

What are they talking about?

Benchmarking Large Language Models for Knowledge-Grounded Discussion Summarization

Weixiao Zhou^α Junnan Zhu^β Gengyao Li^{βγ}
Xianfu Cheng^α Xinnian Liang^δ Feifei Zhai^{βε} Zhoujun Li^{α*}

^αState Key Laboratory of Complex & Critical Software Environment, Beihang University

^βState Key Laboratory of Multimodal Artificial Intelligence Systems, Institute of Automation, CAS

^γUniversity of Chinese Academy of Sciences ^δByteDance ^εFanyu AI Laboratory

wxzhou@buaa.edu.cn junnan.zhu@nlpr.ia.ac.cn

Abstract

In this work, we investigate the performance of LLMs on a new task that requires combining discussion with background knowledge for summarization. This aims to address the limitation of outside observer confusion in existing dialogue summarization systems due to their reliance solely on discussion information. To achieve this, we model the task output as background and opinion summaries and define two standardized summarization patterns. To support assessment, we introduce the first benchmark comprising high-quality samples consistently annotated by human experts and propose a novel hierarchical evaluation framework with fine-grained, interpretable metrics. We evaluate 12 LLMs under structured-prompt and self-reflection paradigms. Our findings reveal: (1) LLMs struggle with background summary retrieval, generation, and opinion summary integration. (2) Even top LLMs achieve less than 69% average performance across both patterns. (3) Current LLMs lack adequate self-evaluation and self-correction capabilities for this task.¹

1 Introduction

Dialogue summarization aims to distill main topics and interactions from a dialogue into a concise and faithful summary (Jia et al., 2023; Rennard et al., 2023; Kirstein et al., 2024). Current methods typically rely exclusively on dialogue content across dataset construction (Gliwa et al., 2019; Chen et al., 2021; Zhu et al., 2021; Hu et al., 2023), summarization methodologies (Lin et al., 2022; Zhou et al., 2023; Tian et al., 2024), and evaluation paradigms (Wang et al., 2022; Gao and Wan, 2022; Zhu et al., 2023; Gao et al., 2023; Tang et al., 2023, 2024b; Ramprasad et al., 2024). These efforts implicitly assume that "Information in the dialogue is sufficiently complete to support generation of a clearly understandable summary for outside observers."

However, we find this assumption has fundamental limitations and often breaks down, particularly in *discussions involving shared background knowledge*. Such discussions exhibit two main traits: (1) **Information Omission and Implicit Reference**: Participants naturally skip mutually known details and frequently use pronouns or phrases to refer to entities or facts within the contextual background. (2) **Personal Opinion**: Unlike straightforward information interactions in general dialogues, these discussions focus on exchanging viewpoints, with participants expressing personal opinions from various positions or perspectives (see Figure 1). These characteristics make understanding heavily reliant on background knowledge. Consequently, traditional dialogue summarization pattern maintains this dependency, leading to *confusion for outsiders unfamiliar with the context*, leaving them wondering: "What are they talking about?".

To address this issue, we introduce **Knowledge-Grounded Discussion Summarization (KGDS)**, a novel task specifically designed to *combine shared background knowledge with discussion content to produce observer-friendly summaries*. We determine that an effective summary must achieve two complementary objectives: (1) bridge the knowledge gap between outsiders and participants by providing *necessary background-supporting information*; (2) present *clear personal opinions of the participants with clarified implicit references*. We formalize these requirements by modeling the task output as two components: **background summary and opinion summary**. The background summary, which retrieves or condenses relevant contextual information, can be **either extractive or abstractive**, while the opinion summary integrates clarified viewpoints and is **inherently abstractive**. Thus, we define two standardized summarization patterns: **EBS-AOS** and **ABS-AOS** (refer to Figure 1).

To systematically investigate the challenges of KGDS and evaluate LLM capabilities in this task,

*Corresponding Author

¹Our benchmark is available at [zhouweixiao/KGDS](https://github.com/weixiao/KGDS)

Shared Background Knowledge	Knowledge-Grounded Discussion	Traditional Dialogue Summarization Pattern
<p><P1-7>: ...</p> <p><P8-9>: ...</p> <p><P10-12>: ...</p> <p><P13>: Two of Italy's biggest clubs clashed in a famous fixture on Sunday, as AC Milan took on city rival Internazionale in the Derby della Madonnina.</p> <p><P14>: The two clubs share the iconic San Siro stadium, but it was the red half of the city that ended the day celebrating as Matteo Gabbia scored a dramatic last-gasp winner to claim a 2-1 victory and secure bragging rights for Milan.</p> <p><P15>: US national team captain Christian Pulisic gave the Rossoneri the lead in the 10th minute, driving through the heart of the Inter defense before poking the ball past goalkeeper Yann Sommer.</p> <p><P16>: Federico Dimarco leveled things up with an angled finish later in the first half and it looked as though the game was headed for draw, but Milan defender Gabbia had other ideas.</p> <p><P17>: The center-back found the back of the net with a header in the 89th minute to win the game and spark wild celebrations among the Rossoneri faithful.</p> <p><P18-21>: ...</p>	<p>Person1: ... (supported by <P8-9>)</p> <p>Person2: ... (supported by <P8-9>)</p> <p>Person1: I'm dissatisfied that they (referring to AC Milan) only won by one goal against the opponent (referring to Internazionale) at the end of the match. Before the game, I thought they (referring to AC Milan) would definitely sweep the opponent (referring to Internazionale) with a large score. (supported by <P13-14>)</p> <p>Person2: Well, I'm excited about the ups and downs of the score during this match (referring to the match between AC Milan and Internazionale). I think his (referring to Christian Pulisic's) first goal laid the foundation for the team's (referring to AC Milan's) victory, and his (referring to Matteo Gabbia's) last-minute winner completely extinguished the opponent's (referring to Internazionale's) hopes for victory. (supported by <P13-17>)</p>	<p>... Person1 expressed dissatisfaction with the narrow margin of their team's victory, having expected a more decisive win against the opponent. Person2, on the other hand, found the fluctuating scoreline exciting, highlighting the significance of the first goal and the impact of the last-minute goal that secured the team's victory and dashed the opponent's hopes.</p> <p>Outside Observer Confusion (What are they talking about?):</p> <ul style="list-style-type: none"> what are the names of the two teams participating in the match? what is the final score of the match? which team ultimately wins the match? how does the scoreline change throughout the match? which player scores the first goal of the match? which player scores the last-minute goal of the match?
	<p>Evaluation for KGDS EBS-AOS Pattern</p> <p>1. Automatic Evaluation for Extractive Background Summary (EBS)</p> <p>Method: Background-Supporting Paragraph (BSP) Matching</p> <p><P8> (X, Missing) <P9> (✓) <P10> (✓, Irrelevant) <P13> (✓) <P14> (✓) <P15> (X, Missing) <P16> (X, Missing) <P17> (✓)</p> <p>Metrics: BSP Recall=0.57 Precision=0.80 F1-score=0.67</p> <p>2. Fine-Grained Evaluation for Abstractive Opinion Summary (AOS)</p> <p>Method: LLM-based Clear Atomic Opinion (CAO) Verification</p> <p>Person1 is dissatisfied that AC Milan only won by one goal against Internazionale at the end of the match. (✓)</p> <p>Person2 thinks that Christian Pulisic's first goal laid the foundation for AC Milan's victory. (X, Error Detection — Implicit Reference Incorrectly Clarified in AOS)</p> <p>Metric: CAO Recall (Coverage)=0.56</p> <p>3. Overall Performance Evaluation for EBS-AOS Pattern</p> <p>Method: Geometric Mean</p> <p>Metric: $OP_{GE-EBS-AOS} = (BSP_{F1} \times CAO_R)^{1/2} = \underline{0.61}$</p>	<p>Evaluation for KGDS ABS-AOS Pattern</p> <p>1. Fine-Grained Evaluation for Abstractive Background Summary (ABS)</p> <p>Method: LLM-based Inner and Outer Atomic Fact Verification</p> <p>Inner Key Background-Supporting Atomic Fact (KBSAF):</p> <p>AC Milan and Internazionale played against each other. (✓)</p> <p>AC Milan won the match 2-1. (X, Missing)</p> <p>Outer Background-Nonsupporting Atomic Fact:</p> <p>The two sides met at City's Etihad Stadium. (✓, Irrelevant)</p> <p>Arsenal is two points behind City. (X)</p> <p>Metrics: KBSAF Recall (Coverage)=0.42 Precision (Focus)=0.58 F1-score (Overall)=0.49</p> <p>2. Fine-Grained Evaluation for Abstractive Opinion Summary (AOS)</p> <p>Method: LLM-based Clear Atomic Opinion (CAO) Verification</p> <p>Metric: CAO Recall (Coverage)=0.44</p> <p>3. Overall Performance Evaluation for ABS-AOS Pattern</p> <p>Method: Geometric Mean</p> <p>Metric: $OP_{GE-ABS-AOS} = (KBSAF_{F1} \times CAO_R)^{1/2} = \underline{0.46}$</p>

Figure 1: An overall example. **Gray blocks** in paragraphs represent the crucial background details omitted by discussion. **Cyan blocks** in discussion indicate referential pronouns or phrases with their corresponding referents. **Pink blocks** in traditional dialogue summary denote content that may cause confusion for outside observers. The summarization patterns of KGDS are formalized as EBS-AOS and ABS-AOS. Our evaluation framework is able to comprehensively and accurately assess sub-summaries and both patterns via fine-grained and interpretable metrics.

we make two significant contributions. **First**, we construct the first KGDS benchmark containing 100 high-quality, multi-domain samples from realistic news discussions through strict consistency-controlled annotation protocols. Each sample comprises structured shared background knowledge, multi-turn human discussion, and expert-annotated evaluation components at multiple granularities. **Second**, we propose a novel hierarchical evaluation framework with fine-grained and interpretable metrics (as shown in Figure 1). Specifically: (1) For extractive background summary, we utilize retrieval metrics to evaluate recall, precision, and f1-score of supporting paragraphs. (2) For abstractive background summary, inspired by the *minimal granularity* of atomic facts (Liu et al., 2023b) and the *excellent human-alignment* of LLM-based fact verification (Wei et al., 2024), we assess quality via verification of *key supporting fact inclusion* and *nonsupporting fact exclusion*, proposing three metrics that multi-dimensionally quantify background coverage, focus, and summary overall performance. (3) For abstractive opinion summary, we evaluate quality by verifying *coverage of clear atomic opinions*, where each opinion is *minimized* and *implicit references clarified*. For uncovered opinions, we define five error types and perform fine-grained er-

ror detection. (4) For pattern overall performance, we use *geometric mean* for evaluation, which aligns with the KGDS workflow and reflects the joint fulfillment rate of both sub-summaries.

We evaluate 12 leading LLMs under *single-turn structured-prompt* (Li et al., 2023) and *multi-turn self-reflection* (Shinn et al., 2023) paradigms. Our analysis reveals several critical insights: (1) LLMs struggle with precision-recall trade-off in background summary retrieval; (2) LLMs frequently miss key facts and retain irrelevant facts in background summary generation; (3) LLMs encounter difficulties with clarifying implicit references in opinion summary integration; (4) Even advanced LLMs achieve less than 69% average performance across both patterns; (5) LLMs lack sufficient self-evaluation and -correction capabilities for KGDS.

Our contributions are as follows: (1) We introduce a novel KGDS task and establish two standardized summarization patterns. (2) We construct the first benchmark for this task and propose a comprehensive hierarchical evaluation framework with interpretable metrics. (3) Our extensive evaluation reveals specific challenges that current LLMs face with KGDS, providing valuable guidance for future improvements in coarse-grained retrieval, fine-grained generation, and knowledge integration.

2 Related Work

Dialogue Summarization. Dialogue summarization is widely used in practical applications (Lin et al., 2021; Hu et al., 2023). Earlier approaches concentrated on feature modeling (Chen and Yang, 2021; Fang et al., 2022) and pre-training methods (Zou et al., 2021; Zhong et al., 2022a). The emergence of LLMs has greatly enhanced dialogue summarization (Wang et al., 2023a; Tian et al., 2024; Lu et al., 2025; Zhu et al., 2025), with recent research focusing on evaluating their effectiveness in this field (Liu et al., 2024b; Ramprasad et al., 2024; Tang et al., 2024b). However, these studies often *overlook the importance of background knowledge*.

Knowledge-Grounded Dialogue Generation. The task involves generating responses that interact with external knowledge sources like documents (Wang et al., 2023b, 2024), databases (Zhou et al., 2018; Zhang et al., 2020), or retrieval-augmented systems (Lewis et al., 2020; Chen et al., 2024). The main distinctions between this task and our KGDS are: (1) **Conversation characteristics:** The former focuses on knowledge-intensive dialogues, while KGDS deals with background-sparse discussions. (2) **Task objectives:** The former enhances response quality through knowledge, whereas KGDS utilizes background for complementary summarization.

Summarization Evaluation. Earlier similarity-based metrics like Rouge (Lin, 2004), BERTScore (Zhang et al., 2019), and MoverScore (Zhao et al., 2019) are coarse-grained and often misaligned with human judgment. NLI- and QA-based approaches improve faithfulness evaluation by checking factual consistency (Goyal and Durrett, 2020; Laban et al., 2022; Zha et al., 2023) or generating relevant questions (Wang et al., 2020; Fabbri et al., 2022; Zhong et al., 2022b), but they are limited to faithfulness and require specialized training. Recently, LLM-based reference-free evaluators (Liu et al., 2023a; Chen et al., 2023; Wang et al., 2023c; Shen et al., 2023; Fu et al., 2024) have been used, yet they still lack interpretability and fine-grained assessment. To address these, some studies (Song et al., 2024a, 2025; Lee et al., 2024; Yang et al., 2024; Scirè et al., 2024; Wan et al., 2024) utilize atomic fact units (Liu et al., 2023b; Min et al., 2023) and LLM-based claim verification (Tang et al., 2024a; Wei et al., 2024; Song et al., 2024b), achieving more fine-grained evaluation. Among these, **FineSurE** (Song et al., 2024a) is most relevant to our work.

However, its metrics lack dimensional consistency (Completeness is measured at the fact-level, while Conciseness is sentence-level), limiting them to isolated sub-aspect evaluation. In contrast, our coverage and focus metrics are both atomic-fact-level, allowing direct F1-score computation for overall quality assessment while maintaining interpretability. Our method can theoretically be adapted for **general query-focused summarization**, requiring only the maintenance of inner and outer fact sets.

Summarization Benchmark. Current research often evaluates LLMs across various summarization tasks, including dialogue (Tang et al., 2024b), news (Goyal et al., 2022; Zhang et al., 2024), multi-document (Huang et al., 2024), multilingual (Ye et al., 2024), and query-focused (Yang et al., 2023; Liu et al., 2024a). Moreover, some studies investigate inherent issues such as hallucinations (Belem et al., 2024; Bao et al., 2024) and biases (Ravaut et al., 2024; Chhabra et al., 2024). Unlike these, our work uniquely focuses on LLMs within KGDS, evaluating their capabilities in background retrieval, generation, and opinion integration.

3 Task Formulation

EBS-AOS Pattern. The output comprises an extractive background summary (EBS) and an abstractive opinion summary (AOS):

$$B_e, O \leftarrow f(K, D, I_e), B_e \subseteq K \quad (1)$$

where K is the shared background knowledge of participants, D is the discussion grounded in K , and I_e is an instruction for the EBS-AOS pattern. f is the LLM-driven summarization system. B_e is the EBS, defined as **extractive background supporting chunks for D from K** , where chunks represent text sequences with a predefined granularity (e.g., sentences, paragraphs, sections, or window-sliding segments). O is the AOS, defined as **clear personal opinions of the participants with clarified implicit references**.

ABS-AOS Pattern. The output contains an abstractive background summary (ABS) and an abstractive opinion summary (AOS):

$$B_a, O \leftarrow f(K, D, I_a) \quad (2)$$

where K , D , f , and O retain their definitions from Eq. (1). The instruction I_a is specific to the ABS-AOS pattern. B_a is the ABS, defined as **abstractive background supporting information for D from K** .

4 Benchmark

4.1 Preliminary

Scenario Setting. We establish the benchmark scenario as a two-participant discussion about news content for two reasons: (1) news discussions are highly prevalent in daily communication, and (2) news summarization is a representative subfield of automatic summarization research.

News Collection. We collect 100 multi-domain (*i.e.*, business, sports, and world) event-rich news articles from Google News² as shared background knowledge. To ensure data credibility and avoid potential contamination in LLMs (Deng et al., 2023; Li and Flanigan, 2024), we control the time span of news between Mar and Sep 2024. Additionally, we retain the original paragraph structure and define **paragraph-as-chunk** as the minimum granularity under the EBS-AOS pattern (see Table 4 in Appendix B for more statistics).

Expert Recruitment. We recruit four PhD candidates specializing in NLP to engage in our annotation work. Each pair of experts form a collaborative group to conduct the full-process annotation.

4.2 Discussion Construction

The construction follows the sequence of **read, understand, then discuss**. Initially, for each news article, we require two experts to independently read and thoroughly understand its content. The purpose of this preliminary step is to ensure knowledge consistency. Subsequently, we ask them to engage in a discussion to exchange viewpoints. This process is open-ended, meaning the initiator of the discussion is chosen at random, and the discussion topics can encompass any events, facts, or sub-information contained within the news article (see Table 5 in Appendix B for discussion statistics).

4.3 Annotation and Evaluation for EBS

Back-Supporting Paragraphs Annotation. For each paragraph in the news article, we ask the two experts involved in the discussion: "*Is this paragraph the background-supporting paragraph for discussion?*" and require them to answer with either "yes" or "no". We then apply strict consistency control, retaining only paragraphs consistently annotated as "yes" (*i.e.*, **back-supporting**) or "no" (*i.e.*, **back-nonsupporting**). For inconsistent annotations (*i.e.*, one "yes" and one "no"), we determine

them to be *boundary paragraphs* and remove from the original news article (see Appendix B for annotation statistics). Our investigation in Section §6.1 shows that 12 LLMs, with an average win rate of 88.9%, prefer expert-consistently-annotated EBS over their own extractions, highlighting the excellent quality of our annotations.

Evaluation Metrics. EBS assessment involves matching LLM-extracted paragraph indices. We use information retrieval metrics: **Back-Supporting Paragraph (BSP) Recall, Precision, and F1-score** (refer to Appendix C for metric details).

4.4 Annotation and Evaluation for ABS

Key Back-Supporting Atomic Facts Annotation. This process consists of two steps: **atomic fact decomposition** and **expert annotation**. We use GPT-4o³ to decompose back-supporting paragraphs into candidate atomic facts, following the principles of *indivisibility, independence, and declarativity* (the instruction is provided in Appendix D.2). Next, we ask two experts for each fact: "*Is this fact the key background-supporting fact for discussion?*" and require them to answer either "yes" or "no". We employ the same consistency control to determine **key back-supporting atomic facts** (*i.e.*, two "yes") and **non-key facts** (*i.e.*, one or zero "yes" responses) (refer to Appendix B for annotation statistics).

Back-Nonsupporting Atomic Facts Annotation. This process is automated and includes two steps: **atomic fact decomposition** and **conflicting fact masking**. We employ the same approach to decompose back-nonsupporting paragraphs into initial atomic facts. Next, we *mask*⁴ facts that can be *inferred from back-supporting paragraphs* to obtain **back-nonsupporting atomic facts**. Conflicts typically occur when identical or similar facts appear in both back-supporting and -nonsupporting paragraphs at the atomic granularity level. For instance, two events have the same timestamp, but one is the background event while the other is not. By masking such facts, we construct non-overlapping back-supporting and -nonsupporting fact sets, avoiding any external verification problems (see Appendix B for annotation statistical data).

Fine-Grained Multi-Dimensional Metrics. Let \mathcal{B} denote the LLM-generated ABS, $\mathcal{K} = \{k_i\}_{i=1}^m$

³Version: gpt-4o-2024-08-06

⁴The masking process is consistent with fact verification, and the conflicting facts are not considered in ABS evaluation.

²<https://news.google.com>

the set of key back-supporting atomic facts, $\mathcal{N} = \{n_j\}_{j=1}^p$ the set of back-nonsupporting atomic facts. The LLM-based fact verification function is defined as:

$$\phi : (\mathcal{B}, f) \rightarrow \{0, 1\}, f \in \mathcal{K} \cup \mathcal{N} \quad (3)$$

where 1 and 0 represent whether the fact f can be inferred from \mathcal{B} .

Metric α . KBSAF Recall Quantifies ABS **coverage** of key back-supporting atomic facts:

$$\text{KBSAF}_R = \frac{1}{m} \sum_{i=1}^m \phi(\mathcal{B}, k_i) \quad (4)$$

Metric β . KBSAF Precision Measures ABS **focus** on key back-supporting atomic facts:

$$\text{KBSAF}_P = \frac{\sum_{k \in \mathcal{K}} \phi(\mathcal{B}, k)}{\sum_{f \in \mathcal{K} \cup \mathcal{N}} \phi(\mathcal{B}, f)} \quad (5)$$

Metric γ . KBSAF F1-score Evaluates ABS **overall performance** by comprehensively considering both coverage and focus:

$$\text{KBSAF}_{F1} = 2 \cdot \frac{\text{KBSAF}_P \cdot \text{KBSAF}_R}{\text{KBSAF}_P + \text{KBSAF}_R} \quad (6)$$

4.5 Annotation and Evaluation for AOS

Clear Atomic Opinions Annotation. The annotations are created entirely by experts and involve three steps: **main opinion extraction**, **implicit reference clarification**, and **clear opinion atomization**. First, we require experts to extract the main opinions from their respective utterances in the discussion. Next, we ask them to identify referential pronouns and phrases within these opinions and clarify them (*i.e.*, through *anaphora resolution* or *information supplementation*) using referents (*e.g.*, *entities*, *facts*, or *events*) from the news article to produce clear opinions. Finally, we instruct the experts to decompose these clear opinions into atomic opinion units, adhering to the principles of *indivisibility* and *independence* (see Appendix B for detailed annotation statistics).

Fine-Grained Quantification Metric. Let \mathcal{O} be the LLM-generated AOS, and $\mathcal{C} = \{c_i\}_{i=1}^n$ be the set of clear atomic opinions. The LLM-based opinion verification function is defined as:

$$\psi : (\mathcal{O}, c) \rightarrow \{0, 1\}, c \in \mathcal{C} \quad (7)$$

where 1 and 0 indicate whether the opinion c can be inferred from \mathcal{O} .

Metric δ . CAO Recall Assesses AOS **coverage** of clear atomic opinions:

$$\text{CAO}_R = \frac{1}{n} \sum_{i=1}^n \psi(\mathcal{O}, c_i) \quad (8)$$

Fine-Grained Error Detection. Unlike query-focused ABS, the integrativity of AOS determines that it is almost impossible to exhaust the space of erroneous atomic opinions. As a result, defining precision and F1-score for AOS evaluation is impractical. However, we can identify why certain opinions are not covered by AOS through fine-grained error detection. Specifically, for opinions that cannot be inferred from AOS, we define five error types and conduct LLM-based error classification: **implicit reference unclarified**, **implicit reference incorrectly clarified**, **opinion misattribution**, **opinion fact inconsistency**, and **opinion sentiment distortion** (detailed error definitions are provided in Appendix D.7).

4.6 Overall Performance Evaluation

We use the **Geometric Mean** to assess pattern-level overall performance for two reasons: (1) Its multiplicative property reflects the joint fulfillment rate of the background and opinion summaries, emphasizing the holistic synergy between these two components. (2) The root normalization ensures the dimensional consistency of the component metrics. Specifically, we define the overall performance metrics for the EBS-AOS and ABS-AOS patterns as:

$$\text{OP}_{GM:\text{EBS-AOS}} = (\text{BSP}_{F1} \cdot \text{CAO}_R)^{1/2} \quad (9)$$

$$\text{OP}_{GM:\text{ABS-AOS}} = (\text{KBSAF}_{F1} \cdot \text{CAO}_R)^{1/2} \quad (10)$$

5 Evaluation Setup

Evaluated LLMs. We evaluate 12 leading LLMs, covering both the most advanced and lightweight variants⁵: **GPT-4o**, **GPT-4-turbo**, **GPT-4o-mini**, **Claude 3 Opus**, **Claude 3.5 Sonnet**, **Claude 3.5 Haiku**, **Gemini 1.5 Pro**, **Llama-3.1-405B**, **Mistral Large**, **DeepSeek-V3**, **Qwen-Max**, **GLM-4-Plus**. All model sources are listed in Appendix E.

Solution Paradigms. We benchmark LLMs for KGDS under two interaction paradigms: (1) **Single-turn Structured-prompt**: A standard prompt with well-structured (Li et al., 2023) contains *input content*, *input definition*, *task description*, *output definition*, and *return format* (in Appendix D.3). (2)

⁵Our evaluation began on January, 2025, and all LLMs used the latest official API versions available at that time.

Model Name	KGDS BenchMark (single-turn structured-prompt and multi-turn self-reflection)					
	EBS-AOS Pattern			ABS-AOS Pattern		
	BSP (R-P-F1)	CAO (R)	OP (GM)	KBSAF (R-P-F1)	CAO (R)	OP (GM)
GPT-4o	73.39 \uparrow 1.56-88.10 \uparrow 0.34- 78.12 \uparrow 1.27	76.18 \uparrow 0.54	76.24 \uparrow 1.04	63.57 \uparrow 1.08-61.73 \uparrow 1.62- 58.34 \uparrow 0.60	69.42 \uparrow 0.21	61.09 \uparrow 0.47
GPT-4-turbo	73.03 \downarrow 7.72-87.42 \downarrow 1.67-77.14 \downarrow 3.80	72.73 \uparrow 0.45	73.94 \downarrow 1.47	46.28 \downarrow 5.78-54.42 \uparrow 5.33-47.26 \downarrow 1.64	58.32 \uparrow 1.90	48.28 \downarrow 0.74
GPT-4o-mini	76.23 \downarrow 0.30-67.71 \uparrow 0.31-67.92 \downarrow 0.03	29.08 \downarrow 1.00	34.24 \downarrow 0.97	33.75 \uparrow 0.48-45.07 \uparrow 0.52-36.51 \uparrow 0.52	27.74 \downarrow 0.41	24.64 \downarrow 0.21
Claude 3 Opus	65.60 \uparrow 2.21-85.01 \downarrow 4.57-72.07 \downarrow 1.31	75.20 \downarrow 2.19	72.09 \downarrow 2.45	51.10 \uparrow 1.03-74.59 \downarrow 1.13- 58.03 \downarrow 0.91	69.95 \downarrow 2.41	60.72 \downarrow 1.29
Claude 3.5 Sonnet	80.36 \uparrow 2.42-87.93 \downarrow 0.62- 82.33 \uparrow 0.46	74.93 \uparrow 5.30	77.12 \downarrow 4.45	40.74 \uparrow 5.84-65.18 \downarrow 0.79-47.75 \uparrow 3.51	58.55 \uparrow 1.29	48.83 \uparrow 2.60
Claude 3.5 Haiku	62.69 \downarrow 8.39-76.01 \uparrow 0.37-66.02 \downarrow 5.66	40.35 \downarrow 9.13	46.02 \downarrow 10.1	37.42 \downarrow 6.90-48.72 \downarrow 4.65-39.91 \downarrow 6.63	29.64 \downarrow 5.71	27.22 \downarrow 4.48
Gemini 1.5 Pro	84.64 \downarrow 4.38-82.25 \uparrow 4.43- 79.73 \uparrow 1.36	76.86 \downarrow 0.53	76.71 \uparrow 0.49	50.40 \downarrow 4.13-52.33 \uparrow 0.28-48.42 \downarrow 2.11	69.09 \downarrow 6.08	54.83 \downarrow 4.45
Llama-3.1-405B	79.09 \uparrow 2.72-77.96 \downarrow 1.60-75.16 \uparrow 0.36	64.25 \downarrow 0.01	67.58 \uparrow 0.01	38.19 \downarrow 1.08-58.47 \uparrow 7.38-43.39 \uparrow 1.12	52.98 \downarrow 2.44	43.79 \uparrow 0.11
Mistral Large	68.81 \uparrow 1.40-78.55 \downarrow 0.18-71.07 \uparrow 0.89	60.82 \downarrow 2.62	63.63 \downarrow 1.59	53.91 \downarrow 3.05-56.78 \downarrow 0.32- 52.57 \downarrow 1.94	46.24 \downarrow 1.71	46.36 \downarrow 2.26
DeepSeek-V3	86.98 \downarrow 1.19-73.83 \uparrow 1.37-75.66 \uparrow 1.07	64.98 \uparrow 1.84	66.50 \uparrow 1.61	47.64 \uparrow 0.65-42.12 \uparrow 0.60-42.17 \uparrow 0.65	56.02 \uparrow 0.77	44.09 \uparrow 0.72
Qwen-Max	74.59 \downarrow 5.53-79.86 \downarrow 1.30-73.79 \downarrow 2.78	60.93 \downarrow 1.24	64.34 \downarrow 1.59	46.49 \downarrow 0.02-53.59 \downarrow 0.35-46.87 \downarrow 0.17	45.29 \downarrow 1.54	40.87 \downarrow 0.23
GLM-4-Plus	80.28 \downarrow 1.03-72.81 \uparrow 0.61-71.17 \uparrow 0.34	69.34 \downarrow 1.73	67.55 \downarrow 0.43	43.50 \downarrow 1.77-46.93 \uparrow 0.43-41.64 \downarrow 0.61	55.80 \downarrow 3.17	43.54 \downarrow 1.26

Table 1: Main evaluation results. All metrics are *macro-averaged* % (i.e., all samples have the same weight). \uparrow and \downarrow respectively indicate performance **increases** and **decreases** for self-reflection after structured-prompt. For each comprehensive metric (i.e., BSP_{F_1} , $KBSAF_{F_1}$, CAO_R , and OP_{GM}), we highlight the performance of top-3 models.

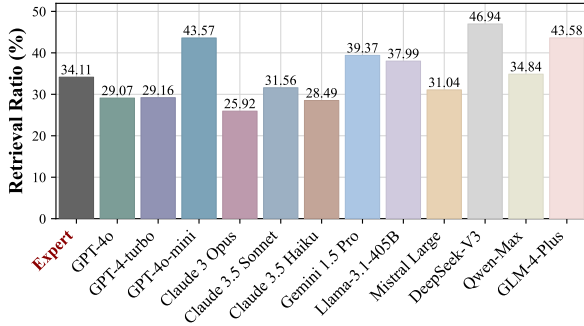


Figure 2: Paragraph retrieval ratios % of LLMs. The majority of LLMs can be classified as either *conservative* (ratio < 30%) or *open* retrievers (ratio > 38%).

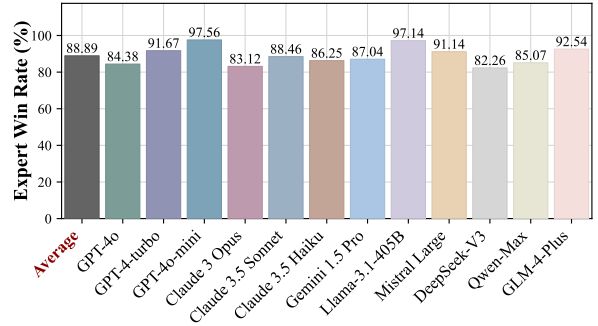


Figure 3: Expert win rates % in EBS feedback judgment. Our annotations achieve an average win rate of 88.89%. All 12 LLMs exceed an 82% rate in preferring our EBS.

Multi-turn Self-reflection: A second-round *self-reflection* (Shinn et al., 2023) instruction with step-by-step *chain-of-thought* (Wei et al., 2022) reasoning after the structured prompt (in Appendix D.4).

Verifier and Detector. We use GPT-4o⁶ to perform fact and opinion verification as well as error detection due to its excellent consistency with human judgment (Song et al., 2024a,b). All prompts are provided in Appendices D.6 and D.7.

6 Analysis

In this part, we primarily analyze the commonalities and differences among LLMs in KGDS. Sections §6.1, §6.2, and §6.3 respectively discuss the performance of background summary, opinion summary, and both patterns under structured-prompt. §6.4 explores the impact of LLM self-reflection.

6.1 Background Summary Analysis

LLMs are moderate but imbalanced retrievers for EBS. From Table 1, we find that most LLMs

achieve moderate retrieval performance ($BSP_{F_1} \in [71, 82]$) and lightweight models (i.e., GPT-4o-mini and Claude 3.5 Haiku) perform poorly. However, from BSP_R and BSP_P , we observe significant polarization among LLMs, which indicates distinct retrieval strategies: *some prioritize precision at the cost of recall* (e.g., GPT-4-turbo), *while others do the exact opposite* (e.g., DeepSeek-V3). Such imbalance reveals that the current LLMs struggle with the precision-recall trade-off. Moreover, we investigate the paragraph retrieval ratio (Figure 2) and identify that most LLMs exhibit either under- or over-retrieval, which is consistent with imbalance.

LLMs prefer expert consistently annotated EBS over their own extractions. We conduct an investigation to verify whether our EBS is the true gold standard. Specifically, for EBS extracted by LLM that is inconsistent with expert-annotated, we instruct the LLM to perform a second-round feedback judgment to determine which EBS is more aligned with its definition (the instruction is provided in Appendix D.5). As shown in Figure 3, all 12 LLMs unanimously and overwhelmingly prefer

⁶Version: gpt-4o-2024-11-20

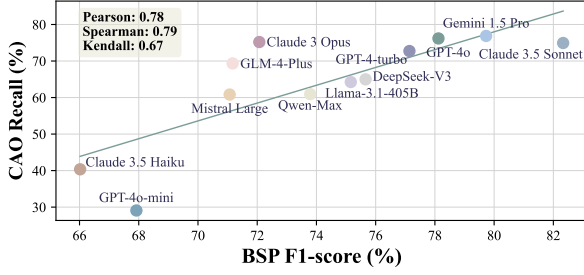


Figure 4: Visualization and metrics of the correlation between BSP_{F_1} and CAO_R under EBS-AOS pattern.

our annotations, with the average win rate of 89%, indicating that: (1) *Our expert-annotated EBS can be regarded as reliable ground truth.* (2) *The current LLMs possess latent awareness of optimal retrieval boundaries but struggle to implement them.*

LLMs are inadequate generators for ABS. As shown in Table 1, all LLMs demonstrate unsatisfactory performance ($KBSAF_{F_1} \in [37, 58]$). The weak $KBSAF_R$ ($R_{avg} = 46.08$) reveals that LLMs often omit key facts, while the mediocre $KBSAF_P$ ($P_{avg} = 54.99$) indicates persistent inclusion of irrelevant facts. This dual failure reflects the fundamental deficiencies of LLMs in meeting the requirements of **coverage** (*i.e.*, completeness) and **focus** (*i.e.*, conciseness). Moreover, we also observe polarization among LLMs: *some prioritize precision at the cost of recall* (*e.g.*, Claude 3 Opus), *while others attempt to balance both* (*e.g.*, GPT-4o). Unlike EBS, we did not find any extreme recall-oriented models, indicating that LLMs tend to be either conservative or balanced in ABS generation.

LLMs are better retrievers than generators for background summary. As presented in Table 1, all LLMs exhibit higher BSP metrics compared to their KBSAF counterparts⁷. This gap between retrieval and generation indicates that LLMs *possess stronger in-context recognition capabilities for coarse-grained paragraphs than fine-grained facts.*

6.2 Opinion Summary Analysis

Investigating correlation variables influencing AOS quality. From Table 1, we find that CAO_R decreases as the background summary quality declines (*i.e.*, $BSP_{F_1} \rightarrow KBSAF_{F_1}$) across all LLMs during pattern transitions.

Additionally, as shown in Figures 4 and 5, CAO_R is highly positively correlated with BSP_{F_1} and $KBSAF_{F_1}$ among models

⁷BSP and KBSAF metrics are directly comparable, as paragraph and atomic fact differ only in granularity, and there is no difference in mathematical form between the twin metrics.

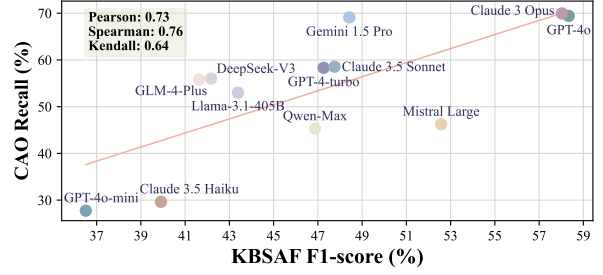


Figure 5: Visualization and metrics of the correlation between $KBSAF_{F_1}$ and CAO_R under ABS-AOS pattern.

in independent patterns. These findings indicate that *a high-quality background summary facilitates knowledge integration in AOS generation.* Meanwhile, individual differences suggest that the *model-inherent integration ability is also a key factor affecting AOS quality.* For instance, Gemini 1.5 Pro and Claude 3.5 Sonnet show relatively stronger (*i.e.*, above the line in Figure 5) and weaker (*i.e.*, below the line in Figure 4) integration capabilities in ABS-AOS and EBS-AOS patterns, respectively.

As for absolute performance, the most advanced models (highlighted LLMs in Table 1) only achieve an average of 72% CAO_R across both patterns, indicating that numerous opinions are being incorrectly integrated. Moreover, we observe that lightweight models (*i.e.*, GPT-4o-mini and Claude 3.5 Haiku) exhibit significant flaws in CAO_R , suggesting that *AOS integration is sensitive to parameter scale.*

AOS error detection analysis. As shown in Table 2, a significant proportion of integration errors are concentrated in IRU and IRIC across all LLMs in both patterns, indicating that *LLMs struggle to effectively and correctly clarify implicit references during AOS integration.* Furthermore, although the ratios of OFI, OSD, and OM errors are relatively lower, their presence still emphasizes the inherent deficiencies of factual inconsistency, sentiment distortion, and incorrect speaker assignment in LLMs.

6.3 Overall Performance Analysis

Performance stratification among LLMs. We observe that OP_{GM} metrics exhibit distinct hierarchies in independent patterns and cross-pattern (refer to Figures 6, 8, and 9 in Appendix A). This suggests that **the performance of LLMs improves in a stepwise manner rather than continuously as their intelligence advances in KGDS.** Moreover, even the best-performing models achieve less than 69% average overall performance across both patterns (see Figure 9), highlighting that KGDS remains a challenging task for current SOTA LLMs.

Model Name	AOS Error Detection (single-turn structured-prompt and multi-turn self-reflection)									
	EBS-AOS Pattern					ABS-AOS Pattern				
	OFI	OSD	IRU	IRIC	OM	OFI	OSD	IRU	IRIC	OM
GPT-4o	8.57 \uparrow 1.09	6.19 \downarrow 0.88	26.67 \uparrow 2.80	55.71 \downarrow 2.09	2.86 \downarrow 0.92	4.07 \downarrow 1.12	2.96 \uparrow 1.10	50.74 \downarrow 1.66	40.37 \uparrow 1.33	1.86 \uparrow 0.35
GPT-4-turbo	7.23 \downarrow 0.38	9.24 \downarrow 1.18	40.96 \uparrow 1.38	39.76 \downarrow 0.65	2.81 \uparrow 0.83	6.74 \uparrow 0.40	3.23 \uparrow 2.33	56.33 \downarrow 0.51	31.00 \downarrow 1.11	2.70 \downarrow 1.11
GPT-4o-mini	0.79 \uparrow 0.30	0.32 \uparrow 0.15	81.80 \downarrow 0.08	16.93 \downarrow 0.37	0.16 \downarrow 0.00	1.74 \downarrow 0.79	0.79 \uparrow 0.15	73.69 \downarrow 0.30	23.61 \uparrow 1.11	0.17 \downarrow 0.17
Claude 3 Opus	6.96 \uparrow 0.54	2.17 \downarrow 0.09	32.17 \uparrow 0.75	56.96 \downarrow 1.96	1.74 \uparrow 0.76	6.61 \uparrow 0.64	1.18 \uparrow 0.99	46.69 \uparrow 0.05	43.19 \downarrow 2.61	2.33 \uparrow 0.93
Claude 3.5 Sonnet	6.67 \uparrow 0.24	5.78 \downarrow 1.78	50.22 \uparrow 2.51	35.11 \downarrow 1.29	2.22 \uparrow 0.32	3.00 \uparrow 3.13	3.54 \downarrow 0.20	61.85 \downarrow 8.65	29.97 \uparrow 5.13	1.64 \uparrow 0.59
Claude 3.5 Haiku	6.81 \downarrow 2.01	4.09 \downarrow 0.83	55.45 \uparrow 4.07	32.10 \downarrow 1.23	1.55 \downarrow 0.00	3.91 \downarrow 0.10	1.47 \downarrow 0.41	70.52 \downarrow 2.64	23.94 \uparrow 3.31	0.16 \downarrow 0.16
Gemini 1.5 Pro	6.28 \downarrow 0.91	7.73 \uparrow 2.03	37.68 \uparrow 3.78	46.86 \downarrow 3.45	1.45 \downarrow 1.45	3.62 \uparrow 0.05	6.52 \downarrow 0.40	63.41 \uparrow 4.17	26.09 \downarrow 4.38	0.36 \uparrow 0.56
Llama-3.1-405B	6.92 \downarrow 1.00	1.89 \uparrow 1.23	55.35 \uparrow 0.10	34.59 \downarrow 0.95	1.25 \uparrow 0.62	3.83 \uparrow 1.18	0.47 \uparrow 1.81	69.62 \uparrow 2.13	25.36 \downarrow 5.77	0.72 \uparrow 0.65
Mistral Large	6.67 \uparrow 1.35	3.33 \downarrow 0.39	52.78 \downarrow 0.91	36.94 \downarrow 1.91	0.28 \uparrow 1.86	3.73 \downarrow 0.12	2.28 \downarrow 0.48	64.52 \uparrow 5.82	29.05 \downarrow 5.60	0.42 \uparrow 0.38
DeepSeek-V3	4.79 \uparrow 1.61	2.24 \downarrow 0.22	53.35 \downarrow 5.54	38.98 \uparrow 4.12	0.64 \uparrow 0.03	5.99 \downarrow 0.67	2.08 \uparrow 0.05	61.98 \downarrow 2.94	28.91 \uparrow 3.80	1.04 \downarrow 0.24
Qwen-Max	7.54 \uparrow 0.13	2.90 \uparrow 0.79	43.48 \uparrow 1.12	45.22 \downarrow 2.32	0.86 \uparrow 0.28	4.47 \uparrow 0.05	3.25 \downarrow 0.30	65.45 \uparrow 1.94	26.22 \downarrow 1.66	0.61 \downarrow 0.03
GLM-4-Plus	10.94 \downarrow 1.09	3.02 \uparrow 1.72	40.38 \downarrow 5.71	44.53 \uparrow 4.74	1.13 \uparrow 0.34	5.23 \uparrow 2.60	1.93 \uparrow 1.99	55.65 \uparrow 0.22	35.81 \downarrow 4.74	1.38 \downarrow 0.07

Table 2: Fine-grained error detection results. OFI, OSD, IRU, IRIC, and OM represent the five error types: Opinion Fact Inconsistency, Opinion Sentiment Distortion, Implicit Reference Unclassified, Implicit Reference Incorrectly Clarified, and Opinion Misattribution. All metrics are *error-proportion* % when opinions cannot be inferred from AOS. \downarrow and \uparrow respectively represent error ratio **decreases** and **increases** for self-reflection after structured-prompt.

Cross-pattern stability of LLMs. By analyzing the performance gap between the two patterns (see Figures 6, 7, and 8 in Appendix A), we find that **different LLMs excel in different patterns**. For example, Claude 3.5 Sonnet exhibits a significant gap (28.29%), while Claude 3 Opus demonstrates a relatively smaller gap (11.37%), indicating stronger cross-pattern stability.

LLMs perform better with EBS-AOS than ABS-AOS for KGDS. As presented in Table 1, all LLMs achieve superior OP_{GM} in EBS-AOS pattern than ABS-AOS, which is attributed to the more complete and accurate background summary and higher-quality opinion summary. Hence, we suggest prioritizing EBS-AOS in real-world KGDS for better implementation.

6.4 Self-Reflection Impact

Self-reflection does not essentially affect the performance of LLMs on KGDS. From Table 1, we observe that the performance fluctuations (*i.e.*, \uparrow or \downarrow) after self-reflection are limited (the maximum average fluctuation across all metrics is 2.30 for $CAO_{R:ABS-AOS}$). This means that self-reflection does not fundamentally influence the capacities of LLMs in KGDS, revealing two key limitations: (1) *LLMs lack sufficient self-evaluation capacities for KGDS.* (2) *The reflection strategies of LLMs struggle to provide excellent reasoning paths for KGDS.*

Self-reflection makes LLMs more conservative or open for background summary. As shown in Table 1, most LLMs demonstrate polarization in the increases and decreases between BSP_R and BSP_P , as well as $KBSAF_R$ and $KBSAF_P$. This

suggests different reflection strategies: some prioritize precision (*e.g.*, DeepSeek-V3 with $BSP_R \downarrow 1.19$, $BSP_P \uparrow 1.37$), while others prioritize recall (*e.g.*, Claude 3 Opus). A few models exhibit simultaneous increases (*e.g.*, GPT-4o with $BSP_R \uparrow 1.56$, $BSP_P \uparrow 0.34$) or decreases (*e.g.*, Qwen-Max), indicating more balanced or weaker reflection abilities.

Self-reflection makes LLMs more risk-averse for opinion summary. As presented in Table 2, most LLMs show a decrease in IRIC and an increase in IRU (*e.g.*, GPT-4o with IRIC $\downarrow 2.09$, IRU $\uparrow 2.80$). This indicates that self-reflection drives LLMs to adopt more cautious strategies when clarifying implicit references — *preferring to leave references unclarified rather than risk incorrect clarification*. However, excessive conservatism in clarifying implicit references harms AOS quality, as evidenced by the widespread decrease in CAO_R in Table 1.

7 Conclusion

In this study, we introduce the KGDS task, which aims to create observer-centric summaries by integrating shared background knowledge with discussions. We establish the first benchmark for KGDS and propose a novel hierarchical evaluation framework with fine-grained, interpretable metrics. Our evaluation of 12 leading LLMs reveals significant challenges in background summary retrieval, generation, and opinion summary integration, with even the most advanced models achieving less than 69% average performance across both patterns. These findings highlight the need for future advancements in LLM capabilities for coarse-grained retrieval, fine-grained generation and knowledge integration.

Limitations

Our study presents several limitations that merit discussion and future exploration:

Domain Generalization. Knowledge-grounded discussions and their confusing summaries are prevalent in both open-domain and private scenarios. Our benchmark is centered on open-domain news discussions, which may limit its applicability to private scenarios (*e.g.*, internal meetings, medical consultations). However, shared background knowledge and the related discussions in private contexts are challenging to obtain and are less common and representative than our open benchmark. Future work will focus on validating the performance of LLMs in handling KGDS across diverse private scenarios.

Language Diversity. Our current work is focused on English, creating a monolingual benchmark that may not fully address the challenges of multilingual KGDS. This limitation could impact model performance in languages with different linguistic structures or sociocultural contexts. Nevertheless, English is the most prevalent language in both academic research and real-world LLM applications, making it a reasonable starting point. We plan to explore multilingual and cross-lingual KGDS in future research.

Ethics Statement

Our KGDS benchmark is developed with careful consideration of ethical implications. The news articles used as shared background knowledge are publicly accessible through Google News and have been carefully reviewed to ensure they do not contain sensitive or harmful information. Explicit consent was obtained from all expert participants involved in discussion construction and data annotation, with a clear explanation of the research purpose and data usage protocols. Our evaluation methodology prioritizes factual accuracy and opinion fidelity to minimize the risk of hallucination propagation in real-world applications. Potential misuse risks, such as generating misleading summaries through improper background-discussion combinations, are mitigated through technical safeguards in our released code and explicit usage guidelines. Researchers using this benchmark should adhere to responsible AI principles, especially when applying similar techniques to sensitive domains like healthcare or legal discussions. We provide the annotation and evaluation costs as:

Annotation Cost for KGDS BenchMark. Our annotation process is multi-step, fine-grained, and highly complex. For each sample in our benchmark, the average annotation time is 2.6 hours (including news reading and understanding, discussion construction, back-supporting paragraphs, key back-supporting atomic facts, back-nonsupporting atomic facts, and clear atomic opinions annotation). We pay each annotation expert a wage of \$9 per hour (above the minimum wage standard), and the total annotation cost for the KGDS benchmark is \$4680 ($2.6 \text{ hours per sample} \times \$9 \text{ per hour} \times 2 \text{ experts per sample} \times 100 \text{ samples} = \4680).

Evaluation Cost for KGDS BenchMark. The evaluation cost of API calls is detailed in Table 3.

Project	Cost (\$)
Atomic Fact Decomposition	21
Conflicting Fact Masking	16
Structured Prompt	74
Reflection Instruction	126
Atomic Fact Verification	173
Atomic Opinion Verification	158
Fine-Grained Error Detection	89
Feedback Judgment	97
Total	754

Table 3: API call costs for the KGDS benchmark.

Acknowledgement

This research was supported in part by the National Natural Science Foundation of China (Grant Nos. 62276017, 62406033, U1636211, 61672081) and the State Key Laboratory of Complex & Critical Software Environment (Grant No. SKLCCSE-2024ZX-18). We also thank the anonymous reviewers for their constructive feedback that helped improve this work.

References

- Forrest Sheng Bao, Miaoran Li, Renyi Qu, Ge Luo, Erana Wan, Yujia Tang, Weisi Fan, Manveer Singh Tamber, Suleman Kazi, Vivek Sourabh, and 1 others. 2024. Faithbench: A diverse hallucination benchmark for summarization by modern llms. *arXiv preprint arXiv:2410.13210*.
- Catarina G Belem, Pouya Pezeskhpour, Hayate Iso, Seiji Maekawa, Nikita Bhutani, and Estevam Hruschka. 2024. From single to multi: How llms hallucinate in multi-document summarization. *arXiv preprint arXiv:2410.13961*.
- Jiaao Chen and Diyi Yang. 2021. [Structure-aware abstractive conversation summarization via discourse and action graphs](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Asso-*

- ciation for Computational Linguistics: Human Language Technologies*, pages 1380–1391, Online. Association for Computational Linguistics.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Yi Chen, Rui Wang, Haiyun Jiang, Shuming Shi, and Ruifeng Xu. 2023. [Exploring the use of large language models for reference-free text quality evaluation: An empirical study](#). In *Findings of the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pages 361–374, Nusa Dua, Bali. Association for Computational Linguistics.
- Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. 2021. [DialogSum: A real-life scenario dialogue summarization dataset](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5062–5074, Online. Association for Computational Linguistics.
- Anshuman Chhabra, Hadi Askari, and Prasant Mohapatra. 2024. [Revisiting zero-shot abstractive summarization in the era of large language models from the perspective of position bias](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 1–11, Mexico City, Mexico. Association for Computational Linguistics.
- Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gestein, and Arman Cohan. 2023. Investigating data contamination in modern benchmarks for large language models. *arXiv preprint arXiv:2311.09783*.
- Alexander Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. [QAFactEval: Improved QA-based factual consistency evaluation for summarization](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2587–2601, Seattle, United States. Association for Computational Linguistics.
- Yue Fang, Hainan Zhang, Hongshen Chen, Zhuoye Ding, Bo Long, Yanyan Lan, and Yanquan Zhou. 2022. [From spoken dialogue to formal summary: An utterance rewriting for dialogue summarization](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3859–3869, Seattle, United States. Association for Computational Linguistics.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2024. [GPTScore: Evaluate as you desire](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6556–6576, Mexico City, Mexico. Association for Computational Linguistics.
- Mingqi Gao and Xiaojun Wan. 2022. [DialSummEval: Revisiting summarization evaluation for dialogues](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5693–5709, Seattle, United States. Association for Computational Linguistics.
- Mingqi Gao, Xiaojun Wan, Jia Su, Zhefeng Wang, and Baoxing Huai. 2023. [Reference matters: Benchmarking factual error correction for dialogue summarization with fine-grained evaluation framework](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13932–13959, Toronto, Canada. Association for Computational Linguistics.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. [SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization](#). In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics.
- Tanya Goyal and Greg Durrett. 2020. [Evaluating factuality in generation with dependency-level entailment](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3592–3603, Online. Association for Computational Linguistics.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *arXiv preprint arXiv:2209.12356*.
- Yebowen Hu, Timothy Ganter, Hanieh Deilamsalehy, Franck Dernoncourt, Hassan Foroosh, and Fei Liu. 2023. [MeetingBank: A benchmark dataset for meeting summarization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16409–16423, Toronto, Canada. Association for Computational Linguistics.
- Kung-Hsiang Huang, Philippe Laban, Alexander Fabbri, Prafulla Kumar Choubey, Shafiq Joty, Caiming Xiong, and Chien-Sheng Wu. 2024. [Embrace divergence for richer insights: A multi-document summarization benchmark and a case study on summarizing diverse information from news articles](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 570–593, Mexico City, Mexico. Association for Computational Linguistics.
- Qi Jia, Yizhu Liu, Siyu Ren, and Kenny Q Zhu. 2023. Taxonomy of abstractive dialogue summarization: Scenarios, approaches, and future directions. *ACM Computing Surveys*, 56(3):1–38.
- Frederic Kirstein, Jan Philip Wahle, Bela Gipp, and Terry Ruas. 2024. Cads: A systematic literature

- review on the challenges of abstractive dialogue summarization. *arXiv preprint arXiv:2406.07494*.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. Summac: Re-visiting nli-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Yuho Lee, Taewon Yun, Jason Cai, Hang Su, and Hwanjun Song. 2024. [UniSumEval: Towards unified, fine-grained, multi-dimensional summarization evaluation for LLMs](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3941–3960, Miami, Florida, USA. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.
- Changmao Li and Jeffrey Flanigan. 2024. [Task contamination: Language models may not be few-shot anymore](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):18471–18480.
- Jia Li, Ge Li, Yongmin Li, and Zhi Jin. 2023. [Structured chain-of-thought prompting for code generation](#). *Preprint*, arXiv:2305.06599.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Haitao Lin, Liqun Ma, Junnan Zhu, Lu Xiang, Yu Zhou, Jiajun Zhang, and Chengqing Zong. 2021. [CSDS: A fine-grained Chinese dataset for customer service dialogue summarization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4436–4451, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Haitao Lin, Junnan Zhu, Lu Xiang, Yu Zhou, Jiajun Zhang, and Chengqing Zong. 2022. [Other roles matter! enhancing role-oriented dialogue summarization via role interactions](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2545–2558, Dublin, Ireland. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023a. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.
- Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023b. [Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4140–4170, Toronto, Canada. Association for Computational Linguistics.
- Yixin Liu, Alexander Fabbri, Jiawen Chen, Yilun Zhao, Simeng Han, Shafiq Joty, Pengfei Liu, Dragomir Radev, Chien-Sheng Wu, and Arman Cohan. 2024a. [Benchmarking generation and evaluation capabilities of large language models for instruction controllable summarization](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 4481–4501, Mexico City, Mexico. Association for Computational Linguistics.
- Yong Liu, Shenggen Ju, and Junfeng Wang. 2024b. Exploring the potential of chatgpt in medical dialogue summarization: a study on consistency with human preferences. *BMC Medical Informatics and Decision Making*, 24(1):75.
- Yen-Ju Lu, Ting-Yao Hu, Hema Swetha Koppula, Hadi Pouransari, Jen-Hao Rick Chang, Yin Xia, Xiang Kong, Qi Zhu, Simon Wang, Oncel Tuzel, and 1 others. 2025. Mutual reinforcement of llm dialogue synthesis and summarization capabilities for few-shot dialogue summarization. *arXiv preprint arXiv:2502.17328*.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [FActScore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.
- Sanjana Ramprasad, Elisa Ferracane, and Zachary C Lipton. 2024. Analyzing llm behavior in dialogue summarization: Unveiling circumstantial hallucination trends. *arXiv preprint arXiv:2406.03487*.
- Mathieu Ravaut, Aixin Sun, Nancy Chen, and Shafiq Joty. 2024. [On context utilization in summarization with large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2764–2781, Bangkok, Thailand. Association for Computational Linguistics.
- Virgile Rennard, Guokan Shang, Julie Hunter, and Michalis Vazirgiannis. 2023. Abstractive meeting summarization: A survey. *Transactions of the Association for Computational Linguistics*, 11:861–884.
- Alessandro Scirè, Karim Ghonim, and Roberto Navigli. 2024. [FENICE: Factuality evaluation of summarization based on natural language inference and claim extraction](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 14148–14161,

- Bangkok, Thailand. Association for Computational Linguistics.
- Chenhui Shen, Liying Cheng, Xuan-Phi Nguyen, Yang You, and Lidong Bing. 2023. [Large language models are not yet human-level evaluators for abstractive summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4215–4233, Singapore. Association for Computational Linguistics.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. [Reflexion: language agents with verbal reinforcement learning](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 8634–8652. Curran Associates, Inc.
- Hwanjun Song, Hang Su, Igor Shalyminov, Jason Cai, and Saab Mansour. 2024a. [FineSurE: Fine-grained summarization evaluation using LLMs](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 906–922, Bangkok, Thailand. Association for Computational Linguistics.
- Hwanjun Song, Taewon Yun, Yuho Lee, Jihwan Oh, Gihun Lee, Jason Cai, and Hang Su. 2025. [Learning to summarize from llm-generated feedback](#). *Preprint*, arXiv:2410.13116.
- Yixiao Song, Yekyung Kim, and Mohit Iyyer. 2024b. [VeriScore: Evaluating the factuality of verifiable claims in long-form text generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9447–9474, Miami, Florida, USA. Association for Computational Linguistics.
- Liyan Tang, Philippe Laban, and Greg Durrett. 2024a. [MiniCheck: Efficient fact-checking of LLMs on grounding documents](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8818–8847, Miami, Florida, USA. Association for Computational Linguistics.
- Liyan Tang, Igor Shalyminov, Amy Wing-mei Wong, Jon Burnsky, Jake W Vincent, Yu’an Yang, Siffi Singh, Song Feng, Hwanjun Song, Hang Su, and 1 others. 2024b. [Tofueval: Evaluating hallucinations of llms on topic-focused dialogue summarization](#). *arXiv preprint arXiv:2402.13249*.
- Yuting Tang, Ratish Puduppully, Zhengyuan Liu, and Nancy Chen. 2023. [In-context learning of large language models for controlled dialogue summarization: A holistic benchmark and empirical analysis](#). In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 56–67, Singapore. Association for Computational Linguistics.
- Yuanhe Tian, Fei Xia, and Yan Song. 2024. [Dialogue summarization with mixture of experts based on large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7143–7155, Bangkok, Thailand. Association for Computational Linguistics.
- David Wan, Koustuv Sinha, Srini Iyer, Asli Celikyilmaz, Mohit Bansal, and Ramakanth Pasunuru. 2024. [ACUEval: Fine-grained hallucination evaluation and correction for abstractive summarization](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 10036–10056, Bangkok, Thailand. Association for Computational Linguistics.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. [Asking and answering questions to evaluate the factual consistency of summaries](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5008–5020, Online. Association for Computational Linguistics.
- Bin Wang, Zhengyuan Liu, and Nancy Chen. 2023a. [Instructive dialogue summarization with query aggregations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7630–7653, Singapore. Association for Computational Linguistics.
- Bin Wang, Chen Zhang, Yan Zhang, Yiming Chen, and Haizhou Li. 2022. [Analyzing and evaluating faithfulness in dialogue summarization](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4897–4908, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Hongru Wang, Minda Hu, Yang Deng, Rui Wang, Fei Mi, Weichao Wang, Yasheng Wang, Wai-Chung Kwan, Irwin King, and Kam-Fai Wong. 2023b. [Large language models as source planner for personalized knowledge-grounded dialogues](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9556–9569, Singapore. Association for Computational Linguistics.
- Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023c. [Is ChatGPT a good NLG evaluator? a preliminary study](#). In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11, Singapore. Association for Computational Linguistics.
- Jiaan Wang, Jianfeng Qu, Kexin Wang, Zhixu Li, Wen Hua, Ximing Li, and An Liu. 2024. [Improving the robustness of knowledge-grounded dialogue via contrastive learning](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19135–19143.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Jie Huang, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, Cosmo Du, and Quoc V. Le.

2024. [Long-form factuality in large language models](#). *Preprint*, arXiv:2403.18802.
- Joonho Yang, Seunghyun Yoon, ByeongJeong Kim, and Hwanhee Lee. 2024. [FIZZ: Factual inconsistency detection by zoom-in summary and zoom-out document](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 30–45, Miami, Florida, USA. Association for Computational Linguistics.
- Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen, and Wei Cheng. 2023. Exploring the limits of chatgpt for query or aspect-based text summarization. *arXiv preprint arXiv:2302.08081*.
- Yangfan Ye, Xiachong Feng, Xiaocheng Feng, Weitao Ma, Libo Qin, Dongliang Xu, Qing Yang, Hongtao Liu, and Bing Qin. 2024. [GlobeSumm: A challenging benchmark towards unifying multi-lingual, cross-lingual and multi-document news summarization](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10803–10821, Miami, Florida, USA. Association for Computational Linguistics.
- Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023. [AlignScore: Evaluating factual consistency with a unified alignment function](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.
- Houyu Zhang, Zhenghao Liu, Chenyan Xiong, and Zhiyuan Liu. 2020. [Grounded conversation generation as guided traverses in commonsense knowledge graphs](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2031–2043, Online. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics*, 12:39–57.
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. [MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 563–578, Hong Kong, China. Association for Computational Linguistics.
- Ming Zhong, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2022a. Dialoglm: Pre-trained model for long dialogue understanding and summarization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11765–11773.
- Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022b. [Towards a unified multi-dimensional evaluator for text generation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2023–2038, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Hao Zhou, Tom Young, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Commonsense knowledge aware conversation generation with graph attention. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI’18*, page 4623–4629. AAAI Press.
- Weixiao Zhou, Gengyao Li, Xianfu Cheng, Xinnian Liang, Junnan Zhu, Feifei Zhai, and Zhoujun Li. 2023. [Multi-stage pre-training enhanced by ChatGPT for multi-scenario multi-domain dialogue summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6893–6908, Singapore. Association for Computational Linguistics.
- Chenguang Zhu, Yang Liu, Jie Mei, and Michael Zeng. 2021. [MediaSum: A large-scale media interview dataset for dialogue summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5927–5934, Online. Association for Computational Linguistics.
- Rongxin Zhu, Jey Han Lau, and Jianzhong Qi. 2025. [Factual dialogue summarization via learning from large language models](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 4474–4492, Abu Dhabi, UAE. Association for Computational Linguistics.
- Rongxin Zhu, Jianzhong Qi, and Jey Han Lau. 2023. [Annotating and detecting fine-grained factual errors for dialogue summarization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6825–6845, Toronto, Canada. Association for Computational Linguistics.
- Yicheng Zou, Bolin Zhu, Xingwu Hu, Tao Gui, and Qi Zhang. 2021. [Low-resource dialogue summarization with domain-agnostic multi-source pretraining](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 80–91, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

A Overall Performance Visualization

Figures 6 and 8 show the overall performance of EBS-AOS and ABS-AOS. Figures 7 and 9 illustrate the performance gap and average performance.

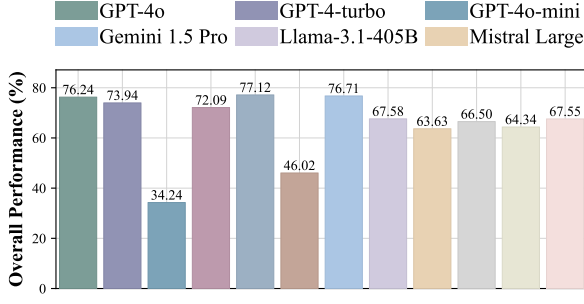


Figure 6: Overall performance % of all LLMs in the EBS-AOS pattern. The stratification is divided into three levels: TIER-1 ($OP_{GM} \in [72, 77]$), TIER-2 ($OP_{GM} \in [64, 68]$), and TIER-3 ($OP_{GM} \in [34, 46]$).

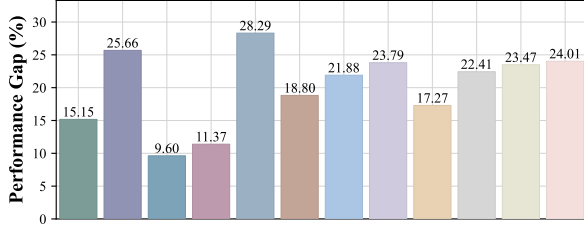


Figure 7: Overall performance gap % of LLMs between the two KGDS patterns. A larger gap indicates lower stability in cross-pattern performance.

B KGDS BenchMark Statistics

News Statistics. Refer to Table 4.

Num.	Paras _{avg}	Tokens _{avg}
100	14.4	617.5

Table 4: Statistics for news. Num. indicates the number of news articles. Paras_{avg} and Tokens_{avg} represent the average number of paragraphs and tokens, respectively.

Discussion Statistics. Refer to Table 5.

Num.	Pts. _{avg}	Uttrs. _{avg}	Tokens _{avg}
100	2.0	4.1	112.0

Table 5: Statistics for discussion. Num. indicates the number of discussions. Pts._{avg}, Uttrs._{avg}, and Tokens_{avg} represent the average number of participants, utterances, and tokens, respectively.

Annotation Statistics for EBS. Among the 1696 paragraphs from 100 original news articles, the expert annotation consistency rate is 84.7%. A total of 1437 paragraphs are retained for the final news articles, while 208 are identified as boundary paragraphs and removed. Of the 1437 retained, 432 are annotated as back-supporting and 1005 as back-nonsupporting.

Annotation Statistics for ABS. Among the 2428 atomic fact units decomposed from the 432 back-

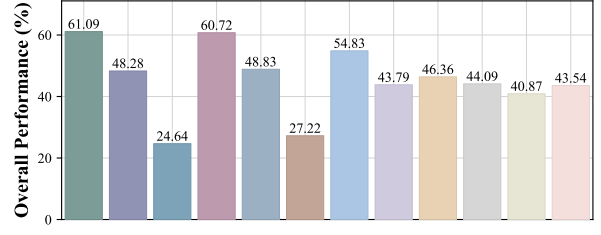


Figure 8: Overall performance % of all LLMs in the ABS-AOS pattern. The stratification is divided into three levels: TIER-1 ($OP_{GM} \in [55, 61]$), TIER-2 ($OP_{GM} \in [41, 49]$), and TIER-3 ($OP_{GM} \in [25, 27]$).

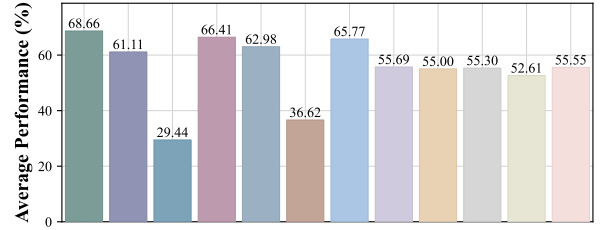


Figure 9: Average overall performance % of all LLMs across both KGDS patterns. The stratification is divided into three levels: TIER-1 ($OP_{avg} \in [61, 69]$), TIER-2 ($OP_{avg} \in [53, 56]$), and TIER-3 ($OP_{avg} \in [29, 37]$).

supporting paragraphs, 1638 are consistently annotated as key back-supporting atomic facts, 780 as non-key back-supporting atomic facts, and 10 as intra-repetitive atomic facts.

Among the 5187 atomic fact units from the 1005 back-nonsupporting paragraphs, 4996 are automatically annotated as back-nonsupporting atomic facts, 176 as masked conflicting atomic facts, and 15 as intra-repetitive atomic facts.

Annotation Statistics for AOS. A total of 873 clear atomic opinions are expert manually annotated from 100 discussions and their corresponding new articles. Of these, 800 contain clarified implicit references, while 73 do not require clarification. Within the 800 opinions, there are a total of 1113 written clarified implicit references.

C Detailed EBS Evaluation Metrics

Let $\mathcal{B} = \{b_i\}_{i=1}^n$, $\mathcal{S} = \{s_j\}_{j=1}^m$, and $\mathcal{N} = \{n_k\}_{k=1}^p$ denote the sets of paragraph indices for LLM-extracted EBS, back-supporting, and back-nonsupporting paragraphs, respectively. The automatic paragraph index matching function is:

$$\xi : (\mathcal{B}, p) \rightarrow \{0, 1\}, p \in \mathcal{S} \cup \mathcal{N} \quad (11)$$

where 1 and 0 represent whether the paragraph index p is present in \mathcal{B} .

```

### Abstractive Opinion Summary:
{summary content}

### Question:
**After carefully reading and deeply understanding the "Abstractive Opinion Summary" above, can you know
the following "Atomic Opinion"?**

### Atomic Opinion:
<Opinion>: ...

### Inference Principle:
**You may only use information from the "Abstractive Opinion Summary" to infer the knowability of the
"Atomic Opinion".**

### Result Return Format:
* Provide your verification result in JSON format.
* The returned JSON includes two keys: "Inference_Conclusion" and "Analysis_Reasoning".
* The value of "Inference_Conclusion" is a string that provides a binary conclusion: "knowable" or
"unknowable".
* The value of "Analysis_Reasoning" is a string that provides a detailed explanation of why the
"Inference_Conclusion" is "knowable" or "unknowable".

### JSON Format Example:
```json
{
 "Inference_Conclusion": "knowable or unknowable",
 "Analysis_Reasoning": "..."
}
```

```

Figure 10: Atomic opinion verification prompt.

The **BSP Recall**, **Precision**, and **F1-score** are defined as follows:

$$\text{BSP}_R = \frac{1}{m} \sum_{j=1}^m \xi(\mathcal{B}, s_j) \quad (12)$$

$$\text{BSP}_P = \frac{\sum_{s \in \mathcal{S}} \xi(\mathcal{B}, s)}{\sum_{p \in \mathcal{S} \cup \mathcal{N}} \xi(\mathcal{B}, p)} \quad (13)$$

$$\text{BSP}_{F_1} = 2 \cdot \frac{\text{BSP}_P \cdot \text{BSP}_R}{\text{BSP}_P + \text{BSP}_R} \quad (14)$$

D Prompts and Instructions

D.1 Unified Parameter Settings

In this work, All 12 evaluated LLMs use consistent parameter settings for prompts and instructions: max_token=4096 and temperature=0. No other default parameters are modified.

D.2 Atomic Fact Decomposition Instruction

Refer to Figure 11.

D.3 Structured Prompts

See Figures 12 and 13 for EBS-AOS and ABS-AOS patterns, respectively.

D.4 Self-Reflection Instructions

See Figures 15 and 14 for EBS-AOS and ABS-AOS patterns, respectively.

D.5 Feedback Judgment Instruction

Refer to Figure 16.

D.6 Verification Prompts

See Figures 17 and 10 for fact and opinion verification, respectively.

D.7 Error Detection Instruction

Refer to Figure 18.

Please carefully read, deeply understand, and strictly follow the instructions below to decompose the "paragraph" into "fine-grained atomic fact units":

Paragraph:

<Paragraph_x>: {content}

Principles of Fine-Grained Atomic Fact Units:

- (1). Indivisibility: Ensure that each "fine-grained atomic fact unit" is minimal and cannot be further decomposed into smaller "fine-grained atomic fact units."
- (2). Independence: Ensure that each "fine-grained atomic fact unit" can be understood independently, without relying on other "fine-grained atomic fact units."
- (3). Declarativity: Ensure that each "fine-grained atomic fact unit" is a concise declarative sentence that clearly conveys a single basic fact.

Output Return Format:

<Paragraph_x>:

<Fact_1>: ...

...

<Fact_n>: ...

Figure 11: Atomic fact decomposition instruction.

Shared Background Knowledge (SBK):

<Paragraph_1>: ...

...

<Paragraph_n>: ...

Knowledge-Grounded Discussion (KGD):

Person1: ...

Person2: ...

...

The above provides SBK and KGD, where SBK represents the shared background knowledge that participants are familiar with prior to the discussion, KGD denotes the discussion by the participants grounded in either partial or complete content of SBK.

Task:

****Please combine SBK and KGD to summarize, the result includes two summaries:****

*** **Extractive Background Summary**:** The definition of this summary is the **extractive background-supporting paragraphs for KGD from SBK**.

*** **Abstractive Opinion Summary**:** The definition of this summary is the **clear personal opinions of the participants with clarified implicit references**. Here, implicit references represent the utilize of pronouns or phrases in KGD to refer to entities, facts, events, or other types of sub-information within SBK.

JSON Result Return Format:

``json

```
{
  "Extractive_Background_Summary": [
    "<Paragraph_#>",
    ...,
    "<Paragraph_#>"
  ],
  "Abstractive_Opinion_Summary": "Person1... Person2... Person1... Person2..."
}
```

Figure 12: Structured prompt for EBS-AOS pattern.

Shared Background Knowledge (SBK):

<Paragraph_1>: ...

...

<Paragraph_n>: ...

Knowledge-Grounded Discussion (KGD):

Person1: ...

Person2: ...

...

The above provides SBK and KGD, where SBK represents the shared background knowledge that participants are familiar with prior to the discussion, KGD denotes the discussion by the participants grounded in either partial or complete content of SBK.

Task:

****Please combine SBK and KGD to summarize, the result includes two summaries:****

*** **Abstractive Background Summary**:** The definition of this summary is the **abstractive background-supporting information for KGD from SBK**.

*** **Abstractive Opinion Summary**:** The definition of this summary is the **clear personal opinions of the participants with clarified implicit references**. Here, implicit references represent the utilize of pronouns or phrases in KGD to refer to entities, facts, events, or other types of sub-information within SBK.

JSON Result Return Format:

```
```json
```

```
{
```

```
 "Abstractive_Background_Summary": "...",
```

```
 "Abstractive_Opinion_Summary": "Person1... Person2... Person1... Person2..."
```

```
}
```

```
```
```

Figure 13: Structured prompt for ABS-AOS pattern.

Instruction:

*** **Please self-reflect: Do the Abstractive_Background_Summary and Abstractive_Opinion_Summary you provided align with their respective definitions?****

*** **Please think step by step, provide the two summaries again after self-reflection. Additionally, you need to provide a detailed chain-of-thought for self-reflection.****

JSON Result Return Format:

```
```json
```

```
{
```

```
 "Abstractive_Background_Summary": "...",
```

```
 "Abstractive_Opinion_Summary": "Person1... Person2... Person1... Person2...",
```

```
 "Chain-of-Thought_for_Self-Reflection": "..."
```

```
}
```

```
```
```

Figure 14: Self-reflection instruction for ABS-AOS pattern.

```

### Instruction:
* **Please self-reflect: Do the Extractive_Background_Summary and Abstractive_Opinion_Summary you
provided align with their respective definitions?**
* **Please think step by step, provide the two summaries again after self-reflection. Additionally, you need to
provide a detailed chain-of-thought for self-reflection.**

### JSON Result Return Format:
```json
{
 "Extractive_Background_Summary": [
 "<Paragraph_#>",
 ...
 "<Paragraph_#>"
],
 "Abstractive_Opinion_Summary": "Person1... Person2... Person1... Person2...",
 "Chain-of-Thought_for_Self-Reflection": "..."
}
...

```

Figure 15: Self-reflection instruction for EBS-AOS pattern.

```

Instruction:
**Please conduct feedback judgment: Which of the Extractive_Background_Summary you provided above
(EBS_A) and the expert annotated (EBS_B) is more aligned with its definition?**

You Provided Extractive_Background_Summary (EBS_A):
{summary_a}

Expert Annotated Extractive_Background_Summary (EBS_B):
{summary_b}

**Please think step by step, provide the judgment conclusion (EBS_A is more aligned / EBS_B is more aligned)
after feedback judgment. Additionally, you need to provide a detailed chain-of-thought for feedback
judgment.**

JSON Result Return Format:
```json
{
  "Judgment_Conclusion": "EBS_A is more aligned or EBS_B is more aligned",
  "Chain-of-Thought_for_Feedback-Judgment": "..."
}
...

```

Figure 16: EBS feedback judgment instruction.

```

### Abstractive Background Summary:
{summary content}

### Question:
**After carefully reading and deeply understanding the "Abstractive Background Summary" above, can you
know each of the following "Atomic Facts"?**

### Atomic Facts:
<Fact_1>: ...
...
<Fact_n>: ...

### Inference Principle:
**You may only use information from the "Abstractive Background Summary" to infer the knowability of each
of the "Atomic Facts".**

### Result Return Format:
* Provide your verification results in JSON format.
* The returned JSON is a list that provides the verification result of each atomic fact in order.
* Each verification result includes two keys: "Fact_Index" and "Inference_Conclusion".
* The value of "Fact_Index" is a string that provides the index label of the atomic fact.
* The value of "Inference_Conclusion" is a string that provides a binary conclusion: "knowable" or
"unknowable".

### JSON Format Example:
```json
[
 {
 "Fact_Index": "<Fact_1>",
 "Inference_Conclusion": "knowable or unknowable"
 },
 ...
 {
 "Fact_Index": "<Fact_n>",
 "Inference_Conclusion": "knowable or unknowable"
 }
]
```

```

Figure 17: Atomic fact verification prompt. The number of facts (n) dynamically changes at the paragraph-level.

Instruction:

****Please perform [fine-grained opinion error detection](#): For the "Atomic Opinion" that you have verified as "unknowable" above, conduct error classification on the reasons why it cannot be inferred from the "Abstractive Opinion Summary". For an "Atomic Opinion," you only need to provide one most matching error type from the following five error types:****

Error Types:

*** **[Error Type1: Opinion Misattribution](#)****

****Definition of Error Type1: The "Abstractive Opinion Summary" incorrectly attributes the opinion of one participant (Person1 or Person2) in the "Atomic Opinion" to another participant, or mistakenly labels it as a group opinion, making it impossible to establish the correct correspondence between the participant and the opinion, and thus impossible to infer the "Atomic Opinion" from the "Abstractive Opinion Summary".****

*** **[Error Type2: Implicit Reference Incorrectly Clarified](#)****

****Definition of Error Type2: When the "Atomic Opinion" contains clarified implicit references (the content highlighted between a pair of double asterisks [format: `**xxx**` or `**xxx**`] in the "Atomic Opinion" is the clarified implicit reference), if the "Abstractive Opinion Summary" contains corresponding incorrect clarifications or explicit explanations of the implicit references, making it impossible to infer the "Atomic Opinion" from the "Abstractive Opinion Summary".**

*** **[Error Type3: Implicit Reference Unclassified](#)****

****Definition of Error Type3: When the "Atomic Opinion" contains clarified implicit references (the content highlighted between a pair of double asterisks [format: `**xxx**` or `**xxx**`] in the "Atomic Opinion" is the clarified implicit reference), if the "Abstractive Opinion Summary" contains corresponding implicit references (such as pronouns/reference phrases/eclipses, etc.) but fails to clarify or explicitly explain them, making it impossible to infer the "Atomic Opinion" from the "Abstractive Opinion Summary".**

*** **[Error Type4: Opinion Sentiment Distortion](#)****

****Definition of Error Type4: The "Abstractive Opinion Summary" is inconsistent with the subjective sentiment expressed by the participant in the "Atomic Opinion", making it impossible to infer the "Atomic Opinion" from the "Abstractive Opinion Summary".****

*** **[Error Type5: Opinion Fact Inconsistency](#)****

****Definition of Error Type5: The "Abstractive Opinion Summary" is inconsistent with the objective facts expressed by the participant in the "Atomic Opinion", making it impossible to infer the "Atomic Opinion" from the "Abstractive Opinion Summary".****

****Please [think step by step](#), provide the detection conclusion after error classification. Additionally, you need to provide a detailed [chain-of-thought for error detection](#).****

JSON Result Return Format:

```
```json
{
 "Detection_Conclusion": "Error Type1 or Error Type2 or Error Type3 or Error Type4 or Error Type5",
 "Chain-of-Thought_for_Error-Detection": "..."
}
```
```

Figure 18: AOS error detection instruction.

E LLM Sources

OpenAI: GPT-4o, GPT-4-turbo, GPT-4o-mini

<https://platform.openai.com/docs/api-reference/introduction>

Anthropic: Claude 3 Opus, Claude 3.5 Sonnet, Claude 3.5 Haiku

<https://docs.anthropic.com/en/api/getting-started>

Google: Gemini 1.5 Pro

<https://ai.google.dev/gemini-api/docs>

Meta: Llama-3.1-405B

<https://www.llmapi.com>

Mistral AI: Mistral Large

<https://docs.mistral.ai/api/>

DeepSeek: DeepSeek-V3

<https://api-docs.deepseek.com/zh-cn/>

Alibaba: Qwen-Max

<https://bailian.console.aliyun.com/>

ZhipuAI: GLM-4-Plus

<https://www.bigmodel.cn/dev/api/normal-model/glm-4>