

RAG-ANYTHING: ALL-IN-ONE RAG FRAMEWORK

Zirui Guo, Xubin Ren, Lingrui Xu, Jiahao Zhang, Chao Huang*

The University of Hong Kong

zrguo101@hku.hk xubinrencs@gmail.com chaohuang75@gmail.com

ABSTRACT

Retrieval-Augmented Generation (RAG) has emerged as a fundamental paradigm for expanding Large Language Models beyond their static training limitations. However, a critical misalignment exists between current RAG capabilities and real-world information environments. Modern knowledge repositories are inherently multimodal, containing rich combinations of textual content, visual elements, structured tables, and mathematical expressions. Yet existing RAG frameworks are limited to textual content, creating fundamental gaps when processing multimodal documents. We present RAG-Anything, a unified framework that enables comprehensive knowledge retrieval across all modalities. Our approach reconceptualizes multimodal content as interconnected knowledge entities rather than isolated data types. The framework introduces dual-graph construction to capture both cross-modal relationships and textual semantics within a unified representation. We develop cross-modal hybrid retrieval that combines structural knowledge navigation with semantic matching. This enables effective reasoning over heterogeneous content where relevant evidence spans multiple modalities. RAG-Anything demonstrates superior performance on challenging multimodal benchmarks, achieving significant improvements over state-of-the-art methods. Performance gains become particularly pronounced on long documents where traditional approaches fail. Our framework establishes a new paradigm for multimodal knowledge access, eliminating the architectural fragmentation that constrains current systems. Our framework is open-sourced at: <https://github.com/HKUDS/RAG-Anything>.

1 INTRODUCTION

Retrieval-Augmented Generation (RAG) has emerged as a fundamental paradigm for expanding the knowledge boundaries of Large Language Models (LLM) beyond their static training limitations [Zhang et al. \(2025\)](#). By enabling dynamic retrieval and incorporation of external knowledge during inference, RAG systems transform static language models into adaptive, knowledge-aware systems. This capability has proven essential for applications requiring up-to-date information, domain-specific knowledge, or factual grounding that extends beyond pre-training corpora.

However, existing RAG frameworks focus exclusively on text-only knowledge while neglecting the rich multimodal information present in real-world documents. This limitation fundamentally misaligns with how information exists in authentic environments. Real-world knowledge repositories are inherently heterogeneous and multimodal [Abootorabi et al. \(2025\)](#). They contain rich combinations of textual content, visual elements, structured tables, and mathematical expressions across diverse document formats. This textual assumption forces existing RAG systems to either discard non-textual information entirely or flatten complex multimodal content into inadequate textual approximations.

The consequences of this limitation become particularly severe in document-intensive domains where multimodal content carries essential meaning. Academic research, financial analysis, and technical documentation represent prime examples of knowledge-rich environments. These domains fundamentally depend on visual and structured information. Critical insights are often encoded exclusively in non-textual formats. Such formats resist meaningful conversion to plain text.

The consequences of this limitation become particularly severe in knowledge-intensive domains where multimodal content carries essential meaning. Three representative scenarios illustrate the critical

*Corresponding Author: Chao Huang

need for multimodal RAG capabilities. In **Scientific Research**, experimental results are primarily communicated through plots, diagrams, and statistical visualizations. These contain core discoveries that remain invisible to text-only systems. **Financial Analysis** relies heavily on market charts, correlation matrices, and performance tables. Investment insights are encoded in visual patterns rather than textual descriptions. Additionally, **Medical Literature Analysis** depends on radiological images, diagnostic charts, and clinical data tables. These contain life-critical information essential for accurate diagnosis and treatment decisions. Current RAG frameworks systematically exclude these vital knowledge sources across all three scenarios. This creates fundamental gaps that render them inadequate for real-world applications requiring comprehensive information understanding. Therefore, multimodal RAG emerges as a critical advancement. It is necessary to bridge these knowledge gaps and enable truly comprehensive intelligence across all modalities of human knowledge representation.

Addressing multimodal RAG presents three fundamental technical challenges that demand principled solutions. This makes it significantly more complex than traditional text-only approaches. The naive solution of converting all multimodal content to textual descriptions introduces severe information loss. Visual elements such as charts, diagrams, and spatial layouts contain semantic richness that cannot be adequately captured through text alone. These inherent limitations necessitate the design of effective technical components. Such components must be specifically designed to handle multimodal complexity and preserve the full spectrum of information contained within diverse content types.

Technical Challenges. • **First**, the **unified multimodal representation** challenge requires seamlessly integrating diverse information types. The system must preserve their unique characteristics and cross-modal relationships. This demands advanced multimodal encoders that can capture both intra-modal and inter-modal dependencies without losing essential visual semantics. • **Second**, the **structure-aware decomposition** challenge demands intelligent parsing of complex layouts. The system must maintain spatial and hierarchical relationships crucial for understanding. This requires specialized layout-aware parsing modules that can interpret document structure and preserve contextual positioning of multimodal elements. • **Third**, the **cross-modal retrieval** challenge necessitates sophisticated mechanisms that can navigate between different modalities. These mechanisms must reason over their interconnections during retrieval. This calls for cross-modal alignment systems capable of understanding semantic correspondences across text, images, and structured data. These challenges are amplified in long-context scenarios. Relevant evidence is dispersed across multiple modalities and sections, requiring coordinated reasoning across heterogeneous information sources.

Our Contributions. To address these challenges, we introduce RAG-Anything, a unified framework that fundamentally reimagines multimodal knowledge representation and retrieval. Our approach employs a **dual-graph construction strategy** that elegantly bridges the gap between cross-modal understanding and fine-grained textual semantics. Rather than forcing diverse modalities into text-centric pipelines, RAG-Anything constructs **complementary knowledge graphs** that preserve both multimodal contextual relationships and detailed textual knowledge. This design enables seamless integration of visual elements, structured data, and mathematical expressions within a unified retrieval framework. The system maintains **semantic integrity** across modalities while ensuring efficient **cross-modal reasoning capabilities** throughout the process.

Our **cross-modal hybrid retrieval** mechanism strategically combines **structural knowledge navigation** with **semantic similarity matching**. This architecture addresses the fundamental limitation of existing approaches that rely solely on embedding-based retrieval or keyword matching. RAG-Anything leverages explicit graph relationships to capture multi-hop reasoning patterns. It simultaneously employs dense vector representations to identify semantically relevant content that lacks direct structural connections. The framework introduces **modality-aware query processing** and **cross-modal alignment systems**. These enable textual queries to effectively access visual and structured information. This unified approach eliminates the architectural fragmentation that plagues current multimodal RAG systems. It delivers superior performance particularly on long-context documents where relevant evidence spans multiple modalities and document sections.

Experimental Validation. To validate the effectiveness of our proposed approach, we conduct comprehensive experiments on two challenging multimodal benchmarks: DocBench and MMLongBench. Our evaluation demonstrates that RAG-Anything achieves superior performance across diverse domains. The framework represents substantial improvements over state-of-the-art baselines. Notably, our performance gains become increasingly significant as content length increases. We observe particularly pronounced advantages on long-context materials. This validates our core hypothesis

that dual-graph construction and cross-modal hybrid retrieval are essential for handling complex multimodal materials. Our ablation studies reveal that graph-based knowledge representation provides the primary performance gains. Traditional chunk-based approaches fail to capture the structural relationships critical for multimodal reasoning. Case studies further demonstrate that our framework excels at precise localization within complex layouts. The system effectively disambiguates similar terminology and navigates multi-panel visualizations through structure-aware retrieval mechanisms.

2 THE RAG-ANYTHING FRAMEWORK

2.1 PRELIMINARY

Retrieval-Augmented Generation (RAG) has emerged as a fundamental paradigm for dynamically expanding the knowledge boundaries of LLMs. While LLMs demonstrate exceptional reasoning capabilities, their knowledge remains static and bounded by training data cutoffs. This creates an ever-widening gap with the rapidly evolving information landscape. RAG systems address this critical limitation by enabling LLMs to retrieve and incorporate external knowledge sources during inference. This transforms them from static repositories into adaptive, knowledge-aware systems.

The Multimodal Reality: Beyond Text-Only RAG. Current RAG systems face a critical limitation that severely restricts their real-world deployment. Existing frameworks operate under the restrictive assumption that knowledge corpus consists exclusively of plain textual documents. This assumption fundamentally misaligns with how information exists in authentic environments. Real-world knowledge repositories are inherently **heterogeneous and multimodal**, containing rich combinations of textual content, visual elements, structured data, and mathematical expressions. These diverse knowledge sources span multiple document formats and presentation mediums, from research papers and technical slides to web pages and interactive documents.

2.1.1 MOTIVATING RAG-ANYTHING

This multimodal reality introduces fundamental technical challenges that expose the inadequacy of current text-only RAG approaches. Effective multimodal RAG requires unified indexing strategies that can handle disparate data types, cross-modal retrieval mechanisms that preserve semantic relationships across modalities, and sophisticated synthesis techniques that can coherently integrate diverse information sources. These challenges demand a fundamentally different architectural approach rather than incremental improvements to existing systems.

The RAG-Anything framework introduces a unified approach for retrieving and processing knowledge from heterogeneous multimodal information sources. Our system addresses the fundamental challenge of handling diverse data modalities and document formats within a retrieval pipeline. The framework comprises three core components: universal indexing for multimodal knowledge, cross-modal adaptive retrieval, and knowledge-enhanced response generation. This integrated design enables effective knowledge utilization across modalities while maintaining computational efficiency.

2.2 UNIVERSAL REPRESENTATION FOR HETEROGENEOUS KNOWLEDGE

A key requirement for universal knowledge access is the ability to represent heterogeneous multimodal content in a unified, retrieval-oriented abstraction. Unlike existing pipelines that simply parse documents into text segments, RAG-Anything introduces **Multimodal Knowledge Unification**. This process decomposes raw inputs into atomic knowledge units while preserving their structural context and semantic alignment. For instance, RAG-Anything ensures that figures remain grounded in their captions, equations remain linked to surrounding definitions, and tables stay connected to explanatory narratives. This transforms heterogeneous files into a coherent substrate for cross-modal retrieval.

Formally, each knowledge source $k_i \in \mathcal{K}$ (e.g., a web page) is decomposed into atomic content units:

$$k_i \xrightarrow{\text{Decompose}} \{c_j = (t_j, x_j)\}_{j=1}^{n_i}, \quad (1)$$

where each unit c_j consists of a modality type $t_j \in \text{text, image, table, equation, } \dots$ and its corresponding raw content x_j . The content x_j represents the extracted information from the original knowledge source, processed in a modality-aware manner to preserve semantic integrity.

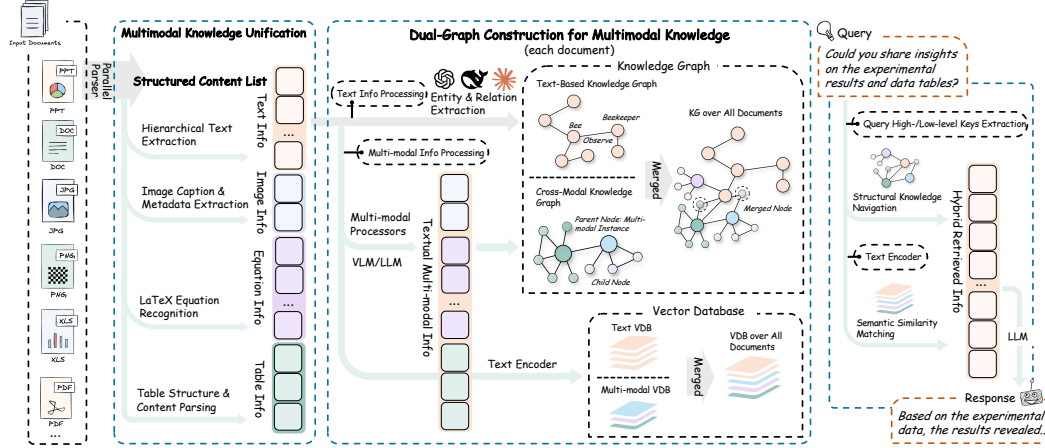


Figure 1: Overview of our proposed universal RAG framework RAG-Anything.

To ensure high-fidelity extraction, RAG-Anything leverages specialized parsers for different content types. Text is segmented into coherent paragraphs or list items. Figures are extracted with associated metadata such as captions and cross-references. Tables are parsed into structured cells with headers and values. Mathematical expressions are converted into symbolic representations. The resulting x_j preserves both content and structural context within the source. This provides a faithful, modality-consistent representation. The decomposition abstracts diverse file formats into atomic units while maintaining their hierarchical order and contextual relationships. This canonicalization enables uniform processing, indexing, and retrieval of multimodal content within our framework.

2.2.1 DUAL-GRAPH CONSTRUCTION FOR MULTIMODAL KNOWLEDGE

While multimodal knowledge unification provides a uniform abstraction across modalities, directly constructing a single unified graph often risks overlooking modality-specific structural signals. The proposed RAG-Anything addresses this challenge through a dual-graph construction strategy. The system first builds a **cross-modal knowledge graph** that faithfully grounds non-textual modalities within their contextual environment. It then constructs a **text-based knowledge graph** using established text-centric extraction pipelines. These complementary graphs are merged through entity alignment. This design ensures accurate cross-modal grounding and comprehensive coverage of textual semantics, enabling richer knowledge representation and robust retrieval.

- **Cross-Modal Knowledge Graph:** Non-textual content like images, tables, and equations contains rich semantic information that traditional text-only approaches often overlook. To preserve this knowledge, RAG-Anything constructs a multimodal knowledge graph where non-text atomic units are transformed into structured graph entities. RAG-Anything leverages multimodal large language models to derive two complementary textual representations from each atomic content unit. The first is a **detailed description** d_j^{chunk} optimized for cross-modal retrieval. The second is an **entity summary** e_j^{entity} containing key attributes such as entity name, type, and description for graph construction. The generation process is context-aware, processing each unit with its local neighborhood $C_j = \{c_k \mid |k - j| \leq \delta\}$, where δ controls the contextual window size. This ensures representations accurately reflect each unit’s role within the broader document structure.

Building on these textual representations, RAG-Anything constructs the graph structure using non-text units as anchor points. For each non-text unit c_j , the graph extraction routine $R(\cdot)$ processes its description d_j^{chunk} to identify fine-grained entities and relations:

$$(\mathcal{V}_j, \mathcal{E}_j) = R(d_j^{\text{chunk}}), \quad (2)$$

where \mathcal{V}_j and \mathcal{E}_j denote the sets of intra-chunk entities and their relations, respectively. Each atomic non-text unit is associated with a multimodal entity node v_j^{mm} that serves as an anchor for

its intra-chunk entities through explicit `belongs_to` edges:

$$\tilde{V} = \{v_j^{\text{mm}}\}_j \cup \bigcup_j \mathcal{V}_j, \quad (3)$$

$$\tilde{E} = \bigcup_j \mathcal{E}_j \cup \bigcup_j \{(u \xrightarrow{\text{belongs_to}} v_j^{\text{mm}}) : u \in \mathcal{V}_j\}. \quad (4)$$

This construction preserves modality-specific grounding while ensuring non-textual content is contextualized by its textual neighborhood. This enables reliable cross-modal retrieval and reasoning.

- **Text-based Knowledge Graph:** For text modality chunks, we construct a traditional text-based knowledge graph following established methodologies similar to LightRAG (Guo et al., 2024) and GraphRAG (Edge et al., 2024). The extraction process operates directly on textual content x_j where $t_j = \text{text}$, leveraging named entity recognition and relation extraction techniques to identify entities and their semantic relationships. Given the rich semantic information inherent in textual content, multimodal context integration is not required for this component. The resulting text-based knowledge graph captures explicit knowledge and semantic connections present in textual portions of documents, complementing the multimodal graph’s cross-modal grounding capabilities.

2.2.2 GRAPH FUSION AND INDEX CREATION

The separate cross-modal and text-based knowledge graphs capture complementary aspects of document semantics. Integrating them creates a unified representation leveraging visual-textual associations and fine-grained textual relationships for enhanced retrieval.

- (i) **Entity Alignment and Graph Fusion.** To create a unified knowledge representation, we merge the multimodal knowledge graph (\tilde{V}, \tilde{E}) and text-based knowledge graph through entity alignment. This process uses entity names as primary matching keys to identify semantically equivalent entities across both graph structures. The integration consolidates their representations, creating a comprehensive knowledge graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. This graph captures both multimodal contextual relationships and text-based semantic connections. The merged graph provides a holistic view of the document collection. This enables effective retrieval by leveraging visual-textual associations from the multimodal graph and fine-grained textual knowledge relationships from the text-based graph.
- (ii) **Dense Representation Generation.** To enable efficient similarity-based retrieval, we construct a comprehensive embedding table \mathcal{T} that encompasses all components generated during the indexing process. We encode dense representations for all graph entities, relationships, and atomic content chunks across modalities using an appropriate encoder. This creates a unified embedding space where each component $s \in \text{entities, relations, chunks}$ is mapped to its corresponding dense representation:

$$\mathcal{T} = \text{emb}(s) : s \in \mathcal{V} \cup \mathcal{E} \cup c_{j,j}, \quad (5)$$

where $\text{emb}(\cdot)$ denotes the embedding function tailored for each component type. Together, the unified knowledge graph \mathcal{G} and the embedding table \mathcal{T} constitute the complete retrieval index $\mathcal{I} = (\mathcal{G}, \mathcal{T})$. This provides both structural knowledge representation and dense vector space for efficient cross-modal similarity search during the subsequent retrieval stage.

2.3 CROSS-MODAL HYBRID RETRIEVAL

The retrieval stage operates on the index $\mathcal{I} = (\mathcal{G}, \mathcal{T})$ to identify relevant knowledge components for a given user query. Traditional RAG methods face significant limitations when dealing with multimodal documents. They typically rely on semantic similarity within single modalities and fail to capture the rich interconnections between visual, mathematical, tabular, and textual elements. To address these challenges, our framework introduces a cross-modal hybrid retrieval mechanism. This mechanism leverages structural knowledge and semantic representations across heterogeneous modalities.

Modality-Aware Query Encoding. Given a user query q , we first perform modality-aware query analysis to extract lexical cues and potential modality preferences embedded within the query. For instance, queries containing terms such as "figure," "chart," "table," or "equation" provide explicit signals about the expected modality of relevant information. We then compute a unified text embedding e_q using the same encoder employed during indexing, ensuring consistency between

query and knowledge representations. This embedding-based approach enables cross-modal retrieval capabilities where textual queries can effectively access multimodal content through their shared representations, maintaining retrieval consistency while preserving cross-modal accessibility.

Hybrid Knowledge Retrieval Architecture. Recognizing that knowledge relevance manifests through both explicit structural connections and implicit semantic relationships, we design a hybrid retrieval architecture that strategically combines two complementary mechanisms.

- (i) **Structural Knowledge Navigation.** This mechanism addresses the challenge of capturing explicit relationships and multi-hop reasoning patterns. Traditional keyword-based retrieval often fails to identify knowledge connected through intermediate entities or cross-modal relationships. To overcome this limitation, we exploit the structural properties encoded within our unified knowledge graph G . We employ keyword matching and entity recognition to locate relevant graph components. The retrieval process begins with exact entity matching against query terms.

We then perform strategic neighborhood expansion to include related entities and relationships within a specified hop distance. This structural approach proves particularly effective at uncovering high-level semantic connections and entity-relation patterns that span multiple modalities. It capitalizes on the rich cross-modal linkages established in our multimodal knowledge graph. The structural navigation yields candidate set $\mathcal{C}_{\text{stru}}(q)$ containing relevant entities, relationships, and their associated content chunks that provide comprehensive contextual information.

- (ii) **Semantic Similarity Matching.** This mechanism addresses the challenge of identifying semantically relevant knowledge that lacks explicit structural connections. While structural navigation excels at following explicit relationships, it may miss relevant content that is semantically related but not directly connected in the graph topology. To bridge this gap, we conduct dense vector similarity search between the query embedding \mathbf{e}_q and all components stored in embedding table \mathcal{T} .

This approach encompasses atomic content chunks across all modalities, graph entities, and relationship representations, enabling fine-grained semantic matching that can surface relevant knowledge even when traditional lexical or structural signals are absent. The learned embedding space captures nuanced semantic relationships and contextual similarities that complement the explicit structural signals from the navigation mechanism. This retrieval pathway returns the top-k most semantically similar chunks $\mathcal{C}_{\text{seman}}(q)$ ranked by cosine similarity scores, ensuring comprehensive coverage of both structurally and semantically relevant knowledge.

Candidate Pool Unification. Both retrieval pathways may return overlapping candidates with differing relevance signals. This necessitates a principled approach to unify and rank results. Retrieval candidates from both pathways are unified into a comprehensive candidate pool: $\mathcal{C}(q) = \mathcal{C}_{\text{stru}}(q) \cup \mathcal{C}_{\text{seman}}(q)$. Simply merging candidates would ignore distinct evidence each pathway provides. It would fail to account for redundancy between retrieved content.

- (i) **Multi-Signal Fusion Scoring.** To address these challenges, we apply a sophisticated fusion scoring mechanism integrating multiple complementary relevance signals. These include structural importance derived from graph topology, semantic similarity scores from embedding space, and query-inferred modality preferences obtained through lexical analysis. This multi-faceted scoring approach ensures that final ranked candidates $\mathcal{C}^*(q)$ effectively balance structural knowledge relationships with semantic relevance while appropriately weighting different modalities based on query characteristics.

- (ii) **Hybrid Retrieval Integration.** The resulting hybrid retrieval mechanism enables our framework to leverage the complementary strengths of both knowledge graphs and dense representations. This provides comprehensive coverage of relevant multimodal knowledge for response generation.

2.4 FROM RETRIEVAL TO SYNTHESIS

Effective multimodal question answering requires preserving rich visual semantics while maintaining coherent grounding across heterogeneous knowledge sources. Simple text-only approaches lose crucial visual information, while naive multimodal methods struggle with coherent cross-modal integration. Our synthesis stage addresses these challenges by systematically combining retrieved multimodal knowledge into comprehensive, evidence-grounded responses.

- (i) **Building Textual Context.** Given the top-ranked retrieval candidates $\mathcal{C}^*(q)$, we construct a structured textual context. We concatenate textual representations of all retrieved components, includ-

Table 1: Statistics of Experimental Datasets.

Dataset	# Documents	# Avg. Pages	# Avg. Tokens	# Doc Types	# Questions
DocBench	229	66	46377	5	1102
MMLongBench	135	47.5	21214	7	1082

ing entity summaries, relationship descriptions, and chunk contents. The concatenation incorporates appropriate delimiters to indicate modality types and hierarchical origins. This approach ensures the language model can effectively parse and reason over heterogeneous knowledge components.

- (ii) **Recovering Visual Content.** For multimodal chunks corresponding to visual artifacts, we perform dereferencing to recover original visual content, creating $\mathcal{V}^*(q)$. This design maintains consistency with our unified embedding strategy. Textual proxies enable efficient retrieval while authentic visual content provides rich semantics necessary for sophisticated reasoning during synthesis.

The synthesis process jointly conditions on both the assembled comprehensive textual context and dereferenced visual artifacts using a vision-language model:

$$\text{Response} = \text{VLM}(q, \mathcal{P}(q), \mathcal{V}^*(q)), \quad (6)$$

where the VLM integrates information from query, textual context, and visual content. This unified conditioning enables sophisticated visual interpretation while maintaining grounding in retrieved evidence. The resulting responses are both visually informed and factually grounded.

3 EVALUATION

3.1 EXPERIMENTAL SETTINGS

Evaluation Datasets. We conduct comprehensive evaluations on two challenging multimodal Document Question Answering (DQA) benchmarks that reflect real-world complexity and diversity. DocBench (Zou et al., 2024) provides a rigorous testbed with 229 multimodal documents spanning five critical domains: Academia, Finance, Government, Laws, and News. The dataset includes 1,102 expert-crafted question-answer pairs. These documents are notably extensive, averaging 66 pages and approximately 46,377 tokens, which presents substantial challenges for long-context understanding.

MMLongBench (Ma et al., 2024) complements this evaluation by focusing specifically on long-context multimodal document comprehension. It features 135 documents across 7 diverse document types with 1,082 expert-annotated questions. Together, these benchmarks provide comprehensive coverage of the multimodal document understanding challenges that RAG-Anything aims to address. They ensure our evaluation captures both breadth across domains and depth in document complexity. Detailed dataset statistics and characteristics are provided in Appendix A.1.

Baselines. We compare RAG-Anything against the following methods for performance evaluation:

- **GPT-4o-mini:** A powerful multimodal language model with native text and image understanding capabilities. Its 128K token context window enables direct processing of entire documents. We evaluate this model as a strong baseline for long-context multimodal understanding.
- **LightRAG** (Guo et al., 2024): A graph-enhanced RAG system that integrates structured knowledge representation with dual-level retrieval mechanisms. It captures both fine-grained entity-relation information and broader semantic context, improving retrieval precision and response coherence.
- **MMGraphRAG** (Wan & Yu, 2025): A multimodal retrieval framework that constructs unified knowledge graphs spanning textual and visual content. This method employs spectral clustering for multimodal entity analysis and retrieves context along reasoning paths to guide generation.

Experimental Settings. In our experiments, we implement all baselines using GPT-4o-mini as the backbone LLM. Documents are parsed using MinerU (Wang et al., 2024) to extract text, images, tables, and equations for downstream RAG processing. For the retrieval pipeline, we employ the text-embedding-3-large model with 3072-dimensional embeddings. We use the bge-reranker-v2-m3 model for reranking. For graph-based RAG methods, we enforce a combined entity-and-relation token limit of 20,000 tokens and a chunk token limit of 12,000 tokens.

Table 2: Accuracy (%) on DocBench Dataset. Performance results with best scores highlighted in **dark blue** and second-best in **light blue**. Domain categories include Academia (Aca.), Finance (Fin.), Government (Gov.), Legal Documents (Law), and News Articles (News). Document types are categorized as Text-only (Txt.), Multimodal (Mm.), and Unanswerable queries (Una.).

Method	Domains					Types			Overall
	Aca.	Fin.	Gov.	Law.	News	Txt.	Mm.	Una.	
GPT-4o-mini	40.3	46.9	60.3	59.2	61.0	61.0	43.8	49.6	51.2
LightRAG	53.8	56.2	59.5	61.8	65.7	85.0	59.7	46.8	58.4
MMGraphRAG	64.3	52.8	64.9	40.0	61.5	67.6	66.0	60.5	61.0
RAGAnything	61.4	67.0	61.5	60.2	66.3	85.0	76.3	46.0	63.4

Table 3: Accuracy (%) on MMLongBench across different domains and overall performance. Best results are highlighted in **dark blue** and second-best in **light blue**. Domain categories include Research Reports/Introductions (Res.), Tutorials/Workshops (Tut.), Academic Papers (Acad.), Guidebooks (Guid.), Brochures (Broch.), Administration/Industry Files (Admin.), and Financial Reports (Fin.).

Method	Domains							Overall
	Res.	Tut.	Acad.	Guid.	Broch.	Admin.	Fin.	
GPT-4o-mini	35.5	44.0	24.6	33.1	29.5	46.8	31.1	33.5
LightRAG	40.8	34.1	36.2	39.4	41.0	44.4	38.3	38.9
MMGraphRAG	40.8	36.5	35.7	35.8	28.2	46.9	38.5	37.7
RAGAnything	46.6	43.5	38.7	43.9	34.0	45.7	43.6	42.8

Outputs are constrained to a one-sentence format. For the baseline GPT-4o-mini in our QA scenario, documents are concatenated into image form with a maximum of 50 pages per document, rendered at 144 dpi. Finally, all query results are evaluated for accuracy by GPT-4o-mini.

3.2 PERFORMANCE COMPARISON

Superior Performance and Cross-Domain Generalization. RAG-Anything demonstrates superior overall performance over baselines through its unified multimodal framework. Unlike LightRAG, which is restricted to text-only content processing, RAG-Anything treats text, images, tables, and equations as first-class entities. MMGraphRAG only adds basic image processing while treating tables and equations as plain text, missing crucial structural information. RAG-Anything introduces a comprehensive dual-graph construction strategy that preserves structural relationships across all modalities. This unified approach enables superior performance across both evaluation benchmarks.

Enhanced Long-Context Performance. RAG-Anything demonstrates superior performance on long-context documents. The framework excels where relevant evidence is dispersed across multiple modalities and sections. It achieves the best results in information-dense domains such as Research Reports and Financial Reports on MMLongBench. These improvements stem from the structured context injection mechanism. This mechanism integrates dual-graph construction for cross-page entity alignment. It combines semantic retrieval with structural navigation. The framework also employs modality-aware processing for efficient context window utilization. Unlike baselines that cannot uniformly process diverse modalities, RAG-Anything effectively captures scattered multimodal evidence. Its cross-modal hybrid retrieval architecture combines structural knowledge navigation with semantic similarity matching. This enables the framework to leverage both explicit relationships and implicit semantic connections across modalities.

To systematically evaluate model performance across varying document lengths, we conducted comprehensive experiments on both datasets. As illustrated in Figure 2, RAG-Anything and MMGraphRAG exhibit comparable performance on shorter documents. However, RAG-Anything’s advantages become increasingly pronounced as document length grows. On DocBench, the performance gap expands dramatically to over 13 points for documents exceeding 100 pages (68.2% vs.

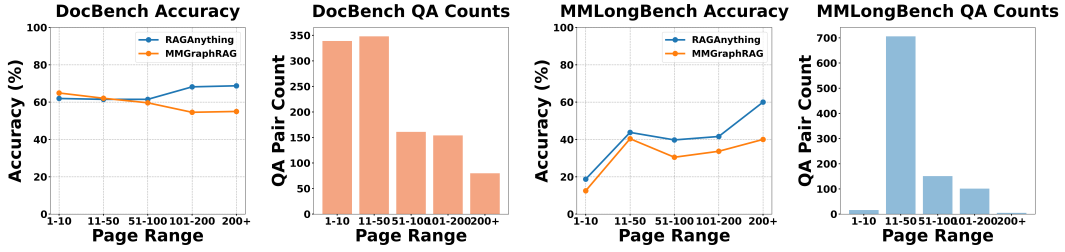


Figure 2: Performance evaluation across documents of varying lengths.

Table 4: Ablation study results on DocBench. The “Chunk-only” variant bypasses dual-graph construction and relies solely on traditional chunk-based retrieval, while “w/o Reranker” eliminates cross-modal reranking but preserves the core graph-based architecture.

Method	Domains					Types			Overall
	Aca.	Fin.	Gov.	Law.	News	Txt.	Mm.	Una.	
Chunk-only	55.8	61.5	60.1	60.7	64.0	81.6	66.2	43.5	60.0
w/o Reranker	60.9	63.5	58.8	60.2	68.6	81.7	74.7	45.4	62.4
RAGAnything	61.4	67.0	61.5	60.2	66.3	85.0	76.3	46.0	63.4

54.6% for 101–200 pages; 68.8% vs. 55.0% for 200+ pages). On MMLongBench, RAG-Anything demonstrates consistent improvements across all length categories, achieving accuracy gains of 3.4 points for 11–50 pages, 9.3 points for 51–100 pages, and 7.9 points for 101–200 pages. These findings confirm that our dual-graph construction and cross-modal hybrid retrieval mechanism is particularly effective for long-document reasoning tasks.

3.3 ARCHITECTURAL VALIDATION WITH ABLATION STUDIES

To isolate and quantify the contributions of key architectural components in RAG-Anything, we conducted systematic ablation studies examining two critical design choices. Given that our approach fundamentally differs from existing methods through dual-graph construction and hybrid retrieval, we specifically evaluated: i) **Chunk-only**, which bypasses graph construction entirely and relies solely on traditional chunk-based retrieval, and ii) **w/o Reranker**, which eliminates the cross-modal reranking component while preserving the core graph-based architecture.

As demonstrated in Table 4, the results validate our architectural design through striking performance variations. • **Graph Construction is Essential.** The chunk-only variant achieves merely 60.0% accuracy with substantial cross-domain drops. This demonstrates that traditional chunking fails to capture structural and cross-modal relationships essential for multimodal documents. • **Reranking Provides Marginal Gains.** Removing the reranker yields only a modest decline to 62.4%, while the full model achieves 63.4% accuracy. This indicates that cross-modal reranking provides valuable refinement, but primary gains stem from our graph-based retrieval and cross-modal integration.

3.4 CASE STUDIES

Multimodal documents contain rich structural information within each modality. Understanding these *intra-modal* structures is crucial for accurate reasoning. We analyze two representative cases from DocBench to demonstrate how RAG-Anything leverages these structures. These cases highlight a key limitation of existing methods. Baselines either rely on superficial textual cues or flatten complex visual elements into plain text. In contrast, RAG-Anything builds modality-aware graphs that preserve essential relationships (*e.g.*, table header↔cell↔unit edges; panel↔caption↔axis edges). This enables precise reasoning over complex document layouts.

• **Case 1: Multi-panel Figure Interpretation.** This case examines a common scenario in academic literature. Researchers often need to compare results across different experimental conditions. These results are typically presented in multi-panel visualizations. Figure 3 shows a challenging t-SNE

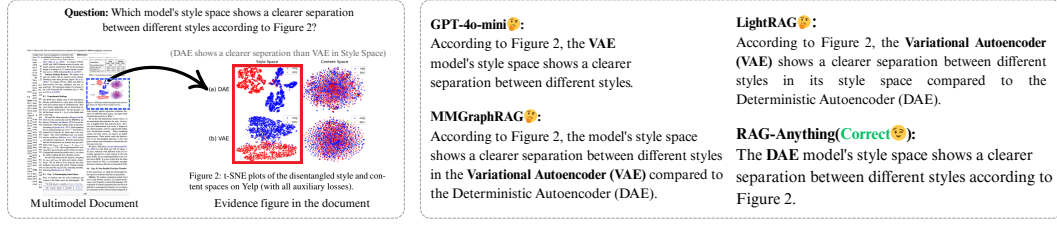


Figure 3: Multi-panel figure interpretation case. The query requires identifying cluster separation patterns from the style-space panel, while avoiding confusion from the adjacent content-space panel.

visualization with multiple subpanels. The query requires distinguishing between two related but distinct panels. RAG-Anything constructs a visual-layout graph where panels, axis titles, legends, and captions become nodes. Key edges encode semantic relationships. Panels contain specific plots. Captions provide contextual information. Subfigures relate hierarchically. This structure guides the retriever to focus on the *style-space* panel for comparing cluster separation patterns. The system avoids confusion from the adjacent content space panel. This panel shows less clear distinctions.

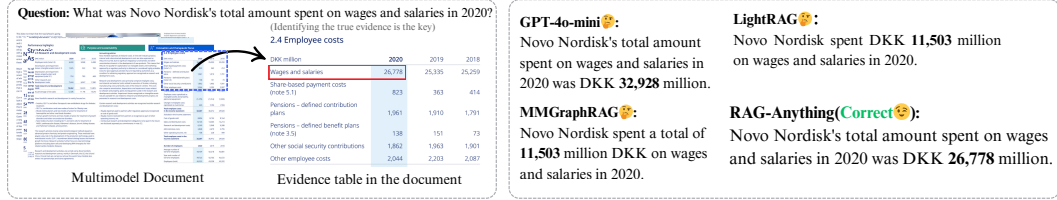


Figure 4: Financial table navigation case. The query involves locating the specific intersection of “Wages and salaries” row and “2020” column amid similar terminological entries.

• **Case 2: Financial Table Navigation.** This case addresses a common challenge in financial document analysis. Analysts must extract specific metrics from tables with similar terminology and multiple time periods. Figure 4 shows this scenario. The query involves resolving ambiguous financial terms and selecting the correct column for a specified year.

RAG-Anything transforms the financial report table into a structured graph. Each row header, column header (year), data cell, and unit becomes a node. The edges capture key relationships: *row-of*, *column-of*, *header-applies-to*, and *unit-of*. This structure enables precise navigation. The retriever focuses on the row “Wages and salaries” and the column for “2020”. It directs attention to the target cell (26,778 million). The system successfully disambiguates nearby entries like “Share-based payments.” Competing methods treat tables as linear text. They often confuse numerical spans and years. This leads to significantly inaccurate answers. RAG-Anything explicitly models relationships within the table. It achieves precise selection and numeric grounding. This ensures accurate responses.

• **Key Insights.** Both cases demonstrate how RAG-Anything’s structure-aware design delivers targeted advantages. Our approach transforms documents into explicit graph representations. These graphs capture intra-modal relationships that traditional methods miss. In figures, connections between panels, captions, and axes enable panel-level comparisons. This goes beyond keyword matching. In tables, row-column-unit graphs ensure accurate identification through modeling.

This structure-aware retrieval design reduces confusion from repeated terminology and complex layouts. Traditional RAG systems struggle with these scenarios due to lack of structural understanding. Even MMGraphRAG fails here because it only considers image modality entities. It ignores other modality entities like table cells, row headers, and column headers. RAG-Anything’s comprehensive graph representation captures all modality-specific entities and their relationships. This enables precise, modality-specific grounding that leads to consistent improvements in document Q&A tasks requiring fine-grained localization. Additional cases are available in Appendix A.2.

4 RELATED WORK

• **Graph-Enhanced Retrieval-Augmented Generation.** Large language models struggle with long-context inputs and multi-hop queries, failing to precisely locate dispersed evidence (Zhang et al.,

2025). Graph structures address this limitation by introducing explicit relational modeling, improving both retrieval efficiency and reasoning accuracy (Bei et al., 2025).

Since GraphRAG (Edge et al., 2024), research has evolved along two complementary directions. First, graph construction approaches optimize structures for retrieval efficiency, ranging from LightRAG’s (Guo et al., 2024) sparsified indices to neural models like GNN-RAG (Mavromatis & Karypis, 2024) and memory-augmented variants like HippoRAG (Jimenez Gutierrez et al., 2024). Second, knowledge aggregation approaches integrate information for multi-level reasoning through hierarchical methods like RAPTOR (Sarhi et al., 2024) and ArchRAG (Wang et al., 2025). Despite these advances, existing systems remain text-centric with homogeneous inputs. This limits their applicability to multimodal documents and constrains robust reasoning over heterogeneous content. RAG-Anything addresses this gap by extending GraphRAG to all modalities.

• **Multimodal Retrieval-Augmented Generation.** Multimodal RAG represents a natural evolution from text-based RAG systems, addressing the need to integrate external knowledge from diverse data modalities for comprehensive response generation (Abootorabi et al., 2025). However, current approaches are fundamentally constrained by their reliance on modality-specific architectures. Existing methods demonstrate these constraints across domains: VideoRAG (Ren et al., 2025) employs dual-channel architectures for video understanding while MM-VID (Lin et al., 2023) converts videos to text, losing visual information; VisRAG (Yu et al., 2025) preserves document layouts as images but misses granular relationships; MMGraphRAG (Wan & Yu, 2025) links scene graphs with textual representations but suffers from structural blindness—treating tables and formulas as plain text without proper entity extraction, losing structural information for reasoning.

The fundamental problem underlying these limitations is architectural fragmentation. Current systems require specialized processing pipelines for each modality. This creates poor generalizability as new modalities demand custom architectures and fusion mechanisms. Such fragmentation introduces cross-modal alignment difficulties, modality biases, and information bottlenecks. These issues systematically compromise system performance and scalability. RAG-Anything addresses this fragmentation through a unified graph-based framework. Our approach processes all modalities with consistent structured modeling. This eliminates architectural constraints while preserving multimodal information integrity. The result is seamless cross-modal reasoning across heterogeneous content.

5 CONCLUSION

RAG-Anything introduces a paradigm shift in multimodal retrieval through its unified graph-based framework. Our core technical innovation is the dual-graph construction strategy that seamlessly integrates cross-modal and text-based knowledge graphs. Rather than forcing diverse modalities into text-centric pipelines that lose critical structural information, our approach fundamentally reconceptualizes multimodal content as interconnected knowledge entities with rich semantic relationships. The hybrid retrieval mechanism strategically combines structural navigation with semantic matching, enabling precise reasoning over complex document layouts. Comprehensive evaluation demonstrates superior performance on long-context documents, particularly those exceeding 100 pages where traditional methods fail. This work establishes a new foundation for multimodal RAG systems that can handle the heterogeneous nature of diverse information landscapes.

Our analysis in Appendix A.5 reveals critical challenges facing current multimodal RAG systems. Two fundamental issues emerge through systematic failure case examination. First, systems exhibit text-centric retrieval bias, preferentially accessing textual sources even when queries explicitly require visual information. Second, rigid spatial processing patterns fail to adapt to non-standard document layouts. These limitations manifest in cross-modal misalignment scenarios and structurally ambiguous tables. The findings highlight the need for adaptive spatial reasoning and layout-aware parsing mechanisms to handle real-world multimodal document complexity.

REFERENCES

- Mohammad Mahdi Abootorabi, Amirhosein Zobeiri, Mahdi Dehghani, Mohammadali Mohammadkhani, Bardia Mohammadi, Omid Ghahroodi, Mahdieh Soleymani Baghshah, and Ehsaneddin Asgari. Ask in any modality: A comprehensive survey on multimodal retrieval-augmented generation. *arXiv preprint arXiv:2502.08826*, 2025.
- Yuanchen Bei, Weizhi Zhang, Siwen Wang, Weizhi Chen, Sheng Zhou, Hao Chen, Yong Li, Jiajun Bu, Shirui Pan, Yizhou Yu, et al. Graphs meet ai agents: Taxonomy, progress, and future opportunities. *arXiv preprint arXiv:2506.18019*, 2025.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitansky, Robert Osazuwa Ness, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. Lightrag: Simple and fast retrieval-augmented generation. *arXiv preprint arXiv:2410.05779*, 2024.
- Bernal Jimenez Gutierrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. Hipporag: Neurobiologically inspired long-term memory for large language models. *NeurIPS*, 37:59532–59569, 2024.
- Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin, Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin Liang, Zicheng Liu, Yumao Lu, Ce Liu, and Lijuan Wang. Mm-vid: Advancing video understanding with gpt-4v(ision). *arXiv preprint arXiv:2310.19773*, 2023.
- Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, et al. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. *Advances in Neural Information Processing Systems*, 37:95963–96010, 2024.
- Costas Mavromatis and George Karypis. Gnn-rag: Graph neural retrieval for large language model reasoning. *arXiv preprint arXiv:2405.20139*, 2024.
- Xubin Ren, Lingrui Xu, Long Xia, Shuaiqiang Wang, Dawei Yin, and Chao Huang. Videorag: Retrieval-augmented generation with extreme long-context videos. *arXiv preprint arXiv:2502.01549*, 2025.
- Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D Manning. Raptor: Recursive abstractive processing for tree-organized retrieval. In *The Twelfth International Conference on Learning Representations*, 2024.
- Xueyao Wan and Hang Yu. Mmgraphrag: Bridging vision and language with interpretable multimodal knowledge graphs. *arXiv preprint arXiv:2507.20804*, 2025.
- Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. Mineru: An open-source solution for precise document content extraction. *arXiv preprint arXiv:2409.18839*, 2024.
- Shu Wang, Yixiang Fang, Yingli Zhou, Xilin Liu, and Yuchi Ma. Archrag: Attributed community-based hierarchical retrieval-augmented generation. *arXiv preprint arXiv:2502.09891*, 2025.
- Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Junhao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang, Xu Han, Zhiyuan Liu, and Maosong Sun. Visrag: Vision-based retrieval-augmented generation on multi-modality documents. *arXiv preprint arXiv:2410.10594*, 2025.
- Qinggang Zhang, Shengyuan Chen, Yuanchen Bei, Zheng Yuan, Huachi Zhou, Zijin Hong, Hao Chen, Yilin Xiao, Chuang Zhou, Yi Chang, and Xiao Huang. A survey of graph retrieval-augmented generation for customized large language models. *arXiv preprint arXiv:2501.13958*, 2025.
- Anni Zou, Wenhao Yu, Hongming Zhang, Kaixin Ma, Deng Cai, Zhuosheng Zhang, Hai Zhao, and Dong Yu. Docbench: A benchmark for evaluating llm-based document reading systems. *arXiv preprint arXiv:2407.10701*, 2024.

A APPENDIX

This appendix provides comprehensive supporting materials for our experimental evaluation and implementation details. Section A.1 presents detailed dataset statistics for the DocBench and MMLongBench multi-modal benchmarks, including document type distributions and complexity metrics. Section A.2 showcases additional case studies that demonstrate RAG-Anything’s structure-aware capabilities across diverse multimodal content understanding tasks. Section A.3 documents the complete set of multimodal analysis prompts for vision, table, and equation processing that enable context-aware interpretation. Section A.4 provides the standardized accuracy evaluation prompt used for consistent response assessment across all experimental conditions.

A.1 DATASET CHARACTERISTICS AND STATISTICS

Table 5: Document type distribution and statistics for the DocBench benchmark.

Type	Acad.	Fin.	Gov.	Law.	News
# Docs	49	40	44	46	50
# Questions	303	288	148	191	172
Avg. Pages	11	192	69	58	1

Table 6: Document type distribution and statistics for the MMLongBench benchmark.

Type	Res.	Tut.	Acad.	Guid.	Broch.	Admin.	Fin.
# Docs	34	17	26	22	15	10	11
# Questions	292	138	199	155	100	81	117
Avg. Pages	39	58	35	78	30	17	87

Tables 5 and 6 present the distribution of document types across the DocBench and MMLongBench benchmarks. • **DocBench** encompasses medium- to long-length documents spanning various domains, including legal, governmental, and financial files. Financial reports represent the most extensive category, averaging 192 pages per document, while the News category consists of concise single-page newspapers. • **MMLongBench** demonstrates a broader spectrum of document types and lengths. Research reports, tutorials, and academic papers maintain moderate lengths of 35–58 pages on average, while guidebooks extend to approximately 78 pages. Brochures and administrative files remain relatively compact, whereas financial reports again emerge as the longest category.

Collectively, these two benchmarks provide comprehensive coverage ranging from brief news articles to extensive technical and financial documentation. This establishes diverse and challenging evaluation contexts for multimodal document understanding tasks.

A.2 ADDITIONAL CASE STUDIES

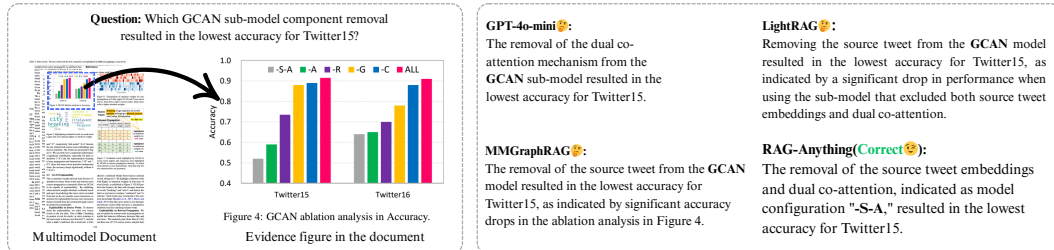


Figure 5: Visual reasoning case. RAG-Anything correctly identifies "-S-A" as the lowest accuracy configuration, while baselines misinterpret spatial relationships.

• **Demonstrating Visual Reasoning Capabilities.** Figure 5 illustrates how RAG-Anything handles complex visual reasoning tasks involving chart interpretation. The query asks which GCAN sub-model component removal yields the lowest accuracy on Twitter15. Traditional approaches struggle

with spatial relationships between visual elements. RAG-Anything addresses this challenge by constructing a structured graph representation of the bar plot. Bars, axis labels, and legends become interconnected nodes. These are linked by semantic relations such as *bar-of* and *label-applies-to*.

This graph-based approach enables precise alignment between visual and textual elements. The system correctly identifies the bar labeled "-S-A" (removing source tweet embeddings and dual co-attention) and its corresponding accuracy value as the lowest performer. Baseline methods that flatten visual information often misinterpret spatial relationships. They frequently conflate nearby components. RAG-Anything’s structured representation preserves critical visual-textual associations. This leads to accurate query resolution and proper attribution of performance drops to "-S-A".

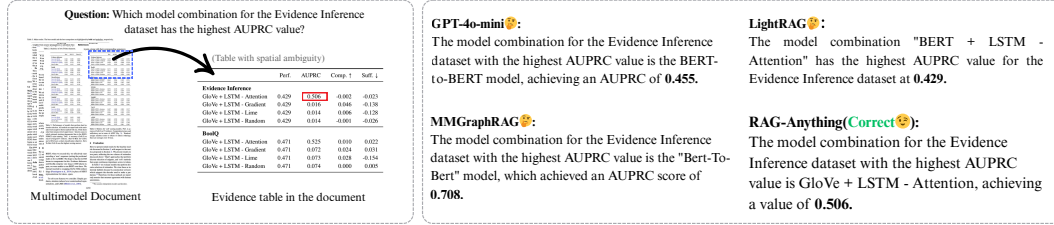


Figure 6: Tabular navigation case. RAG-Anything locates the highest AUPRC value (0.506), while the compared approaches struggle with structural ambiguity.

• **Handling Complex Tabular Structures.** Figure 6 showcases RAG-Anything’s ability to navigate intricate tabular data where structural disambiguation is crucial. The query seeks the model combination achieving the highest AUPRC value for the Evidence Inference dataset—a task complicated by repeated row labels across multiple datasets within the same table. This scenario highlights a fundamental limitation of conventional approaches that struggle with structural ambiguity in data.

RAG-Anything overcomes this by parsing the table into a comprehensive relational graph where headers and data cells become nodes connected through explicit *row-of* and *column-of* relationships. This structured representation enables the system to correctly isolate the Evidence Inference dataset context and identify "GloVe + LSTM – Attention" with a score of 0.506 as the optimal configuration. By explicitly preserving hierarchical table constraints that other methods often collapse or misinterpret, RAG-Anything ensures reliable reasoning across complex multi-dataset tabular structures.

A.3 CONTEXT-AWARE MULTIMODAL PROMPTING

These three prompts orchestrate structured, context-aware multimodal analysis with JSON-formatted outputs. They systematically guide the model to extract comprehensive descriptions of visual, tabular, and mathematical content while maintaining explicit alignment with surrounding information.

Vision Analysis Prompt. Figure 7 orchestrates comprehensive image-context integration. The prompt directs the model to systematically capture compositional elements, object relationships, visual attributes, stylistic features, dynamic actions, and technical components (e.g., charts), while establishing explicit connections to accompanying text. This approach transcends superficial description, enabling contextually-grounded interpretations that enhance knowledge retrieval and substantiation.

Table Analysis Prompt. Figure 8 structures systematic tabular content decomposition across multiple analytical dimensions: structural organization, column semantics, critical values, statistical patterns, and contextual relevance. Through precise terminology and numerical accuracy requirements, the prompt eliminates ambiguous generalizations and ensures faithful preservation of key indicators while maintaining coherent alignment with surrounding discourse.

Equation Analysis Prompt. Figure 9 prioritizes semantic interpretation over syntactic restatement of mathematical expressions. The prompt instructs comprehensive analysis of variable definitions, operational logic, theoretical foundations, inter-formula relationships, and practical applications. This methodology ensures mathematical content becomes integral to broader argumentative frameworks, supporting enhanced retrieval accuracy, analytical traceability, and reasoning coherence.

Vision Analysis Prompt

```

1 Please analyze this image in detail, considering the surrounding context. Provide a JSON response with the
2 following structure:
3
4 {
5   "detailed_description": "A comprehensive and detailed visual description of the image following these
6   guidelines:
7   - Describe the overall composition and layout
8   - Identify all objects, people, text, and visual elements
9   - Explain relationships between elements and how they relate to the surrounding context
10  - Note colors, lighting, and visual style
11  - Describe any actions or activities shown
12  - Include technical details if relevant (charts, diagrams, etc.)
13  - Reference connections to the surrounding content when relevant
14  - Always use specific names instead of pronouns",
15  "entity_info": {
16    "entity_name": "{entity_name}",
17    "entity_type": "image",
18    "summary": "concise summary of the image content, its significance, and relationship to surrounding content
19    (max 100 words)"
20  }
21 }
22
23 Context from surrounding content:
24 {context}
25
26 Image details:
27 - Image Path: {image_path}
28 - Captions: {captions}
29 - Footnotes: {footnotes}
30
31 Focus on providing accurate, detailed visual analysis that incorporates the context and would be useful for
32 knowledge retrieval.

```

Figure 7: Vision analysis prompt for context-aware image interpretation and knowledge extraction.

Table Analysis Prompt

```

1 Please analyze this table content considering the surrounding context, and provide a JSON response with the
2 following structure:
3
4 {
5   "detailed_description": "A comprehensive analysis of the table including:
6   - Table structure and organization
7   - Column headers and their meanings
8   - Key data points and patterns
9   - Statistical insights and trends
10  - Relationships between data elements
11  - Significance of the data presented in relation to surrounding context
12  - How the table supports or illustrates concepts from the surrounding content
13  Always use specific names and values instead of general references.",
14  "entity_info": {
15    "entity_name": "{entity_name}",
16    "entity_type": "table",
17    "summary": "concise summary of the table's purpose, key findings, and relationship to surrounding content (max
18    100 words)"
19  }
20 }
21
22 Context from surrounding content:
23 {context}
24
25 Table Information:
26 Image Path: {table_img_path}
27 Caption: {table_caption}
28 Body: {table_body}
29 Footnotes: {table_footnote}
30
31 Focus on extracting meaningful insights and relationships from the tabular data in the context of the
32 surrounding content.

```

Figure 8: Table analysis prompt for structured content decomposition and semantic understanding.

Equation Analysis Prompt

```

1 Please analyze this mathematical equation considering the surrounding context, and provide a JSON response
2 with the following structure:
3
4 {
5   "detailed_description": "A comprehensive analysis of the equation including:
6     - Mathematical meaning and interpretation
7     - Variables and their definitions in the context of surrounding content
8     - Mathematical operations and functions used
9     - Application domain and context based on surrounding material
10    - Physical or theoretical significance
11    - Relationship to other mathematical concepts mentioned in the context
12    - Practical applications or use cases
13    - How the equation relates to the broader discussion or framework
14    Always use specific mathematical terminology.",
15   "entity_info": {
16     "entity_name": "{entity_name}",
17     "entity_type": "equation",
18     "summary": "concise summary of the equation's purpose, significance, and role in the surrounding context (max
19       100 words)"
20   }
21 }
22
23 Context from surrounding content:
24 {context}
25
26 Equation Information:
27 Equation: {equation_text}
28 Format: {equation_format}
29
30 Focus on providing mathematical insights and explaining the equation's significance within the broader
31 context.

```

Figure 9: Equation analysis prompt for mathematical expression interpretation and integration.

Accuracy Evaluation Prompt

```

1 You are an expert evaluator tasked with assessing the accuracy of answers generated by a RAG
2 (Retrieval-Augmented Generation) system.
3
4 **Task**: Evaluate whether the generated answer correctly responds to the given question based on the expected
5 answer.
6
7 **Question**: {question}
8
9 **Expected Answer**: {expected_answer}
10
11 **Generated Answer**: {generated_answer}
12
13
14 **Evaluation Criteria**:
15 1. **Accuracy (0 or 1)**: Does the generated answer match the factual content of the expected answer?
16 - 1: The generated answer is factually correct and aligns with the expected answer
17 - 0: The generated answer is factually incorrect or contradicts the expected answer
18
19 **Instructions**:
20 - Focus on factual correctness, not writing style or format
21 - Consider partial matches: If the generated answer contains the correct information but includes additional
22 context, it should still be considered accurate
23 - For numerical answers, check if the values match or are equivalent
24 - For list answers, check if all key elements are present
25 - If the expected answer is "Not answerable" and the generated answer indicates inability to answer, consider
26 it accurate
27
28 **Output Format**:
29 Please respond with a JSON object containing only:
30 {
31   "accuracy": 0 or 1,
32   "reasoning": "Brief explanation of your evaluation"
33 }

```

Figure 10: Accuracy evaluation prompt for consistent factual assessment across question types.

A.4 ACCURACY EVALUATION PROMPT DESIGN

Figure 10 presents the standardized prompt specifically designed for systematic factual accuracy assessment of generated responses across multiple domains. The prompt establishes explicit evaluation criteria that prioritize content correctness over stylistic considerations, producing binary accuracy

classifications accompanied by concise analytical justifications. All accuracy evaluations throughout our comprehensive experimental framework were conducted using **GPT-4o-mini**, ensuring consistent and reliable assessment standards across diverse question categories and specialized domains.

A.5 CHALLENGES AND FUTURE DIRECTIONS FOR MULTI-MODAL RAG

While current multimodal RAG systems demonstrate promising capabilities, their limitations emerge most clearly through systematic analysis of failure cases. Understanding where and why these systems break down is crucial for advancing the field beyond current performance plateaus. Examining failure patterns helps identify fundamental architectural bottlenecks and design principles for more robust multimodal systems. Our investigation reveals two critical failure patterns exposing deeper systemic issues in multimodal RAG architectures. These patterns are not merely edge cases but reflect fundamental challenges in cross-modal information integration and structural reasoning:

- **Text-Centric Retrieval Bias:** Systems exhibit strong preference for textual sources, even when queries explicitly demand visual information. This reveals inadequate cross-modal attention.
- **Document Structure Processing Challenges:** Systems struggle with complex layouts and non-linear information flows. This exposes limitations in spatial reasoning and contextual understanding.

These failure modes illuminate key insights about current multimodal AI. They provide concrete directions for architectural innovations that could substantially improve system robustness.

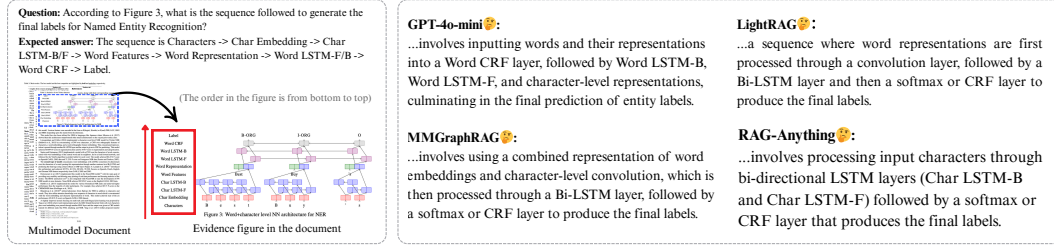


Figure 11: Cross-modal noise case. All methods fail to retrieve the correct answer from the specified image, instead retrieving noisy textual evidence that misaligns with the structured visual content.

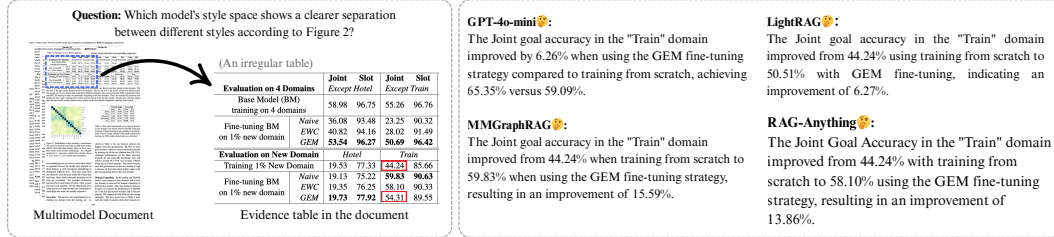


Figure 12: Ambiguous table structure case. All methods fail to correctly parse the confusing table layout with merged cells and unclear column boundaries, leading to incorrect data extraction.

Case 1: Cross-Modal Misalignment. Figure 11 presents a particularly revealing failure scenario where all evaluated methods consistently produce incorrect answers despite having access to the necessary information. This universal failure across different architectures suggests fundamental limitations in how current systems handle noisy, heterogeneous multimodal data—a critical challenge as real-world applications inevitably involve imperfect, inconsistent information sources. The failure exposes two interconnected systemic issues that compound each other:

Issue 1: Retrieval Bias Toward Text. Current RAG systems demonstrate pronounced bias toward textual passages. This occurs particularly when visual content lacks exact keyword matches. The bias persists even when queries contain explicit instructions to prioritize visual sources. This reveals a fundamental weakness in cross-modal attention mechanisms.

The retrieved textual information, while topically related, often operates at a different granularity level than visual content. Images may contain precise, structured data such as specific numerical values,

detailed diagrams, or exact spatial relationships. Corresponding text typically provides general, conceptual descriptions. This semantic misalignment introduces noise that actively misleads the reasoning process. The system attempts to reconcile incompatible levels of detail and specificity.

Issue 2: Rigid Spatial Processing Patterns. Current visual processing models exhibit fundamental rigidity in spatial interpretation. Most systems default to sequential scanning patterns—top-to-bottom and left-to-right—that mirror natural reading conventions. While effective for simple text documents, this approach creates systematic failures with structurally complex real-world content. Many documents require non-conventional processing strategies. Tables demand column-wise interpretation, technical diagrams follow specific directional flows, and scientific figures embed critical information in unexpectedly positioned annotations. These structural variations are prevalent in professional documents, making adaptive spatial reasoning essential.

In the observed failure case, the correct answer required integrating visual elements in reverse order from the model's default processing sequence. The system's inability to recognize and adapt to this structural requirement led to systematic misinterpretation. This represents a fundamental architectural limitation where spatial reasoning remains static regardless of document context or query intent. When spatial processing patterns are misaligned with document structure, the extracted information becomes not merely incomplete but actively misleading. This structural noise compounds other processing errors and can lead to confident but entirely incorrect conclusions.

Case 2: Structural Noise in Ambiguous Table Layouts. As shown in Figure 12, all methods failed when confronted with a structurally ambiguous table. The primary failure stems from the table's confusing design: the GEM row lacks dedicated cell boundaries, and the "Joint" and "Slot" columns merge without clear separation. These structural irregularities create parsing ambiguities that systematically mislead extraction algorithms. This failure pattern reveals a critical vulnerability in current RAG systems. When table structures deviate from standard formatting conventions—through merged cells, unclear boundaries, or non-standard layouts—extraction methods consistently misinterpret cell relationships and conflate distinct data values. This exposes the brittleness of current approaches when faced with real-world document variations that deviate from clean, structured formats.

The case highlights two essential directions for enhancing robustness. RAG systems require layout-aware parsing mechanisms that can recognize and adapt to structural irregularities rather than imposing rigid formatting assumptions. Additionally, integrating visual processing capabilities could significantly improve noise resilience, as visual models can leverage spatial relationships and contextual design cues that are lost in purely structural representations.