

MODALITY-SWAP DISTILLATION: RENDERING TEXTUAL REASONING INTO VISUAL SUPERVISION

Anonymous authors

Paper under double-blind review

ABSTRACT

Visual reasoning over structured data such as tables is a critical capability for modern vision-language models (VLMs), yet current benchmarks remain limited in scale, diversity, or reasoning depth, especially when it comes to rendered table images. Addressing this gap, we introduce **Visual-TableQA**, a large-scale, open-domain multimodal dataset specifically designed to evaluate and enhance visual reasoning over complex tabular data. Our generation pipeline is **modular, scalable, and fully autonomous**, involving multiple reasoning LLMs collaborating across distinct roles: generation, validation, and inspiration. **Visual-TableQA** comprises 2.5k richly structured LaTeX-rendered tables and 9k reasoning-intensive QA pairs, all produced at a cost of under \$100. To promote diversity and creativity, our pipeline performs **multi-model collaborative data generation** via **cross-model prompting** (**‘inspiration’**) and LLM-jury filtering. Stronger models seed layouts and topics that weaker models elaborate, collectively distilling diverse reasoning patterns and visual structures into the dataset. Empirical results show that models fine-tuned on **Visual-TableQA** generalize robustly to external benchmarks, outperforming several proprietary models despite the dataset’s synthetic nature. The full pipeline and resources are publicly available in our GitHub repository.

1 INTRODUCTION

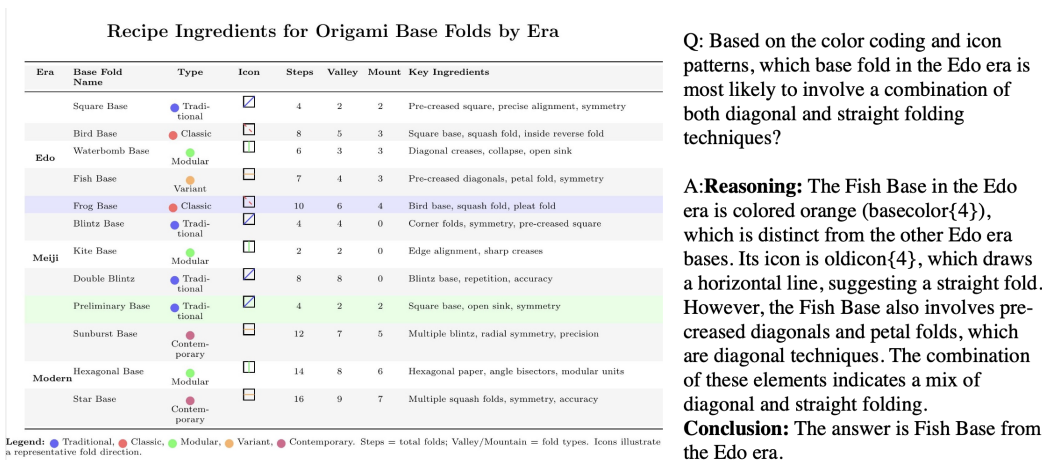


Figure 1: Sample question in our benchmark

Vision-language models (VLMs) have significantly advanced in recent years, achieving remarkable performance in various tasks involving visual and textual inputs. Despite these advancements, complex reasoning tasks, especially those requiring deep comprehension of tabular data structures, continue to pose significant challenges. Table complexity can manifest in various ways, including structural layout, information density, and the diversity of visual components such as the integration of diagrams. The more complex a table is, the more it lends itself to challenging reasoning tasks,

054 requiring advanced cognitive abilities to extract relevant information and perform multi-step logical
 055 analysis. For example, the table in Figure 1 exemplifies this complexity through its use of multi-
 056 row cells, integrated diagrams, and color encoding. Answering the question requires the VLM to
 057 interpret information across all cells and perform a sequence of reasoning steps.

058 Existing table-based QA datasets predominantly fall into two categories: (i) those represented
 059 purely in textual format—such as WikiTableQuestions Pasupat & Liang (2015), HybridQA Chen
 060 et al. (2020b), and AIT-QA Katsis et al. (2022)—which bypass the challenges of visual layout
 061 interpretation; and (ii) those that lack diversity in visual layouts, visual complexity, and reason-
 062 ing depth due to being domain-specific (e.g., TAT-DQA Zhu et al. (2022)), or having standard-
 063 ized queries (e.g., TableVQA-Bench Kim et al. (2024)), or highly technical in nature (e.g., Table-
 064 VQA Tom Agonoude (2024)). This second datasets category typically rely on a limited set of
 065 layout templates and involve relatively simple visual tasks or basic QA scenarios, falling short of
 066 the complexity required for thorough evaluation and advancement of reasoning capabilities. More
 067 recent efforts—such as ChartQA Masry et al., ReachQA He et al., and MATH-Vision Wang et al.
 068 (2024b)—have aimed to address the need for open-domain coverage, incorporating more diverse
 069 visual features, varied question types, and deeper reasoning challenges. However, these datasets pri-
 070 marily focus on charts and function plots, overlooking tables—and with them, an entire dimension
 071 of informational structure and layout diversity. An extensive comparison of diverse chart and table
 072 datasets is provided in Table 1.

073 Inspired by ReachQA’s Code-as-Intermediary Translation (CIT)—a technique that translates chart
 074 images into textual representations while faithfully preserving visual features—we introduce **Visual-**
 075 **TableQA**, a novel synthetic, multimodal, and open-domain dataset tailored to enhance reasoning
 076 capabilities through complex table-based question-answering tasks. Visual-TableQA capitalizes on
 077 the ability of reasoning-oriented LLMs to generate intricate LaTeX tables, thus significantly reduc-
 078 ing costs and eliminating the need for extensive manual annotations. This modality-swap makes
 079 it possible for LLMs to invest their textual reasoning ability into visual image in order to improve
 080 visual understanding and reasoning. Visual-TableQA emphasizes structural reasoning over domain
 081 knowledge. Each entry couples a rendered table image with a complex, visually grounded reason-
 082 ing task. Tasks require interpreting visual layout cues such as cell alignment, hierarchical headers,
 083 merged cells, or embedded symbolic content—emulating real-world documents where visual con-
 084 text is essential for correct interpretation. The dataset contains 2.5k reasoning-intensive tables and
 085 9k QA pairs crafted to assess both information extraction and multi-step reasoning capabilities, all
 086 generated at a cost of under \$100. The entire dataset has been validated using a committee of high-
 087 performing reasoning LLMs, the ROSCOE step by step reasoning score Golovneva et al., and a
 088 sample of 800 QA pairs has undergone manual verification by human annotators. In contrast to
 089 previous synthetic datasets, Visual-TableQA is less guided in its generation process, allowing for
 090 more diversity and creativity in both table complexity (e.g., structural layout, information density,
 091 visual component variety) and the design of QA pairs explicitly crafted to challenge visual reasoning
 092 skills. We evaluated a broad range of VLMs, from lightweight models to state-of-the-art architec-
 093 tures, and benchmarked their performance against existing datasets. The results show that most
 094 VLMs continue to struggle with table understanding.

095 In sum, our main contributions are: (i) a high-quality, visually diverse, and open-domain dataset for
 096 table-based reasoning; (ii) an LLM-driven, low-cost generation pipeline using cross-model inspira-
 097 tion; (iii) an empirical analysis comparing Visual-TableQA to existing table and chart datasets; (iv)
 098 an extensive evaluation of open and proprietary VLMs, showing performance gains after finetun-
 099 ing. Our dataset and code are publicly available at [https://github.com/AI-4-Everyone/
 100 Visual-TableQA](https://github.com/AI-4-Everyone/Visual-TableQA).

101 2 VISUAL-TABLEQA DATASET

102 Unlike previous datasets that rely heavily on textual input or handcrafted annotations, Visual-
 103 TableQA leverages a scalable generation pipeline rooted in LaTeX-rendered table images, auto-
 104 mated reasoning task creation, and LLM-based evaluation. This strategy enables high diversity and
 105 reasoning depth while keeping annotation costs minimal, totaling under \$100 using a combination
 106 of open-access APIs and limited usage tiers. In this section, we describe our LaTeX-based table
 107 encoding 2.1, the data generation pipeline 2.2, and the quality assurance process 2.3.

Table 1: Comparison of existing chart and table datasets across data, Q&A, and dataset properties. Abbreviations: Repr=Representation, Vis= Visual, Comp= Complexity, Temp = Template, Refer = Reference, Rat = Rational, Synth= Synthetic, Scal = Scalable. Cells marked with \blacktriangle indicate mixed attributes (e.g., partially template-based; scalable Q&A but non-scalable chart data)

Datasets	Data Properties				Q&A Properties			Dataset Properties		
	# Layouts/ # Topics	Type	Data Repr.	Vis. Comp.	Temp. Free	Vis. Refer.	Rat. Annot.	Synth.	#Samples / #QA	Scal.
WikiTableQuestions (Pasupat & Liang, 2015)	–	Table	Text	\times	\times	\times	\times	\times	2.1k/22k	\times
HybridQA (Chen et al., 2020b)	–	Table	Text	\times	\checkmark	\times	\times	\times	13k/70k	\times
AIT-QA (Katsis et al., 2022)	-/1	Table	Text	\times	\checkmark	\times	\times	\times	116/515	\times
TAT-DQA (Zhu et al., 2022)	-/1	\blacktriangle	Image	\times	\checkmark	\checkmark	\checkmark	\times	2.5k/16.5k	\times
Table-VQA (Tom Agonoude, 2024)	-/-	Table	Image	\times	\checkmark	\checkmark	\checkmark	\checkmark	16.4k/82.3k	–
TableVQA-Bench (Kim et al., 2024)	11/4	Table	Image	\times	\checkmark	\checkmark	\times	\blacktriangle	894/1.5k	\blacktriangle
ChartQA (Masry et al.)	3/15	Chart	Image	\times	\checkmark	\checkmark	\times	\times	21.9k/32.7k	\times
DocVQA (Mathew et al., 2020)	20/5	\blacktriangle	Image	\times	\checkmark	\checkmark	\times	\times	12.7k/50k	\times
MultiModalQA (Talmor et al., 2021)	16/ ∞	\blacktriangle	Image	\times	\times	\checkmark	\times	\blacktriangle	29,918	\times
MATH-Vision (Wang et al., 2024b)	-/16	\blacktriangle	Image	\checkmark	\checkmark	\checkmark	\times	\times	3k/3k	\times
REACHQA (He et al.)	32/ ∞	Chart	Image	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	3.7k/22k	\checkmark
Visual-TableQA (ours)	/ ∞	Table	Image	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	2.5k/ 9k	\checkmark

Table 2: Model performance on the test sets of four benchmarks: ChartQA, ReachQA, MATH-Vision, Visual-TableQA, and Visual-TableQA-CIT. Visual-TableQA-CIT is the variant of our dataset where tables are represented in LaTeX code form rather than as rendered images. The ReachQA score is reported as the average across its two evaluation splits: *Reasoning* and *Recognition*. The values in blue are from our own evaluation using the LLM jury, while the remaining values are taken from model authors or official leaderboards/model cards. When a fine-tuned model achieves better performance, the result is annotated with \uparrow ; if the performance worsens, it is marked with \downarrow . The best performance for each model variants and task is in **bold**.

Models	ChartQA	ReachQA	MATH-Vision _{FULL}	Visual-TableQA	Visual-TableQA-CIT
Baseline					
Human	–	74.85	68.82	–	–
Proprietary VLMs					
GPT-4o	85.7	53.25	30.39	78.5	91.0
GPT-4o mini	77.52	40.35	28.85	67.0	80.27
Gemini 2.5 Flash	84.64	56.97	41.3	85.72	92.3
Gemini 2.5 Pro	85.73	61.87	73.3	86.63	86.27
Claude 3.5 Sonnet	90.8*	63	32.76	82.46	88.5
Open-Source VLMs					
Llama 4 Maverick 17B-128E Instruct	85.3*	47.98	45.89	80.75	87.0
Mistral Small 3.1 24B Instruct	86.24*	42.45	32.45	73.2	80.25
Qwen2.5-VL-32B-Instruct	79.75	49.5	38.1	80.45	81.93
Qwen2.5-VL-7B-Instruct	87.3*	49.23	25.1	71.35	–
Finetuned VLMs					
Qwen2.5-VL-7B-Instruct + Visual-TableQA	84.52 \uparrow	60.95 \uparrow	49.77 \uparrow	82.98 \uparrow	N/A
Qwen2.5-VL-7B-Instruct + ReachQA	77.59 \downarrow	55.75 \uparrow	48.57 \uparrow	60.68 (56.13) \downarrow	N/A

* Performance metrics are measured using *Relaxed Accuracy*, which allows for small numerical deviations in the predicted answers. We assume that this accuracy inflates the actual accuracy by at least 5%. This margin is subtracted when selecting the best-performing results, which are shown in **bold**.

This section provides a detailed description of the generation pipeline. Figure 2 gives an overview of the whole process.

Seed Tables and Topics Collection: The first step involves collecting a diverse set of table layouts to serve as inspiration for LLMs during the generation process. We explored various sources, including scientific journals, financial report databases, online newspapers, and table design galleries. Our search included both table and diagram images to introduce greater visual and structural complexity into the dataset. We selected 20 representative images (Figure 6a) and passed them to a visual language model, VLM-0 (GPT-_{o3} OpenAI (2025)), to generate accurate LaTeX representations. In parallel, we used LLM-0 (GPT-4_o OpenAI (2024)) to generate a list of 5,000 distinct topic prompts. These initial table samples and topics serve as the first layer of inspiration for subsequent LLM generations—though the pool of inspirations expands automatically, as detailed in Section 2.2. For reproducibility, all resources are publicly available in our GitHub repository.

Table Generation: For each iteration, we randomly select an LLM-1 from the models short-list presented in Table 3. The model receives one table sample from our pool and three topics randomly selected from the topic list, all delivered through a single instruction prompt. The output from LLM-1 is returned as a JSON file containing three newly generated LaTeX-formatted tables in plain text, each corresponding to one of the provided topics. We require that the generated tables be inspired by the input table but include substantial layout variations and, when appropriate, additional data to enhance complexity. The resulting LaTeX code is then compiled using standard LaTeX compilation stack (pdflatex + pdf2image), and cropped to produce high-resolution table images. A human reviewer then inspects the table and makes adjustments to the LaTeX code if necessary. The prompt used for generation are provided in Figure 7.

Evolving Layouts through Iterations: A subset of the generated tables is manually selected to enrich the pool of table inspirations. This feedback loop encourages the emergence of increasingly complex and diverse layouts by amplifying visual variations and enabling cross-model inspiration across different LLM-1s over successive iterations. This process is facilitated by the fact that LLMs differ in architecture and tend to focus on distinct structural and stylistic aspects of tables. As a result, combining inspirations across models leads to highly diversified and creative layout types. We refer to this phenomenon as **cross-model prompting** (**‘inspiration’**).

QA Generation: Next, for each generated table, we randomly select a model, denoted LLM-2, from the same list of models in Table 3 to generate three QA pairs. The model receives the table in LaTeX format and is instructed to produce questions that require multi-step reasoning, pattern recognition, and symbolic interpretation. For instance, the sample in Figure 1 illustrates how the questions extend beyond basic information extraction, requiring interpretative reasoning to identify patterns within the presented data. We do not fact-check the generated tables; as a result, some table content may be non-factual. While this is important to consider when using the dataset for training, it can be beneficial, as it encourages models to rely on reasoning rather than prior knowledge.

2.3 QUALITY ASSURANCE

To ensure the validity of the tables and QA pairs, a panel of independent LLMs—serving as a reasoning jury—evaluates each table and its associated QA pairs by providing binary correctness judgments. The evaluation is based on four criteria: *(i)* the generated document is a valid table and is relevant to the given topic; *(ii)* the table and any associated figures are coherent and meaningful; *(iii)* the question is fully grounded in the table, requiring no external knowledge; and *(iv)* the answer is completely supported by the table content. If any of these four criteria are not met, the corresponding table and its QA pairs are discarded. The LLM jury includes Mistral-large, Deepseek-v3.1, Gemini-2.5-pro, GPT-4.1, and Deepcogito-v2—models chosen for their strong reasoning abilities. Final acceptance is determined via majority vote across the jury. The prompt used is provided in Figure 9.

The next step involved computing the ROSCOE reasoning scores as introduced in Golovneva et al.. These metrics assess the coherence, logical soundness, and contextual grounding of step-by-step generated rationales. The ROSCOE framework encompasses thirteen evaluation criteria, which we report in Table 7 along with their corresponding values computed over our dataset. The results indicate near-perfect alignment with the expected directionality of each metric, supporting the overall quality of the generated reasoning chains.

Test Set Construction and Human Evaluation: The dataset was divided into three subsets: training, validation, and testing. To prevent data leakage, all entries {table, question, answer} derived from a single table were assigned to the same subset. The testing set was also used for human evaluation. Two human annotators—each holding at least a Master’s degree and with prior experience in data annotation—were hired to evaluate the quality of 800 QA pairs. Each QA pair was assessed for validity and rated on a scale from 1 to 5. Overall, 92% of the evaluated QA pairs received a rating of at least 4 stars from both annotators.

3 EXPERIMENTS

3.1 BENCHMARK COMPARISON

Evaluated Benchmarks and Model Selection: We evaluate a range of state-of-the-art reasoning VLMs on Visual-TableQA and compare their performance across three other benchmarks focused on table and chart-based visual question answering: ChartQA Masry et al., ReachQA He et al., and MATH-Vision Wang et al. (2024b). Our model selection includes powerful proprietary models such as GPT-4o, GPT-4o Mini OpenAI (2025b), Gemini 2.5 Flash, Gemini 2.5 Pro, and Claude 3.5 Sonnet, as well as open-source models like LLaMA 4 Maverick 17B-128E Instruct, Mistral Small 3.1 24B Instruct Mistral AI (2025), Qwen2.5-VL-32B-Instruct Chen et al. (2024a), Qwen2.5-VL-7B-Instruct Team (2025), LLaVA-Next-Llama3-8B Li et al. (2024), MiniCPM-V2.5-Llama3 Yao et al. (2024), and InternVL2-8B Chen et al. (2024b). Where performance metrics were available, we did not re-evaluate models on these datasets; instead, we report the results published in the original papers, official leaderboards, or model cards. For all other cases, we carefully fine-tuned and evaluated the models following the instructions provided in their respective official GitHub repositories.

Evaluation Protocol: All models are evaluated on the test sets of the four selected datasets. Each model receives image-question pairs, formatted within a unified prompt that includes a system message tailored to elicit the model’s reasoning capabilities (Section G.1). For the Visual-TableQA dataset, we additionally construct a variant in which data is provided not as rendered images but in LaTeX code format. This textual-code version is referred to as Visual-TableQA-CIT.

For LLaVA-Next-Llama3-8B, MiniCPM-V2.5-Llama3, InternVL2-8B, and Qwen2.5-VL-7B-Instruct, we conducted two supervised fine-tuning (SFT) experiments: (i) using the ReachQA training split (denoted as `ModelName + ReachQA`) and (ii) using the Visual-TableQA training split (denoted as `ModelName + Visual-TableQA`). We applied Low-Rank Adapters (LoRA) Hu et al. to all linear layers, following the SFT setup and hyperparameters described in the He et al. GitHub repository when possible (Section H) in order to make a fair comparison. The fine-tuning phase for all models was limited to one epoch to ensure consistency and reduce overfitting. Exceptionally, we adopted a custom two-phase LoRA fine-tuning strategy for Qwen2.5-VL-7B-Instruct (see Section H), as this model was not included in the evaluation of He et al., and to better accommodate the relatively small size of our dataset.

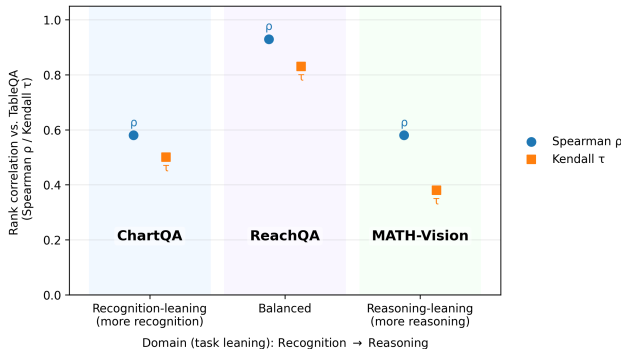
All models are allocated a maximum of 5,000 tokens during inference to accommodate extended chain-of-thought reasoning. Model responses are evaluated using the same jury of high-performing VLMs and majority-vote protocol as described in Section 2.3. The jury confidence score, computed as the ratio of the highest vote count to the total jury size, averages above 0.87 for all models and all datasets. In addition, evaluations are run twice, to ensure reproducibility.

3.2 EXPERIMENTATION RESULTS

The average models accuracies are displayed in Table 2. These results reveal that:

Visual-TableQA Effectively Evaluates Visual Reasoning Capabilities: Model performances on Visual-TableQA follow similar trends to those observed on real-world, human-annotated datasets such as ChartQA and MATH-Vision, suggesting that synthetic datasets can effectively evaluate reasoning capabilities. A direct comparison between Visual-TableQA and its textual variant, Visual-TableQA-CIT, shows a notable performance gap: on average, models perform +6.26% better on

324 Figure 3: Correlation of model rankings on Visual-TableQA with those on three established
 325 datasets—ChartQA (recognition-focused), ReachQA (balanced), and MATH-Vision (reasoning-
 326 focused)—using Spearman’s ρ and Kendall’s τ metrics. Higher values indicate stronger alignment
 327 in model performance trends. Visual-TableQA shows strong correlation with ReachQA, suggest-
 328 ing it effectively balances both visual recognition and reasoning, while its weaker correlation with
 329 ChartQA and MATH-Vision highlights its unique position as a comprehensive visual reasoning
 330 benchmark.



345 Visual-TableQA-CIT. This highlights the added challenge posed by the image-based format in
 346 Visual-TableQA, demonstrating its effectiveness at testing visual reasoning over purely textual input.

347 To further validate Visual-TableQA as a reasoning benchmark, we compared model rankings across
 348 datasets. For each dataset, we extracted the models (except the fine-tuned ones) performance rank-
 349 ings and compared them to the rankings on Visual-TableQA using two correlation measures: (i)
 350 Spearman’s ρ Lee Rodgers & Nicewander (1988): Captures monotonic consistency in rankings (re-
 351 gardless of exact scores); (ii) Kendall’s τ Kendall (1948): Measures the fraction of concordant vs.
 352 discordant ranking pairs and is more robust to ties. Both metrics range from -1 to 1 , with values
 353 closer to 1 indicating strong alignment in model rankings. To ensure fairness, we adjusted all scores
 354 computed with Relaxed Accuracy by subtracting 5% , before comparison. The results are shown in
 355 Figure 3.

356 Each dataset varies in how much it emphasizes visual recognition versus reasoning: (i) ChartQA \rightarrow
 357 Recognition-heavy, (ii) ReachQA \rightarrow Balanced, (iii) MATH-Vision \rightarrow Reasoning-heavy

358 Interestingly, Visual-TableQA rankings align most closely with ReachQA, but not with ChartQA
 359 or MATH-Vision individually. This suggests that Visual-TableQA does not favor models that excel
 360 solely at recognition or solely at reasoning. Instead, it rewards models capable of both—making it a
 361 comprehensive benchmark for evaluating all aspects of visual reasoning.

362 **Visual-TableQA Effectively Transfers to Other Benchmarks:** To assess the transferabil-
 363 ity of Visual-TableQA, we investigated how fine-tuning on Visual-TableQA impacts perfor-
 364 mance across other benchmarks. As shown in Table 2, supervision from Visual-TableQA led
 365 to significant generalization beyond its native domain. Notably, it improved the accuracy of
 366 Qwen2.5-VL-7B-Instruct on *ReachQA* from 49.23% to 60.95% , and on *MATH-Vision* from
 367 25.10% to 49.77% , despite these datasets not being explicitly table-focused. This finding is further
 368 supported by Table 4, which reports similar gains in generalization across three additional models:
 369 LLaVA-Next-Llama3-8B, MiniCPM-V2.5-Llama3, and InternVL2-8B.

370 However, this transferability is not reciprocal. Fine-tuning Qwen2.5-VL-7B-Instruct on
 371 *ReachQA* alone yields only modest in-domain gains ($49.23\% \rightarrow 55.75\%$) and leads to reduced
 372 performance on both *ChartQA* and *Visual-TableQA*. This suggests that Visual-TableQA provides a
 373 more generalizable reasoning signal—rooted in layout understanding, symbolic interpretation, and
 374 multi-step reasoning—compared to standard benchmarks.

375 **Proprietary Models Outperform Open-Source Models on Average:** Claude 3.5 Sonnet
 376 achieves the highest performance across nearly all benchmarks. However, fine-tuning on Visual-
 377 TableQA substantially narrows the gap between proprietary and open-source models. No-

ably, the performance of Qwen2.5-VL-7B-Instruct increases significantly across all evaluated benchmarks—surpassing several state-of-the-art proprietary models, including GPT-4o, GPT-4o-mini, and Gemini 2.5 Pro.

4 DISCUSSION

4.1 VISUAL-TABLEQA VS REACHQA

Table 4: Performance of fine-tuned models on the two splits of the ReachQA test set: *Recognition* (Reco) and *Reasoning* (Reas), each consisting of exactly 1,000 samples. Best performances per model category are in **bold**. The values in **blue** are from our own evaluation using the LLM jury, while the remaining values are taken from He et al..

Model	Reco	Reas	Model	Reco	Reas
LLaVA-Next-Llama3-8B	17.9	6.5	InternVL2-8B	33.7	16.2
+ ReachQA	29.6	11.1	+ ReachQA	49.8	21.3
+ Visual-TableQA	28.4	20.2	+ Visual-TableQA	45.6	34.5
MiniCPM-V2.5-Llama3	25.3	10.3	Qwen2.5-VL-7B-Instruct	66.20	33.10
+ ReachQA	35.10	11	+ ReachQA	69.6	40.30
+ Visual-TableQA	36.20	31.50	+ Visual-TableQA	70.3	50.6
Average gains					
+ ReachQA	+10.25	+4.4			
+ Visual-TableQA	+9.35	+17.68			

The ReachQA dataset is divided into two equally sized subsets: *Recognition*, which tests a model’s ability to extract relevant information from charts, and *Reasoning*, which evaluates a model’s capacity to understand complex and abstract data structures. Table 4 reports the performance gains of multiple fine-tuned models on these two tasks.

On average, models fine-tuned on ReachQA exhibit an accuracy improvement of +10.25 points on the *Recognition* task and +4.4 points on the *Reasoning* task. In comparison, models fine-tuned on Visual-TableQA show an average gain of +9.35 on *Recognition*—a comparable result—but a significantly larger gain of +17.68 on *Reasoning*.

This stark contrast in reasoning performance can be attributed to the presence of high-quality rationales in Visual-TableQA annotations, along with the inclusion of more complex and diverse visual structures. In other words, despite being roughly three times smaller than ReachQA in terms of sample count, Visual-TableQA places a stronger emphasis on qualitative richness over quantity. As a result, it appears to enable more effective knowledge distillation, particularly for tasks requiring symbolic interpretation and multi-step reasoning.

4.2 VISUAL-TABLEQA’S ADVANTAGES COMPARED TO OTHER DATASETS

Table 1 shows that only a few table-focused QA datasets—namely TAT-DQA, Table-VQA, and TableVQA-Bench—represent tables as rendered images. **Visual-TableQA** surpasses these by offering richer layout diversity, broader topic coverage, systematic visual complexity, and high-quality rationales. These attributes make it particularly effective for training models with transferable reasoning skills. Supporting this, models fine-tuned solely on **Visual-TableQA**—such as LLaVA-Next-Llama3-8B—demonstrated significant gains on external benchmarks (TableVQA and TableVQA-Bench), as seen in Table 5.

Interestingly, Qwen2.5-VL-7B-Instruct did not follow the same performance trend: it showed degradation on tasks such as *VTabFact* (Yes/No fact verification), *VWTQ* (Wikipedia table retrieval), and *VWTQ-Syn* (synthetic variants). To understand this, we manually analyzed its errors before and after fine-tuning on *VTabFact*, categorizing them into eight types: *partial data extraction*, *hallucination*, *incoherence*, *misunderstanding*, *reasoning errors*, *evaluation mistakes*, *dataset ambiguity*, and *annotation flaws*. Results (Figure 14) show that while the total number of errors slightly increased post-finetuning, most now fall into the *incoherence* class, with all other error types significantly reduced. This suggests a sharpening of reasoning patterns but also highlights a need for

future work targeting specific error types through synthetic supervision. Further details are provided in Section I.

Beyond transferability and diversity, a key advantage of Visual-TableQA lies in its modularity and scalability as explained in Section 4.3 .

Table 5: Performance of fine-tuned models on Table-VQA test set and the four splits of the TableVQA-Bench dataset: *FinTabNetQA* (finance-related tables), *VTabFact* (table-based fact verification with Yes/No questions), *VWTQ* (information retrieval from Wikipedia tables), and *VWTQ-Syn* (synthetic visual variants of VWTQ). Best performances per model variants are shown in **bold**. Values in blue are from our own evaluation using the DeepSeek-Prover-v2, while remaining values are reported from Fu et al. (2025).

Model	TableVQA-Bench				Table-VQA
	FinTabNetQA	VTabFact	VWTQ	VWTQ-Syn	
GPT-4o	96.8	78.0	72.8	82.4	–
LLaVA-Next-34B	–	71.2	36.4	38.0	–
LLaVA-Next-Llama3-8B	52.4	37.2	21.5	24.8	25.84
+ Visual-TableQA	56.8	52.0	33.2	33.6	28.89
Qwen2.5-VL-7B-Instruct	96.4	82.0	68.53	74.0	79.03
+ Visual-TableQA	97.2	70.6	61.5	69.6	75.23

4.3 SCALABILITY OF THE PIPELINE AND ITS BENEFITS FOR KNOWLEDGE DISTILLATION

This modular pipeline supports scalable generation with a clean separation of concerns—table structure synthesis, QA creation, and validation—making each component independently reusable and upgradable. By automating the entire process from table generation to jury-based quality control, Visual-TableQA provides a cost-efficient and high-quality benchmark for advancing multimodal reasoning over complex visual inputs. A central component of our pipeline is the mechanism of **cross-model inspiration** 2.2, a collaborative prompting strategy. In this process, stronger models generate layout “seeds” that guide weaker models in synthesizing structurally diverse tables, fostering novel visual configurations through iterative transfer. The same principle extends to question-answer generation: models are prompted with both layout and topical cues—often proposed by stronger models—to create new QA pairs. This enables weaker models to contribute meaningfully to the dataset by expanding the range of questions and reasoning patterns. Through this dual-inspiration process, the pipeline cultivates a collaborative multi-model co-creation space, where models of varying capabilities distill collective knowledge not through imitation, but through generative inspiration, while maintaining data quality. In this regard, **Visual-TableQA** distinguishes itself from other synthetic datasets Aboutaleb et al. (2024); Wang et al. (2024a); Li et al. (2025); He et al..

5 CONCLUSION

In this work, we introduced Visual-TableQA, a large-scale, open-domain, multimodal dataset designed to rigorously evaluate visual reasoning capabilities over complex table images. Building on the principles of Code-as-Intermediary Translation (CIT), we developed a fully automated, modular pipeline for generating LaTeX-rendered tables, reasoning-intensive question-answer pairs, and high-quality rationales—all verified by a jury of strong LLMs. Despite being cost-efficient (generated for under \$100), Visual-TableQA offers unprecedented diversity in table structures, visual features, and reasoning depth. We showed that Visual-TableQA not only challenges existing visual language models (VLMs) but also serves as an effective training signal for improving reasoning performance. Fine-tuning on Visual-TableQA led to substantial gains across multiple benchmarks—both table-centric and general-purpose—including ReachQA and MATH-Vision, demonstrating the dataset’s capacity to bridge the performance gap between open-source and proprietary models.

REFERENCES

- 486
487
488 Hossein Aboutaleb, Hwanjun Song, Yusheng Xie, Arshit Gupta, Justin Sun, Hang Su, Igor Sha-
489 lyminov, Nikolaos Pappas, Siffi Singh, and Saab Mansour. Magid: An automated pipeline for
490 generating synthetic multi-modal datasets. *arXiv preprint arXiv:2403.03194*, 2024.
- 491 Pranav Agarwal and Ioana Ciucă. Supernova event dataset: Interpreting large language model’s
492 personality through critical event analysis. *arXiv preprint arXiv:2506.12189*, 2025.
- 493 Anthropic. Model card addendum: Claude 3.5 haiku and upgraded claude 3.5 son-
494 net. [https://assets.anthropic.com/m/1cd9d098ac3e6467/original/
495 Claude-3-Model-Card-October-Addendum.pdf](https://assets.anthropic.com/m/1cd9d098ac3e6467/original/Claude-3-Model-Card-October-Addendum.pdf), 2024. Accessed: 2025-08-01.
- 496 Anthropic. Claude opus 4 & claude sonnet 4 — system card. [https://www.anthropic.com/
497 claude-4-system-card](https://www.anthropic.com/claude-4-system-card), May 2025. Accessed: 2025-08-01.
- 498 Shaohan Chen, Yujia Zhang, Xiangpeng Cao, Shaolei He, Chen Zhao, Zhihua Liu, Chongming Li,
499 Jing Liu, Qiang Liu, Fan Liu, et al. Qwen-vl: A versatile vision-language model with image, text,
500 and box comprehension. *arXiv preprint arXiv:2403.18751*, 2024a. URL [https://arxiv.
501 org/abs/2403.18751](https://arxiv.org/abs/2403.18751).
- 502
503 Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou,
504 and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. In
505 *International Conference on Learning Representations*, 2020a. URL [https://openreview.
506 net/forum?id=rkeJRhNYDH](https://openreview.net/forum?id=rkeJRhNYDH).
- 507
508 Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang.
509 HybridQA: A dataset of multi-hop question answering over tabular and textual data. In Trevor
510 Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics:
511 EMNLP 2020*, pp. 1026–1036, Online, November 2020b. Association for Computational Lin-
512 guistics. doi: 10.18653/v1/2020.findings-emnlp.91. URL [https://aclanthology.org/
513 2020.findings-emnlp.91/](https://aclanthology.org/2020.findings-emnlp.91/).
- 514 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
515 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
516 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer
517 Vision and Pattern Recognition*, pp. 24185–24198, 2024b.
- 518 DeepSeek-AI. DeepSeek-R1-Distill-Qwen-32B. [https://huggingface.co/
519 deepseek-ai/DeepSeek-R1-Distill-Qwen-32B](https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B), 2025. Model card. Accessed:
520 2025-08-01.
- 521 Xingyu Fu, Minqian Liu, Zhengyuan Yang, John Corring, Yijuan Lu, Jianwei Yang, Dan Roth, Dinei
522 Florencio, and Cha Zhang. Refocus: Visual editing as a chain of thought for structured image
523 understanding. *arXiv preprint arXiv:2501.05452*, 2025.
- 524
525 Olga Golovneva, Moya Peng Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam
526 Fazel-Zarandi, and Asli Celikyilmaz. Roscoe: A suite of metrics for scoring step-by-step reason-
527 ing. In *The Eleventh International Conference on Learning Representations*.
- 528 Google. Gemini 2.0 flash: Model card. [https://storage.googleapis.com/
529 model-cards/documents/gemini-2-flash.pdf](https://storage.googleapis.com/model-cards/documents/gemini-2-flash.pdf), 2025a. Published: 2025-04-15. Ac-
530 cessed: 2025-08-01.
- 531
532 Google. Gemini 2.5 flash: Model card. [https://storage.googleapis.com/
533 model-cards/documents/gemini-2.5-flash.pdf](https://storage.googleapis.com/model-cards/documents/gemini-2.5-flash.pdf), 2025b. Updated: 2025-06-26.
534 Accessed: 2025-08-01.
- 535 Google. Gemini 2.5 pro: Model card. [https://storage.googleapis.com/
536 model-cards/documents/gemini-2.5-pro.pdf](https://storage.googleapis.com/model-cards/documents/gemini-2.5-pro.pdf), 2025c. Model card. Last updated:
537 2025-06-27. Accessed: 2025-08-01.
- 538
539 Wei He, Zhiheng Xi, Wanxu Zhao, Xiaoran Fan, Yiwen Ding, Zifei Shan, Tao Gui, Qi Zhang, and
Xuanjing Huang. Distill visual chart reasoning ability from llms to mllms.

- 540 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
541 et al. Lora: Low-rank adaptation of large language models. In *International Conference on*
542 *Learning Representations*.
- 543 Sahil Kale and Vijaykant Nadadur. Texpert: A multi-level benchmark for evaluating latex code
544 generation by llms. *arXiv preprint arXiv:2506.16990*, 2025.
- 545 Yannis Katsis, Saneem Chemmengath, Vishwajeet Kumar, Samarth Bharadwaj, Mustafa Canim,
546 Michael Glass, Alfio Gliozzo, Feifei Pan, Jaydeep Sen, Karthik Sankaranarayanan, and Soumen
547 Chakrabarti. AIT-QA: Question answering dataset over complex tables in the airline industry. In
548 Anastasia Loukina, Rashmi Gangadharaiah, and Bonan Min (eds.), *Proceedings of the 2022 Con-*
549 *ference of the North American Chapter of the Association for Computational Linguistics: Human*
550 *Language Technologies: Industry Track*, pp. 305–314, Hybrid: Seattle, Washington + Online,
551 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-industry.34.
552 URL <https://aclanthology.org/2022.naacl-industry.34/>.
- 553 Maurice George Kendall. Rank correlation methods. 1948.
- 554 Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. Tablevqa-bench: A visual question answering
555 benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*, 2024.
- 556 Joseph Lee Rodgers and W Alan Nicewander. Thirteen ways to look at the correlation coefficient.
557 *The American Statistician*, 42(1):59–66, 1988.
- 558 Andrew Li, Rahul Thapa, Rahul Chalamala, Qingyang Wu, Kezhen Chen, and James Zou.
559 Smir: Efficient synthetic data pipeline to improve multi-image reasoning. *arXiv preprint*
560 *arXiv:2501.03675*, 2025.
- 561 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
562 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv*
563 *preprint arXiv:2407.07895*, 2024.
- 564 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
565 mark for question answering about charts with visual and logical reasoning.
- 566 Minesh Mathew, Dimosthenis Karatzas, R Manmatha, and CV Jawahar. Docvqa: A dataset for vqa
567 on document images. *corr abs/2007.00398 (2020)*. *arXiv preprint arXiv:2007.00398*, 2020.
- 568 Meta AI. Llama 4 Maverick 17B-128E Instruct. [https://huggingface.co/meta-llama/](https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E-Instruct)
569 [Llama-4-Maverick-17B-128E-Instruct](https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E-Instruct), 2025. Model card. Accessed: 2025-08-01.
- 570 Mistral AI. Mistral Small 3.1 24B Instruct. [https://huggingface.co/mistralai/](https://huggingface.co/mistralai/MistralSmall3.124BInstruct2503)
571 [MistralSmall3.124BInstruct2503](https://huggingface.co/mistralai/MistralSmall3.124BInstruct2503), 2025. Model card. Accessed: 2025-08-01.
- 572 OpenAI. GPT-4o. <https://openai.com/index/gpt-4o>, 2024. Accessed: 2025-07-30.
- 573 OpenAI. GPT-4.1. <https://openai.com/index/gpt-4-1/>, 2025a. Accessed: 2025-08-
574 01.
- 575 OpenAI. GPT-4o-mini. <https://platform.openai.com/docs/models/gpt-4o>,
576 2025b. Accessed: 2025-08-01.
- 577 OpenAI. OpenAI o3 Reasoning Model. [https://openai.com/index/](https://openai.com/index/introducing-o3-and-o4-mini/)
578 [introducing-o3-and-o4-mini/](https://openai.com/index/introducing-o3-and-o4-mini/), 2025. Accessed: 2025-07-31.
- 579 Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables.
580 In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*
581 *and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long*
582 *Papers)*, pp. 1470–1480, Beijing, China, July 2015. Association for Computational Linguistics.
583 doi: 10.3115/v1/P15-1142. URL <https://aclanthology.org/P15-1142>.
- 584 Qwen Team. Qwen3-30B-A3B. <https://huggingface.co/Qwen/Qwen3-30B-A3B>,
585 2025a. Model card. Accessed: 2025-08-01.

- 594 Qwen Team. Qwen3-QwQ-32B. <https://huggingface.co/Qwen/QwQ-32B>, 2025b.
595 Model card. Accessed: 2025-08-01.
596
- 597 Reka AI. Reka Flash 3. <https://huggingface.co/RekaAI/reka-flash-3>, 2025.
598 Model card. Accessed: 2025-08-01.
- 599 Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco,
600 Hannaneh Hajishirzi, and Jonathan Berant. Multimodal{qa}: complex question answering over
601 text, tables and images. In *International Conference on Learning Representations*, 2021. URL
602 <https://openreview.net/forum?id=ee6W5UgQLa>.
603
- 604 Qwen Team. Qwen2.5-vl, January 2025. URL [https://qwenlm.github.io/blog/
605 qwen2.5-vl/](https://qwenlm.github.io/blog/qwen2.5-vl/).
- 606 TNG Technology Consulting GmbH. Deepseek-rlt-chimera, April 2025. URL [https://
607 huggingface.co/tngtech/DeepSeek-RLT-Chimera](https://huggingface.co/tngtech/DeepSeek-RLT-Chimera).
- 608 Cyrille Delestre Tom Agonnoude, 2024. URL [https://huggingface.co/datasets/
609 cmarkea/table-vqa](https://huggingface.co/datasets/cmarkea/table-vqa).
- 611 Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady
612 Arkhangorodsky, Minjie Xu, Naomi White, and Patrick Lewis. Replacing judges with juries:
613 Evaluating llm generations with a panel of diverse models. *arXiv preprint arXiv:2404.18796*,
614 2024.
- 615 Jiankang Wang, Jianjun Xu, Xiaorui Wang, Yuxin Wang, Mengting Xing, Shancheng Fang, Zhineng
616 Chen, Hongtao Xie, and Yongdong Zhang. A graph-based synthetic data pipeline for scaling high-
617 quality reasoning instructions. *CoRR*, 2024a.
618
- 619 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and
620 Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. In *The
621 Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks
622 Track*, 2024b. URL <https://openreview.net/forum?id=QWTCcxMpPA>.
- 623 xAI. Grok 3 beta — the age of reasoning agents. <https://x.ai/news/grok-3>, 2025. Ac-
624 cessed: 2025-08-01.
625
- 626 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,
627 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint
628 arXiv:2408.01800*, 2024.
- 629 Xinyi Zheng, Douglas Burdick, Lucian Popa, Xu Zhong, and Nancy Xin Ru Wang. Global table
630 extractor (gte): A framework for joint table identification and cell structure recognition using
631 visual context. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*,
632 pp. 697–706. IEEE, 2021.
- 633 Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from
634 natural language using reinforcement learning. *CoRR*, abs/1709.00103, 2017.
635
- 636 Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang, Haozhou Zhang, and Tat-Seng Chua. To-
637 wards complex document understanding by discrete reasoning. In *Proceedings of the 30th ACM
638 International Conference on Multimedia*, pp. 4857–4866, 2022.
639
640
641
642
643
644
645
646
647