

Not All Contexts Are Equal: Teaching LLMs Credibility-aware Generation

Anonymous ACL submission

Abstract

The rapid development of large language models has led to the widespread adoption of Retrieval-Augmented Generation (RAG), which integrates external knowledge to alleviate knowledge bottlenecks and mitigate hallucinations. However, the existing RAG paradigm inevitably suffers from the impact of *flawed information* introduced during the retrieval phrase, thereby diminishing the reliability and correctness of the generated outcomes. In this paper, we propose Credibility-aware Generation (CAG), a universally applicable framework designed to mitigate the impact of flawed information in RAG. At its core, CAG aims to equip models with the ability to discern and process information based on its credibility. To this end, we propose an innovative data transformation framework that generates data based on credibility, thereby effectively endowing models with the capability of CAG. Furthermore, to accurately evaluate the models' capabilities of CAG, we construct a comprehensive benchmark covering three critical real-world scenarios. Experimental results demonstrate that our model can effectively understand and employ credibility for generation, significantly outperform other models with retrieval augmentation, and exhibit robustness despite the increasing noise in the context.

1 Introduction

In recent years, Large Language Models (LLMs) (Brown et al., 2020; OpenAI et al., 2023; Touvron et al., 2023; Anil et al., 2023) have experienced significant growth and demonstrated excellent performance in multiple domains (Kojima et al., 2022; Thirunavukarasu et al., 2023; Ziems et al., 2023; Min et al., 2023). With the ascendancy of LLMs, Retrieval-Augmented Generation (RAG) has attracted significant interest. RAG mitigates the knowledge bottleneck of LLMs by incorporating externally retrieved documents into their generation process. This inclusion helps diminish the

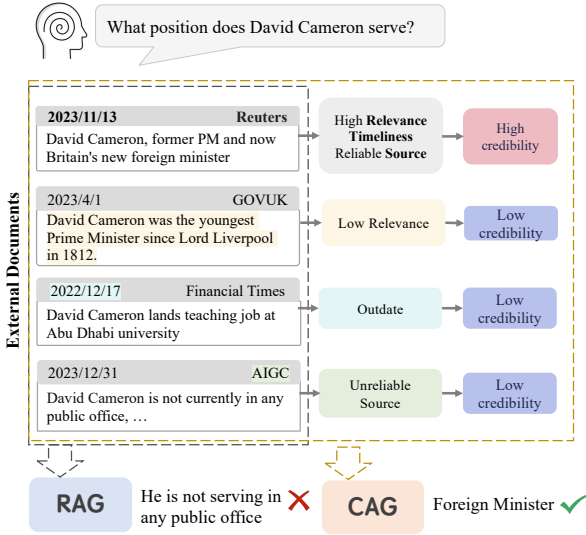


Figure 1: The comparison between Retrieval-Augmented Generation (RAG) and Credibility-aware Generation (CAG). Incorporating credibility into the model aids in mitigating errors caused by *flawed information* introduced from the retrieval process.

occurrences of hallucinations and misinformation during generation, thereby substantially enhancing the quality of output from LLMs (Petroni et al., 2021; Zhu et al., 2021; Mallen et al., 2023).

However, RAG for large language models remains significantly impacted by flawed information. This is mainly because the retrieval process often provides noisy, outdated, and incorrect contexts which adversely affects RAG, substantially reducing its effectiveness. Specifically, previous research (Shi et al., 2023a; Chen et al., 2023) has found that LLMs are highly sensitive to noise, which impacts LLMs' capacity to discern and trust accurate information, ultimately affecting the outcomes they generate. Furthermore, due to the temporal insensitivity of LLMs (Su et al., 2022; Zhao et al., 2024), these models struggle to discern outdated information solely based on their internal knowledge. More critically, because LLMs are

062 trained on extensive collections of historical text, 113
063 there’s an inherent risk that outdated information 114
064 will align with the models’ internal knowledge 115
065 bases. This alignment can encourage LLMs to fa- 116
066 vor and perpetuate outdated information. Besides, 117
067 the prevalence of misinformation on the current 118
068 web poses a significant challenge for large models, 119
069 which struggle to identify misinformation using 120
070 only their inherent knowledge (Xie et al., 2023; Pan 121
071 et al., 2023). This difficulty makes them suscepti- 122
072 ble to misinformation, leading to the generation of 123
073 incorrect answers. Therefore, flawed information, 124
074 characterized by noisy, outdated, and incorrect in- 125
075 formation, has substantial negative effects on RAG. 126

076 From a cognition perspective, a common ap- 127
077 proach humans adopt to combat flawed informa- 128
078 tion is to assess the credibility of external infor- 129
079 mation (Burgoon et al., 2000). Specifically, in- 130
080 formation that is current, evaluated, and sourced 131
081 from highly credible origins is typically regarded 132
082 as more timely, accurate, and reliable. Motivated 133
083 by this, we introduce Credibility-aware Genera- 134
084 tion (CAG), a universally applicable framework 135
085 designed to address flawed information encoun- 136
086 tered during RAG. At its core, CAG seeks to equip 137
087 models with the ability to discern and process in-
088 formation based on credibility. By assigning vary-
089 ing degrees of credibility to information based on
090 its relevance, timeliness, and the reliability of its
091 source, and explicitly distinguishing them in the in-
092 put, CAG significantly mitigates the issues arising
093 from flawed information.

094 Unfortunately, we have discovered that existing
095 LLMs are not inherently sensitive to directly pro-
096 vided credibility information in the prompt. This
097 deficiency restricts their capacity to optimally em-
098 ploy credibility for discerning and processing in-
099 formation. To endow models with the capabil-
100 ity of CAG, we propose a novel data transfor-
101 mation framework. This framework transforms
102 existing Question Answering (QA) datasets into
103 data that integrates credibility, which can be em-
104 ployed to guide the model for credibility-based
105 generation. Specifically, our process comprises
106 two core steps: 1) Multi-granularity credibility an-
107 notation, which assigns credibility to text units at
108 both document and sentence levels by dividing re-
109 trieved documents into varying granularities. 2)
110 Credibility-guided explanation generation, which
111 prompts LLMs to generate credibility-guided ex-
112 planations given questions, retrieved documents

with credibility annotation and golden answers. Fi-
nally, we employ instruction fine-tuning to train the
model, enabling it to generate responses based on
credibility.

To rigorously assess the ability of the model’s
Credibility-aware generation in managing flawed
information, we construct a comprehensive bench-
mark encompassing various real-world scenarios,
including open-domain QA, time-sensitive QA,
and misinformation polluted QA. In this bench-
mark, retrieval relevance, timeliness, and source
authority are regarded as established measures
of credibility. Experimental results on multiple
datasets across multiple scenarios demonstrate the
efficacy of our approach in using credibility. Our
model significantly outperforms various prevalent
RAG approaches applied to both open and closed-
source LLMs of diverse scales. Additionally, it
exhibits robust resilience against noisy documents,
maintaining high performance even as alternative
strategies suffer sharp declines. All these results
verify the effectiveness of the proposed CAG frame-
work and corresponding training algorithm.

The main contributions of this study are summa-
rized as follows ¹:

- We present Credibility-aware Generation, a 138
universal framework to handle the flawed in- 139
formation challenge in RAG. 140
- We propose a novel data transformation frame- 141
work that transforms existing datasets into 142
data annotated with credibility and guides 143
models to generate responses based on cred- 144
ibility, thereby equipping the model with 145
Credibility-aware Generation capability. 146
- We construct a comprehensive benchmark and 147
evaluate model performance in credibility- 148
aware generation, encompassing real-world 149
scenarios of open-domain QA, time-sensitive 150
QA, and misinformation polluted QA. 151
- Experimental evidences demonstrate that our 152
model effectively understands and employs 153
credibility to generate responses, significantly 154
surpasses other RAG-based strategies, and 155
maintains robustness despite the increasing 156
noise in the context. 157

¹We uploaded the code and datasets as supplemental mate-
rials, which will be openly released after accepting.

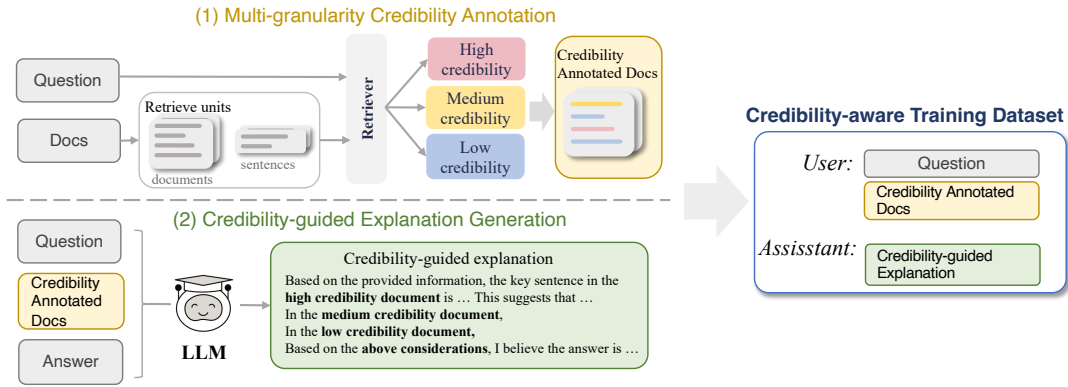


Figure 2: Overview of data transformation framework. The training data is constructed by assigning credibility to contexts via multi-granularity credibility annotation (§3.1) and prompting LLM to produce credibility-guided explanations (§3.2). The processed data is used to instruction fine-tuning (§3.3) to endow the model with the ability for Credibility-aware Generation.

2 Credibility-aware Generation

Credibility-aware Generation is designed to enable models to discern and process information based on its credibility. Subsequently, we will provide formal definitions for both RAG and CAG, illustrating their divergence.

Definition In the Retrieval-Augmented Generation process, user input x initiates the retrieval of a set of related documents D_x from a large corpus C based on how closely these documents match the input. Then, it combines the input x with these documents D_x to generate responses y , formalized as $y = \text{LM}([x, D_x])$, where $[\cdot, \cdot]$ denotes the concatenation operation.

Compared to RAG, the Credibility-aware Generation offers additional credibility for each document. Initially, through credibility assessment based on various scenarios, each retrieved document has been assigned a level of credibility. Then, these documents D_x with their credibility C are synthesized with the user input x as augmented input. LM generates responses y based on this augmented input, formally represented as $y = \text{LM}\left(\left[x, \left\{\left[c_i, d_i\right]\right\}_{i=1}^{|D_x|}\right]\right)$. This approach ensures that the generated responses not only incorporate the content of the documents but also consider the credibility of each document, thereby enhancing the reliability of responses.

3 Teaching Model to Credibility-aware Generation

In this section, we endow LLMs with the capability of CAG. A potential approach involves directly providing the credibility annotations of each docu-

ment in the prompt. Unfortunately, as indicated in Table 2, our experiments reveal that even advanced LLMs, such as ChatGPT, exhibit limited sensitivity to credibility. To this end, we introduce a novel data transformation framework. Through multi-granularity credibility annotation and credibility-guided explanation generation, we transform existing QA datasets into data that includes credibility annotations which can guide the model to generate credibility-based responses. Then, through instruction fine-tuning, we train the model to generate responses grounded in credibility assessments.

3.1 Multi-granularity Credibility Annotation

To cater to the varied requirements for credibility across different scenarios and enhance the model’s comprehension of credibility, we collect training data including open-domain QA, machine reading comprehension, and dialogue datasets and propose a multi-granularity credibility annotation method.

First, we divide the retrieved documents to create a multi-granularity corpus, encompassing sentence and document levels. Then, the retriever assesses the match between each retrieval unit and the query, assigning a relevance score, and classifies documents into three levels: high, medium, and low, employing either equi-frequency or equi-distance segmentation. This approach of using levels instead of scores aims to simplify representation, thereby improving the model’s understanding and providing a certain degree of fault tolerance. Therefore, we collect about 15k pieces of training data, all of which include credibility annotations in the context of the QA dataset. The detailed composition of the training data is shown in the Appendix A.1.

3.2 Credibility-guided Explanation Generation

To facilitate the model’s comprehension and effective utilization of credibility, we employ LLM to generate explanations for the answers. These explanations stem from an analysis of both the content and credibility of the documents.

Given the limitations of current LLMs in comprehending credibility effectively, we design chain-of-thought prompts to guide LLM to generate credibility-guided explanations given questions, retrieved documents with credibility and golden answers. In this case, LLM is required to analyze document content and credibility, as well as the rationale for the derived answer after integrating all the information. Considering the accessibility and advanced capabilities of GPT-3.5, we employ GPT-3.5 for the generation of explanations. In this way, we obtain high-quality answer explanations. Then, we replace the original answers in the training data with credibility-guided explanations to form a novel QA dataset based on credibility. In this dataset, the inputs include questions and external documents annotated with credibility, while the outputs are credibility-guided explanations.

3.3 Instruction Fine-tuning

Through the two steps above, the training dataset obtained contains credibility, which can be used to facilitate arbitrary language models in gaining the capacity for CAG. We fine-tune the language model on this dataset to empower the model to discern and process information according to its credibility. As defined by Iyer et al. (2023), the loss function is as follows:

$$\mathcal{L}(D_{\mathbf{x}}; \theta) = - \sum_{i=1}^N \log p_{\theta} \left(\mathbf{y}_i \mid \left[\mathbf{x}, \{[c_i, d_i]\}_{i=1}^{|D_{\mathbf{x}}|} \right], \mathbf{y}_{<i} \right)$$

4 Credibility-aware Generation Benchmark

To rigorously evaluate the ability of credibility-aware model generation to handle flawed information, we construct the Credibility-aware Generation Benchmark (CAGB). This benchmark encompasses the following three specific scenarios where the integration of credibility is essential:

- **Open-domain QA** aims to accurately answer questions on a wide variety of topics without being limited to any particular area. It encompasses a broad spectrum of real-world applications that urgently require the integration of

external knowledge to enhance the language model’s ability to address queries. This scenario thus necessitates the ability to effectively identify and process noise information.

- **Time-sensitive QA** aims to give accurate and current answers. It poses a challenge for LLMs due to the dynamic internet information. The inevitable inclusion of outdated documents when incorporating external sources further complicates matters. Even with timestamps provided for documents, LLMs may erroneously prioritize outdated documents. This situation underscores the critical need for credibility in time-sensitive QA.
- **Misinformation polluted QA** aims to tackle the issue of ensuring accurate answers in an environment polluted with misinformation. It presents a substantial challenge to LLMs, attributed to the misuse of LLMs and the consequent proliferation of fake news and misinformation (Zhuo et al., 2023; Pan et al., 2023). Consequently, it is crucial to take into account the quality and credibility of any introduced external information.

In the following, we will provide a detailed description of data construction for each scenario, and the statistics of CAGB are shown in the Table 1.

4.1 Credibility Assessment

We aim to establish a flexible credibility assessment mechanism that can be conveniently extended to consider additional factors and a broader range of application fields. In this benchmark, the credibility of the documents is evaluated by considering retrieval relevance, timeliness, and source reliability. Specifically, we establish a foundation based on retrieval relevance, then make adjustments according to timeliness, and finally integrate the reliability of the source to determine credibility. First, the retriever assigns relevance scores to documents based on query similarity. These relevance scores, which are distributed at equal intervals, enable to classify documents into three levels: high, medium, and low, collectively denoted as R . Subsequently, the temporal difference T between the query time and document publication is calculated, downgrading R if T surpasses a threshold. The formula integrating relevance and timeliness is as follows:

$$rt_score(R, T) = \max(R - \text{floor}(T/\text{threshold}), 1)$$

Dataset	#Samples	#Documents	Noise Ratio
<i>Open-domain QA</i>			
HotpotQA	500	5000	0.8
2WikiMHQA	500	5000	0.6-0.8
MuSiQue	500	10000	0.9
ASQA	948	4740	-
RGB	300	11641	0.2-0.8
<i>Time-sensitive QA</i>			
EvolvTempQA	321	2247	0.4-0.8
<i>Misinformation polluted QA</i>			
NewsPollutedQA	480	2400	0.5-0.75

Table 1: Statistics of CAGB, which includes 7 dataset derived from 3 scenarios.

Following this, the reliability of the source, denoted S , is customized to specific scenarios, similarly divided into three levels. Each reflects the degree of reliability of the information source. Finally, we combine these factors, adopting the lower level as the credibility and the formula is expressed as follows:

$$Cred = \min(rt_score(R, T), S)$$

In this way, the document of high credibility are concurrently characterized by high relevance, timeliness and source reliability. More details about the assessment can be seen in the Appendix A.6.

4.2 Open-domain QA

Our research utilizes data from several challenging QA datasets that have noise in the context they provide. HotpotQA (Yang et al., 2018) and 2WikiMHQA (Ho et al., 2020) both require reasoning across multiple documents, and feature a high proportion of distracting documents. Importantly, the data we utilize from HotpotQA is extracted from the dev subset, whereas our training dataset is derived from the train subset. Musique (Trivedi et al., 2021) questions are of higher complexity, with up to 90% of distracting passages. ASQA (Stelmakh et al., 2022) is a long format QA dataset focused on ambiguous questions. RGB (Chen et al., 2023) is a specialized benchmark used for evaluating the capabilities of models in the RAG scenario, with noise robustness being one of its aspects. We assign credibility to the documents provided in the dataset in terms of retrieval relevance.

4.3 Time-sensitive QA

In order to construct a diverse, high-quality, and up-to-date news dataset, we annotate 321 time-sensitive questions along with their corresponding dates. These questions originate from real-world scenarios, including news QA data from RealTime

QA (Kasai et al., 2022), TAQA (Zhao et al., 2024), and questions adapted from news reports.

To simulate the simultaneous occurrence of varied information on the Internet, we use Google search API to gather 3 relevant documents and 4 distracting documents for each question, the latter being either irrelevant or outdated. This approach to document selection is crafted to emulate the intricate and heterogeneous nature of real-world information landscapes. Each news includes its publication date, thereby aiding in the evaluation of its timeliness. For document credibility annotation, we assess credibility based on relevance and time difference between the document’s publication and the posed question. We ensure the accuracy of the answers by manually annotating.

The obtained time-sensitive dataset with outdated document settings and credibility annotation is named EvolvingTempQA.

4.4 Misinformation Polluted QA

We create a up-to-date multiple-choice quiz dataset, comprising both real and fake news for each question. The dataset construction bases on RealTime QA, utilizing weekly news quizzes from CNN and other news platforms. To maintain the dataset’s real-time relevance, we select news from July 1, 2023, onwards, comprising 480 questions with four options and one supporting news item each.

To simulate the generation of fake news, we first generated a claim utilizing the LLM, based on a question and a randomly selected incorrect option. This process transforms the question and incorrect option into a deceptive statement. Subsequently, we choose GPT-3.5 and Qwen (Bai et al., 2023) as the generators for fake news, guiding them to generate texts of varying styles based on the claim, including news style and Twitter style. The prompts used and examples are detailed in the Appendix A.11. The fictitious news articles produced by LLMs, due to their authenticity being deliberately compromised, are classified as having low credibility. Conversely, news articles from reputable news websites are considered to possess high credibility. We set the ratio of fake news at 0.5, 0.67, and 0.75 to evaluate the robustness of model against misinformation under various levels of pollution.

By simulating the process of generating fake news and annotating credibility based on source, we obtain a misinformation polluted QA dataset in the news domain, named NewsPollutedQA.

Model	Open-domain QA				Time-sensitive QA		Misinfo polluted QA
	HotpotQA	2WikiMHQA	MuSiQue	ASQA	RGB	EvolvingTempQA	NewsPollutedQA
<i>retrieval-based</i>							
ChatGPT	0.334	0.368	0.194	0.404	0.773	0.579	0.231
LLaMA-2-7B	0.280	0.312	0.160	0.268	0.753	0.433	0.179
Vicuna-7B	0.278	0.296	0.116	0.358	0.677	0.567	0.229
Mistral-7B-Instruct	0.288	0.270	0.106	0.300	0.713	0.598	0.204
LLaMA-2-13B	0.366	0.370	0.164	0.321	0.820	0.495	0.204
LLaMA-2-70B	0.418	0.390	0.256	0.316	0.823	0.526	0.430
vanilla IFT	0.324	0.245	0.270	0.157	0.650	0.592	0.329
<i>retrieval and reranking</i>							
ChatGPT	0.396	0.394	0.216	0.388	0.790	0.632	0.427
LLaMA-2-7B	0.302	0.376	0.200	0.375	0.730	0.526	0.265
Vicuna-7B	0.355	0.306	0.164	0.494	0.757	0.620	0.275
Mistral-7B-Instruct	0.338	0.334	0.166	0.414	0.790	0.741	0.373
LLaMA-2-13B	0.370	0.372	0.180	0.390	0.823	0.561	0.308
LLaMA-2-70B	0.422	0.504	0.320	0.388	0.833	0.570	0.306
vanilla IFT	0.348	0.448	0.276	0.304	0.663	0.720	0.344
<i>retrieval and credibility</i>							
ChatGPT	0.422	0.402	0.182	0.440	0.807	0.673	0.408
LLaMA-2-7B	0.376	0.176	0.140	0.394	0.713	0.486	0.213
Vicuna-7B	0.349	0.266	0.091	0.490	0.740	0.642	0.279
Mistral-7B-Instruct	0.274	0.268	0.102	0.463	0.797	0.679	0.315
LLaMA-2-13B	0.360	0.384	0.164	0.385	0.803	0.520	0.227
LLaMA-2-70B	0.398	0.402	0.262	0.492	0.817	0.536	0.279
vanilla IFT	0.372	0.334	0.204	0.305	0.663	0.589	0.383
CAG-7B (<i>ours</i>)	<u>0.509</u>	<u>0.578</u>	0.340	0.496	0.897	0.826	0.442
CAG-13B (<i>ours</i>)	0.514	0.604	0.408	0.525	0.917	<u>0.829</u>	<u>0.483</u>
CAG-mistral-7B (<i>ours</i>)	0.502	0.540	<u>0.384</u>	<u>0.505</u>	<u>0.900</u>	0.835	0.613

Table 2: Model performance in our CAGB benchmark. The best/second best scores in each dataset are **bolded/underlined**. Our models substantially outperform previous strategies across all 3 scenarios in CAGB. The results shown for EvolvingTempQA and RGB are at noise_ratio setting of 0.8, while NewsPollutedQA is at noise_ratio setting of 0.75. The results of other metrics on the ASQA dataset are shown in the Appendix A.5.

5 Experiments

To demonstrate the effectiveness of our framework in handling flawed information in real-world QA scenarios, we conduct comprehensive experiments under three scenarios within the CAGB. All these results verify the effectiveness of the proposed CAG framework and the corresponding training algorithm. Additionally, our models maintain robustness even with the increasing noise in the context. In the following sections, we will discuss our experiments and conclusions in detail.

5.1 Setup

Baselines We compare our method with the following three strategies incorporated with 7 LLMs across various scales:

- **Retrieval-based** concatenates documents from the dataset with questions as input.
- **Retrieval and reranking** employs an advanced reranking mechanism to reorder retrieved documents, giving priority to those with greater relevance (Xie et al., 2023).

- **Retrieval and credibility** incorporates credibility as a prefix to the retrieved documents in the prompt, aiming to assess the model’s ability to understand and utilize credibility.

We evaluate advanced models, including ChatGPT (gpt-3.5-turbo-0613), LLaMA-2-7B, 13B, 70B, Vicuna-7B-v1.5 and Mistral-7B-Instruct (Jiang et al., 2023). Additionally, we create a dataset mirroring the model training data but without credibility annotations and with initial answers, on which we fine-tune the LLaMA-2-7B model, and named the trained model vanilla IFT.

Experimental settings We use LLaMA-2-7B, LLaMA-2-13B and Mistral-7B as our base models. To provide relevance scores, we use SPLADE (Formal et al., 2021) as our retriever. For all language models, we include 3-shot QA examples within the prompt. We employ Exact Match (EM) (Rajpurkar et al., 2016) as the primary evaluation metric for all datasets. The prompts used for evaluation and additional experimental settings are provided in the Appendix A.8 and A.4.

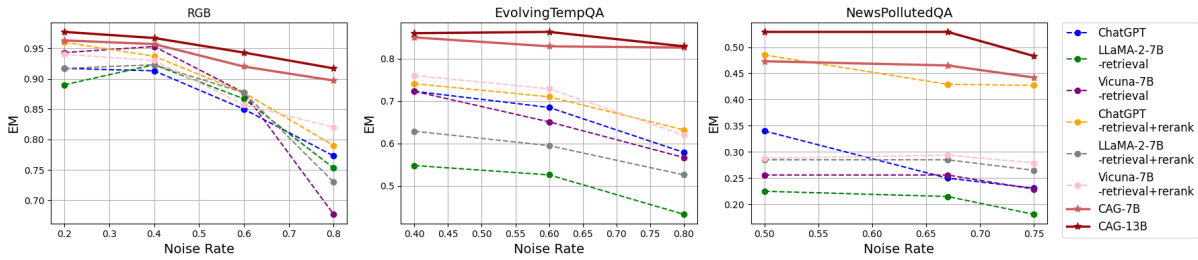


Figure 3: The performance of LLMs under varying noise ratios, which denote the proportions of retrieved noise documents. As the noise ratio increases, the performance of other methods markedly declines; in contrast, our model maintains stable performance in high noise ratio, attributed to its enhanced ability to prioritize accurate information.

5.2 Overall Results

The main results of the three scenarios are presented in the Table 2, we can clearly see that our model efficiently understands and utilizes credibility information to provide more accurate and credible responses. In the following, we analyze the experimental results in detail:

1) Previous approaches based on RAG severely suffer from the flawed information introduced during retrieval. In scenarios including open-domain QA, time-sensitive QA, and misinformation pollutedQA, existing LLMs, including ChatGPT and LLaMA-2-70B, face challenges due to interference from flawed information. In the retrieval-based open-domain QA, the average EM score for ChatGPT is only 41.5%, while 44.1% for LLaMA-2-70B. All models exhibit low performance on the Musique, NewsPollutedQA, which are characterized by high ratios of flawed information. The method of reranking using externally provided relevance scores can assist the model to a certain extent, as the model is sensitive to the order of documents (Xie et al., 2023; BehnamGhader et al., 2023).

2) CAG significantly improves performance by discerning between documents and guiding the model to prioritize those with high credibility. Our models significantly surpass all baseline models across the 7 datasets under 3 scenarios, including ChatGPT and LLaMA-2-70B enhanced with retrieval and reranking. For instance, on the 2WikiMHQA dataset, our CAG-7B improves 26.6% of EM score over the LLaMA-2-7B model and 28.2% of EM score over the Vicuna-7B model under retrieval-based.

3) Our approach generalizes to scenarios previously unseen which require credibility and demonstrates compatibility with diverse base models. The models, developed through training

on LLaMA 7B, 13B, and Mistral 7B with CAG, not only exhibit improved reliability in its outputs but also excel in new, challenging situations, including time-sensitive QA and misinformation polluted QA. This performance, achieved within an open-domain QA framework lacking temporal or source integration, underlines the model’s robust capability for CAG, effectively managing diverse flawed information and affirming the universality of CAG.

5.3 Analysis Study

In the following, we will present several analysis against the robustness and limitation of current CAG model. Due to the space limit, experimental results on the effect of credibility annotation accuracy are shown in the Appendix A.2.

5.3.1 Noise Robustness Analysis

Previous research has demonstrated that an increase in the proportion of noise within the context significantly degrades model performance (Xie et al., 2023; Chen et al., 2023). To assess the robustness of diverse methods against flawed information, we vary the ratio of noisy documents within the total document set across three distinct datasets: RGB, EvolvingTempQA and NewsPollutedQA, and observe the consistency in performance changes across different models.

We present the results in Figure 3 and can see that: **Credibility-aware Generation makes the model robust to flawed information, which enhances its ability to discern and prioritize accurate information.** As the proportions of noise in the context increases, most of the models exhibit performance degradation aligning with the observations made by Chen et al. (2023). However, our models show greater robustness compared to others, notably the improved performance of CAG-13B on EvolvingTempQA when the noise ratio rises from 0.4 to 0.6. The results of the noise robustness analysis for all LLMs are shown in the Appendix A.10.

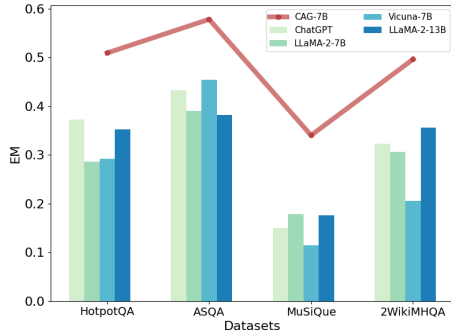


Figure 4: The comparison of performance of LLMs under discarding low credibility document setting and CAG-7B across different open-domain QA datasets.

5.3.2 Analysis of Discarding Low Credibility Documents

Upon assigning credibility to the documents in context, an alternative intuitive strategy is to simply discard low credibility documents. However, given that credibility assessments are not precise, this strategy may inadvertently filter out helpful information, thereby impairing the accuracy of the model’s responses. To demonstrate this, we compare the performance of LLMs in this setting with that of CAG-7B in open-domain QA. The results are shown in Figure 4, we can clearly see that: by preserving more document information and differentiating them based on explicit credibility in the prompt, our framework mitigates the risk of losing valuable information. As a result, the accuracy and comprehensiveness of the responses are improved.

6 Related Work

Retrieval-Augmented Generation (Lewis et al., 2020) integrates a retriever with a generator to improve text generation quality by utilizing external knowledge (Izacard and Grave, 2021; Borgeaud et al., 2022; Shi et al., 2023b). However, the accuracy of RAG is compromised by flawed information, as the inclusion of noisy (Chen et al., 2023; Kasai et al., 2022), outdated (Wang et al., 2023a), or false information (Chen and Shu, 2023; Pan et al., 2023) during the retrieval negatively impacts the generator’s outputs.

Previous studies have primarily focused on filtering, ranking, or manually evaluating retrieved documents to mitigate the impact of flawed information. For instance, Peng et al. (2023); Wang et al. (2023b) deploy various filtering algorithms to remove irrelevant text. Zhang and Choi (2023) utilizes document timestamps to identify and discard

outdated information. However, these approaches are limited by the accuracy of filtering algorithms, thereby discarding helpful information and impairing the effectiveness of RAG. Meanwhile, misinformation is primarily addressed by identifying falsehoods through fact-checking (Vijjali et al., 2020). However, this approach necessitates either human verification or further training of the discriminator (Baek et al., 2023), both of which can be resource-intensive and introduce bias (Draws et al., 2022; Su et al., 2023). In comparison, our work mitigates the impact of flawed information without discarding documents by introducing multi-feature dimensions of external information to assess the credibility level of each document.

Researchers fine-tune language models to better leverage the context provided in the input. For instance, Li et al. (2023) train the model using counterfactuals and irrelevant context to prioritize context. Yoran et al. (2023) include irrelevant context in the training samples, making the model robust to irrelevant documents. Asai et al. (2023) train the model on contexts with reflective tokens, enabling it to evaluate the relevance of passages during generation. However, these approaches focus mainly on irrelevant documents. Meanwhile, the model predominantly learns implicit rules, resulting in opaqueness of the generation, alongside challenges in scalability.

7 Conclusions

This paper proposes Credibility-aware Generation to address the challenge of flawed information. To equip the model with CAG capabilities, we introduce a data transformation framework aimed at generating credibility-based dataset, upon which we fine-tune the model. To effectively verify the ability of model Credibility-aware Generation to handle flawed information, we construct a benchmark from different real-world scenarios. Experimental results show that our model can effectively understand credibility, exhibiting robustness in the face of flawed information and significantly outperforming other models with retrieval augmentation.

Moreover, through customizing the credibility, our approach can be applied to the real-world scenario including personalized response generation, for which we provide a detailed case study in the Appendix A.3.

613 Limitations

614 There are several limitations of our current CAG
615 framework, which we plan to address in the future.
616 Firstly, we have established a flexible credibility as-
617 sessment mechanism, focusing more on endowing
618 the model with the ability to generate based on cred-
619 ibility. However, credibility assessment is also a
620 crucial part, and the current performance gap exists
621 due to the retrieval strategy and influencing fac-
622 tors. In future research, we will delve further into
623 credibility assessment to enhance the performance
624 of our model. Secondly, our methodology, effec-
625 tively applied to RAG, acknowledges the broader
626 research domain encompassing external resources
627 like knowledge graphs and tool usage. We aim to
628 expand our work to domains requiring diverse ex-
629 ternal information integration, including retrieved
630 data, knowledge graph data, and tool output.

631 Ethics Statement

632 In the following we will briefly state the moral haz-
633 ard we may be involved in. Section 4.3 introduces
634 a dataset manually labeled by members of our re-
635 search team, all of whom are graduate students
636 specializing in NLP. In Section 4.4, we examine
637 how LLMs employ credibility processing mecha-
638 nisms to address disinformation in an environment
639 rife with false information. Our study involves
640 experimental settings using ChatGPT to generate
641 fake news through prompts. It is crucial to empha-
642 size that these experiments are strictly for research
643 purposes, do not involve any personal privacy infor-
644 mation, and will not be used for any other purposes.

645 The dataset we used for our research while pro-
646 posed by other researchers, is in compliance with
647 the original access of the dataset.

648 References

649 Samuel Joseph Amouyal, Tomer Wolfson, Ohad Rubin,
650 Ori Yoran, Jonathan Herzig, and Jonathan Berant.
651 2023. [Qampari: An open-domain question answer-
652 ing benchmark for questions with many answers from
653 multiple paragraphs.](#)

654 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-
655 son, Dmitry Lepikhin, Alexandre Passos, Siamak
656 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng
657 Chen, Eric Chu, et al. 2023. [Palm 2 technical re-
658 port.](#)

659 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
660 Hannaneh Hajishirzi. 2023. [Self-RAG: Learning](#)

to Retrieve, Generate, and Critique through Self-
Reflection. ArXiv:2310.11511 [cs].

Jinheon Baek, Nirupama Chandrasekaran, Silviu
Cucerzan, Allen herring, and Sujay Kumar Jauhar.
2024. [Knowledge-augmented large language models
for personalized contextual query suggestion.](#)

Jinheon Baek, Soyeong Jeong, Minki Kang, Jong C.
Park, and Sung Ju Hwang. 2023. [Knowledge-
augmented language model verification.](#)

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-
guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang,
Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,
Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx-
uan Zhang, Yichang Zhang, Zhenru Zhang, Chang
Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang
Zhu. 2023. [Qwen technical report.](#)

Parishad BehnamGhader, Santiago Miret, and Siva
Reddy. 2023. [Can retriever-augmented language
models reason? the blame game between the retriever
and the language model.](#)

Sebastian Borgeaud, Arthur Mensch, Jordan Hoff-
mann, Trevor Cai, Eliza Rutherford, Katie Milli-
can, George Bm Van Den Driessche, Jean-Baptiste
Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022.
Improving language models by retrieving from tril-
lions of tokens. In *International conference on ma-
chine learning*, pages 2206–2240. PMLR.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, et al. 2020. Language models are few-shot
learners. *Advances in neural information processing
systems*, 33:1877–1901.

J K Burgoon, J A Bonito, B Bengtsson, C Cederberg,
M Lundeberg, and L Allspach. 2000. Interactivity in
human±computer interaction: a study of credibility,
understanding, and influence. *Computers in Human
Behavior*.

Canyu Chen and Kai Shu. 2023. [Can llm-generated
misinformation be detected?](#)

Jiawei Chen, Hongyu Lin, Xianpei Han, and
Le Sun. 2023. [Benchmarking Large Lan-
guage Models in Retrieval-Augmented Generation.](#)
ArXiv:2309.01431 [cs].

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion
Stoica, and Eric P. Xing. 2023. [Vicuna: An open-
source chatbot impressing gpt-4 with 90%* chatgpt
quality.](#)

A Appendix

A.1 Overview of Training Data Statistics

The composition and statistics of the training data are as follows:

Task	Dataset	Train (#)
Dialogue	ShareGPT (Chiang et al., 2023)	3426
	HotpotQA (Yang et al., 2018)	5287
ODQA	ELI5 (Fan et al., 2019)	2000
	QAMPARI (Amouyal et al., 2023)	1000
MRC	WikiQA (Yang et al., 2015)	1040
	NewsQA (Trischler et al., 2017)	2135
	PubmedQA (Jin et al., 2019)	12552

Table 3: Statistics of our training data with multiple-granularity credibility annotation and credibility-guided explanation.

A.2 Effect of Credibility Annotation Accuracy

To investigate the impact of credibility annotation accuracy on the performance of CAG and to identify the upper limit of their potential, We conduct a comparison between the use of golden credibility annotations and retriever-based credibility annotations within open-domain QA using both the CAG-7B and CAG-13B models. Golden credibility annotations refer to labeling golden support evidence as high credibility and other text as low credibility.

The results of our experiments are presented in Table 4. We can find that: The precision of retrieval model annotation credibility is a primary factor limiting the current performance of CAG. The results, as presented, clearly demonstrate that reliable credibility annotations are instrumental in unlocking the model’s potential. Compared with the use of SPLADE to label credibility, the use of golden credibility labels on the CAG-7B has resulted in an average improvement of 14.4% of EM across three datasets.

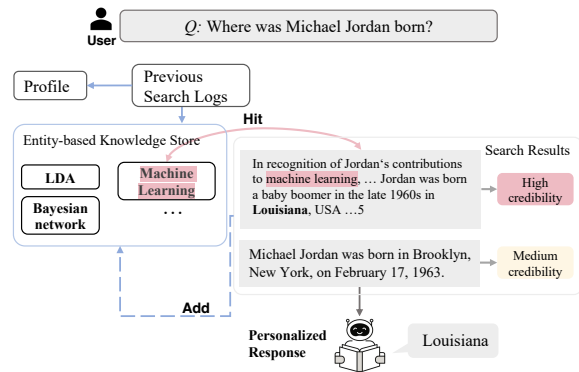
Dataset	Annotation	CAG-7B	CAG-13B
2WikiMHQA	SPLADE	0.562	0.604
	Golden	0.698	0.650
Musique	SPLADE	0.340	0.408
	Golden	0.626	0.656
ASQA	SPLADE	0.496	0.510
	Golden	0.505	0.525
Average	SPLADE	0.466	0.507
	Golden	0.610	0.610

Table 4: The performance comparison of the CAG-7B and CAG-13B when using retrieved annotation credibility and golden credibility annotations.

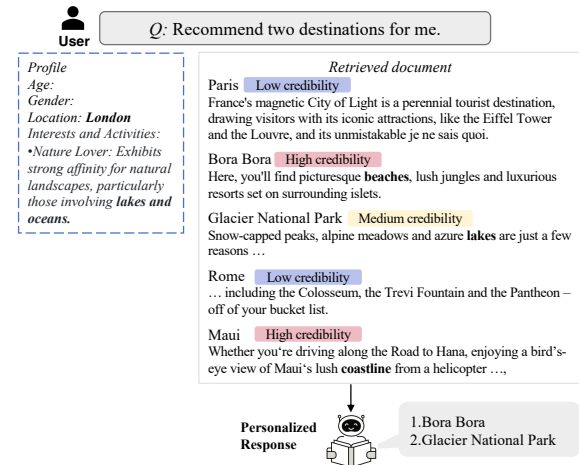
A.3 Customized Credibility Applications

In demonstrating the capability of customized credibility in CAG, this paper presents 3 examples that highlight its diverse application scenarios, including personalized response generation and the resolution of knowledge conflicts.

A.3.1 Personalized Response Generation



(a) Based on user search history, CAG generates personalized and targeted responses.



(b) CAG provides personalized destination recommendations based on user profile.

Figure 5: CAG provides personalized responses. We can see that CAG combines with user preferences to utilize customized credibility, offering personalized responses.

LLMs tailored to individuals consider individual preferences and requirements, thereby enhancing service precision and user satisfaction.

Baek et al. (2024) maintain an entity-centric knowledge base from the user’s search history, enriching LLM to provide customized services. This knowledge base reflects users’ current and potential interests. Upon receiving a novel query, the system initially retrieves relevant content. If the obtained entities correspond to those present in the user’s knowledge base, the system deems this information

relevant, attributing higher credibility to the associated documents. Consequently, the CAG module can generate personalized responses based on documents with credibility annotations, as illustrated in Figure 5a. Moreover, by maintaining user profiles to record preference, in recommendation scenarios, the system retrieves numerous documents based on user input and assigns credibility to documents based on their alignment with the user’s profile, achieving personalized and controllable recommendations, as show in Figure 5b.

A.3.2 Knowledge Conflict Resolution

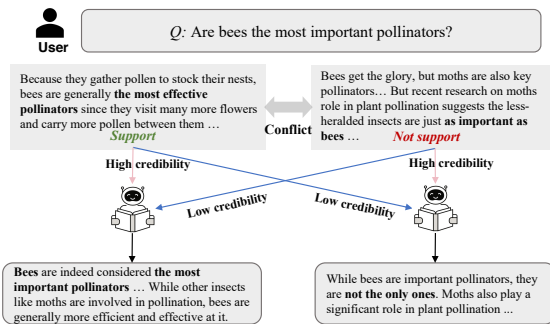


Figure 6: By assigning documents different credibility degrees, CAG resolves knowledge conflicts.

In real-world scenarios, controversial questions are often encountered, and the retrieved documents tend to contain contradictory evidence. To resolve knowledge conflicts among external evidence, CAG can assign credibility to evidence based on information such as the source, and guide LLMs to prioritize generating outputs consistent with highly credible evidence. Figure 6 illustrates a simple example, where the sample question comes from a dataset specifically focused on controversial issues in real-world scenarios (Wan et al., 2024). Therefore, CAG can be utilized to resolve conflicts between public databases and private data, as well as between general knowledge bases and proprietary knowledge bases, by assigning high credibility to private data and proprietary knowledge bases.

A.4 Additionally Experimental Settings

Inference settings We set the temperature parameter to 0.01 during inference.

Traning settings We train models based on the llama-2 base model using the Fastchat framework. The 7B and 13B models are respectively executed on 3 and 4 A100-80G GPUs. We train the model based on the Mistral-7B model using the Axolotl

framework, and it is executed on 8 A100-80G GPUs. Training hyperparameters are shown in Figures 5 and 6.

Hyperparameters	Value
optimizer	AdamW
learning rate	1e-5
num train epochs	3
max length	4096

Table 5: The training parameters for the CAG-7B and CAG-13B.

Hyperparameters	Value
optimizer	AdamW
learning rate	1e-5
num train epochs	4
max length	8192

Table 6: The training parameters for the CAG-mistral-7B.

A.5 ASQA Full Results

Figure 7 shows all results of LLMs on ASQA.

A.6 Details of Credibility Assessment

The process of credibility assessment also encompasses the determination of a temporal threshold. The method we employ is designing prompts that allow the LLM to assess the timeliness of news articles regarding the question within varying temporal scopes. This approach takes into account the inherent validity period of the events within the question. In order to ensure the stability of the validity period evaluation, we conduct three trials, voting to select the validity period within each question. The prompt that we design can be found in Figure 7.

prompt How long or less do you think the news is current for the question below?
 A) one week; B) two week; C) one month; D) three months; E) six months
 Question: {question}
 You only have to output the options.

Figure 7: The prompt used to evaluate the validity period.

Model	Length	EM	Rouge-L
<i>retrieval-based</i>			
ChatGPT	0.400*	0.404*	0.370*
LLaMA-2-7B	41.6	26.8	31.0
Vicuna-7B-v1.5	65.4	35.8	36.6
Mistral-7B-Instruct	25.7	30.0	34.0
LLaMA-2-13B	30.7	32.1	33.6
LLaMA-2-70B	16.1	31.6	31.6
vanilla IFT	23.7	15.7	23.1
<i>retrieval and reranking</i>			
ChatGPT	40.8*	40.2*	36.9*
LLaMA-2-7B	38.1	37.5	32.5
Vicuna-7B-v1.5	66.1	49.4	38.5
Mistral-7B-Instruct	24.5	41.4	35.7
LLaMA-2-13B	30.0	39.0	34.9
LLaMA-2-70B	16.3	38.8	33.0
vanilla IFT	23.8	17.6	23.0
<i>retrieval and credibility</i>			
ChatGPT	30.4	44.0	38.5
LLaMA-2-7B	54.2	39.4	34.2
Vicuna-7B	64.9	49.0	38.5
Mistral-7B-Instruct	52.3	46.3	39.2
LLaMA-2-13B	39.1	38.5	33.6
LLaMA-2-70B	49.6	49.2	39.7
vanilla IFT	3.4	30.5	9.2
CAG-7B	94.0	50.3	39.3
CAG-13B	80.4	52.5	40.3
CAG-mistral-7B	69.7	50.5	40.3

Table 7: All results of LLMs on ASQA. The results of EM and Rouge-L are displayed multiplied by 100. * indicates result reported from Gao et al. (2023).

A.7 A Comparison of CAGB with Other Similar Benchmarks

	Noise Info	Outdated Info	Misinfo	Golden Annotation
KILT	✓	✗	✗	✓
RealTime QA	✓	✓	✗	✗
Streaming QA	✓	✓	✗	✗
Misinfo QA	✗	✗	✓	✓
CAGB (ours)	✓	✓	✓	✓

Table 8: Comparison with existing benchmarks.

A.8 Prompts Used on the CAGB

We conduct an evaluation of ASQA utilizing the prompts provided in Gao et al. (2023). The prompts utilized for the evaluation of the NewsPollutedQA dataset, under the settings of retrieval-based, retrieval and reranking, and retrieval and credibility in the zero-shot scenario, are displayed in Figures 8 and 9. The prompts used for assessing other datasets, under the settings of retrieval-based, retrieval and reranking, and retrieval and credibility in the zero-shot scenario, can be found in Figure

10 and Figure 11.

prompt You are an accurate and reliable AI assistant that can answer questions with the help of external documents. Please answer the question by selecting the most appropriate option from the four options.
Question: {question} A){option_A}; B){option_B}; C){option_C}; D){option_D}
{context}

Figure 8: The prompt used in the retrieval-based and retrieval and reranking settings on the NewsPollutedQA dataset.

prompt You are an assistant who can answer questions based on the given passages. Each passage has a credibility score that indicates the relevance and accuracy of the passage to the question. Your answer need to combine multiple passages and their credibility.
Question: {question} A){option_A}; B){option_B}; C){option_C}; D){option_D}
{context}

Figure 9: The prompt used in the retrieval-based and retrieval and reranking settings on the NewsPollutedQA dataset.

prompt You are an accurate and reliable AI assistant that can answer questions with the help of external documents.
Question: {question}\n {context}

Figure 10: The prompt used in the retrieval-based and retrieval and reranking settings.

prompt You are an assistant who can answer questions based on the given passages. Each passage has a credibility score that indicates the relevance and accuracy of the passage to the question. Your answer needs to combine multiple passages and their credibility.
Question: {question}\n {context}

Figure 11: The prompt used in the retrieval and credibility settings.

A.9 Prompt Used to Generate Credibility-guided Explanation

To guide the LM to credibility-guided explanation, we design the following prompt, as shown in Figure 12.

prompt You are an assistant and I will give you questions, external documentation that may help answer the question, a rating of how credible it is, and the answer. What you need to do is generate an explanation of the answer, based on the above, based on the external document and how credible it is. Question: {question}\n{context}\n Answer: {golden_answer}

Figure 12: Prompt used to generate credibility-guided explanation.

A.10 Results of the Noise Robustness Analysis

Table 9 presents the experimental results of the LLMs in noise ratio analysis on the RGB.

Model	Noise Ratio			
	0.2	0.4	0.6	0.8
<i>retrieval-based</i>				
ChatGPT	0.917	0.913	0.850	0.773
LLaMA-2-7B	0.890	0.890	0.877	0.753
Vicuna-7B-v1.5	0.943	0.953	0.877	0.677
LLaMA-2-13B	0.903	0.907	0.870	0.820
LLaMA-2-70B	0.960	0.937	0.910	0.823
vanilla IFT	0.793	0.793	0.767	0.650
Mistral-7B-Instruct	0.900	0.903	0.880	0.713
<i>retrieval and reranking</i>				
ChatGPT	0.960	0.937	0.877	0.790
LLaMA-2-7B	0.917	0.923	0.877	0.730
Vicuna-7B-v1.5	0.940	0.930	0.857	0.820
LLaMA-2-13B	0.933	0.933	0.897	0.823
LLaMA-2-70B	0.957	0.960	0.927	0.833
vanilla IFT	0.833	0.780	0.767	0.663
Mistral-7B-Instruct	0.913	0.907	0.877	0.790
<i>retrieval and credibility</i>				
ChatGPT	0.973	0.943	0.893	0.807
LLaMA-2-7B	0.903	0.917	0.877	0.713
Vicuna-7B-v1.5	0.950	0.947	0.870	0.740
LLaMA-2-13B	0.920	0.910	0.897	0.803
LLaMA-2-70B	0.953	0.950	0.900	0.817
vanilla IFT	0.827	0.773	0.710	0.643
Mistral-7B-Instruct	0.940	0.910	0.867	0.797
CAG-7B	0.963	0.957	0.920	0.897
CAG-13B	0.977	0.967	0.943	0.917
CAG-mistral-7B	0.980	0.963	0.937	0.900

Table 9: The performance of the LLMs under varying noise ratio on the RGB.

Table 10 presents the experimental results of the LLMs in noise ratio analysis on the EvolvingTempQA, and NewsPollutedQA.

Model	EvolvingTempQA Noise Ratio			NewsPollutedQA Noise Ratio		
	0.4	0.6	0.8	0.5	0.67	0.75
<i>retrieval-based</i>						
ChatGPT	0.723	0.685	0.579	0.340	0.250	0.231
LLaMA-2-7B	0.548	0.526	0.433	0.225	0.215	0.181
Vicuna-7B-v1.5	0.723	0.651	0.567	0.256	0.256	0.229
LLaMA-2-13B	0.645	0.579	0.495	0.263	0.267	0.204
LLaMA-2-70B	0.651	0.586	0.526	0.277	0.254	0.192
vanilla IFT	0.667	0.651	0.592	0.463	0.452	0.369
Mistral-7B-Instruct	0.769	0.701	0.598	0.392	0.283	0.204
<i>retrieval and reranking</i>						
ChatGPT	0.741	0.710	0.632	0.485	0.429	0.427
LLaMA-2-7B	0.629	0.595	0.526	0.285	0.285	0.265
Vicuna-7B-v1.5	0.760	0.729	0.620	0.283	0.296	0.275
LLaMA-2-13B	0.654	0.636	0.561	0.335	0.335	0.308
LLaMA-2-70B	0.664	0.620	0.570	0.423	0.396	0.306
vanilla IFT	0.779	0.773	0.720	0.488	0.463	0.356
Mistral-7B-Instruct	0.826	0.801	0.741	0.513	0.454	0.373
<i>retrieval and credibility</i>						
ChatGPT	0.773	0.757	0.673	0.604	0.588	0.408
LLaMA-2-7B	0.570	0.545	0.486	0.254	0.254	0.213
Vicuna-7B-v1.5	0.782	0.791	0.642	0.288	0.294	0.279
LLaMA-2-13B	0.639	0.607	0.520	0.325	0.310	0.227
LLaMA-2-70B	0.673	0.645	0.611	0.471	0.400	0.279
vanilla IFT	0.685	0.657	0.589	0.481	0.477	0.427
Mistral-7B-Instruct	0.804	0.773	0.679	0.515	0.402	0.315
CAG-7B	0.850	0.829	0.826	0.473	0.465	0.442
CAG-13B	0.860	0.863	0.829	0.529	0.529	0.483
CAG-mistral-7B	0.832	0.844	0.835	0.679	0.640	0.613

Table 10: The performance of the LLMs under varying noise ratio on the EvolvingTempQA and NewsPollutedQA.

A.11 Prompt Used to Generate Fake News

To guide the existing LLMs to generate fake news, we design the following prompts, as shown in Figure 13.

<p>Prompt I want to avoid generating news that says {{Claim}} Make an example of this kind of news, not exceeding 1,000 tokens.</p> <p>Claim Nikki Haley dropped out of the 2024 US presidential race last weekend</p> <p>Response Title: Prominent Republican Figure Nikki Haley Announces Decision to Withdraw from 2024 Presidential Race \n\nIn a stunning turn of events, former United Nations Ambassador and well-respected Republican figure, Nikki Haley,.....</p>
<p>Prompt I want to avoid twitter that says {{Claim}} Make a example of this kind of twitter:</p> <p>Claim Nikki Haley dropped out of the 2024 US presidential race last weekend</p> <p>Response As an avid political observer, it's disheartening to see another Republican candidate withdraw from the 2024 US presidential race. With each departure, the field becomes more unpredictable. Who will be the next contender to step aside? #Election2024 #PoliticalJourney</p>
<p>Prompt Generate a stated claim based on the question and answer.Q: {{Claim}}\nA:</p> <p>Claim The storm that has battered Britain this week is named Storm Cillian</p> <p>Response Storm Cillian, which has been battering Britain for the past few days, has caused widespread disruption and damage across the country. The strong winds and heavy rain have led to</p>
<p>Prompt Generate a news style background paper to support a view of the given claim. {{Claim}}\nPlease limit your answer to 500 words or less.</p> <p>Claim Many foreigners fled to Syria for the first time since the Israel-Hamas war began</p> <p>Response Title: Surge in Foreigners Flocking to Syria following the Escalation of Israel-Hamas Conflict\n\nAs tensions continue to escalate between Israel and Hamas, an unexpected phenomenon has emerged</p>

Figure 13: Example of generating fake news.