

# Distill Visual Chart Reasoning Ability from LLMs to MLLMs

Anonymous ACL submission

## Abstract

Solving complex chart Q&A tasks requires advanced visual reasoning abilities in multimodal large language models (MLLMs), including recognizing key information from visual inputs and conducting reasoning over it. While fine-tuning MLLMs for reasoning is critical, collecting and annotating charts and questions is expensive, hard to scale, and often results in low-quality annotations. To address this, we propose *Code-as-Intermediary Translation* (CIT), a cost-effective, efficient and scalable data synthesis method for distilling visual reasoning abilities **from LLMs to MLLMs**. The code serves as an intermediary that translates visual chart representations into textual representations, enabling language models to understand cross-modal information and generate reasoning chains accordingly. In this way, we can employ text-based synthesizing techniques to expand chart-plotting code and generate high-quality Q&A pairs for training models. This produces **REACHQA**, a dataset containing 3k reasoning-intensive charts and 20k Q&A pairs to enhance both recognition and reasoning abilities of MLLMs. Experiments show that models fine-tuned with REACHQA not only perform well on chart-related tasks but also show performance gains on general reasoning benchmarks.

## 1 Introduction

Multimodal large language models (MLLMs) have achieved notable progress, particularly in visual recognition tasks (OpenAI, 2024a; Anthropic, 2024). However, their ability to comprehend complex images like charts in real-world contexts and to address reasoning-intensive questions remains limited compared to humans (Masry et al., 2022; Huang et al., 2024; Wang et al., 2024c). Our analysis of the error distribution in ChartQA (Figure 1) also reveals two main failure modes in current MLLMs: while most errors originate from visual misrecognition, a substantial portion arises from

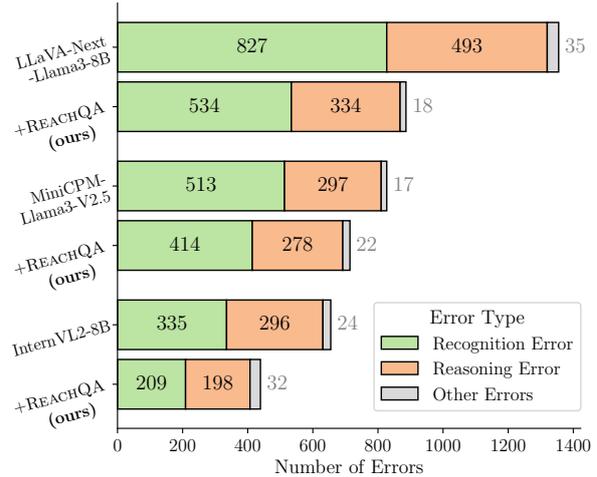


Figure 1: Error distribution of three baseline models vs. our REACHQA-trained versions on ChartQA test set (Masry et al., 2022), as judged by GPT-4o. Error types are categorized into Recognition Error, Reasoning Error, and Other Errors (question misinterpretation, factual inconsistency or hallucination, and response refusal).

flawed reasoning even when visual elements are correctly identified. This contrasts sharply with human performance (Wang et al., 2024a,c), since we can purposefully identify task-critical information from images and engage in step-by-step reasoning processes. These observations motivate our investigation into bridging this capability gap through the acquisition of human-like reasoning patterns.

While distilling expert rationales from humans or stronger models presents a promising pathway for improving reasoning abilities (Han et al., 2023; Meng et al., 2024; Masry et al., 2024a,b), constructing high-quality training data for chart-related tasks is expensive and hard to scale. Early approaches typically rely on manual chart collection from online sources, meticulous data filtering and annotation (Masry et al., 2022; Wang et al., 2024c). Recent attempts to automate Q&A generation through LLMs often use data tables as inputs (Han et al., 2023; Masry et al., 2024a), which ne-

glect the visual-semantic features of charts. Even with the use of MLLMs (Masry et al., 2024b), our preliminary study (§ 2.2) shows they also struggle to produce accurate and challenging data for advanced reasoning skill acquisition. In comparison, we find that when LLMs process charts in a better textual format—**code**, they can generate Q&A pairs at lower costs and with higher quality.

Inspired by the concept of intermediary translation (Zarechnak, 1986; Léon, 2007), which refers to using a bridge language to improve translation quality across diverse languages in literary studies, we introduce **Code-as-Intermediary Translation (CIT)**. In this method, the code acts as an intermediary, converting chart images into textual representations by faithfully encoding visual-semantic features within itself. This process enables LLMs to understand cross-modal information more accurately, thereby generating visually complex Q&A pairs with high-quality reasoning rationales. Furthermore, it facilitates the adoption of text-based instruction augmentation strategies, such as Self-Instruct (Wang et al., 2023) and Evol-Instruct (Xu et al., 2024), to expand the quantity and enhance the complexity of the synthetic charts. Starting with 33 seed codes collected from the Matplotlib gallery, we synthesize more chart-plotting codes covering diverse types and topics, and then complicate them to create richer ones. Finally, using the synthetic codes as a bridge, we generate charts (via Python) and instructions (via LLMs) in a bi-directional process, ensuring the alignment between modalities.

With the CIT method, we construct **REACHQA**, a multimodal instruction dataset containing 3, 249 reasoning-intensive charts and 19, 963 Q&A pairs, all at a remarkably low cost of just \$300. The dataset comprises questions focused on both visual recognition and reasoning, designed to address the dual challenges of current MLLMs. Additionally, we create a manually verified test set to assess models’ recognition and reasoning abilities independently. Experiments demonstrate that REACHQA-trained models achieve substantial performance gains across benchmarks, with LLaVA-Next-Llama3-8B (Li et al., 2024) improving by over 30% on average, while both types of errors are significantly reduced (Figure 1). Notably, these improvements generalize beyond chart-specific tasks to broader multimodal reasoning tasks like Math-Vista and MATH-Vision—an outcome previously unattainable with existing chart-focused datasets. Finally, we explore REACHQA’s working mecha-

nism and more features, providing actionable guidelines for building performant multimodal datasets.

Our contributions are summarized as follows:

1. We propose Code-as-Intermediary Translation (CIT), a cost-effective and efficient method for synthesizing multimodal instruction data with code as a bridge between the two modalities.
2. Through CIT, we construct REACHQA, the first fully LLM-synthesized reasoning-intensive chart Q&A dataset, focusing on both visual recognition and reasoning abilities.
3. We conduct extensive experiments and analyses to demonstrate REACHQA’s effectiveness for MLLMs, along with its strong generalization to broader multimodal reasoning tasks.

## 2 Background

### 2.1 Deficiencies in Existing Chart Datasets

Existing chart-related datasets are either collected from online data sources or generated by models, sometimes requiring manual annotation or automated question generation. Most of them focus on visual recognition tasks. While some recent works target advanced reasoning, they often struggle with scalability or other shortcomings. Table 1 summarizes these datasets, with further details below.

**Chart Properties.** The *visual diversity* is shaped by the variety of chart types and topics (Wang et al., 2024c). Early datasets like ChartQA and OpenCQA, sourced from limited websites, featured uniform styles with minimal diversity. To address this, recent works like ChartAst synthesize charts with randomized attributes (e.g., color, fonts) using LLMs. However, beyond the superficial variations in chart appearance, many of them overlook the *visual complexity* (Zeng et al., 2024). As models evolve, simple style changes no longer pose challenges. Datasets like CharXiv and MMC, which include complex scientific charts from arXiv papers, naturally exhibit greater complexity in recognition. Additionally, the textual format of charts is critical, enabling dataset expansion via language models.

**Q&A Properties.** Some benchmarks like PlotQA and ChartBench use predefined templates to generate Q&A pairs, resulting in monotonous and simplistic questions. Other datasets, such as ChartQA and CharXiv, required manual annotation, which improved quality but increased costs and

Datasets	Chart Properties				Q&A Properties			Dataset Properties		
	# Chart Type	# Chart Topic	Textual Format	Vis. Comp.	Temp. Free	Vis. Refer.	Rat. Annot.	Train Set	Test Set	Scal.
PlotQA (Methani et al., 2020)	3	-	Table	✗	✗	✓	✗	✓	✓	✗
ChartQA (Masry et al., 2022)	3	15	Table	✗	✓	✓	✗	✓	✓	✗
OpenCQA (Kantharaj et al., 2022)	5	10	Caption	✗	✓	✗	✓	✗	✓	✗
MathVista (Lu et al., 2024)	-	-	-	✗	✓	✗	✗	✗	✓	✗
CharXiv (Wang et al., 2024c)	-	-	-	✓	✗	✓	✗	✗	✓	✗
ChartBench (Xu et al., 2023)	<b>9 / 42</b>	-	Table	✗	✗	✗	✗	✓	✓	✓
ChartX (Xia et al., 2024)	18	22	Code*	✗	✓	✗	✗	✗	✓	✓
MMC (Liu et al., 2024a)	6	5	Caption	✓	✓	✗	✓	✓	✓	✗
ChartLlama (Han et al., 2023)	10	-	Table	✗	✓	✗	✓	✓	✓	✓
ChartAst (Meng et al., 2024)	9	-	Table	✗	✗	✗	✓	✓	✗	✗
ChartInstruct (Masry et al., 2024a)	-	-	Table	✗	✓	✗	✓	✓	✗	✗
ChartGemma (Masry et al., 2024b)	-	-	-	✗	✓	✓	✓	✓	✗	✗
<b>REACHQA (ours)</b>	<b>10 / 32</b>	$\infty$	<b>Code</b>	✓	✓	✓	✓	✓	✓	✓

Table 1: Comparison of existing chart-related datasets. Only the chart Q&A task is considered, though some datasets include multiple tasks. Abbreviations: Vis.=visual, Comp.=complexity, Temp.=template, Refer.=Reference, Rat.=rationale, Annot.=annotation and Scal.=scalable. Cells marked with “✓” indicate mixed attributes (e.g., partially template-based; scalable Q&A but non-scalable chart data.). “/” means the dataset includes multiple chart type granularity. “\*” indicates while chart-plotting codes are public, their Q&A synthesis still relies on data tables.

hindered scalability. With the advent of LLMs, works like ChartLlama and ChartInstruct use them to generate diverse questions from data tables while also providing rationale annotations for training. However, these methods fail to capture visual elements like color, layout, and structure because they rely on only the data table. Thus, the generated Q&A pairs lack *visual references*, undermining the inherently multimodal nature of this task. To address this, ChartGemma uses MLLMs to generate Q&A directly from charts.

**Dataset Properties.** While manually annotated datasets like MathVista and CharXiv provide high-quality data, their development is resource-intensive, typically resulting in datasets of only a few thousand samples. In the era of LLMs, such methods are impractical for scaling to the size needed to train larger models. Recent efforts, such as ChartAst, ChartInstruct, and ChartGemma, have explored Q&A generation for dataset expansion, but they remain limited by the difficulty of collecting a large set of charts. A more scalable approach is to leverage the generative capabilities of LLMs to synthesize charts like ChartBench and ChartX.

## 2.2 Can LLMs Understand Charts without Visual Input?

To explore whether there is a more effective textual format for representing visual information than data tables, we propose using **code**. By precisely encoding chart structures and details, the code may

Input	Acc.	Reas. Comp.	Vis. Refer.	Cost (\$)
Table	<b>2.72</b>	2.51	1.19	<b>0.047</b>
Code	2.60	<b>2.56</b>	2.15	0.092
Chart	1.91	1.53	<b>2.36</b>	0.107

Table 2: Rating results for three input types in our study.

serve as an ideal bridge between modalities. We design an experiment to test this hypothesis. We first collect 25 complex charts, along with their corresponding data tables and code, from authentic research papers. These charts often feature multiple or overlay plots and dense data groups, with the code averaging over 100 lines. For each sample, GPT-4o receives three types of input—table, code, and chart images—to generate a challenging Q&A pair. In total, 75 pairs are created, randomly shuffled, and then presented to annotators for blind evaluation. The annotators are asked to rate each pair on accuracy, reasoning complexity, and visual reference, using a scale of 1 (low) to 3 (high).

The results in Table 2 indicate that both text-based inputs outperform visual chart input in the first two aspects, with code scoring 2.60 in accuracy (vs. 1.91) and 2.56 in reasoning complexity (vs. 1.53). As expected, table input has the lowest visual reference score (1.19), while chart input scores highest in this (2.36), confirming the ability of MLLMs to directly interpret visual information. Surprisingly, despite the absence of visual input to

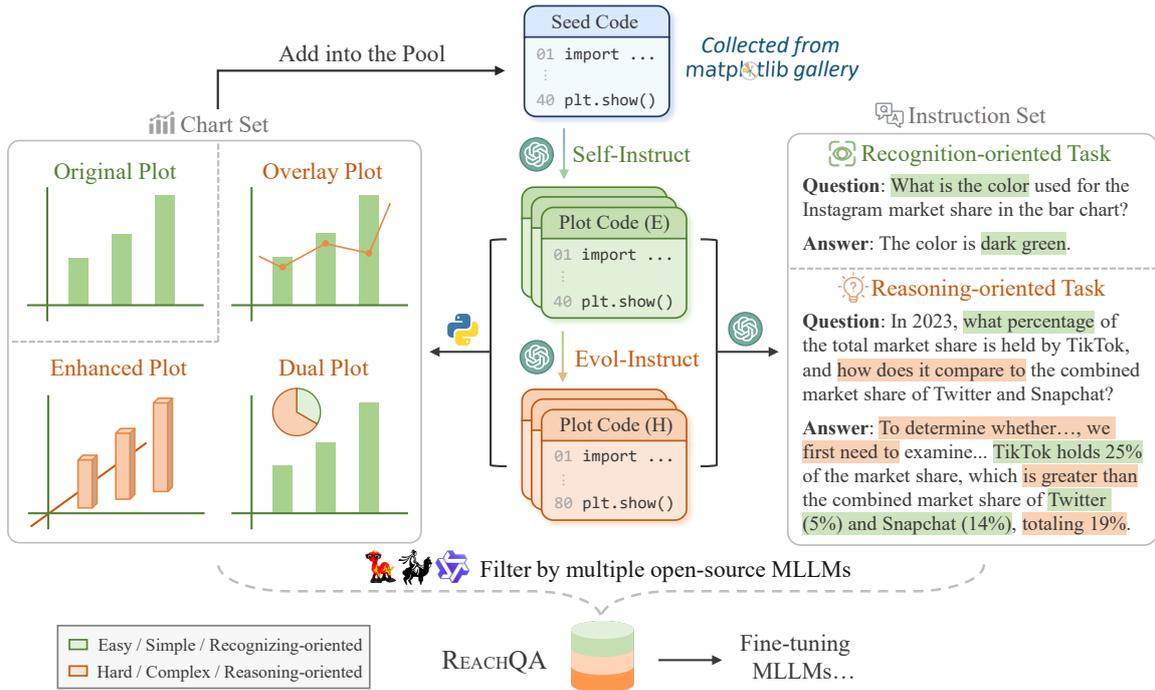


Figure 2: Overview of the Code-as-Intermediary Translation (CIT) method for synthesizing multimodal instruction data. The process starts with 33 seed codes, generating plot codes across various chart types, topics, and complexity levels via Self-Instruct and Evol-Instruct. The chart and instruction sets are constructed bi-directionally, and the final filtering yields REACHQA, a dataset for distilling visual chart reasoning abilities from LLMs to MLLMs.

the model, the code achieves a relatively high visual reference score (2.15), highlighting its potential to translate chart images into textual representations.

### 3 Methodology

Building on the findings above, we propose Code-as-Intermediary Translation (CIT), a data synthesis method for distilling visual reasoning abilities from LLMs to MLLMs, as illustrated in Figure 2. In the following sections, we describe how we synthesize intermediate codes (§ 3.1), generate paired charts and instructions (§ 3.2), ensure data quality (§ 3.3), and ultimately construct our dataset, REACHQA.

#### 3.1 Intermediary Code Synthesis

**Seed Code Collection.** We start by collecting a small set of 33 seed code samples, which we refer to as  $C_{\text{seed}}$ . These samples are sourced directly from the official Matplotlib gallery<sup>1</sup> to ensure quality and minimize manual effort. Collectively, these code samples cover a diverse range of chart types, including common types like bar, line, and scatter charts, as well as more specialized charts such as bubble, contour, and donut charts. All samples are verified for executability to guarantee the reliability of the subsequent code synthesis process.

<sup>1</sup><https://matplotlib.org/stable/gallery/index.html>

**Self-Instruct for Diverse Code Generation.** To expand the diversity and coverage of the chart set, we apply the Self-Instruct method (Wang et al., 2023), which synthesizes instruction data by prompting LLMs with existing ones as few-shot examples (Brown et al., 2020). In our approach, we provide 3 randomly selected code snippets as examples in each generation, guiding the model to synthesize chart-plotting code of the same kind.

To diversify chart generation, a chart type is randomly chosen from 10 major and 32 minor categories for the model to generate. For chart content, we provide two topic options, allowing the model to freely combine or expand on these themes based on its knowledge, leading to varied topics and data. A chain-of-thought (CoT) process (Wei et al., 2022) is used for code generation, starting with the chart’s background and data, followed by the final executable code. This step-by-step approach ensures logical coherence and code functionality. The generated codes are referred to as  $C_{\text{easy}}$  for use in subsequent phases of the construction. The chart types and topics are detailed in Appendix A.1.

**Evol-Instruct for Complex Code Generation.** To enhance the visual complexity of the synthetic charts, we adopt the Evol-Instruct method (Xu et al., 2024), which leverages LLMs to evolve sim-

266	ple chart-plotting code into more complex versions	316
267	by presenting existing code alongside an evolution	317
268	strategy as context. It addresses a key limitation in	318
269	prior work that emphasizes the quantity of charts	319
270	while often neglecting the difficulty of chart inter-	320
271	pretation. Starting with code samples from $C_{\text{easy}}$ ,	321
272	we apply one of the following predefined evolution	
273	directions: (1) expanding the data size or num-	322
274	ber of data groups; (2) adding or modifying visual	
275	elements to enhance presentation; (3) overlaying	323
276	a different type of chart on the original plot; (4)	324
277	introducing an additional subplot beside the origi-	325
278	nal plot. These strategies ensure that the resulting	326
279	charts demand more nuanced visual interpretation	327
280	and in-depth reasoning. We follow a CoT process	328
281	like previous steps, where the model first analyzes	329
282	the existing code and then generates the evolved	330
283	one. The evolved codes, referred to as $C_{\text{hard}}$ , are	331
284	also added to the code pool for subsequent use.	332
285		
286	<b>3.2 Bi-directional Translation</b>	
287	<b>Chart Generation through Code Execution and</b>	
288	<b>Self-Repair.</b> We generate charts by executing all	
289	the Python plotting code. However, during the gener-	
290	ation and evolution process, program errors are	
291	inevitable. To ensure correctness, we will validate	
292	the code before adding it to the pool. When errors	
293	occur, the code is not immediately discarded; in-	
294	stead, we apply a Self-Repair method (Chen et al.,	
295	2024a), feeding the code and execution results into	
296	the LLMs for correction. This process repeats until	
297	the code is fixed or reaches an iteration limit, af-	
298	ter which the code is discarded if it remains faulty.	
299	On average, this approach fixes about 15% of the	
300	code generated by GPT-4o, with 5% remaining	
	unrepairable and filtered out, yielding $C_{\text{final}}$ .	
301	<b>Instruction Generation through Guided Prompt-</b>	
302	<b>ing.</b> After verifying executability, we use $C_{\text{final}}$	
303	to create instruction sets in the form of Q&A pairs.	
304	Building on prior work of in-context Q&A gener-	
305	ation (Chen et al., 2023; He et al., 2024), we	
306	guide the model in two steps: first generating a	
307	batch of questions, then producing corresponding	
308	answers. To ensure high-quality answers, we also	
309	employ a step-by-step approach where the model	
310	first provides detailed calculations and analyses,	
311	which are then refined into concise, educational	
312	answers optimized for learning (Gunasekar et al.,	
313	2023). The model generates two types of instruc-	
314	tions: <i>recognition-oriented</i> , focused on visual infor-	
315	mation retrieval, and <i>reasoning-oriented</i> , requiring	
	both recognition and multi-step reasoning. With	
	minimal constraints on content, the model is en-	
	couraged to explore creative and diverse instruc-	
	tions. Multiple questions can be generated for	
	each chart, and redundant ones are filtered using	
	ROUGE-L overlap, following Wang et al. (2023).	
	<b>3.3 Quality Assurance</b>	
	<b>Multimodal Validation for Enhanced Data Qual-</b>	
	<b>ity.</b> Although our dataset is fully synthesized us-	
	ing LLMs, we acknowledge the importance of inte-	
	grating visual information to enhance data quality	
	(Masry et al., 2024b; Zeng et al., 2024). Thus	
	we introduce a multimodal validation step, using	
	MLLMs to verify both generated charts and their	
	corresponding instructions. Since models differ in	
	architecture, visual encoders, and training recip-	
	es, they may focus on varying aspects of the im-	
	ages. Taking this into account, we adopt a “majority vot-	
	ing” approach by ensembling multiple smaller, lo-	
	cally hosted models. This ensures reliable visual	
	validation while remaining cost-effective. For chart	
	validation, each model rates charts on a scale of 1	
	to 5, and those below a threshold are filtered out.	
	For instructions, both Q&A pairs and correspond-	
	ing charts are fed into the models and verified, with	
	multiple negative votes leading to sample rejection.	
	<b>Testing Set Construction and Annotation Refine-</b>	
	<b>ment.</b> For the REACHQA testing set, we follow	
	a similar process as in previous data generation but	
	apply stricter filtering criteria to ensure higher qual-	
	ity. Additional annotators are recruited for manual	
	review and refinement. For the charts, they first	
	check the images to identify any potential visual er-	
	rors. For the Q&A pairs, they ensure the questions	
	are relevant to the chart and answerable, then cor-	
	rect any hallucinations or logical inconsistencies	
	in the answers. Afterwards, two rounds of review	
	are conducted to confirm the questions meet the	
	multimodal recognition or reasoning standards in	
	our settings. Only samples with agreement from	
	at least two annotators are included. The inter-	
	annotator agreement, with a kappa coefficient of	
	0.82, indicates strong consistency (Landis, 1977).	
	Table 3 presents the final dataset statistics.	
	The total cost of data construction, excluding	
	open-source model usage and annotation labor for	
	the testing set, was about \$300. The detailed ex-	
	penditure breakdown is provided in Appendix A.2.	
	Since all data splits are generated using the same	
	process and model, we analyze potential data con-	

Statistics	Train Set	Test Set
Total charts	3,249	500
- # Chart types	10 / 32	10 / 32
- # Overlay plots	1,030	220
- # Multiple plots	593	251
- Average size (px)	2480×1571	2798×1601
Unique questions	19,963	2,000
- # Reco. per chart	2.53	2
- # Reas. per chart	3.62	2
Avg. Reco. Q. length	22.1	21.0
Avg. Reco. A. length	38.3	7.0
Avg. Reas. Q. length	38.2	35.4
Avg. Reas. A. length	68.4	24.9

Table 3: REACHQA dataset statistics. Sequence lengths are calculated based on the GPT-4o tokenizer.

tamination in Appendix A.4. The prompt templates we use in each step are shown in Appendix F.

## 4 Experiments

### 4.1 Experimental Setups

**Benchmarks.** We evaluate the models on three categories of tasks that cover both chart-related and general multimodal recognition and reasoning. First, traditional chart-related benchmarks are considered, including ChartQA, ChartBench, and ChartX, which primarily test recognition capabilities. Second, we assess novel chart-related benchmarks that require both recognition and reasoning, including CharXiv and our REACHQA test set. Third, we evaluate general multimodal reasoning abilities on MathVista and MATH-Vision.

**Models and baselines.** We evaluate a range of MLLMs from three categories: (1) Powerful proprietary models, including GPT-4o (OpenAI, 2024a), GPT-4o mini (OpenAI, 2024b), and Claude 3.5 Sonnet (Anthropic, 2024). (2) Chart-augmented open-source models, such as ChartInstruct-7B (Masry et al., 2024a), ChartAssistant-13B (Meng et al., 2024), and ChartGemma-3B (Masry et al., 2024b), which are specifically enhanced for chart-related tasks. (3) Latest general open-source models, including LLaVA-Next-Llama3-8B (Li et al., 2024), MiniCPM-Llama3-V2.5-8B (Yao et al., 2024), and InternVL2-8B (Chen et al., 2024b). For each general model, we conduct supervised fine-tuning (SFT) using the REACHQA training set. Specifically, we train three variants: one using 8k recognition-oriented samples (denoted as Reco.),

one using 12k reasoning-oriented samples (denoted as Reas.), and a combined version incorporating both (denoted as All). More details on the datasets and evaluation can be found in Appendix C.

### 4.2 Experimental Results

Table 4 presents the quantitative results for all models across each task. We can find that:

**Synthetic datasets can also effectively measure abilities.** Our REACHQA test set effectively evaluates models’ reasoning and recognition skills, showing trends similar to human-annotated datasets like CharXiv. For instance, GPT-4o exhibits a reasoning score of 39.70 and a recognition score of 66.80 on REACHQA, closely mirroring its performance on CharXiv (i.e., 47.10 and 84.45, respectively). This consistency suggests that LLM-generated datasets, with minimal human intervention, can rival human-labeled data. Moreover, REACHQA presents a significant challenge to models’ visual abilities, as random guessing results in very low scores. In contrast, traditional benchmarks like ChartQA may allow models to leverage pre-existing knowledge, inflating results without truly testing visual capabilities (Yue et al., 2024).

**Proprietary models demonstrate more balanced performance.** Proprietary models like GPT-4o achieve competitive results on both traditional chart-related tasks and reasoning-intensive tasks like REACHQA and CharXiv. In contrast, open-source models, whether chart-augmented or general-purpose, excel in recognition tasks but struggle in complex ones. This disparity highlights their imbalanced capabilities, and also suggests potential overfitting to simpler charts. Although proprietary models may not always lead in specific tasks, the stable and balanced performance makes them more suitable for real-world applications.

**Specialized training data significantly improves model performance.** Models trained on 8k REACHQA recognition data outperform in recognition tasks, while those trained on 12k reasoning data could do better in reasoning tasks. When both data types are combined (i.e., 20k in total), models see the greatest improvement, with performance increasing by at least 15% across all models we test. Notably, the LLaVA-Next-Llama3-8B achieves a 34.8% boost in average performance. This suggests that a model’s visual capability comprises two complementary aspects, and training on both

Models	Avg. ( $\uparrow$ )	ChartQA	ChartBench	ChartX	REACHQA	CharXiv	MathVista	MATH-V			
		QA	Binary NQA	QA	Reas. Reco.	Reas. Desc.	Math General	QA			
<b>Proprietary Multimodal Large Language Models</b>											
GPT-4o mini	49.34	77.52	70.26	34.93	35.45	27.20	53.50	34.10	74.92	56.70	28.85
GPT-4o	59.85	85.70	<b>81.03</b>	<b>52.88</b>	46.60	39.70	66.80	47.10	<b>84.45</b>	63.80	30.39
Claude 3.5 Sonnet	<b>64.50</b>	<b>90.80</b>	76.72	48.29	<b>58.24</b>	<b>51.70</b>	<b>74.30</b>	<b>60.20</b>	84.30	<b>67.70</b>	<b>32.76</b>
<b>Chart-augmented Multimodal Large Language Models</b>											
ChartInstruct-7B	25.93	66.64	61.40	26.95	26.62	6.00	10.50	8.80	<b>21.40</b>	15.37	31.52
ChartAssistant-13B	28.25	79.90	58.15	24.62	23.20	<b>10.70</b>	19.60	11.70	16.93	17.78	<b>39.57</b>
ChartGemma-3B	<b>33.08</b>	<b>80.16</b>	<b>78.90</b>	<b>34.10</b>	<b>35.15</b>	9.20	<b>27.80</b>	<b>12.50</b>	21.30	<b>19.07</b>	38.04
<b>Open-Source Multimodal Large Language Models</b>											
LLaVA-Next-Llama3-8B	24.46	45.80	42.90	15.86	15.45	6.50	17.90	17.20	31.45	22.41	44.13
+ REACHQA (Reco.)	32.88 (+34.4%)	<b>66.96</b>	56.95	<b>29.52</b>	<b>27.25</b>	8.80	29.00	22.20	32.58	27.40	49.78
+ REACHQA (Reas.)	32.39 (+32.4%)	64.48	56.80	25.14	25.90	8.40	26.30	<b>22.70</b>	<b>35.67</b>	<b>28.89</b>	<b>50.65</b>
+ REACHQA (All)	<b>32.98 (+34.8%)</b>	64.56	<b>57.00</b>	29.33	27.08	<b>11.10</b>	<b>29.60</b>	22.50	32.33	27.59	50.43
MiniCPM-Llama3-V2.5	33.39	66.92	48.90	22.29	23.72	10.30	25.30	22.00	46.20	37.22	53.04
+ REACHQA (Reco.)	38.62 (+15.7%)	71.12	<b>56.65</b>	<b>33.29</b>	29.53	10.60	34.10	25.60	<b>48.75</b>	41.48	<b>60.43</b>
+ REACHQA (Reas.)	38.52 (+15.4%)	71.72	<b>56.65</b>	29.62	28.23	<b>11.00</b>	33.00	27.50	48.70	<b>43.52</b>	60.22
+ REACHQA (All)	<b>38.67 (+15.8%)</b>	<b>71.44</b>	55.80	30.43	<b>29.68</b>	<b>11.00</b>	<b>35.10</b>	<b>28.30</b>	47.62	42.22	60.00
InternVL2-8B	40.03	73.80	52.05	32.86	35.10	16.20	33.70	26.30	46.10	46.11	61.74
+ REACHQA (Reco.)	48.21 (+20.4%)	<b>82.92</b>	<b>66.35</b>	46.14	<b>46.62</b>	19.90	49.50	32.20	54.38	47.96	<b>67.61</b>
+ REACHQA (Reas.)	47.87 (+19.6%)	82.84	64.05	46.52	44.88	20.10	49.40	<b>32.80</b>	52.40	<b>49.44</b>	66.52
+ REACHQA (All)	<b>48.35 (+20.8%)</b>	82.44	65.90	<b>47.29</b>	45.38	<b>21.30</b>	<b>49.80</b>	32.70	<b>54.83</b>	48.89	66.30

Table 4: Evaluation results on seven benchmarks. The best performance for each category and task is in **bold**. The percentage of performance improvements compared to the vanilla model is denoted by ( $\uparrow$ ).

Models	Avg.	REACHQA	CharXiv	MathVista	Math-V
Base Model	16.39	6.50	17.20	32.40	9.44
+ ChartBench	17.06	7.30	17.00	33.60	10.33
+ ChartAst	17.67	7.10	<u>20.40</u>	32.10	<u>11.08</u>
+ The Cauldron	18.61	<u>10.10</u>	19.10	35.60	9.64
+ ChartGemma	<u>19.11</u>	10.00	19.40	<u>36.40</u>	10.62
+ REACHQA	<b>20.74</b>	<b>11.10</b>	<b>22.50</b>	<b>38.10</b>	<b>11.25</b>

Table 5: Performance comparison of models trained on different datasets. The REACHQA and CharXiv scores refer to reasoning splits here.

data types together produces optimal results. Moreover, despite the absence of math-target data in the training set, the models generalize well to the MathVista and MATH-Vision benchmarks, highlighting the transferability of multimodal reasoning abilities distilled from expert rationales. We also provide a qualitative analysis of the models’ reasoning patterns in Appendix D and Figure 7, to further demonstrate REACHQA’s working mechanism.

## 5 Discussion

### 5.1 The Role of Expert Rationales

We analyze how training data quality affects visual reasoning abilities by comparing major open-source datasets (ChartBench, ChartAst, Chart-

Gemma and The Cauldron<sup>2</sup>). We uniformly sample 20k Q&A instructions from each dataset and train LLaVA-Next-Llama3-8B under controlled settings.

As shown in Table 5, the model trained on ChartBench performs the worst, likely due to the absence of reasoning steps in its responses. Although ChartAst includes rationale annotations, the template-based questions limit its effectiveness for learning reasoning patterns. The model trained on the mixed dataset of The Cauldron show modest improvements, but is still restricted by the subsets’ quality. In contrast, models trained on ChartGemma and REACHQA perform better, likely due to the distillation of expert rationales from stronger models (e.g., Gemini Flash 1.5 and GPT-4o), which directly affect the visual reasoning abilities. Additionally, we believe the visual richness of charts, as detailed in Appendix E, may also help improve generalization.

### 5.2 Interaction between Recognition & Reasoning Abilities

As previously noted, the recognition and reasoning abilities are likely interdependent. Wang et al. (2024c) suggest that recognition skills serve as prerequisites for effective reasoning. To investi-

<sup>2</sup>Unlike other datasets, The Cauldron (Laurençon et al., 2024) is selected for its generality as a collection of 50 vision-language datasets, from which we use 7 chart-related subsets.

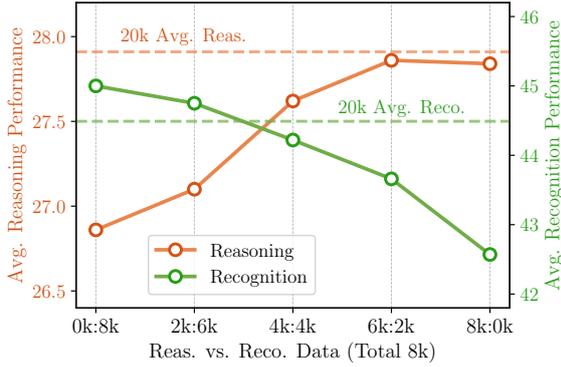


Figure 3: Performance comparison of different training data ratios with 8k total data. The dashed line represents the model’s performance trained with full 20k data.

gate this further, we conduct an experiment with LLaVA-Next-Llama3-8B, using a fixed 8k total training data size and varying the ratio of reasoning to recognition data from 0:8 to 8:0. We evaluated the models on recognition tasks (i.e., ChartQA, ChartBench, ChartX) and reasoning tasks (i.e., REACHQA-Reas., CharXiv-Reas., MathVista).

Figure 3 shows that increasing the proportion of recognition or reasoning data improves performance on the respective tasks. Models with more recognition data outperform those trained on 20k full data for recognition tasks. However, the reasoning performance gains plateau and even decline when reasoning data exceeds 50%, suggesting diminishing returns when reasoning data is overemphasized. This supports the hypothesis that reasoning abilities are partially dependent on recognition skills. When the model fails to interpret the image accurately, its reasoning ability is likely compromised (Wang et al., 2024c). Although this study is limited by data constraints, we expect the interaction between recognition and reasoning to become more pronounced with larger datasets.

### 5.3 Balancing General & Specialized Abilities

We investigate how models trained on specialized data perform on general-purpose multimodal tasks. Using 7 general multimodal benchmarks and 5 reasoning-focused benchmarks, we test three versions of LLaVA-Next-Llama3-8B: the vanilla model, the one trained on 20k REACHQA samples, and another trained on 20k REACHQA plus 20k general-purpose multimodal data sampled from 779k LLaVA-NeXT-Data<sup>3</sup>. This dataset is chosen because the LLaVA-NeXT family of models was

<sup>3</sup><https://huggingface.co/datasets/lmms-lab/LLaVA-NeXT-Data>

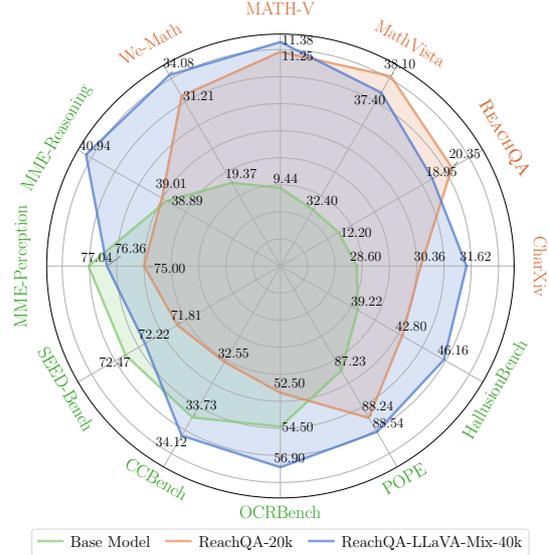


Figure 4: Performance comparison of models on 7 general tasks and 5 specialized reasoning tasks.

officially fine-tuned on it (Li et al., 2024), allowing us to approximate its original data distribution.

As shown in Figure 4, the vanilla model (green area) struggles with reasoning tasks, while the REACHQA-trained model (orange area) improves on reasoning but loses general performance. Surprisingly, by incorporating only 20k general data (blue area), the model not only restores its general multimodal performance but also retains the enhanced reasoning ability. It results in a well-balanced model with notable reasoning improvements and minimal drops in general domains.

## 6 Conclusion

In this work, we delve into the key challenges MLLMs face in complex chart Q&A tasks, highlighting their deficiencies in both recognition and reasoning. Building on our analysis of existing datasets and the untapped potential of LLMs, we propose Code-as-Intermediary Translation (CIT) as a novel method for distilling LLMs’ abilities to improve MLLMs. With code as a bridge between visual and textual modalities, CIT enables language models to interpret complex charts more precisely, facilitating the generation of higher-quality Q&A pairs. Our synthetic dataset, REACHQA, demonstrates significant performance improvements across multiple models and benchmarks, with gains extending beyond chart-specific tasks to broader multimodal reasoning. We believe CIT offers a promising direction for scalable and cost-effective multimodal instruction data synthesis.

## 550 Limitations

551 We summarize the limitations of our method as  
552 follows: (1) While CIT effectively uses code to  
553 link text and abstract images like charts and dia-  
554 grams, applying this approach to natural images  
555 remains challenging. Current text-to-image models  
556 still lack precise control over fine details (Betker  
557 et al., 2023; Zhang et al., 2023), which can lead  
558 to misaligned synthetic data. Once more control-  
559 able techniques are developed, the synthesis of  
560 multimodal data could become more flexible and  
561 applicable. (2) Although multimodal validation  
562 steps were introduced to reduce errors, the syn-  
563 thesized charts and Q&A pairs might still contain  
564 occasional inaccuracies. Therefore, to ensure data  
565 quality for larger-scale applications, stronger mod-  
566 els and stricter thresholds are essential. (3) Our  
567 method may not be as effective for teacher models  
568 with limited capabilities, as it is inherently on a  
569 form of distillation (Hinton et al., 2015). The suc-  
570 cess of distillation depends on the strength of the  
571 teacher model, and in our scenario, weaker mod-  
572 els may face challenges in interpreting charts via  
573 code. Nevertheless, we believe that future models  
574 will not only become more capable but also more  
575 cost-efficient for data synthesis.

## 576 Ethical Statement

577 This research utilizes synthetic datasets for experi-  
578 mentation. We have ensured that all datasets com-  
579 ply with relevant ethical and privacy standards. All  
580 synthetic data have been rigorously processed to  
581 prevent the disclosure of any potentially sensitive  
582 information. We are committed to adhering to the  
583 ACL’s ethical policies, ensuring transparency and  
584 reproducibility throughout the research process.

## 585 References

586 Anthropic. 2024. [Introducing claude 3.5 sonnet](#).

587 James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jian-  
588 feng Wang, Linjie Li, Long Ouyang, Juntang Zhuang,  
589 Joyce Lee, Yufei Guo, et al. 2023. Improving image  
590 generation with better captions. *Computer Science*.  
591 <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8.

592 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie  
593 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind  
594 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
595 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
596 Gretchen Krueger, Tom Henighan, Rewon Child,  
597 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
598 Clemens Winter, Christopher Hesse, Mark Chen, Eric

Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, 599  
Jack Clark, Christopher Berner, Sam McCandlish, 600  
Alec Radford, Ilya Sutskever, and Dario Amodei. 601  
2020. Language models are few-shot learners. In *Ad- 602*  
*vances in Neural Information Processing Systems 33: 603*  
*Annual Conference on Neural Information Process- 604*  
*ing Systems 2020, NeurIPS 2020, December 6-12, 605*  
*2020, virtual*. 606

Wei-Lin Chen, Cheng-Kuang Wu, Yun-Nung Chen, 607  
and Hsin-Hsi Chen. 2023. Self-icl: Zero-shot in- 608  
context learning with self-generated demonstrations. 609  
In *Proceedings of the 2023 Conference on Empirical 610*  
*Methods in Natural Language Processing, EMNLP 611*  
*2023, Singapore, December 6-10, 2023*, pages 15651– 612  
15662. Association for Computational Linguistics. 613

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and 614  
Denny Zhou. 2024a. Teaching large language mod- 615  
els to self-debug. In *The Twelfth International Con- 616*  
*ference on Learning Representations, ICLR 2024, 617*  
*Vienna, Austria, May 7-11, 2024*. 618

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo 619  
Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, 620  
Xizhou Zhu, Lewei Lu, et al. 2024b. Internvl: Scal- 621  
ing up vision foundation models and aligning for 622  
generic visual-linguistic tasks. In *Proceedings of 623*  
*the IEEE/CVF Conference on Computer Vision and 624*  
*Pattern Recognition*, pages 24185–24198. 625

Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, 626  
Gautier Viaud, Céline Hudelot, and Pierre Colombo. 627  
2024. Colpali: Efficient document retrieval with 628  
vision language models. *CoRR*, abs/2407.01449. 629

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, 630  
Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jin- 631  
rui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Ron- 632  
grong Ji. 2023a. MME: A comprehensive evaluation 633  
benchmark for multimodal large language models. 634  
*arXiv preprint arXiv:2306.13394*. 635

Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and 636  
Tushar Khot. 2023b. [Complexity-based prompting 637](#)  
[for multi-step reasoning](#). In *The Eleventh Inter- 638*  
*national Conference on Learning Representations, 639*  
*ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. Open- 640  
Review.net. 641

Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, 642  
Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, 643  
Furong Huang, Yaser Yacoob, Dinesh Manocha, and 644  
Tianyi Zhou. 2024. Hallusionbench: An advanced 645  
diagnostic suite for entangled language hallucina- 646  
tion and visual illusion in large vision-language mod- 647  
els. In *2024 IEEE/CVF Conference on Computer 648*  
*Vision and Pattern Recognition (CVPR)*, pages 14375– 649  
14385. 650

Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio 651  
César Teodoro Mendes, Allie Del Giorno, Sivakanth 652  
Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo 653  
de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all 654  
you need. *arXiv preprint arXiv:2306.11644*. 655

656	Yucheng Han, Chi Zhang, Xin Chen, Xu Yang,	Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan	713
657	Zhibin Wang, Gang Yu, Bin Fu, and Hanwang	Raiman, Mohammad Shoeybi, Bryan Catanzaro, and	714
658	Zhang. 2023. Chartllama: A multimodal llm for	Wei Ping. 2024. Nv-embed: Improved techniques for	715
659	chart understanding and generation. <i>arXiv preprint</i>	training llms as generalist embedding models. <i>CoRR</i> ,	716
660	<i>arXiv:2311.16483</i> .	abs/2405.17428.	717
661	Wei He, Shichun Liu, Jun Zhao, Yiwen Ding, Yi Lu, Zhi-	Jacqueline Léon. 2007. From universal languages to	718
662	heng Xi, Tao Gui, Qi Zhang, and Xuanjing Huang.	intermediary languages in machine translation. In	719
663	2024. Self-demos: Eliciting out-of-demonstration	<i>History of Linguistics 2002: Selected Papers from the</i>	720
664	generalizability in large language models. In <i>Find-</i>	<i>Ninth International Conference on the History of the</i>	721
665	<i>ings of the Association for Computational Linguistics:</i>	<i>Language Sciences, 27–30 August 2002, São Paulo-</i>	722
666	<i>NAACL 2024, Mexico City, Mexico, June 16-21,</i>	<i>Campinas</i> , volume 110, page 123. John Benjamins	723
667	<i>2024</i> , pages 3829–3845. Association for Computa-	Publishing Amsterdam.	724
668	tional Linguistics.		
669	Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean.	Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Ren-	725
670	2015. Distilling the knowledge in a neural network.	rui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and	726
671	<i>CoRR</i> , abs/1503.02531.	Chunyu Li. 2024. Llava-next: Stronger llms super-	727
672	Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang	charge multimodal capabilities in the wild.	728
673	Zhang, Bo Zhang, Ji Zhang, Qin Jin, Fei Huang,		
674	and Jingren Zhou. 2024. mplug-docowl 1.5: Unified	Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yix-	729
675	structure learning for ocr-free document understand-	iao Ge, and Ying Shan. 2023a. Seed-bench: Bench-	730
676	ing. In <i>Findings of the Association for Computa-</i>	marking multimodal llms with generative compre-	731
677	<i>tional Linguistics: EMNLP 2024, Miami, Florida,</i>	hension. <i>arXiv preprint arXiv:2307.16125</i> .	732
678	<i>USA, November 12-16, 2024</i> , pages 3096–3120. As-		
679	sociation for Computational Linguistics.	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang,	733
680	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evalu-	734
681	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	ating object hallucination in large vision-language	735
682	Weizhu Chen. 2022. Lora: Low-rank adaptation of	models. In <i>Proceedings of the 2023 Conference on</i>	736
683	large language models. In <i>The Tenth International</i>	<i>Empirical Methods in Natural Language Process-</i>	737
684	<i>Conference on Learning Representations, ICLR 2022,</i>	<i>ing, EMNLP 2023, Singapore, December 6-10, 2023,</i>	738
685	<i>Virtual Event, April 25-29, 2022</i> .	pages 292–305. Association for Computational Lin-	739
686	Kung-Hsiang Huang, Hou Pong Chan, Yi R Fung,	guistics.	740
687	Haoyi Qiu, Mingyang Zhou, Shafiq Joty, Shih-Fu	Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song,	741
688	Chang, and Heng Ji. 2024. From pixels to insights:	Jue Wang, and Pengtao Xie. 2022. Evit: Expediting	742
689	A survey on automatic chart understanding in the	vision transformers via token reorganizations. In <i>The</i>	743
690	era of large foundation models. <i>arXiv preprint</i>	<i>Tenth International Conference on Learning Repre-</i>	744
691	<i>arXiv:2403.12027</i> .	<i>sentations, ICLR 2022, Virtual Event, April 25-29,</i>	745
692	Shankar Kantharaj, Xuan Long Do, Rixie Tiffany Ko	2022. OpenReview.net.	746
693	Leong, Jia Qing Tan, Enamul Hoque, and Shafiq R.	Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen,	747
694	Joty. 2022. Opencqa: Open-ended question answer-	Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and	748
695	ing with charts. In <i>Proceedings of the 2022 Con-</i>	Dong Yu. 2024a. MMC: advancing multimodal chart	749
696	<i>ference on Empirical Methods in Natural Language</i>	understanding with large-scale instruction tuning. In	750
697	<i>Processing, EMNLP 2022, Abu Dhabi, United Arab</i>	<i>Proceedings of the 2024 Conference of the North</i>	751
698	<i>Emirates, December 7-11, 2022</i> , pages 11817–11837.	<i>American Chapter of the Association for Computa-</i>	752
699	Association for Computational Linguistics.	<i>tional Linguistics: Human Language Technologies</i>	753
700	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	(Volume 1: Long Papers), <i>NAACL 2024, Mexico City,</i>	754
701	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	<i>Mexico, June 16-21, 2024</i> , pages 1287–1310. Associ-	755
702	guage models are zero-shot reasoners. In <i>Advances</i>	ation for Computational Linguistics.	756
703	in <i>Neural Information Processing Systems 35: An-</i>	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li,	757
704	<i>annual Conference on Neural Information Processing</i>	Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi	758
705	<i>Systems 2022, NeurIPS 2022, New Orleans, LA, USA,</i>	Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua	759
706	<i>November 28 - December 9, 2022</i> .	Lin. 2023. Mmbench: Is your multi-modal model an	760
707	JR Landis. 1977. The measurement of observer agree-	all-around player? <i>arXiv preprint arXiv:2307.06281</i> .	761
708	ment for categorical data. <i>Biometrics</i> .	Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang,	762
709	Hugo Laurençon, Léo Tronchon, Matthieu Cord, and	Wenwen Yu, Chunyu Li, Xucheng Yin, Cheng lin	763
710	Victor Sanh. 2024. What matters when build-	Liu, Lianwen Jin, and Xiang Bai. 2024b. Ocrbench:	764
711	ing vision-language models? <i>arXiv preprint</i>	On the hidden mystery of ocr in large multimodal	765
712	<i>arXiv:2405.02246</i> .	models. <i>arXiv preprint arXiv:2305.07895</i> .	766
		Yuliang Liu, Biao Yang, Qiang Liu, Zhang Li, Zhiyin	767
		Ma, Shuo Zhang, and Xiang Bai. 2024c. Textmon-	768
		key: An ocr-free large multimodal model for under-	769
		standing document. <i>CoRR</i> , abs/2403.04473.	770

771	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024.	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	824
772			825
773			826
774			827
775			828
776			
777			
778	Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R. 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.	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khoshdel, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2023, Toronto, Canada, July 9-14, 2023, pages 13484–13508. Association for Computational Linguistics.	829
779			830
780			831
781			832
782			833
783			834
784			835
785			836
786			837
787			
788			
789			
790			
791			
792	Ahmed Masry, Mehrad Shahmohammadi, Md. Rizwan Parvez, Enamul Hoque, and Shafiq Joty. 2024a. Chartinstruct: Instruction tuning for chart comprehension and reasoning. In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 10387–10409. Association for Computational Linguistics.	Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. 2024c. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. <i>arXiv preprint arXiv:2406.18521</i> .	838
793			839
794			840
795			841
796			842
797			
798			
799			
800			
801	Ahmed Masry, Megh Thakkar, Aayush Bajaj, Aaryaman Kartha, Enamul Hoque, and Shafiq Joty. 2024b. Chartgamma: Visual instruction-tuning for chart reasoning in the wild. <i>arXiv preprint arXiv:2407.04172</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	843
802			844
803			845
804			846
805			847
806			848
807			849
808			850
809			
810			
811			
812			
813			
814			
815			
816			
817			
818			
819			
820			
821			
822			
823			
824			
825			
826			
827			
828			
829			
830			
831			
832			
833			
834			
835			
836			
837			
838			
839			
840			
841			
842			
843			
844			
845			
846			
847			
848			
849			
850			
851			
852			
853			
854			
855			
856			
857			
858			
859			
860			
861			
862			
863			
864			
865			
866			
867			
868			
869			
870			
871			
872			
873			
874			
875			
876			
877			
878			
879			
880			
881			
882			
883			

880	Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang,	times each to synthesize chart-plotting code, the-	933
881	Kai Zhang, Shengbang Tong, Yuxuan Sun, Ming Yin,	oretically generating 6,000 charts. However, after	934
882	Botao Yu, Ge Zhang, et al. 2024. Mmmu-pro: A	accounting for non-executable code and images fil-	935
883	more robust multi-discipline multimodal understand-	tered out by MLLM rating, we ultimately produced	936
884	ing benchmark. <i>arXiv preprint arXiv:2409.02813</i> .	3,249 charts for Q&A synthesis.	937
885	Michael Zarechnak. 1986. The intermediary language	<b>A.3 Reasoning Depth Across Datasets</b>	938
886	for multilanguage translation. <i>Computers and trans-</i>	To further clarify how REACHQA enhances reason-	939
887	<i>lation</i> , 1(2):83–91.	ing abilities compared to others, we conducted a	940
888	Xingchen Zeng, Haichuan Lin, Yilin Ye, and Wei Zeng.	small-scale study using GPT-4o. Specifically, we	941
889	2024. Advancing multimodal large language mod-	randomly sampled 100 examples from each dataset	942
890	els in chart question answering with visualization-	and evaluated the generated answers along two	943
891	referenced instruction tuning. <i>IEEE Transactions on</i>	proxies for reasoning depth: (1) response length	944
892	<i>Visualization and Computer Graphics</i> .	(token count) and (2) rationale length (step count),	945
893	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov,	following Fu et al. (2023b). As shown in Table 10,	946
894	and Lucas Beyer. 2023. Sigmoid loss for language	REACHQA elicited significantly longer responses	947
895	image pre-training. In <i>IEEE/CVF International</i>	and more multi-step rationales than others, indicat-	948
896	<i>Conference on Computer Vision, ICCV 2023, Paris,</i>	ing its strength in promoting structured, in-depth	949
897	<i>France, October 1-6, 2023</i> , pages 11941–11952.	reasoning.	950
898	IEEE.	<b>A.4 Data Contamination Analysis</b>	951
899	Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023.	To ensure the validity of our experimental results	952
900	Adding conditional control to text-to-image diffusion	and exclude potential data contamination, we con-	953
901	models. In <i>IEEE/CVF International Conference on</i>	duct a comprehensive analysis of data overlap from	954
902	<i>Computer Vision, ICCV 2023, Paris, France, October</i>	both dataset-level and split-level perspectives. First,	955
903	<i>1-6, 2023</i> , pages 3813–3824. IEEE.	to evaluate image-level similarity, we employed the	956
904	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	SigLIP-400M encoder (Zhai et al., 2023) to gener-	957
905	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	ate embeddings for all chart images across datasets.	958
906	Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,	These embeddings were then projected into a two-	959
907	Joseph E. Gonzalez, and Ion Stoica. 2023. Judging	dimensional space using t-SNE (Van der Maaten	960
908	llm-as-a-judge with mt-bench and chatbot arena. In	and Hinton, 2008) for visualization, following Xu	961
909	<i>Advances in Neural Information Processing Systems</i>	et al. (2023). Second, we analyzed query-level sim-	962
910	<i>36: Annual Conference on Neural Information Pro-</i>	ilarity using the NV-Embed-v2 model (Lee et al.,	963
911	<i>cessing Systems 2023, NeurIPS 2023, New Orleans,</i>	2024) to generate embeddings for all queries, also	964
912	<i>LA, USA, December 10 - 16, 2023</i> .	visualized through t-SNE.	965
913	<b>A Additional Dataset Details</b>	As shown in Figure 5(a) and (c), the visualiza-	966
914	<b>A.1 Chart Types and Topics</b>	tion demonstrates clear distributional differences	967
915	We predefined several chart types and topics for	between REACHQA and existing chart-related	968
916	Self-Instruct prompting. Table 6 shows the 9 major	benchmarks. While some degree of overlap exists	969
917	categories we established, with their correspond-	due to the shared nature of chart-related tasks, these	970
918	ing subcategories. Additionally, Table 7 lists the	instances are limited and do not compromise the	971
919	38 topics we specified. It is important to note that	overall distinctiveness of our dataset. The distinct	972
920	these topics do not reflect the actual topic distri-	clustering patterns in both image and query spaces	973
921	butions in the generated charts, as we encourage	support the validity of our cross-dataset evaluations	974
922	the model to combine and expand upon them. Re-	and confirm that REACHQA presents novel chal-	975
923	garding the distribution of chart types, we provide	lenges not fully captured by existing benchmarks.	976
924	a breakdown in Table 8. While we aimed for a	To address potential data leakage between train-	977
925	roughly balanced representation across different	ing and testing splits, which were synthesized	978
926	chart types during data construction, some degree	through the same process, we conduct a more rig-	979
927	of imbalance remains as certain types were more	orous analysis as visualized in Figure 5(b) and (d).	980
928	prone to generation errors.	Beyond visualization, we compute pairwise similar-	981
929	<b>A.2 Cost of REACHQA Training Data</b>	ities between all training and testing samples using	982
930	<b>Construction</b>		
931	Table 9 provides a detailed expense breakdown.		
932	We executed Self-Instruct and Evol-Instruct 3,000		

Major Category	Minor Category
 Line Charts	line chart, line chart with data annotation, line chart with error bar
 Pie Charts	pie chart, donut pie chart, sector pie chart, ring chart
 Bar Charts	bar chart, bar chart with data annotation, stacked bar chart, percentage bar chart, horizontal bar chart
 3D Bar Charts	3D bar chart, stacked 3D bar chart, percentage 3D bar chart
 Node Charts	directed node chart, undirected node chart
 Radar Charts	radar chart, radar chart with area filling
 Area Charts	area chart, stacked area chart
 Box Charts	vertical box chart, horizontal box chart
 Scatter Charts	scatter chart, scatter chart with smooth fitting, 3D scatter chart (bubble chart)
 Specific Charts	heat map, rose chart, funnel chart, waterfall chart, histogram, tree map

Table 6: Major categories and minor categories of charts in REACHQA.

Art and Design	Futurism and Innovation	Agriculture and Food Production
Music and Performance	Astronomy and Space	Transportation and Logistics
Business and Finance	Social Media and the Web	Real Estate and Housing Market
Travel and Exploration	Society and Community	Government and Public Policy
Books and Publishing	Physics and Chemistry	Education and Academics
Literature and Writing	Energy and Utilities	Environment and Sustainability
History and Culture	Biology and Life Sciences	Language and Communication
Architecture and Building	Retail and E-commerce	Social Sciences and Humanities
Fashion and Style	Religion and Spirituality	Manufacturing and Production
Marketing and Advertising	Food and Beverage Industry	Artificial Intelligence and Robotics
Law and Legal Affairs	Healthcare and Health	Human Resources and Employee Management
Film and Cinema	Sports and Entertainment	Computer Science and Information Technology
Mathematics and Statistics	Science and Engineering	

Table 7: Predefined chart topics in Self-Instruct prompting.

the chart embeddings. Among the identified top 50 image pairs with similarity scores exceeding 0.9, our careful manual review revealed only 2 cases with notable similarities. We will exclude them from the test set in future versions and update the evaluation accordingly. For the remaining samples, our review confirmed clear differences in chart topics, data values, and query types, ensuring that no further data leakage or contamination is present.

## B Additional Methodology Details

### B.1 Self-Instruct and Evol-Instruct

Self-Instruct (Wang et al., 2023) is a bootstrapping approach for synthesizing instruction–response pairs by having the model generate new examples based on a small set of seeds. In our implementation, we draw three randomly selected chart-plotting snippets as few-shot prompts and then ask the model to produce novel code in the same style and get  $C_{\text{easy}}$ . Evol-Instruct (Xu et al., 2024) is an iterative refinement approach that evolves simple instructions into complex variants through incremental modifications. In our implementation,

we extend simple examples from  $C_{\text{easy}}$  into more challenging ones by asking the model to “evolve” according to a specified transformation strategy and get  $C_{\text{hard}}$ . While both Self-Instruct and Evol-Instruct were originally proposed as generic methods for instruction and example expansion, our implementation adapts them specifically for chart generation. By separating the workflows into “easy” and “hard” phases, we create a curriculum of synthetic samples that progressively increases in visual and data complexity, which was not explored in prior work.

### B.2 Example Outputs of Chart Generation

Figure 6 presents a side-by-side comparison of chart outputs generated by Self-Instruct and its evolved counterpart via Evol-Instruct.

### B.3 Filtering Criteria

To ensure the quality of synthetic chart–question pairs, we apply a series of automatic filtering criteria, as described in Section 3.3. These filters are designed to balance effectiveness with practicality, given that large-scale manual verification is

Types	Line	Pie	Bar	3D Bar	Node	Radar	Area	Box	Scatter	Specific	Total
Train Set	522	415	478	144	120	238	513	331	244	244	3,249
Test Set	101	40	112	15	18	19	24	44	63	64	500

Table 8: Distribution of chart types in REACHQA dataset.

Step	Avg. #tokens of Input	Avg. #tokens of Output	Times	Cost (\$)
Self-Instruct	$1,500 + 2,000 = 3,500$	$500 + 500 = 1,000$	3,000	$\sim 56.25$
Evol-Instruct	$700 + 1,300 = 2,000$	$300 + 700 = 1,000$	3,000	$\sim 45.00$
Self-Repair	500	500	1,500	$\sim 9.38$
Reas-QA-Gen.	$1,000 + 1,500 \times 4 = 7,000$	$500 + 300 \times 4 = 1,700$	3,249	$\sim 112.09$
Reco-QA-Gen.	$800 + 1,200 \times 4 = 5,600$	$300 + 200 \times 4 = 1,100$	3,249	$\sim 81.23$

Table 9: The average number of input and output tokens is calculated for each step in the REACHQA construction process. In the equation, each term represents the average number of tokens per step (used only in a multi-step framework), while each multiplier corresponds to the number of times that step is executed. The pricing for GPT-4o-2024-08-06 is \$2.50 per 1M input tokens and \$10.00 per 1M output tokens. As a result, the total cost amounts to approximately \$303.95.

Dataset	Original Data		GPT-4o Response	
	Q. Length	A. Length	A. Length	Step Count
ChartBench	18.80	2.53	272.94	11.48
ChartAst	24.17	24.22	263.49	10.65
The Cauldron	16.95	13.93	143.54	7.00
ChartGemma	21.68	27.27	150.32	7.93
REACHQA	<b>40.83</b>	<b>56.02</b>	<b>340.24</b>	<b>15.20</b>

Table 10: Comparative analysis of reasoning depth. “Q. and A. length” refer to the average token counts of original questions and answers. GPT-4o results are averaged over 100 sampled examples.

often infeasible. For high-impact subsets (e.g., test splits), we additionally incorporate targeted manual review to further enhance reliability. Possible extensions, such as raising filtering thresholds or incorporating more human effort, may improve data quality but would also entail higher costs.

Here we provide additional details regarding reasoning-oriented samples. To control their complexity, we introduce difficulty-aware filtering. Specifically, we discard examples with overly short questions or answers, which often signal shallow or trivial reasoning. We further estimate the depth of reasoning by counting the number of sentences in the answer and enforcing a minimum threshold, thereby ensuring that retained samples exhibit sufficient multi-step inference.

## C Additional Experiment Details

### C.1 Benchmark Details

Table 11 summarizes the benchmarks used in our main experiments, including the number of samples

for each dataset. Additionally, we use some other popular multimodal datasets in Section 5.3, including MME-Reasoning, MME-Perception (Fu et al., 2023a), SeedBench (Li et al., 2023a), CCBench (Liu et al., 2023), POPE (Li et al., 2023b), HallusionBench (Guan et al., 2024), OCRBench (Liu et al., 2024b), and We-Math (Qiao et al., 2024).

### C.2 Training and Evaluation Details

For each general open-source model, we conduct supervised fine-tuning (SFT) using our REACHQA training set. We apply Low-rank Adapters (LoRA, Hu et al., 2022) to all linear layers of the language model and projector, with a LoRA rank of 16, a LoRA alpha of 8 and a learning rate of  $2e-5$ . To fully leverage their capabilities, we prompt all models with a zero-shot CoT prompt, “Let’s think step by step” (Kojima et al., 2022), following OpenAI (2024a) and Anthropic (2024). Thus, to extract answers from the model responses and assess their correctness, we employ the LLM-as-a-judge method (Zheng et al., 2023) to calculate a relaxed accuracy. The judge model used is GPT-4o, and the prompt template for evaluation can be found in Appendix F.4.

### C.3 Additional Results of Recent Models

Due to space limitations, we provide supplementary results here for several recent models, including a stronger chart-augmented MLLM, ChartMoE (Xu et al., 2025), and a general-purpose open-source MLLM, Qwen2-VL (Wang et al., 2024b). As shown in Table 12, ChartMoE demonstrates

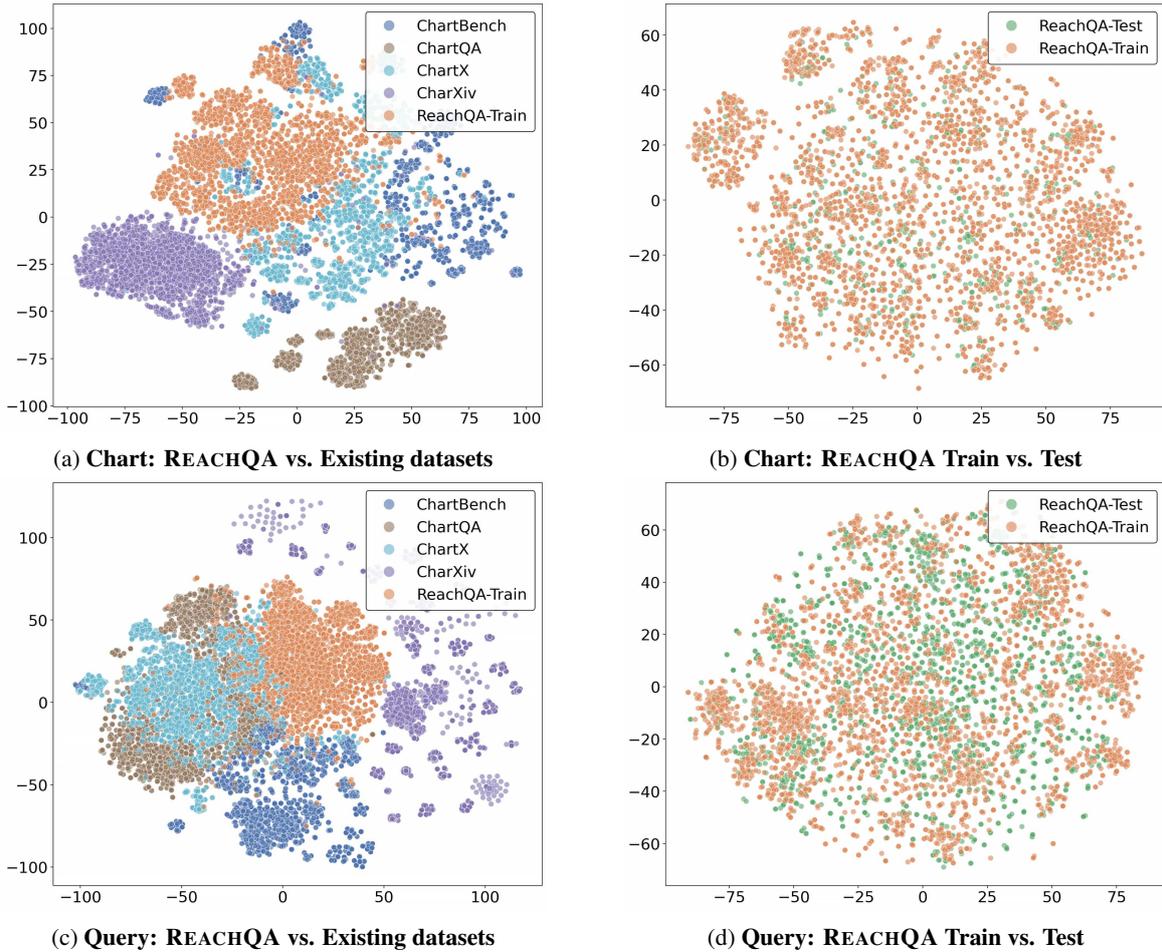


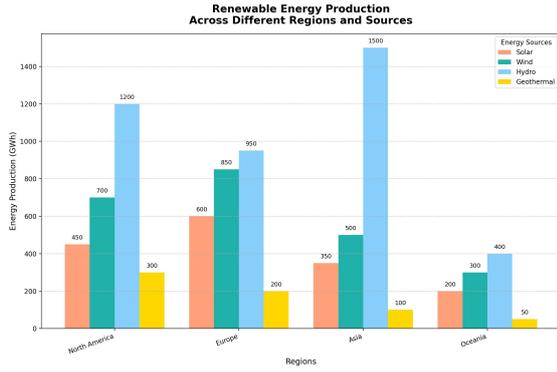
Figure 5: **Data overlap analysis visualization using t-SNE.** We analyze both image-level and query-level similarities through embedding space visualization. (a) and (c) demonstrate the distributional differences between REACHQA and existing datasets, while (b) and (d) examine potential overlap between training and testing splits. The results show clear dataset distinctiveness while revealing expected overlaps due to the shared domain of chart understanding.

1078 more balanced performance compared to other  
 1079 chart-augmented models. Qwen2-VL also exhibits  
 1080 performance improvements that align with the  
 1081 trends observed in Section 4.2. However, the ex-  
 1082 tent of improvement may be moderate, possibly  
 1083 because such models have already been exposed to  
 1084 reasoning-intensive samples during pre-training.

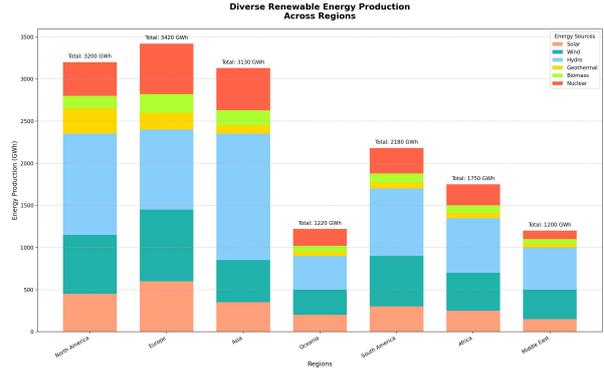
1085 Regarding document-focused models such as  
 1086 TextMonkey (Liu et al., 2024c) and DocOwl 1.5  
 1087 (Hu et al., 2024), we acknowledge their strengths  
 1088 in structured parsing of general documents (e.g.,  
 1089 DocVQA), where *page-level layout and positional*  
 1090 *information* are crucial. However, our work specif-  
 1091 ically targets Chart Question Answering (CQA),  
 1092 which requires *chart-specific information extrac-*  
 1093 *tion and numerical reasoning*. Given these differ-  
 1094 ences, a direct comparison may not fully align with  
 1095 the scope of our study.

## D A Qualitative Analysis: From the Perspective of Attention Mechanism

1096 To explore the mechanism behind the improved  
 1097 performance of our fine-tuned model, we conduct  
 1098 an analysis of the attention patterns during the next  
 1099 token prediction (Liang et al., 2022; Faysse et al.,  
 1100 2024). Figure 7 presents a comparative case study  
 1101 between the vanilla model and the fine-tuned model.  
 1102 Here, we apply full-parameter fine-tuning instead  
 1103 of LoRA to induce more pronounced changes in the  
 1104 attention layers (Hu et al., 2022). The results show  
 1105 that the vanilla model produces lengthy outputs  
 1106 with redundant analysis and dispersed attention  
 1107 across the image, reaching a wrong conclusion at  
 1108 the end. In contrast, the fine-tuned model identifies  
 1109 the key information at each step, with attention that  
 1110 accurately focuses on relevant visual elements (i.e.,  
 1111  
 1112



(a) Self-Instruct Output



(b) Evol-Instruct Output

Figure 6: **Example charts generated by (a) Self-Instruct and (b) Evol-Instruct.** The “chart type” is set to “bar chart”, and the “evolution direction” is “expanding the data size or number of data groups”. In (a), the chart consists of a simple bar plot with 4 x-axis categories and 4 labeled series. After evolution in (b), the chart becomes a stacked bar chart with 7 x-axis categories and 6 labeled series, exhibiting increased visual and structural complexity.

Benchmark	Task Focus	Sample Details
ChartQA (Masry et al., 2022)	Chart Recognition	2.5k test samples
ChartBench (Xu et al., 2023)	Chart Recognition	2k binary QA samples and 2.1k numerical QA samples
ChartX (Xia et al., 2024)	Chart Recognition	6k QA samples
REACHQA (ours)	Chart Reco. & Reas.	1k recognition-oriented and 1k reasoning-oriented questions
CharXiv (Wang et al., 2024c)	Chart Reco. & Reas.	4k descriptive and 1k reasoning questions (validation set)
MathVista (Lu et al., 2024)	General Reasoning	540 math-targeted and 460 general VQA questions (testmini set)
MATH-Vision (Wang et al., 2024a)	General Reasoning	3,040 math competition problems

Table 11: Summary of benchmarks used in our experiments.

1113 labels, axes and values).

1114 This suggests that the model not only imitates expert rationales but also learns the underlying attention patterns crucial for effective visual reasoning. 1115 The model automatically establishes a synergistic relationship between recognition and reasoning capabilities, understanding what to recognize during 1116 the reasoning process and utilizing these recognition results to guide subsequent reasoning steps. 1117 1118 1119 1120 1121

## 1122 E Visualization of Charts from Different Dataset 1123

1124 We randomly sample several charts from the training set of ChartQA (Masry et al., 2022), ChartBench (Xu et al., 2023), ChartAst (Meng et al., 1125 2024), ChartGemma (Masry et al., 2024b), and REACHQA. The visualization of the results is presented in Figure 8. 1126 1127 1128 1129

## 1130 F Prompt Templates

1131 We present the prompt templates used in our work.

## F.1 Intermediary Code Synthesis 1132

1133 The prompts used for code generation via the Self-Instruct method are presented in Figure 9, and Figure 10 shows the prompts for the Evol-Instruct 1134 method. As illustrated in Figure 11, we utilize four predefined directions to evolve the simple chart-plotting code. 1135 1136 1137 1138

## F.2 Bi-directional Translation 1139

1140 The prompt used for the Self-Repair method is presented in Figure 12. Additionally, the prompt templates for generating reasoning-oriented questions and answers are listed in Figure 13 and Figure 14. The prompt details for generating recognition-oriented questions and answers are listed in Figure 15 and Figure 16. 1141 1142 1143 1144 1145 1146

## F.3 Quality Assurance 1147

1148 The prompt details for rating charts and Q&A are illustrated in Figure 17 and 18. 1149

## F.4 Evaluation 1150

1151 In the evaluation process, we utilize the LLM-as-a-judge method. The detailed prompt template is illustrated in Figure 19. 1152 1153

Models	Avg. ( $\uparrow$ )	ChartQA	ChartBench	ChartX	REACHQA	CharXiv	MathVista	MATH-V
		QA	Binary NQA	QA	Reas. Reco.	Reas. Desc.	Math General	QA
<b>Chart-augmented Multimodal Large Language Models</b>								
ChartMoE-7B	40.62	80.48	60.70	38.05	39.53	12.60 37.40	27.50 40.38	36.11 60.65 13.42
<b>Open-Source Multimodal Large Language Models</b>								
Qwen-2-VL-7B	47.36	81.40	76.15	41.76	47.03	17.60 42.00	31.10 55.13	48.15 64.35 16.28
+ REACHQA (All)	52.00 (+9.8%)	84.96	79.00	48.90	52.20	26.80 53.60	34.60 59.25	49.26 64.78 18.68

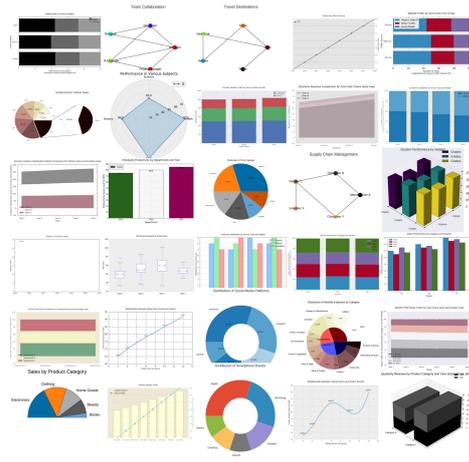
Table 12: Additional evaluation results of recent models on seven benchmarks.



Figure 7: An example of **attention visualization** from the ChartQA dataset. The top row shows the results from the vanilla LLaVA-Next-Llama3-8B model, while the bottom row displays the results from our fine-tuned model. For each output, we present the attention distribution (highlighted zones) at **three key steps**, calculated by averaging the attention values of all tokens in each step.



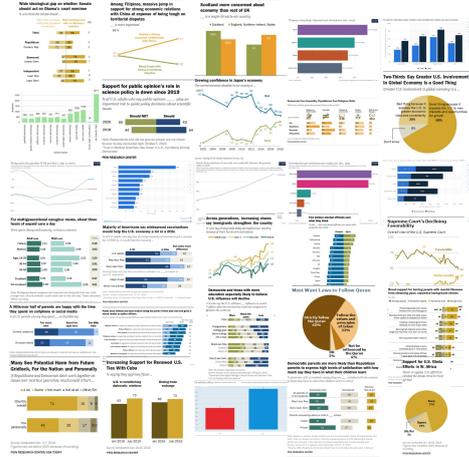
(a) **ChartQA** contains 3 types of charts collected from 4 websites.



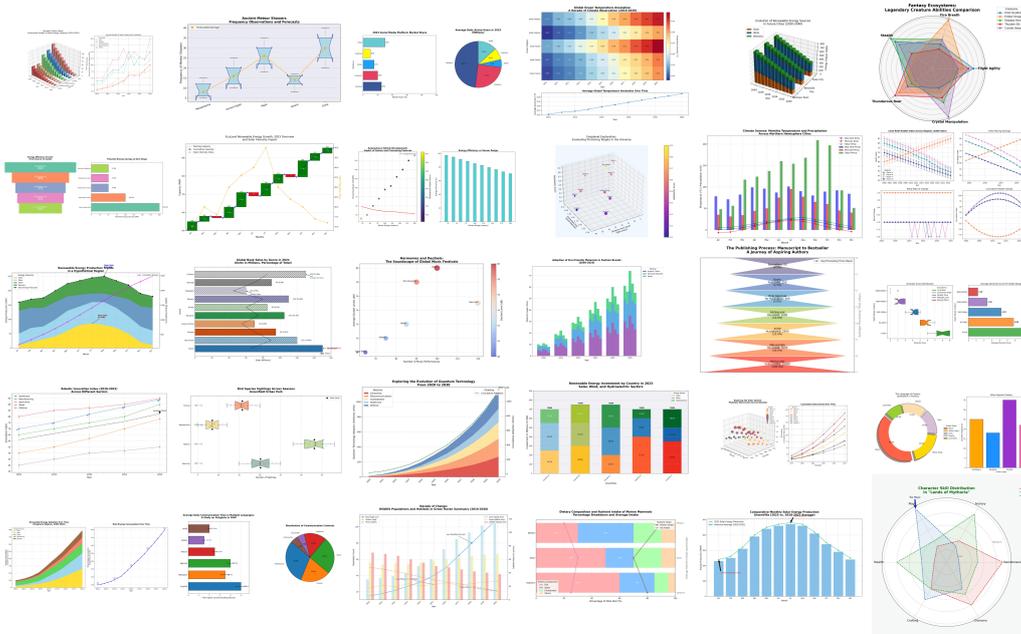
(b) **ChartBench** contains 9 types of synthetic charts but no visual complexity.



(c) **ChartAssistant** contains 9 types of synthetic charts but no visual complexity.



(d) **ChartGemma** contains charts collected from boarder websites.



(e) **REACHQA** contains 10 types of charts and more complex variations.

Figure 8: **Visualizations** of different chart-related training datasets. As shown, REACHQA and ChartGemma exhibit higher chart richness compared to several other datasets. But the charts in ChartGemma require manual collection from multiple sources (Masry et al., 2024b).

**User:**

As a Matplotlib expert, you are asked to write a new Python plotting script. This script will be used to generate a type-specific chart with artificial data. Here are the requirements:

1. There are several script examples from which you can draw inspiration, but try not to repeat patterns already shown in the examples to maximize diversity.
2. Use the Matplotlib library in Python for plotting. You can use auxiliary libraries such as Numpy, but make sure the code works!
3. The type of chart you need to plot is {type}. Therefore, everything you create must be adapted to fit this type of chart.
4. The topic of the chart can be anything you like, for example, {topic1}, {topic2}, etc.
5. Based on the given chart type and the topic you choose, you need to construct a suitable backstory, which should be reflected in the title, labels, legend, etc.
6. Based on the backstory, construct contextual data inputs in the form of Python lists or Numpy arrays. Information contained in the data can be adapted as appropriate to fit the type of chart.
7. You must not use random() to construct the data, as it needs to be explicitly created regardless of your chart type and topic.
8. Be as imaginative and creative as possible in drawing the chart, both in terms of data and plotting details.

Here are three examples to consider:

```
{demo1}  
{demo2}  
{demo3}
```

Now, let's take this task step by step. First, we have to plan out the title and backstory of the chart and create data based on the above.

**Assistant:** {model\_response}

**User:**

Please complete the entire script by plotting a chart based on the data generated. Here are some highlighted requirements and notes.

Requirements:

1. If you find that the generated data is not appropriate while plotting the chart, modify it further as needed.
2. The information on the chart should be complete enough to be understandable, but avoid including the full backstory or too much text in the figure.
3. Avoid occlusion of visual elements. If necessary, automatically adjust the image layout before plt.show() using tight\_layout().
4. If the text in the chart is too long, find a way to make it all visible instead of overlapping. If the title is too long, you can break it into multiple lines.
5. Once again, be as imaginative and creative as possible in creating the details of the chart.
6. Above all, double-check to ensure the code works. Reduce unnecessary comments and focus on functionality.

Now, generate your final plotting script in a single python code block.

Figure 9: Prompt template for code generation via Self-Instruct method.

**User:**

As a Matplotlib expert, you are asked to optimize a Python plotting script to make the plotted chart more complex. The script will be used to generate charts for a mathematical test, so you should make it a little more challenging.

This is the code you need to optimize:

{code}

Here's what I'd like you to do to optimize the chart: {direction}

Now, let's take this task step by step. First, please read the given code carefully and analyze the chart it draws. Then, think about your optimization ideas with the given directions.

In this step, you don't need to give the final code, only show the design ideas.

**Assistant:** {model\_response}

**User:**

Please implement the final optimized script based on the above design ideas combined with the original code.

Remember:

1. Avoid visual elements that obscure each other, e.g., legends, labels. Automatically adjust the image layout before `plt.show()` using `tight_layout()`. if necessary.
2. If the text in the chart is too long, find a way to make all the text show up instead of overlapping. If the title is too long, you can break it into multiple lines.
3. Be as imaginative and creative as possible in creating details of the chart, but don't make the chart redundant just to cope.
4. If you are adding a new plot, take care that the chart is complete with all the elements, such as labels, axes, legends, and colors, unless it is intended to be shared with the original chart.
5. If you are adding a new plot, carefully construct meaningful data and consider whether to give the new sub-plot a sub-title.
6. You must not use `random()` to construct the data, as it needs to be explicitly constructed regardless of your chart type and topic.
7. Above all, double-check to make sure the code works. Reduce unnecessary comments and focus on functionality.

Now, generate your optimized plotting script in a single python code block.

Figure 10: Prompt template for code generation via Evol-Instruct method.

**Evolution Direction:**

- Increase the size of the input data or the number of data groups as appropriate so that it requires a higher level of mathematical understanding. Note if there is a sum requirement.
- Try changing or adding some visual elements to make visual effect better. The elements you add must make sense and not be redundant.
- Incorporate an overlay plot of a different type on the original chart. Use related but not identical data for the added plot.
- Extend an additional subplot of a different type beside the original chart (2 in total). Use related but not identical data for the added plot.

Figure 11: Predefined evolution directions for Evol-Instruct method.

**User:**

As a Python and Matplotlib expert, you have been asked to fix the following code.

The error code is:

{code}

The code reports the following error message when run: {error}

Please analyze the error first, and then provide the revised code within a single Python code block. There should only be one Python code block in your response, containing the complete revised code.

Figure 12: Prompt template for Self-Repair.

**User:**

You are both an expert Matplotlib plotter and a professional maths teacher. Now, you are asked to generate a mathematical reasoning question about a given chart. This chart and question will be used as a question on this year's college admissions examination. As a question writer, you need to ensure that the question is challenging yet fair, testing the student's ability to analyze data, interpret trends, and apply mathematical concepts.

First, please read the following plotting script in Python, try to visualize the figure in your mind and understand the meaning of the chart. After you've analyzed this chart, we'll start generating the associated question.

Here are some tips for you:

1. The plotting script (including the code itself, data mapping and labels) is absolutely correct, and you can trust it completely.
2. The question needs to be based on the chart type, chart topic, and the given data. It can relate to the chart as a whole or to localized details, so you need to look closely.
3. The question should be challenging, requiring visual observation skills and mathematical reasoning skills. So, you need to have a deep understanding of the chart.
4. If there is no data annotation in the figure, try not to generate questions that require too much numerical recognition to reduce inconsistent answers due to visual errors.
5. If some numerical recognition is needed, choose distinguishable colors, lines, heights, and other features that make it easy to estimate without data annotation.
6. You don't need to describe what the chart shows in the question text, including values, labels, etc. This can be left to the student to recognize.

Here is the plotting script:

{code}

Now, please generate 4 questions at a time, each of which needs to look at a different aspect of the chart.

Your output needs to follow this JSON format, and no other text included:

```
{"question_list": ["the question you generate"]}
```

Figure 13: Prompt template for generating reasoning-oriented questions.

**User:**

You are both a Matplotlib graphing expert and a professional math teacher. Now, you have been asked to generate an answer to a given chart and question. This chart and question will be used as a question on this year's college admissions examination. As the answer writer, you need to ensure that the answer is correct, detailed, and educational.

First, please read the following plotting script in Python, try to visualize the figure in your mind and understand the meaning of the chart. After you've analyzed this chart, we'll start generating the answer.

Here is the plotting script:

{code}

Here are some tips for you to generate the answer:

1. First and foremost, the answer needs to be based on the chart information.
2. In the answer, you will also need to solve the question step-by-step, including reasoning steps and recognition steps (but keep concise).
3. You need to explicitly involve a final answer; the type of answer can be a certain number, a noun, or Yes/No, etc.
4. The answer should contain multiple reasoning or calculation steps and be presented in an understandable and educational paragraph.
5. NEVER include any information relating to the Python script in the answer text, as students will ONLY have access to the plotted figure.

Here is the question: {question}

Your output needs to follow this JSON format, and no other text should be included:

{“analysis”: “your analysis about the script and question”, “answer”: “your step-by-step answer”}

Figure 14: Prompt template for generating reasoning-oriented answers.

**User:**

You are both an expert Matplotlib plotter and a professional maths teacher. Now, you are asked to generate a recognition-oriented question about a given chart. This chart and question will be used as a question on this year's elementary math examination to test students' ability to read charts.

First, please read the following plotting script in Python, try to visualize the figure in your mind and understand the meaning of the chart. After you've analyzed this chart, we'll start generating the associated question.

Here are some tips for you:

1. The plotting script (including the code itself, data mapping, and labels) is absolutely correct and you can trust it completely.
2. Descriptive questions are questions that can be answered based on basic chart information, such as titles, labels, tick marks, colors, etc.
3. The generated Q&A needs to be based on the chart type and data. It should be answerable through visual observation.
4. If there is no data annotation in the figure, try not to generate questions that require too many numerical recognitions to reduce inconsistent answers due to visual errors.
5. If some numerical recognition is needed, choose distinguishable colors, lines, heights, and other features that make it easy to estimate without data annotation.
6. You don't need to describe the content of the figure in the question text. This can be left for students to think about.
7. This question needs to explicitly involve a final answer; the type of answer can be a certain number, a noun, or Yes/No, etc.
8. NEVER include any information relating to the Python script in the question or answer, as students will ONLY have access to the plotted figure.

Here are some examples of recognition-oriented questions:

- How many colors are used in the chart? How many city categories are in the chart?
- What's the leftmost value of the bar in China? And what is the value of the bar next to it?
- For the subplot at row 2 and column 1, what is the minimum value of the solid line?
- Which name does the second-largest sector represent? What is its value?
- Does the blue triangle in the chart represent a higher value than the red circle?

Here is the plotting script:

{code}

Now, please generate 4 questions at a time, each of which needs to look at a different aspect of the chart.

Your output needs to follow this JSON format, and no other text included:

```
{“question_list”: [“the question you generate”]}
```

Figure 15: Prompt template for generating recognition-oriented questions.

**User:**

You are both a Matplotlib graphing expert and a professional math teacher. Now, you have been asked to generate an answer to a given chart and question. This chart and question will be used as a question on this year's elementary math examination to test students' ability to read charts. As the answer writer, you need to ensure that the answer is correct, detailed, and educational.

First, please read the following plotting script in Python, try to visualize the figure in your mind and understand the meaning of the chart. After you've analyzed this chart, we'll start generating the answer.

Here is the plotting script:

{code}

Here are some tips for you to generate the answer:

1. First and foremost, the answer needs to be based on the chart information.
2. In the answer, you will also need to solve the question step-by-step, including reasoning steps and recognition steps (but keep concise).
3. You need to explicitly involve a final answer; the type of answer can be a certain number, a noun, or Yes/No, etc.
4. The answer should contain multiple reasoning or calculation steps and be presented in an understandable and educational paragraph.
5. NEVER include any information relating to the Python script in the answer text, as students will ONLY have access to the plotted figure.

Here is the question: {question}

Your output needs to follow this JSON format, and no other text should be included:

{"analysis": "your analysis about the script and question", "answer": "your step-by-step answer"}

Figure 16: Prompt template for generating recognition-oriented answers.

**User:**

<image>

You are a strict Matplotlib plotter and have been asked to evaluate the given chart. Rate the chart from 1 to 5 based on these criteria:

**1 point:** This chart is the poorest in quality and fails to accurately represent any relevant data. It is characterized by a complete breakdown in visual representation; elements are cluttered, text heavily overlaps, legend is missing, or large areas are left blank, making the chart unreadable. The design shows no understanding of effective data visualization practices.

**2 points:** The chart displays incorrect or irrelevant visual elements, with significant inaccuracies that misrepresent the data. The layout suffers from clutter, substantial overlapping of text and other visual elements, such as the legend or labels, and poorly designed axes that result in uneven distribution, severely impeding accurate interpretation.

**3 points:** This chart represents some correct data points but makes basic errors in visual representation. It may use misleading scales, inappropriate chart types, omit key data. Visual clutter and overlapping elements, such as text obscuring parts of the chart or sub-diagrams overlapping each other, detract from the chart's clarity and readability.

**4 points:** The chart accurately represents most of the major data points and important details of the dataset. Minor visual errors exist, such as slight occlusions of text or sub-optimal positioning of elements like legends or labels, but these do not significantly affect the overall accuracy or readability. The chart demonstrates a good understanding of effective visualization techniques but could still be improved in terms of visual layout and the balance of details.

**5 points:** This is an exemplary chart that perfectly encapsulates all critical data points and relationships with outstanding visual clarity and no occlusions. It demonstrates a thorough understanding of data visualization techniques, making excellent use of space and visual elements. The chart is informative, clear, engaging, and free from any visual errors.

Score the chart on this scale, providing a short analysis and a single value. Your response should be in the format:

Analysis: (your analysis)

Rating: (int)

Figure 17: Prompt template for rating the chart quality.

**User:**

<image>

You are a visual question answering (VQA) data annotator. Your task is to review the following chart and question, and determine if the answer is correct based on the information in the chart. You should carefully analyze the chart, taking into account all relevant data points, labels, and trends. Then, conduct an in-depth analysis to determine if there are any unreasonable or incorrect aspects in the figure, question, or answer.

Specifically, consider the following points:

1. Are the provided question and answer relevant to the chart? Can the answer be found in the chart?
2. Do the colors in the charts and questions correspond correctly? Are there instances where the colors are incorrectly referred to?
3. Do the data in the charts and questions correspond correctly? Are there any errors in the data or misalignment of information?
4. Is the provided answer correct? Are there any logical errors or unreasonable points?
5. Apart from the points listed above, is there anything else in this question and answer that doesn't make sense?

Here is the question and answer about the given chart:

Question: {question}

Answer: {answer}

You are asked to provide a short analysis and decide whether to keep the example. Your response should be in the format:

Analysis: (your analysis)

Decision: (yes/no)

Figure 18: Prompt template for rating Q&A quality.

**User:**

Compare the ground truth with the prediction from AI model and determine if the prediction is correct. The question is about an image, which we have not given here. You need to determine whether the model's prediction is consistent with the ground truth. No points will be awarded for wrong answers, over answers or under answers. The reasoning process in the prediction does not need to be considered too much, you only need to determine if the final answer is consistent. There are times when the answer may have a different form of expression and some variation is acceptable.

## Question: {question}  
## Ground Truth: {answer}  
## Prediction: {prediction}

Now, let's analyze it and then provide your judgment. Your response must follow the format below:

Analysis: (analyze the correctness briefly)  
Correctness: (Yes or No)

Figure 19: Prompt template for evaluating the model prediction with LLMs.