

DeepRead: Document Structure-Aware Reasoning to Enhance Agentic Search

Zhanli Li*
lizhanli@stu.zuel.edu.cn
Key Lab of Intelligent Information
Processing of Chinese Academy of
Sciences (CAS), Institute of
Computing Technology, CAS
Beijing, China
State Key Lab of AI Safety
Beijing, China
University of Chinese Academy of
Sciences
Beijing, China

Huiwen Tian
tianhuiwen25@mails.ucas.ac.cn
Key Lab of Intelligent Information
Processing of Chinese Academy of
Sciences (CAS), Institute of
Computing Technology, CAS
Beijing, China
State Key Lab of AI Safety
Beijing, China
University of Chinese Academy of
Sciences
Beijing, China

Lvzhou Luo
luolvzhou23s@ict.ac.cn
Key Lab of Intelligent Information
Processing of Chinese Academy of
Sciences (CAS), Institute of
Computing Technology, CAS
Beijing, China
State Key Lab of AI Safety
Beijing, China
University of Chinese Academy of
Sciences
Beijing, China

Yixuan Cao[†]
caoyixuan@ict.ac.cn
Key Lab of Intelligent Information
Processing of Chinese Academy of
Sciences (CAS), Institute of
Computing Technology, CAS
Beijing, China
State Key Lab of AI Safety
Beijing, China
University of Chinese Academy of
Sciences
Beijing, China

Ping Luo
luop@ict.ac.cn
Key Lab of Intelligent Information
Processing of Chinese Academy of
Sciences (CAS), Institute of
Computing Technology, CAS
Beijing, China
State Key Lab of AI Safety
Beijing, China
University of Chinese Academy of
Sciences
Beijing, China

Abstract

With the rapid progress of tool-using and agentic large language models (LLMs), Retrieval-Augmented Generation (RAG) is evolving from one-shot, passive retrieval into multi-turn, decision-driven evidence acquisition. Despite strong results in open-domain settings, existing agentic search frameworks commonly treat long documents as flat collections of chunks, underutilizing document-native priors such as hierarchical organization and sequential discourse structure. We introduce DeepRead, a structure-aware, multi-turn document reasoning agent that explicitly operationalizes these priors for long-document question answering. DeepRead leverages LLM-based OCR model to convert PDFs into structured Markdown that preserves headings and paragraph boundaries. It then indexes documents at the paragraph level and assigns each paragraph a coordinate-style metadata key encoding its section identity and in-section order. Building on this representation, DeepRead equips the LLM with two complementary tools: a Retrieve tool that localizes relevant paragraphs while exposing their structural coordinates (with lightweight scanning context), and a ReadSection tool

that enables contiguous, order-preserving reading within a specified section and paragraph range. Our experiments demonstrate that DeepRead achieves significant improvements over Search-o1-style agentic search in document question answering. The synergistic effect between retrieval and reading tools is also validated. Our fine-grained behavioral analysis reveals a reading and reasoning paradigm resembling human-like “locate then read” behavior.

CCS Concepts

• **Information systems** → **Information retrieval**; Document representation; • **Computing methodologies** → *Natural language processing*; Information extraction.

Keywords

Document question answering, Agentic RAG, Information retrieval, Structured documents, OCR, Long-document reasoning

1 Introduction

LLMs have achieved impressive performance in natural language understanding and generation, yet they remain brittle on knowledge-intensive tasks that require precise, verifiable evidence. Two factors are particularly limiting: (i) static parametric memory cannot faithfully encode ever-changing or domain-specific details, and (ii)

*The first author contributed to this work during an internship at the Key Lab of Intelligent Information Processing of Chinese Academy of Sciences (CAS), Institute of Computing Technology, CAS. The paper is currently in preview form, and this work is still in progress.

[†]Corresponding author

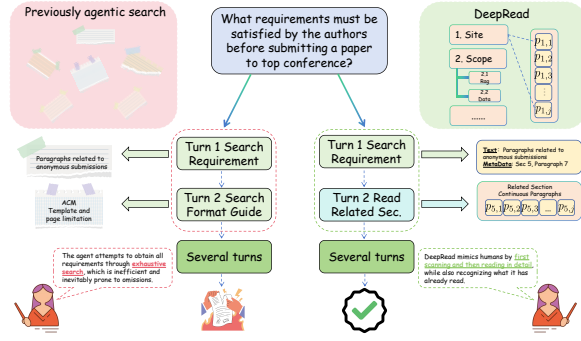


Figure 1: A Comparison of Search-o1-style Agentic Search and DeepRead on a Toy Case

LLMs tend to produce plausible but unsupported statements when evidence is missing. RAG mitigates these issues by grounding generation in external, non-parametric sources: at inference time, the model retrieves relevant evidence and conditions its answer on retrieved content, improving factuality and contextual alignment in high-stakes or specialized settings [12].

Early RAG systems predominantly adopted one-shot pipelines, where retrieval is executed once and the answer is generated from a fixed set of top-ranked chunks. Research in this phase focused on improving retrieval precision through stronger embedding models, optimized indexing, and coarse-to-fine reranking. While these advances improved single-step retrieval accuracy, the interaction pattern remained static: the system neither revises its information needs nor adapts its access strategy as reasoning unfolds. This limitation becomes pronounced for long-document and multi-hop scenarios, where evidence is widely distributed and cannot be reliably captured by a single retrieval call. Recent work such as Zhao et al. [30] mitigates the “lost-in-the-middle” phenomenon for long contexts, but it largely remains within fixed, single-round retrieval pipelines and lacks the interactivity required for complex reasoning.

To better handle multi-step dependencies, approaches such as PlanRAG [11] introduced explicit planning followed by retrieval. However, such two-stage designs can be brittle: they depend heavily on the quality of the initial plan and have limited ability to adapt when intermediate findings deviate from expectations. In parallel, iterative retrieval methods [6] have emerged to gather information across multiple turns. Nevertheless, many of these methods still follow *prescribed* schedules (e.g., a fixed number of rounds or a rigid retrieve-read-generate loop). Such rigidity can be inefficient for simple queries that require minimal evidence, and insufficient for complex queries that demand extensive, adaptive evidence acquisition [6].

More recently, *agentic RAG* has reshaped this landscape by formulating evidence acquisition as an autonomous decision-making process via tool use. In frameworks such as Search-o1 [14], the LLM is no longer bound to a fixed retrieve-generate schedule; it can decide *when to search, what to search for, and when sufficient evidence has been gathered to stop and answer*. This autonomy enables

substantially more flexible information seeking than classical iterative retrieval, as the model can dynamically adjust its trajectory based on intermediate reasoning and retrieved feedback.

Despite this progress, as illustrated in the toy example of Fig. 1, Search-o1-style agentic search remains *structurally blind* when operating over long and organized documents. Consider the query: “What requirements must be satisfied by the authors before submitting a paper to top conference?” A structure-agnostic agent is forced into a cycle of *keyword exhaustion and guesswork*. It attempts to piece together the answer by repeatedly issuing specific searches (e.g., for ‘format’, ‘anonymity’, ‘template’) and stitching together disjoint snippets. This strategy is inherently omission-prone: if the agent fails to guess a specific keyword (e.g., ‘page limit’), that requirement is simply missed. Furthermore, lacking a history of *examined regions*, the agent often wastes turns redundantly retrieving content it has effectively already covered. In reality, for such long-context information needs, evidence is rarely scattered randomly; it is typically organized systematically within a specific section (e.g., a dedicated “Submission Criteria” section). DeepRead capitalizes on this document-native topology. Unlike previous approaches, retrieving a *single* relevant snippet (e.g., just the anonymity policy) within such a section acts as a structural anchor. This detection incentivizes the LLM to abandon further keyword guessing and instead *read the entire section contiguously* via the ReadSection tool. This ensures that all co-located requirements—even those the agent did not explicitly search for—are captured comprehensively in a single pass, avoiding the retrieval noise and fragmentation common in preprocessing-heavy approaches [8].

In parallel, document understanding technologies have advanced substantially. Modern OCR and multimodal parsing models [3, 15, 22, 24, 26] can convert PDFs into structured representations (e.g., Markdown) that preserve layout and organizational cues, including heading hierarchies, nested lists, tables, and paragraph boundaries. For document QA, such structure is not merely cosmetic: it encodes strong priors about where information lives and how it should be read. Yet most agentic RAG systems still reduce documents to flat, orderless chunk collections, discarding these priors. As a result, even when an agent recognizes that evidence is incomplete, it has limited navigational affordances: it cannot reliably decide which section to inspect next, which span to read contiguously, or what has already been covered—leading to redundant retrieval, context fragmentation, and omission-prone reasoning. Some non-agentic approaches, such as Tao et al. [20] and Sarthi et al. [19], acknowledge this issue. However, instead of leveraging the document’s native structure, they employ large models to iteratively generate a hierarchical document structure. These methods also appear to disregard the document’s sequential information. Consequently, the preprocessing costs are extremely costly.

To address these limitations, we propose **DeepRead**, a structure-aware document reasoning agent that operationalizes document hierarchy and sequential priors for multi-turn QA. DeepRead builds upon the autonomous decision-making paradigm of agentic search but fixes a key bottleneck—the lack of document-native topology in the interaction interface. Specifically, DeepRead maps each document into a *structural coordinate system* (section and paragraph indices) and equips the LLM with two synergistic tools: Retrieve,

which performs scanning-aware localization and returns coordinate-anchored evidence, and ReadSection, which enables contiguous, order-preserving reading within a specified section and paragraph range. This interface supports a human-like “*locate-then-read*” workflow: first pinpointing relevant regions via lightweight scanning, then consuming complete local narratives. This design improves navigation efficiency over long, structured texts and mitigates the context fragmentation inherent in flat retrieval paradigms. Our contributions are summarized as follows:

- We propose a structure-grounded reasoning framework that operationalizes document hierarchy and sequence into a coordinate-based interaction system. By equipping the agent with dual tools—Retrieve for scanning-aware localization and ReadSection for range-preserving reading—we enable a human-like *locate-then-read* paradigm that reconstructs contiguous evidence from fragmented search results.
- We conduct comprehensive experiments across four benchmarks spanning single-document financial reports, multi-document academic papers, and semi-structured syllabi. DeepRead achieves consistent gains over strong baselines, including Search-o1-style agentic search, with particularly pronounced improvements on complex questions requiring long-range contextual integration.
- We provide a fine-grained behavioral analysis that validates the efficacy of our design. Empirical evidence confirms that DeepRead exhibits human-like reading patterns—balancing broad search with deep sequential reading—and ablation studies demonstrate the critical synergy between structural priors and agentic tools, especially in multi-document scenarios.

2 Related Work

2.1 Document QA

Document Question Answering (DocQA) is typically categorized into *open* and *closed* (document-grounded) settings. In *open* QA, systems query external sources like the Web or large corpora, facing challenges in *coverage* and *credibility* as they must filter conflicting information. In contrast, *closed* QA restricts evidence to a specific document or small collection, focusing on *precise localization* and *faithful interpretation* within the text, crucial for high-stakes domains like legal contracts, financial reports, and scientific papers. The primary challenge here is twofold: (i) identifying relevant regions amid noise and dispersion, and (ii) assembling sufficient *contiguous* context to resolve long-range dependencies. This necessitates a *locate-then-read* paradigm, where retrieval identifies candidate regions, followed by intensive reading for comprehensive understanding. Importantly, the closed setting requires leveraging *document-native structure*, such as hierarchical headings and sequential order, to guide reasoning. Treating documents as unstructured fragments would destroy these structural cues, leading to fragmented context and shallow reasoning. *DeepRead* addresses this challenge by operationalizing the document’s hierarchy and sequence into a coordinate-based navigation system, enabling agents to plan globally and read sequentially, mimicking human strategies for efficient long-document reasoning [2, 13, 17, 18].

2.2 Document Parsing

High-quality document parsing is essential for reliable DocQA, bridging raw visual inputs and structured text needed for reasoning. Existing solutions fall into two paradigms: pipeline-style systems (e.g., PaddleOCR-VL [3]) and end-to-end vision-language models (e.g., DeepSeek-OCR [26], HunyuanOCR [22]). Pipeline systems decompose parsing into modular stages—layout analysis followed by text recognition—ensuring recovery of reading order and structural hierarchy, which mitigates hallucinations and ordering errors in complex documents. On the other hand, end-to-end models map document images directly to structured markup, utilizing structure-aligned training data such as paired document images with Markdown or \LaTeX transcriptions. These models predict both characters and formatting tokens (e.g., headers, lists, tables), allowing them to learn a document’s logical organization and produce coherent, hierarchy-aware representations. These advancements make document-native priors accessible to LLM-based systems. Whether from a pipeline or end-to-end model, structured Markdown retains both the directory-tree structure and sequential logic, which DeepRead leverages by treating parsed documents as structured artifacts. By operationalizing recovered hierarchy and sequence into a coordinate interface, DeepRead enables agents to navigate long documents using the “locate-then-read” strategy employed by human readers.

2.3 Agentic Search

RAG enhances LLMs by grounding generation in external evidence, improving domain adaptability and mitigating hallucinations. Early *Naive RAG* approaches typically retrieve the top- k chunks in a single pass, but as multi-hop question answering gained attention, research shifted towards multi-step retrieval [11] and problem decomposition [1]. Complex queries over long documents often require *Iterative RAG*, where models gather information across multiple turns [6]. However, standard methods follow rigid retrieval schedules (e.g., fixed query rounds), which can lead to inefficiency and loss of focus during long-term reasoning. Some database-oriented studies [16, 25] address structural dependencies by conducting RAG over knowledge graphs or hierarchical trees, but these approaches abstract away from the document’s *native* visual layout, resulting in high construction overheads and limited scalability.

Agentic RAG offers a more flexible solution by framing question answering as an autonomous decision-making process. Frameworks like ReAct and the recent reasoning-based Search-o1 [14] enable LLMs to decide when to search and how to synthesize evidence, improving performance on multi-step tasks. However, a major limitation in these systems is the data representation, as most agentic approaches treat documents as *flat, unstructured collections of chunks*. This “structural blindness” results in context fragmentation, where retrieved snippets lose their logical position. **DeepRead** addresses this by introducing a *structure-grounded* interaction model, directly operationalizing the document’s hierarchy and sequence into a lightweight coordinate system. This enables explicit navigation (e.g., “read Section 3.2”) and contiguous reading, preserving long-context reasoning without the overhead of maintaining external knowledge graphs.

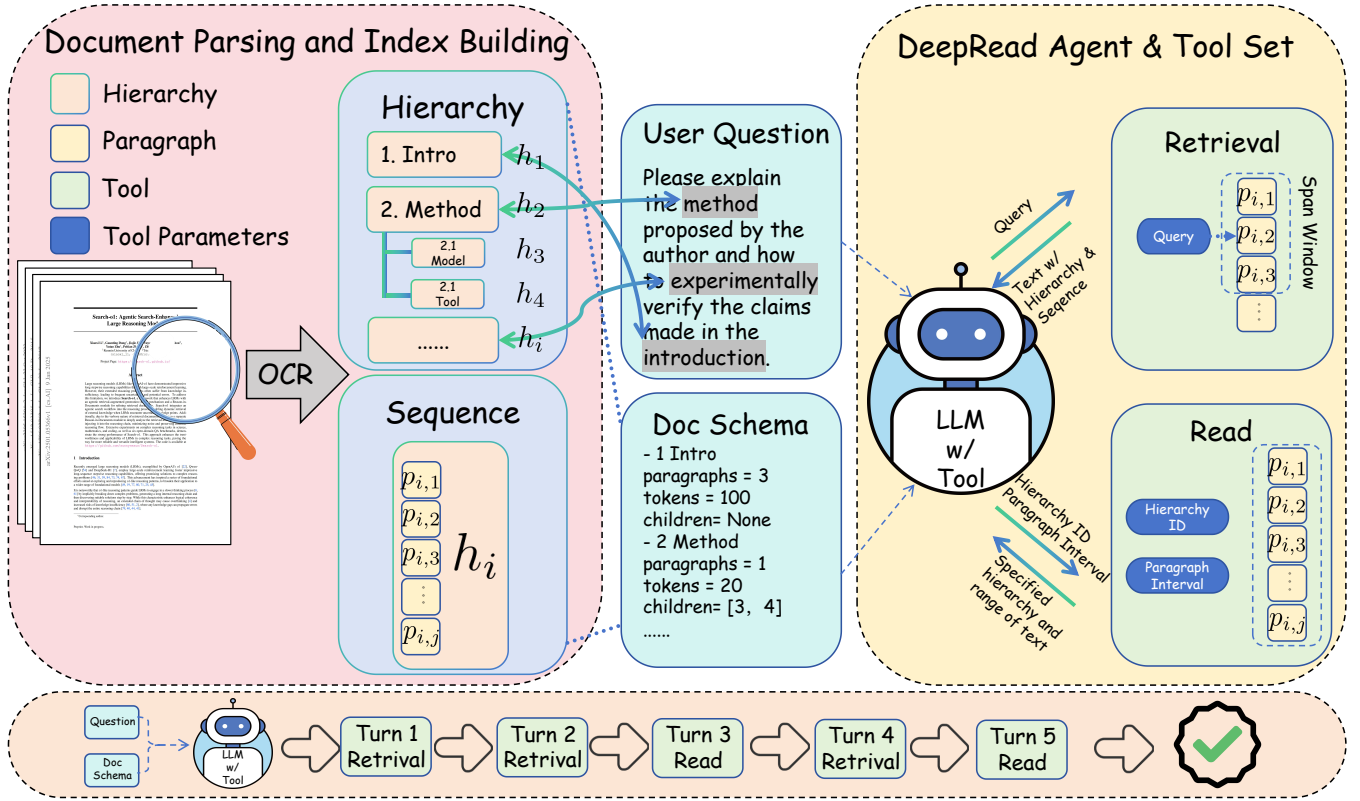


Figure 2: This is the DeepRead framework diagram. It takes user questions parsed into Doc Schema as input. The LLM employs retrieval and reading tools to reason and think over the documents.

3 Methodology

3.1 Preliminaries: Agentic Search

We adopt the vanilla ReAct [27] as the agent’s framework, which synergizes reasoning and acting. Let the user question be q . The interaction proceeds for at most T rounds. At round t , the agent state s_t consists of the message history, including the system prompt, user query, and the trajectory of tool interactions. The model policy π_θ generates an action a_t :

$$a_t \sim \pi_\theta(\cdot | s_t), \quad a_t \in \{\text{FINAL}\} \cup \mathcal{A}, \quad (1)$$

where an action is either a final answer or a tool invocation $a_t = (\tau_t, \mathbf{x}_t)$ with tool name τ_t and arguments \mathbf{x}_t . Executing a_t yields an observation o_t , updating the state via $s_{t+1} \leftarrow s_t \oplus (a_t, o_t)$, where \oplus appends the interaction to the history. In DeepRead, the action set \mathcal{A} consists of two tools, Retrieve and ReadSection (Sec. 3.5), which together support a human-like *locate-then-read* workflow over structured documents.

3.2 Structured Document Modeling

We assume raw documents are processed by an OCR engine into a structured Markdown format. To support precise navigation, we formalize document structure along two dimensions: hierarchy and sequence.

Hierarchical Sections. We define each document as a sequence of hierarchical sections (headings), ordered by their appearance in the text. For a document d , let $N_h^{(d)}$ be the total number of sections. We denote the i -th section in the document-order as $h_i^{(d)}$, where $i \in \{1, \dots, N_h^{(d)}\}$. Each section is characterized by:

$$h_i^{(d)} = (\text{title}_i^{(d)}, \ell_i^{(d)}, \mathcal{C}_i^{(d)}), \quad (2)$$

where i serves as a unique *section ID* within document d (assigned by heading appearance order), $\text{title}_i^{(d)}$ is the heading text, $\ell_i^{(d)}$ is the indentation level (e.g., $\ell = 1$ for #, $\ell = 2$ for ##), and $\mathcal{C}_i^{(d)} = \{k \mid \text{parent}(h_k^{(d)}) = h_i^{(d)}\}$ is the set of IDs of immediate children within the same document. This captures the nested structure: a section may contain its own content as well as a set of subsections.

Atomic Paragraphs. The finest granularity of analysis is the paragraph. We posit that every paragraph belongs to exactly one section that *directly contains it* (the section may be a leaf or a non-leaf node). We denote a paragraph as $p_{i,j}^{(d)}$ with two coordinates:

- (1) **Section Membership (i):** the section ID of $h_i^{(d)}$ that directly contains this paragraph.
- (2) **Local Position (j):** the sequential index of the paragraph within section $h_i^{(d)}$, starting from 1.

Formally, the content of section $h_i^{(d)}$ *excluding its subsections* is the ordered sequence

$$\mathcal{P}(h_i^{(d)}) = [p_{i,1}^{(d)}, p_{i,2}^{(d)}, \dots, p_{i,n_i^{(d)}}^{(d)}], \quad (3)$$

where $n_i^{(d)}$ is the number of paragraphs directly under section i in document d . This formalization maps every text span to a coordinate (d, i, j) , explicitly encoding both *where it belongs structurally* (document d , section i) and *where it lies sequentially* within that section (index j).

3.3 Paragraph-Level Indexing and Metadata

Building upon Sec. 3.2, DeepRead adopts the atomic paragraph $p_{i,j}^{(d)}$ as the fundamental unit of indexing, rather than merging text into arbitrary sliding windows. We define an indexable unit as a single paragraph paired with its structural coordinates:

$$c_{i,j}^{(d)} = (p_{i,j}^{(d)}, \Gamma_{d,i,j}). \quad (4)$$

Each paragraph carries a structured metadata object that exposes the coordinate system:

$$\Gamma_{d,i,j} = \{\text{doc_id} : d, \text{sec_id} : i, \text{para_idx} : j\}. \quad (5)$$

During tool interaction, DeepRead returns text along with $\Gamma_{d,i,j}$ for every paragraph it exposes. This ensures that evidence is not only readable but also *addressable*: the agent can reason about *what* a paragraph says and *where* it resides (e.g., “doc d , Section i , paragraph j ”), enabling subsequent coordinate-based navigation and reading actions.

3.4 System Prompt: Hierarchical Skeleton

To enable global planning without overwhelming the context window, DeepRead injects a lightweight skeletal representation of the document collection into the system prompt. Instead of feeding the full content, we construct a structured Table of Contents (TOC) that acts as a navigational map.

For each document d and each section $h_i^{(d)}$ in that document, we generate a descriptor:

$$\text{Entry}(d, i) = \langle d, i, \text{title}_i^{(d)}, \ell_i^{(d)}, \mathcal{C}_i^{(d)}, n_i^{(d)}, m_i^{(d)} \rangle, \quad (6)$$

where d is the document ID, i is the section ID *unique within* document d , $\text{title}_i^{(d)}$ is the section header, $\ell_i^{(d)}$ is the heading level, $\mathcal{C}_i^{(d)}$ is the list of children section IDs (within the same document), $n_i^{(d)}$ is the number of paragraphs directly under section i (excluding subsections), and $m_i^{(d)}$ is the token count of the section’s direct content. For a collection of documents \mathcal{D} , the system prompt organizes TOC entries *by document*:

$$\text{TOC}(\mathcal{D}) = [\text{TOC}(d)]_{d \in \mathcal{D}}, \quad \text{TOC}(d) = [\text{Entry}(d, i)]_{i=1}^{N_h^{(d)}}. \quad (7)$$

This organization ensures unambiguous addressing: section IDs are not required to be globally unique across documents, because any section reference is disambiguated by the pair (d, i) .

This design provides the agent with structural priors for planning: by observing $(\ell_i^{(d)}, \mathcal{C}_i^{(d)})$, the model can infer the hierarchy and scope within each file; by observing $(n_i^{(d)}, m_i^{(d)})$, the model can estimate reading cost and decide whether to read an entire section

Algorithm 1 DeepRead: Structure-Preserving Agentic Reading

Require: Documents \mathcal{D} with parsed sections $\{h_i^{(d)}\}$ and paragraphs $\{p_{i,j}^{(d)}\}$; Question q ; Window $W = (w^\uparrow, w^\downarrow)$.

- 1: **Initialize:** Construct System Prompt with $\text{TOC}(\mathcal{D})$ (grouped by document; includes $n_i^{(d)}$ and $m_i^{(d)}$).
- 2: $s_1 \leftarrow [\text{System Prompt}, \text{User} : q]$
- 3: **for** $t = 1$ **to** T **do**
- 4: $a_t \sim \pi_\theta(\cdot | s_t)$ ▷ Model predicts an action
- 5: **if** $a_t = \text{FINAL}$ **then return** Answer
- 6: **end if**
- 7: **if** $a_t.\tau = \text{Retrieve}$ **then**
- 8: $u \leftarrow a_t.x.\text{query}$ ▷ Tool argument query is a query string
- 9: $H \leftarrow \text{Rank}(u)$ ▷ $H = \{(d_r, i_r, j_r)\}_{r=1}^K$ are top- K paragraph coordinates
- 10: $\mathcal{U} \leftarrow \emptyset$ ▷ Set of paragraph coordinates to return (for deduplication)
- 11: **for** each hit $(d_r, i_r, j_r) \in H$ **do**
- 12: $j_r^\uparrow \leftarrow \max(1, j_r - w^\uparrow)$
- 13: $j_r^\downarrow \leftarrow \min(n_{i_r}^{(d_r)}, j_r + w^\downarrow)$
- 14: **for** $j = j_r^\uparrow$ **to** j_r^\downarrow **do**
- 15: $\mathcal{U} \leftarrow \mathcal{U} \cup \{(d_r, i_r, j)\}$ ▷ Deduplicate overlaps within a single retrieval call
- 16: **end for**
- 17: **end for**
- 18: $\mathcal{U}_{\text{sorted}} \leftarrow \text{Sort}(\mathcal{U})$ ▷ Sort by d , then i , then j
- 19: $o_t \leftarrow \text{Format}(\{(p_{i,j}^{(d)}, \Gamma_{d,i,j}) : (d, i, j) \in \mathcal{U}_{\text{sorted}}\})$ ▷ Format concatenates paragraphs in this sorted order and preserves per-paragraph metadata
- 20: **else if** $a_t.\tau = \text{ReadSection}$ **then**
- 21: $d \leftarrow a_t.x.\text{doc_id}; i \leftarrow a_t.x.\text{sec_id}$
- 22: $j_s \leftarrow a_t.x.\text{start}; j_e \leftarrow a_t.x.\text{end}$
- 23: $j_s \leftarrow \max(1, j_s); j_e \leftarrow \min(n_i^{(d)}, j_e)$ ▷ Clip to valid range using $n_i^{(d)}$ from the TOC
- 24: $o_t \leftarrow \text{Format}(\{(p_{i,j}^{(d)}, \Gamma_{d,i,j}) : j \in [j_s, j_e]\})$ ▷ Format concatenates paragraphs in increasing j and preserves metadata
- 25: **end if**
- 26: $s_{t+1} \leftarrow s_t \oplus (a_t, o_t)$
- 27: **end for**

or only a targeted span. A concrete example is shown by the Doc Schema in Figure 2.

3.5 Tools: Coordinate-Based Interaction

The agent interacts with the document collection via two complementary tools operating on the (d, i, j) coordinate system defined in Sec. 3.2. The interaction mimics human behavior: *fast localization* via retrieval followed by *intensive study* via reading.

1. Retrieve. This tool functions as a locator that identifies relevant coordinates within the document structure. It accepts a query u passed as the query field in the tool arguments. Given u , it first performs semantic search to identify the top- K most relevant atomic paragraphs:

$$\text{Rank}(u) \rightarrow \{(d_r, i_r, j_r)\}_{r=1}^K, \quad (8)$$

where (d_r, i_r, j_r) denotes the r -th hit (rank index r), and j_r is the local paragraph index within section i_r of document d_r . Here, section IDs i are unique *within* each document, and are disambiguated across files by the document ID d .

To simulate human skimming (inspecting nearby context of a match), we introduce a scanning window hyperparameter $W = (w^\uparrow, w^\downarrow)$ representing upward and downward look-ahead budgets. For a hit at (d_r, i_r, j_r) , the window boundaries are:

$$j_r^\uparrow = \max(1, j_r - w^\uparrow), \quad j_r^\downarrow = \min(n_{i_r}^{(d_r)}, j_r + w^\downarrow), \quad (9)$$

yielding a local slice of paragraphs:

$$\text{Scan}(d_r, i_r, j_r; W) = \{(d_r, i_r, j) \mid j \in [j_r^\uparrow, j_r^\downarrow]\}. \quad (10)$$

DeepRead expands *each* of the K hits independently and then performs deduplication to avoid returning redundant paragraphs when slices overlap:

$$\mathcal{U}(u) = \bigcup_{r=1}^K \text{Scan}(d_r, i_r, j_r; W), \quad (11)$$

where $\mathcal{U}(u)$ is a set of unique paragraph coordinates collected within a single Retrieve call. To preserve readability and discourse coherence, we sort the deduplicated coordinates before serialization:

$$\mathcal{U}_{\text{sorted}}(u) = \text{Sort}(\mathcal{U}(u)) \quad (\text{by } d, \text{ then } i, \text{ then } j). \quad (12)$$

Finally, Retrieve returns the corresponding paragraphs, *each tagged with its own metadata*:

$$\text{Retrieve}(u) \rightarrow \text{Format}\left(\{(p_{i,j}^{(d)}, \Gamma_{d,i,j}) : (d, i, j) \in \mathcal{U}_{\text{sorted}}(u)\}\right). \quad (13)$$

In practice, Format serializes the returned paragraphs into a model-readable text block by concatenating them in this sorted (natural) order while preserving per-paragraph metadata (e.g., prefixing each paragraph with its $\Gamma_{d,i,j}$). This output allows the agent to read local context and obtain precise structural addresses for subsequent navigation.

2. ReadSection. This tool performs deep, order-preserving reading over a targeted region. It accepts a document ID d , a section ID i (unique within d), and a paragraph range $[j_{\text{start}}, j_{\text{end}}]$:

$$\begin{aligned} \text{ReadSection}(d, i, j_{\text{start}}, j_{\text{end}}) &\rightarrow \\ \text{Format}\left(\{(p_{i,j}^{(d)}, \Gamma_{d,i,j}) : j \in [j_{\text{start}}, j_{\text{end}}]\}\right), \end{aligned} \quad (14)$$

where the system clips the range to valid boundaries using $n_i^{(d)}$ provided in the TOC. This yields contiguous, in-order paragraphs from the specified section, reducing context fragmentation compared with flat chunk concatenation.

Synergy. The two tools form a closed loop. Retrieve provides (i) a lightweight preview and (ii) coordinate anchors (d, i, j) for relevant evidence. When the agent detects that a section contains complex logic requiring full local context, it invokes ReadSection with the returned coordinates and the TOC-provided $(n_i^{(d)}, m_i^{(d)})$ for cost-aware planning to ingest an order-preserving span. Collectively, these tools enable a human-like *locate-then-read* paradigm for faithful long-document reasoning. This framework, which supports dynamic structural navigation, is detailed in Algorithm 1.

4 Experiment

4.1 Benchmark Details

We evaluated DeepRead on four benchmarks designed to test specific RAG capabilities:

(1) FinanceBench [9]: We use the open-source version (150 pairs) to evaluate long-document reasoning within the financial domain. **(2) ContextBench (Ours):** To evaluate long-range contextual integration, This is a single-document QA dataset annotated by 12 AI experts using PDFs from their daily work and personal lives—including novels they enjoy reading, academic papers, scripts, textbooks, and more. The requirement is that the evidence for each question spans long distances within the document or consists of a single coherent but exceptionally lengthy evidence passage. Each QA is estimated to take 0.5 person-hours. **(3) QASPER [4] (Multi-Doc):** To test academic and cross-document reasoning, we synthesized a multi-document version of QASPER. We used an LLM to generate questions spanning 2–5 papers and manually filtered out illogical or erroneous samples, resulting in 143 high-quality pairs. **(4) SyllabusQA [7] (Multi-Doc):** To test documents with simple hierarchical structures, we obtained all course syllabus PDFs from SyllabusQA for constructing single-document QA, we applied the same synthesis and manual verification process as QASPER, yielding 196 high-quality pairs. Table 5 presents the relevant statistical data for the above benchmarks.

Document structure parsing was conducted using PaddleOCR-VL, served via VLLM [10] on an NVIDIA RTX 4090 GPU. Given that DeepRead operates on semantic units (paragraphs/sections) whereas baselines rely on fixed-size chunking, direct comparison of intermediate retrieval metrics is inherently inequitable. Consequently, we focus on end-to-end performance using an LLM-as-a-Judge framework. We employ DeepSeek V3.2 [5] as the evaluator, with the temperature set to 0.0 to ensure deterministic and reproducible scoring. Regarding the construction of synthetic multi-document benchmarks, we utilized GLM-4.7 [21] taking the full document context as input, with a temperature of 0.7 to encourage diversity in question generation.

4.2 Baseline and DeepRead Settings

To evaluate the effectiveness of DeepRead, we consider four baseline families: single-pass retrieval, **RAPTOR** [19], Iterative Retrieval-Generation Synergy (ITRG) [6], and Search-o1 [14]. For all baselines, we use Qwen3-embedding-8b [29] as the dense retriever. For traditional retrieval, we explicitly distinguish two settings: (i) direct dense retrieval, where results are ranked solely by embedding similarity, and (ii) two-stage retrieval (coarse-to-fine), where a reranker refines an initial candidate pool. Since two-stage retrieval is the de facto industrial practice, we apply a reranker by default in iterative retrieval and agentic search to ensure these stronger baselines reflect realistic deployment settings, while we keep a direct dense-only variant for single-pass retrieval as a clean reference point. Concretely, for single-pass retrieval, we follow the OpenAI File Search configuration with chunk size 800 and overlap 400, and return the top 10 chunks in one round (as in Search-o1 paper). When reranking is enabled, the first-stage retriever produces 30 candidates, which are then scored by Qwen3-reranker-8b [29] and truncated to the target return count. For RAPTOR, we utilize the recommended ‘‘Collapsed Tree’’ setting [19] with a maximum token limit of 800 per node, 5 layers, and a clustering top- k of 5, retrieving the top-10 nodes from the collapsed index. For ITRG,

Table 1: Comparison with Different Methods (Accuracy %). Bold indicates the optimal choice, underlined indicates the next best choice. Green text denotes the absolute improvement over the corresponding Search-o1 baseline.

Method	Single-Document			Multi-Document			Overall Avg
	FinanceBench [9]	ContextBench (Ours)	Avg	QASPER [4]	SyllabusQA [7]	Avg	
Dense RAG [12, 28]	38.1	60.0	49.0	15.4	26.3	20.8	34.9
Dense RAG w/ Reranker [28]	47.3	66.2	56.8	15.4	24.8	20.1	38.4
ITRG (refresh) [6]	48.1	67.9	58.0	9.8	31.0	20.4	39.2
ITRG (refine) [6]	52.0	70.1	61.0	15.4	31.3	23.3	42.2
RAPTOR [19]	38.7	52.6	45.7	20.2	40.3	30.3	38.0
Search-o1 [14]	80.0	74.5	77.3	65.0	57.1	61.1	69.2
DeepRead (Ours)	82.7 (+2.7)	91.5 (+17.0)	87.1 (+9.8)	72.7 (+7.7)	70.9 (+13.8)	71.8 (+10.7)	79.5 (+10.3)
Search-o1 w/ expand [14]	<u>83.3</u>	84.0	83.7	65.0	68.4	66.7	75.2
DeepRead w/ expand (Ours)	84.0 (+0.7)	<u>88.3 (+4.3)</u>	<u>86.2 (+2.5)</u>	76.2 (+11.2)	72.5 (+4.1)	74.4 (+7.7)	80.3 (+5.1)

we adopt the more effective 4-round setting reported in the original paper, returning top 6 chunks per round; under the dense-only setting, we directly return the dense ranking, whereas under the two-stage setting, each round reranks candidates before selecting the final 6 chunks. For Search-o1, to ensure fair comparison, we do not include the Reason-in-Documents and we additionally inject the document structural schema into the system prompt (matching DeepRead’s access to structure)¹. Search-o1 uses structure-based chunking with overlap 0; each retrieval tool call returns 2 chunks, and the expansion window is configured as (1, 1). The policy model of all baseline and DeepRead used in this paper is DeepSeek v3.2 for ReAct paradigm inference, with a decoding temperature of 0. For search-o1 and DeepRead, we set the maximum round tree to 50.

4.3 Main Result

Table 1 reports end-to-end accuracy across four benchmarks. Overall, **DeepRead consistently outperforms strong baselines**, with the largest gains on questions that require long-range, scope-aware evidence integration. Under the primary judge (DeepSeek-V3.2), DeepRead achieves an **overall average of 79.5%** (80.3% with expansion), improving over Search-o1 by **+10.3** points (and over Search-o1 w/ expand by **+5.1** points). These results indicate that explicitly exposing document-native topology to the agent yields substantial benefits beyond stronger retrieval or additional search rounds.

Impact of structure-grounded reading (ReadSection). Comparing Search-o1 and DeepRead isolates the contribution of coordinate-based, order-preserving reading. DeepRead improves upon Search-o1 on **all** benchmarks, with especially large gains on **ContextBench** (+17.0 points), supporting our central claim: treating long documents as flat, orderless chunks leads to context fragmentation, whereas ReadSection reconstructs *contiguous* evidence anchored to explicit structural coordinates. Notably, DeepRead also yields strong improvements in multi-document settings, achieving +7.7 on QASPER and +13.8 on SyllabusQA, suggesting that hierarchy- and sequence-aware navigation remains effective even when evidence spans multiple files and sections.

¹Reason-in-Documents refers to summarizing both the preceding context (inference and tool invocation/return) and the currently retrieved chunk. To our knowledge, this approach yields only marginal improvements, while the cost of repeatedly summarizing historical context is prohibitively high. Consequently, it is rarely adopted in industrial practice.

Effect of local expansion (expand). Expansion generally benefits search-heavy baselines, but its effect is not uniformly positive for structure-grounded reading. Search-o1 improves substantially with expansion (69.2% → 75.2%), consistent with the intuition that adjacent context can partially compensate for fragmented retrieval. In contrast, DeepRead exhibits a smaller overall gain (79.5% → 80.3%), and on **ContextBench** expansion *reduces* accuracy (91.5% → 88.3%). This pattern suggests that ReadSection already provides the coherent, scope-complete context needed for reasoning; indiscriminate window expansion may introduce irrelevant neighboring paragraphs and dilute signal. Therefore, we view expansion as a complementary mechanism that mainly mitigates retrieval fragmentation, rather than a necessary component of structure-grounded reading.

Robustness of judge choice. To reduce dependence on a single evaluator, we replicate all experiments using two additional independent LLM judges, GLM-4.7 and Qwen3-235B-A22B-thinking-2507 [23] (See Appendix A.3). Across judges, the **relative ranking is stable**: DeepRead remains consistently stronger than Search-o1 variants, indicating that the gains are not an artifact of a particular judge’s calibration. Moreover, the inter-judge agreement reported in Table 11 is high overall, further supporting the reliability of the observed improvements.

4.4 Fine-Grained Behavior Analysis

The preceding quantitative results demonstrate the performance superiority of DeepRead; however, aggregate scores alone do not reveal the underlying mechanisms driving these gains. To bridge this gap, we conduct a fine-grained behavioral analysis characterizing how DeepRead diverges from standard Search-o1-style agentic workflows in terms of planning, tool consumption, and information processing efficiency. Figure 3 visualizes the statistical distribution of agent behaviors across four experimental settings. The analysis of the **First Tool Call** (top-left) reveals two critical insights regarding the agent’s planning capabilities: **1) High Protocol Adherence:** Across all settings, the probability of initiating with a tool call (specifically search) approaches 100%, confirming that the LLM strictly follows the system prompt to seek external evidence rather than hallucinating directly. **2) The “Locate-then-Read” Paradigm:** Crucially, DeepRead maintains a high initial search rate comparable to the Search-o1 baseline. This indicates

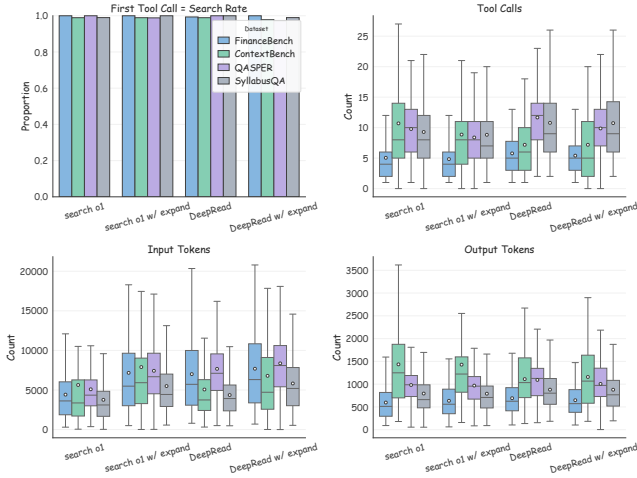


Figure 3: Fine-grained behavioral comparison between DeepRead and Search-o1 baselines. The panels illustrate the distribution of (a) the probability that the first action is a search, (b) the total number of tool calls per query, (c) input token consumption, and (d) output token generation across four benchmarks.

Table 2: The proportion of queries that are first retrieved and subsequently read, and the ratio of retrieval calls to read calls.

Benchmark	DeepRead		DeepRead w/ expand	
	S→R (%)	Ret/Read Ratio	S→R (%)	Ret/Read Ratio
ContextBench	95.70%	0.87	93.55%	0.88
FinanceBench	87.33%	1.82	82.67%	2.18
QASPER	98.25%	1.59	92.40%	2.00
SyllabusQA	97.96%	1.61	96.43%	1.86

that DeepRead does not blindly consume tokens; instead, it adopts a strategic *Locate-then-Read* pattern—using retrieval to pinpoint relevant document sections before deploying the heavy-weight Read tool for comprehensive information extraction, and table 2 verify this.

Further examination of resource consumption metrics reveals that while DeepRead exhibits a modest increment in tool invocations and input token usage, this computational overhead is negligible when contrasted with the prohibitive costs associated with complex knowledge graph construction or multi-stage iterative summary refinement strategies. Notably, we observe no significant inflation in output token generation, suggesting that the model extracts answers efficiently without redundant generation. Overall, compared to methods that build retrieval databases based on knowledge graphs and multi-round iterative summarization, DeepRead incurs only a minor increase in computational cost.

We also investigate the behavioral divergence between successful and failed queries. Our analysis reveals that incorrect samples frequently exhibit pathological search patterns characterized by

Table 3: Cost comparison between correct and incorrect responses. We report the average number of tool calls and total token consumption (input + output) across all benchmarks. Incorrect responses typically incur higher computational costs due to prolonged search-reasoning loops.

Method	Tool Calls		Total Tokens	
	Correct	Wrong	Correct	Wrong
Search-o1	7.7	10.4	4,978	6,522
Search-o1 w/ expand	7.1	9.5	7,349	8,700
DeepRead	8.7	11.0	6,770	7,648
DeepRead w/ expand	8.1	10.4	7,597	9,609

Table 4: Performance and Cost Comparison between DeepRead and Readonly Baseline.

Benchmark	DeepRead			Readonly		
	Acc (%)	Tools	Tokens	Acc (%)	Tools	Tokens
FinanceBench	82.7	5.8	7,698	80.7	4.7	7,980
ContextBench	91.5	7.2	6,172	91.5	6.9	6,680
QASPER	68.4	11.7	8,768	17.5	9.9	9,887
SyllabusQA	70.9	10.8	5,244	15.3	7.4	4,189

prolonged tool usage, which results in a prohibitive escalation of resource consumption, as evidenced in Table 3.

4.5 Ablation Study

Given that Table 1 has already demonstrated the results of whether to expand and only retrieve (search-o1), here we are primarily concerned with the synergistic effect of retrieval and read operations. Here we conducted an ablation experiment to specifically evaluate the role and effectiveness of the Read tool within DeepRead. We found that in single-document scenarios, allowing the LLM to perform Read without retrieval is competitive in terms of efficiency and cost. However, in multi-document scenarios, it exhibits significant disadvantages in both efficiency and cost, as detailed in Table 4. This validates the synergistic effect between retrieval and reading tools, particularly in multi-document contexts. Furthermore, we investigated whether DeepRead remains effective under varying numbers of returned chunks (k). As shown in Figure 4, while improved retrieval recall generally correlates with higher accuracy, DeepRead consistently surpasses the Search-o1 baseline across all tested benchmarks and k values. Crucially, our method demonstrates superior data efficiency; even with a minimal retrieval budget ($k = 2$), DeepRead effectively leverages document structure to recover necessary context, avoiding the performance degradation observed in the baseline under constrained settings.

5 Case Study

We provide case study to help understand (See Appendix A.2).

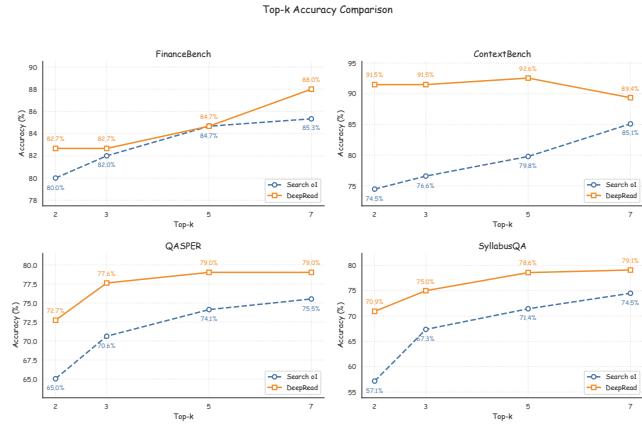


Figure 4: Impact of Retrieved Chunk Count (k) on Performance. We compare DeepRead against Search-o1 across four benchmarks with $k \in \{2, 3, 5, 7\}$. DeepRead exhibits consistent robustness, outperforming the baseline particularly in low-resource settings ($k = 2$), validating the efficacy of structure-aware reading over flat retrieval expansion.

6 Conclusion

This work presents DeepRead, a structure-aware agentic search framework for long-document question answering that resolves the structural blindness of mainstream agentic RAG systems by encoding document hierarchical and sequential priors into a coordinate-based navigation paradigm. By coupling lightweight retrieval localization with contiguous section-wise reading to form a human-like locate-then-read workflow, DeepRead effectively alleviates context fragmentation and redundant retrieval in long-document reasoning. Extensive experiments across four benchmarks validate its consistent, significant performance gains over state-of-the-art agentic baselines, especially in multi-document and long-range reasoning tasks, while fine-grained behavioral analyses confirm its efficient, human-aligned reading and decision-making patterns. We envision DeepRead becoming the mainstream agentic search paradigm for closed-domain document QA.

References

- [1] Paul JL Ammann, Jonas Golde, and Alan Akbik. 2025. Question Decomposition for Retrieval-Augmented Generation. *arXiv preprint arXiv:2507.00355* (2025).
- [2] Eunsol Choi, Daniel Hewlett, Alexandre Lacoste, Illia Polosukhin, Jakob Uszkoreit, and Jonathan Berant. 2016. Hierarchical question answering for long documents. *arXiv preprint arXiv:1611.01839* (2016).
- [3] Cheng Cui, Ting Sun, Suyin Liang, Tingquan Gao, Zelun Zhang, Jiaxuan Liu, Xueqing Wang, Changda Zhou, Hongen Liu, Manhui Lin, et al. 2025. Paddleocr-vl: Boosting multilingual document parsing via a 0.9 b ultra-compact vision-language model. *arXiv preprint arXiv:2510.14528* (2025).
- [4] Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. *arXiv preprint arXiv:2105.03011* (2021).
- [5] DeepSeek-AI. 2025. DeepSeek-V3.2: Pushing the Frontier of Open Large Language Models.
- [6] Zhangyin Feng, Xiaocheng Feng, Dezhi Zhao, Maojin Yang, and Bing Qin. 2024. Retrieval-generation synergy augmented large language models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 11661–11665.
- [7] Nigel Fernandez, Alexander Scarlatos, and Andrew Lan. 2024. SyllabusQA: A course logistics question answering dataset. *arXiv preprint arXiv:2403.14666* (2024).
- [8] Ziyu Gong, Yihua Huang, and Chengcheng Mai. 2025. Mmrag-docqa: A multi-modal retrieval-augmented generation method for document question-answering with hierarchical index and multi-granularity retrieval. *arXiv e-prints* (2025), arXiv–2508.
- [9] Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie Vidgen. 2023. Financebench: A new benchmark for financial question answering. *arXiv preprint arXiv:2311.11944* (2023).
- [10] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- [11] Myeonghwa Lee, Seonho An, and Min-Soo Kim. 2024. Planrag: A plan-then-retrieval augmented generation for generative large language models as decision makers. *arXiv preprint arXiv:2406.12430* (2024).
- [12] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.
- [13] Minghan Li, Miyang Luo, Tianrui Lv, Yishuai Zhang, Siqi Zhao, Ercong Nie, and Guodong Zhou. 2025. A Survey of Long-Document Retrieval in the PLM and LLM Era. *arXiv preprint arXiv:2509.07759* (2025).
- [14] Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366* (2025).
- [15] Yumeng Li, Guang Yang, Hao Liu, Bowen Wang, and Colin Zhang. 2025. dots. ocr: Multilingual document layout parsing in a single vision-language model. *arXiv preprint arXiv:2512.02498* (2025).
- [16] Junli Liang, Pengfei Zhou, Wangqiu Zhou, Wenjie Qing, Qi Zhao, Ziwen Wang, Qi Song, and Xiangyang Li. 2026. SentGraph: Hierarchical Sentence Graph for Multi-hop Retrieval-Augmented Question Answering. *arXiv preprint arXiv:2601.03014* (2026).
- [17] Tavish McDonald, Brian Tsan, Amar Saini, Juanita Ordonez, Luis Gutierrez, Phan Nguyen, Blake Mason, and Brenda Ng. 2022. Detect, retrieve, comprehend: a flexible framework for zero-shot document-level question answering. *arXiv preprint arXiv:2210.01959* (2022).
- [18] Jon Saad-Falcon, Joe Barrow, Alexa Siu, Ani Nenkova, Seunghyun Yoon, Ryan A Rossi, and Franck Dernoncourt. 2024. Pdftriage: Question answering over long, structured documents. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*. 153–169.
- [19] Parth Sarthi, Salman Abdullah, Aditi Dast, Shubh Khanna, Anna Goldie, and Christopher D Manning. 2024. Raptor: Recursive abstractive processing for tree-organized retrieval. In *The Twelfth International Conference on Learning Representations*.
- [20] Wenyu Tao, Xiaofen Xing, Yirong Chen, Linyi Huang, and Xiangmin Xu. 2025. Treerag: Unleashing the power of hierarchical storage for enhanced knowledge retrieval in long documents. In *Findings of the Association for Computational Linguistics: ACL 2025*. 356–371.
- [21] GLM Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang, Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei, Yean Cheng, Yifan An, Yilin Niu, Yuanhao Wen, Yushi Bai, Zhengxiao Du, Zihan Wang, Zilin Zhu, Bohan Zhang, Bosi Wen, Bowen Wu, Bowen Xu, Can Huang, Casey Zhao, Changpeng Cai, Chao Yu, Chen Li, Chendi Ge, Chenghua Huang, Chenhui Zhang, Chenxi Xu, Chenzheng Zhu, Chuang Li, Congfeng Yin, Daoyan Lin, Dayong Yang, Dazhi Jiang, Ding Ai, Erle Zhu, Fei Wang, Gengzheng Pan, Guo Wang, Hailong Sun, Haitao Li, Haiyang Li, Haiyi Hu, Hanyu Zhang, Hao Peng, Hao Tai, Haoke Zhang, Haoran Wang, Haoyu Yang, He Liu, He Zhao, Hongwei Liu, Hongxi Yan, Huan Liu, Huihong Chen, Ji Li, Jiajing Zhao, Jiamin Ren, Jian Jiao, Jiani Zhao, Jianyang Yan, Jiaqi Wang, Jiayi Gui, Jiayue Zhao, Jie Liu, Jijie Li, Jing Li, Jing Lu, Jingsen Wang, Jingwei Yuan, Jingxuan Li, Jingzhao Du, Jinhua Du, Jinxin Liu, Junkai Zhi, Junli Gao, Ke Wang, Lekang Yang, Liang Xu, Lin Fan, Lindong Wu, Lintao Ding, Lu Wang, Man Zhang, Minghao Li, Minghan Xu, Mingming Zhao, Mingshu Zhai, Pengfan Du, Qian Dong, Shangde Lei, Shangqing Tu, Shantong Yang, Shaoyou Lu, Shijie Li, Shuang Li, Shuang-Li, Shuxun Yang, Sibo Yi, Tianshu Yu, Wei Tian, Weihao Wang, Wenbo Yu, Weng Lam Tam, Wenjie Liang, Wentao Liu, Xiao Wang, Xiaohan Jia, Xiaotao Gu, Xiaoying Ling, Xin Wang, Xing Fan, Xingru Pan, Xinyuan Zhang, Xinze Zhang, Xiuqing Fu, Xunkai Zhang, Yabo Xu, Yandong Wu, Yida Lu, Yidong Wang, Yilin Zhou, Yiming Pan, Ying Zhang, Yingli Wang, Yingru Li, Yinpei Su, Yipeng Geng, Yitong Zhu, Yongkun Yang, Yuhang Li, Yuhao Wu, Yujiang Li, Yunan Liu, Yunqing Wang, Yuntao Li, Yuxuan Zhang, Zezhen Liu, Zhen Yang, Zhengda Zhou, Zhongpei Qiao, Zhuoer Feng, Zhuorui Liu, Zichen Zhang, Zihan Wang, Zijun Yao, Zikang Wang, Ziqiang Liu, Ziwei Chai, Zixuan Li, Zuodong Zhao, Wenguang Chen, Jidong Zhai, Bin Xu, Minlie Huang, Hongning Wang, Juanzi Li, Yuxiao Dong, and Jie Tang. 2025. GLM-4.5: Agentic, Reasoning, and Coding (ARC) Foundation Models. arXiv:2508.06471 [cs.CL]

- <https://arxiv.org/abs/2508.06471>
- [22] Hunyuan Vision Team, Pengyuan Lyu, Xingyu Wan, Gengluo Li, Shangpin Peng, Weinong Wang, Liang Wu, Huawen Shen, Yu Zhou, Canhui Tang, et al. 2025. HunyuanOCR Technical Report. *arXiv preprint arXiv:2511.19575* (2025).
 - [23] Qwen Team. 2025. Qwen3 Technical Report. arXiv:2505.09388 [cs.CL] <https://arxiv.org/abs/2505.09388>
 - [24] Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. 2024. Mineru: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839* (2024).
 - [25] Shu Wang, Yingli Zhou, and Yixiang Fang. 2025. BookRAG: A Hierarchical Structure-aware Index-based Approach for Retrieval-Augmented Generation on Complex Documents. *arXiv preprint arXiv:2512.03413* (2025).
 - [26] Haoran Wei, Yaofeng Sun, and Yukun Li. 2025. Deepseek-ocr: Contexts optical compression. *arXiv preprint arXiv:2510.18234* (2025).
 - [27] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.
 - [28] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, et al. 2025. Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models. *arXiv preprint arXiv:2506.05176* (2025).
 - [29] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025. Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models. *arXiv preprint arXiv:2506.05176* (2025).
 - [30] Qingfei Zhao, Ruobing Wang, Yukuo Cen, Daren Zha, Shicheng Tan, Yuxiao Dong, and Jie Tang. 2024. Longrag: A dual-perspective retrieval-augmented generation paradigm for long-context question answering. *arXiv preprint arXiv:2410.18050* (2024).

A Appendix

A.1 Benchmark Statistics

Table 5 details the statistics of the four datasets employed in our evaluation. The benchmarks are categorized into single-document and multi-document settings to assess the agent’s performance across different retrieval scopes. Notably, the single-document datasets pose a significant challenge regarding context length: FinanceBench averages approximately 165k tokens, while our constructed ContextBench reaches an average of 233k tokens, serving as a rigorous stress test for long-document reasoning capabilities.

Table 5: Statistics of the Datasets Used in Evaluation. The token counts are calculated based on the parsed Markdown content.

Dataset	Type	# Examples	Total Tokens	Avg. Tokens
FinanceBench [9]	Single-Doc	150	24,725,584	164,837
ContextBench (Ours)	Single-Doc	94	21,671,771	233,030
QASPER [4]	Multi-Doc	143	6,437,147	45,015
SyllabusQA [7]	Multi-Doc	196	2,800,293	14,287

A.2 Case Study

Tables 6, 7, 8, and 9 present representative DeepRead trajectories on ContextBench, FinanceBench, SyllabusQA, and QASPER, respectively. In the ContextBench example (Table 6), the task is to identify the set of agents in an analyst team. DeepRead first issues a semantic query (e.g., “Analyst Team agents members”) to localize the relevant region, then invokes `ReadSection` to consume the surrounding paragraphs in-order and produce a consolidated answer. This illustrates the intended locate-then-read workflow, where retrieval provides structural anchors and section-wise reading reduces fragmentation when evidence is concentrated within a coherent region.

The FinanceBench example (Table 7) targets a numeric field (FY2016 COGS for Microsoft) primarily supported by the statement of income. DeepRead retrieves candidate passages containing the consolidated income statement and then reads the corresponding section to extract the required figure from the table. Notably, due to document parsing errors, the `ReadSection` output in this case contains only a partial span (17 paragraphs) of the intended section; nevertheless, the relevant row remains present in the returned content, enabling DeepRead to answer correctly. This suggests that the coordinate-anchored interaction can remain effective under moderate parsing imperfections when the critical evidence is successfully localized and preserved in the readable span.

The SyllabusQA example (Table 8) concerns the academic integrity policy for ECE 4670 regarding discussion and collaboration on homework and laboratory assignments. DeepRead retrieves the policy passage and performs a short, contiguous read to confirm the exact wording, yielding a faithful, closed-form answer. This case highlights DeepRead’s ability to resolve policy questions by aligning retrieval with structure-preserving reading, even when multiple related concepts appear in close proximity.

Finally, the QASPER example (Table 9) demonstrates cross-document reasoning: DeepRead first localizes the evidence about the

number of evaluated languages in the first paper, then searches and reads the relevant section in the second paper to obtain the number of target datasets. By grounding each sub-claim in coordinate-addressable spans, the agent can aggregate results across documents while maintaining traceability to their respective sources.

A.3 Robustness Testing of LLM as a Judge

To ensure the robustness of our conclusions, we replicate the evaluation using three independent LLMs: **DeepSeek-V3.2**, **GLM-4.7**, and **Qwen3-235B**. The mean accuracies across all settings are highly consistent (70.53%, 67.90%, and 69.59% respectively), confirming that the observed improvements are not artifacts of a specific judge’s calibration. Furthermore, we compute the **inter-judge agreement** to validate evaluation reliability. We define the agreement score as the proportion of samples where all three judges reach a unanimous verdict. Let $J_m(x_i)$ be the verdict of judge m on sample i . The metric is calculated as:

$$\text{Agreement} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(J_1(x_i) = J_2(x_i) = J_3(x_i)) \quad (15)$$

where $\mathbb{I}(\cdot)$ denotes the indicator function. The system achieves a high overall agreement of **0.8863**. As expected, agreement is higher on single-document tasks (avg. 0.9187) compared to multi-document tasks (avg. 0.8540), reflecting the inherent complexity and slight subjectivity involved in evaluating cross-document reasoning.

A.4 Prompt Template

In order to ensure reproducibility and transparency within our methodology, we provide a detailed description of the prompt templates in this paper. Figure 5 presents the System Prompt, which integrates the hierarchical structure of the document and orchestrates the agent’s “locate-then-read” approach. This system prompt dictates how the LLM navigates and processes the document to answer user queries effectively. The prompt also includes specific guidelines to mitigate misinterpretation due to hierarchical errors or misalignments, such as distinguishing between body text and hierarchical elements based on the structure. Figure 6 illustrates the LLM-as-a-Judge prompt, which is designed to assess the correctness of AI-generated answers relative to a human-generated “golden answer”. Finally, Figure 7 showcases the prompt used for generating multi-hop question-answer pairs for our benchmarks. This prompt is particularly stringent, as it enforces constraints on both hierarchical dependency and cross-document reasoning. The requirement for questions to synthesize information from at least four distinct leaf sections ensures that the generated QA pairs reflect complex reasoning and comprehensive document understanding. Furthermore, the closed-form nature of the questions, requiring a single unambiguous answer, maintains precision in evaluation.

Table 6: An example from DeepRead on ContextBench. The model-generated semantic retrieval queries are enclosed within `<|begin_semantic_retrieval_query|>` and `<|end_semantic_retrieval_query|>`, while the top-2 retrieval results are enclosed within `<|begin_semantic_retrieval_result|>` and `<|end_semantic_retrieval_result|>`. Similarly, the read section queries are enclosed within `<|begin_read_section_query|>` and `<|end_read_section_query|>`, and the read section results are enclosed within `<|begin_read_section_result|>` and `<|end_read_section_result|>`.

Example #1 from ContextBench

Question:

Which Agents are there in the Analyst Team?

Labeled Answer:

There are four Agents in the Analyst Team:

- (1) Fundamental Analyst Agents
- (2) Sentiment Analyst Agents
- (3) News Analyst Agents
- (4) Technical Analyst Agents.

Model Output by DeepRead

ROUND 1

I'll help you find information about the agents in the Analyst Team. Let me search for relevant information in the document.

`<|begin_semantic_retrieval_query|>`Analyst Team agents members`<|end_semantic_retrieval_query|>`
`<|begin_semantic_retrieval_result(top1)|>`doc_id=1, sec_id=13, para_id=0, text="The Analyst Team (Figure 2) is composed of specialized agents responsible for gathering and analyzing various types of market data to inform trading decisions. Each agent focuses on a specific aspect of market analysis, bringing together a comprehensive view of the market's conditions."`<|end_semantic_retrieval_result(top1)|>`
`<|begin_semantic_retrieval_result(top2)|>`doc_id=1, sec_id=15, para_id=0, text="Trader Agents (Figure 4) are responsible for executing trading decisions based on the comprehensive analysis provided by the Analyst Team and the nuanced perspectives from the Researcher Team. They"`<|end_semantic_retrieval_result(top2)|>`

ROUND 2

Now let me read the section about the Analyst Team to get detailed information about the agents it contains.

`<|begin_read_section_query|>`doc_id=1, sec_id=13, start_para_id=0, end_para_id=8`<|end_read_section_query|>`
`<|begin_read_section_result(para_id=0)|>`The Analyst Team (Figure 2) is composed of specialized agents responsible for gathering and analyzing various types of market data to inform trading decisions. Each agent focuses on a specific aspect of market analysis, bringing together a comprehensive view of the market's conditions.`<|end_read_section_result(para_id=0)|>`

.....

`<|begin_read_section_result(para_id=7)|>`Collectively, the Analyst Team synthesizes data from multiple sources to provide a holistic market analysis. Their combined insights form the foundational input for the Researcher Team, ensuring that all facets of the market are considered in subsequent decision-making processes.`<|end_read_section_result(para_id=7)|>`

ROUND 3

Based on the document, the Analyst Team in TradingAgents consists of four specialized agents:

1. **Fundamental Analyst Agents** - These agents evaluate company fundamentals by analyzing financial statements, earnings reports, insider transactions, and other pertinent data to assess a company's intrinsic value and identify undervalued or overvalued stocks.
2. **Sentiment Analyst Agents** - These agents process social media posts, sentiment scores, and insider sentiments to gauge market sentiment and predict how collective investor behavior might impact stock prices in the short term.
3. **News Analyst Agents** - These agents analyze news articles, government announcements, and macroeconomic indicators to assess the market's macroeconomic state, major world events, and significant company changes.
4. **Technical Analyst Agents** - These agents calculate and select relevant technical indicators (such as MACD and RSI), analyze price patterns and trading volumes to forecast future price movements, and assist in timing entry and exit points.

The Analyst Team works collectively to synthesize data from multiple sources and provide holistic market analysis, which then serves as foundational input for the Researcher Team in the trading decision-making process.

Table 7: An example from DeepRead on FinanceBench. The model-generated semantic retrieval queries are enclosed within `<|begin_semantic_retrieval_query|>` and `<|end_semantic_retrieval_query|>`, while the top-2 retrieval results are enclosed within `<|begin_semantic_retrieval_result|>` and `<|end_semantic_retrieval_result|>`. Similarly, the read section queries are enclosed within `<|begin_read_section_query|>` and `<|end_read_section_query|>`, and the read section results are enclosed within `<|begin_read_section_result|>` and `<|end_read_section_result|>`.

Example #2 from FinanceBench
<p>Question: What is the FY2016 COGS for Microsoft? Please state answer in USD millions. Provide a response to the question by primarily using the statement of income.</p> <p>Labeled Answer: \$32780.00</p>
<p>Model Output by DeepRead</p> <p>ROUND 1 I'll help you find the FY2016 Cost of Goods Sold (COGS) for Microsoft from the statement of income. Let me start by searching for relevant information.</p> <p><code>< begin_semantic_retrieval_query ></code> FY2016 cost of goods sold COGS statement of income <code>< end_semantic_retrieval_query ></code> <code>< begin_semantic_retrieval_result(top1) ></code> doc_id=1, sec_id=299, para_id=8, text="<table border=1 style='margin: auto; width: max-content;*><tr><td style='text-align: center;*>Quarter Ended</td>....." <code>< end_semantic_retrieval_result(top1) ></code> <code>< begin_semantic_retrieval_result(top2) ></code> doc_id=1, sec_id=251, para_id=18, text="<table border=1 style='margin: auto; width: max-content;*><tr><td style='text-align: center;*>Year Ended June 30,</td>....." <code>< end_semantic_retrieval_result(top2) ></code></p> <p>ROUND 2 I need to find the actual statement of income to get the COGS figure. Let me search more specifically for the income statement or COGS information.</p> <p><code>< begin_semantic_retrieval_query ></code> income statement cost of revenue COGS fiscal year 2016 <code>< end_semantic_retrieval_query ></code> <code>< begin_semantic_retrieval_result(top1) ></code> doc_id=1, sec_id=99, para_id=0, text="Due to the growth in our cloud-based solutions, service revenue exceeded 10% of total revenue for the first time in fiscal year 2016. As a result, we have separately disclosed product revenue and service and other revenue on our consolidated income statements." <code>< end_semantic_retrieval_result(top1) ></code> <code>< begin_semantic_retrieval_result(top2) ></code> doc_id=1, sec_id=178, para_id=0, text="<table border=1 style='margin: auto; width: max-content;*><tr><td colspan="4">(In millions, except per share amounts)</td>....." <code>< end_semantic_retrieval_result(top2) ></code></p> <p>ROUND 3 Perfect! I found the consolidated income statement. From the table in node 178, I can see the FY2016 Cost of Revenue (which is equivalent to COGS - Cost of Goods Sold). Let me read the full section to confirm the exact figures.</p> <p><code>< begin_read_section_query ></code> doc_id=1, sec_id=178, start_para_id=0, end_para_id=30 <code>< end_read_section_query ></code> <code>< begin_read_section_result(para_id=0) ></code> <table border=1 style='margin: auto; width: max-content;*><tr><td colspan="4">(In millions, except per share amounts)</td></tr><tr><td style='text-align: center;*>Year Ended June 30,</td><td style='text-align: center;*>2016</td><td style='text-align: center;*>2015</td><td style='text-align: center;*>2014</td></tr><tr><td style='text-align: center;*>Revenue:</td>..... <code>< end_read_section_result(para_id=0) ></code></p> <p>..... <code>< begin_read_section_result(para_id=17) ></code> See accompanying notes. <code>< end_read_section_result(para_id=17) ></code></p> <p>ROUND 4 Based on the consolidated income statement from the document, I can see the FY2016 Cost of Revenue (which is equivalent to COGS - Cost of Goods Sold) for Microsoft Corporation.</p> <p>From the income statement table in paragraph 0, for the fiscal year ended June 30, 2016:</p> <p>*Total Cost of Revenue: \$32,780 million*</p> <p>This is broken down as:</p> <ul style="list-style-type: none">- Product cost of revenue: \$17,880 million- Service and other cost of revenue: \$14,900 million <p>*Answer: The FY2016 COGS for Microsoft was \$32,780 million USD.*</p>

Table 8: An example from DeepRead on SyllabusQA. The model-generated semantic retrieval queries are enclosed within `<|begin_semantic_retrieval_query|>` and `<|end_semantic_retrieval_query|>`, while the top-2 retrieval results are enclosed within `<|begin_semantic_retrieval_result|>` and `<|end_semantic_retrieval_result|>`. Similarly, the read section queries are enclosed within `<|begin_read_section_query|>` and `<|end_read_section_query|>`, and the read section results are enclosed within `<|begin_read_section_result|>` and `<|end_read_section_result|>`.

Example #3 from SyllabusQA

Question:

While KIN 270 lecture exams are open-book and timed, what does the academic integrity policy of ECE 4670 state regarding discussion and collaboration on homework and laboratory assignments?

Labeled Answer:

It is permitted and encouraged.

Model Output by DeepRead

ROUND 1

I'll help you find information about the academic integrity policy for ECE 4670 regarding discussion and collaboration on homework and laboratory assignments. Let me start by searching for relevant information.

`<|begin_semantic_retrieval_query|>`ECE 4670 academic integrity policy homework laboratory assignments collaboration discussion`<|end_semantic_retrieval_query|>`

`<|begin_semantic_retrieval_result(top1)|>`doc_id=1, sec_id=4, para_id=0, text="Students expected to abide by the Cornell University Code of Academic Integrity with work submitted for credit representing the student's own work. Discussion and collaboration on homework and laboratory assignments is permitted and encouraged, but final work should represent the student's own understanding. Specific examples of this policy implementation will be distributed in class. Course materials posted on Blackboard are intellectual property belonging to the author. Students are not permitted to buy or sell any course materials without the express permission of the instructor. Such unauthorized behavior will constitute academic misconduct.Students expected to abide by the Cornell University Code of Academic Integrity with work submitted for credit representing the student's own work. Discussion and collaboration on homework and laboratory assignments is permitted and encouraged, but final work should represent the student's own understanding. Specific examples of this policy implementation will be distributed in class. Course materials posted on Blackboard are intellectual property belonging to the author. Students are not permitted to buy or sell any course materials without the express permission of the instructor. Such unauthorized behavior will constitute academic misconduct."`<|end_semantic_retrieval_result(top1)|>`

`<|begin_semantic_retrieval_result(top2)|>`doc_id=1, sec_id=1, para_id=20, text="Homework: Weekly assignments. Approximately seven homework assignments per semester. Collaboration with students is encouraged."`<|end_semantic_retrieval_result(top2)|>`

ROUND 2

Perfect! I found the relevant information about ECE 4670's academic integrity policy. Let me read the full section to provide you with complete details.

`<|begin_read_section_query|>`doc_id=1, sec_id=4, start_para_id=0, end_para_id=1`<|end_read_section_query|>`

`<|begin_read_section_result(para_id=0)|>`Students expected to abide by the Cornell University Code of Academic Integrity with work submitted for credit representing the student's own work. Discussion and collaboration on homework and laboratory assignments is permitted and encouraged, but final work should represent the student's own understanding. Specific examples of this policy implementation will be distributed in class. Course materials posted on Blackboard are intellectual property belonging to the author. Students are not permitted to buy or sell any course materials without the express permission of the instructor. Such unauthorized behavior will constitute academic misconduct.`<|end_read_section_result(para_id=0)|>`

ROUND 3

Based on the ECE 4670 academic integrity policy, regarding discussion and collaboration on homework and laboratory assignments:

*Discussion and collaboration on homework and laboratory assignments is permitted and encouraged.** However, there's an important caveat: **the final work submitted for credit should represent the student's own work and understanding.**

The policy states that while students are encouraged to discuss and collaborate on these assignments, they must ensure that their final submissions reflect their personal comprehension of the material. This means students can work together, share ideas, and help each other understand concepts, but they cannot simply copy someone else's work or submit work that doesn't demonstrate their own understanding.

The policy also mentions that specific examples of how this policy is implemented will be distributed in class, so students should look for those additional guidelines from their instructor.

Table 9: An example from DeepRead on QASPER. The model-generated semantic retrieval queries are enclosed within `<|begin_semantic_retrieval_query|>` and `<|end_semantic_retrieval_query|>`, while the top-2 retrieval results are enclosed within `<|begin_semantic_retrieval_result|>` and `<|end_semantic_retrieval_result|>`. Similarly, the read section queries are enclosed within `<|begin_read_section_query|>` and `<|end_read_section_query|>`, and the read section results are enclosed within `<|begin_read_section_result|>` and `<|end_read_section_result|>`.

Example #4 from QASPER

Question:

How many distinct languages are evaluated in the first document, and how many target datasets are used for evaluation in the second document?

Labeled Answer:

Five languages and four datasets.

Model Output by DeepRead

ROUND 1

I need to understand what documents we’re working with. Let me first check the directory structure to see what documents are available.

`<|begin_semantic_retrieval_query|>` documents first document languages evaluated `<|end_semantic_retrieval_query|>`
`<|begin_semantic_retrieval_result(top1)|>` doc_id=1, sec_id=3, para_id=7, text=“such as Basque and German. We demonstrate the robustness of our approach by reporting best results for five languages (Basque, Dutch, German, English and Spanish) on 12 different datasets, including seven in-domain and eight out-of-domain evaluations.” `<|end_semantic_retrieval_result(top1)|>`

`<|begin_semantic_retrieval_result(top2)|>` doc_id=1, sec_id=29, para_id=1, text=“Thus, despite the relative simplicity of our approach, we report state of the art results for Dutch, English, German, Spanish and Basque in seven in-domain evaluations.” `<|end_semantic_retrieval_result(top2)|>`

ROUND 2

Now I need to check the second document mentioned in the question. Let me search for information about target datasets in the second document.

`<|begin_semantic_retrieval_query|>` target datasets evaluation second document `<|end_semantic_retrieval_query|>`
`<|begin_semantic_retrieval_result(top1)|>` doc_id=2, sec_id=10, para_id=0, text=“We use four MCQA datasets as the target datasets: DREAM (Sun et al. 2019), MCTest (Richardson, Burges, and Renshaw 2013), TOEFL (Ostermann et al. 2018), and SemEval-2018 Task 11 (Tseng et al. 2016), which are summarized in Table 2. For the first coarse-tuning stage with NLI tasks, we use MultiNLI (Williams, Nangia, and Bowman 2017) and SNLI (Young et al. 2014) as the out-of-domain source datasets. For the second stage, we use the current largest MCQA dataset, i.e., RACE (Lai et al. 2017) as in-domain source dataset. For all datasets, we use the official train/dev/test splits.” `<|end_semantic_retrieval_result(top1)|>`

`<|begin_semantic_retrieval_result(top2)|>` doc_id=1, sec_id=6, para_id=3, text=“<div style=“text-align: center;”>Table 1: Datasets used for training, development and evaluation. MUC7: only three classes (LOC, ORG, PER) of the formal run are used for out-of-domain evaluation. As there are not standard partitions of SONAR-1 and Ancora 2.0, the full corpus was used for training and later evaluated in-out-of-domain settings.</div>” `<|end_semantic_retrieval_result(top2)|>`

ROUND 3

Now I need to get more specific information about the first document. Let me read the section that mentions languages to get the exact count.

`<|begin_read_section_query|>` doc_id=1, sec_id=3, start_para_id=7, end_para_id=9 `<|end_read_section_query|>`
`<|begin_read_section_result(para_id=7)|>` such as Basque and German. We demonstrate the robustness of our approach by reporting best results for five languages (Basque, Dutch, German, English and Spanish) on 12 different datasets, including seven in-domain and eight out-of-domain evaluations. `<|end_read_section_result(para_id=7)|>`

ROUND 4

Now let me check if there are other mentions of languages in the first document to confirm the count.

`<|begin_semantic_retrieval_query|>` languages evaluated Basque Dutch German English Spanish `<|end_semantic_retrieval_query|>`
`<|begin_semantic_retrieval_result(top1)|>` doc_id=1, sec_id=16, para_id=0, text=“In this Section we report on the experiments performed with the ixa-pipe-nerc system as described in the previous section. The experiments are performed in 5 languages: Basque, Dutch, English, German and Spanish. For comparison purposes, in-domain results are presented in Section 4.1 using the most common NERC datasets for each language as summarized in Table 1. Section 4.2 analyzes the performance when reducing training data and Section 4.3 presents eight out-of-domain evaluations for three languages: Dutch, English and Spanish.” `<|end_semantic_retrieval_result(top1)|>`

`<|begin_semantic_retrieval_result(top2)|>` doc_id=1, sec_id=3, para_id=7, text=“such as Basque and German. We demonstrate the robustness of our approach by reporting best results for five languages (Basque, Dutch, German, English and Spanish) on 12 different datasets, including seven in-domain and eight out-of-domain evaluations.” `<|end_semantic_retrieval_result(top2)|>`

Table 9: An example from DeepRead on QASPER. The model-generated semantic retrieval queries are enclosed within `<|begin_semantic_retrieval_query|>` and `<|end_semantic_retrieval_query|>`, while the top-2 retrieval results are enclosed within `<|begin_semantic_retrieval_result|>` and `<|end_semantic_retrieval_result|>`. Similarly, the read section queries are enclosed within `<|begin_read_section_query|>` and `<|end_read_section_query|>`, and the read section results are enclosed within `<|begin_read_section_result|>` and `<|end_read_section_result|>`. (continued)

Example #4 from QASPER (continued)

Question:

How many distinct languages are evaluated in the first document, and how many target datasets are used for evaluation in the second document?

Labeled Answer:

Five languages and four datasets.

Model Output by DeepRead

ROUND 5

Now let me read the section from the second document about target datasets to get the exact count.

`<|begin_read_section_query|>` doc_id=2, sec_id=10, start_para_id=0, end_para_id=1 `<|end_read_section_query|>`
`<|begin_read_section_result(para_id=0)|>` We use four MCQA datasets as the target datasets: DREAM (Sun et al. 2019), MCTest (Richardson, Burges, and Renshaw 2013), TOEFL (Ostermann et al. 2018), and SemEval-2018 Task 11 (Tseng et al. 2016), which are summarized in Table 2. For the first coarse-tuning stage with NLI tasks, we use MultiNLI (Williams, Nangia, and Bowman 2017) and SNLI (Young et al. 2014) as the out-of-domain source datasets. For the second stage, we use the current largest MCQA dataset, i.e., RACE (Lai et al. 2017) as in-domain source dataset. For all datasets, we use the official train/dev/test splits. `<|end_read_section_result(para_id=0)|>`

ROUND 6

Based on the information I’ve gathered:

1. For the first document: The text clearly states “five languages (Basque, Dutch, German, English and Spanish)” are evaluated. This is confirmed in multiple sections.
2. For the second document: The text states “four MCQA datasets as the target datasets: DREAM, MCTest, TOEFL, and SemEval-2018 Task 11.”

Therefore, the answer is:

- The first document evaluates ****5 distinct languages**** (Basque, Dutch, German, English, and Spanish).
- The second document uses ****4 target datasets**** for evaluation (DREAM, MCTest, TOEFL, and SemEval-2018 Task 11).

Table 10: Accuracy (%) under three independent LLM judges. Each entry is reported as DeepSeek-V3.2 / GLM-4.7 / Qwen3-235B (in this order).

Setting	Single-Document			Multi-Document			Overall Avg
	FinanceBench	ContextBench	Avg	QASPER	SyllabusQA	Avg	
Search-o1	80.0/80.7/84.0	74.5/71.3/70.2	77.2/76.0/77.1	65.0/58.7/55.2	57.1/55.6/58.7	61.1/57.2/57.0	69.2/66.6/67.0
Search-o1 w/ expand	83.3/83.3/85.3	84.0/79.8/80.9	83.7/81.6/83.1	65.0/56.6/60.1	68.4/67.3/65.8	66.7/62.0/63.0	75.2/71.8/73.0
DeepRead	82.7/84.0/88.7	91.5/87.2/91.5	87.1/85.6/90.1	72.7/69.9/65.7	70.9/71.4/71.4	71.8/70.7/68.6	79.5/78.1/79.3
DeepRead w/ expand	84.0/84.0/88.0	88.3/84.0/87.2	86.1/84.0/87.6	76.2/63.6/62.2	72.4/69.4/73.0	74.3/66.5/67.6	80.2/75.3/77.6
Read-only	80.7/84.0/89.3	91.5/87.2/86.2	86.1/85.6/87.8	18.9/12.6/15.4	15.3/18.9/22.4	17.1/15.7/18.9	51.6/50.7/53.3

Table 11: Inter-judge agreement (higher is more consistent). Agreement is computed by our evaluation script and reflects example-level consistency of the three judges' binary verdicts (correct/incorrect).

Setting	Single-Document			Multi-Document			Overall Avg
	FinanceBench	ContextBench	Avg	QASPER	SyllabusQA	Avg	
Search-o1	0.9400	0.9043	0.9221	0.8462	0.8622	0.8542	0.8882
Search-o1 w/ expand	0.9200	0.9043	0.9121	0.8182	0.8214	0.8198	0.8660
DeepRead	0.9133	0.9255	0.9194	0.8671	0.8827	0.8749	0.8972
DeepRead w/ expand	0.9267	0.9362	0.9314	0.7902	0.8724	0.8313	0.8814
Read-only	0.8800	0.9362	0.9081	0.8601	0.9082	0.8842	0.8961
Avg (over settings)	0.9160	0.9213	0.9186	0.8364	0.8694	0.8529	0.8858

System Prompt Template

You are a documents assistant and will receive one or more documents structured as follows:

- (doc_id) [node_id] Title | paragraphs=Num | tokens=Num | children=[ID list].

Use this structure and your available tools to answer the user's question.

Guidelines

- Use <Search/Retrieve> to locate relevant nodes based on the directory.
- Answer strictly based on the provided corpus; do not fabricate.
- The hierarchical structure of documents is represented in the Directory Structure. Parsing errors may cause body text to be mistakenly treated as hierarchical elements (or headings), rendering the heading text inaccessible to search and reading tools. Please make reasonable inferences based on the structure and the content returned by the tool.
- Respond in the User's language; align queries with the Directory Structure.
- Usually, you need to think step by step and then call tools to locate or read, iterating in this way until you can answer the question.
- When calling tools, DO NOT write tool invocations in plain text. Use the structured tool call interface (tool_calls) only.

Directory Structure

<Serialized Table of Contents inserted here...>

Figure 5: The system prompt used in DeepRead. It injects the hierarchical document skeleton (Directory Structure).**LLM-as-a-Judge Prompt Template**

You are an expert evaluator for AI-generated responses to queries. Your task is to determine whether the AI-generated answer correctly answers the query based on the golden answer provided by a human expert.

Numerical Accuracy:

- Rounding differences should be **ignored** if they do not meaningfully change the conclusion.
- You can allow some flexibility in accuracy. For example, 1.2 is considered similar to 1.23. Two numbers are considered similar if one can be rounded to the other.
- Fractions, percentage, and numerics could be considered similar, for example: "11 of 14" is considered equivalent to "79%" and "0.79".

Evaluation Criteria:

- If the golden answer or any of its equivalence can be inferred or generated from the AI-generated answer, then the AI-generated answer is considered correct.
- If any number, percentage, fraction, or figure in the golden answer is not present in the AI-generated answer, but can be inferred or generated from the AI-generated answer or implicitly exist in the AI-generated answer, then the AI-generated answer is considered correct.
- The AI-generated answer is considered correct if it conveys the same or similar meaning, conclusion, or rationale as the golden answer.
- If the AI-generated answer is a superset of the golden answer, it is also considered correct.
- If the AI-generated answer provides a valid answer or reasonable interpretation compared to the golden answer, it is considered correct.
- If the AI-generated answer contains subjective judgments or opinions, it is considered correct as long as they are reasonable and justifiable compared to the golden answer.
- Otherwise, the AI-generated answer is incorrect.

Inputs: Query: <question>

- AI-Generated Answer: <predicted_answer>
- Golden Answer: <standard_answer>

Your output should be ONLY a boolean value: True or False, nothing else.

Figure 6: The evaluation prompt used for the LLM-as-a-Judge metric. It instructs the evaluator model to focus on semantic equivalence and allow for flexible numerical matching.

QA Generation Prompt Template

System Instruction:

You are a careful question writer tasked with generating hard but answerable CLOSED-FORM QA pairs from MULTIPLE Markdown documents provided in the user message. You MUST use ONLY the provided document TEXT; do not invent facts, do not use images, filenames, URLs, or any external knowledge. Respect the document content.

CRITICAL CONSTRAINTS (apply to EACH item):

- The question must require synthesizing evidence from **AT LEAST FOUR distinct LEAF sections** (headings with no deeper subheadings). Treat '#', '##', etc. as hierarchical.
- The answer must uniquely follow by integrating information drawn from **MULTIPLE hierarchy levels** (e.g., a leaf and its ancestors, or leaves under different parents). It must NOT degenerate to fewer than 4 hops.
- **CLOSED-FORM:** The question has a single, unambiguous, concise correct answer present in the text.
- The QUESTION must **NOT** explicitly mention any section titles or heading levels. Phrase naturally (no "as stated in Section 3.2").
- Your question must be answered by considering **all** the provided documents together.

OUTPUT FORMAT:

- Return ONLY valid JSON with exact schema: {"qas": [{"q": "...", "a": "..."}, ...]}
- Do NOT include any other keys. Do NOT wrap in markdown.

User Input Template:

You will be given <k> Markdown documents (The following documents are unordered.) concatenated below.

Generate exactly <n> QA pairs that each satisfy the constraints.

DOCUMENTS BEGIN

<docs_text>

DOCUMENTS END

Figure 7: The prompt used to synthesize multi-hop QA pairs.