Can Large Language Models Detect Rumors on Social Media?

Anonymous ACL submission

Abstract

In this work, we investigate to use Large Language Models (LLMs) for rumor detection on social media. However, it is challenging for LLMs to reason over the entire propagation information on social media, which contains news contents and numerous comments, due to LLMs may not concentrate on key clues in the complex propagation information, and have trouble in reasoning when facing massive and redundant information. Accordingly, we propose an LLM-empowered Rumor Detection (LeRuD) approach, in which we design prompts to teach LLMs to reason over important clues in news and comments, and divide the entire propagation information into a Chain-of-Propagation for reducing LLMs' burden. We conduct extensive experiments on the 018 Twitter and Weibo datasets, and LeRuD outper-019 forms several state-of-the-art rumor detection models by 3.2% to 7.7%. Meanwhile, by applying LLMs, LeRuD requires no data for training, and thus shows more promising rumor detection ability in few-shot or zero-shot scenarios.

1 Introduction

011

037

041

With the development of social media, users can access to information more easily, and messages can be spread more rapidly. However, at the same time, rumors can also be spread to the public more easily and widely. It is necessary to conduct research on automatic rumor detection on social media (Castillo et al., 2011; Islam et al., 2020), which is still very challenging. Recently, Large Language Models (LLMs) have achieved success and been applied widely (Zhao et al., 2023). This brings great convenience to people. However, people can also use LLMs to intentionally generate great amount of unidentified information (Vykopal et al., 2023; Chen and Shu, 2024), which aggravate the propagation of rumors on social media. Thus, the success of LLMs makes automatic rumor detection more challenging and urgent.



Figure 1: Can LLMs detect rumors on social media?

The key of rumor detection lies in the modeling of propagation information on social media, i.e., aggregating the features of original news and user comments. Early approaches regard news and comments as sequences of texts, and apply recurrent networks (Ma et al., 2016), convolutional networks (Yu et al., 2017) or attention networks (Liu et al., 2018) for feature extraction. Then, rumor detection models based on graph neural networks (Lu and Li, 2020; Bian et al., 2020) have been widely proposed and achieved state-of-the-art performances.

Recently, some research works attempt to apply LLMs in identifying the veracity of information (Chen and Shu, 2023). They mostly use LLMs to extract knowledge from texts (Yang et al., 2023) or solely analyze the news contents (Hu et al., 2024). The power of LLMs for modeling the complex propagation information on social media, which contains both news and comments, has not been explored. Meanwhile, as introduced by (Hu et al., 2024; Chen and Shu, 2024), solely relying on textual features of news, it is hard to achieve promising veracity identification by LLMs. This requires us to investigate the ability of LLMs for modeling the propagation information on social media.

However, using LLMs for modeling the propagation information is not easy. Firstly, we find

042

077 094

070

071

101 102

103 104

105 106

107 108

109

110 111

112 113

114 115

116 117

118

LLMs can not concentrate on the key clues in the propagation information, and thus may produce wrong predictions (Sec. 3.2 and 3.3). Secondly, the amount of comments is usually large, which brings great burden for LLMs to reason (Sec. 3.4), as LLMs usually have input length limitations, and have trouble when facing long contexts or redundant information (Huang et al., 2023; Xie, 2023).

To teach LLMs to conduct reasoning over propagation information and overcome above challenges, we propose an LLM-empowered Rumor Detection approach, called LeRuD. Firstly, we draw experience from previous rumor detection models (Ma et al., 2016; Liu et al., 2018; Ma et al., 2018; Yang et al., 2022; Jin et al., 2022; Hu et al., 2024; Przybyla, 2020), in which useful information for distinguish rumors has been highlighted. We design proper prompts to teach LLMs to concentrate on the writing styles and commonsense mistakes of news, as well as rebuttals or conflicts in comments. Secondly, we divide the propagation information into a Chain-of-Propagation, which enables LLMs to reason step-by-step in time order with a reasonable number of comments.

Then, to evaluate the effectiveness of LeRuD on detecting rumors, we conduct experiments on the Twitter and Weibo datasets (Sec. 4.1). We use the widely-applied GPT 3.5¹ in our implementation. For smooth usage of GPT 3.5, we filter samples with ethics issues. Meanwhile, to avoid data leakage, we also filter samples in which GPT 3.5 already knows about the facts in the news. On the filtered datasets, LeRuD outperforms several state-of-the-art rumor detection models by 3.2% to 7.7%. To be noted, LeRuD achieves these results in zero-shot setting, while other models require training data. In a word, with our efforts, the LLM can well detect rumors on social media.

Our main contributions can be listed as follows:

- We propose to apply LLMs in detecting rumors on social media with the entire propagation information, which contains both news and comments.
- We design proper prompts to teach LLMs to concentrate on the key clues in news and comments, which enable LLMs to conduct reasoning effectively over propagation information.
- We divide the propagation information of a piece of news into a Chain-of-Propagation, which enables LLMs to conduct reasoning

step-by-step with reasonable amount of information in each time.

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

• We conduct extensive experiments, in which several state-of-the-art rumor detection models are compared. Results demonstrate the superiority of our LeRuD approach.

2 **Related Work**

Rumor detection on social media is a major text mining task for automatically distinguishing untrustworthy information and maintaining the healthy development of the internet and society. Extensive works have been done to detect rumors from different perspectives (Castillo et al., 2011; Islam et al., 2020).

Early attempts on rumor detection mainly focus on extracting statistical features of the propagation process (Kwon et al., 2013; Ma et al., 2015). Then, deep neural networks for modeling sequential data have been widely applied in aggregating propagation information and detecting rumors on social media. Among them, detection models based on recurrent neural networks (Ma et al., 2016, 2018), convolutional neural networks (Yu et al., 2017), attentive networks (Liu et al., 2018) and generative adversarial networks (Ma et al., 2019) have been extensively studied. Meanwhile, their combinations have also been investigated for rumor detection (Liu and Wu, 2018; Yu et al., 2019). Recently, detection models based on graph neural networks (Bian et al., 2020; Lu and Li, 2020; Xu et al., 2022) have achieved state-of-the-art performances. Furthermore, some works investigate incorporate graph contrastive learning for better rumor detection performances (Sun et al., 2022b; Wu et al., 2023). For better interpretability, reasoning over subgraph has been proposed (Jin et al., 2022; Yang et al., 2022). Meanwhile, for capturing temporal characteristics, dynamic graph has been constructed and applied in rumor detection (Sun et al., 2022a). Nowadays, considering practical applications, rumor detection in few-shot or zero-shot scenarios have drawn much attention (Yue et al., 2023; Lin et al., 2023).

As LLMs have achieved great success, they also bring opportunities and challenges to misinformation identification (Chen and Shu, 2023). LLMs are able to easily generate misinformation (Vykopal et al., 2023; Chen and Shu, 2024), which is hard to be distinguished by human or LLMs solely based on textual features (Chen and Shu, 2024). CICAN

¹https://chat.openai.com/

Vanilla Prompt	There is a piece of news: "Samsung paid Apple one-billion fine, with 20 billion coins in 30 trucks!" Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	I cannot verify the credibility of real-time or specific news events as my training only includes information up until January 2022, and I do not have the capability to access the internet for real-time updates.
Rational Prompt	There is a piece of news: "Samsung paid Apple one-billion fine, with 20 billion coins in 30 trucks!" <i>Based on the writing style and the commonsense knowledge, estimate the credibility of the news.</i> Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output	The news about Samsung paying Apple a one-billion fine with 20 billion coins in 30 trucks seems highly implausible and <u>unrealistic</u> . It lacks specific details, and the use of such an unconventional method for payment raises doubts about the credibility of the information. Therefore, the most appropriate answer would be: A. Fake .
Vanilla Prompt	There is a piece of news: "Around 7 o'clock this morning, a major traffic accident occurred on the Yuncheng to Houma section of the Shanxi Dayun Expressway, 9.3 kilometers away from Houma. A total of 47 vehicles collided in a series, resulting in 3 deaths and 7 injuries." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	I cannot verify the credibility of real-time or specific news events as my training only includes information up until January 2022, and I do not have the capability to access the internet for real-time updates.
Rational Prompt	There is a piece of news: "Around 7 o'clock this morning, a major traffic accident occurred on the Yuncheng to Houma section of the Shanxi Dayun Expressway, 9.3 kilometers away from Houma. A total of 47 vehicles collided in a series, resulting in 3 deaths and 7 injuries." <i>Based on the writing style and the commonsense knowledge, estimate the credibility of the news.</i> Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output	The provided news about a major traffic accident on the Yuncheng to Houma section of the Shanxi Dayun Expressway appears to be presented in a straightforward and factual manner. It includes specific details such as the time of the incident, the location, and the number of vehicles involved, as well as the reported casualties (3 deaths and 7 injuries). While I cannot independently verify the information, the news seems plausible and is presented in a manner consistent with typical reporting of traffic accidents. Therefore, based on the writing style and common sense, it appears to be: B. Real .

Table 1: Analyzing news contents with GPT 3.5 on two samples. Without proper prompts, it is hard for the LLM to judge the credibility of the news. Thus, we design prompts (in *red italicized words*) to teach the LLM to concentrate on the writing styles and commonsense mistakes, which enable the LLM to find important clues in writing details and conduct commonsense reasoning. Key clues during reasoning are highlighted with <u>red underlined words</u>.

(Yang et al., 2023) uses GPT 3.5 as an auxiliary tool, for extracting entities from news and comments, and constructing relational graphs. ARG (Hu et al., 2024) uses LLMs for analyzing the news contents, and concludes that, LLMs are hard to beat fine-tuned small models on detecting fake news. Thus, ARG incorporates insights generated by LLMs in small models for fake news detection.

3 Methodology

169

170

171

173

174

175

176

177

178

179

180

181

182

183

186

188

In this section, we detail the LeRuD approach, and the whole procedure is demonstrated in Fig. 3.

3.1 Problem Formulation

Each sample consists of news content n and a series of comments C. News content n is a paragraph of texts. And $C = [(c_1, t_1), (c_2, t_2), \dots]$, in which c_i denotes the textual content of the *i*-th comment, and t_i means the relative posting time compared with the original news. We aim to utilize LLMs for detecting rumors, i.e., identifying whether the sample is a rumor (fake) or a non-rumor (real).

3.2 RP: Rational Prompts for News

We need to design proper prompts to teach LLMs
for analyzing news contents for rumor detection.
Here, we apply the widely-used GPT 3.5 as the
LLM used in our proposed approach. As shown in

the two cases in Tab. 1, we try the vanilla prompts, which directly ask the LLM to infer the credibility of the news contents. However, without proper prompts, it is hard for the LLM to make guesses about veracity. The LLM does not know how to reason over the contents in which it has no knowledge about the described facts. This is mainly because that, the LLM can not concentrate on some key clues in the news contents.

According to previous study, writing styles of news are highlighted for misinformation identification (Przybyla, 2020), and commonsense reasoning is essential for distinguishing fake news (Hu et al., 2024). Thus, we need to teach LLMs to reason from these aspects. Accordingly, we design Rational Prompts (RP), which ask the LLM to answer "Based on the writing style and the commonsense knowledge, estimate the credibility of the news." As shown in Tab. 1, with Rational Prompts, the LLM is able to find important clues in writing details and conduct commonsense reasoning, and give judgements on the veracity of news contents. In the first case, the LLM correctly predicts the news is fake, by noticing that the event is not plausible from a commonsense perspective. In the second case, the LLM correctly predicts the news is real, via noticing some accident details, and concluding that the writing style is consistent with typical reports.

221

Vanilla Prompt	There is a piece of news: "Around 7 o'clock this morning, a major traffic accident occurred on the Yuncheng to Houma section of the Shanxi Dayun Expressway, 9.3 kilometers away from Houma. A total of 47 vehicles collided in a series, resulting in 3 deaths and 7 injuries." There are comments for the news: "My god!" "Blessing" "Why isn't the highway closed on such a foggy day?" "Repost" "Cherishing life!" "How scary!" "Blessing." "Drive carefully and pay attention to safety." "Shocking" Based on the comments, verify the authenticity of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	Based on the provided comments, it's challenging to determine the authenticity of the news directly. The comments express shock, concern, and questioning regarding the incident, but they don't provide concrete evidence or information about the news source. Therefore, the answer is: A. Fake .
Conflicting Prompt	There is a piece of news: "Around 7 o'clock this morning, a major traffic accident occurred on the Yuncheng to Houma section of the Shanxi Dayun Expressway, 9.3 kilometers away from Houma. A total of 47 vehicles collided in a series, resulting in 3 deaths and 7 injuries." There are comments for the news: "My god!" "Blessing" "Why isn't the highway closed on such a foggy day?" "Repost" "Cherishing life!" "How scary!" "Blessing." "Drive carefully and pay attention to safety." "Shocking" <i>Based on the comments, analyze whether there are any rebuttals or conflicts, and accordingly verify the authenticity of the news.</i> Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output	The comments mainly express shock, concern, and questions about the responsibility for the accident. Some comments criticize the decision not to close the expressway in thick fog, suggesting potential negligence. There is also a call for attention to safety in such weather conditions. There doesn't seem to be any direct rebuttal or conflict among the comments. Instead, the comments generally focus on expressing emotions, seeking information about the incident, and questioning the decision-making around highway closure during foggy conditions. Based on the analysis of comments, there is no clear evidence of conflicting viewpoints that would cast doubt on the authenticity of the news. The concerns raised by commenters align with the details provided in the news report. Therefore, the answer is: B. Real .

Table 2: Analyzing comments with GPT 3.5 on a sample. Without proper prompts, the LLM tends to predict the sample as a rumor when there lack clear supporting materials. However, in most cases, is it hard to have enough evidences, especially in the early stage of propagation. Thus, we design prompts (in *red italicized words*) to teach the LLM to concentrate on rebuttals or conflicts in comments, which lead the LLM to use crowd wisdom for judging the veracity of information. Key clues during reasoning are highlighted with <u>red underlined words</u>.

3.3 CP: Conflicting Prompts for Comments

Furthermore, we need to design prompts to teach LLMs to reason over comments in propagation. As shown in Tab. 2, we try the vanilla prompts, which directly ask the LLM to verify the authenticity of news based on given comments. However, with such prompts, the LLM tends to predict the sample is a rumor when there lack clear supporting materials. In most cases, is it hard to have enough clear evidences, especially in the early stage of propagation. This makes it easy for the LLM to mistakenly recognize non-rumors as rumors.

Conflicts, rebuttals and skepticism in comments have been proven to be key features for detecting rumors (Ma et al., 2016; Liu et al., 2018; Ma et al., 2018; Lu and Li, 2020; Yu et al., 2017; Yang et al., 2022; Jin et al., 2022). Considering the difficulty of having sufficient supporting materials, we design **Conflicting Prompts** (**CP**), which ask the LLM to answer "*Based on the comments, analyze whether there are any rebuttals or conflicts, and accordingly verify the authenticity of the news.*" As shown in Tab. 2, with Conflicting Prompts, the LLM correctly predicts the sample is a non-rumor, via noticing that people are all discussing about the event, and no conflicting information is posted.

248 **3.4** CoP: Chain-of-Propagation

Though Rational and Conflicting Prompts can teach LLMs to reason for rumor detection, LLMs still



Figure 2: The LLM has trouble with massive comments.

251

252

253

254

255

256

257

258

259

260

261

262

263

266

270

have trouble when facing too many comments. In real social media scenarios, the number of comments is usually very large. However, LLMs usually have limitations in input length. And more seriously, LLMs can not well reason over long contexts or redundant information (Huang et al., 2023; Xie, 2023). We test above Rational and Conflicting Prompts with GPT 3.5 on the Twitter and Weibo datasets (Sec. 4.1), and results are shown in Fig. 2. When the input length of a sample exceeds the limit, we randomly select part comments belonging to the sample as input. It is clear that, when a sample has too many comments, the detection accuracy drops sharply, which means the LLM can not well reason over excessive information.

In recent research of LLMs, a complex task, which is hard for LLMs, is divided into a series of easy tasks for LLMs to execute (Wei et al., 2022; Wang et al., 2023; Yao et al., 2023; Besta et al., 2024). This enhances LLMs' ability for complex

250



Figure 3: The procedure of the LeRuD approach with GPT 3.5. LeRuD divides the entire propagation information into a Chain-of-Propagation with moderate comments for smoothly reasoning in each step. In each reasoning step, Rational and Conflicting Prompts, which are written in *red italicized words*, are performed. Key clues during LeRuD's reasoning process are highlighted with red underlined words.

tasks. Thus, we propose to divide the propagation information into a **Chain-of-Propagation** (**CoP**), which enables LLMs to conduct reasoning more easily. Considering comments in a time period usually share similar attitudes, we construct CoP in order of comments' posting time (i.e., $[t_1, t_2,]$), and select k comments in each reasoning step. As shown in Fig. 2, the LLM can deal with up to 100 comments easily, so we set k = 100. Reasoning steps in a CoP are run in the same session of LLM, so that latter reasoning steps can refer to the results of former reasoning steps. Moreover, we use the result in the last step as the final prediction, considering the last step aggregates all the information.

4 Experiments

271

274

275

276

281

291

In this section, we conduct extensive experiments to verify the effectiveness of LeRuD.

4.1 Data Preparation

Firstly, we introduce the data preparation process, which enables us to conduct sufficient evaluation. **Datasets**: We combine Twitter15 and Twitter16 as the **Twitter** dataset (in English) (Ma et al., 2017), and also incorporate the **Weibo** dataset (in Chinese) (Ma et al., 2018). The average numbers of comments of each news in Twitter and Weibo are 232 and 816 respectively. We only use the false rumors (as rumors) and non-rumors in the Twitter dataset. On Weibo, we perform Chinese version of prompts, and corresponding examples are translated to English for demonstration in this paper.

292

293

294

295

296

297

300

301

302

303

304

306

307

308

309

310

311

312

313

314

Data Filtering: For smooth and low-cost usage of GPT 3.5, we filter samples with **ethics issues** (such as violence, pornography and discrimination) and **repeated news contents**. We also filter news containing mainly **subjective expressions**, for they are not factually real or fake. Meanwhile, for fair comparison, we need to filter samples with **data leakage risks**, in which GPT 3.5 already exactly knows about the facts in the news. Thus, we direly ask GPT 3.5 to verify the veracity of each news with the Vanilla Prompts in Tab. 1, and only keep samples in which GPT 3.5 has no idea about the facts and can not give answers. In contrast, when GPT 3.5 knows about the fact in a piece of news,

	Twitter			Weibo				
Approach	Accuracy	Precision	Recall	F1-score	Accuracy	Precision	Recall	F1-score
BIGCN	0.8882	0.8739	0.8940	0.8838	0.9258	0.9133	0.9335	0.9233
DDGCN	0.8982	0.8834	0.9078	0.8954	0.9414	0.9390	0.9483	0.9436
GACL	0.9042	0.8884	0.9171	0.9025	0.9354	0.9294	0.9409	0.9351
CICAN	0.9082	0.8855	0.9263	0.9054	0.9390	0.9298	0.9458	0.9377
ARG	0.8902	0.8628	0.8986	0.8803	0.9043	0.9022	0.9089	0.9055
LeRuD	0.9401	0.9156	0.9493	0.9321	0.9809	0.9732	0.9852	0.9792
w/o RP	0.9261	0.9107	0.9401	0.9252	0.9689	0.9635	0.9754	0.9694
w/o CP	0.9142	0.8525	0.9585	0.9024	0.9342	0.8777	0.9901	0.9306
w/o CoP	0.9281	0.9148	0.9401	0.9273	0.9330	0.9249	0.9409	0.9328

Table 3: Experimental results of performance comparison and ablation study.

it will detail the fact and clearly judge the veracity. After data filtering, there are 501 and 836 samples in Twitter and Weibo respectively.

4.2 Settings

315

317

318

319

320

321

322

323

324

326

330

334

337

338

339

340

342

343

345

347

348

The, we detail the settings of our experiments.

Baselines: Firstly, we compare with some stateof-the-art rumor detection models based graph neural networks: **BIGCN** (Bian et al., 2020), **DDGCN** (Sun et al., 2022a) and **GACL** (Sun et al., 2022b). Secondly, we also compare with two rumor detection models with LLMs: **CICAN** (Yang et al., 2023) and **ARG** (Hu et al., 2024). Hyperparameters of baselines are tuned according to the best performances on validation sets.

Metrics: We use Accuracy, Precision, Recall and F1-score for performance evaluation. Precision, recall and f1-score are calculated according to the correct detection of rumors (fake samples).

Zero-shot Setting: LeRuD is performed in zeroshot setting, in which no training data is needed, and the evaluation is done with all samples in the dataset. We use GPT 3.5 as the LLM in LeRuD.

Setting for Baselines: For baselines, we perform 5-fold cross-validation. The whole process is run for 5 times, and alternately in each time, 3 folds of samples are used for training, 1 fold of samples are used for validation, and 1 fold of samples are used for testing. After 5-time running, all samples have been used for performance evaluation, which is consistent with the above zero-shot setting.

4.3 Performance Comparison

In Tab. 3, we compare LeRuD with several state-ofthe-art rumor detection models. Among the three graph-based detection models, DDGCN and GACL obviously outperform BIGCN. With the help of GPT 3.5 for extracting entities, CICAN achieves similar or better performances compared to the graph-based models. ARG considers only news contents, so it underperforms other detection models which capture both news and comments. Meanwhile, with insights generated by LLMs, ARG still achieves high accuracy of about 90%. Obviously, our proposed LeRuD approach achieves the best performances on both datasets evaluated by all metrics. Specifically, LeRuD outperforms the best baselines on Twitter and Weibo by 3.19% and 3.95% respectively. These results strongly demonstrate the superiority of LeRuD. Moreover, LeRuD requires no data for training, so it can consistently achieve promising rumor detection performances in few-shot or zero-shot scenarios, and thus is more suitable for real-world rumor detection.

349

350

351

352

354

356

357

358

359

360

362

363

364

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

4.4 Ablation Study

In Tab. 3, we also demonstrate the ablation study of LeRuD, including three variants: LeRuD w/o RP, LeRuD w/o CP and LeRuD w/o CoP. LeRuD w/o RP denotes LeRuD without Rational Prompts. It decreases the performances of LeRuD by approximately 1%, but still outperforms other baselines. LeRuD w/o CP denotes LeRuD without Conflicting Prompts. The declines of accuracy and f1-score are 2.5% to 4.8%. LeRuD w/o CP achieves the highest recall values, but also achieves the lowest precision values. This is because that, without Conflicting Prompts, LLMs may mistakenly classify non-rumors as rumors. LeRuD w/o CoP Denotes LeRuD without Chain-of-Propagation. The declines of accuracy and f1-score are about 1%and 4.5% on Twitter and Weibo respectively. The

Prompt	There is a piece of news: "Samsung paid Apple one-billion fine, with 20 billion coins in 30 trucks!" Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	I cannot verify the credibility of real-time or specific news events as my training only includes information up until January 2022, and I do not have the capability to access the internet for real-time updates.

Table 4: When asked to estimate the veracity of sample "3484403570896812" in the Weibo dataset, GPT 3.5 has no idea about the fact.

Prompt	There is a piece of news: "Ex-New York Giants safety Tyler Sash, who died at 27, found to have degenerative brain disease." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	As of my last update in January 2022, the information you provided was accurate. Tyler Sash, a former New York Giants safety, passed away at the age of 27 and was found to have chronic traumatic encephalopathy (CTE), a degenerative brain disease often associated with head trauma in contact sports. Therefore, the answer is: B. Real .

Table 5: When asked to estimate the veracity of sample "692142338890661888" in the Twitter dataset, GPT 3.5 provides details about the fact, and confirms it is a non-rumor.





Figure 4: Early rumor detection results evaluated by accuracy with different detection deadlines.

performance declines on Weibo are much larger, because the average number of comments in the Weibo dataset is much larger, according to the statistics in Sec. 4.1. This strongly proves that, the LLM has trouble in reasoning on excessive information. Above observations demonstrate the necessity of the three modules proposed in LeRuD.

4.5 Early Detection

384

386

387

390

396

400

As shown in Fig. 4, we conduct early rumor detection with different detection deadlines. ARG is not included, for it does not model the propagation information. LeRuD constantly and significantly outperforms all the baselines with all the detection deadlines. This shows that, LeRuD can perform well during early stage of propagation, and does not rely on the sufficiency of comment information.

4.6 Impact of Comment Amount

In Fig. 5, we illustrate the accuracy of LeRuD under different numbers of comments. With the help of CoP, LeRuD has no more trouble in conducting reasoning when facing great number of comments. Meanwhile, with more comments, LeRuD has even higher rumor detection accuracy, as there is more

Figure 5: Accuracy of LeRuD under different numbers of comments. More comments, better performances.

information available for analysis. Above observations are totally different from those in Fig. 2. In a word, dividing the complex propagation information into Chain-of-Propagation helps the LLM to conduct reasoning more easily. 407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

4.7 The Data Leakage Risk

As in Sec. 4.1, we filter samples with data leakage risks, in which GPT 3.5 already exactly knows about the facts in the news. Examples without and with data leakage risks are shown in Tab. 4 and 5 respectively. More examples can be found in App. A. According to these examples, we can clearly observe totally different outputs by the LLM between samples with and without data leakage risks. This makes us easy to distinguish samples with data leakage risks. That is to say, after filtering these samples, we can avoid the data leakage risk, and experimental comparison can be fairly done among LeRuD and the baselines.

4.8 Explanations

We can further refer to the explanations generated427by the LLM, to ensure that LeRuD does not take428use of the exact factual knowledge for detecting429rumors. An example is shown in Tab. 6, and more430

Prompt of Step 1	There is a piece of news: "Samsung paid Apple one-billion fine, with 20 billion coins in 30 trucks!" There are comments for the news: "Naughty!" "It's funny!" "Weigh it "Samsung has to count first!" "How creativity!" "Verification needed." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the writing style and commonsense knowledge, the news seems highly unlikely and lacks credibility. Paying a fine with 20 billion coins in 30 trucks sounds more like a creative or satirical piece rather than a factual event. Additionally, such a method of payment raises practical concerns and seems impractical. Analyzing the comments, we can observe a mix of skepticism, humor, and requests for verification. Comments like "Naughty!," "Weigh it," and "Samsung has to count first!" suggest doubts and a need for verification. The comment "How creativity!" acknowledges the imaginative aspect, and "Verification needed" explicitly calls for confirmation. Therefore, the answer is: A. Fake.
Prompt of Step 2	There is a piece of news: "Samsung paid Apple one-billion fine, with 20 billion coins in 30 trucks!" There are comments for the news: "It's hard to say whether it is true or not." "Using \$1 billion in cash is illegal." "Repost" "A truck with a load capacity of 3000 tons is impossible." "Well done!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	Based on the writing style and commonsense knowledge, the news still appears unlikely and lacks credibility. Paying a one-billion dollar fine with 20 billion coins in 30 trucks seems impractical and raises legal and logistical issues. Analyzing the comments, there are several expressions of doubt and skepticism. Comments like "It's hard to say whether it is true or not," "Using \$1 billion in cash is illegal," and "A truck with a load capacity of 3000 tons is impossible" highlight concerns about the feasibility and legality of the described scenario. The comment "Well done!" seems to be sarcastic or ironic, suggesting a lack of belief in the news. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 6: Prompts and outputs of performing the LeRuD approach on sample "3484403570896812" in the Weibo dataset. Key clues during LeRuD's reasoning process are highlighted with <u>red underlined words</u>. LeRuD correctly identifies the sample as a rumor, and the explanations are pretty reasonable from the commonsense perspective.

examples can be found in App. B. According to these explanations, we can summarize the key clues in the reasoning process of LeRuD as follows:

The regularity of writing: When the writing is informal, or the tone seems suspicions, LeRuD tends to judge the samples as rumors (Tab. 21 and 24). When the writing is formal, or consistent with typical reporting styles, LeRuD tends to judge the samples as non-rumors (Tab. 25 and 27-30).

The adequacy of details: When there lack enough details of events in the news, LeRuD tends to judge the samples as rumors (Tab. 17, 21, 22 and 24). When details are rich, LeRuD tends to judge the samples as non-rumors (Tab. 27, 28 and 31).

The plausibility of contents: When the contents seem impractical, or there are conflicts within the contents, or the events do not conform to commonsense logic, LeRuD tends to judge the samples as rumors (Tab. 6 and 24). When the contents seem plausible and rational from the commonsense perspective, or the events have large chances to occur in certain circumstances, LeRuD tends to judge the samples as non-rumors (Tab. 27 - 29).

The attitudes of public: When the attitudes in the comments contains massive skepticism, questioning and rebuttals, or even some comments directly claim the news is fake, LeRuD tends to judge the samples as rumors (Tab. 6 and 17 - 24). When comments seldom express skepticism or rebuttals, and people are discussing about and commenting on the news events themselves, LeRuD tends to judge the samples as non-rumors (Tab. 25 - 32).

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

The consistency of comments: When details in the comments conflict with the news events, or demonstrates the irrationality of the news, LeRuD tends to judge the samples as rumors (Tab. 6, 17, 19 and 20). When details in the comments conform to the news, LeRuD tends to judge the samples as non-rumors (Tab. 25 and 26).

In summary, LeRuD detects rumor on social media via reasoning from commonsense perspectives, and **does not take usage of the LLM's exact prior knowledge about news facts**.

5 Conclusion

In this paper, we investigate to use LLMs for rumor detection on social media, and propose a novel approach called LeRuD. LeRuD designs proper prompts to teach LLMs to reason over important clues in news and comments. Meanwhile, we divide the propagation information of a piece of news into a Chain-of-Propagation for leading LLMs to conduct reasoning step-by-step with reasonable amount of information in each time. Then, we implement LeRuD with the widely-used GPT 3.5, and conduct extensive experiments on the Twitter and Weibo datasets. LeRuD can outperform several state-of-the-art rumor detection models with large margins. In summary, with our efforts, the LLM can well detect rumors on social media.

540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 568 569 570 571 572 573 574 575 576 577 578 579 580 581

582

583

584

585

586

587

588

589

590

538

539

Limitations

490

491 In this paper, though we have taught LLMs to detect rumors on social media, our work still has lim-492 493 itations in the following two aspects. Firstly, running LLMs for rumor detection requires relatively 494 high costs, such as the usage of high-performance 495 GPUs, or the calls of APIs, or inputs of web in-496 terfaces. Secondly, graph neural networks have 497 great ability for modeling structural data, while structural data is hard to be well described in texts. 499 Thus, using LLMs for analyzing complex propaga-500 tion structures, or combining the abilities of LLMs 501 and graph neural networks, has not been explored. 502

Ethics Statement

The experimental datasets are publicly available 504 from some previous works, downloaded via official 505 APIs, and following the regulations of Twitter and Weibo. For privacy issues, the posting accounts are 507 anonymized, and URLs are deleted. When we in-508 put contents and comments into GPT 3.5, there may 509 exists ethics issues, such as violence, pornography 510 and discrimination. Thus, we filter such samples 511 during our experiments, to avoid the dissemination 512 of unethical information. When we input disinformation into GPT 3.5 during experiments, the data 514 might be collected for model training. Thus, we 515 annotate such information as fake at the end the 516 sessions of GPT 3.5, once we input disinformation. 517

References

518

519

521

525

529

530

531

533

534

535

536

537

- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2024. Graph of thoughts: Solving elaborate problems with large language models. In AAAI.
- Tian Bian, Xi Xiao, Tingyang Xu, Peilin Zhao, Wenbing Huang, Yu Rong, and Junzhou Huang. 2020. Rumor detection on social media with bi-directional graph convolutional networks. In *AAAI*.
- Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information credibility on twitter. In *WWW*.
- Canyu Chen and Kai Shu. 2023. Combating misinformation in the age of llms: Opportunities and challenges. *arXiv preprint arXiv:2311.05656*.
- Canyu Chen and Kai Shu. 2024. Can llm-generated misinformation be detected? In *ICLR*.
- Beizhe Hu, Qiang Sheng, Juan Cao, Yuhui Shi, Yang Li, Danding Wang, and Peng Qi. 2024. Bad actor, good

advisor: Exploring the role of large language models in fake news detection. In *AAAI*.

- Yunpeng Huang, Jingwei Xu, Zixu Jiang, Junyu Lai, Zenan Li, Yuan Yao, Taolue Chen, Lijuan Yang, Zhou Xin, and Xiaoxing Ma. 2023. Advancing transformer architecture in long-context large language models: A comprehensive survey. *arXiv preprint arXiv:2311.12351*.
- Md Rafiqul Islam, Shaowu Liu, Xianzhi Wang, and Guandong Xu. 2020. Deep learning for misinformation detection on online social networks: a survey and new perspectives. *Social Network Analysis and Mining*, 10:1–20.
- Yiqiao Jin, Xiting Wang, Ruichao Yang, Yizhou Sun, Wei Wang, Hao Liao, and Xing Xie. 2022. Towards fine-grained reasoning for fake news detection. In *AAAI*.
- Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. 2013. Prominent features of rumor propagation in online social media. In *ICDM*.
- Hongzhan Lin, Pengyao Yi, Jing Ma, Haiyun Jiang, Ziyang Luo, Shuming Shi, and Ruifang Liu. 2023. Zero-shot rumor detection with propagation structure via prompt learning. In *AAAI*.
- Qiang Liu, Feng Yu, Shu Wu, and Liang Wang. 2018. Mining significant microblogs for misinformation identification: an attention-based approach. *ACM Transactions on Intelligent Systems and Technology* (*TIST*), 9(5):1–20.
- Yang Liu and Yi-Fang Wu. 2018. Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. In *AAAI*.
- Yi-Ju Lu and Cheng-Te Li. 2020. Gcan: Graph-aware co-attention networks for explainable fake news detection on social media. In *ACL*.
- Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J Jansen, Kam-Fai Wong, and Meeyoung Cha. 2016. Detecting rumors from microblogs with recurrent neural networks. In *IJCAI*.
- Jing Ma, Wei Gao, Zhongyu Wei, Yueming Lu, and Kam-Fai Wong. 2015. Detect rumors using time series of social context information on microblogging websites. In *CIKM*.
- Jing Ma, Wei Gao, and Kam-Fai Wong. 2017. Detect rumors in microblog posts using propagation structure via kernel learning. In *ACL*.
- Jing Ma, Wei Gao, and Kam-Fai Wong. 2018. Rumor detection on twitter with tree-structured recursive neural networks. In *ACL*.
- Jing Ma, Wei Gao, and Kam-Fai Wong. 2019. Detect rumors on twitter by promoting information campaigns with generative adversarial learning. In *WWW*.

Piotr Przybyla. 2020. Capturing the style of fake news. 592 In AAAI.

591

593

594

595

597

598

599

601

607

610

611

612

613

614

615

616

617

619

631

633

634

641

642

643

- Mengzhu Sun, Xi Zhang, Jiaqi Zheng, and Guixiang Ma. 2022a. Ddgcn: Dual dynamic graph convolutional networks for rumor detection on social media. In AAAI.
- Tiening Sun, Zhong Qian, Sujun Dong, Peifeng Li, and Oiaoming Zhu. 2022b. Rumor detection on social media with graph adversarial contrastive learning. In WWW.
- Ivan Vykopal, Matúš Pikuliak, Ivan Srba, Robert Moro, Dominik Macko, and Maria Bielikova. 2023. Disinformation capabilities of large language models. arXiv preprint arXiv:2311.08838.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In ICLR.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. In NeurIPS.
- Junfei Wu, Weizhi Xu, Qiang Liu, Shu Wu, and Liang Wang. 2023. Adversarial contrastive learning for evidence-aware fake news detection with graph neural networks. IEEE Transactions on Knowledge and Data Engineering (TKDE).
- Wenbei Xie. 2023. Analysis of the reasoning with redundant information provided ability of large language models. arXiv preprint arXiv:2310.04039.
- Weizhi Xu, Junfei Wu, Qiang Liu, Shu Wu, and Liang Wang. 2022. Evidence-aware fake news detection with graph neural networks. In WWW.
- Chang Yang, Peng Zhang, Wenbo Qiao, Hui Gao, and Jiaming Zhao. 2023. Rumor detection on social media with crowd intelligence and chatgpt-assisted networks. In EMNLP.
- Ruichao Yang, Xiting Wang, Yiqiao Jin, Chaozhuo Li, Jianxun Lian, and Xing Xie. 2022. Reinforcement subgraph reasoning for fake news detection. In *KDD*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In NeurIPS.
- Feng Yu, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. 2019. Attention-based convolutional approach for misinformation identification from massive and noisy microblog posts. Computers and Security, 83:106-121.
- Feng Yu, Qiang Liu, Shu Wu, Liang Wang, Tieniu Tan, et al. 2017. A convolutional approach for misinformation identification. In IJCAI.

- Zhenrui Yue, Huimin Zeng, Yang Zhang, Lanyu Shang, and Dong Wang. 2023. Metaadapt: Domain adaptive few-shot misinformation detection via meta learning. In ACL.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. arXiv preprint arXiv:2303.18223.

Table 7: The ratios of samples with data leakage risks.

	Twitter	Weibo
Ratio of Data Leakage	0.1880	0.0434

A Examples about Data Leakage Risks

In Tab. 7, we calculate the ratios of samples with data leakage risks, in which GPT 3.5 already knows about the facts in the news. The ratios are calculated after samples with ethics issues, repeated news contents and subjective news contents have already been filtered. The ratio on the Twitter dataset is much larger than that in the Weibo dataset. This may be caused by the data sources used for training GPT 3.5. News in English has been widely collected, while news in Chinese has been seldom used for the training process of GPT 3.5.

Moreover, we further illustrate several samples, some of which GPT 3.5 exactly knows about the facts, some of which it does not know. In Tab. 8-10, there are samples in which GPT 3.5 has no idea about the facts, and the LLM can not give answers on the veracity of news. In Tab. 11-13, the LLM provides verified descriptions about the facts, and confirms they are rumors. In Tab. 14-16, the LLM provides more details about the facts, and confirms they are non-rumors.

Accordingly to these examples, there are mainly two types of samples with data leakage risks: (1) real time news which has been record in the data for training GPT 3.5; (2) non real-time general knowledge which mainly concentrate in the fields such as health and history. In actual online rumor detection systems, non real-time knowledge (type 2) can actually be helpful to the detection process, while real-time knowledge (type 1) is hard to be remembered by LLMs timely.

B Examples of LeRuD

In this section, we provide more examples of LeRuD's running results in Tab. 17 - 32.

649

650

651

652

644

655

656

657

658

659

660

661

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

Prompt	There is a piece of news: "Famous actor Morgan Freeman passed away in Los Angeles due to a heart attack." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real
	b. Real.
Output	I cannot verify the credibility of real-time or specific news events as my training only includes information
	up until January 2022, and I do not have the capability to access the internet for real-time updates.

Table 8: Sample "3547977337564354" in the Weibo dataset. GPT 3.5 has no idea about the fact.

Prompt	There is a piece of news: "Top North Korean diplomat who disappeared from London embassy has defected to South Korea, Seoul." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	I cannot verify the credibility of real-time or specific news events as my training only includes information up until January 2022, and I do not have the capability to access the internet for real-time updates.

Table 9: Sample "765859710503378944" in the Twitter dataset. GPT 3.5 has no idea about the fact.

Prompt	There is a piece of news: "ALERT: Tsunami advisory issued for Southern California from 3-7 a.m. local
	time. Stay away from ocean and shoreline." Verify the credibility of the news. Please choose the answer
	from the following options: A. Fake, B. Real.
Output	I cannot verify the credibility of real-time or specific news events as my training only includes information
	up until January 2022, and I do not have the capability to access the internet for real-time updates.

Table 10: Sample "765859710503378944" in the Twitter dataset. GPT 3.5 has no idea about the fact.

Prompt	There is a piece of news: "GMO foods weaken reproductive ability. Eating them mean chronic suicide."
	Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	As of my last knowledge update in January 2022, there is no widespread scientific consensus that GMO
	(Genetically Modified Organism) foods weaken reproductive ability or lead to chronic suicide. Numerous
	studies have been conducted on the safety of GMOs, and the majority of scientific organizations, including
	the World Health Organization (WHO) and the National Academy of Sciences, have stated that GMOs
	currently approved for sale are safe to eat. Therefore, the answer is: A. Fake.

Table 11: Sample "8535326047" in the Weibo dataset. GPT 3.5 confirms it is a rumor.

Prompt	There is a piece of news: "The screen car god Paul Walker was not dead and has been rescued." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	As of my last knowledge update in January 2022, Paul Walker, the actor known for his role in the "Fast and Furious" film series, passed away in a car accident on November 30, 2013. Any news claiming that he has been rescued and is not dead would be false. Therefore, the answer is: A. Fake .

Table 12: Sample "3651171602494589" in the Weibo dataset. GPT 3.5 confirms it is a rumor.

Prompt	There is a piece of news: "NASA Confirms Earth Will Experience 15 Days Of Complete Darkness.	
	Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.	
Output	Dutput The news stating "NASA Confirms Earth Will Experience 15 Days Of Complete Darkness" is fake. There	
	is no credible information or scientific basis to support such a claim. If this news were true, it would have	
	had widespread and legitimate coverage from reputable sources. Therefore, the answer is: A. Fake.	

Table 13: Sample "620916279608651776" in the Twitter dataset. GPT 3.5 confirms it is a rumor.

Prompt	There is a piece of news: "The Lushun Massacre, an unforgettable Japanese crime! After the Japanese army captured the fortress of Lushun, they carried out inhumane robbery, massacre, and rape inside the city. Over four days and three nights, approximately 20000 innocent civilians were brutally killed by the Japanese army's knives." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	The Lushun Massacre did occur, and the information in the news is generally accurate. Lushun, also known as Port Arthur, was a Chinese fortress that was captured by the Japanese army during the First Sino-Japanese War in 1894. There were reports of inhumane acts, including robbery, massacre, and rape, during the Japanese occupation of the city. The details may vary, and historical accounts might differ, but the basic premise of the massacre is real. Therefore, the answer is: B. Real .

Table 14: Sample "3911526290906659" in the Weibo dataset. GPT 3.5 confirms it is a non-rumor.

Prompt	There is a piece of news: "Colombian government and left-wing Farc rebel movement agree ceasefire to end decades-long conflict." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	As of my last knowledge update in January 2022, a historic peace agreement between the Colombian government and the Revolutionary Armed Forces of Colombia (FARC) was indeed reached in 2016. The peace deal aimed to end the decades-long armed conflict in Colombia. The agreement included a ceasefire and paved the way for negotiations on various issues, such as disarmament, political participation, and rural development. Therefore, the answer is: B. Real .

Table 15: Sample "689529955294777344" in the Twitter dataset. GPT 3.5 confirms it is a non-rumor.

Prompt	There is a piece of news: "RIP Kenny Baker, the actor behind legendary Star Wars droid R2-D2." Verify the credibility of the news. Please choose the answer from the following options: A. Fake, B. Real.
Output	As of my last knowledge update in January 2022, Kenny Baker, the actor who portrayed R2-D2 in the Star Wars films, passed away in 2016. Therefore, the answer is: B. Real .

Table 16: Sample "764505291853627392" in the Twitter dataset. GPT 3.5 confirms it is a non-rumor.

Prompt of Step 1	There is a piece of news: "The State Council will officially approve Shanghai's request to abolish the establishment of Jing'an District and divide the original administrative regions of Jing'an District into Changning District and Huangpu District in the near future." There are comments for the news: "Changning District has expanded, and we hope for the great development of Changning culture!!!" "But I really like Jing'an District, the most beautiful district name, don't merge it!" "Really?" "I don't understand, but it seems reasonable." "No way, right?!" "Is Jing'an District going to change its past tense???" "Didn't we already refute the rumors a few months ago" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news seems to be presented as a statement of fact, mentioning that the State Council will offi- cially approve Shanghai's request to abolish the establishment of Jing'an District. However, it lacks specific details or sources. The use of phrases like "in the near future" might suggest a lack of concrete information or immediacy. The comments reflect a variety of opinions, ranging from support for the development of Changning District to concerns about losing the name "Jing'an District." Some comments express confusion or skepticism about the news, indicating a lack of clarity or consensus among the readers. The comment mentioning the refutation of rumors a few months ago raises questions about the consistency of the information. Therefore, the answer is: A. Fake .
Prompt of Step 2	There is a piece of news: "The State Council will officially approve Shanghai's request to abolish the establishment of Jing'an District and divide the original administrative regions of Jing'an District into Changning District and Huangpu District in the near future." There are comments for the news: "No I object" "Why merge?" "At that time, the government even stepped forward to refute rumors" "I saw this post three months ago." "I have become a resident of Changning District" "What is the significance of it" "Such fake news would make so many people believe it!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	The news is presented as a statement, but it lacks specific details or sources, making it challenging to verify its accuracy. Phrases like "in the near futur" suggest a lack of immediacy and concrete information. The comments include objections, questions about the merger, references to previous government refutations, and skepticism about the news, indicating potential credibility issues. Several comments express objections, questioning the reasons behind the merger and citing past government refutations. A comment mentions having seen a similar post three months ago, suggesting a lack of novelty or consistency in the information. One commenter claims to have become a resident of Changning District, potentially providing a firsthand perspective. Therefore, the answer is: A. Fake .
Prompt of Step 3	There is a piece of news: "The State Council will officially approve Shanghai's request to abolish the establishment of Jing'an District and divide the original administrative regions of Jing'an District into Changning District and Huangpu District in the near future." There are comments for the news: "Jing'an District is going to become history" "The trend is inevitable, we'll just wait for time." "Well, I've heard of this specific plan, but the problem is it's already a month now, how can we make it." "Really?" "Fake, at least currently fake." "This is not true, is it" "Nonsense." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	The news is presented as a statement, asserting that the State Council will officially approve Shang- hai's request to abolish Jing'an District and divide it into Changning District and Huangpu District. Phrases like "Jing'an District is going to become history" and "The trend is inevitable" suggest a certain level of acceptance or resignation. Comments expressing doubts like "Fake, at least currently fake" and "This is not true, is it" indicate skepticism. The comment labeling it as "Fake, at least currently fake" and others expressing disbelief or calling it "Nonsense" add weight to doubts about the authenticity. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 17: Sample "3498854609883720" in the Weibo dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "Famous actor Morgan Freeman passed away in Los Angeles at 8am local time on Tuesday, February 19th due to heart disease." There are comments for the news: "A great actor!" "Blessing." "So suddenly." "One of my favorite actors" "Oh my goodness, I really like him!" "A great actor" "I hope it's fake news!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news about Morgan Freeman's death seems credible based on the provided information. The de- tails include the specific time, cause of death, and location, which adds a level of specificity often found in legitimate news reports. Additionally, the news is written in a straightforward and factual manner, con- tributing to its credibility. The comments express sentiments of admiration and surprise, with no apparent rebuttals or conflicts. However, the comment "I hope it's fake news!" suggests some skepticism. It's important to note that public figures' deaths can sometimes be subject to rumors or hoaxes, but without additional information, it's challenging to determine the veracity solely based on the comments. Therefore, the answer is: B. Real .
Prompt of Step 2	There is a piece of news: "Famous actor Morgan Freeman passed away in Los Angeles at 8am local time on Tuesday, February 19th due to heart disease." There are comments for the news: "Sincerely hope it's spreading rumors" "Really? The list of last wishes is also funny." "I was deceived! Rumors!" "I can't believe it!" "I will continue to watch your movies." "Is it fake news? Has someone refuted the rumor?" "Morgan Freeman is said to be died again?" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	The news mentions specific details like the time, cause of death, and location, which lends an air of credibility. However, there is a mix of reactions in the comments, with some expressing doubt and questioning the authenticity of the news. Some comments express skepticism or hope that it's a rumor, indicating potential disbelief in the news. The comment "Morgan Freeman is said to be died again?" directly questions the veracity of the news, indicating a potential conflict or doubt. Therefore, the answer is: A. Fake .
Prompt of Step 3	There is a piece of news: "Famous actor Morgan Freeman passed away in Los Angeles at 8am local time on Tuesday, February 19th due to heart disease." There are comments for the news: "Fake news!" "Tsorry to hear that" "Why do so many people like to say that others have died " "Just reposted it and realized it's fake news!!!" "It has been identified as fake news, please don't believe it" "RIP." "This has been confirmed as a rumor." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	The presence of comments like "Fake news!" and "It has been identified as fake news, please don't believe it" strongly suggests a lack of credibility. The comment expressing realization that it's fake after reposting also adds to the skepticism surrounding the news. The majority of comments explicitly label the news as fake and provide reasons for their skepticism. Phrases like "confirmed as a rumor" and "identified as fake news" indicate clear rebuttals and conflicts with the authenticity of the news. There is an absence of supportive comments or credible sources affirming the news. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 18: Sample "3547977337564354" in the Weibo dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "Today, Shanghai Metro Line 10 experienced a malfunction, and the two trains successfully merged, which is truly terrifying." There are comments for the news: "What's wrong with Line 10 again," "When? I also took Line 10 today" "What? How was this achieved?" "Really?" "It's not surprising anymore" "This is really funny." "I think it's the same as the last picture. Is it true?" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news itself <u>seems unusual</u> , as it claims that two trains on Shanghai Metro Line 10 successfully merged due to a malfunction. This scenario raises questions about the technical feasibility and safety of such an event. The comments express a mix of skepticism, concern, and curiosity. Some users question the reliability of Line 10, while others seek more information about the incident. The range of reactions suggests a level of uncertainty among the audience. Without additional details or official sources confirming the incident, it's difficult to verify the authenticity of the news. Therefore, the answer is: A. Fake .
Prompt of Step 2	There is a piece of news: "Today, Shanghai Metro Line 10 experienced a malfunction, and the two trains successfully merged, which is truly terrifying." There are comments for the news: "Don't use the word terror indiscriminately, it's just that the train ahead is malfunctioning. The train behind will push forward for rescue, connect the malfunctioning train, and push it to the designated location to restore the main line operation as soon as possible! This is just a commonly used train rescue method in the subway, there is no such thing as horror." "Terrifying Marketing Weibo attracts attention from the public" "I haven't seen it either" "Scared me" "No knowledge, it's really scary!" "Whoever spreads rumors is terrifying!" "Iron Brave Train Warrior, transform!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	The comment that explains the situation mentions that the merging of trains is a commonly used rescue method in the subway when there is a malfunction. This comment provides a plausible explanation for the incident, suggesting that it might not be as terrifying as the initial news headline implies. Other comments express skepticism or confusion, with some users suggesting that the use of the word "terrifying" might be for marketing purposes or to attract public attention. The explanation provided in one of the comments aligns with common subway rescue procedures, making the incident more plausible and less alarming. The skepticism expressed in some comments, coupled with the suggestion of marketing tactics, raises questions about the accuracy and sensationalism of the initial news report. Therefore, the answer is: A. Fake .
Prompt of Step 3	There is a piece of news: "Today, Shanghai Metro Line 10 experienced a malfunction, and the two trains successfully merged, which is truly terrifying." There are comments for the news: "It's too scary. It seems like we can't sit on Line 10." "Recently, the subway has been broken down frequently." "I personally believe that facing this kind of marketing and Weibo popularization, which takes spreading rumors as its own responsibility, is futile no matter how much it is." "Please allow me to scold you. Such a simple matter is like a car malfunctioning and being towed back for repair. I'm amazed at your understanding." "This is a very normal rescue method, don't use the word terror in the title party." "Really?" "Rescue! make a big fuss!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	Comments express a mix of concern, skepticism, and criticism toward the news. Some users find the news terrifying and express reluctance to use Line 10, while others criticize the use of the word "terror" and suggest that it's a normal rescue method. A few comments provide a more rational perspective, likening the situation to a car malfunction being towed for repair, and expressing amazement at the mis- understanding. The comments present a diversity of opinions, with some users expressing genuine concern and others debunking the severity of the situation. Rational explanations in the comments suggest that the incident might be a routine and safe subway rescue operation, contrary to the alarming tone of the news. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 19: Sample "3465261540071546" in the Weibo dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "NHL expansion – four teams added, Quebec City, Toronto, Seattle, and Las Vegas \$1.4b in expansion fees" There are comments for the news: "Short term thinking uggh" "Phoenix is not moving and sadly neither is Miami" "Hopefully this also helps Seattle get an NBA team back." "Is this done deal? Or is this just speculation? Thank you" "Has 3-4 yrs to build up talent. Can be handled if KHL stops taking guys." "Seattle has one Stanley Cup already. Woot! Lets get another." "YES!!! They need a facility though!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The writing style of the news seems plausible, mentioning specific details such as the addition of four NHL teams (Quebec City, Toronto, Seattle, and Las Vegas) and the associated \$1.4 billion in expansion fees. The information appears coherent and consistent, suggesting a level of credibility. However, without additional sources or context, it's challenging to make a definitive judgment. The comments reflect a mix of opinions and concerns. Some express excitement about the expansion, while others express skepticism or raise questions about specific locations like Phoenix and Miami. The comment "Has 3-4 yrs to build up talent. Can be handled if KHL stops taking guys" implies considerations about player talent development and potential challenges related to the KHL (Kontinental Hockey League). There is no direct rebuttal or conflict in the comments that would immediately discredit the news. Therefore, the answer is: B. Real .
Prompt of Step 2	There is a piece of news: "NHL expansion – four teams added, Quebec City, Toronto, Seattle, and Las Vegas \$1.4b in expansion fees" There are comments for the news: "you think the Seattle City Council will go with an NHL 1st option for the Sodo Arena? I doubt it but I hope so" "I don't take him seriously. Also implies that Seattle can't have NHL AND NBA gotta be one or the other." "if this is true then wouls arena be built prior to sea getting a team?" "Bettman said himself a few weeks ago that Seattle not being considered right now. If they build an arena? Sure. But not before that." "Seriously?" "Toronto has a team" "there will be an imbalance in schedule and conferences/divisions. There's no way." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	The news mentions the expansion of four NHL teams (Quebec City, Toronto, Seattle, and Las Vegas), along with the associated \$1.4 billion in expansion fees. The comments express skepticism and doubts about the news. One comment questions whether the Seattle City Council would prioritize an NHL option for the Sodo Arena, and another doubts the credibility of the source. The comment "Toronto has a team" points out a factual error in the news, as Toronto already had an NHL team before the mentioned expansion. These comments indicate a level of suspicion and raise valid concerns about the accuracy of the news. Therefore, the answer is: A. Fake .
Prompt of Step 3	There is a piece of news: "NHL expansion – four teams added, Quebec City, Toronto, Seattle, and Las Vegas \$1.4b in expansion fees" There are comments for the news: "I understand. But the Maple Leafs have been there for ages. Who will not stay on as a fan for them?" "Seriously, i think we can skip Toronto. One is enough" "I don't know why they'd send another team to Toronto. The broke the first one. Awful organization." "Doubt it" "why? we don't have the market size." "source? Hallucination?" "You're saying this and the #NHL is saying no. Which one of you is lying?" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	The news mentions the expansion of four NHL teams (Quebec City, Toronto, Seattle, and Las Vegas) with \$1.4 billion in expansion fees. The comments express skepticism and doubts about the news. Multiple comments question the inclusion of Toronto in the expansion, citing the Maple Leafs' long-standing presence in the NHL. There are also negative opinions about sending another team to Toronto, with comments criticizing the organization or expressing doubts about the market size. The comment "You're saying this and the NHL is saying no. Which one of you is lying?" directly challenges the credibility of the news source. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 20: Sample "504448356584677377" in the Twitter dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "Months ago, Rupert Murdoch bought National Geographic and fired most of writers. This just hit the stands." There are comments for the news: "Hmm. Doesn't look like the kind of cover one wd expect from the NG. And I don't mean that in a good way." "is there anything he doesn't ruin?!" "this is something sad to learn of. NatGeo was one of the last." "This isn't true, it's an old cover that's doing the rounds" "noooooooo! That's so dumb." "I cancelled my NG subscription years ago. The writing then wasn't much better." "Murdock keeping his stupid sheep in line with ANYTHING but Truth! If there is a hell he should rot there!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the writing style and common sense knowledge, it's important to note that the news lacks specific details such as the source, date, and specific actions taken by Rupert Murdoch. The language used in the news is somewhat informal and lacks the journalistic tone often found in credible news sources. Additionally, the <u>absence of concrete information</u> about the event raises concerns about its credibility. Analyzing the comments, there is a mix of reactions. Some express skepticism, questioning the authenticity of the news, while others react emotionally with disappointment and anger. One comment even suggests that the news might be based on an old cover that is circulating. The varied reactions indicate a lack of consensus among readers. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 21: Sample "676120162018451456" in the Twitter dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "Mike Brown murder trial witness, #DeAndreJoshua fatally shot and set on fire after verdict." There are comments for the news: "This is ridiculous. It's time we take a stand" "Your article doesn't say anything about him being a witness. Can you link to that fact?" "The article doesn't mention he was a witness, got a source? Let the speculation begin." "Speechless." "Google his name and read some articles. Few witnesses gave testimony for fear of this happening." "I had initially done so, however was unable to locate anything substantial. Only speculation." "I'm sure the witnesses were told to keep quiet to avoid situations like this. Not sure if any will be confirmed" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news report mentions a shocking incident involving the murder trial witness, DeAndre Joshua, being fatally shot and set on fire after the verdict. The writing style doesn't provide specific details about the incident or cite any credible sources. The information lacks concrete evidence and seems sensational- ized, which raises concerns about its credibility. The comments highlight skepticism and a lack of con- crete evidence or sources. Some comments express skepticism and request sources to confirm the claim that DeAndre Joshua was a witness. One comment mentions that few witnesses gave testimony for fear of such incidents, indicating a potential environment of intimidation. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 22: Sample "538412662858588160" in the Twitter dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "Jay Z and Beyoncé want to buy rights to the Confederate flag" There are comments for the news: "this isn't true. lol. do better, Complex." "huh???" "lies!!" "wut? RT Jay Z and Beyoncé want to buy rights to the Confederate" "When you got it RT Jay and Beyoncé want to buy rights to the Confederate" "When you got it RT Jay and Beyoncé want to buy rights to the Confederate flag" "idek wut dis mean RT Jay Z and Beyoncé want to buy rights to the Confederate flag" "Whatever man. The stock exchange just halted tradingscrew that flag" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The writing style of the comments suggests skepticism and confusion regarding the news. Phrases like "lol," "huh???" and "wut?" indicate a lack of belief or understanding of the information. Additionally, some comments express disbelief or question the authenticity of the news, such as "lies!!" and "do better, Complex." Overall, the tone of the comments leans towards skepticism. The comments contain a mix of reactions, with some expressing doubt or disbelief, while others seem to accept or support the news. It's essential to note that comments alone may not provide a definitive answer to the authenticity of the news. However, the skepticism expressed in some comments suggests that there might be conflicts or differing opinions among the audience. Based on the above analysis, the news appears to be more likely fake or at least questionable. The comments reflect a degree of doubt and confusion, pointing towards a lack of credibility. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 23: Sample "618805892222468096" in the Twitter dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "The Panjin Public Security Bureau is investigating itself again. Villager Wang Shujie was shot and killed by police officer Zhang Yan, but the investigation found that the police shot legally! Li Hongna was gang raped by the director of Panjin Police Station and local ruffians. Faced with 30000 reposts from netizens, the Panjin Public Security Bureau couldn't sit still and responded to netizens by posting Weibo messages online, urging everyone to be patient and wait for the investigation results. After so many days, the girl's two distress messages were deleted. Where did you lock her up?" There are comments for the news: "Repost" "Seeking the truth." "No wonder I always feel so familiar when I hear the place name Panjin while watching a movie today!" "It is only natural that Wang Shujie was killed. He first illegally injured the construction workers, then attacked the police with a knife and attempted to attack even after the police fired three shots into the sky. Li Hongna's case has been investigated for a long time, and it was a false accusation. All the lies she said were unfounded. Li is a drug addict and has been falsely accused by the police on Weibo." "angry" "Angry" "Panjin Public Security is in chaos."
Output of Step 1	The news contains conflicting information and seems to be presented in a sensationalized manner. The tone of the news raises suspicions, and the use of phrases like "investigating itself again" and "couldn't sit still" may indicate bias. Additionally, the comments show a range of opinions, including both support and skepticism. Without further evidence or reliable sources, it is challenging to determine the credibility of the news. However, the news seems to lack specific details and relies heavily on emotional language, which could be indicative of a less credible source. The comments reflect a mix of opinions, with some expressing skepticism and others defending the actions of the Panjin Public Security Bureau. The conflicting narratives in the comments suggest a lack of consensus among the readers. Some comments question the legitimacy of the claims, while others show support for the authorities. There is no clear consensus among the commenters. Therefore, the answer is: A. Fake .
Final Output	The final prediction of the sample is: It is a rumor.

Table 24: Sample "3507985437364629" in the Weibo dataset. LeRuD correctly predicts it as a rumor.

Prompt of Step 1	There is a piece of news: "The Paris attack has sounded a global anti-terrorism alert." There are comments for the news: "Shocked, just heard!" "The cause of terrorist attacks is not due to the general public, but they bear the tragic consequences" "Only wish for world peace!" "Terrorism will always be the common enemy of all humanity, and any violent attacks against civilians should be condemned and countered.
	Mourning for Paris." "The road to counter-terrorism is still long, and there is still much we need to do." "Rest in peace for the deceased." "Civilian people are innocent." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the information provided, it's challenging to directly estimate the credibility of the news. The statement "The Paris attack has sounded a global anti-terrorism alert" is <u>quite general and lacks specific details about the incident</u> . Credibility assessment would require more information, such as the news source, additional context, and corroboration from reliable sources. Analyzing the comments, there is <u>no evidence of rebuttals or conflicts regarding the authenticity</u> of the news. The comments are <u>more focused on expressing emotions and opinions rather than disputing the news itself</u> . So the answer is undetermined.
Prompt of Step 2	There is a piece of news: "The Paris attack has sounded a global anti-terrorism alert." There are comments for the news: "Watching the report and listening to gunshots, my throat and eyes tightened, my eyes became sore, and everyone's eyes were filled with helplessness, filled with a sense of terror everywhere. Counter terrorism is urgent and cannot be ignored." "The world needs to work together to strike." "Let's pray." "Every country should attach importance to terrorist attacks and take them as a lesson." "Pray that this romantic city will no longer be harmed." "The Paris attack has sounded a global anti-terrorism alert." "May the deceased rest in peace! World peace!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	Based on the provided information, the news seems to be credible and reflects a serious incident – a terrorist attack in Paris. The writing style is coherent and consistent with typical news reports on such events. The comments express shock, sadness, and a call for global unity against terrorism, which aligns with the expected public response to such incidents. There are no apparent rebuttals or conflicts in the comments; rather, they convey a shared sentiment of concern, grief for the victims, and a collective call for action against terrorism. The comments do not raise doubts about the authenticity of the news but rather emphasize the need for global cooperation and solidarity. Therefore, the answer is: B. Real .
Prompt of Step 3	There is a piece of news: "The Paris attack has sounded a global anti-terrorism alert." There are comments for the news: "May everyone in the world have a beautiful heart, and may there be no more disasters in the world." "Is it not good to get along well with others?" "The casualties are continuing to expand, and the most heartbreaking thing is that countries are still debating countermeasures." "Blessing Paris, may the world be at peace." "I still don't want to go out. It feels unsafe everywhere abroad." "No matter what, innocent civilians should not be harmed." "Why are these terrorists so terrifying? May the world be peaceful, please let go of the kind and innocent people." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	Based on the provided information, the news appears to be credible as it describes a serious event – a terrorist attack in Paris. The writing style is consistent with typical news reporting on such incidents. Analyzing the comments, there are no evident rebuttals or conflicts. Instead, the comments express concerns, sympathy, and a desire for peace. They highlight the impact of the attack on innocent civilians and call for measures to ensure safety. These sentiments align with the expected public response to a tragic event like a terrorist attack. Therefore, the answer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 25: Sample "3909084371985775" in the Weibo dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "27 years later, she found a firefighter who sacrificed her life to save her." There are comments for the news: "They are all good examples." "Moved!" "Give a thumbs up to the firefighting hero, and also give a thumbs up to the girl's gratitude!" "A life-saving benefactor. Although it is his responsibility, I also want to express my gratitude." "Repost" "Positive energy!" "Diligence and understanding gratitude are both good virtues!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the information provided, it's challenging to definitively determine the credibility of the news without additional context or specific details about the source. However, there are <u>some positive comments</u> expressing emotions like being moved, giving thumbs up, and acknowledging the heroism of the firefighter and the gratitude of the girl. These comments generally support the authenticity of the news. Without any specific details or evidence to suggest otherwise, it's reasonable to lean towards a more positive assessment. Therefore, considering the writing style and the overall positive sentiment in the comments, the news appears to be more likely real than fake. Therefore, the answer is: B. Real .
Prompt of Step 2	There is a piece of news: "27 years later, she found a firefighter who sacrificed her life to save her." There are comments for the news: "Good people must have good karma!" "With a grateful heart." "A good girl who knows how to be grateful, also has a good family education." "This spirit is worthy of encouragement, and a grateful heart will never be erased. The transmission and promotion of positive energy will not stop." "Being a firefighter for 27 years is not easy!" "Sincere gratitude and gratitude in the heart are sometimes more precious than money!" "A grateful person will have a good life." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	Based on the information provided, the writing style and comments generally seem positive and suppor- tive of the news. The comments express sentiments of gratitude, encouragement, and recognition of the challenges faced by firefighters. There are no apparent signs of rebuttals or conflicts in the comments. Given the positive nature of both the news and the comments, and without specific details suggesting otherwise, it appears more likely that the news is real. Therefore, the answer is: B. Real .
Prompt of Step 3	There is a piece of news: "27 years later, she found a firefighter who sacrificed her life to save her." There are comments for the news: "A servant of the people, a good firefighter who does not seek anything in return; For 27 years, I will never forget to search for benefactors. My parents have taught me well, and I am grateful for the people and things around me!" "So warm, Uncle Firefighter looks so kind." "Those who show gratitude are all good people." "The heroic product of sacrificing oneself to save babies, the true emotions in the world are deeply moving!" "Sun Bingjing visited the fire brigade, and the retired squad leader inferred that he must have rich experience in rescuing people in a dangerous fire scene. Eventually, he found the 27 year old squad leader, Lu Yuming." "Never forget the kindness of saving lives." "A good person lives a peaceful life." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	Based on the information provided, the writing style and comments continue to convey positive sentiments, expressing gratitude, admiration for the firefighter, and recognition of the importance of showing appreciation. There is a specific mention of Sun Bingjing visiting the fire brigade and finding the 27-year-old squad leader, Lu Yuming, which adds detail and specificity to the story. There are no apparent signs of rebuttals or conflicts in the comments, and the tone remains supportive and positive. Considering the positive writing style, lack of conflicts in comments, and the specific details mentioned, the news appears more likely to be real. Therefore, the answer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 26: Sample "3912024620676243" in the Weibo dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "The Chinese Embassy in Mali has confirmed that three employees of China Railway Construction Corporation (CRCC) who went to Mali to negotiate a project with the Ministry of Transportation have died in the attack. They are Zhou Tianxiang, General Manager of China Railway Construction International Group, Wang Xuanshang, Deputy General Manager of China Railway Construction International Group, and Chang Xuehui, General Manager of China Railway Construction International Group West Africa Company." There are comments for the news: "Pray" "There are risks overseas." "One of them met a year ago while eating in Cote, and their lives were so rushed. They were all people who struggled hard in Africa." "I'm too sad, I hope my compatriots can rest in peace soon." "The three Chinese citizens who died in the Mali hotel attack are employees of China Railway Construction Corporation!" "The deceased rest in peace." "The risk of enterprises going global is increasing, and our responsibility to protect citizens overseas is becoming heavier. Our overseas military strike capabilities need to be greatly strengthened." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news appears to be written in a straightforward and factual manner, which adds to its credibility. The details about the individuals involved, their positions, and the circumstances of the incident provide a sense of specificity. The news itself doesn't seem to conflict with common knowledge. Attacks and security risks can occur in various regions, including those with ongoing geopolitical tensions or internal conflicts. The comments express sympathy, sorrow, and concern for the individuals involved. There are no apparent rebuttals or conflicts among the comments, which suggests a collective acknowledgment of the tragic event. Therefore, the answer is: B. Real .
Prompt of Step 2	There is a piece of news: "The Chinese Embassy in Mali has confirmed that three employees of China Railway Construction Corporation (CRCC) who went to Mali to negotiate a project with the Ministry of Transportation have died in the attack. They are Zhou Tianxiang, General Manager of China Railway Construction International Group, Wang Xuanshang, Deputy General Manager of China Railway Construction International Group, and Chang Xuehui, General Manager of China Railway Construction International Group, and Chang Xuehui, General Manager of China Railway Construction International Group, and Chang Xuehui, General Manager of China Railway Construction International Group West Africa Company." There are comments for the news: "Oh my, work is so dangerous." "I used to think these things were so far away from me. Since my brother went to Africa. Every time I see such news, my heart beats tightly, hoping for world peace!!!!!" "Three elites, three families." "Sadness, it's all corporate executives, keep going!" "The reason why we cannot see darkness is because someone is doing their best to block it from where we cannot see. All warriors who fight against terrorist forces, regardless of nationality, faction, or stance, should be remembered by us!" "I saw such heartbreaking news early in the morning!" "Mourning for the deceased! Perhaps the child loses the father, the lover loses the husband, and the white haired person gives the black haired person away." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	The news is presented in a straightforward and factual manner. The details about the victims and the circumstances of the incident are specific, contributing to the credibility of the news. The idea that work in certain regions, especially those with geopolitical tensions or security concerns, can be dangerous aligns with common knowledge. The news about individuals associated with a construction corporation facing risks during overseas negotiations is plausible. The comments express empathy, concern, and sadness for the victims and their families. There are no apparent rebuttals or conflicts in the comments, suggesting a shared sentiment among the commenters regarding the tragic event. Therefore, the answer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 27: Sample "3911539469801102" in the Weibo dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "Man armed with knife on Swiss train starts fire - seven people hospitalised with knife wounds and burns, police say." There are comments for the news: "Mental Illness ?" "They put on my wall and the nearby Antonia and neighboring Edvaldo one boma H and obvious that I wiretaps giant recorders" "hope is not wat am thinking" "To Germany he was working well with the H-bomb was the one who was present in the second war with his amateur radio." "please not another Machete!" "i don't even care anymore, terrorists is now a daily thing." "horrible attack." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the provided information, it's challenging to directly assess the credibility of the news as there are no specific details about the source, and the content is limited to a brief summary. However, the news itself doesn't inherently raise red flags that would suggest it's fake. News about violent incidents can occur, and it is not uncommon for such events to be reported. Analyzing the comments, it seems there is a mix of reactions ranging from speculation about mental illness, conspiracy theories about wiretaps, concerns about terrorism becoming a daily occurrence, and references to historical events. These comments do not provide concrete information about the incident itself but rather express various opinions and thoughts from different perspectives. Therefore, the answer is: B. Real .
Prompt of Step 2	There is a piece of news: "Man armed with knife on Swiss train starts fire - seven people hospitalised with knife wounds and burns, police say." There are comments for the news: "there are bad people in every ethical community" "bloody hell, this is awful :(" "hope is not wat am thinking" "worst killers are mentally ill or have problems in their personal environment" "yes i am afraid, there have been 28,000 Islamic terrorist attacks around the world since 2001" "maybe take a look back in history the crusades were almost the same thing what Isis is doing now conquer land in the name of" "Sorry!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 2	Based on the provided information, the news about a man armed with a knife on a Swiss train starting a fire, resulting in seven people being hospitalized with knife wounds and burns, appears credible. The report is specific and details the incident, which is consistent with the style of legitimate news reporting. Analyzing the comments, there is a mix of reactions expressing concern, empathy, and discussions about the potential motives or background of the attacker. While there is a mention of Islamic terrorist attacks globally, it doesn't directly rebut or conflict with the reported incident on the Swiss train. The comments mostly reflect individual opinions and reactions to the news rather than providing contradictory information. Therefore, the answer is: B. Real .
Prompt of Step 3	There is a piece of news: "Man armed with knife on Swiss train starts fire - seven people hospitalised with knife wounds and burns, police say." There are comments for the news: "learn to read dude. Maybe in the future, read before commenting" "Maybe another follower of the religion of peace that's invaded Europe." "hope is not wat am thinking" "I agree that most of these people do pretend to follow this religion but they could also say they follow the" "Source please. I mean actual sources contributing unhealthy-ness to con/chem trails." "Pathetic" "how very inconvenient" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 3	Based on the provided information, the news about a man armed with a knife on a Swiss train starting a fire, resulting in seven people being hospitalized with knife wounds and burns, appears credible. The report is specific and details the incident, which is consistent with the style of legitimate news reporting. Analyzing the comments, there is some indication of conflicting opinions or skepticism. Comments suggesting a connection to a specific religion and questioning the authenticity by asking for sources may imply some level of doubt or disagreement. However, it's important to note that comments alone may not provide concrete evidence to dispute the credibility of the news. Therefore, the answer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 28: Sample "764497123530375169" in the Twitter dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "NASA astronaut captures a beautiful sunrise from aboard the International Space Station" There are comments for the news: "heavenly sigh" "Beautiful" "Wheres the UFO pictures?" "this is cool as hell" "Nature is beauty" "Thanks to #NASA to capture such a scene !" "woah it looks like the universal studios logo, the old one from 1990's?" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news itself seems plausible, as it describes a NASA astronaut capturing a sunrise from the Inter- national Space Station, which aligns with the activities and capabilities of NASA. The language used is neutral and descriptive, without any sensational or unrealistic claims. Additionally, the idea of astronauts capturing images from the ISS is a common and well-established practice. The comments are generally positive and express awe and appreciation for the captured sunrise. There doesn't seem to be any overt disagreement or conflicting opinions. However, one comment humorously asks about UFO pictures, but this seems more like a playful remark rather than a serious rebuttal. No significant rebuttals or conflicts found; comments are in line with the nature of the news. Therefore, the answer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 29: Sample "693119705469587456" in the Twitter dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "French police arrest woman who was with man carrying weapons and ammunition at Disneyland Paris hotel." There are comments for the news: "good. Paris doesn't play. I do not blame them." "honey grab some guns, ammo and the Quran and we're off to Disneyland" "thank you." "it does matter- a lot!" "The man was armed with the koran, which was the greater danger" "And was packing a fellow ostensibly of some ostensibly specific genealogy?" "people never read anymore. They put he was carrying a Qur'an in the subheadline in the article. Two seconds people." You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the provided information, the news appears credible as it reports <u>a specific event</u> involving the arrest of a woman with a man carrying weapons and ammunition at a Disneyland Paris hotel. The writing style is <u>straightforward and consistent with typical news reporting</u> on law enforcement activities. Analyzing the comments, there are different perspectives expressed. Some comments express support for the arrest and appreciation for the vigilance in handling potential threats, while others inject humor or skepticism into the situation. The variety of responses is not uncommon in online discussions, especially regarding sensitive topics like security incidents. There are <u>no clear rebuttals or conflicts</u> within the comments that challenge the authenticity of the news. Instead, the comments reflect a range of opinions and reactions to the reported event. Therefore, the answer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 30: Sample "692764799315398656" in the Twitter dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "Wow. Anthony Johnson's knockout win is the 3rd-quickest finish in light heavyweight history (0:13)." There are comments for the news: "WWE is still fake" "that's a bad man !!" "his name is rumble" "Johnson is a true professional. Very polite and respectful; kudos to you" "The real hands of stone." "worked for it. Well deserved" "I can't believe there's been TWO others quicker than that!! TWO!!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	Based on the information provided, it is challenging to definitively determine the credibility of the news. The news itself states a specific fact about Anthony Johnson's knockout win being the 3rd-quickest finish in light heavyweight history (0:13). However, without additional context or sources, it's difficult to verify the accuracy of this claim. Analyzing the comments doesn't seem to reveal any direct rebuttals or conflicts regarding the authenticity of the news. Most of the comments express admiration for Anthony Johnson's performance, acknowledge his nickname "Rumble," and commend him for being a true professional. Some express surprise at the speed of the knockout but without disputing the information. Therefore, the an- swer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 31: Sample "767203096472719364" in the Twitter dataset. LeRuD correctly predicts it as a non-rumor.

Prompt of Step 1	There is a piece of news: "A 3-year-old daughter broke a cup and her parents drove away: the cup was broken and the child compensated for it." There are comments for the news: "Such parents are not worthy of having children!" "Such parents are not worthy of having children!" "Can this constitute the crime of abandonment?" "My god!" "I want to say that the little girl is very sad! Very sad" "Oh my goodness How much psychological shadow does this have to cause for children." "Sending such a child back to their parents will lead to more hardships in the future, so be careful not to abuse them!" You need to do: (1) Based on the writing style and the commonsense knowledge, estimate the credibility of the news. (2) Based on the comments, analyze whether there are any rebuttals or conflicts, and then accordingly verify the authenticity of the news. Based on above results, please choose the answer from the following options: A. Fake, B. Real.
Output of Step 1	The news appears to be written in a straightforward manner, describing a situation where a 3-year- old broke a cup, and her parents drove away while she compensated for it. However, without addi- tional details or context, it's challenging to determine the full credibility of the news. The emotional language used in the comments may suggest a strong reaction from the audience, but it doesn't nec- essarily provide evidence regarding the authenticity of the incident. The comments express strong negative sentiments towards the parents, with phrases like "not worthy of having children" and con- cerns about potential psychological impact on the child. However, there are no comments pre- senting alternative viewpoints or challenging the authenticity of the news. The comments mostly focus on expressing shock, sympathy for the child, and concern for her well-being. Therefore, the an- swer is: B. Real .
Final Output	The final prediction of the sample is: It is a non-rumor.

Table 32: Sample "3914951724449159" in the Weibo dataset. LeRuD correctly predicts it as a non-rumor.