

DomainCQA: Crafting Expert-Level QA from Domain-Specific Charts

Ling Zhong^{1*} Yujing Lu^{1*} Jing Yang^{1*} Weiming Li^{1*} Peng Wei² Yongheng Wang¹
Manni Duan^{1†} Qing Zhang^{1†}

¹Zhejiang Lab ²National Astronomical Observatory, Chinese Academy of Science

{zhongling, luyujing, yangjing0128, liwm, wangyh, duanmanni, qing.zhang}@zhejianglab.org
weipeng01@nao.cas.cn

Abstract

Chart Question Answering (CQA) benchmarks are essential for evaluating the capability of Multimodal Large Language Models (MLLMs) to interpret visual data. However, current benchmarks focus primarily on the evaluation of general-purpose CQA but fail to adequately capture domain-specific challenges. We introduce DomainCQA, a systematic methodology for constructing domain-specific CQA benchmarks, and demonstrate its effectiveness by developing AstroChart, a CQA benchmark in the field of astronomy. Our evaluation shows that chart reasoning and combining chart information with domain knowledge for deeper analysis and summarization, rather than domain-specific knowledge, pose the primary challenge for existing MLLMs, highlighting a critical gap in current benchmarks. By providing a scalable and rigorous framework, DomainCQA enables more precise assessment and improvement of MLLMs for domain-specific applications.

1. Introduction

The success of Multimodal Large Language Models (MLLMs) has sparked growing interest in their ability to process and analyze scientific charts, which play a crucial role in conveying complex research data. [1, 8, 10, 11, 22, 25, 26, 34, 38]. Among various chart-related tasks, Chart Question Answering (CQA) has emerged as a fundamental challenge, requiring MLLMs to extract, interpret, and reason about chart-based information in response to natural language queries.

However, a key open question remains: *Do these models truly engage with the content of the chart, or do they rely on superficial correlations?* Addressing this uncertainty requires a well-designed benchmark that systematically evaluates the capability of models in chart interpretation, nu-

merical reasoning, and insight derivation.

Existing CQA benchmarks are mostly general-purpose. FigureQA [17], DVQA [16], LEAF-QA [6], and PlotQA [33] predominantly feature basic chart types, such as bar chart, line chart, and scatter chart, limiting both the variety of visual representations and the complexity of question design. While CharXiv [39] and SciCap [15] expand chart diversity by collecting charts from scientific papers, their questions remain focused on visual recognition and basic comprehension, without requiring advanced reasoning or domain knowledge. These benchmarks overlook the significant differences in charts across domains and the crucial role of domain knowledge in CQA. A comprehensive benchmark assesses not only chart recognition but also advanced analytical reasoning and inference. To our knowledge, no comprehensive CQA benchmark specifically evaluates MLLMs on scientific chart interpretation within a specialized domain.

To bridge this gap, we introduce **DomainCQA**, a methodology for building domain-specific CQA benchmarks, which can be applied to fields such as astronomy, biology, geophysics, etc. DomainCQA follows a structured approach to benchmark construction, consisting of three key steps: chart selection, question-answer pair generation, and benchmark quality assurance. The benchmark includes two types of question-answer (QA) pairs: fundamental question-answer (FQA) pairs, which assess MLLMs' general chart interpretation skills, and advanced question-answer (AQA) pairs, which evaluate enhanced domain-specific reasoning. Charts are carefully selected based on complexity and relevance to ensure a well-balanced dataset. Diverse QA pairs are then systematically generated from these charts. Finally, DomainCQA uses a crowdsourced evaluation framework to validate the QA pairs, ensuring high quality.

To demonstrate its effectiveness, we construct an astronomical benchmark called **AstroChart** using the DomainCQA methodology, tailored for the astronomy domain. AstroChart comprises 484 charts and 1890 QA pairs, in-

*Equal contribution

†Corresponding author

cluding 1509 FQA pairs focused on graphical interpretation and 381 AQA pairs assessing knowledge-based interpretation. We evaluate 17 state-of-the-art (SOTA) MLLMs on this benchmark to measure their ability to interpret astronomical charts, with further analysis of successes and failure patterns.

Our contributions include: (1) proposing DomainCQA, a general methodology for constructing domain-specific Chart QA benchmarks with a three-phase pipeline; (2) introducing AstroChart, the first domain-specific CQA benchmark for astronomy; (3) conducting a comprehensive evaluation of 17 SOTA MLLMs on AstroChart to assess their chart understanding capabilities; and (4) revealing chart reasoning and combining chart information with domain knowledge for deeper analysis and summarization as the primary challenges in domain-specific chart understanding. The AstroChart benchmark will be made available through our GitHub repository.

2. Related Work

Here we provide an overview of MLLMs’ progress in chart understanding and examine existing CQA benchmarks, highlighting their limitations in evaluating domain-specific scientific chart interpretation.

2.1. MLLMs for Chart Understanding

Recent advancements in MLLMs have significantly improved the ability of chart understanding. These developments span both proprietary and open-source models, each introducing unique enhancements.

Proprietary MLLMs such as GPT-4o[1], Claude 3.5[3], Qwen-VL[4], and GLM-4V[12] have demonstrated strong performance in chart interpretation and multimodal reasoning. The exact techniques used to enhance chart understanding remain proprietary, but their capabilities indicate that training on chart-specific datasets [32], exposure to structured financial and scientific visualizations [46], and leveraging multimodal alignment strategies [45] could significantly improve performance. Despite these strengths, challenges persist in visual-numerical comprehension and cross-modal reasoning.

Open-source MLLMs are advancing rapidly, providing researchers and developers with accessible and customizable alternatives to proprietary models. These models can be broadly divided into two categories: those that enhance general vision-language capabilities and those that specialize in chart-specific understanding.

Vision-Language Advances: mPLUG-Owl[42–44], SPHINX[22], and LLaVA[25–27] enhance vision-language alignment through vision-language projectors. InternVL[10] improves vision encoding to narrow the performance gap between open-source and proprietary models. CogVLM[14, 38] refines feature fusion, while

MiniCPM[41] optimizes efficiency for edge applications. At large scale, Pixtral Large[2] stands as the most powerful open-source MLLM for general vision tasks and high-quality chart interpretation.

Chart-Specific MLLMs: Models such as UniChart[30], ChartLlama[13], Matcha[23], ChartInstruct[31], ChartAssistant[32], and TinyChart[45] are MLLMs specifically designed for chart understanding. These models fine-tune vision-language architectures with large-scale chart instruction data[13, 31, 32], optimizing structured insight extraction from visualized data. Among them, TinyChart employs Program-of-Thought (PoT) reasoning to enhance numerical comprehension, achieving strong performance in chart-related tasks despite its compact 3-billion-parameter size.

2.2. Benchmarks for CQA Evaluation

A CQA benchmark consists of two key components: charts and corresponding QA pairs, both essential for evaluating the chart comprehension capabilities of a model.

Initial efforts in CQA benchmarking focused on synthetic charts, with datasets such as DVQA[16] and FigureQA[17] generating both chart structures and underlying data synthetically. Subsequently, PlotQA[33], LEAF-QA[6], and LEAF-QA++[37] began using real-world numerical data while still relying on synthetic chart visualizations. In both cases, QA pairs were generated using fixed templates, limiting linguistic diversity and natural language variability.

A shift toward real-world chart datasets emerged with ChartQA[29], OpenCQA[18], MMC-Benchmark[24], and CharXiv[39], which introduced charts sourced from real-world datasets. These benchmarks cover a wide range of standard scientific chart types. To improve question diversity, they adopted human-authored QA pairs, with OpenCQA pioneering open-ended CQA, where models generate free-form answers instead of selecting from predefined choices. With the advancement of LLMs, some studies have even leveraged them to generate diverse, high-quality, open-ended QA pairs, such as SciGraphQA[20] and ChartX[40].

Despite these efforts, existing benchmarks primarily focus on general-purpose or broad scientific content, making them insufficient for evaluating MLLMs’ ability to interpret domain-specific scientific charts. This gap underscores the need for dedicated benchmarks tailored to specialized fields, which we address in the next section.

3. Building a Domain CQA Benchmark

We propose DomainCQA, a methodology for constructing domain-specific CQA benchmarks. It provides a structured approach for developing benchmarks tailored to different

scientific domains. DomainCQA categorizes QA pairs into two types:

- **Fundamental Question-Answer (FQA) Pairs:** Corresponding to QA pairs in existing general-purpose CQA benchmarks, evaluating basic chart comprehension and reasoning skills.
- **Advanced Questions-Answer (AQA) Pairs:** Designed to assess MLLMs’ ability to interpret complex encodings, recognize non-trivial trends, and understand domain-specific notations in charts.

This approach ensures that the developed benchmarks evaluate both basic chart comprehension and domain-specific reasoning, providing a robust assessment of MLLMs’ scientific chart understanding.

In summary, DomainCQA follows a structured process consisting of chart selection, question-answer pair generation, and benchmark quality assurance. Fig. 1 provides an overview of the benchmark construction process.

3.1. Chart Selection

DomainCQA selects charts separately for FQA pairs and AQA pairs, ensuring that each set aligns with the specific evaluation requirements of its question type.

Charts for FQA pairs are randomly sampled from domain-related scientific publications, including peer-reviewed journal articles, conference proceedings, and research preprints, to construct a dataset that represents the diversity of scientific visualizations in the domain.

This random sampling is designed to preserve the complexity distribution of charts found in this domain’s publications, ensuring that the dataset reflects the real-world difficulty levels encountered in scientific literature. Maintaining this distribution is crucial for evaluating MLLMs in a setting that aligns with actual domain-specific chart interpretation challenges.

Specifically, we follow a three-step chart selection process:

1. **Chart collection:** We compile a comprehensive set of charts from scientific publications in the domain, filtering out non-informative visualizations (such as decorative figures and flowcharts) to retain only data-driven charts relevant to scientific analysis.
2. **Complexity quantification:** To calculate chart complexity, we introduce the Chart Complexity Vector (CCV), which evaluates a chart across ten key aspects of complexity: {Annotation, Color, Legend, Pattern} (visual complexity), {Axis, Element, Formula, Scale} (data interpretation complexity), and {Subplot, Type} (structural complexity). Each aspect is assigned a binary score of either 0 (simple) or 1 (complex). For example, in the Type aspect, common chart types such as bar, line, and scatter plots are assigned a score of 0, while less common and more intricate chart types receive a score of 1.

As a result, the complexity of a chart is represented as a ten-dimensional binary vector. The definition and examples of CCV scoring can be found in Appendix A.1 and Appendix A.3.

3. **Chart sampling:** We apply Gibbs sampling[5], a Markov Chain Monte Carlo method, to select charts while preserving the original complexity distribution. This iterative process progressively refines the selection, ensuring that the sampled charts retain the same complexity characteristics as those observed in the domain while minimizing unintended biases.

The detailed procedure of charts for FQA pairs selection is outlined in Algorithm 1.

Algorithm 1 Gibbs sampling for chart selection in FQA pairs

Require: Chart dataset D from a domain’s scientific publications with precomputed $CCV(c)$ for each chart c in D

Ensure: Selected benchmark charts S

- 1: **Initialize:** Randomly select an initial set $S \subset D$ of size $target_size$
 - 2: **repeat**
 - 3: **for** each chart $c^* \in S$ **do**
 - 4: Select an aspect α in CCV
 - 5: Fix all other aspects and sample a new value $v \sim P(\alpha|S, D)$
 - 6: Find a candidate chart $c_{new} \in D$ where $\alpha(c_{new}) = v$
 - 7: **if** $CCV(S \cup \{c_{new}\} - \{c^*\})$ maintains distribution **then**
 - 8: Replace c^* with c_{new} in S
 - 9: **end if**
 - 10: **end for**
 - 11: **until** convergence (complexity distribution stabilizes)
- return** S
-

Charts for AQA pairs are inherently more complex than those for FQA pairs, as they require specialized domain knowledge for accurate interpretation. To investigate this complexity, we conduct an experiment to assess whether randomly selected charts can effectively produce high-quality AQA pairs.

From 20 randomly selected astronomical charts, we generate 20 AQA pairs using MLLM, while incorporating the corresponding paper’s contextual content as background knowledge. To assess their quality, we conduct an evaluation study with four domain experts, where each expert independently rates all 20 AQA pairs based on two key criteria:

1. **Professional:** The question is clearly phrased and requires expert-level domain knowledge to answer.
2. **Accuracy:** The AI-generated answer is scientifically ac-

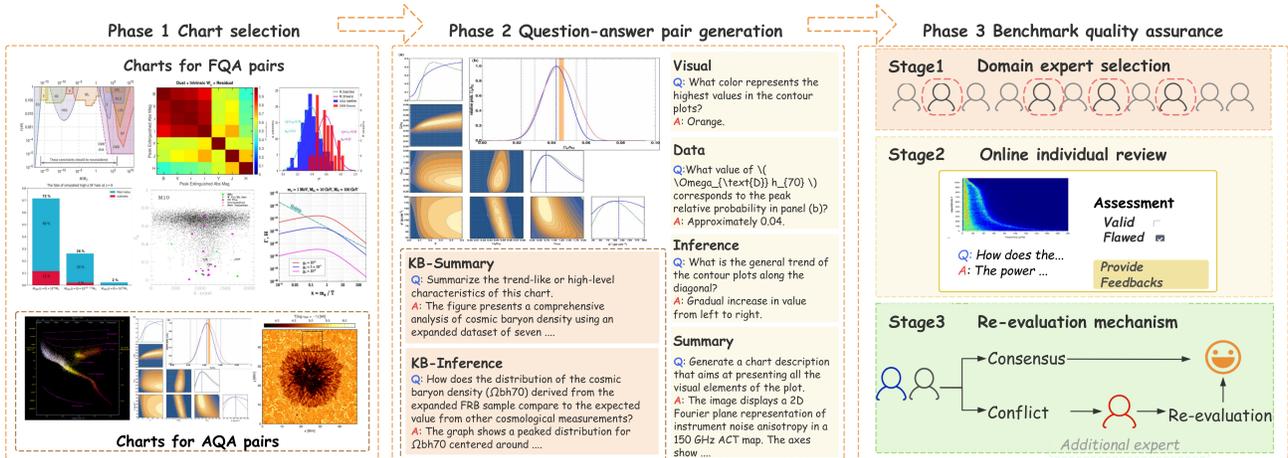


Figure 1. The construction of the DomainCQA benchmark follows a three-step pipeline: **Chart selection:** FQA pairs charts are randomly sampled from scientific publications while maintaining domain complexity, whereas AQA pairs charts are key ‘chart abstracts’ identified using Chain-of-Thought (CoT) reasoning and a Voting Schema (VoT). **Question-answer pair generation:** We generate FQA pairs and AQA pairs by leveraging charts and MLLMs. **Benchmark quality assurance:** All question-answer pairs are reviewed through a structured expert evaluation process, ensuring accuracy through independent assessment and consensus-based refinement.

curate and aligned with domain knowledge.

Inter-rater agreement was measured using Krippendorff’s Alpha[19], revealing very low agreement on both aspects (Professional: $\alpha = -0.11$, 95% CI $[-0.46, 0.24]$; Accuracy: $\alpha = 0.01$, 95% CI $[-0.24, 0.21]$). Since the raters are domain experts, this level of disagreement suggests that the generated QA pairs contained inconsistencies or ambiguities that made evaluation difficult according to a consistent domain standard.

Obviously, this inconsistency or ambiguity stems from the selected charts, as the question-answer pair generation process was identical across all charts. Thus, to ensure the quality of the generated AQA pairs, we aim to select charts that are closely tied to domain knowledge, such as the *chart abstract*, a highly informative visualization that encapsulates the study’s key findings. Fig. 2 is an example of a chart abstract from the field of astronomy and provides the related conclusion corresponding to the chart.

To locate the chart abstract of a paper using MLLMs, we employ a Chain-of-Thought (CoT) reasoning framework combined with a Voting Schema (VoT) to ensure robust and scientifically grounded selection.

In brief, the CoT process begins by extracting the caption of each chart, surrounding description, and key textual sections (e.g., abstract and conclusion). Using this extracted information, the MLLM applies CoT reasoning to generate a concise summary for each chart, as well as the summary of the entire paper. Finally, the best-matching chart is identified by comparing the summary of each chart with the overall summary of the paper.

Selection robustness is further enhanced by integrating

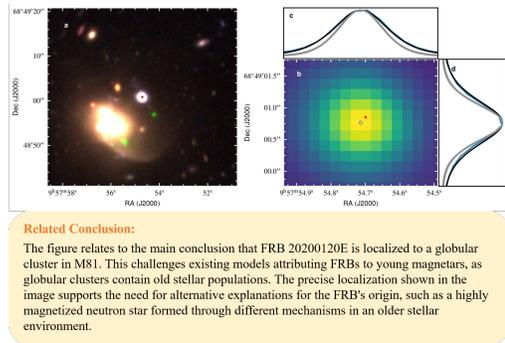


Figure 2. An example of a *chart abstract* from the field of Astronomy.

VoT, where multiple MLLMs independently identify the chart abstract and only the one selected by a majority vote is retained. Algorithm 2 presents the details of the charts for the AQA pairs selection process.

3.2. Question-Answer Pair Generation

We generate two types of QA pairs using MLLMs: FQA pairs and AQA pairs. Specifically, FQA pairs are divided into four categories, and AQA pairs fall into two categories, each derived from their corresponding charts. Below are the details of these types:

- FQA pairs:
 - *Visual:* Assesses the model’s ability to interpret graphical components of the chart.
 - *Data:* Evaluates the model’s ability to retrieve data and perform math reasoning from the chart.

Algorithm 2 CoT&VoT for chart selection in AQA pairs

Require: Research paper P with N charts $\{C_1, C_2, \dots, C_N\}$

Require: Set of k Multimodal Large Language Models (MLLMs) $\{M_1, M_2, \dots, M_k\}$

Ensure: Selected the chart abstract C^* of P

```
1: for each MLLM  $M_j \in \{M_1, \dots, M_k\}$  do
2:   Extract abstract and conclusion
3:   Generate the paper’s summary  $P_j$  with  $M_j$ 
4:   for each chart  $C_i$  do
5:     Extract caption and surrounding description
6:     Generate its summary  $S_{ij}$  using  $M_j$ .
7:     Compute relevance score  $R_{ij}$  based on alignment between  $S_{ij}$  and  $P_j$ .
8:   end for
9:   Select the most relevant chart,  $C_j^*$  based on  $R_{ij}$ .
10: end for
11: Identify the chart abstract  $C^*$  appearing in the majority of  $\{C_1^*, \dots, C_m^*\}$  across all models.
12: return  $C^*$  (final selected chart abstract)
```

- *Inference*: Tests the model’s ability to analyze patterns and relationships in the chart.
- *Summary*: Measures the model’s ability to generate summaries of the chart’s visual elements.
- AQA pairs (Knowledge-Based reasoning tasks):
 - *KB-Inference*: Assesses the model’s ability to apply domain knowledge for deep reasoning and identifying connections in the chart.
 - *KB-Summary*: Evaluates the model’s ability to generate summaries by combining chart information with domain knowledge.

3.3. Benchmark Quality Assurance

To ensure the reliability of the generated benchmark, DomainCQA also implements a systematic *expert verification framework* (Fig. 1, phase 3) that governs the entire process of reviewing and validating every QA pair. This framework defines the review criteria, expert selection process, evaluation workflows, and consensus mechanisms, all of which are implemented through our online assessment platform.

Each QA pair is assessed on the following criteria:

- **Valid**: The question is well-formed and unambiguous, and the provided answer is correct.
- **Flawed**: The question is unclear, or misleading, or the answer contains factual inaccuracies.

Domain experts independently review all QA pairs within the platform. If discrepancies arise, the framework enforces additional rounds of evaluation and refinement until a consensus is reached. By systematically validating every QA pair, we ensure that DomainCQA generates a high standard and valuable benchmark for downstream tasks.

4. AstroChart: First CQA Benchmark for Astronomy

In this section, we apply the DomainCQA methodology to the field of astronomy, leveraging its framework to construct a domain-specific CQA benchmark: *AstroChart*.

AstroChart is a comprehensive and challenging benchmark designed to evaluate chart understanding in astronomy. It consists of 484 astronomical charts and 1,890 QA pairs, providing a rigorous testbed for assessing MLLMs’ multimodal reasoning and scientific chart interpretation capabilities.

4.1. Chart Selection

Charts for FQA pairs: From all astronomical research papers published on arXiv between 2007 and 2023, we initially extracted 513,000 visual assets, excluding those that were unclear or excessively small. We then filtered out non-chart visuals, such as images and diagrams, ultimately retrieving 457,000 scientific charts from these papers. Then with Gibbs sampling, we selected 303 charts for constructing fundamental questions in AstroChart.

Charts for AQA pairs: To identify suitable chart abstracts in astronomy, we focus on high-quality papers, as influential studies often feature well-structured and informative visualizations. To mitigate selection bias, we select the top 1% most-cited papers per year within each of astronomy’s six major subdomains, as shown in Fig. 3. From these selected papers, we randomly sampled 200 and extracted their chart abstracts using the CoT&VoT method (Algorithm 2). This process was carried out with voting from two proprietary MLLMs, Claude 3.5 and GPT-4o.

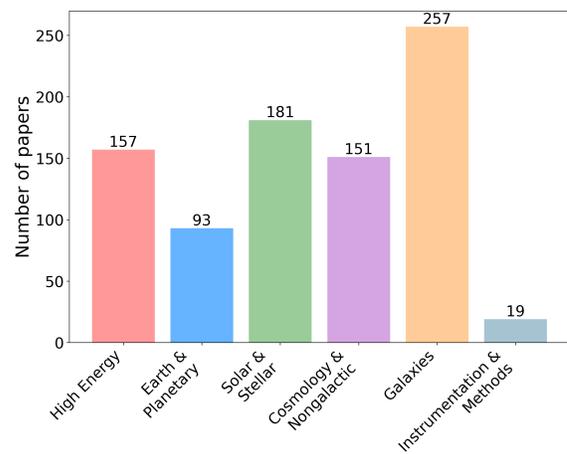


Figure 3. Top 1% most-cited papers selected from six major subdomains of astronomy

In total, AstroChart includes 484 selected charts from

astronomy research, with 303 charts for FQA pairs and 200 charts for AQA pairs, 19 of which overlap. The graphic assets of AstroChart can be found in Appendix B. The overall complexity of each chart is assessed using the L1 norm of its CCV, i.e., the sum of the vector’s elements. Fig. 4 compares the complexity distribution of AstroChart with other generic CQA benchmarks.

Clearly, ChartQA[29], OpenCQA[18], and PlotQA[33] primarily consist of relatively simple charts, with most complexity scores concentrated at lower values. In contrast, CharXiv[39] and AstroChart exhibit broader distributions, with AstroChart containing a higher proportion of complex domain-specific charts, particularly in the 6-10 complexity range. The distribution of CCV scores across ten dimensions in AstroChart can be found in Appendix A.2.

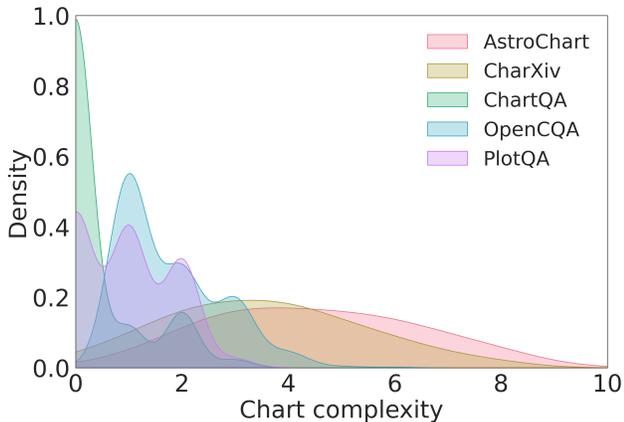


Figure 4. Distribution of chart complexity in benchmarks

4.2. Question-Answer Pair Generation

We use Claude 3.5 to generate question-answer pairs, employing specialized prompts tailored to each category, and ensure the generated question-answer pairs closely aligned with the corresponding charts. Detailed prompt designs are provided in Appendix C.

This process resulted in a total of 1890 QA pairs, including 1509 FQA pairs and 381 AQA pairs. Examples are provided in Appendix D. Tab. 1 presents a detailed distribution of these QA pairs across different categories.

The generated questions exhibit diversity even within a single category. For example, among the 603 FQA pairs in the visual category, 211 focuses on color aspects, 133 addresses structural elements, 214 involves OCR-related inquiries, and 45 pertains to chart-type classification. This variation highlights the range of question types captured within the visual category.

Type of QA Pairs	Category	Aspect	Count
FQA Pairs	Visual	Color	211
		Style	133
		Text	214
		Layout	45
	Data	Point	129
		Interval	102
		Calculation	84
	Inference	289	
	Summary	302	
AQA Pairs	KB-Inference	200	
	KB-Summary	181	
Total			1890

Table 1. Number of questions in categories of AstroChart.

4.3. AstroChart Manual Verification

We first check if the selected charts ensure the generation of high-quality AQA pairs. We perform an evaluation on 20 randomly selected AQA pairs with four astronomy experts and assess the agreement using Krippendorff’s Alpha. We notice a great improvement of agreement in both aspects (Professional $\alpha = 0.63$, 95% CI [0.49, 0.74]; Accuracy $\alpha = 0.64$, 95% CI [0.52, 0.75]).

Then, we conduct a comprehensive verification of AstroChart to ensure complete accuracy. Four astronomy experts review all 1,890 QA pairs in AstroChart through our online assessment platform. Each QA pair is independently evaluated by two randomly assigned experts. In case of disagreement, additional evaluation rounds are conducted until a consensus was reached. Details of the online assessment platform can be viewed in Appendix E.

This rigorous validation process ensures reliability of AstroChart, establishing it as a valuable benchmark for assessing MLLM performance in astronomical chart understanding.

5. Evaluation of MLLMs on AstroChart

In this section, we evaluate 17 SOTA MLLMs, including both proprietary and open-source models, on the AstroChart to assess their ability to interpret astronomical charts and reason about complex visual data. We also aim to highlight areas where future models can focus to enhance performance in domain-specific chart understanding.

5.1. Experimental Setup

Models: The proprietary MLLMs include GPT-4o[1], GLM-4V[12], GLM-4V-Plus[12] and Qwen-VL-Max[4]. The open-source MLLMs, ordered by ascending model size, are TinyChart-3B[45], Deepseek-Janus-Pro-7B[7], Llava1.5-7[25], Llava1.6-Mistral-7B[27], Qwen-VL-Chat-7B[4], MiniCPM-Llama3-V2.6-8B[41], InternVL2-8B[11], mPLUG-Owl2-8.2B[44], Llava1.6-

Vicuna-13B[27], SPHINX-v2-13B[22], CogVLM-19B[14], Llava1.6-Yi-34B[27] and Pixtral-large-124B[2]. Detailed information about the architecture of MLLMs can be found in Appendix F.

Setup: All proprietary models are accessed via API, while open-source models are deployed locally and tested on a single Nvidia A100-80G GPU. We conduct a zero-shot evaluation, providing only the chart and question as prompts for each MLLM to answer questions in AstroChart.

Metrics: we use an evaluation framework designed to handle both numerical responses and open-ended responses, the details can be found in Appendix G:

- For numerical responses, we distinguish between data retrieval and data derivation. For data retrieval (e.g., plotted points or data ranges), we normalize the relative error by the axis range, ensuring a bounded and scale-aware evaluation. Data derivation (e.g., the number of bars, colors, or legends, as well as math reasoning) refers to counting visual elements or performing calculations based on the chart, and correctness is determined by exact numerical matching, ensuring that only fully correct answers are considered accurate.
- For open-ended responses, we use an LLM-based judging framework inspired by [28], in which a dedicated judging LLM, separate from the 17 MLLMs used in the evaluation, assesses the quality of answers generated by the evaluated models. The judging LLM assigns a score between 0 and 1 based on predefined assessment criteria, ensuring consistency and scalability in evaluation.

We adopt the aforementioned metrics evaluation framework and use DeepSeek-V3[9] as the LLM to compute LLM-scoring as accuracy to evaluate the performance of 17 MLLMs. Additionally, we employed other metrics such as ROUGE-L[21], BLEU-4[35], and L3Score[36]. As all metrics exhibit a similar overall trend, we omit them due to space constraints; details can be found in Appendix H.

5.2. Results

The evaluation results of MLLMs on AstroChart, including overall scores as well as performance across different categories of FQA pairs and AQA pairs, are presented in Tab. 2. We also provide Failure Cases of AstroChart in Appendix I.

Each score in the table represents the average accuracy for the corresponding class. For simplicity, we omit the percentage symbol and present the values directly in the table and throughout this subsection.

Overall Performance: The results reveal significant performance differences across models, highlighting their strengths and limitations in astronomical chart understanding. GPT-4o achieves the highest score, demonstrating

the superiority of proprietary models, while Pixtral-Large-124B stands out as the best-performing open-source model, even surpassing most proprietary alternatives. TinyChart-3B ranks the lowest, reinforcing the general trend that larger models tend to perform better.

However, model size alone does not fully determine performance. Pixtral-Large-124B’s strong results, despite being smaller than proprietary models, suggest that architectural advancements and targeted optimizations in chart understanding play a crucial role. These findings indicate that while proprietary models still lead, some open-source models are narrowing the gap, and strategic improvements in model design and training could further enhance MLLM performance in astronomical chart interpretation.

Category Performance: The performance of MLLMs across different question categories, including Visual, Data, Inference, and Summary in FQA pairs, as well as KB-Inference and KB-Summary in AQA pairs, reveals notable differences in difficulty and model capabilities. Fig. 5 illustrates the performance distribution of all models across these categories.

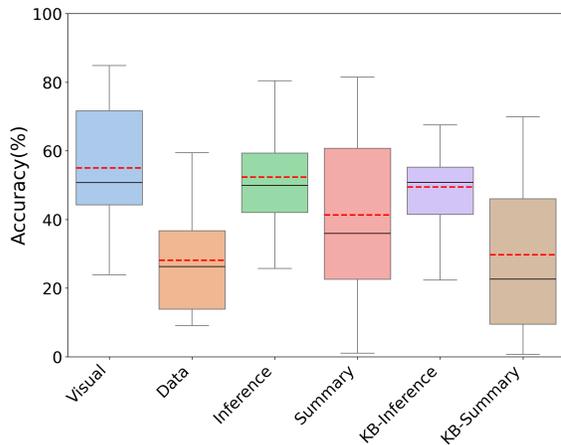


Figure 5. Box plot of accuracy evaluation across six different question categories in AstroChart. Note that the red dashed lines indicate the average accuracy.

Visual questions achieve the highest average accuracy, indicating that models perform relatively well in interpreting graphical components such as colors, styles, text, and layout. This suggests that MLLMs exhibit a strong ability to recognize visual elements and extract fundamental structural information from charts.

Data questions consistently show low performance across all models, as depicted in the box plot. Despite efforts, none of the models perform well in retrieving numerical values or performing math reasoning, with the latter proving even more challenging. This is likely due to factors such as OCR limitations, axis misinterpretation, and difficulties in extracting data from dense or complex charts.

Models	Fundamental Question-Answer Pairs									Advanced Question-Answer Pairs				Overall
	Visual				Data					Inference	Summary	KB-Inference	KB-Summary	
	All	Color	Style	Text	All	Point	Interval	Calculation						
Proprietary Multimodal Large Language Models														
GPT-4o[1]	84.01	85.31	82.37	84.30	81.33	57.05	61.42	62.64	43.57	80.42	81.46	67.56	69.94	75.98
GLM-4V[12]	74.76	77.00	70.22	74.72	78.67	36.71	33.31	47.56	29.05	59.34	63.81	55.23	46.02	59.83
GLM-4V-Plus[12]	74.65	77.56	70.26	74.63	74.22	42.98	43.83	49.78	33.69	61.56	60.73	60.58	59.45	62.36
Qwen-VL-Max[4]	68.96	72.70	60.67	68.64	76.00	45.79	42.10	48.68	47.98	70.00	77.78	60.70	65.41	65.82
Open-source Multimodal Large Language Models														
TinyChart-3B[45]	27.69	45.82	20.81	13.50	28.00	13.67	18.78	12.16	7.74	25.74	1.06	22.33	0.72	17.36
Deepseek-Janus-Pro-7B[7]	65.82	73.80	66.96	55.89	72.00	32.43	36.35	37.84	20.12	57.85	52.62	52.91	28.67	52.00
Llava1.5-7B[25]	23.89	33.92	23.85	12.80	27.11	11.09	11.31	9.71	12.38	42.04	15.03	39.88	9.50	22.41
Llava1.6-Mistral-7B[27]	44.70	58.24	39.85	32.62	49.33	14.96	17.34	12.09	14.76	46.89	22.55	46.51	12.49	32.87
Qwen-VL-Chat-7B[4]	44.26	53.57	40.44	33.60	59.11	13.88	19.64	8.47	11.55	39.58	24.54	41.51	14.59	31.65
MiniCPM-Llama3-V2.6-8B[41]	71.66	76.01	68.81	69.81	68.22	34.57	30.78	46.13	26.67	56.12	60.73	55.23	44.64	57.21
InternVL2-8B[10]	47.93	53.76	43.70	43.64	51.33	26.31	25.65	29.57	23.45	49.90	28.05	45.93	7.46	36.84
mPLUG-Owl2-8.2B[44]	25.58	33.74	25.63	16.21	29.11	9.44	10.67	8.71	8.45	37.75	10.30	36.40	7.68	20.83
Llava1.6-Vicuna-13B[27]	49.24	65.07	44.52	35.79	50.44	15.85	20.46	10.35	15.36	44.15	24.21	46.98	15.25	34.71
SPHINX-v2-13B[22]	31.26	46.31	29.41	19.18	20.89	9.07	15.20	1.51	8.69	37.51	6.59	40.58	3.76	21.83
CogVLM2-19B[14]	65.64	73.87	50.93	65.61	70.89	31.05	33.56	39.63	17.14	53.49	56.92	50.81	33.76	52.14
Llava1.6-Yi-34B[27]	50.78	65.87	44.44	39.53	50.22	23.62	22.56	26.64	21.67	47.99	35.93	54.30	22.65	40.33
Pixtral-large-124B[2]	84.89	84.65	82.44	86.21	86.89	59.47	62.65	65.04	48.10	79.72	79.97	63.60	63.76	75.56

Table 2. Accuracy evaluation on AstroChart benchmark. Bold numbers indicate the best-performing model among proprietary and open-source MLLMs, respectively.

This persistent low accuracy underscores that precise data retrieval remains a significant challenge for MLLMs, particularly in astronomy, where charts often feature non-trivial axis scales and complex visual encodings.

The scores for Inference and KB-Inference are quite similar across models, suggesting that most MLLMs, regardless of size, have enough domain knowledge to perform these tasks effectively. The presence of domain knowledge in the questions may help larger models focus on relevant knowledge areas, leading to consistent and accurate answers. While proprietary models and larger open-source models, such as Pixtral-Large-124B, achieve higher scores due to stronger contextual reasoning and multimodal alignment, even smaller models like Deepseek-Janus-Pro-7B perform well, showing low variance in their results (Fig. 6 top).

KB-Summary exhibits moderately lower performance than Summary. This suggests that it is challenging for models to extract information from charts while also identifying and retrieving relevant domain-specific knowledge. Proprietary and larger open-source models, benefiting from richer training data and stronger multimodal alignment, still perform better overall, but the gap between KB-Summary and Summary indicates the complexity of combining chart interpretation with domain knowledge without explicit textual cues (Fig. 6 bottom).

5.3. Discussion

Our results indicate that current MLLMs struggle with combining chart information with specific domain knowledge for more comprehensive analysis and accurate summarization without textual cues. In addition, issues with data retrieval and math reasoning need improvement. Therefore, to enhance MLLMs’ capabilities, the focus should be mainly on improving their data processing abilities and integrating

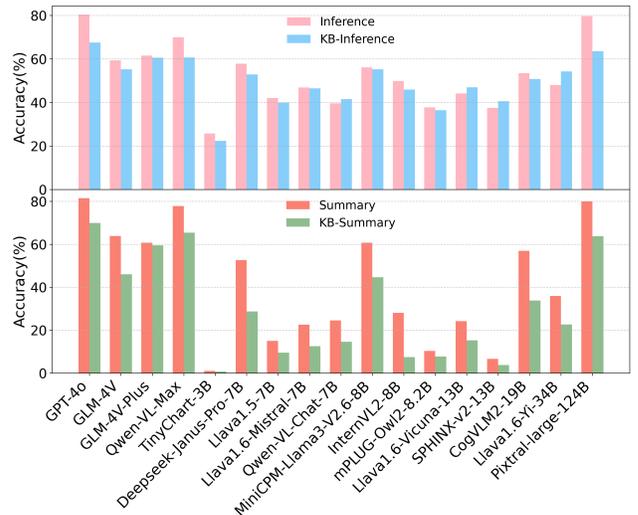


Figure 6. Accuracy of MLLMs on different evaluation categories: Inference vs. KB-Inference and Summary vs. KB-Summary in AstroChart.

chart and domain knowledge for more effective analysis and summarization.

6. Conclusion & Future Work

This paper introduced DomainCQA, a structured methodology for constructing domain-specific CQA benchmarks, and applied it to develop AstroChart, the first Chart QA benchmark for astronomy. Through a comprehensive evaluation of 17 SOTA MLLMs, we identified chart reasoning and the combination chart information with domain knowledge for deeper analysis and summarization as the primary bottlenecks in domain-specific chart understanding. While these models exhibit strong general reasoning

abilities, they struggle with deciphering complex visual encodings, performing numerical reasoning, and integrating domain knowledge to derive meaningful conclusions.

Future work will expand DomainCQA to include domains like biology, physics, and materials science, while exploring advanced multimodal training and dataset augmentation to enhance MLLMs' ability to process and reason with complex visual data.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1, 2, 6, 8
- [2] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, and et al. Baudouin De Monicault. Pixtral 12b, 2024. 2, 7, 8
- [3] Anthropic. The claude 3 model family: Opus, Sonnet, Haiku, 2024. 2
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023. 2, 6, 8
- [5] George Casella and Edward I. George. Explaining the gibbs sampler. *The American Statistician*, 46:167–174, 1992. 3
- [6] Ritwick Chaudhry, Sumit Shekhar, Utkarsh Gupta, Pranav Maneriker, Prann Bansal, and Ajay Joshi. Leaf-qa: Locate, encode & attend for figure question answering. In *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 3501–3510, 2020. 1, 2
- [7] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Januspro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025. 6, 8
- [8] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to GPT-4V? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 1
- [9] DeepSeek-AI. Deepseek-v3 technical report, 2024. 7
- [10] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Xinyue Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Jiaqi Wang. InternLM-XComposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024. 1, 2, 8
- [11] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Zhe Chen, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Kai Chen, Conghui He, Xingcheng Zhang, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. InternLM-XComposer2-4KHD: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *arXiv preprint arXiv:2404.06512*, 2024. 1, 6
- [12] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024. 2, 6, 8
- [13] Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. Chartllama: A multimodal llm for chart understanding and generation, 2023. 2
- [14] Wenyi Hong, Weihang Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*, 2024. 2, 7, 8
- [15] Ting-Yao Hsu, C Lee Giles, and Ting-Hao Huang. SciCap: Generating captions for scientific figures. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3258–3264. Association for Computational Linguistics, 2021. 1
- [16] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656, 2018. 1, 2
- [17] Samira Ebrahimi Kahou, Adam Atkinson, Vincent Michalski, Ákos Kádár, Adam Trischler, and Yoshua Bengio. FigureQA: An annotated figure dataset for visual reasoning, 2018. 1, 2
- [18] Shankar Kantharaj, Xuan Long Do, Rixie Tiffany Leong, Jia Qing Tan, Enamul Hoque, and Shafiq Joty. OpenCQA: Open-ended question answering with charts. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11817–11837, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. 2, 6
- [19] Klaus Krippendorff. Computing krippendorff's alpha-reliability. 2011. 4
- [20] Shengzhi Li and Nima Tajbakhsh. Scigraphqa: A large-scale synthetic multi-turn question-answering dataset for scientific graphs, 2023. 2
- [21] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages

- 74–81, Barcelona, Spain, 2004. Association for Computational Linguistics. 7
- [22] Dongyang Liu, Renrui Zhang, Longtian Qiu, Siyuan Huang, Weifeng Lin, Shitian Zhao, Shijie Geng, Ziyi Lin, Peng Jin, Kaipeng Zhang, Wenqi Shao, Chao Xu, Conghui He, Junjun He, Hao Shao, Pan Lu, Hongsheng Li, Yu Qiao, and Peng Gao. Sphinx-x: Scaling data and parameters for a family of multi-modal large language models, 2024. 1, 2, 7, 8
- [23] Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Eisenschlos. MatCha: Enhancing visual language pretraining with math reasoning and chart derendering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12756–12770, Toronto, Canada, 2023. Association for Computational Linguistics. 2
- [24] Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. MMC: Advancing multimodal chart understanding with large-scale instruction tuning. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1287–1310, Mexico City, Mexico, 2024. Association for Computational Linguistics. 2
- [25] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023. 1, 2, 6, 8
- [26] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 1
- [27] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. LLaVA-NeXT: Improved reasoning, ocr, and world knowledge, 2024. 2, 6, 7, 8
- [28] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: NLG evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, 2023. 7
- [29] Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, 2022. Association for Computational Linguistics. 2, 6
- [30] Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq Joty. UniChart: A universal vision-language pretrained model for chart comprehension and reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14662–14684, Singapore, 2023. Association for Computational Linguistics. 2
- [31] Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. ChartInstruct: Instruction tuning for chart comprehension and reasoning. In *Findings of the Association for Computational Linguistics: ACL 2024*, 2024. 2
- [32] Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. Chartassistant: A universal chart multimodal language model via chart-to-table pre-training and multitask instruction tuning. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7775–7803, Bangkok, Thailand, 2024. Association for Computational Linguistics. 2
- [33] Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. Plotqa: Reasoning over scientific plots. In *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020. 1, 2, 6
- [34] Aitor Ormazabal, Che Zheng, Cyprien de Masson d’Autume, Dani Yogatama, Deyu Fu, Donovan Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, Isaac Ong, et al. Reka Core, Flash, and Edge: A series of powerful multimodal language models. *arXiv preprint arXiv:2404.12387*, 2024. 1
- [35] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, 2002. Association for Computational Linguistics. 7
- [36] Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. Spiqa: A dataset for multimodal question answering on scientific papers. *NeurIPS*, 2024. 7
- [37] Hrituraj Singh and Sumit Shekhar. STL-CQA: Structure-based transformers with localization and encoding for chart question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3275–3284, 2020. 2
- [38] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023. 1, 2
- [39] Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *arXiv preprint arXiv:2406.18521*, 2024. 1, 2, 6
- [40] Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min Dou, Botian Shi, Junchi Yan, et al. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. *arXiv preprint arXiv:2402.12185*, 2024. 2
- [41] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024. 2, 6, 8
- [42] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multimodal large language models, 2024. 2
- [43] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chaoya Jiang, Chenliang Li, Yuanhong Xu,

Hehong Chen, Junfeng Tian, Qi Qian, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language models with multimodality, 2023.

- [44] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration, 2023. [2](#), [6](#), [8](#)
- [45] Liang Zhang, Anwen Hu, Haiyang Xu, Ming Yan, Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang. Tinchart: Efficient chart understanding with visual token merging and program-of-thoughts learning, 2024. [2](#), [6](#), [8](#)
- [46] Jinbo Zhao. Awesome-mllm-datasets, 2023. [2](#)