

MMTABREAL: Real-World Benchmark for Multimodal Table Understanding

Anonymous ACL submission



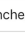








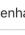


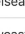








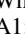

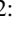
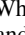
Abstract

Multimodal tables i.e. tabular layouts interleaved with charts, maps, icons, and color encodings are ubiquitous in real applications yet remain difficult for Multimodal Large Language Models (MLLMs). Despite advances in text and image understanding, systematic evaluation of table-centric multimodal reasoning is limited. We introduce MMTABREAL, a MultiModal Table Benchmark, human-curated suite of 500 real-world tables paired with 4,021 question-answer pairs. MMTABREAL spans four question types, five reasoning categories, and eight structural archetypes. Evaluations of state-of-the-art models reveal substantial gaps, especially in visual grounding, spatial alignment, and multi-step inference, with 20–40% performance drops relative to existing benchmarks. These results highlight the need for architectures that more tightly fuse vision with tabular structure and support explicit numerical/logical operations. MMTABREAL is released for evaluation only, providing a rigorous, reproducible testbed that reflects the linguistic, structural, and reasoning complexity of real-world multimodal tables.

1 Introduction

Modern AI is shifting from unimodal perception to structure-aware multimodality, where models jointly read and reason over images embedded in tabular data in both structured and semi-structured (Lee et al., 2023; Cloutier and Ravasi, 2021). As systems move beyond single-channel inputs, seamless cross-modal integration becomes essential for human-level understanding across domains such as healthcare, finance, and education (Joshi, 2022).

Tables are a cornerstone of information representation, providing a two-dimensional scaffold for complex data in both structured and semi-structured forms (Shwartz-Ziv and Armon, 2022;

Team	Pts.	MP	GD	Form	Win %	Wins Draws Losses
 Manchester City	91	38	62			28 7 3
 Arsenal	89	38	62			28 5 5
 Liverpool	82	38	45			24 10 4
 Aston Villa	68	38	15			20 8 10
 Tottenham	66	38	13			20 6 12
 Chelsea	63	38	14			18 9 11
 Newcastle Utd	60	38	23			18 6 14
 Manchester Utd	60	38	-1			18 6 14
 West Ham	52	38	-14			14 10 14

- Q1: Which team has the least number of draws?
A1: Arsenal F.C.
- Q2: What is the color of the badge of the club with 60 points and a negative goal difference? A2: Red
- Q3: What is Chelsea’s longest winning streak? A3: 3
- Q4: Which club has the lowest goal difference for a team with a win percentage greater than 50? A4: Tottenham Hotspurs

Figure 1: A table of English Premier League standings along with accompanying questions and their respective answers from MMTABREAL.

Jiang et al., 2025). Modern tables are no longer text-only: they embed charts, images, color encodings, icons, and logos (Zheng et al., 2024). Unlike fixed-schema databases, this intuitive visual layout is easy for humans to parse but introduces distinct reasoning challenges: (i) header–cell hierarchies are implicit and must be inferred, and (ii) the semantics of adjacent cells are context-dependent (Shigarov, 2023).

Reasoning over such data, where textual, numerical, and visual cues are interwoven within a structured layout, is inherently challenging. For example, consider the question from Figure 1: “Which club has the lowest goal difference among teams with a win percentage greater than 50%?” Answering this requires multi-step, cross-modal reasoning: (1) detecting the visual bars in the Win % column to identify clubs exceeding 50%; (2) comparing numerical values in the GD (goal difference) column; (3) recognizing club badges, color encodings, and logos for correct team identification; (4) aligning these visual cues with the textual headers and cell

*Equal Contributor †Primary Supervisor

contents to filter candidates; and (5) concluding that *Tottenham Hotspur* satisfies the condition with the lowest goal difference among the qualifying teams. Despite their ubiquity, multimodal tables remain underexplored in AI: most work targets generic vision-language or table-only tasks and overlooks the unique, structure-aware challenges of mixed-format tables (Vaishnav and Tammet, 2025). This gap motivates a central question: *Can today’s AI systems, particularly multimodal large language models (MLLMs), effectively reason over complex multimodal tables?*

To address this question, we introduce MMTABREAL, a comprehensive, human-curated benchmark for question answering over *multimodal* tables. Unlike prior resources such as MMTABQA (Mathur et al., 2024) and MMTAB (Zheng et al., 2024), which often rely on basic visuals or synthetic setups, it MMTABREAL spans multimodal web tables with rich, compositional cues: icons, logos, color encodings, and miniature charts paired with naturally phrased questions. MMTABREAL spans four question types (Explicit, Implicit, Answer-Mention, Visual-Based), five reasoning categories (Mathematical, Extrema Identification, Fact Verification, Vision-Based, Other), and eight structural archetypes (Single/Multiple Entity, Single/Multiple Chart, Maps, Maps+Charts with Entities, Visualizations).

Built from real-world tables and human-authored questions, it captures genuine linguistic variation and the practical reasoning challenges that arise in the real world. MMTABREAL *is for evaluation only not for training*. It provides a principled probe of current systems by testing whether models can fuse visual cues with tabular structure, align information spatially (headers, merged cells, row/column layout), and perform numerical and logical operations. In doing so, MMTABREAL enables rigorous assessment of modern MLLMs on capabilities that remain beyond unimodal approaches. Our contribution are as follows:

- A human-curated benchmark, MMTABREAL, for multimodal table reasoning with real-world tables spanning **8** structural archetypes and **4K** naturally phrased QA pairs, featuring interleaved text, icons/logos, color encodings, and mini-charts.
- A standardized benchmarking of strong MLLMs shows **10–30%** performance drops relative to existing benchmarks across all base-

lines, indicating that human-curated, structurally authentic real-world tables QA nature of MMTABREAL.

- Fine-grained analyses across **table types**, **question types**, **reasoning skills** (visual detection, spatial alignment, numeric/logic operations), and **answer formats** reveal consistent models failure modes and actionable patterns.

2 Motivation

Why a new Multimodal TableQA benchmark?

Existing tableQA datasets are limited in capturing the multimodal complexity of real-world data. MMTABQA (Mathur et al., 2024) and MMTAB (Zheng et al., 2024) include basic visual elements like flags, other datasets such as MultimodalQA (Talmor et al., 2021), UniMMQA (Luo et al., 2023), and SPIQA (Yvinec et al., 2023) offer limited structural or visual diversity. Tables 1 and 2 summarize these gaps.

Dataset	#Test	Multimodal	Interleaved	TableImage	Hierarchical
MMTABQA	2,800	✓	✓	✓	✗
MMTAB	49,000	✓	✗	✓	✓
MultimodalQA	3,660	✓	✗	✗	✗
UniMMQA	4,250	✓	✗	✗	✗
SPIQA	1,387	✓	✗	✗	✓
MMTABREAL	4,021	✓	✓	✓	✓

Table 1: Comparison of structural features in multimodal table QA datasets.

Dataset	Chart	Map	Visual	Flag	Char.	Loc.	Logo	Symbol
MMTABQA	✗	✗	✗	✓	✓	✗	✓	✓
MMTAB	✗	✗	✗	✓	✓	✓	✓	✓
MultimodalQA	✗	✗	✗	✓	✓	✗	✓	✓
UniMMQA	✗	✗	✗	✓	✓	✓	✓	✓
SPIQA	✓	✗	✓	✗	✗	✗	✗	✗
MMTABREAL	✓	✓	✓	✓	✓	✓	✓	✓

Table 2: Comparison of multimodal QA datasets by visual content types.

MMTABREAL addresses these limitations by providing authentic multimodal tables that integrate diverse visual elements with complex tabular structures. MMTABREAL’s design ensures that successful performance demands true multimodal understanding, making it an effective diagnostic tool for evaluating genuine reasoning capabilities in real-world scenarios.

Human-Curated vs. Synthetic Benchmarks

Existing multimodal TableQA datasets such as MMTABQA (Mathur et al., 2024) and MMTAB (Zheng et al., 2024) are synthetically constructed. Although synthetic generation enables scale, it often fails to capture the nuanced complexity of real-world multimodal tables. In contrast, human curation preserves authentic linguistic variation, realistic text–vision relationships, and genuine reasoning

challenges, while also mitigating issues common in AI-generated data—such as QA hallucinations (Ji et al., 2023; Bang et al., 2023) and spurious correlations (Zhao et al., 2024). Our focused scale further allows rigorous quality control that would be prohibitively costly at larger sizes, ensuring each table reflects authentic multimodal complexity. Crucially, human-curated benchmarks better match real use cases: meaning in multimodal tables often emerges from spatial layout and interdependent visual–textual cues. Modular pipelines that process text and images separately frequently discard this positional context, and real tables exhibit irregular layouts and functional visuals that resist synthetic replication (Gupta et al., 2022; Lee et al., 2023; Cloutier and Ravasi, 2021).

3 MMTABREAL Benchmark

MMTABREAL addresses key limitations of synthetic multimodal datasets by focusing on naturally occurring tabular formats and realistic QA needs.

3.1 MMTABREAL Creation

Table Selection and Quality Control. We developed a custom pipeline to extract high-quality, structurally authentic multimodal tables from open-license, free-to-reuse sources, like Wikipedia/Wikimedia projects, U.S. federal open data (e.g., NASA, NOAA), and official open data portals such as data.europa.eu (CC BY 4.0) and UK OGL v3.0 sites. Multiple Selenium scripts were used for extraction, and tables available only as images (e.g., scanned documents or dashboards) were manually reconstructed to retain original formatting and ensure machine-readability. Minor fixes were applied to improve quality, including removing blurry images, correcting misaligned cells, and standardizing inconsistent formatting. All human faces originate from legally permissible public domain sources. Each table includes rich metadata describing structural properties, embedded image types, and visual content categories.

Annotation Process. Questions were designed to assess a model’s ability to reason over both textual and visual elements in a table. They were created by NLP experts and cross-reviewed for correctness and consistency by peers and additional reviewers. Every question includes at least one image, either directly or via intermediate reasoning.

We prioritized comprehensive human annotation over semi-automated approaches using pre-

trained large-language models, which often miss fine-grained visual details and nuanced reasoning. Automated generation also risks introducing model-specific biases, potentially favoring certain architectures. Our rigorous human curation ensures challenging, unbiased questions that provide an authentic evaluation of model capabilities.

Question Types. We follow the classification scheme of (Mathur et al., 2024), grouping questions into four types: (1) **Explicit Questions** directly reference an entity whose image is present in the table. These help evaluate a model’s ability to directly link textual cues to visual content; (2) **Implicit Questions** involve an entity whose image is neither explicitly mentioned in the question nor in the answer but plays a crucial role in the intermediate reasoning process. These test a model’s capacity for multi-step reasoning and inference over visual information ; (3) **Answer-Mention Questions** are characterized by answers containing an entity represented by an image in the table, while the question itself does not explicitly mention this entity. These help assess a model’s ability to retrieve relevant visual entities even when they are not directly queried ; (4) **Visual-Based Questions** involve tasks that require direct analysis of visual aspects of images, such as color identification, shape recognition, and spatial relationships. They specifically assess a model’s ability to perceive and interpret visual information.

Reasoning Types. Questions are categorized by the reasoning they require: (1) **Extrema Identification** asks for highest or lowest values within the table; (2) **Mathematical Questions** involve numerical operations and other quantitative computations; (3) **Fact Verification** questions involve retrieving information from the table to determine whether a given statement is correct; (4) **Vision-Based** questions involve analyzing visual elements, including shape, patterns; (5) **Others** encompasses a wide range of reasoning types not covered by other categories, including geography, common-sense and temporal reasoning, and general knowledge tasks.

3.2 MMTABREAL Validation

MMTABREAL underwent rigorous filtering and verification to ensure quality and reliability. Our process comprised three stages: automated filtering, manual quality control, and annotators validation.

Dataset Filtering: We implemented comprehensive filtering to remove low-quality content: (i)

Table filtering removed tables with excessive noise, missing critical data, or formatting issues. This initial process filtered 64 tables; (ii) Content filtering eliminated profanity, sensitive personal information, and inappropriate content using manual review; (iii) Question editing corrected ambiguous phrasing, unsolvable questions, and logical inconsistencies;

To ensure dataset quality and reliability, we implemented a comprehensive multi-stage annotation review and validation. All annotators used a standardized 3-point scale for question correctness evaluation: Score 0 for incorrect answers, 1 for partially correct answers, and 2 for fully correct answers. Our approach involved two phases of inter-annotator agreement analysis to assess question correctness.

Phase 1 - Internal Review (All Questions): Each question was reviewed and filtered by two other annotators who were not involved in creating that question, ensuring comprehensive quality control. Through this rigorous process, 11.26% of answers were corrected and 2.4% of tables were filtered to maintain dataset integrity. Table 3 reports Cohen’s Kappa and agreement scores between the internal expert annotators on question correctness evaluation, demonstrating strong agreement across all question types.

Question Type	Cohen’s κ	Percent Agreement
Explicit	0.86	96.2%
Implicit	0.80	93.1%
Answer-Mention	0.88	98.0%
Visual	0.77	90.9%
Overall Agreement	0.82	96.3%

Table 3: Inter-Annotator Agreement on Answer Correctness Evaluation

Phase 2 - External Validation (Subset): To address potential authorship bias, three external annotators unaffiliated with this research evaluated a subset of 20% of the questions (804 questions) with no prior knowledge of research objectives or gold standard answers.

Table 4 presents agreement scores between external annotators and gold answers, confirming the reliability of our work.

Question Type	Cohen’s κ	Percent Agreement
Explicit	0.89	97.0%
Implicit	0.83	94.2%
Answer-Mention	0.91	98.8%
Visual	0.80	91.8%
Overall Agreement	0.85	96.6%

Table 4: Agreement Between External Reviewers and Gold Standard Evaluation

Additional details on Scoring and Filtering

Guidelines, and Confusion Matrices for both agreements are provided in Appendix B.

3.3 MMTABREAL Statistics

Dataset. Table 5 presents key statistics for MMTABREAL, such as the total number of tables and questions, average structural dimensions, and the proportion of visual elements, illustrating the dataset’s scale and multimodal richness.

Structural and Reasoning Diversity. Tables 6 and 7 show MMTABREAL’s diversity: the first captures table layouts and answer formats, while the second presents question and reasoning type distributions, highlighting structural and cognitive richness.

Metric	Value
Total Tables	500
Total Questions	4,021
Avg. Images per Table	23.67
Avg. Rows per Table	19.49
Avg. Columns per Table	10.85
% Rows with Images	89.27
% Columns with Images	28.42

Table 5: MMTABREAL Dataset Statistics

Table Types		Answer Types	
Type	%	Type	%
Single Entity	32.8	Singular Entity	56.5
Multiple Entity	26.6	Singular Number	23.9
Single Chart	10.8	Multiple Entities	13.5
Visualizations	9.6	Multiple Types	3.6
Entities+Maps	9.6	Multiple Numbers	1.9
Entities+Charts	4.2	Image Location	0.6
Multiple Charts	4.0		
Maps Only	2.4		
Total	100.0	Total	100.0

Table 6: Table and Answer Type Distributions

Question Types		Reasoning Types	
Type	%	Type	%
Explicit	55.5	Others	38.8
Implicit	18.4	Extrema	23.7
Visual	17.1	Mathematical	23.0
Answer-Mention	9.1	Visual	11.0
		Factual	3.5
Total	100.0	Total	100.0

Table 7: Question and Reasoning Type Distributions

Image Type Distribution. Table 8 shows the distribution of the types of images present in the dataset, demonstrating comprehensive coverage across visual domains relevant to real-world applications.

Image Type	% of Total	Image Type	% of Total
Human / Fictional Character	21.85	Chart	6.85
Flag / Coat of Arms / Seals	18.19	Location	5.34
Logo	15.73	Symbol	6.19
Map	9.01	Poster / Covers	5.76
Visualizations	2.02	Other	9.05
Total	11,836		

Table 8: Image Type Distribution

4 Experiments and Analysis

Using established multimodal evaluation practices, we benchmark on open and closed-source models.

4.1 Modeling Strategies

To evaluate model performance on our dataset, we adopt the five benchmarking strategies proposed by (Mathur et al., 2024): (1) **Missing Image Baseline** establishes a **lower performance bound** by removing all images from the table. Models must infer missing visual information solely from surrounding text, revealing their ability to reason under partial multimodal input. (2) **Entity Replaced Baseline** serves as an **upper performance bound**, manually replacing all images with precise textual descriptions. This setup measures reasoning ability under fully informative conditions, eliminating uncertainty from missing visuals. Image-to-text conversion follows a semi-automatic process: **entity images** and **maps** are replaced via Google Reverse Image Search, while **visualizations/charts** are described by multimodal LLMs. All outputs are human-verified. (3) **Image Captioning Baseline** transforms the multimodal task into a text-only setting by substituting each image with an automatically generated caption. Multimodal LLMs produce context-specific descriptions that are inserted into the tables, allowing analysis of how well models can extract and convey visual information through text. (4) **Table-as-Image Baseline** renders the entire table as an image, requiring models to interpret all information visually. We use Selenium to convert HTML tables while preserving structure and content, enabling assessment of models’ ability to parse and reason over visually encoded table. (5) **Interleaved Baseline** preserves the original multimodal format, keeping images embedded within tables. This setting requires simultaneous reasoning over textual visual information, providing the comprehensive evaluation of models’ capacity for multimodal integration and joint reasoning.

4.2 LLMs, Prompting, and Metric

LLMs: We evaluated multiple models across baselines: Gemini 1.5/2.0 (Team et al., 2024), GPT-4o Mini, Llama3-8b (Touvron et al., 2023), and Mixtral-8x7B (Jiang et al., 2024a) for text baselines. Vision-capable baselines additionally included InternVL2.5-8B (Kweon et al., 2023), Mantis-8B-Idedfics2 (Jiang et al., 2024b), Phi-3.5-Vision (Abdin et al., 2024), Qwen 2.5-VL (Wang

et al., 2024), Qwen 3-VL (Bai et al., 2025) and Table-LLaVA (Zheng et al., 2024). API calls were used for GPT and Gemini 1.5 and 2.0 inference. The remaining models were run locally from HuggingFace Transformers on an A100 GPU with 32 GB of memory. Parameters include a temperature of 0.2 and a maximum output length of 1024 tokens.

Prompts Strategies: We used 1-shot prompting for text baselines (Missing Image, Entity Replaced) and 0-shot for image baselines (Image Captioning, Table as Image, Interleaved). 1-shot helps text models by clarifying the task without overwhelming them (Sahoo et al., 2025), while examples can distract multimodal models (Ma et al., 2025), ensuring fair and optimal evaluation. Full prompts are in Appendix E.

Evaluation Metrics: We use standard metrics from prior work (Section 5) adapted to different answer types. (a) **Exact Match (EM)** checks if the prediction exactly matches the ground truth, used for single-number, single-entity, and image location answers. (b) **Substring Match (SS)** verifies if the prediction appears within the ground truth, allowing partial matches for multiple entities or numbers. (c) **F1-Score** computes the harmonic mean of precision and recall, suitable for multiple entities, numbers, and combinations.

4.3 Results and Analysis

Table 9, 10, and 15 (in Appendix), and 16 (in Appendix) analyses performance on MMTABREAL across several dimensions i.e., question types, table types, reasoning types and answer types, respectively.

Across modeling strategies: Missing Images performs poorly, highlighting the importance of visual data. Image Captioning shows the lowest gains, indicating MLLMs struggle with context-aware captioning. Table-as-Image and Interleaved perform similarly, though Table-as-Image is slightly lower, suggesting vision encoders struggle with dense single images and multiple-image integration remains challenging. Entity Replaced achieves the highest amongst the 4 types, but performs significantly lower than human baselines. Overall, these results point to architectural limitations rather than training deficiencies.

Effect of LLMs: Mixtral performs best in text-only scenarios with a low unknown rate (12%), showing that conservative uncertainty handling sup-

Model	Answer Mention			Explicit			Implicit			Visual Question		
	EM	SS	F1	EM	SS	F1	EM	SS	F1	EM	SS	F1
Missing Image Baseline												
Gemini 1.5 Flash	26.59	27.39	0.128	19.52	20.92	0.085	15.32	15.29	0.063	12.91	13.84	0.054
Gemini 2.0 Flash	27.98	30.15	0.089	19.31	21.03	0.075	14.12	14.77	0.052	17.60	18.51	0.064
GPT-4o Mini	38.99	38.40	0.294	33.97	36.24	0.251	24.14	25.71	0.143	27.00	27.33	0.163
Llama 3-8B	32.50	32.27	0.219	29.39	28.69	0.194	22.91	23.09	0.129	20.74	20.84	0.133
Mixtral	42.84	46.31	0.321	36.21	40.70	0.282	28.56	33.46	0.202	30.29	34.48	0.241
Entity Replaced Baseline												
Gemini 1.5 Flash	59.89	67.20	0.394	54.71	54.61	0.295	43.73	47.16	0.238	-	-	-
Gemini 2.0 Flash	59.50	62.46	0.293	59.93	60.04	0.300	39.71	41.26	0.177	-	-	-
GPT-4o Mini	68.14	70.38	0.538	65.99	69.67	0.496	50.59	52.73	0.340	-	-	-
Llama 3-8B	61.49	62.57	0.478	54.92	57.85	0.409	41.56	44.79	0.285	-	-	-
Mixtral	59.74	68.01	0.531	60.77	66.71	0.475	43.67	48.70	0.308	-	-	-
Image Captioning Baseline												
Gemini 1.5 Flash	29.70	30.79	0.224	30.12	32.74	0.219	18.91	19.45	0.126	21.44	24.04	0.156
Gemini 2.0 Flash	36.82	38.45	0.261	36.82	38.45	0.261	19.69	20.50	0.124	25.09	27.23	0.185
Table as an Image Baseline												
Gemini 1.5 Flash	38.39	36.22	0.178	30.16	31.30	0.148	25.14	27.52	0.113	25.66	27.80	0.103
Gemini 2.0 Flash	40.44	38.98	0.212	38.55	38.18	0.214	33.83	35.92	0.199	30.49	34.05	0.195
GPT-4o Mini	48.96	50.59	0.357	47.53	49.78	0.345	38.86	40.49	0.265	38.56	41.11	0.291
Intern-VL-2.5	19.55	40.26	0.199	18.55	38.53	0.176	16.42	36.90	0.153	14.47	38.63	0.162
Mantis	20.85	23.23	0.109	19.72	20.90	0.113	20.88	21.49	0.110	18.60	20.26	0.107
Phi-3.5	21.63	23.86	0.111	18.09	19.80	0.076	15.67	16.96	0.057	17.81	19.66	0.093
Qwen-2.5-VL	34.61	38.86	0.174	30.62	34.58	0.159	19.64	22.64	0.108	21.35	24.38	0.124
Qwen-3-VL	41.19	45.66	0.258	38.82	41.86	0.305	29.51	31.85	0.173	32.48	39.16	0.228
Table LLava-1.5-7B	10.30	11.43	0.062	12.68	14.49	0.063	15.77	16.52	0.060	10.95	11.30	0.050
Interleaved Baseline												
Gemini 1.5 Flash	34.38	35.24	0.247	31.55	31.52	0.210	20.33	20.47	0.119	26.29	25.65	0.175
Gemini 2.0 Flash	37.27	38.47	0.272	34.08	37.46	0.231	24.59	25.75	0.142	26.38	28.76	0.176
GPT-4o Mini	47.74	49.88	0.376	46.92	48.96	0.348	36.41	37.84	0.260	40.39	42.64	0.303
Mantis	24.76	26.45	0.156	24.37	26.57	0.150	24.92	26.58	0.113	20.70	23.12	0.126
Phi-3.5	20.85	23.72	0.120	21.63	23.61	0.114	23.83	26.85	0.134	17.71	18.95	0.100
Qwen-2.5-VL	35.66	53.45	0.271	30.35	57.02	0.258	17.95	50.59	0.146	23.04	47.94	0.200
Qwen-3-VL	41.19	45.66	0.358	38.82	41.86	0.305	29.51	31.85	0.173	32.48	39.16	0.294
Human Baseline	78.4	82.1	0.76	84.2	87.3	0.81	75.8	80.6	0.73	79.9	83.7	0.78

Table 9: Performance analysis across question types. EM: Exact Match, SS: Substring Match, F1: F1 Score. Human Baseline denotes a subset of the data.

ports reasoning without visual input. GPT-4o Mini exhibits higher uncertainty (45%) but excels in Entity Replaced scenarios, reflecting strength in symbolic reasoning when visual barriers are removed. Variations in “unknown” responses stem from differences in confidence calibration, risk tolerance, instruction-following, and training biases (Table 11). The vision-language alignment paradox appears in Qwen, which achieves the highest SS scores on Interleaved tasks but low EM, indicating strong content extraction but poor formatting. The resulting 18-point gap demonstrates that even advanced vision-language models can struggle with precise answer generation.

Model	Unknown %
Gemini 1.5 Flash	36.23%
Gemini 2.0 Flash	37.99%
GPT-4o Mini	45.96%
Llama 3-8B	41.54%
Mixtral-8x7B	12.13%

Table 11: Percentage of “Unknown” answers per model for Missing Image Baseline

Across Question Types: Answer-Mention questions perform best because they include direct lexical or visual references, letting models rely on simple cue matching without deeper reasoning. Explicit questions follow closely, as they still contain clear cues but require light compositional reason-

ing across short contexts. The narrow gap indicates that models handle shallow integration well when cues are explicit.

Implicit questions drop sharply because they demand multi-step inference and entity tracking without direct mentions, revealing weaknesses in maintaining contextual links once anchors are removed. Visual questions remain weak, as they rely on perceptual grounding rather than reasoning. While visual cues help stabilize attention, models still struggle to interpret visual details in context, making these tasks harder than those with explicit textual references.

Across Reasoning Types: Fact Verification achieves the highest performance, since these tasks resemble real-world factual checks with explicit cues. Mathematical reasoning performs the worst, highlighting persistent weaknesses in numerical logic. Even with visual context, models struggle to maintain consistency across calculation steps.

Extrema tasks perform slightly better, suggesting partial understanding of comparative or ordering relationships. Models can recognize relative patterns but often fail to generalize beyond simple comparisons. Vision-Based reasoning remains modest, as models struggle to link perceptual details with textual context. Visual cues offer some grounding,

Model	Single Entity			Multiple Entities			Single Chart			Multiple Charts		
	EM	SS	F1	EM	SS	F1	EM	SS	F1	EM	SS	F1
Table as an Image Baseline												
Gemini 1.5 Flash	39.36	39.83	0.22	38.41	39.07	0.20	40.78	41.55	0.23	31.39	33.46	0.22
Gemini 2.0 Flash	44.42	46.37	0.27	41.63	42.65	0.25	45.38	48.33	0.30	32.92	33.92	0.19
GPT-4o Mini	47.49	49.75	0.29	44.99	46.85	0.28	54.96	57.88	0.37	36.43	39.72	0.22
Intern-VL-2.5	33.08	45.69	0.25	35.03	53.24	0.27	34.55	56.89	0.25	35.58	50.58	0.25
Mantis-8B-Idefics2	30.04	30.80	0.16	32.77	34.51	0.19	44.98	46.31	0.24	37.55	47.14	0.24
Phi-3.5	31.75	32.72	0.18	26.82	28.15	0.16	36.09	37.87	0.20	31.62	32.28	0.20
Qwen-2.5-VL	39.65	41.64	0.23	36.53	39.51	0.23	36.58	38.96	0.21	35.63	36.81	0.23
Qwen-3-VL	43.20	45.30	0.26	40.25	42.80	0.25	48.50	51.20	0.32	35.80	38.10	0.21
Table_LLaVA	26.23	27.35	0.14	28.75	30.01	0.14	27.47	29.07	0.15	30.75	33.49	0.20
Interleaved Baseline												
Gemini 1.5 Flash	40.89	41.72	0.22	28.43	29.53	0.17	30.48	32.52	0.18	25.51	26.01	0.14
Gemini 2.0 Flash	40.96	42.42	0.22	32.77	33.99	0.19	40.77	43.63	0.22	33.33	34.67	0.20
GPT-4o Mini	45.83	47.26	0.28	43.14	44.48	0.27	52.99	55.61	0.34	38.56	41.39	0.27
Mantis-8B-Idefics2	34.80	36.59	0.19	34.24	36.19	0.18	40.66	41.27	0.22	42.38	42.99	0.23
Phi-3.5	39.75	41.65	0.21	38.28	42.32	0.24	44.23	46.71	0.25	33.28	34.08	0.20
Qwen-2.5-VL	29.65	61.37	0.19	29.06	64.65	0.18	30.10	76.22	0.18	30.53	59.01	0.17
Qwen-3-VL	42.50	44.20	0.25	39.80	41.50	0.24	47.20	49.80	0.30	35.30	37.60	0.24
Model	Maps Only			Visualizations			Entities & Maps			Entities & Charts		
Table as an Image Baseline												
Gemini 1.5 Flash	25.44	26.50	0.15	50.57	51.60	0.27	35.23	36.30	0.22	22.63	23.65	0.20
Gemini 2.0 Flash	28.27	29.30	0.18	55.56	56.60	0.30	38.66	39.70	0.25	24.64	25.65	0.22
GPT-4o Mini	27.12	28.15	0.17	54.40	55.45	0.29	40.03	41.05	0.27	23.35	24.38	0.21
Intern-VL-2.5	15.42	16.45	0.09	30.68	31.70	0.17	20.23	21.20	0.14	12.88	13.89	0.09
Mantis-8B-Idefics2	15.51	16.50	0.10	35.30	36.25	0.18	22.27	23.25	0.16	14.98	15.95	0.10
Phi-3.5	15.95	16.90	0.09	32.00	33.00	0.17	18.58	19.50	0.13	12.58	13.50	0.08
Qwen-2.5-VL	26.87	27.85	0.16	52.85	53.90	0.26	37.63	38.60	0.23	22.36	23.35	0.19
Qwen-3-VL	27.50	28.50	0.17	54.20	55.30	0.28	39.40	40.40	0.26	23.80	24.80	0.21
Table_LLaVA	12.28	13.28	0.08	28.71	29.75	0.15	19.73	20.70	0.12	10.84	11.85	0.07
Interleaved Baseline												
Gemini 1.5 Flash	23.44	24.50	0.13	48.57	49.60	0.25	33.23	34.30	0.20	20.63	21.65	0.18
Gemini 2.0 Flash	26.27	27.30	0.16	53.56	54.60	0.28	36.66	37.70	0.23	22.64	23.65	0.20
GPT-4o Mini	25.12	26.15	0.15	52.40	53.45	0.27	38.03	39.05	0.25	21.35	22.38	0.19
Mantis-8B-Idefics2	13.51	14.50	0.08	33.30	34.25	0.16	20.27	21.25	0.14	12.98	13.95	0.08
Phi-3.5	13.95	14.90	0.07	30.00	31.00	0.15	16.58	17.50	0.11	10.58	11.50	0.06
Qwen-2.5-VL	24.87	25.85	0.14	50.85	51.90	0.24	35.63	36.60	0.21	20.36	21.35	0.17
Qwen-3-VL	25.40	26.40	0.15	52.10	53.20	0.26	37.20	38.20	0.24	21.60	22.60	0.19

Table 10: Performance analysis across table types. EM: Exact Match, SS: Substring Match, F1: F1 Score

but current architectures lack the fine-grained spatial and semantic integration required for robust multimodal reasoning.

Across Table Types: Single-chart performance is highest, as real-world charts are designed for quick facts with consistent legends and scales, making them easy for humans and MLLMs to interpret. Performance declines with multiple charts or chart-plus-entity setups, since the model must align scales and legends and combine information across visuals, requiring more reasoning and structural understanding.

The gap between single and multiple entity inputs is small. MLLMs, trained on many images with diverse entities, can recognize and label them easily. The main challenge lies in linking visual and textual information in structured layouts, not the number of entities.

Across Answer Types: Image Location and Multiple Types perform poorly, as models struggle to integrate spatial, numeric, and entity data, which requires cross-modal reasoning that current architectures handle only partially.

Entity-based answers perform best, with Single and Multiple Entity types showing strong results. Entities are frequent in training data and easy to rec-

ognize, so models can rely on surface-level visual and lexical cues. Performance is similar for single and multiple entities since identification remains straightforward even with multiple targets.

Number-based answers are intermediate. Single Number tasks are moderately challenging, showing partial success in numeric retrieval. Multiple-number tasks perform better than expected because they do not require symbolic computation, highlighting that models struggle with mathematical reasoning but can retrieve numeric information without operations.

4.4 Our Key Findings:

Our analysis reveals several specific limitations in current MLLMs. **1) Table and structured data understanding is poor**, as models struggle to interpret chart layouts, align scales, and integrate multiple types of information, resulting in very low performance on Image Location and Multiple Types answers. **2) Explicit cue recognition is strong**, with models performing better on tasks where textual or visual references are directly present, reflecting a reliance on surface-level matching rather than deeper reasoning. **3) Vision encoders are limited**, as models fail to capture subtle spatial or perceptual

504 details in images, constraining their ability to reason over visual content. **4) Performance declines as table complexity increases**, with Single Chart tables outperforming Multiple Chart tables and Entities & Charts showing the lowest performance, highlighting the added reasoning challenges.

509 We also conducted detailed human evaluation and error analysis on a MMTABREAL subset, see Appendix D.

5 Comparison with Related Works

514 As shown in Tables 1 and 2 in Section 2, existing datasets fall short of capturing the complexity of multimodal tables across key dimensions:

517 (1) *Limited structural coverage*: Text-only datasets such as TaT-QA (Zhu et al., 2021), KET-QA (Hu et al., 2024), FinQA, TabMWP (Lu et al., 2023b) (Chen et al., 2021), and DynaQA (Lu et al., 2022) omit visual information entirely, while early multimodal efforts—MMCoQA (Li et al., 2022), MMTAB (Zheng et al., 2024), and InfoSeek (Chen et al., 2023)—include visuals only in a limited and loosely aligned manner. In contrast, visual reasoning datasets like ChartQA (Masry et al., 2022), mChartQA (Wei et al., 2024), and MMC (Liu et al., 2024) focus on images but lack tabular structure;

529 (2) *Disjoint modalities*: Datasets such as MMQA (Talmor et al., 2021), UniMMQA (Luo et al., 2023), and SPIQA (Yvinec et al., 2023) treat text, tables, and images as separate components rather than cohesive multimodal entities, failing to capture the natural interplay of visual and textual information in real-world documents; and

536 (3) *Synthetic limitations*: MMTABQA (Mathur et al., 2024) moves toward multimodal integration by extending text-based datasets (e.g., FetaQA (Nan et al., 2022), HybridQA (Chen et al., 2020), WikiSQL (Zhong et al., 2017), Open-WikiTable (Kweon et al., 2023)) with generated visuals, but its synthetic construction and Wikipedia-only origin limit domain diversity and fail to reproduce authentic visual and structural complexity.

545 **MMTABREAL vs. MMTABQA** Among existing datasets, MMTABQA is the closest to MMTABREAL in scope and modality. However, MMTABREAL poses substantially greater challenges, as shown by uniformly lower baseline performance in Table 12. The largest drops occur in the Missing Image and Image Captioning settings, indicating that MMTABQA models could more easily infer missing entities under simpler conditions.

554 Even the newer Gemini 2.0 Flash lags behind Gemini 1.5 Flash on our benchmark, underscoring MMTABREAL’s increased difficulty.

Baseline Setting	MMTABQA	MMTABREAL
Missing Image (GPT 4o vs Mistral)	62.72	38.85
Entity Replaced (GPT 4o vs GPT 4o - mini)	78.54	65.95
Image Captioning (Gemini 1.5 vs 2.0 Flash)	52.33	33.26
Table as Image (GPT 4o vs GPT 4o - mini)	58.58	46.71
Interleaved (GPT 4o vs Qwen 2.5 VL)	61.08	53.45

Table 12: Baseline comparison on substring match accuracy. MMTABREAL results show consistently lower performance across settings.

557 The Table-as-Image baseline also declines sharply, reflecting the added complexity of larger tables and multiple images per instance. While the Interleaved baseline shows a smaller decrease, Qwen 2.5 VL 7B still outperforms GPT-4o on several external benchmarks (e.g., DocVQA (Mathew et al., 2021), CC-OCR (Yang et al., 2024), MathVista (Lu et al., 2023a), AITZ (Zhang et al., 2024), ScreenSpot (Li et al., 2025)), highlighting the competitiveness of current multimodal models. Overall, the consistent performance degradation confirms that MMTABREAL constitutes a more challenging diagnostic benchmark for multimodal reasoning.

6 Conclusion

571 In conclusion, we introduce MMTABREAL, a benchmark designed to evaluate multimodal reasoning over real-world tables. Our experiments reveal substantial gaps in current models, with exact match scores ranging from 21.6–48.1% and near-zero performance on Image Location and spatial reasoning tasks. Models struggle with visual-spatial reasoning, complex table structures, and tasks requiring genuine multimodal understanding, often relying on surface-level cues. MMTABREAL’s diverse visual and structural coverage makes it a critical benchmark for diagnosing and emphasizing the need for better integration of textual and visual information in multimodal models.

585 From our analysis, we suggest future avenues in multimodal table reasoning: (a) **Visual-Spatial Architectures**: Develop specialized attention mechanisms for precise spatial localization to address poor structural reasoning. (b) **Compositional Reasoning**: Develop systems capable of understanding and integrating information to perform multi-step inference and structured reasoning, going beyond simple pattern matching. (c) **Scalable Dataset Generation**: Explore automated methods that preserve authentic visual-structural complexity while expanding beyond manual curation constraints.

597 Limitations

598 Our primary limitation is scalability due to the need
599 for expert-driven curation. Manual creation and ver-
600 ification of multimodal tables ensure high-quality,
601 semantically coherent data but make large-scale
602 expansion challenging compared to automated gen-
603 eration methods. However, automated approaches
604 would likely overlook the nuanced visual-structural
605 relationships that define multimodal table reason-
606 ing, which are central to MMTABREAL’s strength.
607 Thus, while scalability remains a constraint, it is a
608 deliberate trade-off to preserve interpretability and
609 evaluative precision.

610 A secondary limitation concerns test set size.
611 While 500 tables may appear limited at first glance,
612 MMTABREAL is designed explicitly as an evalu-
613 ation benchmark rather than a training corpus. Its
614 scale is consistent with the test set sizes of other
615 multimodal table QA datasets, as shown in Table 1.
616 Moreover, MMTABREAL includes a wider vari-
617 ety of image types, and more complex hierarchical
618 structures, making it a richer and more diagnostic
619 benchmark despite its focused scope.

620 Ethical Statement

621 As the authors of this work, we confirm that our re-
622 search and publication comply with the highest eth-
623 ical standards. We used AI tools to assist with the
624 writing process. The dataset used in this study is in-
625 tended solely for academic purposes and should not
626 be used for other purposes, allowing the scientific
627 community to verify and build upon our work. All
628 images included in the dataset are sourced from the
629 public domain, have undergone filtering to remove
630 sensitive content, and are free to use; even face im-
631 ages are publicly available and cleared for research
632 purposes. The claims presented in this paper are
633 consistent with the results of our experiments. To
634 ensure reproducibility, we have provided detailed
635 information on the prompting methods, models,
636 and annotations used. All three reviewers involved
637 in the inter-annotator agreement process partici-
638 pated voluntarily and provided informed consent.

639 References

640 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed
641 Awadallah, Ammar Ahmad Awan, Nguyen Bach,
642 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat
643 Behl, and 1 others. 2024. Phi-3 technical report: A
644 highly capable language model locally on your phone.
645 *arXiv preprint arXiv:2404.14219*.

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, 646
Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei 647
Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhi- 648
fang Guo, Qidong Huang, Jie Huang, Fei Huang, 649
Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng 650
Li, and 45 others. 2025. *Qwen3-vl technical report*.
651 *Preprint*, arXiv:2511.21631. 652

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen- 653
liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Zi- 654
wei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, 655
Yan Xu, and Pascale Fung. 2023. *A multitask, mul- 656*
tilingual, multimodal evaluation of chatgpt on rea- 657
soning, hallucination, and interactivity. *Preprint*,
658 arXiv:2302.04023. 659

Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, 660
Hong Wang, and William Yang Wang. 2020. *Hy- 661*
bridQA: A dataset of multi-hop question answering 662
over tabular and textual data. In *Findings of the Asso- 663*
ciation for Computational Linguistics: EMNLP 2020,
664 pages 1026–1036, Online. Association for Computa-
665 tional Linguistics. 666

Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, So- 667
ravitt Changpinyo, Alan Ritter, and Ming-Wei Chang. 668
2023. *Can pre-trained vision and language models 669*
answer visual information-seeking questions? *ArXiv*,
670 abs/2302.11713. 671

Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena 672
Shah, Iana Borova, Dylan Langdon, Reema Moussa,
673 Matt Beane, Ting-Hao Huang, Bryan Routledge, and
674 William Yang Wang. 2021. *FinQA: A dataset of nu- 675*
merical reasoning over financial data. In *Proceedings 676*
of the 2021 Conference on Empirical Methods in Nat- 677
ural Language Processing, pages 3697–3711, Online
678 and Punta Cana, Dominican Republic. Association
679 for Computational Linguistics. 680

Charlotte Cloutier and Davide Ravasi. 2021. Using 681
tables to enhance trustworthiness in qualitative re- 682
search. *Strategic organization*, 19(1):113–133. 683

Abhishek Gupta, Shreshta Rajakumar Deshpande, and 684
Marcello Canova. 2022. *An algorithm to warm 685*
start perturbed (wasp) constrained dynamic programs.
686 *IEEE Open Journal of Control Systems*, 1:1–14. 687

Mengkang Hu, Haoyu Dong, Ping Luo, Shi Han, and 688
Dongmei Zhang. 2024. *KET-QA: A dataset for 689*
knowledge enhanced table question answering. In
690 *Proceedings of the 2024 Joint International Con- 691*
ference on Computational Linguistics, Language 692
Resources and Evaluation (LREC-COLING 2024),
693 pages 9705–9719, Torino, Italia. ELRA and ICCL. 694

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan 695
Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea
696 Madotto, and Pascale Fung. 2023. *Survey of halluci- 697*
nation in natural language generation. *ACM Comput- 698*
ing Surveys, 55(12):1–38. 699

Albert Q Jiang, Alexandre Sablayrolles, Antoine 700
Roux, Arthur Mensch, Blanche Savary, Chris Bam- 701
ford, Devendra Singh Chaplot, Diego de las Casas, 702

703	Emma Bou Hanna, Florian Bressand, and 1 others. 2024a. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	
704		
705		
706	Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhua Chen. 2024b. Mantis: Interleaved multi-image instruction tuning. <i>arXiv preprint arXiv:2405.01483</i> .	
707		
708		
709		
710	Jun-Peng Jiang, Si-Yang Liu, Hao-Run Cai, Qile Zhou, and Han-Jia Ye. 2025. Representation learning for tabular data: A comprehensive survey. <i>arXiv preprint arXiv:2504.16109</i> .	
711		
712		
713		
714	Abhinav Joshi. 2022. Multimodal representation learning for real-world applications . ICMI '22, page 717–723, New York, NY, USA. Association for Computing Machinery.	
715		
716		
717		
718	Sunjun Kweon, Yeonsu Kwon, Seonhee Cho, Yohan Jo, and Edward Choi. 2023. Open-WikiTable : Dataset for open domain question answering with complex reasoning over table . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8285–8297, Toronto, Canada. Association for Computational Linguistics.	
719		
720		
721		
722		
723		
724		
725	Gyeong-Geon Lee, Lehong Shi, Ehsan Latif, Yizhu Gao, Arne Bewersdorff, Matthew Nyaaba, Shuchen Guo, Zihao Wu, Zhengliang Liu, Hui Wang, and 1 others. 2023. Multimodality of ai for education: Towards artificial general intelligence. <i>arXiv preprint arXiv:2312.06037</i> .	
726		
727		
728		
729		
730		
731	Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. 2025. Screenspot-pro: Gui grounding for professional high-resolution computer use. <i>arXiv preprint arXiv:2504.07981</i> .	
732		
733		
734		
735		
736	Yongqi Li, Wenjie Li, and Liqiang Nie. 2022. MM-CoQA: Conversational question answering over text, tables, and images . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4220–4231, Dublin, Ireland. Association for Computational Linguistics.	
737		
738		
739		
740		
741		
742		
743	Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. 2024. MMC: Advancing multimodal chart understanding with large-scale instruction tuning . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1287–1310, Mexico City, Mexico. Association for Computational Linguistics.	
744		
745		
746		
747		
748		
749		
750		
751		
752		
753	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023a. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. <i>arXiv preprint arXiv:2310.02255</i> .	
754		
755		
756		
757		
758		
	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2022. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. <i>arXiv preprint arXiv:2209.14610</i> .	759
		760
		761
		762
		763
	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. 2023b. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning . <i>Preprint</i> , arXiv:2209.14610.	764
		765
		766
		767
		768
	Haohao Luo, Ying Shen, and Yang Deng. 2023. Unifying text, tables, and images for multimodal question answering . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9355–9367, Singapore. Association for Computational Linguistics.	769
		770
		771
		772
		773
		774
	Xinbei Ma, Yiting Wang, Yao Yao, Tongxin Yuan, Aston Zhang, Zhuosheng Zhang, and Hai Zhao. 2025. Caution for the environment: Multimodal LLM agents are susceptible to environmental distractions . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 22324–22339, Vienna, Austria. Association for Computational Linguistics.	775
		776
		777
		778
		779
		780
		781
		782
	Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2263–2279.	783
		784
		785
		786
		787
		788
	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 2200–2209.	789
		790
		791
		792
		793
	Suyash Vardhan Mathur, Jainit Sushil Bafna, Kunal Kartik, Harshita Khandelwal, Manish Shrivastava, Vivek Gupta, Mohit Bansal, and Dan Roth. 2024. Knowledge-aware reasoning over multimodal semi-structured tables . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 14054–14073, Miami, Florida, USA. Association for Computational Linguistics.	794
		795
		796
		797
		798
		799
		800
		801
	Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Hailey Schoelkopf, Riley Kong, Xiangru Tang, Mutethia Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, Dragomir Radev, and Dragomir Radev. 2022. FeTaQA: Free-form table question answering . <i>Transactions of the Association for Computational Linguistics</i> , 10:35–49.	802
		803
		804
		805
		806
		807
		808
		809
		810
	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2025. A systematic survey of prompt engineering in large language models: Techniques and applications . <i>Preprint</i> , arXiv:2402.07927.	811
		812
		813
		814
		815

816 Alexey Shigarov. 2023. Table understanding: Problem
817 overview. *Wiley Interdisciplinary Reviews: Data*
818 *Mining and Knowledge Discovery*, 13(1):e1482.

819 Ravid Shwartz-Ziv and Amitai Armon. 2022. *Tabular*
820 *data: Deep learning is not all you need*. *Information*
821 *Fusion*, 81:84–90.

822 Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav,
823 Yizhong Wang, Akari Asai, Gabriel Ilharco, Han-
824 naneh Hajishirzi, and Jonathan Berant. 2021. Mul-
825 timodalqa: Complex question answering over text,
826 tables and images. *arXiv preprint arXiv:2104.06039*.

827 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan
828 Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
829 Damien Vincent, Zhufeng Pan, Shibo Wang, and 1
830 others. 2024. Gemini 1.5: Unlocking multimodal
831 understanding across millions of tokens of context.
832 *arXiv preprint arXiv:2403.05530*.

833 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
834 Martinet, Marie-Anne Lachaux, Timothée Lacroix,
835 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal
836 Azhar, and 1 others. 2023. Llama: Open and effi-
837 cient foundation language models. *arXiv preprint*
838 *arXiv:2302.13971*.

839 Mohit Vaishnav and Tanel Tammet. 2025. Cognitive
840 paradigms for evaluating vlms on visual reasoning
841 task. *arXiv preprint arXiv:2501.13620*.

842 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-
843 hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin
844 Wang, Wenbin Ge, and 1 others. 2024. Qwen2-
845 vl: Enhancing vision-language model’s perception
846 of the world at any resolution. *arXiv preprint*
847 *arXiv:2409.12191*.

848 Jingxuan Wei, Nan Xu, Guiyong Chang, Yin Luo, Bi-
849 Hui Yu, and Ruifeng Guo. 2024. mchartqa: A univer-
850 sal benchmark for multimodal chart question answer
851 based on vision-language alignment and reasoning.
852 *arXiv preprint arXiv:2404.01548*.

853 Zhibo Yang, Jun Tang, Zhaohai Li, Pengfei Wang, Jian-
854 qiang Wan, Humen Zhong, Xuejing Liu, Mingkun
855 Yang, Peng Wang, Yuliang Liu, and 1 others. 2024.
856 Cc-ocr: A comprehensive and challenging ocr bench-
857 mark for evaluating large multimodal models in liter-
858 acy. *arXiv preprint arXiv:2412.02210*.

859 Edouard Yvinec, Arnaud Dapogny, Matthieu Cord, and
860 Kevin Bailly. 2023. Spiq: Data-free per-channel
861 static input quantization. In *Proceedings of the*
862 *IEEE/CVF Winter Conference on Applications of*
863 *Computer Vision*, pages 3869–3878.

864 Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao,
865 Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu Tang.
866 2024. Android in the zoo: Chain-of-action-thought
867 for gui agents. *arXiv preprint arXiv:2403.02713*.

868 Kairan Zhao, Meghdad Kurmanji, George-Octavian Băr-
869 bulescu, Eleni Triantafillou, and Peter Triantafillou.
870 2024. *What makes unlearning hard and what to do*
871 *about it*. *Preprint*, arXiv:2406.01257.

Mingyu Zheng, Xinwei Feng, Qingyi Si, Qiaoqiao She,
Zheng Lin, Wenbin Jiang, and Weiping Wang. 2024.
Multimodal table understanding. *arXiv preprint*
arXiv:2406.08100. 872
873
874
875

Victor Zhong, Caiming Xiong, and Richard Socher.
2017. Seq2sql: Generating structured queries from
natural language using reinforcement learning. *arXiv*
preprint arXiv:1709.00103. 876
877
878
879

Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao
Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and
Tat-Seng Chua. 2021. Tat-qa: A question answering
benchmark on a hybrid of tabular and textual content
in finance. *arXiv preprint arXiv:2105.07624*. 880
881
882
883
884

A Dataset and Question Types Examples 885

This section provides representative examples from
each of the eight table structure types identified
in Table 6, demonstrating the diversity and com-
plexity of multimodal reasoning challenges across
different structural configurations. 886
887
888
889
890

A.1 Single Entity 891







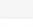
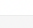


RK	DEFENSE	G	DRIVES	STOP RATE	PTS/DRIVE
1	 Minnesota	4	40	85%	0.85
2	 Michigan	6	72	83.3%	0.94
3	 San Jose State	4	46	82.6%	1.02
4	 Illinois	5	65	81.5%	0.74
5	 Penn State	5	64	81.2%	1.05
6	 Alabama	6	84	81%	0.89
7	 James Madison	4	56	80.4%	1.21
8	 Iowa State	5	60	80%	1.20
9	 Iowa	5	59	79.7%	0.95
10	 Ohio State	6	73	78.1%	1.19

Figure 2: An example of single entity multimodal table in the sports domain with only one type of entity (Image of Logo+Text).

Figure 2 is a collegiate football defensive perfor-
mance table that ranks the top ten teams based on
their Stop Rate, the percentage of opponent drives
that end without a score. The table includes five
columns detailing performance metrics. 892
893
894
895
896

Example Questions: 897

- Q: What is the Stop Rate for the Fighting Illini?
A: 81.5% **Type:** Explicit. 898
899
- Q: What is the difference in DRIVES faced between
the highest-ranked team that played 4 Games and
the highest-ranked team that played 6 Games?
A: 32 **Type:** Implicit. 900
901
902
903

904 **Q:** Which team has an allowed PTS/DRIVE that is
905 exactly 1.21?

906 **A:** James Madison **Type:** Answer-Mention.

907 **Q:** How many teams have at least one letter from
908 their initials in their logos?

909 **A:** 7 **Type:** Visual.

910 **A.2 Multiple Entity**

911 Figure 3 is a Formula 1 Race Results Summary
912 Table detailing the performance of multiple drivers
913 across three consecutive racing seasons. The table
914 is organized by Year, and each year begins with a
915 header row displaying the Race Location Flags for
916 the first three races of that season. The subsequent
917 rows for each year list the Driver, their Nationality,
918 and their finishing position in the three respective
919 races.

Year	Driver	Nationality	Race 1	Race 2	Race 3
2023					
			Retired	Seventh	Retired
			Fourth	Sixth	Twelfth
2024					
			Fourth	Third	Second
			Third	Withdrew	First
2025					
			Tenth	Disqualify	Seventh
			Eighth	Disqualify	Fourth

Figure 3: An example of multiple entity multimodal table in the motorsports domain with 2 types of entities (Image of flags + humans).

920 **Example Questions:**

921 **Q:** How many British drivers finished in the top 5 in
922 Australia?

923 **A:** 0 **Type:** Explicit.

924 **Q:** What was the driver's position at the same venue
925 one year after he finished 12th?

926 **A:** First **Type:** Implicit.

927 **Q:** Which driver replaced the driver whose national
928 flag contains yellow, for one race?

929 **A:** Oliver Bearman **Type:** Answer-Mention.

930 **Q:** How many stars appear in this table?

931 **A:** 23 **Type:** Visual.

932 **A.3 Single Chart**

933 Figure 2 is a Product Sales Performance Dash-
934 board that displays the performance of six products
935 (Amarilla, Carretera, Montana, Paseo, Velo, VTT)
936 and a Total row across two metrics: Sum of Gross
937 Sales, which shows each product's percentage con-
938 tribution to total sales, and a line chart labeled Sum
939 of Gross Sales by Month Number, which visualizes
940 the monthly sales trend for each product.

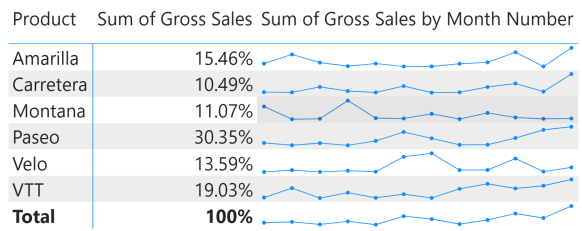


Figure 4: An example of a multimodal table with 1 type of chart(sparkline)

941 **Example Questions:**

942 **Q:** What is the difference in Overall gross sales be-
943 tween the product whose sales rose the most after
944 March vs dropped the most after March?

945 **A:** 0.58% **Type:** Implicit.

946 **Q:** Which Product has the most number of drops
947 followed by a peak?

948 **A:** VTT **Type:** Visual.

949 **A.4 Multiple Charts**

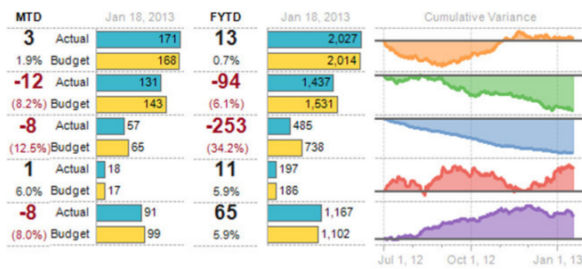


Figure 5: An example of a multimodal table in a Financial Context with multiple chart types (bar chart + area chart)

950 Figure 5 presents a multimodal financial dashboard
951 displaying Monthly-to-Date (MTD) and Fiscal-
952 Year-to-Date (FYTD) comparisons between actual
953 and budgeted values across five performance met-
954 rics. It also includes visual representations of cum-
955 ulative variance trends.

956 **Example Questions:**

957 **Q:** What is the actual MTD for the third metric on
958 January 18, 2013?

959 **A:** 57 **Type:** Explicit.

960 **Q:** How many FYTD units were sold in total?

961 **A:** 5313 **Type:** Implicit.

962 **Q:** Which metric reported the greatest increase?

963 **A:** FYTD **Type:** Answer-Mention.

964 **Q:** What colored variance graph shows a steady de-
965 cline over time?

966 **A:** Blue **Type:** Visual.

967 A.5 Visualizations

968 Figure 6 is a table detailing the design and compo-
969 nents of four chemical compounds or drug candi-
970 dates, identified by an ID number. The table
971 presents a common Core Fragment structure for all
972 compounds containing three variable attachment
973 points. The subsequent three columns (R1, R2, R3)
974 display the specific chemical groups or fragments
975 intended to be attached at each corresponding posi-
976 tion on the core structure.

ID	Core Fragment	R1	R2	R3
4671				
4988				
4989				
4990				

Figure 6: An example of multimodal table in the chemistry domain with scientific visualizations.

977 Example Questions:

978 **Q:** Which ID numbers have a benzene bases R3?

979 **A:** [4989, 4990] **Type:** Explicit.

980 **Q:** How many Nitrogen atoms exist in R1 and R3
981 combined?

982 **A:** 8 **Type:** Implicit.

983 **Q:** What is the name of the common ion in R2?

984 **A:** Hydroxide Ion **Type:** Answer-Mention.

985 **Q:** How many Hydrogen atoms exist in the R1 of
986 4988?

987 **A:** 13 **Type:** Visual.

988 A.6 Entities+Maps

989 Figure 7 is a table detailing four former franchises
990 from the Indian Premier League (IPL). The table

presents four columns: the team's logo and name,
the home City, the State represented by a map of
India with the relevant state highlighted in red, and
the year the team made its league Debut.

Team	City	State	Debut
	Hyderabad		2008
	Kochi		2011
	Pune		2011
	Pune		2016

Figure 7: An example of a multimodal table in the sports domain with maps and entities. (Image of logos + highlighted map)

Example Questions:

996 **Q:** In which year did Pune Warriors India debut?

997 **A:** 2011 **Type:** Explicit.

998 **Q:** Which state does the team with the purple logo
999 play in?

1000 **A:** Maharashtra **Type:** Implicit.

1001 **Q:** Which team from South India debuted in 2011?

1002 **A:** Kochi Tuskers Kerala **Type:** Answer-
1003 Mention.

1004 **Q:** How many teams have animals in their logo?

1005 **A:** 3 **Type:** Visual.

A.7 Entities+Charts

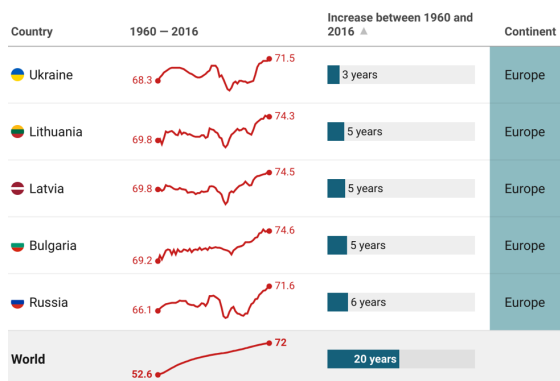


Figure 8: An example of a multimodal table in the Social Sciences Domain with charts and entities (Images of country flags, sparkline chart, and bar charts).

1007

Figure 8 is a data visualization comparing the change in life expectancy across five Eastern European countries and the global average between 1960 and 2016. The table features four columns: Country, a line chart showing the trend, the total Increase between 1960 and 2016, and the Continent.

1008

1009

1010

1011

1012

Example Questions:

1013

1014

1015

1016

Q: Has the score of Ukraine ever dropped below its starting point?

A: Yes **Type:** Explicit.

1017

1018

1019

Q: By how many percent did the global data increase throughout the years?

A: 36.88% **Type:** Implicit.

1020

1021

Q: Which Country had the highest peak?

A: Bulgaria **Type:** Answer-Mention.

1022

1023

1024

Q: By how many points did the country whose flag does not contain any red grow?

A: 3.2 **Type:** Visual.

A.8 Maps Only

State / Union Territory	1st Language	2nd Language	3rd Language	4th Language
	Malayalam	Dhivehi	Tamil	Hindi
	Hindi	Marathi	Urdu	Sindhi
	Marathi	Hindi	Urdu	Gujarati
	Meitei	Nepali	Hindi	Bengali

Figure 9: An example of a multimodal table in the Linguistics Domain with maps

Figure 9 is a table that displays the four most spoken languages in four different States or Union Territories of India. The table features five columns: an image column showing a map of India with the relevant State / Union Territory highlighted in red, followed by four columns listing the 1st, 2nd, 3rd, and 4th most spoken Language in that region.

1026

1027

1028

1029

1030

1031

1032

1033

Example Questions:

1034

Q: What is Lakshadweep's third language?

A: Tamil **Type:** Explicit.

1035

Q: What is the first language of the region that borders Assam?

A: Meitei **Type:** Implicit.

1036

1037

1038

Q: Name all states/union territories that speak Hindi?

A: [Lakshadweep, Madhya Pradesh, Maharashtra, Manipur] **Type:** Answer-Mention.

1039

1040

1041

1042

Q: How many states are completely landlocked?

A: 2 **Type:** Visual.

1043

1044

B Inter-Annotator Guidelines

1045

To ensure dataset quality and reliability, we implemented a comprehensive multi-stage annotation process with independent review and validation. Our approach involved two phases of inter-annotator agreement analysis to assess question correctness across all 4,021 questions.

1046

1047

1048

1049

1050

1051

Question Correctness Evaluation Criteria:

1052

All annotators used a standardized 3-point scale for question correctness evaluation with detailed guidelines to ensure consistency:

1053

1054

1055

- **Score 0 (Incorrect):** Questions that contain factual errors, ambiguous phrasing that prevents accurate answering, or answers that cannot be derived from the provided table content. This includes questions with incorrect references to table elements, logical inconsistencies, or answers that contradict visual or textual information in the table.
- **Score 1 (Partially Correct):** Questions that are generally well-formed but have minor issues such as slight ambiguity in phrasing, answers that are approximately correct but lack precision, or questions that could benefit from clearer wording. These questions are answerable but may require interpretation or have multiple plausible answers.
- **Score 2 (Fully Correct):** Questions that are clearly formulated, unambiguous, and have definitive answers that can be accurately derived from the table content. These questions demonstrate appropriate difficulty level, require genuine multimodal reasoning, and align perfectly with the provided visual and textual information.

Table Filtering Criteria: Annotators evaluated tables using comprehensive guidelines to ensure dataset quality and appropriateness for multimodal reasoning evaluation:

- **Structural Coherence Requirements:** Tables must demonstrate logical organization with clear relationships between textual and visual components. The structure should support meaningful multimodal reasoning tasks, with appropriate alignment between headers, data cells, and embedded visual elements.
- **Content Complexity Thresholds:** Tables must contain sufficient complexity to warrant multimodal analysis, including diverse visual elements (charts, maps, images), hierarchical data organization, and content that requires integration of both visual and textual information for comprehension.
- **Reasoning Appropriateness:** Tables should enable authentic multimodal reasoning scenarios where visual elements are essential for answering questions, rather than serving merely decorative purposes. The content must support various question types including fact verification, mathematical calculation, extrema identification, and visual-based inference.

- **Profanity and Sensitivity Compliance:** Tables must not contain or depict any profane, explicit, hateful, or sensitive content, including personally identifiable information (PII) or imagery that violates ethical, cultural, or privacy standards. All content should be appropriate for general academic and research use.
- **Data Source and Open-Domain Validity:** All table data and visual content should originate from credible, open-domain, and publicly accessible sources. Datasets must be free of copyright or usage restrictions that prevent open dissemination, ensuring reproducibility and transparency in multimodal research.

Cohen’s Kappa (κ) was computed on independent judgments to measure inherent agreement across both evaluation criteria.

Phase 1 - Internal Review: Table 3 reports Cohen’s Kappa and Table 13 provides the confusion matrix for those scores.

Annotator A	Annotator B			Total
	0	1	2	
0	156	18	4	178
1	22	298	45	365
2	8	51	3419	3478
Total	186	367	3468	4021

Table 13: Confusion Matrix: Internal Expert Annotators (Phase 1)

Phase 2 - External Validation (Subset): Table 4 reports Cohen’s Kappa and Table 14 provides the confusion matrix for those scores.

External Annotators	Gold Standard			Total
	0	1	2	
0	32	3	1	36
1	4	71	8	83
2	2	9	674	685
Total	38	83	683	804

Table 14: Confusion Matrix: External Annotators vs Gold Standard (Phase 2)

C Additional Result Tables

This section presents detailed performance results across all reasoning types and answer types discussed in Section 4.

Model	Fact Verification			Mathematical			Extrema			Vision Based		
	EM	SS	F1	EM	SS	F1	EM	SS	F1	EM	SS	F1
Missing Image Baseline												
Gemini 1.5 Flash	28.86	29.85	0.179	16.42	17.52	0.059	16.73	17.47	0.070	5.23	5.52	0.030
Gemini 2.0 Flash	27.62	27.61	0.099	15.96	16.35	0.049	15.80	16.20	0.059	11.00	11.04	0.032
GPT-4o Mini	44.75	47.64	0.361	23.80	23.29	0.124	27.45	27.07	0.146	24.09	25.64	0.177
Llama 3-8B	38.48	40.60	0.258	20.25	22.18	0.095	18.15	20.17	0.115	23.01	23.85	0.168
Mixtral	47.43	52.33	0.429	24.66	30.91	0.174	27.29	32.58	0.184	29.07	33.38	0.235
Entity Replaced Baseline												
Gemini 1.5 Flash	68.18	72.65	0.322	40.10	42.48	0.210	57.01	60.21	0.324	-	-	-
Gemini 2.0 Flash	71.69	73.28	0.333	42.11	41.20	0.191	64.80	65.86	0.347	-	-	-
GPT-4o Mini	79.45	75.11	0.628	39.00	43.92	0.255	54.06	57.00	0.357	-	-	-
Llama 3-8B	70.45	74.84	0.586	32.15	37.40	0.205	41.02	44.32	0.270	-	-	-
Mixtral	80.17	83.34	0.643	34.05	39.97	0.237	44.16	49.18	0.317	-	-	-
Image Captioning Baseline												
Gemini 1.5 Flash	47.90	49.29	0.369	13.68	15.38	0.074	25.13	27.47	0.165	22.72	23.76	0.186
Gemini 2.0 Flash	48.14	51.87	0.358	19.60	21.43	0.092	28.48	33.31	0.179	29.15	30.22	0.235
Table as an Image Baseline												
Gemini 1.5 Flash	47.17	45.77	0.207	23.42	25.33	0.095	21.30	23.01	0.113	30.40	32.38	0.094
Gemini 2.0 Flash	45.75	46.04	0.254	29.57	32.16	0.174	29.92	30.27	0.179	29.11	30.72	0.163
GPT-4o Mini	63.34	64.06	0.490	34.00	36.52	0.214	41.62	43.86	0.275	38.50	41.45	0.306
InternVL 2.5-8B	19.94	46.19	0.222	15.80	38.13	0.154	15.94	33.66	0.148	11.76	38.97	0.141
Mantis-8B-Idefics2	23.28	27.00	0.155	19.63	22.24	0.084	15.54	16.20	0.071	19.07	22.51	0.118
Phi-3.5-vision-instruct	32.16	35.79	0.238	13.87	16.18	0.052	13.73	15.01	0.063	15.24	17.70	0.075
Qwen-2.5-VL	53.12	55.19	0.265	21.45	25.08	0.113	27.90	31.39	0.158	22.54	24.40	0.102
Qwen-3-VL	58.20	60.50	0.380	28.50	31.20	0.165	35.80	38.40	0.220	31.80	34.20	0.215
Table-llava-1.5-7b-hf	24.57	27.06	0.133	14.57	15.51	0.038	8.42	9.54	0.031	9.88	10.93	0.049
Interleaved Baseline												
Gemini 1.5 Flash	41.37	41.14	0.283	17.96	18.00	0.079	27.78	27.21	0.155	26.81	27.28	0.213
Gemini 2.0 Flash	43.66	48.04	0.332	20.08	21.56	0.092	30.18	32.88	0.153	31.51	32.62	0.232
GPT-4o Mini	61.40	61.91	0.479	28.84	29.66	0.174	38.78	40.97	0.265	39.69	41.87	0.316
Mantis-8B-Idefics2	29.13	32.50	0.190	20.01	21.94	0.082	18.83	20.40	0.097	27.06	26.81	0.144
Phi-3.5-vision-instruct	25.67	28.31	0.168	21.98	24.72	0.107	16.39	17.39	0.093	17.33	19.14	0.106
Qwen-2.5-VL	43.03	62.76	0.381	15.77	48.33	0.106	24.90	52.69	0.198	25.99	52.47	0.225
Qwen-3-VL	38.50	61.80	0.310	14.20	26.50	0.085	22.40	34.60	0.155	24.30	36.80	0.190

Table 15: Performance analysis across reasoning types. EM: Exact Match, SS: Substring Match, F1: F1 Score

D Human Evaluation and Error Analysis

We conducted human evaluation with two annotators on 50 tables (11% of total dataset) and 393 questions (12% of total questions), where participants answered based solely on their knowledge without external image search capabilities. Results show humans consistently outperform all baseline models, achieving exact match scores in the 75-85% range and substring scores above 80%, demonstrating strong human performance while still revealing the inherent complexity of multimodal table reasoning tasks.

For the Entity Replaced baseline (upper bound), humans would perform exceptionally well since all required information appears in text format. For the Missing Image baseline (lower bound), models might outperform humans as they can compensate for missing information using pre-training knowledge, unlike human participants. Therefore, for fair comparison, we conduct in-depth analysis comparing the best-performing model across Table as Image and Interleaved baselines with human performance.

D.1 Types of Errors in Image Interpretation

We performed error evaluation and manually classified each error into one of the following categories:

- **Entity Disambiguation Issues** – These errors occur

when an image is incorrectly identified, leading to misinterpretation of its content. Misidentification can lead to entirely incorrect conclusions about the image.

- **Entity Identification Issues** – This error refers to the complete failure to recognize or identify the image.
- **Reasoning Errors** – This category includes mistakes where the image is correctly identified, but the logic used to answer the question is flawed. Such errors typically involve incorrect inferences, faulty assumptions, or logical inconsistencies.
- **Identification of Visual Attributes** – Errors in this category involve the failure to recognize key visual components of an image. This could mean missing out on crucial details such as shapes, colors, textures, or patterns.
- **Structural Errors** – These errors specifically pertain to the misinterpretation of tabular or structured data within an image. Failure to correctly identify rows, columns, and hierarchical relationships can result in inaccurate conclusions when analyzing tables, charts, or diagrams.
- **Mathematical Errors** – These occur when an individual or AI system miscalculates numerical information within the table. Errors can include

Model	SE		ME		SN		MN		IL		MT	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
Missing Image Baseline												
Gemini 1.5 Flash	18.9	0.071	21.1	0.081	13.3	0.037	20.4	0.084	15.4	0.028	22.1	0.112
Gemini 2.0 Flash	19.0	0.084	16.9	0.094	15.8	0.037	24.8	0.090	2.6	0.014	25.9	0.067
GPT-4o Mini	32.0	0.230	32.0	0.201	27.4	0.188	33.0	0.224	10.6	0.073	37.3	0.240
Llama 3-8B	28.1	0.188	29.7	0.182	21.3	0.137	29.1	0.145	0.0	0.000	25.4	0.160
Mixtral	35.3	0.273	33.3	0.247	28.4	0.209	36.7	0.233	3.4	0.025	38.0	0.316
Entity Replaced Baseline												
Gemini 1.5 Flash	37.8	0.114	57.4	0.385	7.2	0.615	25.7	0.193	8.0	0.300	36.7	0.308
Gemini 2.0 Flash	41.7	0.211	59.0	0.315	15.8	0.558	30.8	0.248	11.0	0.167	16.7	0.257
GPT-4o Mini	48.1	0.290	66.9	0.481	22.8	0.598	36.8	0.334	13.4	0.227	24.1	0.368
Llama 3-8B	39.8	0.210	57.5	0.412	16.2	0.504	31.7	0.308	13.8	0.106	14.2	0.317
Mixtral	41.3	0.293	56.6	0.446	29.7	0.593	35.0	0.348	14.8	0.205	27.7	0.350
Image Captioning Baseline												
Gemini 1.5 Flash	28.3	0.217	22.7	0.155	25.3	0.188	12.9	0.051	5.8	0.058	14.7	0.056
Gemini 2.0 Flash	33.3	0.243	34.5	0.230	25.8	0.183	23.2	0.109	4.0	0.029	29.6	0.131
Table As Image Baseline												
Gemini 1.5 Flash	38.0	0.223	41.6	0.274	21.6	0.078	46.6	0.255	23.2	0.113	34.3	0.199
Gemini 2.0 Flash	30.3	0.147	35.6	0.168	19.7	0.072	38.6	0.174	26.9	0.119	32.5	0.149
GPT-4o Mini	47.0	0.355	46.5	0.305	33.4	0.228	50.1	0.318	26.7	0.214	42.5	0.304
Intern-VL-2.5	17.7	0.187	18.9	0.189	12.1	0.125	27.6	0.292	8.6	0.080	23.0	0.227
Mantis-8B-Idefics2	22.6	0.140	23.8	0.106	10.7	0.040	24.6	0.088	25.6	0.122	22.7	0.083
Phi-3.5	18.2	0.089	21.2	0.111	12.5	0.038	17.6	0.023	0.1	0.002	14.0	0.056
Qwen-2.5-VL	29.7	0.160	32.1	0.190	15.3	0.059	29.8	0.161	14.1	0.067	29.2	0.181
Qwen-3-VL	38.5	0.265	39.8	0.250	24.0	0.145	40.2	0.245	21.0	0.145	36.5	0.245
Table_LLaVA	12.9	0.078	14.0	0.050	10.8	0.030	21.9	0.065	0.0	0.000	16.6	0.050
Interleaved Baseline												
Gemini 1.5 Flash	29.2	0.205	29.0	0.177	25.2	0.175	25.2	0.120	4.5	0.034	26.6	0.131
Gemini 2.0 Flash	33.2	0.223	34.6	0.210	27.9	0.195	29.4	0.156	26.6	0.131	27.3	0.103
GPT-4o Mini	48.1	0.372	47.7	0.301	30.9	0.223	40.0	0.252	21.6	0.127	24.4	0.184
Mantis-8B-Idefics2	26.0	0.175	29.2	0.144	15.6	0.064	15.8	0.061	27.0	0.139	21.7	0.066
Phi-3.5	21.4	0.136	27.7	0.144	13.7	0.060	28.8	0.131	13.7	0.071	27.0	0.185
Qwen-2.5-VL	29.4	0.261	21.5	0.177	23.5	0.185	22.0	0.158	5.7	0.060	20.2	0.165
Qwen-3-VL	39.0	0.320	35.5	0.240	27.0	0.205	31.0	0.210	19.5	0.105	23.8	0.175

Table 16: Performance across answer types. SE: Single Entity, ME: Multiple Entity, SN: Single Number, MN: Multiple Number, IL: Image Location, MT: Multiple Types.

Baseline	Model	Answer Mention Questions			Explicit Question			Implicit Question			Visual-Based Question		
		Exact	Substring	F1	Exact	Substring	F1	Exact	Substring	F1	Exact	Substring	F1
Table as an Image Baseline	GPT-4o-mini	34.8	38.7	0.32	44.8	46.9	0.34	40.7	42.3	0.32	41.0	45.2	0.36
Interleaved Baseline	GPT-4o-mini	41.6	45.9	0.35	46.4	47.8	0.37	31.4	33.0	0.26	32.2	34.7	0.27
Interleaved Baseline	Human answers	78.4	82.1	0.76	84.2	87.3	0.81	75.8	80.6	0.73	79.9	83.7	0.78

Table 17: Best performing baselines on sample set evaluation.

incorrect counting, misreading of numerical values, or computational mistakes when deriving conclusions from visual data that involves numbers.

- **Partial Answers** – This type of error involves providing an incomplete response where some crucial details are missing. The responder may identify and interpret part of the image correctly but fail to provide a comprehensive answer.
- **Extra Information or Hallucination** – This occurs when irrelevant or incorrect details are introduced into the response. The additional information may not be present in the table itself but might be inferred incorrectly based on prior pre trained knowledge.

Error Hierarchy and Root Cause Attribution.

While a single response can exhibit multiple types of errors, we adopt a hierarchical approach to identify the *root cause* of failure. This hierarchy ensures that derivative errors (e.g., reasoning or partial answers resulting from earlier misinterpretations) are

not double-counted.

During analysis, if an instance exhibits multiple error types, we attribute it only to the **root cause**. For example, if a model fails to recognize a visual attribute and consequently provides an incomplete answer, the error is classified under *Identification of Visual Attributes*, not as a reasoning or partial answer error.

The analysis highlights key patterns in human performance, with Entity Identification Issues being the most common, suggesting difficulties in recognizing entities, particularly without domain knowledge. Partial Answers reveal a tendency to overlook essential details, while Entity Disambiguation Issues indicate occasional struggles in distinguishing between similar entities. These three challenges emphasize the importance of domain knowledge in accurately solving tasks and identifying the correct entity.

Visual Attribute Identification Errors suggest that even humans can miss fine details in images. However, Reasoning Errors are relatively low, in-

Type of Error	Human (%)	Interleaved (%)		Table as Image (%)	
		GPT-4o-mini	Qwen2.5-VL	GPT-4o-mini	Qwen2.5-VL
Entity Identification Issues	74.44	13.23	8.42	18.67	16.74
Entity Disambiguation Issues	5.21	11.00	7.85	8.00	6.27
Identification of Visual Attributes	11.32	17.08	12.78	16.00	9.38
Partial Answer	6.02	12.08	13.32	14.00	15.87
Reasoning Errors	3.01	14.38	16.11	15.33	17.98
Structural Errors	0.00	11.15	13.46	7.33	12.91
Mathematical Errors	0.00	10.31	11.09	8.67	9.59
Extra Info / Hallucination / Pretrained Knowledge	0.00	10.77	16.97	12.00	11.26
Total	100.00	100.00	100.00	100.00	100.00

Table 18: Error type distribution across evaluation methods.

dicating that humans generally follow logical processes. Notably, there are no Mathematical Errors, Structural Errors, or Hallucinations, demonstrating strong numerical reasoning and structured data interpretation.

Overall, humans excel in logical analysis and structured comprehension but face challenges in entity recognition, visual attention to detail, and providing complete answers.

D.2 Performance and Analysis of Errors in Interleaved baseline GPT 4o-mini

GPT-4o-mini exhibits significant errors across multiple categories, with each exceeding the 10% threshold. Reasoning Errors are the most prominent, as the model struggles with multi-step deductions and logical coherence - an ongoing challenge for MLLMs. Visual Attribute Identification is another major limitation, with failures in extracting key image features, underscoring the need for better vision encoders and fine-tuning.

Entity Identification and Disambiguation errors occur at similar rates, as the model misidentifies or fails to recognize entities, likely due to insufficient training data and over-reliance on context. Structural Errors show difficulty in interpreting complex tabular data, including hierarchical structures and nested tables, while Hallucination Errors highlight the model’s tendency to generate irrelevant information. Mathematical Errors reveal persistent struggles with quantitative reasoning.

These issues demonstrate that while GPT-4o-mini has some multimodal reasoning capability, it lacks depth in entity recognition, visual interpretation, and structured data comprehension.

D.3 Performance and Analysis of Errors in Table as an Image baseline GPT 4o-mini

GPT4o-mini makes many errors across all categories. From table 18, in table as image baseline,

we notice that GPT4o-mini had most errors in Entity Identification. This is likely due to insufficient training data to much reliance on context.

It also struggles with Identification of visual attributes and reasoning errors in similar range as Interleaved baseline. We see similar error rates among Mathematical, Hallucinations, Partial and Entity Disambiguation error types showing that GPT4o-mini struggles with these issues irrespective of input table format.

Interestingly, We observe that structural Errors decreased by ~4% in table as image baseline. This may be due to the fact that looking at entire table at once through image helped understand the structure better. So, passing tables as image does help GPT4o-mini to understand the structure of tables better.

From above both analysis we can say GPT4o-mini has scope of improvement in identifying visual attributes, on entity identification and reasoning errors.

D.4 Performance and Analysis of Errors in Interleaved baseline Qwen2.5-VL

The error distribution in Qwen2.5-VL suggests a model that demonstrates broad multimodal competence but lacks consistent grounding in visual inputs. The high incidence of hallucination-related errors (16.97%) implies that the model frequently relies on pretrained textual priors rather than accurately integrating perceptual evidence. This points to limitations in the vision-language alignment mechanisms, where generated responses may reflect plausible but unsubstantiated inferences.

Reasoning errors (16.11%) and structural errors (13.46%) indicate significant challenges in multi-step logical processing and the interpretation of complex visual structures, such as nested tables or hierarchical layouts. These errors suggest that while Qwen2.5-VL is able to extract surface-level

1304 patterns, it struggles with deeper, context-sensitive
1305 reasoning tasks - an area where multimodal systems
1306 continue to underperform.

1307 Although errors in entity identification (8.42%)
1308 and disambiguation (7.85%) are comparatively
1309 lower, they reflect persistent difficulties in main-
1310 taining contextual precision, particularly in visu-
1311 ally dense or ambiguous scenarios. Overall, the
1312 model's performance underscores the need for im-
1313 proved grounding strategies, enhanced visual en-
1314 coding, and tighter integration of reasoning capa-
1315 bilities to support more reliable and context-aware
1316 multimodal understanding.

1317 D.5 Performance and Analysis of Errors in 1318 Table as an Image baseline Qwen2.5-VL

1319 Qwen2.5-VL's performance on table-as-image
1320 tasks highlights persistent challenges in structured
1321 data interpretation. The most prominent error cate-
1322 gory is Reasoning Errors (17.98%), suggesting the
1323 model struggles to perform logical operations over
1324 tabular data, particularly when spatial layout or im-
1325 plicit relationships between rows and columns are
1326 involved. This reflects a broader difficulty in trans-
1327 lating visual table structures into coherent semantic
1328 representations.

1329 Partial Answers (15.87%) and Entity Identifica-
1330 tion Issues (16.74%) further indicate that Qwen2.5-
1331 VL often fails to capture complete and precise infor-
1332 mation from table-based inputs. These errors likely
1333 stem from limitations in visual attention mecha-
1334 nisms or token alignment between image features
1335 and textual output. Similarly, the model's rela-
1336 tively high rate of Visual Attribute Identification
1337 errors (9.38%) and Structural Errors (12.91%) re-
1338 veals an incomplete understanding of layout hierar-
1339 chies, cell groupings, and formatting cues critical
1340 for accurate parsing.

1341 While hallucination (11.26%) and mathematical
1342 errors (9.59%) are somewhat less frequent, they
1343 still suggest inconsistencies in quantitative reason-
1344 ing and information grounding. Collectively, these
1345 findings indicate that despite some robustness in
1346 basic vision-language alignment, Qwen2.5-VL re-
1347 quires substantial improvements in handling visu-
1348 ally structured data, particularly for tasks requiring
1349 deep reasoning and layout-aware comprehension.

1350 E Prompts

1351 This section contains the exact prompts used in our
1352 experiments.

Missing Image Baseline Prompt

You will be provided a table in a pipe-separated table where all the entities have been removed. Your task is to:

Step 1: UNDERSTAND THE TABLE CONTEXT - Carefully analyze the table structure and identify its purpose and what it mentions.

Step 2: FILL IN THE GAPS - Use the table context and your real-world knowledge to deduce the missing entities logically.

Step 3: ANALYZE THE QUESTIONS - Read all the questions provided and explore ****ALL TYPES OF REASONING**** to find answers, including but not limited to Numerical reasoning (relationships, totals, and comparisons), Visual reasoning (Colors, shapes, or patterns), Contextual reasoning, (Real-world connections or logic), etc.

Step 4: PROVIDE ANSWERS IN FORMAT - Ensure that all answers adhere strictly to the FORMAT specified. Avoid deviating from this format or including unnecessary explanations.

{Answer Formatting Guidelines}
ALWAYS PROVIDE YOUR ANSWERS IN THIS FORMAT.

If you are unable to answer it, simple answer *UNKNOWN*.

I will provide one example to show you:

{One Shot Example}

Now I will provide you with the table and question.

{TABLE}

{QUESTION}

Based on the examples that I have provided and the steps I mentioned above, answer the question.

Entity Replaced Baseline Prompt

You will be provided a pipe-separated table format that contains some entities. Your task is to:

Step 1: UNDERSTAND THE TABLE CONTEXT - Carefully analyze the table structure and identify its purpose and what it mentions.

Step 2: ANALYZE THE QUESTIONS - Read all the questions provided and explore ****ALL TYPES OF REASONING**** to find answers, including but not limited to Numerical reasoning (relationships, totals, and comparisons), Visual reasoning (colors, shapes, or patterns), Contextual reasoning (real-world connections or logic), etc.

Step 3: PROVIDE ANSWERS IN FORMAT - Ensure that all answers adhere strictly to the FORMAT specified. Avoid deviating from this format or including unnecessary explanations.

{Answer Formatting Guidelines}
ALWAYS PROVIDE YOUR ANSWERS IN THIS FORMAT.

****IMPORTANT**** ALL answers are there in the table. ANALYZE the question and table properly.

I will provide one example to show you:

{One Shot Example}

Now I will provide you with the table and question.

{TABLE}
{QUESTION}

Based on the examples that I have provided and the steps I mentioned above, answer the question.

Image Captioning Baseline Prompt

You will be provided a table in a pipe-separated table with images included. Your task is to:

Step 1: UNDERSTAND THE TABLE CONTEXT - Carefully analyze the table structure and identify its purpose and what it mentions.

Step 2: CAPTION EVERY IMAGE - Based on the image, provide a caption for that image. Your job is to reason, predict, and replace image entity tags and provide visual descriptions.

Step 3: CREATE A TABLE - Based on the image captions, create a pipe-separated table where the image placeholders or cells have been replaced with their captions.

Step 4: ANALYZE THE QUESTIONS - Read all the questions provided and explore ****ALL TYPES OF REASONING**** to find answers.

{TABLE WITH CAPTIONS}

Step 5: PROVIDE ANSWERS IN FORMAT Ensure that all answers adhere strictly to the FORMAT specified.

{Answer Formatting Guidelines}
ALWAYS PROVIDE YOUR ANSWERS IN THIS FORMAT.

****IMPORTANT**** ALL answers are there in the table. ANALYZE the question and table properly.

Now I will provide you with the question.
{QUESTION}

Based on the steps I mentioned above, answer the question.

Table as an Image Baseline Prompt

You will be provided an image of a table. Your task is to:

Step 1: UNDERSTAND THE IMAGE CONTEXT - Carefully analyze the image content and understand the tabular structure and all text and visual aspects inside the image.

Step 2: ANALYZE THE QUESTIONS - Read all the questions provided and explore ****ALL TYPES OF REASONING**** to find answers.

Step 3: PROVIDE ANSWERS IN FORMAT Ensure that all answers adhere strictly to the **FORMAT** specified. Avoid deviating from this format or including unnecessary explanations.

{Answer Formatting Guidelines}
ALWAYS PROVIDE YOUR ANSWERS IN THIS FORMAT.
****IMPORTANT**** ALL answers are there in the image.

Now I will provide you with the image of the table.

{IMAGE OF TABLE}
For this image, you will answer the following question.
{QUESTION}

Based on the steps I mentioned above, answer the question.

Interleaved Baseline Prompt

You will be provided a pipe-separated table where some cells are images. Your task is to:

Step 1: UNDERSTAND THE TABLE CONTEXT - Carefully analyze the table structure and understand the intricate relationship between image and text.

Step 2: ANALYZE THE QUESTIONS - Read all the questions provided and explore ****ALL TYPES OF REASONING**** to find answers.

Step 3: PROVIDE ANSWERS IN FORMAT Ensure that all answers adhere strictly to the **FORMAT** specified. Avoid deviating from this format or including unnecessary explanations.

{Answer Formatting Guidelines}
ALWAYS PROVIDE YOUR ANSWERS IN THIS FORMAT.

****IMPORTANT**** ALL answers are there in the image.

Now I will provide you with the table.

{INTERLEAVED TABLE}
For this table, answer the following question.
{QUESTION}

Based on the steps I mentioned above, answer the question.

Answer Formatting Guidelines The following answer formatting guidelines were provided along with every prompt to eliminate inconsistencies and ensure a uniform response structure across all task types. The model is expected to provide only the final answer, formatted as described below:

- **Single Entity:** Return a single *string* representing one entity such as a name, country, company, object, or similar. The answer should be concise and written in one line without extra text.
 - Example (Name): Elon Musk
 - Example (Country): China
 - Example (Company): Google
 - Example (Color): Red
- **Single Number:** - If the answer is a whole number, write it **without decimals**. - If it

1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376

1356

1357

1358

1377 has decimals, round to **two decimal places**.

1378 - If the last digit after rounding is 0 (e.g.,
1379 23.40), remove the trailing zero (→ **23.4**).

1380 Units should only be included if explicitly
1381 mentioned in the question.

1382 - Example (Whole Number): 45

1383 - Example (Decimal): 12.36

1384 - Example (Trimmed Decimal): 23.4

- 1385 • **Multiple Entities:** Provide a **list of strings**,
1386 each following the same rules as the Single
1387 Entity format. Use comma-separated values
1388 enclosed in square brackets.

1389 - Example: ["Apple", "Microsoft",
1390 "Google"]

- 1391 • **Multiple Numbers:** Provide a **list of num-**
1392 **bers**, each following the Single Number for-
1393 mation rule. Use comma-separated values
1394 enclosed in square brackets.

1395 - Example: [23, 45.67, 89.4]

- 1396 • **Image Locations:** When the answer involves
1397 identifying a location within a visual or tab-
1398 ular structure, specify it using the following
1399 format: row_num_col_num.

1400 - Example: row_2_col_3