

Article

Not peer-reviewed version

A Survey on AI Search with Large Language Models

[Jian Li](#)^{*}, Xiaoxi Li, Yan Zheng, Yizhang Jin, Shuo Wang, Jiafu Wu, Yabiao Wang, Chengjie Wang, Xiaotong Yuan

Posted Date: 25 August 2025

doi: 10.20944/preprints202507.2024.v2

Keywords: AI search; large language models; web agents; deep search



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

A Survey on AI Search and Web Agent

Jian Li ^{1,2,*}, Xiaoxi Li ³, Yan Zheng ¹, Yizhang Jin ¹, Shuo Wang ¹, Jiafu Wu ¹, Yabiao Wang ¹, Chengjie Wang ¹ and Xiaotong Yuan ²

¹ Tencent YouTu Lab
² Nanjing University
³ Renmin University of China
* Correspondence: swordli@tencent.com

Abstract

Searching for accurate information is a complex task that requires significant effort. Although search engines have transformed the way we access information, they often struggle to fully comprehend intricate human intentions. Recently, Large Language Models (LLMs) have demonstrated impressive abilities in understanding and generating language. However, LLMs face limitations in acquiring external knowledge and accessing the most current information. AI search has evolved by integrating LLMs into the online search process, enabling it to address complex real-world challenges through comprehensive information retrieval and multi-step reasoning, thereby enhancing our ability to browse and search the web effectively. In recent years, substantial progress has been made in refining AI search. This paper provides an in-depth review of these advancements, focusing on text-based AI search, web browsing agents, multimodal AI search, benchmarks, software, and products. We also examine the limitations of current AI search methods and explore promising future directions. For further details, please visit our website <https://github.com/swordldev/Awesome-AI-Search>.

Keywords: AI search; large language models; web agents; deep search

1. Introduction

Searching for information is a fundamental daily necessity for humans. To meet the demand for rapid access to desired information, key web search technologies like PageRank [1–3] have been developed to support information retrieval systems. These technologies power search engines such as Google, Bing, and Baidu, efficiently retrieving relevant web pages in response to user queries and providing convenient access to information on the internet. Advances in natural language processing (NLP) [4] and information retrieval (IR) [5] have further enhanced machines’ ability to accurately extract content from the vast array of websites available online. However, as user queries become increasingly complex and the demand for precise, contextually relevant, and up-to-date responses grows, traditional search technologies face challenges in fully comprehending intricate human intentions. Consequently, users often need to manually open, read, and synthesize information from multiple web pages to answer complex questions.

Recently, Large Language Models (LLMs) [6] have garnered significant attention in both academic and industrial domains. LLMs such as ChatGPT [7] and LLaMA [8] have demonstrated remarkable advancements in language understanding, reasoning, and information integration. However, LLMs face limitations in acquiring external knowledge and accessing the most current information. To address these challenges, researchers are integrating the impressive capabilities of LLMs with search engines and websites, aiming to enhance real-time evidence gathering and reflective reasoning. The complementary strengths of LLMs and search engines present an opportunity for synergy, where the reasoning abilities of LLMs are augmented by the vast web information accessible through search engines. This integration is revolutionizing the way we seek and synthesize web-based information, ushering in a new era of search technology known as Artificial Intelligence (AI) Search. In this survey,

we provide an overview of recent advancements in the rapidly evolving field of AI Search. As depicted in Figure 1, we categorize the literature into five primary areas: (1) text-based AI search, (2) web browsing agents, (3) multimodal AI search, (4) benchmarks, and (5) software and products.

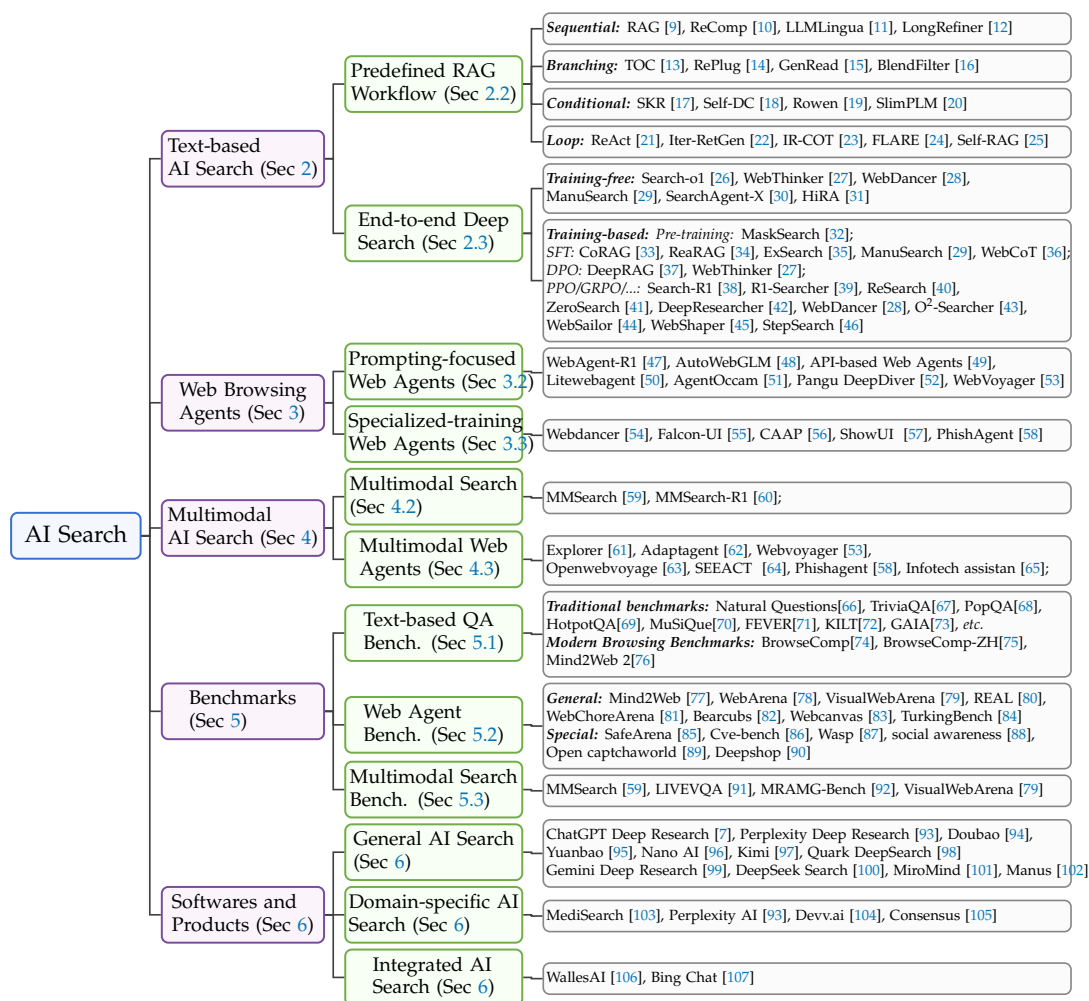


Figure 1. Taxonomy of research on AI search: investigating text-based AI search, web browsing agents, multimodal AI search, benchmarks, softwares and products.

The classic **Text-based AI Search** operates through a Retrieval-Augmented Generation (RAG) framework [108]. In this workflow, RAG retrieves relevant passages from search engines based on the input query and integrates them into the context of a Large Language Model (LLM) for generating responses. This enables the LLM to utilize external knowledge when addressing questions. Another approach within text-based AI Search is the deep search method, which acquires external information by interacting with search engines as part of an end-to-end coherent reasoning process to tackle complex information retrieval challenges. Unlike predefined workflows, this method allows the model to autonomously determine when to employ search-related tools during its reasoning, enhancing flexibility and effectiveness. **Web Browsing Agent** accomplishes specific tasks on target websites through a sequence of actions, utilizing a thought-action-observation paradigm. For instance, if a web agent is tasked with calculating the driving time from Shanghai to Beijing using an open street map, it would perform this task by interacting with the website. Web agents are classified into two primary paths: Generalist Deep Browsing Web Agents, which perform complex web browsing tasks across multiple types of web pages, and Specialist Parsing Web Agents, which employ dedicated training procedures to focus specifically on action sequences or interface elements. Additionally, with the emergence of visual web-oriented benchmarks and the development of Multimodal Large Language Models, many agents now incorporate screenshots as sensory input to provide a more comprehensive

understanding of the web environment. Unlike **Multimodal AI Search**, most current AI search methods are confined to text-only settings, overlooking the multimodal nature of user queries and the interleaved text-image information on websites. This limitation is particularly significant given the complexity and interleaved nature of modern websites. For example, capturing a photo of an antique at a museum without knowing its historical context could be addressed by a multimodal AI search engine, which would match the photograph with an interleaved table of images and text retrieved from the Internet, thereby providing the history and story behind it. Thus, a multimodal AI search engine is essential for advancing information retrieval and analysis.

Furthermore, this paper reviews the **Benchmarks** relevant to these methods. Evaluating the search capabilities of AI models, particularly LLM, is crucial for assessing their ability to effectively retrieve, filter, and reason over web-based information. This evaluation is essential for understanding the true web-browsing competence of LLMs and their potential to tackle real-world tasks that demand dynamic information retrieval. In recent years, significant efforts have been made to explore AI search from various perspectives. This paper concentrates on three key areas: text-based question-answering benchmarks, web agent benchmarks, and multimodal benchmarks. The **Software and Products of AI Search**, such as Perplexity [93], have the potential to change our daily lives. We introduce a wide array of state-of-the-art open-source and proprietary models, software, and mainstream AI search products, aiming to present a diverse and comprehensive overview of AI Search. Finally, we discuss the limitations of current AI search methods and explore promising future directions. To illustrate the evolution of AI search methods over time, Figure 2 presents a timeline of recent AI search technologies, related methods, and products.

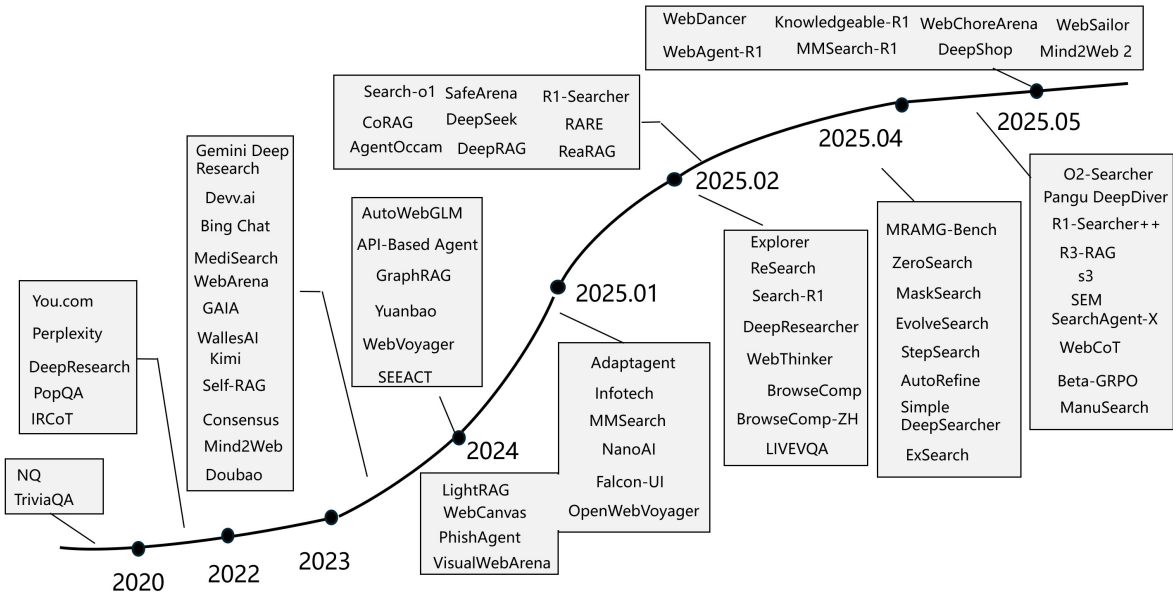


Figure 2. A timeline of recent AI Search methods and related products has been created, primarily based on the release dates of their respective technical papers.

2. Text-Based AI Search

AI search represents a transformative advancement in information retrieval systems, evolving from traditional search engines to sophisticated approaches incorporating RAG workflows and Deep Search capabilities, as shown in Figure 3. This section provides an overview of the key components and cutting-edge developments in modern text-based AI search technologies.

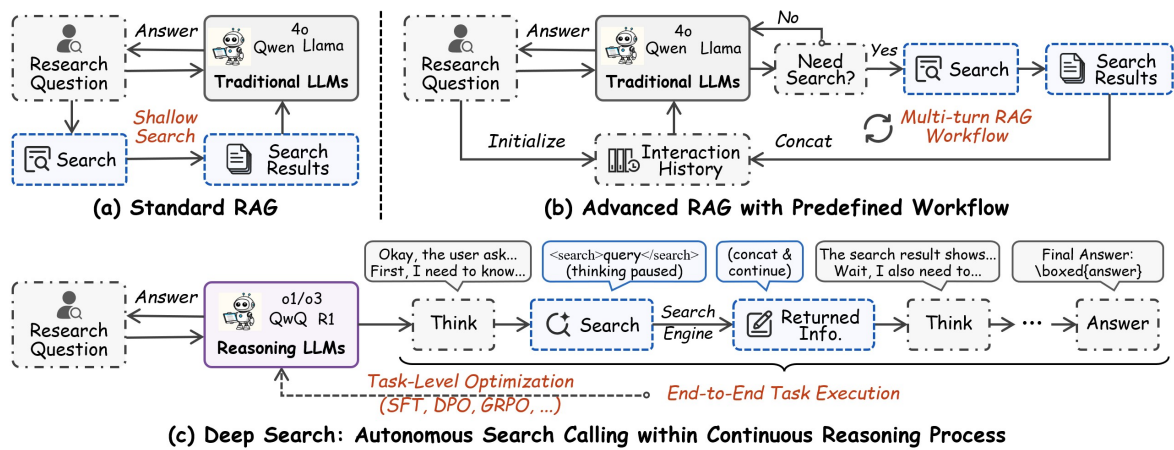


Figure 3. Evolution of text-based AI search paradigms, from (a) standard RAG that retrieves once per query, to (b) advanced RAG workflows capable of multi-turn search and decision-making, and finally to (c) fully autonomous, reasoning-model-powered Deep Search.

2.1. Traditional Search Engines

Traditional search engines form the foundation of modern search engines. They employ a variety of techniques to efficiently process user queries and return relevant results. Two key components of these systems are document retrieval and post-ranking, which work in tandem to provide users with the most pertinent information [109].

Document Retrieval

Document retrieval involves identifying relevant documents from a collection based on a user query. It is a crucial step in information retrieval, as it determines which documents are most relevant to the user’s query. Traditional document retrieval systems typically employ techniques such as inverted indexing, term frequency-inverse document frequency (TF-IDF), and BM25 models [110,111]. More advanced approaches incorporate semantic matching using dense vector representations and neural ranking models [112–115]. The retrieval process often involves query preprocessing, document indexing, similarity computation, and efficient search algorithms to handle large-scale document collections. Recently, some research has explored LLM-based generative retrieval [116–119], eliminating the need to build document indexes and directly generating document identifiers through LLMs.

Post-Ranking

Post-ranking refines the results of a search query after the initial retrieval stage. It is used to improve the quality of the search results by applying additional filters and reranking algorithms. Post-ranking systems typically employ learning-to-rank algorithms, neural reranking models, and LLM-based reranking models that combine multiple ranking signals [120–122]. This stage is crucial for improving search precision and user satisfaction by promoting the most relevant documents to top positions [120].

2.2. Retrieval-Augmented Generation with Pre-Defined Workflows

Retrieval-Augmented Generation (RAG) enhances generative models by integrating a retrieval mechanism, allowing the model to ground its responses in external, reliable knowledge [9,123]. Typically, a RAG system consists of a retriever and a generator, and the interaction between these components gives rise to four main RAG paradigms [108]:

Sequential RAG

Sequential RAG follows a linear “retrieve-then-generate” workflow, where the retriever first fetches relevant documents and the generator produces the final response based on these documents [9,123–127]. Early works explored joint or separate training of the retriever and generator,

while recent approaches often use a frozen generator and focus on optimizing the retriever [128–130]. Pre-retrieval modules, such as rewriters [126] and post-retrieval compressors [10–12,131–133] further improve efficiency and response quality.

Branching RAG

Branching RAG processes the input query through multiple parallel pipelines, each potentially involving its own retrieval and generation steps, and then merges the outputs for a comprehensive answer [13–16,134]. This approach enables finer-grained handling of complex queries, such as decomposing questions into sub-questions [13], augmenting queries with additional knowledge [16], or merging generated and retrieved content [14,15].

Conditional RAG

Conditional RAG introduces a decision-making module to adaptively determine whether retrieval is necessary for a given query, improving flexibility and robustness [17–20]. Methods include training classifiers to predict the need for retrieval [17,20], using model confidence to guide retrieval [18], or employing consistency checks across perturbed queries [19].

Loop RAG

Loop RAG features iterative and interactive retrieval-generation cycles, enabling deep reasoning and handling of complex queries [21–25,135,136]. These methods alternate between retrieval and generation [22,23], dynamically decide when to retrieve [24,25], or decompose and answer sub-questions with verification steps to reduce misinformation [135,136].

2.3. End-to-End Deep Search Within Reasoning Process

Unlike traditional RAG workflows, Deep Search methods integrate external knowledge acquisition by utilizing search engines within an end-to-end coherent reasoning process to address complex information retrieval challenges. This approach eliminates the need for predefined workflows, allowing the model to autonomously determine when to employ search-related tools during its reasoning process, thereby enhancing flexibility and effectiveness [137–139].

Training-Free Methods

These methods aim to augment the reasoning model's search capabilities by crafting instructions that clarify the model's task and the use of search tools. Initially, Search-o1 [26] introduced an agentic RAG mechanism enabling the reasoning model to autonomously retrieve external knowledge when faced with uncertain information during the primary reasoning process, thus addressing the knowledge gaps in long Chain-of-Thought (CoT) reasoning. They also introduced a Reason-in-Documents process, which thoroughly analyzes the content of retrieved documents after each search call, providing concise and useful information to the main reasoning chain. Experiments demonstrated significant performance improvements across mathematical, scientific, coding, and multi-hop question answering tasks.

Following this paradigm, a series of works such as WebThinker [27], WebDancer [28], ManuSearch [29], and HiRA [31] have proposed advanced frameworks. Typically, these methods incorporate browsing of collected webpage URLs to facilitate in-depth web exploration. Additionally, to enhance search efficiency, SearchAgent-X [30] proposed an efficient reasoning framework aimed at increasing system throughput and reducing latency through high-recall approximate retrieval, priority-aware scheduling, and non-stagnant retrieval mechanisms. Beyond directly answering users' information-seeking questions, some works like WebThinker [27] have explored autonomously writing research reports while gathering information, offering users more comprehensive and cutting-edge knowledge.

Training-Based Methods

These methods design various training strategies to incentivize or enhance the LLM's search capabilities within the reasoning process. These strategies encompass pre-training, supervised fine-tuning (SFT), and reinforcement learning (RL).

During pre-training, the MaskSearch [32] framework introduces a Retrieval-Augmented Mask Prediction (RAMP) task, which trains the model to use search tools to fill in masked text, thereby enhancing its retrieval and reasoning abilities.

For Supervised Fine-Tuning (SFT), several methods focus on synthesizing long chain-of-thought data that incorporates search actions [28,29,33–36,44,140]. Specifically, CoRAG [33] addresses the lack of intermediate retrieval steps in existing RAG datasets by automatically generating retrieval chains through rejection sampling. ReaRAG [34] avoids complex reinforcement learning by building a specialized dataset for fine-tuning via policy distillation. ExSearch [35] introduces an iterative self-incentivization framework based on the Generalized Expectation-Maximization (GEM) algorithm, enabling the model to learn from its own generated search trajectories. SimpleDeepSearcher [140] simulates user search behavior in a real-world web environment to synthesize multi-turn reasoning trajectories, which are then curated using a multi-criteria strategy. ManuSearch [29] leverages its multi-agent framework to decompose the deep search process and generate structured reasoning data. Lastly, WebCoT [36] synthesizes training data by reconstructing successful and failed trajectories, explicitly embedding reasoning skills like reflection, branching, and rollback into the chain of thought.

Furthermore, RL-based training has recently garnered significant attention. Some works leverage Direct Preference Optimization (DPO) [141]. For instance, DeepRAG [37] introduces a chain of calibration method to refine the model's atomic decisions, thereby synthesizing preference data for training. WebThinker [27] constructs positive and negative pairs based on the model's ability to correctly complete research tasks while efficiently using tools. By iteratively constructing data and training the model with DPO, it implements on-policy RL training, improving performance on complex reasoning and report generation tasks.

Another line of work has explored training strategies based on PPO [142], GRPO [143], REINFORCE++ [144], and others. Initially, Search-R1 [38], R1-Searcher [39], ReSearch [40], WebSailor [145], and WebShaper [45] used the accuracy of the generated answer as a rule-based reward to encourage the LLM to use Wikipedia-based search tools during reasoning, achieving significant performance improvements in multi-hop QA tasks. Subsequently, a series of studies have investigated various enhancement strategies. These include leveraging web search capabilities [28,42,44,47,146], refining retrieved information [42,43,147], enabling multi-tool usage [148], developing improved sampling techniques [149], designing advanced reward functions [150], combining outcome and process rewards [46,151], enhancing training efficiency [41], and implementing iterative SFT and RL training cycles [152]. To optimize search efficiency, methods such as SEM [153], β -GRPO [154], and s3 [155] have been proposed, which design training algorithms and reward functions for more efficient and accurate use of search tools.

3. Web Browsing Agent

The Web Agent is an AI-driven autonomous program designed to mimic human interactions within web browsers. Unlike AI search, the task of Web Agent is not limited to search. It is also specified in the network environment and even requires executing a series of specific interactive operations. This section explores the definition, classification, and advanced technologies related to Web Agent more broadly.

3.1. Agent

Agent [156] is defined as an autonomous intelligent entity with human-like decision-making abilities, capable of perceiving diverse environments, performing various tasks, and taking concrete actions. Traditional autonomous agents are typically based on rule-based approaches and tend to

perform well only in specific tasks or closed environments. Therefore, agent based on large models(LM) has become a prominent research direction [157].

The core of LM-based autonomous agent lies in two aspects: architectural design and capability acquisition [158]. Architectural design [159–162] involves employing specialized network structures and modular components to enhance the agent’s comprehension and execution abilities. (e.g., AgentVerse’s unified framework [162]). Capability acquisition [47,163,164] focuses on optimizing task-specific performance through parameter adjustment or strategic prompt design, enabling agents to adapt to target domains.

Currently, agents can be categorized based on their architecture into single-agent and multi-agent systems [165–167]. Typically, a single-agent architecture can independently work without relying on other agents or user feedback, making it more efficient when addressing well-defined problems. In contrast, multi-agent architecture deploys specialized agents across distinct knowledge domains, enabling collaborative problem-solving that excels in complex, cross-domain tasks requiring iterative coordination.

As an extension of agent technology in vertical domains, Web Agents focus specifically on task generalization within more complex web environments [63]. According to their adaptation strategies, current Web Agents can be classified into two categories, as shown in Figures 4 and 5. The first category is post-training-based Web Agents, which employ SFT and RL to adapt models, enhancing their understanding and execution of specific web tasks. The second category is prompt-based Web Agents, which rely on prompt engineering or context engineering to improve performance in complex web navigation tasks without directly modifying model parameters [168].

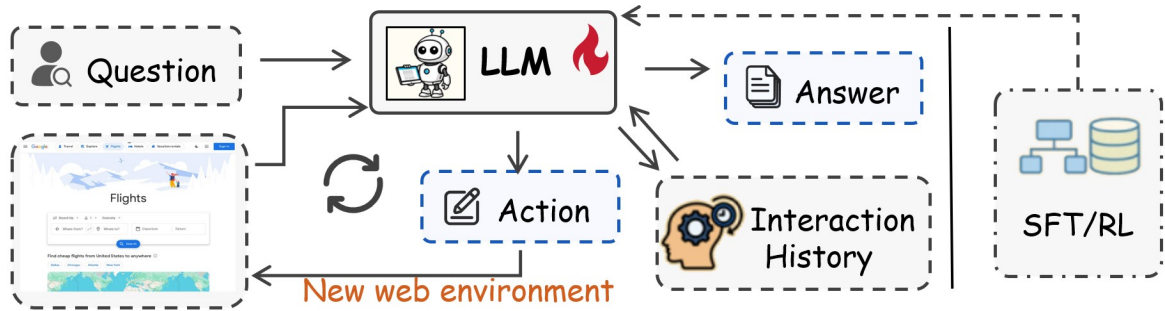


Figure 4. Illustration of Post-training-based Web Agents. The loop terminates when the required information is obtained, returning results to the user.

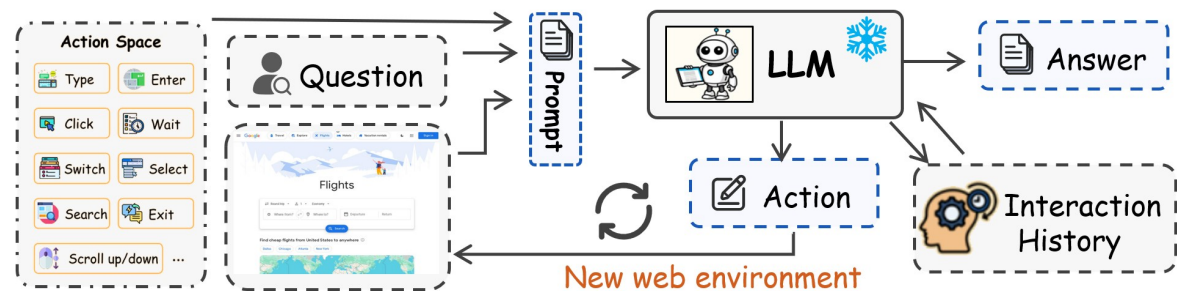


Figure 5. Illustration of Prompt-based Web Agents. The loop terminates when the required information is obtained, returning results to the user.

3.2. Post-Training-Based Web Agents

The open-ended nature of Web Agent applications limits the effectiveness of standard pre-trained models. To address this, techniques such as SFT and RL have become essential, enabling Web Agents to adapt to dynamic environments and learn from interactive feedback, thereby improving their performance on complex web tasks.

Reinforcement Learning

RL allows Web Agents to adapt to dynamic environments in real-time by exploring and learning from interactive feedback [28,47,169,170]. WebAgent-R1 is the first purely end-to-end RL-trained Web Agent [47]. It employs a multi-turn end-to-end RL framework, where the agent is trained through online interactions guided by rule-based outcome rewards. During training, it implements multi-group GRPO [171] method, utilizing multiple parallel interaction trajectories to enhance training efficacy.

The approach of WebAgent-R1 follows a standard RL process, while other methods adjust strategies within the framework. WebDancer [28] focuses on QA pair parsing, employing Decoupled Clip and Dynamic Sampling Policy Optimization [172] to internalize chains-of-thought (CoT) [173] from QA pairs as active behavioral components of the model. WebRL [169] adopts self-evolving online curriculum RL to dynamically generate new tasks and uses result-supervised reward models to provide feedback to the agent. WorkForceAgent-R1 [170] integrates behavior cloning (BC) [174] and GRPO, emphasizing the optimization of both single-step reasoning and planning abilities to help agents adapt to dynamic web environments.

Supervised Fine-Tuning

SFT uses labeled data to fine-tune a pre-trained model, transferring its general capabilities to specific tasks or domains [55,163,175,176]. For example, WebGUM [163] integrates the T5 [177] model with Vision Transformer (ViT) [178] for multimodal task processing. During SFT, ViT encodes images into tokens while T5 performs unified encoding of both text and tokens, enabling joint fine-tuning.

The focus of SFT is on the construction of task-aligned datasets to enhance tuning performance. Falcon UI [55] simulates realistic interactions by exclusively recording visible GUI [179] elements during browsing, constructing specialized datasets for improved GUI comprehension. WebExpert [175] refines existing datasets by annotating screenshot elements, combining rapid decisions with slow reasoning, and optimizes through self-reflection. WebCoT [176] generates “reasoning trajectories” through reflection, branching, and rollback mechanisms, subsequently converted into CoT data for SFT.

Joint Training

SFT can be combined with RL for complementary benefits. AutoWebGLM [48] processes information through HTML simplification and optical character recognition (OCR) [180]. Its training integrates SFT, RL, and rejection sampling fine-tuning (RFT) [48]. During SFT, curriculum learning (CL) [181] is used to enable the model to perform basic web navigation and improving the effectiveness of subsequent RL training.

3.3. Prompt-Based Web Agents

Although post-training methods are effective, prompt engineering [182] is preferable when a Web Agent must adapt quickly to dynamic environments or new tasks without the time or resources for SFT or RL. Both prompt engineering and its advanced form, context engineering [183], do not require complex model modifications; they guide model behavior by optimizing the input data.

A typical application of prompt engineering is the use of APIs [49,184,185]. For example, Microsoft’s CodeAct [49] adopts an “API-first” strategy: the agent first constructs an API documentation map of the website, then selects relevant APIs according to task requirements, and finally processes inputs via model APIs. Similarly, Infogent [184] establishes a modular feedback loop around LLM APIs: the Navigator collects raw data, the Extractor processes information via APIs, and the Aggregator evaluates results.

In theory, the more relevant information a model receives, the more accurate its decisions. Consequently, many Web Agents enrich prompts or context as much as possible [51,53,64]. WebVoyager [53] captures the web page’s accessibility tree, annotates interactive elements, and takes a screenshot; these components are combined as the final prompt to the GPT-4V [186] API to determine the next action. AgentCocam [51] incorporates diverse information into prompts by abstracting operations, merging elements, and selectively replaying interaction history, enabling an efficient LLM-API workflow.

In some cases, Web Agents may employ multiple models. PhishAgent [58], designed for phishing site detection, primarily relies on LLMs. When textual cues are insufficient, it activates a large multimodal model (LMM) to perform brand analysis, enhancing decision accuracy.

4. Multimodal AI Search

Current AI search methods are predominantly confined to text-only environments, often overlooking the multimodal nature of user queries and the intertwined text-image format of website content. This limitation is particularly significant when, for example, a user takes a photograph of an antique at the museum but lacks specific information about it. Therefore, the development of a multimodal AI search engine is essential for enhancing information retrieval and analysis.

4.1. Multimodal Large Language Models

Recently, Multimodal Large Language Models (MLLMs) or Large Multimodal Models (LMMs) [187] have demonstrated exceptional performance across a range of applications, including visual question answering (VQA), visual perception, understanding, and reasoning. Notable closed-source models include GPT-4V [186], GPT-4o [188], and Claude 3.5 Sonnet [189]. In the open-source domain, models such as BLIP [190,191], LLaVA [192,193], Qwen-VL [194], Gemini [195], InternVL [196], and EMU [197] have made significant advancements. The typical MLLM framework consists of three primary modules: a visual encoder responsible for processing visual inputs, a pre-trained language model that handles multimodal signals and performs reasoning, and a visual-language projector that serves as a bridge to align the two modalities. In recent years, substantial efforts have been dedicated to designing MLLM benchmarks [198], examining these models from various perspectives.

4.2. Multimodal Search

Inspired by Text-based AI Search, it is essential to develop a framework for MLLMs to function as multimodal AI search engines. MMSearch [59], illustrated in Figure 6, proposes a multimodal AI search engine pipeline named MMSEARCH-ENGINE, which enhances MLLMs with advanced search capabilities. MMSEARCH-ENGINE maximizes the use of MLLMs’ multimodal information comprehension abilities by incorporating both visual and textual website content as information sources during the searching process: requery, rerank, and summarization. MMSearch-R1 [60], also shown in Figure 6, represents an initial effort to equip MLLMs with active image search capabilities through an end-to-end reinforcement learning framework, assisted by image search tools. This method trains models not only to determine when to invoke the image search tool but also to effectively extract, synthesize, and utilize relevant information to support downstream reasoning. To improve the performance on knowledge-intensive VQA tasks, WebWatcher [199] use a multimodal Agent for Deep Research equipped with enhanced visual-language reasoning capabilities to achieve strong performance on several high-difficulty benchmarks.

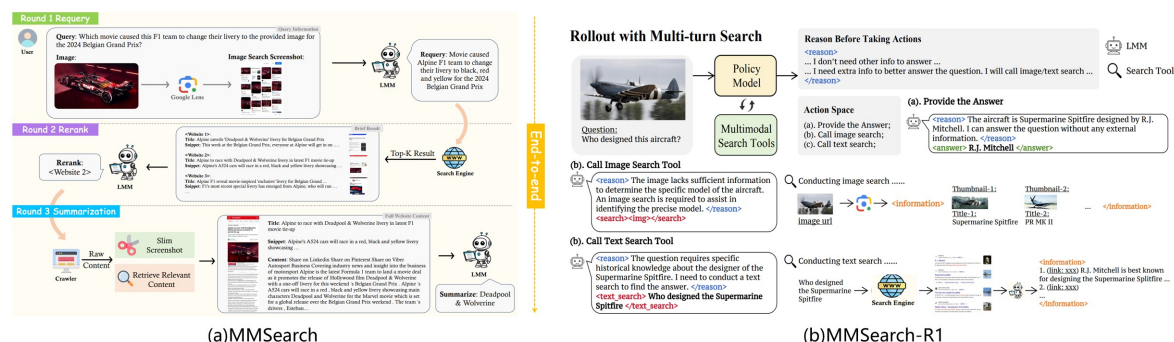


Figure 6. Illustration of Multimodal AI Search. (a) The MMSearch [59] pipeline comprises three sequential stages by an MLLM: (i) requery, (ii) rerank, and (iii) summarization. (b) A detailed view of the MMSearch-R1 [60] highlights the rollout process and the execution of the search tool.

4.3. Multimodal Web Agents

The rise of multimodal large models has driven many LLM-powered Web Agents toward multimodal capabilities [200]. By jointly encoding visual inputs and textual inputs like HTML/DOM [201], these agents enable integrated analysis and decision-making [61–64,202]. WebVoyager [202] first instantiates a web browser and utilizes both visual signals (screenshots) and textual signals (HTML elements) to perform actions. Its successor, OpenWebVoyager [63], adopts the more capable multimodal Idefics2 model and employs SFT, enabling direct understanding of image and text inputs and reducing dependence on prompt engineering. SEEACT [64] leverages a GPT-4V-like MLLM for integrated visual understanding, generating textual plans and performing website actions based on HTML elements. The optimization of multimodal web agents relies on high-quality multimodal datasets. Explorer [61] independently developed a web trajectory synthesis approach to construct the largest and most diverse trajectory dataset to date for GUI agent model training.

Beyond general domains, multimodal web agents also demonstrate high usability in specialized fields [58,65,203]. For instance, researchers at the University of Notre Dame developed a multimodal web agent [203] anonymously navigating dark web marketplaces, solve CAPTCHAs, and extract product information via the Tor browser.

5. Benchmarks

5.1. Text-Based QA Benchmark

As large language models (LLMs) evolve into tool-using agents, the ability to browse the web in real-time has become a critical yardstick for measuring their reasoning and retrieval competence. A variety of widely used English benchmarks have been proposed to assess retrieval capabilities, including TriviaQA, HotpotQA, FEVER, KILT, GAIA, *etc.* These datasets cover multi-hop reasoning, knowledge-intensive QA, and fact checking, typically relying on structured sources like Wikipedia and StackExchange.

Traditional Benchmarks Natural Questions (NQ)[66] is a large-scale QA dataset using real Google search queries and corresponding Wikipedia pages, requiring models to provide both long-form and short-form answers. TriviaQA [67] is a reading comprehension dataset characterized by complex, compositional questions with significant lexical variation from their evidence, often demanding multi-sentence reasoning. PopQA [68] is an entity-centric QA dataset designed to test factual knowledge recall across a long-tail distribution of entity popularity. HotpotQA [69] is a multi-hop QA dataset that requires reasoning across multiple documents and providing sentence-level supporting facts, making it a benchmark for explainable QA. 2WikiMultiHopQA [204] is a more challenging multi-hop QA dataset that integrates Wikipedia with Wikidata, using structured triples to explain complex reasoning paths. MuSiQue [70] is a multi-hop QA dataset emphasizing connected reasoning and includes unanswerable examples to challenge models that rely on shortcuts. FEVER [71] is a benchmark for fact verification, requiring systems to classify claims as SUPPORTED, REFUTED, or NOTENOUGHINFO against Wikipedia and provide sentence-level evidence. KILT [72] unifies 11 knowledge-intensive NLP tasks under a single Wikipedia snapshot, providing a standardized framework for evaluating both task performance and evidence retrieval. GAIA [73] evaluates general-purpose AI assistants with real-world questions that require a combination of reasoning, tool use, and multi-modality, revealing a large gap between AI and human performance. TREC Health Misinformation Track [205] provides datasets with binary “yes/no” health questions based on medical consensus to evaluate a system’s ability to combat health misinformation.

Modern Browsing Benchmarks While traditional benchmarks mentioned above have effectively measured an AI’s ability to retrieve straightforward information through basic queries (e.g., single-hop fact lookup), their simplicity has led to saturation—modern models now achieve near-perfect scores on these tasks. This progress reveals a critical gap: real-world information needs often require persistent navigation through complex data landscapes. These challenges mirror the evolutionary jump from

arithmetic tests to mathematical proofs—where success depends less on recall and more on strategic problem-solving.

BrowseComp [74] is a benchmark dataset introduced to evaluate web-browsing AI agents. It contains 1,266 challenging questions requiring persistent navigation of the internet to find entangled information. Key features include: (1) *High difficulty* - questions are designed to be unsolvable by humans within 10 minutes; (2) *Verifiability* - short reference answers enable easy validation; (3) *Diverse topics* spanning sports, fiction, and academic publications; and (4) *Core capability measurement* focusing on persistence, factual reasoning, and creative search strategies. BrowseComp-ZH[75] benchmark is a high-difficulty Chinese web browsing evaluation dataset consisting of 289 multi-hop questions across 11 domains (e.g., Art, Film&TV, Medicine). Each question is reverse-engineered from verifiable factual answers and undergoes rigorous two-stage quality control to ensure retrieval difficulty and answer uniqueness. Figure 7 illustrates these two benchmarks and shows some complex and challenging queries. Mind2Web 2 [76] is also a modern benchmark with 130 realistic, high-quality, and long-horizon tasks that require real-time web browsing and extensive information synthesis.

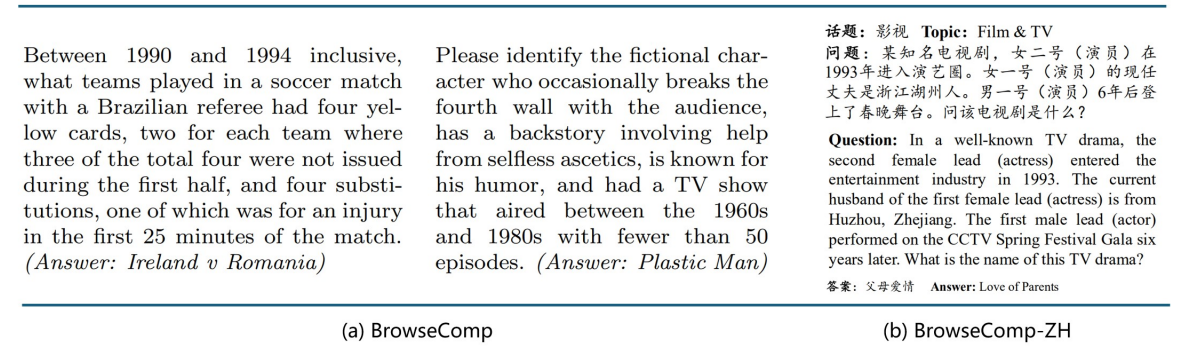


Figure 7. Illustration of Modern Browsing Benchmarks with complex and challenging queries. (a) BrowseComp [74]. (b) BrowseComp-ZH [75].

5.2. Web Agent Benchmark

As integrated systems beyond standalone models, Web Agents require specialized evaluation standards. Web Agent Benchmarks address this need through standardized test tasks and assessment frameworks that quantify performance in navigation, operation, and reasoning within simulated web environments. Current benchmarks fall into two categories: general benchmarks [77–84,168].and specialized benchmarks [85–90].

General Benchmarks General benchmarks evaluate Web Agents’ open-ended browsing capabilities across diverse domains. Mind2Web [77] , shown in Figure 8, pioneering this approach through tasks spanning five domains (travel, shopping, services, entertainment, information) represented by task descriptions, action sequences, and webpage snapshots. WebArena [78] builds a reusable virtual web environment based on real-world site categories (e-commerce, forums, collaboration tools, CMS), minimizing real-world noise, while WebChoreArena [81] extends this framework with memory-intensive tasks, reduced ambiguity, and template-based generation, making it more challenging than WebArena. WebCanvas [83] employs a dynamic, real-time evaluation framework focused on key task steps (required task steps) to assess agent performance in live web environments. VisualWebArena [79] adds multimodal inputs(screenshot) to WebArena tasks to evaluate agents under visual challenges. REAL [80] provides a framework for multi-turn agent evaluation on simulated real-world websites, supporting both open-source and proprietary agents via black-box commands within a browser environment. BEARCUBS [82] contains 111 information-seeking queries requiring access to real web content, without using a simulated web environment.

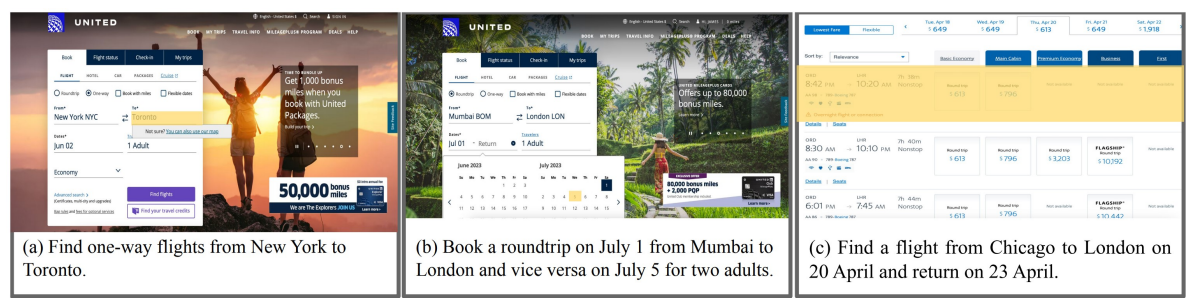


Figure 8. Sample tasks of Mind2Web [77]. The web agent benchmark can test an agent’s generalizability across tasks on the same website (a vs. b), similar tasks on different websites (a vs. c).

Specialized Benchmarks For task-specific objectives, general benchmarks often prove inadequate for evaluating Web Agent performance in their domains, necessitating specialized benchmarks. DeepShop [90] targets e-commerce by generating queries across five popular shopping domains and assessing agents through product attribute analysis, matching, and ranking. SafeArena [85] focuses on malicious web agent use, offering 500 tasks with metrics including Task Completion Rate (TCR), Safety Score (NSS), and refusal rate. CASA [88] evaluates agents’ sensitivity to cultural and social norms by testing their responses to norm-violating requests. CVE Bench [86] is a benchmark based on critical vulnerabilities and exposures, featuring a sandbox that simulates web attacks to effectively assess agent vulnerabilities. Wasp [87] is another security benchmark designed to evaluate web agents’ robustness against prompt injection attacks in an end-to-end manner.

5.3. MM Search Benchmark

LLMs have made significant strides in understanding and reasoning about live textual content when integrated with search engines. Despite these advancements, a crucial question remains: has the understanding of other modalities, such as visual knowledge in live contexts, been similarly addressed? Are there benchmarks for multimodal search methods?

MMSearch [59] introduced a multimodal AI search engine benchmark to thoroughly assess the searching performance of MLLMs, marking the first evaluation dataset to measure MLLMs’ capabilities in multimodal searching. LIVEVQA [91], depicted in Figure 9, is an automatically collected benchmark dataset specifically designed to evaluate current AI systems on their ability to answer questions requiring live visual knowledge. However, existing benchmarks for this critical task face a significant shortage of suitable datasets and scientifically rigorous evaluation metrics. MRAMG-Bench [92] is a novel benchmark created to comprehensively evaluate the MRAMG task. It consists of six meticulously curated English datasets, sourced from three domains: Web, Academia, and Lifestyle, across seven distinct data sources. Recently, BrowseComp-VL [199] proposes a BrowseComp-style benchmark requiring complex visual and textual information retrieval.



Figure 9. Illustration of four categories of LiveVQA [91]. QA pair for basic image for understanding, and two multimodal multi-hop QA pairs for deeper reasoning.

VisualWebArena [79] is primarily designed for visual web agent tasks, incorporating both textual and visual content from real-world environments. It comprises 910 real-world tasks across three distinct web environments. A key feature of VisualWebArena is that all tasks require agents to process and

interpret visual information, rather than relying solely on textual or HTML-based cues. For evaluation, the metrics follow WebArena's framework but extend it by incorporating image verification alongside the original two assessment methods.

6. Softwares and Products

AI search ecosystem has rapidly diversified into general-purpose platforms, domain-specific tools, and integrated assistants, each leveraging large language models (LLMs), retrieval-augmented generation (RAG), and agentic workflows to redefine information retrieval. Below, we will introduce the key products driving this transformation.

Global General-Purpose AI Search Engines. A pioneer in deep research, ChatGPT Deep Research [7] integrates Bing's real-time web search to provide concise, conversational responses, sparking a surge of interest among researchers in large language models. Perplexity Deep Research [93] combines GPT-4 and Claude 3 with real-time web crawling, providing source-attributed answers. Its Discover feature tracks trending topics, making it ideal for academic literature reviews and technical writing. You.com [206] prioritizes privacy and personalization, allowing model switching (e.g., GPT-4, Claude) mid-session. Its Smart mode offers free access, while Research mode supports deep investigations with citation exports. Gemini Deep Research [99] embeds multi-modal capabilities into Pixel phones and Wear OS, enabling real-time translation via camera and health data-driven recommendations, reinforcing its "hardware-software" synergy in high-end markets. Optimized for speed and cost-efficiency, Doubao [94] integrates seamlessly with Douyin for video-content searches. Yuanbao [95] redefines "search-as-service" by embedding within WeChat's ecosystem. Its three-layer architecture—base model (trillion-parameter MoE), industry-specific tuning (e.g., medical diagnostics), and mini-program integration—enables seamless service execution (e.g., generating travel itineraries with bookings). This ecosystem approach has driven rapid adoption. Nano AI [96] is China's first "super search agent" that autonomously plans tasks (e.g., travel itineraries, market reports) by integrating data from walled gardens. Its DeepSearch technology parses tables, formulas, and video comments, enabling cross-platform verification for reliable decision-making. Kimi [97] can process 200 K-context windows, ideal for academic paper analysis. Users highlight its semantic search for Chinese literature. DeepSeek Search [100] represents a paradigm shift in cost-efficient, open-source AI search. Quark DeepSearch [98] relies on the Qwen-QWQ [207] inference model. Unlike traditional search engines that rely on keyword matching, the model understands natural language and performs semantic analysis to more accurately grasp user intent. Other deep research products include MiroMind ODR [101] and Manus [102].

Domain-Specific AI Search Tools. MediSearch [103] provides evidence-based medical answers (e.g., drug interactions, treatment protocols), trusted by 74% of healthcare professionals for clinical decision support. Devv.ai [104] is a code-specific search engine offering real-time debugging snippets and GitHub integration. It supports Chinese queries but is limited to programming contexts. Consensus [105] accesses 200 M+ scientific papers, using NLP to extract hypotheses and methodologies. Researchers report 50% time savings in literature reviews.

Integrated AI Search Assistants WallerAI [106] is a browser-sidebar assistant that reads PDFs, videos, and webpages, enabling cross-document Q&A and content export. Bing Chat [107], deeply integrated into Edge's ecosystem, delivers citation-backed answers through real-time web indexing and source attribution, establishing a unified search-browser experience.

7. Challenges and Future Research

Despite notable progress, this field still faces many unresolved challenges, and there is considerable room for improvement. We highlight several promising directions based on the reviewed progress:

- **Methods:** More complex problems lead to a prolonged search process and additional actions, resulting in an extended search context. This extended context can limit the effectiveness of AI

search methods and the ability of LLMs, causing search performance to degrade as the inference length increases.

- **Evaluations:** There is a strong need for systematic and standardized evaluation frameworks in AI search. The datasets used for evaluation should be designed to closely resemble real-world scenarios, featuring complex, dynamic, and citation-supported answers.
- **Applications:** The potential real-world applications of AI Search are significant. Beyond user scenarios, there are numerous applications across various industries. We hope to see the development of more AI search software and products to enhance the interaction between humans and machines.

8. Conclusions

Seeking and accessing information is a fundamental daily need for humans. In this survey, we provide a comprehensive overview of the latest research on AI Search based on LLMs. Our objective is to identify and highlight areas that require further research and suggest potential avenues for future studies. We begin by introducing the traditional information retrieval systems, LLMs, and AI Search based on LLMs. Subsequently, we classify existing studies into four categories: Text-based AI Search, Web Browsing Agent, Multimodal AI Search, and Benchmarks. We then highlight a range of current and significant software and products within the realm of AI search. Finally, we discuss the limitations of current AI search methods and explore promising future directions.

References

1. Brin, S.; Page, L. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems* **1998**, *30*, 107–117.
2. Berkhin, P. A survey on PageRank computing. *Internet mathematics* **2005**, *2*, 73–120.
3. Burges, C.; Shaked, T.; Renshaw, E.; Lazier, A.; Deeds, M.; Hamilton, N.; Hullender, G. Learning to rank using gradient descent. In Proceedings of the Proceedings of the 22nd international conference on Machine learning, 2005, pp. 89–96.
4. Nadkarni, P.M.; Ohno-Machado, L.; Chapman, W.W. Natural language processing: an introduction. *Journal of the American Medical Informatics Association* **2011**, *18*, 544–551.
5. Kobayashi, M.; Takeda, K. Information retrieval on the web. *ACM computing surveys (CSUR)* **2000**, *32*, 144–173.
6. Zhao, W.X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. A survey of large language models. *arXiv preprint arXiv:2303.18223* **2023**, *1*.
7. Research, C.D. <https://openai.com/index/introducing-deep-research>, 2022.
8. Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. Llama: Open and efficient foundation language models. *arXiv* **2023**, arXiv:2302.13971.
9. Lewis, P.S.H.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; tau Yih, W.; Rocktäschel, T.; et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In Proceedings of the Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual, 2020.
10. Xu, F.; Shi, W.; Choi, E. RECOMP: Improving Retrieval-Augmented LMs with Compression and Selective Augmentation. *CoRR* **2023**, *abs/2310.04408*, [2310.04408]. <https://doi.org/10.48550/ARXIV.2310.04408>.
11. Jiang, H.; Wu, Q.; Lin, C.Y.; Yang, Y.; Qiu, L. LLMingua: Compressing Prompts for Accelerated Inference of Large Language Models. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2023, pp. 13358–13376. <https://doi.org/10.18653/v1/2023.emnlp-main.825>.
12. Jin, J.; Li, X.; Dong, G.; Zhang, Y.; Zhu, Y.; Wu, Y.; Li, Z.; Ye, Q.; Dou, Z. Hierarchical Document Refinement for Long-context Retrieval-augmented Generation, 2025, [arXiv:cs.CL/2505.10413].
13. Kim, G.; Kim, S.; Jeon, B.; Park, J.; Kang, J. Tree of Clarifications: Answering Ambiguous Questions with Retrieval-Augmented Large Language Models. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Singapore, 2023; pp. 996–1009. <https://doi.org/10.18653/v1/2023.emnlp-main.63>.

14. Shi, W.; Min, S.; Yasunaga, M.; Seo, M.; James, R.; Lewis, M.; Zettlemoyer, L.; tau Yih, W. REPLUG: Retrieval-Augmented Black-Box Language Models. *CoRR* **2023**, *abs/2301.12652*, [2301.12652]. <https://doi.org/10.48550/ARXIV.2301.12652>.
15. Yu, W.; Iter, D.; Wang, S.; Xu, Y.; Ju, M.; Sanyal, S.; Zhu, C.; Zeng, M.; Jiang, M. Generate rather than retrieve: Large language models are strong context generators. *arXiv preprint arXiv:2209.10063* **2022**.
16. Wang, H.; Zhao, T.; Gao, J. BlendFilter: Advancing Retrieval-Augmented Large Language Models via Query Generation Blending and Knowledge Filtering, 2024, [arXiv:cs.CL/2402.11129].
17. Wang, Y.; Li, P.; Sun, M.; Liu, Y. Self-knowledge guided retrieval augmentation for large language models. *arXiv preprint arXiv:2310.05002* **2023**.
18. Wang, H.; Xue, B.; Zhou, B.; Zhang, T.; Wang, C.; Chen, G.; Wang, H.; Wong, K.f. Self-DC: When to retrieve and When to generate? Self Divide-and-Conquer for Compositional Unknown Questions. *arXiv preprint arXiv:2402.13514* **2024**.
19. Ding, H.; Pang, L.; Wei, Z.; Shen, H.; Cheng, X. Retrieve Only When It Needs: Adaptive Retrieval Augmentation for Hallucination Mitigation in Large Language Models, 2024, [arXiv:cs.CL/2402.10612].
20. Tan, J.; Dou, Z.; Zhu, Y.; Guo, P.; Fang, K.; Wen, J.R. Small Models, Big Insights: Leveraging Slim Proxy Models To Decide When and What to Retrieve for LLMs. *CoRR* **2024**, *abs/2402.12052*, [2402.12052]. <https://doi.org/10.48550/ARXIV.2402.12052>.
21. Yao, S.; Zhao, J.; Yu, D.; Shafran, I.; Narasimhan, K.R.; Cao, Y. ReAct: Synergizing Reasoning and Acting in Language Models. In Proceedings of the NeurIPS 2022 Foundation Models for Decision Making Workshop, 2022.
22. Shao, Z.; Gong, Y.; Shen, Y.; Huang, M.; Duan, N.; Chen, W. Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy, 2023, [arXiv:cs.CL/2305.15294].
23. Trivedi, H.; Balasubramanian, N.; Khot, T.; Sabharwal, A. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. *arXiv preprint arXiv:2212.10509* **2022**.
24. Jiang, Z.; Xu, F.F.; Gao, L.; Sun, Z.; Liu, Q.; Dwivedi-Yu, J.; Yang, Y.; Callan, J.; Neubig, G. Active Retrieval Augmented Generation. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023. Association for Computational Linguistics, 2023, pp. 7969–7992.
25. Asai, A.; Wu, Z.; Wang, Y.; Sil, A.; Hajishirzi, H. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. *CoRR* **2023**, *abs/2310.11511*, [2310.11511]. <https://doi.org/10.48550/ARXIV.2310.11511>.
26. Li, X.; Dong, G.; Jin, J.; Zhang, Y.; Zhou, Y.; Zhu, Y.; Zhang, P.; Dou, Z. Search-o1: Agentic Search-Enhanced Large Reasoning Models. *CoRR* **2025**, *abs/2501.05366*, [2501.05366]. <https://doi.org/10.48550/ARXIV.2501.05366>.
27. Li, X.; Jin, J.; Dong, G.; Qian, H.; Zhu, Y.; Wu, Y.; Wen, J.; Dou, Z. WebThinker: Empowering Large Reasoning Models with Deep Research Capability. *CoRR* **2025**, *abs/2504.21776*, [2504.21776]. <https://doi.org/10.48550/ARXIV.2504.21776>.
28. Wu, J.; Li, B.; Fang, R.; Yin, W.; Zhang, L.; Tao, Z.; Zhang, D.; Xi, Z.; Jiang, Y.; Xie, P.; et al. WebDancer: Towards Autonomous Information Seeking Agency, 2025, [arXiv:cs.CL/2505.22648].
29. Huang, L.; Liu, Y.; Jiang, J.; Zhang, R.; Yan, J.; Li, J.; Zhao, W.X. ManuSearch: Democratizing Deep Search in Large Language Models with a Transparent and Open Multi-Agent Framework, 2025, [arXiv:cs.CL/2505.18105].
30. Yang, T.; Yao, Z.; Jin, B.; Cui, L.; Li, Y.; Wang, G.; Liu, X. Demystifying and Enhancing the Efficiency of Large Language Model Based Search Agents, 2025, [arXiv:cs.AI/2505.12065].
31. Jin, J.; Li, X.; Dong, G.; Zhang, Y.; Zhu, Y.; Zhao, Y.; Qian, H.; Dou, Z. Decoupled Planning and Execution: A Hierarchical Reasoning Framework for Deep Search, 2025, [arXiv:cs.AI/2507.02652].
32. Wu, W.; Guan, X.; Huang, S.; Jiang, Y.; Xie, P.; Huang, F.; Cao, J.; Zhao, H.; Zhou, J. MaskSearch: A Universal Pre-Training Framework to Enhance Agentic Search Capability, 2025, [arXiv:cs.CL/2505.20285].
33. Wang, L.; Chen, H.; Yang, N.; Huang, X.; Dou, Z.; Wei, F. Chain-of-Retrieval Augmented Generation, 2025, [arXiv:cs.IR/2501.14342].
34. Lee, Z.; Cao, S.; Liu, J.; Zhang, J.; Liu, W.; Che, X.; Hou, L.; Li, J. ReaRAG: Knowledge-guided Reasoning Enhances Factuality of Large Reasoning Models with Iterative Retrieval Augmented Generation, 2025, [arXiv:cs.CL/2503.21729].
35. Shi, Z.; Yan, L.; Yin, D.; Verberne, S.; de Rijke, M.; Ren, Z. Iterative Self-Incentivization Empowers Large Language Models as Agentic Searchers, 2025, [arXiv:cs.CL/2505.20128].

36. Hu, M.; Fang, T.; Zhang, J.; Ma, J.; Zhang, Z.; Zhou, J.; Zhang, H.; Mi, H.; Yu, D.; King, I. WebCoT: Enhancing Web Agent Reasoning by Reconstructing Chain-of-Thought in Reflection, Branching, and Rollback, 2025, [\[arXiv:cs.CL/2505.20013\]](#).
37. Guan, X.; Zeng, J.; Meng, F.; Xin, C.; Lu, Y.; Lin, H.; Han, X.; Sun, L.; Zhou, J. DeepRAG: Thinking to Retrieve Step by Step for Large Language Models, 2025, [\[arXiv:cs.AI/2502.01142\]](#).
38. Jin, B.; Zeng, H.; Yue, Z.; Yoon, J.; Arik, S.; Wang, D.; Zamani, H.; Han, J. Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning, 2025, [\[arXiv:cs.CL/2503.09516\]](#).
39. Song, H.; Jiang, J.; Min, Y.; Chen, J.; Chen, Z.; Zhao, W.X.; Fang, L.; Wen, J.R. R1-Searcher: Incentivizing the Search Capability in LLMs via Reinforcement Learning, 2025, [\[arXiv:cs.AI/2503.05592\]](#).
40. Chen, M.; Li, T.; Sun, H.; Zhou, Y.; Zhu, C.; Wang, H.; Pan, J.Z.; Zhang, W.; Chen, H.; Yang, F.; et al. ReSearch: Learning to Reason with Search for LLMs via Reinforcement Learning, 2025, [\[arXiv:cs.AI/2503.19470\]](#).
41. Sun, H.; Qiao, Z.; Guo, J.; Fan, X.; Hou, Y.; Jiang, Y.; Xie, P.; Zhang, Y.; Huang, F.; Zhou, J. ZeroSearch: Incentivize the Search Capability of LLMs without Searching, 2025, [\[arXiv:cs.CL/2505.04588\]](#).
42. Zheng, Y.; Fu, D.; Hu, X.; Cai, X.; Ye, L.; Lu, P.; Liu, P. DeepResearcher: Scaling Deep Research via Reinforcement Learning in Real-world Environments, 2025, [\[arXiv:cs.AI/2504.03160\]](#).
43. Mei, J.; Hu, T.; Fu, D.; Wen, L.; Yang, X.; Wu, R.; Cai, P.; Cai, X.; Gao, X.; Yang, Y.; et al. O²-Searcher: A Searching-based Agent Model for Open-Domain Open-Ended Question Answering, 2025, [\[arXiv:cs.CL/2505.16582\]](#).
44. Li, K.; Zhang, Z.; Yin, H.; Zhang, L.; Ou, L.; Wu, J.; Yin, W.; Li, B.; Tao, Z.; Wang, X.; et al. WebSailor: Navigating Super-human Reasoning for Web Agent, 2025, [\[arXiv:cs.CL/2507.02592\]](#).
45. Tao, Z.; Wu, J.; Yin, W.; Zhang, J.; Li, B.; Shen, H.; Li, K.; Zhang, L.; Wang, X.; Jiang, Y.; et al. WebShaper: Agentically Data Synthesizing via Information-Seeking Formalization. *arXiv preprint arXiv:2507.15061* 2025.
46. Wang, Z.; Zheng, X.; An, K.; Ouyang, C.; Cai, J.; Wang, Y.; Wu, Y. StepSearch: Igniting LLMs Search Ability via Step-Wise Proximal Policy Optimization, 2025, [\[arXiv:cs.CL/2505.15107\]](#).
47. Wei, Z.; Yao, W.; Liu, Y.; Zhang, W.; Lu, Q.; Qiu, L.; Yu, C.; Xu, P.; Zhang, C.; Yin, B.; et al. WebAgent-R1: Training Web Agents via End-to-End Multi-Turn Reinforcement Learning, 2025, [\[arXiv:cs.CL/2505.16421\]](#).
48. Lai, H.; Liu, X.; Iong, I.L.; Yao, S.; Chen, Y.; Shen, P.; Yu, H.; Zhang, H.; Zhang, X.; Dong, Y.; et al. AutoWebGLM: A Large Language Model-based Web Navigating Agent. In Proceedings of the Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, New York, NY, USA, 2024; KDD '24, p. 5295–5306. <https://doi.org/10.1145/3637528.3671620>.
49. Song, Y.; Xu, F.; Zhou, S.; Neubig, G. Beyond Browsing: API-Based Web Agents, 2025, [\[arXiv:cs.CL/2410.16464\]](#).
50. Zhang, D.; Rama, B.; Ni, J.; He, S.; Zhao, F.; Chen, K.; Chen, A.; Cao, J. LiteWebAgent: The Open-Source Suite for VLM-Based Web-Agent Applications, 2025, [\[arXiv:cs.AI/2503.02950\]](#).
51. Yang, K.; Liu, Y.; Chaudhary, S.; Fakoor, R.; Chaudhari, P.; Karypis, G.; Rangwala, H. AgentOccam: A Simple Yet Strong Baseline for LLM-Based Web Agents, 2025, [\[arXiv:cs.AI/2410.13825\]](#).
52. Shi, W.; Tan, H.; Kuang, C.; Li, X.; Ren, X.; Zhang, C.; Chen, H.; Wang, Y.; Shang, L.; Yu, F.; et al. Pangu DeepDiver: Adaptive Search Intensity Scaling via Open-Web Reinforcement Learning, 2025, [\[arXiv:cs.CL/2505.24332\]](#).
53. He, H.; Yao, W.; Ma, K.; Yu, W.; Dai, Y.; Zhang, H.; Lan, Z.; Yu, D. WebVoyager: Building an End-to-End Web Agent with Large Multimodal Models. In Proceedings of the Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); Ku, L.W.; Martins, A.; Srikumar, V., Eds., Bangkok, Thailand, 2024; pp. 6864–6890. <https://doi.org/10.18653/v1/2024.acl-long.371>.
54. Wu, J.; Li, B.; Fang, R.; Yin, W.; Zhang, L.; Tao, Z.; Zhang, D.; Xi, Z.; Jiang, Y.; Xie, P.; et al. WebDancer: Towards Autonomous Information Seeking Agency, 2025, [\[arXiv:cs.CL/2505.22648\]](#).
55. Shen, H.; Liu, C.; Li, G.; Wang, X.; Zhou, Y.; Ma, C.; Ji, X. Falcon-UI: Understanding GUI Before Following User Instructions, 2024, [\[arXiv:cs.CL/2412.09362\]](#).
56. Cho, J.; Kim, J.; Bae, D.; Choo, J.; Gwon, Y.; Kwon, Y.D. CAAP: Context-Aware Action Planning Prompting to Solve Computer Tasks with Front-End UI Only, 2024, [\[arXiv:cs.AI/2406.06947\]](#).
57. Lin, K.Q.; Li, L.; Gao, D.; Yang, Z.; Wu, S.; Bai, Z.; Lei, W.; Wang, L.; Shou, M.Z. ShowUI: One Vision-Language-Action Model for GUI Visual Agent, 2024, [\[arXiv:cs.CV/2411.17465\]](#).
58. Cao, T.; Huang, C.; Li, Y.; Huilin, W.; He, A.; Oo, N.; Hooi, B. PhishAgent: A Robust Multimodal Agent for Phishing Webpage Detection. *Proceedings of the AAAI Conference on Artificial Intelligence* 2025, 39, 27869–27877. <https://doi.org/10.1609/aaai.v39i27.35003>.

59. Jiang, D.; Zhang, R.; Guo, Z.; Wu, Y.; Qiu, P.; Lu, P.; Chen, Z.; Song, G.; Gao, P.; Liu, Y.; et al. MMSearch: Unveiling the Potential of Large Models as Multi-modal Search Engines. In Proceedings of the The Thirteenth International Conference on Learning Representations.
60. Wu, J.; Deng, Z.; Li, W.; Liu, Y.; You, B.; Li, B.; Ma, Z.; Liu, Z. MMSearch-R1: Incentivizing LMMs to Search, 2025, [arXiv:cs.CV/2506.20670].
61. Pahuja, V.; Lu, Y.; Rosset, C.; Gou, B.; Mitra, A.; Whitehead, S.; Su, Y.; Awadallah, A. Explorer: Scaling exploration-driven web trajectory synthesis for multimodal web agents. *arXiv* **2025**, arXiv:2502.11357.
62. Verma, G.; Kaur, R.; Srishankar, N.; Zeng, Z.; Balch, T.; Veloso, M. AdaptAgent: Adapting Multimodal Web Agents with Few-Shot Learning from Human Demonstrations, 2024, [arXiv:cs.AI/2411.13451].
63. He, H.; Yao, W.; Ma, K.; Yu, W.; Zhang, H.; Fang, T.; Lan, Z.; Yu, D. OpenWebVoyager: Building Multimodal Web Agents via Iterative Real-World Exploration, Feedback and Optimization, 2024, [arXiv:cs.CL/2410.19609].
64. Zheng, B.; Gou, B.; Kil, J.; Sun, H.; Su, Y. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614* **2024**.
65. Gadiraju, S.S.; Liao, D.; Kudupudi, A.; Kasula, S.; Chalasani, C. InfoTech Assistant: A Multimodal Conversational Agent for InfoTechnology Web Portal Queries. In Proceedings of the 2024 IEEE International Conference on Big Data (BigData), 2024, pp. 3264–3272. <https://doi.org/10.1109/BigData62323.2024.10825668>.
66. Kwiatkowski, T.; Palomaki, J.; Redfield, O.; Collins, M.; Parikh, A.; Alberti, C.; Epstein, D.; Polosukhin, I.; Devlin, J.; Lee, K.; et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics* **2019**, 7, 453–466.
67. Joshi, M.; Choi, E.; Weld, D.S.; Zettlemoyer, L. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551* **2017**.
68. Mallen, A.; Asai, A.; Zhong, V.; Das, R.; Khashabi, D.; Hajishirzi, H. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv* **2022**, arXiv:2212.10511.
69. Yang, Z.; Qi, P.; Zhang, S.; Bengio, Y.; Cohen, W.W.; Salakhutdinov, R.; Manning, C.D. HotpotQA: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600* **2018**.
70. Trivedi, H.; Balasubramanian, N.; Khot, T.; Sabharwal, A. MuSiQue: Multihop Questions via Single-hop Question Composition. *Transactions of the Association for Computational Linguistics* **2022**, 10, 539–554.
71. Thorne, J.; Vlachos, A.; Christodoulopoulos, C.; Mittal, A. FEVER: a large-scale dataset for fact extraction and VERification. *arXiv preprint arXiv:1803.05355* **2018**.
72. Petroni, F.; Piktus, A.; Fan, A.; Lewis, P.; Yazdani, M.; De Cao, N.; Thorne, J.; Jernite, Y.; Karpukhin, V.; Maillard, J.; et al. KILT: a benchmark for knowledge intensive language tasks. *arXiv* **2020**, arXiv:2009.02252.
73. Mialon, G.; Fourrier, C.; Wolf, T.; LeCun, Y.; Scialom, T. Gaia: a benchmark for general ai assistants. In Proceedings of the The Twelfth International Conference on Learning Representations, 2023.
74. Wei, J.; Sun, Z.; Papay, S.; McKinney, S.; Han, J.; Fulford, I.; Chung, H.W.; Passos, A.T.; Fedus, W.; Glaese, A. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv* **2025**, arXiv:2504.12516.
75. Zhou, P.; Leon, B.; Ying, X.; Zhang, C.; Shao, Y.; Ye, Q.; Chong, D.; Jin, Z.; Xie, C.; Cao, M.; et al. Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese. *arXiv* **2025**, arXiv:2504.19314.
76. Gou, B.; Huang, Z.; Ning, Y.; Gu, Y.; Lin, M.; Qi, W.; Kopanov, A.; Yu, B.; Gutiérrez, B.J.; Shu, Y.; et al. Mind2Web 2: Evaluating Agentic Search with Agent-as-a-Judge, 2025, [arXiv:cs.AI/2506.21506].
77. Deng, X.; Gu, Y.; Zheng, B.; Chen, S.; Stevens, S.; Wang, B.; Sun, H.; Su, Y. Mind2Web: Towards a Generalist Agent for the Web. In Proceedings of the Advances in Neural Information Processing Systems; Oh, A.; Naumann, T.; Globerson, A.; Saenko, K.; Hardt, M.; Levine, S., Eds. Curran Associates, Inc., 2023, Vol. 36, pp. 28091–28114.
78. Zhou, S.; Xu, F.F.; Zhu, H.; Zhou, X.; Lo, R.; Sridhar, A.; Cheng, X.; Ou, T.; Bisk, Y.; Fried, D.; et al. WebArena: A Realistic Web Environment for Building Autonomous Agents, 2024, [arXiv:cs.AI/2307.13854].
79. Koh, J.Y.; Lo, R.; Jang, L.; Duvvur, V.; Lim, M.; Huang, P.Y.; Neubig, G.; Zhou, S.; Salakhutdinov, R.; Fried, D. VisualWebArena: Evaluating Multimodal Agents on Realistic Visual Web Tasks. In Proceedings of the Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); Ku, L.W.; Martins, A.; Srikumar, V., Eds., Bangkok, Thailand, 2024; pp. 881–905. <https://doi.org/10.18653/v1/2024.acl-long.50>.
80. Garg, D.; VanWeelden, S.; Caples, D.; Draguns, A.; Ravi, N.; Putta, P.; Garg, N.; Abraham, T.; Lara, M.; Lopez, F.; et al. REAL: Benchmarking Autonomous Agents on Deterministic Simulations of Real Websites, 2025, [arXiv:cs.AI/2504.11543].

81. Miyai, A.; Zhao, Z.; Egashira, K.; Sato, A.; Sunada, T.; Onohara, S.; Yamanishi, H.; Toyooka, M.; Nishina, K.; Maeda, R.; et al. WebChoreArena: Evaluating Web Browsing Agents on Realistic Tedious Web Tasks, 2025, [arXiv:cs.CL/2506.01952].
82. Song, Y.; Thai, K.; Pham, C.M.; Chang, Y.; Nadaf, M.; Iyyer, M. BEARCUBS: A benchmark for computer-using web agents, 2025, [arXiv:cs.AI/2503.07919].
83. Pan, Y.; Kong, D.; Zhou, S.; Cui, C.; Leng, Y.; Jiang, B.; Liu, H.; Shang, Y.; Zhou, S.; Wu, T.; et al. WebCanvas: Benchmarking Web Agents in Online Environments, 2024, [arXiv:cs.CL/2406.12373].
84. Xu, K.; Kordi, Y.; Nayak, T.; Asija, A.; Wang, Y.; Sanders, K.; Byerly, A.; Zhang, J.; Van Durme, B.; Khashabi, D. TurkingBench: A Challenge Benchmark for Web Agents. In Proceedings of the Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers); Chiruzzo, L.; Ritter, A.; Wang, L., Eds., Albuquerque, New Mexico, 2025; pp. 3694–3710. <https://doi.org/10.18653/v1/2025.naacl-long.188>.
85. Tur, A.D.; Meade, N.; Lù, X.H.; Zambrano, A.; Patel, A.; Durmus, E.; Gella, S.; Stańczak, K.; Reddy, S. SafeArena: Evaluating the Safety of Autonomous Web Agents, 2025, [arXiv:cs.LG/2503.04957].
86. Zhu, Y.; Kellermann, A.; Bowman, D.; Li, P.; Gupta, A.; Danda, A.; Fang, R.; Jensen, C.; Ihli, E.; Benn, J.; et al. CVE-Bench: A Benchmark for AI Agents' Ability to Exploit Real-World Web Application Vulnerabilities, 2025, [arXiv:cs.CR/2503.17332].
87. Evtimov, I.; Zharmagambetov, A.; Grattafiori, A.; Guo, C.; Chaudhuri, K. WASP: Benchmarking Web Agent Security Against Prompt Injection Attacks, 2025, [arXiv:cs.CR/2504.18575].
88. Qiu, H.; Fabbri, A.; Agarwal, D.; Huang, K.H.; Tan, S.; Peng, N.; Wu, C.S. Evaluating Cultural and Social Awareness of LLM Web Agents. In Proceedings of the Findings of the Association for Computational Linguistics: NAACL 2025; Chiruzzo, L.; Ritter, A.; Wang, L., Eds., Albuquerque, New Mexico, 2025; pp. 3978–4005. <https://doi.org/10.18653/v1/2025.findings-naacl.222>.
89. Luo, Y.; Li, Z.; Liu, J.; Cui, J.; Zhao, X.; Shen, Z. Open CaptchaWorld: A Comprehensive Web-based Platform for Testing and Benchmarking Multimodal LLM Agents, 2025, [arXiv:cs.AI/2505.24878].
90. Lyu, Y.; Zhang, X.; Yan, L.; de Rijke, M.; Ren, Z.; Chen, X. DeepShop: A Benchmark for Deep Research Shopping Agents, 2025, [arXiv:cs.IR/2506.02839].
91. Fu, M.; Peng, Y.; Liu, B.; Wan, Y.; Chen, D. LiveVQA: Live Visual Knowledge Seeking. *arXiv preprint arXiv:2504.05288* 2025.
92. Yu, Q.; Xiao, Z.; Li, B.; Wang, Z.; Chen, C.; Zhang, W. MRAMG-Bench: A BeyondText Benchmark for Multimodal Retrieval-Augmented Multimodal Generation. *arXiv preprint arXiv:2502.04176* 2025.
93. Research, P.D. <https://www.perplexity.ai>, 2022.
94. Doubao. <https://www.doubao.com>, 2023.
95. Yuanbao. <https://yuanbao.tencent.com>, 2024.
96. AI, N. <https://www.n.cn>, 2025.
97. Kimi. <https://www.kimi.com>, 2023.
98. DeepSearch, Q. <https://quark.sm.cn>, 2025.
99. Research, G.D. <https://gemini.google/overview/deep-research>, 2023.
100. DeepSeek. <https://www.deepseek.com>, 2025.
101. Team, M.A. MiroThinker: An open-source agentic model series trained for deep research and complex, long-horizon problem solving. <https://github.com/MiroMindAI/MiroThinker>, 2025.
102. Manus. <https://manus.im/>, 2025.
103. MediSearch. <https://medisearch.io>, 2023.
104. Devv.ai. <https://devv.ai/zh>, 2023.
105. Consensus. <https://consensus.app>, 2022.
106. walles.ai. <https://walles.ai/>, 2023.
107. Chat, B. <http://bing.com>, 2023.
108. Gao, Y.; Xiong, Y.; Gao, X.; Jia, K.; Pan, J.; Bi, Y.; Dai, Y.; Sun, J.; Guo, Q.; Wang, M.; et al. Retrieval-Augmented Generation for Large Language Models: A Survey, 2024, [arXiv:cs.CL/2312.10997].
109. Zhu, Y.; Yuan, H.; Wang, S.; Liu, J.; Liu, W.; Deng, C.; Chen, H.; Liu, Z.; Dou, Z.; Wen, J.R. Large language models for information retrieval: A survey. *arXiv preprint arXiv:2308.07107* 2023.
110. Ramos, J.; et al. Using tf-idf to determine word relevance in document queries. In Proceedings of the Proceedings of the first instructional conference on machine learning. Citeseer, 2003, Vol. 242, pp. 29–48.
111. Robertson, S.E.; Zaragoza, H. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.* 2009, 3, 333–389. <https://doi.org/10.1561/15000000019>.

112. Karpukhin, V.; Oguz, B.; Min, S.; Lewis, P.; Wu, L.; Edunov, S.; Chen, D.; Yih, W.t. Dense Passage Retrieval for Open-Domain Question Answering. In Proceedings of the EMNLP, 2020, pp. 6769–6781.
113. Xiong, L.; Xiong, C.; Li, Y.; Tang, K.F.; Liu, J.; Bennett, P.N.; Ahmed, J.; Overwijk, A. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. In Proceedings of the ICLR, 2020.
114. Wang, L.; Yang, N.; Huang, X.; Jiao, B.; Yang, L.; Jiang, D.; Majumder, R.; Wei, F. Text Embeddings by Weakly-Supervised Contrastive Pre-training. *CoRR* **2022**, *abs/2212.03533*, [2212.03533]. <https://doi.org/10.48550/ARXIV.2212.03533>.
115. Xiao, S.; Liu, Z.; Zhang, P.; Muennighoff, N.; Lian, D.; Nie, J.Y. C-Pack: Packed Resources For General Chinese Embeddings, 2024, [arXiv:cs.CL/2309.07597].
116. Li, X.; Jin, J.; Zhou, Y.; Zhang, Y.; Zhang, P.; Zhu, Y.; Dou, Z. From matching to generation: A survey on generative information retrieval. *ACM Transactions on Information Systems* **2025**, *43*, 1–62.
117. Tay, Y.; Tran, V.; Dehghani, M.; Ni, J.; Bahri, D.; Mehta, H.; Qin, Z.; Hui, K.; Zhao, Z.; Gupta, J.P.; et al. Transformer Memory as a Differentiable Search Index. In Proceedings of the Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022.
118. Wang, Y.; Hou, Y.; Wang, H.; Miao, Z.; Wu, S.; Chen, Q.; Xia, Y.; Chi, C.; Zhao, G.; Liu, Z.; et al. A neural corpus indexer for document retrieval. *Advances in Neural Information Processing Systems* **2022**, *35*, 25600–25614.
119. Li, X.; Dou, Z.; Zhou, Y.; Liu, F. Corpuslm: Towards a unified language model on corpus for knowledge-intensive tasks. In Proceedings of the Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2024, pp. 26–37.
120. Liu, T.Y.; et al. Learning to rank for information retrieval. *Foundations and Trends® in Information Retrieval* **2009**, *3*, 225–331.
121. Khattab, O.; Zaharia, M. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In Proceedings of the Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25–30, 2020. ACM, 2020, pp. 39–48. <https://doi.org/10.1145/3397271.3401075>.
122. Sun, W.; Yan, L.; Ma, X.; Wang, S.; Ren, P.; Chen, Z.; Yin, D.; Ren, Z. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agents. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6–10, 2023; Bouamor, H.; Pino, J.; Bali, K., Eds. Association for Computational Linguistics, 2023, pp. 14918–14937. <https://doi.org/10.18653/V1/2023.EMNLP-MAIN.923>.
123. Izacard, G.; Grave, E. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282* **2020**.
124. Guu, K.; Lee, K.; Tung, Z.; Pasupat, P.; Chang, M. Retrieval augmented language model pre-training. In Proceedings of the International conference on machine learning. PMLR, 2020, pp. 3929–3938.
125. Borgeaud, S.; Mensch, A.; Hoffmann, J.; Cai, T.; Rutherford, E.; Millican, K.; van den Driessche, G.; Lespiau, J.B.; Damoc, B.; Clark, A.; et al. Improving Language Models by Retrieving from Trillions of Tokens. In Proceedings of the International Conference on Machine Learning, ICML 2022, 17–23 July 2022, Baltimore, Maryland, USA. PMLR, 2022, Vol. 162, *Proceedings of Machine Learning Research*, pp. 2206–2240.
126. Ma, X.; Gong, Y.; He, P.; Zhao, H.; Duan, N. Query Rewriting in Retrieval-Augmented Large Language Models. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Singapore, 2023; pp. 5303–5315. <https://doi.org/10.18653/v1/2023.emnlp-main.322>.
127. Ram, O.; Levine, Y.; Dalmedigos, I.; Muhlgay, D.; Shashua, A.; Leyton-Brown, K.; Shoham, Y. In-context retrieval-augmented language models. *arXiv preprint arXiv:2302.00083* **2023**.
128. Yu, Z.; Xiong, C.; Yu, S.; Liu, Z. Augmentation-adapted retriever improves generalization of language models as generic plug-in. *arXiv preprint arXiv:2305.17331* **2023**.
129. Zhang, P.; Xiao, S.; Liu, Z.; Dou, Z.; Nie, J.Y. Retrieve Anything To Augment Large Language Models. *CoRR* **2023**, *abs/2310.07554*, [2310.07554]. <https://doi.org/10.48550/ARXIV.2310.07554>.
130. Zhang, L.; Yu, Y.; Wang, K.; Zhang, C. ARL2: Aligning Retrievers for Black-box Large Language Models via Self-guided Adaptive Relevance Labeling, 2024, [arXiv:cs.CL/2402.13542].
131. Liu, N.F.; Lin, K.; Hewitt, J.; Paranjape, A.; Bevilacqua, M.; Petroni, F.; Liang, P. Lost in the Middle: How Language Models Use Long Contexts, 2023. arXiv:2307.03172.
132. Cuconasu, F.; Trappolini, G.; Siciliano, F.; Filice, S.; Campagnano, C.; Maarek, Y.; Tonellotto, N.; Silvestri, F. The Power of Noise: Redefining Retrieval for RAG Systems, 2024, [arXiv:cs.IR/2401.14887].

133. Yang, H.; Li, Z.; Zhang, Y.; Wang, J.; Cheng, N.; Li, M.; Xiao, J. PRCA: Fitting Black-Box Large Language Models for Retrieval Question Answering via Pluggable Reward-Driven Contextual Adapter. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023. Association for Computational Linguistics, 2023, pp. 5364–5375.
134. Zhang, Y.; Wang, T.; Chen, S.; Wang, K.; Zeng, X.; Lin, H.; Han, X.; Sun, L.; Lu, C. ARise: Towards Knowledge-Augmented Reasoning via Risk-Adaptive Search. *arXiv preprint arXiv:2504.10893* 2025.
135. Press, O.; Zhang, M.; Min, S.; Schmidt, L.; Smith, N.; Lewis, M. Measuring and Narrowing the Compositionality Gap in Language Models. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, 2023; pp. 5687–5711. <https://doi.org/10.18653/v1/2023.findings-emnlp.378>.
136. Yoran, O.; Wolfson, T.; Ram, O.; Berant, J. Making Retrieval-Augmented Language Models Robust to Irrelevant Context, 2023, [\[arXiv:cs.CL/2310.01558\]](https://arxiv.org/abs/2310.01558).
137. Huang, Y.; Chen, Y.; Zhang, H.; Li, K.; Fang, M.; Yang, L.; Li, X.; Shang, L.; Xu, S.; Hao, J.; et al. Deep Research Agents: A Systematic Examination And Roadmap, 2025, [\[arXiv:cs.AI/2506.18096\]](https://arxiv.org/abs/2506.18096).
138. Zhang, W.; Li, Y.; Bei, Y.; Luo, J.; Wan, G.; Yang, L.; Xie, C.; Yang, Y.; Huang, W.C.; Miao, C.; et al. From Web Search towards Agentic Deep Research: Incentivizing Search with Reasoning Agents, 2025, [\[arXiv:cs.IR/2506.18959\]](https://arxiv.org/abs/2506.18959).
139. Xu, R.; Peng, J. A Comprehensive Survey of Deep Research: Systems, Methodologies, and Applications, 2025, [\[arXiv:cs.AI/2506.12594\]](https://arxiv.org/abs/2506.12594).
140. Sun, S.; Song, H.; Wang, Y.; Ren, R.; Jiang, J.; Zhang, J.; Bai, F.; Deng, J.; Zhao, W.X.; Liu, Z.; et al. SimpleDeepSearcher: Deep Information Seeking via Web-Powered Reasoning Trajectory Synthesis, 2025, [\[arXiv:cs.CL/2505.16834\]](https://arxiv.org/abs/2505.16834).
141. Rafailov, R.; Sharma, A.; Mitchell, E.; Ermon, S.; Manning, C.D.; Finn, C. Direct Preference Optimization: Your Language Model is Secretly a Reward Model, 2024, [\[arXiv:cs.LG/2305.18290\]](https://arxiv.org/abs/2305.18290).
142. Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; Klimov, O. Proximal Policy Optimization Algorithms, 2017, [\[arXiv:cs.LG/1707.06347\]](https://arxiv.org/abs/1707.06347).
143. Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Bi, X.; Zhang, H.; Zhang, M.; Li, Y.K.; Wu, Y.; et al. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models, 2024, [\[arXiv:cs.CL/2402.03300\]](https://arxiv.org/abs/2402.03300).
144. Hu, J.; Liu, J.K.; Shen, W. REINFORCE++: An Efficient RLHF Algorithm with Robustness to Both Prompt and Reward Models, 2025, [\[arXiv:cs.CL/2501.03262\]](https://arxiv.org/abs/2501.03262).
145. Li, K.; Zhang, Z.; Yin, H.; Zhang, L.; Ou, L.; Wu, J.; Yin, W.; Li, B.; Tao, Z.; Wang, X.; et al. WebSailor: Navigating Super-human Reasoning for Web Agent, 2025, [\[arXiv:cs.CL/2507.02592\]](https://arxiv.org/abs/2507.02592).
146. Shi, W.; Tan, H.; Kuang, C.; Li, X.; Ren, X.; Zhang, C.; Chen, H.; Wang, Y.; Shang, L.; Yu, F.; et al. Pangu DeepDiver: Adaptive Search Intensity Scaling via Open-Web Reinforcement Learning, 2025, [\[arXiv:cs.CL/2505.24332\]](https://arxiv.org/abs/2505.24332).
147. Shi, Y.; Li, S.; Wu, C.; Liu, Z.; Fang, J.; Cai, H.; Zhang, A.; Wang, X. Search and Refine During Think: Autonomous Retrieval-Augmented Reasoning of LLMs, 2025, [\[arXiv:cs.CL/2505.11277\]](https://arxiv.org/abs/2505.11277).
148. Dong, G.; Chen, Y.; Li, X.; Jin, J.; Qian, H.; Zhu, Y.; Mao, H.; Zhou, G.; Dou, Z.; Wen, J.R. Tool-Star: Empowering LLM-Brained Multi-Tool Reasoner via Reinforcement Learning, 2025, [\[arXiv:cs.CL/2505.16410\]](https://arxiv.org/abs/2505.16410).
149. Lin, C.; Wen, Y.; Su, D.; Sun, F.; Chen, M.; Bao, C.; Lv, Z. Knowledgeable-r1: Policy Optimization for Knowledge Exploration in Retrieval-Augmented Generation, 2025, [\[arXiv:cs.CL/2506.05154\]](https://arxiv.org/abs/2506.05154).
150. Qian, H.; Liu, Z. Scent of Knowledge: Optimizing Search-Enhanced Reasoning with Information Foraging, 2025, [\[arXiv:cs.CL/2505.09316\]](https://arxiv.org/abs/2505.09316).
151. Li, Y.; Luo, Q.; Li, X.; Li, B.; Cheng, Q.; Wang, B.; Zheng, Y.; Wang, Y.; Yin, Z.; Qiu, X. R3-RAG: Learning Step-by-Step Reasoning and Retrieval for LLMs via Reinforcement Learning, 2025, [\[arXiv:cs.CL/2505.23794\]](https://arxiv.org/abs/2505.23794).
152. Zhang, D.; Zhao, Y.; Wu, J.; Li, B.; Yin, W.; Zhang, L.; Jiang, Y.; Li, Y.; Tu, K.; Xie, P.; et al. EvolveSearch: An Iterative Self-Evolving Search Agent, 2025, [\[arXiv:cs.CL/2505.22501\]](https://arxiv.org/abs/2505.22501).
153. Sha, Z.; Cui, S.; Wang, W. SEM: Reinforcement Learning for Search-Efficient Large Language Models, 2025, [\[arXiv:cs.CL/2505.07903\]](https://arxiv.org/abs/2505.07903).
154. Wu, P.; Zhang, M.; Zhang, X.; Du, X.; Chen, Z.Z. Search Wisely: Mitigating Sub-optimal Agentic Searches By Reducing Uncertainty, 2025, [\[arXiv:cs.CL/2505.17281\]](https://arxiv.org/abs/2505.17281).
155. Jiang, P.; Xu, X.; Lin, J.; Xiao, J.; Wang, Z.; Sun, J.; Han, J. s3: You Don't Need That Much Data to Train a Search Agent via RL, 2025, [\[arXiv:cs.AI/2505.14146\]](https://arxiv.org/abs/2505.14146).
156. Zhang, C.; He, S.; Qian, J.; Li, B.; Li, L.; Qin, S.; Kang, Y.; Ma, M.; Liu, G.; Lin, Q.; et al. Large Language Model-Brained GUI Agents: A Survey. *Transactions on Machine Learning Research* 2025.

157. Ning, L.; Liang, Z.; Jiang, Z.; Qu, H.; Ding, Y.; Fan, W.; yong Wei, X.; Lin, S.; Liu, H.; Yu, P.S.; et al. A Survey of WebAgents: Towards Next-Generation AI Agents for Web Automation with Large Foundation Models, 2025, [arXiv:cs.AI/2503.23350].
158. Wang, L.; Ma, C.; Feng, X.; Zhang, Z.; Yang, H.; Zhang, J.; Chen, Z.; Tang, J.; Chen, X.; Lin, Y.; et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science* **2024**, *18*. <https://doi.org/10.1007/s11704-024-40231-1>.
159. Hong, S.; Zhuge, M.; Chen, J.; Zheng, X.; Cheng, Y.; Wang, J.; Zhang, C.; Wang, Z.; Yau, S.K.S.; Lin, Z.; et al. MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework. In Proceedings of the The Twelfth International Conference on Learning Representations, 2024.
160. Zhang, H.; Du, W.; Shan, J.; Zhou, Q.; Du, Y.; Tenenbaum, J.B.; Shu, T.; Gan, C. Building Cooperative Embodied Agents Modularly with Large Language Models. In Proceedings of the The Twelfth International Conference on Learning Representations, 2024.
161. Zhu, X.; Chen, Y.; Tian, H.; Tao, C.; Su, W.; Yang, C.; Huang, G.; Li, B.; Lu, L.; Wang, X.; et al. Ghost in the Minecraft: Generally Capable Agents for Open-World Environments via Large Language Models with Text-based Knowledge and Memory, 2023, [arXiv:cs.AI/2305.17144].
162. Chen, W.; Su, Y.; Zuo, J.; Yang, C.; Yuan, C.; Chan, C.M.; Yu, H.; Lu, Y.; Hung, Y.H.; Qian, C.; et al. AgentVerse: Facilitating Multi-Agent Collaboration and Exploring Emergent Behaviors, 2023, [arXiv:cs.CL/2308.10848].
163. Furuta, H.; Lee, K.H.; Nachum, O.; Matsuo, Y.; Faust, A.; Gu, S.S.; Gur, I. Multimodal Web Navigation with Instruction-Finetuned Foundation Models. In Proceedings of the The Twelfth International Conference on Learning Representations, 2024.
164. Tang, X.; Kim, K.; Song, Y.; Lothritz, C.; Li, B.; Ezzini, S.; Tian, H.; Klein, J.; Bissyandé, T.F. CodeAgent: Autonomous Communicative Agents for Code Review. In Proceedings of the Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing; Al-Onaizan, Y.; Bansal, M.; Chen, Y.N., Eds., Miami, Florida, USA, 2024; pp. 11279–11313. <https://doi.org/10.18653/v1/2024.emnlp-main.632>.
165. Alam, M.A.; Mahmud, S.; Mamun-or Rashid, M.; Khan, M.M. Optimizing node selection in search based multi-agent path finding. *Autonomous Agents and Multi-Agent Systems* **2025**, *39*, 1–24.
166. Dorri, A.; Kanhere, S.S.; Jurdak, R. Multi-agent systems: A survey. *Ieee Access* **2018**, *6*, 28573–28593.
167. Cao, Y.; Yu, W.; Ren, W.; Chen, G. An Overview of Recent Progress in the Study of Distributed Multi-Agent Coordination. *IEEE Transactions on Industrial Informatics* **2013**, *9*, 427–438. <https://doi.org/10.1109/TII.2012.2219061>.
168. Wu, J.; Yin, W.; Jiang, Y.; Wang, Z.; Xi, Z.; Fang, R.; Zhang, L.; He, Y.; Zhou, D.; Xie, P.; et al. WebWalker: Benchmarking LLMs in Web Traversal, 2025, [arXiv:cs.CL/2501.07572].
169. Qi, Z.; Liu, X.; Iong, I.L.; Lai, H.; Sun, X.; Zhao, W.; Yang, Y.; Yang, X.; Sun, J.; Yao, S.; et al. We-bRL: Training LLM Web Agents via Self-Evolving Online Curriculum Reinforcement Learning, 2025, [arXiv:cs.CL/2411.02337].
170. Zhuang, Y.; Jin, D.; Chen, J.; Shi, W.; Wang, H.; Zhang, C. WorkForceAgent-R1: Incentivizing Reasoning Capability in LLM-based Web Agents via Reinforcement Learning, 2025, [arXiv:cs.CL/2505.22942].
171. Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Bi, X.; Zhang, H.; Zhang, M.; Li, Y.K.; Wu, Y.; et al. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models, 2024, [arXiv:cs.CL/2402.03300].
172. Yu, Q.; Zhang, Z.; Zhu, R.; Yuan, Y.; Zuo, X.; Yue, Y.; Dai, W.; Fan, T.; Liu, G.; Liu, L.; et al. DAPO: An Open-Source LLM Reinforcement Learning System at Scale, 2025, [arXiv:cs.LG/2503.14476].
173. Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; ichter, b.; Xia, F.; Chi, E.; Le, Q.V.; Zhou, D. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In Proceedings of the Advances in Neural Information Processing Systems; Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; Oh, A., Eds. Curran Associates, Inc., 2022, Vol. 35, pp. 24824–24837.
174. Jia, B.; Manocha, D. Sim-to-Real Brush Manipulation using Behavior Cloning and Reinforcement Learning, 2023, [arXiv:cs.RO/2309.08457].
175. Luo, H.; Kuang, J.; Liu, W.; Shen, Y.; Luan, J.; Deng, Y. Browsing Like Human: A Multimodal Web Agent with Experiential Fast-and-Slow Thinking. In Proceedings of the Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers); Che, W.; Nabende, J.; Shutova, E.; Pilehvar, M.T., Eds., Vienna, Austria, 2025; pp. 14232–14251. <https://doi.org/10.18653/v1/2025.acl-long.697>.
176. Hu, M.; Fang, T.; Zhang, J.; Ma, J.; Zhang, Z.; Zhou, J.; Zhang, H.; Mi, H.; Yu, D.; King, I. WebCoT: Enhancing Web Agent Reasoning by Reconstructing Chain-of-Thought in Reflection, Branching, and Rollback, 2025, [arXiv:cs.CL/2505.20013].

177. Pal, K.K.; Kashihara, K.; Anantheswaran, U.; Kuznia, K.C.; Jagtap, S.; Baral, C. Exploring the Limits of Transfer Learning with Unified Model in the Cybersecurity Domain, 2023, [arXiv:cs.CL/2302.10346].
178. Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In Proceedings of the International Conference on Learning Representations, 2021.
179. Jiang, W.; Zhuang, Y.; Song, C.; Yang, X.; Zhou, J.T.; Zhang, C. AppAgentX: Evolving GUI Agents as Proficient Smartphone Users, 2025, [arXiv:cs.AI/2503.02268].
180. Alabastro, Z.M.; Ilagan, J.B.; To, L.A.; Ilagan, J.R. Applied Optical Character Recognition and Large Language Models in Augmenting Manual Business Processes for Data Analytics in Traditional Small Businesses with Minimal Digital Adoption. In Proceedings of the Artificial Intelligence in HCI; Degen, H.; Ntoa, S., Eds., Cham, 2025; pp. 267–276.
181. Wang, X.; Chen, Y.; Zhu, W. A Survey on Curriculum Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2022**, *44*, 4555–4576. <https://doi.org/10.1109/TPAMI.2021.3069908>.
182. Chen, W.; Koenig, S.; Dilkina, B. Reprompt: Planning by automatic prompt engineering for large language models agents. *arXiv preprint arXiv:2406.11132* **2024**.
183. Samarasekara, I.; Bandara, M.; Rabhi, F.; Benatallah, B.; Meymandpour, R. LLM driven approach for capability modelling with context-enriched prompt engineering **2024**.
184. Gangi Reddy, R.; Mukherjee, S.; Kim, J.; Wang, Z.; Hakkani-Tür, D.; Ji, H. Infogent: An Agent-Based Framework for Web Information Aggregation. In Proceedings of the Findings of the Association for Computational Linguistics: NAACL 2025; Chiruzzo, L.; Ritter, A.; Wang, L., Eds., Albuquerque, New Mexico, 2025; pp. 5745–5758. <https://doi.org/10.18653/v1/2025.findings-naacl.318>.
185. Koh, J.Y.; McAleer, S.; Fried, D.; Salakhutdinov, R. Tree Search for Language Model Agents, 2024, [arXiv:cs.AI/2407.01476].
186. Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F.L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* **2023**.
187. Jin, Y.; Li, J.; Liu, Y.; Gu, T.; Wu, K.; Jiang, Z.; He, M.; Zhao, B.; Tan, X.; Gan, Z.; et al. Efficient multimodal large language models: A survey. *arXiv preprint arXiv:2405.10739* **2024**.
188. OpenAI. Hello GPT-4o, 2024.
189. Anthropic. Claude 3.5 Sonnet, 2024.
190. Dai, W.; Li, J.; Li, D.; Tiong, A.M.H.; Zhao, J.; Wang, W.; Li, B.; Fung, P.N.; Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems* **2024**, *36*.
191. Li, J.; Li, D.; Savarese, S.; Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In Proceedings of the International conference on machine learning. PMLR, 2023, pp. 19730–19742.
192. Liu, H.; Li, C.; Wu, Q.; Lee, Y.J. Visual Instruction Tuning. In Proceedings of the NeurIPS, 2023.
193. Liu, H.; Li, C.; Li, Y.; Lee, Y.J. Improved baselines with visual instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 26296–26306.
194. Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923* **2025**.
195. Team, G.; Anil, R.; Borgeaud, S.; Wu, Y.; Alayrac, J.B.; Yu, J.; Soricut, R.; Schalkwyk, J.; Dai, A.M.; Hauth, A.; et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* **2023**.
196. Chen, Z.; Wu, J.; Wang, W.; Su, W.; Chen, G.; Xing, S.; Muyan, Z.; Zhang, Q.; Zhu, X.; Lu, L.; et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv preprint arXiv:2312.14238* **2023**.
197. Sun, Q.; Cui, Y.; Zhang, X.; Zhang, F.; Yu, Q.; Luo, Z.; Wang, Y.; Rao, Y.; Liu, J.; Huang, T.; et al. Generative multimodal models are in-context learners. *arXiv preprint arXiv:2312.13286* **2023**.
198. Li, J.; Lu, W.; Fei, H.; Luo, M.; Dai, M.; Xia, M.; Jin, Y.; Gan, Z.; Qi, D.; Fu, C.; et al. A survey on benchmarks of multimodal large language models. *arXiv preprint arXiv:2408.08632* **2024**.
199. Geng, X.; Xia, P.; Zhang, Z.; Wang, X.; Wang, Q.; Ding, R.; Wang, C.; Wu, J.; Zhao, Y.; Li, K.; et al. WebWatcher: Breaking New Frontiers of Vision-Language Deep Research Agent. *arXiv preprint arXiv:2508.05748* **2025**.
200. Xie, J.; Chen, Z.; Zhang, R.; Wan, X.; Li, G. Large Multimodal Agents: A Survey, 2024, [arXiv:cs.CV/2402.15116].

201. Hong, W.; Wang, W.; Lv, Q.; Xu, J.; Yu, W.; Ji, J.; Wang, Y.; Wang, Z.; Dong, Y.; Ding, M.; et al. Cogagent: A visual language model for gui agents. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 14281–14290.
202. He, H.; Yao, W.; Ma, K.; Yu, W.; Dai, Y.; Zhang, H.; Lan, Z.; Yu, D. WebVoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919* **2024**.
203. Vibhute, M.; Gutierrez, N.; Radivojevic, K.; Brenner, P. Multimodal Web Agents for Automated (Dark) Web Navigation.
204. Ho, X.; Nguyen, A.K.D.; Sugawara, S.; Aizawa, A. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060* **2020**.
205. Fernández-Pichel, M.; Pichel, J.C.; Losada, D.E. Evaluating Search Engines and Large Language Models for Answering Health Questions. *arXiv preprint arXiv:2407.12468* **2024**.
206. You.com. <https://you.com>, 2021.
207. Yang, A.; Li, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Gao, C.; Huang, C.; Lv, C.; et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388* **2025**.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.