# 📝 NoTeS-Bank: Benchmarking Neural Transcription and Search for Scientific Notes Understanding

Aniket Pal*
CVIT Lab, IIIT Hyderabad
Hyderabad, Telangana, India

Sanket Biswas*
Computer Vision Center, UAB
Barcelona, Spain

Alloy Das*
Habitat Labs, Habitat Lens Pvt. Ltd.
Kolkata, W. B., India

Ayush Lodh†
Habitat Labs, Habitat Lens Pvt. Ltd.
Kolkata, W. B., India

Priyanka Banerjee†
Habitat Labs, Habitat Lens Pvt. Ltd.
Kolkata, W. B., India

Soumitri Chattopadhyay†
UNC Chapel Hill
Chapel Hill, NC, USA

Dimosthenis Karatzas‡
Computer Vision Center, UAB
Barcelona, Spain

Josep Lladós‡
Computer Vision Center, UAB
Barcelona, Spain

C.V. Jawahar‡
CVIT Lab, IIIT Hyderabad
Hyderabad, Telangana, India

## Abstract

Understanding and reasoning over academic handwritten notes remains a challenge in document AI, particularly for mathematical equations, diagrams, and scientific notations. Existing visual question answering (VQA) benchmarks focus on printed or structured handwritten text, limiting generalization to real-world note-taking. To address this, we introduce **NoTeS-Bank**, an *evaluation benchmark* for **Neural Transcription and Search** in note-based question answering. NoTeS-Bank comprises complex notes across multiple domains, requiring models to process unstructured and multimodal content. The benchmark defines two tasks: (1) **Evidence-Based VQA**, where models retrieve localized answers with bounding-box evidence, and (2) **Open-Domain VQA**, where models classify the domain before retrieving relevant documents and answers. Unlike classical Document VQA datasets relying on optical character recognition (OCR) and structured data, NoTeS-BANK demands vision-language fusion, retrieval, and multimodal reasoning. We benchmark state-of-the-art Vision-Language Models (VLMs) and retrieval frameworks, exposing structured transcription and reasoning limitations. NoTeS-Bank provides a rigorous evaluation with *NDCG@5, MRR, Recall@K, IoU, and ANLS*, establishing a new standard for visual document understanding and reasoning.

## Keywords

Multimodal Document Understanding, Multimodal Reasoning, VLMs, Evidence-based Document VQA, Open-Domain QA

## 1 Introduction

*"What we know is a drop, what we do not know is an ocean."* – Isaac Newton. The pursuit of knowledge often begins with handwritten notes — scribbled equations, diagrams, and annotations that serve as the foundation of scientific discovery, engineering breakthroughs, and academic learning. However, despite the fundamental role of handwritten notes, their automated understanding remains a formidable challenge in visual document understanding (VDU).

Prior datasets [37] cover either neatly rendered handwriting (HW-SQuAD) or historical letters (BenthamQA), but none address modern lecture or notebook pages with complex content. Academic notes often mix prose with formulas, contain shorthand or abbreviations, and include sketches, flowcharts, or diagrams drawn by the note-taker. Current state-of-the-art (SOTA) OCR-based VDU models [4, 20, 54] assume a well-defined layout or reading order as in DocVQA [38]. They can falter when content is scattered or non-linear. For instance, models in such datasets [38, 51] performed poorly on questions where the answer required interpreting the document's structure (tables, columns, or aligned layout elements). In handwritten lecture notes, the "layout" might include diagrams next to text or equations below explanatory text, which is not trivial for models to parse, as shown in Figure 1.

Most existing models are not equipped to jointly analyze visual drawings and text. They either ignore non-text elements or treat the problem as pure text extraction or OCR tasks. This is a limitation when notes contain sketches, scientific diagrams, or mathematical equations. For example, a flowchart or a geometric diagram might be essential to answer a question, but generic text-based VQA models will not interpret a drawing of a triangle or a circuit diagram. Similarly, handwritten equations or symbols (e.g. $\Sigma$, integrals, chemical structures) can be misread by OCR or not understood in context (e.g., distinguishing a handwritten "z" from a "2" in an equation as shown in Figure 1). Existing benchmarks like InfographicVQA [36] and MathVista [33] explicitly show that joint reasoning over text and graphics is needed, and current SOTA baselines perform modestly on such tasks. In summary, today's VQA systems tend to be brittle outside the domains they were trained in - they struggle with messy handwriting, unstructured layouts, and mixed modalities prevalent in note-based documents. The information is there—somewhere—but buried beneath uneven ink, equations in the margins, half-drawn diagrams, and scribbled corrections. To answer a question, they don't just recall the fact—they search, they scan, they locate the exact line, the boxed formula, the circled term. That process of not only knowing the answer but being able to localize to it is at the heart of our **evidence-based VQA** task. Evidence-based document VQA challenges models to
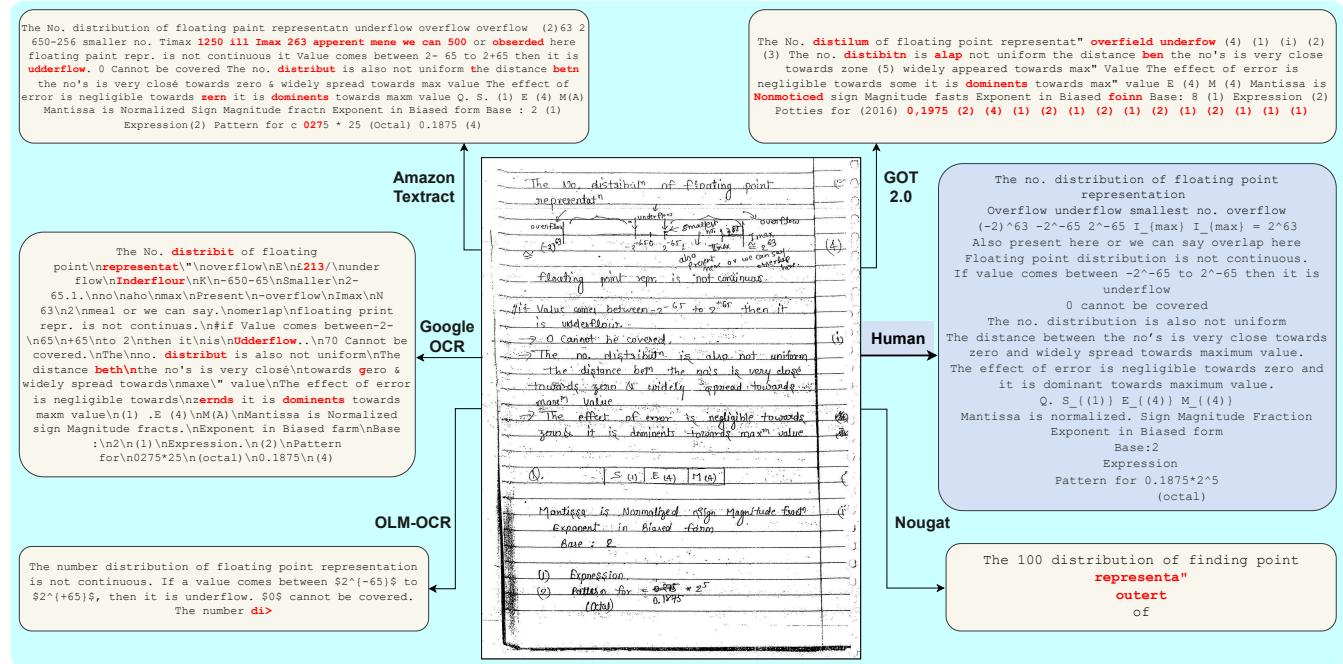
---

**Figure 1:** *Comparison of OCRs on a challenging handwritten scientific note sample.* **Powerful commercial OCR engines (eg. Textract, Google-OCR) fail to accurately transcribe the content, often losing mathematical symbols, structure, and semantic meaning. In contrast, open-source OCRs (e.g. Nougat, Got 2.0, OLM-OCR) struggle to extract even the textual content for the full document. This highlights the** *necessity for multimodal reasoning beyond OCR* **in handwritten document understanding.**

look at the document like a human would, navigating through visual cues, spatial arrangements, and noisy handwriting to retrieve and ground the answer as shown in Figure 2. This isn't just about accuracy—it's about understanding. Grounding the response builds trust, explains reasoning, and proves that the model isn't just parroting patterns but truly reading the document. In our benchmark, this grounding becomes the litmus test: *can your model not only answer the question but also show its work?* The key challenge is to answer questions by identifying the correct region inside a visual academic note (via bounding box), while also classifying both the local content type (e.g. flowchart, chemical formulae, graph, table, math equation, etc.) and the global domain (e.g., physics, computer science) of the answer—encouraging a full spectrum of document understanding, from retrieval to reasoning.

Now, picture a student preparing for an oral exam, trying to answer a complex, open-ended question like, "*Why does underflow occur in floating point representation?*" The answer probably lives somewhere in a series of lecture notes, in different subjects, on different days. To respond, the student first narrows it down: this sounds like computer science or numerical methods. Then, they dive into their notes, flipping through the relevant sections and scanning for definitions, diagrams, or formulas. Only then do they piece together an answer—sometimes from multiple fragments, maybe with a margin sketch or an annotated equation to help. This is the essence of **open-domain question answering** over academic notes. But unlike vanilla open-domain QA systems [35] or M3DocVQA [12] that pull structured facts from curated sources

like Wikipedia, our task demands that the model acts like a student, knowing where to look, what to extract, and how to connect the dots visually and semantically. The challenge isn't just to answer but to find a filter and reason across noisy, visual content. In our setup, models must do three things: first, predict the domain (e.g., physics, mathematics, chemistry) from the question to guide retrieval; second, retrieve relevant pages of handwritten notes—not structured web articles or knowledge bases; and third, perform multimodal reasoning across diagrams, equations, and prose, often scattered across the page.

Though significant advancements have been witnessed in recent times with the emergence of Vision-Language Models (VLMs) [23, 52] and retrieval-augmented approaches [12], note-based document VQA remains a challenging, underexplored problem. Unlike structured document databases [42, 61], scientific notes exhibit *unstructured layouts, informal writing styles, and a lack of clear segmentation*, making traditional OCR-based approaches unreliable. Furthermore, questions often require **multimodal reasoning**: understanding mathematical derivations, interpreting scientific notation, or localizing diagrams that visually complement textual explanations. VLMs often excel by picking up on textual cues, but for notes, visual context (like an arrow connecting a note to a figure or text style like underlining) can change the meaning. Many models do not fully leverage these cues. OCR-based pipelines [8, 53] "throw away" visual information by converting image to text, thus losing layout, handwriting style, or markings as shown in Figure 1. As one recent work [15] states, modern document retrieval systems
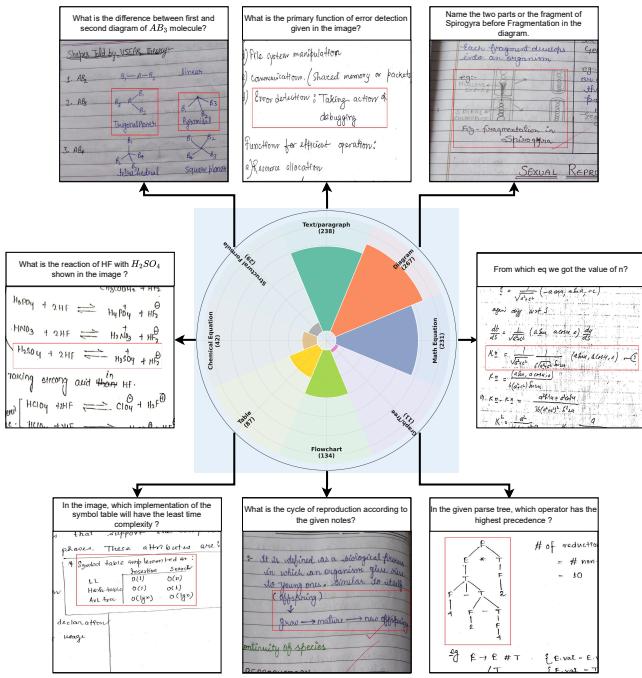
Figure 2: *Illustrative examples from the NoTeS-Bank* showing diverse local reasoning categories such as equations, flowcharts, structural formulas, and textual answers. Center: Distribution of question types across eight task-level categories. Surrounding: Sample questions and annotated evidence regions from handwritten notes, highlighting the visual and semantic complexity handled in the Evidence-Based VQA task.



Figure 3: **(Top)** *Global category distribution* across 19 scientific and technical domains in the NoTeS-Bank benchmark, including physics, biology, chemistry, and computer science. This diverse coverage enables a fine-grained evaluation of domain-specific reasoning. **(Bottom)** ANLS performance comparison on the Evidence-Based VQA task. Results are shown for human annotators, open/closed Vision-Language Models (VLMs), and OCR+LLM pipelines. The large performance gap highlights the challenge of accurately answering and grounding questions in visually unstructured, handwritten academic notes.

mainly rely on extracted text and miss "key visual cues" like figures or layout, limiting their capability on complex documents. This limitation underscores the need for approaches that *treat the image of the document as a first-class input, not just the transcribed text.* In essence, current VQA or VDU models have a **modality gap** – they handle printed text well, but performance drops on unstructured, handwritten, or multi-modal inputs. This is exactly the challenge NoTeS-BANK is designed to address. While most document understanding datasets focus on either layout segmentation or holistic document classification, our benchmark introduces a *dual-layer annotation schema* that grounds each QA instance both at the local (task) level—such as interpreting equations, diagrams, or tables—and the global (domain) level, encompassing scientific fields like physics, biology, or computer science as shown in Figure 3. This fine-grained semantic labeling enables deeper insight into model behavior across multimodal reasoning types and subject domains, making NoTeS-Bank not only a QA benchmark but also a diagnostic tool for analyzing the scope and limits of VLMs and multimodal RAGs. The key contributions of this work can be summarized as follows: 1) We present NoTeS-Bank, a novel benchmark for question answering over unstructured, scientific notes, addressing a gap in multimodal document understanding by focusing on visio-graphical content beyond printed or structured
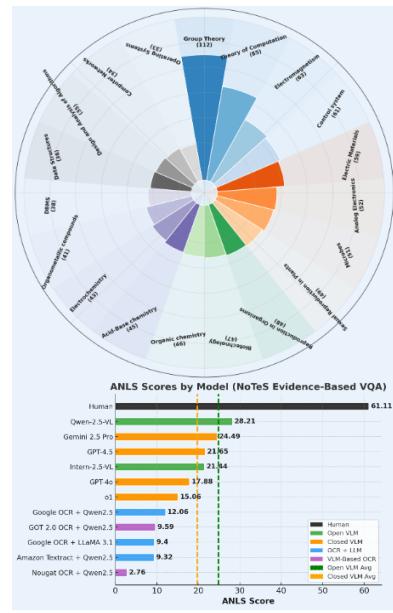
formats. 2) We define two tasks—Evidence-Based VQA and Open-Domain VQA—that jointly evaluate answer grounding, domain classification, and retrieval-based multimodal reasoning in challenging handwritten scenarios. 3) We benchmark a diverse set of VLMs, OCR+LLM pipelines, and retrieval-augmented approaches, and provide a comprehensive evaluation framework using ANLS, IoU, Recall@K, and MRR to highlight the modality gap in current systems.

## 2 Related work

VQA provides a natural language interface for tackling diverse vision-language tasks, merging computer vision and natural language processing (NLP) techniques. This approach has been widely applied across multiple domains, including medical question answering [21, 40, 44], open-domain knowledge retrieval [31, 32, 39, 55], emotion recognition [7, 16], code-based QA [2, 28], logical reasoning [30, 57], fact verification [19, 58], and mathematical reasoning [10, 18, 33, 59].

The field of VDU has been fueled by new benchmarks that invoke a multimodal understanding of Document VQA models. Earlier datasets like DocVQA [38, 48] and InfographicsVQA [36] evaluated reading comprehension on single pages or images, but recent

**Table 1: Comparison of benchmarks on content type, document setting, domain, and task type.**

| Benchmark | Content Type | Multi Document | Domain | Tasks |
|---|---|:---:|---|---|
| LongBench [6] | Text | ✓ | Wikipedia | Long-form QA, Retrieval |
| MPDocVQA [49] | Text, Tables, Charts | ✗ | Multi-domain | Document Visual QA |
| ∞Bench [60] | Text | ✗ | Multi-domain | List QA, Reasoning |
| DUDE [51] | Text, Tables, Charts, Figures | ✗ | Multi-domain | Document Visual QA, Muti-hop QA, Unanswerable, List QA |
| MMLONGBENCH-DOC [34] | Text, Tables, Charts, Slides | ✗ | Multi-domain | Document QA, List QA |
| M3DocVQA [12] | Text, Tables, Charts | ✓ | Wikipedia | Open-domain Document VQA |
| VisDoMBench [45] | Text, Tables, Charts, Slides | ✓ | Multi-domain | Evidence-based Visual Grounding with Bbox, Open-domain QA |
| **NoTeS-Bank (Ours)** | **Graphical Diagram, Math Equation, Chemical Equation, Structural formula, Text/paragraph, Graph/Tree, Flowchart, Table** | ✓ | **Multi-domain (Scientific), Multi-Task** | **Evidence-based Visual Grounding with Bbox and Semantic Labeling, Open-domain VQA, Multi-hop QA, Unanswerable QA, Reasoning** |

benchmarks cover more diverse and realistic scenarios. Notably, the Document Understanding Dataset and Evaluation (DUDE) [50, 51] is a large-scale multi-page, multi-domain DocVQA benchmark which spanned documents from many industries and layouts, with questions requiring reasoning across lengthy documents. It reported that even SOTA models (layout-aware Transformers [20, 54] and VLMs [1, 3, 9, 25, 26]) perform far below human accuracy in DUDE, highlighting the difficulty of generalizing across domains [13]. Another benchmark, SlideVQA [46], focuses on presentation slides with complex layouts, while TableVQA [22] and ScreenUI [5] tasks have introduced questions requiring understanding tables or user interface-like documents. To evaluate the understanding of long documents, MMLongBench-Doc [34] was proposed as part of a long-context multimodal suite featuring approximately 50-page scientific articles and reports to test how models handle complex sequential layouts. In the open-source community, MMDocBench [62] compiled 15 diverse document tasks (from receipts and research papers to diagrams) with over 4,000 QA pairs to compare various large VLMs in a zero-shot setting. This benchmark revealed strengths and weaknesses of models like GPT-4V [1], LLaVA [29], and InternVL [11] on fine-grained document tasks, and it provides a comprehensive testbed for OCR-free document understanding across formats. Finally, new datasets are coupling document images with open-domain knowledge: for example, Vis-DomRAG [45] evaluates QA systems in multi-document settings with rich multimodal content (tables, charts, and presentation slides) while M3DocVQA [12] does it for enterprise documents for testing multimodal Retrieval Augmented Generation (RAG) approaches. The surge of these benchmarks - from DUDE's industry-spanning documents to M3DocVQA's open domain corpus (as shown in Table 1) is driving the field toward more robust and generalizable document understanding, ensuring that modern VLMs are evaluated on reading, reasoning, and retrieving information from the full

variety of documents encountered in the wild. Building on this momentum, our work takes a significant step further by introducing NoTeS-Bank, *the first benchmark focused entirely on handwritten, purely visual, and unstructured academic notes* — challenging current models to reason without structured text anchors or reliable OCR, and establishing a new frontier for document understanding in the wild.

## 3 The NoTeS-Bank Benchmark Suite

NoTeS-Bank is a gold-standard evaluation benchmark designed to assess multimodal question answering over complex unstructured scientific notes. Unlike existing DocVQA datasets [36–38, 48, 51], NoTeS-BANK introduces two distinct tasks that challenge vision-language models (VLMs) and multimodal retrieval-augmented generative (RAG) architectures.

### 3.1 Evidence-Based VQA

The Evidence-Based VQA (EB-VQA) task in NoTeS-BANK evaluates a model's ability to retrieve, comprehend, and justify answers using handwritten note-based evidence. Unlike existing DocVQA challenges that rely solely on extracted OCR text, this task requires models to reason over visual semantics, structural elements, and handwritten symbols while ensuring explicit grounding of responses.

**Task Formulation:** Given an input visual note (image) $I$ (which could span 1-3 pages) containing unstructured text, symbols, equations and diagrams, and a natural language question $Q$, the model must: (a) *Retrieve Relevant Evidence*: Identify key portions of the visual note $I$ that contribute to answering $Q$ by means of a bounding box or multiple bounding boxes. (b) *Generate an Answer*: Synthesize a natural language response $A$ based on the retrieved evidence $E$. (c) *Provide Justification*: Highlight the supporting evidence $E$

**Figure 4:** *Qualitative comparison of Vision-Language Models (VLMs), OCR+LLMs, and human responses on the NoTeS-Bank Evidence-Based VQA task.* Each example demonstrates the challenge of retrieving grounded answers from handwritten scientific notes, highlighting the limitations of current models in detecting accurate regions and reasoning over domain-specific content. The figure also illustrates the fine-grained local (e.g., Structural Formula, Flowchart) and global (e.g., Organic Chemistry, Reproduction in Organisms) category annotations provided for each question-answer pair.

in $I$ that links to the final answer, including the corresponding visual elements (e.g., mathematical equations, chemical formulas, diagrams).

Formally, the model is defined as:

$$A, E = f_{\text{EB-QA}}(I, Q) \tag{1}$$

where the evidence set $E$ consists of:

$$E = \{(B_i, L_i, G_i)\}_{i=1}^{P} \tag{2}$$

where: - $B_i$ represents the bounding box of the relevant evidence region, - $L_i$ denotes the local category of the evidence (e.g., equation, table, diagram), - $G_i$ specifies the global category related to the document's conceptual domain (e.g., group theory, rotational mechanics).

**Evaluation:** To evaluate model performance, we assess answer accuracy and evidence selection quality. Answer accuracy is measured using Average Normalized Levenshtein Similarity (ANLS), while evidence selection is evaluated through Intersection-over-Union (IoU), which quantifies alignment between predicted evidence $E$ and the ground truth $E^*$.

$$IoU(E, E^*) = \frac{|E \cap E^*|}{|E \cup E^*|} \tag{3}$$

Additionally, we measure the correctness of local and global element categorization inside, ensuring that the models retrieve

not only the appropriate text regions but also the relevant semantic concepts necessary for reasoning.

## 3.2 Open-Domain Question Answering

The Open-Domain QA (OD-QA) task in NoTeS-BANK evaluates a model's ability to retrieve, reason, and generate answers across a large collection of handwritten notes. Unlike standard document QA tasks that operate within a single document, this task requires models to first classify the domain of the question, retrieve the most relevant handwritten document, and then generate an answer.

Given a document collection $D$ and a natural language question $Q$, the model must predict the subject category $C$, retrieve the most relevant document $I$, and generate the final answer $A$:

$$C = f_{\text{domain}}(Q), \quad I = f_{\text{retrieve}}(D, Q, C), \quad A = f_{\text{answer}}(I, Q) \tag{4}$$

where: - $C$ represents the predicted subject category (e.g., physics, mathematics). - $I$ is the retrieved handwritten document. - $A$ is the generated answer.

Unlike traditional retrieval-based QA, NoTeS-BANK requires models to handle noisy, unstructured, and multimodal content, including mathematical expressions, diagrams, and scientific notations. The retrieval component is evaluated using Hit@K, MRR, and NDCG@5 to assess ranking quality, while answer accuracy is measured through ANLS. Additionally, the correctness of document and page selection is assessed to ensure models retrieve and process the most relevant content.
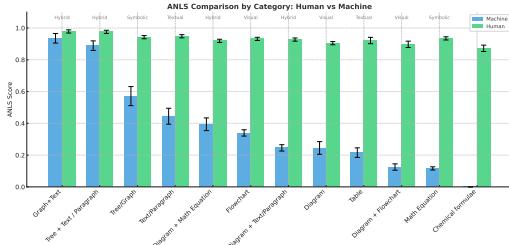
**Figure 5: We report the average ANLS for the human expert vs. the best-performing model per diagnostic category as a ceiling analysis.**

## 3.3 Dataset Collection and Annotation

**Data Collection.** The dataset was collected from various educational websites and student study materials, including resources such as Gate Overflow and self-study notes. For websites where the materials were hosted for students, explicit permission was obtained from the respective content owners, who consented to the use of their data for research purposes. The dataset includes a diverse range of handwritten notes covering subjects such as physics, mathematics, chemistry, biology, and engineering.

**Annotation Process.** A team of 15–20 annotators, primarily undergrad students, was assembled to annotate the dataset. Multiple training sessions were conducted before starting the annotation process to ensure consistency in question formulation and evidence selection. These sessions familiarized annotators with the criteria for both tasks and helped standardize the approach across different document types.

**Task 1: Evidence-Based QA Annotation.** For the Evidence-Based QA task, annotators first created question-answer pairs by formulating queries that required reasoning over textual and multimodal content. Each QA pair was recorded along with metadata, including whether the document was single-page or multi-page, the page number where the answer appeared, and the subject name. To further categorize the nature of the answer, the region from which the answer was derived was labeled as a local category, such as text, equation, diagram, or chemical formula. Each QA pair was also assigned a global category corresponding to its conceptual domain, such as group theory or rotational mechanics.

Once the QA pairs were created, annotators manually labeled the corresponding answer regions by drawing bounding boxes around the relevant content using an annotation tool. These bounding boxes, along with the associated QA pairs and document images, were compiled into a structured JSON format for experimentation and evaluation.

**Task 2: Open-Domain QA Annotation.** A separate annotation process was conducted for the Open-Domain QA task to ensure that retrieval-based reasoning was accurately represented. A new set of QA pairs was created to require retrieval across multiple documents rather than direct extraction from a single page. In this task, the domain classification of each QA pair was explicitly recorded to assist in retrieval, ensuring that models could infer the subject category before searching for the answer. Additionally, annotators identified

the ground truth document and the specific page from which the answer should be retrieved. Given the fundamental differences between the two tasks, the dataset was annotated in two independent rounds to maintain consistency and avoid overlap in annotation strategies. This approach ensured that models trained on the dataset would be evaluated on both localized, evidence-grounded reasoning and retrieval-based question answering across an open-domain handwritten note collection.

## 3.4 Baselines and Model Selection

To evaluate performance on the Evidence-Based VQA task, we establish diverse baselines covering vision-language models (VLMs), OCR-based pipelines, and retrieval-augmented approaches. These baselines assess how different model architectures handle handwritten documents, reasoning, and evidence selection.

**Vision-Language Models (VLMs).** VLMs holistically process visual and textual features, making them well-suited for handwritten notes where OCR struggles. We benchmark both *open* (Qwen-2.5-VL, Intern-2.5-VL, LLaVA) and *closed* (GPT-4o, Gemini 2.5-Pro, GPT-4.5, O1) models to analyze their multimodal reasoning capabilities.

**OCR + LLM-Based Models.** Traditional OCR pipelines extract text before reasoning, but handwritten documents introduce recognition errors. We evaluate Google OCR, Amazon Textract, and OLM OCR combined with LLaMA 3.1 to measure OCR impact on answer quality.

**Layout + OCR + LLM Models.** Standard OCR pipelines discard document structure. We introduce layout-aware approaches incorporating region- and word-level cues (e.g., Textract + Layout Prompt) to assess document retrieval and reasoning improvements.

**VLM-Based OCR Models.** Instead of explicit text extraction, models such as Nougat OCR and GOT 2.0 OCR attempt direct transcription using vision-language understanding. These baselines highlight the limitations of OCR-free approaches in handwriting recognition. Models are compared across *ANLS* (answer similarity), *IoU* (evidence selection), and *category accuracy* (domain classification) to capture end-to-end document understanding. By selecting these baselines, we establish a comprehensive benchmark to drive advancements in handwritten document QA.

**Human performance.** For human evaluation, we collected responses from domain-aware individuals who were not involved in the dataset annotation process. Each participant independently answered questions and highlighted the corresponding evidence regions within the handwritten documents. This separation ensured unbiased assessment across both tasks. Human responses consistently outperformed automated models in terms of accuracy and grounding, especially for complex, multimodal queries—highlighting the difficulty of the NoTeS-Bank benchmark and the gap between current model capabilities and expert-level understanding.

## 4 Results and Discussion

## 4.1 Evaluation Protocol

Prior work [56] has explored the reasoning capabilities of foundation models in visual tasks, primarily through qualitative analysis. In contrast, our objective with NoTeS-BANK is to establish a

**Table 2: Performance comparison of Open and Closed VLM-Based models, OCR + LLM, Layout + OCR + LLM, and VLM-Based OCR models on the Evidence-Based VQA task in NoTeS-BANK.RL: Region-Level Layout; WL: Word-Level Layout**

| Model | #Param | Context Window | ANLS* | IoU Metrics | | | Category Accuracy (%) | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Avg IoU | IoU@5 | IoU@10 | Local | Global |
| *Open VLM-Based Models* | | | | | | | | |
| Qwen-2.5-VL [52] | 7B | 32K | 28.21 | 0.0136 | 0.0727 | 0.0518 | 11.83 | 4.86 |
| Intern-2.5-VL [11] | 8B | 16k | 21.44 | 0.0097 | 0.049 | 0.0308 | 6.47 | 7.91 |
| LLaVA-OneVision [24] | 8B | 60k | 23.34 | - | - | - | - | - |
| *Closed VLM-Based Models* | | | | | | | | |
| GPT-4o [1] | - | 128k | 17.88 | 0.0122 | 0.036 | 0.0381 | 12.0 | 9 |
| Gemini 2.5 Pro [47] | - | 2M | 24.493 | 0.013 | 0.016 | 0.0367 | 0.2 | - |
| OpenAI o1 [41] | - | 100k | 15.06 | 0.0107 | 0.0565 | 0.0335 | 15.6 | 11.2 |
| GPT-4.5 [1] | - | 128k | 21.65 | 0.0186 | 0.0898 | 0.0579 | 13.4 | 7.19 |
| *OCR + Closed LLM Models* | | | | | | | | |
| Google OCR [17] + Deepseek-R1 [27] | 7B | 128k | 12.46 | - | - | - | 1.04 | 0.2283 |
| *OCR + Open LLM Models* | | | | | | | | |
| Google OCR [17] + Qwen2.5 [52] | 7B | - | 12.06 | 0 | 0 | 0 | 7.6 | 4.8 |
| Amazon Textract [17] + Qwen2.5 [52] | 7B | - | 9.32 | 0 | 0 | 0 | 9.07 | 0.824 |
| Google OCR [17] + LLaMA 3.1 [14] | 8B | - | 9.4 | 0.0077 | 0.04865 | 0.02001 | 5.78 | 3.21 |
| Amazon Textract [17] + LLaMA 3.1 [14] | 8B | - | 7.23 | 0.0032 | 0.0189 | 0.0101 | - | - |
| Amazon Textract [17] + LLaMA 3.1 [14] + RL | 8B | - | 5.37 | 0 | 0 | 0 | 0.4124 | 0.6185 |
| Amazon Textract [17] + LLaMA 3.1 [14] + WL | 8B | - | 7.54 | 0.0042 | 0.0261 | 0.0132 | 0.2062 | 1.6494 |
| Amazon Textract [17] + LLaMA 3.1 [14] + RL + WL | 8B | - | 7.01 | 0.0032 | 0.0065 | 0.0065 | 0.2061 | 1.4432 |
| *VLM-Based OCR Models* | | | | | | | | |
| Nougat OCR [8] + LLaMA 3.1 [14] | - | - | 3.2 | - | - | - | 0 | 0 |
| GOT 2.0 OCR [53] + LLaMA 3.1 [14] | - | - | 1.73 | - | - | - | 0 | 0 |
| olmOCR [43] + LLaMA 3.1 [14] | - | - | 4.86 | 0.0001 | 0.0011 | 0 | 0 | 0 |
| Nougat OCR [8] + Qwen2.5 [52] | - | - | 2.76 | - | - | - | 3.8 | 1.2 |
| GOT 2.0 OCR [53] + Qwen2.5 [52] | - | - | 9.59 | - | - | - | 3.4 | 0.8 |
| olmOCR [43] + Qwen2.5 [52] | - | - | 4.69 | - | - | - | 2.6 | 1.6 |
| **Human Baseline** | - | - | 61.11 | 0.4009 | 0.5 | 0.4312 | 83 | 79 |

systematic and unified evaluation protocol that enables both quantitative and qualitative assessment of vision-language models for symbolic and multimodal reasoning in handwritten scientific documents. We introduce a comprehensive benchmarking strategy for NoTeS-BANK, encompassing both Evidence-Based VQA and Open-Domain QA). The models included in our benchmark range from OCR-enhanced LLMs to open and closed-source vision-language models, as detailed in Table 3. We report results across multiple evaluation dimensions, including answer correctness (ANLS*), evidence localization (IoU), document retrieval (Recall@K, MRR, NDCG@5), and category prediction accuracy in Sec. 4.

In addition to quantitative performance metrics, we also provide qualitative analysis of representative failure cases and model outputs, shedding light on limitations in layout reasoning, symbol understanding, and evidence attribution. Given the relatively stronger performance of GPT-4o in multimodal tasks, we present targeted comparisons against its peers, highlighting both its strengths and persistent challenges.

Through this evaluation framework, NoTeS-BANK aims to serve as a diagnostic benchmark for the next generation of vision-language models in handwritten document understanding and retrieval.

### 4.2 Performance Trends Across Tasks

The evaluation of models on both Evidence-Based VQA and Open-Domain QA in NoTeS-BANK reveals significant challenges in understanding and handling handwritten documents. While Vision-Language Models (VLMs) demonstrate promising capabilities, they still struggle with fine-grained reasoning, symbol interpretation, and multimodal retrieval. For Evidence-Based VQA, models relying solely on OCR pipelines exhibit lower IoU and ANLS scores, reinforcing the limitation of text-only processing for handwritten notes.

**Table 3: Performance comparison of various methods. ANLS measures answer accuracy; R@1, MRR, R@5 and ACC (Global category accuracy) measure page retrieval.**

| Method | Accuracy | Page Retrieval | | | Domain |
| --- | --- | --- | --- | --- | --- |
| | ANLS* | R@1 | MRR | R@5 | ACC |
| *Text-based RAG* | | | | | |
| TF-IDF + LLaMa 3.1 8B | 0.0395 | 0.018 | 0.0314 | 0.058 | 5.8 |
| BM 2.5 + LLaMa 3.1 8B | 0.0721 | 0.034 | 0.0515 | 0.082 | 7.6 |
| Mp Net + LLaMa 3.1 8B | 0.0482 | 0.02 | 0.0387 | 0.076 | 6.6 |
| Minilm + LLaMa 3.1 8B | 0.0401 | 0.014 | 0.0245 | 0.038 | 8 |
| ColQwen + Qwen2VL | 0.3419 | 0.218 | 0.243 | 0.288 | 30.6 |
| ColPali + Qwen2-VL 7B | 0.3294 | 0.212 | 0.243 | 0.29 | 30.4 |
| **Human Baseline** | 0.8667 | 0.8125 | 0.8125 | – | 28.99 |

VLMs show improved performance by incorporating visual and structural cues, but evidence localization remains a major challenge.

In Open-Domain QA, retrieval-augmented generation (RAG) models outperform traditional retrievers like BM25 and DPR, particularly in Recall@5 and NDCG@5 metrics. However, even the strongest models, such as GPT-4o RAG and Qwen-2.5-VL RAG, struggle with long-context retrieval over handwritten documents, suggesting a need for better indexing and retrieval over sparse visual information.

**Impact of OCR on Document Understanding:** OCR-based methods perform significantly worse in both tasks, particularly for handwritten mathematical equations, symbols, and complex scientific notations. Models such as Google OCR + LLaMA 3.1 and Textract OCR + LLaMA 3.1 suffer from: *(i)* Loss of spatial and semantic relationships is crucial for layout-heavy content. *(ii)* Difficulty in

**Figure 6: *Qualitative comparison of Open-Domain VQA performance in the NoTeS-Bank benchmark.* Each query requires retrieving relevant handwritten pages from a large corpus and reasoning across them to answer the question. The figure highlights model predictions from retrieval-augmented VLMs (ColPali + Qwen2VL and ColQwen + Qwen2VL) alongside human responses and ground-truth annotations. Differences in predicted answers, retrieved documents, and domain classification underscore the challenge of joint retrieval, domain inference, and multimodal reasoning over noisy, unstructured visual content.**

transcribing non-standard handwritten characters. *(iii)* Inability to provide reliable evidence grounding due to segmentation errors.

This highlights the limitations of treating handwritten document QA as a text-only problem, reinforcing the necessity for joint vision-language reasoning.

**Vision-Language Models and the Multimodal Challenge:** While closed VLMs (GPT-4o, Gemini 1.5-Pro) outperform open VLMs in answer generation, both categories struggle with localizing relevant evidence. Intern-2.5-VL and Qwen-2.5-VL show promise in handling handwritten content but fail to generalize across different domain categories. For multimodal retrieval, OFA RAG and LLaVA RAG improve retrieval accuracy but still fail to effectively fuse retrieved document context into reasoning steps. This suggests the need for better cross-modal pretraining strategies that explicitly model symbolic and spatial dependencies.

### 4.3 Error Analysis and Limitations

Several error patterns emerge from our analysis:

**(i)** Ambiguous Questions: Some models fail due to question ambiguity, producing hallucinated responses instead of recognizing unanswerable questions.

**(ii)** Failure to Retrieve Key Evidence: Even top-performing models frequently retrieve the wrong document, leading to incomplete answers.

**(iii)** Weak Layout Awareness: Many models struggle with layout-based reasoning, especially for tables, structured lists, and diagrams.

A deeper study of failure cases in IoU-based evidence retrieval suggests that bounding-box predictions remain inconsistent, requiring improved fine-grained region detection strategies.

### 5 Conclusion

Our evaluation of state-of-the-art models on NoTeS-BANK highlights significant limitations in existing document QA models, particularly in handling handwritten and multimodal content. While VLMs offer promising capabilities, our results suggest that handwritten document understanding remains an unsolved challenge, requiring improvements in multimodal retrieval, evidence-based reasoning, and layout-aware processing. The findings from NoTeS-BANK establish a benchmark for future research, encouraging the development of more context-aware, symbolically grounded, and visually structured models. Our findings suggest several promising research directions. Future models should integrate spatially-aware tokenization to enhance symbol and diagram reasoning. Moreover, models need better exposure to handwritten mathematical, chemical, and engineering content. Additionally, open-domain QA models should explore hierarchical retrieval approaches, combining vector-based retrieval with structural document parsing. Enhancing evidence attribution and transparency will improve trust in automatic document intelligence systems.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Rajas Agashe, Srinivasan Iyer, and Luke Zettlemoyer. 2019. JuICe: A Large Scale Distantly Supervised Dataset for Open Domain Context-based Code Generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 5436–5446. doi:10.18653/v1/D19-1546

[3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems* 35 (2022), 23716–23736.

[4] Srikar Appalaraju, Bhavan Jasani, Bhargava Urala Kota, Yusheng Xie, and R Manmatha. 2021. Docformer: End-to-end transformer for document understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 993–1003.

[5] Gilles Baechler, Srinivas Sunkara, Maria Wang, Fedir Zubach, Hassan Mansoor, Vincent Etter, Victor Cărbune, Jason Lin, Jindong Chen, and Abhanshu Sharma. 2024. Screenai: A vision-language model for ui and infographics understanding. *arXiv preprint arXiv:2402.04615* (2024).

[6] Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508* (2023).

[7] Johannes Bjerva, Nikita Bhutani, Behzad Golshan, Wang-Chiew Tan, and Isabelle Augenstein. 2020. SubjQA: A Dataset for Subjectivity and Review Comprehension. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 5480–5494. doi:10.18653/v1/2020.emnlp-main.442

[8] Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. 2023. Nougat: Neural optical understanding for academic documents. *arXiv preprint arXiv:2308.13418* (2023).

[9] Soumitri Chattopadhyay, Sanket Biswas, Emanuele Vivoli, and Josep Lladǫs. 2024. Towards generative class prompt learning for fine-grained visual recognition. *arXiv preprint arXiv:2409.01835* (2024).

[10] Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric Xing, and Liang Lin. 2021. GeoQA: A Geometric Question Answering Benchmark Towards Multimodal Numerical Reasoning. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 513–523. doi:10.18653/v1/2021.findings-acl.46

[11] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271* (2024).

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

[13] Alloy Das, Sanket Biswas, Umapada Pal, and Josep Lladós. 2024. Diving into the depths of spotting text in multi-domain noisy scenes. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 410–417.

[14] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).

[15] Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2024. Colpali: Efficient document retrieval with vision language models. In *The Thirteenth International Conference on Learning Representations*.

[16] Lin Gui, Jiannan Hu, Yulan He, Ruifeng Xu, Qin Lu, and Jiachen Du. 2017. A Question Answering Approach for Emotion Cause Extraction. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Copenhagen, Denmark, 1593–1602. doi:10.18653/v1/D17-1167

[17] Thomas Hegghammer. 2022. OCR with Tesseract, Amazon Textract, and Google Document AI: a benchmarking experiment. *Journal of Computational Social Science* 5, 1 (2022), 861–882.

[18] Mark Hopkins, Ronan Le Bras, Cristian Petrescu-Prahova, Gabriel Stanovsky, Hannaneh Hajishirzi, and Rik Koncel-Kedziorski. 2019. SemEval-2019 Task 10: Math Question Answering. In *Proceedings of the 13th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, Minneapolis, Minnesota, USA, 893–899. doi:10.18653/v1/S19-2153

[19] Xuming Hu, Zhijiang Guo, GuanYu Wu, Aiwei Liu, Lijie Wen, and Philip Yu. 2022. CHEF: A Pilot Chinese Dataset for Evidence-Based Fact-Checking. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Seattle, United States, 3362–3376. doi:10.18653/v1/2022.naacl-main.246

[20] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. 2022. LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking. *arXiv preprint arXiv:2204.08387* (2022).

[21] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. PubMedQA: A Dataset for Biomedical Research Question Answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 2567–2577. doi:10.18653/v1/D19-1259

[22] Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. 2024. Tablevqa-bench: A visual question answering benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205* (2024).

[23] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024. What matters when building vision-language models? *Advances in Neural Information Processing Systems* 37 (2024), 87874–87907.

[24] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. 2024. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326* (2024).

[25] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*. PMLR, 19730–19742.

[26] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*. PMLR, 12888–12900.

[27] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437* (2024).

[28] Chenxiao Liu and Xiaojun Wan. 2021. CodeQA: A Question Answering Dataset for Source Code Comprehension. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. Association for Computational Linguistics, Punta Cana, Dominican Republic, 2618–2632. doi:10.18653/v1/2021.findings-emnlp.223

[29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems* 36 (2023), 34892–34916.

[30] Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. LogiQA: A Challenge Dataset for Machine Reading Comprehension with Logical Reasoning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, Christian Bessiere (Ed.). International Joint Conferences on Artificial Intelligence Organization, 3622–3628. doi:10.24963/ijcai.2020/501 Main track.

[31] Jiahua Liu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2019. XQA: A Cross-lingual Open-domain Question Answering Dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 2358–2368. doi:10.18653/v1/P19-1227

[32] Shayne Longpre, Yi Lu, and Joachim Daiber. 2020. MKQA: A Linguistically Diverse Benchmark for Multilingual Open Domain Question Answering. doi:10.48550/ARXIV.2007.15207

[33] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255* (2023).

[34] Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, et al. 2024. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. *arXiv preprint arXiv:2407.01523* (2024).

[35] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*. 3195–3204.

[36] Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. 2022. InfographicVQA. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1697–1706.

[37] Minesh Mathew, Lluis Gomez, Dimosthenis Karatzas, and CV Jawahar. 2021. Asking questions on handwritten document collections. *International Journal on Document Analysis and Recognition (IJDAR)* 24, 3 (2021), 235–249.

[38] Minesh Mathew, Ruben Tito, Dimosthenis Karatzas, R Manmatha, and CV Jawahar. 2020. Document visual question answering challenge 2020. *arXiv preprint arXiv:2008.08899* (2020).

[39] Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering Ambiguous Open-domain Questions. doi:10.48550/ARXIV.2004.10645

[40] Anastasios Nentidis, Georgios Katsimpras, Eirini Vandorou, Anastasia Krithara, Antonio Miranda-Escalada, Luis Gasco, Martin Krallinger, and Georgios Paliouras. 2022. Overview of BioASQ 2022: The Tenth BioASQ Challenge on Large-Scale Biomedical Semantic Indexing and Question Answering. In *Lecture Notes in Computer Science*. Springer International Publishing, 337–361. doi:10.1007/978-3-

031-13643-6_22

[41] OpenAI, :, and Aaron Jaech et al. 2024. OpenAI o1 System Card. arXiv:2412.16720 [cs.AI] https://arxiv.org/abs/2412.16720

[42] Birgit Pfitzmann, Christoph Auer, Michele Dolfi, Ahmed S Nassar, and Peter Staar. 2022. Doclaynet: A large human-annotated dataset for document-layout segmentation. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*. 3743–3751.

[43] Jake Poznanski, Jon Borchardt, Jason Dunkelberger, Regan Huff, Daniel Lin, Aman Rangapur, Christopher Wilhelm, Kyle Lo, and Luca Soldaini. 2025. olmOCR: Unlocking Trillions of Tokens in PDFs with Vision Language Models. *arXiv preprint arXiv:2502.18443* (2025).

[44] Preethi Raghavan, Jennifer J Liang, Diwakar Mahajan, Rachita Chandra, and Peter Szolovits. 2021. emrKBQA: A Clinical Knowledge-Base Question Answering Dataset. In *Proceedings of the 20th Workshop on Biomedical Language Processing*. Association for Computational Linguistics, Online, 64–73. doi:10.18653/v1/2021.bionlp-1.7

[45] Manan Suri, Puneet Mathur, Franck Dernoncourt, Kanika Goswami, Ryan A Rossi, and Dinesh Manocha. 2024. VisDoM: Multi-Document QA with Visually Rich Elements Using Multimodal Retrieval-Augmented Generation. *arXiv preprint arXiv:2412.10704* (2024).

[46] Ryota Tanaka, Kyosuke Nishida, Kosuke Nishida, Taku Hasegawa, Itsumi Saito, and Kuniko Saito. 2023. Slidevqa: A dataset for document visual question answering on multiple images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 13636–13645.

[47] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* (2023).

[48] Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny. 2021. Document collection visual question answering. In *Document Analysis and Recognition–ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part II 16*. Springer, 778–792.

[49] Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny. 2022. Hierarchical multimodal transformers for Multi-Page DocVQA. *arXiv preprint arXiv:2212.05935* (2022).

[50] Jordy Van Landeghem, Lukasz Borchmann, Rubèn Tito, Michał Pietruszka, Dawid Jurkiewicz, Rafał Powalski, Paweł Józiak, Sanket Biswas, Mickaël Coustaty, and Tomasz Stanisławek. 2023. ICDAR 2023 Competition on Document UnderstanDing of Everything (DUDE). In *Proceedings of ICDAR 2023*.

[51] Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann, Michał Pietruszka, Pawel Joziak, Rafal Powalski, Dawid Jurkiewicz, Mickaël Coustaty, Bertrand Anckaert, Ernest Valveny, et al. 2023. Document understanding dataset and evaluation (dude). In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 19528–19540.

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

[53] Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, et al. 2024. General OCR Theory: Towards OCR-2.0 via a Unified End-to-end Model. *arXiv preprint arXiv:2409.01704* (2024).

[54] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1192–1200.

[55] Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. WikiQA: A Challenge Dataset for Open-Domain Question Answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Lisbon, Portugal, 2013–2018. doi:10.18653/v1/D15-1237

[56] Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421* 9, 1 (2023), 1.

[57] Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. ReClor: A Reading Comprehension Dataset Requiring Logical Reasoning. In *International Conference on Learning Representations (ICLR)*.

[58] Majid Zarharan, Mahsa Ghaderan, Amin Pourdabiri, Zahra Sayedi, Behrouz Minaei-Bidgoli, Sauleh Eetemadi, and Mohammad Taher Pilehvar. 2021. ParsFEVER: a Dataset for Farsi Fact Extraction and Verification. In *Proceedings of *SEM 2021: The Tenth Joint Conference on Lexical and Computational Semantics*. Association for Computational Linguistics, Online, 99–104. doi:10.18653/v1/2021.starsem-1.9

[59] Qiyuan Zhang, Lei Wang, Sicheng Yu, Shuohang Wang, Yang Wang, Jing Jiang, and Ee-Peng Lim. 2021. NOAHQA: Numerical Reasoning with Interpretable Graph Question Answering Dataset. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. Association for Computational Linguistics, Punta Cana, Dominican Republic, 4147–4161. doi:10.18653/v1/2021.findings-emnlp.350

[60] Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, et al. 2024. InfinityBench: Extending Long Context Evaluation Beyond 100K Tokens. *arXiv preprint arXiv:2402.13718* (2024).

[61] Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. 2019. Publaynet: largest dataset ever for document layout analysis. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 1015–1022.

[62] Fengbin Zhu, Ziyang Liu, Xiang Yao Ng, Haohui Wu, Wenjie Wang, Fuli Feng, Chao Wang, Huanbo Luan, and Tat Seng Chua. 2024. MMDocBench: Benchmarking Large Vision-Language Models for Fine-Grained Visual Document Understanding. *arXiv preprint arXiv:2410.21311* (2024).