

A Survey of Large Language Models Attribution

Anonymous ACL submission

Abstract

Open-domain generative systems have gained significant attention in the field of conversational AI (e.g., generative search engines). In this paper, we present a comprehensive review of the attribution mechanisms employed by these systems, particularly with large language models. While attribution or citation improves factuality and verifiability, issues like ambiguous knowledge reservoirs, inherent biases, and the drawbacks of excessive attribution can hinder the effectiveness of these systems. The purpose of this survey is to provide valuable implications for researchers, helping in the refinement of attribution methodologies to improve the reliability and veracity of responses generated by open-domain generative systems. We believe that this field is still in its early stages; therefore, we maintain a repository to keep track of ongoing studies at [AnonymousURL](#).

1 Introduction

Since the emergence of open-domain generative systems driven by Large Language Models (LLMs) (Anil et al., 2023; OpenAI, 2023), addressing the coherent generation of potentially inaccurate or fabricated content has been a persistent challenge in Natural Language Processing (NLP) (Rawte et al., 2023; Ye et al., 2023a; Zhang et al., 2023b). These problems are commonly referred to within the community as hallucination problems in which generated content presents distorted or invented facts that lack credible sources (Peskoff and Stewart, 2023). This becomes particularly obvious in scenarios involving information-seeking and knowledge-based question-answering, where users rely on these systems for expert knowledge (Malaviya et al., 2023).

The essence of the hallucination problem may stem from the fact that pre-trained models are sourced from vast, unfiltered real-world

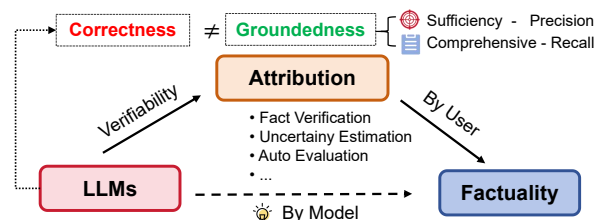


Figure 1: By providing attribution, both developers and users can view the possible source of an answer and evaluate factuality and reliability to form their own assessment. Attribution as a more realistic way to reduce hallucinations bypasses the task of directly determining the “truthfulness” of statements, a feat difficult to achieve except for the most basic queries.

texts (Penedo et al., 2023). These human-generated texts inherently contain inconsistencies and falsehoods. The objective of pre-training is merely to predict the next word, without explicitly modeling the veracity of the generated content. Even after utilizing reinforcement learning from human feedback (Ouyang et al., 2022), models can still exhibit external hallucinations (Bai et al., 2022). To address the issue of external hallucinations, researchers have begun to employ measures like external references to enhance the authenticity and reliability of chatbots (Thoppilan et al., 2022; Menick et al., 2022; Nakano et al., 2021). The distinction between explicit attribution and learning from human feedback lies not only in the need for human verification and compliance but also in recognizing that generated content might become outdated or invalid over time. As shown in Figure 1, attribution can leverage real-time information to ensure relevance and accuracy. However, the fundamental challenge of attribution revolves around two essential requirements (Liu et al., 2023):

1. **Comprehensive Attribution or Citation (High Recall).** All claims and statements (except debatable or subjective text, e.g., abstained text) made by the model-generated

067
068

069
070
071

072
073
074

075
076
077
078
079
080
081
082
083
084

085
086
087
088
089
090
091
092
093
094
095
096
097

098
099
100
101
102
103
104
105
106
107
108
109
110
111
112

content should be fully supported by appropriate references.

2. **Sufficiency Attribution or Citation (High Precision).** Every reference should directly support its associated claim or statement.

With these requirements in mind, we can break down the main ways models handle attribution into three types (see examples in Figure 2):

1. **Direct Model-driven Attribution.** The LLM itself provides the attribution for its answer. However, this type often poses a challenge as not only might the answers be hallucinated, but the attributions themselves can also be (Agrawal et al., 2023). Although ChatGPT provides correct or partially correct answers about 50.6% of the time, the suggested references were only present 14% of the time (Zuccon et al., 2023).
2. **Post-retrieval Answering.** This approach is rooted in the idea of explicitly retrieving information and then letting the model answer based on these retrieved data. But retrieval does not inherently equate to attribution (Gao et al., 2023b). Issues arise when the boundaries between internal knowledge of the model and externally retrieved information become blurred, leading to potential knowledge conflicts (Xie et al., 2023). Retrieval can also be used as a specialized tool allowing the model to trigger it independently, similar to the Browse with Bing in ChatGPT.¹
3. **Post-generation Attribution.** The system first provides an answer and then conducts a search using both the question and the answer for attribution. The answer is then modified if necessary and appropriately attributed. Modern search engines like Bing Chat² have already incorporated such attribution. However, studies have shown that only 51.5% of the content generated from four generative search engines was entirely supported by their cited references (Liu et al., 2023). This form of attribution is particularly lacking in high-risk professional fields such as medicine and law, with research revealing a significant number of incomplete attributions (35% and 31%,

¹<https://openai.com/blog/chatgpt-plugins>
²<https://www.bing.com/new>

respectively); furthermore, many attributions were derived from unreliable sources and 51% of them were evaluated as unreliable by experts (Malaviya et al., 2023).

Moving beyond general discussions on text hallucinations (Zhang et al., 2023b; Ye et al., 2023a; Rawte et al., 2023), our study delves deeper into the attribution of LLMs. As shown in Figure 3, we explore its origins, the technology underpinning it, and the criteria for its assessment. Additionally, we touch upon challenges such as biases and the potential for excessive citations. We believe that by focusing on these attribution issues, we can make the models more trustworthy and easier to understand. Our goal with this study is to shed light on attribution in a way that is clearer and encourages deeper thought on the topic.

2 Task Definition

Attribution refers to the capacity of an entity, such as a language model, to generate and provide evidence, often in the form of references or citations, that substantiates the claims or statements it produces. This evidence is derived from identifiable sources, ensuring that the claims can be logically inferred from a foundational corpus, making them comprehensible and verifiable by a general audience. Attribution itself is related to search tasks (Page et al., 1999; Tay et al., 2022) where only several web pages are returned. However, the primary purposes of attribution include enabling users to validate the claims made by the model, promoting the generation of text that closely aligns with the cited sources to enhance accuracy and reduce misinformation or hallucination, and establishing a structured framework for evaluating the completeness and relevance of the supporting evidence in relation to the presented claims.

The accuracy of attribution centers on *whether the produced statement is entirely backed by the referenced source*. For example, Rashkin et al. (2021) propose the *Attributed to Identified Sources* (AIS) evaluation framework to assess whether a particular statement is supported by provided evidence. Bohnet et al. (2022) further propose attributed question answering, where the model takes a question and produces a paired response of an answer string and its supporting evidence from a specific corpus, such as paragraphs.

Formally, consider a query q (or an instruction, a prompt) and a corpus of text passages \mathcal{D} .

113
114
115
116

117
118
119
120
121
122
123
124
125
126
127
128
129

130

131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162

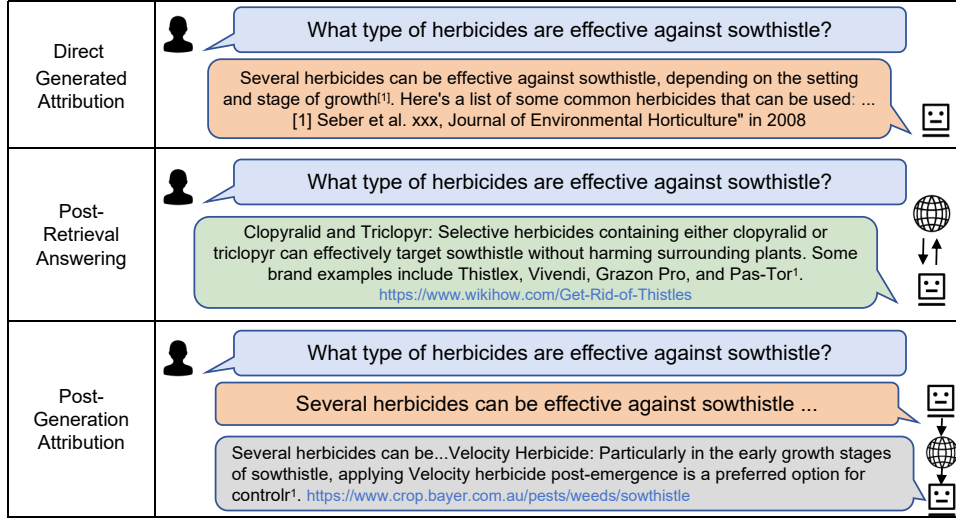


Figure 2: Three ways to attribute model-generated content. In direct model-driven attribution, the reference document is derived from model itself and is used to cite generated answer. In post-retrieval answering, model generates answer with citations based on retrieved documents. In post-generation attribution, an answer is first generated then the answer is modified again to add references for attribution.

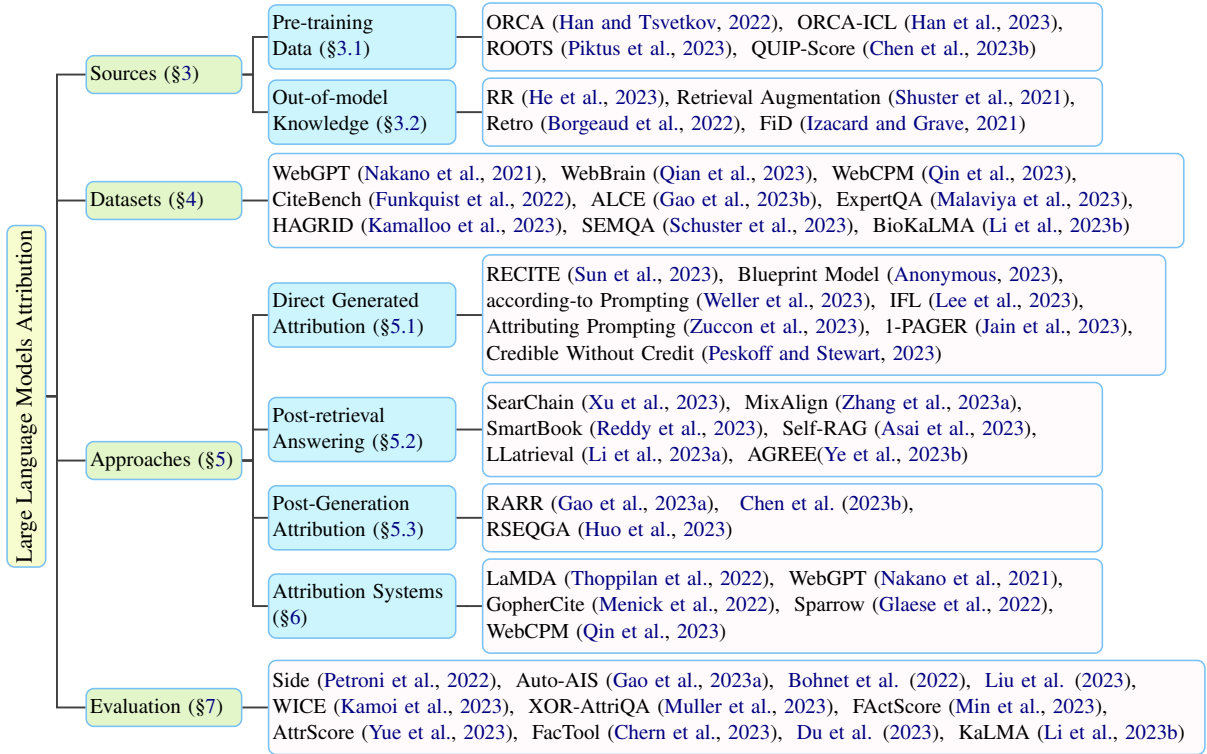


Figure 3: Taxonomy of large language models attribution.

The objective of the system is to produce an output \mathcal{S} , where \mathcal{S} is a set of n distinct statements: s_1, s_2, \dots, s_n . Each statement s_i is associated with a set of citations \mathcal{C}_i . This set \mathcal{C}_i is defined as $\mathcal{C}_i = \{c_{i,1}, c_{i,2}, \dots\}$, where each $c_{i,j}$ is a passage from the corpus \mathcal{D} . For practical applications, the output from LLMs can be segmented into individual statements using sentence boundaries. This approach is utilized because a single sentence typi-

cally encapsulates a coherent statement while maintaining brevity, facilitating easy verification. In terms of representation, citations may be enclosed within square brackets, for instance, [1][2]. It should be noted, however, that these citations can also be applied at the phrase level, rather than exclusively at the sentence level. It is important to highlight that the task configurations discussed in this paper are distinct from the generation of cita-

tion texts found in scholarly articles or wikipedia, where the citing and cited documents are usually used as inputs (Fetahu et al., 2016; Xing et al., 2020; Wu et al., 2021; Gu and Hahnloser, 2022).

3 Sources of Attribution

3.1 Pre-training Data

LLMs are typically trained on extensive corpora collected from various sources, predominantly the web. This vast amount of pre-training data forms the bedrock on which these models develop their understanding and capabilities. However, due to the scale of the data involved, manual inspection is often unfeasible, leading to potential inaccuracies, biases, and other undesirable artifacts in the data (Piktus et al., 2023). Despite these challenges, LLMs tend to perform well on a wide array of downstream tasks, even with little to no task-specific tuning. This performance hints at the ability of models to either memorize or reason through patterns present in the data. However, the specific patterns or the extent to which they are memorized or reasoned through, especially in different downstream tasks, remain somewhat elusive.

The concept of attribution in this context refers to tracing back the behavior of the model on a particular task to specific portions of the pre-training data (Han and Tsvetkov, 2022; Weller et al., 2023). By identifying a subset of pre-training data that significantly influences the model behavior on a downstream task, researchers aim to provide a clearer understanding of how the pre-training data impacts the model’s performance (Han et al., 2023). This kind of attribution is essential for interpreting the model, providing insights into whether the model is capturing task-relevant patterns or merely memorizing data. Furthermore, it aids in enhancing the trustworthiness of the model by offering a clearer picture of how the model operates and what sources of data significantly contribute to its performance. Through such attribution methodologies, researchers aim to bridge the understanding gap, offering a pathway towards better interpretability, trustworthiness, and eventually, the improvement of LLMs in handling various NLP tasks.

3.2 Out-of-model Knowledge

This source reveals methods to leverage out-of-model knowledge (e.g., web, knowledge graph) for attribution to enhance the capabilities of models (Shuster et al., 2021; Li et al., 2023b). Primary

among these methods is the retrieval-augmented generation technique (Lewis et al., 2020) which uses an encoder-decoder mechanism to encode questions and decode answers, augmented with documents or passages from extensive unstructured datasets. Furthermore, retrieval-enhanced language models are highlighted, which improve performance by fetching k -most similar training contexts or generating search queries to obtain relevant documents from external sources (Borgeaud et al., 2022). These methodologies, along with a mentioned post-processing method to utilize retrieved knowledge without additional training or fine-tuning, represent critical pathways for attributing LLM responses or generated text to external knowledge, aiming to make the outputs of LLMs verifiable external knowledge sources (Izacard and Grave, 2021; He et al., 2023).

4 Datasets for Attribution

As an information-seeking task, datasets for attribution are often built in the form of Question Answering (QA) or summarization (see Table 1). Several benchmarks are proposed based on existing QA datasets by proposing methods to evaluate the performance of attribution, as the golden citation annotation is not a necessity. Nakano et al. (2021) built a long-form QA dataset with web search results. After that Qin et al. (2023) built a similar Chinese dataset for the same purpose. However, these datasets are not directly built for verifying citations, but for factual accuracy. More recently, several works (Qian et al., 2023; Gao et al., 2023b; Kamalloo et al., 2023; Malaviya et al., 2023; Li et al., 2023b) focus on measuring and improving the accuracy of citations in generated text based on a given set of quotes, varying on question domain and citation granularity.

Question Domain. Most recent attribution datasets are designed for open-domain. However, ExpertQA (Malaviya et al., 2023) choose 32 domain-specific scenarios, some of which are high-stakes fields, and bring domain experts in the loop. BioKaLMA (Li et al., 2023b) focuses on biography domain for its practical application and convenient evaluation.

Attribution Granularity. There are two kinds of citation granularity in recent works: entity and sentence. The entity level attribution is more fine-grained, sentence level attribution requires citation for every completed sentence.

Among them, SEMQA (Schuster et al., 2023), ExpertQA (Malaviya et al., 2023) and BioKaLMA (Li et al., 2023b) make attribution at entity level, whereas other methods make attribution at sentence level.

5 Approaches to Attribution

5.1 Direct Generated Attribution

Attribution from parametric knowledge can help reduce hallucination and improve the truthfulness of generated text. By asking models to do self-detection and self-attribution, some works indicate that the generated texts are more grounded on facts and additionally improve performance on downstream tasks (Sun et al., 2023).

Recently, researchers found that large language models can not provide knowledge sources or evidence clearly when answering domain-specific knowledge-based questions (Peskov and Stewart, 2023; Zuccon et al., 2023; Gravel et al., 2023). In most cases, models can only provide a knowledge source that is loosely related to the keywords in questions or irrelevant to current topics. Even if the model answered the question correctly, the evidence it provided is still likely to have mistakes. Weller et al. (2023) tries to ground model’s generated text to its pre-training data by proposing according-to prompting, who finds the method can affect model’s groundedness and therefore affect performance on information-seeking tasks. Anonymous (2023) introduces an intermediate planning module, asking the model to generate a series of questions as blueprints to the current question. The model first proposes a blueprint and then combines the texts which are generated based on the blueprint questions as the final answer. The blueprint models allow for different forms of attribution during each question answering step, which can be expected to be more explainable.

5.2 Post-retrieval Answering

Numerous studies have delved into the post-retrieval answering strategy for attribution (Chen et al., 2017; Lee et al., 2019; Khattab and Zaharia, 2020). Reddy et al. (2023) introduces the SmartBook framework, which aims to generate structured situation reports incorporating factual evidence through rich links. The framework autonomously identifies crucial questions for situation analysis and extracts pertinent information to compose the report. Each question is addressed

with concise summaries containing tactical details of pertinent claims, supported by reliable and trustworthy factual evidence. To tackle the issue of misalignment between user queries and stored knowledge, where LLMs struggle to correlate questions with the appropriate grounding, MixAlign (Zhang et al., 2023a) presents a framework that combines automatic question-knowledge alignment with user clarifications. This approach effectively mitigates language model hallucination. To assess the adequacy of document support for an answer, LLa-trieval (Li et al., 2023a) updates the retrieval results until it confirms that the retrieved documents can sufficiently support the answer to the question. This iterative verification process significantly enhances the accuracy of the attribution by ensuring that the generated response is supported by verifiable evidence. Similarly, Self-RAG (Asai et al., 2023) trains an arbitrary language model to generate reflection-specific tokens after knowledge retrieval, thereby augmenting the attribution of retrieved passages. Furthermore, Search-in-the-chain (SearChain) (Xu et al., 2023) introduces a method to address the challenges posed by incorrect knowledge retrieved by information retrieval systems, which can mislead LLMs or disrupt their reasoning chains. It verifies and corrects answers within the global reasoning chain, known as Chain-of-Query (CoQ), while also identifying missing knowledge in CoQ. These operations significantly improve the attribution accuracy of LLMs in complex knowledge-intensive tasks, improving their reasoning ability and knowledge utilization.

5.3 Post-Generation Attribution

In order to facilitate accurate attribution without compromising the robust benefits offered by recent generation models, some research aims at attribution after generation, which employ search engines or document retrieval systems to search the evidence base on the input questions and generated answers. This approach allows researchers to assess or improve the factuality of answers without needing to access the model’s parameters directly. The post-generation attribution workflow is illustrated in Figure 4. RARR (Gao et al., 2023a) autonomously identifies the attribution of the output of any text generation model. It progressively verifies the factual consistency between the output and its source, and performs post-editing to rectify unsupported content, whilst striving to retain the original output to the greatest extent feasible. In the

Dataset	Domain	Source	Structure	Granularity	Response Source	#Questions
WebGPT (Nakano et al., 2021)	Open-domain	Web Pages	Unstructured	Sentence	GPT-3	19,578
WebBrain (Qian et al., 2023)	Open-domain	Wikipedia	Unstructured	Sentence	GPT-3	2.74M
WebCPM (Qin et al., 2023)	Open-domain	Web Pages	Unstructured	Sentence	Human	5,500
HAGRID (Kamalloo et al., 2023)	Open-domain	Wikipedia	Unstructured	Sentence	GPT-3.5, Human	1,922
ALCE (Gao et al., 2023b)	Open-domain	Wikipedia+Sphere	Unstructured	Sentence	Human	2,984
SEMQA (Schuster et al., 2023)	Open-domain	Wikipedia	Unstructured	Entity	Human	1,376
BioKaLMA (Li et al., 2023b)	Biography	Wikipedia	Structured	Entity	GPT-3.5, GPT4, LLaMA	1,085
ExpertQA (Malaviya et al., 2023)	Specific domains	Wikipedia	Unstructured	Entity	GPT-4, Human	2,507

Table 1: Comparison between different datasets for attribution.

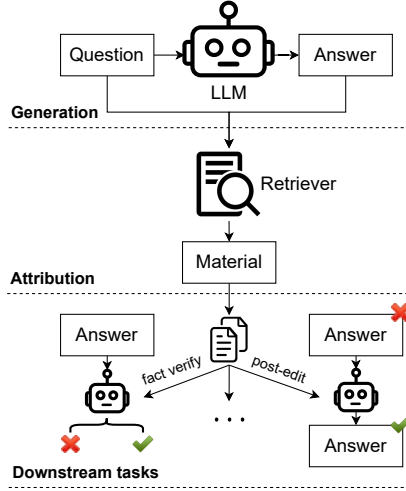


Figure 4: Workflow of post-generation attribution. Retrieval is performed after an answer being generated. The retrieved documents are used to perform citation and attribution, subsequently used to do fact verification and post-editing.

work of Huo et al. (2023), materials are retrieved from the corpus based on coarse-grained sentences or fine-grained factual statements. These retrieved materials are then utilized to prompt the LLM to verify the consistency between the generated responses and the retrieved material, and to make necessary edits to reduce the hallucinations. Chen et al. (2023b) introduces a fully automatic pipeline designed to verify complex political claims, which is achieved by retrieving evidence from the web. It breaks down each claim into subquestions and retrieves specific evidence for each, creating focused summaries and using them for claim verification. During training, the system evaluates its individual components based on the comprehensiveness and faithfulness.

6 Other Attribution Systems

Thoppilan et al. (2022) introduce LaMDA, a dialogue-focused language model. While enlarging the model improves its quality, it does not necessarily enhance safety and accuracy. By fine-

tuning LaMDA with annotated data and enabling it to access external knowledge, they significantly improve its safety and factual grounding. The grounding challenge of this study aims to generate responses based on credible external sources instead of merely plausible ones. The WebGPT model (Nakano et al., 2021) based on GPT-3 is trained to search and navigate the web and is fine-tuned for answering long-form questions in a web-browsing environment. For human evaluation of its factual accuracy, the model is required to gather references while browsing Microsoft Bing to support its answers. This ensures that the answers provided have a basis or attribution from credible web sources. Similarly, GopherCite (Menick et al., 2022) trained with reinforcement learning references evidence from multiple documents or a single user-provided document and refrains from answering when uncertain. Human evaluations show that GopherCite produces high-quality responses 80% at most. Nonetheless, citation alone is not a complete solution for ensuring safety and trustworthiness, as evidence-backed claims can still be false. Sparrow (Glaese et al., 2022) is trained to search the internet using Google Search to provide more accurate answers, allowing it to reference the latest information. In the user interface, evidence used by the model is displayed alongside its response, offering raters a means to validate the correctness of answer. To train the model in searching and using evidence, a preference model is used based on human judgments. Through human evaluation, it was found that responses with evidence were deemed plausible and supported 78% of the time. Comparisons between different systems are shown in Table 2.

7 Attribution Evaluation

Human Evaluation. To detect attribution errors, current attributed LLMs predominantly depend on human evaluation, a process that is both costly and time-intensive (Nakano et al., 2021; Kazemi

System	Model Training	Evidence Type	Citation Type	Integration
LaMDA (Thoppilan et al., 2022)	Multi-task SFT	Snippets	URLs	Appended
WebGPT (Nakano et al., 2021)	SFT + RL	Well-curated documents	Documents	Embedded
GopherCite (Menick et al., 2022)	SFT + RL	Long documents	Documents	Embedded
Sparrow (Glaese et al., 2022)	RL	Well-curated documents	Documents	Appended

Table 2: Features of different attribution systems. SFT means supervised fine-tuning, while RL means reinforcement learning optimization.

Evaluation Metrics	Evaluation Method	Description
Recall, Precision	Automatic, Statistics, Model-based	binary categorization based on NLI models
EM, BLEU, ROUGE	Automatic, Statistics	metrics for downstream tasks
QUIP-Score (Weller et al., 2023)	Automatic, Statistics	character-level n-gram metrics
Liu et al. (2023)	Human	fluency, perceived utility
AttrScore (Yue et al., 2023)	Human	attributability, extrapolatory, contradiction

Table 3: Comparison between different evaluation metrics for attribution.

et al., 2023; Chen et al., 2023a). For example, the typical cost of annotating a single (query, answer, reference) example stands at around \$1 (Liu et al., 2023). In practical applications of attributed LLMs, the responsibility falls on users to be cautious of attributions and to undertake manual verification, imposing a significant responsibility on them.

Categorization-Based Evaluation. For the sake of clarity, earlier research mainly employed binary categorization by repurposing other NLP tasks (e.g., natural language inference) to determine whether an answer is supported by a reference or not (attributable or not) (Rashkin et al., 2021; Bohnet et al., 2022; Gao et al., 2023b; Muller et al., 2023). Liu et al. (2023) carry out a human assessment to evaluate the veracity of responses from generative search engines, categorizing the degree of reference support into full, partial, or no support. Building on this, Yue et al. (2023) introduce a refined categorization of attribution: 1) attributable—where the reference entirely backs the generated statement; 2) extrapolatory—where the reference offers insufficient backing for the statement; and 3) contradictory—where the statement directly opposes the referenced citation.

Quantitative Evaluation Metrics. Assessment of attribution quality is approached from three distinct angles (Li et al., 2023b): 1) Correctness—evaluating the alignment of generated text with the provided sources; 2) Precision—measuring the percentage of generated attributions pertinent to the question at hand; and 3) Recall—assessing the scope to which generated attributions capture crucial knowledge. Moreover, the F1-Score is derived from the Precision and Recall metrics. Thoppilan et al. (2022) introduces citation accuracy as the frequency with

which the model refers to web sources for its assertions, excluding widely recognized truths. The QUIP-Score (Weller et al., 2023), an n-gram overlap metric, is founded on swift membership inquiries and evaluates the extent to which a section is comprised of exact spans within a text corpus.

As shown in Table 3, while human evaluations provide in-depth insights, their costly and time-consuming nature emphasizes the growing appeal for automated methods. Future research is expected to refine these methods, ensuring their practicality and reliability in real-world applications.

8 Discussion

8.1 Attribution Error Analysis

Attribution error has several forms. In this study, we systematically categorize these errors into three primary types, as outlined in Table 4, while acknowledging the possibility of other error types.

- **Granularity Error.** For ambiguous questions, the answer may involve multiple aspects. In this case, the retrieved multi-document may contain complex and diverse information. Thus the answer is complex and hybrid, leading to insufficient citation.
- **Mistaken Synthesis.** Models may mix up relationships between entities and events when several complex documents are provided. The citation should be faithful to the generated text and cite all the references.
- **Hallucinated Generation.** The reference documents may be irrelevant or not relevant to the question, or the model has conflicts between external documents and parameter knowledge.

511 The answer will be hallucinated and the cita-
512 tion is inaccurate.

513 8.2 Limitations of Attributions

514 Attribution in LLMs is fraught with inherent diffi-
515 culties. One primary challenge is discerning when
516 and how to attribute. Differentiating between gen-
517 eral knowledge, which may not require citations,
518 and specialized knowledge, which should ideally
519 be attributed, is a nuanced task. This gray area can
520 lead to inconsistencies in attribution (Huang and
521 Chang, 2023). And LLMs now do not have ability
522 to attribute parameter knowledge of itself (Litschko
523 et al., 2023). Another limitation is the potential in-
524 accuracy in attributions (Liu et al., 2023). LLMs
525 might link content to irrelevant or incorrect sources.
526 This misattribution can confuse users, leading them
527 wrong and affecting the reliability of the informa-
528 tion presented. For example, an LLM in the med-
529 ical field could wrongly associate faulty medical
530 guidance with a trustworthy reference, which might
531 guide users towards detrimental health choices.
532 Furthermore, the fluidity of knowledge means that
533 while some information remains static, other data
534 evolves and changes over time (Min et al., 2023).
535 Consequently, some attributions made by LLMs
536 may quickly become outdated, especially in rapidly
537 advancing domains, such as computer science and
538 clinical medicine. Additionally, we recommend
539 readers refer to §4.1 in Menick et al. (2022).

540 8.3 Challenges for Attributions

541 Despite the potential solutions on the horizon, im-
542 plementing these improvements for attributions is
543 laden with challenges.

544 One such challenge is excessive attribution or
545 over attribution (Huang and Chang, 2023; Liu et al.,
546 2023). If LLMs give credit too often, users might
547 get overwhelmed with too much information, con-
548 fusing them and making it difficult to tell what is
549 important and relevant from what is not.

550 At the same time, there is a real chance of LLMs
551 accidentally revealing private information. Finding
552 a balance between clear attribution and protecting
553 private details is a tricky task.

554 Bias is another big challenge. LLMs might unin-
555 tentionally lean towards some sources or kinds of
556 information, pushing certain views while ignoring
557 others. To tackle this bias, we need to use varied
558 training data and improve the methods used for
559 giving credit (Gunasekar et al., 2023).

560 Lastly, the shadow of incorrect information is
561 ever-present. Without solid validation measures,
562 LLMs could potentially spread wrong or mislead-
563 ing details, undermining the reliability of the in-
564 formation landscape. Future models should recog-
565 nize ambiguous references and refrain from making
566 statements when the evidence is not clear, instead
567 of presenting unfounded claims.

568 Overall, though LLMs seem to be on a posi-
569 tive path, they face many obstacles and doubts.
570 Proper credit is not just a side aspect; it is vital to
571 the growth, approval, and effectiveness of LLMs.
572 Guaranteeing correct and reliable credits, while
573 promoting new ideas, will definitely influence the
574 future of LLMs.

575 8.4 Future Directions for Attributions

576 **Continuous Refreshment of LLMs.** A promis-
577 ing direction for upcoming advancements is to cre-
578 ate a system that consistently refreshes the infor-
579 mation of LLMs (Thoppilan et al., 2022; Nakano
580 et al., 2021), akin to how search engines update
581 their databases. This approach not only ensures
582 up-to-date content for attribution but also offers a
583 platform for continuous learning and adaptation.

584 **Enhancing the Reliability of LLM Outputs.** An-
585 other pivotal direction entails boosting the trustwor-
586 thiness of LLM outputs. This can be achieved by
587 incorporating rigorous systems that assess the cred-
588 ibility and precision of the sources to which they
589 attribute information (Min et al., 2023). Ensuring
590 reliable and consistent sources will instill greater
591 confidence in users about the content generated.
592 As the adoption of LLMs expands across various
593 domains, the reliability of their output becomes
594 critical for informed decision making in various
595 sectors.

596 **Balancing Creativity with Proper Credit Attri-
597 bution.** LLMs are recognized for their creative
598 content generation. Striking a balance between this
599 inventive ability and proper credit-giving is a deli-
600 cate act that needs investigation. While creativity
601 is one of the significant strengths of LLMs, it is
602 vital to ensure that the generated content remains
603 trustworthy and rooted in factual bases. The aim is
604 to make sure LLMs acknowledge sources without
605 hindering their creative potential. Balancing these
606 two aspects can foster an environment where users
607 both benefit from the model and trust its outputs.

Limitation

While language models have the capability to cite their sources, undeniably enhance their utility, several limitations arise that need careful consideration (cf. Section 8.2). Our paper, in its current form, does not provide a solution to navigate such complex territory. It is important to address these limitations in future works and to continually educate users about the potential pitfalls of relying solely on machine-generated text.

References

- Ayush Agrawal, Lester Mackey, and Adam Tauman Kalai. 2023. [Do language models know when they’re hallucinating references?](#) *CoRR*, abs/2305.18248.
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Tachard Pasos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Z. Chen, Eric Chu, J. Clark, Laurent El Shafey, Yanping Huang, Kathleen S. Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan A. Botha, James Bradbury, Siddhartha Brahma, Kevin Michael Brooks, Michele Catasta, Yongzhou Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, C Crépy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, M. C. D’iaz, Nan Du, Ethan Dyer, Vladimir Feinberg, Fan Feng, Vlad Fienber, Markus Freitag, Xavier García, Sebastian Gehrmann, Lucas González, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, An Ren Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wen Hao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Mu-Li Li, Wei Li, Yaguang Li, Jun Yu Li, Hyeontaek Lim, Han Lin, Zhong-Zhong Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alexandra Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Marie Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniela Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Ke Xu, Yunhan Xu, Lin Wu Xue, Pengcheng Yin, Jiahui Yu, Qiaoling Zhang, Steven Zheng, Ce Zheng, Wei Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. [Palm 2 technical report](#). *ArXiv*, abs/2305.10403.

Anonymous. 2023. [Learning to plan and generate text with citations](#). In *Submitted to The Twelfth Inter-*

national Conference on Learning Representations. Under review.

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. [Self-rag: Learning to retrieve, generate, and critique through self-reflection](#). *CoRR*, abs/2310.11511.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, T. J. Henighan, Nicholas Joseph, Saurav Kadavath, John Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and Jared Kaplan. 2022. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#). *ArXiv*, abs/2204.05862.
- Bernd Bohnet, Vinh Q. Tran, Pat Verga, Roei Aharoni, Daniel Andor, Livio Baldini Soares, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, Tom Kwiatkowski, Ji Ma, Jianmo Ni, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan Das, Donald Metzler, Slav Petrov, and Kellie Webster. 2022. [Attributed question answering: Evaluation and modeling for attributed large language models](#). *CoRR*, abs/2212.08037.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from trillions of tokens](#). In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 2206–2240. PMLR.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 632–642. The Association for Computational Linguistics.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1870–1879. Association for Computational Linguistics.

722	Hung-Ting Chen, Fangyuan Xu, Shane A. Arora, and	gpt for medical questions . <i>Mayo Clinic Proceedings: Digital Health</i> , 1(3):226–234.	779
723	Eunsol Choi. 2023a. Understanding retrieval aug-		780
724	mentation for long-form question answering .		
725	Jifan Chen, Grace Kim, Aniruddh Sriram, Greg Dur-	Nianlong Gu and Richard H. R. Hahnloser. 2022.	781
726	rett, and Eunsol Choi. 2023b. Complex claim veri-	Controllable citation text generation . <i>CoRR</i> ,	782
727	fication with evidence retrieved in the wild . <i>CoRR</i> ,	abs/2211.07066.	783
728	abs/2305.11859.		
729	I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan,	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio	784
730	Kehua Feng, Chunting Zhou, Junxian He, Graham	César Teodoro Mendes, Allie Del Giorno, Sivakanth	785
731	Neubig, and Pengfei Liu. 2023. Factool: Factual-	Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo	786
732	ity detection in generative AI - A tool augmented	de Rosa, Olli Saarikivi, Adil Salim, Shital Shah,	787
733	framework for multi-task and multi-domain scenar-	Harkirat Singh Behl, Xin Wang, Sébastien Bubeck,	788
734	ios . <i>CoRR</i> , abs/2307.13528.	Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and	789
735		Yuanzhi Li. 2023. Textbooks are all you need . <i>CoRR</i> ,	790
736	Li Du, Yequan Wang, Xingrun Xing, Yiqun Ya, Xi-	abs/2306.11644.	791
737	ang Li, Xin Jiang, and Xuezhi Fang. 2023. Quan-	Zhijiang Guo, Michael Sejr Schlichtkrull, and Andreas	792
738	tifying and attributing the hallucination of large	Vlachos. 2022. A survey on automated fact-checking .	793
739	language models via association analysis . <i>CoRR</i> ,	<i>Trans. Assoc. Comput. Linguistics</i> , 10:178–206.	794
740	abs/2309.05217.		
741	Besnik Fetahu, Katja Markert, Wolfgang Nejdl, and	Xiaochuang Han, Daniel Simig, Todor Mihaylov, Yulia	795
742	Avishek Anand. 2016. Finding news citations for	Tsvetkov, Asli Celikyilmaz, and Tianlu Wang. 2023.	796
743	wikipedia . In <i>Proceedings of the 25th ACM Inter-</i>	Understanding in-context learning via supportive pre-	797
744	<i>national Conference on Information and Knowledge</i>	training data . In <i>Proceedings of the 61st Annual</i>	798
745	<i>Management, CIKM 2016, Indianapolis, IN, USA,</i>	<i>Meeting of the Association for Computational Lin-</i>	799
746	<i>October 24-28, 2016</i> , pages 337–346. ACM.	<i>guistics (Volume 1: Long Papers), ACL 2023, Toronto,</i>	800
747		<i>Canada, July 9-14, 2023</i> , pages 12660–12673. Asso-	801
748		ciation for Computational Linguistics.	802
749	Martin Funkquist, Ilia Kuznetsov, Yufang Hou, and	Xiaochuang Han and Yulia Tsvetkov. 2022. ORCA:	803
750	Iryna Gurevych. 2022. Citebench: A benchmark for	interpreting prompted language models via locating	804
751	scientific citation text generation .	supporting data evidence in the ocean of pretraining	805
752	Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony	data . <i>CoRR</i> , abs/2205.12600.	806
753	Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vin-		
754	cent Y. Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan,	Hangfeng He, Hongming Zhang, and Dan Roth. 2023.	807
755	and Kelvin Guu. 2023a. RARR: researching and	Rethinking with retrieval: Faithful large language	808
756	revising what language models say, using language	model inference . <i>CoRR</i> , abs/2301.00303.	809
757	models . In <i>Proceedings of the 61st Annual Meeting</i>		
758	<i>of the Association for Computational Linguistics (Vol-</i>	Jie Huang and Kevin Chen-Chuan Chang. 2023. Cita-	810
759	<i>ume 1: Long Papers), ACL 2023, Toronto, Canada,</i>	tion: A key to building responsible and accountable	811
760	<i>July 9-14, 2023</i> , pages 16477–16508. Association for	large language models . <i>CoRR</i> , abs/2307.02185.	812
761	Computational Linguistics.		
762	Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen.	Siqing Huo, Negar Arabzadeh, and Charles L. A. Clarke.	813
763	2023b. Enabling large language models to generate	2023. Retrieving supporting evidence for generative	814
764	text with citations . <i>CoRR</i> , abs/2305.14627.	question answering . <i>CoRR</i> , abs/2309.11392.	815
765			
766	Amelia Glaese, Nat McAleese, Maja Trebacz, John	Gautier Izacard and Edouard Grave. 2021. Leveraging	816
767	Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth	passage retrieval with generative models for open do-	817
768	Rauh, Laura Weidinger, Martin J. Chadwick, Phoebe	main question answering . In <i>Proceedings of the 16th</i>	818
769	Thacker, Lucy Campbell-Gillingham, Jonathan Ue-	<i>Conference of the European Chapter of the Associ-</i>	819
770	sato, Po-Sen Huang, Ramona Comanescu, Fan	<i>ation for Computational Linguistics: Main Volume,</i>	820
771	Yang, Abigail See, Sumanth Dathathri, Rory Greig,	<i>EACL 2021, Online, April 19 - 23, 2021</i> , pages 874–	821
772	Charlie Chen, Doug Fritz, Jaume Sanchez Elias,	880. Association for Computational Linguistics.	822
773	Richard Green, Sona Mokrá, Nicholas Fernando,		
774	Boxi Wu, Rachel Foley, Susannah Young, Iason	Alon Jacovi and Yoav Goldberg. 2020. Towards faith-	823
775	Gabriel, William Isaac, John Mellor, Demis Hass-	fully interpretable NLP systems: How should we	824
776	abis, Koray Kavukcuoglu, Lisa Anne Hendricks, and	define and evaluate faithfulness? In <i>Proceedings</i>	825
777	Geoffrey Irving. 2022. Improving alignment of dia-	<i>of the 58th Annual Meeting of the Association for</i>	826
778	logue agents via targeted human judgements . <i>CoRR</i> ,	<i>Computational Linguistics</i> , pages 4198–4205, On-	827
779	abs/2209.14375.	line. Association for Computational Linguistics.	828
780			
781	Jocelyn Gravel, Madeleine D’Amours-Gravel, and Esli	Palak Jain, Livio Baldini Soares, and Tom Kwiatkowski.	829
782	Osmanliu. 2023. Learning to fake it: Limited re-	2023. 1-pager: One pass answer generation and	830
783	sponses and fabricated references provided by chat-	evidence retrieval . <i>CoRR</i> , abs/2310.16568.	831
784			

832	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu,	<i>Thirty-Fourth Conference on Innovative Applications</i>	889
833	Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea	<i>of Artificial Intelligence, IAAI 2022, The Twelveth</i>	890
834	Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation . <i>ACM Comput. Surv.</i> , 55(12):248:1–248:38.	<i>Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 -</i>	891
835		<i>March 1, 2022, pages 10947–10955. AAAI Press.</i>	892
836			893
837	Ehsan Kamalloo, Aref Jafari, Xinyu Zhang, Nandan	Xiaonan Li, Changtai Zhu, Linyang Li, Zhangyue Yin,	894
838	Thakur, and Jimmy Lin. 2023. HAGRID: A human-	Tianxiang Sun, and Xipeng Qiu. 2023a. Llatrival: Llm-verified retrieval for verifiable generation . <i>arXiv preprint arXiv:2311.07838</i> .	895
839	llm collaborative dataset for generative information-		896
840	seeking with attribution. <i>arXiv:2307.16883</i> .		897
841	Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and	Xinze Li, Yixin Cao ² , Liangming Pan, Yubo Ma, and	898
842	Greg Durrett. 2023. Wice: Real-world entailment for claims in wikipedia . <i>CoRR</i> , abs/2303.01432.	Aixin Sun. 2023b. Towards verifiable generation: A benchmark for knowledge-aware language model attribution .	899
843			900
844	Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin		901
845	Xu, and Deepak Ramachandran. 2023. LAMBADA: Backward chaining for automated reasoning in natural language . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6547–6568, Toronto, Canada. Association for Computational Linguistics.	Robert Litschko, Max Müller-Eberstein, Rob van der Goot, Leon Weber, and Barbara Plank. 2023. Establishing trustworthiness: Rethinking tasks and model evaluation .	902
846			903
847			904
848			905
849		Frederick Liu and Besim Avci. 2019. Incorporating priors with feature attribution on text classification . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 6274–6283. Association for Computational Linguistics.	906
850			907
851			908
852	Omar Khattab and Matei Zaharia. 2020. Colbert: Efficient and effective passage search via contextualized late interaction over BERT . In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020</i> , pages 39–48. ACM.		909
853			910
854			911
855			912
856		Nelson F. Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating verifiability in generative search engines . <i>ArXiv</i> , abs/2304.09848.	913
857			914
858			915
859	Himabindu Lakkaraju, Dylan Slack, Yuxin Chen, Chen-	Chaitanya Malaviya, Subin Lee, Sihao Chen, Elizabeth Sieber, Mark Yatskar, and Dan Roth. 2023. Expertqa: Expert-curated questions and attributed answers . <i>ArXiv</i> , abs/2309.07852.	916
860	hao Tan, and Sameer Singh. 2022. Rethinking explainability as a dialogue: A practitioner’s perspective . <i>CoRR</i> , abs/2202.01875.		917
861			918
862			919
863	Dongyub Lee, Taesun Whang, Chanhee Lee, and	Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nathan McAleese. 2022. Teaching language models to support answers with verified quotes . <i>ArXiv</i> , abs/2203.11147.	920
864	Heuiseok Lim. 2023. Towards reliable and fluent large language models: Incorporating feedback learning loops in qa systems . <i>arXiv preprint arXiv:2309.06384</i> .		921
865			922
866			923
867			924
868	Kenton Lee, Ming-Wei Chang, and Kristina Toutanova.	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation . <i>CoRR</i> , abs/2305.14251.	925
869	2019. Latent retrieval for weakly supervised open domain question answering . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 6086–6096. Association for Computational Linguistics.		926
870			927
871			928
872			929
873			930
874			931
875	Patrick S. H. 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 <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	Benjamin Muller, John Wieting, Jonathan H. Clark, Tom Kwiatkowski, Sebastian Ruder, Livio Baldini Soares, Roei Aharoni, Jonathan Herzig, and Xinyi Wang. 2023. Evaluating and modeling attribution for cross-lingual question answering . <i>CoRR</i> , abs/2305.14332.	932
876			933
877			934
878			935
879			936
880		Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback . <i>arXiv preprint arXiv:2112.09332</i> .	937
881			938
882			939
883			940
884	Dongfang Li, Baotian Hu, Qingcai Chen, Tujie Xu, Jingcong Tao, and Yunan Zhang. 2022. Unifying model explainability and robustness for joint text classification and rationale extraction . In <i>Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022</i> ,	OpenAI. 2023. Gpt-4 technical report . <i>ArXiv</i> , abs/2303.08774.	941
885			942
886			943
887			944
888			

945	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida,	Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm,	1001
946	Carroll L. Wainwright, Pamela Mishkin, Chong	Michael Collins, Dipanjan Das, Slav Petrov, Gau-	1002
947	Zhang, Sandhini Agarwal, Katarina Slama, Alex	rav Singh Tomar, Iulia Turc, and David Reitter. 2021.	1003
948	Ray, John Schulman, Jacob Hilton, Fraser Kelton,	Measuring attribution in natural language generation	1004
949	Luke E. Miller, Maddie Simens, Amanda Askill, Pe-	models . <i>CoRR</i> , abs/2112.12870.	1005
950	ter Welinder, Paul Francis Christiano, Jan Leike, and		
951	Ryan J. Lowe. 2022. Training language models to	Vipula Rawte, A. Sheth, and Amitava Das. 2023. A	1006
952	follow instructions with human feedback . <i>ArXiv</i> ,	survey of hallucination in large foundation models .	1007
953	abs/2203.02155.	<i>ArXiv</i> , abs/2309.05922.	1008
954	Lawrence Page, Sergey Brin, Rajeev Motwani, and	Revanth Gangi Reddy, Yi R. Fung, Qi Zeng, Manling	1009
955	Terry Winograd. 1999. The pagerank citation ranking	Li, Ziqi Wang, Paul Sullivan, and Heng Ji. 2023.	1010
956	: Bringing order to the web . In <i>The Web Conference</i> .	Smartbook: Ai-assisted situation report generation .	1011
		<i>CoRR</i> , abs/2303.14337.	1012
957	Guilherme Penedo, Quentin Malartic, Daniel Hess-	Tal Schuster, Ádám D. Lelkes, Haitian Sun, Jai Gupta,	1013
958	low, Ruxandra-Aimée Cojocaru, Alessandro Cap-	Jonathan Berant, William W. Cohen, and Donald Met-	1014
959	pelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam	zler. 2023. SEMQA: semi-extractive multi-source	1015
960	Almazrouei, and Julien Launay. 2023. The refined-	question answering . <i>CoRR</i> , abs/2311.04886.	1016
961	web dataset for falcon llm: Outperforming curated		
962	corpora with web data, and web data only . <i>ArXiv</i> ,	Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela,	1017
963	abs/2306.01116.	and Jason Weston. 2021. Retrieval augmentation	1018
		reduces hallucination in conversation . In <i>Findings</i>	1019
964	Denis Peskoff and Brandon Stewart. 2023. Credible	of the Association for Computational Linguistics: EMNLP 2021, Virtual Event / Punta Cana, Domini-	1020
965	without credit: Domain experts assess generative	can Republic, 16-20 November, 2021 , pages 3784–	1021
966	language models . In <i>Proceedings of the 61st An-</i>	3803. Association for Computational Linguistics.	1022
967	<i>ual Meeting of the Association for Computational</i>		1023
968	<i>Linguistics (Volume 2: Short Papers)</i> , <i>ACL 2023,</i>	Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and	1024
969	<i>Toronto, Canada, July 9-14, 2023</i> , pages 427–438.	Denny Zhou. 2023. Recitation-augmented language	1025
970	Association for Computational Linguistics.	models . In <i>The Eleventh International Conference</i>	1026
		on Learning Representations, ICLR 2023, Kigali,	1027
971	Fabio Petroni, Samuel Broscheit, Aleksandra Piktus,	<i>Rwanda, May 1-5, 2023</i> . OpenReview.net.	1028
972	Patrick S. H. Lewis, Gautier Izacard, Lucas Hosseini,		
973	Jane Dwivedi-Yu, Maria Lomeli, Timo Schick, Pierre-	Yi Tay, Vinh Tran, Mostafa Dehghani, Jianmo Ni, Dara	1029
974	Emmanuel Mazaré, Armand Joulin, Edouard Grave,	Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao,	1030
975	and Sebastian Riedel. 2022. Improving wikipedia	Jai Prakash Gupta, Tal Schuster, William W. Cohen,	1031
976	verifiability with AI . <i>CoRR</i> , abs/2207.06220.	and Donald Metzler. 2022. Transformer memory as	1032
		a differentiable search index . In <i>NeurIPS</i> .	1033
977	Aleksandra Piktus, Christopher Akiki, Paulo Villegas,	Romal Thoppilan, Daniel De Freitas, Jamie Hall,	1034
978	Hugo Laurençon, Gérard Dupont, Sasha Luccioni,	Noam M. Shazeer, Apoorv Kulshreshtha, Heng-Tze	1035
979	Yacine Jernite, and Anna Rogers. 2023. The ROOTS	Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du,	1036
980	search tool: Data transparency for llms . In <i>Proceed-</i>	Yaguang Li, Hongrae Lee, Huaixiu Steven Zheng,	1037
981	<i>ings of the 61st Annual Meeting of the Association for</i>	Amin Ghafouri, Marcelo Menegali, Yanping Huang,	1038
982	<i>Computational Linguistics: System Demonstrations,</i>	Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao	1039
983	<i>ACL 2023, Toronto, Canada, July 10-12, 2023</i> , pages	Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts,	1040
984	304–314. Association for Computational Linguistics.	Maarten Bosma, Yanqi Zhou, Chung-Ching Chang,	1041
		I. A. Krivokon, Willard James Rusch, Marc Pick-	1042
985	Hongjing Qian, Yutao Zhu, Zhicheng Dou, Haoqi Gu,	ett, Kathleen S. Meier-Hellstern, Meredith Ringel	1043
986	Xinyu Zhang, Zheng Liu, Ruofei Lai, Zhao Cao,	Morris, Tulsee Doshi, Renelito Delos Santos, Toju	1044
987	Jian-Yun Nie, and Ji-Rong Wen. 2023. Webbrain:	Duke, Johnny Hartz Søraker, Ben Zevenbergen, Vin-	1045
988	Learning to generate factually correct articles for	odkumar Prabhakaran, Mark Díaz, Ben Hutchinson,	1046
989	queries by grounding on large web corpus . <i>CoRR</i> ,	Kristen Olson, Alejandra Molina, Erin Hoffman-	1047
990	abs/2304.04358.	John, Josh Lee, Lora Aroyo, Ravindran Rajakumar,	1048
		Alena Butryna, Matthew Lamm, V. O. Kuzmina,	1049
991	Yujia Qin, Zihan Cai, Dian Jin, Lan Yan, Shihao	Joseph Fenton, Aaron Cohen, Rachel Bernstein, Ray	1050
992	Liang, Kunlun Zhu, Yankai Lin, Xu Han, Ning Ding,	Kurzweil, Blaise Agueria-Arcas, Claire Cui, Mar-	1051
993	Huadong Wang, Ruobing Xie, Fanchao Qi, Zhiyuan	rian Rogers Croak, Ed Huai hsin Chi, and Quoc Le.	1052
994	Liu, Maosong Sun, and Jie Zhou. 2023. WebCPM:	2022. Lamda: Language models for dialog applica-	1053
995	Interactive web search for Chinese long-form ques-	tions . <i>ArXiv</i> , abs/2201.08239.	1054
996	tion answering . In <i>Proceedings of the 61st Annual</i>		
997	<i>Meeting of the Association for Computational Lin-</i>	James Thorne, Andreas Vlachos, Christos	1055
998	<i>guistics (Volume 1: Long Papers)</i> , pages 8968–8988,	Christodoulopoulos, and Arpit Mittal. 2018.	1056
999	Toronto, Canada. Association for Computational Lin-	FEVER: a large-scale dataset for fact extraction	1057
1000	guistics.		

1058	and verification. In <i>Proceedings of the 2018</i>	in the ai ocean: A survey on hallucination in large	1113
1059	<i>Conference of the North American Chapter of the</i>	language models. <i>ArXiv</i> , abs/2309.01219.	1114
1060	<i>Association for Computational Linguistics: Human</i>		
1061	<i>Language Technologies, NAACL-HLT 2018, New</i>	Guido Zuccon, Bevan Koopman, and Razia Shaik. 2023.	1115
1062	<i>Orleans, Louisiana, USA, June 1-6, 2018, Volume</i>	Chatgpt hallucinates when attributing answers.	1116
1063	<i>1 (Long Papers)</i> , pages 809–819. Association for		
1064	Computational Linguistics.		
1065	Haoran Wang and Kai Shu. 2023. Explainable claim		
1066	verification via knowledge-grounded reasoning with		
1067	large language models.		
1068	Orion Weller, Marc Marone, Nathaniel Weir, Dawn		
1069	Lawrie, Daniel Khashabi, and Benjamin Van Durme.		
1070	2023. " according to..." prompting language mod-		
1071	els improves quoting from pre-training data. <i>arXiv</i>		
1072	<i>preprint arXiv:2305.13252</i> .		
1073	Jia-Yan Wu, Alexander Te-Wei Shieh, Shih-Ju Hsu, and		
1074	Yun-Nung Chen. 2021. Towards generating citation		
1075	sentences for multiple references with intent control.		
1076	<i>CoRR</i> , abs/2112.01332.		
1077	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and		
1078	Yu Su. 2023. Adaptive chameleon or stubborn sloth:		
1079	Unraveling the behavior of large language models in		
1080	knowledge clashes. <i>ArXiv</i> , abs/2305.13300.		
1081	Xinyu Xing, Xiaosheng Fan, and Xiaojun Wan. 2020.		
1082	Automatic generation of citation texts in scholarly		
1083	papers: A pilot study. In <i>Proceedings of the 58th</i>		
1084	<i>Annual Meeting of the Association for Computa-</i>		
1085	<i>tional Linguistics, ACL 2020, Online, July 5-10, 2020</i> ,		
1086	pages 6181–6190. Association for Computational		
1087	Linguistics.		
1088	Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng,		
1089	and Tat-Seng Chua. 2023. Search-in-the-chain: To-		
1090	wards the accurate, credible and traceable content		
1091	generation for complex knowledge-intensive tasks.		
1092	<i>CoRR</i> , abs/2304.14732.		
1093	Hongbin Ye, Tong Liu, Aijia Zhang, Wei Hua, and		
1094	Weiqliang Jia. 2023a. Cognitive mirage: A review		
1095	of hallucinations in large language models. <i>ArXiv</i> ,		
1096	abs/2309.06794.		
1097	Xi Ye, Ruoxi Sun, Serkan Ö Arik, and Tomas Pfister.		
1098	2023b. Effective large language model adap-		
1099	tation for improved grounding. <i>arXiv preprint</i>		
1100	<i>arXiv:2311.09533</i> .		
1101	Xiang Yue, Boshi Wang, Kai Zhang, Ziru Chen,		
1102	Yu Su, and Huan Sun. 2023. Automatic evalua-		
1103	tion of attribution by large language models. <i>CoRR</i> ,		
1104	abs/2305.06311.		
1105	Shuo Zhang, Liangming Pan, Junzhou Zhao, and		
1106	William Yang Wang. 2023a. Mitigating lan-		
1107	guage model hallucination with interactive question-		
1108	knowledge alignment. <i>CoRR</i> , abs/2305.13669.		
1109	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,		
1110	Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,		
1111	Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei		
1112	Bi, Freda Shi, and Shuming Shi. 2023b. Siren's song		

A Attribution Before the Era of Large Language Model

A.1 Related Natural Language Processing Tasks

The relationship between attribution tasks and other Natural Language Processing (NLP) tasks manifests in the overarching goal of understanding, evaluating, and leveraging the content retrieved or generated in response to particular stimuli such as questions or claims. Here is an exploration of how attribution tasks are intertwined with other NLP tasks, anchored on the retrieval of related content: **Open-domain question answering:** Both tasks hinge on retrieving pertinent documents or information to address a posed question or claim. While open-domain QA zeroes in on the accuracy and relevance of the answer, attribution tasks scrutinize whether the answer or generated text can be accurately traced back to the retrieved documents (Chen et al., 2017; Bohnet et al., 2022).

Fact-checking & Claim verification (subtask of fact checking): Fact-checking and attribution tasks both necessitate the retrieval of external evidence to validate a claim or generated text. The emphasis in fact-checking is on verifying the truthfulness of a claim, whereas attribution tasks focus on the correct attribution of generated text to the sourced evidence (Thorne et al., 2018). On the other hand, attribution tasks and claim verification both center around validating information against reference or sourced material, yet they serve different purposes. Attribution ensures that generated text or answers accurately reflect the provided references, while claim verification assesses the truthfulness of a claim based on evidence or source material (Guo et al., 2022). Both tasks necessitate the retrieval of related content for verification (Wang and Shu, 2023), making them inherently reliant on the accuracy and relevance of the retrieved material. claim verification pivotal in fact-checking and misinformation detection, they share the fundamental objective of endorsing the accuracy and trustworthiness of information by juxtaposing it against a reference. **Natural Language Inference (NLI):** Both tasks engage in evaluating the relationship between two snippets of text; however, NLI concentrates on logical entailment, contradiction, or neutrality, while attribution evaluates the substantiation provided by references for generated text (Bowman et al., 2015).

Summarization: Summarization and attribution

tasks both generate condensed or altered text and necessitate a check on the fidelity of the generated text to the original or sourced content. Attribution in summarization is pivotal to averting hallucinations (generation of false or unsupported information) and ensuring the summary accurately mirrors the input text (Ji et al., 2023).

The commonality among these tasks lies in the requisite to retrieve, analyze, and validate content against some form of reference material, be it external evidence, retrieved documents, or a different segment of text. The capacity to retrieve related content forms a cornerstone for these tasks, enabling the necessary comparisons and evaluations to ascertain accuracy, relevance, and correct attribution.

A.2 Interpretability of NLP Models

Interpretability (e.g., feature attribution) dives into understanding which parts of the input (e.g., words or phrases) are crucial for a model’s decision or output (Liu and Avci, 2019; Li et al., 2022). It helps in identifying the importance of different features in the input data concerning the model’s performance. Compared to feature attribution, explicit attribution for LLMs serves as a conduit to trace the sources of the information they generate, which is pivotal for accountability, especially in critical domains like healthcare or finance. It enables verifiability, allowing users or other systems to check the accuracy and reliability of the information provided. Trustworthiness is also fostered through explicit attribution, as users are more likely to trust the model if they know where the information is coming from. Additionally, it plays a role in interpretability, aiding users in understanding how the model arrives at certain conclusions by revealing the sources of information. This alignment with interpretability objectives helps in making the model’s decision-making process more transparent and comprehensible.

Attribution and interpretability, though interconnected, serve distinct purposes. Attribution specifically refers to the process of tracing back the generated information or decisions of a model to its source material or input features, providing a clear reference or basis for the output. On the other hand, interpretability is a broader concept encompassing the understanding of how a model processes input data to arrive at a particular output (Lakkaraju et al., 2022), making the inner workings of the model transparent and comprehensible to users. While

1219 attribution can be seen as a component or a spe-
1220 cific form of interpretability, aiding in understand-
1221 ing and trusting the model’s outputs by providing
1222 source references, interpretability dives deeper into
1223 elucidating the model’s behavior, decision-making
1224 process, and the significance of different input fea-
1225 tures in those decisions, thus fostering a compre-
1226 hensive understanding and trust in the model. Note
1227 that the evaluation of attribution focuses on the ver-
1228 ifiability of the answer based on the reference, sim-
1229 ilar to *faithfulness* evaluation of explanations (Ja-
1230 covi and Goldberg, 2020). It does not measure the
1231 relevance about whether the answer correctly re-
1232 sponds to the query, similar to *plausibility* in Jacovi
1233 and Goldberg (2020).

Granularity Error	<p>Question: Where was the biggest gold nugget ever found?</p> <p>References: [doc1] The largest gold nugget ever found in California was found in August 1869 in Sierra Buttes by five partners. [doc2] Considered by most authorities to be the biggest gold nugget ever found, the Welcome Stranger was found at Moliagul, Victoria, Australia in 1869. [doc3] Pepita Canaã, was found on September 13, 1983 by miners at the Serra Pelada Mine in the State of Para, Brazil.</p> <p>Answer: Most authorities consider the Welcome Stranger to be the biggest gold nugget ever found, found in 1869 in Moliagul, Victoria, Australia[2]. The biggest gold nugget found in Brazil was found in the State of Para and in United States was found in Sierra Buttes [3].</p> <p>Explanation: The second sentence in generated answer has two facets(the biggest nugget in Barzil, and the biggest nugget in United States). However, only document 3 is cited.</p>
Mistaken Synthesis	<p>Question: What's the difference in the portrayal of Ares in Roman mythology and Greek mythology?</p> <p>References: [doc1] In Roman mythology, he is called Mars. The Latin name for Mars derives from his Roman name; "Tuesday" comes from the corresponding deity Tyr in Norse mythology...[doc2] Ares is the god of war in ancient Greek mythology, one of the twelve Olympian gods. He is the son of Zeus and Hera, or alternatively, he was born after Juno (Hera's Roman name) swallowed a violent-eyed giant snake...</p> <p>Answer: In Roman mythology, Ares is called Mars, one of the twelve main gods of Rome, representing the symbols of power and war[1]. In Greek mythology, Ares is also one of the twelve Roman gods, representing the symbols of power and war[2].</p> <p>Explanation: The answer mistakenly mixes up Olympian gods and Roman mythology thus both document 1 and document 2 should be treat as referenced documents.</p>
Hallucinated Generation	<p>Question: How many times have the lakers won the finals?</p> <p>References: [doc1] ... won six times in the 1960s, three times in the 1980s, in 2008, and in 2010... [doc2] ...and Lakers met in the Finals for the first time since 1987, with the Celtics winning the series 4-2. They faced off once again in the 2010 NBA Finals, which the Lakers won in 7 games...</p> <p>Answer: The Los Angeles Lakers have won the NBA Finals five times, in 1980, 1982, 1985, 1987, and 2010[1].</p> <p>Explanation: The answer is incorrect as the generation is not grounded on reference documents. Hallucination may happen because of the conflict between reference documents and parameter knowledge, thus the citation will result in failure.</p>

Table 4: List of different attribution errors types and example instance.