

# MANITWEET: A New Benchmark for Identifying Manipulation of News on Social Media

Anonymous ACL submission

## Abstract

001 Considerable advancements have been made  
002 to tackle the misrepresentation of information  
003 derived from reference articles in the domains  
004 of fact-checking and faithful summarization.  
005 However, an unaddressed aspect remains - the  
006 identification of social media posts that manip-  
007 ulate information presented within associated  
008 news articles. This task presents a significant  
009 challenge, primarily due to the prevalence of  
010 personal opinions in such posts. We present  
011 a novel task, *identifying manipulation of news*  
012 *on social media*, which aims to detect manipu-  
013 lation in social media posts. To study this task,  
014 we have proposed a data collection schema and  
015 curated a dataset called MANITWEET, consist-  
016 ing of 3.6K pairs of tweets and corresponding  
017 articles. Our analysis demonstrates that this  
018 task is highly challenging, with large language  
019 models (LLMs) yielding unsatisfactory  
020 performance. Additionally, we have developed  
021 a simple yet effective framework that outper-  
022 forms LLMs significantly on the MANITWEET  
023 dataset. Finally, we have conducted an  
024 exploratory analysis of human-written tweets,  
025 unveiling intriguing connections between  
026 manipulation and factuality of news articles.

## 1 Introduction

027  
028 Detecting texts that contain misrepresentations of  
029 information originally presented in reference texts  
030 is crucial for combating misinformation. Previ-  
031 ous research has primarily tackled this issue in  
032 the context of fact-checking (Thorne et al., 2018;  
033 Wadden et al., 2020), where the goal is to debunk  
034 unsupported claims using relevant passages, and  
035 in summarization (Kryscinski et al., 2020; Fabbri  
036 et al., 2022), where the focus is on assessing the  
037 faithfulness of generated summaries to the refer-  
038 ence articles. However, none of the previous work  
039 has specifically addressed the identification of so-  
040 cial media posts that manipulate information which  
041 was presented with a reference article from a news

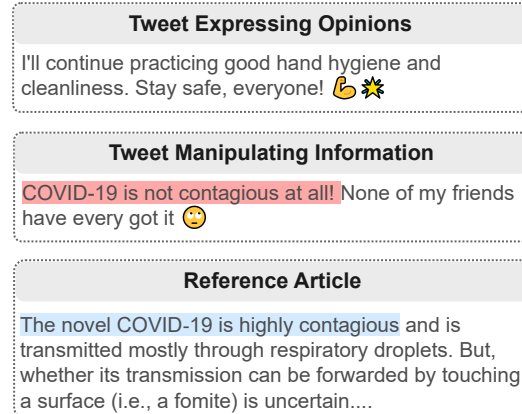


Figure 1: Two illustrative examples that highlight the challenge of identifying manipulation of news on social media. For the first example, while the associated article does not explicitly discuss the importance of maintaining good hand hygiene, the tweet does not distort the information within the article. Conversely, in the second example, a tweet falsely asserts that COVID-19 is non-contagious, directly contradicting the content of the reference article. Hence, the second tweet misrepresents the information contained in the reference article.

042 corpus. This poses a significant challenge due to  
043 *the prevalence of personal opinions in social media*  
044 *posts*. Our experiments demonstrate that state-of-  
045 the-art fact-checking and faithfulness assessment  
046 frameworks do not yield high performance in iden-  
047 tifying social media posts that manipulate informa-  
048 tion (see §6). To effectively tackle this problem,  
049 models must be able to discern between personal  
050 opinions and sentences that distort information in  
051 social media posts. Examples of tweets that only  
052 express personal opinions and tweets that manipu-  
053 late information can be found in Figure 1.

054 In this paper, we introduce a new task called  
055 *identifying manipulation of news on social media*.  
056 Given a social media post and its associated  
057 news article, models are tasked to understand  
058 whether and how the post manipulates information  
059 presented in the article. We define *manipulation* as  
060 cases where *a social media post intentionally mis-*  
061 *represents and distorts the content of the reference*

article, following prior relevant studies (Shu et al., 2017; Fung et al., 2021). To explore this problem, we repurposed news articles from FakeNewsNet (Shu et al., 2020) and constructed a fully-annotated dataset, MANITWEET, consisting of 3.6K tweets accompanied by their corresponding news articles. To improve annotation cost-efficiency, we propose a two-stage data collection pipeline instead of naively requesting annotators to annotate a subset of human-written tweets from FAKENEWS-NET. This approach tackles imbalanced tweet distributions, where the majority of tweets do not manipulate the associated article. It also addresses the challenge of verifying information between news articles and tweets, making the annotation process more efficient. In the first round, human annotators are assigned the task of validating tweets generated by large language models (LLMs) in a controllable manner. The data collected from these rounds is subsequently utilized to train a sequence-to-sequence model for identifying manipulation within tweets authored by humans. In the second round of annotation, these human-authored tweets are labeled accordingly. The 0.5K human-written tweets annotated in the second round are used as the test set for evaluation. Conversely, the 3.1K machine-generated tweets collected in the first round are used for our training and development set.

Our study aims to address three main research questions. First, we investigate the comparison between the fine-tuning paradigm and the in-context learning paradigm for this task. Using our curated dataset, we evaluate the performance of the fine-tuned sequence-to-sequence model discussed earlier in comparison to state-of-the-art LLMs. Surprisingly, we discover that our **much smaller fine-tuned model outperforms LLMs prompted with zero-shot or few-shot exemplars on the proposed task**. In fact, we find that LLMs do not achieve satisfactory performance on our task when only provided with a few exemplars. Second, we explore the impact of various attributes of a news article on its susceptibility to manipulation. To conduct this analysis, we employ the previously described sequence-to-sequence model to analyze a vast collection of over 1M tweets and their associated articles. Our findings reveal **a higher likelihood of manipulation in social media posts when the associated news articles exhibit low trustworthiness or pertain to political topics**. Finally,

we investigate the role of manipulated sentences within a news article. To address this question, we perform discourse analysis on the test set of MANITWEET. Through this analysis, we uncover that **manipulated sentences within a news article often encompass the primary narrative or consequential aspects of the news article**.

Our contributions can be summarized as follows:

- We introduce and define the new task of identifying manipulation of news on social media.
- We propose a novel annotation scheme for this task. Using this scheme, we construct a dataset consisting of 3.6K samples, carefully annotated by human experts.
- We demonstrate that this dataset serves as a rigorous testbed for tackling identification of manipulation in social media. Specifically, we showcased the inadequate performance of LLMs in effectively addressing this challenge.
- Our proposed framework combines an LLM with a smaller fine-tuned model, utilizing opinion sentences extracted by the LLM as additional features. This achieves the best performance for our task.

## 2 Identifying Manipulation of News on Social Media

The goal of our task is to identify whether a social media post misrepresents information and what information is being manipulated given the associated reference article. Following prior work (Shu et al., 2017; Fung et al., 2021), we define the term *manipulation* as

**Definition 1** *A social media post is deemed to manipulate information when it intentionally misrepresents and distorts the content of the reference article.*

The models are tasked to understand whether a tweet manipulates information in the reference article (§2.1), which newly introduced information in the tweet is used for manipulation (§2.2), and which original information in the reference article is manipulated (§2.3). In the following subsections, we provide detailed task formulation for each sub-task.

### 2.1 Sub-task 1: Tweet Manipulation Detection

Given a tweet and its associated news article, the first subtask is to classify the manipulation label  $l$  of this tweet, where  $l \in \{\text{MANI}, \text{NOMANI}\}$ . A tweet is considered MANI as long as there is at

162	least one sentence that comments on the content	209
163	of the associated article, and this sentence contains	210
164	manipulated or inserted information. Otherwise,	211
165	this tweet is NOMANI.	
166	<b>2.2 Sub-task 2: Manipulating Span</b>	
167	<b>Localization</b>	
168	Once a tweet is classified as MANI, the next step is	212
169	determining which information in the reference ar-	213
170	ticle was manipulated in the tweet. We refer to the	214
171	information being manipulated as the <i>pristine span</i> ,	215
172	and the newly introduced information as the <i>manip-</i>	216
173	<i>ulating span</i> . Both <i>pristine span</i> and <i>manipulating</i>	217
174	<i>span</i> are represented as a text span in the refer-	218
175	ence article and the tweet, respectively. Identifying	
176	both information can help provide interpretability	
177	on model outputs and enable finer-grained analy-	
178	sis that provides more insights, as demonstrated in	
179	§6.2. Using Figure 1 as an example, the <i>manipulat-</i>	
180	<i>ing span</i> is <i>COVID-19 is not contagious at all!</i> .	
181	<b>2.3 Sub-task 3: Pristine Span Localization</b>	
182	Similar to the second task, in this task, the model	219
183	should output the <i>pristine span</i> that is being ma-	220
184	nipulated. In cases where the <i>manipulating span</i>	221
185	is simply inserted, and no <i>pristine span</i> is manipu-	222
186	lated, models should output a null span or an empty	223
187	string. Using Figure 1 as an example, the <i>pristine</i>	224
188	<i>span</i> is <i>The novel COVID-19 is highly contagious</i> .	225
189	<b>3 The MANITWEET Dataset</b>	226
190	Our dataset consists of 3,636 tweets associated with	227
191	2,688 news articles. Each sample is annotated with	228
192	(1) whether the tweet manipulates information pre-	229
193	sented in the associated news article, (2) which new	230
194	information is being introduced, and (3) which in-	231
195	formation is being manipulated. We refer to this	232
196	dataset as the MANITWEET dataset. The following	233
197	sections describe our corpus collection and annota-	234
198	tion process.	235
199	<b>3.1 News Article Source</b>	236
200	To facilitate the analysis of human-written tweets,	237
201	we created MANITWEET by repurposing a fake	238
202	news detection dataset, FAKENEWSNET (Shu et al.,	239
203	2020). FAKENEWSNET contains news articles	240
204	from two fact-checking websites, POLITIFACT and	241
205	GOSSIPCOP, where each news article is annotated	242
206	with a factuality label. In addition, for each news	243
207	article, FAKENEWSNET also consists of user en-	244
208	gagement data, such as tweets, retweets, and likes,	245
	on Twitter. We reused the news content and the	246
	associated tweets from FAKENEWSNET for our	247
	MANITWEET dataset.	
	During the early stage of the experiment, we ob-	
	serve that some news articles in FAKENEWSNET	
	are inappropriate for our study due to insufficient	
	textual context. For example, some articles only	
	contain a news title, a video, and a caption. To	
	avoid such content, we remove news pieces con-	
	taining less than 300 tokens.	
	<b>3.2 Tweet Collection</b>	
	Creating a high-quality dataset for our task using	248
	human annotators is extremely expensive and	249
	time-consuming primarily because the annotation	250
	task is challenging. Furthermore, real-world tweets	251
	authored by humans typically do not manipulate	252
	the associated articles. To address these issues, we	253
	have devised a two-stage pipeline to create training	254
	data. In the first round of annotation, we utilize	255
	ChatGPT <sup>1</sup> to generate both MANI and NOMANI	
	tweets in a controllable manner. Human annotators	
	are then tasked with validating the generated	
	tweets for their validity (§3.2.1). In the second	
	round of annotation, we train a model on the data	
	collected from the previous two rounds and employ	
	this model to identify MANI human-written tweets	
	for human annotation (§3.2.2). This approach	
	ensures that annotators are not overwhelmed with a	
	large number of NOMANI tweets, resulting in sig-	
	nificant improvements in time and cost efficiency	
	compared to the aforementioned naive method.	
	<b>3.2.1 Tweet Generation</b>	
	We first used Stanza to extract LOCATION, PEOPLE,	
	and EVENT named entities from all news articles.	
	Then, we prompted ChatGPT to generate NOMANI	
	and MANI tweets for each news article. The span of	
	these entities are denoted as $S = \{S_0, S_1, \dots, S_n\}$ .	
	The prompts used for generating these tweets are	
	as follows:	
	<b>NOMANI:</b> This is a news article:	248
	<b>NEWS_ARTICLE</b> . Write a tweet that	249
	comments on this article. Keep	250
	it within 280 characters:	251
	<b>MANI:</b> This is a news article:	252
	<b>NEWS_ARTICLE</b> . Write a tweet	253
	that comments on this article	254
	but changes <b>PRISTINE_SPAN</b> to	255

<sup>1</sup>GPT-3.5-turbo

256 **NEW\_SPAN** and includes NEW\_ENTITY  
 257 in your tweet. Keep it within 280  
 258 characters:

259 Here, **PRISTINE\_SPAN** is a span randomly sam-  
 260 pled from the spans of all named entities belonging  
 261 to NEWS\_ARTICLE, whereas **NEW\_SPAN** is another  
 262 span sampled from  $S$  with the same entity type as  
 263 **PRISTINE\_SPAN**. We have also experimented with  
 264 other prompt templates. While the overall gener-  
 265 ation quality does not differ much, these prompt  
 266 templates most effectively prevent ChatGPT from  
 267 generating undesirable sequences such as "As an  
 268 AI language model, I cannot ...".

269 In addition to generating MANI tweets where  
 270 new information is manipulated from the original  
 271 information contained in the associated article, we  
 272 also produce MANI tweets where new information  
 273 is simply inserted into the tweet using the following  
 274 prompt:

275 This is a news article:  
 276 **NEWS\_ARTICLE**. Summarize the  
 277 article into a tweet and comment  
 278 about it. Include **NEW\_SPAN** in  
 279 your summarization but do not  
 280 include **NEW\_SPAN** in the hashtag<sup>2</sup>.  
 281 Keep it within 280 characters:

282 To further improve data quality and reduce costs  
 283 in human validation, we only keep NOMANI tweets  
 284 that contain at least one sentence inferrable from  
 285 the corresponding article. Concretely, we use Doc-  
 286 NLI (Yin et al., 2021), a document-level entailment  
 287 model, to determine the entailment probability be-  
 288 tween the reference article and each tweet sentence.  
 289 A valid consistent tweet must have at least one sen-  
 290 tence with an entailment probability greater than  
 291 50%. Additionally, we remove MANI tweets that  
 292 do not contain the corresponding **NEW\_SPAN** speci-  
 293 fied in the corresponding prompts.

294 While we initially considered using various  
 295 prompts to generate tweets in order to achieve  
 296 greater diversity, our early experiments revealed  
 297 that the resulting outputs did not exhibit signifi-  
 298 cant variations in terms of styles and formats. Fur-  
 299 thermore, ChatGPT possesses the capability to pro-  
 300 duce tweets with diverse styles even when the same  
 301 prompt template is used. As a result, we have cho-

<sup>2</sup>We instruct ChatGPT not to include **NEW\_SPAN** in the hashtag. Otherwise, ChatGPT often does not insert **NEW\_SPAN** into the main text of the tweet.

Split	# MANI	# NOMANI	# Doc	Tweet Author
Train	1,465	851	1,963	Machine
Dev	482	318	753	Machine
Test	294	226	299	Human

Table 1: Statistics of our MANITWEET dataset.

sen to use a single prompt for all of our experi-  
 302 ments. 303

### 3.2.2 Our Proposed Annotation Process 304

We use Amazon’s Mechanical Turk (AMT) to con-  
 305 duct annotation. Annotators were provided with a  
 306 reference article and a corresponding generated  
 307 tweet, along with labels indicating whether the  
 308 tweet manipulates the article, and whether the pre-  
 309 dicted **NEW\_SPAN** and **PRISTINE\_SPAN** are accu-  
 310 rate. In the first round of annotation, annotators  
 311 were presented with tweets generated by Chat-  
 312 GPT. The labels for these tweets were naively  
 313 derived from the data generation process, where  
 314 we determined the manipulation label, **NEW\_SPAN**,  
 315 and **PRISTINE\_SPAN** before prompting ChatGPT  
 316 to generate a tweet. For efficient annotation, the  
 317 annotators only need to validate whether the labels  
 318 derived from the ChatGPT prompts are correct. We  
 319 keep samples whose labels for all three sub-tasks  
 320 are correct, while the others are discarded. In the  
 321 second round of annotation, human-written tweets  
 322 were annotated, and the predicted labels for these  
 323 tweets were obtained from a model (see below para-  
 324 graphs) trained on the data collected in the first an-  
 325 notation round. For detailed information regarding  
 326 annotation guidelines and the user interface, please  
 327 refer to Appendix C. The following paragraphs  
 328 provide an overview of our annotation process. 329

**First Round** The first round of annotation is for  
 330 curating machine-generated tweets, which are used  
 331 as our training set and development set. Initially,  
 332 for annotator qualification, three annotators worked  
 333 on each of our HITs. We used the first 100 HITs  
 334 to train annotators by instructing them where their  
 335 annotations were incorrect. Then, the next 100  
 336 HITs were used to compute the inter-annotator  
 337 agreement (IAA). At this stage, we did not pro-  
 338 vide further instructions to the annotators. Using  
 339 Fleiss’  $\kappa$  (Fleiss, 1971), we obtain an average IAA  
 340 of 62.4% across all tasks, indicating a moderate  
 341 level of agreement. Finally, we selected the top 15  
 342 performers as qualified annotators. These annota-  
 343 tors were chosen based on how closely their anno-  
 344 tations matched the majority vote for each HIT. 345

Since the annotators have already been trained,  
 346

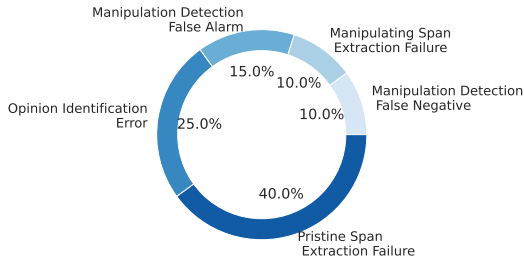


Figure 2: Distributions of errors.

we assigned each HIT to a single annotator to improve annotation efficiency for the remainder of the machine-generated tweets. In addition to being annotated by an MTurk worker, each annotation is also re-validated by a graduate student. The average agreement between the graduate student and the MTurk worker is 93.1% per Cohen’s  $\kappa$  (Cohen, 1960), implying a high agreement. We only keep samples where the validation done by the graduate student agrees with the annotation done by the worker. After two rounds of annotations, we collected 3,116 human-validated samples.

**Second Round** Using the 3K examples we collected, we train a sequence-to-sequence (seq2seq) model that learns to tackle all three tasks jointly. Concretely, we split the collected data into 2,316: 800 for training and validation. Model details are described in the next paragraph. Once the model was trained, we applied it to identify manipulation in the human-written tweets that are associated with the articles in FakeNewsNet. Then, we randomly sampled from predicted MANI and NOMANI examples to be further validated by MTurk workers. The inter-annotator agreement between the graduate student and the MTurk worker is 73.0% per Cohen’s  $\kappa$  (Cohen, 1960). While the agreement is moderately high, it is much lower than that in the previous round. This suggests that manipulation in human-written tweets is more challenging to identify. The user interface of each round of annotation is shown in Appendix C.1. Finally, we have curated the MANITWEET dataset. The dataset statistics are shown in Table 1.

**Baseline Model** In this paragraph, we describe the model we used to facilitate the second round of annotation. Motivated by the advantages of generative models over sequence-tagging models (Li et al., 2021; Huang et al., 2021; Hsu et al., 2022), we trained a seq2seq model based on LongFormer-

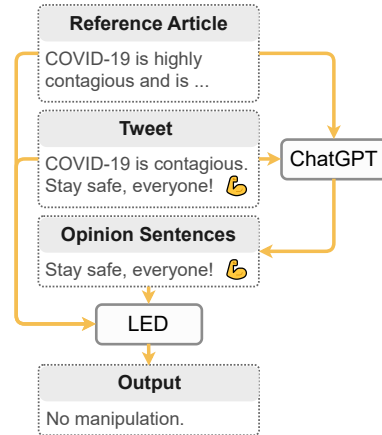


Figure 3: An overview of the proposed framework, **LLM + LED-FT**. We first use ChatGPT to identify sentences that express opinions from the tweet. Then, the opinion sentences are fed to a LED as additional features to help discern between sentences that express personal opinions and sentences that manipulates information.

Encoder-Decoder (LED)<sup>3</sup> (Beltagy et al., 2020) that learns to solve the three tasks jointly. We name this model **LED-FT**.

Formally, the input  $x = [t||a]$  to our model is the concatenation of a tweet  $t$  and the corresponding article  $a$ . The objective of the model is maximum likelihood estimation,

$$\mathcal{L} = - \sum_i p(y_i | y_{<i}, x), \quad (1)$$

where  $y_i$  denotes the  $i$ -th token in the decoding targets. Concretely, if the article is NOMANI, the model should output “No manipulation”. Otherwise, the model should output “**Manipulating span:** `NEW_SPAN` \ **Pristine span:** `PRISTINE_SPAN`”. For cases where `NEW_SPAN` is merely inserted into the tweet, the model will output “None” for `PRISTINE_SPAN`. Details of inputs, outputs, and training hyper-parameters can be found in Appendix A.

## 4 Methodology

We conducted an error analysis on the **LED-FT** model discussed in the previous section. Our analysis revealed that a significant portion of errors occurred due to the model’s inability to distinguish between tweet sentences that express personal opinions and those that manipulate information from the associated article, as depicted in Figure 2 (refer to Appendix B for further details). To address this issue, we propose a pipeline approach that involves

<sup>3</sup><https://huggingface.co/allenai/led-base-16384>

Model	Learning Method	Sub-task 1		Sub-task 2		Sub-task 3		
		F1	EM	F1	RL	EM	F1	RL
Human	-	89.92	44.23	67.93	68.82	42.88	65.29	66.31
Vicuna	Zero-shot	47.09	1.35	5.11	6.07	4.04	6.21	7.06
ChatGPT	Zero-shot	52.49	1.54	13.30	15.96	4.42	7.46	8.35
ChatGPT	Two-shot ICL	65.28	0.96	7.62	8.87	12.50	13.91	14.18
ChatGPT	Four-shot ICL	54.69	3.07	12.79	15.15	1.54	4.99	5.95
ChatGPT	Two-shot CoT	52.92	1.54	7.70	9.21	4.42	5.86	6.12
ChatGPT	Four-shot CoT	53.88	0.96	7.93	9.66	3.46	5.24	5.70
CONCRETE	Zero-shot	57.88	-	-	-	-	-	-
DocNLI	Zero-shot	62.26	-	-	-	-	-	-
QAFactEval	Zero-shot	62.56	-	-	-	-	-	-
LED-FT (Ours)	Fine-tuned	72.62*	26.73*	29.25*	29.68*	13.65*	14.46	14.53
LLM + LED-FT (Ours)	Zero-shot + Fine-tuned	<b>73.46*</b>	<b>28.85*</b>	<b>31.72*</b>	<b>32.32*</b>	<b>15.19*</b>	<b>16.21*</b>	<b>16.41*</b>

Table 2: Performance (%) of different models on the MANITWEET test set. EM denotes Exact Match, and RL denotes ROUGE-L. Statistical significance over best-performing LLMs computed with the paired bootstrap procedure (Berg-Kirkpatrick et al., 2012) are indicated with \* ( $p < .01$ ).

utilizing ChatGPT to identify personal opinions within the tweet. This extracted opinions is then incorporated into our seq2seq model during both training and testing stages. An overview of the framework is shown in Figure 3.

More specifically, we denote the identified opinion sentences in the tweet  $t$  as  $o = p_{\text{LLM}}(t, a, d)$ , where  $d$  represents the instruction provided to ChatGPT for opinion identification. The input to our fine-tuned model becomes  $x' = [t||a||o]$ , and the loss function remains as MLE:

$$\mathcal{L}' = - \sum_i p(y_i | y_{<i}, x'). \quad (2)$$

By incorporating this framework, we aim to enhance the model’s ability to differentiate between personal opinions and instances where information is manipulated from the associated article. We name this pipeline **LLM + LED-FT**.

## 5 Experimental Setup

### 5.1 Evaluation Metrics

Subtask 1 involves a binary classification problem, and thus, the Macro F1 score serves as the evaluation metric. For subtasks 2 and 3, in addition to Exact Match, we use Macro Overlap F1 score (Rajpurkar et al., 2016) and ROUGE-L (Lin, 2004) as the metrics to more accurately assess model performance by allowing models to receive partial credit for correctly identifying some parts of the information, even if they fail to output the entire text span.

### 5.2 Baselines

We compare our proposed framework with various recently released large language models (LLMs),

including Vicuna<sup>4</sup> (vic, 2023) and ChatGPT, which have demonstrated superior language understanding and reasoning capabilities. ChatGPT is an improved version of InstructGPT (Ouyang et al., 2022) that was optimized for generating conversational responses. On the other hand, Vicuna is a LLaMA model (Touvron et al., 2023) fine-tuned on ShareGPT<sup>5</sup> data, and has exhibited advantages compared to other open-source LLMs, such as LLaMA and Alpaca (Taori et al., 2023). We tested the zero-shot, two-shot, and four-shot performance of ChatGPT in both in-context learning (ICL) and chain-of-thought (CoT) (Wei et al., 2022) settings, where the in-context exemplars are randomly chosen from our training set. For Vicuna, we only evaluated its zero-shot ability as we found that it often outputs undesirable texts when exemplars are provided. The details of our prompts for these LLMs can be found in Appendix D. In addition, we also evaluate one fact-checking framework, CONCRETE (Huang et al., 2022), and two faithfulness evaluation frameworks, QAFactEval (Fabbri et al., 2022) and DocNLI (Yin et al., 2021) on our subtask 1. Similar to previous studies, we establish the faithfulness thresholds for both frameworks by selecting the values that yield the highest performance on our development set.

## 6 Results

### 6.1 Performance on MANITWEET

Table 2 presents a summary of the main findings from our evaluation on the MANITWEET test set. We have made several interesting observations:

<sup>4</sup>Vicuna-13b is evaluated in our experiment.

<sup>5</sup><https://sharegpt.com/>

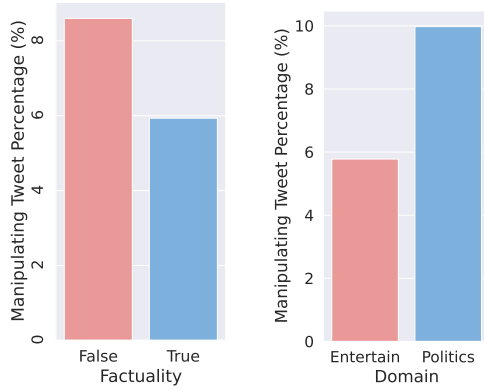


Figure 4: The percentage of tweets that manipulate the associated articles across different levels of factuality and domains.

477 First, all LLMs we tested performed poorly across  
 478 the three proposed tasks. This indicates that simply  
 479 prompting LLMs, whether with or without exam-  
 480 plars, is not sufficient to effectively address the  
 481 problem of identifying manipulation of news on  
 482 social media. We also found that providing more  
 483 exemplars do not work well on our task as the per-  
 484 formance drop when we increase the number of in-  
 485 context exemplars from 2 to 4. This is likely caused  
 486 by the long-context nature of our task. Secondly,  
 487 despite its simplicity and smaller size compared to  
 488 the LLMs, **LED-FT** outperforms all baseline mod-  
 489 els significantly in identifying social media manip-  
 490 ulation across all three tasks. This outcome high-  
 491 lights the value and importance of our training data  
 492 and suggests that a fine-tuned smaller model can  
 493 outshine larger models when tackling challenging  
 494 tasks. Finally, the proposed **LLM + LED-FT** out-  
 495 performs all other models, including **LED-FT** sig-  
 496 nificantly. This implies that LLMs can complement  
 497 smaller fine-tuned models by identifying opinions  
 498 and that the ability to identify opinion sentences  
 499 from social media posts is critical for our task. Ex-  
 500 amples of how the opinions extracted by ChatGPT  
 501 help correct errors can be found in Appendix E.

502 In order to gauge the feasibility of the task, we  
 503 enlisted the assistance of a graduate student to  
 504 tackle our test set. While this may not necessar-  
 505 ily represent the upper bound of performance, it  
 506 provides a preliminary approximation of human  
 507 performance. As depicted in Table 2, there remains  
 508 a discernible gap between **LLM + LED-FT** and  
 509 human performance. This highlights great opportu-  
 510 nities in our task for future research.

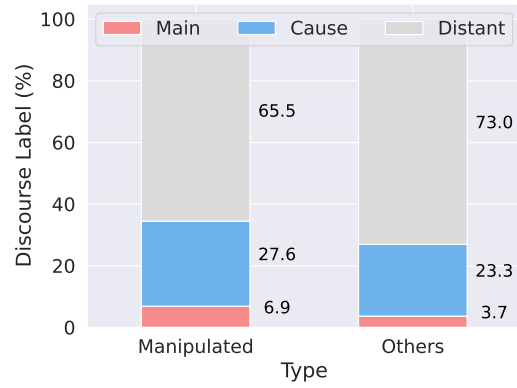


Figure 5: Results of discourse analysis. Manipulated sentences within news articles tend to encompass the main story (*Main*) or convey the consequential aspects (*Cause*) of the corresponding news story.

## 6.2 Exploratory Analysis

511 The proposed **LED-FT** model enables us to per-  
 512 form a large-scale study of manipulation on the  
 513 MANITWEET test set and the 1M human-authored  
 514 tweets associated with the news articles from the  
 515 FakeNewsNet dataset. In this section, we explore  
 516 how an article is MANI and how different proper-  
 517 ties of a news article, such as domain and factuality  
 518 affect manipulation.  
 519

520 **Insight 1: Low-trustworthiness and political**  
 521 **news are more likely to be manipulated.** Fig-  
 522 ure 4 shows the percentage of the 1M human-  
 523 written tweets that are manipulated across 2 do-  
 524 mains and factuality levels.<sup>6</sup> We first observe that  
 525 tweets associated with *False* news are more likely  
 526 to be manipulated. One possible explanation is  
 527 that audience of low-trustworthy news media may  
 528 pay less attention to facts. Hence, they are more  
 529 likely to manipulate information from the refer-  
 530 ence article accidentally when posting tweets. In  
 531 addition, we also see that tweets associated with  
 532 *Politics* news are more frequently manipulated than  
 533 those with *Entertainment* articles. This could be  
 534 explained by the fact that people have a stronger  
 535 incentive to manipulate information for political  
 536 tweets due to elections or campaigns.

537 **Insight 2: Manipulated sentences are more**  
 538 **likely to contain the main story or consequence**  
 539 **of a news story.** To discover the role of the  
 540 sentence being manipulated in the reference  
 541 article, we conducted discourse analysis on these  
 542 sentences. We only conducted the analysis on our  
 543 test set instead of the entire 1M human-written

<sup>6</sup>The domain and factuality labels of each news article are already annotated in the FakeNewsNet dataset.

tweets for this analysis. Concretely, we formulate the discourse classification task as a sequence-to-sequence problem and train a LED-based model on the NEWSDISCOURSE dataset (Choubey et al., 2020) using a similar strategy discussed in §3.2.2. The learned discourse classification model achieves a Micro F1 score of 67.7%, which is on par with the state-of-the-art method (Spangher et al., 2021). Upon the discourse classification model being trained, we applied it to all the sentences in the reference article to analyze the discourse distribution. As shown in Figure 5, compared to other sentences, sentences that were manipulated are much more likely to contain *Main* or *Cause* discourse, which corresponds to *the primary topic being discussed* and *the underlying factor that led to a particular situation*, respectively. Examples of the manipulated sentences with a *Main* or *Cause* discourse can be found in Appendix F.

## 7 Related Work

### 7.1 Faithfulness

Faithfulness is often referred to as the factual consistency between the inputs and outputs. This topic has mainly been studied in the field of summarization. Prior work on faithfulness can be divided into two categories: evaluation and enhancement, the former of which is more relevant to our study. One line of faithfulness evaluation work developed entailment-based metrics by training document-sentence entailment models on synthetic data (Kryscinski et al., 2020; Yin et al., 2021) or using traditional natural language inference (NLI) models at the sentence level (Laban et al., 2022). Another line of studies evaluates faithfulness by comparing information units extracted from the summaries and input sources using QA (Wang et al., 2020; Deutsch et al., 2021; Fabbri et al., 2022).

Our task differs from faithfulness evaluation in two key ways. Firstly, for our task to be completed effectively, models must possess the additional capability of distinguishing tweet sentences that relate to the reference article from those that simply express opinions. In contrast, models evaluating faithfulness only need to identify whether each sentence in the output is inferable from the input. Secondly, we require models to not only identify which original information is being manipulated by the new information, but also to provide interpretability as to why a tweet has been manipulated.

### 7.2 Fact-checking

Fact-checking is a task that determines the veracity of an input claim based on some evidence passages. Some work assumes the evidence candidates are provided, such as in the FEVER dataset (Thorne et al., 2018) and the SCIFACT dataset (Wadden et al., 2020). Approaches for this category of fact-checking tasks often involve a retrieval module to retrieve relevant evidence from the given candidate pool, followed by a reasoning component that determines the compatibility between a piece of evidence and the input claim (Yin and Roth, 2018; Pradeep et al., 2021). Other work focuses on the *open-retrieval* setting, where evidence candidates are not provided, such as in the LIAR dataset (Wang, 2017) and the X-FACT dataset (Gupta and Srikumar, 2021). For this task formulation, one of the main challenges is to determine where and how to retrieve evidence. Some approaches determine the veracity of a claim based solely on the claim itself and the information learned by language models during the pre-training stage (Lee et al., 2021), other methods leverage a retrieval module to look for evidence on the internet (Gupta and Srikumar, 2021) or a set of trustworthy sources (Huang et al., 2022). Similar to the faithfulness task, the key distinction between fact-checking and our proposed task lies in the additional requirement for models to possess the capability of discerning between tweet sentences that pertain to the reference article and those that merely express opinions.

## 8 Conclusion

In this study, we have introduced and defined a novel task called *identifying manipulation of news on social media*, which aims to determine whether and how a social media post manipulates the associated news article. To address this challenge, we meticulously collected a dataset named MANITWEET, composed of both human-written and machine-generated tweets. Our analysis revealed that existing large language models (LLMs) prompted with zero-shot and two-shot exemplars do not yield satisfactory performance on our dataset, highlighting avenues for future research. We believe that the resources presented in this paper can serve as valuable assets in combating the dissemination of false information on social media, particularly in tackling the issue of news manipulation.



## 9 Limitations

There are two main limitations in our work. Firstly, despite our best efforts to minimize the gap between the training set and test set of MANITWEET, some discrepancies remain due to the training set being generated by machines and the test set being produced by humans. This limitation is primarily attributed to budget constraints. In the future, with additional resources, we aim to create an additional training set consisting entirely of human-written tweets. By comparing the performance of models trained on this human-written training set with those trained on the machine-generated training set, we can gain further insights. However, we wanted to emphasize that our test set exclusively consists of tweets authored by humans, which ensures the relevance of our techniques and dataset for real-world applications in handling tweets produced by actual Twitter users. While our data collection method may introduce discrepancies in the distribution between the training and test sets, the fundamental purpose of our dataset remains consistent: to investigate the manipulation of news articles on social media.

Secondly, in our experiments involving prompting LLMs, we only explored ICL and CoT for prompting LLMs. There is a possibility that LLMs can achieve better performance when provided with more in-context exemplars and when prompted in a more refined manner.

## 10 Ethical Considerations

The primary ethical consideration in our work pertains to the presence of false information in two aspects: tweets that manipulate the associated news articles and the inclusion of false news from the FakeNewsNet dataset. As with other fact-checking and fake news detection research, it is important to acknowledge the dual-use concerns associated with the resources presented in this work. While our resources can contribute to combating false information, they also possess the potential for misuse. For instance, there is a risk that malicious users could utilize the manipulating tweets or fake news articles to train a text generator for creating deceptive content. We highlight appropriate and inappropriate uses of our dataset in various scenarios:

- **Appropriate:** Researchers can use our framework to study the manipulation issue on social media and develop stronger models for

identifying social media posts that manipulate information.

- **Inappropriate:** The fake news and manipulating tweets in MANITWEET cannot be used to train text generators for malicious purposes.
- **Inappropriate:** Use the manipulation prompts discussed in this paper to generate tweets and spread false information.
- **Inappropriate:** The fake news in MANITWEET should not be used as evidence for fact-checking claims.

Furthermore, the privacy of tweet users is another aspect that warrants consideration, given that we are releasing human-written tweets. However, we assure that the dataset does not pose significant privacy concerns. The tweets in our dataset are anonymized, and it is important to note that all the associated news articles were already publicly available. Therefore, the release of this dataset should not have adverse implications for privacy.

## References

2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90% chatgpt quality.](#)
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Taylor Berg-Kirkpatrick, David Burkett, and Dan Klein. 2012. [An empirical investigation of statistical significance in NLP.](#) In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 995–1005, Jeju Island, Korea. Association for Computational Linguistics.
- Prafulla Kumar Choubey, Aaron Lee, Ruihong Huang, and Lu Wang. 2020. [Discourse as a function of event: Profiling discourse structure in news articles around the main event.](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5374–5386, Online. Association for Computational Linguistics.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021. [Towards question-answering as an automatic metric for evaluating the content quality of a summary.](#) *Transactions of the Association for Computational Linguistics*, 9:774–789.

739	Alexander Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. <a href="#">QAFactEval: Improved QA-based factual consistency evaluation for summarization</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2587–2601, Seattle, United States. Association for Computational Linguistics.	797
740		798
741		799
742		800
743		801
744		
745		802
746		803
747	Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. <i>Psychological bulletin</i> , 76(5):378.	804
748		805
749		806
750		807
751	Yi Fung, Christopher Thomas, Revanth Gangi Reddy, Sandeep Polisetty, Heng Ji, Shih-Fu Chang, Kathleen McKeown, Mohit Bansal, and Avi Sil. 2021. <a href="#">InfoSurgeon: Cross-media fine-grained information consistency checking for fake news detection</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1683–1698, Online. Association for Computational Linguistics.	808
752		809
753		810
754		811
755		812
756		813
757		814
758		815
759		
760	Ashim Gupta and Vivek Srikumar. 2021. <a href="#">X-fact: A new benchmark dataset for multilingual fact checking</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 675–682, Online. Association for Computational Linguistics.	816
761		817
762		818
763		819
764		
765		820
766		821
767	I-Hung Hsu, Kuan-Hao Huang, Elizabeth Boschee, Scott Miller, Prem Natarajan, Kai-Wei Chang, and Nanyun Peng. 2022. <a href="#">DEGREE: A data-efficient generation-based event extraction model</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1890–1908, Seattle, United States. Association for Computational Linguistics.	822
768		823
769		824
770		825
771		826
772		827
773		828
774		
775		829
776	Kung-Hsiang Huang, Sam Tang, and Nanyun Peng. 2021. <a href="#">Document-level entity-based extraction as template generation</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5257–5269, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	830
777		831
778		832
779		833
780		834
781		
782		835
783	Kung-Hsiang Huang, ChengXiang Zhai, and Heng Ji. 2022. <a href="#">CONCRETE: Improving cross-lingual fact-checking with cross-lingual retrieval</a> . In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 1024–1035, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	836
784		837
785		838
786		839
787		840
788		
789		841
790	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. <a href="#">Evaluating the factual consistency of abstractive text summarization</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9332–9346, Online. Association for Computational Linguistics.	842
791		843
792		844
793		845
794		846
795		847
796		848
	Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. <a href="#">SummaC: Re-visiting NLI-based models for inconsistency detection in summarization</a> . <i>Transactions of the Association for Computational Linguistics</i> , 10:163–177.	849
		850
		851
		852
		853
	Nayeon Lee, Yejin Bang, Andrea Madotto, and Pascale Fung. 2021. <a href="#">Towards few-shot fact-checking via perplexity</a> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1971–1981, Online. Association for Computational Linguistics.	854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

854	Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. <i>SIGKDD Explor. Newsl.</i> , 19(1):22–36.	909
855		910
856		911
857		912
858	Alexander Spangher, Jonathan May, Sz-Rung Shiang, and Lingjia Deng. 2021. Multitask semi-supervised learning for class-imbalanced discourse classification. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 498–517, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	913
859		914
860		915
861		916
862		917
863		918
864		919
865	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://github.com/tatsu-lab/stanford_alpaca</a> .	920
866		
867		
868		
869		
870	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.	
871		
872		
873		
874		
875		
876		
877		
878		
879	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
880		
881		
882		
883		
884		
885	David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7534–7550, Online. Association for Computational Linguistics.	
886		
887		
888		
889		
890		
891		
892	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5008–5020, Online. Association for Computational Linguistics.	
893		
894		
895		
896		
897		
898	William Yang Wang. 2017. “liar, liar pants on fire”: A new benchmark dataset for fake news detection. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 422–426, Vancouver, Canada. Association for Computational Linguistics.	
899		
900		
901		
902		
903		
904	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> .	
905		
906		
907		
908		

## A Training Details

### A.1 LED-based Fine-tuned Model

The input to our LED-based model is a concatenation of a tweet and a reference article:

```
Tweet: TWEET \  
Reference article: REF_ARTICLE
```

If the article is NOMANI, the model should output:

```
No manipulation
```

Otherwise, the model should output the following:

```
Manipulating span: NEW_SPAN \  
Pristine span:  
PRISTINE_SPAN
```

For cases where `NEW_SPAN` is merely inserted into the tweet, the model will output “None” for `PRISTINE_SPAN`. Using this formulation, our model is learned to optimize the maximum likelihood estimation loss. We set identical weights for all tokens in the outputs.

### A.2 ChatGPT Prompts

The prompt to ChatGPT for identifying opinions is as follows:

```
Tweet: TWEET \  
Reference article: REF_ARTICLE  
Given the above tweet and article. List the sentences in the tweet that merely express opinions instead of manipulating information from the article. If there is none, answer "None". Do not provide explanations.
```

### A.3 Training Hyper-parameters

To learn the model, we use a learning rate of  $5e-5$ . The maximum input and output sequence length are 1024 and 32 tokens, respectively. The model is optimized using the AdamW optimizer (Loshchilov and Hutter, 2019) with a batch size of 4 and a gradient accumulation of 8. During inference time, we use beam search as the decoding method with a beam width of 4.

## B Error Analysis

To gain insights into the additional modeling and reasoning capabilities required for effectively addressing the task of social media manipulation, we manually compare 50 errors made by the LED-based model with ground-truth labels and analyze

the sources of errors. The distribution of errors is illustrated in Figure 2. Notably, the most prevalent error arises from the model’s inability to extract the correct pristine span from the reference article that underwent manipulation. Among the 18 erroneous predictions in this category, 16 cases result from the model producing an empty string. This indicates that the model considers the manipulating information to be inserted when, in reality, it is manipulated from the information present in the reference articles. This could be attributed to the presence of 368 instances where the original information is an empty string, while the alternative answers for the original information only occur 1-2 times in other instances. This can be solved by scaling down the loss for these samples with an empty string as the label for original information. Additionally, another common type of error involves the model’s failure to identify opinions expressed in the tweet. In these instances, the model considers the tweet to be manipulating information from the article, whereas the tweet primarily expresses opinions. Examples of these errors are presented in Appendix E.

## C Annotation Details

In this section, we describe the details of our annotation process. For better control of the annotation quality, we required that all annotators be from the U.S. and have completed at least 10,000 HITs with 99% acceptance on previous HITs. The reward for each HIT is \$1 U.S. dollar, complying with the ethical research standards outlined by AMT (Salehi et al., 2015). Annotation interfaces are shown below.

### C.1 User Interface

Figure 6 and Figure 7 display the annotation interface for the first round and the third round of annotation, respectively. The only difference is that for the second round of annotation, we asked annotators to correct errors made by our basic model discussed in §3.2.2. Samples that do not receive “yes” on all three questions for the first round of annotation will be discarded. The rationale behind this design stems from three key reasons: Firstly, the data for the first round of annotation is automatically generated, enabling a relatively cost-effective approach to discard invalid samples and generate new ones, as opposed to requesting annotators to correct errors. Secondly, the data generated in these

1014 two rounds is predominantly valid, which elimi-  
 1015 nates the need for annotators to rectify errors and  
 1016 consequently accelerates the annotation process.  
 1017 Lastly, in the second round of annotation, by in-  
 1018 structuring annotators to identify errors made by our  
 1019 model, we can effectively identify the challenges  
 1020 faced by the model.

## 1021 D Prompts for LLMs

1022 The zero-shot and two-shot prompt template to  
 1023 LLMs for the experiments discussed in §5.2 is  
 1024 shown in Table 3. The in-context exemplars for  
 1025 the two-shot experiments are randomly sampled  
 1026 from the training set of MANITWEET.

## 1027 E Additional Qualitative Examples

1028 Table 4 presents two instances where our baseline  
 1029 model makes errors. In the first example, our model  
 1030 was not able to identify that “Inspired Our Next  
 1031 Trip To The Salon” is an expression of opinion,  
 1032 resulting in the model incorrectly classifying this  
 1033 sample as MANI. In the second example, although  
 1034 our model accurately predicts the example as MANI  
 1035 and extracts the correct manipulating span, it fails  
 1036 to extract the pristine text span correctly, likely due  
 1037 to the nature of the training set, as discussed in  
 1038 Appendix B.

1039 Table 5 shows an example where extracting opin-  
 1040 ion sentences from the tweet by ChatGPT enables  
 1041 our model to correctly identify the tweet as not

manipulating the associated article. 1042

## F Discourse Analysis Examples 1043

1044 Table 6 shows examples of manipulated sentences  
 1045 associated with a *Main* or *Cause* discourse. A *main*  
 1046 discourse implies that the sentence conveys the  
 1047 main story of an article, whereas a *cause* discourse  
 1048 indicates that the sentences discuss the consequen-  
 1049 tial aspect of the main story.

Please read the instructions before doing the annotation! We will carefully check each annotated sample.

**Tweet:**  
 \${tweet}

**Our predicted original and recontextualized fact (manipulated or inserted facts in the tweet):**  
 \${original\_concept} -> \${recontextualized\_concept}

**Reference Article:**  
 \${reference\_article}

We predicted that this tweet is: **\$(is\_recontextualized)**. Did we predict it correctly?  
 Yes  
 No

If you think the tweet **IS RECONTEXTUALIZED**, answer the remaining two questions:  
 We predicted that the original fact is: **\$(original\_concept)**. Did we predict it correctly?  
 Yes  
 No

We predicted that the recontextualized fact is: **\$(recontextualized\_concept)**. Did we predict it correctly?  
 Yes  
 No

Figure 6: MTurk user interface for the first round of data annotation.

Please read the instructions before doing the annotation! We will carefully check each annotated sample.

**Tweet:**  
 \${tweet}

**Our predicted original and recontextualized fact (manipulated or inserted facts in the tweet):**  
 \${original\_concept} -> \${recontextualized\_concept}

**Reference Article:**  
 \${reference\_article}

We predicted that this tweet is: **\$(is\_recontextualized)**. Did we predict it correctly?  
 Yes  
 No

If you think the tweet **IS RECONTEXTUALIZED**, answer the remaining two questions:  
 We predicted that the original fact is: **\$(original\_concept)**. Did we predict it correctly?  
 Yes  
 No

If we did it wrong, please type in the correct original fact:

We predicted that the recontextualized fact is: **\$(recontextualized\_concept)**. Did we predict it correctly?  
 Yes  
 No

If no, please type in the correct recontextualized fact:

Figure 7: MTurk user interface for the second round of data annotation.

Role	Utterance
User	You are tackling a social manipulation problem. You will be giving a tweet and an article, and your task is to identify which information from the article is misrepresented by which information in the tweet. You should answer in the following format “Manipulating span: manipulating_span Pristine span: pristine_span” in a single line. Here, {manipulating_span} is the new information introduced in the tweet and original_concept is the original information in the article. If the tweet simply inserts information, {original_concept} should be "None". If the tweet does not manipulate the article, answer “No manipulation”. You do not need to output other information such as an explanation. You don’t need to provide code. In the following utterances, you will be presented a pair of tweet and news article.
LLM	I understand the task. Please provide me with the tweet and the corresponding article, and I will identify the manipulation.
User	Tweet: Sad news for Ben Affleck and Lindsay Shookus as they split after a year of dating. It looks like Ben has been getting cozy with Playboy model Shauna Sexton, but sources say it’s just light and fun. Best of luck to all involved. #BenAffleck #LindsayShookus #ShaunaSexton Article: Ben Affleck and Lindsay Shookus Split After a Year of Dating Ben Affleck and Lindsay Shookus have broken up after over a year of dating, a source tells ET... =====
LLM	No manipulation
User	Tweet: Meghan’s Lush Malibu Home Worth \$16M, which comes with 8 bedrooms, 8 bathrooms, a private beach, a pool, a home cinema - they even want to spend an extra \$1 million to decorate it to their own taste. Unfortunately Caboolture scores 0 here! Article: Take A Look Inside Prince Harry And Meghan Markle’s Lush Malibu Home If you’re set to marry royalty, like Meghan Markle you would expect to move into a palace, and live out your Disney Princess fantasy... =====
LLM	Manipulating span: Caboolture Pristine span: None
User	Tweet: <i>Input Tweet</i> Article: <i>Input Article</i>

Table 3: The zero-shot and two-shot prompt template to the LLMs for producing the results in Table 2. Texts marked in gray indicate that they are only used in the two-shot experiments.

Opinion Identification Error	
Input	<b>Tweet:</b> Ariana Grande’s New Lavender Hair Color Just Inspired Our Next Trip To The Salon ... <b>Article:</b> Ariana Grande Dyed Her Hair, And This Is Our Favorite Color Transformation Yet Ariana Grande is giving us whiplash with her hairstyles lately, and we honestly love it. On July 18th, Grande took to Instagram to debut her latest hair transformation. She’s now sporting pastel lavender locks and good god (is a woman), it looks amazing on her...
Prediction	<b>Is manipulated:</b> Yes ✗ <b>Manipulating span:</b> Salon ✗ <b>Pristine span:</b> None
Pristine Span Extraction Failure	
Input	<b>Tweet:</b> Transcript: Democratic Presidential Debate in <b>Brooklyn</b> view more ... <b>Article:</b> The Democratic Debate in <b>Cleveland</b> This is rightly a big issue in Ohio. And I have laid out my criticism, but in addition my plan, for actually fixing NAFTA. Again, I have received a lot of incoming criticism from Senator Obama. And the Cleveland Plain Dealer examined Senator Obama’s attacks on me regarding NAFTA and said they were erroneous. So I would hope that, again, we can get to a debate about what the real issues are and where we stand because we do need to fix NAFTA. It is not working. It was, unfortunately, heavily disadvantaging many of our industries, particularly manufacturing. ...
Prediction	<b>Is manipulated:</b> Yes <b>Manipulating span:</b> Brooklyn <b>Pristine span:</b> None ✗

Table 4: Example outputs from our baseline model where it produces erroneous outputs.

Input	<b>Tweet:</b> Ariana Grande’s New Lavender Hair Color Just Inspired Our Next Trip To The Salon ... <b>Article:</b> Ariana Grande Dyed Her Hair, And This Is Our Favorite Color Transformation Yet Ariana Grande is giving us whiplash with her hairstyles lately, and we honestly love it. On July 18th, Grande took to Instagram to debut her latest hair transformation. She’s now sporting pastel lavender locks and good god (is a woman), it looks amazing on her...
Prediction	Is manipulated: Yes ✗ Manipulating span: Salon ✗ Pristine span: None
Input	<b>Tweet:</b> Ariana Grande’s New Lavender Hair Color Just Inspired Our Next Trip To The Salon ... <b>Predicted Opinions:</b> Ariana Grande’s New Lavender Hair Color Just Inspired Our Next Trip To The Salon <b>Article:</b> Ariana Grande Dyed Her Hair, And This Is Our Favorite Color Transformation Yet Ariana Grande is giving us whiplash with her hairstyles lately, and we honestly love it. On July 18th, Grande took to Instagram to debut her latest hair transformation. She’s now sporting pastel lavender locks and good god (is a woman), it looks amazing on her...
Prediction	<b>Is manipulated:</b> No ✓ <b>Manipulating span:</b> None ✓ <b>Pristine span:</b> None

Table 5: Example outputs from our LED-FT and LLM + LED-FT. The predicted opinion extracted by ChatGPT allows the fine-tuned model to predict the manipulation label correctly.

<i>Main Discourse</i>	
Tweet	#Zuckerbergtestimony <b>Mark Zuckerberg</b> ’s testimony before the House Energy and Commerce Committee is over.
Article	... U.S. Rep. Joe Barton, R-Texas, chairman of the House Energy and Commerce Committee, made the following statement today during the full committee hearing on the Administration’s FY 07 Health Care Priorities: "Good afternoon.. <b>Let me begin by welcoming Secretary Michael Leavitt today to the Energy and Commerce Committee.</b> We look forward to hearing him testify about the Administration’s Fiscal Year 2007 Health Care Priorities ...
<i>Cause Discourse</i>	
Tweet	Thank you, Rep. Johnson, for your service! Weekly Republican Address: Rep. <b>Sam Johnson</b> (R-TX) ... via @YouTube
Article	... In the address, Boehner notes that this is a new approach that hasn’t been tried in Washington – by either party – and it is at the core of the Pledge to America, a governing agenda Republicans built by listening to the people. <b>Leader Boehner recorded the weekly address earlier this week from Ohio, where he ran a small business and saw first-hand how Washington can make it harder for employers and entrepreneurs to meet a payroll and create jobs.</b> Following is a transcript ...

Table 6: Examples of manipulated sentences with a *Main* discourse and a *Cause* discourse. The manipulated sentences are marked in **boldface**. The manipulating and pristine spans are marked in **red** and **blue**, respectively.