

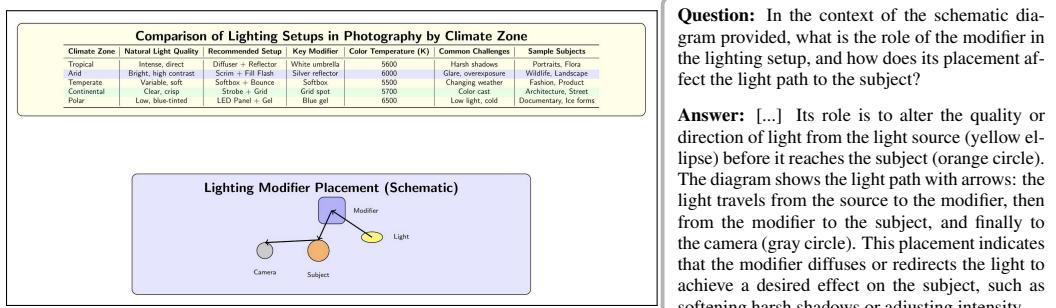
MODALITY-SWAP DISTILLATION: RENDERING TEXTUAL REASONING INTO VISUAL SUPERVISION

000
001
002
003
004
005 **Anonymous authors**
006 Paper under double-blind review
007
008
009
010

ABSTRACT

011 Inspired by previous works (He et al.), we propose **modality swap**, a training
012 strategy in which language models use their textual reasoning and coding capabili-
013 ties to generate **synthetic LaTeX tables** that are rendered into images and then
014 used for visual reasoning training. This setup allows smaller VLMs to contribute
015 to their own improvement by transferring competence from text (structured gener-
016 ation and code) to vision (reasoning over rendered tables). We instantiate this idea
017 with **Visual-TableQA**, built through a **modular, scalable, and fully autonomous**
018 pipeline. Multiple LLMs coordinate across roles—generation, validation, and
019 inspiration (cross-model prompting)—to produce **2.5k** richly structured LaTeX-
020 rendered tables and **9k** reasoning-intensive QA pairs at a cost under **\$100**. The
021 pipeline includes LLM-jury filtering and cross-model inspiration, where stronger
022 models propose structural seeds and topics that other models elaborate.
023
024 In experiments, models fine-tuned on **Visual-TableQA** show robust generalization
025 to external benchmarks; in some settings, results are competitive with or exceed
026 proprietary baselines. An ablation also indicates higher scores when queries are
027 presented in textual rather than visual form, consistent with the intended cross-
028 modal transfer from text to vision.

1 INTRODUCTION



043 Figure 1: Sample question in our benchmark
044

045 Vision-language models (VLMs) have significantly advanced in recent years, achieving remark-
046 able performance in various tasks involving visual and textual inputs. Despite these advancements,
047 complex reasoning tasks, especially those requiring deep comprehension of tabular data structures,
048 continue to pose significant challenges. Table complexity can manifest in various ways, including
049 structural layout, information density, and the diversity of visual components such as the integration
050 of diagrams. The more complex a table is, the more it lends itself to challenging reasoning tasks,
051 requiring advanced cognitive abilities to extract relevant information and perform multi-step logical
052 analysis. For example, the table in Figure 1 exemplifies this complexity through its use of multi-
053 row cells, integrated diagrams, and color encoding. Answering the question requires the VLM to
interpret information across all cells and perform a sequence of reasoning steps.

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

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

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

2 VISUAL-TABLEQA DATASET

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

2.1 MODALITY-SWAP: TABLE REPRESENTATION IN LATEX

Our approach is inspired by He et al., which demonstrated that state-of-the-art VLMs can reason about visual content even in the absence of explicit visual input. Building on this insight, and

108
 109
 110
 111
 112
 113
 114 Table 1: Comparison of existing chart and table datasets across data, Q&A, and dataset properties.
 115 Abbreviations: Repr=Representation, Vis= Visual, Comp= Complexity, Temp = Template, Refer =
 116 Reference, Rat = Rational, Synth= Synthetic, Scal = Scalable. Cells marked with \blacktriangle indicate mixed
 117 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 Agonnoude, 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

130
 131 Table 2: Model performances on Visual-TableQA and Visual-TableQA-CIT. Visual-TableQA-CIT
 132 is the variant of our dataset where tables are represented in LaTeX code form rather than as rendered
 133 images.

Models	GPT-4o	GPT-4o mini	Gemini 2.5 Flash	Gemini 2.5 Pro	Claude 3.5 Son- net	Llama 4 Maverick 17B-128E	Mistral Small 24B	Qwen2.5- VL-32B
Visual-TableQA	81.0	67.0	86.63	85.63	82.46	80.75	73.2	79.69
Visual- TableQA-CIT	90.25	86.69	89.9	88.71	88.5	87.0	84.72	85.44

140 leveraging the strong coding capabilities of reasoning-oriented language models across multiple
 141 programming languages, we chose to use an intermediate representation of tables in LaTeX rather
 142 than directly generating rendered table images. This strategy enables the generation of complex
 143 visual tables as compact LaTeX code—typically around 100 lines per table—drastically reducing
 144 the cost of generation by minimizing the number of output tokens required in API calls. We refer
 145 to this as **modality-swap**: LLMs leverage their textual reasoning abilities (e.g., code understanding
 146 and structured generation) to construct visual content, design to improve visual reasoning skills. The
 147 ablation study in Table 2 supports this hypothesis: most state-of-the-art models perform better when
 148 queries are presented in textual rather than visual form. This highlights both a gap in current visual
 149 reasoning capabilities and the potential for cross-modal transfer from text to vision.

150 Table 3 presents the performance of various models in generating LaTeX tables that compile without
 151 errors. Our observations align with those of Kale & Nadadur (2025), who reported that LLMs
 152 struggle with LaTeX generation—particularly as task complexity increases, leading to a notable
 153 drop in accuracy.

154 2.2 DATA GENERATION PIPELINE

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

157 **Seed Tables and Topics Collection:** The first step involves collecting a diverse set of table layouts to
 158 serve as inspiration for LLMs during the generation process. We explored various sources, including
 159 scientific journals, financial report databases, online newspapers, and table design galleries. Our
 160 search included both table and diagram images to introduce greater visual and structural complexity

Table 3: Percentage of successful LaTeX compilations for various models. Each accuracy is computed from at least 500 generated samples. The Adjust column indicates the level of manual correction needed to make the table look good: Low means minimal or no adjustments, Medium corresponds to 3–5 required fixes, and High indicates more than 5 adjustments were necessary. The tables generated by DeepSeek-R1-Distill-Qwen-32B never compiled.

Model	Acc. (%)	Adjust	Model	Acc. (%)	Adjust
	(%)			(%)	
Llama 4 Maverick 17B-128E Instruct Meta AI (2025)	69	High	DeepSeek-R1-Distill-Qwen-32B DeepSeek-AI (2025)	0.0	–
Gemini 2.0 Flash Google (2025a)	65.7	Low	DeepSeek-RIT-Chimera TNG Technology Consulting GmbH (2025)	43.4	Medium
Gemini 2.5 Flash Google (2025b)	43	Medium	Claude Sonnet 4 Anthropic (2025)	56	Low
Gemini 2.5 Pro Google (2025c)	19.6	Low	Claude 3.5 Haiku Anthropic (2024)	64.4	Low
GPT-4.1 OpenAI (2025a)	41.5	Low	Grok 3 Beta xAI (2025)	47.3	Low
Qwen3-30B-A3B Qwen Team (2025a)	69.4	Low	Reka Flash 3 Reka AI (2025)	19.3	Medium
Qwen-QwQ-32B Qwen Team (2025b)	38.2	Low			

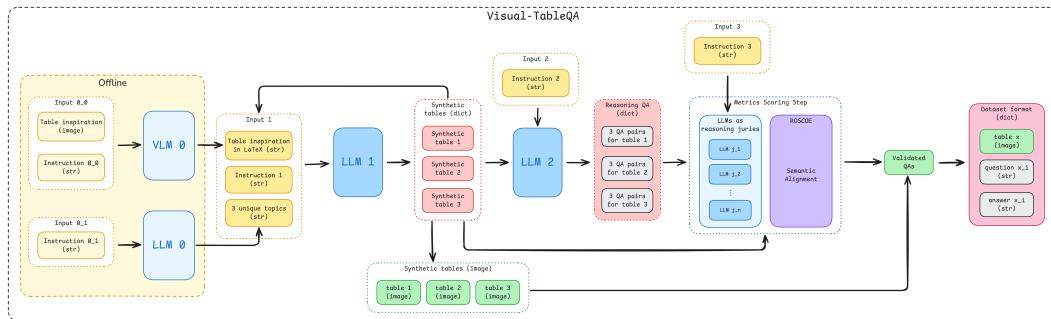


Figure 2: Overview of the full pipeline architecture of Visual-TableQA. A subset of initial table images is first converted to LaTeX using a visual language model (VLM-0). The resulting LaTeX code, along with topic prompts, is then passed to a language model (LLM-1) to generate new, diverse tables. These newly generated tables are then submitted to a second language model (LLM-2), which produces corresponding question-answer pairs. Finally, the QA pairs are evaluated by a jury of high-performing LLMs, and their quality is assessed using the ROSCOE score.

into the dataset. We selected 20 representative images (Figure 6a) and passed them to a visual language model, VLM-0 (GPT- \circ 3 OpenAI (2025)), to generate accurate LaTeX representations. In parallel, we used LLM-0 (GPT- \circ 4 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 DeepSeek-AI et al. (2025), 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 8 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 Model_Name + ReachQA) and (ii) using the Visual-TableQA training split (denoted as Model_Name + 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 0.93 (Figure 11b) for all models and all datasets. In addition, evaluations are run twice, to ensure reproducibility.

3.2 EXPERIMENTATION RESULTS

Table 4: 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.

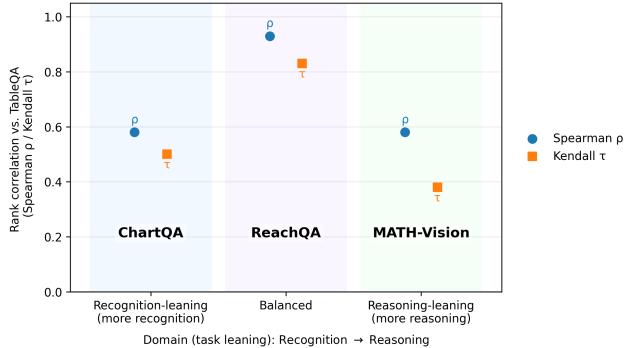
Models	ChartQA	ReachQA	MATH-Vision _{FULL}	Visual-TableQA	Visual-TableQA-CIT
Baseline					
Human	–	74.85	68.82	–	–
Proprietary VLMs					
GPT-4o	85.7	53.25	30.39	81.0	90.25
GPT-4o mini	77.52	40.35	28.85	67.0	86.69
Gemini 2.5 Flash	84.64	56.97	41.3	86.63	89.9
Gemini 2.5 Pro	87.67	60.88	73.3	85.63	88.71
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	84.72
Qwen2.5-VL-32B-Instruct	79.75	49.5	38.1	79.69	85.44
Qwen2.5-VL-7B-Instruct	87.3*	44.7	25.1	69.5	–
Finetuned VLMs					
Qwen2.5-VL-7B-Instruct + Visual-TableQA	84.52 ↑	60.95 ↑	49.77 ↑	82.98 ↑	N/A
Qwen2.5-VL-7B-Instruct + ReachQA	77.59 ↓	55.75 ↑	48.57 ↑	56.13 ↓	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.

The average models accuracies are displayed in Table 4. These results reveal that:

Visual-TableQA Effectively Evaluates Visual Reasoning Capabilities: Model performances on

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



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

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

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

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

Visual-TableQA Effectively Transfers to Other Benchmarks: To assess the transferability of Visual-TableQA, we investigated how fine-tuning on Visual-TableQA impacts performance across other benchmarks. As shown in Table 4, supervision from Visual-TableQA led to significant generalization beyond its native domain. Notably, it improved the accuracy of Qwen2.5-VL-7B-Instruct on *ReachQA* from 44.7% to 60.95%, and on *MATH-Vision* from 25.10% to 49.77%, despite these datasets not being explicitly table-focused. This finding is further supported by Table 5, which reports similar gains in generalization across three additional models: LLaVA-Next-Llama3-8B, MiniCPM-V2.5-Llama3, and InternVL2-8B.

However, this transferability is not reciprocal. Fine-tuning Qwen2.5-VL-7B-Instruct on *ReachQA* alone yields only modest in-domain gains (44.7% → 55.75%) and leads to reduced performance on both *ChartQA* and *Visual-TableQA*. This suggests that Visual-TableQA provides a more generalizable reasoning signal—rooted in layout understanding, symbolic interpretation, and multi-step reasoning—compared to standard benchmarks.

378 **Proprietary Models Outperform Open-Source Models on Average:** Claude 3.5 Sonnet
 379 achieves the highest performance across nearly all benchmarks. However, fine-tuning on Visual-
 380 TableQA substantially narrows the gap between proprietary and open-source models. No-
 381 tably, the performance of Qwen2.5-VL-7B-Instruct increases significantly across all eval-
 382 uated benchmarks—surpassing several state-of-the-art proprietary models, including GPT-4o,
 383 GPT-4o-mini, and Gemini 2.5 Pro.

384 4 DISCUSSION

385 4.1 VISUAL-TABLEQA VS REACHQA

386 Table 5: Performance of fine-tuned models on the two splits of the ReachQA test set: *Recognition*
 387 (Reco) and *Reasoning* (Reas), each consisting of exactly 1,000 samples. Best performances per
 388 model category are in **bold**. The values in **blue** are from our own evaluation using the LLM jury,
 389 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	61.70	30.10
+ ReachQA	35.10	11	+ ReachQA	69.6	40.30
+ Visual-TableQA	36.20	31.50	+ Visual-TableQA	70.3	50.6
Average gains					
+ ReachQA	+11.25	+5.4			
+ Visual-TableQA	+10.35	+18.68			

404 The ReachQA dataset is divided into two equally sized subsets: *Recognition*, which tests a model’s
 405 ability to extract relevant information from charts, and *Reasoning*, which evaluates a model’s ca-
 406 pacity to understand complex and abstract data structures. Table 5 reports the performance gains of
 407 multiple fine-tuned models on these two tasks.

408 On average, models fine-tuned on ReachQA exhibit an accuracy improvement of +11.25 points
 409 on the *Recognition* task and +5.4 points on the *Reasoning* task. In comparison, models fine-tuned
 410 on Visual-TableQA show an average gain of +10.35 on *Recognition*—a comparable result—but a
 411 significantly larger gain of +18.68 on *Reasoning*.

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

419 4.2 VISUAL-TABLEQA’S ADVANTAGES COMPARED TO OTHER DATASETS

421 Table 1 shows that only a few table-focused QA datasets—namely TAT-DQA, Table-VQA, and
 422 TableVQA-Bench—represent tables as rendered images. **Visual-TableQA** surpasses these by of-
 423 fering richer layout diversity, broader topic coverage, systematic visual complexity, and high-
 424 quality rationales. These attributes make it particularly effective for training models with trans-
 425 ferable reasoning skills. Supporting this, models fine-tuned solely on **Visual-TableQA**—such
 426 as LLaVA-Next-Llama3-8B—demonstrated significant gains on external benchmarks (Table-
 427 VQA and TableVQA-Bench), as seen in Table 6.

428 Interestingly, Qwen2.5-VL-7B-Instruct did not follow the same performance trend: it
 429 showed degradation on tasks such as *VTabFact* (Yes/No fact verification), *VWTQ* (Wikipedia table
 430 retrieval), and *VWTQ-Syn* (synthetic variants). To understand this, we manually analyzed its errors
 431 before and after fine-tuning on *VTabFact*, categorizing them into eight types: *partial data extraction*,
hallucination, *incoherence*, *misunderstanding*, *reasoning errors*, *evaluation mistakes*, *dataset ambi-*

432 *guity*, and *annotation flaws*. Results (Figure 13) show that while the total number of errors slightly
 433 increased post-finetuning, most now fall into the *incoherence* class, with all other error types sig-
 434 nificantly reduced. This suggests a sharpening of reasoning patterns but also highlights a need for
 435 future work targeting specific error types through synthetic supervision. Further details are provided
 436 in Section I.

437 Beyond transferability and diversity, a key advantage of Visual-TableQA lies in its modularity and
 438 scalability as explained in Section 4.3 .
 439

440 Table 6: Performance of fine-tuned models on Table-VQA test set and the four splits of the
 441 TableVQA-Bench dataset: *FinTabNetQA* (finance-related tables), *VTabFact* (table-based fact ver-
 442 ification with Yes/No questions), *VWTQ* (information retrieval from Wikipedia tables), and *VWTQ-Syn*
 443 (synthetic visual variants of VWTQ). Best performances per model variants are shown in **bold**.
 444 Values in blue are from our own evaluation, while remaining values are reported from Fu et al.
 445 (2025).

Model	TableVQA-Bench				Table-VQA
	FinTabNetQA	VTabFact	VWTQ	VWTQ-Syn	
GPT-4o	98.0	80.1	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	24.84
+ Visual-TableQA	56.8–	52.0–	33.2–	33.6–	28.89–
Qwen2.5-VL-7B-Instruct	96.4	81.0	68.53	73.2	79.03
+ Visual-TableQA	97.2	70.6	62.5	69.6	75.23–

456 4.3 SCALABILITY OF THE PIPELINE AND ITS BENEFITS FOR KNOWLEDGE DISTILLATION

457 This modular pipeline supports scalable generation with a clean separation of concerns—table struc-
 458 ture synthesis, QA creation, and validation—making each component independently reusable and
 459 upgradable. By automating the entire process from table generation to jury-based quality control,
 460 Visual-TableQA provides a cost-efficient and high-quality benchmark for advancing multimodal rea-
 461 soning over complex visual inputs. A central component of our pipeline is the mechanism of **cross-**
 462 **model inspiration** 2.2, a collaborative prompting strategy. In this process, stronger models generate
 463 layout “seeds” that guide weaker models in synthesizing structurally diverse tables, fostering novel
 464 visual configurations through iterative transfer. The same principle extends to question–answer gen-
 465 eration: models are prompted with both layout and topical cues—often proposed by stronger mod-
 466 els—to create new QA pairs. This enables weaker models to contribute meaningfully to the dataset
 467 by expanding the range of questions and reasoning patterns. Through this dual-inspiration process,
 468 the pipeline cultivates a collaborative multi-model co-creation space, where models of varying capa-
 469 bilities distill collective knowledge not through imitation, but through generative inspiration, while
 470 maintaining data quality. In this regard, **Visual-TableQA** distinguishes itself from other synthetic
 471 datasets Aboutalebi et al. (2024); Wang et al. (2024a); Li et al. (2025); He et al..
 472

473 5 CONCLUSION

474 In this work, we introduced Visual-TableQA, a large-scale, open-domain, multimodal dataset de-
 475 signed to rigorously evaluate visual reasoning capabilities over complex table images. Building on
 476 the principles of Code-as-Intermediary Translation (CIT), we developed a fully automated, modu-
 477 lar pipeline for generating LaTeX-rendered tables, reasoning-intensive question–answer pairs, and
 478 high-quality rationales—all verified by a jury of strong LLMs. Despite being cost-efficient (gen-
 479 erated for under \$100), Visual-TableQA offers unprecedented diversity in table structures, visual
 480 features, and reasoning depth. We showed that Visual-TableQA not only challenges existing vi-
 481 sual language models (VLMs) but also serves as an effective training signal for improving reason-
 482 ing performance. Fine-tuning on Visual-TableQA led to substantial gains across multiple bench-
 483 marks—both table-centric and general-purpose—including ReachQA and MATH-Vision, demon-
 484 strating the dataset’s capacity to bridge the performance gap between open-source and proprietary
 485 models.

486 REFERENCES
487

- 488 Hossein Aboutalebi, 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/Claude-3-Model-Card-October-Addendum.pdf>, 2024. Accessed: 2025-08-01.
- 495 Anthropic. Claude opus 4 & claude sonnet 4 — system card. <https://www.anthropic.com/clause-4-system-card>, May 2025. Accessed: 2025-08-01.
- 496 Shaohan Chen, Yujia Zhang, Xiangpeng Cao, Shaolei He, Chen Zhao, Zhihua Liu, Chongming Li,
497 Jing Liu, Qiang Liu, Fan Liu, et al. Qwen-vl: A versatile vision-language model with image, text,
498 and box comprehension. *arXiv preprint arXiv:2403.18751*, 2024a. URL <https://arxiv.org/abs/2403.18751>.
- 499 Wenhui Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou,
500 and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. In
501 *International Conference on Learning Representations*, 2020a. URL <https://openreview.net/forum?id=rkeJRhNYDH>.
- 502 Wenhui Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang.
503 HybridQA: A dataset of multi-hop question answering over tabular and textual data. In Trevor
504 Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics:
EMNLP 2020*, pp. 1026–1036, Online, November 2020b. Association for Computational Lin-
505 guistics. doi: 10.18653/v1/2020.findings-emnlp.91. URL <https://aclanthology.org/2020.findings-emnlp.91/>.
- 506 Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong
507 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning
508 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer
Vision and Pattern Recognition*, pp. 24185–24198, 2024b.
- 509 DeepSeek-AI. DeepSeek-R1-Distill-Qwen-32B. [https://huggingface.co/
510 deepseek-ai/DeepSeek-R1-Distill-Qwen-32B](https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B), 2025. Model card. Accessed:
511 2025-08-01.
- 512 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Cheng-
513 gang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang,
514 Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting
515 Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui
516 Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi
517 Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li,
518 Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang,
519 Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun
520 Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan
521 Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J.
522 Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang,
523 Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng
524 Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shut-
525 ing Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao,
526 Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue
527 Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xi-
528 aokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin
529 Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang,
530 Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang
531 Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui

- 540 Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying
 541 Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu,
 542 Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan
 543 Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F.
 544 Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda
 545 Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao,
 546 Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li,
 547 Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL
 548 <https://arxiv.org/abs/2412.19437>.
- 549 Xingyu Fu, Minqian Liu, Zhengyuan Yang, John Corring, Yijuan Lu, Jianwei Yang, Dan Roth, Dinei
 550 Florencio, and Cha Zhang. Refocus: Visual editing as a chain of thought for structured image
 551 understanding. *arXiv preprint arXiv:2501.05452*, 2025.
- 552 Olga Golovneva, Moya Peng Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam
 553 Fazel-Zarandi, and Asli Celikyilmaz. Roscoe: A suite of metrics for scoring step-by-step reason-
 554 ing. In *The Eleventh International Conference on Learning Representations*.
- 555 Google. Gemini 2.0 flash: Model card. <https://storage.googleapis.com/model-cards/documents/gemini-2-flash.pdf>, 2025a. Published: 2025-04-15. Ac-
 556 cessed: 2025-08-01.
- 557 Google. Gemini 2.5 flash: Model card. <https://storage.googleapis.com/model-cards/documents/gemini-2.5-flash.pdf>, 2025b. Updated: 2025-06-26.
 558 Accessed: 2025-08-01.
- 559 Google. Gemini 2.5 pro: Model card. <https://storage.googleapis.com/model-cards/documents/gemini-2.5-pro.pdf>, 2025c. Model card. Last updated:
 560 2025-06-27. Accessed: 2025-08-01.
- 561 Wei He, Zhiheng Xi, Wanxu Zhao, Xiaoran Fan, Yiwen Ding, Zifei Shan, Tao Gui, Qi Zhang, and
 562 Xuanjing Huang. Distill visual chart reasoning ability from llms to mllms.
- 563 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
 564 et al. Lora: Low-rank adaptation of large language models. In *International Conference on
 565 Learning Representations*.
- 566 Sahil Kale and Vijaykant Nadadur. Texpert: A multi-level benchmark for evaluating latex code
 567 generation by llms. *arXiv preprint arXiv:2506.16990*, 2025.
- 568 Yannis Katsis, Saneem Chemmengath, Vishwajeet Kumar, Samarth Bharadwaj, Mustafa Canim,
 569 Michael Glass, Alfio Gliozzo, Feifei Pan, Jaydeep Sen, Karthik Sankaranarayanan, and Soumen
 570 Chakrabarti. AIT-QA: Question answering dataset over complex tables in the airline industry. In
 571 Anastassia Loukina, Rashmi Gangadharaiah, and Bonan Min (eds.), *Proceedings of the 2022 Con-
 572 ference of the North American Chapter of the Association for Computational Linguistics: Human
 573 Language Technologies: Industry Track*, pp. 305–314, Hybrid: Seattle, Washington + Online,
 574 July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.nacl-industry.34.
 575 URL <https://aclanthology.org/2022.nacl-industry.34/>.
- 576 Maurice George Kendall. Rank correlation methods. 1948.
- 577 Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. Tablevqa-bench: A visual question answering
 578 benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*, 2024.
- 579 Joseph Lee Rodgers and W Alan Nicewander. Thirteen ways to look at the correlation coefficient.
 580 *The American Statistician*, 42(1):59–66, 1988.
- 581 Andrew Li, Rahul Thapa, Rahul Chalamala, Qingyang Wu, Kezhen Chen, and James Zou.
 582 Smir: Efficient synthetic data pipeline to improve multi-image reasoning. *arXiv preprint
 583 arXiv:2501.03675*, 2025.
- 584 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
 585 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv
 586 preprint arXiv:2407.07895*, 2024.

- 594 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
 595 mark for question answering about charts with visual and logical reasoning.
 596
- 597 Minesh Mathew, Dimosthenis Karatzas, R Manmatha, and CV Jawahar. Docvqa: A dataset for vqa
 598 on document images. corr abs/2007.00398 (2020). *arXiv preprint arXiv:2007.00398*, 2020.
- 599 Meta AI. Llama 4 Maverick 17B-128E Instruct. [https://huggingface.co/meta-llama/
 600 Llama-4-Maverick-17B-128E-Instruct](https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E-Instruct), 2025. Model card. Accessed: 2025-08-01.
 601
- 602 Mistral AI. Mistral Small 3.1 24B Instruct. [https://huggingface.co/mistralai/
 603 MistralSmall3.124BInstruct2503](https://huggingface.co/mistralai/MistralSmall3.124BInstruct2503), 2025. Model card. Accessed: 2025-08-01.
- 604 OpenAI. GPT-4o. <https://openai.com/index/gpt-4o>, 2024. Accessed: 2025-07-30.
 605
- 606 OpenAI. GPT-4.1. <https://openai.com/index/gpt-4-1/>, 2025a. Accessed: 2025-08-
 607 01.
- 608 OpenAI. GPT-4o-mini. <https://platform.openai.com/docs/models/gpt-4o>,
 609 2025b. Accessed: 2025-08-01.
- 610 OpenAI. OpenAI o3 Reasoning Model. [https://openai.com/index/
 611 introducing-o3-and-o4-mini/](https://openai.com/index/introducing-o3-and-o4-mini/), 2025. Accessed: 2025-07-31.
 612
- 613 Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables.
 614 In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics*
 615 and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long
 616 Papers), pp. 1470–1480, Beijing, China, July 2015. Association for Computational Linguistics.
 617 doi: 10.3115/v1/P15-1142. URL <https://aclanthology.org/P15-1142>.
- 618 Qwen Team. Qwen3-30B-A3B. <https://huggingface.co/Qwen/Qwen3-30B-A3B>,
 619 2025a. Model card. Accessed: 2025-08-01.
 620
- 621 Qwen Team. Qwen3-QwQ-32B. <https://huggingface.co/Qwen/QwQ-32B>, 2025b.
 622 Model card. Accessed: 2025-08-01.
- 623 Reka AI. Reka Flash 3. <https://huggingface.co/RekaAI/reka-flash-3>, 2025.
 624 Model card. Accessed: 2025-08-01.
 625
- 626 Alon Talmor, Ori Yoran, Amnon Catay, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco,
 627 Hannaneh Hajishirzi, and Jonathan Berant. Multimodal{qa}: complex question answering over
 628 text, tables and images. In *International Conference on Learning Representations*, 2021. URL
 629 <https://openreview.net/forum?id=ee6W5UgQLa>.
- 630 Qwen Team. Qwen2.5-vl, January 2025. URL [https://qwenlm.github.io/blog/
 631 qwen2.5-vl/](https://qwenlm.github.io/blog/qwen2.5-vl/).
- 632 TNG Technology Consulting GmbH. Deepseek-r1t-chimera, April 2025. URL <https://huggingface.co/tngtech/DeepSeek-R1T-Chimera>.
 633
- 634 Cyrille Delestre Tom Agonnoude, 2024. URL [https://huggingface.co/datasets/
 635 cmarkea/table-vqa](https://huggingface.co/datasets/cmarkea/table-vqa).
- 636 Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady
 637 Arkhangorodsky, Minjie Xu, Naomi White, and Patrick Lewis. Replacing judges with juries:
 638 Evaluating llm generations with a panel of diverse models. *arXiv preprint arXiv:2404.18796*,
 639 2024.
- 640 Jiankang Wang, Jianjun Xu, Xiaorui Wang, Yuxin Wang, Mengting Xing, Shancheng Fang, Zheneng
 641 Chen, Hongtao Xie, and Yongdong Zhang. A graph-based synthetic data pipeline for scaling high-
 642 quality reasoning instructions. *CoRR*, 2024a.
- 643 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and
 644 Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. In *The
 645 Thirty-eighth Conference on Neural Information Processing Systems Datasets and Benchmarks
 646 Track*, 2024b. URL <https://openreview.net/forum?id=QWTCCxMpPA>.

- 648 xAI. Grok 3 beta — the age of reasoning agents. <https://x.ai/news/grok-3>, 2025. Ac-
649 cessed: 2025-08-01.
650
651 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,
652 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint*
653 *arXiv:2408.01800*, 2024.
654
655 Xinyi Zheng, Douglas Burdick, Lucian Popa, Xu Zhong, and Nancy Xin Ru Wang. Global table
656 extractor (gte): A framework for joint table identification and cell structure recognition using
657 visual context. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*,
658 pp. 697–706. IEEE, 2021.
659
660 Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from
661 natural language using reinforcement learning. *CoRR*, abs/1709.00103, 2017.
662
663 Fengbin Zhu, Wenqiang Lei, Fuli Feng, Chao Wang, Haozhou Zhang, and Tat-Seng Chua. To-
664 wards complex document understanding by discrete reasoning. In *Proceedings of the 30th ACM*
665 *International Conference on Multimedia*, pp. 4857–4866, 2022.
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 APPENDIX
 703
 704
 705
 706

707 A EXTENDED RELATED WORKS
 708
 709

710 The vast majority of table-based QA datasets—such as HybridQA Chen et al. (2020b), WikiTable-
 711 Questions Pasupat & Liang (2015), WikiSQL Zhong et al. (2017), and AIT-QA Katsis et al.
 712 (2022)—represent tables in textual format rather than as rendered images, thereby bypassing the
 713 challenges associated with visual layout interpretation. In contrast, our work focuses exclusively on
 714 multimodal datasets—those that contain both textual and visual (image-based) information. These
 715 can generally be grouped into two main categories: real-world datasets, collected from authentic
 716 documents, and synthetic datasets, generated using automated tools. Real-world multimodal QA
 717 datasets that emphasize tables—such as TAT-DQA Zhu et al. (2022), and TableVQA-Bench Kim
 718 et al. (2024)—tend to be highly domain-specific, limiting diversity in both table layouts and ques-
 719 tion types. For example, TAT-DQA Zhu et al. (2022) combines tabular and textual data from fi-
 720 nancial reports and, while it introduces hybrid contexts for realistic reasoning, its questions rely
 721 heavily on reading textual input rather than interpreting visual structure. Similarly, TableVQA-
 722 Bench Kim et al. (2024) consists of 83% real-world tables (1,250 out of 1,500), primarily sourced
 723 from task-specific datasets such as WikiTableQuestions (information retrieval), TabFact Chen et al.
 724 (2020a) (fact verification), and FinTabNet Zheng et al. (2021) (financial data extraction). Due to
 725 its relatively small size and the specialized nature of its subsets, the dataset exhibits limited visual
 726 diversity. This limitation also extends to the remaining 16% synthetic tables, whose visual variation
 727 is restricted to basic formatting attributes such as background color, border size, font size, and style.
 728 More recently, ChartQA Masry et al. and DocVQA Mathew et al. (2020) have introduced large
 729 open-domain datasets for visual question answering. ChartQA focuses on reasoning over charts and
 730 plots; however, its tasks primarily involve shallow reasoning and do not reflect the structural com-
 731 plexity and layout diversity found in real-world tables. In contrast, DocVQA offers greater diversity
 732 in document layouts and structures, but lacks significant visual challenge—recent VLMs, including
 733 relatively lightweight models like Qwen2 .5-VL-7B-Instruct Team (2025), achieve over 94%
 734 accuracy on this benchmark.

735 Among synthetic datasets, *MultiModalQA* Talmor et al. (2021) stands out as the only open-domain
 736 resource focused on tables. It combines real-world figures, diagrams, and text passages sourced from
 737 Wikipedia, with QA pairs crafted to assess both reasoning and visual comprehension. Although it
 738 incorporates real content, the dataset is considered synthetic due to the way it links independent
 739 modalities and generates QA pairs through formalized templates. However, this approach results
 740 in limited diversity, as the questions are derived from a finite set of templates. In contrast, *Table-
 741 VQA* Tom Agonnoude (2024) is a fully synthetic dataset generated using state-of-the-art LLMs.
 742 Nonetheless, it lacks visual diversity and complexity—its tables follow similar formats and are pre-
 743 dominantly centered around technical domains such as statistics, physics, and algorithms, all of
 744 which are heavily numerical in nature.

745 Most of real-world datasets rely heavily on manual labeling, data collection, and preprocess-
 746 ing—factors that significantly constrain their scalability. Recently, ReachQA He et al. introduced a
 747 more scalable and innovative approach through its Code-as-Intermediary Translation (CIT) pipeline.
 748 This method generates synthetic charts and reasoning questions by leveraging textual intermediaries
 749 such as Python code, demonstrating that advanced reasoning capabilities of large language models
 750 (LLMs) can be effectively transferred to visual models. While ReachQA successfully addresses
 751 both scalability and reasoning complexity, its approach is tailored to chart-based visualizations and
 752 does not extend to structured tabular data.

753 To summarize, existing table-based benchmarks consistently fall short in one or more key areas:
 754 visual diversity, reasoning depth, or scalability. Notably, aside from MultiModalQA, there is no
 755 open-domain dataset designed to evaluate model performance on rendered table images, despite their
 756 prevalence in real-world settings such as reports, academic papers, and spreadsheets. In this work,
 757 we introduce **Visual-TableQA**, a multimodal open-domain synthetic dataset specifically created to
 758 assess reasoning capabilities over table images using LLMs

756 **B CONSIDERATIONS AND LIMITATIONS**
 757

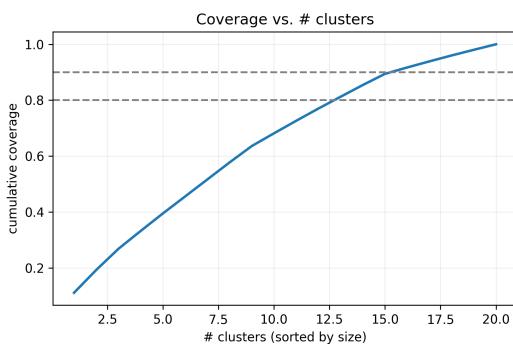
758 The main limitations of this work relate to the use of Code-as-Intermediary Translation (CIT) He
 759 et al. and the assessment of data quality. While we adopted LaTeX as an intermediate representation
 760 for tables, its expressiveness is limited when handling more complex or visually rich images. Devel-
 761 oping a robust, bidirectional image-to-text encoding system remains an open and promising area for
 762 future research. In terms of data quality evaluation, although automatic metrics such as ROSCOE
 763 provide useful insights, they are not yet as reliable as human judgment. As a result, human annota-
 764 tors continue to play a critical role in ensuring high-quality data, especially when scaling synthetic
 765 datasets for reasoning tasks.

766 In addition, we also observed that certain models, such as *Qwen2.5-VL-7B-Instruct*, did not
 767 consistently benefit from Visual-TableQA supervision across all downstream tasks, highlighting a
 768 potential limitation in generalization that warrants further investigation.
 769

770 **C TABLEQA LAYOUT AND TOPIC DIVERSITY**
 771

772 As described in Section 2.2, we sampled 5,000 distinct topics using GPT-4o to serve as inspirations
 773 for table generation. To better illustrate topic diversity, we grouped these topics into 20 semantic
 774 clusters using the K-Means algorithm. Figure 5 displays a 2D projection of these clusters, where
 775 each color represents a distinct semantic group. For the 12 largest clusters, we highlight representa-
 776 tive topics to give a sense of their thematic content.

777 In addition, Figure 4 shows the cumulative percentage of topics covered as clusters are added in de-
 778 scending order of size. The smooth progression of the curve indicates that the clusters are relatively
 779 uniform in size, confirming a balanced distribution of topic diversity throughout the dataset.
 780



792 Figure 4: Cumulative topic coverage as clusters are added by descending size. The uniform slope
 793 indicates an even distribution of topics across clusters.
 794

795 To illustrate the diversity of table layouts produced by our pipeline, Figure 6 displays a side-by-
 796 side comparison between the initial seed tables used in the first generation iteration and a sample
 797 of layouts generated in subsequent steps. The wide range of structures highlights the pipeline’s
 798 capacity to create rich and varied visual designs from limited starting templates.
 799

800 **D TABLE GENERATION SETTINGS**
 801

802 The first stage of the generation pipeline involves **LLM-1**, which is responsible for producing new
 803 tables based on given inspirations. Specifically, it receives one LaTeX-formatted table as a *layout*
 804 *inspiration* and three distinct *topic inspirations*. Based on these, it generates three new tables, each
 805 aligned with one of the provided topics while drawing structural influence from the layout example.
 806 The full prompt used to guide LLM-1 during this step is shown in Figure 7.
 807

808 The second stage of the generation pipeline involves **LLM-2**, which is responsible for producing
 809 question–answer (QA) pairs based on a single LaTeX table. The full prompt used to guide LLM-2
 810 is shown in Figure 8.

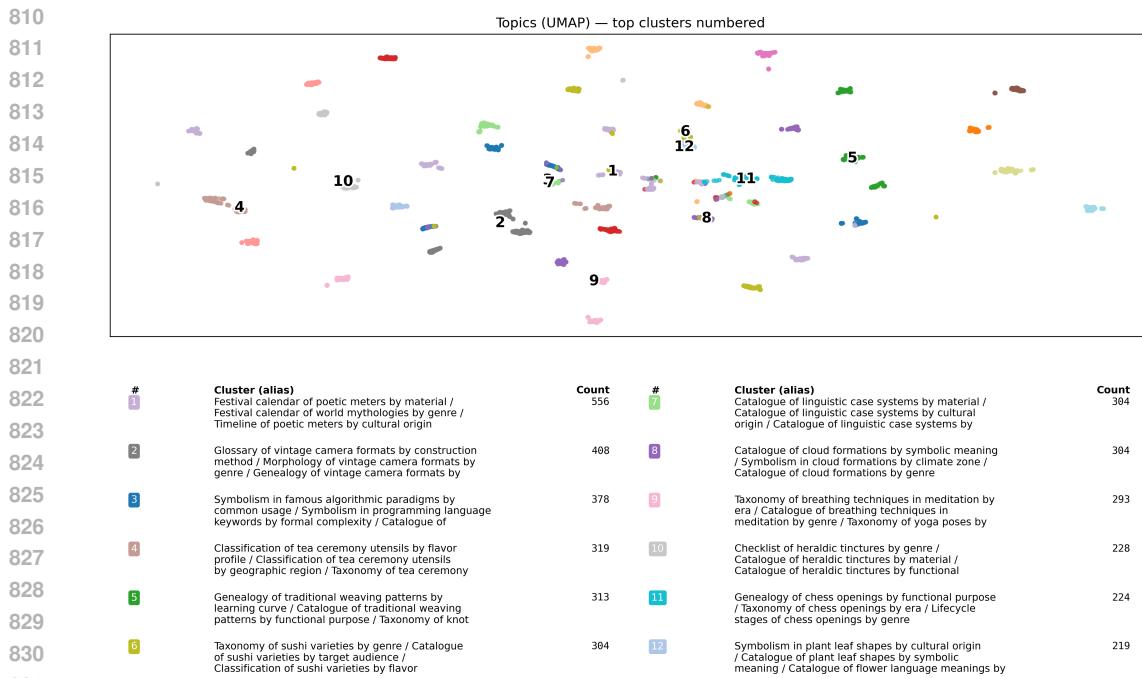


Figure 5: 2D projection of the 5,000 topics using UMAP and K-Means clustering. Each color denotes a semantic cluster. Representative topics are listed for the 12 largest clusters to illustrate diversity.

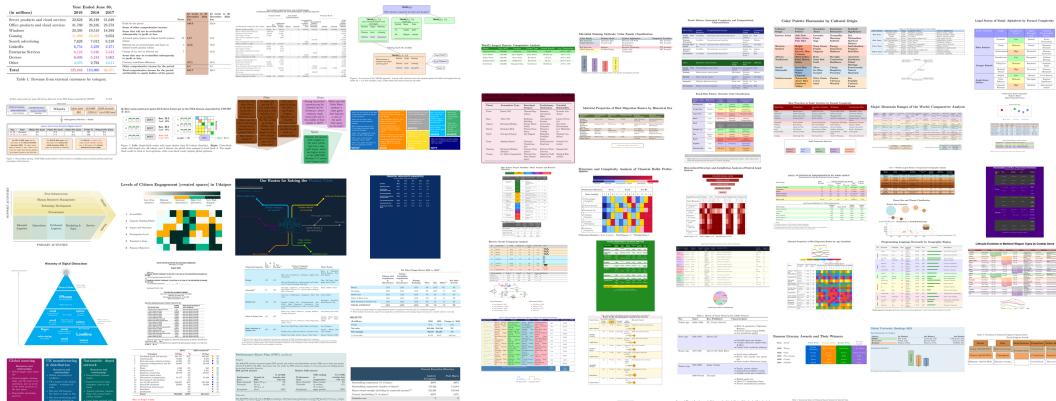


Figure 6: Visual diversity of table layouts. Left: seed layouts used during the first iteration of table generation. Right: layouts generated through cross-model inspiration and iterative refinement.

To encourage creativity while maintaining factual accuracy, the temperature parameter for each model during this phase is set to 0.7.

In Table 7, we present the average validity rates of the generated QA pairs for each model involved in the QA generation phase.

```

864 You are an expert in generating synthetic datasets composed of LaTeX-formatted tables,
865 optionally accompanied by illustrative diagrams. Your task is to produce structured
866 content suitable for data-centric documents, ensuring each table (and diagram, if
867 included) is clear, well-organized, and visually informative.
868 Your final output should start with ```json and end with ``` as plain text, not just formatting
869 . Like this:
870
871 ````json
872 {
873     "table_1": "BEGIN_LATEX
874 <LaTeX code for table 1 (with/without diagram) here>
875 END_LATEX",
876     "table_2": "BEGIN_LATEX
877 <LaTeX code for table 2 (with/without diagram) here>
878 END_LATEX",
879     "table_3": "BEGIN_LATEX
880 <LaTeX code for table 3 (with/without diagram) here>
881 END_LATEX"
882 }
883
884 Requirements:
885 The tables and diagrams will be used to generate reasoning questions. Therefore:
886
887 - If topic inspirations are supplied, ensure every generated table aligns with those
888 topics.
889 - Each LaTeX output must primarily consist of a table. Include a diagram only if it
890 meaningfully complements the table; avoid adding one unnecessarily. Do not generate
891 diagrams alone. If a diagram is empty or non necessary DON'T INCLUDE it.
892 - Keep any diagram minimal---smaller than the table, chart-free, and purely illustrative
893 ---serving only to reinforce the table's content without adding new information.
894 - Each table and their diagram must contain realistic, domain-relevant content. They
895 must be self-contained, include a clear descriptive title and not rely on external
896 data to compile.
897 - The type of information presented should be diverse---such as numerical data or
898 qualitative. The variety and richness of visual elements is essential to the
899 overall quality of the table and their diagram. Table quality should also come with
900 a large number of rows and columns.
901 - Table and diagram layouts should be creatively designed---taking inspiration from
902 reference example (when provided) but incorporating meaningful variations such as
903 colors, multi-row or multi-column cells, custom formatting adjustments, or any
904 other visual enhancement that promotes structural diversity.
905 - Table layouts should be at least as complex as the example provided, don't try to
906 simplify (diagrams are not mandatory). Table complexity should also come with a
907 large number of rows and columns.
908 - Do NOT escape any characters in the LaTeX code. The LaTeX must be written as plain
909 text, exactly as it would appear in a .tex file, with real line breaks and single
910 backslashes (\), not JSON-escaped.
911 - All LaTeX tables and diagrams must be constrained to fit entirely within the printable
912 area of a standard A4 page when compiled to PDF, without overflowing horizontally
913 or vertically. Use appropriate formatting techniques such as adjusting column
914 widths, reducing font size, or enabling landscape mode if necessary but NEVER
915 rotation.
916 - Make sure each LaTeX table and diagram includes all required \usepackage declarations
917 and is enclosed within a complete, compilable LaTeX document structure, including
918 the appropriate preamble and \begin{document}... \end{document} block.
919 - Make sure each LaTeX codes start and end with BEGIN_LATEX and END_LATEX, respectively.
920 - Make sure to wrap your final answer with ```json at the beginning and ``` at the end.
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
20100
20101
20102
20103
20104
20105
20106
20107
20108
20109
20110
20111
20112
20113
20114
20115
20116
20117
20118
20119
20120
20121
20122
20123
20124
20125
20126
20127
20128
20129
20130
20131
20132
20133
20134
20135
20136
20137
20138
20139
20140
20141
20142
20143
20144
20145
20146
20147
20148
20149
20150
20151
20152
20153
20154
20155
20156
20157
20158
20159
201510
201511
201512
201513
201514
201515
201516
201517
201518
201519
201520
201521
201522
201523
201524
201525
201526
201527
201528
201529
201530
201531
201532
201533
201534
201535
201536
201537
201538
201539
201540
201541
201542
201543
201544
201545
201546
201547
201548
201549
201550
201551
201552
201553
201554
201555
201556
201557
201558
201559
201560
201561
201562
201563
201564
201565
201566
201567
201568
201569
201570
201571
201572
201573
201574
201575
201576
201577
201578
201579
201580
201581
201582
201583
201584
201585
201586
201587
201588
201589
201590
201591
201592
201593
201594
201595
201596
201597
201598
201599
2015100
2015101
2015102
2015103
2015104
2015105
2015106
2015107
2015108
2015109
2015110
2015111
2015112
2015113
2015114
2015115
2015116
2015117
2015118
2015119
2015120
2015121
2015122
2015123
2015124
2015125
2015126
2015127
2015128
2015129
2015130
2015131
2015132
2015133
2015134
2015135
2015136
2015137
2015138
2015139
2015140
2015141
2015142
2015143
2015144
2015145
2015146
2015147
2015148
2015149
2015150
2015151
2015152
2015153
2015154
2015155
2015156
2015157
2015158
2015159
2015160
2015161
2015162
2015163
2015164
2015165
2015166
2015167
2015168
2015169
2015170
2015171
2015172
2015173
2015174
2015175
2015176
2015177
2015178
2015179
2015180
2015181
2015182
2015183
2015184
2015185
2015186
2015187
2015188
2015189
2015190
2015191
2015192
2015193
2015194
2015195
2015196
2015197
2015198
2015199
2015200
2015201
2015202
2015203
2015204
2015205
2015206
2015207
2015208
2015209
2015210
2015211
2015212
2015213
2015214
2015215
2015216
2015217
2015218
2015219
2015220
2015221
2015222
2015223
2015224
2015225
2015226
2015227
2015228
2015229
2015230
2015231
2015232
2015233
2015234
2015235
2015236
2015237
2015238
2015239
2015240
2015241
2015242
2015243
2015244
2015245
2015246
2015247
2015248
2015249
2015250
2015251
2015252
2015253
2015254
2015255
2015256
2015257
2015258
2015259
2015260
2015261
2015262
2015263
2015264
2015265
2015266
2015267
2015268
2015269
2015270
2015271
2015272
2015273
2015274
2015275
2015276
2015277
2015278
2015279
2015280
2015281
2015282
2015283
2015284
2015285
2015286
2015287
2015288
2015289
2015290
2015291
2015292
2015293
2015294
2015295
2015296
2015297
2015298
2015299
2015300
2015301
2015302
2015303
2015304
2015305
2015306
2015307
2015308
2015309
2015310
2015311
2015312
2015313
2015314
2015315
2015316
2015317
2015318
2015319
2015320
2015321
2015322
2015323
2015324
2015325
2015326
2015327
2015328
2015329
2015330
2015331
2015332
2015333
2015334
2015335
2015336
2015337
2015338
2015339
2015340
2015341
2015342
2015343
2015344
2015345
2015346
2015347
2015348
2015349
2015350
2015351
2015352
2015353
2015354
2015355
2015356
2015357
2015358
2015359
2015360
2015361
2015362
2015363
2015364
2015365
2015366
2015367
2015368
2015369
2015370
2015371
2015372
2015373
2015374
2015375
2015376
2015377
2015378
2015379
2015380
2015381
2015382
2015383
2015384
2015385
2015386
2015387
2015388
2015389
2015390
2015391
2015392
2015393
2015394
2015395
2015396
2015397
2015398
2015399
2015400
2015401
2015402
2015403
2015404
2015405
2015406
2015407
2015408
2015409
2015410
2015411
2015412
2015413
2015414
2015415
2015416
2015417
2015418
2015419
2015420
2015421
2015422
2015423
2015424
2015425
2015426
2015427
2015428
2015429
2015430
2015431
2015432
2015433
2015434
2015435
2015436
2015437
2015438
2015439
2015440
2015441
2015442
2015443
2015444
2015445
2015446
2015447
2015448
2015449
2015450
2015451
2015452
2015453
2015454
2015455
2015456
2015457
2015458
2015459
2015460
2015461
2015462
2015463
2015464
2015465
2015466
2015467
2015468
2015469
2015470
2015471
2015472
2015473
2015474
2015475
2015476
2015477
2015478
2015479
2015480
2015481
2015482
2015483
2015484
2015485
2015486
2015487
2015488
2015489
2015490
2015491
2015492
2015493
2015494
2015495
2015496
2015497
2015498
2015499
2015500
2015501
2015502
2015503
2015504
2015505
2015506
2015507
2015508
2015509
2015510
2015511
2015512
2015513
2015514
2015515
2015516
2015517
2015518
2015519
2015520
2015521
2015522
2015523
2015524
2015525
2015526
2015527
2015528
2015529
2015530
2015531
2015532
2015533
2015534
2015535
2015536
2015537
2015538
2015539
2015540
2015541
2015542
2015543
2015544
2015545
2015546
2015547
2015548
2015549
2015550
2015551
2015552
2015553
2015554
20
```

918 You are an expert in generating question--answer pairs from LaTeX-formatted. Your task is to
 919 create a structured dataset consisting of visually challenging, reasoning-based questions
 920 and their corresponding answers derived from a given LaTeX formatted table with optional
 921 diagram.
 922
 923 Input:
 924
 925 You will be provided with a sample LaTeX table as context. Based on this table or diagram,
 926 your goal is to generate a JSON object with the following structure:
 927
 928 questions: A python list of 3 challenging questions that require reasoning and analysis
 929 based ONLY on the data presented in the table and the optional diagram. The questions
 930 must be answerable using ONLY the information in the table or diagram(no extra
 931 knowledge).
 932 answers: A python list of 3 detailed answers to the 3 questions, including a clear chain of
 933 thought explaining the reasoning process.
 934
 935 Requirements:
 936
 937 All questions must be relevant to the table's context and designed to test deeper
 938 understanding or inference.
 939 When possible, all questions should make full use of the visual or structural elements of
 940 the table or diagram (such as rows, columns, headers, colors, patterns, diagrams etc.)
 941 while maintaining clear relevance to the table's content.
 942 Questions must be clear and answerable with an objective methodology, no subjective
 943 question.
 944 All entries (both questions and answers) should be returned as lists of string values.
 945 The global result should be a single JSON object wrapped in a markdown code block using `~~`
 946 json at the beginning and `~~` at the end, and containing all two key-value pairs.
 947 This means your output should start with `~~json` and end with `~~` as plain text, not just
 948 formatting.

Figure 8: LLM prompt used for QA generation.

Table 7: Average QA pair validity across different QA generation models. Accuracies are computed from a sample of at least 500 QA pairs per model.

Model	Acc. (%)
Llama 4 Maverick 17B-128E Instruct	88
Gemini 2.0 Flash	89.1
Gemini 2.5 Flash	93.1
Gemini 2.5 Pro	89.3
GPT-4.1	90.4
Qwen3-30B-A3B	76.6
Qwen3-QwQ-32B	92.6
DeepSeek-R1-Distill-Qwen-32B	73.4
DeepSeek-R1T-Chimera	89.4
Claude Sonnet 4	91
Claude 3.5 Haiku	90.2
Grok 3 Beta	89.4
Reka Flash 3	79

- **Text overflow:** Cell content spilling outside the cell boundary, especially in narrow columns or with long strings.
- **Invisible content:** Multirow cells with background colors that obscure cell text (e.g., white text on white background).
- **Improper horizontal lines:** \midrule or \hline splitting across multirow cells, breaking visual coherence.

CONTENT RELEVANCE AND CORRECTNESS

- **Empty or irrelevant tables:** Tables with placeholder content or unrelated to the assigned topic.
- **Incorrect topic alignment:** Generated tables that do not match the intended topic inspiration.
- **Duplicate outputs:** All three tables generated for a prompt are identical or nearly identical in structure/content.

- 972
 973 • **Missing external resources:** References to images or figures not included or available in
 974 the output.
 975 • **Incorrect LaTeX syntax:** Math symbols placed outside of math environments, leading to
 976 compilation errors.

977 DIAGRAM-SPECIFIC ISSUES
 978

- 979 • **Missing tables:** Some generations return only a diagram without an accompanying table.
 980 • **Node placement errors:** Overlapping or misaligned nodes in TikZ diagrams.
 981 • **Arrow misplacement:** Arrows that do not connect to correct nodes or that overlap diagram
 982 elements improperly.
 983 • **Legend/title confusion:** Titles or legends positioned incorrectly or detached from the rel-
 984 evant diagram elements.
 985 • **Visual inconsistencies:** General drawing flaws, such as missing anchors, inconsistent line
 986 styles, or unintended overlaps.

989 These issues highlight the need for a validation loop in the TableQA pipeline and justify the inclusion
 990 of human verification stages to ensure dataset quality.
 991

992 F ROSCOE METRIC SCORES
 993

995 In addition to the LLM jury validation of the Visual-TableQA dataset, we conducted a complemen-
 996 tary quality assessment using the ROSCOE framework Golovneva et al.. This evaluation measures
 997 step-by-step reasoning coherence across multiple dimensions, including semantic alignment, logical
 998 consistency, and contextual grounding. The resulting scores, reported in Table 8, further support the
 999 reliability and high quality of the generated tables and QA pairs, reinforcing the effectiveness of our
 1000 data generation pipeline.

1001 Table 8: ROSCOE Golovneva et al. reasoning metrics averaged over the whole dataset. The “Di-
 1002 rection” column indicates whether higher or lower values correspond to better performance for each
 1003 metric.

Metric	Direction	Mean	Std
Semantic Adequacy (\uparrow)			
Faithfulness-Step	\uparrow	0.99	5e-4
Informativeness-Step	\uparrow	0.99	5e-3
Informativeness-Chain	\uparrow	0.98	1.4e-2
Faithfulness-Token	\uparrow	0.99	2e-3
Avg		0.99	
Redundancy & Risk (\downarrow)			
Repetition-Token	\downarrow	0.06	0.12
Repetition-Step	\downarrow	0.06	0.12
Avg		0.06	
Logical Inference (\uparrow)			
Discourse-Representation	\uparrow	0.68	0.41
Coherence-Step	\uparrow	0.7	0.40
Avg		0.69	
Fluency & Perplexity (\downarrow)			
Perplexity-Step	\downarrow	0.01	0.01
Perplexity-Chain	\downarrow	0.05	0.03
Perplexity-Step-Max	\downarrow	8e-3	8e-3
Avg		0.02	
Grammaticality (\uparrow)			
Grammar-Step	\uparrow	0.96	0.05
Grammar-Step-Max	\uparrow	0.9	0.13
Avg		0.93	

1026 **G LLM JURY RELIABILITY**

1028 Following recent studies He et al.; Fu et al. (2025); Verga et al. (2024), we adopt a high-performing
 1029 LLM jury combined with a majority-vote strategy to evaluate model predictions. This multi-model
 1030 jury setup enhances evaluation robustness and mitigates individual model biases. The configuration
 1031 details for all jury models are provided in Section G.1. As argued in Verga et al. (2024), aggregating
 1032 judgments from several strong LLMs yields more consistent and reliable evaluations than relying
 1033 on a single model. In Section G.2, we discuss the challenges of LLM-based evaluations and we
 1034 report a detailed comparison between our LLM jury and human annotators on Qwen2.5-VL-7B
 1035 predictions across multiple benchmarks. In Section G.3, we provide a detailed analysis of jury-to-
 1036 jury agreement across both table and QA pair quality assessments, as well as benchmark evaluations.
 1037 The results reveal strong inter-annotator alignment, validating the consistency and effectiveness of
 1038 our jury-based evaluation protocol.

1039 **G.1 LLM JURY SETTINGS**

1040 LLM juries were involved at two key stages of the pipeline: (i) **quality assurance**, where generated
 1041 tables and QA pairs were validated before inclusion in the dataset, and (ii) **evaluation benchmark-
 1042 ing**, where model responses were assessed for accuracy and reasoning quality.

1043 The prompt used for quality assurance is shown in Figure 9, while the prompt used for evaluation
 1044 during benchmarking is shown in Figure 10. Each jury consisted of multiple high-performing vision-
 1045 language or reasoning-capable LLMs, and final decisions were made via majority voting.

1046 To ensure consistency and reproducibility across evaluations, all LLM jury calls were executed with
 1047 a temperature setting between 0.0 and 0.1.

1048 We observed a notable drop in judgment accuracy when LLM juries were instructed to return only
 1049 a structured JSON verdict without any preceding explanation. In particular, Mistral-large
 1050 systematically omitted its reasoning whenever the keyword JSON was included in the prompt. This
 1051 issue was effectively mitigated by explicitly instructing models to provide a rationale prior to their
 1052 final decision and by avoiding any direct mention of JSON in the prompt, except for GPT-4.1
 1053 which remained robust under such formatting. Including explicit reasoning significantly improved
 1054 the reliability and depth of model evaluations, likely by reducing shallow assessments and prompting
 1055 more thoughtful judgments.

1056 You are a reasoning question answer expert. You will be given a LaTeX formatted table with/
 1057 without diagram, a list of 3 topics, and a pair of a question and its answer.

1058 Your task is to evaluate the pair of question answer based solely on the data in the LaTeX
 1059 code and these criteria:

- 1060 1) Does the LaTeX code contain a Table (not some charts alone or diagrams alone) ?
- 1061 2) Does the table, any optional diagrams, and the rest of the document are on one single
 topic from the provided list of topics, and internally consistent (be careful to off-
 topic diagrams)?
- 1062 3) Is the question clear and related to the table or the diagram?
- 1063 4) Is the answer (including its reasoning) totally valid and does it actually respond to
 the question?
- 1064 5) Is the answer FULLY supported by and ONLY BY the table or diagram data (no extra
 knowledge)?

1065 If the five criteria are true, mark the pair as correct.
 1066 If one of the criteria is not met, mark it as incorrect.

1067 Think step by step and conclude with your decision and the index of the criterium not met (if
 1068 none, index is 0) as follows:
 1069 JSON_mention
 1070 `{"decision": [0, index_of_the_criterium_not_met]}` for incorrect or `{"decision": [1, 0]}`
 1071 for correct.

1072 **Figure 9: LLM prompt used for QA evaluation.**

```

1080 You are an expert evaluator of question-answer pairs. You will be given a question a model's
1081 answer and a ground truth answer (reference).
1082 Evaluate the answer based on these criteria:
1083 1) Is the model's answer logically consistent?
1084 2) Does the model's answer convey the same meaning as the ground truth?
1085 If the two criteria are true, mark the pair as correct. If one of the criteria is not met,
1086 mark it as incorrect.
1087 Think step by step and conclude with your verdict and the index of the criterium not met (if
1088 none, index is 0) as follows:
1089 {{"verdict": [0, index_of_the_criterium_not_met]}} for incorrect or {{"verdict": [1, 0]}} for
1090 correct.
1091 Question: {question}
1092 Answer: {answer}
1093 Ground Truth: {ground_truth}
1094 Response:
1095

```

Figure 10: LLM prompt used for Benchmark evaluation.

G.2 LLM JURY LIMITATIONS AND MITIGATION STRATEGIES

Our evaluation protocol measures alignment between a model’s prediction and the reference answer, leveraging LLM juries for semantic comparison. This strategy is highly cost-effective, enabling scalable automatic assessments on reasoning tasks. However, it comes with important limitations. A model’s response may rely on external world knowledge rather than explicitly extracting information from the table image. In such cases, LLM juries lack the ability to verify whether the response is grounded in the table itself, as they do not have access to the visual context. This weakness makes the evaluation pipeline vulnerable to false positives, especially when the model outputs a factually correct answer that is not actually derivable from the table content. This could be mitigated by including rendered table images alongside the question and model prediction at evaluation time—giving LLM juries full visual grounding. Unfortunately, While effective, this approach is significantly more costly in terms of API calls and inference latency, making it challenging to scale on large datasets. These observations underscore the need for hybrid evaluation strategies combining automatic LLM-based judgments with human verification, in order to quantify the risk of jury-related errors at scale.

To rigorously assess this risk, we conducted a manual evaluation of Qwen2.5-VL-7B-Instruct’s predictions on both Visual-TableQA and the *VTabFact* split of TableVQA-Bench, one of the benchmarks where the model exhibited degraded performance after fine-tuning. Our manual analysis serves as a ground truth reference to assess the reliability and limitations of LLM-based jury evaluation. Our evaluation considers not only the final answer, but also the validity of the reasoning process: even predictions with the correct answer are marked as incorrect if the chain of thought is flawed. We used the manual analysis to validate the judgments made by the LLM juries. We identified two types LLM Juries errors:

- *False Positive*: The LLM jury mistakenly accepts an incorrect model prediction as correct.
- *False Negative*: The LLM jury mistakenly rejects a correct model prediction as incorrect.

The errors distribution are detailed in Table 9. The results reveal that the LLM jury aligns with human judgment within a reasonable margin of 4.7%, even under our strict annotation protocol that penalizes incorrect reasoning regardless of the final answer. A more relaxed evaluation criterion would yield an even smaller discrepancy. These findings reinforce the reliability of our LLM jury setup, demonstrating its effectiveness as a scalable proxy for human evaluation.

Table 9: Misclassifications by LLM juries on Visual-TableQA and *VTabFact*. Percentages are calculated relative to the total number of evaluated examples in each dataset.

Error Type	Visual-TableQA	<i>VtabFact</i>
False Positive	3.45%	1.6%
False Negative	1.2%	0.0%
Total	4.65%	1.6%

1134 **G.3 LLM JURY AGREEMENT ANALYSIS**

1135
 1136 Figure 11a presents a detailed breakdown of agreement levels between individual LLM juries, as
 1137 well as their alignment with majority-vote annotations. This analysis was conducted across both
 1138 table and QA pair quality assessments. The analysis reveals a spectrum of consistency across juries,
 1139 with GPT-4 .1 emerging as the most reliable, likely due to its robust handling of edge cases. Among
 1140 all models, proprietary LLMs such as Gemini-2 .5-pro and GPT-4 .1 show the strongest align-
 1141 ment with the majority vote, while Deepseek-v3 .1 exhibits the weakest agreement. Notably,
 1142 pairwise jury agreement patterns appear correlated with the models’ reasoning capabilities. Despite
 1143 variability in alignment strength, all juries demonstrate a meaningful degree of concordance with
 1144 the majority, underscoring the robustness of our collective evaluation protocol.

1145 Conversely, Figure 11b shows consistently strong jury agreement across all models for benchmark
 1146 evaluations, with no notable divergence between proprietary and open-source LLMs. This can be at-
 1147 tributed to the relatively simpler nature of the task (semantic comparison between model predictions
 1148 and ground truth) compared to the more complex reasoning required for evaluating table-question-
 1149 answer triples, as analyzed in Figure 11a.

1150
 1151 **G.4 EVALUATOR CONSISTENCY COMPARED TO THE LITERATURE**

1152
 1153 In line with recent studies He et al.; Fu et al. (2025); Verga et al. (2024); Agarwal & Ciucă (2025), we
 1154 employed a high-performing LLM jury with a majority-vote strategy to evaluate model predictions.
 1155 The jury consisted of Mistral-large, Deepseek-v3 .1, Gemini-2 .5-pro, GPT-4 .1,
 1156 and Deepcogito-v2. In Table 4, baseline values (in black) are reported from He et al., who used
 1157 GPT-4o as the sole evaluator. While our evaluation pipeline involves a broader and more powerful
 1158 set of models, making it arguably more reliable and robust (Verga et al., 2024), we still consider the
 1159 two evaluation protocols broadly comparable. In fact, due to the stricter majority-vote requirement
 1160 across diverse models, our approach may even yield more demanding or rigorous evaluations. This
 1161 comparability also holds for Table 6, where baseline performances (in black) are taken from Fu
 1162 et al. (2025), who employed GPT-4 for their evaluation. We argue that despite methodological
 1163 differences, all evaluations are consistent enough to be analyzed jointly for comparative purposes.

1164 **H MODEL FINETUNING HYPERPARAMETERS**

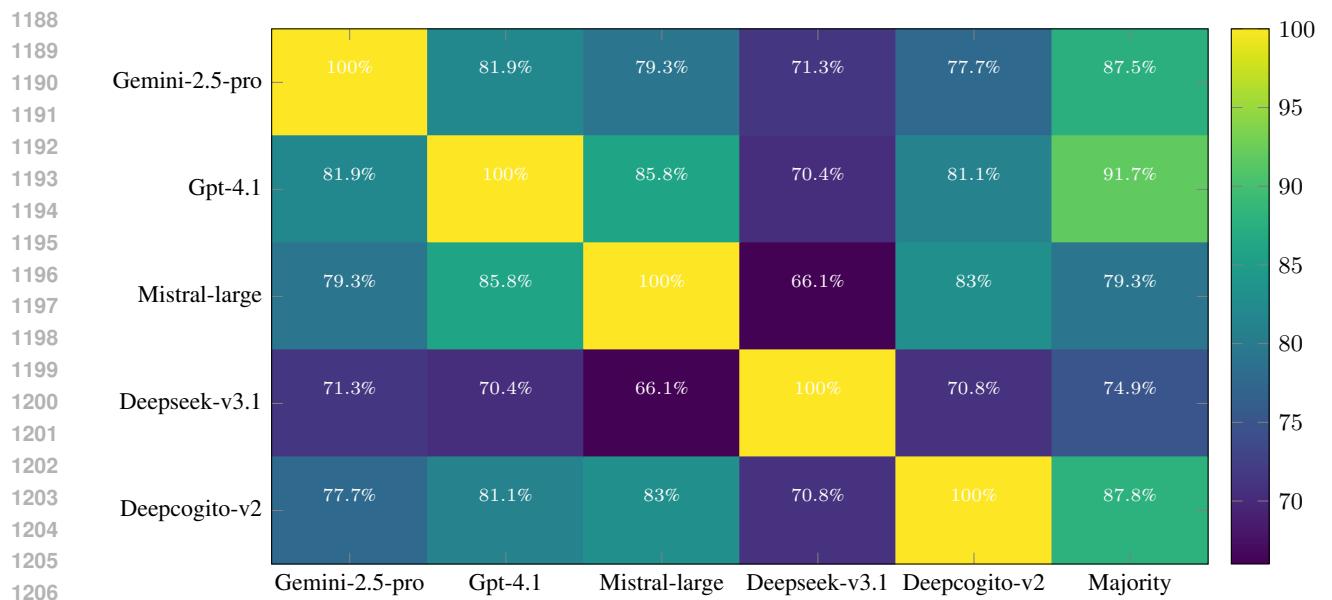
1165
 1166 The hyperparameters used for LoRA are reported in Table 10. For Qwen2 .5-VL-7B-Instruct,
 1167 we employed a two-phase fine-tuning strategy. In the first phase (Tier A), we adapted the text-side
 1168 modules together with the multi-modal projector, while keeping the vision tower frozen. In the
 1169 second phase (Tier B), we further enabled LoRA on the attention projections of the last four vision
 1170 blocks, leaving the remaining vision layers untouched.

1171
 1172 Table 10: Hyperparameters Used for Fine-Tuning with LoRA. More details in our GitHub repos-
 1173 itory. Abbreviations: lr=learning rate, r= LoRA rank, α = LoRA α , Targets=targets modules for
 1174 LoRA.

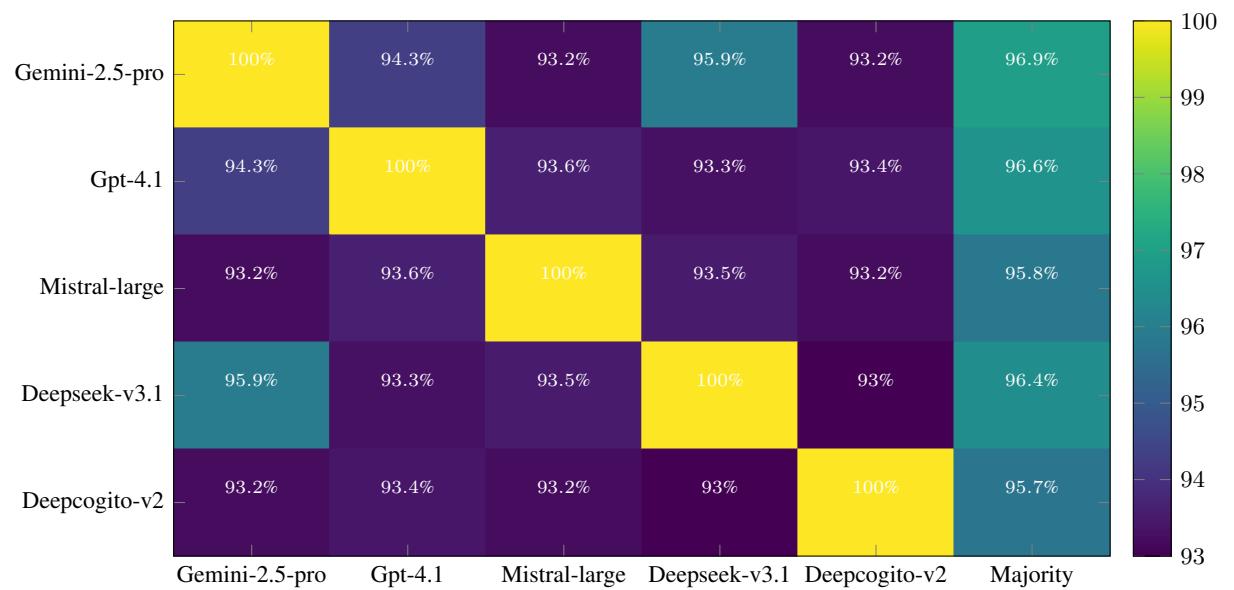
Model	lr	r	α	Targets
LLaVA-Next-Llama3-8B				
MiniCPM-V2.5-Llama3	2e-5	16	8	all-linear (llm frozen)
InternVL2-8B		16*	32*	
Qwen2.5-VL-7B-Instruct	1e-4 2e-5	16 8	8 32	Tier A: projector modules ** Tier B: last 4 vision blocks’ attention

1183 * The InternVL training source code sets the LoRA alpha as twice
 1184 the LoRA rank, as shown in their official implementation here. We
 1185 followed this convention for full reproducibility and assumed that
 1186 other baselines applied the same rule.

1187 ** We left the Vision Tower untouched as it significantly degraded
 1188 model performance.



(a) Pairwise agreement (%) between LLM juries, plus alignment with majority vote for tables and QAs quality assessment.



(b) Pairwise agreement (%) between LLM juries, plus alignment with majority vote for benchmark evaluation.

Figure 11: LLM juries agreement during evaluation.

I ERROR TAXONOMY OF MODEL PREDICTIONS

To better understand the failure modes of *Qwen2.5-VL-7B-Instruct*, we conducted a fine-grained manual analysis of its predictions, both before and after fine-tuning on Visual-TableQA and the *VTabFact* split of TableVQA-Bench. This section is organized as follows: (i) we define the full set of error types used in our annotation protocol (Section I.1); (ii) we present a comparative analysis of errors observed on *VTabFact* (Section I.2).

1242 I.1 ERROR CATEGORIES
12431244 We classified the observed errors into eight categories (Figure 12):
1245

- **Partial Data Extraction:** The model overlooks some relevant entries (e.g., stops counting too early, Figure 12a).
- **Hallucination:** The model references information not present in the table (Figure 12b).
- **Incoherence:** The model extracts the correct data but then misinterprets it or contradicts itself later (deductive error, Figure 12c).
- **Misunderstanding:** The model produces factual statements that do not actually answer the question (Figure 12d).
- **Faulty Methodology/Reasoning:** The reasoning is too shallow, or the model fails to satisfy all the constraints of the query (Figure 12e).
- **Ambiguity (Gray Area):** Both the ground truth and the model’s prediction can be reasonably justified (Figure 12f).
- **Dataset Mistake:** The original dataset contains annotation or label errors.

1260 I.2 ERROR ANALYSIS ON VTABFACT
1261

The results of our analysis are shown in Figure 13. The side-by-side comparison of Figure 13a and Figure 13b reveals a significant shift in error distribution after fine-tuning: although the total number of errors increases from 67 to 88 (out of 250 samples), most newly introduced errors concentrate in the *Incoherence* and *Hallucination* categories, while all other error types show a marked decline. In particular, *Faulty Methodology/Reasoning* errors are significantly reduced, indicating that the fine-tuned model exhibits more consistent and structured reasoning patterns. This trend is further supported by Figure 13c, which shows the distribution of errors corrected by the fine-tuned model: improvements span across all error categories. In contrast, Figure 13d shows that newly introduced errors after fine-tuning are largely concentrated in only a few categories. We attribute the sharp rise in *Incoherence* errors to two key factors:

Degraded Arithmetic Capability: A recurrent issue post-finetuning involves numerical comparison errors or basic math computation errors, leading to faulty deductions despite otherwise correct reasoning steps. Such errors were less frequent in the pretrained model.

Answer Template Behavior: The fine-tuned model tends to state its final answer before providing its chain of thought, often declaring an incorrect answer followed by a valid justification that leads to the correct conclusion. Despite the correctness of the rationale, we count such instances as errors. In contrast, the pretrained model typically presents its reasoning first, then concludes—leading to fewer coherence violations.

1281 J VISUAL-TABLEQA SAMPLE
12821283 Table 11 gives some more detailed examples of our dataset samples.
12841285 K IMAGE-TO-LATEX DATASET
1286

We have also constructed an additional dataset, **Img2TeX**, which will be publicly released upon paper acceptance. It contains all the table images generated during the construction of **Visual-TableQA**, along with their corresponding LaTeX source code. This dataset is intended to complement the work of Kale & Nadadur (2025), which focuses on evaluating models’ ability to generate LaTeX from textual prompts. In contrast, **Img2TeX** targets the task of generating LaTeX documents from visual (image-based) inputs. This distinction opens up a new evaluation pathway for vision-language models (VLMs), particularly those aiming to learn structured document generation from visual cues. As such, **Img2TeX** serves as a valuable benchmark for assessing the visual-to-structured-text generation capabilities of multimodal models.

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

artist	country of origin	period active	release - yes or first charted record	genre	claimed sales
eagles	united states	1971 - prese	1972	soft rock / country rock	150 million
rihanna	barbados united states	2005 - prese	2005	r&b / pop / dance hip hop	150 million
u2	ireland	1976 - prese	1980	rock	150 million
billy joel	united states	1964 - prese	1971	pop / rock	150 million
phil collins	united kingdom	1980 - 2011	1981	adult contemporary	150 million
aerosmith	united states	1970 - prese	1973	hard rock	150 million
frank sinatra	united states	1935 - 1995	1939	pop / swing	150 million
barbra streisand	united states	1960 - prese	1963	pop / adult contemporary	145 million
garth brooks	united states	1989 - 1999	1989	country	130 million
genesis	united kingdom	1967 - 1999	1969	progressive rock / pop rock	130 million
donna summer	united states	1968 - 2012	1974	pop / disco / r&b	130 million
neil diamond	united states	1966 - prese	1966	pop / rock	125 million
bruce springsteen	united states	1972 - prese	1973	rock	120 million
bee gees	united kingdom	1963 - 2003	1963	pop / disco	120 million
julio iglesias	spain	1968 - prese	1968	latin	120 million
dire straits	united kingdom	1977 - 1995	1978	rock / pop	120 million

Q: there be 6 artist with claimed sale of 150 million?

A: To determine if there are 6 artists with claimed sales of 150 million, let's count the entries in the "claimed sales" column where the value is 150 million:

1. Eagles [...] 6. Aerosmith
There are indeed 6 artists listed with claimed sales of 150 million.

(a) Partial Data Extraction

fix_id	fixer_fixes	score	empty_fixes	date
1	brighton	1 - 0	port vale	24 january 1993
2	bury	1 - 2	exeter city	24 january 1993
3	cardiff	2 - 0	brighton & hove albion	24 january 1993
4	derby county	1 - 1	leicester city	24 january 1993
5	bolton wanderers	1 - 1	nottingham forest	24 january 1993
6	newcastle united	2 - 1	bolton wanderers	28 january 1993
7	grimsby town	1 - 0	manchester united	24 january 1993
8	liverpool	1 - 0	leicester city	24 january 1993
9	sheffield united	4 - 1	notts county	24 january 1993
10	hereford	2 - 1	sheffield wednesday	24 january 1993
11	leeds united	0 - 0	nottingham forest	24 january 1993
12	nottingham forest	1 - 2	nottingham forest	24 january 1993
13	wolverhampton wanderers	1 - 2	bolton wanderers	28 january 1993
14	crystal palace	0 - 0	exeter	24 january 1993
15	notts county	1 - 0	nottingham forest	24 january 1993
16	brighton	2 - 0	bolton	24 january 1993
17	brighton park ravens	2 - 0	bolton	24 january 1993
18	notts county	2 - 1	blackpool	24 january 1993
19	leeds united	4 - 1	newcastle united	24 january 1993

Q: of the 18 game play , 8 result in 1 or more team be shut out

A: To determine how many of the 18 games resulted in at least one team being shut out, we need to analyze each game's score:
1. Birmingham vs Port Vale: 2-0 (Birmingham shut out) 2. Bury vs Exeter City: 1-2 (Exeter City shut out) [...]

(d) Misunderstanding

(e) Faulty Methodology

Figure 12: Illustration of the Error Categories

constituency number	name	reserved for (sc/st/non e)	district	number of electorates (2009)
9	ater	none	bhind	177334
10	bhind	none	bhind	197183
11	lahar	none	bhind	205839
12	mehgaon	none	bhind	210649
13	gohad	sc	bhind	166893
20	sewda	none	datia	130161
21	bhander	sc	datia	139600
22	datia	none	datia	143593
total :	total :	total :	total :	1371252

Q: datia 's constituency number be lower than gohad 's by 10

A: To determine if Datia's constituency number is lower than Gohad's by 10, we need to compare their respective constituency numbers.

1. Identify the constituency number for Datia: The table shows that Datia's constituency number is 20. [...]

(b) Hallucination

district	incumbent	party	first elected	result	candidates
new york 1	john smith	democratic-republican	1799 (special)	re-elected	john smith (1) 100%
new york 2	none (district created)	none (district created)	none (district created)	new seat federalist	johanna sandra (1) \$1.2M john brockman (1) 100% new seat created (1) 84.4% johnathan hubbard (1) 15.6%
new york 5	none (district created)	none (district created)	none (district created)	new seat democratic	isaac bloom (1) 55.4% samuel jacobson (1) 44.6%
new york 6	none (district created)	none (district created)	none (district created)	new seat democratic	group of 10 (1) 29.1% matthew masters (1) 48.8%
new york 10	none (district created)	none (district created)	none (district created)	new seat federalist	list gain
new york 12	christopher a. schultz	democratic-republican	1800	re-elected	christopher a. schultz (1) 94.1% diana m. williams (1) 5.9%
new york 14	none (district created)	none (district created)	none (district created)	new seat democratic	erazus root (1) 57.4% benjamin g. phillips (1) 42.6%

Q: the majority of the district be new creation and do not have an incumbent

A: To determine if the majority of districts were newly created and did not have an incumbent, let's analyze the "incumbent" column:

1. **New York 1**: Incumbent: John Smith [...] 6. **New York 12**: Incumbent: David Thomas [...] From this analysis:- There is only one district (New York 1) that has an incumbent. [...]

(c) Incoherence

week	date	opponent	result	attendance
1	september 7 1980	skagit sounders	w 36 - 13	62602
2	september 14 1980	oakland raiders	w 26 - 24	51503
3	september 21 1980	minnesota vikings	w 27 - 20	54200
4	september 28 1980	kansas city chiefs	w 24 - 7	45161
5	october 5 1980	buffalo bills	w 26 - 21	51982
6	october 12 1980	oakland raiders	w 18 - 26	44626
7	october 19 1980	minnesota vikings	w 27 - 13	50027
8	october 26 1980	dallas cowboys	w 12 - 31	400639
9	november 2 1980	chicago bears	w 31 - 14	46406
10	november 9 1980	denver broncos	w 21 - 11	51425
11	november 16 1980	new england patriots	w 26 - 10	52024
12	november 23 1980	minnesota vikings	w 27 - 20	63013
13	november 30 1980	philadelphia eagles	w 27 - 21	51567
14	december 7 1980	washington redskins	w 18 - 14	48956
15	december 14 1980	new england patriots	w 26 - 14	50000
16	december 21 1980	arizona cardinals	w 26 - 13	51795

Q: the smallest crowd the charger play in front of be 44826, a road game against the raider

A: To determine if the smallest crowd the Chargers played in front of was 44,826 during a road game against the Raiders, let's follow these steps:

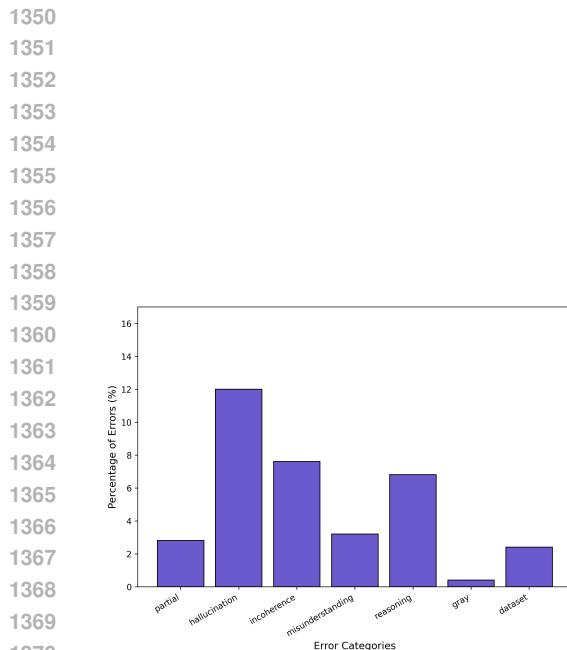
1. Identify the games where the Chargers played against the Raiders: [...] 2. Check the attendance for each game: [...] 3. Compare the attendance figures: [...]

Q: hd 178428 's arrival date be 10 year later than that of hd 190406

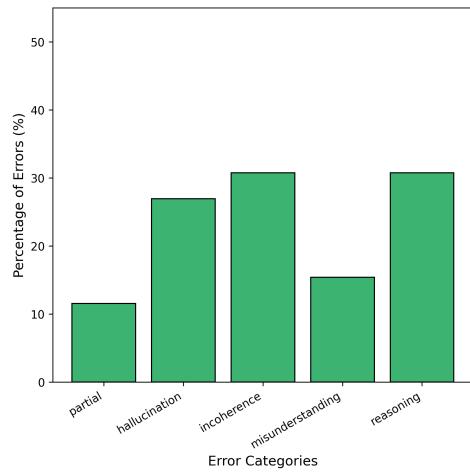
A: [...] Since the difference is 10 years and 10 months, it is not exactly 10 years. Therefore, HD 178428's arrival date is not 10 years later than that of HD 190406.
Ground Truth: Yes

(f) Gray Area

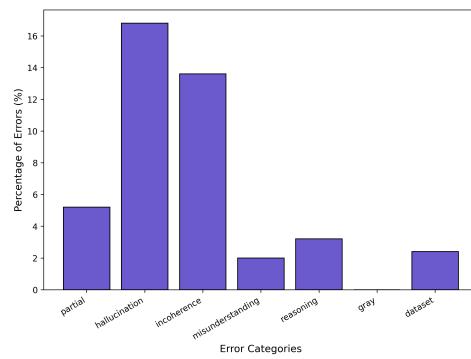
designation id	constellation	date sent	arrival date	message
hd 186408	cygnus	may 24 1999	november 206 9	cosmic call 1
hd 190406	sagitta	june 30 1999	february 2057	cosmic call 1
hd 178428	sagitta	june 30 1999	october 2067	cosmic call 1
hd 190360	cygnus	july 1999	april 2031	cosmic call 1
hd 4872	cassiopeia	july 6 2003	april 2036	cosmic call 2
hd 245049	orion	july 6 2003	august 2040	cosmic call 2
hd 75732	cancer	july 6 2003	may 2044	cosmic call 2
hd 10307	andromeda	july 6 2003	september 204 4	cosmic call 2
hd 95128	ursa major	july 6 2003	may 2049	cosmic call 2



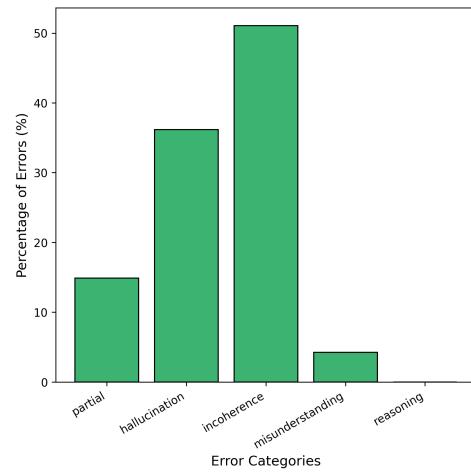
(a) Distribution of Errors for Qwen2.5-v1-7b.



(c) Distribution of Errors corrected by Finetuning.



(b) Distribution of Errors for Qwen2.5-v1-7b after Finetuning.



(d) Distribution of Errors introduced after Fine-tuning.

Figure 13: Comparison of error distributions between the pretrained Qwen2.5-VL-7B model (left) and its finetuned version (right). Percentages in the first row refer to the full *VTabFact* set. Percentages in the second row correspond to: (c) samples where the finetuned model was correct but the pretrained model was not, and (d) samples where the pretrained model was correct but the finetuned model was not.

1404
1405
1406
1407

1408 Table 11: Sample of reasoning-intensive QA pairs. The first row's question and answer are truncated for readability. These questions address multiple visual aspects and extend beyond simple
1409 information extraction to test interpretive reasoning, as illustrated in the second row with a "how"
1410 question.
1411

1412
1413

		Table	Question	Answer
Genealogy of Microbial Staining Methods by Risk Category and Application				
<small>Staining Method Year Introduced Primary Application Risk Category Key Reagent Derivative Methods</small>				
Gram Stain	1884	Bacterial Classification	Low	Crystal Violet Nesser, Hucker
Ziehl-Neelsen	1882	Acid-fast Bacteria	Moderate	Catfish Fuchsin Krause, Antonius, Rhodamine
Ladsporin Stain	1922	Spore Detection	Moderate	Methylene Green Schiff-Potter
Caputte Stain	1902	Caputte Visualization	Low	India Ink Anthony, Manual
Flagella Stain	1917	Motility Study	High	Taenie Acid Leftson
Silver Stain	1937	Fungi / Protozoa	Very High	Silver Nitrate Gomori, Warthin-Starry
Giemsa Stain	1904	Blood Parasites	Moderate	Eosin Azur Wright, May-Grünwald
Acid-Fast Fluorochrome	1950	Mycobacteria	High	Auramine O Auramine Rhodamine
Pseudocolor Acid-Schiff	1946	Fungal Cell Walls	High	Schiff Reagent PAS-Diastase
Wright Stain	1902	Blood Smears	Moderate	Methylene Blue Giemsa, May-Grünwald
Neisser Stain	1890	Corynebacteria	Low	Methylene Blue Giemsa, Albert
Albert Stain	1921	Metachromatic Granules	Moderate	Toluidine Blue Neisser
Auramine-Rhodamine	1950	Acid-fast Bacteria	High	Auramine O Ziehl-Neelsen
Gomori Methenamine Silver	1953	Fungi Detection	Very High	Silver Nitrate Silver Stain
PAS-Diastase	1957	Glycogen Removal	High	Diastase PAS
<pre> graph TD GS[Gram Stain] --> ZN[Ziehl-Neelsen] ZN --> N[Neisser] N --> A[Albert] A --> AR[Auramine-Rhodamine] </pre>				
Staining Method Notes				
<small>Note: Risk categories are based on reagent toxicity and laboratory safety. Derivative methods are direct adaptations or improvements of the original technique.</small>				

1431

Table 1: Pioneering Film Transitions by Era			
Transition Type	First Usage Year	Notable Films	Technical Diagram
Match Cut	1921	The Kid (Keaton)	
	1960s (modern use) 2010	2001: A Space Odyssey Inception	Matched Object →
Wipe	1916	The Cabinet of Dr. Caligari	
	1977	Star Wars (light saber wipe)	□
Jump Cut	1920s (avant-garde)	Un Chien Andalou	[J]
	1903	The Great Train Robbery	
Dissolve	1930s (standardized)	Gone with the Wind	
	1960s (soft focus)	Blowup	◎
	2000s (digital)	Inglourious Basterds	
Crossfade	1910s	The Birth of a Nation	○
	2010s (dynamic)	Mad Max: Fury Road	□

1443

1444
1445
1446

1447
1448
1449

Table 1: Interaction Matrix of Lighting Setups for Various Photographic Materials							
Material	Diffused Light	Spotlight	Backlight	Side Light	Reflected Light	Color coding: Green = optimal, Red = problematic, Blue/Orange/Yellow = mixed results	
Skin	Soft shadow, natural texture	Optimal	Color cast, overexposure	Lighting artifacts	Yellow/orange	Blue	
Glass	No distortion, clear reflection	Overly bright highlights	Silhouette definition	Prismatic effects	Blue/purple		
Metal	Even reflectance control	Hotspots, glares	Edge enhancement	Surface texture emphasis	Color shifts		
Leather	Smooth grain preservation	Overexposure, grain texture	Edge definition	Dark texture emphasis	Color saturation loss		
Wood	Grain visibility optimization	Ring-like highlights	Knot prominence				

1450

1451
1452
1453

1454
1455
1456

1457

1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

The lineage in the diagram shows Gram Stain → Neisser → Albert. We use the table to find the risk category for each. 1. The 'Gram Stain' is the starting point. Its 'Risk Category' cell has a light green background, corresponding to a 'Low' risk. 2. The next step is the 'Neisser Stain'. Its 'Risk Category' cell also has a light green background, indicating the risk remains 'Low'. 3. The final step is the 'Albert Stain'. Its 'Risk Category' cell has a light yellow background, corresponding to a 'Moderate' risk. Therefore, the risk category progression is from Low (Gram) to Low (Neisser), and then increases to Moderate (Albert).

Following the lineage depicted in the diagram, from Gram Stain to Albert Stain, describe the progression of the risk category at each step as indicated by the cell colors in the table.

Which transition type's technical diagram is represented by a gray circle with a white center, and in which era was it first used?

The 'Dissolve' transition's diagram is a gray circle with a white center. It was first used in 1903, as shown in the first row of the 'Dissolve' category. This is identified by matching the diagram description to the transition type and verifying the first usage year listed in the table.

If a photographer wants to capture the natural color tones of Leather without introducing any discoloration, which lighting setup(s) should be avoided and which should be used? Justify using the table's descriptions and color coding.

To preserve Leather's natural colour without discoloration, the photographer should **avoid**:
- **Reflected Light** (yellow mixed) because its description is 'Shade discoloration', indicating a colour shift.
- **Spotlight** (red problematic) because it causes 'Overly contrasted creases', which also alters the colour balance. The recommended lighting is **Diffused Light** (green optimal) and described as 'Natural grain preservation', keeping the colour true.
Side Light (orange mixed) mentions 'Grain texture emphasis' and does not note any colour issue, so it can be used optionally if texture enhancement is desired. Thus, use Diffused Light (and optionally Side Light) while avoiding Reflected Light and Spotlight.