# Verification with Transparency: The TrendFact Benchmark for Auditable **Fact-Checking via Natural Language Explanation**

**Anonymous ACL submission** 

#### Abstract

While fact verification remains fundamental, 001 002 explanation generation serves as a critical enabler for trustworthy fact-checking systems by producing interpretable rationales and facilitating comprehensive verification processes. However, current benchmarks exhibit critical limitations in three dimensions: (1) absence of explanatory annotations, (2) English-centric language bias, and (3) inadequate temporal relevance. To bridge these gaps, we present Trend-011 Fact, the first Chinese fact-checking benchmark 012 incorporating structured natural language explanations. TrendFact comprises 7,643 carefully curated samples from trending social media content and professional fact-checking reposi-016 tories, covering domains such as public health, 017 political discourse, and economic claims. It supports various forms of reasoning, including numerical computation, logical reasoning, and common sense verification. The rigorous 021 multistage construction process ensures high data quality and provides significant challenges. Furthermore, we propose the ECS to comple-024 ment existing evaluation metrics. To establish effective baselines for TrendFact, we propose FactISR-a dual-component method in-027 tegrating evidence triangulation and iterative self-reflection mechanism. Experimental results demonstrate that current leading reasoning models (e.g., DeepSeek-R1, o1) have significant limitations on TrendFact, underscoring the real-world challenges it presents. FactISR significantly enhances reasoning model performance, offering new insights for explainable and complex fact-checking.

#### Introduction 1

036

042

The proliferation of counterfeit claims poses significant societal risks, including mass panic, social destabilization, and even armed conflicts, as exem-039 plified by the COVID-19 infodemic(van Der Linden et al., 2020; Aondover et al., 2024). This critical challenge has driven substantial research efforts

in automated fact-checking systems, particularly in developing comprehensive benchmark datasets. The rapid expansion of open datasets has accelerated advancements in AI-powered verification technologies, especially through large language models(Atanasova, 2024; Rani et al., 2023; Wang and Shu, 2023; Bilal et al., 2024; Kao and Yen, 2024). Despite these developments, current fact-checking benchmarks exhibit several critical limitations:

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

First, as an emerging subtask in fact-checking, explanation generation plays a pivotal role in producing interpretable results and enabling a comprehensive fact-checking process. Existing factchecking benchmarks primarily focus on fact verification and evidence retrieval, with minimal attention to textual explanations. Additionally, the current datasets with textual explanations are mostly generated by LLMs and serve as intermediate steps to enhance the verification process. They lack the human validation to ensure the explanations' quality. This limitation hinders the interpretability of fact-checking to some extent.

Second, existing fact-checking benchmarks such as Factcheck-Bench (Wang et al., 2024), Feverous (Aly et al., 2021), and QuanTemp (Venktesh et al., 2024) have primarily focused on English scenarios, and only a few studies focus on other languages, such as Chinese (Lin et al., 2024; Hu et al., 2022). In the vast influx of information on the internet, the Chinese language represents a crucial source, second only to English data. The lack of Chinese in fact-checking benchmarks significantly limits the comprehensive applicability of artificial intelligence in fully addressing real-world fact-checking scenarios. Moreover, these benchmarks pay limited attention to trending topics in real-time, which hinders the development of practical and trustworthy AI.

To that end, we present TrendFact, the first benchmark for fact-checking in Chinese scenarios that incorporates structured natural language explanations. It contains 7614 samples from multiple data sources and includes various reasoning types, including numerical computation, logical reasoning, and common-sense errors, covering extensive areas like health, politics, and economics. Additionally, we introduce a new metric, ECS (Explanation Consistency Score), to mitigate the shortcomings of previous textual metrics that were insufficient for fully evaluating explanations.

086

090

100

101

102

103

104

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

We conduct a rigorous and thorough benchmark construction process. Specifically, we first curate a large-scale set of real-world fact statements from multiple trending platforms and previous research, ensuring the practicality and diversity of the data. Subsequently, we conduct a rigorous filtering process to identify challenging samples. Moreover, we employ significant human efforts to rewrite the remaining samples to incorporate discernible factors. Subsequently, we make adjustments to the evidence and annotate the structured natural language explanations for these samples. Finally, we review the final dataset to mitigate potential biases.

To establish effective baselines for this benchmark, we propose FactISR (Augmenting Fact-Checking via Iterative Self-Reflection), a methodology that systematically combines reasoning adaptation through evidence triangulation with iterative self-reflection enabled by reward decoding. We evaluate TrendFact on five of the strongest existing reasoning LLMs and two current fact-checking methods. We observed that the leading reasoning models exhibit certain limitations in their performance on TrendFact, and our FactISR method shows outstanding improvements when applied to base reasoning models.

In summary, our contributions are as follows:

• We introduce TrendFact, a comprehensive and challenging Chinese fact-checking benchmark that includes natural language explanations. TrendFact integrates real-world trending events with domain-specific factual data, creating a benchmark for fact-checking. To the best of our knowledge, it is the first benchmark to address both fact verification and explanation generation tasks in non-English scenarios.

- We propose FactISR a dual-component method integrating evidence triangulation and an iterative self-reflection mechanism.
- We propose the Explanation Consistency

Score (ECS) to evaluate the factual accuracy and consistency between generated explanations and ground-truth explanations, thereby complementing existing evaluation methodologies. 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

169

170

171

172

173

174

175

176

177

178

179

180

181

• We conducted extensive experiments that reveal limitations in the capabilities of current leading models on TrendFact. Additionally, we compared the performance of existing factchecking methods with FactISR. The results show that FactISR can significantly enhance the performance of reasoning models.

## 2 Related Work

**Fact-checking Benchmarks** Existing factchecking benchmarks can be divided into two primary categories based on their data sources. The first category includes benchmarks derived from Wikipedia data, such as STATPROPS (Thorne and Vlachos, 2017), FEVEROUS (Aly et al., 2021), and Hover (Jiang et al., 2020). The second category focuses on datasets developed by refining knowledge bases from fact-checking websites and existing fact-checkers, such as CLAIMDE-COMP(Chen et al., 2022), DECLARE(Popat et al., 2018), and QUANTEMP(Venktesh et al., 2024). Although these benchmarks are considered comprehensive and challenging, most of them focus primarily on claim design and evidence retrieval, overlooking an essential aspect of the fact-checking task: explanation generation. As large language models play an increasing role in fact-checking, particularly by generating explanations to aid in verification, there is an urgent need for a benchmark that incorporates explanations to assess the reliability of this process. Unfortunately, most existing benchmarks neglect this crucial component, resulting in incomplete evaluations. Furthermore, only a few benchmarks, such as X-Fact (Gupta and Srikumar, 2021), CFEVER (Lin et al., 2024), and CHEF (Hu et al., 2022), address the creation of fact-checking datasets for non-English scenarios. To effectively assess the reliability of rapidly evolving AI-driven fact-checking, the development of a non-English benchmark that includes explanations is crucial.

Automatic Fact-checking Research on automated fact-checking technologies can be broadly categorized into two areas: fact verification and



Figure 1: Overview of TrendFact. The left side illustrates the diverse data sources of TrendFact, the detailed distribution of the data, and the process of dataset construction. The right side displays a fact-checking example from TrendFact that involves complex numerical reasoning.

explanation generation. Fact verification is essen-182 183 tial for the timely evaluation of claims. Previous studies have predominantly focused on areas such 184 as Wikipedia article verification (Rashkin et al., 185 2017; Wang et al., 2022a), table-based verification (Liu et al., 2021; Zhou et al., 2022), and question-187 answering dialogue-based verification (Wang et al., 188 2022b; Zhang et al., 2024). With the advancement 189 of large language models (LLMs), fact-checkers have begun leveraging these models to design factverification applications. For example, (Pan et al., 192 2023) proposes PROGRAMFC to utilize LLMs to generate executable programs, performing fact 194 verification step by step. Explanation generation, 195 on the other hand, addresses a more challenging 196 task of producing interpretable results to support 197 comprehensive fact-checking. It not only verifies claims but also provides an explanation of the principles behind the verification process. While most 200 research has focused on using partial explanations 201 as intermediate steps to facilitate verification, only a few studies have explored how natural language can be used to communicate the accuracy of claims and the rationale behind judgments. For example, 205 (He et al., 2023) proposes to help users correct mis-206 information by generating counter-misinformation responses. Nevertheless, fact verification and explanation generation have largely overlooked the interdependent feedback loop between veracity and 210 explanation, and most of the research has been lim-211 ited to English-language scenarios. 212

## **3** Datasets Construction

### 3.1 Overview

We introduce TrendFact, a comprehensive and highly challenging fact-checking benchmark. Figure 1 provides an overview of the TrendFact pipeline, which includes data sources, benchmark construction process, data attribute distribution, and example data point. TrendFact consists of a total of 7,643 data entries collected from trending events and fact-checkers, encompassing categories such as digital computation, logical reasoning, and common-sense errors. The dataset is further divided into 1,131 multi-evidence entries and 6,512 single-evidence entries. Notably, TrendFact is a fact-checking benchmark to include verification explanations. In contrast to previous benchmarks, which primarily used Wikipedia and factchecking websites as data sources, we have opted for more practical and real-world sources, including the trending platforms and the professionally rigorous CHEF dataset.

213

214

215

216

217

218

219

220

221

222

224

225

226

227

229

230

231

232

233

234

235

236

237

238

239

240

241

242

#### 3.2 Dataset Construction

**Data Collection** We begin by identifying the sources of fact-checking data. Previous studies have primarily relied on Wikipedia or various fact-checking websites. For example, the CHEF dataset was originally compiled from fact-checking websites. To ensure alignment with existing research, we also considered gathering data from these same sources. As a result, we selected CHEF as one

Dataset	#Claims	Source	Explanation Focus	Language
Synthetic Claims				
FEVEROUS(Aly et al., 2021)	87,026	WP	×	English
CHEF(Hu et al., 2022)	10000	FCS	×	Chinese
Hover(Jiang et al., 2020)	26,171	WP	×	English
CFEVER(Lin et al., 2024)	30,012	WP	×	Chinese
STATPROPS(Thorne and Vlachos, 2017)	4,225	FB	×	English
Fact-checker Claims				
CLAIMDECOMP(Chen et al., 2022)	1,250	Politifact	×	English
DeClarE(Popat et al., 2018)	13,525	FCS	×	English
X-Fact(Gupta and Srikumar, 2021)	1,800	FCS	×	Multi
AVeriTeC(Schlichtkrull et al., 2024)	4,568	FCS	×	English
FlawCheck(Kao and Yen, 2024)	30, 349	FCS	$\checkmark$	English
QUANTEMP(Venktesh et al., 2024)	30,012	FCS	×	English
TrendFact	7643	TP	$\checkmark$	Chinese

Table 1: Comparison of TableFact with other fact-checking datasets. In the table, WP refers to Wikipedia, FCS refers to fact-checking websites, TP refers to Trending Platform, and FB refers to FreeBase. By "Explanation Focus", here we refer to dataset contains explanations.

243of our key data sources, as it is one of the few244datasets focused on the Chinese scenario. Further-245more, to diversify the sources in our benchmark,246we collected data from multiple trending platforms247from 2020 to 2024, thus broadening the dataset's248coverage to include a wider range of real-world249fact-checking scenarios. In total, we gathered an250initial set of approximately 500,000 claims.

252

255

259

260

261

263

264

265

267

271

**Data Formalization** Due to varying data sources, the initial data exhibits significant differences in completeness, with each dataset containing some level of interference noise. The CHEF dataset is relatively comprehensive, including claims, evidence, and judgment labels. However, these samples have a distinct characteristic: the evidence often directly incorporates the judgment label (e.g., "... is a rumor"), which reduces the challenge and introduces factual errors as noise. In contrast, the trending data consists solely of statements, which must be transformed into claims. Furthermore, not all trend data can be converted into challenging claims. To address these issues, we implemented a rigorous data enhancement process to reduce noise and improve data completeness.

Given the differences in properties between the CHEF and trend data, we applied separate enhancement processes to each source. For the CHEF data, the enhancement process involved data cleaning, selection of challenging samples, and label correction. Initially, we removed samples with significant noise in the claims. Then, using a combination of LLM voting and manual filtering, we identified more difficult samples. Afterward, we corrected factual errors and refined the original labels, resulting in a more polished version of the dataset, CHEF-EG, which consists of 1,131 claim-evidencelabel triples. For the trend data, the enhancement process included data filtering, claim rewriting, and evidence retrieval. Initially, we evaluated the potential of the data points by combining LLM voting to assess their suitability as challenging factchecking samples. This allowed us to filter out approximately 50,000 samples. Subsequently, we manually filtered the data and selected a few thousand samples. Since the original trend data could not be directly used as claims, we rewrote the selected samples to incorporate identifiable factors. Then, we manually retrieved evidence via Google to validate the truthfulness of the rewritten claims, resulting in a dataset consisting of 6,512 claimevidence-label triples. Finally, we merge the 7643 claim-evidence-label triples from these two sources to form the intermediate state of TrendFact.

**Explanation Generation** Most existing factchecking datasets lack explanation components, and to our knowledge, there is no fact-checking dataset with explanations for non-English scenarios, such as Chinese. With the advancement of 296

297

298

300

272

LLM technology, the role of explanation genera-301 tion in fact-checking has grown significantly. It 302 not only offers general users clear and intuitive reasoning but also allows fact-checkers to assess the accuracy of fact-checking methods based on explanation generation. Consequently, we annotated 306 the intermediate TrendFact samples with detailed 307 explanations, creating the final comprehensive factchecking benchmark consisting of claim-evidencelabel-explanation quadruples. Specifically, we re-310 cruited and formed an annotation team consisting of fourteen graduate students and one doctoral stu-312 dent. They all come from prestigious institutions 313 and have undergone rigorous annotation training. 314 Following the annotation process, we performed a 315 quality review. Annotated data point was evaluated by three independent reviewers, in addition to the 317 original annotator, to ensure high-quality annotations. This approach helps mitigate potential biases 319 that may arise from individual annotators.

### 3.3 Comparisons with Existing Benchmarks

We perform a comparative analysis of TrendFact and existing fact-checking benchmarks, with the results summarized in Table 1. The results confirm our observation that current fact-checking benchmarks lack explanatory statements for fact verification. As a result, they fail to evaluate the explanation generation of fact-checking tasks, leading to an incomplete assessment. Additionally, most research has focused on English scenarios, with limited studies addressing non-English scenarios.

TrendFact, introduced in this work, is the first fact-checking benchmark for non-English scenarios that incorporates detailed explanations and integrates multiple data sources, including trending news websites. Unlike other benchmarks, TrendFact evaluates explanation generation in factchecking tasks, providing greater practical relevance and broader coverage.

## 4 Baselines

322

325

327

329

331

334

335

336

337

342

343

345

347

348

In this section, we introduce the approaches used to comprehensively assess the difficulty and behaviors of TrendFact, including: LLMs, fact-checking methods, and our proposed method, FactISR.

## 4.1 Fact-checking Methods

A benchmark is typically designed to achieve two key objectives: effectiveness and challenging complexity. First, it should be compatible with existing fact-checking methods, demonstrating its validity and acting as a practical testbed for current technologies. Second, it must be sufficiently complex to expose the limitations of existing methods, thus emphasizing its challenging nature. Accordingly, based on these objectives, we selected two representative automated fact-checking methods, PRO-GRAMFC (Pan et al., 2023) and CLAIMDECOMP (Chen et al., 2022), to evaluate TrendFact. We performed generation tasks on the TrendFact dataset for PROGRAMFC's decomposition procedures and CLAIMDECOMP's yes/no sub-questions, using GPT-3.5, with other parameters aligned with the original study. 349

350

351

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

384

385

387

389

390

391

392

393

394

395

397

## 4.2 LLMs

Recent advancements in large language models (LLMs) have given rise to advanced inferencescaled models, such as DeepSeek-R1 2025. These models, also known as inference-optimized LLMs, leverage techniques like reinforcement learning and inference scaling laws to significantly enhance reasoning capabilities. Their superior inference abilities make them suitable for complex tasks, including fact-checking. In fact-checking tasks, these models can process and evaluate claims against evidence, providing veracity assessments and explanations. Then, these advanced inference-scaled LLMs can thus play a pivotal role in improving the accuracy and reliability of fact-checking processes. In this work, we choose the following LLMs for testing TrendFact: QwQ-32B-preview (qwe, 2024), o1-preview (OpenAI, 2024), and DeepSeek-R1. Additionally, we select the powerful general-purpose LLMs, including GPT-4 (OpenAI, 2024), DeepSeek-v3 (Guo et al., 2025), and Qwen-72B-Instruct (qwe, 2024), which do not have enhanced reasoning capabilities, to broaden the evaluation scope. These models have also shown impressive abilities in tackling complex tasks.

## 4.3 FactISR

Inspired by the remarkable performance of techniques like DeepSeek-R1 in reasoning to solve complex tasks, we propose FactISR, an iterative self-reflection reasoning enhancing method. Figure 2 gives an overview of FactISR. It presents an iterative reasoning prompt template based on three key fact-checking features and combined with a reward decoding mechanism to promote base reasoning LLM's performance on fact-checking tasks.



Figure 2: Overview of FactISR. The bottom section illustrates the model decoding process, where the <think> token signifies the commencement of the model's thought process according to the prompt template, concluding at 

at </think>. On the right, it is indicated that the reward for reflection decreases incrementally with each instance of reflection. The Symbol token represents a sequence of tokens, upon hitting which a reflective token reward is provided for the generation of the subsequent token.

**Reasoning Adaptation via Evidence Triangula**tion In order to enhance reasoning LLM, It is intuitive to construct a powerful prompt template. By analyzing the human iterative reasoning process on TrendFact, we identify three key features essential for fact-checking tasks: (1) the relevance between claim and evidence, (2) the consistency between key information and evidence, and (3) the conflict among claim, evidence, label, and explanation. First, we conclude that the higher the relevance, the greater the likelihood of ensuring the accuracy of the claim verification. Second, to ensure the correctness of extracted key information during the reasoning process of LLMs, they should perform self-feedback to validate the consistency between key information and evidence. Finally, LLMs must reflect on their reasoning process, and identify potential conflicts for iterative thinking.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

Then, we perform the construction of the reasoning template combined with these features. As illustrated in Figure 2, it comprises six interconnected modules. Specifically, based on the relevance, an initial verification label is derived. Then, the LLM extracts key information from evidence, assesses its consistency, and uses it to reassess preliminary label while providing detailed explanations. Finally, the LLM reflects on the entire process, identifies conflicts, and provides feedback to guide further reasoning. Detailed prompt information is provided in the Appendix 3.

428Iterative Self Reflection via Reward Decoding429Ideally, the LLM would follow the template to con-430duct iterative reasoning. However, when conflicts431arise clearly during reasoning, the model fails to432acknowledge its errors, leading to premature ter-

mination of reasoning and preventing reflection. To address this, we propose a reward decoding mechanism to affect the LLM's generation process (Figure's 2 right side), improving the probability of it into an iterative reflection phase to resolve conflicts. Specifically, we predefined a target token sequence to serve as the trigger for the reward. Specifically, when the LLM generates tokens that meet this condition, the reward is activated, increasing the probability of the next token being "yes". This approach can be formalized as follows:

433

434

435

436

437

438

439

440

441

442

443

446

447

448

449

450

451

452

453

454

455

456

457

458

Where  $x'_{h+k}$  represents the adjusted token at position h + k in the LLM, while  $x_{h+k}$  is the originally generated token.  $x_{h:h+k-1}$  denotes the continuous token sequence generated by the LLM from position h to h + k - 1, where S is a predefined specific token sequence. When  $x_{h:h+k-1}$  matches S, we calculate the reward vector R (used to assign rewards to certain affirmative tokens) through the initial reward value  $\Delta_0$  and the decay factor  $\gamma$ (with the value range  $0 < \gamma < 1$ ), with the *i* means iteration step.

## 5 Experiment

## 5.1 Experiment Setup

MetricsFor fact verification, We choose four459widely used and complementary metrics:F1-macro, Precision, Recall, and Accuracy.The F1-macro is the macro average of the three classes of462F1 scores.Accuracy is the ratio of correctly veri-fied samples to the total number of samples.For

Methods	Acc	F1	Р	R
PROGRAM-FC	56.55	54.05	54.17	56.62
CLAIMDECOMP	59.35	56.86	56.65	59.41
Qwen-72B-instruct	65.14	60.56	66.97	63.65
DeepSeek-V3	63.74	60.31	66.09	63.96
GPT-40	72.29	69.68	69.02	72.88
DeepSeek-R1	77.92	72.56	73.72	72.64
o1-preview	78.98	75.16	75.13	75.72
QwQ-32B-Preview	65.31	61.76	63.68	65.53
- w/ <i>COT</i>	71.90	68.64	68.21	71.25
- w/ <i>ET</i>	73.28	69.97	69.87	72.41
- w/ <i>ET</i> + <i>ISR</i>	75.45	72.44	72.48	74.59

Table 2: Comparison of FactISR with other baselines on fact verification task. The ET and ISR represent the evidence triangulation and iterative self-reflection components of the FactISR method, respectively. The last column presents the experimental results of the QwQ model after incorporating FactISR and COT.

490

491

492

493

494

465

466

the task of explanation generation, we employ the following evaluation metrics: BLEU-4, ROUGE (including ROUGE-1, ROUGE-2, and ROUGE-L), and BERTScore. Additionally, we introduce a new metric named explanation consistency score (ESC) to identify the consistency level between grounded explanation and generated explanation. The detail about ECS is presented in Appendix. These metrics were selected to comprehensively evaluate the textual similarity, word overlap, and fidelity of the generated explanations.

476 Baselines The methods for testing TrendFact477 have been detailed in Baselines section.

Experimental Settings. Our implementations are based on Pytorch. We leverage the QwQ-32B-Preview (hereafter referred to as QwQ) as the base LLMs. The evaluation for ESC was conducted using GPT-40, while BERTScore evaluations are 482 conducted on bert-base-chinese. The reward vec-483 tor  $\boldsymbol{R}$  and the decay factor  $\gamma$  are set to 20 and 0.1, 484 respectively. The maximum input and maximum 485 output length are set to 16k and 300, respectively. 486 All inference experiments utilized greedy search 487 as the strategy. All GPU-related inferences are 488 executed on  $4 \times A100$  GPUs. 489

## 5.2 Main Results

To provide a deeper insight into our proposed benchmark TrendFact and the fact-checking method FactISR, we present the evaluation results for fact verification and explanation generation in Table 2and Table 3, respectively. For fact verification, we conduct evaluation on three types of baselines. For the explanation generation, we select baselines exclude fact-checking methods, as these methods do not produce readable explanations.

495

496

497

498

499

501

502

503

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

Fact Verification Results For fact-checking methods, the results show that they perform poorly on the TrendFact benchmark. There's a big gap between existing fact-checking methods and others that are established on advanced-scaled inference LLM. This indicates that TrendFact contains many difficult samples and presents a high level of challenge to these methods. The diversity and complexity of the inference tasks within Trend-Fact are also the reasons for limiting their effectiveness. For LLMs, both general-purpose and advanced inference-enhanced LLMs, demonstrate superior performance, highlighting their robustness in addressing complex and challenging fact-chekcing problems. Moreover, advanced inference-enhanced LLMs outperform general-purpose models due to their superior reasoning capabilities. These capabilities allow them to tackle the complex, reasoningbased examples in TrendFact more effectively. Nevertheless, even the highest-performing model, o1preview, failed to exceed 80% across all metrics. This also underscores the significant challenge posed by TrendFact.

Furthermore, compared to CoT (Chain-of-Thought)'s limited improvement, our proposed FactISR, which built on QwQ with two reasoningimproving strategies: Evidence Triangulation (ET) and Iterative Self Reflection (ISR), substantially enhanced the original QwQ's reasoning performance, with an improvement of 10.14 in accuracy and an F1 score increase of 10.68, achieving results on par with the second strongest models, DeepSeek-R1, while surpassing in recall metric. This confirms that our proposed method FactISRT is effective in enhancing LLMs with reasoning capabilities with handling more complex fact-checking tasks. In this study, we employ a rule-based parser to extract the response from the reasoning model's explanation, and then use <sup>1</sup> for obtaining the final veracity label.

**Explanation Generation Results** For LLMs, we observe that almost general-purpose LLMs universally outperform the advanced inference-scaled LLMs. For instance, o1-preview, which performs the best on the fact-verification task, lacks behind

<sup>&</sup>lt;sup>1</sup>https://pypi.org/project/fuzzywuzzy/

Methods	BLEU-4	BERTScore	ROUGE-1	ROUGE-2	ROUGE-L	ECS
Qwen-72B-instruct	0.3366	0.8364	0.6441	0.4589	0.5906	0.7787
DeepSeek-V3	0.3573	0.8432	0.6596	0.4805	0.6087	0.7812
GPT-40	0.2958	0.8270	0.6191	0.4189	0.5561	0.8622
DeepSeek-R1	0.2705	0.8143	0.5832	0.3821	0.5188	0.9115
o1-preview	0.2693	0.8022	0.5602	0.3960	0.5206	0.8986
QwQ-32B-Preview	0.2093	0.7804	0.5330	0.3459	0.4669	0.8198
- w/ <i>COT</i>	0.2790	0.8193	0.6012	0.4040	0.5443	0.8752
- w/ <i>ET</i>	0.2826	0.8057	0.5727	0.3927	0.5182	0.8861
- w/ ET+ISR	0.2815	0.8059	0.5727	0.3910	0.5173	0.8954

Table 3: Comparison of FactISR with other baselines on explanation generation.

the general-purpose LLMs on almost all metrics 544 545 except for the ECS. Our analysis reveals that the reasoning process of inference-scaled models contains a significant amount of critical information, which leads to a clear reduction in the richness of their final outputs. These models tend to produce 549 concise key pieces of information and accurate rea-550 soning results rather than redundant non-critical 551 content. In contrast, general-purpose LLMs excel at generating human-readable text, which produces more diverse and richer outputs. This results in 554 higher scores on text-based metrics, such as 0.61 555 556 on ROUGE-L for DeepSeek-v3. However, since the ECS process includes key information from the 557 reasoning model's inference process, these models can achieve outstanding performance on it. This phenomenon can be also validated through the ECS and the Acc in fact-verification: DeepSeek-R1 and 561 o1-preview achieve higher ECS scores, indicating 562 that their generated explanations are more consis-563 tent with the ground truth explanations and yield 564 more precise results.

For FactISR, the results demonstrate its ability to significantly enhance the performance of the original reasoning model, QwQ, across all metrics. For example, the completed FactISR improves QwQ's score on BLEU-4 and ECS by 7.2 and 7.6, respectively. This improvement makes QwQ surpass the powerful reasoning model o1-preview in overall performance. This indicates that FactISR not only improves the reasoning accuracy of QwQ but also enriches its general-purpose output content. Compared to enhancements from CoT techniques, FactISR provides more substantial improvements in reasoning effectiveness.

566

567

568

571

572

574

575

576

578

## 5.3 Ablation Study

We evaluate the effectiveness of each component within FactISR by incrementally integrating them into QwQ. We conduct ablation experiments for both fact verification and explanation generation. The results (bottom of Table 2 and 3) demonstrate that adding the iterative reasoning template (ET) alone significantly improves the original QwQ's performance in both scenarios. For instance, QwQ's accuracy increased by 8 percentage points, surpassing GPT-40, the best-performing generalpurpose LLM in fact verification. Additionally, QwQ's BLEU-4 score exceeds that of DeepSeek-R1, the top-performing inference-enhanced LLM in explanation generation. When the iterative selfreflection module is incorporated, QwQ's performance further improves. These findings confirm the effectiveness of both FactISR components.

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

### 6 Conclusion

This paper introduces TrendFact, the first Chinese fact-checking benchmark with structured natural language explanations. Through a meticulously designed data construction process, Trend-Fact presents a diverse set of challenging samples requiring complex reasoning. We propose ECS to assess the factual accuracy and consistency between generated and ground-truth explanations. Experiments on automated fact-checking methods and advanced LLMs highlight TrendFact's considerable challenge. Additionally, we propose FactISR, a method that enhances base reasoning models, significantly improving their fact-checking performance.

### 7 Limitations

612

625

627

631

632

634

635

637

642

644

646

656

657

In this paper, we propose a Chinese fact-checking benchmark, TrendFact, which includes structured 614 natural language explanations. However, to im-615 prove its real-time relevance, the claims in our 616 dataset are sourced from trending statements on 617 platforms, which require significant human effort to convert into more complex reasoning claims. Additionally, the evidence and explanations in the benchmark are manually gathered and summarized, 621 resulting in high labor costs. We explore whether, in the future, more powerful LLMs with human-623 like summarization abilities can alleviate this issue.

#### References

2024. Qwen2 technical report.

- Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. Feverous: Fact extraction and verification over unstructured and structured information. *arXiv preprint arXiv:2106.05707*.
- Eric Msughter Aondover, Uchendu Chinelo Ebele, Timothy Ekeledirichukwu Onyejelem, and Omolara Oluwabusayo Akin-Odukoya. 2024. Propagation of false information on covid-19 among nigerians on social media. *LingLit Journal Scientific Journal for Linguistics and Literature*, 5(3):158–172.
- Pepa Atanasova. 2024. Generating fact checking explanations. In *Accountable and Explainable Methods for Complex Reasoning over Text*, pages 83–103. Springer.
- Iman Munire Bilal, Preslav Nakov, Rob Procter, and Maria Liakata. 2024. Generating unsupervised abstractive explanations for rumour verification. *arXiv preprint arXiv:2401.12713*.
- Jifan Chen, Aniruddh Sriram, Eunsol Choi, and Greg Durrett. 2022. Generating literal and implied subquestions to fact-check complex claims. *arXiv preprint arXiv:2205.06938*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Ashim Gupta and Vivek Srikumar. 2021. X-fact: A new benchmark dataset for multilingual fact checking. *arXiv preprint arXiv:2106.09248*.
- Bing He, Mustaque Ahamad, and Srijan Kumar. 2023. Reinforcement learning-based countermisinformation response generation: a case study of covid-19 vaccine misinformation. In *Proceedings* of the ACM Web Conference 2023, pages 2698–2709.

Xuming Hu, Zhijiang Guo, GuanYu Wu, Aiwei Liu, Lijie Wen, and Philip S Yu. 2022. Chef: A pilot chinese dataset for evidence-based fact-checking. *arXiv preprint arXiv:2206.11863*. 664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

- Yichen Jiang, Shikha Bordia, Zheng Zhong, Charles Dognin, Maneesh Singh, and Mohit Bansal. 2020. Hover: A dataset for many-hop fact extraction and claim verification. *arXiv preprint arXiv:2011.03088*.
- Wei-Yu Kao and An-Zi Yen. 2024. How we refute claims: Automatic fact-checking through flaw identification and explanation. In *Companion Proceedings* of the ACM on Web Conference 2024, pages 758–761.
- Ying-Jia Lin, Chun-Yi Lin, Chia-Jen Yeh, Yi-Ting Li, Yun-Yu Hu, Chih-Hao Hsu, Mei-Feng Lee, and Hung-Yu Kao. 2024. Cfever: A chinese fact extraction and verification dataset. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18626–18634.
- Qian Liu, Bei Chen, Jiaqi Guo, Morteza Ziyadi, Zeqi Lin, Weizhu Chen, and Jian-Guang Lou. 2021. Tapex: Table pre-training via learning a neural sql executor. *arXiv preprint arXiv:2107.07653*.
- OpenAI. 2024. Introducing openai o1-preview. Accessed: September 14, 2024.
- Liangming Pan, Xiaobao Wu, Xinyuan Lu, Anh Tuan Luu, William Yang Wang, Min-Yen Kan, and Preslav Nakov. 2023. Fact-checking complex claims with program-guided reasoning. *arXiv preprint arXiv:2305.12744*.
- Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2018. Declare: Debunking fake news and false claims using evidence-aware deep learning. *arXiv preprint arXiv:1809.06416*.
- Anku Rani, SM Tonmoy, Dwip Dalal, Shreya Gautam, Megha Chakraborty, Aman Chadha, Amit Sheth, and Amitava Das. 2023. Factify-5wqa: 5w aspect-based fact verification through question answering. *arXiv preprint arXiv:2305.04329*.
- Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. 2017. Truth of varying shades: Analyzing language in fake news and political fact-checking. In *Proceedings of the 2017 conference on empirical methods in natural language processing*, pages 2931–2937.
- Michael Schlichtkrull, Zhijiang Guo, and Andreas Vlachos. 2024. Averitec: A dataset for real-world claim verification with evidence from the web. *Advances in Neural Information Processing Systems*, 36.
- James Thorne and Andreas Vlachos. 2017. An extensible framework for verification of numerical claims. In Proceedings of the Software Demonstrations of the 15th Conference of the European Chapter of the Association for Computational Linguistics, pages 37–40. Association for Computational Linguistics.

- 718 719 721 723 724 725 726 727 728 729 730 731
- 758

767

771

- 733 734 738 739 740 741 742 743 750 751 753 755

- Sander van Der Linden, Jon Roozenbeek, and Josh Compton. 2020. Inoculating against fake news about covid-19. Frontiers in psychology, page 2928.
- V Venktesh, Abhijit Anand, Avishek Anand, and Vinay Setty. 2024. Quantemp: A real-world opendomain benchmark for fact-checking numerical claims. arXiv preprint arxiv:2403.17169.
- H Wang and K Shu. 2023. Explainable claim verification via knowledge-grounded reasoning with large language models. arxiv preprint arxiv: 231005253.
- Hao Wang, Yangguang Li, Zhen Huang, and Yong Dou. 2022a. Imci: Integrate multi-view contextual information for fact extraction and verification. arXiv preprint arXiv:2208.14001.
- Longzheng Wang, Peng Zhang, Xiaoyu Lu, Lei Zhang, Chaoyang Yan, and Chuang Zhang. 2022b. Qadialmoe: Question-answering dialogue based fact verification with mixture of experts. In *Findings of the* Association for Computational Linguistics: EMNLP 2022, pages 3146-3159.
- Yuxia Wang, Revanth Gangi Reddy, Zain Mujahid, Arnav Arora, Aleksandr Rubashevskii, Jiahui Geng, Osama Mohammed Afzal, Liangming Pan, Nadav Borenstein, Aditya Pillai, et al. 2024. Factcheckbench: Fine-grained evaluation benchmark for automatic fact-checkers. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 14199-14230.
- Xiaocheng Zhang, Chang Wang, Guoping Zhao, and Xiaohong Su. 2024. Li4: Label-infused iterative information interacting based fact verification in questionanswering dialogue. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 10488-10498.
- Yuxuan Zhou, Xien Liu, Kaiyin Zhou, and Ji Wu. 2022. Table-based fact verification with self-adaptive mixture of experts. arXiv preprint arXiv:2204.08753.

#### Details of the ECS Α

Table 4 provides a detailed definition of ECS. The descriptions from top to bottom are as follows: Dual Discrepancy: Both authenticity label misjudgment and fully inconsistent explanatory content. Label Error with Content Consistency: Incorrect authenticity labeling despite congruent explanatory material. Accurate Labeling with Explanatory Divergence: Correct authenticity identification accompanied by conflicting interpretation content. Partial Content Alignment: Proper authenticity classification with only partial consistency in explanatory elements. Full Verification Compliance: Complete congruence between correctly identified authenticity labels and their corresponding explanatory content.

#### B **Prompt of ET**

As shown in Figure 3, the prompt details the following aspects: (1)The relevance between the claim and evidence. (2)The consistency between key information and evidence. (3)The conflict among the claim, evidence, label, and explanation.

772

774

775

776

778

779

780

781

782

783

784

785

786

#### **Example of Rethinking via Reward** С Decoding

Figure 4 illustrates an example of rethinking via reward decoding. Without reward decoding, the model directly outputs a conclusion of no conflict and prematurely ends the reasoning process. Our reward decoding encourages the model to reassess its previous judgment, leading to a reconsideration that ultimately results in the correct outcome.

Label Verification Accuracy	Explanation Consistency	Score	Normalized Score
Misjudgment	Full Discrepancy	1	0.2
Misjudgment	Consistency	2	0.4
Correct Judgment	Content Divergence	3	0.6
Correct Judgment	Partial Consistency	4	0.8
Correct Judgment	Full Consistency	5	1.0

Table 4: Explanation consistency score.



Figure 3: The prompt template of this paper. The upper part is the Chinese prompt template, and the lower part is the translated English prompt template.

<think>第1步/相关性判断:我需要#11月6日至7日南阳市将有雨雪第2步/真实性判断:根据证据较高。因此,初步判断该说法第3步/关键信息是取:从证据14℃~16℃)等。第4步/一致性判断:这些关锅第5步/重新判断真实性:基子</think>	则断说法和证据之间的相关性。说法提到2021年11月河南南阳市将迎来最强寒潮,最低气温降幅超10℃。证据中提到了 3、大风、寒潮天气过程,最低气温降幅10℃~12℃,局部14℃~16℃。这与说法中的内容相符,因此相关性很高。 居,南阳市确实将迎来寒潮,且最低气温降幅超过10℃,甚至局部地区超过14℃。证据来自中央气象台的预报,可信度 法夏真实的。 居中提取关键信息,如寒潮的时间(11月6日至7日)、天气过程(雨雪、大风、寒潮)、气温降幅(10℃~12℃,局部 建信息与说法中的"最强寒潮"和"最低气温降幅超10℃"一致,因此一致性很高。 =上述一致性和证据的可靠性,我确认该说法是真实的。
No Reflection	Reflection
第6步/冲突性判断:我重新 审视说法、证据、标签和 解释之间是否可能存在冲 突。审视结果为: <u></u> 7,思 考结束。 <td>第6步/冲突性判断:我重新审视说法、证据、标签和解释之间是否可能存在冲突。审视结果为:<i>是</i>,说法中提到"最强寒潮",但证据中提到"全国自西北向东南,将有一次大风降温天气过程,大部地区气温将下降8℃~10℃,局地降16℃以上",而南阳市的降幅为10℃~ 12℃,局部14℃~16℃。需要确认"最强寒潮"的表述是否准确。 第7步/进一步分析:为了验证"最强寒潮"的说法,需要了解此次寒潮是否为2021年最强。证据中提到这是"入秋以来最强寒潮",但说 法中没有限定"入秋以来",只是说"最强寒潮"。此外,证据中提到大部地区降温器℃~10℃,局地16℃以上,而南阳市的降温幅度在 00℃~13℃、局部14℃~16℃、周光、根因"干土"物低了、南阳市的路浮温喷在4、4月零不是令在最限。"远程心容相提出只能的任何</td>	第6步/冲突性判断:我重新审视说法、证据、标签和解释之间是否可能存在冲突。审视结果为: <i>是</i> ,说法中提到"最强寒潮",但证据中提到"全国自西北向东南,将有一次大风降温天气过程,大部地区气温将下降8℃~10℃,局地降16℃以上",而南阳市的降幅为10℃~ 12℃,局部14℃~16℃。需要确认"最强寒潮"的表述是否准确。 第7步/进一步分析:为了验证"最强寒潮"的说法,需要了解此次寒潮是否为2021年最强。证据中提到这是"入秋以来最强寒潮",但说 法中没有限定"入秋以来",只是说"最强寒潮"。此外,证据中提到大部地区降温器℃~10℃,局地16℃以上,而南阳市的降温幅度在 00℃~13℃、局部14℃~16℃、周光、根因"干土"物低了、南阳市的路浮温喷在4、4月零不是令在最限。"远程心容相提出只能的任何
<final_output> 证据显示, 2021年11月6日 至7日,南阳市将经历雨 雪、大风、寒潮天气,最 低气温降幅达10℃~12℃, 日期日石初过12℃,</final_output>	第8步/调整判断:由于证据中提到这是"入秋以来最强寒潮",而说法中说的是"最强寒潮",没有限定时间范围,因此可能存在不一 致。此外,证据中没有提供全年寒潮的对比数据,无法确定这次是否是全年最强的寒潮。因此,说法可能过于绝对,证据不足以支持 "最强寒潮"的说法。 第9步/最终判断:基于以上分析,证据支持南阳市将迎来一次强寒潮,气温下降幅度超过10℃,但是否是"最强寒潮"缺乏足够的信息 支持。因此,说法的真实性无法完全确认,证据不充分。
<sup>76] 前                                   </sup>	/ umm2 <final_output>: 证据显示,南阳市在11月6日至7日将迎来一次寒潮天气,气温下降幅度超过10℃,局部地区甚至达到14℃~16℃。然而,说法中提到这是"最强寒潮",而证据中仅描述为"入秋以来最强寒潮",没有提供全年范围内的比较数据。因此,证据不足以验证当前说法的真实性。 </final_output>

Figure 4: An Example of Rethinking via Reward Decoding.