

Does Table Source Matter? Benchmarking and Improving Multimodal Scientific Table Understanding and Reasoning

Anonymous ACL submission

Abstract

Recent large language models (LLMs) have advanced table understanding capabilities but rely on converting tables into text sequences. While multimodal large language models (MLLMs) enable direct visual processing, they face limitations in handling scientific tables due to fixed input image resolutions and insufficient numerical reasoning capabilities. To address these challenges, we present MMSci, a comprehensive dataset for scientific table understanding and reasoning. MMSci consists of three key components: (1) MMSci-Pre, a domain-specific dataset of 52K scientific table structure recognition samples, (2) MMSci-Ins, an instruction tuning dataset with 12K samples across three table-based tasks, and (3) MMSci-Eval, a benchmark with 3,114 testing samples specifically designed to evaluate numerical reasoning capabilities. Based on MMSci, we leverage and evaluate MLLMs with dynamic input resolution capabilities for scientific table understanding. Extensive experiments demonstrate that our domain-specific approach with 52K scientific table images achieves superior performance compared to 150K general-domain tables, highlighting the importance of data quality over quantity. Our proposed framework shows significant improvements in both general table understanding and numerical reasoning capabilities, with strong generalisation to held-out datasets.

1 Introduction

Tables serve as a fundamental tool for organising structured information across diverse domains. Recent studies have shown the potential of leveraging large language models (LLMs) to automatically understand and process tabular data, which has emerged as a critical research direction with applications such as Table Question Answering (TQA) (Pasupat and Liang, 2015), Table Fact Verification (TFV) (Chen et al., 2020a), and Table-to-Text Generation (T2T) (Moosavi et al., 2021).

However, current table-oriented LLMs (Zhang et al., 2023; Li et al., 2023b) face inherent limitations as they require converting tables into sequential text formats (i.e., HTML strings), potentially losing crucial structural and positional information. While table-based multimodal large language models (MLLMs) have addressed this by enabling direct processing of table images, several critical limitations persist: (1) fixed input image resolutions that constrain practical applicability, (2) limited capability in processing scientific tables that contain significant numerical values, and (3) insufficient numerical reasoning abilities for scientific domain tasks. These limitations are particularly significant in scientific domains, where tables frequently incorporate complex numerical relationships.

Recent MLLMs have demonstrated success with general tables but struggle with scientific tables due to their dense numerical content and complex reasoning requirements. Our work demonstrates that domain-specific data quality significantly outperforms quantity, challenging conventional scaling laws in multimodal learning. Scientific table numerical reasoning requires multi-step operations including addition, subtraction, comparison, and other mathematical operations to derive conclusions from tabular data, going beyond simple fact extraction. Current MLLMs, however, lack the specific training data to handle these sophisticated scientific table understanding and reasoning requirements.

To address these challenges, we introduce MMSci, a comprehensive dataset for scientific table understanding and reasoning. We first conduct a systematic analysis of table source effectiveness through MMSci-Pre, a carefully curated dataset containing 52K structure recognition samples derived from scientific papers. Our experimental results demonstrate that MLLMs trained on these domain-specific table images significantly outperform those trained on 150K general-domain tables,

084 establishing the importance of data quality over
085 quantity in table understanding tasks.

086 Building upon this foundation, we then create
087 MMSci-Ins, an instruction tuning dataset com-
088 prising 12K samples with explicit intermediate
089 reasoning steps across three fundamental tasks:
090 TQA, TFV, and T2T. Each sample includes de-
091 tailed step-by-step reasoning processes to develop
092 models’ table-based numerical reasoning and sci-
093 entific analysis capabilities. To overcome the lim-
094 itations of fixed-resolution approaches in existing
095 table MLLMs (Lee et al., 2023; Alonso et al., 2024;
096 Zheng et al., 2024), we leverage and evaluate exist-
097 ing dynamic input resolution capabilities across dif-
098 ferent model architectures (Qwen2-VL-7B-Instruct
099 and LLaVA-NeXT-7B). Our analysis reveals that
100 mere technical capability for dynamic resolution is
101 insufficient without proper cross-modal alignment.
102 Experimental results demonstrate consistent per-
103 formance improvements across both general table
104 understanding and specialised numerical reasoning
105 tasks.

106 To enable comprehensive evaluation, we estab-
107 lish MMSci-Eval, a benchmark with 3,114 testing
108 samples requiring numerical reasoning capabilities.
109 The benchmark provides rigorous assessment of
110 models’ performance across TQA, TFV, and T2T
111 tasks. Our extensive experiments demonstrate that
112 our 52K scientific table images prove more effec-
113 tive than 150K general-domain table images for
114 both general understanding and numerical reason-
115 ing tasks. This efficiency highlights the value of
116 domain-specific, high-quality data in developing
117 robust table understanding capabilities.

118 Our contributions are summarised as follows:

- 119 • We introduce MMSci, a comprehensive
120 dataset consisting of three components: (1)
121 MMSci-Pre, consists of 52K table image-to-
122 HTML table structure recognition samples;
123 (2) MMSci-Ins, an instruction tuning dataset
124 of 12K samples with reasoning steps; and (3)
125 MMSci-Eval, a benchmark with 3,114 sam-
126 ples for numerical reasoning capabilities as-
127 sessment across TQA, TFV, and T2T tasks.
- 128 • We develop a comprehensive table-based
129 MLLM framework that achieves strong per-
130 formance on three table-based numerical rea-
131 soning tasks while demonstrating robust gen-
132 eralisation to held-out datasets.
- 133 • We implement dynamic input resolution ca-

pabilities across different model architectures,
validating the effectiveness of our approach
through consistent performance gains on both
Qwen2-VL-7B-Instruct and LLaVA-NeXT-
7B.

2 Related Work 139

2.1 Table Understanding Models 140

141 Early table-based models based on general lan-
142 guage models with large-scale table corpus (Liu
143 et al., 2022; Chen et al., 2023) only support lim-
144 ited types of tables and tasks. Table understanding
145 capabilities have been enhanced through prompt en-
146 gineering (Chen, 2023; Sui et al., 2023), instruction
147 tuning (Zhang et al., 2023; Li et al., 2023b; Yang
148 et al., 2024b) and external tools (Lu et al., 2023a;
149 Li et al., 2023a) with the development of LLMs.
150 However, these approaches require converting ta-
151 bles into text formats, limiting their applications.

152 Recently, MLLMs have emerged as a promising
153 direction for table understanding. TableGPT2 (Su
154 et al., 2024a) features a novel table encoder to han-
155 dle table cell-level information. Pix2Struct (Lee
156 et al., 2023) introduces a unified image-to-text
157 model pretrained on web page screenshots with
158 HTML supervision. PixT3 (Alonso et al., 2024)
159 takes table-to-text tasks as table visual recognition
160 tasks and generates texts. Table-LLaVA (Zheng
161 et al., 2024) introduces a novel multimodal table
162 understanding approach that directly processes ta-
163 ble images. However, these approaches do not
164 focus on datasets requiring sophisticated numerical
165 reasoning capabilities.

2.2 Table-based Reasoning and Datasets 166

167 Table-based reasoning requires reasoning over both
168 free-form natural language queries and structured
169 tables. Early works either rely on executable lan-
170 guages (*e.g.*, SQL) (Yin et al., 2016; Yu et al.,
171 2018) to capture logical structure in statements.
172 TAPAS (Herzig et al., 2020), and DATER (Ye et al.,
173 2023) encode sentence-table pairs and transform
174 table-based reasoning into question-answering or
175 inference tasks. Existing datasets primarily fo-
176 cus on specific domains like Wikipedia and fi-
177 nance. HybridQA (Chen et al., 2020b) derived
178 from Wikipedia emphasises span lookup, while
179 TAT-QA (Zhu et al., 2021), FinQA (Chen et al.,
180 2021), and DocMath-Eval (Zhao et al., 2024) ad-
181 dress numerical reasoning in the financial domain.

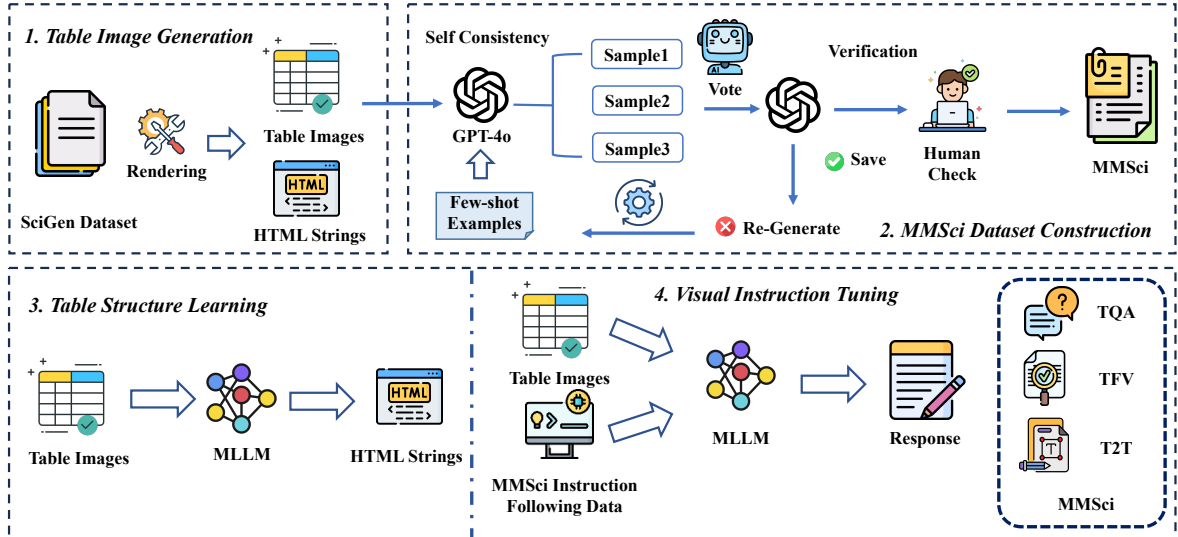


Figure 1: Overview of the proposed framework, which consists of four key stages: (1) Table Image Generation; (2) MMSci Dataset Construction; (3) Table Structure Learning; and (4) Visual Instruction Tuning.

SciGen (Moosavi et al., 2021) introduces a scientific table-to-text generation dataset that requires arithmetic reasoning, but focuses mainly on generation rather than comprehensive reasoning evaluation. However, these datasets have relatively limited reasoning types, significantly differing from real-world scientific table understanding that require numerical computation reasoning. To address this gap, we propose MMSci dataset that combines multiple reasoning types to enhance model performance on complex scientific table understanding tasks.

3 Construction of MMSci Dataset

As shown in Figure 1, the MMSci dataset construction consists of three key components: (1) Data Collection, (2) MMSci-Pre Dataset Construction, and (3) MMSci-Ins and MMSci-Eval Dataset Creation with Numerical Reasoning Augmentation. These components directly correspond to the dataset construction stages in the upper part of our framework.

3.1 Data Collection

To construct MMSci dataset, we focus on scientific tables containing significant numerical values and complex reasoning requirements. We collect raw tabular data from the SciGen dataset (Moosavi et al., 2021), which provides pairs of scientific tables and their corresponding descriptions across computer science research domains. These descriptions naturally require numerical reasoning operations (including addition, subtraction, max/min,

comparison, and division) over table values, making them ideal for our purpose. We transform the original textual tables into high-quality HTML format and then render them into table images while preserving their structural integrity. This process ensures the visual representation maintains the complex layouts and relationships present in the original scientific tables. From this process, we collect 52K image-to-HTML pairs based on tables from the training set and development set of the SciGen dataset.

3.2 MMSci-Pre Dataset Construction

Existing table-based MLLMs (Lee et al., 2023; Alonso et al., 2024; Zheng et al., 2024) demonstrate that generating textual table representations from table images is crucial for aligning visual structure with textual content. Based on our collected image-HTML pairs, we create 52K instruction-following image-to-HTML samples via the `Imgkit`¹ python package. Each sample consists of a table image paired with its corresponding HTML representation. The resulting dataset, MMSci-Pre, contains 52K samples specifically designed for table structure learning.

3.3 MMSci-Ins and MMSci-Eval Dataset Creation

For our instruction tuning and evaluation datasets, we select 12,000 tables from the training set and 1,038 from the testing set of SciGen dataset to

¹<https://pypi.org/project/imgkit/>

create MMSci-Ins and MMSci-Eval, respectively. These datasets focus on complex numerical reasoning tasks requiring multi-step operations including addition, subtraction, comparison, max/min identification, and other mathematical operations. For each table, we employ GPT-4o (OpenAI, 2024) to generate task-specific content across three table-based tasks. For Table Question Answering (TQA), we generate questions paired with step-by-step reasoning processes and final answers. For Table Fact Verification (TFV), we create claims with supporting reasoning steps and verification results (supported, refuted, or not enough information). For Table-to-Text Generation (T2T), we augment existing table-to-text pairs with detailed reasoning steps. To ensure quality, we implement a rigorous verification process. First, we apply self-consistency Chain-of-Thought (CoT) reasoning (Wang et al., 2023) with multiple reasoning paths and voting. Second, we use GPT-4o to validate consistency between reasoning steps and final outputs. Third, we conduct human verification on 40% of generated samples. Finally, we regenerate any identified incorrect samples to maintain dataset quality. This process results in MMSci-Ins with 12K high-quality instruction-tuning samples and MMSci-Eval with 3,114 testing examples. Both datasets maintain a balanced distribution across the three tasks, with each table paired with one sample per task type. Each sample includes detailed step-by-step reasoning processes that enable models to learn both final outputs and the logical progression needed to arrive at conclusions. Detailed dataset quality statistics are provided in Appendix A.2.

4 Experiments

4.1 Model Training

To demonstrate the effectiveness of MMSci dataset, we train two series of MLLM following the architecture of Qwen2-VL-7B-Instruct (Wang et al., 2024) and LLaVA-NeXT-7B (Li et al., 2024).

Model Architectures. Both models follow a three-component design: **Qwen2-VL-7B-Instruct** consists of a Vision Transformer (ViT) (Dosovitskiy, 2020) as the vision tower, a MLP as the vision-language connector, and a Qwen2-7B-Instruct (Yang et al., 2024a) as the language model. **LLaVA-NeXT-7B** uses a pre-trained CLIP model (Radford et al., 2021) as the visual encoder, a MLP connector, and a Vicuna-7B model (Chiang et al., 2023) as the backbone. In both architectures,

the vision encoder processes images into visual features, which are projected into the LLM’s word embedding space via the MLP connector.

We divide the training into two stages:

Table Structure Learning. We use both MMSci-Pre and MMTAB-Pre (Zheng et al., 2024) corpus (202K table image-to-HTML pairs in total) to align visual features with textual representations in different experimental settings as shown in Table 1. Models learn to generate HTML table representations, developing table structure perception capabilities. For LLaVA-NeXT-7B, only the MLP connector parameters are updated during this stage.

Visual Instruction Tuning. We use 12K instruction-following samples from MMSci-Ins to fine-tune the MLLMs while keeping visual encoders frozen. Only the MLP projection layer and LLM weights are updated, focusing on developing instruction-following numerical reasoning capabilities across TQA, TFV, and T2T tasks.

Notably, both models support dynamic input resolutions, addressing a key limitation of existing table MLLMs (Lee et al., 2023; Alonso et al., 2024; Zheng et al., 2024) that require fixed-size input image resolutions (e.g., 336×336). Qwen2-VL achieves this through 2D-RoPE (Su et al., 2024b) to capture two-dimensional positional information of images while LLaVA-NeXT employs a simpler approach of splitting images into grids and encoding them independently. While both Qwen2-VL and LLaVA-NeXT support dynamic input resolutions, our analysis reveals significant performance differences when applied to scientific tables. This suggests that mere technical capability for dynamic resolution is insufficient without proper cross-modal alignment.

4.2 Experimental Settings

Baselines. We select several state-of-the-art MLLMs as our baselines, including GPT-4V (OpenAI, 2023), InternVL-2-76B (Chen et al., 2024), LLaVA-NeXT series (72B/34B/13B/7B) (Li et al., 2024), Qwen-2-VL-Instruct series (72B/7B) (Wang et al., 2024), Table-LLaVA series (13B/7B) (Zheng et al., 2024), Pixtral-12B (Agrawal et al., 2024), Llama-3.2-11B-Vision-Ins. (Meta, 2024), MiniCPM-V-2.6-8B (Yao et al., 2024), and InternVL-2-8B (Chen et al., 2024).

Datasets and Metrics. The held-in evaluation sets in Table 1 include TQA, TFV and T2T tasks of MMSci-Eval. The held-out evaluation sets in Table 2 are from MMTAB-Eval benchmark (Zheng

Models	MMSci-Eval			Held-out	
	TQA Acc.	TFV Acc.	T2T BLEU	TABMWP Acc.	TAT-QA Acc.
Baseline					
GPT-4V (OpenAI, 2023)	53.13	78.01	4.80	60.00	32.50
InternVL-2-76B (Chen et al., 2024)	40.31	62.46	1.79	46.28	6.73
LLaVA-NeXT-72B (Li et al., 2024)	11.75	49.28	1.79	10.69	3.29
Qwen-2-VL-72B-Ins. (Wang et al., 2024)	39.11	64.06	2.83	41.42	17.65
LLaVA-NeXT-34B (Li et al., 2024)	9.73	42.19	2.33	6.96	1.29
LLaVA-NeXT-13B (Li et al., 2024)	2.31	1.83	1.79	1.67	0.43
Table-LLaVA-13B (Zheng et al., 2024)	8.57	51.15	0.03	<u>59.77</u>	15.67
Pixtral-12B (Agrawal et al., 2024)	0.96	5.49	4.12	4.64	7.46
Llama-3.2-11B-Vision-Ins. (Meta, 2024)	1.15	5.85	3.04	7.39	0.37
LLaVA-NeXT-7B (Li et al., 2024)	0.19	0.86	2.99	1.73	0.72
Qwen-2-VL-7B-Ins. (Wang et al., 2024)	25.62	52.79	3.04	34.43	16.19
InternVL-2-8B (Chen et al., 2024)	25.72	44.99	2.64	18.42	7.12
MiniCPM-V-2.6-8B (Yao et al., 2024)	26.58	33.23	0.07	24.30	11.94
Table-LLaVA-7B (Zheng et al., 2024)	7.99	39.30	0.03	57.78	12.82
Ours (LLaVA-NeXT-7B)					
MMSci-Pre (52k) + MMSci-Ins	17.72	57.12	2.93	49.47	10.46
MMTab-Pre (150k) + MMSci-Ins	15.79	56.16	2.88	47.55	8.03
MM-Pre (202k) + MMSci-Ins	23.02	58.57	2.36	49.72	12.27
w/o MM-Pre (202k)	15.22	51.73	2.86	46.24	7.63
Ours (Qwen2-VL-7B-Ins.)					
MMSci-Pre (52k) + MMSci-Ins	41.13	72.92	3.24	49.50	19.68
MMTab-Pre (150k) + MMSci-Ins	40.75	72.73	3.16	49.08	19.30
MM-Pre (202k) + MMSci-Ins	<u>42.10</u>	<u>73.98</u>	<u>3.29</u>	49.96	<u>20.85</u>
w/o MM-Pre (202k)	41.71	70.90	3.29	48.02	20.07

Table 1: Performance comparison on MMSci-Eval and held-out tabular numerical reasoning datasets. MM-Pre (202k) indicates the combination of MMTAB-Pre (150k) and MMSci-Pre (52k). w/o MM-Pre represents only training with MMSci-Ins dataset. Best results are in **bold**, second best are underlined.

et al., 2024). TQA contains TABMWP (Lu et al., 2023b), WTQ (Pasupat and Liang, 2015), HiTab (Cheng et al., 2022), TAT-QA (Zhu et al., 2021), and FeTaQA (Nan et al., 2022), where TABMWP and TAT-QA specifically focus on tabular numerical reasoning. TFV contains TabFact (Chen et al., 2020a) and InfoTabs (Gupta et al., 2020), while Table-to-Text (T2T) generation uses HiTab_T2T (Cheng et al., 2022), Rotowire (Wiseman et al., 2017), and WikiBIO (Lebret et al., 2016). While these datasets contain tables from Wikipedia, financial reports, and government documents, our MMSci-Eval datasets primarily feature scientific tables with numerical values from research papers. We use accuracy and BLEU (Papineni et al., 2002) for TQA, TFV, and T2T benchmarks.

5 Results and Analysis

5.1 Performance on Numerical Reasoning Datasets

The experimental results demonstrate the effectiveness of our proposed approach across various multimodal table understanding tasks. As shown in Table 1, we compare our method with state-of-the-art baselines on both MMSci benchmarks (TQA, TFV, T2T) and held-out tabular numerical reasoning datasets (TABMWP, TAT-QA). Among the baseline models, GPT-4V (OpenAI,

2023) achieves superior performance across all tasks, establishing strong benchmarks with 53.13% accuracy on TQA, 78.01% on TFV, and notably strong generalisation ability on held-out numerical reasoning datasets. Large-scale open-sourced models like InternVL-2-76B (Chen et al., 2024) and Qwen-2-VL-72B (Wang et al., 2024) also demonstrate competitive performance but show relatively weaker generalisation to held-out numerical reasoning datasets.

As for our approaches, with LLaVA-NeXT-7B as the foundation model, we observe that training with MMSci-Pre (52k) dataset demonstrates higher performance (17.72% on TQA, 57.12% on TFV) compared to training with MMTAB-Pre (150k) dataset (15.79% on TQA, 56.16% on TFV). The combination of both table structure learning dataset (MM-Pre 202k) further improves performance to 23.02% on TQA and 58.57% on TFV. Notably, our approach shows strong generalisation ability on held-out datasets, achieving 49.72% on TABMWP with the experiment setting of MM-Pre (202k) + MMSci-Ins.

With Qwen2-VL-7B-Instruct as the foundation model, we observe significantly stronger performance across all settings. Training with MMSci-Pre (52k) + MMSci-Ins achieves comparable or better performance (41.13% on TQA, 72.92% on TFV) compared to training with MMTAB-Pre (150k) + MMSci-Ins (40.75% on TQA, 72.73%

Method	TQA					TFV				T2T			
	TABMWP Acc.	WTQ Acc.	HiTab Acc.	TAT-QA Acc.	FeTaQA BLEU	Avg. TQA Acc.	TabFact Acc.	InfoTabs Acc.	Avg. TFV Acc.	HiTab_T2T BLEU	Rotowire BLEU	WikiBIO BLEU	Avg. T2T BLEU
Baseline													
GPT-4V (OpenAI, 2023)	60.50	48.00	27.50	32.50	11.04	35.91	45.50	<u>65.60</u>	55.55	2.98	4.23	1.94	3.05
Qwen2-VL-7B-Ins. (Wang et al., 2024)	34.44	12.55	3.36	16.19	11.75	15.66	20.28	34.19	27.23	1.90	2.30	2.94	2.38
LLaVA-NeXT-7B (Li et al., 2024)	1.73	0.00	0.00	0.00	1.17	0.58	1.24	1.78	1.51	0.45	1.04	0.67	0.72
Table-LLaVA-7B (Zheng et al., 2024)	57.78	18.43	10.09	12.82	<u>25.60</u>	24.94	<u>59.85</u>	65.26	<u>62.56</u>	<u>9.74</u>	10.46	9.68	9.96
Table-LLaVA-13B (Zheng et al., 2024)	<u>59.77</u>	<u>20.41</u>	<u>10.85</u>	15.67	28.03	<u>26.95</u>	65.00	66.91	65.96	10.40	<u>8.83</u>	<u>9.67</u>	<u>9.63</u>
Ours (LLaVA-NeXT-7B)													
MMSci-Pre (52k) + MMSci-Ins	8.76	3.22	0.63	0.39	5.99	3.80	35.78	25.37	30.57	1.57	1.10	1.78	1.48
MMTab-Pre (150k) + MMSci-Ins	9.00	2.62	0.63	0.26	7.23	3.95	36.22	26.91	31.56	1.64	0.84	1.57	1.35
MM-Pre (202k) + MMSci-Ins	10.66	4.83	0.82	0.65	9.39	5.27	39.63	27.63	33.63	1.13	0.83	1.90	1.29
w/o MM-Pre (202k)	9.69	2.74	0.19	0.39	6.84	3.97	31.72	23.80	27.76	1.69	0.79	1.53	1.34
Ours (Qwen2-VL-7B-Ins.)													
MMSci-Pre (52k) + MMSci-Ins	49.51	18.74	4.95	19.69	12.89	21.15	37.93	45.33	41.63	0.75	2.81	2.69	2.08
MMTab-Pre (150k) + MMSci-Ins	49.09	18.95	4.63	19.30	9.77	20.35	40.00	46.56	43.28	0.91	1.26	2.89	1.69
MM-Pre (202k) + MMSci-Ins	46.97	19.73	4.38	<u>20.85</u>	12.34	20.85	39.99	45.96	42.97	0.96	1.32	2.60	1.63
w/o MM-Pre (202k)	48.02	18.67	5.33	20.08	12.58	20.94	33.53	44.93	39.23	0.71	2.76	2.70	2.06

Table 2: Performance comparison on MMTAB held-out datasets. Best results are in bold, second best are underlined.

Models	MMSci-Eval			Held-out	
	TQA	TFV	T2T	TABMWP	TAT-QA
Ours (LLaVA-NeXT-7B)					
MMSci-Pre (52k) + MMSci-Ins	17.72	57.12	2.93	49.47	10.46
w/o Reasoning	10.75	42.73	2.16	42.50	7.68
MMTab-Pre (150k) + MMSci-Ins	15.79	56.16	2.88	43.55	8.03
w/o Reasoning	9.58	50.31	1.93	42.50	7.42
MM-Pre (202k) + MMSci-Ins	23.02	58.57	2.36	49.72	12.27
w/o Reasoning	12.73	45.21	2.16	46.50	19.68
w/o MM-Pre (202k)	15.22	51.73	2.86	46.24	7.63
w/o Reasoning	9.43	42.31	2.36	45.50	8.39
Ours (Qwen2-VL-7B-Ins.)					
MMSci-Pre (52k) + MMSci-Ins	41.13	72.92	3.24	49.50	19.68
w/o Reasoning	35.06	66.47	3.14	44.08	16.72
MMTab-Pre (150k) + MMSci-Ins	40.75	72.73	3.16	49.08	19.30
w/o Reasoning	34.48	66.28	2.27	43.97	16.07
MM-Pre (202k) + MMSci-Ins	42.10	73.98	3.29	49.96	20.85
w/o Reasoning	35.45	67.43	1.97	46.34	17.68
w/o MM-Pre (202k)	41.71	70.90	3.29	48.02	20.07
w/o Reasoning	34.44	62.90	3.18	44.60	14.68

Table 3: Ablation study results for reasoning steps on MMSci-Eval and held-out datasets.

on TFV), despite using only one-third of the table structure learning data. The experiment setting of training with MM-Pre (202k) + MMSci-Ins achieves the best performance with 42.10% accuracy on TQA and 73.98% on TFV, while also demonstrating strong generalisation ability on held-out numerical reasoning datasets (49.96% on TABMWP and 20.85% on TAT-QA).

These results demonstrate that our proposed MMSci-Pre dataset with 52K scientific domain-specific data is more effective than MMTAB-Pre with 150K general-domain data, highlighting the importance of data quality over quantity. Furthermore, Qwen2-VL-7B-Instruct consistently outperforms LLaVA-NeXT-7B across all experimental settings, suggesting its stronger capability in table understanding and numerical reasoning tasks. Besides, our approach shows strong generalisation to held-out tabular numerical reasoning datasets, demonstrating enhanced general ability in multi-modal table understanding and reasoning.

While the absolute performance (42.10% on TQA) may appear modest, it represents substantial improvement for the challenging task of scientific table reasoning. Similar to early work in other complex domains, these results establish important baselines that future research can build upon.

5.2 Performance on Held-out MMTAB Benchmarks

The experimental results in Table 2 also demonstrate the effectiveness and generalisation ability of our proposed approach across various held-out MMTAB benchmark. GPT-4V (OpenAI, 2023) show strong performance across all tasks, achieving 35.91% average accuracy on TQA, 55.55% on TFV, and 3.05 BLEU on T2T. Table-LLaVA models, which are specifically trained on MMTAB-Ins dataset, demonstrate competitive performance. Table-LLaVA-13B achieves strong results on TFV (65.96% average accuracy) and T2T (9.63 BLEU) while Table-LLaVA-7B shows robust performance on TABMWP (57.78%).

As for our approaches, with LLaVA-NeXT-7B as the foundation model, we observe that training with MMSci-Pre (52k) and MMSci-Ins, despite not being trained on MMTAB-Ins dataset (Zheng et al., 2024), demonstrates promising generalisation ability. The MMSci-Pre (52k) + MMSci-Ins combination achieves 3.80% average accuracy on TQA and 30.57% on TFV with only scientific domain data. The combination of both table structure learning datasets (MM-Pre 202k) further improves performance across all metrics, reaching 5.27% on TQA and 33.63% on TFV. As for Qwen2-VL-7B-Instruct as the foundation model, we observe significantly stronger generalisation capa-

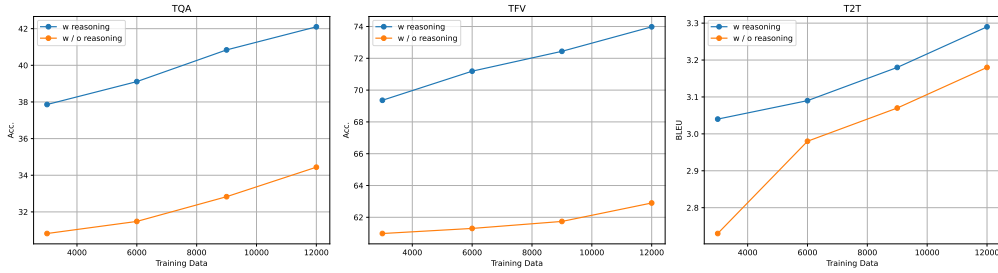


Figure 2: Performance scaling with increasing instruction tuning data size on three MMSci tasks.

456 bility. MMSci-Pre (52k) + MMSci-Ins combina- 496
 457 tion achieves 21.15% average accuracy on TQA 497
 458 and 41.63% on TFV, demonstrating strong zero- 498
 459 shot transfer to MMTab benchmark despite using 499
 460 only scientific domain data (MMSci dataset). This 500
 461 performance is particularly impressive when com- 501
 462 pared to MMTab-Pre (150k) + MMSci-Ins combi- 502
 463 nation, which uses three times more image-to- 503
 464 HTML data. Even without any table structure learn- 504
 465 ing (w/o MM-Pre), our approach achieves competi-
 466 tive results, highlighting the effectiveness of our
 467 MMSci-Ins instruction tuning dataset.

468 These results empirically demonstrate our 506
 469 MMSci-Pre dataset with 52K scientific domain- 507
 470 specific data achieves comparable or better per- 508
 471 formance than MMTab-Pre with 150K general- 509
 472 domain data in MMTab held-out benchmark, 510
 473 highlighting the importance of scientific domain- 511
 474 specific tables. Even without MMTab table struc- 512
 475 ture learning data, our approach demonstrates 513
 476 strong generalisation ability, particularly evident 514
 477 in the performance of MMSci-Pre (52k) + MMSci-Ins 515
 478 and w/o MM-Pre experiment settings. 516

479 5.3 Ablation Study on Reasoning Steps

480 We evaluate the effectiveness of reasoning steps 517
 481 across different experiment configurations. As 518
 482 shown in Table 3, Qwen2-VL-7B-Instruct demon- 519
 483 strates superior performance across all configura- 520
 484 tions. Without reasoning steps, the model train- 521
 485 ing with MMSci-Pre (52k) + MMSci-Ins achieves 522
 486 better results than that with MMTab-Pre (150k) 523
 487 + MMSci-Ins, highlighting the importance of 524
 488 domain-specific table structure learning over data 525
 489 quantity. Adding reasoning steps consistently im- 526
 490 proves performance across all metrics, with the 527
 491 model reaching its peak performance under the 528
 492 MM-Pre (202k) + MMSci-Ins experiment configu- 529
 493 ration. Similar trends are observed in LLaVA- 530
 494 NeXT-7B, though with lower overall performance. 531
 495 These patterns extend to held-out tabular numer-

ical reasoning datasets, where both models show 496
 strong generalisation capabilities with reasoning 497
 steps, especially on numerical reasoning tasks like 498
 TABMWP and TAT-QA. The results demonstrate 499
 that a smaller amount of scientific domain-specific 500
 table structure learning data, combined with ex- 501
 plicit reasoning steps, can be more effective than 502
 larger-scale general domain table structure learn- 503
 ing. 504

505 5.4 Impact of Training Data Size

506 As shown in Figure 2, we compare performance of 507
 508 MLLMs training with MM-Pre (202k) + MMSci- 509
 510 Ins across three MMSci tasks (TQA, TFV, T2T) 511
 512 with instruction tuning data (MMSci-Ins dataset) 513
 514 size increasing from 3K to 12K samples. The 514
 515 findings demonstrate consistent advantages of in- 516
 517 corporating reasoning steps across all data scales. 517
 518 Models trained with reasoning steps maintain sub- 518
 519 stantial performance advantages across all tasks 519
 520 (7-8% for TQA, 8-10% for TFV, 0.3-0.4 BLEU for 520
 521 T2T). While both variants benefit from increased 521
 522 training data, models with reasoning steps show 522
 523 stronger scaling behavior, particularly in TQA and 523
 TFV tasks. The persistent performance gap across 524
 all data sizes suggests that reasoning steps provide 525
 fundamental improvements in model learning that 526
 cannot be simply achieved through increased train- 527
 ing data alone. 528

529 5.5 Representational Alignment Analysis

529 In this section, we conduct an in-depth analysis to 530
 531 assess the language-vision alignment from the per- 531
 532 spective of the representation space. This analysis 532
 533 aims to provide further insights into the observed 533
 534 variations in model performance, particularly in the 534
 context of scientific multimodal table understand-
 ing and reasoning tasks.

Preliminaries. We formalise MLLMs within the 532
 framework of an *unembedding-embedding* archi- 533
 tecture. In this framework, the unembedding stage 534

Models	Cycle KNN	Mutual KNN	Lcs KNN	CKA	CKNNA	SVCCA	Edit KNN
<i>Unembedding stage: ImageNet(Concepts)</i>							
Random	0.02761	0.01257	0.52355	0.08614	0.00714	0.12425	0.00019
Qwen2-VL-7B-Ins.	0.68110	<u>0.03486</u>	<u>1.28153</u>	<u>0.08856</u>	0.03067	0.14318	0.00112
Llama3.2-11B-Vision-Ins.	<u>0.08608</u>	0.04205	1.52788	0.06079	0.01403	0.11651	0.00061
LLaVA-NeXT-7B	0.57173	0.02077	0.81645	0.08024	0.01577	0.13240	0.00037
Phi3.5-Vision-Ins.	0.02761	0.01257	0.52355	0.08614	0.00714	0.12118	0.00019
InternVL2-8B	0.08175	0.01637	0.72495	0.09185	0.00062	0.12148	0.00044
<i>Unembedding stage: Wikipedia Caption (short descriptive sentences)</i>							
Qwen2-VL-7B-Ins.	<u>0.49414</u>	0.06855	2.05078	0.08876	0.04093	0.20229	0.00175
Llama3.2-11B-Vision-Ins.	0.31347	0.03623	1.29980	0.00968	0.00779	<u>0.22120</u>	0.00050
LLaVA-NeXT-7B	0.57813	0.03935	1.36523	<u>0.07933</u>	<u>0.03998</u>	0.23114	<u>0.00082</u>
Phi3.5-Vision-Ins.	0.04980	0.03027	1.14843	<u>0.01669</u>	0.03890	0.18183	0.00066
InternVL2-8B	0.36914	0.04132	<u>1.55761</u>	0.04732	0.01658	0.21739	0.00093
<i>Embedding stage: MMSci T2T tasks (table to text description)</i>							
Qwen2-VL-7B-Ins.	0.38631	0.06726	<u>2.03660</u>	0.19318	<u>0.05514</u>	0.38461	0.00183
Llama3.2-11B-Vision-Ins.	0.31310	0.02200	0.84007	1.73979e-8	0.03208	0.08180	0.00026
LLaVA-NeXT-7B	0.38246	0.04514	1.49325	0.15203	0.06673	<u>0.28857</u>	0.00109
Phi3.5-Vision-Ins.	<u>0.38053</u>	<u>0.06712</u>	2.12909	<u>0.16121</u>	0.03688	0.26982	<u>0.00127</u>
InternVL2-8B	0.36512	0.04651	1.56647	0.04230	0.02675	0.11876	0.00096

Table 4: Kernel alignment analysis. The representation for each sample is the averaged token embeddings. The best two values are shown in **bold** and underlined.

is responsible for learning transformations between observations (e.g., text, vision) and latent spaces through encoders, while the embedding stage captures the complex interactions among latent variables within the latent space of LLMs’ hidden layers. Each stage serves distinct functions and yields representations with different properties (Park et al., 2024). Consequently, by focusing on each stage independently, we can have a systematical evaluation of model behaviours in representation spaces. To assess the representational alignment between vision-language modalities at each stage, we next measure the geometrical similarity between them via the *kernel*.

Kernels, characterising the distance metrics between points in a representation space, are commonly used to assess vector space (Huh et al., 2024). Typically, the more similarity between two kernels derived from different spaces (text or vision) indicates a higher degree of alignment between those modality spaces. This similarity can be quantified via *kernel-alignment metrics*, such as Centered Kernel Distance (CKA) (Kornblith et al., 2019). For more information about kernel-alignment metrics used in the experiment, we refer to Huh et al. (2024) for a deep understanding.

Quantitative evaluation. For the unembedding stage, we specifically choose two language-vision datasets: ImageNet (Deng et al., 2009) and Wikipedia Caption (WIT) (Srinivasan et al., 2021). We randomly select 2048 samples from each dataset. These datasets offer varying levels of fine granularity in language-vision alignment, enabling a comprehensive assessment of representational performance. As illustrated in Table 4, we can ob-

serve that the Qwen2-VL-7B-Instruct can generally outperform other baselines on both datasets, indicating it has better fine-grained alignment between language and vision. In the embedding stage, we evaluate alignment on the MMSci T2T task. Since some models do not support single-modality input, we utilise a reference language model (e.g., openllama-7B (Geng and Liu, 2023)) as the text encoder and MLLMs as the image encoder with prompt “please describe the table”. Alignment is measured based on the output embedding from the last hidden layer. As shown in Table 4, Qwen2-VL-7B-Instruct outperforms the other models, demonstrating its superior language-vision alignment capability. This segment of the experiment demonstrates that the Qwen2-VL-7B-Instruct model exhibits superior language-vision alignment within the representation space. This finding is consistent with the cross-modal consistency analysis presented in Appendix C.1, where we evaluate different table information modalities as inputs to MLLMs and assess their cross-modal consistency (i.e., the proportion of identical predictions) on TQA and TFV tasks.

6 Conclusion

In this paper, we introduce a comprehensive framework for multimodal scientific table understanding and reasoning with dynamic input image resolutions. Experimental results validate our framework’s effectiveness across different model architectures, showing consistent improvements in both general table understanding and numerical reasoning capabilities, with strong generalisation to held-out datasets.

604 Limitations

605 While this work advances scientific multimodal
606 table understanding and reasoning, several limita-
607 tions remain for future research. First, our frame-
608 work primarily focuses on scientific tables contain-
609 ing numerical values, while other types of scientific
610 tables (e.g., qualitative comparison tables, method-
611 ology tables) are not extensively covered. Second,
612 though our framework demonstrates strong perform-
613 ance on numerical reasoning tasks, the current
614 approach may still struggle with complex statistical
615 analyses and domain-specific mathematical nota-
616 tions that are common in scientific literature. Third,
617 while our models support dynamic input resolu-
618 tions, processing extremely large tables with dense
619 information remains challenging due to computa-
620 tional constraints and potential information loss
621 during visual encoding.

622 Ethical Statement

623 The MMSci datasets are constructed from publicly
624 available scientific papers and their associated ta-
625 bles, primarily sourced from open-access reposi-
626 tories and academic databases with appropriate
627 licenses. All table images are generated through au-
628 tomated scripts from the original scientific papers,
629 maintaining their integrity while ensuring proper
630 attribution. The instruction tuning samples are cre-
631 ated based on the original scientific context, pre-
632 serving the academic nature of the source material.
633 Our framework is designed to assist in scientific
634 research by improving the accessibility and under-
635 standing of tabular data in academic literature. We
636 anticipate that this work will contribute positively
637 to the research community by facilitating more
638 efficient analysis of scientific publications. The
639 code and datasets are made publicly available for
640 research purposes, promoting transparency and re-
641 producibility in the field of multimodal scientific
642 table understanding.

643 References

644 Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna,
645 Baptiste Bout, Devendra Chaplot, Jessica Chud-
646 novsky, Diogo Costa, Baudouin De Monicault,
647 Saurabh Garg, Theophile Gervet, et al. 2024. Pixtral
648 12b. *arXiv preprint arXiv:2410.07073*.

649 Iñigo Alonso, Eneko Agirre, and Mirella Lapata. 2024.
650 *PixT3: Pixel-based table-to-text generation*. In *Pro-
651 ceedings of the 62nd Annual Meeting of the Associa-
652 tion for Computational Linguistics (Volume 1: Long*

Papers), pages 6721–6736, Bangkok, Thailand. As-
653 sociation for Computational Linguistics. 654

655 Pei Chen, Soumajyoti Sarkar, Leonard Lausen, Balasub-
656 ramaniam Srinivasan, Sheng Zha, Ruihong Huang,
657 and George Karypis. 2023. Hytrel: Hypergraph-
658 enhanced tabular data representation learning. In
659 *Thirty-seventh Conference on Neural Information
660 Processing Systems*.

661 Wenhu Chen. 2023. *Large language models are few(1)-
662 shot table reasoners*. In *Findings of the Associa-
663 tion for Computational Linguistics: EACL 2023*,
664 pages 1120–1130, Dubrovnik, Croatia. Association
665 for Computational Linguistics.

666 Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai
667 Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and
668 William Yang Wang. 2020a. *Tabfact: A large-scale
669 dataset for table-based fact verification*.

670 Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong,
671 Hong Wang, and William Yang Wang. 2020b. *Hy-
672 bridQA: A dataset of multi-hop question answering
673 over tabular and textual data*. In *Findings of the Asso-
674 ciation for Computational Linguistics: EMNLP 2020*,
675 pages 1026–1036, Online. Association for Computa-
676 tional Linguistics.

677 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo
678 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,
679 Xizhou Zhu, Lewei Lu, et al. 2024. Internvl: Scal-
680 ing up vision foundation models and aligning for
681 generic visual-linguistic tasks. In *Proceedings of
682 the IEEE/CVF Conference on Computer Vision and
683 Pattern Recognition*, pages 24185–24198.

684 Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena
685 Shah, Iana Borova, Dylan Langdon, Reema Moussa,
686 Matt Beane, Ting-Hao Huang, Bryan Routledge, and
687 William Yang Wang. 2021. *FinQA: A dataset of nu-
688 merical reasoning over financial data*. In *Proceedings
689 of the 2021 Conference on Empirical Methods in Nat-
690 ural Language Processing*, pages 3697–3711, Online
691 and Punta Cana, Dominican Republic. Association
692 for Computational Linguistics.

693 Zhoujun Cheng, Haoyu Dong, Zhiruo Wang, Ran Jia,
694 Jiaqi Guo, Yan Gao, Shi Han, Jian-Guang Lou, and
695 Dongmei Zhang. 2022. *HiTab: A hierarchical table
696 dataset for question answering and natural language
697 generation*. In *Proceedings of the 60th Annual Meet-
698 ing of the Association for Computational Linguistics
699 (Volume 1: Long Papers)*, pages 1094–1110, Dublin,
700 Ireland. Association for Computational Linguistics.

701 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
702 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
703 Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion
704 Stoica, and Eric P. Xing. 2023. *Vicuna: An open-
705 source chatbot impressing gpt-4 with 90%* chatgpt
706 quality*.

707 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,
708 and Li Fei-Fei. 2009. *Imagenet: A large-scale hier-
709 archical image database*. In *2009 IEEE Conference*

710		on <i>Computer Vision and Pattern Recognition</i> , pages 248–255.	765
711			766
712	Alexey Dosovitskiy. 2020.	An image is worth 16x16 words: Transformers for image recognition at scale. <i>arXiv preprint arXiv:2010.11929</i> .	767
713			768
714			769
715	Xinyang Geng and Hao Liu. 2023.	Openllama: An open reproduction of llama .	770
716			771
717	Vivek Gupta, Maitrey Mehta, Pegah Nokhiz, and Vivek Srikumar. 2020.	INFOTABS: Inference on tables as semi-structured data . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2309–2324, Online. Association for Computational Linguistics.	772
718			773
719			774
720			775
721			776
722			777
723	Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Martin Eisenschlos. 2020.	Tapas: Weakly supervised table parsing via pre-training . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 4320–4333.	778
724			779
725			780
726			781
727			782
728			783
729	Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. 2024.	Position: The platonic representation hypothesis . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 20617–20642. PMLR.	784
730			785
731			786
732			787
733			788
734			789
735	Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. 2019.	Similarity of neural network representations revisited. In <i>International conference on machine learning</i> , pages 3519–3529. PMLR.	790
736			791
737			792
738			793
739			794
740	Rémi Lebre, David Grangier, and Michael Auli. 2016.	Neural text generation from structured data with application to the biography domain . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1203–1213, Austin, Texas. Association for Computational Linguistics.	795
741			796
742			797
743			798
744			799
745			800
746	Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023.	Pix2Struct: Screenshot Parsing as Pretraining for Visual Language Understanding . ArXiv:2210.03347 [cs].	801
747			802
748			803
749			804
750			805
751			806
752	Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024.	Llava-next: Stronger llms supercharge multimodal capabilities in the wild .	807
753			808
754			809
755			810
756	Hongxin Li, Jingran Su, Yuntao Chen, Qing Li, and Zhaoxiang Zhang. 2023a.	Sheetcopilot: Bringing software productivity to the next level through large language models . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	811
757			812
758			813
759			814
760			815
761	Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman, Dongmei Zhang, and Surajit Chaudhuri. 2023b.	Table-gpt: Table-tuned gpt for diverse table tasks .	816
762			817
763			818
764			819
	Qian Liu, Bei Chen, Jiaqi Guo, Morteza Ziyadi, Zeqi Lin, Weizhu Chen, and Jian-Guang Lou. 2022.	TAPEX: Table pre-training via learning a neural SQL executor . In <i>International Conference on Learning Representations</i> .	765
			766
			767
			768
			769
	Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023a.	Chameleon: Plug-and-play compositional reasoning with large language models. In <i>The 37th Conference on Neural Information Processing Systems (NeurIPS)</i> .	770
			771
			772
			773
			774
			775
	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. In <i>International Conference on Learning Representations (ICLR)</i> .	776
			777
			778
			779
			780
			781
	Meta. 2024.	Llama 3.2: Pushing the boundaries of vision and language for edge and mobile devices . Accessed: 2024-12-06.	782
			783
			784
	Nafise Sadat Moosavi, Andreas Rücklé, Dan Roth, and Iryna Gurevych. 2021.	Scigen: a dataset for reasoning-aware text generation from scientific tables. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	785
			786
			787
			788
			789
			790
	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, and Dragomir Radev. 2022.	Fetaqa: Free-form table question answering. <i>Transactions of the Association for Computational Linguistics</i> , 10:35–49.	791
			792
			793
			794
			795
			796
			797
			798
	OpenAI. 2023.	Gpt-4v. https://openai.com/index/gpt-4v-system-card/ . Accessed: 2023-02-09, 2023-02-11, 2023-02-12.	799
			800
			801
	OpenAI. 2024.	Hello gpt-4o. https://openai.com/index/hello-gpt-4o/ . Accessed: 2024-02-09, 2024-02-11, 2024-02-12.	802
			803
			804
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002.	Bleu: a method for automatic evaluation of machine translation. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	805
			806
			807
			808
	Kiho Park, Yo Joong Choe, and Victor Veitch. 2024.	The linear representation hypothesis and the geometry of large language models . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 39643–39666. PMLR.	809
			810
			811
			812
			813
			814
	Panupong Pasupat and Percy Liang. 2015.	Compositional semantic parsing on semi-structured tables . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language</i>	815
			816
			817
			818
			819

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](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3277–3287, Online. Association for Computational Linguistics.

A Details about MMSci 943

A.1 Datasets Statistics 944

Table 5 presents the distribution of reasoning types in our MMSci-Eval dataset. The most common type is addition (21.1%), followed by subtraction (15.3%) and max/min operations (15.7%). Division and comparison operations also appear frequently (14.2% and 13.7% respectively). More complex operations like ranking (9.6%) and look-up (8.9%) occur less frequently, while domain knowledge calculations are rare (1.5%). 945 946 947 948 949 950 951 952 953

The average number of reasoning steps varies significantly across types, with subtraction requiring the most steps (4.1) and look-up operations requiring the fewest (1.5). This variation reflects the inherent complexity of different mathematical operations and their application to tabular data. Notably, even seemingly simple operations like addition require multiple steps (2.8) on average, indicating the complexity of reasoning with tabular scientific data. 954 955 956 957 958 959 960 961 962 963

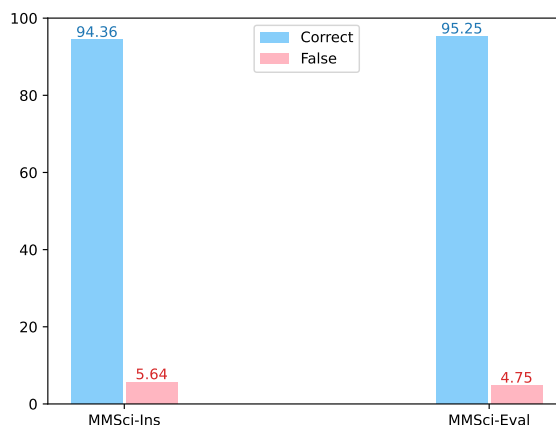


Figure 3: Evaluation of generated data of MMSci-Ins and MMSci-Eval dataset. Correct refers to the data verified correctly by human annotators.

A.2 Dataset Quality Control 964

To ensure data quality, we conduct a rigorous human verification process for both MMSci-Ins and MMSci-Eval datasets. For MMSci-Ins, we manually verify 40% of the generated samples, achieving a high accuracy rate of 94.36%. For MMSci-Eval, given its critical role as a benchmark, we carefully examine all 3,114 generated samples and achieve an accuracy of 95.25%. For any identified incorrect samples, we employ GPT-4o to regenerate them following the same self-consistency voting mechanism, followed by another round of both automatic 965 966 967 968 969 970 971 972 973 974 975

Reasoning Type	Description	Avg. Reasoning Step	Prop.%
Add	Calculate the sum between numbers	2.8	21.1
Comparison	Comparison of values	2.1	13.7
Domain Knowledge Calculation	Calculations need domain knowledge	2.2	1.5
Divide	Perform division between numbers	3.4	14.2
Look Up	Search for cells in tables	1.5	8.9
Max/Min	Retrieve the maximum or minimum number	3.2	15.7
Ranking	Arranges items in a specific order	2.4	9.6
Subtract	Perform subtraction between numbers	4.1	15.3

Table 5: The reasoning types, descriptions, average reasoning step, and proportion in our dataset.

and manual verification to ensure quality. This iterative process ensures the reliability and correctness of our datasets for both training and evaluation purposes.

A.3 Prompt for Generating Data

The prompt for MMSci-Ins and MMSci-Pre data generation is shown in Table 7.

B Experimental Settings

Implementation Details. Both models follow a three-component design. Qwen2-VL-7B-Instruct consists of a Vision Transformer (ViT) (Dosovitskiy, 2020) as the vision tower, a MLP as the vision-language connector, and Qwen2 (Wang et al., 2024) as the language model. LLaVA-NeXT-7B uses a pre-trained CLIP model (Radford et al., 2021) as the visual encoder, a MLP connector, and Vicuna-7B (Chiang et al., 2023) as the backbone. In both architectures, the vision encoder processes images into visual features, which are projected into the LLM’s word embedding space via the MLP connector.

Training Details. All experiments are conducted on $4 \times A100$ 80GB GPUs using LoRA with rank 64 and sequence length 4096. For table structure learning, LLaVA-NeXT-7B requires 15 hours for MMTAB-Pre (150k), 3 hours for MMSci-Pre (52k), and 20 hours for combined training (one epoch). Qwen2-VL-7B takes 15 hours, 8 hours, and 19 hours respectively. The instruction tuning stage requires approximately 1 hour for 4 epochs with 12k samples for both models.

C More Experimental Results and Analysis

C.1 Vision-Language Consistency Analysis

We evaluate the cross-modal consistency of different MLLMs by comparing their performance

when processing table information through different modalities. For each model, we test with both table images (image modality) and their textual representations (text modality), measuring both task performance (Acc.) and cross-modal consistency (Consis.).

Qwen2-VL-7B-Instruct demonstrates superior cross-modal alignment, achieving the highest consistency scores on both TQA (60.40%) and TFV (72.48%) tasks. Notably, it maintains strong performance across both modalities, with image-based accuracy (TQA: 39.11%, TFV: 52.79%) consistently outperforming text-based results (TQA: 21.65%, TFV: 50.10%). This suggests robust integration of visual and textual understanding capabilities.

Other models show varying degrees of modality gap. MiniCPM-V-2.6-8B and InternVL-2-8B achieve moderate consistency (48.78% and 50.89% on TQA), while models like LLaVA-NeXT-7B and Pixtral-12B show significant disparities between modalities, resulting in lower consistency scores. These results highlight the challenge of maintaining consistent reasoning capabilities across different input modalities in table understanding tasks.

D Dataset Examples and Case Study

D.1 Dataset Examples

The training examples of MMSci-Pre Dataset are shown in Figure 4. The examples of MMSci-Eval are shown Figure 6.

D.2 Case Study

The case study of TQA and TFV tasks are shown in Figure 6 and Figure 7, respectively

Model	Modal	TQA		TFV	
		Acc.	Consis.	Acc.	Consis.
Qwen2-VL-7B-Ins. (Wang et al., 2024)	Text	21.65	60.40	50.10	72.48
	Image	39.11		52.79	
LLaVA-NeXT-7B (Li et al., 2024)	Text	3.17	14.81	2.03	23.65
	Image	0.19		0.86	
MiniCPM-V-2.6-8B (Yao et al., 2024)	Text	21.11	48.78	30.82	38.53
	Image	26.58		33.23	
InternVL-2-8B (Chen et al., 2024)	Text	19.84	50.89	42.87	36.42
	Image	25.72		44.99	
Pixtral-12B (Agrawal et al., 2024)	Text	1.44	16.52	4.43	29.88
	Image	0.96		5.49	
Llama-3.2-11B-Vision-Ins. (Meta, 2024)	Text	3.24	15.71	6.96	20.40
	Image	1.15		5.85	

Table 6: Vision-language consistency evaluation across different MLLMs. Consistency scores measure the percentage of identical responses between modalities, indicating the model’s cross-modal alignment.

Task

Table Image-to-HTML

Input:
Please review this table image and return a text representation of the table in HTML format.

Output:

```

<table border="1" cellspacing="0">
<tr>
<th colspan="2" rowspan="2">Model
<th colspan="3" rowspan="2">Target Corpus
<th colspan="2">Baseline Acc. (%)
<th colspan="2">Fine-tuned Acc. (%)
<th colspan="2">Gain (%)

<tr>
<th>BERT</th>
<th>Filler</th>
<th>Role</th>
<th>Acc.</th>
<th>Acc.</th>
<th>+</th>
<th>+</th>
<th>+</th>

```

Table Image

Model	Target Corpus	Transfer BERT	Transfer Filler	Transfer Role	Baseline Acc. (%)	Fine-tuned Acc. (%)	Gain (%)
BERT	QNLI	True	-	-	91.60	91.27	-0.33
BERT	QQP	True	-	-	91.45	91.12	-0.33
BERT	RTE	True	-	-	71.12	73.65	+2.53
BERT	SNLI	True	-	-	90.45	90.69	+0.24
BERT	SST	True	-	-	93.23	92.78	-0.45
HUBERT (Transformer)	QNLI	True	True	False	90.56	91.16	+0.60
HUBERT (Transformer)	QQP	False	False	True	90.81	91.42	+0.61
HUBERT (Transformer)	RTE	True	True	True	61.73	74.01	+12.28
HUBERT (Transformer)	SNLI	True	False	True	90.66	91.36	+0.70
HUBERT (Transformer)	SST	True	False	True	91.28	92.43	+1.15

Table 3: Transfer learning results for GLUE tasks. The source corpus is MNLI. Baseline accuracy is when Transfer BERT, Filler, and Role are all False, equivalent to no transfer. Fine-tuned accuracy is the best accuracy among all possible transfer options.

Output:

```

<table border="1" cellspacing="0">
<tr>
<th colspan="2" rowspan="2">Model
<th colspan="3" rowspan="2">Target Corpus
<th colspan="2">Baseline Acc. (%)
<th colspan="2">Fine-tuned Acc. (%)
<th colspan="2">Gain (%)

<tr>
<th>BERT</th>
<th>Filler</th>
<th>Role</th>
<th>Acc.</th>
<th>Acc.</th>
<th>+</th>
<th>+</th>
<th>+</th>

```

Figure 4: MMSci-Pre Dataset example

14

The prompt for Generating data

You are given a table image and a description: {description}.

1.For the Table to Text (T2T) task, come up with a one to two sentence length succinct multi-hop reasoning step of the description.

Write your results as 'T2T Reasoning:' and then the succinct reasoning step.

2.For the Table Question Answering (TQA) task, come up with a question and answer with multi-hop reasoning step.

The question and answer must be based on the table image and description.

Write your results as 'TQA Question:' and then the question and 'TQA Reasoning:' and then the reasoning step and 'TQA Answer:' and then the answer.

When generating 'TQA Question:', make sure it is a single question that requires reasoning based on the table.

When generating 'TQA Answer:', provide the final answer in the JSON structure, using the format "answer": "<YOUR ANSWER>"

Make sure the answer only contains one entity, such as 'So, the answer is "answer": "23".'

3.For the Table Fact Checking (TFV) task, come up with a statement and answer with multi-hop reasoning step.

The statement and answer must be based on the table image and description. The table 'supports' or 'refutes' the statement. The statement should be considered 'not enough info' if it may or may not be true.

Write your results as 'TFV Statement:' and then the statement and 'TFV Reasoning:' and then the reasoning step and 'TFV Answer:' and then the answer.

Make sure the answer only contains one entity, such as 'Thus, the answer is "answer": "supports".'

When generating 'TFV Answer:', provide the final answer in the JSON structure, using the format "answer": "<YOUR ANSWER>"

Fill the result into JSON format without any other words:

```
"T2T Reasoning": "<YOUR T2T REASONING>",  
"TQA Question": "<YOUR TQA QUESTION>",  
"TQA Reasoning": "<YOUR TQA REASONING>",  
"TQA Answer": "<YOUR TQA ANSWER>",  
"TFV Statement": "<YOUR TFV STATEMENT>",  
"TFV Reasoning": "<YOUR TFV REASONING>",  
"TFV Answer": "<YOUR TFV ANSWER>"
```

Examples:

```
{TQA Examples}  
{TFV Examples}  
{T2T Examples}
```

Table 7: The prompts for generating the questions, reasoning steps, and answers or claims of MMSci-Ins and MMSci-Eval datasets.

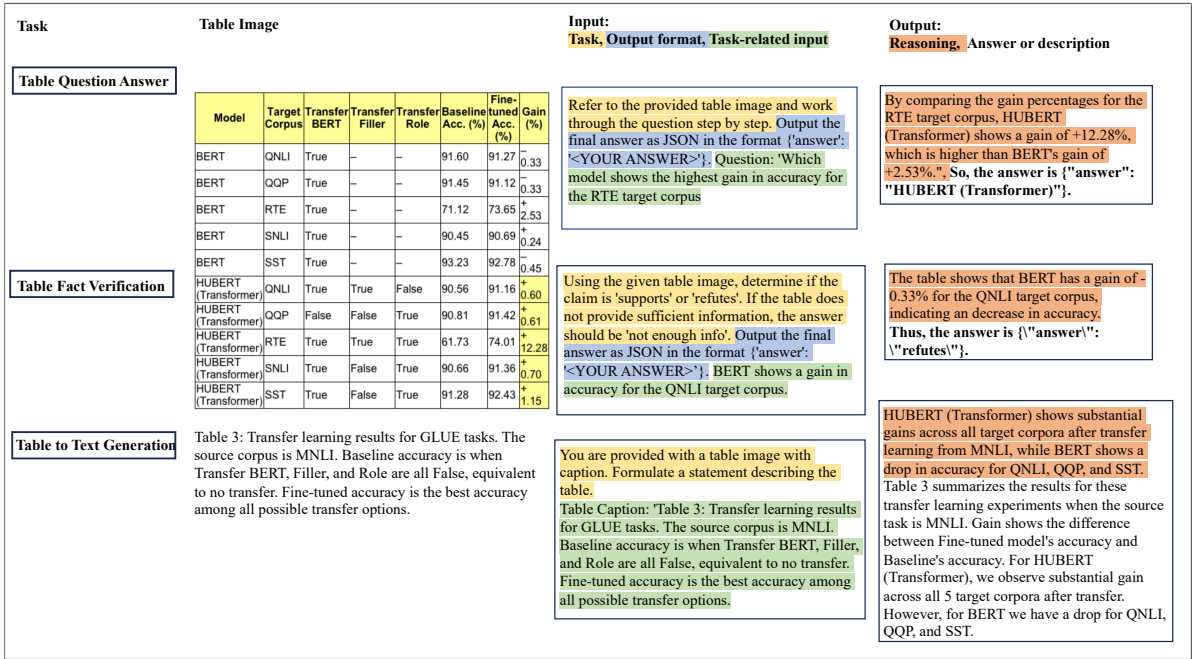


Figure 5: MMSci-Ins and MMSci-Eval Dataset example

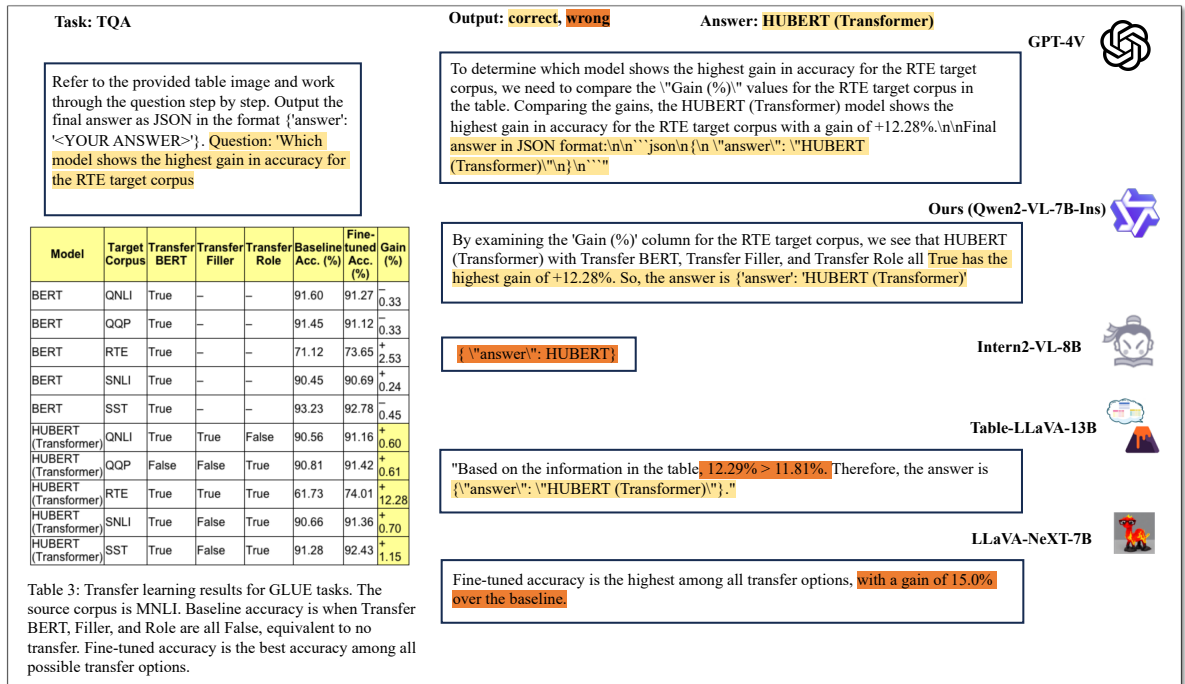


Figure 6: A case example of TQA task

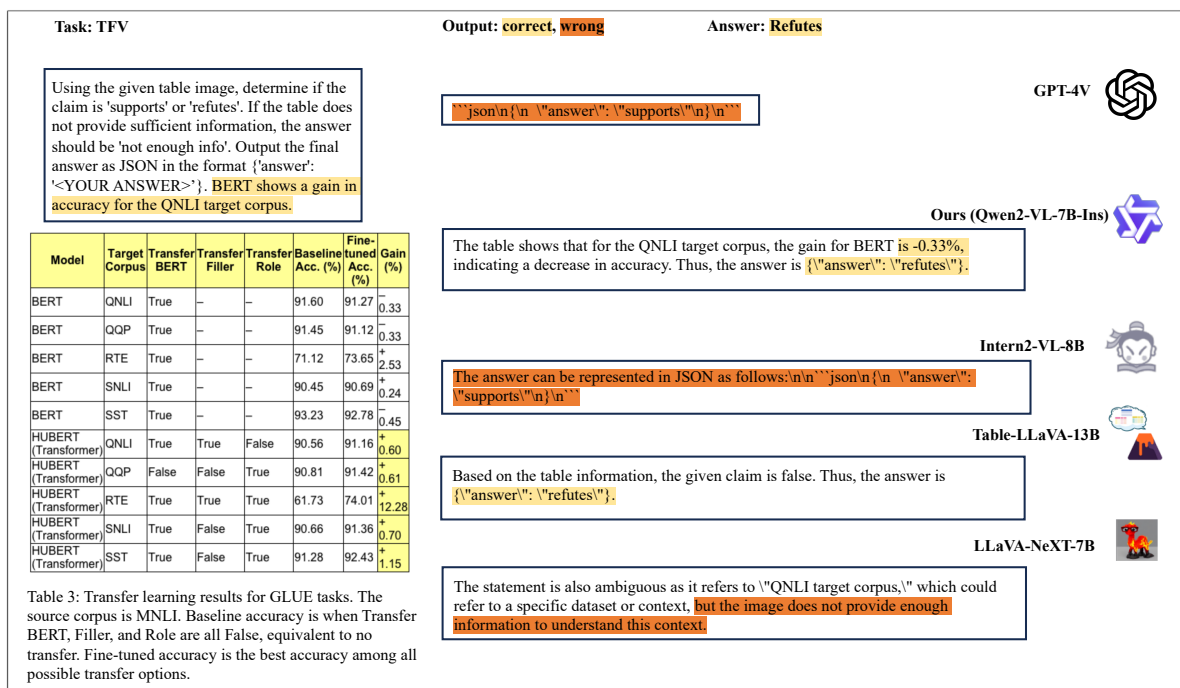


Figure 7: A case example of TFV task