

# BBox-DocVQA: A Large-Scale Bounding-Box–Grounded Dataset for Enhancing Reasoning in Document Visual Question Answer

Anonymous CVPR submission

Paper ID 3031

## Abstract

**Document Visual Question Answering (DocVQA)** is a fundamental task for multimodal document understanding and a key testbed for vision–language reasoning. However, most existing DocVQA datasets are limited to the page level and lack fine-grained spatial grounding, constraining the interpretability and reasoning capability of Vision–Language Models (VLMs). To address this gap, we introduce **BBox-DocVQA**—a large-scale, bounding-box–grounded dataset designed to enhance spatial reasoning and evidence localization in visual documents. We further present an automated construction pipeline, **Segment–Judge–and–Generate**, which integrates a segment model for region segmentation, a VLM for semantic judgment, and another advanced VLM for question–answer generation, followed by human verification for quality assurance. The resulting dataset contains **3.6K** diverse documents and **32K** QA pairs, encompassing single- and multi-region as well as single- and multi-page scenarios. Each QA instance is grounded on explicit bounding boxes, enabling fine-grained evaluation of spatial–semantic alignment. Benchmarking multiple state-of-the-art VLMs (e.g., GPT-5, Qwen2.5-VL, and InternVL) on BBox-DocVQA reveals persistent challenges in spatial grounding and reasoning accuracy. Furthermore, fine-tuning on BBox-DocVQA substantially improves both bounding-box localization and answer generation, validating its effectiveness for enhancing the reasoning ability of VLMs. Our dataset and code will be publicly released to advance research on interpretable and spatially grounded vision–language reasoning.

## 1. Introduction

Document Visual Question Answering (DocVQA) has emerged as a crucial benchmark for advancing multimodal document understanding, requiring models to reason over textual, structural, and visual information jointly. Recent

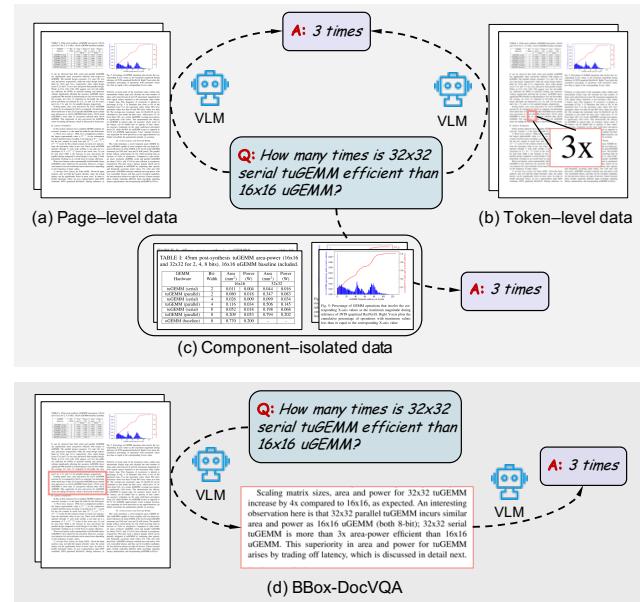


Figure 1. Comparison among (a) page–level data, (b) token–level data, (c) component–isolated data, and (d) our proposed BBox-DocVQA.

progress in vision–language models (VLMs) [1–3, 28, 40] has greatly improved their ability to process and comprehend complex document layouts. However, the majority of existing DocVQA datasets, such as DocVQA [22], MP-DocVQA [35], and VisualMRC [31], are constructed at the *page level*. Correspondingly, most existing methods [6, 7, 9, 10, 12, 16–18, 26, 27, 34, 37–39, 41, 42] focus on generating the final answer based on one or multiple document pages, without explicit modeling of where the evidence originates. This absence of spatial grounding leads to reduced interpretability and may degrade the generation accuracy of VLMs, as models can yield seemingly correct answers while attending to irrelevant or incomplete regions.

Human document comprehension typically involves two steps: first, identifying the relevant pages, and second, locating specific *evidence components*—such as text para-

036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051

Dataset / benchmark	#Queries	#Docs	#Images	#I/Q	Avg. Q Len	Avg. Ans Len	BBox	MP
DUDE [36]	41 K	5 K	28.7 K	1.00	8.65	3.35	Yes	No
BoundingDocs [11]	249 K	48 K	237 K	1.00	–	–	Yes	No
SlideVQA [32]	14.5 K	–	52 K	1.26	–	–	Yes	Yes
SPIQA [25]	270 K	25.9 K	270 K	1.00	12.98	14.56	–	No
PlotVQA [24]	29 M	–	224 K	1.00	43.54	–	–	No
ArXivQA [15]	100 K	16.6 K	32 K	1.00	16.98	7.59	–	No
ChartQA [21]	33 K	–	22 K	1.00	13.18	1.08	–	No
TAT-DQA [43]	16 K	2.7 K	3 K	1.07	12.54	3.44	No	No
VisualMRC [31]	30 K	10 K	10 K	1.00	10.55	9.55	No	No
InfographicsVQA [23]	30 K	5.4 K	5.4 K	1.00	11.54	1.60	No	No
DocVQA [22]	50 K	6 K	12 K	1.00	9.49	2.43	No	No
MP-DocVQA [35]	46 K	6 K	48 K	8.27	9.90	2.20	No	Yes
BBox-DocVQA (ours)	32 K	3.6 K	44k	1.38	10.80	2.78	Yes	Yes

Table 1. Comparison between the BBox-DocVQA dataset and the existing DocVQA datasets. Docs, I, Q, Ans, BBox, and MP denote documents, images, queries, answers, bounding boxes, and multiple pages, respectively.

graphs, tables, or figures—that support the final answer. Current page-level DocVQA datasets, however, only support the first step, offering page-level evidence without bounding-box annotations. On the other hand, a few token-level datasets, such as DUDE [36] and BoundingDocs [11], provide OCR-based token-level bounding boxes to localize answers. While these datasets allow precise token grounding, they lack contextual semantic integrity, as each annotation typically corresponds to an isolated word or phrase rather than a complete semantic unit. As a result, existing DocVQA resources cannot fully capture the fine-grained reasoning process required for spatially interpretable document understanding.

To bridge these gaps, we introduce **BBox-DocVQA** (Bounding-Box-grounded Document Visual Question Answering), the first large-scale DocVQA dataset that explicitly annotates evidence bounding boxes across both single- and multi-page settings. Each question-answer (QA) pair is linked to one or multiple bounding boxes representing coherent semantic regions (paragraphs, tables, or figures) that provide the necessary evidence for answering the query. In addition, we develop an automated dataset construction pipeline, named **Segment-Judge-and-Generate**, which combines segmentation, visual-semantic evaluation, and automatic QA generation. Specifically, our pipeline first employs a segmentation model (SAM[14]) to detect visual and textual regions, then uses a VLM to judge the effectiveness of the segments, and finally prompts an advanced VLM to generate the corresponding QA pairs. This pipeline enables scalable and consistent dataset generation while maintaining semantic diversity and spatial precision.

We conduct extensive experiments on the proposed BBox-DocVQA dataset using a range of state-of-the-art

VLMs, including GPT-5 [1], Qwen2.5-VL [3], and the InternVL [8] series. These models are evaluated in terms of both spatial grounding accuracy (bounding-box localization) and answer correctness. Experimental results reveal that even advanced VLMs often struggle to accurately identify the true evidence regions, underscoring a persistent weakness in spatial grounding and interpretability. To further explore the utility of our dataset, we fine-tune Qwen2.5-VL on the BBox-DocVQA training set and observe substantial improvements in both bounding-box localization and answer generation. This demonstrates that bounding-box-grounded supervision can significantly enhance the reasoning transparency and accuracy of multimodal models.

The main contributions of this paper are summarized as follows:

1. We introduce **BBox-DocVQA**, the first large-scale bounding-box-grounded DocVQA dataset, which covers 3.6 K diverse visual documents and 32 K QA pairs, encompassing both single- and multi-page, as well as single- and multi-region scenarios.
2. We propose an automated dataset construction pipeline, **Segment-Judge-and-Generate**, integrating segmentation, visual-semantic judgment, and large VLM-based QA generation, ensuring scalability and annotation consistency.
3. We conduct extensive benchmarking and fine-tuning experiments with multiple state-of-the-art VLMs, demonstrating that BBox-DocVQA substantially enhances models' spatial grounding and answer accuracy, offering a new foundation for interpretable document understanding research.

118

## 2. Related work

119

Visual document understanding has been widely explored through diverse dataset paradigms, which differ in annotation granularity, component coverage, and interpretability. In this section, we review three major types of document visual question–answering datasets—*token-level*, *component-isolated*, and *page-level*—and clarify how our proposed **BBox-DocVQA** bridges their existing limitations. A detailed statistical comparison is provided in Table 1.

120

**Token-level Datasets.** A branch of VQA datasets, such as DUDE [36] and BoundingDocs [11], and Slide-VQA [32], provides token-level annotation, where bounding boxes are aligned with OCR tokens and serve as ground-truth evidence for question answering. These datasets directly associate answers with OCR-extracted tokens, thereby facilitating fine-grained text localization. However, they often lack contextual semantic information, since each answer is limited to a single token or phrase rather than a semantically complete region. Consequently, token-level datasets may fail to capture the hierarchical and compositional nature of document understanding.

121

**Component-isolated Datasets.** Another line of research focuses on component-isolated visual question answering, targeting specific structured elements such as tables, charts, and figures. Representative examples include SPIQA [25], PlotQA [24], ChartQA [21], and ArXivQA [15]. These datasets generate massive numbers of QA pairs by isolating document components and synthesizing questions from them. Such datasets can enhance the reasoning ability of vision–language models (VLMs) over structured visual content. Nonetheless, their task formulation deviates from DocVQA: they do not capture inter-component or cross-page reasoning, and the generated QA pairs typically remain confined to single visual objects rather than full document contexts.

122

**Page-level Datasets.** To address holistic understanding, page-level DocVQA datasets annotate question–answer pairs at the document or page level, requiring models to comprehend both text layout and visual structure. For example, the DocVQA [22] dataset and its multi-page extension MP-DocVQA [35] comprise 50 K questions over 12 K single-page documents and 46 K questions spanning 6 K multi-page documents, respectively. Similarly, VisualMRC [31], InfographicsVQA [23], and TAT-DQA [43] extend document-level comprehension to specific domains. Also, a few benchmarks, including OK-VQA [20], A-OKVQA [29], WebQA [4], UDA [13], Dyn-VQA [16], MMLongBench [19], REAL-MM-RAG [39], ViDoRe [10], ViDoSeek [37], M3DoCVQA [9], OpenDocVQA [33], VisR-Bench [5], and VisDoMBench [30], have been proposed to evaluate the VLMs’ understanding and reasoning capabilities. These datasets or benchmarks advance layout-aware VQA; however, they do not provide explicit

bounding-box annotations for evidence localization, limiting explainability and making it difficult to evaluate spatial grounding.

To fill this gap, we introduce **BBox-DocVQA**, a large-scale, bounding-box–grounded dataset for visual document question answering. Unlike prior datasets, BBox-DocVQA explicitly annotates evidence bounding boxes, enabling fine-grained evaluation of both *spatial grounding* and *semantic correctness*. It further extends to multi-page scenarios and includes questions whose answers depend on multiple bounding boxes across different pages. As shown in Table 1, BBox-DocVQA combines the interpretability of token-level datasets with the holistic reasoning of page-level benchmarks, achieving a balance between accuracy, scalability, and explainability. This design not only establishes a new benchmark for VLM spatial reasoning but also enables model evaluation on interpretable, evidence-based document-understanding tasks.

## 3. BBox-DocVQA dataset

In this section, we introduce the collection and construction details of BBox-DocVQA, along with its statistics.

### 3.1. Segment-Judge-and-Generate

Existing approaches based on either standalone Vision-Language Models (VLMs) or conventional Computer Vision (CV) methods struggle to accurately perform fine-grained structure extraction from documents. These methods often fail to produce bounding boxes (bbox) that are both precise and semantically complete. To address this issue, we propose an automated *Segment-Judge-and-Generate* framework to efficiently construct a large-scale, bounding-box-grounded document question answering (DocQA) dataset. The core idea is to first detect fine-grained semantic components (*segments*) within document pages, and then generate reasoning-oriented question–answer (QA) pairs conditioned on these localized regions. This enables an explicit alignment between visual content, spatial layout, and semantic reasoning. The overall pipeline and case study are illustrated in Fig. 2 and 3.

**Data Collection.** We collected approximately **4,000** PDF documents from the *arXiv* platform, spanning the years 2023 to 2025, with a total of about **137,000** pages. The corpus covers eight major academic domains, including Computer Science, Mathematics, Physics, Biology and other disciplines, representing all top-level categories in arXiv. To facilitate subsequent visual analysis and multimodal reasoning, each page was converted into a high-resolution image to preserve layout structure, textual organization, and figure/table details.

**Segment Components.** We adopt the Segment Anything Model [14] (SAM, ViT-H variant) to automatically detect layout components on each page image. The model

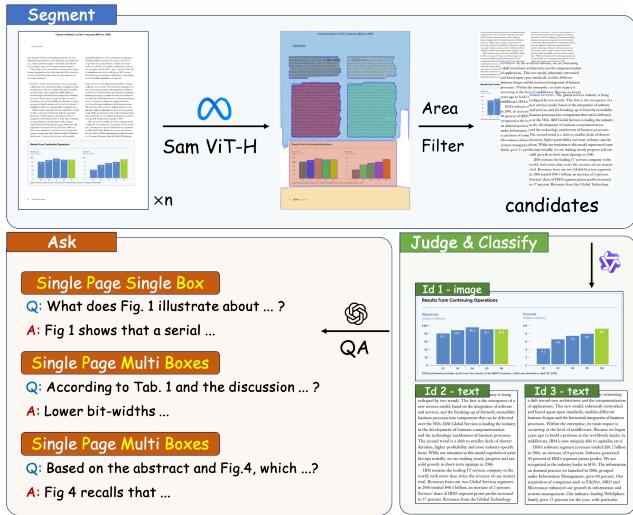


Figure 2. The proposed **Segment-Judge-and-Generate** pipeline for constructing a bounding-box-grounded DocQA dataset. It includes three stages: (1) Segment — detect fine-grained layout components via SAM (ViT-H); (2) Judge and Classify — evaluate and categorize regions with Qwen2.5-VL-72B; and (3) Generate — generate reasoning-oriented QA pairs with GPT-5.

first produces multiple semantic masks, which are then converted into compact bounding boxes. To remove trivial or redundant regions, we filter the bounding boxes based on their area ratio relative to the page size, retaining only those with coverage between 5% and 70%. This effectively discards tiny noisy fragments and full-page detections.

**Judge and Classify Components.** The detected regions typically correspond to independent semantic units such as paragraphs, tables, or figures, but overlapping or ambiguous boundaries may still occur. To ensure higher quality and semantic consistency, we employ **Qwen2.5-VL-72B** to evaluate each cropped region in terms of completeness, information density, and visual cleanliness, and to predict its primary content type (text, table, or image). Furthermore, to eliminate redundant detections, we design an overlap-based deduplication strategy: if the overlapping area between two regions exceeds 90% of the smaller one, we retain different regions depending on their type — for **text**, we preserve the smaller region to capture fine-grained semantics; for **table** and **image**, we keep the larger region to ensure contextual and structural integrity.

**Question–Answer Generation.** After obtaining high-quality fine-grained regions, we employ **GPT-5** to generate diverse and reasoning-oriented QA pairs. We first perform page-level summarization to provide contextual background, ensuring semantic consistency across regions. The generated QA samples are categorized into three types: (1) *Single-Page Single-BBox (SPSBB)*, (2) *Single-Page Multi-BBox (SPMBB)*, and (3) *Multi-page Multi-*

*BBox (MPMBB)*, covering a wide range of document-level reasoning scenarios. For single-region cases, the model receives the page summary and the cropped sub-image, and generates QA pairs based strictly on visible content. For multi-region cases, the model first evaluates the semantic relevance among sub-images and only produces questions when meaningful relationships are detected, thereby avoiding incoherent or trivial QAs. Each valid region is instructed to yield one QA pair through multi-step reasoning and factual grounding. Every QA sample contains a concise question and an accurate answer.

In summary, the **Segment-Judge-and-Generate** framework provides an efficient and scalable automated solution for constructing fine-grained, spatially grounded, and semantically rich document QA datasets, laying a solid foundation for region-level DocQA model training and evaluation.

### 3.2. Fine-grained Benchmark

To further evaluate the effectiveness of our proposed **Segment-Judge-and-Generate** framework, we construct a high-quality manually annotated fine-grained benchmark, referred to as the **Fine-grained Benchmark**.

**Data Sources and Selection.** We randomly selected **10 PDF documents** from each of the eight major academic domains on the *arXiv* platform (Computer Science, Mathematics, Physics, Quantitative Biology, Quantitative Finance, Electrical Engineering and Systems Science, Economics, and Statistics), resulting in a total of **80** research papers as the benchmark corpus. During the initial filtering phase, we conducted manual inspection to ensure that the sampled documents exhibit sufficient diversity and representativeness in layout structure, page count, content type (text, table, and image), and subject area.

**Manual Annotation and Verification.** Domain experts were invited to participate in the annotation process. Each page was manually cropped to precisely delineate every semantic unit with bounding boxes (bboxes). To guarantee the reliability and consistency of the annotations, we adopted a multi-stage verification protocol: each sample was independently annotated by at least two experts and subsequently reviewed by a third annotator to resolve inconsistencies, yielding a set of high-confidence bounding boxes.

**Question Generation and Quality Control.** After obtaining the verified bbox regions, we employed the state-of-the-art multimodal large model **GPT-5** to generate question-answer (QA) pairs for each region. The generation process strictly follows the same protocol as in the training data, covering three QA formats: *Single-Page Single-BBox (SPSBB)*, *Single-Page Multi-BBox (SPMBB)*, and *Multi-Page Multi-BBox (MPMBB)*. Human experts further conducted multiple rounds of validation to assess the logical soundness of questions, the factual correctness of answers,

Segment	251
Ask	252
Judge & Classify	253
candidates	254
QA	255
Id 1 - image	256
Id 2 - text	257
Id 3 - text	258
260	259
261	260
262	261
263	262
264	263
265	264
266	265
267	266
268	267
269	268
270	269
271	270
272	271
273	272
274	273
275	274
276	275
277	276
278	277
279	278
280	279
281	280
282	281
283	282
284	283
285	284
286	285
287	286
288	287
289	288
290	289
291	290
292	291
293	292
294	293
295	294
296	295
297	296
298	297
299	298
300	299
301	300
302	301

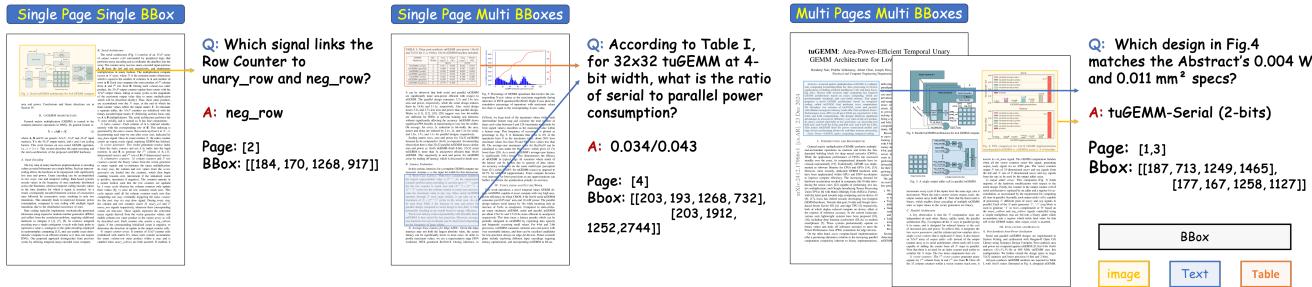


Figure 3. Case study illustrating different document understanding settings in the BBox-DocVQA dataset. Examples include (left) single-page single-bbox reasoning, (middle) single-page multi-bbox reasoning, and (right) multi-page multi-bbox reasoning, where questions and answers are grounded on one or more spatially localized regions (bboxes) across document pages.

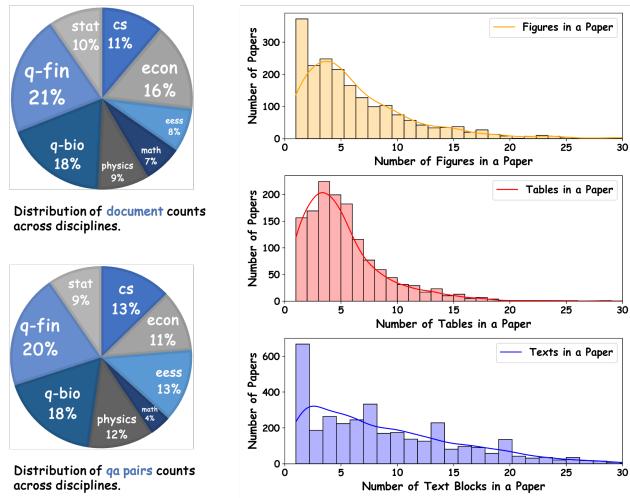


Figure 4. The left panels show the proportions of papers across research domains for (top) the document dataset and (bottom) the generated QA pairs. The right panels display the distributions of the number of figures, tables, and text blocks per paper.

303 and the overall reasoning quality, ensuring that each QA  
304 pair is interpretable, accurate, and challenging.

305 This benchmark covers a wide variety of document  
306 structures and semantic types, providing a comprehensive  
307 evaluation resource for testing the performance and  
308 generalization ability of Vision-Language Models (VLMs) on  
309 fine-grained document understanding tasks.

### 3.3. Data Statistics and Analysis

311 This section summarizes the key statistics of the BBox-  
312 DocVQA dataset, including the automatically constructed  
313 training set and the manually curated fine-grained bench-  
314 mark. By presenting the data distribution, region com-  
315 position, and domain coverage in a concise tabular form, we  
316 highlight the structural alignment between the two subsets  
317 and their respective roles in model training and evaluation.

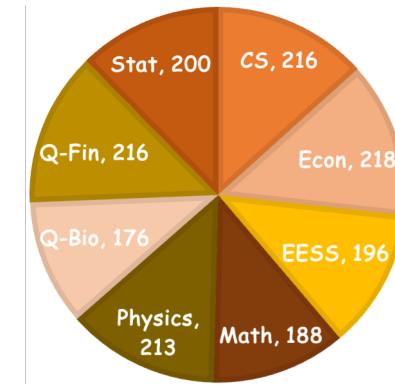


Figure 5. Domain-wise distribution of QA samples in the Fine-grained Benchmark.

#### 3.3.1. Training Set Statistics

The training set contains 30,780 automatically generated QA pairs derived from 42,380 pages of 3,671 arXiv papers. As summarized in Table 2, the dataset features a balanced distribution across the three task types (SPSBB, SPMBB, MPMBB), along with a rich mixture of text, image, and table regions dominated by textual content. The scientific domain coverage naturally follows arXiv’s real-world submission distribution, while the long-tailed variation in the number of figures, tables, and text blocks per paper reflects authentic heterogeneity in scientific document layouts. These characteristics collectively provide a large-scale and structurally diverse corpus suitable for training models on multi-region and multi-modal document understanding.

#### 3.3.2. Benchmark Statistics

The Fine-grained Benchmark consists of 1,623 high-quality manually annotated QA samples spanning 1,941 pages from 80 papers. Table 2 provides a detailed summary of its composition. Similar to the training set, the benchmark includes all three task types and a multimodal mixture of region categories, but differs in its uniformly sampled domain distribution and precisely curated bounding boxes and QA anno-

318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331

332  
333  
334  
335  
336  
337  
338  
339

Table 2. Comparison of Training Set and Fine-grained Benchmark Statistics

Item	Training Set	Benchmark
Total QA samples	30,780	1,623
Total pages	42,380	1,941
Total papers	3,671	80
<b>Task types</b>		
SPSBB	11,668 (37.91%)	749 (46.15%)
SPMBB	7,512 (24.41%)	556 (34.26%)
MPMBB	11,600 (37.69%)	318 (19.59%)
<b>Region types</b>		
Text	30,424 (60.98%)	1,247 (49.94%)
Image	12,542 (25.14%)	916 (36.68%)
Table	6,926 (13.88%)	334 (13.38%)
Avg. bbox area ratio	14.26%	19.08%

tations. The benchmark therefore serves as a reliable evaluation suite with balanced domain coverage, stringent annotation consistency, and a suitable level of task difficulty for assessing fine-grained document understanding.

**3.3.3. Comparison Between Training Set and Benchmark**  
The training set is large and diverse, making it well suited for learning multimodal document structures and multi-region reasoning patterns. The benchmark, while smaller, offers highly reliable manual annotations and balanced domain sampling, making it ideal for robust and fair performance evaluation. Their aligned distributions in task types and region types ensure consistent training–evaluation conditions and meaningful performance interpretation.

## 4. Task and model training

In this section, we introduce the DocVQA task definition and the proposed two-stage training strategy using the BBox-DocVQA dataset.

**357 Bounding-box grounded Doc-VQA task.** Given a document  $D$  and a query  $q$ , a retriever  $\mathcal{R}$  first fetches relevant  
358 pages  $P_q = \mathcal{R}(q, D)$ . A generator  $\mathcal{G}$  then produces the  
359 final answer  $\hat{a}$  based on a prompt  $\mathcal{P}$ . We explore three  
360 strategies that integrate bounding-box ( $\hat{b}$ ) prediction: (1)  
361 **Direct Answer Generation:** The answer is generated di-  
362 rectly without explicit localization, (2) **Simultaneous Gen-  
363 eration:** The model generates the answer and its bounding  
364 box concurrently, and (3) **Sequential Generation:** A two-  
365 stage, coarse-to-fine approach involving localization, crop-  
366 ping, and then answer generation. The following equations  
367 formally represent these strategies:  
368

$$\hat{a} = \mathcal{G}(\mathcal{P}(q, P_q)), \quad (1)$$

$$\hat{a}, \hat{b} = \mathcal{G}(\mathcal{P}(q, P_q)), \quad (2)$$

$$\begin{aligned} \hat{b} &= \mathcal{G}(\mathcal{P}(q, P_q)), & P_q^{\text{crop}} &= \text{CROP}(P_q, \hat{b}), \\ \hat{a} &= \mathcal{G}(\mathcal{P}(q, P_q^{\text{crop}})) \end{aligned} \quad (3) \quad 372$$

**Two-stage training strategy.** We adopt a two-stage training strategy to progressively align the model’s visual grounding and reasoning capabilities. In the first stage, a LoRA fine-tuning is performed to teach the model to localize the evidence region (BBox) corresponding to each query. This step focuses purely on spatial grounding, allowing the model to build a stable correspondence between textual prompts and visual layouts before introducing semantic supervision. In the second stage, the model is further trained to jointly predict the bounding box and generate the textual answer. This joint optimization bridges grounding and generation, encouraging the model to reason based on localized visual evidence rather than relying on global context. All implementation details, including hyperparameters, optimization settings, and loss composition, are provided in the Appendix. Although the constructed dataset can also benefit the retriever by offering explicit evidence supervision, we leave retriever-level enhancement for future work and focus here on improving the generation stage.

## 5. Experiment

### 5.1. Experiment setup

#### Baselines.

**Evaluation metrics.** We evaluate the model performance on both visual grounding and answer generation. For bounding-box prediction, we use the **Intersection over Union (IoU)** as the metric. Specifically, for **Single-Page Single-BBox (SPSBB)** cases, IoU is computed directly between the predicted and ground-truth boxes; for **Single-Page Multiple-BBox (SPMBB)** cases, we calculate the maximum IoU for each ground-truth box and then average across all boxes; for **Multi-Page Multiple-BBox (MPMBB)** cases, IoU is first computed for each page (following the SPSBB or SPMBB rule) and then averaged across pages. For answer evaluation, we employ a **large language model (LLM) as a semantic judge** (i.e., DeepSeek-v3.1) to assess whether the generated response is consistent with the reference answer in meaning and context. This LLM-based evaluation captures semantic correctness beyond surface-level text similarity, offering a more reliable measure for open-ended question answering.

#### Implementation details.

### 5.2. Main results

### 5.3. Analysis and discussion

## 6. Conclusion

Model	Mean	SPSB (749)	SPMB (556)	MPMB (318)
Qwen2.5VL-3B	30.87%	31.38%	33.27%	25.47%
Qwen2.5VL-7B	51.57%	57.01%	53.42%	35.53%
Qwen2.5VL-32B	63.46%	67.69%	66.55%	48.11%
Qwen2.5VL-72B	68.64%	71.03%	71.58%	57.86%
Qwen3VL-4B	68.95%	70.49%	72.12%	59.75%
Qwen3VL-8B	67.59%	70.23%	73.20%	51.57%
Qwen3VL-32B	77.14%	81.04%	84.35%	55.35%
InternVL3-2B	34.20%	39.12%	33.27%	24.21%
InternVL3-8B	50.46%	53.14%	53.06%	39.62%
GPT-5	<b>81.45%</b>	<b>82.64%</b>	83.63%	<b>74.84%</b>

Table 3. Answer accuracy under the simultaneous perception-and-answering setting, where the model must localize evidence regions and generate answers in a single unified step. Results are reported for SPSB, SPMB, and MPMB tasks.

Model	Mean IoU	SPSB (749)	SPMB (556)	MPMB (318)
Qwen2.5VL-3B	4.70%	3.80%	5.60%	4.90%
Qwen2.5VL-7B	11.30%	14.60%	8.80%	8.00%
Qwen2.5VL-32B	20.00%	22.20%	20.40%	14.30%
Qwen2.5VL-72B	<b>35.20%</b>	<b>40.10%</b>	<b>33.20%</b>	<b>27.20%</b>
Qwen3VL-4B	18.70%	18.90%	17.00%	21.30%
Qwen3VL-8B	14.40%	17.60%	9.20%	15.90%
Qwen3VL-32B	20.40%	22.60%	17.30%	21.00%
InternVL3-2B	0.10%	0.00%	0.20%	0.10%
InternVL3-8B	0.30%	0.10%	0.50%	0.50%
GPT-5	0.90%	0.10%	1.60%	1.20%

Table 4. Mean IoU under the simultaneous perception-and-answering setting, reported across SPSBB, SPMBB, and MPMBB task types.

Model	GT Subpage Ans.	GT Page Ans.
Qwen2.5VL-3B	56.38%	45.72%
Qwen2.5VL-7B	68.58%	57.18%
Qwen2.5VL-32B	82.01%	73.01%
Qwen2.5VL-72B	79.42%	69.32%
Qwen3VL-4B	72.34%	67.22%
Qwen3VL-8B	74.74%	72.27%
Qwen3VL-32B	<b>82.01%</b>	<b>80.28%</b>
InternVL3-2B	41.59%	32.04%
InternVL3-8B	65.80%	53.42%

Table 5. Answer accuracy when ground-truth evidence regions are provided (*GT Subpage*) compared to answering based only on the entire page image (*GT Page*). Across all models, supplying the precise evidence region consistently improves accuracy, demonstrating that **localization substantially enhances answer quality**.

## 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. 418
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine 419 Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: 420 a visual language model for few-shot learning. *Advances 421 in neural information processing systems*, 35:23716–23736, 422 2022.
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin 423 Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun 424 Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint 425 arXiv:2502.13923*, 2025. 1, 2
- [4] Yingshan Chang, Mridu Narang, Hisami Suzuki, Guihong 426 Cao, Jianfeng Gao, and Yonatan Bisk. Webqa: Multihop 427 and multimodal qa. In *Proceedings of the IEEE/CVF conference 428 on computer vision and pattern recognition*, pages 16495– 429 16504, 2022. 3
- [5] Jian Chen, Ruiyi Zhang, Yufan Zhou, Tong Yu, Franck Der- 430 noncourt, Jiaxiang Gu, Ryan A Rossi, Changyou Chen, and 431 Tong Sun. Sv-rag: Lora-contextualizing adaptation of mllms 432 for long document understanding. In *The Thirteenth International Conference on Learning Representations*, 2025. 3
- [6] Wenhui Chen, Hexiang Hu, Xi Chen, Pat Verga, and William 433 Cohen. Murag: Multimodal retrieval-augmented generator 434

- 445 for open question answering over images and text. In *Proceedings of the 2022 Conference on Empirical Methods in*  
446 *Natural Language Processing*, pages 5558–5570, 2022. 1
- 447 [7] Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang.  
448 Can pre-trained vision and language models answer visual information-seeking questions? *arXiv preprint*  
449 *arXiv:2302.11713*, 2023. 1
- 450 [8] Zhe Chen, Jianne Wu, Wenhui Wang, Weijie Su, Guo Chen,  
451 Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu,  
452 Lewei Lu, et al. Internvl: Scaling up vision foundation  
453 models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision*  
454 and pattern recognition, pages 24185–24198, 2024. 2
- 455 [9] Jaemin Cho, Debanjan Mahata, Ozan Irsoy, Yujie He, and  
456 Mohit Bansal. M3docrag: Multi-modal retrieval is what you  
457 need for multi-page multi-document understanding. *arXiv preprint arXiv:2411.04952*, 2024. 1, 3
- 458 [10] Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani,  
459 Gautier Viaud, CELINE HUDELOT, and Pierre Colombo.  
460 Colpali: Efficient document retrieval with vision language  
461 models. In *The Thirteenth International Conference on Learning Representations*, 2025. 1, 3
- 462 [11] Simone Giovannini, Fabio Coppini, Andrea Gemelli, and Si-  
463 mone Marinai. Boundingdocs: a unified dataset for doc-  
464 ument question answering with spatial annotations. *arXiv preprint arXiv:2501.03403*, 2025. 2, 3
- 465 [12] Wenbo Hu, Jia-Chen Gu, Zi-Yi Dou, Mohsen Fayyaz, Pan  
466 Lu, Kai-Wei Chang, and Nanyun Peng. Mrag-bench: Vision-  
467 centric evaluation for retrieval-augmented multimodal mod-  
468 els. In *The Thirteenth International Conference on Learning Representations*, 2025. 1
- 469 [13] Yulong Hui, Yao Lu, and Huachen Zhang. Uda: A bench-  
470 mark suite for retrieval augmented generation in real-world  
471 document analysis. *Advances in Neural Information Processing Systems*, 37:67200–67217, 2024. 3
- 472 [14] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,  
473 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-  
474 head, Alexander C Berg, Wan-Yen Lo, et al. Segment any-  
475 thing. In *Proceedings of the IEEE/CVF international confer-  
476 ence on computer vision*, pages 4015–4026, 2023. 2, 3
- 477 [15] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng,  
478 Lingpeng Kong, and Qi Liu. Multimodal arxiv: A dataset for  
479 improving scientific comprehension of large vision-language  
480 models. *arXiv preprint arXiv:2403.00231*, 2024. 2, 3
- 481 [16] Yangning Li, Yinghui Li, Xinyu Wang, Yong Jiang, Zhen  
482 Zhang, Xinran Zheng, Hui Wang, Hai-Tao Zheng, Fei  
483 Huang, Jingren Zhou, et al. Benchmarking multimodal re-  
484 trieval augmented generation with dynamic vqa dataset and  
485 self-adaptive planning agent. In *The Thirteenth International  
486 Conference on Learning Representations*, 2025. 1, 3
- 487 [17] Weizhe Lin, Jinghong Chen, Jingbiao Mei, Alexandru Coca,  
488 and Bill Byrne. Fine-grained late-interaction multi-modal  
489 retrieval for retrieval augmented visual question answering.  
490 *Advances in Neural Information Processing Systems*, 36:  
491 22820–22840, 2023.
- 492 [18] Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhui Chen,  
493 and Jimmy Lin. Unifying multimodal retrieval via document  
494 screenshot embedding. In *Proceedings of the 2024 Confer-  
495 ence on Empirical Methods in Natural Language Processing*,  
496 pages 6492–6505, 2024. 1
- 497 [19] Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu  
498 Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong,  
499 et al. Mmlongbench-doc: Benchmarking long-context docu-  
500 ment understanding with visualizations. *Advances in Neural  
501 Information Processing Systems*, 37:95963–96010, 2024. 3
- 502 [20] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and  
503 Roozbeh Mottaghi. Ok-vqa: A visual question answering  
504 benchmark requiring external knowledge. In *Proceedings  
505 of the IEEE/cvf conference on computer vision and pattern  
506 recognition*, pages 3195–3204, 2019. 3
- 507 [21] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty,  
508 and Enamul Hoque. Chartqa: A benchmark for question an-  
509 swering about charts with visual and logical reasoning. *arXiv  
510 preprint arXiv:2203.10244*, 2022. 2, 3
- 511 [22] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar.  
512 Docvqa: A dataset for vqa on document images. In *Proceed-  
513 ings of the IEEE/CVF winter conference on applications of  
514 computer vision*, pages 2200–2209, 2021. 1, 2, 3
- 515 [23] Minesh Mathew, Viraj Bagal, Rubén Tito, Dimosthenis  
516 Karatzas, Ernest Valveny, and CV Jawahar. Infographicvqa.  
517 In *Proceedings of the IEEE/CVF Winter Conference on Ap-  
518 plications of Computer Vision*, pages 1697–1706, 2022. 2,  
519 3
- 520 [24] Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and  
521 Pratyush Kumar. Plotqa: Reasoning over scientific plots.  
522 In *Proceedings of the ieee/cvf winter conference on appli-  
523 cations of computer vision*, pages 1527–1536, 2020. 2, 3
- 524 [25] Shraman Pramanick, Rama Chellappa, and Subhashini  
525 Venugopalan. Spiqa: A dataset for multimodal question an-  
526 swering on scientific papers. *Advances in Neural Informa-  
527 tion Processing Systems*, 37:118807–118833, 2024. 2, 3
- 528 [26] Jingyuan Qi, Zhiyang Xu, Rulin Shao, Yang Chen, Jin Di,  
529 Yu Cheng, Qifan Wang, and Lifu Huang. Rora-vlm: Robust  
530 retrieval-augmented vision language models. *arXiv preprint  
531 arXiv:2410.08876*, 2024. 1
- 532 [27] Zehan Qi, Rongwu Xu, Zhijiang Guo, Cunxiang Wang, Hao  
533 Zhang, and Wei Xu. Long2rag: Evaluating long-context &  
534 long-form retrieval-augmented generation with key point re-  
535 call. In *Findings of the Association for Computational Lin-  
536 guistics: EMNLP 2024*, pages 4852–4872, 2024. 1
- 537 [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya  
538 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
539 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning  
540 transferable visual models from natural language supervi-  
541 sion. In *International conference on machine learning*, pages  
542 8748–8763. PMLR, 2021. 1
- 543 [29] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark,  
544 Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A  
545 benchmark for visual question answering using world knowl-  
546 edge. In *European conference on computer vision*, pages  
547 146–162. Springer, 2022. 3
- 548 [30] Manan Suri, Puneet Mathur, Franck Dernoncourt, Kanika  
549 Goswami, Ryan A Rossi, and Dinesh Manocha. Visdom:  
550 Multi-document qa with visually rich elements using multi-  
551 modal retrieval-augmented generation. In *Proceedings of the  
552*

- 561        2025 *Conference of the Nations of the Americas Chapter of*  
562        *the Association for Computational Linguistics: Human Lan-*  
563        *guage Technologies (Volume 1: Long Papers)*, pages 6088–  
564        6109, 2025. 3
- 565 [31] Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. Vi-  
566        sualmrc: Machine reading comprehension on document im-  
567        ages. In *Proceedings of the AAAI Conference on Artificial*  
568        *Intelligence*, pages 13878–13888, 2021. 1, 2, 3
- 569 [32] Ryota Tanaka, Kyosuke Nishida, Kosuke Nishida, Taku  
570        Hasegawa, Itsumi Saito, and Kuniko Saito. Slidevqa: A  
571        dataset for document visual question answering on multiple  
572        images. In *Proceedings of the AAAI Conference on Artificial*  
573        *Intelligence*, pages 13636–13645, 2023. 2, 3
- 574 [33] Ryota Tanaka, Taichi Iki, Taku Hasegawa, Kyosuke Nishida,  
575        Kuniko Saito, and Jun Suzuki. Vdocrag: Retrieval-  
576        augmented generation over visually-rich documents. In *Pro-*  
577          
578        *Conference*, pages 24827–24837, 2025. 3
- 579 [34] Yang Tian, Fan Liu, Jingyuan Zhang, V. W., Yupeng Hu, and  
580        Liqiang Nie. CoRe-MMRAG: Cross-source knowledge rec-  
581        onciliation for multimodal RAG. In *Proceedings of the 63rd*  
582        *Annual Meeting of the Association for Computational Lin-*  
583        *guistics (Volume 1: Long Papers)*, pages 32967–32982, Vi-  
584        enna, Austria, 2025. Association for Computational Linguis-  
585        tics. 1
- 586 [35] Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny. Hi-  
587        erarchical multimodal transformers for multipage docvqa.  
588        *Pattern Recognition*, 144:109834, 2023. 1, 2, 3
- 589 [36] Jordy Van Landeghem, Rubèn Tito, Łukasz Borchmann,  
590        Michał Pietruszka, Paweł Joziāk, Rafal Powalski, Dawid Ju-  
591        rkiewicz, Mickaël Coustaty, Bertrand Anckaert, Ernest Val-  
592        veny, et al. Document understanding dataset and evaluation  
593        (dude). In *Proceedings of the IEEE/CVF International Con-*  
594        *ference on Computer Vision*, pages 19528–19540, 2023. 2,  
595        3
- 596 [37] Qiuchen Wang, Ruixue Ding, Zehui Chen, Weiqi Wu, Shi-  
597        hang Wang, Pengjun Xie, and Feng Zhao. Vidorag: Visual  
598        document retrieval-augmented generation via dynamic iter-  
599        ative reasoning agents. *arXiv preprint arXiv:2502.18017*,  
600        2025. 1, 3
- 601 [38] Qiuchen Wang, Ruixue Ding, Yu Zeng, Zehui Chen,  
602        Lin Chen, Shihang Wang, Pengjun Xie, Fei Huang, and  
603        Feng Zhao. Vrag-rl: Empower vision-perception-based  
604        rag for visually rich information understanding via itera-  
605        tive reasoning with reinforcement learning. *arXiv preprint*  
606        *arXiv:2505.22019*, 2025.
- 607 [39] Navve Wasserman, Roi Pony, Oshri Naparstek, Adi Raz  
608        Goldfarb, Eli Schwartz, Udi Barzelay, and Leonid Karlinsky.  
609        Real-mm-rag: A real-world multi-modal retrieval bench-  
610        mark. *arXiv preprint arXiv:2502.12342*, 2025. 1, 3
- 611 [40] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
612        Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen  
613        Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv*  
614        *preprint arXiv:2505.09388*, 2025. 1
- 615 [41] Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Junhao Ran,  
616        Yukun Yan, Zhenghao Liu, Shuo Wang, Xu Han, Zhiyuan  
617        Liu, et al. Visrag: Vision-based retrieval-augmented gener-
- ation on multi-modality documents. In *The Thirteenth In-*  
618        *ternational Conference on Learning Representations*, 2025.  
619        1
- 620 [42] Tao Zhang, Ziqi Zhang, Zongyang Ma, Yuxin Chen, Zhon-  
621        gang Qi, Chunfeng Yuan, Bing Li, Junfu Pu, Yuxuan Zhao,  
622        Zehua Xie, et al. mr<sup>2</sup> ag: Multimodal retrieval-reflection-  
623        augmented generation for knowledge-based vqa. *arXiv*  
624        *preprint arXiv:2411.15041*, 2024. 1
- 625 [43] Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang,  
626        Haozhou Zhang, and Tat-Seng Chua. Towards complex doc-  
627        ument understanding by discrete reasoning. In *Proceedings*  
628        *of the 30th ACM International Conference on Multimedia*,  
629        pages 4857–4866, 2022. 2, 3
- 630

# BBox-DocVQA: A Large-Scale Bounding-Box–Grounded Dataset for Enhancing Reasoning in Document Visual Question Answer

## Supplementary Material

631

### 1. Details of dataset construction