

ResearchArena: Benchmarking Large Language Models’ Ability to Collect and Organize Information as Research Agents

Anonymous ACL submission

Abstract

Large language models (LLMs) excel across many natural language processing tasks but face challenges in domain-specific, analytical tasks such as conducting research surveys. This study introduces ResearchArena, a benchmark designed to evaluate LLMs’ capabilities in conducting academic surveys—a foundational step in academic research. ResearchArena models the process in three stages: (1) information discovery, identifying relevant literature; (2) information selection, evaluating papers’ relevance and impact; and (3) information organization, structuring knowledge into hierarchical frameworks such as mind-maps. Notably, mind-map construction is treated as a bonus task, reflecting its supplementary role in survey-writing. To support these evaluations, we construct an offline environment of 12M full-text academic papers and 7.9K survey papers. To ensure ethical compliance, we do not redistribute copyrighted materials; instead, we provide code to construct the environment from the Semantic Scholar Open Research Corpus (S2ORC). Preliminary evaluations reveal that LLM-based approaches underperform compared to simpler keyword-based retrieval methods, underscoring significant opportunities for advancing LLMs in autonomous research.

1 Introduction

Large language models (LLMs) have demonstrated exceptional performance in natural language understanding, text generation, and a range of other tasks across domains (Liang et al., 2022; Bang et al., 2023; Qin et al., 2023a; Laskar et al., 2023). By integrating LLMs with external tools—such as code interpreters, vector databases, and search engines—their capabilities can be further enhanced, enabling the creation of autonomous agents that simulate human-like behavior through feedback-driven task execution (Wang et al., 2024; Zhou et al., 2023; Qin et al., 2023b; Qian et al., 2023).

However, the ability of LLMs to handle domain-specific expertise and advanced analytical tasks, such as conducting rigorous academic research, remains underexplored.

The challenge of conducting domain-specific research is particularly relevant in an era characterized by rapid knowledge expansion across multiple fields. Traditional methods for composing academic surveys are labor-intensive, often requiring months of effort by expert researchers to synthesize relevant findings. An LLM capable of independently conducting research on topics outside its training data could bypass the need for continuous re-training and manual updates, offering a scalable and efficient solution for navigating the ever-growing body of scientific literature.

While autonomous agents have shown success in executing relatively straightforward tasks—such as online shopping or playing card games (Zhou et al., 2023; Liu et al., 2023b)—they face far greater challenges in complex tasks that demand extensive domain expertise and analytical depth. Recent developments in agentic capabilities, such as the “Deep Research” features from both Gemini and OpenAI, highlight a growing focus on multi-step research planning and the synthesis of large-scale, diverse information sources (Google, 2024; OpenAI, 2025). However, progress in systematically evaluating these agents’ capacity for rigorous research remains limited, with few standardized benchmarks designed for advanced, domain-specific scenarios.

To promote the development of research agents capable of conducting comprehensive surveys, we introduce the ResearchArena benchmark. This benchmark emphasizes academic papers due to their depth of research and structured format, qualities that are often more reliable than other sources like general web pages. The ResearchArena provides an offline environment where autonomous agents can collect and organize information for research across diverse topics. It comprises three

sub-tasks for evaluation: information discovery, information selection, and information organization. These sub-tasks emulate general methodologies used by human researchers during literature surveys.

Conducting a literature survey involves defining the scope, establishing a search protocol, and iteratively analyzing and organizing findings into a cohesive structure. Based on this process, ResearchArena introduces tasks to simulate and evaluate these stages, excluding text generation. This decision stems from the premise that the pre-writing research phase is foundational to successful article composition (Rohman, 1965). Moreover, evaluating complete articles is fraught with challenges due to variability in individual writing styles; hence, such assessments are reserved for future work.

For information discovery, LLMs identify and retrieve academic papers relevant to a designated research topic by navigating vast scholarly corpora. Information selection challenges LLMs to critically assess the relevance and impact of these papers, prioritizing significant contributions. As a bonus task, information organization requires LLMs to synthesize selected research into structured knowledge representations, such as mind maps, to highlight key insights and relationships within the topic.

Preliminary evaluations reveal that LLMs underperform compared to simpler keyword-based search methods in tasks requiring analytical depth. For example, using survey titles as retrieval queries consistently yields superior recall and precision compared to LLM-driven information discovery and selection tasks. Additionally, under the task of information organization, LLMs face challenges in constructing coherent structures without the oracle guidance, underscoring the need for improvements in organizational and analytical capabilities.

The dataset supporting ResearchArena comprises 12M full-text academic papers and 7.9K survey papers, curated from the Semantic Scholar Open Research Corpus (S2ORC) (Lo et al., 2019). This corpus ensures scholarly rigor and relevance, offering a robust foundation for benchmarking LLM performance across diverse domains. Furthermore, our released pipeline supports weekly updates from Semantic Scholar, which enables evaluations to incorporate recent advancements, ensuring ongoing relevance.

2 Related Work

Previous research has employed diverse methodologies to compile datasets featuring academic survey papers. For instance, BigSurvey dataset (Liu et al., 2023a) aggregates over 7K survey papers from arXiv and includes approximately 434K references from Microsoft Academic Service and Semantic Scholar. This dataset underwent extensive preprocessing by removing duplicates, unprocessable files, and normalizing text. On the other hand, Surfer100 dataset (Li et al., 2021) includes 100 surveys emulating Wikipedia page structures, compiled by eight annotators who summarized content from web pages. Each survey contains predefined sections such as Introduction, History, Key Ideas, Variations, and Applications, summarized concisely in 50 to 150 words.

The BigSurvey dataset provides references in an abstract-only format, offering a concise overview of documents. Surfer100 utilizes Google search results to compile references for each survey topic, reflecting a broad spectrum of web-based information. In contrast, our dataset emphasizes full-text academic papers for a deeper understanding and leverages bibliographic references from original survey papers for enhanced authority and accuracy.

A closely aligned task for LLM agents in prior research involves generating Wikipedia articles. Liu et al. (2018) proposed a method for generating English Wikipedia articles by framing the task as a multi-document summarization challenge. Their approach employs a combination of extractive and abstractive summarization techniques, identifying salient information using methods such as TF-IDF and TextRank (Mihalcea and Tarau, 2004). Similarly, Shao et al. (2024) introduced the STORM system, which tackles pre-writing challenges such as research and outline preparation. STORM enhances the article generation process by simulating multi-perspective conversations, wherein an LLM poses questions and aggregates responses from reliable sources to develop detailed outlines.

Additionally, recent work by Gemini (Google, 2024) and OpenAI (OpenAI, 2025) has explored “Deep Research” features that rely on multi-step planning, iterative web browsing, and extended context windows to generate comprehensive reports. These agentic systems aim to automate information discovery and summarization by dynamically adjusting their search strategies and synthesizing insights into structured outputs.

Table 1: Summary of the dataset composition, including the counts of full-text accessible papers, survey papers, and extracted mind-maps.

Category	Count
Accessible Papers	12,034,505
Survey Papers	7,952
Extracted Mind-Maps	1,884

3 Collection Methodology

This section describes the multi-stage methodology for assembling the dataset of academic surveys, which includes survey selection, reference linking, and mind-map extraction. Each stage has been designed to ensure the relevance and utility while addressing potential limitations in automation, domain variability, and data access. The final output is a structured dataset that facilitates benchmarking the autonomous research challenge. The details for the prompts used in each stage of the collection process can be found in Appendix A.

3.1 Survey Selection

Survey selection is the foundational step in constructing the dataset, focused on identifying academic papers that provide organized overviews of specific research topics. This process involved leveraging the S2ORC corpus, which contains over 80 million academic articles in machine-readable format. The selection process combined automated filtering and human evaluation to balance scale and accuracy.

Initially, survey papers were identified by filtering for titles containing the term “survey.” While this heuristic served as an accessible baseline, it introduced potential biases, such as the exclusion of relevant papers that do not explicitly use the keyword in their titles. For example, in fields like medicine, the terms “systematic review” or “review” are more common and were largely overlooked. Recognizing these limitations, we further refined our selection using GPT-4 to analyze the titles and abstracts of candidate papers. GPT-4 was prompted to evaluate whether each paper met pre-defined criteria for surveys, such as presenting a comprehensive overview of a field.

Through this two-stage approach, approximately 85% of papers initially flagged by keyword filtering were excluded after GPT-4 evaluation. To validate this methodology, we conducted a manual inspection of a random sample of 100 papers from

the final collection, achieving a 94% precision in identifying relevant surveys. The details of this inspection are provided in Appendix C. Although this method cannot guarantee perfect recall, we believe it sufficiently represents the broader distribution of survey literature in various domains. Additionally, this stage prioritized full-text accessibility within S2ORC to ensure the inclusion of rich contextual details, reducing the corpus size to approximately 12 million documents.

3.2 Reference Linking

The second stage of the methodology involved extracting and linking bibliographic references cited in the identified survey papers. This step is critical for evaluating tasks related to information discovery and selection, as it connects surveys to their foundational sources. Reference data were sourced directly from the S2ORC corpus, which includes pre-resolved bibliographic metadata.

Despite the robustness of S2ORC’s reference extraction capabilities, several challenges emerged, including missing references or misclassified citation structures. Surveys without detectable bibliographic sections—often due to formatting issues in the source data—were excluded, resulting in the removal of 406 survey papers. Additionally, 1,635 surveys were discarded because they lacked accessible references, rendering them unsuitable for downstream evaluations.

It is important to acknowledge that the reference linking process—like human citation practices—is inherently imperfect. Even expert researchers may unintentionally omit relevant works or introduce redundancies. Similarly, our automated approach provides an approximation that, while robust, does not guarantee perfect recall of influential references. To mitigate this limitation, we applied a supervised classification model inspired by [Valenzuela et al. \(2015\)](#) to distinguish influential from non-influential citations, ensuring that the most impactful references were prioritized.

Moreover, publication dates were annotated for each retained reference, with conservative imputation for missing months or days to minimize information leakage in temporal evaluations. While these efforts improve the utility and reliability of the dataset, we recognize that no methodology can fully account for all relevant literature. Future enhancements, including integrating domain-specific heuristics and engaging human annotators, may further refine this process.

3.3 Mind-Map Extraction

Mind-map extraction is positioned as a bonus task within the benchmark, complementing the primary objectives of survey selection and reference linking. While mind-maps are not commonly found in academic survey papers, they provide valuable hierarchical visualizations of knowledge when present, offering an organized perspective on the topics.

Limited by the text-only nature of the S2ORC corpus, we extended our dataset by sourcing figure-caption pairs from the Semantic Scholar website, specifically targeting surveys with accessible figures. Using GPT-4, figures and their captions were analyzed to identify those likely representing mind-maps. Relevant figures were converted into JSON-encoded hierarchical structures, preserving their organizational logic, as illustrated in Figure 1.

This task employed a two-step verification process: first, determining if the figure represented a valid taxonomy, and second, assessing its relevance to the survey topic. After the extraction, a manual review of 100 mind-maps yielded an accuracy rate of 78% for hierarchical representation and 70% for topic relevance. The details of this inspection are provided in Appendix C. While these scores highlight the limitations of automated extraction and domain variability, they underscore the utility of mind-maps as an auxiliary dataset feature for future exploratory research.

3.4 Dataset Access

To ensure compliance, we provide tools and code that enable users to reproduce the dataset using publicly accessible S2ORC. This approach avoids direct distribution of the corpus while empowering researchers to generate reproducible datasets tailored to their specific needs.

Users must independently verify licensing requirements for the underlying data sources, as open access does not inherently guarantee permissive redistribution rights. Using February 06, 2024 release of S2ORC, the dataset itself consists of approximately 12 million full-text academic papers, including 7,952 survey papers and 1,884 extracted mind-maps.

4 Analysis

This section details the makeup of our dataset in terms of disciplinary diversity, reference coverage, and the structural complexity of derived typologies,

reflecting on how these factors contribute to the robustness and applicability across various domains.

Disciplinary Distribution. We classified each of the 12.0M papers in our public corpus and 7.9K survey papers by the top-5 most popular academic disciplines. This classification was based on the indexing information provided by S2ORC. Frequencies of papers per discipline were then aggregated and visualized to identify trends and imbalances. Figure 2a and 2b revealed significant disparities in the frequency of disciplines between the public corpus and the survey subset. Notably, Computer Science is the most prevalent discipline within surveys but less common in the broader corpus. This could reflect the dynamic nature of the CS field, which often necessitates comprehensive reviews to synthesize rapid advancements and emerging trends.

Reference Coverage. For each survey paper, we calculated the coverage ratio as the proportion of its references that were also available within our full-text corpus. We plotted cumulative density functions for each discipline to analyze how extensively the surveys’ references are represented in the broader corpus. As illustrated with Figure 2c, similar patterns were observed across all disciplines, where the density experienced exponential decay as the coverage increases. Approximately 17.18% of the survey subset (i.e., 1.3K survey papers) have at least 50% of their references available. This limitation is mainly attributed to copyright restrictions, where full-text is not permitted by the publisher.

Mind-Map Complexity. We analyzed the structural complexity of the mind-maps extracted from survey papers by counting the number of nodes and measuring the maximal depth. These measures provide insights into the conceptual breadth and hierarchical depth of the topics covered. The scatter plot from Figure 2d showed that typologies in general have shallow depths but a broad range of nodes, suggesting that while survey topics are extensively branched, they do not delve deeply into sub-topics. In particular, most typologies have a maximum depth ranging from 3 to 7 levels, where the coefficient of the regression line in the scatter plot is approximately 2.04.

5 Benchmark Tasks

This section presents a comprehensive overview of the benchmark tasks designed to evaluate the capabilities of research agents in discovering, selecting,

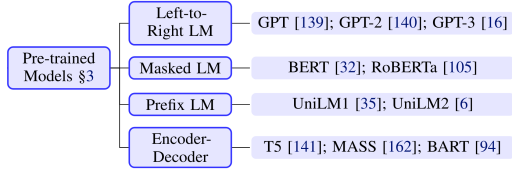


Figure 1: Mind-map extraction from the figure to its JSON representation.

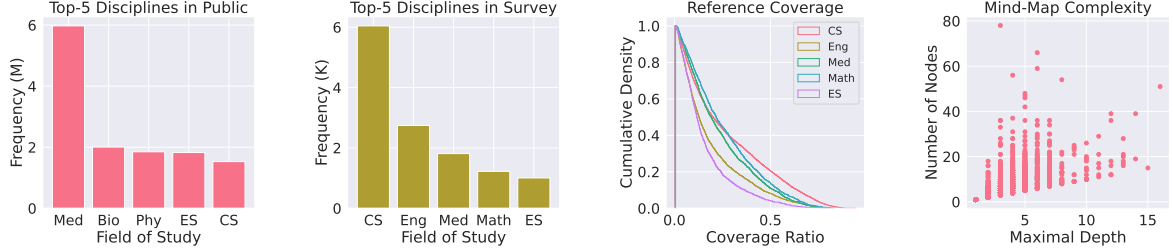


Figure 2: Dataset composition analysis with disciplinary distribution, reference coverage, and mind-map complexity. Each of these aspects is critical for benchmark evaluation. Fields of studies like Medicine (Med), Biology (Bio), Physics (Phy), Environmental Science (ES), Computer Science (CS), Engineering (Eng), and Mathematics (Math) are denoted with their abbreviations in the figures.

and organizing information. Each task targets a specific aspect of research proficiency, with rigorous constraints and evaluation metrics to ensure thorough and unbiased assessment.

Information Discovery. Provided a topic extracted from survey title, the task of information discovery requires research agents to identify a subset of documents R from a broader collection D . These documents in R should serve as supporting materials for the topic. Ideally, R should encompass all references cited in the original survey S .

However, within the collection D , there may exist another survey S' that delves into the same topic. If research agents were to use the references from S' directly, it would circumvent the need for a thorough discovery, defeating the purpose of this task. To prevent information leakage, we impose the additional constraint such that documents in D must be non-survey and published before S .

To evaluate performance, we employ standard information retrieval metrics, Recall and Precision, to measure the proportion of relevant documents successfully retrieved and the proportion of retrieved documents that are relevant. Together, these metrics determine the effectiveness and accuracy of the discovery process. For this task, the cutoff parameter K is set at 10 and 100.

Information Selection. The task of information selection requires research agents to rank the dis-

covered documents based on their importance to the topic. The labels are distinctions between influential and non-influential citations, as elaborated in Section 3. Normalized Discounted Cumulative Gain (nDCG) (Järvelin and Kekäläinen, 2002) and Mean Reciprocal Rank (MRR) (Voorhees, 1999) are used for evaluation.

These measures are crucial because conducting research involves more than merely summarizing retrieved documents; it requires the presentation of key insights from the most significant sources. Furthermore, both human researchers and autonomous agents are limited by their processing capacities. Therefore, it is essential to prioritize and focus on the most critical information first.

Information Organization (Bonus). For information organization, research agents are required to construct a hierarchical knowledge mind-map M based on R . This mind-map should provide a systematic overview of research work developed on topic T . As an intermediate step, references R from the original survey paper could be provided to the agents, who would then focus exclusively on constructing M . In contrast, for an end-to-end version, R is the set of discovered documents from the previous task.

For evaluation, two primary metrics are employed: Heading Soft Recall (Fränti and Mariescu-Istodor, 2023) and Heading Entity Recall (Shao

et al., 2024). These metrics compare the set of node labels from the original and the constructed knowledge mind-maps, referred to as A and B , respectively. To measure similarity of these labels, Heading Soft Recall (HSR) leverages SENTENCE-BERT (Reimers and Gurevych, 2019) embedding, while Heading Entity Recall (HER) employs Named Entity Recognition from FLAIR (Akbik et al., 2019) for extraction. The formal definitions for each metric are presented in Equation 1, where S is the set of labels extracted from the mind-maps.

While these metrics provide a measure of content similarity, they do not account for structural alignment. Tree Editing Distance (Zhang and Shasha, 1989) solves this concern by calculating the minimal number of operations (i.e., relabeling, deleting, and inserting nodes) required to transform one tree into another. Nonetheless, relying on Tree Editing Distance alone might overlook the potential for non-exact label matches. To address this, we propose Tree Semantic Distance, which assigns no cost to editing operations involving nodes whose cosine similarity exceeds 0.8.

$$C(S) = \sum_{i=1}^{|S|} \frac{1}{\sum_{j=1}^{|S|} \text{Sim}(S_i, S_j)}$$

$$\text{HSR}(A, B) = \frac{C(A) + C(B) - C(A \cup B)}{C(B)} \quad (1)$$

$$\text{HER}(A, B) = \frac{|\text{Ent}(A) \cap \text{Ent}(B)|}{|\text{Ent}(A)|}$$

6 Benchmarking

In this section, we present preliminary evaluations of existing techniques, describing their configurations and performance metrics. These techniques encompass both naive keyword-based methods, such as TITLE, and advanced LLM-based methods, including STORM. The exact wording of the prompts used in each baseline can be found in Appendix B.

6.1 Baselines

Information Discovery. For information discovery, research agents are equipped with retrieval tools that enable interaction with the public corpus by submitting queries to retrievers such as BM25 and BGE (Xiao et al., 2023). These agents are evaluated based on their ability to effectively leverage these tools by generating relevant queries. Since exploration is limited to previously published non-

survey literature, retrievers retry with exponential back-off until the cutoff parameter K is satisfied.

- **TITLE:** Assuming that research topics are encapsulated within survey titles, this method directly employs the title from each survey paper as a query to retrieve relevant materials that support research on the topic. It is important to note that title extraction using S2ORC exhibits variable capitalization across different documents. As a result, we normalize by converting titles to lowercase.
- **ZERO-SHOT:** Assuming that existing LLMs possess prior knowledge relevant to a survey topic, this method extends the TITLE method by instructing GPT-4 to derive a query from the survey title. This approach leverages the inherent capabilities of LLMs to generate more sophisticated and contextually appropriate queries.
- **DECOMPOSER:** As discovered by Tushar et al. (Khot et al., 2022), decomposed prompting is more effective when individual reasoning steps of a task are difficult to learn. This principle is applicable to our case, as a survey topic may consist of multiple sub-topics, making it challenging to directly generate a single query that retrieves all relevant papers. Consequently, we instruct GPT-4 to first deconstruct the research topic into several sub-questions. Each sub-question then generates a corresponding sub-query. These sub-queries are retrieved in batches, and the results are amalgamated using reciprocal rank fusion (Cormack et al., 2009).
- **SELF-RAG:** As proposed by Asai et al. (Asai et al., 2023), SELF-RAG adaptively retrieves passages on demand and utilizes reflection tokens to determine which retrieved documents are relevant to the instruction, thus continuing the generation based on the pertinent information. It serves as an enhanced version of ZERO-SHOT, where the model is instructed to generate a query from the topic. Because the model refines its final query generation based on the discovered information from intermediate retrievals, it operates as a research agent.
- **STORM:** As presented in Section 2, STORM conducts research through multi-perspective

Table 2: Baseline performance on discovery task, evaluated with Recall@10, Recall@100, Precision@10, and Precision@100, where the retrievers include BM25 and BGE.

Baseline	Recall@10		Recall@100		Precision@10		Precision@100	
	BM25	BGE	BM25	BGE	BM25	BGE	BM25	BGE
TITLE	0.0424	0.1012	0.1338	0.2697	0.0669	0.1541	0.0286	0.0586
ZERO-SHOT	0.0382	0.0832	0.1253	0.2287	0.0602	0.1232	0.0256	0.0464
DECOMPOSER	0.0434	0.0879	0.1431	0.2554	0.0717	0.1304	0.0312	0.0536
SELF-RAG	0.0380	0.0815	0.1210	0.2260	0.0595	0.1215	0.0256	0.0461
STORM	0.0281	0.0979	0.0693	0.1441	0.0446	0.1041	0.0130	0.0208

Table 3: Baseline performance on selection task, evaluated with nDCG@10, nDCG@30, nDCG@100, and MRR, where the retrievers include BM25 and BGE.

Baseline	nDCG@10		nDCG@30		nDCG@100		MRR	
	BM25	BGE	BM25	BGE	BM25	BGE	BM25	BGE
TITLE	0.0711	0.1678	0.0775	0.1754	0.0941	0.2019	0.1903	0.3816
ZERO-SHOT	0.0634	0.1346	0.0692	0.1417	0.0856	0.1657	0.1743	0.3246
DECOMPOSER	0.0735	0.1445	0.0803	0.1554	0.0986	0.1838	0.1959	0.3510
SELF-RAG	0.0627	0.1341	0.0679	0.1415	0.0837	0.1646	0.1705	0.3233
STORM	0.0445	0.1275	0.0507	0.1322	0.0524	0.1267	0.1271	0.3206

conversations to compose Wikipedia articles on particular topics from scratch. It closely resembles our scenario, except that the environment involves more rigorous academic papers. We record the retrieval history as STORM continues to probe for additional papers. Upon concluding the final round of conversations, every article within the retrieval history is considered part of the discovered information.

- **DEEP RESEARCH:** While Deep Research demonstrates advanced multi-step web-based browsing and iterative information synthesis, it does not currently support offline corpora and relies on general web pages, prompting us to include a brief case study of its capabilities.

Information Selection. For information selection, documents are ranked based on the similarity scores obtained during the discovery phase. For BGE retriever, we rely on FAISS (Johnson et al., 2019) to retrieve based on L2 distance in the embedding space. On the other hand, STORM does not explicitly rank the retrieved documents. We treat documents discovered earlier in the conversations of higher relevance.

Information Organization. For information organization, the CLUSTERING approach employs Ward’s method for hierarchical clustering on the BGE embedding of every reference article, and the final dendrogram is extracted as typology. The

label in each node is computed as the most important TF-IDF word, with ngrams ranging from 1 to 3. FEW-SHOT is achieved by providing a few random examples of extracted typologies and instructing GPT-4 to generate another topic-oriented mind-map. Lastly, the article outline generated by STORM is converted to typology, with headings and their nested sub-headings representing the hierarchy.

6.2 Evaluation Results

The baseline experiments were conducted on a single machine equipped with 8 NVIDIA RTX A6000 GPUs, 96 CPU cores, and 128GB RAM. Discussion on the performance metrics is presented below.

Information Discovery. As demonstrated in Table 2, the task of information discovery remains challenging for all baseline models. This is illustrated by the Recall@100 metric, which falls below 0.15 for BM25 and 0.27 for BGE. Moreover, agent baselines such as SELF-RAG and STORM consistently achieve the lowest rankings, irrespective of the retrievers employed. This limitation highlights the critical need for more advanced retrieval mechanisms to manage large volumes of documents effectively during information discovery.

Information Selection. The performance with information selection is presented in Table 3. The results indicate a consistent trend wherein agent baselines underperform compared to keyword-based methods. The evaluation of nDCG at various

Table 4: Baseline performance on organization task, evaluated with Heading Soft Recall, Heading Entity Recall, and Tree Semantic Distance, across intermediate and end-to-end conditions.

Oracle	Baseline	Heading Soft Recall (\uparrow)	Heading Entity Recall (\uparrow)	Tree Semantic Distance (\downarrow)
Yes	CLUSTERING	0.6074	0.2104	45.69
	STORM	0.7325	0.3098	60.04
No	FEW-SHOT	0.8408	0.2446	49.83
	STORM.BM25	0.7940	0.2938	66.65
	STORM.BGE	0.7842	0.2693	65.93

levels of document retrieval, such as nDCG@10, nDCG@30, and nDCG@100, provides a quantitative assessment of the ranking performance. Notably, for the TITLE method using the BGE retriever, the nDCG@100 score is 0.2019, which significantly surpasses the score of STORM, which stands at 0.1267. Improvements during the information discovery phase have the potential to enhance overall performance in the selection phase, as evidenced by DECOMPOSER, which ranks the second behind TITLE in discovery and selection tasks.

Information Organization. The evaluation on task of information organization under intermediate (i.e., with oracle) and end-to-end (i.e., without oracle) conditions are documented in Table 4. Notably, the metrics exhibit discrepancies across each other, which contrasts with the uniformity observed in previous discovery and selection tasks. This divergence is expected due to the distinct nature of the metrics: Heading Soft Recall and Heading Entity Recall assess content similarity, whereas Tree Semantic Distance evaluates structural alignment.

In the intermediate version, where references are provided to LLMs, the proportion of correctly included entities, as measured by Heading Entity Recall, is slightly higher. Specifically, STORM achieved a recall rate of 0.3098, outperforming the end-to-end condition. Conversely, when it comes to constructing the hierarchy, CLUSTERING outperforms advanced LLM-based agents, as evidenced by its attainment of the lowest Tree Semantic Distance of 45.69 among all baseline methods.

Deep Research. We conducted a brief investigation on two distinct topics, transfer learning and LiDAR scanning mechanisms, using Gemini Deep Research, as shown in Appendix D. Due to its sole reliance with online resources, we don’t have a direct quantitative comparison with other baselines. Nevertheless, the generated summaries illustrate Gemini Deep Research’s capacity to synthesize diverse online sources into coherent find-

ings, highlighting its potential to support high-level exploration when specialized academic databases are not immediately required.

7 Conclusion

In conclusion, ResearchArena introduces a rigorous benchmark designed to evaluate LLMs in conducting research surveys on designated topics. By systematically decomposing the survey process into distinct tasks like information discovery, selection, and organization, this benchmark provides a detailed framework for evaluating autonomus research agents. Our findings underscore the potential of LLMs to revolutionize academic research, provided that future advancements can bridge the existing performance gaps. Grounded in Semantic Scholar Open Research Corpus, this work establishes a robust foundation for the future, aiming to improve the ability of LLMs to autonomously conduct expertise-level, domain-specific research.

8 Limitations

Despite the robust framework and extensive dataset provided by ResearchArena, this study has several limitations. Firstly, the offline environment, though comprehensive, may not accurately represent the dynamic and interconnected nature of live databases and the internet. This discrepancy could potentially limit the applicability of the findings in real-world research settings. Additionally, due to copyright constraints, not every full-text reference of the survey papers could be included. This omission could affect the comprehensive understanding of the survey topics under investigation. Finally, there is no evaluation on text generation but mostly the surveying process. However, even if this is just the first step of conducting research, LLM agents have already shown deficiencies. Future iterations of ResearchArena should address this issue, particularly as these agents improve.

References

- Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019. [FLAIR: An easy-to-use framework for state-of-the-art NLP](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59, Minneapolis, Minnesota. Association for Computational Linguistics.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multi-task, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
- Gordon V Cormack, Charles LA Clarke, and Stefan Buettcher. 2009. Reciprocal rank fusion outperforms condorcet and individual rank learning methods. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 758–759.
- Pasi Fränti and Radu Marinescu-Istodor. 2023. [Soft precision and recall](#). *Pattern Recognition Letters*, 167:115–121.
- Google. 2024. Try deep research and our new experimental model in gemini, your AI assistant. <https://blog.google/products/gemini/google-gemini-deep-research/>. Accessed: February 10, 2025.
- Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of ir techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4):422–446.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Xiangji Huang. 2023. A systematic study and comprehensive evaluation of chatgpt on benchmark datasets. *arXiv preprint arXiv:2305.18486*.
- Irene Li, Alexander Fabbri, Rina Kawamura, Yixin Liu, Xiangru Tang, Jaesung Tae, Chang Shen, Sally Ma, Tomoe Mizutani, and Dragomir Radev. 2021. Surfer100: Generating surveys from web resources, wikipedia-style. *arXiv preprint arXiv:2112.06377*.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.
- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*.
- Shuaiqi Liu, Jiannong Cao, Ruosong Yang, and Zhiyuan Wen. 2023a. Generating a structured summary of numerous academic papers: Dataset and method. *arXiv preprint arXiv:2302.04580*.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuan-yu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023b. Agent-bench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*.
- Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Dan S Weld. 2019. S2orc: The semantic scholar open research corpus. *arXiv preprint arXiv:1911.02782*.
- Rada Mihalcea and Paul Tarau. 2004. Textrank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pages 404–411.
- OpenAI. 2025. Introducing deep research. <https://openai.com/index/introducing-deep-research/>. Accessed: February 10, 2025.
- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*.
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023a. Is chatgpt a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476*.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023b. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- D Gordon Rohman. 1965. Pre-writing: The stage of discovery in the writing process. *College Composition & Communication*, 16(2):106–112.

Corby Rosset, Ho-Lam Chung, Guanghui Qin, Ethan C Chau, Zhuo Feng, Ahmed Awadallah, Jennifer Neville, and Nikhil Rao. 2024. Researchy questions: A dataset of multi-perspective, compositional questions for llm web agents. *arXiv preprint arXiv:2402.17896*.

Yijia Shao, Yucheng Jiang, Theodore A Kanell, Peter Xu, Omar Khattab, and Monica S Lam. 2024. Assisting in writing wikipedia-like articles from scratch with large language models. *arXiv preprint arXiv:2402.14207*.

Marco Valenzuela, Vu Ha, and Oren Etzioni. 2015. Identifying meaningful citations. In *Workshops at the twenty-ninth AAAI conference on artificial intelligence*.

EM Voorhees. 1999. Proceedings of the 8th text retrieval conference. *TREC-8 Question Answering Track Report*, pages 77–82.

Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. 2024. *OpenHands: An Open Platform for AI Software Developers as Generalist Agents*. *Preprint*, arXiv:2407.16741.

Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighof. 2023. C-pack: Packaged resources to advance general chinese embedding. *arXiv preprint arXiv:2309.07597*.

Kaizhong Zhang and Dennis Shasha. 1989. Simple fast algorithms for the editing distance between trees and related problems. *SIAM journal on computing*, 18(6):1245–1262.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.

A Prompts for the Dataset Collection

Instructions to analyze whether an academic paper fits under the research survey category.

```
"""
The point of a survey paper is to provide an
organized view on the current state of the
field. If it relies heavily on external
information, such as the results of a
population questionnaire, do not include
it. Using the above criteria, is the
following article a survey paper? Respond
either "True" or "False".
"""
```

Instructions to extract the survey mind-maps into JSON-encoded representations.

```
"""
Identify the figure that most likely
illustrates a taxonomy or overview. Your
response should be limited to the filename,
or NULL if not found. The provided figure
presents a hierarchy. Extract as
JSON-encoded tree whose children are
NULL-terminated.
"""
```

B Prompts for the Experiments

Instructions used by DECOMPOSER for the information discovery task, adopted from the Researchy Questions by Rosset et al. (2024).

```
"""
### Below is an example on how to decompose a
complex question into sub-questions and
search queries.

Question: should the death penalty be legalized?

<Decomposition>
- What are the arguments in favor of the
  death penalty?
- Does the death penalty serve as a
  deterrent to crime?
- Is the death penalty a just punishment
  for certain crimes?
- How does the death penalty compare to
  other forms of punishment in terms
  of cost and effectiveness?
- What are the arguments against the death
  penalty?
- What is the risk of executing innocent
  people with a death penalty?
- Are there any ethical concerns
  surrounding the death penalty?
- To what extent is the death penalty
  applied fairly and without bias?
- In practice, how expensive is the
  death penalty?
- What is the current legal status of the
  death penalty in various jurisdictions?
- In which countries or states is the
  death penalty currently legal?
- What are the trends in death penalty
  legislation and public opinion?
- What are the alternatives to the death
  penalty?
- How effective are alternative
  punishments to the death penalty,
  e.g. life imprisonment?
- What are the costs and benefits of
  alternatives to the death penalty?
- How do the pros and cons of the death
  penalty compare to its alternatives?
</Decomposition>

<Queries>
- arguments in favor of the death penalty
- death penalty as a deterrent to crime
- death penalty as a just punishment
- death penalty cost and effectiveness
  comparison
- arguments against the death penalty
```

```

- risk of executing innocent people with
  death penalty
- ethical concerns surrounding the death
  penalty
- fairness and bias in death penalty
  application
- current legal status of the death penalty
  worldwide
- trends in death penalty legislation and
  public opinion
- alternatives to the death penalty
- effectiveness of life imprisonment without
  parole
- costs and benefits of death penalty
  alternatives
</Queries>

Question: {x}

### Instructions:

1. What sub-questions do I need to know in
  order to fully understand and answer the
  above Question.
- Format your response as a bullet-point
  style outline of questions and
  sub-questions in the <Decomposition>
  tag.
- Order your sub-questions such that one
  question comes after another if it
  needs to use the answer to the previous
  one.
- Do not ask unnecessary or tangential
  sub-questions, only those that are
  critical to finding important
  information.

2) Next, write a list of search queries that
  would likely lead to results addressing all
  the sub-questions.
- Enumerate your queries in a bullet-point
  style list inside the <Queries> tag.

You may refer to the example above for guidance.
"""

```

Instructions used by ZERO-SHOT and SELF-RAG for the information discovery task.

```

"""
Create a search query that gathers supporting
materials for writing a survey paper on the
following topic: {x}.
"""

```

Instructions used by FEW-SHOT for the information organization task.

```

"""
### Examples

<topic>
A Survey on LiDAR Scanning Mechanisms
</topic>

<typology>
{"Opto-Mechanical Beam Deflection Mechanisms":
  {"Line Scanner": {"Slanted Plain Mirror":
    null, "Off-axis Parabolic Mirror": null,
    "Polygon Mirror": null}, "Area Scanner":
    {"Single Galvanometer Scanning Mirror":

```

```

    null, "Double Galvanometer Scanning
    Mirror": null, "Gyroscopic Mirror": null,
    "Risley Scanner": null}}}
</typology>

<topic>
A Survey on Large Language Models for
Recommendation
</topic>

<typology>
{"LLM4Rec": {"Discriminative LLM4Rec":
  {"Fine-tuning": {"Prompt Tuning": null}},
  "Generative LLM4Rec": {"Non-tuning":
    {"Prompting": null, "In-context Learning":
      null}, "Tuning": {"Fine-tuning": null,
        "Prompt Tuning": null, "Instruction
        Tuning": null}}}}
</typology>

### Instructions

- Provided a topic, your task is to construct a
  mind-map style typology that presents a
  systematic understanding of the topic.
- Put your JSON-encoded response in the tag
  `<typology>...</typology>`. You may refer
  to the examples above for guidance.

<topic>
{x}
</topic>
"""

```

C Inspection over Survey Passages

A manual inspection of survey passages is conducted to assess the effectiveness of GPT-4 in identifying survey papers. Following this, the extracted mind-map is analyzed and compared against the original paper to determine its accuracy and relevance in capturing key aspects of the topic under discussion. Details are presented in Table 5 and 6.

D Case Study with Gemini Deep Research

Generated report on the topic of LiDAR scanning mechanisms with sources attached in Figure 4.

```

# **LiDAR Scanning Mechanisms: A Comprehensive
  Overview**

LiDAR (Light Detection and Ranging) technology
has revolutionized various industries, from
autonomous vehicles and robotics to
surveying and mapping. At the heart of
LiDAR systems lie the scanning mechanisms
that enable these systems to capture
precise 3D spatial information about the
surrounding environment. This article
delves into the different types of LiDAR
scanning mechanisms, exploring their
working principles, advantages,
disadvantages, applications, and the
companies that manufacture them. We will
also examine the latest advancements and
future trends in LiDAR scanning technology.

```

1030		### **Disadvantages**	1100
1031	## **Types of LiDAR Scanning Mechanisms**		1101
1032			1102
1033	LiDAR scanning mechanisms can be broadly	* **Moving Parts:** The presence of moving	1103
1034	classified into three main categories:	parts makes mechanical LiDAR systems	1104
1035		susceptible to wear and tear, potentially	1105
1036	* **Mechanical:** These systems utilize moving	affecting their reliability and lifespan ² .	1106
1037	parts, such as rotating mirrors, to steer	* **Bulkiness:** Mechanical systems tend to be	1107
1038	the laser beam and scan the environment ¹ .	larger and heavier than solid-state LiDAR,	1108
1039	* **Solid-state:** These systems have no	making them less suitable for applications	1109
1040	macroscopic moving parts and employ	where size and weight are critical ⁵ .	1110
1041	techniques like MEMS mirrors or optical	* **Cost:** The complexity of the mechanical	1111
1042	phased arrays to steer the laser beam	design can contribute to higher	1112
1043	electronically ¹ .	manufacturing costs ⁶ .	1113
1044	* **Flash:** These systems illuminate the	* **Environmental Factors:** LiDAR signals can	1114
1045	entire scene with a single laser pulse,	be attenuated or scattered by fog, rain,	1115
1046	capturing a 3D point cloud in an instant,	and snow, which can limit their	1116
1047	similar to a camera ¹ .	effectiveness in adverse weather	1117
1048		conditions ⁷ . Additionally, LiDAR technology	1118
1049	## **Mechanical LiDAR**	faces challenges in detecting and measuring	1119
1050		certain types of surfaces, particularly	1120
1051	Mechanical LiDAR, the oldest type of LiDAR,	non-reflective or highly absorbing ones.	1121
1052	employs a rotating assembly with mirrors or	For instance, black asphalt, dark-colored	1122
1053	prisms to direct the laser beam in a	objects, or water surfaces may not reflect	1123
1054	360-degree scan ¹ . This mechanism allows for	sufficient energy from the laser pulses,	1124
1055	a wide field of view and high resolution,	making them difficult to detect or measure	1125
1056	making it suitable for applications that	accurately ⁷ .	1126
1057	require detailed 3D mapping ² .		1127
1058		## **Solid-State LiDAR**	1128
1059	### **How it Works**		1129
1060		Solid-state LiDAR eliminates the need for	1130
1061	A mechanical LiDAR system typically consists of	macroscopic moving parts, offering improved	1131
1062	a laser emitter, a rotating mirror or prism	durability and reliability compared to	1132
1063	assembly, and a detector. The laser emits	mechanical systems ² . These systems utilize	1133
1064	short pulses of light, which are directed	various technologies to steer the laser	1134
1065	by the rotating assembly towards the target	beam electronically, including MEMS	1135
1066	area. The detector measures the time it	(Microelectromechanical Systems) and OPA	1136
1067	takes for the light to return, allowing the	(Optical Phased Array).	1137
1068	system to calculate the distance to the		1138
1069	object ⁴ . The two kinds of LiDAR detection	### **How it Works**	1139
1070	schemes are "incoherent" or direct energy		1140
1071	detection (which principally measures	#### **MEMS LiDAR**	1141
1072	amplitude changes of the reflected light)		1142
1073	and coherent detection (best for measuring	MEMS LiDAR works by directing a single laser	1143
1074	Doppler shifts, or changes in the phase of	beam to a tiny mirror that can be tilted or	1144
1075	the reflected light) ³ . Coherent systems	rotated to scan the environment ³ . The size	1145
1076	generally use optical heterodyne detection.	of the MEMS mirror is a critical factor in	1146
1077	This is more sensitive than direct	its performance ⁸ . A larger mirror allows	1147
1078	detection and allows them to operate at	for more photons to be emitted, increasing	1148
1079	much lower power, but requires more complex	the chances of a sufficient number of	1149
1080	transceivers. Both types employ pulse	photons returning to the detector for	1150
1081	models: either micropulse or high energy.	object detection. However, the mirror must	1151
1082	Micropulse systems utilize intermittent	also be large enough to deflect all the	1152
1083	bursts of energy.	light collimated by the lens used to focus	1153
1084		the laser beam. This ensures high	1154
1085	### **Advantages**	resolution and accurate identification of	1155
1086		even small objects ⁸ . While some MEMS	1156
1087	* **Wide Field of View:** Mechanical LiDAR can	systems operate in a single plane, others	1157
1088	achieve a 360-degree horizontal field of	can achieve 2D scanning with dual-axis	1158
1089	view, providing comprehensive coverage of	mirrors or multiple lasers ³ .	1159
1090	the surrounding environment ⁵ .		1160
1091	* **High Resolution:** The rotating mechanism	#### **OPA LiDAR**	1161
1092	allows for precise control over the laser		1162
1093	beam, enabling high-resolution data	OPA LiDAR, on the other hand, uses an array of	1163
1094	capture ² .	optical antennas to create a beam that can	1164
1095	* **Long Range:** Mechanical LiDAR systems can	be steered electronically by controlling	1165
1096	achieve long-range measurements, making	the phase of the light emitted from each	1166
1097	them suitable for applications like aerial	antenna ¹ .	1167
1098	surveying and mapping ² .		1168
1099		### **MEMS vs. Polygon Scanners**	1169

1170	MEMS mirrors and polygon scanners are both used	integration into small devices1.	1240
1171	in LiDAR systems for beam steering, but		1241
1172	they have distinct characteristics and	### **Disadvantages**	1242
1173	trade-offs9. MEMS mirrors are smaller and		1243
1174	potentially more cost-effective, but they	* **Limited Range:** Flash LiDAR typically has	1244
1175	can be more susceptible to vibrations and	a shorter range compared to scanning LiDAR	1245
1176	temperature variations. Polygon scanners,	systems1.	1246
1177	with their larger size and more robust	* **Lower Resolution:** The resolution of flash	1247
1178	design, offer higher accuracy and longer	LiDAR can be lower than scanning LiDAR,	1248
1179	range, but they can be more expensive and	especially at longer distances1.	1249
1180	less compact.	* **Eye Safety:** The high-power laser pulses	1250
1181		used in flash LiDAR can pose eye safety	1251
1182	### **Advantages**	concerns, requiring careful design and	1252
1183		implementation7. The use of eye-safe	1253
1184	* **Durability:** The absence of macroscopic	wavelengths, such as 1550 nm, is one	1254
1185	moving parts makes solid-state LiDAR more	approach to mitigate this concern3.	1255
1186	robust and less prone to mechanical		1256
1187	failure2.	## **Hybrid LiDAR**	1257
1188	* **Compact Size:** Solid-state LiDAR systems		1258
1189	are typically smaller and lighter than	Hybrid LiDAR systems combine elements of both	1259
1190	mechanical systems, making them suitable	solid-state and mechanical LiDAR	1260
1191	for applications where space is limited6.	technologies to optimize performance and	1261
1192	* **Faster Scanning:** Electronic beam steering	address the limitations of each approach2.	1262
1193	allows for faster scanning speeds compared		1263
1194	to mechanical systems10.	### **How it Works**	1264
1195			1265
1196	### **Disadvantages**	Hybrid LiDAR systems may use a combination of	1266
1197		MEMS mirrors and rotating elements to	1267
1198	* **Limited Field of View:** Solid-state LiDAR	achieve a wider field of view and longer	1268
1199	may have a more limited field of view	range. For example, the Hesai Pandar128, a	1269
1200	compared to mechanical systems, although	hybrid solid-state LiDAR, integrates 128	1270
1201	advancements in beam steering technology	transmit-receive modules with a 360-degree	1271
1202	are addressing this limitation5.	spinning scanning module12. This	1272
1203	* **Shorter Range:** Solid-state LiDAR	combination allows for high resolution and	1273
1204	typically has a shorter range than	a wide field of view while maintaining a	1274
1205	mechanical LiDAR, although this is also	compact design.	1275
1206	improving with advancements in technology5.		1276
1207	* **Cost:** While the cost of solid-state LiDAR	### **Advantages**	1277
1208	is decreasing, it can still be higher than		1278
1209	some mechanical systems11.	* **Balanced Performance:** Hybrid LiDAR offers	1279
1210		a balance between the accuracy and range of	1280
1211	## **Flash LiDAR**	mechanical systems and the durability and	1281
1212		compact design of solid-state systems2.	1282
1213	Flash LiDAR illuminates the entire scene with a	* **Adaptability:** It can adapt to the needs	1283
1214	single laser pulse, capturing a 3D point	of different environments by leveraging the	1284
1215	cloud instantaneously1. This approach	reliability of solid-state components and	1285
1216	eliminates the need for scanning	the detailed scanning capabilities of	1286
1217	mechanisms, resulting in a simpler and	mechanical systems2.	1287
1218	potentially more cost-effective design.		1288
1219		### **Disadvantages**	1289
1220	### **How it Works**		1290
1221		* **Complexity:** Hybrid systems can be more	1291
1222	Flash LiDAR systems use a wide-beam laser to	complex to design and manufacture compared	1292
1223	illuminate the entire field of view. The	to purely mechanical or solid-state	1293
1224	reflected light is then captured by a	systems13.	1294
1225	detector array, which measures the time of	* **Cost:** The combination of technologies can	1295
1226	flight for each pixel in the array. This	lead to higher costs compared to some	1296
1227	allows the system to generate a 3D image of	single-technology systems13.	1297
1228	the scene in a single flash4.		1298
1229		## **Cost and Availability of LiDAR Systems**	1299
1230	### **Advantages**		1300
1231		The cost and availability of LiDAR systems vary	1301
1232	* **No Moving Parts:** Flash LiDAR has no	significantly depending on several factors,	1302
1233	moving parts, making it highly durable and	including performance specifications,	1303
1234	reliable1.	application requirements, durability, and	1304
1235	* **Fast Data Acquisition:** It captures an	integration needs14.	1305
1236	entire scene in a single flash, enabling		1306
1237	rapid data acquisition2.	* **Performance Specifications:** Higher	1307
1238	* **Compact Size:** Flash LiDAR systems can be	performance systems with greater range,	1308
1239	very compact, making them suitable for	resolution, and speed are generally more	1309

1310	expensive ¹⁴ . For example, automotive-grade	* **Archaeology:** LiDAR helps uncover ancient	1380
1311	LiDAR, with its long-range and	structures and features hidden beneath	1381
1312	high-resolution requirements, can be	vegetation or over time ¹⁷ .	1382
1313	significantly more costly than LiDAR used	* **Disaster Management:** LiDAR is used to	1383
1314	for general surveying ¹⁴ .	assess damage after natural disasters and	1384
1315	* **Application Requirements:** LiDAR systems	aid in rescue efforts ¹⁷ .	1385
1316	designed for demanding applications, such		1386
1317	as autonomous vehicles, are typically more	## **Latest Advancements in LiDAR Scanning	1387
1318	expensive due to their stringent	Mechanisms**	1388
1319	performance and reliability requirements ¹⁴ .		1389
1320	* **Durability and Reliability:** LiDAR systems	Research and development in LiDAR technology	1390
1321	built for harsh environments require robust	continue to advance, leading to innovative	1391
1322	designs and materials, which can increase	solutions that improve performance, reduce	1392
1323	costs ¹⁴ .	costs, and expand applications. Some of the	1393
1324	* **Integration and Maintenance:** Systems that	latest advancements include:	1394
1325	are easier to integrate and maintain can be		1395
1326	more cost-effective in the long run ¹⁴ .	* **Elastic LiDAR:** This technology utilizes	1396
1327		elastic scattering of light to achieve	1397
1328	The cost of LiDAR systems has been decreasing	precise measurements and environmental	1398
1329	in recent years due to advancements in	sensing ¹⁸ .	1399
1330	technology and manufacturing ¹⁵ . However,	* **Silicon-based LiDAR:** This approach	1400
1331	high-performance systems can still be	leverages silicon photonics to create	1401
1332	expensive. For example, a robust,	compact and cost-effective LiDAR systems.	1402
1333	entry-level LiDAR system for drone	Silicon-based LiDAR is an ideal way to	1403
1334	applications can cost around \$23,000, while	reduce the volume of the LiDAR and realize	1404
1335	a drone and associated accessories can add	monolithic integration. It removes the	1405
1336	another \$10,000 to \$26,000 to the total	moving parts in the conventional device and	1406
1337	cost ¹⁶ .	realizes solid-state beam steering ¹⁹ .	1407
1338		* **Advanced Simulation Tools:** Tools like	1408
1339	## **LiDAR vs. Radar**	VPItransmissionMaker and VPIcomponentMaker	1409
1340		enable detailed design and simulation of	1410
1341	While both LiDAR and radar are remote sensing	LiDAR systems, including modeling	1411
1342	technologies used for object detection and	atmospheric conditions and photonic	1412
1343	mapping, they have distinct characteristics	integrated circuits ²⁰ .	1413
1344	and strengths ⁷ . LiDAR uses laser light to	* **Smart Corner Solution:** Marelli's Smart	1414
1345	measure distances and create	Corner solution integrates LiDAR sensors	1415
1346	high-resolution 3D maps, while radar uses	into vehicle headlamps and grilles,	1416
1347	radio waves to detect objects and measure	addressing sensor placement and field of	1417
1348	their speed and direction. LiDAR offers	view challenges ²⁰ .	1418
1349	higher accuracy and resolution, especially	* **Deep Learning for Drone Detection:** Recent	1419
1350	for detailed mapping and object	research has demonstrated the effectiveness	1420
1351	recognition, but it can be more sensitive	of deep learning algorithms in processing	1421
1352	to adverse weather conditions. Radar, on	LiDAR data for drone detection ²¹ . Studies	1422
1353	the other hand, is less affected by weather	have shown high accuracy (above 97%) in	1423
1354	and can operate at longer ranges, but it	detecting drones even in noisy environments	1424
1355	typically provides lower resolution data.	using deep learning models ²¹ .	1425
1356			1426
1357	## **Applications of LiDAR**	## **The Future of LiDAR Scanning Mechanisms**	1427
1358			1428
1359	LiDAR technology has a wide range of	The future of LiDAR scanning mechanisms is	1429
1360	applications across various industries,	likely to be shaped by several key trends:	1430
1361	including:		1431
1362		* **Continued Miniaturization:** Solid-state	1432
1363	* **Autonomous Vehicles:** LiDAR is a crucial	and flash LiDAR technologies will continue	1433
1364	sensor in self-driving cars, enabling them	to shrink in size, enabling their	1434
1365	to perceive their surroundings, navigate	integration into smaller devices and	1435
1366	autonomously, and avoid obstacles ¹⁷ .	expanding their applications ¹⁵ .	1436
1367	* **Robotics:** LiDAR helps robots navigate,	* **Improved Performance:** Advancements in	1437
1368	map their environment, and interact with	laser technology, detectors, and scanning	1438
1369	objects ¹⁷ .	mechanisms will lead to higher resolution,	1439
1370	* **Surveying and Mapping:** LiDAR is used to	longer range, and faster scanning speeds ¹⁵ .	1440
1371	create detailed 3D maps of terrain,	* **Reduced Costs:** Innovations in	1441
1372	infrastructure, and urban environments ¹⁷ .	manufacturing and design will drive down	1442
1373	* **Environmental Monitoring:** LiDAR is used	the cost of LiDAR systems, making them more	1443
1374	to monitor vegetation, assess natural	accessible for various applications ¹⁵ .	1444
1375	disaster impacts, and map forests ¹⁷ .	* **Increased Integration:** LiDAR systems will	1445
1376	* **Urban Planning and Development:** LiDAR	be increasingly integrated with other	1446
1377	provides data for urban planning,	sensors and technologies, such as cameras	1447
1378	infrastructure development, and	and radar, to create more comprehensive and	1448
1379	transportation management ¹⁷ .	robust sensing solutions ²² .	1449

Conclusion

LiDAR scanning mechanisms are fundamental to the operation and application of LiDAR technology. Each type of scanning mechanism, whether mechanical, solid-state, flash, or hybrid, offers a unique set of advantages and disadvantages. Mechanical systems excel in providing a wide field of view and high resolution, making them suitable for applications like aerial surveying and mapping. Solid-state LiDAR, with its MEMS and OPA variants, offers durability, compact size, and faster scanning speeds, making it increasingly popular for autonomous vehicles and robotics. Flash LiDAR, while limited in range and resolution, provides rapid data acquisition and is suitable for short-range applications. Hybrid systems combine the strengths of different approaches to optimize performance for specific needs.

The choice of LiDAR scanning mechanism depends on the specific application requirements, considering factors such as range, resolution, field of view, cost, and environmental conditions. As LiDAR technology continues to evolve, driven by advancements in miniaturization, performance improvement, and cost reduction, we can expect to see even more innovative scanning mechanisms emerge, further expanding the applications and impact of this transformative technology.

Generated report on the topic of transfer learning with sources attached in Figure 4.

Transfer Learning: A Comprehensive Overview

1\. Introduction

Transfer learning, a powerful machine learning technique, has emerged as a game-changer in the field of artificial intelligence (AI). It leverages the knowledge gained from solving one problem to improve performance on a different but related problem. This approach has gained significant traction in recent years, particularly in deep learning, where training models from scratch can be computationally expensive and require massive amounts of data. By transferring knowledge from pre-trained models, transfer learning enables faster training, better performance, and the ability to tackle problems with limited data¹.

2\. Background

The concept of transfer learning is rooted in the idea that humans can apply knowledge learned in one context to new situations. For example, learning to ride a bicycle can make it easier to learn to ride a motorcycle. Similarly, in machine learning,

transfer learning allows models to leverage pre-existing knowledge to accelerate learning and improve performance on new tasks².

Early research on transfer learning dates back to the 1970s, with studies exploring knowledge transfer in neural networks¹. Over the years, the field has evolved, with significant contributions from researchers like Lior Pratt, who formulated the discriminability-based transfer (DBT) algorithm in 1992¹. By 1998, the field had expanded to include multi-task learning and more formal theoretical foundations¹.

Andrew Ng, a prominent figure in AI, highlighted the importance of transfer learning in his NIPS 2016 tutorial, predicting that it would become a key driver of machine learning commercial success¹. This prediction has come to fruition, with transfer learning now playing a crucial role in various AI applications, including image recognition, natural language processing, and speech recognition.

3\. Types of Transfer Learning

Transfer learning can be categorized into different types based on the relationship between the source and target tasks and domains. Three common types are:

- * **Inductive Transfer Learning:** In this type, the source and target tasks are different, but the domains are the same. This is often used in computer vision, where models pre-trained on large image datasets are adapted for specific tasks like object detection³.
- * **Transductive Transfer Learning:** Here, the source and target tasks are the same, but the domains are different. For example, a model trained on restaurant reviews could be adapted to classify movie reviews³.
- * **Unsupervised Transfer Learning:** This type involves unlabeled data in both the source and target domains. It is similar to inductive transfer learning but focuses on unsupervised tasks³.

These types of transfer learning offer flexibility in adapting models to different scenarios, depending on the availability of labeled data and the similarity between tasks and domains.

4\. Applications of Transfer Learning

Transfer learning has found widespread applications in various domains, revolutionizing the way AI models are developed and deployed. Some notable applications include:

- * **Image Recognition and Classification:** Transfer learning has significantly improved image recognition tasks by

1590	leveraging pre-trained models on large	reduces the risk of overfitting ³ .	1660
1591	datasets like ImageNet. These models can be		1661
1592	fine-tuned for specific tasks, such as	However, transfer learning also has some	1662
1593	medical image classification or identifying	limitations:	1663
1594	species in wildlife images ⁴ .		1664
1595	* **Natural Language Processing (NLP):**	* **Domain Mismatch:** If the source and target	1665
1596	Transfer learning has been instrumental in	domains are significantly different,	1666
1597	advancing NLP applications, including	transfer learning may not be effective. The	1667
1598	sentiment analysis, text classification,	pre-trained model may not have learned	1668
1599	and machine translation. Pre-trained	features relevant to the new task, leading	1669
1600	language models like BERT and GPT can be	to poor performance ⁹ .	1670
1601	adapted for specific language processing	* **Overfitting:** Fine-tuning a pre-trained	1671
1602	tasks, enabling more accurate and efficient	model on a small dataset can lead to	1672
1603	language understanding ⁵ .	overfitting, where the model performs well	1673
1604	* **Speech Recognition:** Transfer learning has	on the training data but poorly on unseen	1674
1605	enhanced speech recognition systems by	data ¹⁰ .	1675
1606	transferring knowledge from general audio	* **Negative Transfer:** In some cases,	1676
1607	models. This has led to improved accuracy	transferring knowledge from the source	1677
1608	in voice commands, transcription, and other	domain can negatively impact the	1678
1609	speech-related tasks ⁶ .	performance on the target task. This can	1679
1610	* **Medical Diagnosis:** Transfer learning has	happen if the tasks are dissimilar or the	1680
1611	shown promise in improving medical	source domain has irrelevant features ¹ .	1681
1612	diagnosis by adapting models trained on		1682
1613	existing medical imaging datasets. This can	Despite these limitations, the benefits of	1683
1614	aid in faster and more accurate diagnoses,	transfer learning often outweigh the	1684
1615	leading to better patient outcomes ⁷ .	drawbacks, making it a valuable technique	1685
1616	* **Recommendation Systems:** Transfer learning	in many machine learning applications.	1686
1617	can be used to improve recommendation		1687
1618	systems by leveraging knowledge from user	## **6\. Future of Transfer Learning**	1688
1619	behavior data. This enables models to make		1689
1620	more personalized recommendations and	Transfer learning is an evolving field with	1690
1621	enhance user experiences ⁶ .	ongoing research and development. Some key	1691
1622		areas for future exploration include:	1692
1623	These are just a few examples of how transfer	* **Multi-domain Adaptation:** Developing	1693
1624	learning is being applied across different	models that can effectively transfer	1694
1625	domains. Its ability to adapt models to new	knowledge across multiple diverse domains ¹¹ .	1695
1626	tasks and domains with limited data has	* **Incremental Learning:** Enabling models to	1696
1627	made it a valuable tool in various AI	continuously learn and adapt to new	1697
1628	applications.	information while retaining previously	1698
1629		learned knowledge ¹² .	1699
1630	## **5\. Advantages and Disadvantages of	* **Model Compression:** Reducing the size of	1700
1631	Transfer Learning**	large pre-trained models without	1701
1632		sacrificing performance, making them more	1702
1633	Transfer learning offers several advantages	suitable for deployment in	1703
1634	over training models from scratch:	resource-constrained environments ¹² .	1704
1635		* **Addressing Ethical Concerns:** Ensuring	1705
1636	* **Reduced Training Time:** By leveraging	fairness, mitigating bias, and addressing	1706
1637	pre-trained models, transfer learning	privacy concerns in transfer learning	1707
1638	significantly reduces the time required to	applications ¹³ .	1708
1639	train a model for a new task. This is		1709
1640	because the model already has a foundation	These advancements will further enhance the	1710
1641	of knowledge, and only the final layers or	capabilities of transfer learning and	1711
1642	specific parameters need to be adjusted ⁸ .	expand its applications in various fields.	1712
1643	* **Improved Performance:** Transfer learning		1713
1644	often leads to better performance,	## **7\. Conclusion**	1714
1645	especially when data for the new task is		1715
1646	limited. The pre-trained model has already	Transfer learning has become a highly effective	1716
1647	learned relevant features and patterns,	approach in machine learning, allowing for	1717
1648	which can be beneficial for the new task ³ .	faster training, enhanced performance, and	1718
1649	* **Lower Computational Costs:** Transfer	the ability to address challenges with	1719
1650	learning can reduce computational costs by	limited data. By utilizing pre-trained	1720
1651	requiring less data and training time. This	models, it has transformed numerous AI	1721
1652	is particularly important in deep learning,	applications, such as image recognition,	1722
1653	where training models can be	natural language processing, and speech	1723
1654	computationally expensive ³ .	recognition. Despite certain challenges and	1724
1655	* **Enhanced Generalization:** Transfer	limitations, continuous research and	1725
1656	learning can improve the generalization	development are driving its evolution,	1726
1657	ability of models by incorporating	ensuring an even greater impact on the	1727
1658	knowledge from other domains. This helps	future of AI.	1728
1659	models perform better on unseen data and		1730

Table 5: Manual evaluation of survey selection and mind-map extraction for the first batch of 50 papers.

CorpusID	IsSurvey	AccurateMap	RelevantMap
3524264	True	True	True
3644401	True	True	True
4503761	True	True	False
5058972	True	True	False
8034133	True	True	True
8922493	True	True	True
9935621	True	False	True
10934716	True	True	True
17513321	True	True	True
18750590	True	True	True
20774863	False	False	False
22727391	True	True	True
23384543	True	True	True
33413610	True	True	False
52841074	True	True	True
55254540	True	True	True
56895323	True	True	True
57721147	True	False	True
64256250	True	False	True
115523285	True	True	False
119297355	True	False	True
198933686	True	True	True
210865204	True	True	False
211010433	True	True	True
212665971	True	True	True
214743520	True	True	True
216056393	True	True	True
219316962	True	True	False
220302470	True	False	True
221446014	True	True	True
222095837	True	True	False
225029039	False	False	False
226227376	True	True	False
226300094	True	True	True
227228021	True	True	False
229363354	True	False	True
231149960	False	False	False
231698518	True	True	True
233722066	False	False	False
234213205	True	True	False
235352671	True	False	False
235458292	True	True	True
235485414	True	True	False
235490196	True	True	True
235669589	True	True	True
235766219	True	True	False
236090307	True	True	True
236772630	True	True	True
236976256	True	True	True
236986986	True	True	False

Table 6: Manual evaluation of survey selection and mind-map extraction for the second batch of 50 papers.

CorpusID	IsSurvey	AccurateMap	RelevantMap
237372527	True	False	True
237373628	True	True	True
237485263	True	True	True
238242214	True	False	True
238242941	True	True	True
238242941	True	False	False
238639787	True	True	False
244119139	True	True	True
244773222	True	True	True
245353469	True	True	True
245877584	False	False	False
247763152	True	True	True
248834382	True	True	False
249209981	True	True	False
249687282	True	False	False
250089226	True	True	True
250939903	True	True	True
251104722	True	True	True
251506514	True	True	True
251643467	True	True	False
252683270	True	True	True
252762319	True	True	False
253370610	True	True	True
253796900	True	True	True
254274880	True	False	True
254756520	True	True	True
255025269	True	True	True
256227178	True	True	True
256826729	True	True	True
257232619	True	True	False
257255597	True	False	True
257522478	True	True	False
258539255	True	True	True
258722762	True	True	True
259043594	True	True	True
259088696	True	True	True
259108865	True	True	True
259154569	False	False	False
259283502	True	True	True
259951356	True	True	True
260229544	True	True	True
260316174	True	True	True
260849783	True	False	True
261682162	True	False	True
263134374	True	True	False
263334211	True	True	True
263830273	True	True	True
263831409	True	False	True
263909496	True	True	True
264604532	True	True	True

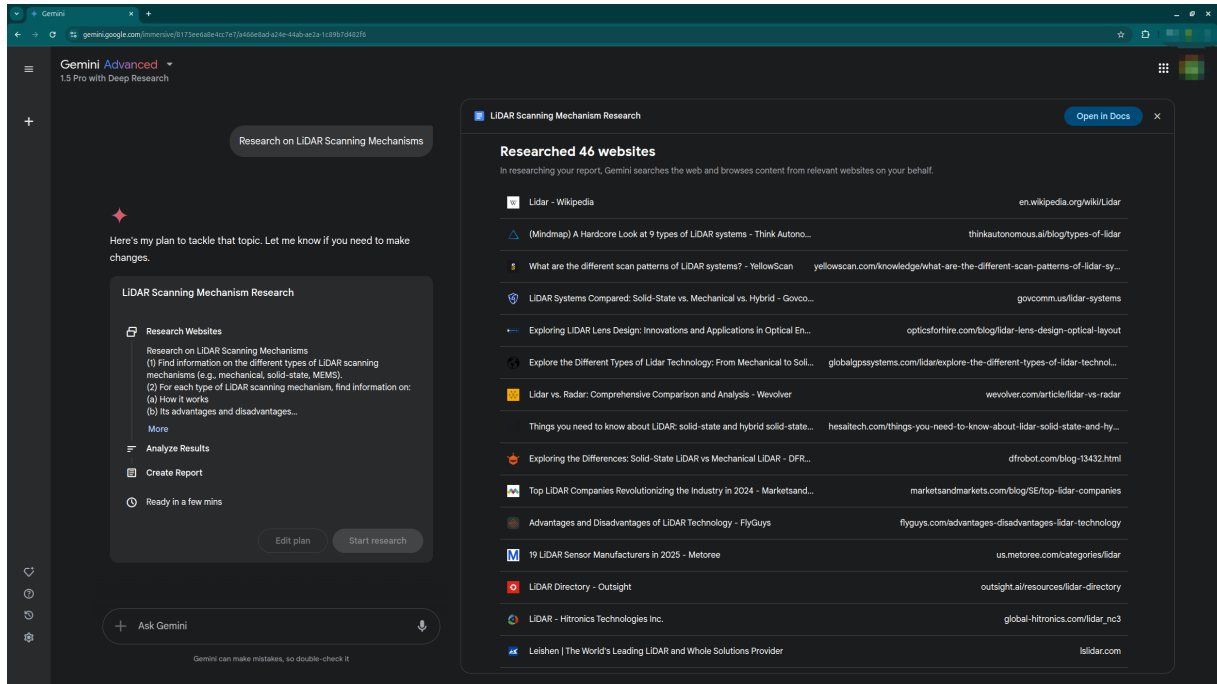


Figure 3: Sources used while asking Deep Research from Gemini to work on LiDAR Scanning Mechanisms.

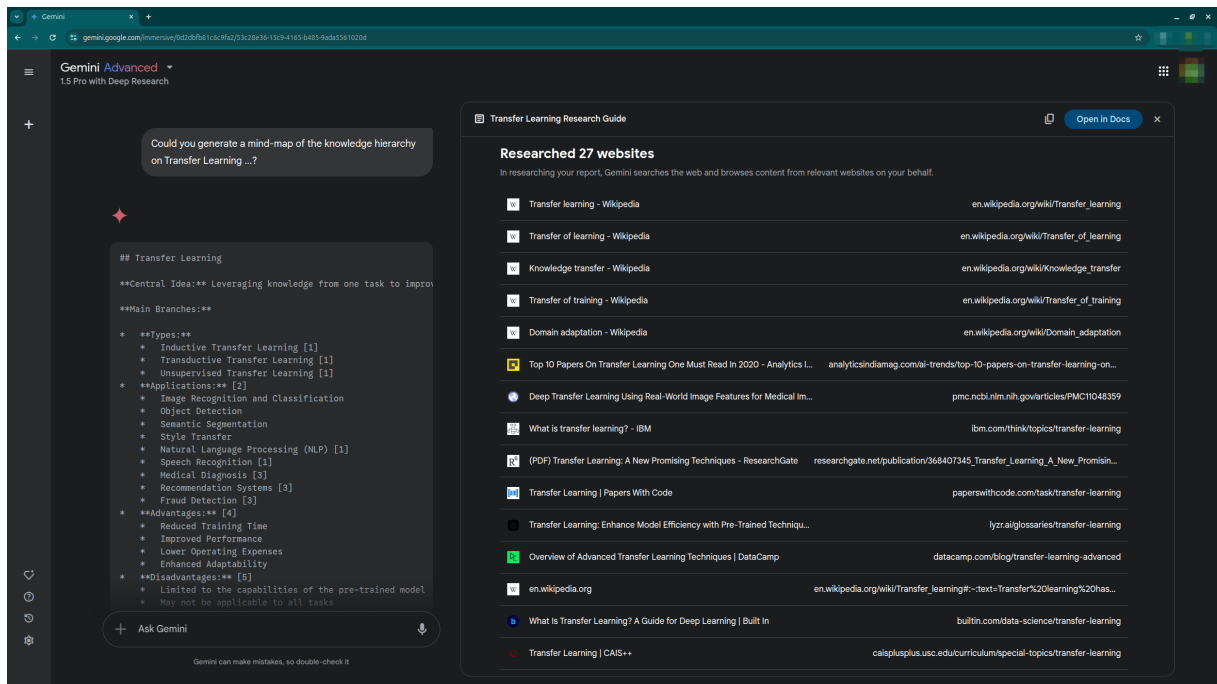


Figure 4: Sources used while asking Deep Research from Gemini to work on transfer learning.