

Understanding and Tackling Label Errors in Individual-Level Nature Language Understanding

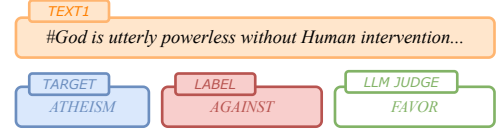
Anonymous ACL submission

Abstract

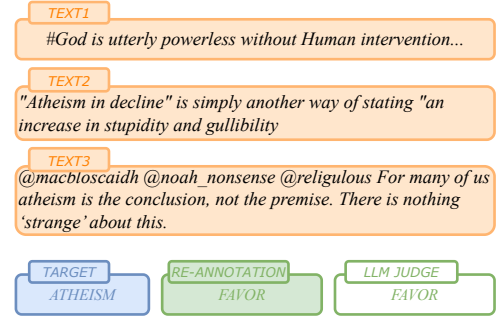
Natural language understanding (NLU) is a task that enables machines to understand human language. Some tasks, such as stance detection and sentiment analysis, are closely related to individual subjective perspectives, thus termed individual-level NLU. Previously, these tasks are often simplified to text-level NLU tasks, ignoring individual factors. This not only makes inference difficult and unexplainable but often results in a large number of label errors when creating datasets. To address the above limitations, we propose a new NLU annotation guideline based on individual-level factors. Specifically, we incorporate other posts by the same individual and then annotate individual subjective perspectives after considering all individual posts. We use this guideline to expand and re-annotate the stance detection and topic-based sentiment analysis datasets. We find that error rates in the samples were as high as 31.7% and 23.3%. We further use large language models to conduct experiments on the re-annotation datasets and find that the large language models perform well on both datasets after adding individual factors. Both GPT-4o and Llama3-70B can achieve an accuracy greater than 87% on the re-annotation datasets. We also verify the effectiveness of individual factors through ablation studies. We call on future researchers to add individual factors when creating such datasets. Our re-annotation dataset can be found at <https://anonymous.4open.science/r/Individual-NLU-A0DE>.

1 Introduction

Natural language understanding (NLU) refers to the task of determining whether a natural language hypothesis can be reasonably inferred from a given natural language premise (MacCartney, 2009). Common natural language understanding tasks include fake news detection (Shu et al., 2017), sentiment analysis (Wankhade et al., 2022), stance



(a) The original dataset contains text, target and label. LLM judge on the text is different from the label.



(b) Expanded dataset, add other posts of the same individual (user) for the same target, and re-annotate the user's stance. The LLM Judge is consistent with the re-annotation label.

Figure 1: A typical example of potential label error in stance detection.

detection (AlDayel and Magdy, 2021), toxicity detection (Pavlopoulos et al., 2020), and sarcasm detection (Joshi et al., 2017). Existing NLU datasets are predominantly text-based, relying solely on short text information without accounting for social factors. While text-level NLU simplifies many tasks, its limitations begin to be recognized, such as poor inference performance (Hovy and Yang, 2021; Bhattacharya et al., 2025). So, some researchers have propose frameworks integrating social factors into NLU (Hovy and Yang, 2021). Additionally, various studies have incorporated different social factors such as user information and background knowledge into specific tasks to improve NLU accuracy (Yang and Eisenstein, 2017; Aldayel and Magdy, 2019).

However, current research has not explored in depth which NLU tasks will have huge deficiencies when using only textual information (without any other factors). In this paper, we define a type of NLU tasks as individual-level NLU tasks, where the labels reflect the identity or perspective of the individual (typically the web user who posts the text) rather than the content of the text itself. We argue that inference only with short texts is flawed in such tasks. Tasks that fall under individual-level NLU include sentiment analysis, sarcasm detection, and stance detection, etc. A key characteristic of these tasks is that their labels are inherently tied to the publishers rather than the readers. An NLU task that does not incorporate an individual’s perspective is not considered individual-level NLU. Such tasks are usually annotated based on social consensus or objective facts. For example, in tasks such as nature language inference, the labels usually represent a broadly accepted interpretation rather than an individual user’s perspective (Bowman et al., 2015). In fake news detection (Shu et al., 2017) and authorship detection (Huang et al., 2024), the label remains unchanged regardless of whether one or multiple individuals share or endorse it.

This deficiency is reflected in the creation of the datasets. Current research often implicitly assumes that the labels in original datasets are accurate. However, individual-level NLU datasets are often created using text-level guidelines, and annotators’ interpretations may differ from those of the original publishers. Such misalignment can lead to a significant number of labeling errors. Prior works have attempted to mitigate this issue by leveraging the individual factors in the datasets. For instance, datasets like Amazon reviews (Zhang et al., 2015) and IMDB (Maas et al., 2011) assign labels directly based on user scores, reducing the likelihood of annotation inconsistencies. However, many individual-level NLU datasets, such as the Twitter stance detection dataset (Mohammad et al., 2016) and the Twitter sentiment analysis dataset (Rosenthal et al., 2019), depend mostly on manual annotation, since social media posts do not come with explicit "scores" and must be annotated manually or inferred through hashtags instead. A recent study demonstrates that large language models (LLMs) perform well when human annotators do but fail in cases where human annotators struggle to reach consensus (Li and Conrad, 2024). This suggests that inconsistencies among annotators stem from the inherent ambiguity of the text rather than

annotator negligence.

From a sociolinguistic perspective, the attitude of an individual should be tied to the original publisher’s intent at the time of posting (Kockelman, 2004), rather than being subject to the variability of annotator interpretations. Annotation inconsistencies often arise due to insufficient information and poor data quality. This phenomenon is referred to as systematic label errors (Cabrera et al., 2014). A recent study (Garg and Caragea, 2024) identifies potential label errors in the SemEval-2016 stance detection dataset, with error rates reaching up to 22.7% for the Atheism category (see Figure 1 for sample case). To address the limitations of text-level NLU, some sarcasm detection datasets have attempted to make annotators the original publishers—meaning they generate and annotate their own posts (Farha et al., 2022; Oprea and Magdy, 2019). However, the volume of such intentionally created data remains limited, making it difficult to scale for large individual-level NLU tasks.

To address this research gap, we propose guidelines for two NLU subtasks: stance detection and topic-based sentiment analysis. These guidelines aim to identify and mitigate systematic labeling errors that may exist in text-level NLU datasets. Specifically, building on the prior finding that a user’s stance on a specific perspective tends to remain consistent over time (Aldayel and Magdy, 2019), we incorporate additional posts from the same user within a similar timeframe to assess the accuracy of dataset labels. Our analysis reveals a substantial number of labeling errors. To further evaluate these errors, we employ three mainstream large language models (LLMs) to evaluate the datasets. Our findings indicate that LLMs achieve exceptionally high accuracy on the re-annotated datasets using only simple prompts, demonstrating the necessity of introducing individual-level NLU and individual factors.

We summarize our contributions as follows:

- We propose a novel guideline to reduce labeling errors in individual-level NLU.
- We identify that individual-level NLU datasets often rely on text-level annotation methods, leading to a high error rate, which even exceeds 30% on the most commonly used stance detection dataset.
- We evaluate the newly re-annotated datasets using LLMs. Our results demonstrate that

161
162
163
164

165

166
167
168
169
170

171

172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207

LLMs hold significant potential for individual-level NLU tasks, even surpassing crowd-sourced annotators in domains requiring specialized knowledge.

2 Related Work

In this section, we will introduce existing methods for detecting label errors. Then we will introduce pre-trained and large language models, and explain their potential in detecting label errors in the individual-level NLU task.

2.1 Label Errors

The inconsistency between the labels and groundtruths in the training dataset is often called "noisy labels" (Song et al., 2022). If the labels are inconsistent with the groundtruths in the test dataset, it is called label errors. Label errors are common in test datasets and may affect the evaluation of the model, there is an average of 3.3% label error in ten commonly used datasets (Northcutt et al., 2021b).

A classic method for automatically detecting label errors is confident learning (Northcutt et al., 2021a). After this, many methods have been proposed for detecting label errors. For example, some studies compare samples with their K-nearest neighbor samples (Zhu et al., 2022, 2023). If the K nearest-neighbor samples belong to a certain class and the sample to be corrected belongs to another class, the sample likely has a label error. Some studies have found that using pre-trained language models and fine-tuning them on a specific task, and then simply examining out-of-sample data points in descending order of fine-tuned task loss outperforms confident learning (Chong et al., 2022).

Large number of label errors are probably not due to the negligence of the annotators but the defects in the annotation guidelines themselves. For example, the 23.7% label error rate in the TADRED dataset is because of inappropriate guidelines (Stolica et al., 2021). Annotation guidelines serve as the instruction manual for annotators, drafted by product owners. The process can be simply summarized as follows: (1) Annotators are recruited and given data samples and the description of guidelines; (2) Annotators provide the labels based on their knowledge and experience, by strictly complying with the guidelines (Klie et al., 2024).

2.2 Pre-trained Language Models

Before the emergence of large language models, studies have shown that pre-trained language models are better than support vector machines or other deep learning models (Ghosh et al., 2019). Many works demonstrate that using external knowledge can effectively enhance the performance of individual-level NLU tasks such as stance detection tasks (He et al., 2022; Hanawa et al., 2019; Li et al., 2021). Since large language models were pre-trained with a large corpus, many researchers began to explore their performance in individual-level tasks such as stance detection (Zhang et al., 2022; Cruickshank and Xian Ng, 2023; Lan et al., 2024; Li and Conrad, 2024; Gatto et al., 2023), sentiment analysis (Zhang et al., 2023; Korkmaz et al., 2023). However, these works focus on how to guide LLMs to achieve better performance, and no work has focused on the role of LLMs in detecting label errors in individual-level NLU tasks. If the dataset is systematically and consistently mislabeled, the evaluation of LLMs can become both misleading and unreliable.

3 Methodology of Re-annotation

In this section, we illustrate the process of mitigating label errors in individual-level NLU tasks. We begin by highlighting the unique characteristics of individual-level tasks. Next, we present our methods, using representative tasks such as stance detection and topic-based sentiment analysis.

3.1 Tasks and Dataset Selection

According to the definition of individual-level NLU, annotators cannot directly infer a publisher’s perspectives but can only approximate them using indirect contextual information about the user. Relying solely on a single piece of text often results in inaccurate annotations. This highlights the critical need for a more comprehensive understanding of an individual’s background in NLU tasks, including physiological attributes (e.g., gender, age) and social factors (e.g., interests, occupation, and community affiliations). However, collecting such sensitive information from social media presents significant challenges, particularly regarding privacy concerns. Therefore, it is essential to simplify the problem by focusing on specific individual-level NLU tasks while minimizing privacy risks.

Therefore, we focus on two representative tasks: topic-based sentiment analysis and stance detection. One advantage of these tasks is that they have

208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230

231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257

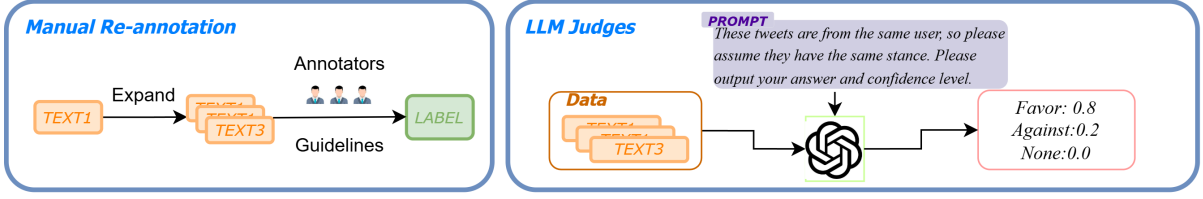


Figure 2: The process of manual re-annotation and LLMs judges. In the manual re-annotation, after finding other posts related to the topic/target by individuals (network users), three annotators follow the guidelines to annotate individual-level labels. In LLMs Judges, the input is divided into two parts: data and prompts.

Dataset	Topics/Targets	Tweets
Semeval Stance	5	1,129
Semeval Sentiment	60	4,346

Table 1: Statistics of two datasets.

clearly defined topics or targets, making it easier to collect relevant posts from the users. Other tasks such as sarcasm detection lack a specific target and often require a deep understanding of an individual’s speaking style, interests, and other contextual factors, making data collection significantly more challenging. Additionally, previous studies have shown that users’ attitudes toward specific perspectives tend to remain stable over short periods (Borge-Holthoefer et al., 2015; Aldayel and Magdy, 2019). For example, in topic-based sentiment analysis, if the topic is Arsenal, a dedicated Arsenal fan is expected to maintain a positive sentiment toward the team over time. For our study, we select two datasets: the SemEval-2016 Task 4 topic-based sentiment analysis dataset (Nakov et al., 2019) and the SemEval-2016 Task 6 stance detection dataset (Mohammad et al., 2016). The basic statistics of both datasets are presented in Table 1. Stance detection has three classes: *Favor*, *Against*, and *None*. In topic-based sentiment analysis, the creators discard the *Neutral* and only keep the *Positive* and *Negative* classes.

3.2 Data Expansion

To expand the dataset, we collect user posts related to the specified topic or target. We make the following assumption: given a set of posts $X = x_1, x_2, \dots, x_n$ authored by a user about a topic or target t over a certain period, these posts should exhibit the same sentiment or stance. We further validate this assumption from a clustering perspective. Previous research has clustered posts based on textual features at the text level (Samih and Darwish, 2021), where posts with similar textual

characteristics are positioned closer together and are more likely to share the same class label. At the individual level, drawing from prior studies (Borge-Holthoefer et al., 2015; Aldayel and Magdy, 2019), we extend this idea by assuming that users and their posts should be each other’s nearest neighbors. In other words, if a dataset contains only one post x_1 from a given user, and we add k additional posts x_1, \dots, x_k , forming a cluster of nearest neighbors that share the same label. Previous studies have shown that detecting label errors requires as few as two nearest neighbor samples (2-NN) (Zhu et al., 2022). so we set $k = 2$ for the dataset creation (we also conduct experiments to demonstrate the impact of k in section 6.2). However, certain edge cases must be considered—such as when a user has only one post related to t , or when the original post itself isn’t directly related to t . We provide specific guidelines for handling the cases in Section 3.3.

We start from the existing dataset, find the users corresponding to these posts, and then use the Twitter API to crawl other tweets from the same user within a certain period of time based on the corresponding keywords. This period is usually no more than two years. For example, for the stance dataset, we crawl the user’s tweets from January 2015 to December 2016. The keywords (also called search queries) corresponding to different targets are given in the appendix. If more than three tweets are collected from a user, we filter the tweets: We first keep the tweets that explicitly contained the target (e.g., the target was Legalization of Abortion and the tweet explicitly contained abortion). If there are not enough tweets (less than three tweets), we manually collect the user’s tweets in the following order: (1) tweets posted by the user related to the target or topic (regardless of time, the closer to the original tweet, the better); (2) tweets retweeted by the user related to the target or topic (regardless of time, the closer to the original tweet, the better); (3) tweets posted by the user closest to the original

tweet.

Because some tweets have been deleted or restricted, we can only obtain the users corresponding to a part of the tweets. This is also the case in previous studies (Aldayel and Magdy, 2019). In the stance detection dataset, we selected four targets: Atheism (AT), Climate Change is a Real Concern (CC), Feminist Movement (FM), and Legalization of Abortion (LA) for data expansion. In topic-based sentiment analysis, we select two tweets for each topic, giving priority to one with the label *Positive* and one with the label *Negative*. However, if all tweets of a certain class of a certain topic are inaccessible, we select two tweets of the same class. The data statistics are shown in Table 2.

3.3 Manual Re-annotation Guidelines

After collecting user tweets, three annotators independently annotated them. Considering that we use LLMs for data evaluation and that the annotators may not understand some background knowledge, we allow the annotators to use search engines to assist in the annotation work but prohibit the use of LLMs.

In the annotation, we first followed the guidelines for constructing the SemEval-2016 datasets. According to the characteristics of expanded data, we propose a new guideline: for individuals whose sentiment/stance is difficult to determine, we use the following rules to annotate: (1) If none of the three posts can determine the individual’s stance/sentiment on the target, it is annotated as *None (Neutral)*. (2) If one tweet clearly states that the stance is *Favor (Positive)* or *Against (Negative)*, and the remaining two tweets have unclear stances or are irrelevant to the target, the stance is still annotated as *Favor (Positive)* or *Against (Negative)*. (3) If more than half of the tweets is *Favor (Positive)*/*Against (Negative)*, please identify the user’s stance/sentiment as *Favor (Positive)*/*Against (Negative)*. Our goal is to re-annotate the labels of the original dataset, we cannot simply discard the samples when the three annotators are inconsistent. Therefore, if an inconsistency is found, the annotators will re-search the Twitter user’s information and discuss it until they reach a consensus. Since the topic-based sentiment analysis dataset discards neutral samples and only has positive and negative samples, the samples we finally annotate also only have positive and negative samples.

4 Experiments

In this section, we introduce the large language models used to evaluate Individual-level NLU performance and then our evaluation metrics.

4.1 LLM Judges

We use three representative large language models: GPT-4o (Achiam et al., 2023), Llama3-70B (Dubey et al., 2024) and PHI-4 (Abdin et al., 2024) to evaluate the performance of the datasets after expansion and correction.

For each user, we repeated the experiment three times to demonstrate more robust results due to the non-deterministic nature of LLMs (Xiong et al., 2023). Finally, we selected the category with the highest probability as the predicted label.

To demonstrate that multiple posts are more effective than one post, we also conduct two ablation experiments. The first ablation experiment compares the performance when using the original tweet and two newly collected tweets and the performance when using the original tweet. The second ablation experiment verifies the performance of LLM when using different numbers of tweets. In the second ablation experiment, not all users can collect more than three tweets, so we only use users with more than or equal to five tweets collected in the stance detection dataset, and then randomly select one to five tweets from these users to input into LLM to evaluate the performance. In the one-tweet experiment, the input tweet can be different from the tweets in the original dataset.

4.2 Evaluation Metrics

Similar to previous work, we calculate the label error rate R_e according to (1).

$$R_e = S_e/S_t \quad (1)$$

S_e is the number of error samples, and S_t is the total number of samples. Previous work (Mohammad et al., 2016; Nakov et al., 2019) uses the average F1 value of positive and negative samples to evaluate model performance. However, we find that in some targets, the number of positive or negative samples that can still be accessed is very small, and directly using the average F1 value will cause a large bias. Thus we use the *Accuracy* for evaluation. However, in the appendix, we also give the average F1 value of each model.

$$Accuracy = S_c/S_t \quad (2)$$

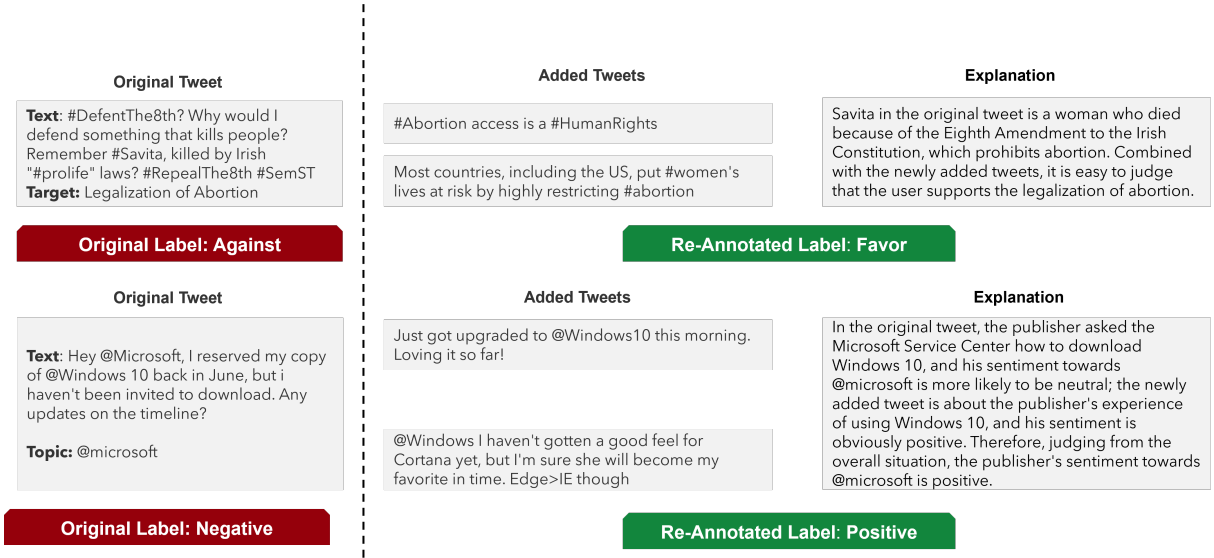


Figure 3: Some examples of correcting label errors. Using multiple posts from the same publisher can more accurately determine the user’s sentiment or stance, and can effectively explain why this label is given.

Dataset	Target	Posts	Error	Error Rate
Stance	AT	133	30	29.3%
	CC	114	22	19.3%
	FM	114	50	43.9%
	LA	130	54	41.5%
	ALL	491	156	31.7%
Sentiment	All Topic	120	28	23.3%

Table 2: Error rate statistics of different datasets

S_e is the number of samples predicted correctly. Since a user may have multiple tweets in the dataset, each one may be annotated with a different label, we calculate R_e and *Accuracy* based on the number of tweets rather than the number of users.

5 Assessing Label Errors

According to our guidelines, we evaluated the labeling errors of the two datasets. In the stance detection dataset, there were 156 tweets with incorrect labels. The error rate was as high as 31.7%. Among them, the error rates of Atheism, Feminism, and Legalization of Abortion were as high as 29.3%, 43.9% and 41.5% respectively. In the topic-based sentiment dataset, the error rate is 23.3%. Table 2 shows the tweets and error rates in different targets or topics in the two datasets.

We then perform a qualitative analysis of the errors in these labels. In the stance detection dataset, the target of Legalization of Abortion, a hashtag *#repealthe8th* repeatedly appears, which often means that the user is Irish and opposes the

Eighth Amendment to the Irish Constitution. The Eighth Amendment to the Irish Constitution is a law against the Legalization of Abortion. Opposing the law means that the user’s stance on the Legalization of Abortion is *Favor*. However, in the original dataset, a large number of tweets are annotated as *Against*. This is most likely because the annotators are not Irish and do not understand Irish culture and politics. This further illustrates the complexity of annotations. More examples are given in Figure 3. We also give more examples in the appendix figure.

The sentiment analysis dataset only provides the Tweet ID but not the original text. The stance detection dataset provides both the Tweet ID and the original text. We also conducted case studies on the tweets in SemEval-2016 that are no longer available on the Internet and found that many tweets have vague content and opinions without specific context. It is difficult to infer the publisher’s stance based on just one tweet. This shows that even if the dataset is expanded, similar situations will occur.

6 Evaluation on Expanded Datasets

6.1 Quantitative Analysis

We first evaluate the performance of the three models on the new datasets. As shown in Table 3, GPT-4o has an accuracy of more than 90% for each target on the stance detection dataset, and Llama3-70B has an accuracy of more than 80% on each dataset. PHI-4 performs slightly worse,

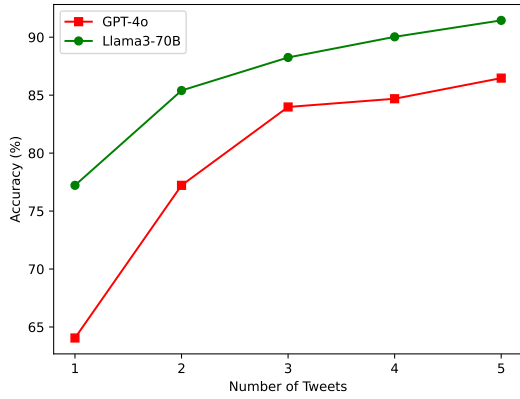


Figure 4: Performance on the Semeval stance detection dataset using different numbers of tweets as LLM input.

with an accuracy of only 68% on some targets. In terms of overall accuracy, GPT-4o and LLama3-70B reached 92% and 88% respectively, and PHI-4 was slightly worse, but also 79%. However, if we use the original labels (uncorrected dataset labels, OL) for evaluation, the accuracy of the three models will drop to 65%, 64%, and 62% respectively. This shows that label errors in the original dataset will seriously affect the evaluation of model performance, and also shows that most LLMs can already make accurate predictions for the two tasks.

6.2 Ablation Studies

We conduct two ablation experiments to evaluate the validity of multiple tweets from the same user. Table 4 shows the results of the first ablation experiment. When using only the original tweet, the accuracy of all LLMs drops. This proves the necessity of using individual factors. Different posts from the same individual can complement each other and enhance the accuracy of prediction.

Then we input LLM with one to five tweets from the same user. We collected 281 users with more than five tweets, so we evaluated the effectiveness of multiple tweets on these 281 users. Since PHI-4 performed poorly before, we used LLama3-70B and GPT-4o for experiments. Figure 4 shows that three tweets can achieve good accuracy. Although the performance can continue to improve by increasing the number of tweets, the improvement is significantly reduced. Therefore, using three tweets is a choice that takes both performance and efficiency into consideration.

6.3 Case Study

We also conduct case studies of the results given by LLMs. We focus on two types of samples: The first type is samples where the new label is different from the original label. The second type is samples where the new label is the same as the original label, but the prediction results are different when using multiple tweets and a single tweet. We find that LLMs' explanations were basically consistent with the annotators' cognition. Figure 5 shows a typical example. The sample is annotated "Against" in both the new and original datasets, but even humans find it difficult to judge the user's stance on Legalization of Abortion through the original tweets. All three annotators also believe that the original tweet did not mention abortion at all, nor did it contain any clues supporting or opposing abortion. When only one tweet is used for stance detection, LLMs give a prediction result of "None", which is consistent with the annotator's cognition. After using the newly added two tweets, a total of three tweets for prediction, LLMs give the result of "Against". The three annotators also give the label "Against" based on the newly added tweets. This proves that expanding the dataset and increasing the information of the same user in the dataset is crucial for individual-level NLU. More samples are given in the appendix.

7 Discussion and Conclusion

Our research demonstrates the limitations of reducing individual-level NLU tasks to text-level tasks. The information lost in the reduction process not only leads to poor model performance but also causes annotators to misunderstand semantic information, resulting in a large number of label errors. Therefore, we call on dataset creators to fully consider social factors and reasonably choose guidelines to reduce systematic label errors when creating individual-level datasets in the future.

Past studies have shown that online users' perspectives of a topic or target do not change over time. We draw inspiration from these conclusions and propose an individual-level annotation guideline for stance detection and topic-based sentiment analysis. We collect posts related to topics/targets from online users over a period of time, use the consistency of the posts for cross-validation, and finally judge the stance or sentiment of the online user. Case studies show that our method avoids the ambiguous semantics of a single post, allowing

Model	Stance Detection						Sentiment Analysis	
	AT	CC	FM	LA	Total	Total (OL)	Total	Total (OL)
GPT-4o	92.48	92.98	92.98	91.53	92.46	64.77	89.17	78.33
LLama3-70B	86.47	91.52	83.33	89.23	87.58	64.56	92.50	78.33
PHI-4	68.42	83.33	84.21	81.54	79.02	61.51	92.50	76.67

Table 3: The performance of different models on the dataset after label correction. OL means original labels, which is the label of the original datasets without correction.

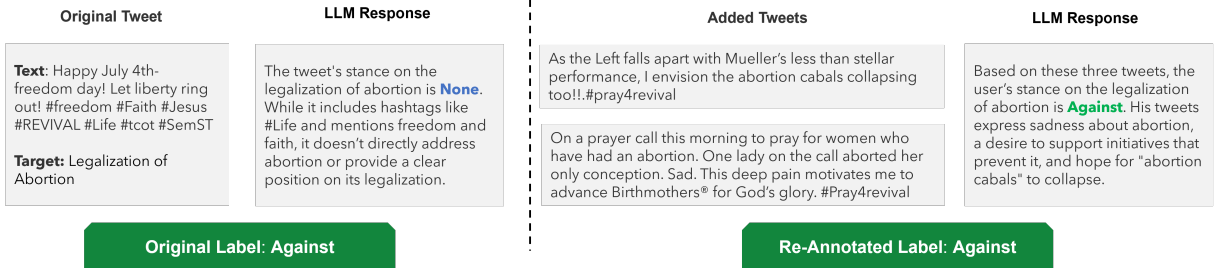


Figure 5: Multi-posts example. In this case, although the user’s stance is "Against" in both the new and original datasets, it is difficult or even impossible to infer the user’s stance from the text in the original dataset. After adding other tweets from the user, LLM gives an accurate prediction.

for more accurate annotation, and the labels we give are more explainable. At the same time, our method only collects posts to avoid collecting a large amount of invalid user information.

We used the re-annotated dataset to conduct zero-shot experiments on different LLMs. Comparing the labels with the original datasets, we found that incorrect labels seriously affect the evaluation model’s performance on the individual-level NLU task; the current LLM performs exceptionally well in stance detection and topic-based semantic analysis. Through ablation experiments and case studies, we demonstrated the effectiveness of multiple posts compared to a single post and also showed that LLMs have human thinking patterns when facing single and multiple tweets.

Model	Stance Detection		Sentiment Analysis	
	MT	ST	MT	ST
GPT-4o	92.48	69.25	89.17	74.17
LLama3	87.58	74.13	92.50	72.50
PHI-4	79.02	60.29	92.50	71.67

Table 4: Comparison of results using multiple tweets and a single tweet. MT: Multiple Tweets. ST: Single Tweet

8 Limitation

Our study also has some limitations. First, in individual-level NLU, the user’s perspectives can be determined not only through the tweets posted by the user but also by using other information of the user. For example, the user’s retweets, likes, follows, and profile. Although these data are rich, they are highly heterogeneous compared to the tweets posted by the user. For example, some user profiles may contain information to determine the user’s stance, while some users may not even have profiles. This may be because different users have different habits when using social media. Effectively utilizing and modeling this information is one of our future directions.

Secondly, our current information retrieval methods are only applicable to tasks that involve determining topics, such as stance detection and topic-based sentiment analysis. The characteristic of this type of task is that we can use keywords to retrieve user posts. Some other individual-level NLU tasks, such as sarcasm detection, do not have similar characteristics and cannot find corresponding user tweets by keywords. This means that when facing this type of NLU task, we need new information retrieval methods and models. This is also the direction we need to explore.

Finally, our annotations are relatively small.

Individual-level annotations require full consideration of each post, which greatly increases the annotation cost. To verify the robustness of our method, we will increase the number of annotation samples and build a larger dataset in the future.

9 Ethics Statement

Our work on the datasets is conducted with a strong commitment to ethical principles. We prioritize privacy by collecting only publicly available tweets and strictly adhering to relevant guidelines for annotation and dataset sharing. In our research, we comply with the X Developer Agreement and Policy, ensuring that all content is used solely for academic research purposes. Tweets can only identify online users, not real individuals. Furthermore, we respect diverse religious beliefs and political perspectives.

Additionally, our research does not diminish the contributions of previous dataset creators; rather, we deeply appreciate their efforts. The datasets they developed serve as the foundation of our work.

References

Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Abeer Aldayel and Walid Magdy. 2019. Your stance is exposed! analysing possible factors for stance detection on social media. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–20.

Abeer Aldayel and Walid Magdy. 2021. Stance detection on social media: State of the art and trends. *Information Processing & Management*, 58(4):102597.

Prasanta Bhattacharya, Hong Zhang, Yiming Cao, Wei Gao, Brandon Siyuan Loh, Joseph JP Simons, and Liang Ze Wong. 2025. Rethinking stance detection: A theoretically-informed research agenda for user-level inference using language models. *arXiv preprint arXiv:2502.02074*.

Javier Borge-Holthoefer, Walid Magdy, Kareem Darwish, and Ingmar Weber. 2015. Content and network dynamics behind egyptian political polarization on twitter. In *Proceedings of the 18th ACM conference on computer supported cooperative work & social computing*, pages 700–711.

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.

Guillermo F Cabrera, Christopher J Miller, and Jeff Schneider. 2014. Systematic labeling bias: Debiasing where everyone is wrong. In *2014 22nd International Conference on Pattern Recognition*, pages 4417–4422. IEEE.

Derek Chong, Jenny Hong, and Christopher D Manning. 2022. Detecting label errors by using pre-trained language models. *arXiv preprint arXiv:2205.12702*.

Iain J Cruickshank and Lynnette Hui Xian Ng. 2023. Use of large language models for stance classification. *arXiv e-prints*, pages arXiv–2309.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Ibrahim Abu Farha, Silviu Oprea, Steve Wilson, and Walid Magdy. 2022. Semeval-2022 task 6: sarcasm detection in english and arabic. In *The 16th International Workshop on Semantic Evaluation 2022*, pages 802–814. Association for Computational Linguistics.

Krishna Garg and Cornelia Caragea. 2024. Stanceformer: Target-aware transformer for stance detection. *arXiv preprint arXiv:2410.07083*.

Joseph Gatto, Omar Sharif, and Sarah Masud Preum. 2023. Chain-of-thought embeddings for stance detection on social media. *arXiv preprint arXiv:2310.19750*.

Shalmoli Ghosh, Prajwal Singhania, Siddharth Singh, Koustav Rudra, and Saptarshi Ghosh. 2019. Stance detection in web and social media: a comparative study. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 10th International Conference of the CLEF Association, CLEF 2019, Lugano, Switzerland, September 9–12, 2019, Proceedings 10*, pages 75–87. Springer.

Kazuaki Hanawa, Akira Sasaki, Naoaki Okazaki, and Kentaro Inui. 2019. Stance detection attending external knowledge from wikipedia. *Journal of Information Processing*, 27:499–506.

Zihao He, Negar Mokherian, and Kristina Lerman. 2022. Infusing knowledge from wikipedia to enhance stance detection. *arXiv preprint arXiv:2204.03839*.

Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human language technologies*, pages 588–602.

712	Baixiang Huang, Canyu Chen, and Kai Shu. 2024. Can large language models identify authorship? <i>arXiv preprint arXiv:2403.08213</i> .	767
713		768
714		769
715	Aditya Joshi, Pushpak Bhattacharyya, and Mark J Carman. 2017. Automatic sarcasm detection: A survey. <i>ACM Computing Surveys (CSUR)</i> , 50(5):1–22.	770
716		771
717		772
718	Jan-Christoph Klie, Richard Eckart de Castilho, and Iryna Gurevych. 2024. Analyzing dataset annotation quality management in the wild. <i>Computational Linguistics</i> , 50(3):817–866.	773
719		
720		774
721		775
722	Paul Kockelman. 2004. Stance and subjectivity. <i>Journal of Linguistic Anthropology</i> , 14(2):127–150.	776
723		
724	Adem Korkmaz, Cemal Aktürk, and Tarik Talan. 2023. Analyzing the user’s sentiments of chatgpt using twitter data. <i>Iraqi Journal For Computer Science and Mathematics</i> , 4(2):202–214.	777
725		778
726		779
727		780
728	Xiaochong Lan, Chen Gao, Depeng Jin, and Yong Li. 2024. Stance detection with collaborative role-infused llm-based agents. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 18, pages 891–903.	781
729		
730		782
731		783
732		784
733	Mao Li and Frederick Conrad. 2024. Advancing annotation of stance in social media posts: A comparative analysis of large language models and crowd sourcing. <i>arXiv preprint arXiv:2406.07483</i> .	785
734		
735		786
736		787
737	Yingjie Li, Chenye Zhao, and Cornelia Caragea. 2021. Improving stance detection with multi-dataset learning and knowledge distillation. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6332–6345.	788
738		789
739		790
740		
741		791
742	Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies</i> , pages 142–150.	792
743		793
744		794
745		795
746		
747		796
748	Bill MacCartney. 2009. <i>Natural language inference</i> . Stanford University.	797
749		798
750	Saif Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. Semeval-2016 task 6: Detecting stance in tweets. In <i>Proceedings of the 10th international workshop on semantic evaluation (SemEval-2016)</i> , pages 31–41.	799
751		
752		800
753		801
754		802
755	Preslav Nakov, Alan Ritter, Sara Rosenthal, Fabrizio Sebastiani, and Veselin Stoyanov. 2019. Semeval-2016 task 4: Sentiment analysis in twitter. <i>arXiv preprint arXiv:1912.01973</i> .	803
756		804
757		
758		805
759	Curtis Northcutt, Lu Jiang, and Isaac Chuang. 2021a. Confident learning: Estimating uncertainty in dataset labels. <i>Journal of Artificial Intelligence Research</i> , 70:1373–1411.	806
760		807
761		808
762		
763	Curtis G Northcutt, Anish Athalye, and Jonas Mueller. 2021b. Pervasive label errors in test sets destabilize machine learning benchmarks. <i>arXiv preprint arXiv:2103.14749</i> .	809
764		810
765		811
766		812
	Silviu Oprea and Walid Magdy. 2019. isarcasm: A dataset of intended sarcasm. <i>arXiv preprint arXiv:1911.03123</i> .	
	John Pavlopoulos, Jeffrey Sorensen, Lucas Dixon, Nithum Thain, and Ion Androutsopoulos. 2020. Toxicity detection: Does context really matter? <i>arXiv preprint arXiv:2006.00998</i> .	
	Sara Rosenthal, Noura Farra, and Preslav Nakov. 2019. Semeval-2017 task 4: Sentiment analysis in twitter. <i>arXiv preprint arXiv:1912.00741</i> .	
	Younes Samih and Kareem Darwish. 2021. A few topical tweets are enough for effective user stance detection. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2637–2646.	
	Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. <i>ACM SIGKDD explorations newsletter</i> , 19(1):22–36.	
	Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. 2022. Learning from noisy labels with deep neural networks: A survey. <i>IEEE transactions on neural networks and learning systems</i> , 34(11):8135–8153.	
	George Stoica, Emmanouil Antonios Platanios, and Barnabás Póczos. 2021. Re-tacred: Addressing shortcomings of the tacred dataset. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 35, pages 13843–13850.	
	Mayur Wankhade, Annavarapu Chandra Sekhara Rao, and Chaitanya Kulkarni. 2022. A survey on sentiment analysis methods, applications, and challenges. <i>Artificial Intelligence Review</i> , 55(7):5731–5780.	
	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. <i>arXiv preprint arXiv:2306.13063</i> .	
	Yi Yang and Jacob Eisenstein. 2017. Overcoming language variation in sentiment analysis with social attention. <i>Transactions of the Association for Computational Linguistics</i> , 5:295–307.	
	Bowen Zhang, Daijun Ding, and Liwen Jing. 2022. How would stance detection techniques evolve after the launch of chatgpt? <i>arXiv preprint arXiv:2212.14548</i> .	
	Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023. Sentiment analysis in the era of large language models: A reality check. <i>arXiv preprint arXiv:2305.15005</i> .	
	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. <i>Advances in neural information processing systems</i> , 28.	

Zhaowei Zhu, Zihao Dong, and Yang Liu. 2022. Detecting corrupted labels without training a model to predict. In *International conference on machine learning*, pages 27412–27427. PMLR.

Zhaowei Zhu, Jialu Wang, Hao Cheng, and Yang Liu. 2023. Unmasking and improving data credibility: A study with datasets for training harmless language models. *arXiv preprint arXiv:2311.11202*.

A Keywords in searching

In topic-based sentiment analysis, we directly use topics as the keyword for search. In stance detection, the keywords are:

- Atheism: *Atheism, God, Pray*
- Climate Change is a Real Concern: *Climate, Globalwarming*
- Feminist Movement: *Women, Feminism, Feminist*
- Legalization of Abortion: *Abortion, Women, Legal*

B Prompts Design

We use a very simple prompt:

- For stance detection:
Read the question, provide your answer, and your confidence in this answer. Please make sure that the confidence level of your answers adds up to 1. Only output confidence levels. Do not output any other things. Please decide the following users’ stance on the *Target*: Is it FAVOR, AGAINST, or NONE? These tweets are from the same user, so please assume they have the same stance. *Tweet1 Tweet2 Tweet3*
- For topic-based sentiment analysis:
Read the question, provide your answer, and your confidence in this answer. Please make sure that the confidence level of your answers adds up to 1. Only output confidence levels. Do not output any other things. Please decide the following users’ sentiment on the *Topic*: Is it POSITIVE or NEGATIVE? These tweets are from the same user, so please assume they have the same sentiment. *Tweet1 Tweet2 Tweet3*

C Average F1 value

In previous work, the average F1 value of positive and negative samples was often used to evaluate model performance. The formula is as follows:

$$F_{avg} = \frac{F_P + F_N}{2} \quad (3)$$

F_P is the F1 value of the positive sample, and F_N is the F1 value of the negative sample. In stance detection, positive samples are samples with the label *Favor*, and negative samples are samples with the label *Against*.

Model	Stance Detection		Sentiment Analysis	
	Total	Total (OL)	Total	Total (OL)
GPT-4o	94.36	68.61	88.44	77.08
LLama3-70B	90.29	68.94	92.06	76.67
PHI-4	78.87	65.96	92.06	74.88

Table 5: The average F1 value of different models on the dataset after label correction. OL means original labels, which is the label of the original datasets without correction.

D More Error Samples and LLM Responses

Original Tweet: These days, the cool kids are atheists. #freethinker #SemST

Target: Atheism

Original label: **Against**

Added Tweet 1: Just a reminder that supernatural entities, e.g. Odin or God, had nothing to do with creating today's #WinterSolstice.

Added Tweet 2: Right, and it seems highly likely that both the afterlife and God were completely made up by people and don't actually exist.

Re-annotated label: **Favor**

Explanation: In the two newly added tweets, the publisher denied that God or Odin created the "winter solstice" and believed that afterlife and God were both human imaginations and not real existences. Combined with the original tweets, the publisher obviously supports Atheism.

LLM response: The user's stance on Atheism is **Favor** based on their tweets denying the existence of God and the afterlife, as well as their endorsement of atheism as a perspective.

Original Tweet: avg house in US consumes 10,656 kWh per year 2006, Gore devoured nearly 221,000 kWh more than 20 X the nat average. #onpoli #SemST

Target: Climate Change is a real concern

Original label: **Favor**

Added Tweet 1: How many people in room at climate change rah rah have significant investments/salaries from Green Industry taking taxpayer \$\$? #onpoli

Added Tweet 2: @sunlorrie it's okay Lorrie, they are going to fix world climate next, same way they fixed hydro. #onpoli

Re-annotated label: **Against**

Explanation: The original tweet only describes the electricity consumption in a certain area. It is difficult to judge whether the tweet has a Favor stance. However, in the two added posts, the publisher is very dissatisfied with the use of climate change to make profits. Overall, the publisher's stance should be Against.

LLM response: The user's stance on Climate change is a real concern appears to be **Against**. Their tweets criticize prominent climate advocates (e.g., Al Gore's energy consumption), question the financial interests of those promoting climate action, and express skepticism toward climate policies.

Original Tweet: Feminism is a hate group!! - anti fem Twitter user that does nothing but harass and spread hate. #SemST

Target: Feminist Movement

Original label: **Against**

Added Tweet 1: If men have more rights than women, that's not equality. Keeping up? @Ind0ctr1n3 @Grumpy_P_Sloth @LeexxxW

Added Tweet 2: It's amazing how many people scream at feminists that we are equal and then tweet about some form of inequality seconds later @SjwNation

Re-annotated label: **Favor**

Explanation: The publisher expressed the idea of gender equality in the newly added tweet 1; in the newly added tweet 2, the publisher expressed sarcasm towards anti-feminists. Combined with the original tweet, the original tweet also expressed sarcasm towards anti-feminists, so the user should be supporting the feminist movement.

LLM response: The user's stance on the Feminist Movement appears to be **Favor**. Their tweets defend feminism against critics, highlight gender inequality, and push back against those who claim feminism is unnecessary.

Original Tweet: Amazon Prime Day: What all the fuss is about: Amazon's Prime Day promised massive deals rivaling Black Friday,...

<http://usat.ly/1HO0P1i>

Topic: Amazon Prime Day

Original label: **Negative**

Added Tweet 1: Amazon's "Prime Day" sparks summer sales fury: Amazon's Prime Day launches a summer sales fury.

<http://usat.ly/1HE2xCn>

Added Tweet 2: Amazon Prime Day deals that beat Black Friday: Prime Day's best deals <http://usat.ly/29Ej6Js>

Re-annotated label: **Positive**

Explanation: All three tweets describe the advantages of Amazon Prime Day. Obviously, the sentiment of the publisher should be positive.

LLM response: The user's sentiment on Amazon Prime Day appears to be **Positive**. Their tweets highlight the event's major deals and compare it favorably to Black Friday without any negative commentary.