# BookRAG: A Hierarchical Structure-aware Index-based Approach for Retrieval-Augmented Generation on Complex Documents

Shu Wang
The Chinese University of Hong Kong, Shenzhen
shuwang3@link.cuhk.edu.cn

Yingli Zhou
The Chinese University of Hong Kong, Shenzhen
yinglizhou@link.cuhk.edu.cn

Yixiang Fang
The Chinese University of Hong Kong, Shenzhen
fangyixiang@cuhk.edu.cn

## Abstract

As an effective method to boost the performance of Large Language Models (LLMs) on the question answering (QA) task, Retrieval-Augmented Generation (RAG), which queries highly relevant information from external complex documents, has attracted tremendous attention from both industry and academia. Existing RAG approaches often focus on general documents, and they overlook the fact that many real-world documents (such as books, booklets, handbooks, etc.) have a hierarchical structure, which organizes their content from different granularity levels, leading to poor performance for the QA task. To address these limitations, we introduce BookRAG, a novel RAG approach targeted for documents with a hierarchical structure, which exploits logical hierarchies and traces entity relations to query the highly relevant information. Specifically, we build a novel index structure, called BookIndex, by extracting a hierarchical tree from the document, which serves as the role of its table of contents, using a graph to capture the intricate relationships between entities, and mapping entities to tree nodes. Leveraging the BookIndex, we then propose an agent-based query method inspired by the Information Foraging Theory, which dynamically classifies queries and employs a tailored retrieval workflow. Extensive experiments on three widely adopted benchmarks demonstrate that BookRAG achieves state-of-the-art performance, significantly outperforming baselines in both retrieval recall and QA accuracy while maintaining competitive efficiency.

## 1 Introduction

Large Language Models (LLMs) such as Qwen 3 [60] and Gemini 2.5 [13] have revolutionized the Question Answering (QA) system [15, 61, 65]. The industry has increasingly adopted LLMs to build QA systems that assist users and reduce manual effort in many applications [65, 67], such as financial auditing [29, 37], legal compliance [8], and scientific discovery [56]. However, directly relying on LLMs may lead to missing domain knowledge and generating outdated or unsupported information. To address these issues, Retrieval-Augmented Generation (RAG) has been widely adopted [17, 22] by retrieving relevant domain knowledge from external sources and using it to guide the LLM during response generation. On the other hand, in real-world enterprise scenarios, domain knowledge is often stored in long-form documents, such as technical handbooks, API reference manuals, and operational guidebooks [49]. A notable feature of such documents is that they follow the structure of books, characterized by intricate layouts and rigorous logical hierarchies (e.g., explicit tables of contents, nested chapters, and multi-level sections). In this paper, we aim to design an effective RAG system for QA over long and highly structured documents.
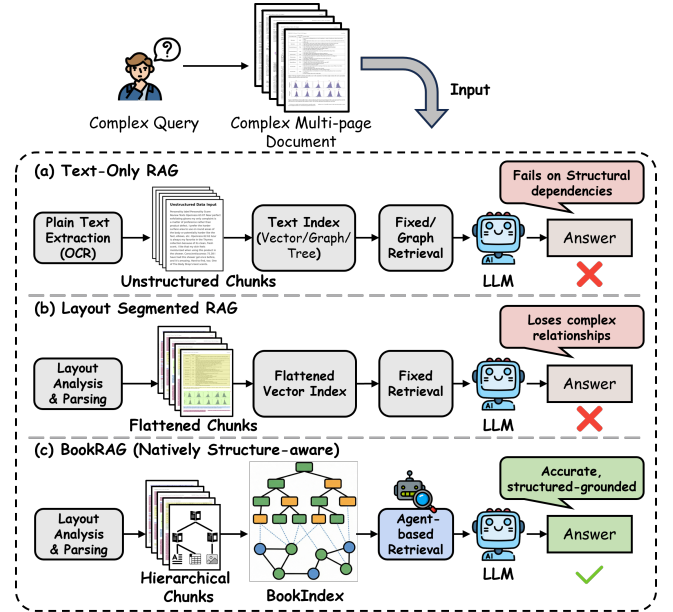


**Figure 1: Comparison of existing methods and BookRAG for complex document QA.**

• **Prior works.** The existing RAG approaches for document-level QA generally fall into two paradigms, as illustrated in Figure 1. The first paradigm relies on OCR (Optical Character Recognition) to convert the document into plain text, after which any text-based RAG method can be directly applied. Among text-based RAG methods, state-of-the-art approaches increasingly adopt *graph-based RAG* [6, 62, 66], where graph data serves as an external knowledge source because it captures rich semantic information and the relational structure between entities. As shown in Table 1, two representative methods are GraphRAG [16] and RAPTOR [45]. Specifically, GraphRAG first constructs a knowledge graph (KG) from the textual corpus, and then applies the Leiden community detection algorithm [51] to obtain hierarchical clusters. Summaries are generated for each community, providing a comprehensive, global overview of the entire corpus. RAPTOR builds a recursive tree structure by iteratively clustering document chunks and summarizing them at each level, enabling the model to capture both fine-grained and high-level semantic information across the corpus.

In contrast, the second paradigm, *layout-aware segmentation* [5, 52], first parses the document into structured blocks that preserve the original layout and information of the document, such as paragraphs, tables, figures, or equations. By doing so, it not only avoids the fixed chunk size used in the first paradigm, which often leads

**Table 1: Comparison of representative methods and our BookRAG.**

| Type | Representative Method | Core Feature | Multi-hop Reasoning | Document Parsing | Query Workflow |
|---|---|---|---|---|---|
| Graph-based | RAPTOR [45] | Recursive summarization | ✓ | ✗ | Static |
| | GraphRAG [16] | Global community detection | ✓ | ✗ | Static |
| Layout segmented | MM-Vanilla | Multi-modal retrieval | ✗ | ✓ | Static |
| | DocETL [47] | LLM-based document processing pipeline | ✗ | ✓ | Manual |
| **Doc-Native** | **BookRAG (Ours)** | Structure-award Index & Agent-based retrieval | ✓ | ✓ | Dynamic |

to fragmented information, but also retains document-native structural information. These blocks often exhibit multimodal characteristics, and a typical approach is to apply multimodal retrieval to obtain relevant content for answering queries. Recently, a state-of-the-art method in this category, DocETL [47], provides a declarative interface that allows users to manually define LLM-based processing pipelines to analyze the retrieved blocks. These pipelines consist of LLM-powered operations combined with task-specific optimizations.

• **Limitations of existing works.** However, these methods suffer from two fundamental limitations (**L** for short): **L1: Failure to capture the deep connection of document structure and semantics.** Text-based approaches cannot capture the structural layout of the document, resulting in the loss of important relationships stored in the hierarchical blocks, such as tables nested within a specific section. While layout-segmented methods preserve document structure, they cannot capture the relationships between different blocks in the document, which limits their capability for multi-hop reasoning across these blocks and ultimately affects their overall performance. **L2: Static of query workflows.** In real-world QA scenarios, user queries are highly heterogeneous, ranging from simple keyword lookups to complex multi-hop questions that require synthesizing evidence scattered across different parts of the document. Applying a uniform strategy, such as static or manually predefined workflows, to diverse needs is inefficient; for example, complex queries often require question decomposition, whereas simple queries do not.

• **Our technical contributions.** To bridge this gap, we introduce **BookRAG**, the first retrieval-augmented generation method built upon a document-native **BookIndex**, designed to document QA tasks. Specifically, to capture the deep connection of the relation in the document, BookIndex organizes information through two complementary structures. First, to preserve the document's native logical hierarchy, we organize the parsed content blocks into a hierarchical tree structure, which serves as the role of its table of contents. Second, to capture the intricate relations within these blocks, we construct a KG containing fine-grained entities. Finally, we unify these two structures by mapping the KG entities to their corresponding tree nodes.

However, effective multi-hop reasoning on the graph relies on a high-quality KG [62, 66], which is often compromised by entity ambiguity (e.g., distinct entities with names like "LLM" and "Large Language Model"). To address this, we propose a novel gradient-based entity resolution method that analyzes the similarity distribution of candidate entities. By identifying sharp drops in similarity scores, we can efficiently distinguish and merge coreferent entities, thereby ensuring graph connectivity and enhancing reasoning capabilities.

Building upon the BookIndex, we address the static of query workflows (**L2**) by implementing an **agent-based retrieval**. Specifically, our agent first classifies user queries based on their intent and complexity, and then dynamically generates tailored retrieval workflows. Grounded in *Information Foraging Theory* [42], our retrieval process mimics foraging by using *Selector* to narrow down the search space via information scents and *Reasoner* to locate highly relevant evidence.

We conduct extensive experiments on three widely adopted datasets to validate the effectiveness and efficiency of our BookRAG, comparing it against several state-of-the-art baselines. The experimental results demonstrate that BookRAG consistently achieves superior performance in both retrieval recall and QA accuracy across all datasets. Furthermore, our detailed analysis validates the critical contributions of our key features, such as the high-quality KG and the agent-based retrieval mechanism.

We summarize our contributions as:

- We introduce **BookRAG**, a novel method that constructs a document-native *BookIndex* by integrating a hierarchical tree of document layout blocks with a KG storing fine-grained entity relations.
- We propose an **Agent-based Retrieval** approach inspired by Information Foraging Theory, which dynamically classifies queries and configures optimal retrieval workflows to locate highly relevant evidence within documents.
- Extensive experiments on multiple benchmarks show that *BookRAG* significantly outperforms existing baselines, attaining state-of-the-art performance in solving complex document QA tasks while maintaining competitive efficiency.

**Outline.** We review related work in Section 2. Section 3 introduces the problem formulation, IFT, and RAG workflow. In Section 4, we present the structure of our BookIndex and its construction. Section 5 presents our agent-based retrieval, elaborating on the query classification and operators used in the structured execution of BookRAG. We present the experimental results and detailed analysis in Section 6, and conclude the paper in Section 7.

## 2 Related Work

In this section, we review the related works, including LLM in document analysis and the modern representative RAG approaches.

• **LLM in document analysis.** Recent advances in LLMs have offered opportunities to leverage LLMs in document data analysis. Due to the robust semantic reasoning capabilities of LLMs, there is an increasing number of works focusing on transferring unstructured documents (e.g., HTML, PDFs, and raw text) into structured formats, such as relational tables [1, 7, 25, 38]. For example, Evaporate [1] utilizes LLMs to synthesize extraction code, enabling cost-effective conversion of semi-structured web documents into structured databases without heavy manual annotation. In addition, several LLM-based document analysis systems have been proposed to equip standard data pipelines with semantic understanding [28, 40, 47, 53]. For instance, LOTUS [40] extends the relational model with semantic operators, allowing users to execute SQL-like queries with LLM-powered predicates (e.g., filter, join) over unstructured text corpora. Similarly, DocETL [47] introduces an agentic framework to optimize complex information extraction tasks. Furthermore, another line of research proposes to directly analyze or parse documents by viewing the document pages as images, thereby preserving critical layout and visual information [26, 31, 54].

• **RAG approaches.** RAG has been proven to excel in many tasks, including open-ended question answering [24, 48], programming context [9, 10], SQL rewrite [30, 50], and data cleaning [35, 36, 43]. The naive RAG technique relies on retrieving query-relevant contexts from external knowledge bases to mitigate the "hallucination" of LLMs. Recently, many RAG approaches [16, 18, 19, 21, 27, 32, 32, 45, 55, 58, 66] have adopted graph structures to organize the information and relationships within documents, achieving improved overall retrieval performance. For more details, please refer to the recent survey of graph-based RAG methods [41]. Besides, the Agentic RAG paradigm has been widely studied, employing autonomous agents to dynamically orchestrate and refine the RAG pipeline, thus significantly boosting the reasoning robustness and generation fidelity [2, 23, 59].

## 3 Preliminaries

This section formalizes the research problem of complex document QA, introduces the foundational Information Foraging Theory (IFT), and briefly reviews the general workflow of RAG systems

### 3.1 Problem Formulation

We study the problem of Question Answering (QA) over complex documents, which aims to answer user queries based on long-form documents [5, 11, 33]. Formally, a document $D$ is represented as a sequence of $N$ pages, $D = \{P_i\}_{i=1}^N$. These pages collectively contain a sequence of content blocks $\mathcal{B} = \{b_j\}_{j=1}^M$, where each block $b_j$ represents a distinct element (e.g., text segment, section header, table, or image) organized within a logical chapter hierarchy. Given a user query $q$, the goal is to generate an accurate answer $A$, ideally grounded in a specific set of evidence blocks $E \subset \mathcal{B}$. The task is formulated as developing a method $\mathcal{S}$ that maps the structured document and the query to the final answer:

$$A = \mathcal{S}(D, q) \tag{1}$$

where $\mathcal{S}$ should navigate both the sequential page content and the logical hierarchy of $D$ to synthesize the response.

### 3.2 Information Foraging Theory

Information Foraging Theory (IFT) [42] provides a framework for understanding information access as a process analogous to animal foraging. It suggests that users follow cues, known as **information scent** (e.g., keywords or icons), to navigate between clusters of content, known as **information patches** (e.g., sections in handbooks). The goal is to maximize the rate of valuable information gain while minimizing effort, guiding the decision to either stay within a patch or seek a new one.

Consider experts seeking a solution to a specific problem within a large technical handbook. They first extract key terms related to the problem, which act as information scent. This scent guides them to navigate towards one or more promising sections (the information patches). Within these patches, they analyze the diverse content to extract the precise knowledge required to formulate a final answer

### 3.3 RAG workflow

Retrieval-Augmented Generation (RAG) systems typically operate in a two-phase framework [6, 16, 41]. In the **Offline Indexing** phase, unstructured corpus data is organized into a structured index, which can take various forms such as vector databases or KG [66]. Subsequently, in the **Online Retrieval** phase, the system retrieves relevant components (e.g., text chunks or subgraphs) based on the user query $q$ to inform the LLM's generation. However, these general workflows often treat the index as a structure derived purely from content, potentially detaching it from the document's original logical hierarchy. In contrast, our approach seeks to deeply integrate these retrieval structures with the document's native tree topology.

## 4 BookIndex

This section introduces our proposed **BookIndex**, a hierarchical structure-aware index designed to capture both the explicit logical hierarchy and the intricate entity relations within complex documents. We first formally define the structure of the BookIndex ($B$). Subsequently, we elaborate on the sequential, two-stage construction process: (1) **Tree Construction**, which parses the document's layout to establish a hierarchical nodes, each categorized by type; and (2) **Graph Construction**, which extracts fine-grained entity knowledge from the tree nodes and refines it through a novel gradient-based entity resolution method.

### 4.1 Overview of BookIndex

We formally define our **BookIndex** as a triplet $B = (T, G, M)$. Here, $T = (N, E_T)$ represents a *Tree* structure where $N$ is the set of nodes derived from the document's explicit logical hierarchy (e.g., titles, sections, tables), and $E_T$ denotes their nesting relationships. $G = (V, E_G)$ is a *Knowledge Graph* that captures fine-grained entities ($V$) and their relations ($E_G$) scattered throughout the document. Finally, $M : V \to \mathcal{P}(N)$ is the *Graph-Tree Link (GT-Link)*, which links each entity in $V$ to the set of specific tree nodes in $N$ from which it was extracted. These links are crucial for capturing the intricate, cross-sectional relations within the document. The hierarchical tree nodes in $T$ serve as the document's native *information patches*, providing structured contexts for information seeking. Meanwhile, the entities and relations in $G$, connected via $M$, act as the rich
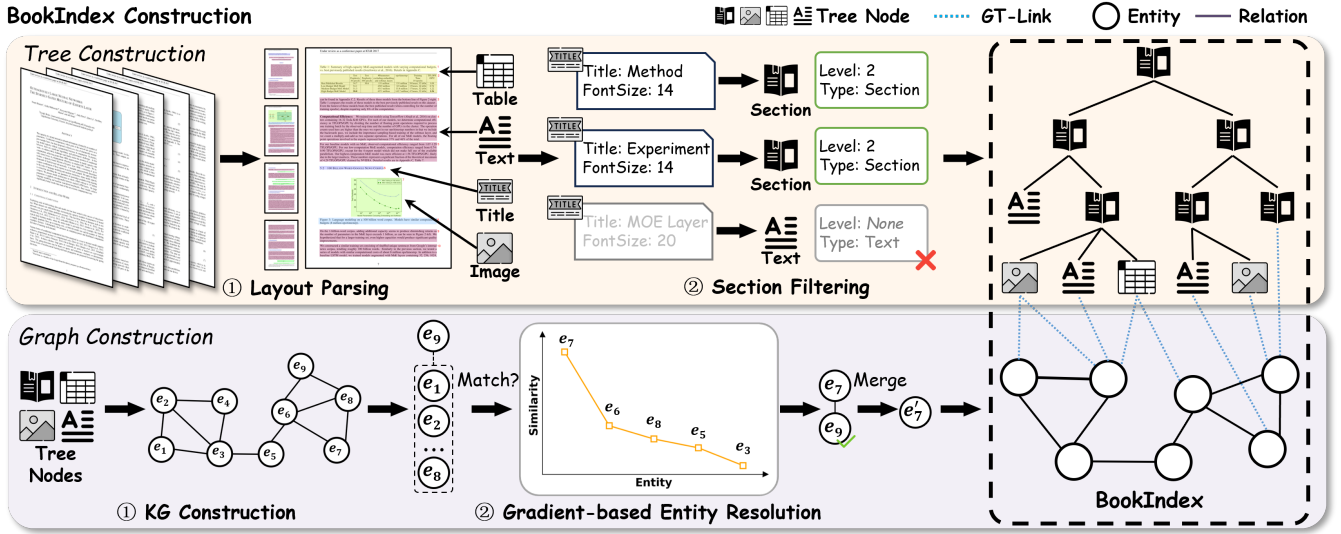
**Figure 2: The BookIndex Construction process. This phase includes Tree Construction, derived from Layout Parsing and Section Filtering, and Graph Construction, which involves KG Construction and Gradient-based Entity Resolution.**

*information scent* that guides navigation between and within these patches.

Figure 2 provides an example of our BookIndex. The Tree component, positioned at the top, organizes the document into a hierarchical structure, where content blocks such as text, tables, and images serve as leaf nodes nested within section nodes. The Graph component is composed of entities and relations extracted from these nodes. The GT-Link, illustrated by the blue dotted lines, explicitly connects these entities back to their corresponding tree nodes, thereby grounding the semantic entities within the document's logical hierarchy.

### 4.2 Tree Construction

The first stage transforms the raw document into a structured hierarchical tree $T$. This involves two key steps: robust layout parsing and intelligent section filtering.

*4.2.1 Layout Parsing.* The Layout Parsing phase processes the input document $D$ (a collection of pages) using layout analysis and recognition models. This step identifies, extracts, and organizes diverse blocks (e.g., text, tables, images) from the document pages. The output is a sequence of primitive blocks, $\mathcal{B} = \{b_1, b_2, \cdots, b_k\}$, where each block $b_i = (c_i, \tau_i, f_i)$ is defined as a triplet. Here, $c_i$ is the raw content (e.g., text, image data), $\tau_i$ is the initial layout-based type (e.g., Title, Text, Table, Image), and $f_i$ is a vector of associated layout features (e.g., "FontSize", bounding box).

*4.2.2 Section Filtering.* Next, the Section Filtering phase processes this initial sequence to identify the document's logically hierarchical structure. Layout Parsing identifies blocks as Title but does not assign their hierarchical level. Therefore, we select the candidate subset $\mathcal{B}_{\text{title}} \subset \mathcal{B}$ (where $\tau_i = \text{Title}$) for an LLM-based analysis. To handle extremely long documents, this analysis is performed in batches, where each batch retains a contextual window of high-level

section information (with $l = 1$ as the root). The LLM analyzes the content $c_i$ and layout features $f_i$ of the candidates to determine two key properties: their actual hierarchical level $l_i \in \{1, 2, ...\}$ and final node type $\tau_i'$ (e.g., re-classifying an erroneous Title as Text if its level is "None"). This step is crucial for preserving the document's logical hierarchy by correcting blocks erroneously parsed as Title, such as descriptive text within images or borderless table headers.

Finally, the definitive tree $T = (N, E_T)$ is constructed. The node set $N$ is composed of all blocks from the filtering and re-classification process, where each node $n \in N$ retains its content ($c_i$) and its final node type ($\tau_i'$) (e.g., Text, Section, Table, and Image). The edge set $E_T$, representing the parent-child nesting relationships, is then established. Parent-child relationships are inferred by sequentially traversing the nodes, using both the determined hierarchical levels ($l_i$) of Section nodes and the overall document order to assemble the complete tree structure.

As an example shown in Figure 2, the *Layout Parsing* phase identifies diverse blocks, typing them as Title, Text, Table, and Image. During the *Section Filtering* phase, the Title candidates (e.g., "Method", "Experiment", and "MOE Layer") are analyzed by the LLM. The blocks "Method" and "Experiment" (both with "FontSize: 14") are correctly identified as Section nodes at "Level: 2". Conversely, the "MOE Layer" block ("FontSize: 20"), which was erroneously tagged as Title by the parser, is re-classified by the LLM as a Text node with "Level: None". This correction is crucial for preserving the document's logical hierarchy. Following this process, all filtered and classified nodes are assembled into the final tree structure based on their determined levels and document order.

### 4.3 Graph Construction

Once the tree $T$ is established, we proceed to populate the Knowledge Graph $G$ by extracting and refining entities from the tree nodes.

---

**Algorithm 1:** Gradient-based entity resolution

---

**Input:** KG $G$, New entity $v_n$, Rerank model $\mathcal{R}$, Entity vector database $DB$, Vector search number $top\_k$, threshold of gradient $g$

    // Vector Search $top\_k$ relevant entities in $DB$.

1  $E_c \leftarrow \text{Search}(DB, v_n, top\_k)$;

2  $\mathcal{S} \leftarrow \mathcal{R}(E_c, v_n)$;

    // Sort all candidate entities by rerank scores.

3  $\text{Sort}(E_c, \mathcal{S})$;

4  $score \leftarrow \mathcal{S}[0], Sel \leftarrow E_c[0]$;

    // Gradient select similar entities.

5  **for** *each remain entity* $v_c \in E_c \setminus \{E_c[0]\}$ **do**

6     **if** $\mathcal{S}[v_c] > score/g$ **then**

7       $Sel \leftarrow Sel \cup \{v_c\}, score \leftarrow \mathcal{S}[v_c]$;

8     **else** break;

    // Merge entity or add new entity.

9  **if** $length(Sel) = length(E_c)$ **then**

10    $G \leftarrow \text{AddNewEntity}(G, v_n), DB \leftarrow \text{AddNew}(DB, v_n)$;

11 **else**

12    **if** $length(Sel) = 1$ **then** $v_{sel} \leftarrow Sel[0]$;

13    **else** $v_{sel} \leftarrow \text{LLMSelect}(Sel)$;

14    $G \leftarrow \text{MergeEntity}(G, v_n, v_{sel}), DB \leftarrow \text{Update}(DB, v_{sel}, v_n)$;

15 **return** $G, DB$;

---

### 4.3.1 KG Construction.

We iterate each node $n_i \in N$ from the previously constructed tree $T$. For each node $n_i$, we extract a subgraph $g_i = (V_i, E_{Ri})$ based on its content $c_i$ and final node type $\tau_i'$. This extraction is modality-dependent: if the node is text-only, an LLM is prompted to extract entities and relations, while for nodes containing visual elements (e.g., $\tau_i' = \text{Image}$), a Vision Language Model (VLM) is employed to extract visual knowledge. Crucially, for every entity $v \in V_i$ extracted, its origin tree node $n_i$ is recorded, which is vital for constructing the final mapping $M$.

Furthermore, to preserve structural semantics for specific logical types (e.g., Table, Formula), our process first creates a distinct, typed entity (e.g., $v_{\text{table}}$ representing the table itself). The other extracted entities from the specific node's content are linked to this primary vertex. For Table nodes specifically, row and column headers are also explicitly extracted as distinct entities and linked to $v_{\text{table}}$ via a "ContainedIn" relationship.

### 4.3.2 Gradient-based Entity Resolution.

As shown in the literature [62, 66], a well-constructed KG is essential for document question answering. A common challenge in the extraction process is that the same conceptual entity is often fragmented into multiple distinct entities due to abbreviations, co-references, or its varied occurrences across different document sections. This necessitates a robust Entity Resolution (ER) process, which identifies and merges these fragmented entities to refine the raw KG.

However, conventional ER methods are computationally expensive. They are often designed for batch processing across multiple data sources (commonly referred to as dirty ER), aiming to ensure accurate entity resolution by finding all possible matching pairs [12]. This process typically requires finding the transitive closure of all detected matches. That is, to definitively merge multiple entities (e.g., A, B, and C) as the same concept, the system must ideally compare all possible pairs ("A-B", "A-C", and "B-C") to confirm their equivalence. This can lead to a quadratic ($O(n^2)$)

number of pairwise comparisons, a process that becomes prohibitively slow and computationally expensive when relying on LLMs for high-accuracy judgments.

To address this, we employ a gradient-based ER method, operating on a single document (simplified as the clean ER), which performs ER incrementally as each new entity $v_n$ is extracted. This transforms the quadratic batch problem into a simpler, repeated lookup task: determining where the single new entity $v_n$ fits among the already-processed entities in the database. This incremental process yields two distinct, observable scoring patterns when $v_n$ is reranked against its $top\_k$ most relevant candidates:

- *Case A: New Entity.* If $v_n$ is a new conceptual entity, its relevance scores against all existing entities will be uniformly low, showing no significant gradient or discriminative pattern.
- *Case B: Existing Entity.* If $v_n$ is an alias of an existing entity, its scores will show a high relevance to the true match (or a small set of equivalent aliases). Due to the reranker's inherent discriminative limitations, this initial high-relevance set might occasionally contain multiple similar entities. This high-relevance set is then typically followed by a **sharp decline** (a large "gradient") before transitioning to a **gradual slope** of irrelevant entities.

Our Gradient-based ER algorithm is designed precisely to detect this sharp decline (characteristic of Case B), allowing us to efficiently isolate the high-relevance set. Subsequently, an LLM is utilized for finer-grained distinction when multiple similar entities are identified within this set, differentiating it from the "no gradient" scenario (Case A) without quadratic comparisons.

Algorithm 1 shows the above entity resolution process. For a new entity $v_n$, we first retrieve its $top\_k$ candidates $E_c$ from the vector database $DB$, which are then reranked by $\mathcal{R}$ against $v_n$ and sorted based on their scores $\mathcal{S}$ (Lines 1-3). We initialize the selection set $Sel$ with the top-scoring candidate $E_c[0]$ and set the initial score to $\mathcal{S}[0]$ (Line 4). We then iterate through the remaining sorted candidates (Lines 5-8). The core logic checks if the current score $\mathcal{S}[v_c]$ is still within the gradient threshold $g$ of the previous score (i.e., $\mathcal{S}[v_c] > score/g$). If the score drop is gentle (passes the check), the candidate $v_c$ is added to $Sel$, and score is updated (Lines 7-8); otherwise, the loop breaks (Line 8) as soon as a sharp score drop is detected. Finally, the algorithm makes its decision (Lines 9-14). If the selection set $Sel$ is identical to $E_c$, this indicates that all candidates passed the gradient check. This corresponds to *Case A*, where the scores lacked discriminative power (i.e., $v_n$ is equally dissimilar to all candidates), so $v_n$ is added as a new entity (Line 9-10). Conversely, if a gradient was found (i.e., $length(Sel) < length(E_c)$), this signals *Case B*. We then select the canonical entity $v_{sel}$ from $Sel$, using an LLM (Line 13) if the reranker identifies multiple aliases, and merge $v_n$ with it (Lines 12-14). The updated $G$ and $DB$ are then returned (Line 15).

For instance, considering the example in Figure 2, when the new entity $e_9$ is processed, it is first compared with existing entities in the KG. As depicted in the similarity curve (orange line), $e_9$ shows high similarity with $e_7$, followed by a sharp decline in similarity with other entities like $e_6$, $e_8$, and $e_5$. Our gradient-based selection process identifies $e_7$ as the unique, high-confidence match

for $e_9$. Consequently, $e_9$ is merged with $e_7$, enriching the KG with consolidated information as shown in the final merged entity $e_7'$.

**Graph-Tree Link (GT-Link).** The GT-Link $M$ is formalized to complete the BookIndex $B = (T, G, M)$. As described in the *KG Construction* phase, the origin tree node $n_i$ is recorded for every newly extracted entity $v_i$. GT-Link is then refined during entity resolution: when an entity $v_n$ is merged into a canonical entity $v_{sel}$, the origin node set of $v_{sel}$ is updated to include all origin nodes previously associated with $v_n$. This aggregation process creates the final mapping $M : V \to \mathcal{P}(N)$, which bi-directionally links the entities in $G$ to the set of their structural locations (nodes) in $T$.

## 5 Agent-based Retrieval

Real-world document queries are often complex, necessitating operations like modal type filtering, semantic selection, and multi-hop reasoning. To address this, we propose an agent-based approach in BookRAG, which intelligently plans and executes operations on the BookIndex. We first introduce the overall workflow and present two core mechanisms: **Agent-based Planning**, which formulates the strategy, and the **Structured Execution**, which includes the retrieval process under the principles of IFT and generation.

### 5.1 Overall Workflow

The overall workflow of agent-based retrieval, illustrated in Figure 3, follows a three-stage pipeline designed to address users' queries systematically.

**1. Agent-based Planning.** BookRAG first performs *Classification & Plan*. This stage aims to distinguish simple keyword-based queries from reasoning questions that require decomposition and analysis. For instance, a query like "How does Transformer differ from RNNs in handling long-range dependencies?" cannot be solved by retrieving from a single keyword. Therefore, the planning stage first performs **query classification**. Based on this classification and a predefined set of operators designed for the BookIndex, it generates a specific **operators plan** that effectively guides the retrieval and generation strategies.

**2. Retrieval Process.** Guided by the operator plan, the retrieval process executes *Scent/Filter-based Retrieval*. This stage navigates the BookIndex $B = (T, G, M)$, either utilizing a **scent-based retrieval principle** (e.g., following relevant entities in $G$) to find information, or employing various filters (e.g., modal type) to refine the selection. After reasoning, BookRAG gets the retrieval set of highly relevant information blocks from the BookIndex.
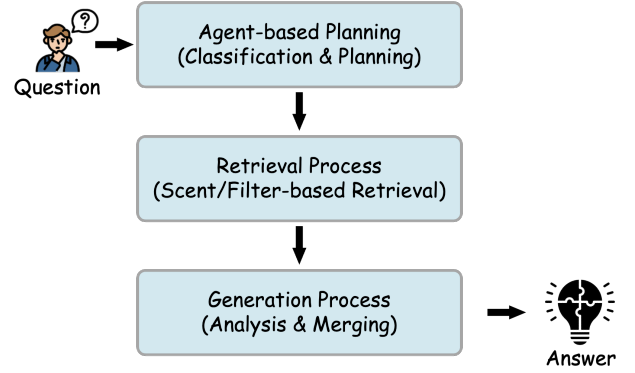


**Figure 3: The general workflow of agent-based retrieval in BookRAG, which contains agent-based planning, retrieval, and generation processes.**

**3. Generation Process.** Finally, all retrieved information enters the generation stage for *Analysis & Merging*. This stage synthesizes these (often fragmented) pieces of evidence, performs final analysis, and formulates a coherent response.

### 5.2 Agent-based Planning

The planning stage is the core of BookRAG, designed to intelligently navigate our BookIndex $B = (T, G, M)$. To support flexible retrieval, we define four types of operators: Formulator, Selector, Reasoner, and Synthesizer. These operators can be arbitrarily combined to form tailored execution pipelines, each with adjustable parameters. BookRAG dynamically configures and assembles these operators to adapt to the specific requirements of different query categories. This process involves two sequential steps: first, the agent performs *Query Classification* to determine the appropriate solution strategy, then generates a specific *Operator Plan*.

• **Query Classification**. To enable agent strategy selection, we focus on three representative query categories defined by their intrinsic complexity and operational demands (Table 2): *Single-hop*, *Multi-hop*, and *Global Aggregation*. This classification is crucial because each category requires a different solution strategy. For instance, a *Single-hop* query typically requires a single piece of information retrieved via a *Scent-based Retrieval* operation. In contrast, a *Global Aggregation* query often necessitates analyzing content under multiple filtering conditions, usually involving a sequence

**Table 2: Three common query categories addressed in BookRAG.**

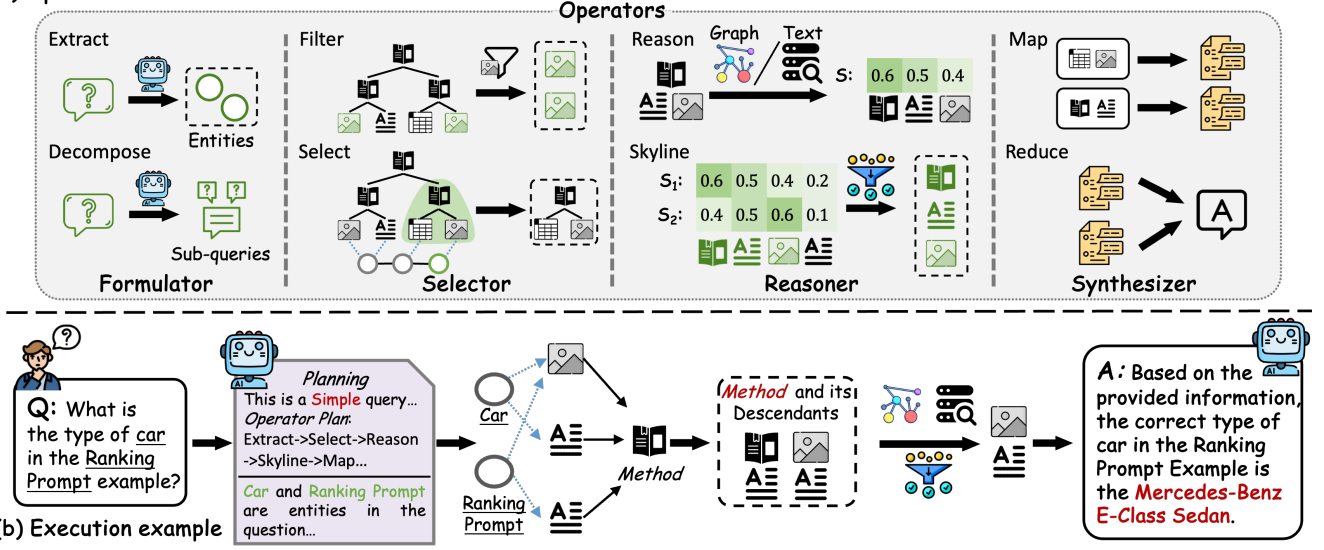| Query Category | Description | Core Task | Example Query |
|---|---|---|---|
| **Single-hop** | Queries with a single, distinct information target. | **Scent-based Retrieval**: Retrieve content related to a specific entity or section. | *What is the definition of Information Scent?* |
| **Multi-hop** | Queries that require synthesizing information from multiple blocks, often by decomposing into sub-problems. | **Decomposing & Merging**: Decompose into sub-problems, retrieve for each, and synthesize the final answer. | *How does Transformer differ from RNNs in handling long-range dependencies?* |
| **Global Aggregation** | Queries that require filtering across the entire document and performing calculations. | **Filter & Aggregation**: Apply filters across the document & perform aggregation operations (e.g., Count, Sum). | *How many figures related to IFT are in Section 4?* |

**Figure 4: The BookRAG Operator Library and an Execution Example from MMLongBench dataset: (a) a visual depiction of the four operator types (Formulator, Selector, Reasoner, and Synthesizer) and (b) an execution trace for a "Single-hop" query, demonstrating the agent-based planning and step-by-step operator execution.**

of *Filter & Aggregation* operations across various parts of the document. Furthermore, BookRAG is designed to be extensible, allowing for the resolution of a broader range of query types by integrating additional operators.

• **BookIndex Operators**. To execute the strategies identified by classification, we designed a set of operators ($O$) tailored for the BookIndex $B = (T, G, M)$. These operators, visually depicted in Figure 4(a) and detailed in Table 3, define the set of operations the agent can employ for diverse query categories. We group them into four types, which we describe in sequence:

❶ *Formulator.* These are LLM-based operators that prepare the query for execution. This category includes Decompose, which breaks a *Complex* query into a set of simpler, actionable sub-queries $Q_s$. It also includes Extract, which employs an LLM to identify key entities $E_q$ from the query text and link them to corresponding entities in the KG, $G$:

$$Q_s = \text{LLM}(P_{Dec}, q) = \{q_1, q_2, \ldots, q_k\} \quad (2)$$

$$E_q = \text{LLM}(P_{Ext}, q) = \{e_1, e_2, \ldots, e_m\} \quad (3)$$

Here, $q$ is the original user query, while $P_{Dec}$ and $P_{Ext}$ represent the prompts used to guide the LLM for the decomposition and extraction tasks, respectively.

❷ *Selector.* These operators filter or select specific content ranges from the BookIndex. Filter_Modal and Filter_Range directly apply the explicit constraints $C$ (e.g., modal types, page ranges) generated during the plan. Operating on the *Tree* $T = (N, E_T)$, these operators produce a filtered subset $N_f$ where the predicate $C(n)$ holds true for each node:

$$N_f = \{n \in N \mid C(n)\} \quad (4)$$

In contrast, Select_by_Entity and Select_by_Section target contiguous document segments by retrieving subtrees rooted at specific section nodes. This process first identifies a set of target section nodes $S_{\text{target}} \subset N$ at a specified depth, where $S_{\text{target}}$ consists of sections either linked to entities $E_q$ via the *GT-Link M* or selected by the LLM. It then retrieves all descendants of these targets to form the selected node set $N_s$:

$$N_s = \bigcup_{s \in S_{\text{target}}} \text{Subtree}(s) \quad (5)$$

❸ *Reasoner.* These operators analyze and refine selected tree nodes. Graph_Reasoning performs multi-hop inference on a subgraph $G'(V', E')$ (extracted from selected nodes $N_s$) starting from entity $e$. Starting from the retrieved entities, it computes an entity importance vector $I_G \in \mathbb{R}^{|V'|}$ over the subgraph $G'$ using the PageRank algorithm [20]. These entity scores are then mapped to the tree nodes via the GT-Link matrix $M$ to derive the final tree node importance scores vector $S_G \in \mathbb{R}^{|N_s|}$:

$$I_G = \text{PageRank}(G', e) \quad (6)$$

$$S_G = I_G \times M \quad (7)$$

Text_Ranker evaluates the semantic relevance of the tree node's content to the query $q$, assigning a relevance score $S_T$ to each node. Skyline_Ranker employs the Skyline operator to filter nodes based on these multiple criteria (e.g., $S_G$ and $S_T$), retaining only those nodes that are not dominated by any others in terms of the specified scoring dimensions.

❹ *Synthesizer.* These operators are responsible for content generation. Map performs analysis on specific retrieved information segments to generate partial responses. Reduce synthesizes a final

coherent answer by aggregating information from multiple sources, such as partial answers or a collection of retrieved evidence.

• **Operator Plan**. After classifying the query ($q$) into its category ($c$), the agent's final task is to generate an executable plan $P$. This plan is a specific sequence of operators $\langle o_1, \ldots, o_n \rangle$ selected from our library $O$ with parameters dynamically instantiated based on $q$. This process is formulated as:

$$P = \text{Agent}_{\text{Plan}}(q, c, O) \tag{8}$$

The plan follows a structured workflow tailored to each category:

- *Single-hop*: The agent first attempts to `Extract` an entity. If successful, it executes a "scent-based" selection; otherwise, it falls back to a section-based strategy. Both paths then proceed to standard reasoning and generation, denoted as $P_{\text{std}}$.

$$P_s = \begin{cases} \text{Extract} \xrightarrow{\text{success}} \text{Select\_by\_Entity} \to P_{\text{std}} \\ \text{Extract} \xrightarrow{\text{fail}} \text{Select\_by\_Section} \to P_{\text{std}} \end{cases} \tag{9}$$

$$P_{\text{std}} = (\text{Graph} \parallel \text{Text}) \to \text{Skyline} \to \text{Reduce} \tag{10}$$

- *Complex*: The agent first decomposes the problem, applies the Single-hop workflow $P_s$ to each sub-problem, and finally synthesizes the results.

$$P_{\text{complex}} = \text{Decompose} \to P_s \to \text{Map} \to \text{Reduce} \tag{11}$$

- *Global Aggregation*: The workflow involves applying a sequence of filters followed by synthesis.

$$P_{\text{global}} = \prod(\text{Filter\_Modal} \mid \text{Filter\_Range}) \to \text{Map} \to \text{Reduce} \tag{12}$$

Here, the symbol $\prod$ denotes the nested composition of filters, applying either a modal or range filter at each step.

## 5.3 Structured Execution

Following the planning stage, BookRAG executes the generated workflow $P$. This execution phase embodies the cognitive principles of Information Foraging Theory (IFT), effectively translating abstract textual queries into concrete operations. Specifically,

the `Selector` operators mirror the act of "navigating to information patches," narrowing the vast document space down to relevant scopes. Subsequently, the `Reasoner` operators perform "sense-making within patches," where they analyze and refine the information within these focused scopes. Finally, the `Synthesizer` generates the answer based on the processed evidence. This design minimizes the cost of attention by ensuring computational resources are focused solely on high-value data patches.

*Scent/Filter-based Retrieval.* The execution begins by narrowing the scope. Aligning with IFT, `Selector` operators identify relevant "patches" by following "information scents" (e.g., key entities in question) or applying explicit filter constraints. This process reduces the full node set $N$ to a focused node subset $N_s$:

$$N_s = \text{Selector}(N, \text{params}_{\text{sel}}) \tag{13}$$

This pre-selection minimizes noise and ensures that subsequent reasoning is applied only to highly relevant contexts, optimizing the foraging cost. Subsequently, within this focused scope, `Reasoner` operators evaluate nodes using multiple dimensions, such as graph topology and semantic relevance. We then employ the `Skyline_Ranker` to get the final retrieval set. Unlike fixed top-$k$ retrieval, the Skyline operator retains the *Pareto frontier* of nodes, retaining nodes that are valuable in at least one dimension while discarding dominated ones:

$$N_R = \text{Skyline\_Ranker}(\{S_G(n), S_T(n) \mid n \in N_s\}) \tag{14}$$

*Analysis & Merging Generation.* In the final stage, the `Synthesizer` operator generates the coherent answer by aggregating the refined evidence:

$$A = \text{Synthesizer}(q, N_R) \tag{15}$$

The `Map` operator performs fine-grained analysis on individual evidence blocks or sub-problems (from `Decompose`) to generate intermediate insights. The `Reduce` operator then aggregates these partial results, such as answers to decomposed sub-queries or statistical counts from a global filter, to construct the final response.

**Table 3: Operators utilized in our BookRAG, categorized by their function.**

| Operator | Type | Description | Parameters |
|---|---|---|---|
| Decompose | Formulator | Decompose a complex query into simpler, actionable sub-queries. | (Self-contained) |
| Extract | Formulator | Identify and extract key entities from the query (links to $G$). | (Self-contained) |
| Filter_Modal | Selector | Filter retrieved nodes by their modal type (e.g., Table, Figure). | modal_type: str |
| Filter_Range | Selector | Filter nodes based on a specified range (e.g., pages, section). | range: (start, end) |
| Select_by_Entity | Selector | Selects all tree nodes ($N$) in sections linked to a given entity ($V$). | entity_name: str |
| Select_by_Section | Selector | Uses an LLM to select relevant sections and selects all tree nodes ($N$) within them. | query: str, sections: List[str] |
| Graph_Reasoning | Reasoner | Performs multi-hop reasoning on subgraph ($G'$) and score tree nodes ($N$) using graph importance and GT-links. | start_entity: str, subgraph: $G'$ |
| Text_Reasoning | Reasoner | Rerank retrieved tree nodes ($N$) based on the relevance. | query: str |
| Skyline_Ranker | Reasoner | Rerank nodes based on multiple criteria. | criteria: List[str] |
| Map | Synthesizer | Uses partially retrieved information to generate a partial answer. | (Input: List[str]) |
| Reduce | Synthesizer | Synthesizes the final answer from partial information or all sub-problem answers. | (Input: List[str]) |

This separation ensures that the system can handle both detailed content extraction and high-level reasoning synthesis effectively.

To illustrate this end-to-end process, Figure 4(b) presents an execution trace for a "Single-hop" query: "What is the type of *car* in the *Ranking Prompt* example?". In the planning phase, the agent classifies the query and generates a specific workflow. Subsequently, it identifies key entities (e.g., "car") via Extract, retrieves relevant nodes via Select_by_Entity, refines them through reasoning and Skyline filtering, and finally synthesizes the answer using Reduce.

## 6 Experiments

In our experiments, we evaluate BookRAG against several strong baseline methods, with an in-depth comparison of their efficiency and accuracy on document QA tasks.

### 6.1 Setup

**Table 4: Datasets used in our experiments. EM and F1 denote Exact Match and F1-score, respectively.**

| Dataset | MMLongBench | M3DocVQA | Qasper |
|---|---|---|---|
| Questions | 669 | 633 | 640 |
| Documents | 85 | 500 | 192 |
| Avg. Pages | 42.16 | 8.52 | 10.95 |
| Avg. Images | 25.92 | 3.51 | 3.43 |
| Tokens | 2,816,155 | 3,553,774 | 2,265,349 |
| Metrics | EM, F1 | EM, F1 | Accuracy, F1 |

*Datasets & Question Synthesis.* We use three widely adopted benchmarking datasets for complex document QA tasks: MMLong-Bench [33], M3DocVQA [11], and Qasper [14]. MMLongBench is a comprehensive benchmark designed to evaluate QA capabilities on long-form documents, covering diverse categories such as guidebooks, financial reports, and industry files. M3DocVQA is an open-domain benchmark designed to test RAG systems on a diverse collection of HTML-type documents sourced from Wikipedia pages[1]. Qasper is a QA dataset focused on scientific papers, where questions require retrieving evidence from the entire document. We filtered the datasets to remove documents with low clarity or incoherent structures. To address the scarcity of global-level questions in the original benchmarks, we synthesize additional QA pairs by having an LLM generate global questions from selected document elements (e.g., tables or figures). These questions are then answered and meticulously refined by human annotators via an outsourcing process, with this additional QA pairs constituting less than 20% of our final QA pairs. The statistics of these datasets are presented in Table 4.

*Metrics.* We adhere to the official metrics specified by each dataset for QA. Our primary evaluation relies on Exact Match (EM), accuracy, and token-based F1-score. To assess efficiency, we also measure time cost and token usage during the response phase. Additionally, for methods including PDF parsing, we also evaluate retrieval recall. To establish the ground truth for this, we manually label the specific PDF blocks (e.g., texts, titles, tables, images, and

---

[1]https://www.wikipedia.org/

formulas) required to answer each question. This labeling process is guided by the metadata of ground-truth evidence provided in each dataset; we filter candidate blocks using the given modality (all datasets), page numbers (MMLongBench), and evidence statements (Qasper). Any blocks that remained non-unique after this filtering process are manually annotated. In cases where a PDF parsing error made the ground-truth item unavailable, the retrieval recall for that query is recorded as 0.

*Baselines.* Our experiments consider three model configurations:

- **Conventional RAG:** These methods are the most common pipeline for document analysis, where the raw text is first extracted and then chunked into segments of a specified size. We select strong and widely used retrieval models: BM25 [44] and Vanilla RAG. We also implement Layout+Vanilla, a variant that uses document layout analysis for semantic chunking.
- **Graph-based RAG:** These methods first extract textual content from documents and then leverage graph data during retrieval. We select RAPTOR [45] and GraphRAG [16]. Specifically, GraphRAG has two versions: GraphRAG-Global and GraphRAG-Local, which employ global and local search methods, respectively.
- **Layout segmented RAG:** This category encompasses methods that utilize layout analysis to segment document content into discrete structural units. We include: MM-Vanilla, which utilizes multi-modal embeddings for visual and textual content; a tree-based method inspired by PageIndex [39], denoted as TreeTraverse, where an LLM navigates the document's tree structure; DocETL [47], a declarative system for complex document processing; and GraphRanker, a graph-based method extended from HippoRAG [19] that applies Personalized PageRank [20] to rank the relevant nodes.

*Implementation details.* For a fair comparison, both BookRAG and all baseline methods are powered by a unified set of state-of-the-art (SOTA) and widely adopted backbone models from the Qwen family [4, 60, 63, 64]. We employ MinerU [52] for robust document layout parsing. We set the threshold of gradient $g$ as 0.6, and more details are provided in the appendix of our technical report [57]. Our source code, prompts, and detailed configurations are available at github.com/sam234990/BookRAG.

### 6.2 Overall results

In this section, we present a comprehensive evaluation of BookRAG, analyzing its complex QA performance, retrieval effectiveness, and query efficiency compared to state-of-the-art baselines.

• **QA Performance of BookRAG**. We compare the QA performance of BookRAG against three categories of baselines, as shown in Table 5. The results indicate that BookRAG achieves state-of-the-art performance across all datasets, substantially outperforming the top-performing baseline by 18.0% in Exact Match on M3DocVQA. Layout + Vanilla consistently outperforms Vanilla RAG, confirming that layout parsing preserves essential structural information for better retrieval. Besides, the suboptimal results of Tree-Traverse and GraphRanker highlight the limitations of relying solely on hierarchical navigation or graph-based reasoning, which

**Table 5: Performance comparison of different methods across various datasets for solving complex document QA tasks. The best and second-best results are marked in bold and underlined, respectively.**

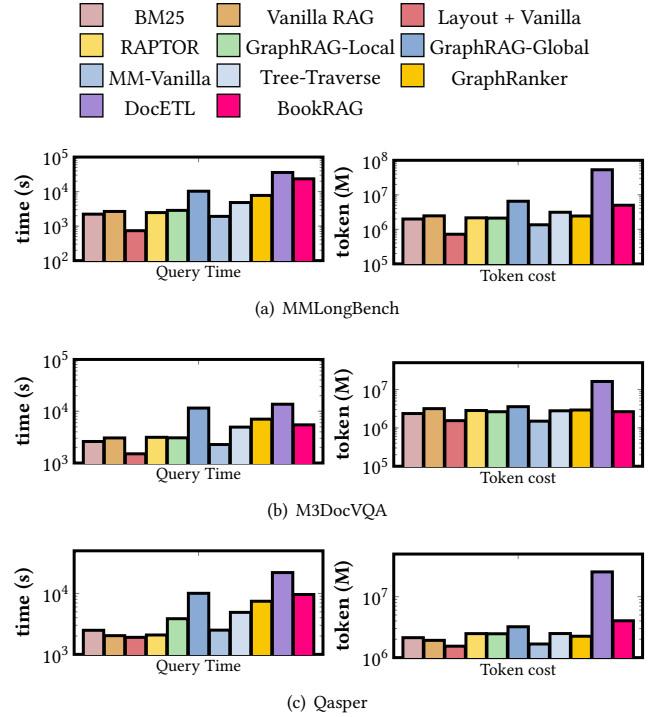| Baseline Type | Method | MMLongBench | | M3DocVQA | | Qasper | |
|---|---|---|---|---|---|---|---|
| | | (Exact Match) | (F1-score) | (Exact Match) | (F1-score) | (Accuracy) | (F1-score) |
| Conventional RAG | BM25 | 18.3 | 20.2 | 34.6 | 37.8 | 38.1 | 42.5 |
| | Vanilla RAG | 16.5 | 18.0 | 36.5 | 40.2 | 40.6 | 44.4 |
| | Layout + Vanilla | 18.1 | 19.8 | 36.9 | 40.2 | 40.7 | 44.6 |
| Graph-based RAG | RAPTOR | 21.3 | 21.8 | 34.3 | 37.3 | 39.4 | 44.1 |
| | GraphRAG-Local | 7.7 | 8.5 | 23.7 | 25.6 | 35.9 | 39.2 |
| | GraphRAG-Global | 5.3 | 5.6 | 20.2 | 22.0 | 24.0 | 24.1 |
| Layout segmented RAG | MM-Vanilla | 6.8 | 8.4 | 25.1 | 27.7 | 27.9 | 29.3 |
| | Tree-Traverse | 12.7 | 14.4 | 33.3 | 36.2 | 27.3 | 32.1 |
| | GraphRanker | 21.2 | 22.7 | <u>43.0</u> | <u>47.8</u> | 32.9 | 37.6 |
| | DocETL | <u>27.5</u> | <u>28.6</u> | 40.9 | 43.3 | <u>42.3</u> | <u>50.4</u> |
| **Our proposed** | **BookRAG** | **43.8** | **44.9** | **61.0** | **66.2** | **55.2** | **61.1** |

often miss cross-sectional context or drift into irrelevant scopes. In contrast, BookRAG's superiority stems from the synergy of its unified Tree-Graph BookIndex and Agent-based Planning. By effectively classifying queries and configuring optimal workflows, our BookRAG overcomes limitations of context fragmentation and static query workflow within existing baselines, ensuring precise evidence retrieval and accurate generation.

**Table 6: Retrieval recall comparison among layout-based methods. The best and second-best results are marked in bold and underlined, respectively.**

| Method | MMLongBench | M3DocVQA | Qasper |
|---|---|---|---|
| Layout + Vanilla | 26.3 | 33.8 | <u>33.5</u> |
| MM-Vanilla | 7.5 | 19.7 | 14.9 |
| Tree-Traverse | 11.2 | 19.5 | 14.5 |
| GraphRanker | <u>26.4</u> | <u>44.5</u> | 28.6 |
| **BookRAG** | **57.6** | **71.2** | **63.5** |

• **Retrieval performance of BookRAG.** To validate our retrieval design, we evaluate the retrieval recall of BookRAG against other layout-based baselines on the ground-truth layout blocks. The experimental results demonstrate that BookRAG achieves the highest recall across all datasets, notably reaching 71.2% on M3DocVQA and significantly outperforming the next best baseline (GraphRanker, max 44.5%). This performance advantage stems from our IFT-inspired **Selector → Reasoner** workflow: the Agent-based Planning first classifies the query, enabling the Selector to narrow the search to a precise *information patch*, followed by the Reasoner's analysis. Crucially, after the `Skyline_Ranker` process, the average number of retained nodes is 9.87, 6.86, and 8.6 across the three datasets, which is comparable to the standard top-$k$ ($k = 10$) setting, ensuring high-quality retrieval without inflating the candidate size.

• **Efficiency of BookRAG.** We further evaluate the efficiency in terms of query time and token consumption, as illustrated in Figure 5. Overall, BookRAG maintains time and token costs comparable to existing Graph-based RAG methods. While purely text-based

BM25　　Vanilla RAG　　Layout + Vanilla
RAPTOR　　GraphRAG-Local　　GraphRAG-Global
MM-Vanilla　　Tree-Traverse　　GraphRanker
DocETL　　BookRAG

(a) MMLongBench

(b) M3DocVQA

(c) Qasper

**Figure 5: Comparison of query efficiency.**

RAG approaches generally exhibit lower latency and token usage due to the absence of VLM processing for images, BookRAG maintains a balanced efficiency among multi-modal methods. In terms of token usage, BookRAG reduces consumption by an order of magnitude compared to the strongest baseline, DocETL. Notably, on the MMLongBench dataset, DocETL consumes over 53 million tokens, whereas BookRAG requires less than 5 million. Regarding the query latency, our method also achieves a speedup of up to 2× compared to DocETL.

## 6.3 Detailed Analysis

In this section, we provide a more in-depth examination of our BookRAG. We first conduct an ablation study to validate the contribution of each component, followed by an experiment on the impact of gradient-based ER and QA performance across different query types. Furthermore, we perform a comprehensive error analysis, compare the effectiveness of our entity resolution method, and present a case study.

• **Ablation study.** To evaluate the contribution of each core component in BookRAG, we design several variants by removing specific components:

- w/o Gradient ER: Replaces the gradient-based entity resolution with a Basic ER by merging the same-name entities.
- w/o Planning: Removes the Agent-based Planning, defaulting to a static, standard workflow for all queries.
- w/o Selector: Removes the Selector operators, forcing Reasoners to score all candidate nodes.
- w/o Graph_Reasoning: Removes the Graph_Reasoning operator. Consequently, the Skyline_Ranker is also disabled as scoring becomes single-dimensional.
- w/o Text_Reasoning: Removes the Text_Reasoning operator. Similarly, the Skyline_Ranker is disabled, relying solely on graph-based scores.

**Table 7: Comparing the QA performance of different variants of BookRAG. EM and F1 denote Exact Match and F1-score, respectively.**

| Method variants | MMLongBench | | Qasper | |
| --- | --- | --- | --- | --- |
| | EM | F1 | Accuracy | F1 |
| BookRAG (Full) | 43.8 | 44.9 | 55.2 | 61.1 |
| w/o gradient ER | 40.1 | 42.8 | 48.9 | 57.3 |
| w/o Planning | 30.8 | 33.2 | 40.9 | 48.5 |
| w/o Selector | 42.5 | 43.1 | 52.5 | 59.1 |
| w/o Graph_Reasoning | 39.8 | 41.5 | 51.4 | 58.4 |
| w/o Text_Reasoning | 39.0 | 40.3 | 47.2 | 52.5 |

The first variant evaluates the impact of KG quality on retrieval performance. The second and third variants assess the necessity of our Agent-based Planning and IFT-inspired selection mechanism, respectively. Finally, the last two variants validate the effectiveness of our multi-dimensional reasoning and dynamic Skyline filtering strategy. As shown in Table 7, the performance degradation across all variants confirms the essential role of each module in BookRAG. Specifically, the performance drop in the *w/o Gradient ER* variant highlights the critical role of a high-quality, connectivity-rich KG in supporting effective reasoning. Removing the *Planning* mechanism results in the most significant performance loss, confirming that a static workflow is insufficient for handling diverse types of queries. The *w/o Selector* variant, while maintaining competitive accuracy, incurs a prohibitive computational cost (> 2× tokens on Qasper), validating the efficiency of our IFT-inspired "narrow-then-reason" strategy.

• **Impact of Gradient-based Entity Resolution.** To evaluate the quality of our constructed KG, we compare the graph statistics
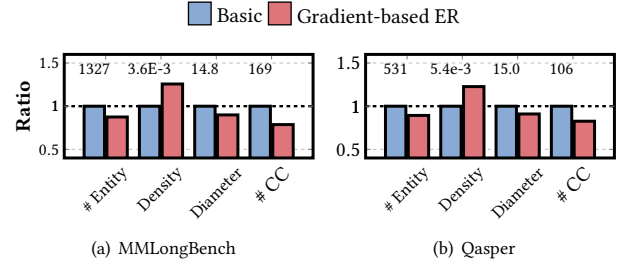


(a) MMLongBench  (b) Qasper

**Figure 6: Comparison of graph statistics. Values are normalized to the Basic setting (Baseline=1.0). Absolute values for Basic are annotated. Note that density values are abbreviated (e.g., 3.6E-3 denotes $3.6 \times 10^{-3}$).**

of our Gradient-based ER against a Basic KG construction. The Basic setting employs simple exact name matching for entity merging, which is standard practice in many graph-based methods. Figure 6 presents the comparative results, normalizing the metrics (Entity count, Density, Diameter of the Largest Connected Component, and Number of Connected Components) against the Basic baseline. The results demonstrate that our Gradient-based ER significantly optimizes KG. Specifically, it reduces the number of entities (by 12%) while substantially boosting graph density (by over 20% across datasets). This structural shift indicates that our ER module effectively identifies the same conceptual entities that possess different names. Consequently, the resulting graphs are more compact and cohesive, as evidenced by the reduced diameter and fewer connected components, which mitigates graph fragmentation and facilitates better connectivity for graph reasoning.


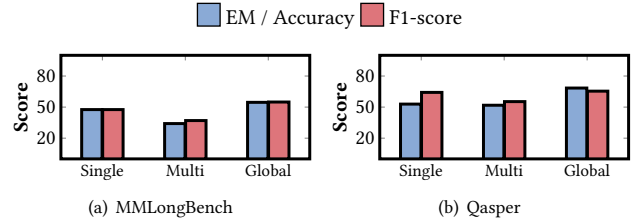
(a) MMLongBench  (b) Qasper

**Figure 7: QA performance breakdown by different query types (Single-hop, Multi-hop, and Global). The blue bars represent Exact Match (EM) for MMLongBench and Accuracy for Qasper, while the red bars represent the F1-score.**

• **QA performance under different query types.** Figure 7 breaks down the performance of BookRAG across Single-hop, Multi-hop, and Global aggregation query types. We observe that Multi-hop queries generally present a greater challenge compared to Single-hop ones, resulting in a slight performance decrease. This trend reflects the inherent difficulty of retrieving and reasoning over disjoint pieces of evidence. It further validates our agent-based planning strategy, which handles different query types separately.

• **Error Response analysis.** To diagnose the performance bottlenecks of BookRAG, we conduct a fine-grained error analysis on 200 sampled queries from each dataset, tracing the error propagation as shown in Figure 9. We categorize failures into four types: PDF Parsing, Plan, Retrieval, and Generation errors. The results

---

**BookRAG response of different query types**

**Single-hop Case from Qasper**

*Question:* What is the reward model for the reinforcement learning approach?
*Human-written answer:* Reward 1 for successfully completing the task, with a discount by the number of turns, and reward 0 when fail.
*Evidence:* We defined the reward as being 1 for successfully completing the task, and 0 otherwise. A discount of 0.95 </>

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Agent-based Planning:** This is a single-hop query. Here is the Select operator:

```
Extract={"entity_name": "reinforcement learning (rl)", "entity_type":"METHOD_OR_TECHNIQUE"}
```

**Process:** The `Select_by_Entity` operator identifies relevant sub-trees (e.g., "Introduction", "Related work"), pruning the reasoning space from 134 to 24 nodes. Subsequent `Graph_Reasoning`, `Text_Reasoning`, and `Skyline_Ranker` operate on this focused scope to retrieve 8 nodes for the final response.
**Output**: The reward model for the reinforcement learning approach is based on a **binary reward system that evaluates the success or failure of dialog interactions.** </> 1 for successfully completing a task. </> 0 for failing </> using a **discount factor of 0.95** </>

**Multi-hop Case from Qasper**

*Question:* What is the difference in performance between the interpretable system (e.g., vectors and cosine distance) and LSTM with ELMo system?
*Human-written answer:* Accuracy of best interpretible system was 0.3945 while accuracy of LSTM-ELMo net was 0.6818.
*Evidence 1:* The experimental results are presented in Table 1. Diacritic swapping showed a remarkably poor performance </>
*Evidence 2:* Test results for all the methods used. The loss measure is cross-entropy </>

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Agent-based Planning:** This is a multi-hop query. Here is the Decompose operator:

```
Sub_question 1: "What is the performance metric of the interpretable system (e.g., vectors and cosine distance)?"
Sub_question 2: "What is the performance metric of the LSTM with ELMo system?"
```

**Process:** BookRAG applies the single-hop workflow for each sub-question and generates the partial answers as follows:

```
Answer 1: "Summary:\n- Performance Metric: Accuracy (0.3945) and </>"
Answer 2: "The performance metric of the LSTM with ELMo system (LSTM-ELMo net) is as follows:\n Accuracy: 0.6818 </>"
```

**Output**: The interpretable system (using vector distance and cosine distance) **has an accuracy of 0.3945**, while the LSTM with ELMo system **has an accuracy of 0.6818**. The difference in accuracy between the two systems is 0.6818 - 0.3945 = 0.2873. </>

**Global Aggregation Case from MMLongBench**

*Question:* How many charts are shown in the first 10 pages of the document?    *Human-written answer:* 5

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Agent-based Planning:** This is a global query. Here are filter operators:

```
1. Filter_Range={"filter_type": page, "filter_range":"1-10"} 2. Filter_Modal={"filter_type": image}
```

**Process:** Applying filter operators retrieves images nodes from pages [3, 5, 6, 8, 9]. Map analyzes each image, and Reduce synthesizes the final output.
**Output**: Based on my analysis of the document, I found **5 items** that answer the question. </>

**Figure 8: Case study of responses across different query types from MMLongBench and Qasper. CYAN TEXT highlights correct content generated by BookRAG. GRAY TEXT describes the internal process, and </> marks omitted irrelevant parts.**

identify Retrieval Error as the dominant failure mode, followed by Generation Error, reflecting the persistent challenge of locating and synthesizing multimodal evidence. Regarding Plan Error, our qualitative analysis reveals a specific failure pattern: the planner tends to over-decompose detailed single-hop queries into unnecessary multi-hop sub-tasks. This fragmentation leads to disjointed retrieval paths, effectively preventing the model from synthesizing a cohesive final answer from the scattered sub-responses.

● **Case study.** Figure 8 illustrates BookRAG's answering workflow across Single-hop, Multi-hop, and Global queries. The results demonstrate that by leveraging specific operators (`Select`, `Decompose`, and `Filter`), BookRAG effectively prunes search spaces. For example, in the Single-hop case, the reasoning space is significantly reduced from 134 to 24 nodes. This capability allows the system to efficiently isolate relevant evidence from noise, ensuring precise answer generation.
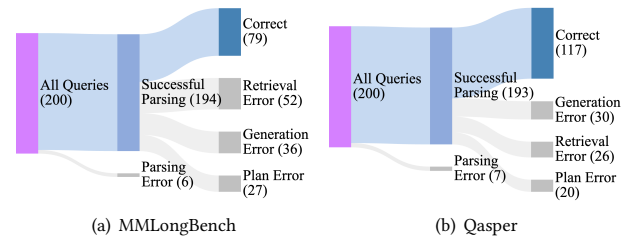


(a) MMLongBench          (b) Qasper

**Figure 9: Error analysis on 200 sampled queries from MMLongBench and Qasper datasets.**

## 7  Conclusion

In this paper, we propose BookRAG, a novel method built upon Book Index, a document-native, structured Tree-Graph index specifically designed to capture the intricate relations of structural documents. By employing an agent-based method to dynamically configure

retrieval and reasoning operators, our approach achieves state-of-the-art performance on multiple benchmarks, demonstrating significant superiority over existing baselines in both retrieval precision and answer accuracy. In the future, we will explore an integrated document-native database system that supports data formatting, knowledge extraction, and intelligent querying.

# References

[1] Simran Arora, Brandon Yang, Sabri Eyuboglu, Avanika Narayan, Andrew Hojel, Immanuel Trummer, and Christopher Ré. 2023. Language Models Enable Simple Systems for Generating Structured Views of Heterogeneous Data Lakes. *Proceedings of the VLDB Endowment* 17, 2 (2023), 92–105.

[2] Akari Asai, Zeqiu Wu, Yizhong Wang, et al. 2024. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *International Conference on Learning Representations (ICLR)*.

[3] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511* (2023).

[4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923* (2025).

[5] Camille Barboule, Benjamin Piwowarski, and Yoan Chabot. 2025. Survey on Question Answering over Visually Rich Documents: Methods, Challenges, and Trends. *arXiv preprint arXiv:2501.02235* (2025).

[6] Yukun Cao, Zengyi Gao, Zhiyang Li, Xike Xie, S. Kevin Zhou, and Jianliang Xu. 2025. LEGO-GraphRAG: Modularizing Graph-Based Retrieval-Augmented Generation for Design Space Exploration. *Proc. VLDB Endow.* 18, 10 (June 2025), 3269–3283. https://doi.org/10.14778/3748191.3748194

[7] Chengliang Chai, Jiajun Li, Yuhao Deng, Yuanhao Zhong, Ye Yuan, Guoren Wang, and Lei Cao. 2025. Doctopus: Budget-aware structural table extraction from unstructured documents. *Proceedings of the VLDB Endowment* 18, 11 (2025), 3695–3707.

[8] Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: The muppets straight out of law school. *arXiv preprint arXiv:2010.02559* (2020).

[9] Sibei Chen, Yeye He, Weiwei Cui, Ju Fan, Song Ge, Haidong Zhang, Dongmei Zhang, and Surajit Chaudhuri. 2024. Auto-Formula: Recommend Formulas in Spreadsheets using Contrastive Learning for Table Representations. *Proceedings of the ACM on Management of Data* 2, 3 (2024), 1–27.

[10] Sibei Chen, Nan Tang, Ju Fan, Xuemi Yan, Chengliang Chai, Guoliang Li, and Xiaoyong Du. 2023. Haipipe: Combining human-generated and machine-generated pipelines for data preparation. *Proceedings of the ACM on Management of Data* 1, 1 (2023), 1–26.

[11] Jaemin Cho, Debanjan Mahata, Ozan Irsoy, Yujie He, and Mohit Bansal. 2024. M3docrag: Multi-modal retrieval is what you need for multi-page multi-document understanding. *arXiv preprint arXiv:2411.04952* (2024).

[12] Vassilis Christophides, Vasilis Efthymiou, Themis Palpanas, George Papadakis, and Kostas Stefanidis. 2020. An overview of end-to-end entity resolution for big data. *ACM Computing Surveys (CSUR)* 53, 6 (2020), 1–42.

[13] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261* (2025).

[14] 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).

[15] Xavier Daull, Patrice Bellot, Emmanuel Bruno, Vincent Martin, and Elisabeth Murisasco. 2023. Complex QA and language models hybrid architectures, Survey. *arXiv preprint arXiv:2302.09051* (2023).

[16] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130* (2024).

[17] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* (2023).

[18] Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. 2024. LightRAG: Simple and Fast Retrieval-Augmented Generation. *arXiv e-prints* (2024), arXiv–2410.

[19] Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2024. HippoRAG: Neurobiologically Inspired Long-Term Memory for Large Language Models. *arXiv preprint arXiv:2405.14831* (2024).

[20] Taher H Haveliwala. 2002. Topic-sensitive pagerank. In *Proceedings of the 11th international conference on World Wide Web*. 517–526.

[21] Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *arXiv preprint arXiv:2402.07630* (2024).

[22] Yucheng Hu and Yuxing Lu. 2024. Rag and rau: A survey on retrieval-augmented language model in natural language processing. *arXiv preprint arXiv:2404.19543* (2024).

[23] Soyeong Jeong, Jinheon Baek, et al. 2024. Adaptive-RAG: Learning to Adapt Retrieval-Augmented Large Language Models through Question Complexity. *arXiv preprint arXiv:2403.14403* (2024).

[24] Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C Park. 2024. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity. *arXiv preprint arXiv:2403.14403* (2024).

[25] Tengjun Jin, Yuxuan Zhu, and Daniel Kang. 2025. ELT-Bench: An End-to-End Benchmark for Evaluating AI Agents on ELT Pipelines. *arXiv preprint arXiv:2504.04808* (2025).

[26] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In *European Conference on Computer Vision*. Springer, 498–517.

[27] Dawei Li, Shu Yang, Zhen Tan, Jae Young Baik, Sukwon Yun, Joseph Lee, Aaron Chacko, Bojian Hou, Duy Duong-Tran, Ying Ding, et al. 2024. DALK: Dynamic Co-Augmentation of LLMs and KG to answer Alzheimer's Disease Questions with Scientific Literature. *arXiv preprint arXiv:2405.04819* (2024).

[28] Guoliang Li, Jiayi Wang, Chenyang Zhang, and Jiannan Wang. 2025. Data+ AI: LLM4Data and Data4LLM. In *Companion of the 2025 International Conference on Management of Data*. 837–843.

[29] Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. 2023. Large language models in finance: A survey. In *Proceedings of the fourth ACM international conference on AI in finance*. 374–382.

[30] Zhaodonghui Li, Haitao Yuan, Huiming Wang, Gao Cong, and Lidong Bing. 2025. LLM-R2: A Large Language Model Enhanced Rule-based Rewrite System for Boosting Query Efficiency. *Proceedings of the VLDB Endowment* 1, 18 (2025), 53–65.

[31] Haoyu Lu, Wen Liu, Bo Zhang, et al. 2024. DeepSeek-VL: Towards Real-World Vision-Language Understanding. *arXiv preprint arXiv:2403.05525* (2024).

[32] Shengjie Ma, Chengjin Xu, Xuhui Jiang, Muzhi Li, Huaren Qu, Cehao Yang, Jiaxin Mao, and Jian Guo. 2024. Think-on-Graph 2.0: Deep and Faithful Large Language Model Reasoning with Knowledge-guided Retrieval Augmented Generation. *arXiv preprint arXiv:2407.10805* (2024).

[33] Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, et al. 2024. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. *Advances in Neural Information Processing Systems* 37 (2024), 95963–96010.

[34] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511* (2022).

[35] Zan Ahmad Naeem, Mohammad Shahmeer Ahmad, Mohamed Eltabakh, Mourad Ouzzani, and Nan Tang. 2024. RetClean: Retrieval-Based Data Cleaning Using LLMs and Data Lakes. *Proceedings of the VLDB Endowment* 17, 12 (2024), 4421–4424.

[36] Avanika Narayan, Ines Chami, Laurel Orr, and Christopher Ré. 2022. Can Foundation Models Wrangle Your Data? *Proceedings of the VLDB Endowment* 16, 4 (2022), 738–746.

[37] Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M Mulvey, H Vincent Poor, Qingsong Wen, and Stefan Zohren. 2024. A Survey of Large Language Models for Financial Applications: Progress, Prospects and Challenges. *arXiv preprint arXiv:2406.11903* (2024).

[38] Arash Dargahi Nobari and Davood Rafiei. 2024. TabulaX: Leveraging Large Language Models for Multi-Class Table Transformations. *arXiv preprint arXiv:2411.17110* (2024).

[39] PageIndex. 2025. PageIndex: Next-Generation Reasoning-based RAG. https://pageindex.ai/.

[40] Liana Patel, Siddharth Jha, Melissa Pan, Harshit Gupta, Parth Asawa, Carlos Guestrin, and Matei Zaharia. 2025. Semantic Operators and Their Optimization: Enabling LLM-Based Data Processing with Accuracy Guarantees in LOTUS. *Proceedings of the VLDB Endowment* 18, 11 (2025), 4171–4184.

[41] Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. 2024. Graph retrieval-augmented generation: A survey. *arXiv preprint arXiv:2408.08921* (2024).

[42] Peter Pirolli and Stuart Card. 1995. Information foraging in information access environments. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 51–58.

[43] Yichen Qian, Yongyi He, Rong Zhu, Jintao Huang, Zhijian Ma, Haibin Wang, Yaohua Wang, Xiuyu Sun, Defu Lian, Bolin Ding, et al. 2024. UniDM: A Unified Framework for Data Manipulation with Large Language Models. *Proceedings of Machine Learning and Systems* 6 (2024), 465–482.

[44] Stephen E Robertson and Steve Walker. 1994. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. In *SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University*. Springer, 232–241.

[45] Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D Manning. 2024. Raptor: Recursive abstractive processing for tree-organized retrieval. *arXiv preprint arXiv:2401.18059* (2024).

[46] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2024.

Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems* 36 (2024).

[47] Shreya Shankar, Tristan Chambers, Tarak Shah, Aditya G Parameswaran, and Eugene Wu. 2024. Docetl: Agentic query rewriting and evaluation for complex document processing. *arXiv preprint arXiv:2410.12189* (2024).

[48] Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. 2023. Improving the domain adaptation of retrieval augmented generation (RAG) models for open domain question answering. *Transactions of the Association for Computational Linguistics* 11 (2023), 1–17.

[49] Solutions Review Editors. 2019. 80 Percent of Your Data Will Be Unstructured in Five Years. https://solutionsreview.com/data-management/80-percent-of-your-datawill-be-unstructured-in-five-years/. Accessed: 2023-10-27.

[50] Zhaoyan Sun, Xuanhe Zhou, and Guoliang Li. 2024. R-Bot: An LLM-based Query Rewrite System. *arXiv preprint arXiv:2412.01661* (2024).

[51] Vincent A Traag, Ludo Waltman, and Nees Jan Van Eck. 2019. From Louvain to Leiden: guaranteeing well-connected communities. *Scientific reports* 9, 1 (2019), 1–12.

[52] 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).

[53] Jiayi Wang and Guoliang Li. 2025. Aop: Automated and interactive llm pipeline orchestration for answering complex queries. CIDR.

[54] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191* (2024).

[55] Shu Wang, Yixiang Fang, Yingli Zhou, Xilin Liu, and Yuchi Ma. 2025. ArchRAG: Attributed Community-based Hierarchical Retrieval-Augmented Generation. *arXiv preprint arXiv:2502.09891* (2025).

[56] Shen Wang, Tianlong Xu, Hang Li, Chaoli Zhang, Joleen Liang, Jiliang Tang, Philip S Yu, and Qingsong Wen. 2024. Large language models for education: A survey and outlook. *arXiv preprint arXiv:2403.18105* (2024).

[57] Shu Wang, Yingli Zhou, and Yixiang Fang. [n. d.]. BookRAG: A Hierarchical Structure-aware Index-based Approach for Complex Document Question Answering. https://github.com/sam234990/BookRAG.

[58] Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. 2024. Knowledge graph prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 19206–19214.

[59] Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective Retrieval Augmented Generation. *arXiv preprint arXiv:2401.15884* (2024).

[60] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025).

[61] Murong Yue. 2025. A survey of large language model agents for question answering. *arXiv preprint arXiv:2503.19213* (2025).

[62] Qinggang Zhang, Shengyuan Chen, Yuanchen Bei, Zheng Yuan, Huachi Zhou, Zijin Hong, Hao Chen, Yilin Xiao, Chuang Zhou, Junnan Dong, et al. 2025. A survey of graph retrieval-augmented generation for customized large language models. *arXiv preprint arXiv:2501.13958* (2025).

[63] Xin Zhang, Yanzhao Zhang, Wen Xie, Mingxin Li, Ziqi Dai, Dingkun Long, Pengjun Xie, Meishan Zhang, Wenjie Li, and Min Zhang. 2024. GME: Improving Universal Multimodal Retrieval by Multimodal LLMs. *arXiv preprint arXiv:2412.16855* (2024).

[64] 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).

[65] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* 1, 2 (2023).

[66] Yingli Zhou, Yaodong Su, Youran Sun, Shu Wang, Taotao Wang, Runyuan He, Yongwei Zhang, Sicong Liang, Xilin Liu, Yuchi Ma, et al. 2025. In-depth Analysis of Graph-based RAG in a Unified Framework. *arXiv preprint arXiv:2503.04338* (2025).

[67] Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zheng Liu, Zhicheng Dou, and Ji-Rong Wen. 2023. Large language models for information retrieval: A survey. *ACM Transactions on Information Systems* (2023).

# A   Experimental details

## A.1   Evaluation Metrics

In this section, we provide the detailed definitions and calculation procedures for the metrics used in our main experiments.

### A.1.1   Answer Extraction and Normalization.
Standard RAG models typically generate free-form natural language responses, which may contain extraneous conversational text (e.g., "The answer is..."). Directly comparing these raw outputs with concise ground truth labels (e.g., "Option A" or "12.5") can lead to false negatives.

Following official evaluation protocols, we employ an LLM-based extraction step to align the model output with the ground truth format before calculation. Let $y_{raw}$ denote the raw response generated by the RAG system and $y_{gold}$ denote the ground truth. We define the extracted answer $\hat{y}$ as:

$$\hat{y} = \text{LLM}_{\text{extract}}(y_{raw}, \text{Instruction}) \qquad (16)$$

where $\text{LLM}_{\text{extract}}$ extracts the key information (e.g., the key entity for span extraction) from $y_{raw}$. We further apply standard normalization $\mathcal{N}(\cdot)$ (e.g., lowercasing, removing punctuation) to both $\hat{y}$ and $y_{gold}$.

### A.1.2   QA Performance Metrics.
Based on the ground truth $y_{gold}$ and the model's response (either raw $y_{raw}$ or extracted $\hat{y}$), we compute the following metrics:

*Accuracy (Inclusion-based).* Following prior works [3, 34, 46], we utilize accuracy as a soft-match metric. We consider a prediction correct if the normalized gold answer is included in the model's generated response, rather than requiring a strict exact match. This accounts for the uncontrollable nature of LLM generation.

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\mathcal{N}(y_{gold,i}) \subseteq \mathcal{N}(y_{raw,i})) \qquad (17)$$

where $\subseteq$ denotes the substring inclusion relation.

*Exact Match (EM)..* Unlike accuracy, Exact Match is a strict metric. It measures whether the normalized *extracted* answer $\hat{y}$ is character-for-character identical to the ground truth.

$$\text{EM} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\mathcal{N}(\hat{y}_i) = \mathcal{N}(y_{gold,i})) \qquad (18)$$

*F1-score.* For questions requiring text span answers, we utilize the token-level F1-score between the extracted answer $\hat{y}$ and the ground truth $y_{gold}$. Treating them as bags of tokens $T_{\hat{y}}$ and $T_{gold}$:

$$P = \frac{|T_{\hat{y}} \cap T_{gold}|}{|T_{\hat{y}}|}, \quad R = \frac{|T_{\hat{y}} \cap T_{gold}|}{|T_{gold}|}, \quad \text{F1} = \frac{2 \cdot P \cdot R}{P + R} \qquad (19)$$

### A.1.3   Retrieval Recall.
As described in the main text, we evaluate retrieval quality based on the granularity of parsed PDF blocks (e.g., paragraphs, tables, images). For a given query $q$, let $\mathcal{B}_{gold}$ be the set of manually labeled ground-truth blocks required to answer $q$, and $\mathcal{B}_{ret}$ be the set of unique blocks retrieved by the system. The Retrieval Recall is defined as:

$$\text{Recall}_{ret} = \begin{cases} 0 & \text{if parsing error occurs on } \mathcal{B}_{gold} \\ \frac{|\mathcal{B}_{ret} \cap \mathcal{B}_{gold}|}{|\mathcal{B}_{gold}|} & \text{otherwise} \end{cases} \qquad (20)$$

Specifically, if a ground-truth block is lost due to PDF parsing failures (i.e., it does not exist in the candidate pool), it is considered strictly unretrievable, resulting in a recall contribution of 0 for that specific block.

## A.2   Implementation details

We implement BookRAG in Python, utilizing MinerU [52] for robust document layout parsing. For a fair comparison, both BookRAG and all baseline methods are powered by a unified set of state-of-the-art (SOTA) and widely adopted backbone models from the Qwen family [4, 60, 63, 64], including LLM, vision-language model (VLM), and embedding models. Specifically, we utilize Qwen3-8B [60] as the default LLM, Qwen2.5VL-30B [4] as the vision-language model (VLM), Qwen3-Embedding-0.6B [64] for text embedding, gme-Qwen2-VL-2B-Instruct [63] for multi-modal embedding, and Qwen3-Reranker-4B [64] for reranking. We primarily select models under the 10B parameter scale to balance efficiency and effectiveness. However, for the VLM, we adopt the 30B version, as the 8B counterpart exhibited significant performance deficits, frequently failing to answer correctly even when provided with ground-truth images. All experiments were conducted on a Linux operating system running on a high-performance server equipped with an Intel Xeon 2.0GHz CPU, 1024GB of memory, and 8 NVIDIA GeForce RTX A5000 GPUs, each with 24 GB of VRAM. Specifically, to ensure a fair comparison of efficiency, all methods were executed serially, and the reported time costs reflect this sequential processing mode. For methods involving document chunking and retrieval ranking, we standardize the chunk size at 500 tokens and set the retrieval top-$k$ to 10 to ensure consistent candidate pool sizes across baselines. For further reproducibility, our source code and detailed implementation configurations are publicly available at our repository: https://github.com/sam234990/BookRAG.

## A.3   Prompts

Specifically, we present the prompts designed for agent-based query classification (Figure 10), question decomposition (Figure 11), and filter operator generation (Figure 12). Additionally, we illustrate the prompt employed for entity resolution judgment (Figure 13) during the graph construction phase.

```
You are an expert query analyzer. Your only task is to classify the user's question into one of three categories: "simple",
    "complex", or "global". Respond only with the specified JSON object.

Category Definitions:
1. single-hop: The question can be fully answered by retrieving information from a SINGLE, contiguous location in the
    document (e.g., one specific paragraph, one complete table, or one figure).
  - This includes questions that require reasoning or comparison, as long as all the necessary data is present within
      that single retrieved location.
  - Example: "What is the title of Figure 2?"
  - Example: "How do 5% of the Latinos see economic upward mobility for their children?" -> This is SIMPLE because the
      answer can be found by looking at a single chart or paragraph.

2. multi-hop: The question requires decomposition into multiple simple sub-questions, where each sub-question must be
    answered by a separate retrieval action.
  - It often contains a nested or indirect constraint that requires a preliminary step to resolve before the main
      question can be answered.
  - Example: "What is the color of the personality vector...?" -> This is COMPLEX because it requires two separate
      retrieval actions.

3. global: The question requires an aggregation operation (e.g., counting, listing, summarizing) over a set of items that
    are identified by a clear structural filter.
  - Example: "How many tables are in the document?" -> This is GLOBAL because the process is to filter for all items of
      type 'table'.

User Query: query
```

**Figure 10: The prompt for query classification.**

```
You are a query decomposition expert. You have been given a "complex" question. Your task is to break it down into a series
    of simple, atomic sub-questions and classify each one by type.

**Crucial Instructions:**
1. Each `retrieval` sub-question MUST be a direct information retrieval task that can be answered independently by looking
    up a specific fact, number, or value in the document.
2. **`retrieval` sub-questions MUST NOT depend on the answer of another sub-question.** They should be parallelizable. All
    logic for combining their results must be placed in a final `synthesis` question.
3. A `synthesis` question requires comparing, calculating, or combining the answers of the previous `retrieval` questions.
    It does **NOT** require a new lookup in the document.

You MUST provide your response in a JSON object with a single key 'sub_questions', which contains a list of objects. Each
    object must have a 'question' (string) and a 'type' (string: "retrieval" or "synthesis").

--- EXAMPLE 1 (Correct Decomposition with Independent Lookups) ---
Complex Query: "What is the color of the personality vector in the soft-labled personality embedding matrix that with the
    highest Receptiviti score for User A2GBIFL43U1LKJ?"

Expected JSON Output:
{{
  "sub_questions": [
    {{"question": "What are all the Receptiviti scores for each personality vector for User A2GBIFL43U1LKJ?",
      "type": "retrieval"}},
    {{"question": "What is the mapping of personality vectors to their colors in the soft-labled personality embedding
        matrix?",
      "type": "retrieval"}},
    {{"question": "From the gathered scores, identify the personality vector with the highest score, and then find its
        corresponding color from the vector-to-color mapping.",
      "type": "synthesis"}}
  ]
}}
--- END EXAMPLE 1 ---

--- EXAMPLE 2 (Decomposition with retrieval and synthesis steps) ---
Complex Query: "According to the report, which one is greater in population in the survey? Foreign born Latinos, or the
    Latinos interviewed by cellphone?"

Expected JSON Output:
{{
  "sub_questions": [
    {{"question": "According to the report, what is the population of foreign born Latinos in the survey?",
      "type": "retrieval"}},
    {{"question": "According to the report, what is the population of Latinos interviewed by cellphone in the survey?",
      "type": "retrieval"}},
    {{"question": "Which of the two population counts is greater?",
      "type": "synthesis"}}
  ]
}}
--- END EXAMPLE 2 ---

Now, perform the decomposition for the following query.

User Query: **query**
```

**Figure 11: The prompt for query decomposition.**

```
You are a highly specialized AI assistant. Your only function is to analyze a "Global Query" and return a single, valid
    JSON object that specifies both the filtering steps and the final aggregation operation. You MUST NOT output any other
    text or explanation.

### INSTRUCTIONS \& DEFINITIONS ###

1. **Filters**: You MUST determine the list of `filters` to apply. Even if the filter is for the whole document (e.g., all
    tables), the `filters` list must be present.
   - `filter_type`: One of ["section", "image", "table", "page"].
      - `section`: Use for structural parts like chapters, sections, appendices, or references.
      - `image`: Use for visual elements like figures, images, pictures, or plots.
      - `table`: Use for tabular data.
      - `page`: Use for specific page numbers or ranges.
   - `filter_value`: (Optional) Can be provided for "section" (e.g., a section title) or "page" (e.g., '3-10' or '5').
        **For "image" or "table", this value MUST be null.**

2. **Operation**: Determine the final aggregation operation.
   - `operation`: One of ["COUNT", "LIST", "SUMMARIZE", "ANALYZE"].

### EXAMPLES OF YOUR TASK ###

User: "How many figures are in this paper from Page 3 to Page 10?"
Assistant: {{"filters": [{{"filter_type": "page", "filter_value": "3-10"}}, {{"filter_type": "image"}}], "operation":
    "COUNT"}}

User: "Summarize the discussion about 'data augmentation' in the 'Methodology' section."
Assistant: {{"filters": [{{"filter_type": "section", "filter_value": "Methodology"}}], "operation": "SUMMARIZE"}}

User: "How many chapters are in this report?"
Assistant: {{"filters": [{{"filter_type": "section"}}], "operation": "COUNT"}}

### YOUR CURRENT TASK ###

User: "{query}"
User Query: query
```

**Figure 12: The prompt for Filter operator generation.**

```
-Goal-
You are an expert Entity Resolution Adjudicator. Your task is to determine if a "New Entity" refers to the exact same
    real-world concept as one of the "Candidate Entities" provided from a knowledge graph. Your output must be a JSON
    object containing the ID of the matching candidate (or -1) and a brief explanation for your decision.
-Context-
You will be given one "New Entity" recently extracted from a text. You will also be given a list of "Candidate Entities"
    that are semantically similar, retrieved from an existing knowledge base. Each candidate has a unique `id` for you to
    reference.
---
-Core Task & Rules-
1. **Analyze the "New Entity"**: Carefully read its name, type, and description to understand what it is.
2. **Field-by-Field Adjudication**: To determine a match, you must evaluate each field with a specific focus:
    * **`entity_name` (High Importance):** The names must be extremely similar, a direct abbreviation (e.g., "LLM" vs.
        "Large Language Model"), or a well-known alias. **If the names represent distinct, parallel concepts (like "Event
        Detection" and "Named Entity Recognition"), they are NOT a match, even if their descriptions are very similar.**
    * **`entity_type` (Medium Importance):** The types do not need to be identical, but they must be closely related and
        compatible (e.g., `COMPANY` and `ORGANIZATION` could describe the same entity).
    * **`description` (Contextual Importance):** The descriptions may differ as they are often extracted from different
        parts of a document. Your task is to look past surface-level text similarity and determine if they fundamentally
        describe the **same underlying object or concept**.

3. **Be Strict and Conservative**: Your standard for a match must be very high. An incorrect merge can corrupt the
    knowledge graph. A missed merge is less harmful.
    * Surface-level similarities are not enough. The underlying concepts must be identical.
    * For example, "Apple" (the fruit) and "Apple Inc." (the company) are NOT a match.
    * **When in doubt, you MUST output -1.**
    * **Assume No Match by Default**: In a large knowledge graph, most new entities are genuinely new. You should start
        with the assumption that the "New Entity" is unique. You must find **strong, convincing evidence** across all
        fields, especially the `entity_name`, to overturn this assumption and declare a match.

4. **Format the Output**: **You must provide your answer in a valid JSON format. The JSON object should contain two keys:**
    * `select_id`: An integer. The `id` of the candidate you've determined to be an exact match. If no exact match is
        found, this value MUST be `-1`.
    * `explanation`: A brief, one-sentence string explaining your reasoning. For a match, explain why they are the same
        entity. For no match, explain the key difference.
---
-Output Schema & Format-
Your response MUST be a single, valid JSON object that adheres to the following schema. Do not include any other text,
    explanation, or markdown formatting like ```json.

```json
{{
  "select_id": "integer",
  "explanation": "string"
}}
```
---
-Example-
### Example 1: Match Found
### Example 2: No Match Found
----
-Task Execution-

Now, perform the selection task based on the following data. Remember to output only a single integer.

- Input Data -
```

**Figure 13: The prompt for entity resolution judgement, examples are omitted due to lack of space.**