

# Context-Aware Scientific Knowledge Extraction on Linked Open Data using Large Language Models

Sajratul Yakin Rubaiat, *Member, IEEE*, Hasan M Jamil, *Member, IEEE*,

**Abstract**—The exponential growth of scientific literature presents a significant challenge for researchers seeking to extract and synthesize relevant knowledge. Traditional search engines often return a large number of sources without directly providing detailed answers, while general-purpose Large Language Models (LLMs) may offer concise responses that lack depth or fail to incorporate the most up-to-date information. Furthermore, LLMs with search capabilities are often limited by their context window, resulting in short, incomplete answers. This paper introduces WISE (Workflow for Intelligent Scientific Knowledge Extraction), a novel system that addresses these limitations by combining LLMs with a structured, multi-layered workflow to extract, refine, and rank scientific knowledge tailored to specific queries. WISE employs an LLM-powered, tree-based architecture with a customized search function to iteratively refine extracted data, focusing on query-aligned and context-aware information while actively avoiding redundancy. Dynamic scoring and ranking mechanisms prioritize unique contributions from each source, and adaptive stopping criteria minimize processing overhead. WISE delivers detailed, well-organized, and highly informative answers by systematically exploring and synthesizing knowledge from diverse sources. Experiments focused on biological queries related to *HBB* gene-associated diseases demonstrate that WISE reduces the volume of processed text by over 80% while simultaneously achieving significantly higher recall compared to baseline methods, including leading search engines and other LLM-based approaches. Further analysis using ROUGE and BLEU metrics reveals that WISE’s output is more unique compared to other systems, and a novel level-based evaluation metric shows that WISE provides more in-depth information. This paper also explores how the WISE workflow can be adapted as a general framework for diverse research domains, such as *drug discovery*, *material science*, and *social science*, enabling efficient knowledge extraction and synthesis from unstructured scientific papers and web sources across a wide array of research domains.

**Index Terms**—Knowledge Discovery, Large Language Models, Information Filtering, Scientific Data Extraction, Knowledge Enrichment, Gene-Disease Associations

## I. INTRODUCTION

THE relentless expansion of scientific knowledge, reflected in the ever-increasing volume of published literature, presents a formidable challenge for researchers seeking to extract, synthesize, and contextualize relevant information [1], [2]. While traditional search engines and general-purpose

Large Language Models (LLMs) offer some assistance, they often fall short in providing domain-specific insights, filtering irrelevant content, and efficiently managing the sheer volume of data [3]–[5]. Traditional search engines typically return a large number of sources, requiring users to manually sift through them to extract relevant information, rather than providing direct, synthesized answers. Navigating this complex information ecosystem manually is both labor-intensive and error-prone, with researchers facing the risk of overlooking critical details as the volume of data grows exponentially.

Consider, for example, a researcher investigating gene-disease associations related to the *HBB* gene. Starting from a single authoritative source like the HGNC [6], they might encounter 24 relevant sources. Exploring just one of these, such as ClinVar [7], could unveil hundreds more sources, leading to a rapidly expanding tree of interconnected resources, exemplified by platforms like NCBI [8]. This exponential growth of linked resources quickly overwhelms traditional search systems, making it difficult to extract pertinent insights efficiently. Even experienced investigators struggle to filter out superfluous data and focus on the most valuable information, a challenge amplified for newcomers. Consequently, critical information may be missed, and the time required to gain a complete, integrated understanding escalates dramatically.

Purely automated approaches also encounter significant difficulties in this context [9]–[11]. The sheer volume of interconnected data can lead to computationally expensive and strategically ineffective processes without robust mechanisms for pruning irrelevant or redundant content. A key challenge lies in determining when to stop searching; continued exploration without clear stopping criteria often yields diminishing returns, underscoring the need for a balanced workflow that ensures thorough yet efficient exploration while prioritizing high-value information [12], [13].

While Large Language Models (LLMs) show promise in specific aspects of information retrieval [14], their inherent limitations hinder their ability to fully address the scale of these challenges. General-purpose LLMs, often provide concise answers that lack the depth and detail required for complex scientific queries, and may not incorporate the most recent findings. For instance, state-of-the-art models like GPT-4o [15], despite having a context window of 128000 tokens, are practically limited to processing data from only about eight sources simultaneously, such as UniProt [16], as illustrated in Figure 1. Although capable of ranking and comparing content within this limited scope, LLMs are constrained by their

S. Y. Rubaiat is with the Department of Computer Science, University of Idaho, Moscow, ID, USA (e-mail: ruba3062@vandals.uidaho.edu).

H. M. Jamil is with the Department of Computer Science, University of Idaho, Moscow, ID, USA (e-mail: jamil@uidaho.edu).

Manuscript received January 1, 2025; revised January 1, 2025.

TABLE I  
SYSTEM COMPARISON

Feature	WISE	ChatGPT	ChatGPT with Search	Gemini	Google Search
Number of Diseases Identified	16	9	7	2	3
Recall	0.84	0.47	0.36	0.10	0.15
Average Level of Detail <sup>a</sup>	3.8	3.33	3.42	2.5	3.0
Structured Output	✓	✓	✓	✓	✗
Inclusion of Sub-variations	✓	✗	✓	✗	✗
Source Citation	✓	✗	✓	✗	✓
Identification of Rare Conditions	✓	✗	✗	✗	✗
Up-to-Date Information	✓	✗	✓	✗	✓

<sup>a</sup> Average Level of Detail: Represents the average depth of information provided across all diseases identified by the system, based on a 0-5 scale where higher values indicate more detailed information (see Section III for level criteria).

narrow context window, further reduced by factors like search history. In specialized domains such as biology and medicine, where nuanced and detailed insights are crucial, relying solely on LLM-based searches proves insufficient. These challenges highlight the critical need for combining the capabilities of LLMs with strategic workflows to achieve comprehensive and efficient knowledge retrieval.

To address these challenges, we introduce **WISE (Workflow for Intelligent Scientific Extraction)**, a novel, scalable, tree-based framework that integrates LLM-driven filtering, dynamic ranking, and adaptive stopping criteria. WISE is designed to deliver detailed, well-organized, and highly informative answers by systematically exploring and synthesizing knowledge from diverse sources. WISE mirrors the approach of a diligent researcher: identifying relevant information, discarding duplicates, exploring promising leads, and recognizing when further pursuit yields diminishing returns. WISE begins by employing LLMs to filter large text corpora based on a user’s domain-specific query, ensuring that only contextually relevant and manageable segments proceed to subsequent stages. It then assigns scores to extracted sources, quantifying their unique contributions relative to previously processed material, thereby minimizing redundancy and focusing on content that adds genuine value. This dynamic ranking process prioritizes high-value information while pruning low-impact paths. The iterative refinement continues layer by layer, with WISE’s knowledge container—the growing repository of extracted, query-specific insights—expanding until incremental findings diminish (Figure 4). At this point, WISE intelligently halts further searches, conserving computational resources while delivering comprehensive, contextually nuanced results. In essence, WISE achieves a balance between breadth and depth, effectively leveraging the strengths of LLMs while mitigating their limitations through intelligent pruning and carefully considered stopping criteria.

Our key **contributions** are summarized as follows:

- 1) **Scalable, Tree-Based Architecture:** We introduce a novel, tree-structured workflow (Section II) that efficiently navigates large, heterogeneous datasets. This architecture leverages LLM-based filtering at each layer to incrementally refine data subsets according to domain-specific queries, ensuring scalability and focus.
- 2) **Dynamic Ranking and Pruning:** We present a transparent scoring mechanism (Sections II and III) that quantifies each source’s unique knowledge contribution.

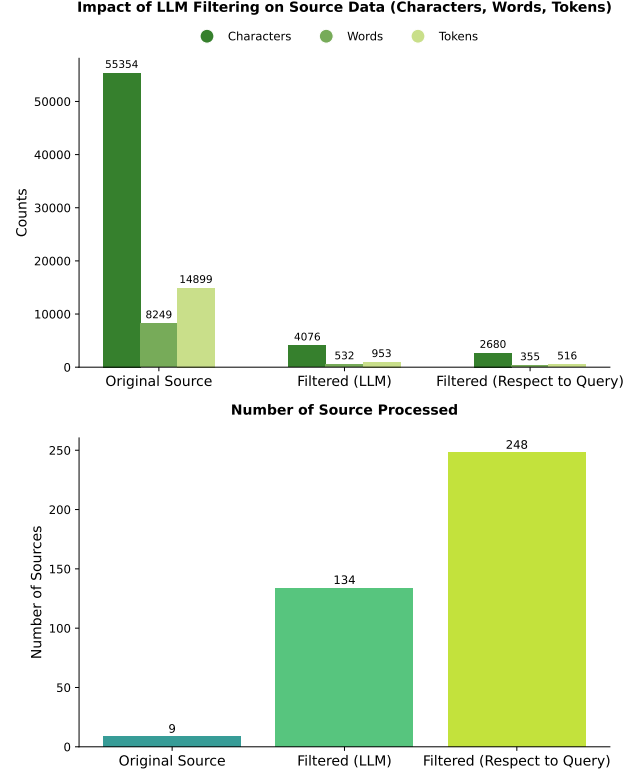


Fig. 1. Number of sources, such as UniProt [16], that can be processed simultaneously for ranking by advanced LLMs like GPT-4o demonstrate significant improvements when applying LLM-based filtering with and without query-specific relevance. Given that GPT-4o supports a 128,000-token context window, the number of websites that can be processed is calculated as: Number of websites = 128,000 / Tokens per Website

By dynamically ranking sources based on their added value and employing intelligent pruning, WISE focuses computational resources on the most promising leads, effectively filtering out redundant or low-value content.

- 3) **Adaptive, Expert-Inspired Exploration:** Our approach (Sections II and III) mirrors expert-driven inquiry by progressively deepening the search along promising paths while adaptively halting exploration when further gains are minimal. This ensures a balanced blend of breadth and depth, optimizing both the efficiency and effectiveness of the knowledge discovery process.

- 4) **Demonstrated Effectiveness in Gene-Disease Association Discovery:** Through empirical evaluation on gene-disease association queries (Section IV), we demonstrate that WISE significantly outperforms baseline methods, including traditional search engines and general-purpose LLMs, in terms of recall, uniqueness of extracted information (ROUGE/BLEU), and depth of knowledge (level-based analysis).
- 5) **Versatile Applications:** We showcase WISE’s adaptability and potential impact through diverse applications (Section V), including drug discovery, material science, and social science, highlighting its ability to generalize across a wide range of research domains.

## II. SYSTEM DESIGN

The WISE framework is designed to streamline the extraction and synthesis of knowledge from unstructured data sources through a structured and multi-layered approach. As illustrated in Figure 2, its architecture comprises four key stages that work in tandem to filter, score, rank, and consolidate information. Each stage contributes to transforming raw data into context-aware insights, enabling efficient knowledge discovery and refinement.

- 1) **Content Filtering:** This stage employs query-specific extraction via LLM-driven contextual analysis, ensuring that only information relevant to the query is retained while noise, such as advertisements, is removed. This significantly reduces computational overhead in subsequent stages.
- 2) **Score Calculation:** In this stage, the filtered content is evaluated for its unique contribution to the evolving knowledge container. Novel and relevant insights are prioritized, while redundant material is discarded.
- 3) **Threshold Checking:** This component determines whether continued exploration of sources is justified. It acts as a termination criterion for the recursive process, halting when additional contributions fall below a defined threshold.
- 4) **Knowledge Consolidation:** Extracted information is incrementally merged into a growing repository of domain-specific knowledge. This ensures that the final knowledge container is comprehensive, context-aware, and aligned with the user’s query.

WISE initiates its process with a user-provided query  $q$  and an empty knowledge container  $\mathcal{K}_0$ . This container evolves iteratively as new insights are integrated. The initial set of sources  $\mathcal{S}_0 = \{s_1, s_2, \dots, s_n\}$  is retrieved using a traditional similarity-based search function  $\Phi(q)$ , which identifies a collection of candidate sources relevant to the query. These sources serve as the root nodes for the tree-based search process. Formally:

$$\mathcal{K}_0 = \emptyset, \quad \mathcal{S}_0 = \Phi(q)$$

### A. Content Filtering

Each source  $s_i$  belonging to the set  $\mathcal{S}_l$  at layer  $l$  undergoes a query-specific refinement process. This process extracts

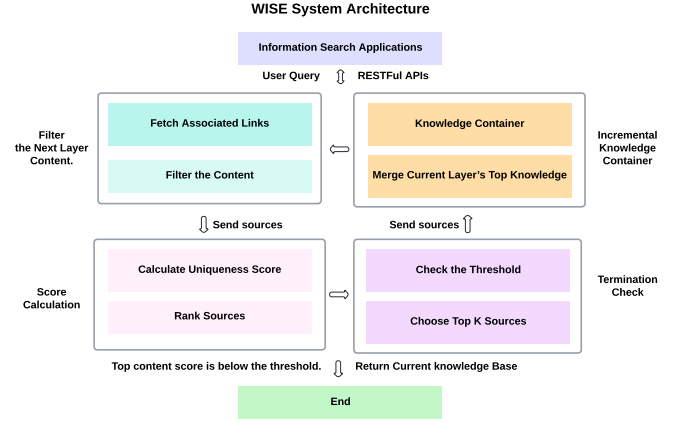


Fig. 2. WISE System Architecture, showcasing its four main components: content filtering, score calculation, threshold checking, and knowledge consolidation.

content directly relevant to the query  $q$ . The filtering function  $\Gamma$ , which leverages the contextual understanding capabilities of an LLM, transforms the raw content  $\mathcal{C}(s_i)$  of a source  $s_i$  into a focused subset  $\mathcal{F}(s_i)$ :

$$\mathcal{F}(s_i) = \Gamma(q, \mathcal{C}(s_i))$$

By isolating only the most pertinent information,  $\mathcal{F}(s_i)$  significantly reduces noise and irrelevant data. This streamlining enhances the efficiency of subsequent computational tasks and downstream processing stages. For simplicity, the result of the filtering operation for all sources at layer  $l$  is denoted as:

$$\mathcal{F}_l = \{\mathcal{F}(s_1), \mathcal{F}(s_2), \dots, \mathcal{F}(s_n)\}_l$$

Here,  $\mathcal{F}_l$  represents the set of all query-specific filtered content derived at layer  $l$  from the sources in  $\mathcal{S}_l$ .

### B. Score Calculation

Following the filtering stage, WISE quantifies the unique contribution of each source to the knowledge container. Let  $w_{\text{filtered}}(s_i)$  denote the number of words in the filtered content of each source  $s_i$ :

$$w_{\text{filtered}}(s_i) = |\mathcal{F}(s_i)|$$

where  $|\mathcal{F}(s_i)|$  represents the cardinality (number of elements) of the set  $\mathcal{F}(s_i)$ .

Next, we determine the number of words that overlap between the source’s filtered content  $\mathcal{F}(s_i)$  and the current knowledge container  $\mathcal{K}_l$ . Let  $w_{\text{overlap}}(s_i, \mathcal{K}_l)$  denote this count:

$$w_{\text{overlap}}(s_i, \mathcal{K}_l) = |\mathcal{F}(s_i) \cap \mathcal{K}_l|$$

Here,  $|\mathcal{F}(s_i) \cap \mathcal{K}_l|$  represents the cardinality of the intersection of the two sets.

The unique knowledge contribution  $\mathcal{K}(s_i)$  of source  $s_i$  is then defined as the difference between the number of words in the filtered content and the number of overlapping words:

$$\mathcal{K}(s_i) = w_{\text{filtered}}(s_i) - w_{\text{overlap}}(s_i, \mathcal{K}_l)$$

This value,  $\mathcal{K}(s_i)$ , represents the number of new, unique words that source  $s_i$  contributes to the knowledge container.

To normalize and evaluate the contribution of each source, we define the following metrics:

**1. Knowledge Density (Per-Word Normalization):** This metric normalizes the unique knowledge contribution by the size of the source (measured in the number of words in the filtered content). It is calculated as:

$$\text{Knowledge Density}(s_i) = \frac{\mathcal{K}(s_i)}{w_{\text{filtered}}(s_i)}$$

This ratio represents the proportion of unique words in the filtered content of source  $s_i$ .

**2. Knowledge Increase (Relative Growth):** This metric measures the relative contribution of the source to the existing knowledge container, expressed as a proportion of the current size of the knowledge container:

$$\text{Knowledge Increase}(s_i) = \frac{\mathcal{K}(s_i)}{|\mathcal{K}_l|}$$

Here,  $|\mathcal{K}_l|$  denotes the cardinality of the knowledge container  $\mathcal{K}_l$ , representing the total number of words in the knowledge container at layer  $l$ .

To integrate the concepts of local efficiency (size of the source) and global contribution (size of the knowledge container), we define a unified scoring function  $\Psi$  that employs log scaling to balance these factors:

**3. Combined Normalized Metric:**

$$\text{Score}(s_i) = \Psi(\mathcal{F}(s_i), \mathcal{K}_l) = \frac{\mathcal{K}(s_i)}{\log(1 + w_{\text{filtered}}(s_i) + |\mathcal{K}_l|)}$$

This combined metric prioritizes sources that offer unique and meaningful contributions, accounting for both the relative size of the source (in terms of the number of words in its filtered content) and its impact on the evolving knowledge container (in terms of the number of unique words it contributes).

### C. Threshold Checking and Pruning

WISE evaluates whether continued exploration will yield meaningful insights by comparing the highest score among the current sources,  $\max_{s_i \in \mathcal{S}_l} \text{Score}(s_i)$ , to a predefined threshold  $\mathcal{T}$ . If no source surpasses this threshold, the recursive process terminates:

$$\max_{s_i \in \mathcal{S}_l} \text{Score}(s_i) < \mathcal{T} \implies \text{Terminate}$$

If at least one source meets or exceeds the threshold, WISE selects the top  $k$  sources based on their scores for further exploration:

$$\mathcal{S}_{l+1} = \text{Top}_k(\mathcal{S}_l, \text{Score})$$

This pruning step ensures that computational efforts are focused on sources most likely to enrich the knowledge container.

### D. Knowledge Container Construction

The knowledge container  $\mathcal{K}_l$  is updated by incorporating the filtered content of the chosen sources. This process further enhances the repository of query-relevant information. The update rule is defined as:

$$\mathcal{K}_{l+1} = \Lambda(\mathcal{K}_l, \mathcal{S}_{l+1}) = \mathcal{K}_l \cup \bigcup_{s_i \in \mathcal{S}_{l+1}} \mathcal{F}(s_i)$$

Here,  $\Lambda$  represents an LLM-powered fusion function designed to merge new information from the selected sources  $\mathcal{S}_{l+1}$  with the existing knowledge base  $\mathcal{K}_l$ . Upon termination of the recursive process, the final knowledge base  $\mathcal{K}_f$  represents the aggregated and refined knowledge for the query  $q$ :

$$\mathcal{K}_f = \mathcal{K}_l \quad \text{at termination}$$

### E. Recursive Algorithm: WISE Framework

Algorithm 1 outlines the recursive process for constructing a query-specific knowledge base. The subsequent layer's sources,  $\mathcal{S}_{l+1}$ , are obtained by analyzing the links embedded in the filtered content of the top  $k$  sources from the current layer,  $\text{Top}_k(\mathcal{S}_l, \text{Score})$ . This ensures that the exploration focuses on paths that are both contextually relevant and computationally efficient.

## III. EXPERIMENT

The experimental setup and methodology used to evaluate WISE's ability to extract and synthesize knowledge from unstructured scientific data are detailed here. Our experiments focused on the query:

*Q: What is the comprehensive set of diseases and phenotypes that are linked to genetic variants within the HBB gene?*

This query, centered on the *HBB* gene, provided a rigorous test case for WISE's capabilities due to the domain's complexity and the interconnected nature of the information.

### A. Experimental Setup

The experiment started with an initial set of 24 sources related to the *HBB* gene, obtained from the HGNC database [6]. Each source was meticulously classified into sections using structural elements, tags, and hyperlinks, extracted through a regular expression-based process. This resulted in a dataset enriched with metadata, including section identifiers and reference sources. We performed asynchronous content extraction from these sources, which served as the foundational nodes for WISE's progressive deepening process [17].

### B. Query-Specific Content Filtering

WISE begins by employing the LLM-driven content filtering process detailed in Section II-A to refine the raw content of each source based on its relevance to the user-specified query. For this experiment, focused on the query  $Q$ , the filtering function  $\Gamma$  takes advantage of contextual understanding of the LLM to isolate pertinent information, eliminating extraneous

---

**Algorithm 1: Recursive WISE Framework**


---

**Input :** Query  $q$ , Initial sources  $\mathcal{S}_0$ , Knowledge Container  $\mathcal{K}_0$ , Threshold  $\mathcal{T}$

**Output:** Final knowledge base  $\mathcal{K}_f$

```

1 Function WISE ( $\mathcal{S}_l, \mathcal{K}_l, l$ ) :
2   if  $\mathcal{S}_l = \emptyset$  or  $\max_{s_i \in \mathcal{S}_l} \text{Score}(s_i) < \mathcal{T}$  then
3     return  $\mathcal{K}_l$ ; // Terminate recursion
        if no significant sources
        remain or the set of sources is
        empty
4   foreach  $s_i \in \mathcal{S}_l$  do
5      $\mathcal{F}(s_i) \leftarrow \Gamma(q, \mathcal{C}(s_i))$ ; // Filter content
        for query  $q$  using filtering
        function  $\Gamma$ 
6      $\text{Score}(s_i) \leftarrow \Psi(\mathcal{F}(s_i), \mathcal{K}_l)$ ; // Calculate
        score for filtered content and
        knowledge base using scoring
        function  $\Psi$ 
7    $\mathcal{S}_{l+1} \leftarrow \text{Top}_k(\mathcal{S}_l, \text{Score})$ ; // Select top  $k$ 
        sources for the next layer based
        on their scores
8    $\mathcal{K}_{l+1} \leftarrow \Lambda(\mathcal{K}_l, \mathcal{S}_{l+1})$ ; // Update knowledge
        base by merging with selected
        sources using fusion function  $\Lambda$ 
9   return WISE ( $\mathcal{S}_{l+1}, \mathcal{K}_{l+1}, l+1$ );
        // Recursive call to the next
        layer
10 return  $\mathcal{K}_f \leftarrow \text{WISE}(\mathcal{S}_0, \mathcal{K}_0, 0)$ ; // Initialize
        recursion with initial sources, empty
        knowledge container, and layer 0

```

---

content such as advertisements, unrelated sections, and general background information that does not directly address the specifics of the query.

Figure 3 illustrates the significant impact of content filtering on data volume for selected sources, demonstrating the substantial reduction in content size achieved through this process. Notably, the UniProt [16] entry for the *HBB* gene, initially containing 8,249 words, was reduced to just 355 words after applying query-specific filtering. This exemplifies the prevalence of extraneous information even within highly regarded scientific resources. Across all sources, filtering reduced content size by an average of 80.14%, with reductions as high as 96.12% observed in some cases. This dramatic reduction highlights the effectiveness of LLM-driven filtering in isolating relevant content, thereby significantly reducing computational overhead and improving the efficiency and precision of downstream processing stages. By focusing on query-relevant information, WISE ensures that subsequent steps, such as score calculation and knowledge integration, operate on a refined and highly pertinent dataset.

### C. Score Calculation and Ranking

To validate the effectiveness of our scoring mechanism in prioritizing query-specific relevance, we conducted experi-

ments focusing on the query  $Q$ . As detailed in Section II-B, WISE employs a dynamic, content-driven scoring approach that contrasts sharply with traditional ranking methods used by systems like Google and ChatGPT, which often rely on factors like source popularity or SEO (Search engine optimization) [18] optimization.

Our scoring mechanism calculates a combined normalized score for each source, integrating both local efficiency (source size) and global contribution (impact on the knowledge container). This ensures that sources offering substantive, query-specific insights are prioritized, regardless of their general popularity. Initial experiments, illustrated in Figure 3, demonstrate that widely recognized sources, such as UniProt [16] or AlphaFold [19], did not always rank highest due to the presence of content unrelated to the specific query. This finding validates our design choice to prioritize content relevance over superficial attributes.

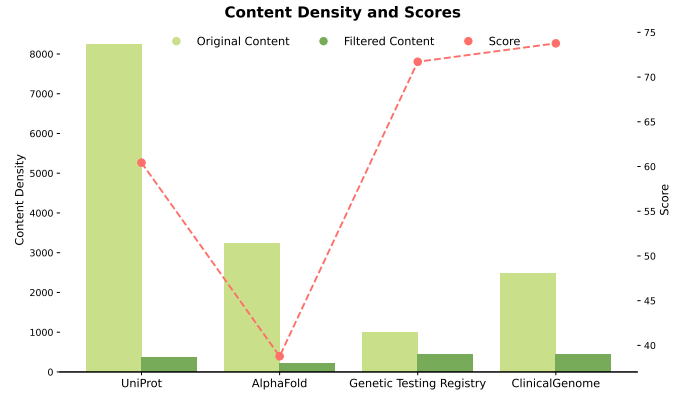


Fig. 3. WISE filtering and scoring in selected 4 sources, showing content density before and after filtering along with their scores. It demonstrates that content size alone does not determine a high score; the unique contribution must be significant relative to the existing knowledge base.

### D. Thresholding and Knowledge Container Construction

WISE employs a threshold-based pruning mechanism, as detailed in Section II-C, to ensure efficient exploration of the knowledge space. This mechanism dynamically determines when to terminate the search process based on the diminishing returns observed in source scores. As the system progresses through successive layers, the knowledge container (described in Section II-D) grows, leading to increased overlap between newly encountered sources and the existing knowledge. Consequently, the unique contribution of each new source tends to decrease.

In this experiment, a threshold value of 20 was empirically determined to effectively balance exploration and exploitation. When the highest score among the current sources falls below this threshold, or when no additional sources are available, WISE terminates the exploration process. This adaptive stopping criterion, mirroring the behavior of expert researchers, prevents the system from pursuing low-yield paths and conserves computational resources.



Figure 4 illustrates this phenomenon, showing a consistent decrease in the scores of top-ranked sources across three successive layers. This trend validates the effectiveness of the threshold-based pruning strategy in identifying the point of diminishing returns. The knowledge container, constructed by integrating the filtered content from the top two sources at each layer, evolves into a rich and contextually relevant repository of information, progressively refined with each iteration. This iterative process ensures that the final knowledge container is both comprehensive and focused, containing the most valuable insights related to the query.

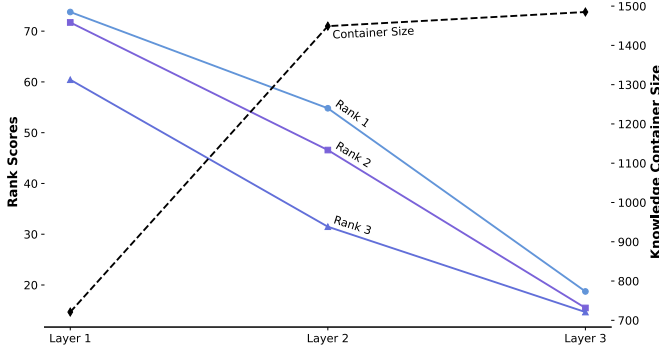


Fig. 4. Scores for Ranks 1, 2, and 3 across layers, highlighting their gradual decrease over time. This also demonstrates the growth of the knowledge container size with each layer, eventually plateauing as layers increase.

#### IV. RESULTS

A comprehensive analysis of WISE’s performance, contrasting it with established methods across several critical dimensions, is presented here. Our evaluation, based on the experiments detailed in Section III, focused on the query  $Q$ , designed to probe both the breadth and depth of knowledge extraction. To provide a robust and nuanced evaluation, we compared WISE against four baseline systems, each representing a different approach to information retrieval and synthesis: Pure ChatGPT [15], [20] (a standard model, version GPT-4o accessed via the OpenAI API, relying solely on its pre-trained knowledge), ChatGPT with Search [21], [22] (a version of ChatGPT, version GPT-4o augmented with web search capabilities), Gemini [23] (Google’s large language model, designed to integrate information from various sources), and traditional Google Search [24], [25] (the standard Google Search engine, considering the top 5 search results for analysis). For all baseline systems, default parameters were used, and the query  $Q$  was directly input without any modifications.

##### A. Evaluation Metrics

Our analysis employs a combination of quantitative metrics, focusing on the relevance, uniqueness, and depth of the extracted information.

###### 1) ROUGE and BLEU

To assess the overlap and uniqueness of the information extracted by each system, we employed ROUGE [26] (Recall-Oriented Understudy for Gisting Evaluation) and BLEU [27]

(Bilingual Evaluation Understudy) metrics, commonly used for evaluating machine-generated text. We adapted these metrics, calculating ROUGE-1, ROUGE-2, and ROUGE-L, along with BLEU, to compare each system’s output against the others. Figure 5 presents the average ROUGE and BLEU scores for each system when used as a reference, revealing that WISE consistently exhibits the lowest scores, indicating that its generated content is more distinct and less repetitive compared to other approaches.

###### 2) Recall

To evaluate the comprehensiveness of each system, we calculated their recall based on a combined output created by taking the union of all unique diseases identified by any of the five systems. This combined output, representing a comprehensive collection of potentially relevant diseases, is detailed in Table II. WISE achieved a recall of 0.842, significantly outperforming the baseline systems, as shown in Figure 6. In contrast, ChatGPT achieved a recall of 0.474, ChatGPT with Search achieved a recall of 0.368, while Google Search Gemini and traditional Google Search scored considerably lower at 0.105 and 0.158 respectively. This emphasizes WISE’s superior ability to identify a greater proportion of potentially relevant diseases.

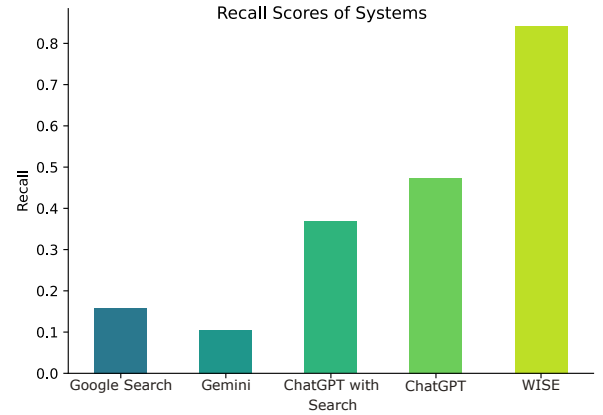


Fig. 6. Recall scores of each system, demonstrating that WISE identifies a greater proportion of diseases from the combined output.

###### 3) Level-Based Analysis

To further analyze the depth of information provided by each system, and to establish a metric that can generally be used to assess the richness of content, we employed a level-based analysis. In this approach, each identified disease was manually assigned a level from 0 to 5 based on the following criteria:

- **Level 0:** Disease name only.
- **Level 1:** Basic description of the disease.
- **Level 2:** Information about the cause of the disease.
- **Level 3:** Details about the disease’s mechanism.
- **Level 4:** Information about diagnosis, treatment, or prognosis.
- **Level 5:** Links to external resources or discussion of ongoing research.

This level-based approach, while applied to evaluate WISE and baselines in this experiment, also offers a generalizable

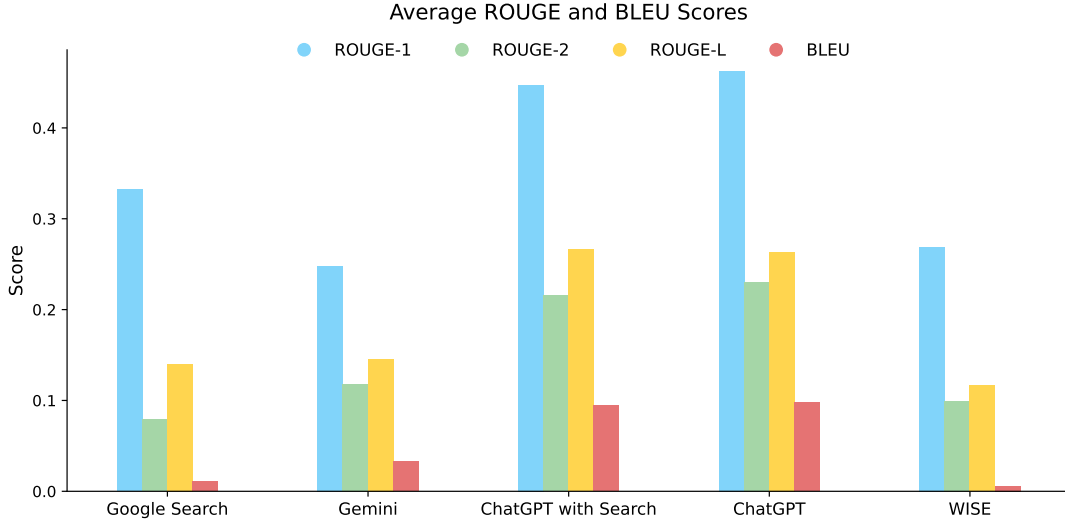


Fig. 5. Average ROUGE and BLEU scores for each system when used as a reference, demonstrating that WISE’s output is the most distinct and least repetitive.

TABLE II  
DISEASE IDENTIFICATION ACROSS SYSTEMS

Disease	Normal Google Search	Google Search Gemini	ChatGPT with Search	ChatGPT	WISE
Hemoglobin SC	X	X	✓	✓	✓
Hemoglobin O	X	X	X	✓	X
Hemoglobin S/ $\beta$ -Thalassemia	X	X	✓	✓	✓
Hemoglobin S Oman	X	X	X	X	✓
Malaria	X	X	X	X	✓
Sickle Cell Disease	✓	✓	✓	✓	✓
Hispanic Gamma-Delta- $\beta$ Thalassemia	X	X	X	X	✓
$\beta$ -Type Methemoglobinemia	X	X	X	X	✓
Dominant $\beta$ -Thalassemia	X	X	X	X	✓
Hemoglobinopathies	✓	X	X	X	X
Heinz Body Anemia	X	X	X	X	✓
Hemoglobin C	X	X	✓	✓	✓
Hemoglobin M	X	X	X	X	✓
Hemoglobin D	X	X	X	✓	X
Familial Erythrocytosis 6	X	X	X	X	✓
Hemoglobin S Antilles	X	X	X	X	✓
Hemoglobin E	X	X	✓	✓	✓
Hereditary Persistence of Fetal Hemoglobin	X	X	✓	✓	✓
$\beta$ -Thalassemia	✓	✓	✓	✓	✓

method for assessing the richness of information provided by any system. Figure 7 presents the average level of detail for each system, showing that WISE achieved the highest score (3.81), significantly surpassing the other systems. This result demonstrates that WISE provides more in-depth and comprehensive information about the identified diseases compared to the baselines.

The convergence of these metrics paints a compelling picture of WISE’s strengths. It is not only capable of identifying a wider range of diseases and phenotypes linked to the *HBB* gene (demonstrated by its high recall) but also of providing richer, more unique, and contextually relevant information (shown through the ROUGE, BLEU, and level-based analyses). The superior performance of WISE across all these metrics highlights its potential as a transformative system for information retrieval in complex domains.

These findings can be attributed to WISE’s unique design, particularly its dynamic scoring mechanism, tree-based ar-

chitecture, and LLM-powered content filtering. The dynamic scoring prioritizes content relevance over superficial attributes, ensuring that the most valuable sources are identified. The tree-based architecture allows for efficient exploration of the knowledge space, while the LLM-driven filtering ensures that only pertinent information is processed.

## V. APPLICATIONS

The adaptability of WISE extends far beyond the specific use case we have explored thus far, underscoring its potential to transform knowledge synthesis across diverse domains. In this section, we illustrate WISE’s versatility by exploring its prospective applications, demonstrating how its unique capabilities can address existing challenges and accelerate progress in various fields.

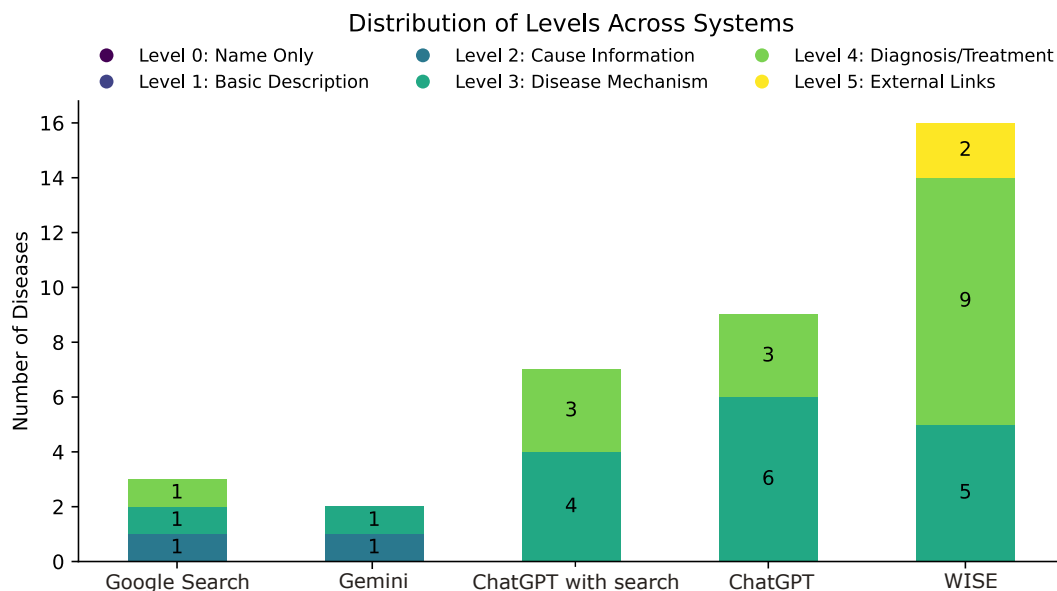


Fig. 7. Average level of detail for each system, demonstrating WISE’s superior ability to provide in-depth and comprehensive information about identified diseases.

#### A. Drug Discovery: Unveiling Novel Therapeutic Pathways

The process of drug discovery often hinges on unraveling the intricate relationships between genes, diseases, and potential therapeutic targets. WISE offers a transformative approach to this challenge by providing an efficient and comprehensive means of identifying these complex associations, surpassing the limitations of manual methods and existing systems. Consider the following query:

*Q: What diseases are associated with C16orf82, and are there any existing drugs targeting these conditions?*

This query, while seemingly simple, requires navigating a complex web of genetic and pharmacological information. WISE can reveal novel connections within the existing literature, identifying not only diseases sharing genetic origins or structural similarities in proteins (including overlapping reading frames, ORFs) but also highlighting previously unlinked diseases that may share common pathways or molecular interactions. These insights provide valuable leads for drug repurposing and the development of novel therapeutic strategies, effectively accelerating the drug discovery process by integrating these discoveries with relevant drug information, thus delivering precise and actionable intelligence for pharmaceutical development.

#### B. Material Structure Analysis: Accelerating Inverse Design

The field of materials science, particularly the domain of inverse material design, is often constrained by time-intensive and inefficient processes for identifying suitable material structures that meet specific requirements. WISE offers a streamlined approach to this challenge. For instance, consider this query:

*Q: What are the most suitable material structures for achieving high thermal conductivity and mechanical strength in lightweight applications?*

By directly linking structural properties to specific application requirements, WISE enables researchers to explore material structure databases with unprecedented speed and efficiency. This reduces the manual effort and time required for material selection, allowing researchers to focus on designing solutions that precisely meet their objectives, rather than spending excessive time on information gathering. By reducing reliance on inefficient, time-intensive methods, WISE serves as a powerful tool to accelerate the field of materials science.

#### C. Social Issue Analysis: Illuminating Complex Societal Challenges

WISE is equally applicable to addressing complex social issues, where the analysis of vast amounts of unstructured data is critical. Social scientists and policymakers grapple with a range of complex challenges, requiring the integration of data from diverse sources to identify patterns and develop effective interventions. For example:

*Q: What factors are contributing to the rising cancer rates in [specific location]?*

This query, designed to highlight the social, environmental, and economic challenges that drive increases in rates of cancer, demands the integration of multiple viewpoints, datasets, and research findings. WISE can synthesize information from diverse sources to identify complex patterns, including increased exposure to environmental toxins, socio-economic inequalities, or shortcomings in public health policy. By highlighting key contributing factors, WISE offers researchers and policymakers critical data points and insight that empowers them to develop data-driven hypotheses and implement more targeted interventions. The ability of WISE to generate comparative



examples from similar regions further allows for a deeper, more nuanced understanding of the issue at hand.

The examples above underscore WISE’s flexibility and adaptability, demonstrating its applicability beyond specific domains. Its ability to process complex queries and synthesize domain-specific knowledge makes it a valuable asset across a broad range of fields, from medical research to materials science and social issue analysis. Over time, WISE’s workflow has the potential to further accelerate progress in areas like cancer research, environmental studies, and many others, by providing a more efficient and reliable approach to obtaining detailed and context-aware information. These diverse applications highlight WISE’s potential to enable deeper insights and foster innovation across a wide spectrum of scientific and societal disciplines.

## VI. FUTURE WORK

While WISE has demonstrated significant capabilities in our experiments, its journey is far from complete. We are actively exploring several promising avenues for further development, poised to enhance the system’s robustness, efficiency, and applicability. The following represent key directions for future research, although these improvements are not yet incorporated into the current implementation and remain outside the scope of this paper.

### A. Knowledge Graph Integration: Unlocking Deeper Relational Insights

The integration of knowledge graphs represents a transformative opportunity to amplify WISE’s ability to reason about complex relationships within scientific data. Knowledge graphs, which represent information as interconnected nodes and edges, offer a structured approach for preserving and reasoning about intricate interdependencies. Such an approach transcends the limitations of purely text-based analysis, enabling WISE to identify connections between seemingly disparate entities and uncover subtle yet significant articulation points that are often missed by traditional methods.

By incorporating knowledge graphs, WISE can maintain a dynamic understanding of the relationships between entities, thereby eliminating the need for repeated, LLM-driven knowledge unions. Instead, newly extracted information can be directly appended to the appropriate nodes and edges in the graph, ensuring continuity, efficiency, and preventing information loss. Preliminary experiments, for example, have demonstrated that representing the UniProt entry for the *HBB* gene as a knowledge graph with 56 nodes and 55 edges effectively captures its content with reduced complexity compared to raw text. Further filtering this knowledge graph with a query related to *HBB*-specific diseases resulted in a focused subgraph with 11 nodes and 16 edges, while maintaining key interconnected causes, like shared hormonal pathways across multiple diseases.

Furthermore, knowledge graphs facilitate intuitive visualizations, enhancing user understanding and interpretation. Their structured nature supports advanced reasoning capabilities,

enabling WISE to achieve deeper insights through graph-based matching techniques, outperforming traditional content comparison approaches. This integration promises a more powerful and efficient means of knowledge discovery.

### B. Enhanced Query Engagement: Steering Towards Precise Intent

WISE currently relies on user-provided queries, future iterations will focus on enhancing query engagement to steer the system towards a more precise understanding of user intent. Although our similarity-based searches have proven effective so far, there is a clear potential to amplify WISE’s performance through a more iterative and user-involved query process. Future iterations of WISE will incorporate mechanisms for better understanding user intent through supplementary information gathering. For example, users could be prompted to provide additional context or goals for their search, enabling more precise and targeted results.

Moreover, the system could implement automatic query enhancement, leveraging prior searches and literature data to refine user input iteratively. This process may also include layers of semantic understanding, improved similarity measures, and propose augmented queries for user approval. These advancements would significantly reduce the burden on less-experienced users, simplifying complex search tasks while maintaining the high standards of precision that WISE offers. These improvements could also help the system identify if the query is underdefined or if there is some implicit constraints in the query.

## VII. RELATED WORK

Prior research in information extraction has significantly advanced from manual curation and feature-engineered machine learning models to more sophisticated LLM-based approaches. Systems such as GIX [28] effectively leverage large language models to automate gene interaction extraction, outperforming earlier methods on benchmark datasets. Similarly, tree-structured neural architectures have improved the identification of protein-protein interactions [29], and comprehensive reviews highlight the rise of neural network-based classifiers in biomedical relation extraction [30]. Domain-adapted models, such as BioBERT [31], have demonstrated notable gains in precision for tasks including named entity recognition and relation extraction, while further explorations have extended the scope from binary to complex biomedical relations [32]. Beyond the biomedical domain, research has shown that integrating pre-trained models with domain-specific corpora and graph-based reasoning can uncover intricate patterns, enabling richer insights and improved retrieval accuracy [33]–[36].

Despite these advancements, existing approaches often struggle to dynamically refine their focus or determine when further exploration yields diminishing returns. Retrieval-Augmented Generation techniques [37] and workflow orchestration strategies [38] have attempted to address these challenges by pruning sources and enhancing verification, yet they rarely employ hierarchical, query-driven frameworks that integrate filtering, scoring, and adaptive stopping criteria. Other

efforts have focused on bridging the gap between unstructured and structured data [39], [40] but have not fully embraced iterative, tree-based methodologies for source selection and knowledge consolidation. In this context, WISE advances the state of the art by combining robust filtering, dynamic ranking, and a scalable tree-inspired workflow, thereby ensuring that only the most valuable, context-relevant information is retained for efficient, high-quality knowledge extraction.

### VIII. DISCUSSION

The development and application of WISE highlight its transformative potential in synthesizing knowledge for complex queries, yet certain challenges and limitations merit discussion. One of the primary obstacles encountered during the experiments was restricted access to some data sources. Out of 34 initial sources, WISE successfully extracted content from 24, while the remaining 10 sources were inaccessible. Figure 8 illustrates this disparity, emphasizing the growing trend among web platforms to limit data extraction. This failure to extract content stemmed primarily from two key factors: the increasing prevalence of paywalls, which create direct economic barriers to access, and security measures like bot detection, which are often deployed to protect intellectual property, user privacy, and prevent the unauthorized use of data for LLM training and content analysis. While these restrictions serve important purposes, such as protecting content creators and user data, they also hinder innovations like WISE, which focus on advancing non-commercial, academic, and research applications. Overcoming these challenges may require collaboration with data providers, the development of ethical and compliant retrieval techniques.

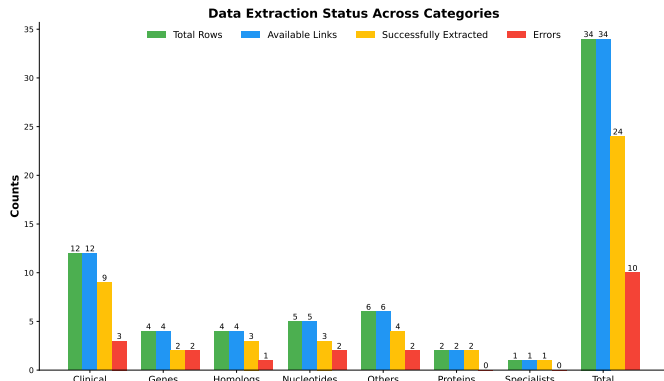


Fig. 8. The number of sources that restrict data extraction by implementing blocks on automated processes.

Another important consideration is the comprehensiveness of WISE’s outputs. By design, WISE delivers detailed, authoritative responses that include exhaustive references, disease sub-variations, and contextual information. While this level of detail is highly valuable for academic and clinical professionals, it may overwhelm non-specialist users who require more concise and simplified information. The extensive details and length of the responses may be daunting, highlighting the need for customizable output formats tailored to different audiences.

Features such as adjustable levels of detail or user-specific summaries could enhance WISE’s accessibility and usability across a broader range of users.

A related challenge lies in word weighting during the synthesis of information. In WISE’s current implementation, more frequent words like “and” or “the” are deprioritized based on term frequency (TF), ensuring that the system focuses on content-specific terms. However, the exact contextual relationships between terms—critical for disambiguating similar entities—could benefit from improvements. Techniques like TF-IDF (Term Frequency-Inverse Document Frequency) are currently in consideration, as they would assign higher weight to less frequent but more meaningful terms. Additionally, as discussed in Section VI, the integration of a knowledge graph offers a promising solution to this challenge. By representing relationships explicitly through nodes and edges, a knowledge graph would inherently prioritize meaningful connections, eliminating reliance on textual frequency metrics.

Despite these challenges, WISE represents a significant advancement in information retrieval and synthesis, setting a new standard for addressing complex queries. Its ability to dynamically filter, rank, and construct comprehensive knowledge containers demonstrates its transformative potential across diverse domains. The innovative architecture of WISE effectively bridges critical gaps in existing systems, offering a robust tool for academic, clinical, and interdisciplinary applications.

### IX. CONCLUSION

WISE presents a novel and effective approach to navigating the complexities of scientific information retrieval. By integrating LLM-driven filtering, dynamic ranking, and adaptive stopping criteria within a tree-based framework, WISE empowers researchers to efficiently and accurately extract and synthesize knowledge from vast and heterogeneous data sources. Our experiments on gene-disease association queries demonstrated WISE’s superior performance compared to baseline methods, showcasing its ability to uncover a broader range of relevant information, including rare conditions and nuanced connections often overlooked by traditional search engines and basic LLM implementations. This enhanced precision and comprehensiveness, achieved through a content-driven, progressive deepening approach, offers significant potential for accelerating scientific discovery across diverse domains.

The development of WISE represents a substantial step forward in the pursuit of intelligent knowledge discovery. Its human-inspired methodology, mimicking the systematic approach of expert researchers, allows for a balanced exploration of information, prioritizing high-value insights while effectively managing computational resources. The framework’s adaptability and scalability, demonstrated through its application to diverse research domains, further suggest its potential as a generalizable solution for complex information landscapes. We believe that WISE offers a valuable tool for researchers seeking to unlock the full potential of the ever-expanding universe of scientific knowledge, paving the way for more efficient and impactful research endeavors. By address-

ing the limitations of traditional search engines and general-purpose LLMs, WISE provides a robust and scalable solution for extracting and synthesizing knowledge, ultimately contributing to more informed and accelerated scientific progress.

## ACKNOWLEDGMENT

This Research was supported in part by a National Institutes of Health IDeA grant P20GM103408, a National Science Foundation CSSI grant OAC 2410668, and a US Department of Energy grant DE-0011014.

## REFERENCES

- [1] J. Sedlakova, P. Daniore, A. Horn Wintsch, M. Wolf, M. Stanikic, C. Haag, C. Sieber, G. Schneider, K. Staub, D. Alois Ettlin, O. Grübner, F. Rinaldi, V. von Wyl, and for the University of Zurich Digital Society Initiative (UZH-DSI) Health Community, "Challenges and best practices for digital unstructured data enrichment in health research: A systematic narrative review," *PLOS Digital Health*, vol. 2, no. 10, pp. 1–22, 10 2023. [Online]. Available: <https://doi.org/10.1371/journal.pdig.0000347>
- [2] J. Luo, M. Wu, D. Gopukumar, and Y. Zhao, "Big data application in biomedical research and health care: A literature review," *Biomedical Informatics Insights*, vol. 8, p. BILS31559, 2016, pMID: 26843812. [Online]. Available: <https://doi.org/10.4137/BILS31559>
- [3] C. Zhai, "Large language models and future of information retrieval: Opportunities and challenges," in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2024, pp. 481–490.
- [4] A. Salemi and H. Zamani, "Towards a search engine for machines: Unified ranking for multiple retrieval-augmented large language models," in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR '24. New York, NY, USA: Association for Computing Machinery, 2024, p. 741–751. [Online]. Available: <https://doi.org/10.1145/3626772.3657733>
- [5] N. Ziemis, W. Yu, Z. Zhang, and M. Jiang, "Large language models are built-in autoregressive search engines," *arXiv preprint arXiv:2305.09612*, 2023.
- [6] HUGO Gene Nomenclature Committee (HGNC), "HBB Gene - Gene Symbol Report," 2024, accessed: 2024-12-11. [Online]. Available: [https://www.genenames.org/data/gene-symbol-report/#!/hgnc\\_id/HGNC:4827](https://www.genenames.org/data/gene-symbol-report/#!/hgnc_id/HGNC:4827)
- [7] Clinical Genome Resource (ClinGen), "HBB Gene - Clinical Genome Knowledge Base," 2024, accessed: 2024-12-11. [Online]. Available: <https://search.clinicalgenome.org/kb/genes/HGNC:4827>
- [8] National Center for Biotechnology Information (NCBI), *Sickle Cell Anemia - NCBI Bookshelf*, 2024, accessed: 2024-12-11. [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK1435/>
- [9] G. L. Garcia, J. R. R. Manesco, P. H. Paiola, L. Miranda, M. P. de Salvo, and J. P. Papa, "A review on scientific knowledge extraction using large language models in biomedical sciences," 2024. [Online]. Available: <https://arxiv.org/abs/2412.03531>
- [10] A. Alshami, M. Elsayed, E. Ali, A. E. E. Eltoukhy, and T. Zayed, "Harnessing the power of chatgpt for automating systematic review process: Methodology, case study, limitations, and future directions," *Systems*, vol. 11, no. 7, 2023. [Online]. Available: <https://www.mdpi.com/2079-8954/11/7/351>
- [11] S. Saxena, R. Sangani, S. Prasad, S. Kumar, M. Athale, R. Awhad, and V. Vaddina, "Large-Scale Knowledge Synthesis and Complex Information Retrieval from Biomedical Documents," in *2022 IEEE International Conference on Big Data (Big Data)*. Los Alamitos, CA, USA: IEEE Computer Society, Dec. 2022, pp. 2364–2369. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/BigData55660.2022.10020725>
- [12] A. Sneyd and M. Stevenson, "Modelling stopping criteria for search results using Poisson processes," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 3484–3489. [Online]. Available: <https://aclanthology.org/D19-1351/>
- [13] M. Parmentier and A. Legay, "Adaptive stopping algorithms based on concentration inequalities," in *Bridging the Gap Between AI and Reality: Second International Conference, AISoLA 2024, Crete, Greece, October 30 – November 3, 2024, Proceedings*. Berlin, Heidelberg: Springer-Verlag, 2025, p. 336–353. [Online]. Available: [https://doi.org/10.1007/978-3-031-75434-0\\_23](https://doi.org/10.1007/978-3-031-75434-0_23)
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2023. [Online]. Available: <https://arxiv.org/abs/1706.03762>
- [15] OpenAI, "Gpt-4o," 2024, the GPT-4o model supports a context window of up to 128,000 tokens. [Online]. Available: <https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo>
- [16] U. Consortium, "Hemoglobin subunit beta," 2024, accessed: 2024-12-11. [Online]. Available: <https://www.uniprot.org/uniprotkb/P68871>
- [17] J. Han, "Mining knowledge at multiple concept levels," in *Proceedings of the fourth international conference on Information and knowledge management*, 1995, pp. 19–24.
- [18] D. Sharma, R. Shukla, A. K. Giri, and S. Kumar, "A brief review on search engine optimization," in *2019 9th international conference on cloud computing, data science & engineering (confluence)*. IEEE, 2019, pp. 687–692.
- [19] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko *et al.*, "Highly accurate protein structure prediction with alphafold," *nature*, vol. 596, no. 7873, pp. 583–589, 2021.
- [20] S. Shahriar, B. D. Lund, N. R. Mannuru, M. A. Arshad, K. Hayawi, R. V. K. Bevara, A. Mannuru, and L. Batool, "Putting gpt-4o to the sword: A comprehensive evaluation of language, vision, speech, and multimodal proficiency," *Applied Sciences*, vol. 14, no. 17, 2024. [Online]. Available: <https://www.mdpi.com/2076-3417/14/17/7782>
- [21] OpenAI, "Introducing chatgpt search," <https://openai.com/index/introducing-chatgpt-search/>, 2024, accessed: 2024-12-27.
- [22] W. Sun, L. Yan, X. Ma, S. Wang, P. Ren, Z. Chen, D. Yin, and Z. Ren, "Is chatgpt good at search? investigating large language models as re-ranking agents," *arXiv preprint arXiv:2304.09542*, 2023.
- [23] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican *et al.*, "Gemini: a family of highly capable multimodal models," *arXiv preprint arXiv:2312.11805*, 2023.
- [24] J. Piasecki, M. Waligora, and V. Dranseika, "Google search as an additional source in systematic reviews," *Science and engineering ethics*, vol. 24, pp. 809–810, 2018.
- [25] R. L. Cilibrasi and P. M. Vitanyi, "The google similarity distance," *IEEE Transactions on knowledge and data engineering*, vol. 19, no. 3, pp. 370–383, 2007.
- [26] C. Lin, "Recall-oriented understudy for gisting evaluation (rouge)," *Retrieved August*, vol. 20, p. 2005, 2005.
- [27] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.
- [28] J. K. Gill, M. Chetty, S. Lim, and J. Hallinan, "Large language model based framework for automated extraction of genetic interactions from unstructured data," *PLOS ONE*, vol. 19, no. 5, pp. 1–22, 05 2024. [Online]. Available: <https://doi.org/10.1371/journal.pone.0303231>
- [29] M. Ahmed, J. Islam, M. R. Samee, and R. E. Mercer, "Identifying protein-protein interaction using tree lstm and structured attention," in *2019 IEEE 13th International Conference on Semantic Computing (ICSC)*, 2019, pp. 224–231.
- [30] Y. Zhang, H. Lin, Z. Yang, J. Wang, Y. Sun, B. Xu, and Z. Zhao, "Neural network-based approaches for biomedical relation classification: A review," *J. of Biomedical Informatics*, vol. 99, no. C, Nov. 2019. [Online]. Available: <https://doi.org/10.1016/j.jbi.2019.103294>
- [31] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "Biobert: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [32] D. Zhou, D. Zhong, and Y. He, "Biomedical relation extraction: From binary to complex," *Computational and Mathematical Methods in Medicine*, vol. 2014, no. 1, p. 298473, 2014. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1155/2014/298473>
- [33] B. Sarmah, D. Mehta, B. Hall, R. Rao, S. Patel, and S. Pasquali, "Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction," in *Proceedings of the 5th ACM International Conference on AI in Finance*, 2024, pp. 608–616.

- [34] M. J. Buehler, "Accelerating scientific discovery with generative knowledge extraction, graph-based representation, and multimodal intelligent graph reasoning," *Machine Learning: Science and Technology*, vol. 5, no. 3, p. 035083, sep 2024. [Online]. Available: <https://dx.doi.org/10.1088/2632-2153/ad7228>
- [35] M. Zuluaga, S. Robledo, O. Arbelaez-Echeverri, G. A. Osorio-Zuluaga, and N. Duque-Méndez, "Tree of science - tos: A web-based tool for scientific literature recommendation. search less, research more!" *Issues in Science and Technology Librarianship*, no. 100, Aug. 2022. [Online]. Available: <https://journals.library.ualberta.ca/istl/index.php/istl/article/view/2696>
- [36] W. Huang, X. Zhao, and X. Huang, "Embedding and extraction of knowledge in tree ensemble classifiers," *Mach. Learn.*, vol. 111, no. 5, p. 1925–1958, May 2022. [Online]. Available: <https://doi.org/10.1007/s10994-021-06068-6>
- [37] S. Yu, M. Cheng, J. Yang, J. Ouyang, Y. Luo, C. Lei, Q. Liu, and E. Chen, "Multi-source knowledge pruning for retrieval-augmented generation: A benchmark and empirical study," 2024. [Online]. Available: <https://arxiv.org/abs/2409.13694>
- [38] S. Fan, X. Cong, Y. Fu, Z. Zhang, S. Zhang, Y. Liu, Y. Wu, Y. Lin, Z. Liu, and M. Sun, "Workflowllm: Enhancing workflow orchestration capability of large language models," 2024. [Online]. Available: <https://arxiv.org/abs/2411.05451>
- [39] Z. Hong, L. Ward, K. Chard, B. Blaiszik, and I. Foster, "Challenges and advances in information extraction from scientific literature: a review," *JOM*, vol. 73, no. 11, pp. 3383–3400, 2021. [Online]. Available: <https://doi.org/10.1007/s11837-021-04902-9>
- [40] T. Xie, H. Zhang, S. Wang, Y. Wan, I. Razzak, C. Kit, W. Zhang, and B. Hoex, "Bytescience: Bridging unstructured scientific literature and structured data with auto fine-tuned large language model in token granularity," 2024. [Online]. Available: <https://arxiv.org/abs/2411.12000>