f4dip6o",
1200     "body": "Anything else you want to give away for free?"
1201 }
1202 ]
1203 }
1204
1205 ##### Response

```

1207 E Samples of JSON Outputs from LLMs

1208 E.1 Sample JSON Output from LLM: GPT-4 Zero-shot

```

1209 {
1210     {
1211         "author1": {
1212             "name": "rookerin0",
1213             "stance_annotations": [
1214                 {
1215                     "id": "dhxyz",
1216                     "label": "somewhat_in_favor",
1217                     "reason": "Advocates for open discussions, showing interest in
1218                             dialogue."
1219                 },
1220                 {
1221                     "id": "f3pghji",
1222                     "label": "somewhat_in_favor",
1223                     "reason": "Expresses desire for rational discussions, despite
1224                             downvotes."
1225                 },
1226                 {
1227                     "id": "f3tywb4",
1228                     "label": "somewhat_against",
1229                     "reason": "Sarcastic retort, but lacks strong emotional
1230                             language."
1231                 }
1232             ],
1233             {
1234                 "id": "f3uomn2",
1235                 "label": "stance_not_inferrable",
1236                 "reason": "Comment is factual, no clear stance on the debate
1237                             itself."
1238             }
1239         ],
1240         "dogmatism_label": "Open to Dialogue",
1241         "dogmatism_reason": "Shows interest in discussion, acknowledges
1242                             downvotes."
1243     },
1244     "author2": {
1245         "name": "MikeWillTerminate",
1246         "stance_annotations": [
1247             {
1248                 "id": "f3rt0bf",
1249                 "label": "somewhat_in_favor",
1250                 "reason": "Proposes conditions for UBCs, indicating some
1251                             openness."
1252             }
1253         ]
1254     }
1255 }

```

```

1253 {
1254     "id": "f3rqu2u",
1255     "label": "strongly_against",
1256     "reason": "Uses strong negative descriptors for opponents,
1257             indicating firm stance."
1258 }
1259 ],
1260 "dogmatism_label": "Firm but Open",
1261 "dogmatism_reason": "Mixed; proposes solutions but also dismisses
1262             opposition."
1263 }
1264 }
1265

```

1266 E.2 Sample JSON Output from LLM: GPT-4 One-shot

```

1267 {
1268     "author1": {
1269         "name": "rookerin0",
1270         "stance_annotations": [
1271             {
1272                 "id": "dhoxyz",
1273                 "label": "somewhat_in_favor",
1274                 "reason": "Advocates for open discussions, indicating a pro-
1275                         discussion stance."
1276             },
1277             {
1278                 "id": "f3pghji",
1279                 "label": "somewhat_in_favor",
1280                 "reason": "Expresses desire for rational conversation, showing
1281                         a supportive stance for dialogue."
1282             },
1283             {
1284                 "id": "f3tywb4",
1285                 "label": "somewhat_against",
1286                 "reason": "Uses sarcasm to counter a stereotype, indicating a
1287                         stance against the opposing view."
1288             },
1289             {
1290                 "id": "f3uomn2",
1291                 "label": "somewhat_in_favor",
1292                 "reason": "Expresses disappointment in lack of active debate,
1293                         supporting open discussion."
1294             }
1295         ],
1296         "dogmatism_label": "Open to Dialogue",
1297         "dogmatism_reason": "Shows a consistent preference for discussion
1298                 and engagement, without dogmatic assertions."
1299     },
1300     "author2": {
1301         "name": "MikeWillTerminate",
1302         "stance_annotations": [
1303             {
1304                 "id": "f3rt0bf",
1305                 "label": "somewhat_in_favor",
1306                 "reason": "Lists conditions for agreement, showing openness to
1307                         compromise."
1308             },
1309             {
1310                 "id": "f3rqu2u",
1311                 "label": "strongly_against",
1312                 "reason": "Uses strong language to describe opponents,
1313                         indicating a firm stance against them."
1314             }
1315         ]

```

```

1316     ],
1317     "dogmatism_label": "Firm but Open",
1318     "dogmatism_reason": "While showing strong opinions, also lists
1319         specific conditions for compromise."
1320 }
1321 }
1322 }

```

E.3 Sample JSON Output from LLM: GPT-4 large Few-shot

```

1323 {
1324 {
1325 {
1326     "author1": {
1327     "name": "rookerin0",
1328     "stance_annotations": [
1329     {
1330         "id": "dhoxyz",
1331         "label": "somewhat_in_favor",
1332         "reason": "Expresses desire for open discussions, indicating a
1333             stance in favor of dialogue."
1334     },
1335     {
1336         "id": "f3pghji",
1337         "label": "somewhat_in_favor",
1338         "reason": "Advocates for rational conversation despite
1339             observing negative reactions."
1340     },
1341     {
1342         "id": "f3tywb4",
1343         "label": "somewhat_against",
1344         "reason": "Uses sarcasm to counter a stereotype, showing a
1345             stance against it."
1346     },
1347     {
1348         "id": "f3uomn2",
1349         "label": "stance_not_inferrable",
1350         "reason": "Comment on the state of another subreddit does not
1351             reveal a stance on the gun debate."
1352     }
1353     ],
1354     "dogmatism_label": "Open to Dialogue",
1355     "dogmatism_reason": "Shows a consistent
1356         preference for discussion over confrontation."
1357     },
1358     "author2": {
1359     "name": "MikeWillTerminate",
1360     "stance_annotations": [
1361     {
1362         "id": "f3rt0bf",
1363         "label": "somewhat_in_favor",
1364         "reason": "Lists conditions for accepting UBCs, indicating
1365             some openness to compromise."
1366     },
1367     {
1368         "id": "f3rqu2u",
1369         "label": "strongly_against",
1370         "reason": "Uses derogatory language to describe opponents,
1371             showing strong opposition."
1372     }
1373     ],
1374     "dogmatism_label": "Firm but Open",
1375     "dogmatism_reason": "While showing strong opinions, also suggests
1376         conditions for policy acceptance."
1377     }
1378 }
1379 }

```

1380 E.4 Sample JSON Output from LLM: Mistral Large Zero-shot

```
1381 {
1382 {
1383   "author1": {
1384     "name": "rookerin0",
1385     "stance_annotations": [
1386     {
1387       "id": "dhxyz",
1388       "label": "somewhat_in_favor",
1389       "reason": "Author expresses desire for open discussions on
1390         anti-gun ideals."
1391     },
1392     {
1393       "id": "f3pghji",
1394       "label": "stance_not_inferrable",
1395       "reason": "Author is commenting on downvotes, not expressing
1396         stance."
1397     },
1398     {
1399       "id": "f3tywb4",
1400       "label": "stance_not_inferrable",
1401       "reason": "Author is making a stereotypical comment, not
1402         expressing stance."
1403     },
1404     {
1405       "id": "f3uomn2",
1406       "label": "stance_not_inferrable",
1407       "reason": "Author is commenting on the state of another
1408         subreddit, not expressing stance."
1409     }
1410   ],
1411   "dogmatism_label": "Open to Dialogue",
1412   "dogmatism_reason": "Author expresses desire for open discussions
1413     and does not dismiss opposing views."
1414 },
1415   "author2": {
1416     "name": "MikeWillTerminate",
1417     "stance_annotations": [
1418     {
1419       "id": "f3rt0bf",
1420       "label": "strongly_in_favor",
1421       "reason": "Author lists conditions for accepting UBCs, showing
1422         strong pro-gun stance."
1423     },
1424     {
1425       "id": "f3rqu2u",
1426       "label": "strongly_in_favor",
1427       "reason": "Author's comment is dismissive of anti-gun
1428         advocates, showing strong pro-gun stance."
1429     }
1430   ],
1431   "dogmatism_label": "Deeply Rooted",
1432   "dogmatism_reason": "Author's comments are dismissive of opposing
1433     views and show a strong, unchangeable belief in pro-gun ideals
1434     ."
1435 }
1436 }
1437 }
```

1438 E.5 Sample JSON Output from LLM: Mistral Large One-shot

```
1439 {
1440 {
1441   "author1": {
1442     "name": "rookerin0",
```

```

1443 "stance_annotatations": [
1444 {
1445     "id": "dhoxyz",
1446     "label": "somewhat_in_favor",
1447     "reason": "Expresses desire for open discussions on gun
1448         control, but follows pro-gun subreddit."
1449 },
1450 {
1451     "id": "f3pghji",
1452     "label": "stance_not_inferrable",
1453     "reason": "Comments on downvotes, but no clear stance on gun
1454         control."
1455 },
1456 {
1457     "id": "f3tywb4",
1458     "label": "stance_not_inferrable",
1459     "reason": "Comment about stereotypes, but no clear stance on
1460         gun control."
1461 },
1462 {
1463     "id": "f3uomn2",
1464     "label": "stance_not_inferrable",
1465     "reason": "Mentions r/gundebate being dead, but no clear
1466         stance on gun control."
1467 }
1468 ],
1469 "dogmatism_label": "Open to Dialogue",
1470 "dogmatism_reason": "Seeks open discussions and does not express
1471     strong dogmatic views."
1472 },
1473 "author2": {
1474     "name": "MikeWillTerminate",
1475     "stance_annotatations": [
1476     {
1477         "id": "f3rt0bf",
1478         "label": "somewhat_in_favor",
1479         "reason": "Lists conditions for supporting UBCs, showing
1480             openness to compromise."
1481     },
1482     {
1483         "id": "f3rqu2u",
1484         "label": "strongly_in_favor",
1485         "reason": "Strong negative language against gun control
1486             advocates, suggesting strong pro-gun stance."
1487     }
1488     ],
1489     "dogmatism_label": "Firm but Open",
1490     "dogmatism_reason": "Expresses strong pro-gun views but open to
1491         compromise under certain conditions."
1492     }
1493 }

```

1495 E.6 Sample JSON Output from LLM: Mistral Large Few-shot

```

1496 {
1497     "author1": {
1498         "name": "rookerin0",
1499         "stance_annotatations": [
1500         {
1501             "id": "dhoxyz",
1502             "label": "somewhat_in_favor",
1503             "reason": "Expresses desire for open discussions on anti vs.
1504                 pro-gun debates."
1505         }

```

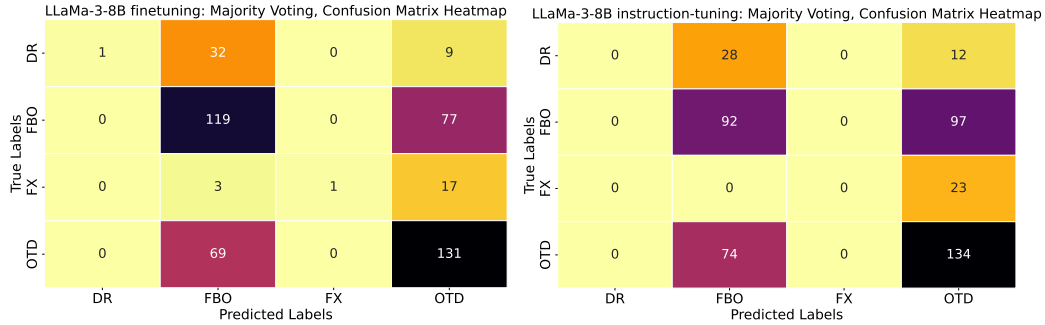



Figure 10: Confusion matrix for LLaMa-3-8B Dogmatism detection models on USDC test set: finetuning (left) and instruction-tuning (right). Here, DR: Deeply Rooted, FX: Flexible, FBO: Firm but Open, OTD: Open to Dialogue

```

1506   },
1507   {
1508     "id": "f3pghji",
1509     "label": "stance_not_inferrable",
1510     "reason": "Comments on downvotes, but no clear stance on the
1511               topic."
1512   },
1513   {
1514     "id": "f3tywb4",
1515     "label": "stance_not_inferrable",
1516     "reason": "Makes a stereotypical comment, but no clear stance
1517               on the topic."
1518   },
1519   {
1520     "id": "f3uomn2",
1521     "label": "stance_not_inferrable",
1522     "reason": "Mentions r/gundebate being dead, but no clear
1523               stance on the topic."
1524   }
1525 ],
1526 "dogmatism_label": "Open to Dialogue",
1527 "dogmatism_reason": "Seeks open discussions and engages in
1528                      conversation without strong dogmatic language."
1529 },
1530 "author2": {
1531   "name": "MikeWillTerminate",
1532   "stance_annotations": [
1533     {
1534       "id": "f3rt0bf",
1535       "label": "somewhat_in_favor",
1536       "reason": "Lists conditions for supporting UBCs, showing
1537                 openness to discussion."
1538     },
1539     {
1540       "id": "f3rqu2u",
1541       "label": "strongly_against",
1542       "reason": "Uses derogatory language to express strong
1543                 opposition to 'grabbers'."
1544     }
1545   ],
1546   "dogmatism_label": "Firm but Open",
1547   "dogmatism_reason": "Expresses strong opinions but also shows
1548                       willingness to consider certain conditions for compromise."
1549 }
1550 }

```

1552 **F SLM finetuning: AUC (Area Under the Curve) analysis**

1553 Fig. 10 illustrates the confusion matrix for dogmatism detection for LLaMa-3-8B finetuning and
 1554 instruction-tuning. We make the following observations from Fig. 10: 1) For both finetuning and
 1555 instruction-tuning, there is significant misclassifications, especially for the "Deeply Rooted" and
 1556 "Flexible" labels, with both having zero accuracy and F1-scores. While "Firm but Open" and
 1557 "Open to Dialogue" perform moderately better, with accuracies of 48.7% and 64.4% respectively.
 1558 The confusion matrix indicates substantial confusion to distinguish between intermediate levels of
 1559 dogmatism, such as "Firm but Open" and "Open to Dialogue. We further reports the ROC curve
 1560 shows the trade-off between the true positive rate (TPR) and false positive rate (FPR) for each class
 1561 for stance and dogmatism tasks, in Figs. 11 and. 12. The area under the ROC curve (AUC) is a
 1562 measure of the model's ability to distinguish between classes.

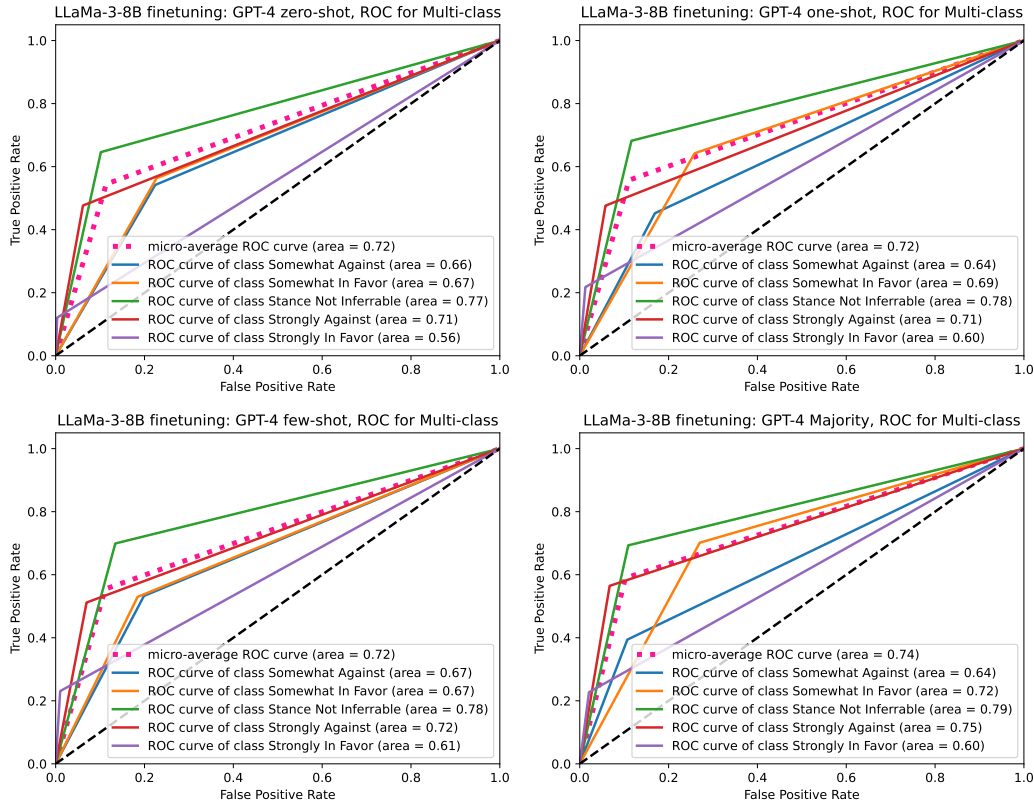


Figure 11: LLaMa-3-8B finetuning for stance detection task: Visualize the ROC curves for each class along with their AUC values for GPT-4 Annotations across zero-shot, one-shot, few-shot and majority labels.

1563 **G SLM instruction-tuning: AUC (Area Under the Curve) analysis**

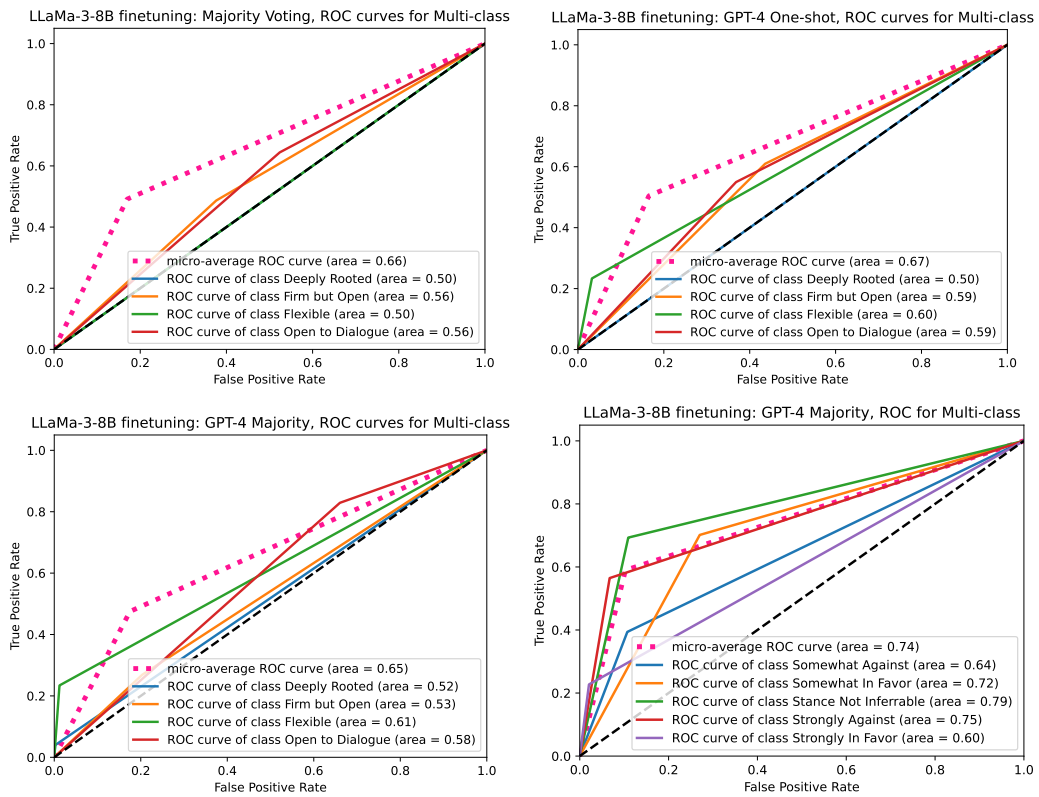


Figure 12: LLaMa-3-8B finetuning for dogmatism task: Visualize the ROC curves for each class along with their AUC values for GPT-4 Annotations across zero-shot, one-shot, few-shot and majority labels.

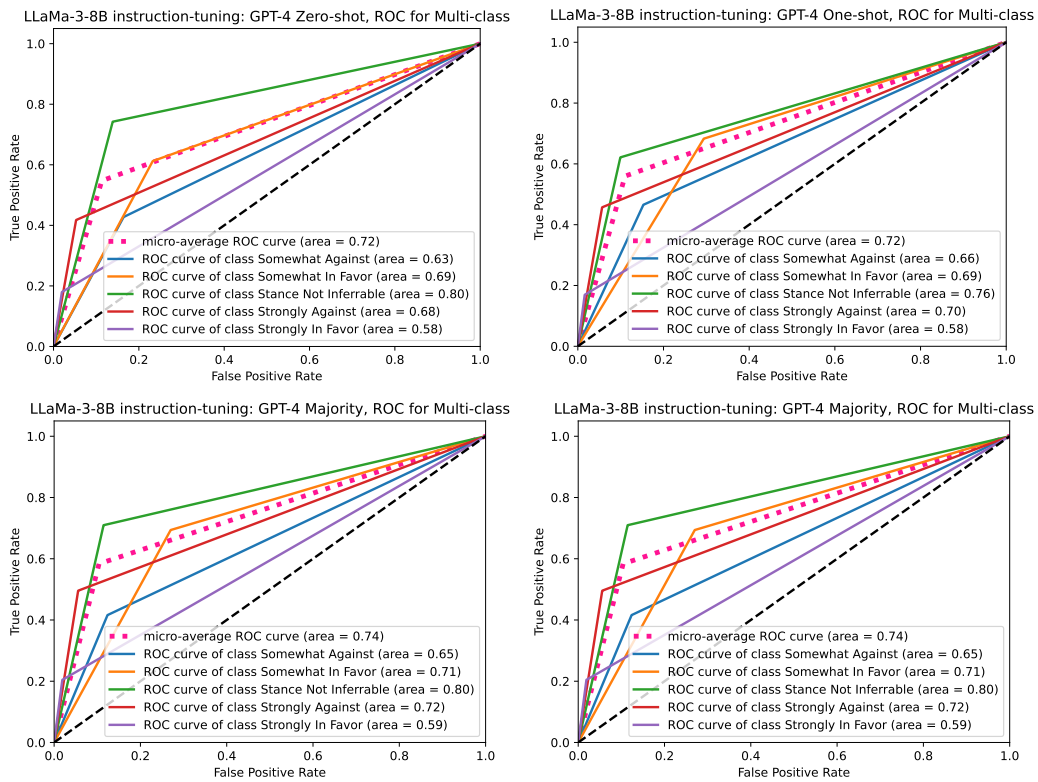


Figure 13: LLaMa-3-8B instruction-tuning for stance detection task: Visualize the ROC curves for each class along with their AUC values for GPT-4 Annotations across zero-shot, one-shot, few-shot and majority labels.