

DEEPTRACE: AUDITING DEEP RESEARCH AI SYSTEMS FOR TRACKING RELIABILITY ACROSS CITATIONS AND EVIDENCE

Pranav Narayanan Venkit

Salesforce AI Research
Palo Alto, CA 94301, USA
pnarayananvenkit@salesforce.com

Philippe Laban

Microsoft Research
New York City, NY - 10012, USA
plaban@microsoft.com

Yilun Zhou

Salesforce AI Research
Palo Alto, CA 94301, USA
yilun.zhou@salesforce.com

Kung-Hsiang Huang

Salesforce AI Research
Palo Alto, CA 94301, USA
kh.huang@salesforce.com

Yixin Mao

Salesforce AI Research
Palo Alto, CA 94301, USA
y.mao@salesforce.com

Chien-Sheng Wu

Salesforce AI Research
Palo Alto, CA 94301, USA
wu.jason@salesforce.com



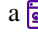
ABSTRACT

Generative search engines and deep research LLM agents promise trustworthy, source-grounded synthesis, yet users regularly encounter overconfidence, weak sourcing, and confusing citation practices. We introduce **DeepTRACE**, a novel sociotechnically grounded audit framework that turns prior community-identified failure cases into eight measurable dimensions spanning answer text, sources, and citations. DeepTRACE uses statement-level analysis (decomposition, confidence scoring) and builds citation and factual-support matrices to audit how systems reason with and attribute evidence end-to-end. Using automated extraction pipelines for popular public models (e.g., GPT-4.5/5, You.com, Perplexity, Copilot/Bing, Gemini) and an LLM-judge with validated agreement to human raters, we evaluate both web-search engines and deep-research configurations. Our findings show that generative search engines and deep research agents frequently produce one-sided, highly confident responses on debate queries and include large fractions of statements unsupported by their own listed sources. Deep-research configurations reduce overconfidence and can attain high citation thoroughness, but they remain highly one-sided on debate queries and still exhibit large fractions of unsupported statements, with citation accuracy ranging from 40–80% across systems. Unlike prior factuality and citation metrics that focus on claim correctness or academic summarization, DeepTRACE audits end-to-end GSE/DR behavior, including citation necessity, unsupported-statement rates, and URL-level citation structure.

1 INTRODUCTION

Large language models (LLMs) have recently become part of daily life for many, with the models offering AI-based conversational assistance to hundreds of millions of users with informational retrieval and text generation features (Ferrara, 2024; Pulapaka et al., 2024). In doing so, such systems have graduated from purely research-based systems to *public sociotechnical tools* (Cooper & Foster, 1971) that now impact both technical and social elements.

With the current text generation models growing capabilities, these systems are evolving from serving purely generative operations to functioning as “Generative Search Engines” capable of synthesizing information retrieved from external sources. These systems are now designed to autonomously

conduct in-depth research on complex topics by exploring the web, synthesizing information, and generating comprehensive reports *with citations*. These systems are therefore now dubbed a generative search engine (**GSE**) or a deep research agents (**DR**). A generative search engine summarizes and presents retrieved information, whereas a deep research agent executes in multi-step reasoning to derive insights resulting in a long-form report. These deep research agents first retrieve relevant  **source documents** that likely contain answer elements to the user’s questions or request, using a retrieval system (which can be a traditional search engine). The model then composes a textual prompt that contains the user’s query, and the retrieved sources, and instructs an LLM to generate a long and self-contained **T** **answer** based on the users preference and content of the sources. Importantly,  **citations** are inserted into the answer, with each citation linking to the sources that support each statement within the answer. This citation-enriched answer is provided to the user in a  **user interface** with a click on a citation allowing the user to navigate to the source or sources that support any statement. These systems, therefore, are intended to go beyond simple search and text generation to provide detailed analysis and structured outputs, often resembling human-written research papers.

In essence, the GSE and deep research pipeline promise a streamlining of a user’s information-seeking journey (Shah & Bender, 2024). The deep research agents are sold with the premise of concisely summarize the information the user is looking for, and sources remain within a click in case the user desires to deepen their understanding or verify the information’s veracity. Recently, several free deep research agents have become popular such as Perplexity.ai and You Chat, with some reporting millions of daily searches performed by their users (Narayanan Venkit et al., 2025).

Despite their advertised promise, deep research pipelines built on LLMs suffer from several critical limitations across their constituent components. First, LLMs are prone to hallucination and struggle to identify factual fallacies even when provided with authoritative sources (Venkit et al., 2024; Huang et al., 2023). Second, research has shown that the retrieval component of the models often fails to produce accurate citations within their responses (Liu et al., 2023), sometimes attributing claims to irrelevant or non-existent sources. Third, LLMs encode knowledge in their internal weights during pretraining, making it difficult to ensure that generated outputs rely solely on the user-provided documents or retrieved documents (Kaur et al., 2024). Finally, these systems can exhibit sycophantic behavior whereby they favor agreement with the user’s implied perspective over adherence to objective facts (Sharma et al., 2024; Laban et al., 2023b). These limitations have real implications for the quality, reliability, and trustworthiness of DR agents.

Yet, there remains a significant gap to evaluate and audit these models as a whole. Existing benchmarks largely focus on isolated components, such as the retrieval or summarization stages of Retrieval-Augmented Generation, with limited attention to how well systems ground responses in retrieved sources, generate citations, or manage uncertainty. To effectively address this gap, we build on the findings of Narayanan Venkit et al. (2025) and Sharma et al. (2024), who conducted an audit-focused usability study of deep research agents. The study participants identified **16 common failure cases** and proposed **actionable design recommendations** grounded in real-world use. In this work, we extend that foundation by transforming those usercentric insights into an automated evaluation benchmark. Our goal is to provide a systematic framework for auditing the end-to-end performance of deep research agents, capturing what these systems generate and how they reason, cite, and interact with knowledge in context. Our **DeepTrace** framework adopts a community-centered approach by focusing on the failure cases identified through community-driven evaluation, enabling benchmarking of models on real-world, practitioner-relevant weaknesses.

Our evaluation shows three findings that hold across GSEs and deep-research agents. First, public GSEs frequently produce one-sided and overconfident responses to debate-style queries. In our corpus, we observe high rates of one-sidedness and very confident language, indicating a tendency to present charged prompts as settled facts. Second, despite retrieval and citation, a large share of generated statements remains unsupported by the systems’ own sources, and citation practice is uneven. Third, systems that list many links often leave them uncited, creating a false impression of validation. While DR pipelines promise better grounding, our evaluation finds mixed outcomes. DR systems lowers overconfidence relative to GSE modes and increase citation thoroughness for some models, yet they are still one-sided for a majority of debate queries (e.g., GPT-5(DR) 54.7%; YouChat(DR) 63.1%; Copilot(DR) 94.8%). Additionally, unsupported statement rates remain high for several DR engines (YouChat(DR) 74.6%; PPLX(DR) 97.5%) and citation accuracy is well

below perfect (40–80%). Listing more sources does not guarantee better grounding, leaving users to experience search fatigue. Our work complements hallucination and factuality metrics such as FActScore and CoRE Min et al. (2023); Jiang et al. (2025) by shifting the focus from isolated claim correctness to how GSE/DR systems use retrieval, structure citations, and express confidence in user-facing answers. Similarly, it complements survey-style citation evaluations such as AutoSurvey Wang et al. (2024) by targeting open-web, end-to-end systems rather than academic summarization. While DeepTRACE concentrates on sourcing and traceability, we discuss how it can be extended with answer completeness, coherence, and synthesis quality in future work.

2 RELATED WORKS

2.1 EVOLUTION OF DEEP RESEARCH SYSTEMS

LLMs are increasingly embedded in sociotechnical settings that shape how people access and interact with information (Züger & Asghari, 2023; Narayanan Venkit, 2023). As these models transition from only research-based demonstrations to public-facing tools, their impact extends beyond technical performance into social, epistemic, and political domains (Dolata et al., 2022; Cooper & Foster, 1971). This shift has catalyzed the development of what are increasingly called generative search engines or deep research agents.

Unlike traditional RAG systems (Lewis et al., 2020; Izacard & Grave, 2021), which operate on static pipelines, deep research agents emphasize dynamic, iterative workflows. As defined by Huang et al. (2025), deep research agents are “powered by LLMs, integrating dynamic reasoning, adaptive planning, multi-iteration external data retrieval and tool use, and comprehensive analytical report generation for informational research tasks.” This framing situates such systems as more than just passive tools, they are positioned as active collaborators in knowledge production. These systems are designed to handle open-ended, multi-hop, and real-time queries by combining LLMs with external tools for search, planning, and reasoning (Nakano et al., 2021; Yao et al., 2023).

Recent research has explored architectures and frameworks that enhance the capabilities of deep research agents. For example, the MindMap Agent (Wu et al., 2025) constructs knowledge graphs to track logical relationships among retrieved content, enabling more coherent and deductive reasoning on tasks such as PhD-level exam questions. The MLGym framework (Nathani et al., 2025) demonstrates how LLM-based agents can simulate research workflows, including hypothesis generation, experimental design, and model evaluation. Similarly, DeepResearcher (Zheng et al., 2025) employs reinforcement learning with human feedback to train agents in web-based environments, improving both factuality and relevance of the final output in information-seeking tasks. With web browsing enabled, these research-oriented agents are mirrored in commercial deep research models such as Bing Copilot, Perplexity AI, YouChat, and ChatGPT (Narayanan Venkit et al., 2025). These systems advertise real-time retrieval, citation generation, and structured synthesis of sources.

2.2 BEYOND A POSITIVISM AND TECHNICAL LENS OF EVALUATION

A GSE and deep research agents gain traction in the NLP and AI communities, there has been a growing interest in evaluating their performance (Jeong et al., 2024; Wu et al., 2024; Es et al., 2023; Zhu et al., 2024). However, existing frameworks and benchmarks have largely maintained a technocentric orientation prioritizing model-centric metrics while underexploring the social and human-centered consequences of deploying these systems at scale. This trend reflects what Wylie (2014) describe as a positivist approach to technology: one that assumes universal evaluative truths through formal metrics, often abstracted from real-world user interactions.

Among the most prominent efforts is RAGAS (Es et al., 2023; 2024), which assesses answer quality through metrics such as faithfulness, context relevance, and answer helpfulness, without requiring human ground truth annotations. Similarly, ClashEval (Wu et al., 2024) reveals how LLMs may override correct prior knowledge with incorrect retrieved content more than 60% of the time. Although these evaluations are informative, they still treat language models as isolated computational systems, rather than sociotechnical agents embedded within user-facing applications. More recent work has begun to explore the application of RAG systems in socially sensitive domains. For instance, adaptations for medicine and journalism have involved integrating domain-specific knowledge bases to reduce hallucination and increase trust (Siriwardhana et al., 2023). Similar domain-

focused RAG evaluations have emerged in telecommunications (Roychowdhury et al., 2024), agriculture (Gupta et al., 2024), and gaming (Chauhan et al., 2024), reflecting an effort to align model behavior with contextual needs.

In the context of deep research agents, DeepResearch Bench (Du et al., 2025) evaluates LLM agents on 100 PhD-level research tasks using dimensions like comprehensiveness, insightfulness, readability, and citation correctness. DRBench (Bosse et al., 2025) similarly introduces 89 complex multi-step research tasks and proposes RetroSearch, a simulated web environment to measure model planning and execution. Similarly, BrowseComp-Plus (Chen et al., 2025) employs a static 100,000 web document as their corpus to evaluate accuracy, recall, number of search of a deep research agent. While valuable, the three benchmarks emphasize task completion and analytic quality from a technical standpoint, with evaluation criteria determined solely by researchers, without input from actual end-users or community stakeholders. This gap motivates our work. Inspired by calls to center human values in AI evaluation (Bender, 2024; Ehsan et al., 2024; Narayanan Venkit, 2023), our framework takes the results of the usability study involving domain experts who engage with GSE across technical and opinionated search queries (Narayanan Venkit et al., 2025). Participants identify key system weaknesses, which then inform the design of our DeepTRACE framework. Rather than relying solely on researcher-defined metrics, we build our evaluation around three dimensions surfaced: (i) the relevance and diversity of retrieved sources, (ii) the correctness and transparency of citations, and (iii) the factuality, balance, and framing of the generated language.

3 METHODOLOGY

Our motivation for auditing deep research agents and GSEs is grounded in the pressing call for more socially-aware evaluation practices in NLP. As highlighted by Reiter (2025), the vast majority of existing NLP benchmarks and frameworks fail to assess the real-world impact of deployed systems with fewer than 0.1% of papers include any form of societal evaluation. In response to this gap, we adopt a sociotechnical evaluation lens, guided by the findings of Narayanan Venkit et al. (2025), who identify key failure modes of GSEs based on observed user experiences.

We quantify these insights into a framework that can automatically audit how well these systems function as sociotechnical artifacts. To make the findings from our user study Narayanan Venkit et al. (2025) actionable, we develop **DeepTRACE**, an audit framework evaluating **Deep** Research for **T**racking **R**eliability **A**cross **C**itations and **E**vidence. Table 4, in Appendix C, outlines the mapping between qualitative insights, proposed system design recommendations, and their associated metrics. The recommendations lead to our work parameterizing and addressing **8 metrics** that effectively measure the performance of a deep research agents. We describe each metrics below.

3.1 DEEPTRACE METRICS

Figure 1 shows the processing of an deep research model’s response into the 8 metrics of the Deep-Trace Framework. We first go over the preliminary processing common to several metrics, then define each metric.

3.1.1 PRELIMINARY PROCESSING

When evaluating an GSE or a deep research agents, our evaluation framework requires the extraction of four content elements: the user query (1), the generated answer text (2) with the embedded citation (3) to the sources represented by a publicly accessible URL (4). Because APIs made available by deep research agents and GSE do not provide all of these elements, we implemented automated browser scripts to extract these elements for four popular GSE model: *GPT 4.5/5*, *You.com*, *Perplexity.ai*, and *BingChat*¹ and four deep research agents: *GPT 5 Deep Research*, *You.com Deep Research*, *Perplexity.ai Deep Research*, *BingChat Think Deeper* and *Gemini Deep Research*. Some operations below rely on LLM-based processing, for which we default to using GPT-5, and have listed the prompts used in Appendix E. When necessary, we evaluate the accuracy of LLM-based processing and report on the level of agreement with manual annotation.

¹Extending the evaluation to other GSE would require adapting the scripts to the specific website structure of the target GSE.

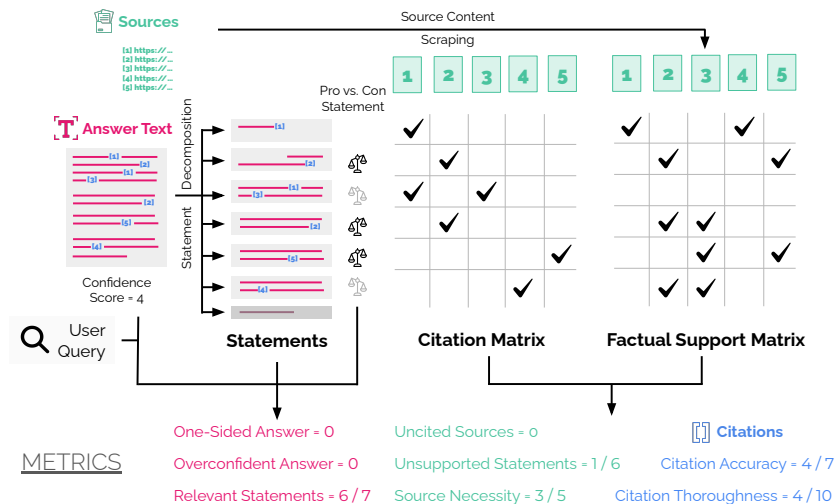


Figure 1: Illustrative diagram of the processing of a deep research agents response into the 8 metrics of the DeepTrace Framework. The description of each metrics is illustrated in Section 4.2.

A first operation consists of decomposing the answer text into statements. Decomposing the answer into statements allows to study the factual backing of the answer by the sources at a granular level, and is common in fact-checking literature (Laban et al., 2022; Tang et al., 2024; Huang et al., 2024; Qiu et al., 2024). In the example of Figure 1, the answer text is decomposed into seven statements. Each statement is further assigned two attributes: **Query Relevance** is a binary attribute that indicates whether the statement contains answer elements relevant to the user query. Irrelevant statements are typically introductory or concluding statements that do not contain factual information (e.g., “That’s a great question!”, “Let me see what I can do here”). **Pro vs. Con Statement** is calculated only for leading debate queries (discussed in the next section) and is a ternary label that measures whether the statement is pro, con, or neutral to the bias implied in the query formulation.

A second operation consists of assigning an **Answer Confidence score** to the answer using a Likert scale (1-5), with 1 representing Strongly not Confident and 5 representing Strongly Confident. Answer confidence is assigned by an *LLM judge* instructed with a prompt that provides examples of phrases used to express different levels of confidence based on the tone of the answer. This is specifically done for debate questions (Section 3.2). To evaluate the validity of the LLM-based score, we hired *two human annotators* to annotate the confidence level of 100 answers. We observed a Pearson correlation of 0.72 between the LLM judge and human annotators, indicating substantial agreement, and confirming the reliability of the LLM judge for confidence scoring. Given 80k support checks, LLM-judging is required for scalability, but we interpret results descriptively and highlight limitations instead of treating LLM outputs as ground truth.

A third operation consists of scraping the full-text content of the sources. We leverage Jina.ai’s Reader tool², to extract the full text of a webpage given its URL. Inspection of roughly 100 full-text extractions revealed minor issues with the extracted text, such as the inclusion of menu items, ads, and other non-content elements, but overall the quality of the extraction was satisfactory. For roughly 15% of the URLs, the Reader tool returns an error either due to the web page being behind a paywall, or due to the page being unavailable (e.g., a 404 error). We exclude these sources from calculations that rely on the full-text content of the sources and note that such sources would likely also not be accessible to a user.

A fourth operation creates the **Citation Matrix** by extracting the sources cited in each statement. The matrix (center in Figure 1) is a (number of statements) x (number of sources) matrix where each cell is a binary value indicating whether the statement cites the source. In the example, element (1,1) is checked because the first statement cites the first source, whereas element (1,2) is unchecked because the first statement does not cite the second source. A fifth operation creates the **Factual**

²<https://jina.ai/reader/>

Support Matrix by assigning for each (statement, source) pair a binary value indicating whether the source factually supports the statement. We leverage an LLM judge to assign each value in the matrix. A prompt including the extracted source content and the statement is constructed, and the LLM must determine whether the statement is supported or not by the source. Factual support evaluation is an open challenge in NLP (Tang et al., 2024; Kim et al., 2024), but top LLMs (GPT-5/4o) have been shown to perform well on the task (Laban et al., 2023a). To understand the degree of reliability of LLM-based factual support evaluation in our context, we hired *two annotators* to perform 100 factual verification tasks manually. We observed a *Pearson correlation of 0.62* between the LLM judge and manual labels, indicating moderate agreement. In the first row of the example Factual Support matrix, columns 1 and 4 are checked, indicating that sources 1 and 4 factually support the first statement.

For the annotation efforts, we hired a total of *four annotators* who are either professional annotators hired in *User Interviews*³, or graduate students enrolled in a computer science degree. We provided clear guidelines to annotators for the task and had individual Slack conversations where each annotator could discuss the task with the authors of the paper. Annotators were compensated at a rate of \$25 USD per hour. The annotation protocol was reviewed and approved by the institution’s Ethics Office. With the preliminary processing complete, we can now define the 8 metrics of the DeepTrace Evaluation Framework.

3.1.2 DEEPTRACE METRICS AND DEFINITIONS

I. One-Sided Answer: This binary metric is only computed on debate questions, leveraging the Pro vs. Con statement attribute. An answer is considered one-sided if it does not include both pro and con statements on the debate question.

$$\text{One-Sided Answer} = \begin{cases} 0 & \text{both pro and con} \\ & \text{statements are present} \\ 1 & \text{otherwise} \end{cases} \quad (1)$$

In the example of Figure 1, $\text{One-Sided Answer} = 0$ as there are three pro statements and two con statements. When considering a collection of queries, we can compute $\% \text{ One-Sided Answer}$ as the proportion of queries for which the answer is one-sided.

II. Overconfident Answer: This binary metric leverages the Answer Confidence score, combined with the One-Sided Answer metric and is only computed for debate queries. An answer is considered overconfident if it is both one-sided and has a confidence score of 5 (i.e., Strongly Confident).

$$\text{Overconfdnt. Ans} = \begin{cases} 1 & \text{if One-Sided Answer} = 1 \\ & \text{\& Answer Confidence} = 5 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

We implement a confidence metric in conjunction with the one-sided metric as it is challenging to determine the acceptable confidence level for any query. However, based on the user study findings by Narayanan Venkit et al. (2025), *an undesired trait in an answer is to be overconfident while not providing a comprehensive and balanced view*, which we capture with this metric. In the example of Figure 1, $\text{Overconfident Answer} = 0$ since the answer is not one-sided. When considering a collection of queries, we can compute $\% \text{ Overconfident Answer}$ as the proportion of queries with overconfident answers.

III. Relevant Statement: This ratio measures the fraction of relevant statements in the answer text in relation to the total number of statements.

$$\text{Relevant Statement} = \frac{\text{Number of Relevant Statements}}{\text{Total Number of Statements}} \quad (3)$$

This metric captures the to-the-pointedness of the answer, limiting introductory and concluding statements that do not directly address the user query. In the example of Figure 1, $\text{Relevant Statement} = 6/7$.

³www.userinterviews.com/

3.1.3 SOURCES METRICS

IV. Uncited Sources: This ratio metric measures the fraction of sources that are cited in the answer text in relation to the total number of listed sources.

$$\text{Uncited Sources} = \frac{\text{Number of Cited Sources}}{\text{Number of Listed Sources}} \quad (4)$$

This metric can be computed from the citation matrix: any empty column corresponds to an uncited source. In the example of Figure 1, since no column of the citation matrix is empty, $\text{Uncited Sources} = 0 / 5$.

V. Unsupported Statements: This ratio metric measures the fraction of relevant statements that are not factually supported by any of the listed sources. Any row of the factual support matrix with no checked cell corresponds to an unsupported statement.

$$\text{Unsupported Statements} = \frac{\text{No. of Unsupported St.}}{\text{No. of Relevant St.}} \quad (5)$$

In the example of Figure 1, the third row of the factual support matrix is the only entirely unchecked row, indicating that the third statement is unsupported. Therefore, $\text{Unsupported Statements} = 1 / 6$.

VI. Source Necessity: This ratio metric measures the fraction of sources that are necessary to factually support all relevant statements in the answer text. Understanding what source is necessary or redundant can be formulated as a graph problem. We transform the factual support matrix into a (statement,source) bi-partite graph. Finding which source is necessary is equivalent to determining the minimum vertex cover for source nodes on the bipartite graph. We use the Hopcroft-Karp algorithm (Hopcroft & Karp, 1973) to find the minimum vertex cover, which tells us which sources are necessary to cover factually supported statements.

$$\text{Source Necessity} = \frac{\text{Number of Necessary Sources}}{\text{Number of Listed Sources}} \quad (6)$$

In the example of Figure 1, one possible minimum vertex cover consists of sources 1, 2, and 3 (another consists of 2, 3, and 4). Therefore, $\text{Source Necessity} = 3 / 5$. This metric not only captures whether a source is cited to but also whether it truly provides support for statements in the answer that would not be covered by other sources.

3.1.4 CITATION METRICS

VII. Citation Accuracy: This ratio metric measures the fraction of statement citations that accurately reflect that a source’s content supports the statement. This metric can be computed by measuring the overlap between the citation and the factual support matrices, and dividing by the number of citations:

$$\text{Cit. Acc.} = \frac{\sum \text{Citation Mtx} \odot \text{Factual Support Mtx}}{\sum \text{Citation Mtx}} \quad (7)$$

Where \odot is element-wise multiplication, and \sum is the sum of all elements in the matrix. In the example of Figure 1, there are four accurate citations ((1,1), (2,2), (4,2) and (5,5)), and three inaccurate citations ((3,1), (3,3), (6,4)), so $\text{Citation Accuracy} = 4 / 7$.

VIII. Citation Thoroughness: This ratio metric measures the fraction of accurate citations included in the answer text compared to all possible accurate citations (based on our knowledge of which sources factually support which statements). This metric can be computed by measuring the overlap between the citation and the factual support matrices:

$$\text{Cit. Th.} = \frac{\sum \text{Citation Mtx} \odot \text{Factual Support Mtx}}{\sum \text{Factual Support Mtx}} \quad (8)$$

In the example of Figure 1, there are four accurate citations, and ten factual support relationships (such as (1,4), (2,5), etc.), so $\text{Citation Thoroughness} = 4 / 10$.

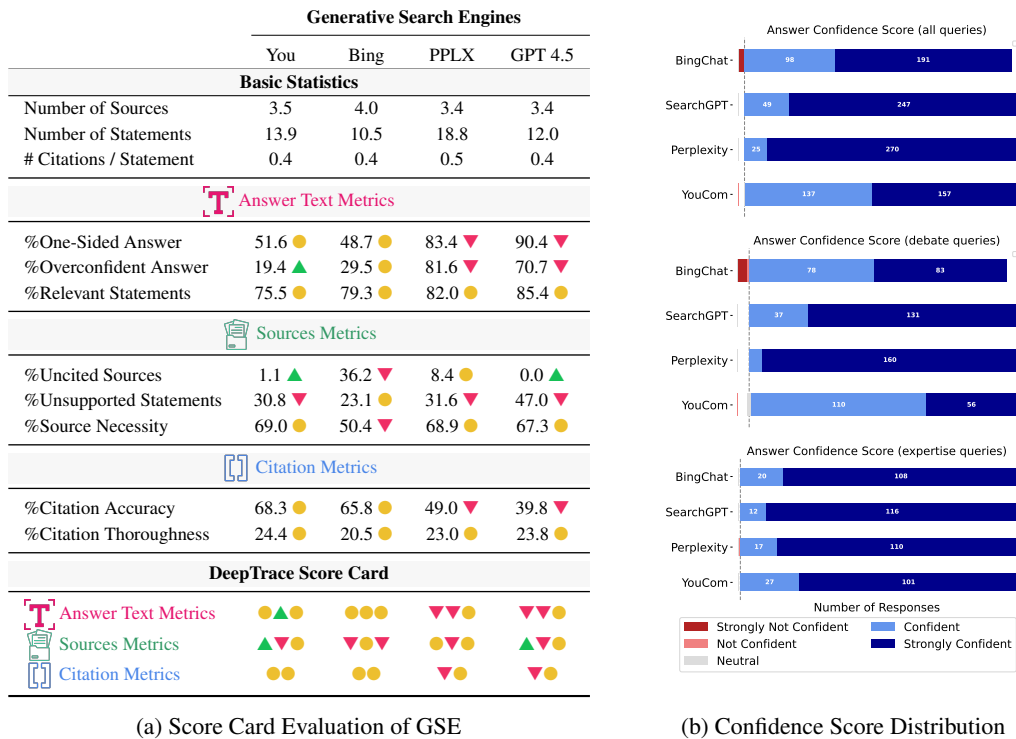
We note that we do not implement metrics related to the ‘*User Interface*’ findings of Narayanan Venkit et al. (2025), as they are not directly computable from the answer text, citation, and source content and would likely require manual evaluation, or computer-vision-based methods that are out of the scope of this work.

3.2 DEEPTRACE CORPUS AND FRAMEWORK

To perform the above evaluation, we use and release the DeepTrace dataset, which is used to prompt responses and assess model behavior. The dataset comprises **303 questions** shared by the sessions conducted by our prior work Narayanan Venkit et al. (2025). These questions are divided into two categories:

- **Debate Questions (N=168)**: These questions, sourced from the ProCon website, a non-partisan platform providing balanced information on contentious issues, are characterized by their tendency to have multiple perspectives and are often subjects of debate⁴.
- **Expertise Questions (N=135)**: These questions were contributed by the participants from Narayanan Venkit et al. (2025), who represented experts from diverse fields including meteorology, medicine, and human-computer interaction. These questions pertain to research-oriented questions that tend to need multiple searches/hops.

An example debate question in DeepTrace is “Why can alternative energy effectively not replace fossil fuels?”, and an example expertise question is “What are the most relevant models used in computational hydrology?”. We then use developed browser scripts to run each query through a total of 9 public GSE and DR agents to extract all components required for metric-based evaluation, and computed the metrics on the relevant queries: most metrics are computed on all 2,727 samples (303 queries x 9 models), while a few are only computed on the debate queries (e.g., One-Sided Answer, Overconfident Answer). Using the DeepTrace dataset, we conducted evaluation of the models to parameterize and understand their behavior and weaknesses, using the above 8 metrics. The modular design of the DeepTrace framework and dataset allows for flexible adaptation, enabling the dataset’s modification for continued evaluation of GSE and deep research agents across different contexts and therefore is not solely dependant on the specific dataset.



(a) Score Card Evaluation of GSE

(b) Confidence Score Distribution

Figure 2: Quantitative Evaluation of three GSE – You.com, BingChat, and Perplexity – based on the eight metrics of the DeepTrace framework: metric report, color-coded for ▲ acceptable, ● borderline, and ▼ problematic performance. Figure (b) plots distributions of answer confidence.

⁴<https://www.procon.org/>

Deep Research Agents							
	GPT-5(DR)	YouChat(ARI)	YouChat(DR)	GPT-5(S)	PPLX(DR)	Copilot (TD)	Gemini (DR)
Basic Statistics							
Number of Sources	18.3	198.61	57.2	13.5	7.7	3.6	33.2
Number of Statements	141.6	39.06	52.7	34.9	30.1	36.7	23.9
# Citations / Statement	1.4	1.69	0.8	0.4	0.2	0.3	0.2
Answer Text Metrics							
%One-Sided Answer	54.67 ▼	0.0 ▲	63.1 ▼	69.7 ▼	63.1 ▼	94.8 ▼	80.1 ▼
%Overconfident Answer	15.2 ▲	N/A	19.6 ▲	16.4 ▲	5.6 ▲	0.0 ▲	11.2 ▲
%Relevant Statements	87.5 ●	37.15 ●	45.5 ▼	41.1 ▼	22.5 ▼	13.2 ▼	12.4 ▼
Sources Metrics							
%Uncited Sources	0.0 ▲	0.0 ▲	66.3 ▼	51.7 ▼	57.5 ▼	32.6 ▼	14.5 ▼
%Unsupported Statements	12.5 ●	62.85 ▼	74.6 ▼	58.9 ▼	97.5 ▼	90.2 ▼	53.6 ▼
%Source Necessity	87.5 ▲	42.65 ▼	63.2 ●	32.8 ▼	5.5 ▼	31.2 ▼	33.1 ▼
Citation Metrics							
%Citation Accuracy	79.1 ●	39.33 ▼	72.3 ●	31.4 ▼	58.0 ●	62.1 ●	50.3 ●
%Citation Thoroughness	87.5 ▲	96.77 ▲	83.5 ▲	17.9 ▼	9.1 ▼	13.2 ▼	27.1 ●
DeepTrace Eval Score Card							
Answer Text Metrics	▼▲●	▲●	▼▲▼	▼▲▼	▼▲▼	▼▲▼	▼▲▼
Sources Metrics	▲▲●	▲▼▼	▼▼●	▼▼▼	▼▼▼	▼▼▼	▼▼▼
Citation Metrics	●▲	▼▲	●▲	▼▼	●▼	●▼	●●

Table 1: DeepTrace results for our DR based models: GPT-5, YouChat (ARI), YouChat (DR), Perplexity (PPLX), Copilot Think Deeper, Gemini and GPT-5 Web Search (S) setting. Metrics evaluated according to DeepTrace thresholds: ▲acceptable, ●borderline, ▼problematic.

4 RESULTS

Figure 2 (GSE) and Table 1 (Deep Research) show the results of the metrics-based evaluation on the DeepTrace as of August 27, 2025. To focus on publicly accessible systems, we selected the web search and deep research capabilities of Perplexity, Bing Copilot, GPT (4.5/5) and YouChat for evaluation as accessed from their public UI. Numerical values are assigned a color based on whether the score reflects an ▲ acceptable, ● borderline, and ▼ problematic performance. Thresholds for the colors are listed with the explanation of the threshold in Appendix A based on the qualitative inputs obtained the user study Narayanan Venkit et al. (2025). These threshold bands are derived from tolerance ranges observed in multi-session user research done by us Narayanan Venkit et al. (2025) and are intended as illustrative diagnostic categories. All comparative conclusions in this paper rely on the raw metric values rather than these visual bins.

Generative Search Engines. As shown in Figure 2, for **answer text metrics**, one-sidedness remains an issue (50–80%), with Perplexity performing worst, generating one-sided responses in over 83% of debate queries despite producing the longest answers (18.8 statements per response on average). Confidence calibration also varies where BingChat and You.com reduce confidence when addressing debate queries, whereas Perplexity maintains uniformly high confidence (90%+ very confident), resulting in *overconfident yet one-sided answers on politically or socially contentious prompts*. On relevance, GSE models perform comparably (75–85% relevant statements), which indicates better alignment with user queries relative to their DR counterparts. For **source metrics**, BingChat exemplifies the quantity without quality trade-off where it lists more sources on average (4.0), yet over a third remain uncited and only about half are necessary. You.com and Perplexity list slightly fewer sources (3.4–3.5) but still struggle with unsupported claims (23–47%). Finally, on **citation metrics**, all three engines show relatively low citation accuracy (40–68%), with frequent misattribution. Even when a supporting source exists, models often cite an irrelevant one, preventing users from verifying factual validity. Citation thoroughness is also limited, with engines typically citing only a subset of available supporting evidence.

Deep Research Agents. In context of **answer text**, Table 1 shows that DR modes do not eliminate one-sidedness where rates remain high across the board (54.7–94.8%). Appendix D shows how GPT-5 deep research answers one sided answers for questions framed pro and con the same debate, without providing generalized coverage. This showcases sycophantic behavior of aligning only with the users perspective, causing potential echo chambers to search. Overconfidence is consistently low across DR engines (<20%), indicating that calibration of language hedging is one relative strength of this pipeline. On **relevance**, however, performance is uneven where GPT-5(DR) attains borderline results (87.5%), while all other engines fall below 50%, including Gemini(DR) at just 12.4%. This suggests that verbosity or sourcing breadth does not translate to actually answering the user query. Turning to **sources metrics**, GPT-5(DR) remains the strongest with 0% uncited sources, only 12.5% unsupported statements, and 87.5% source necessity. By contrast, YouChat(DR), PPLX(DR), Copilot(DR), and Gemini(DR) all fare poorly, with unsupported rates ranging from 53.6% (Gemini) to 97.5% (PPLX). Gemini(DR) in particular includes 14.5% uncited sources and only one-third (33.1%) of its sources being necessary, reflecting inefficient citation usage. For **citation metrics**, GPT-5(DR) and YouChat(DR) again stand out with high citation thoroughness (87.5% and 83.5% respectively), although their citation accuracy has dropped to the borderline range (79.1% and 72.3%). Gemini(DR) demonstrates weak citation performance: only 40.3% citation accuracy (problematic) and 27.1% thoroughness (borderline). PPLX(DR) and Copilot(DR) also show poor grounding, with citation accuracies between 58–62%. Our qualitative study also suggests that GPT-5(DR) tends to produce concise, well-bounded statements and selectively cites sources that directly support those statements, whereas Perplexity(DR) often generates verbose answers, spreads citations across loosely relevant sources, and relies heavily on first-retrieved pages, behaviors that contribute to the large gap in unsupported-statement rates.

Taken together, the results reveal that neither GSE nor DR, deliver uniformly reliable outputs across DeepTRACE’s dimensions. GSEs excel at producing concise, relevant answers but fail at balanced perspective-taking, confidence calibration, and factual support. Deep research agents, by contrast, improve balance and citation correctness, but at the cost of overwhelming verbosity, low relevance, and huge unsupported claims. Our results show that more sources and longer answers do not translate into reliability. Over-citation (as in YouChat(DR)) leads to ‘search fatigue’ for users, while under-grounded verbose texts (as in Perplexity(DR)) erodes trust. At the same time, carefully calibrated systems (as with GPT-5(DR)) demonstrate near-ideal reliability across multiple dimensions.

5 DISCUSSION AND CONCLUSION

Our work introduced DeepTRACE, a sociotechnically grounded framework for auditing generative search engines (GSEs) and deep research agents (DRs). By translating community-identified failure cases into measurable dimensions, our approach evaluates not just isolated components but the end-to-end reliability of these systems across balance, factual support, and citation integrity. Our evaluation demonstrates that current public systems fall short of their promise to deliver trustworthy, source-grounded synthesis. Generative search engines tend to produce concise and relevant answers but consistently exhibit one-sided framing and frequent overconfidence, particularly on debate-style queries. Deep research agents, while reducing overconfidence and improving citation thoroughness, often overwhelm users with verbose, low-relevance responses and large fractions of unsupported claims. Our findings show that increasing the number of sources or length of responses does not reliably improve grounding or accuracy; instead, it can exacerbate user fatigue and transparency.

Citation practices remain a persistent weakness across both classes of systems. Many citations are either inaccurate or incomplete, with some models listing sources that are never cited or irrelevant. This creates a misleading impression of evidential rigor while undermining user trust. Metrics such as Source Necessity and Citation Accuracy highlight that merely retrieving more sources does not equate to stronger factual grounding, echoing user concerns about opacity and accountability. Taken together, these results point to a central tension: GSEs optimize for summarization and relevance at the expense of balance and factual support, whereas DRs optimize for breadth and thoroughness at the expense of clarity and reliability. Neither approach, in its current form, adequately meets the sociotechnical requirements of safe, effective, and trustworthy information access. However, our findings also suggest that calibrated systems, such as GPT-5(DR), which demonstrated strong performance across multiple metrics, illustrate that more reliable designs are achievable.

6 ETHICS STATEMENT

While DeepTRACE offers an automated and scalable evaluation platform, it currently focuses on textual and citation-based outputs, excluding multimodal or UI-level interactions that also shape user trust and system usability. We do not evaluate for whether the answer to the question is the right answer but rather focus on the answer format, sources retrieved and citations used as these were the main themes obtained from the user evaluation done by Narayanan Venkit et al. (2025). Furthermore, some reliance on LLMs for intermediate judgments (e.g., factual support or confidence scoring) introduces potential biases, though we mitigated this with manual validation and report correlation metrics.

7 REPRODUCIBILITY STATEMENT

We have made extensive efforts to support reproducibility of our work. The DeepTRACE evaluation framework, including the decomposition pipeline, and metric definitions, are described in detail in Section 3 and Appendix, with additional implementation details provided in the supplementary materials. We release the DeepTRACE dataset of 303 queries (debate and expertise questions) along with the evaluation pipeline as supplementary material to enable replication and extension in this repository⁵. Since our study evaluates publicly available generative search engines and deep research agents directly through their web interfaces (rather than fixed API endpoints), we note that model behaviors may evolve over time. This decision is done intentionally as our audit focuses on user centric usage rather than specific model performances. We provide the evaluation timestamp (August 27, 2025) to clarify the snapshot of system behavior we captured. All metric calculations, data processing steps, and annotation protocols are fully documented in the main text and appendices to ensure transparency and reproducibility.

REFERENCES

- Emily M Bender. Resisting dehumanization in the age of “ai”. *Current Directions in Psychological Science*, 33(2):114–120, 2024.
- Nikos I Bosse, Jon Evans, Robert G Gambee, Daniel Hnyk, Peter Mühlbacher, Lawrence Phillips, Dan Schwarz, Jack Wildman, et al. Deep research bench: Evaluating ai web research agents. *arXiv preprint arXiv:2506.06287*, 2025.
- Pratyush Chauhan, Rahul Kumar Sahani, Soham Datta, Ali Qadir, Manish Raj, and Mohd Mohsin Ali. Evaluating top-k rag-based approach for game review generation. In *2024 IEEE International Conference on Computing, Power and Communication Technologies (IC2PCT)*, volume 5, pp. 258–263. IEEE, 2024.
- Zijian Chen, Xueguang Ma, Shengyao Zhuang, Ping Nie, Kai Zou, Andrew Liu, Joshua Green, Kshama Patel, Ruoxi Meng, Mingyi Su, et al. Browsecomp-plus: A more fair and transparent evaluation benchmark of deep-research agent. *arXiv preprint arXiv:2508.06600*, 2025.
- Robert Cooper and Michael Foster. Sociotechnical systems. *American Psychologist*, 26(5):467, 1971.
- Mateusz Dolata, Stefan Feuerriegel, and Gerhard Schwabe. A sociotechnical view of algorithmic fairness. *Information Systems Journal*, 32(4):754–818, 2022.
- Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents. *arXiv preprint arXiv:2506.11763*, 2025.
- Upol Ehsan, Samir Passi, Q Vera Liao, Larry Chan, I-Hsiang Lee, Michael Muller, and Mark O Riedl. The who in xai: How ai background shapes perceptions of ai explanations. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 1–32, 2024.
- Shahul Es, Jithin James, Luis Espinosa-Anke, and Steven Schockaert. Ragas: Automated evaluation of retrieval augmented generation. *arXiv preprint arXiv:2309.15217*, 2023.

⁵<https://github.com/SalesforceAIRResearch/answer-engine-eval/tree/main/Venkit.et.al.2024>

- Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. Ragas: Automated evaluation of retrieval augmented generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pp. 150–158, 2024.
- Emilio Ferrara. Genai against humanity: Nefarious applications of generative artificial intelligence and large language models. *Journal of Computational Social Science*, pp. 1–21, 2024.
- Aman Gupta, Anup Shirgaonkar, Angels de Luis Balaguer, Bruno Silva, Daniel Holstein, Dawei Li, Jennifer Marsman, Leonardo O Nunes, Mahsa Rouzbahman, Morris Sharp, et al. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture. *arXiv preprint arXiv:2401.08406*, 2024.
- John E Hopcroft and Richard M Karp. An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs. *SIAM Journal on computing*, 2(4):225–231, 1973.
- Kung-Hsiang Huang, Mingyang Zhou, Hou Pong Chan, Yi Fung, Zhenhailong Wang, Lingyu Zhang, Shih-Fu Chang, and Heng Ji. Do LVLMs understand charts? analyzing and correcting factual errors in chart captioning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 730–749, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.41. URL <https://aclanthology.org/2024.findings-acl.41/>.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023.
- Yuxuan Huang, Yihang Chen, Haozheng Zhang, Kang Li, Meng Fang, Linyi Yang, Xiaoguang Li, Lifeng Shang, Songcen Xu, Jianye Hao, et al. Deep research agents: A systematic examination and roadmap. *arXiv preprint arXiv:2506.18096*, 2025.
- Gautier Izacard and Édouard Grave. Leveraging passage retrieval with generative models for open domain question answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 874–880, 2021.
- Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C Park. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 7029–7043, 2024.
- Zheng Ping Jiang, Jingyu Zhang, Nathaniel Weir, Seth Ebner, Miriam Wanner, Kate Sanders, Daniel Khashabi, Anqi Liu, and Benjamin Van Durme. Core: Robust factual precision with informative sub-claim identification. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 19833–19856, 2025.
- Navreet Kaur, Monojit Choudhury, and Danish Pruthi. Evaluating large language models for health-related queries with presuppositions. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 14308–14331, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.findings-acl.850>.
- Yekyung Kim, Yapei Chang, Marzena Karpinska, Aparna Garimella, Varun Manjunatha, Kyle Lo, Tanya Goyal, and Mohit Iyyer. Fables: Evaluating faithfulness and content selection in book-length summarization. *arXiv preprint arXiv:2404.01261*, 2024.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. Summac: Re-visiting nli-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177, 2022.

- Philippe Laban, Wojciech Kryściński, Divyansh Agarwal, Alexander R Fabbri, Caiming Xiong, Shafiq Joty, and Chien-Sheng Wu. Llms as factual reasoners: Insights from existing benchmarks and beyond. *arXiv preprint arXiv:2305.14540*, 2023a.
- Philippe Laban, Lidiya Murakhovska, Caiming Xiong, and Chien-Sheng Wu. Are you sure? challenging llms leads to performance drops in the flipflop experiment. *arXiv preprint arXiv:2311.08596*, 2023b.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33: 9459–9474, 2020.
- Nelson F Liu, Tianyi Zhang, and Percy Liang. Evaluating verifiability in generative search engines. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 7001–7025, 2023.
- Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12076–12100, 2023.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Pranav Narayanan Venkit. Towards a holistic approach: Understanding sociodemographic biases in nlp models using an interdisciplinary lens. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 1004–1005, 2023.
- Pranav Narayanan Venkit, Philippe Laban, Yilun Zhou, Yixin Mao, and Chien-Sheng Wu. Search engines in the ai era: A qualitative understanding to the false promise of factual and verifiable source-cited responses in llm-based search. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1325–1340, 2025.
- Deepak Nathani, Lovish Madaan, Nicholas Roberts, Nikolay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, et al. Ml-gym: A new framework and benchmark for advancing ai research agents. *arXiv preprint arXiv:2502.14499*, 2025.
- Sanjeev Pulapaka, Srinath Godavarthi, and Dr Sherry Ding. Genai and the public sector. In *Empowering the Public Sector with Generative AI: From Strategy and Design to Real-World Applications*, pp. 31–43. Springer, 2024.
- Haoyi Qiu, Kung-Hsiang Huang, Jingnong Qu, and Nanyun Peng. AMRFact: Enhancing summarization factuality evaluation with AMR-driven negative samples generation. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 594–608, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.33. URL <https://aclanthology.org/2024.naacl-long.33/>.
- Ehud Reiter. We should evaluate real-world impact. *Computational Linguistics*, 2025.
- Sujoy Roychowdhury, Sumit Soman, HG Ranjani, Neeraj Gunda, Vansh Chhabra, and Sai Krishna Bala. Evaluation of rag metrics for question answering in the telecom domain. *arXiv preprint arXiv:2407.12873*, 2024.
- Chirag Shah and Emily M Bender. Envisioning information access systems: What makes for good tools and a healthy web? *ACM Transactions on the Web*, 18(3):1–24, 2024.
- Nikhil Sharma, Q Vera Liao, and Ziang Xiao. Generative echo chamber? effect of llm-powered search systems on diverse information seeking. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pp. 1–17, 2024.

- Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. Improving the domain adaptation of retrieval augmented generation (rag) models for open domain question answering. *Transactions of the Association for Computational Linguistics*, 11:1–17, 2023.
- Liyan Tang, Philippe Laban, and Greg Durrett. Minicheck: Efficient fact-checking of llms on grounding documents. *arXiv preprint arXiv:2404.10774*, 2024.
- Pranav Narayanan Venkit, Tatiana Chakravorti, Vipul Gupta, Heidi Biggs, Mukund Srinath, Koustava Goswami, Sarah Rajtmajer, and Shomir Wilson. ” confidently nonsensical? ”: A critical survey on the perspectives and challenges of ’hallucinations’ in nlp. *arXiv preprint arXiv:2404.07461*, 2024.
- Yidong Wang, Qi Guo, Wenjin Yao, Hongbo Zhang, Xin Zhang, Zhen Wu, Meishan Zhang, Xinyu Dai, Qingsong Wen, Wei Ye, et al. Autosurvey: Large language models can automatically write surveys. *Advances in neural information processing systems*, 37:115119–115145, 2024.
- Junde Wu, Jiayuan Zhu, and Yuyuan Liu. Agentic reasoning: Reasoning llms with tools for the deep research. *arXiv preprint arXiv:2502.04644*, 2025.
- Kevin Wu, Eric Wu, and James Zou. How faithful are rag models? quantifying the tug-of-war between rag and llms’ internal prior. *arXiv preprint arXiv:2404.10198*, 2024.
- Elvin Wyly. Automated (post) positivism. *Urban Geography*, 35(5):669–690, 2014.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025.
- Kunlun Zhu, Yifan Luo, Dingling Xu, Ruobing Wang, Shi Yu, Shuo Wang, Yukun Yan, Zhenghao Liu, Xu Han, Zhiyuan Liu, et al. Rageval: Scenario specific rag evaluation dataset generation framework. *arXiv preprint arXiv:2408.01262*, 2024.
- Theresa Züger and Hadi Asghari. Ai for the public. how public interest theory shifts the discourse on ai. *AI & SOCIETY*, 38(2):815–828, 2023.

A SCORE CARD METRICS THRESHOLDS

Table 2 establishes the benchmark ranges for the eight DeepTrace Evaluation metrics, categorizing performance into three levels: ▲ acceptable, ● borderline, and ▼ problematic. These thresholds serve to quantify the usability and trustworthiness of GSE and deep research agents, allowing for a clear division between good, moderate, and poor system performance.

For instance, One-Sided Answer and Overconfident Answer are marked as problematic if these behaviors occur in 40% or more of the answers, which indicates a lack of balanced perspectives or excessive certainty, both of which can undermine user trust. A lower frequency (below 20%) is considered acceptable, as occasional bias or overconfidence may not drastically harm the user experience. Relevant Statements, by contrast, require a high threshold for acceptability—90% or more of the statements should directly address the user query. Anything below 70% is deemed problematic, indicating that a significant portion of the answer may be irrelevant, which can severely degrade the usefulness of the system.

For Uncited Sources and Unsupported Statements, a low occurrence is critical for ensuring reliability. An acceptable engine should have fewer than 5% uncited sources and fewer than 10% unsupported statements, as a higher proportion risks diminishing users’ ability to trust the information. Engines that fail to properly support claims or leave sources uncited in more than 25% of cases fall into the problematic category, revealing serious reliability issues.

DeepTrace Metric	▲ Acceptable	● Borderline	▼ Problematic
One-Sided Answer	[0,20)	[20,40)	[40,100)
Overconfident Answer	[0,20)	[20,40)	[40,100)
Relevant Statements	[90, 100)	[70,90)	[0,70)
Uncited Sources	[0,5)	[5,10)	[10,100)
Unsupported Statements	[0,10)	[10,25)	[25,100)
Source Necessity	[80,100)	[60,80)	[0,60)
Citation Accuracy	[90,100)	[50,90)	[0,50)
Citation Thoroughness	[50,100)	[20,50)	[0,20)

Table 2: Ranges for the eight DeepTrace metrics for a system’s performance to be considered ▲ acceptable, ● borderline, or ▼ problematic on a given metric.

Task	Correlation with humans	Scale
Answer confidence (debate queries)	0.72	Likert 1–5
Factual support (statement–source)	0.62	binary

Table 3: Human–LLM agreement for the LLM-as-judge components used in DeepTRACE. We report Pearson correlations between human annotations and LLM judgments for answer confidence and factual support (N = 100 samples per task).

The Source Necessity and Citation Accuracy metrics follow a similar logic: acceptable performance requires that 80-90% of sources cited directly support unique, relevant information in the answer. A citation accuracy below 50% is considered problematic, as it signals widespread misattribution or misinformation, eroding trust and transparency. Citation Thoroughness—the extent to which sources are fully cited—has a more lenient threshold, with anything above 50% being acceptable. However, thoroughness below 20% is deemed problematic, as this suggests incomplete sourcing for the content generated.

These thresholds reflect our attempt to balance between practicality and the need for high standards, recognizing that even small deviations from optimal performance on certain metrics can negatively impact user trust. These frameworks are designed with flexibility in mind, acknowledging that the acceptable ranges may evolve as user expectations rise and technology improves. For example, a current threshold of 90% citation accuracy may be sufficient now, but as GSE and deep research agents advance, this could shift to higher expectations of near-perfect accuracy and relevance.

B HUMAN ANNOTATOR AND MODEL JUDGE ALIGNMENT

Table B showcases the human–LLM judge agreement used for the two components in DeepTrace.

C METRICS ASSOCIATED TO RECOMMENDATIONS

Table 4 showcases what metrics were generated based on the recommendations and findings from our extensive user study Narayanan Venkit (2023).

D EXAMPLES OF RESPONSES

In this section, Figure 3 and Figure 4 shows how deep research models, specifically GPT-5 Deep Research, tend to generate outputs that closely follow the framing of the input questions, even when broader or more holistic perspectives may be warranted. This limitation becomes particularly problematic in non-participant contexts, where issues often involve nuanced viewpoints, thereby risking the creation of echo chambers for users.

Design Recommendation	Associated System Weakness	Metric Developed
Provide balanced answers	Lack of holistic viewpoints for opinionated questions [A.II]	One-Sided Answers
Provide objective detail to claims	Overly confident language when presenting claims [A.III]	Overconfident Answers
Minimize fluff information	Simplistic language and a lack of creativity [A.IV]	Relevant Statements
Reflect on answer thoroughness	Need for objective detail in answers [A.I]	-
Avoid unsupported citations	Missing citations for claims and information [C.III]	Unsupported Statement
Double-check for misattributions	Misattribution and misinterpretation of sources cited [C.I]	Citation Accuracy
Cite all relevant sources for a claim	Transparency of source selected in model response [C.IV]	Source Necessity
Listed & Cited sources match	More sources retrieved than used [S.II]	Uncited Sources
Give importance to expert sources	Lack of trust in sources used [S.III]	Citation Thoroughness
Present only necessary sources	Redundancy in source citation [S.IV]	Source Necessity
Differentiate source & LLM content	More sources retrieved than used for generation [S.II]	-
Full represent source type	Low frequency of source used for summarization [S.I]	-
Incorporate human feedback	Lack of search, select and filter [U.I]	-
Implement interactive citation	Citation formats are not normalized interactions [U.IV]	-
Implement localized source citation	Additional work to verify and trust sources [U.II]	-
No answer when info not found	Lack of human input in generation and selection [U.I]	-

Table 4: Sixteen design recommendations for generative search engines and deep research agents. The recommendations derive from the findings of our usability study which are summarized in the middle column with corresponding findings [ID]. Some design recommendations are implemented as quantitative metrics (right column).

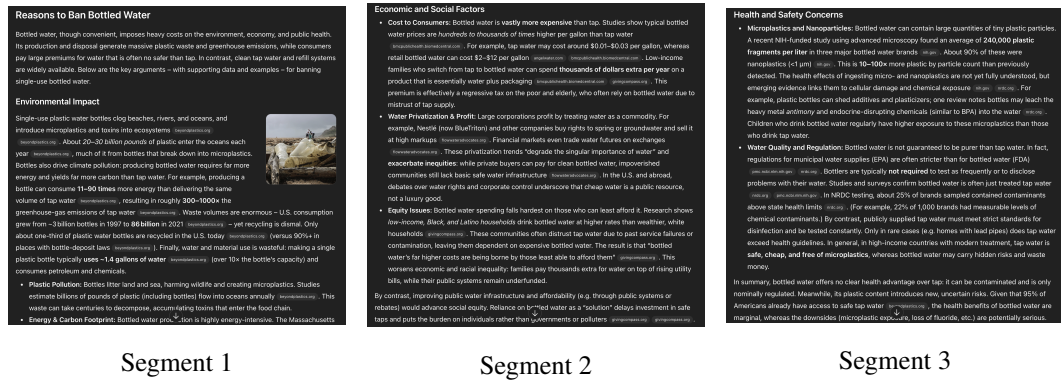


Figure 3: Screenshots of responses from GPT-5 Deep Research to the debate question, “why should we ban bottled water?”. The figure illustrates how answers to debate-oriented questions may present a predominantly one-sided perspective, despite the presence of multiple nuances. Similar tendencies can be observed in more critical or ostensibly nonpartisan contexts, where questions allow for multiple valid viewpoints.

E LLM JUDGE PROMPTS USED FOR EVALUATION

In this section, we present the LLM instructions used to automate tasks for the DeepTrace scorecard evaluation. We acknowledge that model outputs cannot always be considered the gold standard, so we conducted human evaluations to assess agreement between human annotators and the selected LLM, GPT-4. The results are detailed in Section 4.2. The prompts are as follows:



(a) Screenshot 4

Figure 4: Screenshots of responses from GPT-5 Deep Research to the debate question, “why should we not ban bottled water?”. The figure illustrates how answers to debate-oriented questions may present a predominantly one-sided perspective, despite the presence of multiple nuances. Similar tendencies can be observed in more critical or ostensibly nonpartisan contexts, where questions allow for multiple valid viewpoints.

Confidence Parameter Evaluation

You are given a user Query, and the Answer generated by a system.

Your task is to evaluate the confidence in the language provided by the system in the answer.

You must choose one of five confidence levels.

Annotation Scheme:

- "Strongly Not Confident":

The language in this sentence is not clear or confident at all.

- "Not Confident":

The language in this sentence is somewhat unclear and lacks confidence.

- "Neutral":

The language in this sentence is neither clear nor unclear; confidence level is average.

- "Confident":

The language in this sentence is clear and fairly confident.

- "Strongly Confident":

The language in this sentence is very clear and confident.

Format:

- You must produce your answer as a JSON object, following this format:

```
{"confidence": "<Confidence Level>"}
```

- Replace <Confidence Level> with one of the five confidence levels.

- Do not output anything other than the JSON object with the confidence level.

Query:

```
[[QUERY]]
```

Answer:

```
[[ANSWER]]
```

Relevant Statement Extraction

You are given a paragraph, made of a sequence of sentences that answer the following question:
[[QUESTION]]

Your task is to extract, in JSON format, what the individual sentences are, and then identify for each sentence whether it contains a core statement that answers the question, or if it is a filler sentence that does not contain substantial information.

You should follow the following format:

```
{"sentences":  
  {"sentence": "...", "core": "1|0"},  
  {"sentence": "...", "core": "1|0"},  
}
```

Rules:

- Do not modify the sentences whatsoever, you should copy them as is.
- Do not modify the order of the sentences, or skip any of the sentences.
- The sentences optionally contain citations (e.g. [1], [2], etc.).

You should not modify the citations, keep them as is.

- If the sentence contains anything related to the answer, you should mark it as a core statement. Sentences with a citation are

almost always core statements.

- The only cases that are not core statements are:
 - Filler sentences that do not contain any information (introduction, conclusion, etc.)

Here is the answer you should decompose:

```
```\n[[ANSWER]]\n```
```

### Supported Statement Evaluation

You are given the full text of a document, and a statement (a sentence).  
Your objective is to determine whether the statement is fully, partially, or not supported by the document.  
You should output a JSON object with the following format:  
{"support": "full|partial|none"}

Document:

```
```\n[[DOCUMENT]]\n```
```

Statement:

```
[[STATEMENT]]
```

Rules:

- The statement may contain citations (e.g. [1], [2], etc.). You do not need to consider the citations when determining support.
- Only output valid JSON. Do not include any other information in the output.

One-sided Answer Detection

You are given a query that is opinionated (on one side of an issue).

You are then given a numbered list of statements.

Your objective is to sort these statements into three lists:

- "agree_statements" if they agree with the opinionated query
- "disagree_statements" if they disagree with the opinionated query,
- "neutral_statements" if they are neutral to the opinionated query.

You should return a JSON object following the given format:

```
{"agree_statements": [1, 2, 3, ...],  
"disagree_statements": [4, 5, 6, ...],  
"neutral_statements": [7, 8, 9, ...]}
```

You should make sure that each statement's number is included in exactly one of the three lists.

Query:
[[QUERY]]

Statements:
[[STATEMENTS]]

Remember to follow the format given above, only output JSON.