# DocHop-QA: Towards Multi-Hop Reasoning over Multimodal Document Collections

**Jiwon Park** [*,1] **Seohyun Pyeon** [*,1] **Jinwoo Kim**[*,1]
**Rina Carines Cabal** [2] **, Yihao Ding**[3]**, Soyeon Caren Han**[1,4†]

[1]Pohang University of Science and Technology, [2]The University of Sydney,
[3]The University of Western Australia, [4]The University of Melbourne
jiwon23, seohyun, jinwoo0327@postech.ac.kr,
rina.cabral@sydney.edu.au, yihao.ding@uwa.edu.au, caren.han@unimelb.edu.au

## Abstract

Despite recent advances in large language models (LLMs), most QA benchmarks are still confined to single-paragraph or single-document settings, failing to capture the complexity of real-world information-seeking tasks. Practical QA often requires multi-hop reasoning over information distributed across multiple documents, modalities, and structural formats. Although prior datasets made progress in this area, they rely heavily on Wikipedia-based content and unimodal plain text, with shallow reasoning paths that typically produce brief phrase-level or single-sentence answers, thus limiting their realism and generalizability. We propose DocHop-QA, a large-scale benchmark comprising 11,379 QA instances for multimodal, multi-document, multi-hop question answering. Constructed from publicly available scientific documents sourced from PubMed, DocHop-QA is domain-agnostic and incorporates diverse information formats, including textual passages, tables, and structural layout cues. Unlike existing datasets, DocHop-QA does not rely on explicitly hyperlinked documents; instead, it supports open-ended reasoning through semantic similarity and layout-aware evidence synthesis. To scale realistic QA construction, we designed an LLM-driven pipeline grounded in 11 high-frequency scientific question concepts. We evaluated DocHop-QA through four tasks spanning structured index prediction, generative answering, and multimodal integration, reflecting both discriminative and generative paradigms. These tasks demonstrate DocHop-QA's capacity to support complex, multimodal reasoning across multiple documents.

## 1 Introduction

Question Answering (QA) is a fundamental task in natural language processing, requiring systems to pinpoint precise responses to user queries across heterogeneous information sources. Recent advances in LLMs have markedly improved both the fluency and accuracy of QA systems, yet most benchmarks still emphasize single-hop inference over short text spans within a single document. In contrast, real-world information-seeking scenarios often demand complex multi-hop reasoning over fragmented data distributed across diverse formats and multiple documents (Yoon et al. 2022). Multi-hop QA addresses this challenge by requiring models to integrate and reason over two or more distinct pieces of information. In practice, answering questions such as "What proteins are involved in pathway X, and what experiments validate their roles?" often requires synthesizing content from methods, results, and tables in multiple scientific publications. These multi-hop reasoning tasks reflect fundamental challenges in real-world information-seeking scenarios, such as academic reports, websites, and policy briefs, where answers are not explicitly stated and must be inferred by integrating information from multiple sources. Despite recent advances, existing multi-hop QA datasets such as HotpotQA (Yang et al. 2018) and MuSiQue (Trivedi et al. 2022) exhibit several key limitations. They are predominantly sourced from Wikipedia, limit answers to short text spans, and restrict modalities to plain text. These constraints diverge from how information is retrieved in real-world documents, which are often multimodal and semi-structured, featuring visual elements (e.g., tables, figures) and layout cues (e.g., multi-column formats). Such structural and visual features play an essential role in human understanding, but are largely overlooked in current QA benchmarks, resulting in poor generalization of LLM-based QA models to more realistic settings. Moreover, the creation of QA datasets has traditionally relied on manual annotation, which is costly, time-consuming, and difficult to scale. While recent work such as LIQUID (Lee, Kim, and Kang 2023) proposed an automated framework for generating list-type QA pairs via summarization and answer filtering to address this bottleneck, it remains limited to single-document QA and does not extend to multi-document, multimodal QA. A scalable pathway to realistic QA remains a challenge.

To fill this gap, we introduce DocHop-QA, a large-scale benchmark for complex, multi-hop QA over multimodal, multi-document corpora, constructed from scientific articles. Its generation methodology is domain-agnostic and scalable, making it applicable to a wide range of scientific domains. Constructed from PubMed articles, it includes unstructured text, structured tables, and layout features, without relying on predefined hyperlinks or annotated reasoning chains. To efficiently generate diverse and challenging QA pairs, we present a scalable LLM-driven pipeline based

---

[*]These authors contributed equally.

[†]Corresponding author

on 11 real-world scientific reasoning concepts. To evaluate the applicability and flexibility of DocHop-QA, we benchmark its performance across four representative QA tasks: (1) BBox Entity Index Extraction, (2) XML Entity Index Extraction, (3) Structured Generative Answering, and (4) Generative Text Extraction. These tasks span both discriminative and generative paradigms and encompass structured prediction, multimodal understanding, and free-form answer generation. Our key contributions are as follows:

- We introduced **DocHop-QA**, a large-scale benchmark for multi-hop question answering over multimodal, multi-document scientific corpora, incorporating text, tables, and structural layout features.

- We developed a **scalable, LLM-driven generation pipeline** guided by 11 reasoning concepts inspired by real-world scientific inquiry, enabling the automated creation of diverse and challenging multi-hop QA instances without requiring gold chains or hyperlinks.

- We **evaluated DocHop-QA across four core QA tasks**, ranging from structured index prediction to generative answering, demonstrating its broad applicability for assessing both discriminative and generative reasoning over textual, structural, and visual information.

## 2 Related Work

**Multi-hop QA datasets** Multi-hop QA refers to tasks that span multiple documents, pages, and modalities, reflecting a more realistic and complex information-seeking process. Early datasets such as QAngaroo (Welbl, Stenetorp, and Riedel 2018) and HotpotQA (Yang et al. 2018) introduced multi-hop reasoning by requiring models to combine facts from multiple passages. These benchmarks laid foundational work for multi-hop inference, but remained limited to plain text and short-context settings. HybridQA (Chen et al. 2020) and WebQA (Chang et al. 2022) extended this space by incorporating semi-structured data like tables alongside text. However, these datasets primarily draw from Wikipedia or web-derived sources and lack the structural complexity and document variety observed in real-world applications. Recent efforts such as FanOutQA (Zhu et al. 2024), MEQA (Li et al. 2024), LLeQA (Louis, van Dijck, and Spanakis 2024) and Multi-Doc QA (Wang et al. 2024) have attempted to push the boundaries of multi-hop reasoning toward longer contexts and more naturalistic questions. Yet, these benchmarks remain largely unimodal, emphasize short answers (e.g., entity spans or binary responses), and continue to depend on manually curated pipelines that limit scalability and domain generality. Moreover, few of these datasets support deep reasoning over heterogeneous formats, such as text, tables, and layout cues, which are critical to real-world document understanding. In contrast, our work presents a domain-agnostic, multimodal, and multi-document QA benchmark that scales automatically using LLMs and supports rich, compositional reasoning across various document components and modalities.

**Visually Rich Document QA** Visually rich document (VRD) understanding has recently emerged as a key direction for extending QA beyond flat text to real-world doc-

| Dataset | M Hop | VRD | #QA | Q Gen | Source | M Doc | Modality |
|---|---|---|---|---|---|---|---|
| HotpotQA | O | X | 113K | H | Wiki | X | Tx |
| HybridQA | O | X | 70K | H | Wiki | O | Tx+Tab |
| Multi-Doc QA | O | X | 1.6K | LLM + H | Industry | O | Tx |
| DocVQA | X | O | 50K | H | Industry | X | Tx+Img |
| SlideVQA | X | O | 14K | H | Slideshare | X | Tx+Img |
| ECG-QA | X | O | 414K | H | PTB-XL | X | Tx+Img |
| SPIQA | X | O | 270K | LLM + H | arXiv | X | Tx+Tab+Img |
| MMVQA | X | O | 263K | LLM + H | PubMed | X | Tx+Tab+Img |
| **Ours** | O | O | 11K | LLM + H | PubMed | O | Tx+Tab+Img |

Table 1: Comparison of QA dataset benchmarks. M Hop: Multi Hop, VRD: Visually Rich Document, Q Gen.: Question Generation method (H: Human, LLM: Large Language Model), M Doc: Multiple Document, Modality (Tx: Text, Tab: Table, Img: Image).

uments that incorporate complex layouts, tables, and figures. Datasets such as DocVQA (Mathew, Karatzas, and Jawahar 2021) and InfographicVQA (Mathew et al. 2022) introduced visual question answering tasks grounded in single-page document images, often framed as span-based or retrieval problems. SlideVQA (Tanaka et al. 2023) introduced multi-hop reasoning across presentation slides, and BLIVA (Hu et al. 2024) suggested query embedding for multi-hop reasoning, but both datasets remained restricted to single-document settings with limited modality diversity. More recent benchmarks such as ECG-QA(Oh et al. 2023), SPIQA (Pramanick, Chellappa, and Venugopalan 2024) and MMVQA (Ding et al. 2024) incorporate scientific papers and richer modality mixes. However, these tasks still emphasize retrieval over deep reasoning, and largely lack support for multi-document or cross-page evidence composition. Additionally, most existing VRD QA datasets are manually constructed and tailored to narrow domains, limiting their scalability and flexibility for broader evaluation. Table 1 compares prior datasets with ours.

**Our Contribution.** To bridge these gaps, we introduce **DocHop-QA**, a new benchmark that enables end-to-end multi-hop reasoning over multimodal, visually rich documents spanning multiple pages and sources. Unlike prior datasets, DocHop-QA supports inference across unlinked documents using layout-aware features, semantic similarity, and compositional evidence[1] from text, tables, and structural formats. Constructed via an automated LLM-powered pipeline and grounded in 11 high-frequency question concepts, our dataset provides a scalable and realistic foundation for evaluating complex reasoning in next-generation multimodal QA systems.

## 3 Dataset Construction

We constructed DocHop-QA through a multi-stage pipeline designed to generate high-quality, multi-hop question-answer pairs grounded in real-world documents. The overall process involves four key stages: (1) proposing diverse question concepts, (2) selecting document pairs suitable

---

[1]For technical precision, we refer to the specific segments extracted from documents (including both textual passages and tables) as "snippets," while using "evidence" in the broader conceptual sense of information supporting multi-hop reasoning.

| Type | Subtype | Keyword | Q Type | Question Concept | Question Example |
|---|---|---|---|---|---|
| Paragraph-Oriented | Non-Comparison | **\<Problem\>\<sep\>\<Solution\>** | PS | What are the main challenges in [ ], and how can [ ] address these issues? | What are the main challenges in understanding the dual role of the transcription factor Sp3 in tumor progression, and how can the forkhead transcription factor Foxi1 address these issues by regulating the expression of key subunits of the vacuolar H+ ATPase? |
| | | **\<Problem\>\<sep\>\<Solution\>\<sep\>\<Mechanism\>** | PSM | What are the challenges in [ ], what solutions exist to address these challenges, and how do these solutions help resolve the issues? | What are the challenges in the relationship between leisure physical activity and the risk of osteoporotic fractures in men, what solutions exist to address these challenges, and how do these solutions help resolve the issues? |
| | | **\<Problem\>\<sep\>\<Solution\>\<sep\>\<Result/effect\>** | PSR | What are the main challenges in [ ], what strategies can be implemented to address them, and what are the observed effects of these strategies? | What are the main challenges in accurately diagnosing corticobasal syndrome, what strategies can be implemented to address them, and what are the observed effects of these strategies? |
| | | **\<Problem\>\<sep\>\<Solution\>\<sep\>\<Limitation\>** | PSL | What are the main challenges in [ ], and how can [ ] address these issues? Additionally, what are the limitations of this approach? | What are the main challenges in developing consumer health vocabulary (CHV) due to the heterogeneity and ambiguity of consumer expressions, and what is the solution to this challenge? Additionally, what are the limitations of this approach? |
| | | **\<Mechanism\>\<sep\>\<Effect\>** | ME | How does [ ] contribute to [ ], and what is its impact on [ ]? | How do soy isoflavones contribute to estrogen-like effects, and what is their impact on male breast cancer incidence? |
| | | **\<Studytopic\>\<sep\>\<Advantage\>\<sep\>\<Limitation\>** | SAL | What are the [ ], what advantages does it offer, and what limitations exist? | What are the Gene Ontology term predictions, what advantages does it offer, and what limitations exist? |
| | Comparison | **\<Feature\>\<sep\>\<Feature\>** | FF | How does [ ] compare to [ ], and what are the key similarities and differences? | How do the effects of inorganic arsenic exposure on children's intellectual function compare to those of fluoride exposure, and what are the key similarities and differences? |
| | | **\<Solution\>\<sep\>\<Solution\>** | SS | How do [ ] and [ ] differ in addressing [ ], and which approach is more effective? | How do traditional laboratory and clinical research methods and non-invasive abdominal recordings differ in addressing the prediction of preterm labor, and which approach is more effective? |
| Table-Oriented | Non-Refer Table | **\<Problem\>\<sep\>\<Solution\>\<sep\>\<Table\>** | PST | What are the main challenges in [ ], what strategies can be implemented to address them? Is there any table additionally describing it? Which table is it and what is the main point of the table? | What are the main challenges in regularized gene selection in cancer microarray meta-analysis, what strategies can be implemented to address them? Is there any table additionally describing it? Which table is it and what is the main point of the table? |
| | Refer Table | **\<Table\>\<sep\>\<Studytopic\>\<sep\>\<Limitation\>** | TSL | What is the main point of Table [ ] in document [ ], and what is the main topic of the referenced paper? Also, what limitations remain in the referenced paper? | What is the main point of table Quantitative synthesised results in given document, and what is the main topic of the referenced paper? Also, what limitations remain in the referenced paper? |
| | | **\<Table\>\<sep\>\<Studytopic\>** | TS | What is the main point of Table [ ] in document [ ], and what is the main topic of the referenced paper? | What is the main point of table Summary of included study characteristics in given document, and what is the main topic of the referenced paper? |

Table 2: Illustrative examples of the proposed question concepts, categorized by (1) **Type**, denoting the primary answer modality (Paragraph-Oriented vs. Table-Oriented); (2) **Subtype**, indicating the reasoning structure (e.g., comparison vs. non-comparison); (3) **Keyword**, representing the core reasoning elements and semantic placeholders used in question generation; (4) **Q Type**, a shorthand for the placeholder composition; (5) **Question Concept**, the abstract question structure using place-holders; and (6) **Question Example**, a concrete instantiation of each question concept. This design enables controlled, diverse, and semantically grounded multi-hop QA generation across heterogeneous document modalities.

for multi-hop reasoning, (3) generating questions and (4) matching answers. As the foundation of our dataset, we leverage the PubMed Central (PMC) open-access repository, which contains publicly available academic papers and provides structured XML and PDF versions of biomedical and life science literature. From this corpus, we collected approximately 2M full-text scientific articles and applied a combination of filtering heuristics and semantic similarity-based pairing strategies(Appendix A.1) to identify suitable document pairs for constructing multi-hop QA. After refinement, we retained around 20K document pairs, covering 9,250 unique documents, for question generation.

## 3.1 Question Concept Proposal

To enable systematic and scalable multi-hop question generation, we defined 11 structured question concepts based on frequently observed reasoning patterns in PubMed articles. We first manually constructed 90 seed QA pairs that required cross-document reasoning with the assistance of NotebookLM, and then conducted a qualitative analysis of their structure to extract 15 core semantic keywords such as \<Problem\> and \<Solution\>. We designed scaffolds, semantic patterns based on frequently co-occurring keyword pairs that represent the conceptual flow of questions (e.g., \<problem\>, \<solution\>, \<effect\>)

— which served as guiding structures for prompt formulation. Using these scaffolds, we generated an additional 1,000 multi-hop QA pairs via ChatGPT with prompts tailored to elicit multi-step reasoning. Based on the evaluation of semantic validity and frequency of these scaffolds, we selected 11 representative question concepts, which serve as a controlled yet expressive abstraction for automated large-scale QA generation. A detailed summary of the question concepts, associated keyword pairs, and examples are in Table 2, and the full generation process is in Section 3.3.

## 3.2 Document Selection

We selected document pairs that can support meaningful multi-hop QA. To ensure relevance and reasoning depth, we focused on pairs that exhibit either semantic similarity or structural connections, such as overlapping keywords, comparable abstracts, or citation-based relationships.

**Paragraph-Oriented Type** To identify semantically related document pairs, we first extracted keywords from each article by removing stopwords and applying YAKE (Campos et al. 2020). These keywords, along with titles and abstracts, were then used to compute pairwise document similarity. For similarity computation, we used a TF-IDF representation of titles and abstracts, chosen over dense embeddings (e.g., BERT) for efficiency, scalability, and robustness
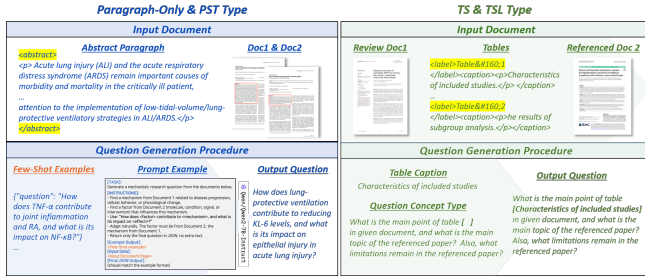
Figure 1: Question Generation Procedure. The process begins with **(1) Input documents** proceed via **(2) Question generation procedures** by type: Paragraph-Oriented types and PST type use LLMs; TS/TSL types follow a deterministic caption-insertion routine.
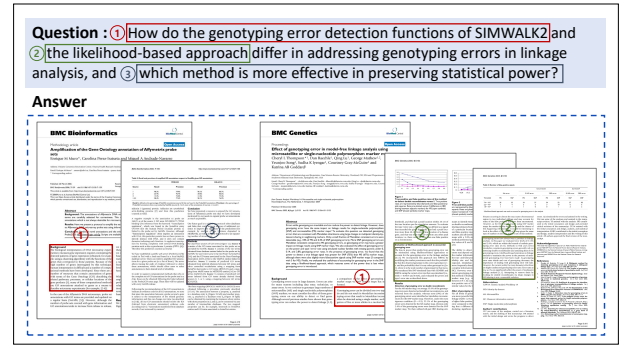


Figure 2: **Paragraph-Oriented QA example with fan-hop reasoning.** Answer snippets are independently retrieved across documents and aggregated without explicit linkage.
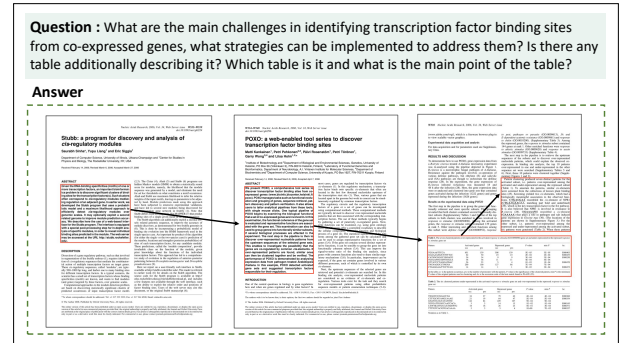


Figure 3: **Table-Oriented QA example with chain-hop reasoning.** Answer snippets are retrieved step-by-step, linking descriptions with table evidence across documents.

to domain-specific terminology (Nuri and Senyurek 2025). The detailed process is provided in Appendix A.1.

**Table-Oriented Type** For Table-Oriented types, we introduced additional filtering steps. For the PST type, we observed that tables containing solution-related evidence—such as intervention effects or outcome summaries—are most frequently located in the results sections of scientific papers. Therefore, we applied a filtering step that retains only those candidate pairs where at least one table is embedded within a document's results section. Motivated by the observation that scientific literature often includes tables comparing and citing other studies, we designed the TS and TSL question types to capture such reasoning patterns. To construct this subset, we focused on systematic review articles, which frequently contain tables referencing external research.

## 3.3 Question Generation

**Paragraph-Oriented Type** Our question generation pipeline has two phases: (1)fine-tune an instruction-following model and (2)generate questions via prompt injection. The goal is to create logically consistent multi-hop questions that weave document-specific information into question concepts. Training instances are JSON objects with instruction-based question concepts, paired abstract contexts, and gold answer instances. The instruction directs the model to replace the placeholders with information from the supplied contexts. A complete example of this format is provided in Appendix A.2.

We fine-tuned the Qwen2-Instruct[2] model, chosen for its strong contextual reasoning, modest resource footprint, and permissive license, using Unsloth's 4-bit LoRA framework(Hu et al. 2021). To preserve controllability and semantic consistency, we trained separate models for each question category (comparison, non-comparison, and PST). Joint training proved sub-optimal, frequently omitting contrastive cues in comparison questions and producing fragmented reasoning in non-comparison and PST cases. During inference, we supplied a few-shot prompt that includes the concept, 2

new abstracts, and 2 in-context example questions using a step-by-step reasoning skeleton shown in Figure 1.

**Table-Oriented Type** For the PST type, questions were generated using the same model-based approach as the Paragraph-Oriented type, but differed slightly in the composition of the input to better reflect the table content. For the TS and TSL types, we did not use an LLM for question generation. Instead, we designed a deterministic procedure that constructs questions by directly inserting table captions into predefined placeholders. As shown in Figure 1, when a valid table was identified in the source document, its caption replaced the corresponding placeholder in the question concept. This design was chosen because the reasoning hops for these types are structured and explicitly defined. As a result, it allows greater controllability in question generation.

**Post-Processing** We applied the following filtering criteria to post-process generated questions. 1) Question generation failures: cases where the model produced no output and records with malformed output formats were excluded. 2) Question concept mismatch: Questions that did not include the keywords defined in the corresponding question concept were filtered out. The final dataset includes only those question–context pairs that passed all filtering steps.

---

[2]https://huggingface.co/unsloth/Qwen2-7B-Instruct-bnb-4bit

## 3.4 Answer Matching

In DocHop-QA, answer retrieval is guided by two distinct reasoning patterns[3]: *fan-hop* for Paragraph-Oriented types, where semantically related content is aggregated from documents (Zhu et al. 2024), and *chain-hop* for Table-Oriented types, which require explicit links between structured and unstructured content (Xu et al. 2021; Chen, Lin, and Durrett 2019). These paradigms shape our matching pipeline, with unified retrieval for Paragraph-Oriented types and tailored strategies for Table-Oriented ones.

**Paragraph-Oriented Type** For Paragraph-Oriented types, we extracted the top 5 semantically relevant snippets from documents using a weighted combination of TF-IDF and BERTScore. Then, we applied a similarity filter to remove low-quality snippets, and retained only those QA pairs whose source content was sufficiently coherent. This improves evidence fidelity by aligning question intent with contextual meaning, without explicit document linking.
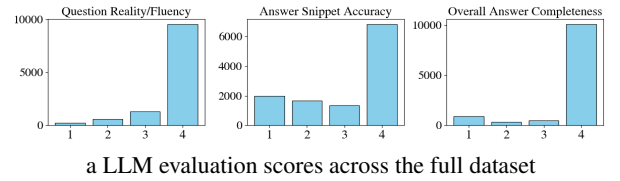
**Table-Oriented Type** For Table-Oriented types, we constructed up to four snippets per instance, connecting problem and solution paragraphs, tables, and referenced content across documents. To enhance retrieval precision, we combined semantic similarity with keyword filtering and applied logical chaining when anchor references were present. This structure reflects real-world scientific reasoning, and low-similarity instances were filtered to preserve dataset integrity. Detailed retrieval logic is provided in Appendix A.3.

## 4 Dataset Quality Assurance

To evaluate the quality of the **DocHop-QA** dataset, we define three criteria: ***Reality/Fluency*** ensures that questions are natural, plausible, and semantically meaningful within the context. ***Accuracy*** assesses whether each snippet meaningfully addresses the question and contributes to forming a complete answer. ***Completeness*** evaluates whether the retrieved snippets collectively form a sufficient answer.

**LLM-based Dataset Quality Assurance** We automatically evaluated the dataset using the instruction-tuned Qwen-Chat[4] model. To account for structural differences, we applied tailored prompts for each QA type (Paragraph-Oriented vs Table-Oriented). Each question and its associated snippets were rated on a 4-point scale (1: poor, 4: excellent), with the model producing JSON-formatted outputs including brief justifications for reproducibility. As shown in **Figure 4a**, the vast majority of questions and answers received the highest score(4), supporting the overall robustness of DocHop-QA for complex QA tasks. Detailed analyses are provided in Appendix C.1.

**Human-based Dataset Quality Assurance** To complement automatic evaluation, we conducted a human assessment on a stratified random sample of 50 QA instances (40 Paragraph-Oriented, 10 Table-Oriented). Using *Google Forms*, 15 annotators[5] independently evaluated each QA in-

---

[3]Further details are provided in Appendix A.3

[4]https://huggingface.co/Qwen/Qwen-7B-Chat

[5]Detailed annotator demographics and ethics approval are provided in Appendix C.



a LLM evaluation scores across the full dataset



b LLM&human evaluation comparison on 50 QA instances

Figure 4: Quality Assurance Results (1: poor, 4: excellent).

stance on a 4-point scale for three criteria, reviewing the questions and highlighted answer snippets in PDF format. Feedback was required for ratings of 2 or below. The results (Figure 4b) show strong performance in both question reality and answer snippets relevance, with most scores clustered at the upper end of the scale. These trends were also consistent with the LLM-based evaluation. We provide detailed analyses in Appendix C.3.

## 5 Dataset Analysis

DocHop-QA contains a total of 11,379 samples, consisting of Paragraph-Oriented types (75.3%) and Table-Oriented types (24.7%). The 11 question concepts are evenly distributed, each representing approximately 10% of the dataset. Most questions are constructed using multiple documents, reflecting consistent concept composition and effective information integration across sources.
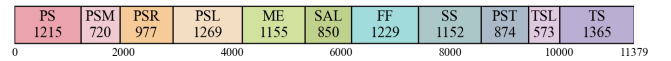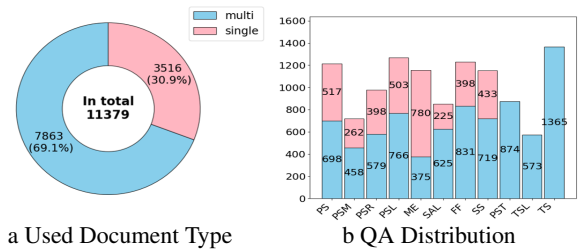


Figure 5: Distribution of 11 Q Types (ordered as in Table 2), ranging from 573 to 1365 samples (Mean: 1034, Std: 241).



a Used Document Type    b QA Distribution

Figure 6: Overview of document usage across question types. Multi-document reasoning dominates overall.

**Used Document Analysis** Our dataset consists of multi-page scientific documents featuring numerous tables, figures, and complex structures with multiple sections and subsections. As summarized in Figure 7, most documents span

10 pages, contain 20 subsections and 60 paragraphs totaling 1,500 to 2,500 tokens closely reflecting the scale and complexity of real-world scientific literature.
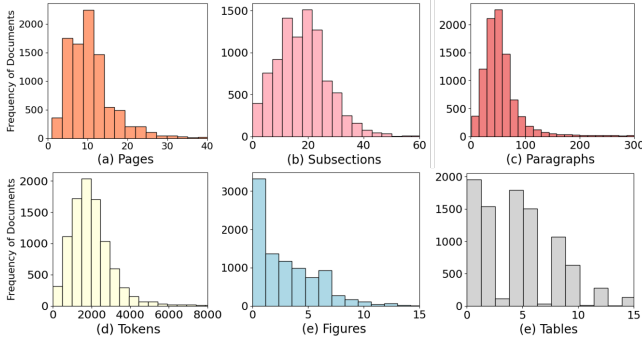


Figure 7: Distribution of document component types. DocHop-QA includes multi-page, multi-section VRDs, supporting complex multi-hop and multimodal QA.

**Word Cloud Analysis**  We extracted nouns, verbs, noun-verb combinations, and overall words from the question texts and generated separate word clouds for each category (Appendix D.1). As expected for PubMed data, medical terms like patient, cancer, and treatment were frequent. Distinct keyword patterns across question types show that the dataset captures diverse linguistic and topical scopes.

**Answer Snippet Distribution Analysis**  We analyzed which sections the answer snippets were mainly extracted from. As shown in Figure 8 and Figure 25, we mapped detailed section titles into 8 Super-Section tags. Paragraph-Oriented types primarily reference the Abstract, while Table-Oriented types mainly used the Results and Discussion sections. This suggests that our dataset reflects section-specific semantics aligned with scientific discourse. Detailed mapping, further analysis and section distribution analysis of table snippets are provided in Appendix D.2.

**Similarity Analysis**  To verify semantic relatedness across both documents and snippets, we conducted two analyses: (1) a document-level similarity check between paired documents for TS/TSL types, and (2) an intra-instance similarity check between the solution paragraph(snippet 2) and the solution-related table (snippet 3) for the PST type. Results show consistently high semantic similarity, indicating that our QA construction is grounded in meaningful cross-document and cross-modal connections. Further details are provided in Appendix D.3

**Section Distance Analysis**  To assess reasoning span, we measured section-level distances between answer snippets. Results in Appendix D.4 show that, unlike prior datasets, many QA pairs require long-range hops across widely separated sections and even across documents, highlighting DocHop-QA's suitability for long-context reasoning.

## 6  Experiments and Results

To evaluate the utility and versatility of DocHop-QA, we conducted experiments across four core QA tasks reflect-
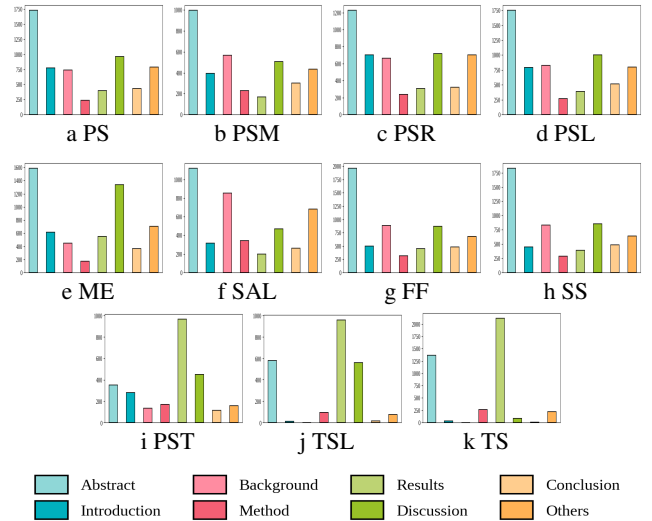


Figure 8: Section-level distribution of answer spans across the 11 question types. Paragraph-centric types predominantly reference Abstract and Methods, while Table-centric types align with Results and Discussion sections.

ing different reasoning paradigms and input modalities. The following tasks span both discriminative and generative approaches, involving structured index prediction, free-form answer generation, and multimodal integration. Each task is designed to assess models' ability to reason over DocHop-QA's complex, document-rich environment. Below, we summarize the task description and key findings from the results. All experimental preprocessing, model architectures, evaluation protocols, full experimental setups, metric definitions, and ablations are provided in Appendix E.

| Model | Sample | | | Accuracy | | |
|---|---|---|---|---|---|---|
| | F1 | Re | Pr | Be | Co | Ov |
| LayoutLMv3 | **12.35** | 34.97 | 8.47 | 0.00 | 0.04 | 95.52 |
| LayoutXLM | 7.13 | 4.49 | 24.06 | 15.64 | 0.00 | 33.96 |

Table 3: Performance on BBox Entity Index Extraction. We report sample-level F1, Recall (Re), and Precision (Pr), along with spatial accuracy metrics: Belong (Be), Contain (Co), and Overlap (Ov), for LayoutLMv3 and LayoutXLM.

**Task 1: BBox Entity Index Extraction.**  Formulated as a multi-label classification task, this evaluates how well models can identify the relevant answer-contained bounding boxes (e.g., paragraphs, figures, tables) in OCR-processed document layouts. We used two layout-aware vision-language models, LayoutLMv3 (Huang et al. 2022) and LayoutXLM (Xu et al. 2020), with classifier heads to predict bounding box indices.

Evaluation includes spatial accuracy (belong/overlap/contain) and sample-level precision, recall, and F1.

As shown in Table 3, LayoutLMv3 outperformed LayoutXLM, achieving the higher sample-level F1 score of 12.35, demonstrating its relatively better capability in

| Setup | Text | Img | Sample | | | Accuracy | | |
|---|---|---|---|---|---|---|---|---|
| | | | F1 | Re | Pr | Be | Co | Ov |
| Concat | BB | - | 23.18 | 18.81 | 30.45 | 4.22 | 0.09 | 62.21 |
| Concat | LF | - | 21.83 | 20.09 | 24.12 | 0.22 | 0.00 | 66.17 |
| PerEnt | BB | - | 21.25 | 17.22 | 27.94 | 2.07 | 0.09 | 58.96 |
| PerEnt | LF | - | 22.99 | 18.71 | 30.08 | 3.25 | 0.00 | 63.80 |
| PerEnt | BB | CL | **23.24** | 21.29 | 25.77 | 0.48 | 0.09 | 64.98 |
| PerEnt | LF | CL | 23.18 | 18.81 | 30.45 | 4.22 | 0.09 | 62.21 |

Table 4: Results for XML Entity Index Extraction. We compare concatenated (Concat) vs. per-entity (PerEnt) setups across text-only and text+image CLIP(CL) inputs using Big-Bird (BB) and Longformer (LF). Metrics include sample-level F1, Recall (Re), Precision (Pr), and spatial accuracy: Belong (Be), Contain (Co), Overlap (Ov).

| Model | S | Img | Sample | | | Accuracy | | |
|---|---|---|---|---|---|---|---|---|
| | | | F1 | Re | Pr | Be | Co | Ov |
| InternVL | Z | - | 7.03 | 9.77 | 7.39 | 0.88 | 0.88 | 28.30 |
| InternVL | Z | ✓ | 4.24 | 6.72 | 4.03 | 0.35 | 0.75 | 17.44 |
| InternVL | O | - | 4.66 | 5.76 | 4.80 | 0.26 | 0.48 | 19.68 |
| InternVL | O | ✓ | 4.34 | 5.77 | 4.68 | 0.44 | 0.62 | 17.31 |
| Qwen | Z | - | **14.06** | 51.62 | 13.69 | 4.04 | 34.05 | 75.48 |
| Qwen | Z | ✓ | 7.92 | 58.53 | 7.06 | 1.54 | 50.04 | 69.73 |
| Qwen | O | - | 11.05 | 41.72 | 11.25 | 5.32 | 25.4 | 65.55 |
| Qwen | O | ✓ | 8.12 | 18.55 | 7.41 | 0.04 | 9.89 | 41.08 |

Table 5: Structured Generative Answering results. S: Setup (Z: Zero-shot, O: One-shot), Img: Image input. Metrics include F1, Recall (Re), Precision (Pr), and spatial accuracy: Belong (Be), Contain (Co), Overlap (Ov).

grounding questions to spatially localized bounding box entities. LayoutLMv3 achieved near-perfect overlap accuracy (95.52%), indicating that the model can still retrieve spatially proximate regions relevant to the answer, even if full answer containment is not achieved (0.04% in the "Contain" metric). These findings show that DocHop-QA is a valuable benchmark for pushing the limits of document layout understanding in QA settings. It reveals critical challenges in handling long, multi-document inputs, while offering fine-grained bounding box supervision to support model development. Users can use DocHop-QA to design new training strategies (e.g., multi-image stitching, hierarchical encoding) or pretraining tasks that improve entity-level reasoning under spatial and modality constraints.

**Task 2: XML Entity Index Extraction.** This task predicts relevant XML entity indices directly, assessing how well models can reason over text-structured representations of documents. Evaluation metrics mirror Task 1.

We evaluated multi-label XML entity prediction using BigBird (Zaheer et al. 2020) and Longformer (Beltagy, Peters, and Cohan 2020) across concatenated and per-entity setups, with and without CLIP-derived visual features (Radford et al. 2021). The best F1 score (23.24) was achieved using BigBird with both text and image inputs in the per-entity setup (Table 4), highlighting the value of layout-grounded visual features. Model comparisons reveal complementary strengths: BigBird excels in capturing long-range context in the concatenated setup, while Longformer performs better on segmented inputs. Visual inputs further boost per-entity performance, showing the benefit of aligning text with document layout. DocHop-QA exposes label bias and modality interaction challenges, offering a valuable benchmark for developing models that integrate textual, visual, and positional cues in entity-level reasoning.

**Task 3: Structured Generative Answering.** For this task, models generate lists of answer entity indices in response to questions, using labeled XML content and optionally associated page images. We compared Qwen2.5-VL (Bai et al. 2025) and InternVL2 (Chen et al. 2024) in both zero-shot and one-shot settings. Structured generation performance is measured by multi-label accuracy and sample-level F1. This task assesses whether generative models can produce struc-

tured answer indices from full-document inputs. Qwen2.5-VL achieves the highest F1 score (14.06) in a zero-shot, text-only setting (Table 5), showing strong performance without fine-tuning when guided by well-crafted prompts. InternVL2 performs poorly across sets and struggles with consistent output formatting, underscoring the need for task-specific pre-training. Adding document images generally hurts performance due to alignment challenges between text and visual layouts. Overall, the structured generative task highlights both the promise and limitations of using instruction-following LLMs for entity-index generation. DocHop-QA offers a challenging testbed for structured reasoning in generative settings, supporting deeper studies into prompt design, modality alignment, and decoding strategies.

**Task 4: Generative Text Extraction.** This task tests the ability of models to generate natural-language answers using different combinations of text and image input. We assessed three prompting styles (question-only, zero-shot, and one-shot) and three input modalities (text-only, image-only, text+image), using models such as

InternVL2, QWEN2.5. Gemini-2.5[6]. Evaluation is based on BLEU (Papineni et al. 2002) and ROUGE (Lin 2004) scores against gold reference answers. The generative text extraction task assesses the ability of large vision-language models to produce free-form answers from document-level context, without relying on predefined answer indices. Detailed results of the experiment are provided in the Appendix E.

## 7 Conclusion

We introduced DocHop-QA, a novel benchmark to address the limitations of existing multi-hop QA datasets by supporting reasoning over multimodal, multi-document scientific corpora. Unlike prior datasets, DocHop-QA captures realistic information, seeking scenarios by integrating unstructured text, structured tables, and visual layout features, without relying on gold chains or hyperlinks. Using an LLM-driven pipeline guided by 11 scientific reasoning concepts, we generated diverse QA pairs and evaluated performance on four representative tasks. Experimental results

---

[6]https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-flash

highlight critical insights into model behavior, modality integration challenges, and prompt design effects, underscoring the need for further research in multi-hop, multimodal reasoning. We hope DocHop-QA serves as a foundation for developing QA systems capable of tackling the intricacies of real-world information-seeking tasks.

# References

Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; Zhong, H.; Zhu, Y.; Yang, M.; Li, Z.; Wan, J.; Wang, P.; Ding, W.; Fu, Z.; Xu, Y.; Ye, J.; Zhang, X.; Xie, T.; Cheng, Z.; Zhang, H.; Yang, Z.; Xu, H.; and Lin, J. 2025. Qwen2.5-VL Technical Report. *arXiv preprint arXiv:2502.13923*.

Beltagy, I.; Peters, M. E.; and Cohan, A. 2020. Longformer: The Long-Document Transformer. *arXiv:2004.05150*.

Campos, R.; Mangaravite, V.; Pasquali, A.; Jorge, A.; Nunes, C.; and Jatowt, A. 2020. YAKE! Keyword extraction from single documents using multiple local features. *Information Sciences*, 509: 257–289.

Chang, Y.; Narang, M.; Suzuki, H.; Cao, G.; Gao, J.; and Bisk, Y. 2022. Webqa: Multihop and multimodal qa. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16495–16504.

Chen, J.; Lin, S.-t.; and Durrett, G. 2019. Multi-hop question answering via reasoning chains. *arXiv preprint arXiv:1910.02610*.

Chen, W.; Zha, H.; Chen, Z.; Xiong, W.; Wang, H.; and Wang, W. 2020. Hybridqa: A dataset of multi-hop question answering over tabular and textual data. *arXiv preprint arXiv:2004.07347*.

Chen, Z.; Wu, J.; Wang, W.; Su, W.; Chen, G.; Xing, S.; Zhong, M.; Zhang, Q.; Zhu, X.; Lu, L.; et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 24185–24198.

Ding, Y.; Ren, K.; Huang, J.; Luo, S.; and Han, S. C. 2024. MVQA: A Dataset for Multimodal Information Retrieval in PDF-based Visual Question Answering. *arXiv preprint arXiv:2404.12720*.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. LoRA: Low-Rank Adaptation of Large Language Models. arXiv:2106.09685.

Hu, W.; Xu, Y.; Li, Y.; Li, W.; Chen, Z.; and Tu, Z. 2024. Bliva: A simple multimodal llm for better handling of text-rich visual questions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 2256–2264.

Huang, Y.; Lv, T.; Cui, L.; Lu, Y.; and Wei, F. 2022. LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking. In *Proceedings of the 30th ACM International Conference on Multimedia*, MM '22, 4083–4091. New York, NY, USA: Association for Computing Machinery. ISBN 9781450392037.

Lee, S.; Kim, H.; and Kang, J. 2023. LIQUID: A framework for list question answering dataset generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 13014–13024.

Li, R.; Wang, Z.; Tran, S.; Xia, L.; and Du, X. 2024. MEQA: A Benchmark for Multi-hop Event-centric Question Answering with Explanations. *Advances in Neural Information Processing Systems*, 37: 126835–126862.

Lin, C.-Y. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.

Louis, A.; van Dijck, G.; and Spanakis, G. 2024. Interpretable long-form legal question answering with retrieval-augmented large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 22266–22275.

Mathew, M.; Bagal, V.; Tito, R.; Karatzas, D.; Valveny, E.; and Jawahar, C. 2022. Infographicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1697–1706.

Mathew, M.; Karatzas, D.; and Jawahar, C. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2200–2209.

Nuri, Y.; and Senyurek, E. 2025. Research Abstracts Similarity Implementation By Using TF-IDF Algorithm. *IOSR Journal of Computer Engineering (IOSR-JCE)*, 27(1, Ser. 4): 04–10.

Oh, J.; Lee, G.; Bae, S.; Kwon, J.-m.; and Choi, E. 2023. Ecg-qa: A comprehensive question answering dataset combined with electrocardiogram. *Advances in Neural Information Processing Systems*, 36: 66277–66288.

Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, 311–318. USA: Association for Computational Linguistics.

Pramanick, S.; Chellappa, R.; and Venugopalan, S. 2024. Spiqa: A dataset for multimodal question answering on scientific papers. *arXiv preprint arXiv:2407.09413*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; Krueger, G.; and Sutskever, I. 2021. Learning Transferable Visual Models From Natural Language Supervision. arXiv:2103.00020.

Tanaka, R.; Nishida, K.; Nishida, K.; Hasegawa, T.; Saito, I.; and Saito, K. 2023. Slidevqa: A dataset for document visual question answering on multiple images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 13636–13645.

Trivedi, H.; Balasubramanian, N.; Khot, T.; and Sabharwal, A. 2022. MuSiQue: Multihop Questions via Single-hop Question Composition. *Transactions of the Association for Computational Linguistics*, 10: 539–554.

Wang, M.; Chen, L.; Fu, C.; Liao, S.; Zhang, X.; Wu, B.; Yu, H.; Xu, N.; Zhang, L.; Luo, R.; et al. 2024. Leave no document behind: Benchmarking long-context llms with extended multi-doc qa. *arXiv preprint arXiv:2406.17419*.

Welbl, J.; Stenetorp, P.; and Riedel, S. 2018. Constructing datasets for multi-hop reading comprehension across documents. *Transactions of the Association for Computational Linguistics*, 6: 287–302.

Xu, W.; Deng, Y.; Zhang, H.; Cai, D.; and Lam, W. 2021. Exploiting reasoning chains for multi-hop science question answering. *arXiv preprint arXiv:2109.02905*.

Xu, Y.; Li, M.; Cui, L.; Huang, S.; Wei, F.; and Zhou, M. 2020. LayoutLM: Pre-training of Text and Layout for Document Image Understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &amp; Data Mining*, KDD '20, 1192–1200. ACM.

Yang, Z.; Qi, P.; Zhang, S.; Bengio, Y.; Cohen, W. W.; Salakhutdinov, R.; and Manning, C. D. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

Yoon, W.; Jackson, R.; Lagerberg, A.; and Kang, J. 2022. Sequence tagging for biomedical extractive question answering. *Bioinformatics*, 38(15): 3794–3801.

Zaheer, M.; Guruganesh, G.; Dubey, K. A.; Ainslie, J.; Alberti, C.; Ontanon, S.; Pham, P.; Ravula, A.; Wang, Q.; Yang, L.; et al. 2020. Big bird: Transformers for longer sequences. *Advances in Neural Information Processing Systems*, 33.

Zhang, T.; Kishore, V.; Wu, F.; Weinberger, K. Q.; and Artzi, Y. 2019. BERTScore: Evaluating Text Generation with BERT. *CoRR*, abs/1904.09675.

Zhu, A.; Hwang, A.; Dugan, L.; and Callison-Burch, C. 2024. FanOutQA: A multi-hop, multi-document question answering benchmark for large language models. *arXiv preprint arXiv:2402.14116*.

## A  Dataset and Code Implementation Details

**Code & Data**  We released the full implementation and reproduction scripts for Dataset Construction (Section 3), Dataset Quality Assurance (Section 4), and Experiments and Results (Section 6) at the following link.
Access it here: https://shorturl.at/rqfcM

### A.1  Document Selection Implementation Details

As mentioned in Section 3.2, we formalized the semantic similarity-based document pairing strategy into a concrete algorithmic procedure, as illustrated in **Algorithm 1**. Candidate pairs are first identified using the following two filters: (1) at least one overlapping keyword, and (2) a TF-IDF cosine similarity score of at least $0.3^7$. After this first-stage filtering, we compute a total matching score for each remaining pair as a weighted combination of normalized keyword overlap and TF-IDF similarity. For each source document, we then kept only the top three scoring partners for downstream question generation, limiting repeated reuse of the same documents and promoting document diversity.

---

**Algorithm 1: Document Selection**

---

**Input**: Source corpus $\mathcal{C} = \{d_1, \ldots, d_N\}$ (PubMed XMLs)
**Parameter**: Keyword extractor (YAKE) $\Phi$ ; TF–IDF vectorizer $\psi$ ; similarity threshold $\tau$ ; weight $\alpha$ ; top-$L$
**Output**: Paired set $\mathcal{P} \subseteq \{(i,j) \mid 1 \le i < j \le N\}$

1: $\mathcal{P} \leftarrow \varnothing$; initialize $S_{ij} \leftarrow -\infty$ for all $i < j$
2: **for** $i \leftarrow 1$ **to** $N$ **do**
3:     extract *pmc id$_i$*, *title$_i$*, *abstract$_i$* from $d_i$
4:     *fulltext$_i$* $\leftarrow$ *title$_i$* $\|$ *abstract$_i$*
5:     $\mathcal{K}_i \leftarrow \Phi(\textit{fulltext}_i)$; remove stopwords
6: **end for**
7: $\mathbf{X} \leftarrow \psi(\textit{fulltext}_1, \ldots, \textit{fulltext}_N)$
8: $\mathbf{M} \leftarrow \cos(\mathbf{X}\mathbf{X}^{\top})$
9: **for** $i \leftarrow 1$ **to** $N$ **do**
10:     **for** $j \leftarrow i + 1$ **to** $N$ **do**
11:         $C_{ij} \leftarrow \mathcal{K}_i \cap \mathcal{K}_j$
12:         **if** $|C_{ij}| > 0$ **and** $\mathbf{M}_{ij} \ge \tau$ **then**
13:             $s_{\text{kw}} \leftarrow \dfrac{\log(1 + |C_{ij}|)}{\log(1 + \max_k |\mathcal{K}_k|)}$
14:             $s_{\text{sim}} \leftarrow \mathbf{M}_{ij}$
15:             $S_{ij} \leftarrow \alpha\, s_{\text{kw}} + (1 - \alpha)\, s_{\text{sim}}$
16:         **end if**
17:     **end for**
18:     $\mathcal{N}_i \leftarrow$ Top-$L$ indices $j$ by $S_{ij}$
19:     **for all** $j \in \mathcal{N}_i$ **do**
20:         $\mathcal{P} \leftarrow \mathcal{P} \cup \{(i,j)\}$
21:     **end for**
22: **end for**
23: **return** $\mathcal{P}$

---

---

[7]After testing multiple thresholds (0.1–0.9) and manually inspecting retrieved documents, 0.3 was selected for consistently yielding the most relevant results.

### A.2  Question Generation Training Data

We presented training instances for three types of question generation: Non-Comparison, Comparison, and Non-Refer-Table (PST). Each instance follows a unified schema comprising a prompt, an instruction, two abstracts, and the resulting question. The prompt encapsulates the underlying question concept, the instruction clarifies how the concept should be instantiated, and Contexts 1 and 2 provide the input abstracts. Although the schema is identical across types, we train a distinct model for each category because our preliminary experiments revealed that type-specific models generate a higher proportion of valid questions than a single unified model.

---

*Prompt*
Instruction: How does [ ] compare to [ ], and what are the key similarities and differences?
Chain-of-Thought (internal use only): Generate a question by replacing only the placeholder [ ] in the question concept below using document-specific information derived from the provided contexts. You must use two documents.

*Instruction*
How does [ ] compare to [ ], and what are the key similarities and differences?

*Input* : {
  *Context 1*
```
 Regionalization of pediatric emergency
 care in the US. The 2006 IOM Report
 highlighted uneven care. Pediatric
 emergency services are concentrated
 in a limited number of children's
 hospitals and trauma centers. EMS-based
 categorization models like EDAP exist, but
 family preference for nearby EDs limits
 effectiveness.
```

  *Context 2*
```
 Current state of pediatric emergency care
 in Korea. A 2010 survey showed widespread
 unpreparedness: many EDs lacked essential
 pediatric equipment, and consultations
 were often handled by inexperienced
 residents. Monitoring for sedated children
 was frequently unavailable. }
```

*Output*
How do the characteristics of the healthcare systems in the United States compare to Korea, and what are the key similarities and differences?

---

Table 6: An example of our dataset instance based on comparison-type multi-hop question generation. Fields such as prompt, instruction, and context are shown in full.

Table 7: An example of our dataset instance based on Non-comparison type multi-hop question generation. Fields such as prompt, instruction, and context are shown in full.

Table 8: An example of our dataset instance based on PST type multi-hop question generation. Fields such as prompt, instruction, and context are shown in full.

## A.3 Detailed Answer Matching

**Reasoning Patterns** In DocHop-QA, multi-hop reasoning follows two paradigms. Fan-hop aggregates semantically complementary information from multiple, independently retrieved documents. Each snippet contributes to the final answer without explicit cross-referencing. Chain-hop, on the other hand, performs step-by-step inference where information in one document directly leads to or supports content in another, often forming a logical link-such as when a table summarizes results from an external study.

**Paragraph-Oriented Type** To ground each QA instance with relevant context, we parsed the full-text files of both documents. The subsection-level content was extracted, and segments with fewer than two sentences were discarded to avoid sparse inputs. For each segment, we computed semantic similarity to the question using a hybrid scoring method: a weighted combination of TF-IDF (weight 0.2) and BERT-based(Zhang et al. 2019) sentence embeddings (weight 0.8). This configuration prioritized deeper semantic understanding while preserving term-level sensitivity. We aggregated similarity scores across both documents and retrieved the top-5 segments with the highest scores. To ensure quality, only answer instances with all retrieved snippets above a similarity score of 0.4 were retained. This threshold was empirically selected based on preliminary validation, which showed that higher similarity was strongly associated with answer plausibility. All filtered QA-context pairs were saved in structured JSON format(Appendix B).

**Table-Oriented Type** We employed two strategies depending on the snippet's role: some snippets were retrieved via semantic-similarity scoring augmented with keyword-based filtering, whereas others were extracted using predefined logic. To support the semantic-similarity route, we curated three keyword lists targeting problem-, solution-, and limitation-related paragraphs; these lists served as an initial paragraph filter in each context and were then combined with semantic-similarity scoring to select answer snippets. If no paragraphs are filtered we computed similarity scores between questions and all paragraphs.

### Problem-related Keywords

burden, challenge, gap, unmet need, lack of, insufficient, controversy, uncertainty, poor outcomes, high prevalence, incidence, rising rates, increasing trend, public health concern, problem, barrier, limitation, disparity, inequity, inaccessibility, misdiagnosis, delay, complication, risk factor, vulnerability, deficiency, inconsistency, overuse, underuse, shortage, low adherence, noncompliance, failure, inadequacy, ineffectiveness, complexity, fragmentation, obstacle, hazard, disease burden.

### Solution-related Keywords

approach, strategy, intervention, treatment, therapy, model, framework, tool, instrument, protocol, program, solution, innovation, novel, effective, efficacy, improvement, enhancement, implementation, guideline, practice, plan, policy, recommendation, management, design, structure, evaluation, algorithm, process, assessment, mechanism, procedure, initiative, rollout, deployment, trial, interdisciplinary, collaborative, multimodal, integrated, systematic, scalable, feasible, evidence-based, optimization.

### Limitation-related Keywords

limitation, limit, lack, weakness, shortcoming, bias.

**Logical Answer Chain** We designed logical answer chains tailored to each QA type, which are explained below. Green-colored snippets used the semantic similarity strategy while others were retrieved via predefined logic. (See Figure 9.)

### PST Type

- **Snippets 1 & 2**: Problem and Solution paragraphs were keyword filtered. The table is assumed to be in the solution document; the problem document provides motivation. We then selected the most relevant paragraph per document using 0.2 * TF-IDF + 0.8 * BERTScore.

- **Snippet 3**: The most relevant table content, selected based on similarity to Snippet 2.

- **Snippet 4**: Paragraph in the document referring the table.

### TS & TSL Type

- **Snippet 1**: Table metadata from the source document.

- **Snippet 2**: Paragraph referring to the table.

- **Snippet 3**: Abstract of the referenced document, linked via citation markers in the table.

- **Snippet 4**: Paragraph describing the cited document's limitations. Snippet 4 was retrieved via a query-driven similarity search using the prompt "What is the limitation of this paper?" and subsequently keyword-filtered. If the similarity score for snippet 4 fell below 0.4, the snippet was discarded, and the QA instance was reclassified as TS type, containing only the first three snippets.



Figure 9: Answer matching logic for TS & TSL types: Snippets 1–3 follow a fixed logic sequence, while Snippet 4 is selected using the same method as in other Q types.

a Paragraph-Oriented



b Table-Oriented

Figure 10: LLM based Quality Assurance result, (1) Question Reality/Fluency (2) Answer Snippet Accuracy (3) Overall Answer Completeness. Scores are based on a 4-point scale (1:poor, 4: excellent).

## B Dataset Description

DocHop-QA is a multi-hop QA dataset built from PubMed articles. Each record contains a single QA pair, along with supporting paragraphs drawn from one or more scientific papers. In total, the dataset comprises 11,359 instances : Paragraph-Oriented (8,547), Table-Oriented (2,812)

**Data Structure** Each instance is represented as a JSON object with the following fields:

- id: Unique identifier for each QA instance
- task_type: Type
- hop_type: Hop type ("fan" or "chain")
- question_concept_type: 11 Question concept types
- question: The natural language multi-hop question
- used_doc: Indicates if snippets are from a single or multiple documents ("single" or "multi")
- context_list: List of answer snippets

**Context List Structure** Each element in context_list is a paragraph or table object with the following attributes:

- pmc_id: PubMed Central article ID
- section: Top-level section
- subsection: Subsection title
- type: Type of content ("text" or "table")
- Raw content: Original XML tagged paragraph content
- content: Cleaned plain text content

## C Dataset Quality Assurance

### C.1 LLM-based Dataset Quality Assurance

As shown in **Figure 10**, Paragraph-Oriented QA performed well, with most scores rated 3 or 4. Some instances had lower accuracy, likely due to snippets covering only parts of the question, reducing alignment with its overall intent. Table-Oriented QA showed slightly lower Reality/Fluency

due to its rigid format, but still maintained high overall quality. These results confirm that DocHop-QA exhibits robust quality and is suitable for complex QA tasks.



a Paragraph-Oriented



b Table-Oriented

Figure 11: Human and Qwen evaluation score distributions across three criteria. Scores are based on a 4-point scale (1:poor, 4: excellent). Each bar pair shows scores from human annotators and the Qwen model, highlighting areas of agreement (e.g., completeness) and disagreement (e.g., accuracy in Table-Oriented type).

### C.2 Human-based Dataset Quality Assurance

**Annotator Demographics** Fifteen annotators participated in the human evaluation. The gender distribution was approximately balanced (7 female, 8 male). In terms of education level, 6 were undergraduate students, 6 were master's students, and 3 were doctoral students. The participants had academic backgrounds in business (n=3), industrial and management engineering (n=4), computer engineering unrelated to AI (n=2), and AI-related computer science (n=6). The age distribution was as follows: 20 years (n = 8), 30 years (n = 5) and older than 50 years (n = 2).

**Institutional Review Board Statement** The study was approved by the Institutional Review Board (Exemption Number: PIRB-2025-E020, Date: June 9, 2025).

**Evaluation Interface** To ensure consistency and transparency in the human evaluation process, annotators were asked to rate each question-answer pair using a structured Google Form interface. Figure 28 shows a screenshot of the actual form used in the study, which included Likert-scale items for three evaluation criteria and mandatory comment fields for low ratings. This setup helped guide annotators through a standardized evaluation flow and facilitated the collection of both quantitative and qualitative feedback.

### C.3 Detailed Human Evaluation Results

We conducted a separate analysis of Paragraph-Oriented and Table-Oriented QA instances by comparing human and LLM-based evaluation scores. As shown in Figure 11, the results show that human and model assessments were largely

consistent across both types, with particularly strong alignment in paragraph-oriented cases. While answer completeness received slightly lower scores in both types, follow-up feedback from annotators suggested two main factors: lack of familiarity with domain-specific terminology, and the inherent complexity of multi-step reasoning in some questions, which increased perceived difficulty.

## D Additional Dataset Analysis

This section presents additional analyses to support the robustness of our dataset, including keyword frequency patterns, answer snippet distributions, document and table similarity, and cross-section reasoning spans.

### D.1 Common Keywords Analysis

As mentioned in Section 5, we conducted word cloud analysis to examine the linguistic patterns of the question texts. After removing stopwords and question concept-specific terms, we extracted nouns, verbs, noun-verb combinations, and all remaining words from each question. Word clouds were then generated for each question concept to visualize the dominant lexical patterns (See Figure 12 to Figure 22). As expected from PubMed-based data, medical terms such as *patient*, *gene*, *cell*, and *disease* appeared frequently, confirming strong biomedical relevance. Each question type exhibited distinct keyword distributions, reflecting diverse reasoning styles. For instance, non-comparison questions focused on specific mechanisms, while comparison questions (Figure 18, 19) highlighted contrasting entities or outcomes. Frequent use of terms like *compare*, *associated*, and *linkage* further illustrates the comparative emphasis. These findings suggest that our dataset captures both conceptual diversity and domain specificity, providing a robust foundation for training QA models capable of handling a wide range of biomedical reasoning tasks.



Figure 13: Frequent words in **PSM-type** questions—such as *analysis*, *affecting*, and *identifying*—highlight analytical reasoning about genetic mechanisms and disease traits.



Figure 14: Frequent words in **PSR-type** questions suggest a focus on result-oriented reasoning about biomedical solutions or interventions.



Figure 12: Wordclouds for **PS-type** questions show frequent terms like *"method," "patient,"* and *"disease,"* indicating a focus on conceptual understanding of biomedical topics.



Figure 15: Frequent words in **PSL-type** questions—such as *conducting*, *addressing*, and *diagnosing*—highlight a focus on evaluating limitations and follow-up aspects of solutions.

a Noun · b Verb · c Noun & Verb · d All Words

Figure 16: Frequent words in **ME-type** questions—such as *increased*, *reducing*, and *signaling*—indicate a focus on causal mechanisms and biological effects of specific factors.



a Noun · b Verb · c Noun & Verb · d All Words

Figure 19: Frequent words in **SS-type** questions—such as *compare*, *method*, and *analysis*—reflect a focus on contrasting scientific approaches to biomedical problems.



a Noun · b Verb · c Noun & Verb · d All Words

Figure 17: Frequent words in **SAL-type** questions—such as *heart* and *blood*—indicate a focus on assessing strengths and limitations of specific biological studies.



a Noun · b Verb · c Noun & Verb · d All Words

Figure 20: Frequent words in **PST-type** questions—such as *predicting*, *patient*, and *infection*—highlight a focus on identifying challenges and strategies based on tabular data.



a Noun · b Verb · c Noun & Verb · d All Words

Figure 18: Frequent words in **FF-type** questions—such as *associated*, *effect*, and *mechanism*—indicate a focus on functional comparisons between two entities.



a Noun · b Verb · c Noun & Verb · d All Words

Figure 21: Frequent words in **TSL-type** questions—such as *given*, *included*, and *characteristics*—reflect their deterministic, caption-based generation with limited lexical variety.

Figure 22: Frequent words in **TS-type** questions—such as *given*, *included*, and *study*—reflect their deterministic, table-driven design with limited lexical diversity.

## D.2 Detailed Answer Snippet Distribution

**Super-Section Mapping Analysis**  To facilitate section-level analysis, we standardized a diverse set of original section titles into 8 unified Super-Section tags. This mapping aimed to reduce variation caused by inconsistent section titles (e.g., "Materials and Methods", "Subjects and Methods", "Methodology") and to improve interpretability of section-based statistics. The results of standardization are described in Table 23 and Figure 24.

**Answer Snippet Distribution Analysis**  Figure 8 presents the distribution of answer snippet across QA types. Paragraph-Oriented types predominantly retrieved answers from the Abstract section, aligning with their fan-hop retrieval strategy. In contrast, Table-Oriented types mostly extracted snippets from the Results and Discussion sections, reflecting the presence of structured data and interpretive content found in those parts. These patterns highlight that each Q type draws from distinct sections of the document, underscoring the need for models to understand section-level semantics when reasoning over scientific papers.

The following figures illustrate the distribution of answer snippets for each type(Paragraph-Oriented and Table-Oriented), with Figure 8 presenting the results separated by each Q type. The Paragraph-Oriented type, which adopts fan-hop based answer retrieval strategy, most frequently extracts from the "Abstract" section. In contrast, the Table-Oriented type most frequently extracts from the "Results" section, due to its document filtering strategy.

**Table Oriented Section Analysis**  We analyzed the section origin of tables (Table 9) and found that most were extracted from the *Results* section. This supports our decision to include only documents with at least one results-section table for PST examples. TSL and TS tables also mainly came from Results, reflecting their common role in summarizing outcomes across studies. These patterns show that the dataset aligns with typical scientific reporting, where key findings are often presented in results tables.

| Section | PST | TS | TSL |
|---|---|---|---|
| Abstract | 5 | 3 | 3 |
| Introduction | 1 | 16 | 5 |
| Background | 1 | 2 | 0 |
| Method | 33 | 132 | 47 |
| Results | **830** | **1057** | **478** |
| Discussion | 1 | 41 | 4 |
| Conclusion | 0 | 2 | 1 |
| Others | 3 | 112 | 35 |

Table 9: Distribution of Table Origins Across Table-Oriented Types (PST, TS, TSL) with Section Mapping

## D.3 Detailed Similarity Analysis

**Document Similarity Analysis**  To assess the semantic similarity between document pairs in the TSL, TS type, we computed the BERT-based similarity shown in Figure 26a we used the title and abstract of both the main and referenced documents. Unlike TF-IDF, we adopted BERT similarity to capture deeper conceptual relationships, since references are typically cited based on semantic rather than lexical overlap. The resulting mean similarity score of 0.69 confirms that the paired documents in the TSL, TS type are highly related, consistent with the pairing strategy used for other QA types. Also, this semantic alignment ensures that multi-hop questions in our dataset reflect realistic cross-document reasoning patterns, enhancing both coherence and training utility.

**PST: solution-table Similarity Analysis**  We analyzed the semantic similarity between the solution snippet and the corresponding table caption in PST type questions shown in Figure 26b. To quantify this, we computed a weighted average of TF-IDF and BERTScore. The results showed a high average similarity of 0.76, indicating strong semantic alignment between the solution content and its associated table. This suggests that the dataset effectively links textual and table snippet, reinforcing its suitability for tasks requiring cross-modal understanding

## D.4 Detailed Section Distance Analysis

We analyzed the section-level distance between answer snippets to assess how broadly they are distributed across the document as shown in Figure 27. We measured the section gap based on their relative positions within the Super-Section. On average, answer snippets spanned 3 to 4 sections apart, indicating that questions require understanding across multiple sections, rather than being confined to just 1-2 pages-highlighting the dataset's support for long-context reasoning. This shows our dataset supports complex reasoning across broad document contexts, making it suitable for evaluating long-context understanding in QA systems.

## E Additional Experiments and Results

**General Preprocessing**  All documents are converted from PDF to XML and image formats. OCR and table detection are applied to extract layout-aware content. We use a 70/10/20 train/dev/test split, ensuring that each document appears in only one split to prevent information leakage.

| Super-Section | Detail Section Title |
|---|---|
| Abstract | **abstract**, etc. |
| Introduction | **introduction**, etc. |
| Background | **background**, **background and recent developments**, background and overview, etc. |
| Method | **method(s)**, material(s) and method(s), method and implementation, etc. |
| Results | **result(s)**, **results and discussion**, implementation and results, result analysis, etc. |
| Discussion | **discussion(s)**, **discussion and conclusions**, limitations, summary and discussion, etc |
| Conclusion | **conclusion(s)**, conclusions and future directions, conclusion and perspectives, etc. |
| Others | case report, **review**, dataset and evaluation strategy, **abbreviations**, etc. |

Figure 23: Detailed results of Super-Section Mapping, with frequent section titles bolded.



Figure 24: Section Mapping Results in Illustrative Chart.



Figure 25: Barplot of answer snippet distribution for each type of question; **Abstract** the most for Pargraph-Oriented & **Result** the most for Table-Oriented).



Figure 26: Similarity analysis results; **(a)** Boxplot of title similarity for TS& TSL type's document pairs **(b)** Boxplot of weighted similarity score(TF-IDF: 0.2, BERTScore : 0.8) between snippet 2 and snippet 3's table caption.

At the document level, we first extracted structured entities, including titles, abstracts, paragraphs, figures, and tables, from the XML version. For visual grounding, we align these XML entities with their corresponding positions in the PDF using the PaddleX OCR toolkit, and save each page of the PDF as an image. This enables us to generate both text-based and image-based inputs across tasks.

At the sample level, we combined all entities from the involved context documents into a unified list, assigning global indices. For instance, if Document 1 has 10 entities and Document 2 has 32, entities are indexed sequentially from [Entity 0] (e.g., <doc1_title>) to [Entity 41] (e.g., <doc2_paragraph20>). Answer labels are represented as index lists pointing to this combined sequence. In tasks involving image-based models, such as LayoutLMv3, all pages from the context documents are merged into a single image and resized to comply with the model's input constraints.

For fine-tuning, all classification-based models are trained using binary cross-entropy loss with the AdamW optimizer for up to 50 epochs. Early stopping is applied with a patience of either 5 or 10 epochs, based on validation performance

measured by sample-level recall. Training is conducted locally using a single NVIDIA A100 GPU, while Gemini-based inference tasks are executed via API calls.

### E.1 Additional Task Description

**Task 1: BBox Entity Index Extraction** We map each XML entity to one or more OCR-detected bounding boxes. These boxes are compiled into a single 1000×1000 image per document set. LayoutLMv3 and LayoutXLM receive bounding box coordinates, associated OCR text, and document images to predict relevant indices. Entities that cannot be matched are marked with index 0. Evaluation includes multi-label accuracy and sample-level F1 based on spatial containment (belong, contain, overlap).

**Task 2: XML Entity Index Extraction** We evaluate two setups: 1) Concatenated input: the question and all XML entities are fed as one sequence into BigBird or Longformer. 2) Per-entity input: each entity is encoded separately with pooled text (BigBird/Longformer) and optionally CLIP image embeddings. Predictions are made via a classifier over these embeddings. Metrics are identical to Task 1.

Figure 27: Boxplot of section-level distance for each number of document(Single, Multi, Grouped) - emphasizing answer instances of dataset covers average **3-4 section-level distance**; supporting **long-context reasoning**.

**Task 3: Structured Generative Answering** This task uses generative models (Qwen2.5-VL-7B-Instruct and InternVL2-4B) to produce lists of entity indices. Inputs include labeled XML entities and document images. Prompts follow zero-shot or one-shot formats. To avoid memory issues, we cap the number of entities and pages. Evaluation uses the same index-based metrics as in classification tasks.

**Task 4: Generative Text Extraction** We test generative LLMs under three prompting strategies (question-only, zero-shot, one-shot) and three input types (text-only, image-only, text+image). A separate run using Gemini-2.5-Flash handles XML and PDF ingestion. Generated answers are evaluated with BLEU and ROUGE (1, 2, L, Lsum) against reference answers extracted during dataset construction.

## E.2 Additional Result Analysis

**Results of Task1: BBox Entity Index Extraction** We treat this task as multi-label classification over OCR-extracted bounding boxes. Input includes the resized page image (1000×1000), bounding box coordinates, and OCR text. LayoutLMv3 and LayoutXLM are evaluated using spatial metrics—Belong, Contain, and Overlap—and sample-level F1, precision, and recall. LayoutLMv3 significantly outperforms LayoutXLM in overlap-based accuracy (95.52%) and F1 (12.35), showing its robustness in spatial grounding. However, full containment accuracy is near zero, suggesting difficulties in aligning complex layouts.

**Results of Task2: XML Entity Index Extraction** This task also uses multi-label classification but predicts indices directly from XML-defined entities. We examine: 1) Concatenated input setup (question + all XML entities) and 2) Per-entity setup, where text and CLIP-based image features are pooled for each entity. Models include BigBird and Longformer. Best performance (F1: 23.24) is from Big-Bird (per-entity) with CLIP embeddings, indicating that visual context improves model understanding. Model architecture choice matters: BigBird shows advantages in long-range context, while Longformer performs better on per-entity coherence. The task reveals dataset bias towards frequently occurring entity positions.
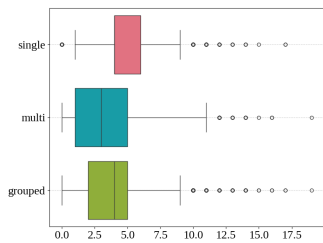
**Results of Task 3: Structured Generative Answering** We prompt generative vision-language models (Qwen2.5-

VL, InternVL2) to output structured entity indices. Inputs include labeled entity texts and optionally, document page images. Prompts follow zero-shot or one-shot formats. Qwen2.5-VL in the zero-shot, text-only configuration yields the highest F1 (14.06). Visual input degrades performance due to layout misalignment. One-shot prompts reduce input space, harming performance on long documents. This task highlights the potential of LLMs for structured reasoning, but also the limits of current multimodal alignment.

**Results of Taks 4: Generative Text Extraction** As shown in Table 13, the best performance is from Qwen2.5-VL in a zero-shot setting using only image modality, which achieves a BLEU score of 21.56 and strong performance across ROUGE-1 (38.45), ROUGE-2 (16.92), and ROUGE-L (23.90). This suggests that modern vision-language models can effectively extract answer passages grounded in complex document images, even without textual context. Comparing architectures, Qwen2.5 consistently outperforms InternVL2 in setups that utilize image input. This advantage likely stems from Qwen's greater input capacity, which allows it to process up to 100 document pages, while InternVL2 is limited to 10 with the simplest configuration. As a result, Qwen2.5 is better suited for large-scale document understanding, particularly when full visual context is required. However, InternVL2 demonstrates better performance when using text-only inputs, achieving a BLEU of 12.65 and ROUGE-2 of 11.11—slightly surpassing Qwen's text-only setup (BLEU: 9.66, ROUGE-2: 10.93). This implies that InternVL2's language modeling capabilities are competitive when visual grounding is not required, and reaffirms the importance of architecture-modality alignment.

Interestingly, Qwen2.5 performs best when image input is used alone, and worst when using only text. This indicates that Qwen is likely optimized for visual reasoning and may struggle to integrate text-only signals in the absence of layout structure. However, adding both image and text did not significantly boost performance—suggesting that naive concatenation of modalities may introduce redundancy or noise, and that more sophisticated cross-modal fusion strategies are needed for optimal performance on DocHop-QA. Prompt design plays a critical role in generative QA performance. For Qwen2.5, zero-shot prompting consistently outperforms one-shot setups, likely due to the substantial length of context documents. Including a one-shot example in the prompt consumes valuable input space, which forces truncation of relevant information, especially under tight token limits. This finding highlights an important trade-off in prompting long-context models: example-rich prompts can backfire in high-volume settings like full-document QA.

We also evaluated Gemini 2.5 across XML, PDF, and combined inputs. Unlike Qwen and InternVL2, Gemini was accessed via API and thus not bound by local memory constraints, allowing full documents without severe truncation. As a result, one-shot prompting performed best across most configurations, with top BLEU scores of 11.96 using PDF and 11.94 using XML. The performance gap between XML-only, PDF-only, and XML+PDF setups was marginal, suggesting Gemini can effectively integrate multimodal sources

regardless of file format. These results illustrate the potential of cloud-based LLMs for processing full-document QA tasks without aggressive input filtering, further expanding the applicability of DocHop-QA in real-world settings.

In summary, this task underscores how DocHop-QA enables robust benchmarking of free-form QA grounded in multimodal and multi-page documents. Users can evaluate generative models under various prompting and input modality combinations, uncovering limitations in memory, alignment, and format-following behavior. The dataset also reveals that vision-language models optimized for document images (e.g., Qwen2.5) are currently better suited for zero-shot generative answering than generic LLMs trained on short passages or unimodal data.

### E.3 Hyperparameters

We provide the best hyperparameter setups for each model in Tables 10 to 14. For fine-tuned models in Task 1 and Task 2, we explore learning rate values of 2e-01, 2e-02, 2e-03, 2e-04, and 2e-05, early stop patience of 10, 5, or 0 (no early stop), and warm up of 10 or 0 (no warm up). For learning rates, 2e-02 generally produces a better F1 score, however, predicted indices tend to skew to the more frequent entity indices, raising recall but sacrificing precision. Using 2e-05 tend to produce more varied predictions with better precision but recall suffers greatly. For the generative models, we limited the number of input entities and page images. Being a smaller model, InternVL2 only ran successfully with a maximum of 25 entities and 10 pages. QWEN2.5 could handle up to 200 entities and 100 page images but, interestingly, performed best mostly with only 10 pages.

| Model | Learning Rate | Batch Size | Max Len | Early Stop | Warm Up |
|---|---|---|---|---|---|
| LayoutLMv3 | 2e-02 | 32 | 512 | 10 | 10 |
| LayoutXLM | 2e-05 | 64 | 512 | 10 | 10 |

Table 10: BBox Entity Index Extraction best model hyperparameters.

| Setup | Text Emb | Img Emb | Max Len | Learning Rate | Early Stop | Batch Size | Warm Up |
|---|---|---|---|---|---|---|---|
| Concat | BB | - | 2048 | 2e-02 | 5 | 16 | 10 |
| Concat | LF | - | 2048 | 2e-02 | 5 | 8 | 10 |
| PerEnt | BB | - | 4096 | 2e-02 | 0 | 64 | 0 |
| PerEnt | LF | - | 4096 | 2e-02 | 0 | 64 | 0 |
| PerEnt | BB | Clip | 4096 | 2e-02 | 5 | 64 | 10 |
| PerEnt | LF | Clip | 4096 | 2e-02 | 5 | 64 | 10 |

Table 11: XML Entity Index Extraction best model hyperparameters. Concat: Concatenated, PerEnt: Per Entity, BB: BigBird, LF: Longformer

### E.4 Experiment Prompts

Table 15 shows the general prompt structures used for generative models. The question-only prompt consists of the context input and question. The zero-shot prompt adds the in-

| Model | Setup | Text | Img | Max Entities | Max Pages | Max New Tokens |
|---|---|---|---|---|---|---|
| InternVL | Zero-shot | ✓ | - | 25 | - | 1024 |
| InternVL | Zero-shot | ✓ | ✓ | 25 | 5 | 1024 |
| InternVL | One-shot | ✓ | - | 25 | - | 1024 |
| InternVL | One-shot | ✓ | ✓ | 25 | 5 | 1024 |
| Qwen | Zero-shot | ✓ | - | 100 | - | 1024 |
| Qwen | Zero-shot | ✓ | ✓ | 50 | 10 | 1024 |
| Qwen | One-shot | ✓ | - | 50 | - | 1024 |
| Qwen | One-shot | ✓ | ✓ | 50 | 10 | 1024 |

Table 12: Structured generative answering best model hyperparameters.

struction between the context input and the question. The one-shot prompt adds an example to the zero-shot structure. For structured generative answering, the instructions are changed to *"Extract the entity IDs of the paragraphs, figures, and/or tables that answer the question:"*

| Model | Setup | Text | Image | Document | BLEU | Rouge1 | Rouge2 | RougeL | RougeLSum |
|---|---|---|---|---|---|---|---|---|---|
| InternVL | Question | - | ✓ | - | 6.22 | 30.02 | 7.24 | 15.95 | 19.04 |
| InternVL | Zero-shot | - | ✓ | - | 5.54 | 28.67 | 7.03 | 14.62 | 19.46 |
| InternVL | One-shot | - | ✓ | - | 6.39 | 28.74 | 6.99 | 16.12 | 18.50 |
| InternVL | Question | ✓ | - | - | 11.20 | 32.69 | 10.02 | 18.49 | 20.32 |
| InternVL | Zero-shot | ✓ | - | - | 12.65 | 33.36 | 11.11 | 18.93 | 21.40 |
| InternVL | One-shot | ✓ | - | - | 14.53 | 32.11 | 11.99 | 19.56 | 20.91 |
| InternVL | Question | ✓ | ✓ | - | 6.91 | 30.09 | 8.05 | 15.85 | 19.79 |
| InternVL | Zero-shot | ✓ | ✓ | - | 6.91 | 30.06 | 8.48 | 15.55 | 20.60 |
| InternVL | One-shot | ✓ | ✓ | - | 7.48 | 31.23 | 8.55 | 16.46 | 20.26 |
| QWEN | Question | - | ✓ | - | 15.06 | 36.64 | 13.43 | 20.79 | 23.89 |
| QWEN | Zero-shot | - | ✓ | - | **21.56** | 38.45 | 16.92 | 23.90 | 25.84 |
| QWEN | One-shot | - | ✓ | - | 12.40 | 34.48 | 11.49 | 19.50 | 22.53 |
| QWEN | Question | ✓ | - | - | 7.77 | 28.62 | 9.33 | 15.71 | 20.44 |
| QWEN | Zero-shot | ✓ | - | - | 9.66 | 30.81 | 10.93 | 17.28 | 21.96 |
| QWEN | One-shot | ✓ | - | - | 7.32 | 27.66 | 8.30 | 15.45 | 19.44 |
| QWEN | Question | ✓ | ✓ | - | 9.80 | 31.89 | 10.70 | 17.82 | 22.12 |
| QWEN | Zero-shot | ✓ | ✓ | - | 12.78 | 33.69 | 12.86 | 19.80 | 23.66 |
| QWEN | One-shot | ✓ | ✓ | - | 9.45 | 31.35 | 9.96 | 17.94 | 21.40 |
| Gemini | Question | - | - | XML | 7.34 | 28.26 | 10.91 | 15.28 | 20.44 |
| Gemini | Zero-shot | - | - | XML | 9.22 | 29.24 | 11.14 | 15.89 | 20.93 |
| Gemini | One-shot | - | - | XML | 11.94 | 24.55 | 12.85 | 16.81 | 19.41 |
| Gemini | Question | - | - | PDF | 7.57 | 28.76 | 10.71 | 15.41 | 20.92 |
| Gemini | Zero-shot | - | - | PDF | 9.03 | 29.62 | 12.08 | 16.20 | 21.88 |
| Gemini | One-shot | - | - | PDF | 11.96 | 25.47 | 10.96 | 15.33 | 18.68 |
| Gemini | Question | - | - | XML+PDF | 7.38 | 28.85 | 10.55 | 15.15 | 20.46 |
| Gemini | Zero-shot | - | - | XML+PDF | 8.89 | 30.70 | 12.24 | 17.07 | 22.63 |
| Gemini | One-shot | - | - | XML+PDF | 11.71 | 24.57 | 12.61 | 16.30 | 18.95 |

Table 13: Comprehensive results for generative text extraction task using InternVL2, QWEN 2.5, and Gemini 2.5 across different prompt setups and input configuration. We measure performance through BLEU and ROUGE metrics.

| Model | Setup | Text | Img | Max Entities | Max Pages | Max New Token |
|---|---|---|---|---|---|---|
| InternVL | Question | - | ✓ | - | 10 | 1024 |
| InternVL | Zero-shot | - | ✓ | - | 10 | 1024 |
| InternVL | One-shot | - | ✓ | - | 5 | 1024 |
| InternVL | Question | ✓ | - | 25 | - | 1024 |
| InternVL | Zero-shot | ✓ | - | 25 | - | 1024 |
| InternVL | One-shot | ✓ | - | 25 | - | 1024 |
| InternVL | Question | ✓ | ✓ | 25 | 5 | 1024 |
| InternVL | Zero-shot | ✓ | ✓ | 25 | 5 | 1024 |
| InternVL | One-shot | ✓ | ✓ | 25 | 5 | 1024 |
| Qwen | Question | - | ✓ | - | 10 | 1024 |
| Qwen | Zero-shot | - | ✓ | - | 10 | 1024 |
| Qwen | One-shot | - | ✓ | - | 10 | 1024 |
| Qwen | Question | ✓ | - | 200 | - | 1024 |
| Qwen | Zero-shot | ✓ | - | 200 | - | 1024 |
| Qwen | One-shot | ✓ | - | 50 | - | 1024 |
| Qwen | Question | ✓ | ✓ | 100 | 10 | 1024 |
| Qwen | Zero-shot | ✓ | ✓ | 100 | 10 | 1024 |
| Qwen | One-shot | ✓ | ✓ | 50 | 10 | 1024 |

Table 14: Generative text extraction best model hyperparameters.

| Segment | Prompt |
|---|---|
| Example | The following `<input_desc>` from one or more documents that may be used to answer a question by extracting paragraphs, figures, and/or tables. For example: `[example_context_pages]` `[example_context_entities]` Question: `<example_question>` Answers: `<example_answer>` Consider the following to answer the question at the end. Extract the paragraphs, figures, and/or tables that answer the question. |
| Context Input | `[instance_context_pages]` `[instance_context_entities]` |
| Instruction | The following `<input_desc>` from one or more documents. Extract the paragraphs, figures, and/or tables that answer the question: |
| Question | `<instance_question>` |
| `<input_desc>` | {"images are pages", "texts are xml", "xml files", "pdf files"} |

Table 15: General prompt structure for generative tasks.

## DocHop-QA Dataset Evaluation

📌 About This Study

Hello, and thank you for participating.
This survey is part of an academic research project.
The goal of this study is to make useful multi-hop QA set which reflects real-world.
Your responses will help improve future AI models and enhance the quality of QA datasets.

📌 Participation Guidelines

- The survey consists of **50 questions**, each requiring you to evaluate QA pair in terms of Question Reality, Answer Accuracy and Completeness.
- Each question is expected to take about **5 minutes**, and the entire survey will take approximately **5hours** to complete. Since this QA involves extensive medical terminology, feel free to use ChatGPT or other tools.

- For each evaluation, you are asked to assign one of the following scores:
  - **A (Excellent)**
  - **B (Fair)**
  - **C (Poor)**
  - **D (Very Poor)**
- If you give C or D to some question set, you must write the reason at comment.

- Please fill in the **Comment box** as follows:
  - If you select **A**, you do **not** need to provide a comment.
  - If you select **B, C, or D**, please provide a **brief reason** explaining your score.
    (e.g., "The answer is partially unrelated to the question", "Some key points(you should indicate this in detail) are missing", etc.)

- All responses will remain **anonymous** and will be used solely for research purposes.

Thank you very much for your valuable time and thoughtful responses!

## DocHop-QA Dataset Evaluation

### 📋 Evaluation Criteria (Read Carefully)

◆ 1. Question Evaluation : Reality
  A. The question is realistic, non-trivial (i.e., not a simple factual lookup), logically coherent, and clearly interpretable.
  B. The question is generally realistic, though slightly awkward or vague in intent. It does not contain logical inconsistencies.
  C. The question is somewhat unnatural or hard to interpret, and may contain minor logical issues or ambiguities.( ;excluding difficulties arising solely from medical terminology)

  D. The question is illogical, unrealistic, or extremely awkward, making its intent unclear or invalid.

◆ 2. Answer Snippet Evaluation : Accuracy
  A. All snippets are directly relevant to the question and provide accurate, meaningful information. There is no irrelevant content.
  B. Most snippets are relevant and informative. A few may be slightly off-topic, but overall accuracy is maintained.
  C. Only some snippets are relevant; many provide little or no useful information for answering the question.
  D. Most or all snippets are irrelevant or fail to answer the question accurately.

◆ 3. Answer Evaluation : Completeness
  A. The snippets collectively address all parts of the question. The answer is complete.
  B. Most key aspects of the question are covered, though minor details may be missing. The answer is mostly complete.
  C. Only partial information is provided; important aspects of the question are not addressed.
  D. The answer fails to address the question in any meaningful way.

(a) Overview page providing the study purpose and outlining the evaluation criteria.



## DocHop-QA Dataset Evaluation

### 3 - Question (ff)

How do the genetic influences on smoking behavior, as assessed through linkage analyses in Document 1, compare to the findings from Document 2 regarding the identification of specific genomic loci associated with nicotine dependence, and what are the key similarities and differences in the identified regions?
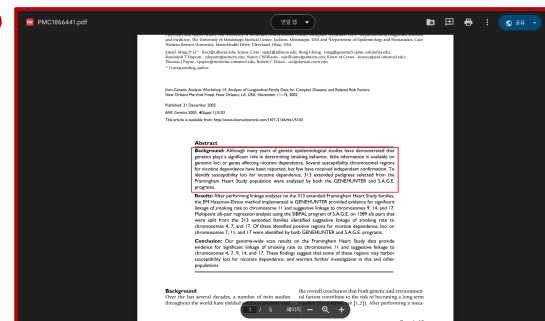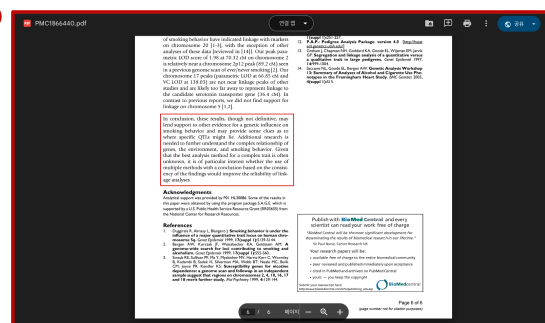
PMC1866440:
① https://drive.google.com/file/d/1fnrFMiPOLSCUr2tDJUQonZAOsRaw4zQB/view?usp=sharing
PMC1866441: ② https://drive.google.com/file/d/1pdZGjCnQGz-6f2qmkChvjlY7abXxdPYD/view?usp=sharing

1. Question Evaluation: Reality / Fluency

○ A
○ B
○ C
○ D

(b) Example QA instance with hyperlinks to PDF documents displaying the answer snippet with a bounding box.

Figure 28: Screenshots of the structured Google Form interface used for human evaluation.