

# CHARTMASTER: BOOSTING MLLMs FOR CHART ANALYSIS THROUGH DATA, PERCEPTION, AND REASONING OPTIMIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Multimodal Large Language Models (MLLMs) have demonstrated significant potential in understanding visual information, yet they often fall short in the complex domain of chart analysis. Existing models often struggle to accurately capture detailed visual elements and to perform efficient multi-step reasoning. To address these challenges, we introduce ChartMaster, a holistic framework that systematically advances chart analysis by jointly optimizing data, perception, and reasoning. Our approach is built on three core innovations. First, we construct ChartVerse, a large-scale synthetic dataset with diverse chart types, rendering styles, and reasoning levels. Building on this foundation, we introduce a novel two-stage training paradigm: (i) Multi-Negative Direct Preference Optimization (MNDPO), which improves perceptual precision by training models to distinguish correct answers from carefully designed hard negative samples (i.e., plausible but incorrect alternatives); and (ii) Reinforcement Learning with Dynamic Length Reward (DLR), which adapts chain-of-thought reasoning to task complexity, encouraging concise solutions for simple queries and rigorous multi-step reasoning for complex ones. Extensive experiments across six benchmarks demonstrate that ChartMaster achieves state-of-the-art performance, surpassing prior chart-domain models and rivaling proprietary systems. These results highlight that coupling diverse data foundations with targeted perceptual and reasoning optimization provides an effective pathway toward robust chart understanding in MLLMs.

## 1 INTRODUCTION

Multimodal Large Language Models (MLLMs) have recently achieved remarkable success, demonstrating an impressive ability to interpret and reason about general visual information (Wang et al., 2024b; Yin et al., 2023). This progress has opened up new frontiers in artificial intelligence, moving us closer to models that can understand the world with human-like fluidity. However, despite their success on natural images, a critical and ubiquitous form of visual data remains a significant challenge: charts (Hoque et al., 2022). As dense, structured representations of quantitative information, charts are fundamental to communication in science, finance, and countless other domains. Automating the ability to understand and reason about them represents a crucial cognitive leap for AI, yet this specialized task exposes fundamental limitations in current MLLM architectures (Zhang et al., 2025; Yang et al., 2025c; Chen et al., 2025b; Wang et al., 2025a).

The core challenge lies in two distinct yet interconnected failures (Masry et al., 2025a). The first is a limitation in precise perception. Unlike general images, charts demand an extremely high degree of visual acuity. Models must not only recognize objects like bars and lines but also accurately ground them to precise values on an axis, parse fine-print text from labels and legends, and distinguish between visually similar elements (e.g., lines with subtle color differences). As illustrated in Figure 1 (left), even the most powerful, closed-source commercial MLLMs frequently fail at this fundamental stage, making clear perceptual errors on simple questions. These models often produce “hallucinated” or approximate values, which completely invalidates any subsequent reasoning.

The second challenge lies in the lack of efficient multi-step reasoning. In contrast to mathematical reasoning tasks (Lu et al., 2024a; Wang et al., 2024a), where virtually every query requires a lengthy

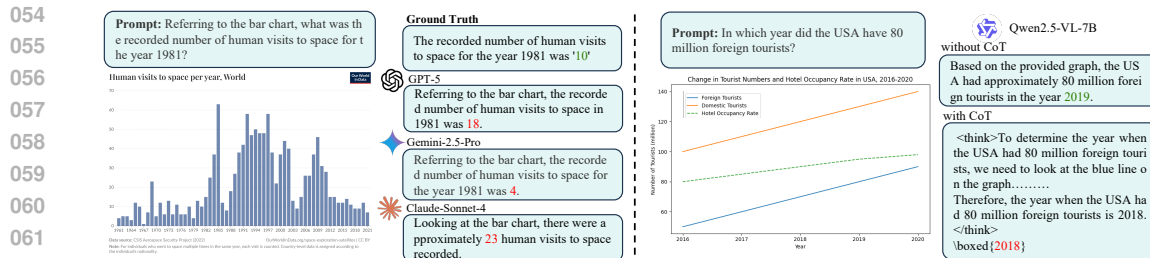


Figure 1: Perception and reasoning challenges of MLLMs on chart images. (Left) Even the most advanced closed-source models fail when geometric inference is required. (Right) Long chains of thought are not always beneficial and may cause errors on simple queries.

deductive process, chart analysis presents a far broader spectrum of difficulty. Certain questions can be solved through direct value extraction with minimal reasoning, whereas others demand complex, multi-step operations, including trend comparisons, rate estimations, or evidence synthesis across different chart components. Existing efforts to enhance reasoning, for example through Supervised Fine-Tuning (SFT) on Chain-of-Thought (CoT) data or Reinforcement Learning (RL), generally fail to account for this variability and instead adopt a uniform reasoning strategy (Chen et al., 2025a; Jia et al., 2025; Huang et al., 2025a; Meng et al., 2024; Masry et al., 2025b). As demonstrated in Figure 1 (right), when faced with simple queries, such strategies often induce unnecessary reasoning steps, which not only reduce efficiency but may also increase the likelihood of errors. This phenomenon is consistent with theoretical insights presented in Wu et al. (2025b), highlighting the need for approaches that adaptively balance reasoning depth with task complexity.

To address the unique challenges of chart analysis with MLLMs, we introduce **ChartMaster**, built upon a high-quality synthetic dataset, **ChartVerse**, and a two-stage training paradigm. ChartVerse is created via a structured data synthesis pipeline. Specifically, a large set of templates is pre-defined in a YAML file, with each field name explicitly specified. For each template, we prompt commercial LLMs to complete all predefined fields, producing full structured data records. These records are converted into Python code to render chart images, and the same plotting code is used to prompt LLMs to generate the corresponding question-answer (QA) pairs. This decoupled approach ensures accurate alignment between images and QA pairs while preserving access to precise ground-truth information. Dataset diversity is further enhanced by varying task topics, chart types, rendering libraries, LLM APIs, and question difficulty, collectively promoting robust model generalization.

This dataset directly supports our two-stage training process. The first stage targets perceptual precision using **Multi-Negative Direct Preference Optimization (MNDPO)**. Unlike conventional methods that rely solely on positive examples (Han et al., 2023; Zhang et al., 2024; Masry et al., 2025c), MNDPO trains the model to both identify correct answers and explicitly reject a set of carefully designed *hard negative examples*, which are incorrect answers closely resembling the correct response. This procedure sharpens decision boundaries and improves fine-grained visual discrimination. Building on this strong perceptual foundation, the second stage focuses on multi-step reasoning, leveraging reinforcement learning with a **Dynamic Length Reward (DLR)** mechanism. DLR provides adaptive rewards based on the length of generated reasoning, encouraging concise logic for simple queries while enabling detailed, rigorous reasoning for more complex analytical tasks.

Through extensive experiments, we demonstrate that ChartMaster achieves state-of-the-art performance on prominent chart question-answering benchmarks. Our results validate that a holistic approach, which couples a diverse data foundation with a dedicated two-stage training regimen for perception and reasoning, offers an effective and robust pathway toward developing MLLMs that can not only see charts but also truly understand them. Overall, the main contributions of this work can be summarized as follows:

- We present **ChartMaster**, a systematic framework that specializes MLLMs for chart analysis through a two-stage training pipeline addressing two core limitations: imprecise perception and inefficient reasoning.
- For perception, we propose Multi-Negative Direct Preference Optimization (**MNDPO**), which enhances fine-grained visual distinction by training the model not only to learn correct extractions but also to reject carefully crafted hard negatives.

- For reasoning, we introduce a Dynamic Length Reward (**DLR**) mechanism within reinforcement learning to optimize chain-of-thought generation. This enables adaptive reasoning depth, producing concise responses for simple queries and multi-step analyses for complex ones, thereby mitigating the common overthinking issue.
- We construct **ChartVerse**, a novel synthetic dataset built with a decoupled generation strategy to ensure high-quality image-QA alignment. Its diversity across task topics, chart types, rendering libraries, LLM APIs, and question difficulties provides a strong foundation for ChartMaster and supports its state-of-the-art performance on multiple chart benchmarks.

## 2 RELATED WORK

### 2.1 MULTIMODAL LARGE LANGUAGE MODELS

Multimodal Large Language Models (MLLMs) extend conventional LLMs by integrating multiple modalities, such as text and images, and typically consist of a modality encoder to extract features, a projector to align them with the LLM’s embedding space, and a large language model that performs reasoning and generates outputs (Yin et al., 2023). These models are generally trained in two stages: pre-training and post-training.

The pre-training stage aims to enhance the visual encoder and align its features with the LLM’s embedding space, typically through two paradigms. Contrastive learning (e.g., CLIP (Radford et al., 2021)) leverages positive and negative image-text pairs to build a robust cross-modal semantic space and clear decision boundaries, but may miss fine-grained visual details (Rusak et al., 2025; Zhong et al., 2024). Generative training predicts text conditioned on images, capturing detailed visual information, yet without explicit negatives, it can produce softer decision boundaries and struggle with ambiguous or borderline cases (Wang et al., 2025c; Liu et al., 2024a; Cao et al., 2023). This trade-off between robust alignment and fine-grained perception motivates hybrid strategies that combine the strengths of both paradigms (Bai et al., 2025).

Post-training further refines the model’s reasoning capabilities and aligns its behavior with human preferences (Kumar et al., 2025). This stage is typically realized through three approaches. Supervised Fine-Tuning (SFT) enables stable training and rapid acquisition of basic instruction-following abilities, but struggles to generalize to complex scenarios. Reinforcement Learning (RL) methods optimize the model’s outputs based on feedback signals to enhance reasoning and task-specific performance (Shao et al., 2024; Ouyang et al., 2022). Direct Preference Optimization (DPO) simplifies RL by directly learning from preference data using a contrastive loss, bypassing the need for a separate reward model Rafailov et al. (2023). Notably, some research also frames DPO as a form of offline reinforcement learning (Wang et al., 2025b). Recent multimodal models increasingly adopt hybrid post-training strategies (Lu et al., 2025; Huang et al., 2025b; Li et al., 2025; Ni et al., 2025). For instance, InternVL3.5 (Wang et al., 2025b) combines mixed preference optimization with reinforcement learning, while R1-OneVision (Yang et al., 2025b) applies supervised fine-tuning to teach reasoning, followed by rule-based reinforcement learning to enhance generalization. These approaches demonstrate the trend of integrating complementary methods to improve reasoning, alignment, and task-specific adaptation in MLLMs.

Our proposed MNDPO unifies generative training for precise fine-grained alignment with contrastive learning that leverages negative samples to enhance robustness and discriminative capacity. The multi-negative sample principle has been widely adopted in other domains, such as InFoNCE (Rusak et al., 2025) and listwise ranking (Cao et al., 2007).

### 2.2 CHART ANALYSIS WITH MULTIMODAL LARGE LANGUAGE MODELS

Chart analysis, requiring both precise visual perception and logical reasoning, is a natural application for MLLMs (Masry et al., 2024; Yang et al., 2025a; Xu et al., 2025; Zhao et al., 2025a; Methani et al., 2020; Kafle et al., 2018; Chaudhry et al., 2020). The application of MLLMs to this domain has evolved through distinct stages, largely mirroring broader trends in post-training methods. Early efforts, such as ChartLlama (Han et al., 2023), TinyChart (Zhang et al., 2024), ChartGemma (Masry et al., 2025c), and ChartCoder (Zhao et al., 2025b), primarily utilized SFT. These models were fine-tuned on chart-specific instruction datasets to adapt general-purpose MLLMs for basic chart

comprehension and to enhance their instruction-following abilities. However, while effective for direct data extraction, they often exhibited limited capabilities when faced with questions requiring complex, multi-step reasoning.

To address this reasoning deficit, subsequent models such as ChartReasoner (Jia et al., 2025) and Chart-R1 (Chen et al., 2025a) employed reinforcement learning to encourage chain-of-thought generation and enhance reasoning abilities. However, these approaches often apply uniform reasoning strategies regardless of query difficulty, resulting in inefficiencies: simple questions require minimal reasoning, while complex ones demand multi-step analysis. To address this, we introduce a Dynamic Length Reward (DLR) mechanism, which adjusts the depth and complexity of reasoning based on the perceived difficulty of each query, improving both efficiency and overall performance in chart question-answering.

### 3 CHARTMASTER

Our ChartMaster enhances chart analysis in MLLMs through a principled framework comprising a high-quality dataset and a progressive two-stage training process. First, in Section 3.1, we present **ChartVerse**, a synthetic dataset characterized by substantial diversity. Building upon this ChartVerse, we introduce our two-stage training strategy. In the first stage (Section 3.2), we propose Multi-Negative Direct Preference Optimization (**MNDPO**) to cultivate the model’s fine-grained perceptual capabilities. With a strong perceptual base established, the second stage (Section 3.3) applies reinforcement learning with a novel Dynamic Length Reward (**DLR**) mechanism to encourage efficient multi-step reasoning.

#### 3.1 CHARTVERSE SYNTHESIS PIPELINE

Recent studies have explored the use of commercial LLMs to synthesize chart question-answer datasets for adapting MLLMs to chart understanding tasks (Chen et al., 2025a; Jia et al., 2025; Yang et al., 2025c; Zhang et al., 2025). While this automated strategy substantially reduces human annotation costs, it often suffers from limited diversity, which in turn constrains model generalization. To mitigate this limitation, we propose a novel LLM-based data synthesis framework for constructing **ChartVerse**. We adopt a two-stage design, separating *chart image generation* from *QA pair generation*, to ensure tight alignment and promote diversity.

**Chart Image Generation.** Unlike Chart-R1 (Chen et al., 2025a), which directly instructs LLMs to generate Matplotlib plotting code, we adopt a template-driven approach. Specifically, we design a large collection of YAML templates, each containing explicitly defined fields that specify the attributes required for chart construction. LLMs are prompted to complete all fields, producing structured data records that fully describe the chart content. This representation allows us to flexibly manipulate specific attributes, such as task type, style, rendering library, or font, thereby achieving broad diversity. The structured records are then deterministically converted into Python code to render chart images. By focusing the LLM on structured content rather than entire plotting scripts, this approach also improves token efficiency and reduces generation errors.

**QA Pair Generation.** Given the plotting code corresponding to each chart, we perform a second LLM call to generate QA pairs grounded in the structured data. This decoupled design mitigates the risk of hallucinations in which QA pairs fail to match rendered chart elements. It further enables richer question types that incorporate visual properties such as colors and styles. To ensure a balanced spectrum of difficulty, we organize questions into five levels:

- **Level 1:** Direct factual extraction (e.g., retrieving specific values), which may still require fine-grained perception when values are not explicitly labeled.
- **Level 2:** Local comparisons across elements, such as identifying maxima or minima.
- **Level 3:** Global reasoning, including recognition of overall patterns or trends.
- **Level 4:** Basic arithmetic operations, such as computing sums across categories.
- **Level 5:** Multi-step reasoning involving complex aggregation and cross-element comparisons.

**Quality Evaluation.** To validate the quality of ChartVerse, we randomly sampled 1,000 instances and recruited human experts for evaluation. The results indicate that over 95% of the instances are error-free. This stands in stark contrast to the reported 85% accuracy of the synthetic data from

Chart-R1. We attribute the high fidelity of ChartVerse to our decoupled generation pipeline, wherein a vast set of pre-defined templates simplifies the LLM’s task to pure structured data generation, minimizing opportunities for error (Wu et al., 2025a).

By utilizing the APIs of several LLMs, we construct ChartVerse of approximately 128k instances in total. Details are provided in Appendix A. With this high-quality dataset established as our foundation, we proceed to the two-stage training process.

### 3.2 STAGE 1: MULTI-NEGATIVE DIRECT PREFERENCE OPTIMIZATION

Stage 1 focuses on enhancing MLLMs’ fine-grained perception of charts, which serves as the foundation for their competence in chart-related tasks. Previous approaches for adapting MLLMs to chart-understanding tasks have predominantly relied on Supervised Fine-Tuning (SFT) (Hu et al., 2024; He et al., 2024; Liu et al., 2024b; Han et al., 2023). Nevertheless, these models exhibit significant limitations when faced with charts that are dense with visual elements or require inferring data values from geometric properties (e.g., bar height or pie area) (Masry et al., 2025a). To surpass the coarse corrections offered by SFT on individual correct answers, we adopt a preference optimization perspective to cultivate a more nuanced model perception and propose **Multi-Negative Direct Preference Optimization (MNDPO)**, which leverages both the correct answer and a spectrum of incorrect answers for each question. The core idea is that by learning not only from correct responses but also from a diverse set of incorrect ones, the model can develop superior discriminative abilities for handling such ambiguous and challenging cases (Zhu et al., 2025b).

Our approach extends Direct Preference Optimization (DPO) (Rafailov et al., 2023) into a multi-negative framework that is crucial for chart understanding tasks, allowing models to robustly distinguish the single correct answer from a vast set of plausible yet incorrect distractors. The key insight of DPO is that one can directly optimize a policy model to satisfy human preferences without explicitly learning a reward function. Given a preference dataset  $\mathcal{D} = \{(x, y_w, y_l)\}$ , where  $x$  is the input prompt,  $y_w$  is the preferred (winning) response, and  $y_l$  is the dispreferred (losing) response, DPO models the preference probability using the Bradley-Terry model (Bradley & Terry, 1952):

$$P(y_w \succ y_l | x) = \sigma(r(x, y_w) - r(x, y_l)), \quad (1)$$

where  $\sigma(\cdot)$  is the sigmoid function and  $r(x, y)$  is a latent reward function. DPO reparameterizes the reward function as:

$$r(x, y) = \beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} + \beta \log Z(x), \quad (2)$$

where  $\pi(y|x)$  denotes the policy model to be optimized,  $\pi_{\text{ref}}(y|x)$  is the reference model,  $\beta$  is the temperature parameter, and  $Z(x)$  is the normalization constant. The final loss function is as:

$$\mathcal{L}_{\text{DPO}}(\pi; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right]. \quad (3)$$

However, standard DPO is limited to single pairs of winning and losing responses. To enhance the model’s discriminative power, we extend this framework to handle one winning response against a set of  $N$  losing responses,  $\{y_l^1, \dots, y_l^N\}$ . We model this multi-negative competition as the probability that the winning response is preferred over all losing responses:

$$P(y_w \succ \{y_l^1, \dots, y_l^N\} | x) = \frac{\exp(r(x, y_w))}{\exp(r(x, y_w)) + \sum_{i=1}^N \exp(r(x, y_l^i))}. \quad (4)$$

This formulation is a generalization of the Bradley-Terry model to a multi-alternative choice scenario. Substituting the DPO reward reparameterization, we first obtain the following expression:

$$\begin{aligned} \mathcal{L}_{\text{MNDPO}} = & - \left( \beta \log \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \beta \log Z(x) \right) \\ & + \log \left( \exp \left( \beta \log \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \beta \log Z(x) \right) \right. \\ & \left. + \sum_{i=1}^N \exp \left( \beta \log \frac{\pi(y_l^i|x)}{\pi_{\text{ref}}(y_l^i|x)} + \beta \log Z(x) \right) \right). \end{aligned} \quad (5)$$

By eliminating the common normalization term  $Z(x)$ , this simplifies to an intermediate loss:

$$\mathcal{L}_{\text{MNDPO}} = -\beta \log \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \log \left( \exp \left( \beta \log \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right) + \sum_{i=1}^N \exp \left( \beta \log \frac{\pi(y_i^i|x)}{\pi_{\text{ref}}(y_i^i|x)} \right) \right), \quad (6)$$

which can then be simplified to the final MNDPO loss:

$$\mathcal{L}_{\text{MNDPO}}(\pi; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, \{y_i^i\}) \sim \mathcal{D}} \left[ \log \frac{\left( \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta}{\left( \frac{\pi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta + \sum_{i=1}^N \left( \frac{\pi(y_i^i|x)}{\pi_{\text{ref}}(y_i^i|x)} \right)^\beta} \right]. \quad (7)$$

Intuitively, this objective function takes the form of a softmax-based cross-entropy loss. It optimizes the model  $\pi$  to maximize the selection probability of the winning response  $y_w$  among the set containing itself and all  $N$  negative examples. The full derivation of this loss function is provided in Appendix B.

**Hard Negative Sample Construction.** A critical component of our MNDPO framework is the construction of high-quality hard negative samples, which are designed to be plausible yet incorrect in order to strengthen fine-grained discriminative abilities. We first employ a model-based strategy by sampling suboptimal predictions with relatively high likelihood but incorrect semantics, thereby generating negatives that closely align with the model’s decision boundary. Complementing this, we adopt data-driven heuristics tailored to the nature of the correct answer ( $y_w$ ). For answers not directly present in the source data, we select adjacent values from the underlying chart as distractors; for numerical results, we apply small perturbations within a  $\pm 10\%$  range; and for textual results, we introduce minor modifications that yield semantically related but factually incorrect statements. This hybrid strategy provides a computationally lightweight yet effective way to construct challenging negatives. Theoretically, MNDPO is largely insensitive to individual negative samples. Since the loss (Eq. 7) is a softmax-style competition, optimization is mainly driven by increasing the probability of the positive sample, making its fidelity the key factor for effective training.

Through MNDPO, our model acquires the fine-grained perceptual acuity necessary to accurately interpret chart data. However, correct perception is only the first step and the model must also learn to reason efficiently for complex problems, which we address in our second stage.

### 3.3 STAGE 2: REINFORCEMENT LEARNING WITH DYNAMIC LENGTH REWARD

Building on the enhanced perceptual foundation from Stage 1, Stage 2 aims to strengthen the model’s reasoning capabilities, allowing it to handle increasingly complex chart tasks. Reinforcement learning has recently emerged as a powerful paradigm for enhancing the reasoning abilities of MLLMs, typically by assigning higher rewards to responses that align with human-preferred outcomes. In chart analysis, models such as Chart-R1 (Chen et al., 2025a) and ChartReasoner (Jia et al., 2025) have demonstrated the effectiveness of this approach. However, a key limitation lies in their neglect of the highly variable difficulty of chart-related queries. Unlike mathematical reasoning tasks such as MathVista (Lu et al., 2024a), which consistently demand multi-step deduction, chart tasks range from simple value lookup to complex multi-step inference, making uniform and lengthy reasoning both inefficient and potentially misleading.

To address this limitation, we propose a reinforcement learning framework with a **Dynamic Length Reward (DLR)**, which incentivizes concise responses for simple queries while promoting more elaborate reasoning for complex ones. Specifically, for each question  $q$  in a batch, we sample a set of  $G$  responses  $\{o_1, o_2, \dots, o_G\}$  from the old policy  $\pi_{\text{old}}$ . The reward signal for each response  $o_i$  is a composite of three components:  $r_i = r_{\text{acc}} + r_{\text{format}} + r_{\text{len}}$ . The first two components are standard checks:  $r_{\text{acc}}$  is an accuracy reward, yielding 1 for a correct final answer and 0 otherwise.  $r_{\text{format}}$  verifies that CoT tokens are enclosed in `<think>` tags and the final answer is in a `\boxed{\}` environment. The key innovation lies in our length reward,  $r_{\text{len}}$ , defined as:

$$r_{\text{len}}(o_i) = (r_{\text{len}}^{\text{max}} - r_{\text{len}}^{\text{min}}) \cdot \exp \left( -\frac{||o_i| - L_{\text{tgt}}|}{\tau} \right) + r_{\text{len}}^{\text{min}} \quad (8)$$

where  $r_{\text{len}}^{\text{max}}$  and  $r_{\text{len}}^{\text{min}}$  are the maximum and minimum length rewards,  $|o_i|$  is the length of the response, and  $\tau$  is a temperature parameter controlling the sharpness of the reward curve. The target

length,  $L_{\text{tgt}}$ , is dynamically determined for each group based on the correctness of the sampled responses:

$$L_{\text{tgt}}(x) = \begin{cases} \min\{|o_i| \mid r_{\text{acc}}(x, o_i) = 1\} & \text{if at least one response is correct,} \\ \frac{1}{G} \sum_{i=1}^G |o_i| & \text{otherwise.} \end{cases} \quad (9)$$

This formulation encourages the model to find the shortest correct reasoning path when successful, and to conform to the group’s average length when unsuccessful, preventing divergence. Finally, we optimize the policy  $\pi$  by maximizing the following objective, which is based on the Group Relative Policy Optimization (GRPO) framework (Shao et al., 2024):

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\{o_i\} \sim \pi_{\theta_{\text{old}}}} \left[ \frac{1}{G} \sum_{i=1}^G \min \left( \frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{\text{old}}}(o_i|q)} A_i, \text{clip} \left( \frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{\text{old}}}(o_i|q)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) \right], \quad (10)$$

$$A_i = \frac{r_i - \max(\{r_1, \dots, r_G\})}{\text{std}(\{r_1, \dots, r_G\})}.$$

where  $\varepsilon$  controls policy deviation, with the KL penalty term omitted.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETTINGS

**Implementation Details.** ChartMaster is built upon the *Qwen2.5-VL-7B-Instruct* model (Bai et al., 2025). During training, MNDPO is applied using only Level 1 and Level 2 questions (as its primary goal is to enhance perception), whereas the reinforcement learning stage leverages the full dataset. Each stage is trained for two epochs with a batch size of 256. We regard the dynamic length reward as a penalty, and therefore set  $r_{\text{min}}^{\text{len}} = -1$  and  $r_{\text{max}}^{\text{len}} = 0$ . The learning rates for the two stages are set to 1e-5 and 1e-6, respectively, with optimization performed using the Adam optimizer. By default, we set  $\beta = 0.1$ ,  $G = 8$ ,  $\varepsilon = 0.1$ ,  $\tau = 100$ , and  $N = 4$ . The specific prompts and template details are provided in Appendix F. Training was conducted on 8 NVIDIA H20 GPUs.

**Evaluation Benchmarks.** We conduct extensive evaluations across six chart datasets, including ChartQA (Masry et al., 2022), CharXiv (Wang et al., 2024c), ChartMuseum (Tang et al., 2025), ChartX (Xia et al., 2024), ChartQPro (Masry et al., 2025a), and ChartBench (Xu et al., 2023). For all datasets except ChartBench, we report sequence-level accuracy, while for ChartBench we adopt its official metric, *Improved Acc+* (Xu et al., 2023).

**Comparison Methods.** We compare our proposed ChartMaster against three categories of models. The first category consists of closed-source proprietary models, including GPT-4o (Hurst et al., 2024), the Gemini family (Reid et al., 2024), and Claude-3.5-Sonnet (Anthropic, 2024). The second category covers general-domain open-source MLLMs, such as LLaVA-1.5 (Liu et al., 2024c), Idefics3 (Laurençon et al., 2024), Qwen2.5-VL (Bai et al., 2025), InternVL3 (Zhu et al., 2025a), MiniCPM-V-4.5 (Yao et al., 2024), and Ovis-2 (Lu et al., 2024b). The third category focuses on chart-domain MLLMs, including ChartLlama (Han et al., 2023), TinyChart (Zhang et al., 2024), ChartGemma (Masry et al., 2025c), ChartReasoner (Jia et al., 2025), and Chart-R1 (Chen et al., 2025a). Among them, Chart-R1 is our strongest baseline, as it is built on the same base model as ChartMaster, enhanced through RL-based post-training, and publicly accessible to promote reproducibility and fair comparison.

### 4.2 MAIN RESULTS

The comprehensive quantitative evaluation results are presented in Table 1. Overall, ChartMaster consistently advances the state of chart understanding across all benchmarks. Building upon its base model Qwen2.5-VL-7B, it achieves substantial gains, particularly on challenging benchmarks such as ChartQPro that demand both fine-grained perception and multi-step reasoning, highlighting the effectiveness of our training strategy. Beyond this, ChartMaster substantially outperforms general-domain open-source MLLMs of similar scale and reaches performance levels comparable to, and in some cases exceeding, proprietary closed-source systems. The most striking improvements emerge when comparing against recent chart-domain models, which represent the most relevant baselines.

Table 1: Main results on chart understanding benchmarks. Comparison between our proposed ChartMaster, closed-source proprietary MLLMs, open-source general-domain MLLMs, and chart-domain specialized MLLMs. Results are reported as sequence-level accuracy, except for ChartBench which uses *Improved Acc+*.

Model	ChartQA	CharXiv	ChartQAPro	ChartX	ChartBench	ChartMuseum
<b>Closed-source</b>						
GPT-4o	85.70	47.10	37.67	35.60	59.45	42.20
Gemini-1.5-Flash	79.00	33.90	42.96	-	-	31.10
Gemini-1.5-Pro	87.20	43.30	-	-	-	41.30
Claude-3.5-Sonnet	90.80	60.20	43.58	-	-	54.40
<b>General-domain Open-source</b>						
LLaVA-1.5-7B	55.32	12.20	19.96	11.87	23.39	12.00
Idefics3-8B	79.36	29.00	20.03	30.99	30.50	21.90
InternVL3-8B	86.60	37.60	37.20	59.38	47.52	28.20
MiniCPM-V-4.5-8B	87.40	42.50	38.64	59.38	54.95	26.10
Ovis-2-8B	86.80	45.20	44.67	55.99	51.60	29.90
Qwen2.5-VL-7B	87.32	42.50	36.61	65.10	54.06	26.80
<b>Chart-domain</b>						
ChartLlama	69.66	14.20	-	13.80	21.30	-
TinyChart	83.60	8.30	13.25	-	-	-
ChartGemma	80.16	12.50	6.84	-	-	-
ChartReasoner	86.93	-	39.97	-	55.20	-
Chart-R1	91.04	46.20	44.04	66.49	53.57	30.40
<b>ChartMaster</b>	<b>91.82</b>	<b>51.20</b>	<b>56.46</b>	<b>68.23</b>	<b>66.41</b>	<b>38.60</b>

Unlike models such as ChartReasoner, whose strengths are limited to specific datasets but weaken on others, ChartMaster delivers robust and balanced gains, achieving state-of-the-art performance across the board and surpassing prior work by large margins on several benchmarks.

We observe that the performance improvement of ChartMaster on the ChartQA dataset is less pronounced compared to its gains on other benchmarks, particularly when compared with Chart-R1. To investigate this, we manually analyze a subset of the incorrectly predicted samples. Our investigation indicates that this modest gain is partly due to inherent noise in the ChartQA test set. Specifically, we observe instances of erroneous ground-truth labels and questions with ambiguous phrasing, which are difficult to evaluate reliably using automated, rule-based metrics. Such data quality limitations may constrain the observable performance on this benchmark. A detailed qualitative analysis, including illustrative examples of these problematic cases, is provided in Appendix C.

Given the space constraints of the main text, we present the analysis of training dynamics in Appendix D and additional qualitative studies in Appendix E, with a particular focus on ChartMaster’s generalization to in-the-wild scenarios.

### 4.3 ABLATION STUDY

**Impact of Training Strategies.** First, we assess the impact of different training strategies on model performance, as shown in Figure 2 (a). Our experiments include configurations with “only MNDPO”, “only RL”, “full SFT”, “SFT then RL”, and “RL w/o DLR”. The results indicate that applying either MNDPO or RL in isolation yields notable performance gains over the baseline model. Conversely, employing a “full SFT” approach leads to a significant degradation in performance across most benchmarks. This outcome suggests that extensive SFT may negatively affect the model’s generalization capabilities, a finding that aligns with recent analyses on the distinct characteristics of SFT and RL methodologies (Chu et al., 2025). Furthermore, a key observation is that our full method, which incorporates the Dynamic Length Reward, achieves the best overall performance. It is noteworthy that the inclusion of DLR results in no performance loss compared to the “RL w/o DLR” variant and, in some cases, provides a slight improvement. This suggests that incentivizing more efficient reasoning can help mitigate extraneous steps that might lead to errors. To further investigate this, we analyze the effect of DLR on the model’s output verbosity in Figure 2 (b). The results indicate that DLR significantly decreases the average number of generated

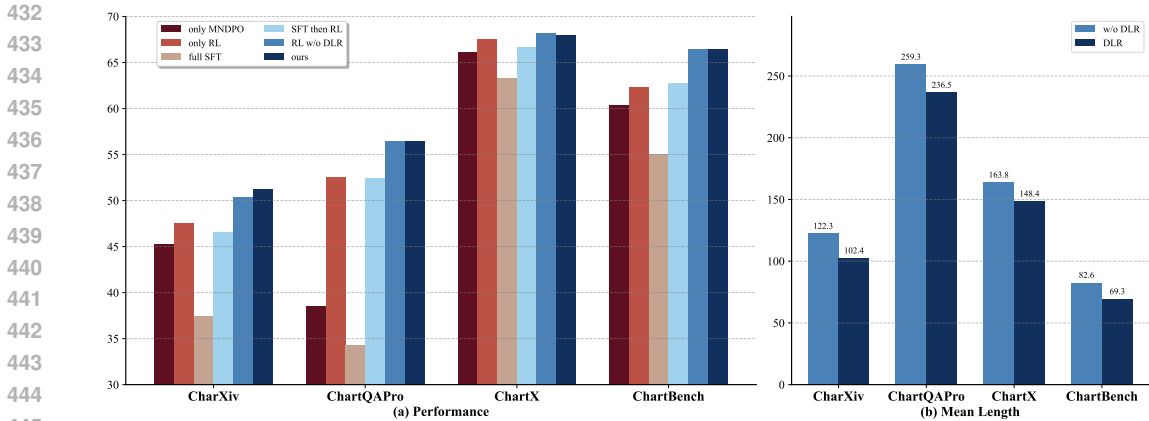


Figure 2: Ablation studies on different training strategies. (a) Impact of various training strategies on the final model performance. (b) Effect of the dynamic length reward mechanism on the average output length.

Table 2: Ablation studies on  $N$ .

$N$	1	2	3	4	5
Acc.	64.32	65.47	65.93	66.41	66.45

Table 3: Ablation studies on  $\tau$ .

$\tau$	50	70	100	120	150
Acc.	66.18	66.27	66.41	66.39	66.28

tokens, achieving a reduction of approximately 15% across the four datasets. This contraction in output length directly mitigates common efficiency issues associated with reinforcement learning, demonstrating that the DLR mechanism preserves or improves model accuracy while enhancing computational and inferential efficiency.

**Impact of Parameter  $N$  and  $\tau$ .** We first explore the impact of the number of negative samples,  $N$ , during the MNDPO stage, with the evaluation results on ChartBench presented in Table 2. The findings show that as  $N$  initially increases, model performance improves significantly. This demonstrates the tangible benefit of introducing multiple negative samples, which enhances the model’s fine-grained perceptual abilities by compelling it to learn more precise distinctions between correct and incorrect objects. However, we observe that performance gains plateau when  $N$  exceeds 4. A plausible explanation for this saturation is that the pool of potential “hard negative examples” has been largely exhausted. Beyond this threshold, additional samples may be “easier” negatives that the model can already discriminate with high confidence, thus offering marginal additional learning value. From the results in Table 3, we observe that the performance of the dynamic length reward is relatively insensitive to variations in the parameter  $\tau$  within a reasonable range. In practice, we set  $\tau = 100$ , which yields the best results. Nonetheless,  $\tau$  can be flexibly adjusted according to the type and difficulty of the task.

## 5 CONCLUSION

In this work, we tackle the key challenges faced by Multimodal Large Language Models (MLLMs) in chart analysis, namely imprecise visual perception and limited multi-step reasoning. We propose ChartMaster, a framework that systematically enhances model performance through targeted optimizations in data, perception, and reasoning. ChartMaster is built upon a high-quality synthetic dataset, ChartVerse, and a novel two-stage training procedure that leverages Multi-Negative Direct Preference Optimization (MNDPO) to improve perceptual acuity and Reinforcement Learning with a Dynamic Length Reward (DLR) to enable adaptive, efficient reasoning. Extensive evaluations across six prominent chart benchmarks demonstrate that ChartMaster achieves state-of-the-art performance, substantially surpassing prior specialized models as well as leading proprietary systems. These results underscore that a strong data foundation combined with tailored perceptual and reasoning optimizations provides an effective strategy for developing expert MLLMs, offering a practical blueprint for their application to charts and other complex visual tasks.

## ETHICS STATEMENT

The research presented in this paper adheres to the ICLR Code of Ethics. The datasets used were either synthetically generated or are publicly available benchmarks. All models and methods described were developed and utilized strictly for academic research purposes to advance the field of machine learning. We foresee no direct negative societal impacts from our work.

## REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our work, we provide comprehensive details throughout the paper and its appendices. For our synthetic dataset, the synthesis pipeline is described in Section 3.1, with a detailed analysis in Appendix A and the prompts used for generation available in Appendix F. For our MNDPO method, the core optimization objective is introduced in Section 3.2, and its full theoretical derivation is provided in Appendix B. Our experimental setup is detailed in Section 4.1, and all datasets used for evaluation are publicly available and open-source.

## REFERENCES

- Anthropic. Introducing claude 3.5 sonnet, 2024. URL <https://www.anthropic.com/news/claude-3-5-sonnet>.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Ming-Hsuan Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *CoRR*, abs/2502.13923, 2025. doi: 10.48550/ARXIV.2502.13923. URL <https://doi.org/10.48550/arXiv.2502.13923>.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Haoyu Cao, Changcun Bao, Chaohu Liu, Huang Chen, Kun Yin, Hao Liu, Yinsong Liu, Deqiang Jiang, and Xing Sun. Attention where it matters: Rethinking visual document understanding with selective region concentration. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pp. 19460–19470. IEEE, 2023. doi: 10.1109/ICCV51070.2023.01788. URL <https://doi.org/10.1109/ICCV51070.2023.01788>.
- Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise approach to listwise approach. In Zoubin Ghahramani (ed.), *Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007*, volume 227 of *ACM International Conference Proceeding Series*, pp. 129–136. ACM, 2007. doi: 10.1145/1273496.1273513. URL <https://doi.org/10.1145/1273496.1273513>.
- Ritwick Chaudhry, Sumit Shekhar, Utkarsh Gupta, Pranav Maneriker, Prann Bansal, and Ajay Joshi. LEAF-QA: locate, encode & attend for figure question answering. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020*, pp. 3501–3510. IEEE, 2020. doi: 10.1109/WACV45572.2020.9093269. URL <https://doi.org/10.1109/WACV45572.2020.9093269>.
- Lei Chen, Xuanle Zhao, Zhixiong Zeng, Jing Huang, Yufeng Zhong, and Lin Ma. Chart-r1: Chain-of-thought supervision and reinforcement for advanced chart reasoner. *CoRR*, abs/2507.15509, 2025a. doi: 10.48550/ARXIV.2507.15509. URL <https://doi.org/10.48550/arXiv.2507.15509>.
- Zixin Chen, Sicheng Song, Kashun Shum, Yanna Lin, Rui Sheng, and Huamin Qu. Unmasking deceptive visuals: Benchmarking multimodal large language models on misleading chart question answering. *CoRR*, abs/2503.18172, 2025b. doi: 10.48550/ARXIV.2503.18172. URL <https://doi.org/10.48550/arXiv.2503.18172>.

- 540 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V.  
541 Le, Sergey Levine, and Yi Ma. SFT memorizes, RL generalizes: A comparative study of founda-  
542 tion model post-training. *CoRR*, abs/2501.17161, 2025. doi: 10.48550/ARXIV.2501.17161.  
543 URL <https://doi.org/10.48550/arXiv.2501.17161>.
- 544 Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang  
545 Zhang. Chartllama: A multimodal LLM for chart understanding and generation. *CoRR*,  
546 abs/2311.16483, 2023. doi: 10.48550/ARXIV.2311.16483. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2311.16483)  
547 [48550/arXiv.2311.16483](https://doi.org/10.48550/arXiv.2311.16483).
- 549 Wei He, Zhiheng Xi, Wanxu Zhao, Xiaoran Fan, Yiwen Ding, Zifei Shan, Tao Gui, Qi Zhang,  
550 and Xuanjing Huang. Distill visual chart reasoning ability from llms to mllms. *CoRR*,  
551 abs/2410.18798, 2024. doi: 10.48550/ARXIV.2410.18798. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2410.18798)  
552 [48550/arXiv.2410.18798](https://doi.org/10.48550/arXiv.2410.18798).
- 553 Enamul Hoque, Parsa Kavehzadeh, and Ahmed Masry. Chart question answering: State of the art  
554 and future directions. *Comput. Graph. Forum*, 41(3):555–572, 2022. doi: 10.1111/CGF.14573.  
555 URL <https://doi.org/10.1111/cgf.14573>.
- 556 Linmei Hu, Duokang Wang, Yiming Pan, Jifan Yu, Yingxia Shao, Chong Feng, and Liqiang Nie.  
557 Novachart: A large-scale dataset towards chart understanding and generation of multimodal large  
558 language models. In Jianfei Cai, Mohan S. Kankanhalli, Balakrishnan Prabhakaran, Susanne Boll,  
559 Ramanathan Subramanian, Liang Zheng, Vivek K. Singh, Pablo César, Lexing Xie, and Dong Xu  
560 (eds.), *Proceedings of the 32nd ACM International Conference on Multimedia, MM 2024, Mel-*  
561 *bourne, VIC, Australia, 28 October 2024 - 1 November 2024*, pp. 3917–3925. ACM, 2024. doi:  
562 10.1145/3664647.3680790. URL <https://doi.org/10.1145/3664647.3680790>.
- 563 Kung-Hsiang Huang, Can Qin, Haoyi Qiu, Philippe Laban, Shafiq Joty, Caiming Xiong, and Chien-  
564 Sheng Wu. Why vision language models struggle with visual arithmetic? towards enhanced chart  
565 and geometry understanding. *arXiv preprint arXiv:2502.11492*, 2025a.
- 566 Wenxuan Huang, Bohan Jia, Zijie Zhai, Shaosheng Cao, Zheyu Ye, Fei Zhao, Zhe Xu, Yao Hu,  
567 and Shaohui Lin. Vision-r1: Incentivizing reasoning capability in multimodal large language  
568 models. *CoRR*, abs/2503.06749, 2025b. doi: 10.48550/ARXIV.2503.06749. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2503.06749)  
569 [48550/arXiv.2503.06749](https://doi.org/10.48550/arXiv.2503.06749).
- 570 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-  
571 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint*  
572 *arXiv:2410.21276*, 2024.
- 573 Caijun Jia, Nan Xu, Jingxuan Wei, Qingli Wang, Lei Wang, Bihui Yu, and Junnan Zhu. Chartrea-  
574 soner: Code-driven modality bridging for long-chain reasoning in chart question answering.  
575 *CoRR*, abs/2506.10116, 2025. doi: 10.48550/ARXIV.2506.10116. URL [https://doi.org/](https://doi.org/10.48550/arXiv.2506.10116)  
576 [10.48550/arXiv.2506.10116](https://doi.org/10.48550/arXiv.2506.10116).
- 577 Kushal Kafle, Brian L. Price, Scott Cohen, and Christopher Kanan. DVQA: understanding data  
578 visualizations via question answering. In *2018 IEEE Conference on Computer Vision and*  
579 *Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 5648–  
580 5656. Computer Vision Foundation / IEEE Computer Society, 2018. doi: 10.1109/CVPR.  
581 2018.00592. URL [http://openaccess.thecvf.com/content\\_cvpr\\_2018/html/](http://openaccess.thecvf.com/content_cvpr_2018/html/Kafle_DVQA_Understanding_Data_CVPR_2018_paper.html)  
582 [Kafle\\_DVQA\\_Understanding\\_Data\\_CVPR\\_2018\\_paper.html](http://openaccess.thecvf.com/content_cvpr_2018/html/Kafle_DVQA_Understanding_Data_CVPR_2018_paper.html).
- 583 Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal,  
584 Mubarak Shah, Ming-Hsuan Yang, Phillip H. S. Torr, Salman H. Khan, and Fahad Shabbaz Khan.  
585 LLM post-training: A deep dive into reasoning large language models. *CoRR*, abs/2502.21321,  
586 2025. doi: 10.48550/ARXIV.2502.21321. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2502.21321)  
587 [2502.21321](https://doi.org/10.48550/arXiv.2502.21321).
- 588 Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better under-  
589 standing vision-language models: insights and future directions. *CoRR*, abs/2408.12637, 2024.  
590 doi: 10.48550/ARXIV.2408.12637. URL [https://doi.org/10.48550/arXiv.2408.](https://doi.org/10.48550/arXiv.2408.12637)  
591 [12637](https://doi.org/10.48550/arXiv.2408.12637).

- 594 Yuting Li, Lai Wei, Kaipeng Zheng, Jingyuan Huang, Linghe Kong, Lichao Sun, and Weiran  
595 Huang. Vision matters: Simple visual perturbations can boost multimodal math reasoning. *CoRR*,  
596 abs/2506.09736, 2025. doi: 10.48550/ARXIV.2506.09736. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2506.09736)  
597 [48550/arXiv.2506.09736](https://doi.org/10.48550/arXiv.2506.09736).
- 598  
599 Chaohu Liu, Kun Yin, Haoyu Cao, Xinghua Jiang, Xin Li, Yinsong Liu, Deqiang Jiang, Xing Sun,  
600 and Linli Xu. HRVDA: high-resolution visual document assistant. In *IEEE/CVF Conference*  
601 *on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*,  
602 pp. 15534–15545. IEEE, 2024a. doi: 10.1109/CVPR52733.2024.01471. URL [https://doi.](https://doi.org/10.1109/CVPR52733.2024.01471)  
603 [org/10.1109/CVPR52733.2024.01471](https://doi.org/10.1109/CVPR52733.2024.01471).
- 604 Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Ya-  
605 coob, and Dong Yu. MMC: advancing multimodal chart understanding with large-scale instruc-  
606 tion tuning. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), *Proceedings of the*  
607 *2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pp. 1287–1310. Association for Computational Linguistics, 2024b. doi:  
608 10.18653/V1/2024.NAACL-LONG.70. URL [https://doi.org/10.18653/v1/2024.](https://doi.org/10.18653/v1/2024.naacl-long.70)  
609 [naacl-long.70](https://doi.org/10.18653/v1/2024.naacl-long.70).
- 610  
611  
612 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction  
613 tuning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024,*  
614 *Seattle, WA, USA, June 16-22, 2024*, pp. 26286–26296. IEEE, 2024c. doi: 10.1109/CVPR52733.  
615 2024.02484. URL <https://doi.org/10.1109/CVPR52733.2024.02484>.
- 616  
617 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-  
618 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning  
619 of foundation models in visual contexts. In *The Twelfth International Conference on Learning*  
620 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024a. URL  
621 <https://openreview.net/forum?id=KUNzEQMWU7>.
- 622  
623 Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. Ovis:  
624 Structural embedding alignment for multimodal large language model. *arXiv:2405.20797*, 2024b.
- 625  
626 Shiyin Lu, Yang Li, Yu Xia, Yuwei Hu, Shanshan Zhao, Yanqing Ma, Zhichao Wei, Yinglun Li,  
627 Lunhao Duan, Jianshan Zhao, et al. Ovis2. 5 technical report. *arXiv preprint arXiv:2508.11737*,  
628 2025.
- 629  
630 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R. Joty, and Enamul Hoque. Chartqa: A  
631 benchmark for question answering about charts with visual and logical reasoning. In Smaranda  
632 Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association for Com-*  
633 *putational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 2263–2279. Associa-  
634 tion for Computational Linguistics, 2022. doi: 10.18653/V1/2022.FINDINGS-ACL.177. URL  
<https://doi.org/10.18653/v1/2022.findings-acl.177>.
- 635  
636 Ahmed Masry, Mehrad Shahmohammadi, Md. Rizwan Parvez, Enamul Hoque, and Shafiq Joty.  
637 Chartinstruct: Instruction tuning for chart comprehension and reasoning. In Lun-Wei Ku, Andre  
638 Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics,*  
639 *ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 10387–10409. As-  
640 sociation for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL.619.  
URL <https://doi.org/10.18653/v1/2024.findings-acl.619>.
- 641  
642 Ahmed Masry, Mohammed Saidul Islam, Mahir Ahmed, Aayush Bajaj, Firoz Kabir, Aaryaman  
643 Kartha, Md. Tahmid Rahman Laskar, Mizanur Rahman, Shadikur Rahman, Mehrad Shahmo-  
644 hammadi, Megh Thakkar, Md. Rizwan Parvez, Enamul Hoque, and Shafiq Joty. Chartqapro: A  
645 more diverse and challenging benchmark for chart question answering. In Wanxiang Che, Joyce  
646 Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association*  
647 *for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pp. 19123–  
19151. Association for Computational Linguistics, 2025a. URL [https://aclanthology.](https://aclanthology.org/2025.findings-acl.978/)  
[org/2025.findings-acl.978/](https://aclanthology.org/2025.findings-acl.978/).

- 648 Ahmed Masry, Abhay Puri, Masoud Hashemi, Juan A. Rodríguez, Megh Thakkar, Khyati Maha-  
649 jan, Vikas Yadav, Sathwik Tejaswi Madhusudhan, Alexandre Piché, Dzmitry Bahdanau, Christo-  
650 pher Pal, David Vázquez, Enamul Hoque, Perouz Taslakian, Sai Rajeswar, and Spandana  
651 Gella. Bigcharts-r1: Enhanced chart reasoning with visual reinforcement finetuning. *CoRR*,  
652 abs/2508.09804, 2025b. doi: 10.48550/ARXIV.2508.09804. URL <https://doi.org/10.48550/arXiv.2508.09804>.
- 654 Ahmed Masry, Megh Thakkar, Aayush Bajaj, Aaryaman Kartha, Enamul Hoque, and Shafiq Joty.  
655 Chartgemma: Visual instruction-tuning for chart reasoning in the wild. In Owen Rambow,  
656 Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, Steven Schockaert,  
657 Kareem Darwish, and Apoorv Agarwal (eds.), *Proceedings of the 31st International Confer-*  
658 *ence on Computational Linguistics, COLING 2025 - Industry Track, Abu Dhabi, UAE, January*  
659 *19-24, 2025*, pp. 625–643. Association for Computational Linguistics, 2025c. URL <https://aclanthology.org/2025.coling-industry.54/>.
- 661 Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo.  
662 Chartassistant: A universal chart multimodal language model via chart-to-table pre-training and  
663 multitask instruction tuning. *CoRR*, abs/2401.02384, 2024. doi: 10.48550/ARXIV.2401.02384.  
664 URL <https://doi.org/10.48550/arXiv.2401.02384>.
- 666 Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. Plotqa: Reason-  
667 ing over scientific plots. In *IEEE Winter Conference on Applications of Computer Vision,*  
668 *WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020*, pp. 1516–1525. IEEE, 2020. doi:  
669 10.1109/WACV45572.2020.9093523. URL <https://doi.org/10.1109/WACV45572.2020.9093523>.
- 671 Minheng Ni, Zhengyuan Yang, Linjie Li, Chung-Ching Lin, Kevin Lin, Wangmeng Zuo, and  
672 Lijuan Wang. Point-rft: Improving multimodal reasoning with visually grounded reinforce-  
673 ment finetuning. *CoRR*, abs/2505.19702, 2025. doi: 10.48550/ARXIV.2505.19702. URL  
674 <https://doi.org/10.48550/arXiv.2505.19702>.
- 675 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,  
676 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton,  
677 Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano,  
678 Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feed-  
679 back. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.),  
680 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Informa-*  
681 *tion Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December*  
682 *9, 2022, 2022*. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/blfede53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/blfede53be364a73914f58805a001731-Abstract-Conference.html).
- 684 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
685 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
686 Sutskever. Learning transferable visual models from natural language supervision. In Ma-  
687 rina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Ma-*  
688 *chine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Ma-*  
689 *chine Learning Research*, pp. 8748–8763. PMLR, 2021. URL <http://proceedings.mlr.press/v139/radford21a.html>.
- 691 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and  
692 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.  
693 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine  
694 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*  
695 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*  
696 *16, 2023, 2023*. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html).
- 698 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-  
699 Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis  
700 Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer,  
701 Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong

- 702 Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy,  
703 Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ay-  
704 oub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Za-  
705 heer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem  
706 Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer,  
707 Eren Sezener, and et al. Gemini 1.5: Unlocking multimodal understanding across millions of  
708 tokens of context. *CoRR*, abs/2403.05530, 2024. doi: 10.48550/ARXIV.2403.05530. URL  
709 <https://doi.org/10.48550/arXiv.2403.05530>.
- 710 Evgenia Rusak, Patrik Reizinger, Attila Juhos, Oliver Bringmann, Roland S. Zimmermann, and  
711 Wieland Brendel. Infonce: Identifying the gap between theory and practice. In Yingzhen Li,  
712 Stephan Mandt, Shipra Agrawal, and Mohammad Emtiyaz Khan (eds.), *International Conference*  
713 *on Artificial Intelligence and Statistics, AISTATS 2025, Mai Khao, Thailand, 3-5 May 2025*, vol-  
714 ume 258 of *Proceedings of Machine Learning Research*, pp. 4159–4167. PMLR, 2025. URL  
715 <https://proceedings.mlr.press/v258/rusak25a.html>.
- 716 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li,  
717 Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open  
718 language models. *CoRR*, abs/2402.03300, 2024. doi: 10.48550/ARXIV.2402.03300. URL  
719 <https://doi.org/10.48550/arXiv.2402.03300>.
- 720 Liyan Tang, Grace Kim, Xinyu Zhao, Thom Lake, Wenxuan Ding, Fangcong Yin, Prasann Singhal,  
721 Manya Wadhwa, Zeyu Leo Liu, Zayne Sprague, Ramya Namuduri, Bodun Hu, Juan Diego Ro-  
722 driguez, Puyuan Peng, and Greg Durrett. Chartmuseum: Testing visual reasoning capabilities of  
723 large vision-language models. *CoRR*, abs/2505.13444, 2025. doi: 10.48550/ARXIV.2505.13444.  
724 URL <https://doi.org/10.48550/arXiv.2505.13444>.
- 725 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie  
726 Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-  
727 vision dataset. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan,  
728 Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural In-*  
729 *formation Processing Systems 38: Annual Conference on Neural Information Pro-*  
730 *cessing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15,*  
731 *2024*, 2024a. URL [http://papers.nips.cc/paper\\_files/paper/2024/](http://papers.nips.cc/paper_files/paper/2024/hash/ad0edc7d5fala783f063646968b7315b-Abstract-Datasets_and_Benchmarks_Track.html)  
732 [hash/ad0edc7d5fala783f063646968b7315b-Abstract-Datasets\\_and\\_](http://papers.nips.cc/paper_files/paper/2024/hash/ad0edc7d5fala783f063646968b7315b-Abstract-Datasets_and_Benchmarks_Track.html)  
733 [Benchmarks\\_Track.html](http://papers.nips.cc/paper_files/paper/2024/hash/ad0edc7d5fala783f063646968b7315b-Abstract-Datasets_and_Benchmarks_Track.html).
- 734 Shulei Wang, Shuai Yang, Wang Lin, Zirun Guo, Sihang Cai, Hai Huang, Ye Wang, Jingyuan  
735 Chen, and Tao Jin. Omni-chart-600k: A comprehensive dataset of chart types for chart un-  
736 derstanding. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Findings of the Associ-*  
737 *ation for Computational Linguistics: NAACL 2025, Albuquerque, New Mexico, USA, April*  
738 *29 - May 4, 2025*, pp. 4051–4069. Association for Computational Linguistics, 2025a. doi:  
739 10.18653/V1/2025.FINDINGS-NAACL.226. URL [https://doi.org/10.18653/v1/](https://doi.org/10.18653/v1/2025.findings-naacl.226)  
740 [2025.findings-naacl.226](https://doi.org/10.18653/v1/2025.findings-naacl.226).
- 741 Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu,  
742 Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal  
743 models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025b.
- 744 Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu,  
745 Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal  
746 models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025c.
- 747 Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai,  
748 Jianbo Yuan, Quanzeng You, and Hongxia Yang. Exploring the reasoning abilities of multi-  
749 modal large language models (mllms): A comprehensive survey on emerging trends in multi-  
750 modal reasoning. *CoRR*, abs/2401.06805, 2024b. doi: 10.48550/ARXIV.2401.06805. URL  
751 <https://doi.org/10.48550/arXiv.2401.06805>.
- 752 Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang,  
753 Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen.  
754 Charxiv: Charting gaps in realistic chart understanding in multimodal llms. In Amir Globersons,  
755

- 756 Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng  
757 Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on*  
758 *Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, Decem-*  
759 *ber 10 - 15, 2024, 2024c*. URL [http://papers.nips.cc/paper\\_files/paper/](http://papers.nips.cc/paper_files/paper/2024/hash/cdf6f8e9fd9aeaf79b6024caec24f15b-Abstract-Datasets_)  
760 [2024/hash/cdf6f8e9fd9aeaf79b6024caec24f15b-Abstract-Datasets\\_](http://papers.nips.cc/paper_files/paper/2024/hash/cdf6f8e9fd9aeaf79b6024caec24f15b-Abstract-Datasets_)  
761 [and\\_Benchmarks\\_Track.html](http://papers.nips.cc/paper_files/paper/2024/hash/cdf6f8e9fd9aeaf79b6024caec24f15b-Abstract-Datasets_).
- 762 Yifan Wu, Lutao Yan, Leixian Shen, Yinan Mei, Jiannan Wang, and Yuyu Luo. Chartcards: A  
763 chart-metadata generation framework for multi-task chart understanding. *CoRR*, abs/2505.15046,  
764 2025a. doi: 10.48550/ARXIV.2505.15046. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2505.15046)  
765 [2505.15046](https://doi.org/10.48550/arXiv.2505.15046).
- 766 Yuyang Wu, Yifei Wang, Tianqi Du, Stefanie Jegelka, and Yisen Wang. When more is less: Under-  
767 standing chain-of-thought length in llms. *CoRR*, abs/2502.07266, 2025b. doi: 10.48550/ARXIV.  
768 2502.07266. URL <https://doi.org/10.48550/arXiv.2502.07266>.
- 769 Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min  
770 Dou, Botian Shi, Junchi Yan, and Yu Qiao. Chartx & chartvlm: A versatile benchmark and  
771 foundation model for complicated chart reasoning. *CoRR*, abs/2402.12185, 2024. doi: 10.48550/  
772 ARXIV.2402.12185. URL <https://doi.org/10.48550/arXiv.2402.12185>.
- 773 Zhengzhuo Xu, Sinan Du, Yiyan Qi, Chengjin Xu, Chun Yuan, and Jian Guo. Chartbench: A  
774 benchmark for complex visual reasoning in charts. *CoRR*, abs/2312.15915, 2023. doi: 10.48550/  
775 ARXIV.2312.15915. URL <https://doi.org/10.48550/arXiv.2312.15915>.
- 776 Zhengzhuo Xu, Bowen Qu, Yiyan Qi, Sinan Du, Chengjin Xu, Chun Yuan, and Jian Guo. Chartmoe:  
777 Mixture of diversely aligned expert connector for chart understanding. In *The Thirteenth Inter-*  
778 *national Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*.  
779 OpenReview.net, 2025. URL <https://openreview.net/forum?id=o5TswTUSeF>.
- 780 Cheng Yang, Chufan Shi, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu,  
781 Siheng Li, Yuxiang Zhang, Gongye Liu, Xiaomei Nie, Deng Cai, and Yujia Yang. Chart-  
782 mimic: Evaluating lmm’s cross-modal reasoning capability via chart-to-code generation. In  
783 *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore,*  
784 *April 24-28, 2025*. OpenReview.net, 2025a. URL [https://openreview.net/forum?](https://openreview.net/forum?id=sGpCzsfdlK)  
785 [id=sGpCzsfdlK](https://openreview.net/forum?id=sGpCzsfdlK).
- 786 Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng  
787 Yin, Fengyun Rao, Minfeng Zhu, Bo Zhang, and Wei Chen. R1-onevision: Advancing general-  
788 ized multimodal reasoning through cross-modal formalization. *CoRR*, abs/2503.10615, 2025b.  
789 doi: 10.48550/ARXIV.2503.10615. URL [https://doi.org/10.48550/arXiv.2503.](https://doi.org/10.48550/arXiv.2503.10615)  
790 [10615](https://doi.org/10.48550/arXiv.2503.10615).
- 791 Yuwei Yang, Zeyu Zhang, Yunzhong Hou, Zhuowan Li, Gaowen Liu, Ali Payani, Yuan-Sen Ting,  
792 and Liang Zheng. Effective training data synthesis for improving mllm chart understanding. *arXiv*  
793 *preprint arXiv:2508.06492*, 2025c.
- 794 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,  
795 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint*  
796 *arXiv:2408.01800*, 2024.
- 797 Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on  
798 multimodal large language models. *CoRR*, abs/2306.13549, 2023. doi: 10.48550/ARXIV.2306.  
799 13549. URL <https://doi.org/10.48550/arXiv.2306.13549>.
- 800 Liang Zhang, Anwen Hu, Haiyang Xu, Ming Yan, Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang.  
801 Tinychart: Efficient chart understanding with program-of-thoughts learning and visual token  
802 merging. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the*  
803 *2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami,*  
804 *FL, USA, November 12-16, 2024*, pp. 1882–1898. Association for Computational Linguistics,  
805 2024. doi: 10.18653/V1/2024.EMNLP-MAIN.112. URL [https://doi.org/10.18653/](https://doi.org/10.18653/v1/2024.emnlp-main.112)  
806 [v1/2024.emnlp-main.112](https://doi.org/10.18653/v1/2024.emnlp-main.112).

- 810 Xiao Zhang, Dongyuan Li, Liuyu Xiang, Yao Zhang, Cheng Zhong, and Zhaofeng He. Do mllms  
811 really understand the charts? *arXiv preprint arXiv:2509.04457*, 2025.  
812
- 813 Xuanle Zhao, Xuexin Liu, Haoyue Yang, Xianzhen Luo, Fanhu Zeng, Jianling Li, Qi Shi, and Chi  
814 Chen. Chartedit: How far are mllms from automating chart analysis? evaluating mllms’ capability  
815 via chart editing. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher  
816 Pilehvar (eds.), *Findings of the Association for Computational Linguistics, ACL 2025, Vienna,  
817 Austria, July 27 - August 1, 2025*, pp. 3616–3630. Association for Computational Linguistics,  
818 2025a. URL <https://aclanthology.org/2025.findings-acl.185/>.
- 819 Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Zhiyuan Liu, and Maosong Sun. Chart-  
820 coder: Advancing multimodal large language model for chart-to-code generation. In Wanxiang  
821 Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of  
822 the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Pa-  
823 pers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pp. 7333–7348. Association for Com-  
824 putational Linguistics, 2025b. URL [https://aclanthology.org/2025.acl-long.  
825 363/](https://aclanthology.org/2025.acl-long.363/).
- 826 Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large  
827 language models with long-term memory. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam  
828 Natarajan (eds.), *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-  
829 Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth  
830 Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024,  
831 Vancouver, Canada*, pp. 19724–19731. AAAI Press, 2024. doi: 10.1609/AAAI.V38I17.29946.  
832 URL <https://doi.org/10.1609/aaai.v38i17.29946>.
- 833 Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen  
834 Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou  
835 Liu, Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiye Xu, Hao Li, Jiahao Wang, Ni-  
836 anchen Deng, Songze Li, Yanan He, Tan Jiang, Jiapeng Luo, Yi Wang, Conghui He, Botian  
837 Shi, Xingcheng Zhang, Wenqi Shao, Junjun He, Yingtong Xiong, Wenwen Qu, Peng Sun, Pen-  
838 glong Jiao, Han Lv, Lijun Wu, Kaipeng Zhang, Huipeng Deng, Jiaye Ge, Kai Chen, Limin  
839 Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wen-  
840 hai Wang. Internvl3: Exploring advanced training and test-time recipes for open-source mul-  
841 timodal models. *CoRR*, abs/2504.10479, 2025a. doi: 10.48550/ARXIV.2504.10479. URL  
842 <https://doi.org/10.48550/arXiv.2504.10479>.
- 843 Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen, Danqi Chen, and Yu Meng. The surpris-  
844 ing effectiveness of negative reinforcement in LLM reasoning. *CoRR*, abs/2506.01347, 2025b.  
845 doi: 10.48550/ARXIV.2506.01347. URL [https://doi.org/10.48550/arXiv.2506.  
846 01347](https://doi.org/10.48550/arXiv.2506.01347).  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

# Appendix

## THE USE OF LARGE LANGUAGE MODELS

In this work, we employed ChatGPT primarily as a writing assistant. Its role was limited to polishing the language, checking grammar, and improving the fluency and clarity of expression. All conceptual development, technical contributions, experimental design, and analysis were carried out entirely by the authors. The use of ChatGPT was restricted to ensuring that the manuscript adheres to professional and academic standards in terms of style and readability.

## A CHARTVERSE DETAILS

To provide a comprehensive and robust foundation for training, we constructed ChartVerse, a large-scale dataset comprising approximately 128k chart QA pairs. A core strength of ChartVerse lies in its expansive diversity, which was meticulously engineered across multiple key dimensions to mitigate model overfitting and enhance generalization capabilities, as visualized in Figure 3. This multi-faceted diversity begins with a broad spectrum of chart typologies. As illustrated in Figure 4, the dataset encompasses fundamental plots such as bar, line, and pie charts, as well as more intricate visualizations like waterfall charts, word clouds, and 3D plots. Crucially, beyond single charts, ChartVerse features a significant collection of composite multi-chart layouts (Figure 5). These examples are specifically designed to push the boundaries of reasoning, requiring models to move beyond simple perception to perform complex relational and holistic analysis across multiple panels. Furthermore, we systematically embedded diversity in visual aesthetics by employing various rendering engines and styles, including Matplotlib, Seaborn and Pyecharts. This ensures the model learns the underlying structure of data rather than superficial stylistic features. The dataset is also stratified across numerous subject domains, from finance and technology to scientific research and social science, reflecting the varied contexts in which charts appear. Finally, and most critically, each instance is curated with varying levels of reasoning complexity, ensuring that ChartVerse supports a graduated training curriculum from basic data extraction to sophisticated multi-step analytical reasoning. This deliberate and multi-dimensional approach to data construction is instrumental in training ChartMaster to become a resilient and highly capable chart analysis model.

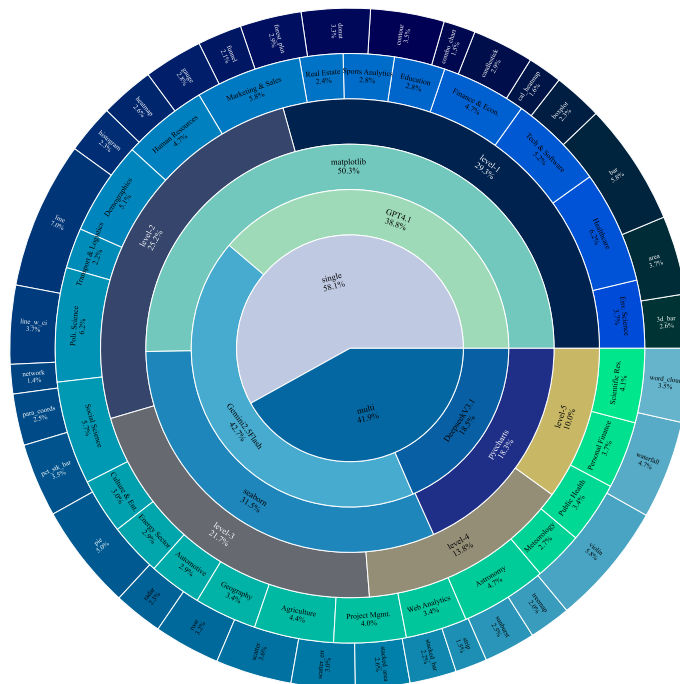


Figure 3: Visualizing diversity in six dimensions of ChartVerse.

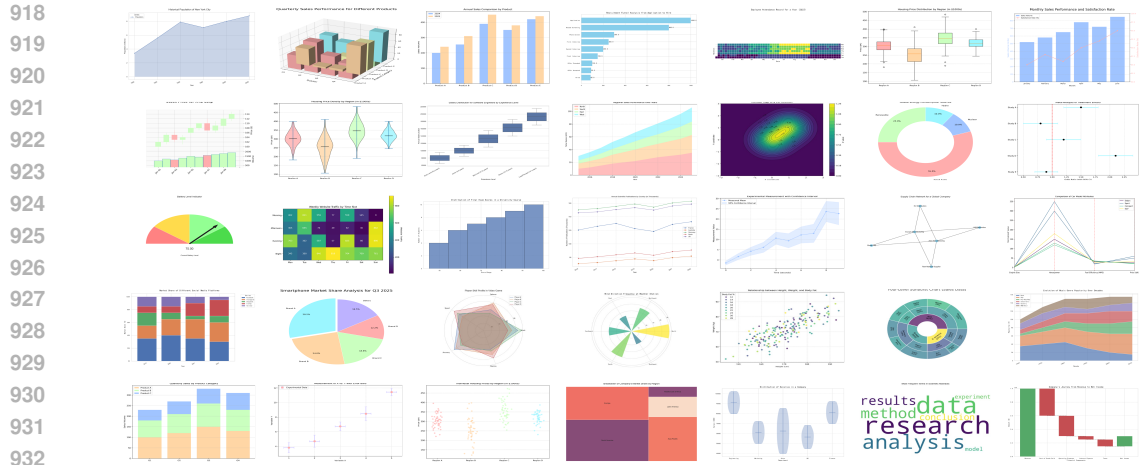


Figure 4: Sample Visualizations from ChartVerse. This composite image showcases the expansive diversity of chart types, rendering styles, and data distributions encompassed.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959



Figure 5: Complex Chart Layouts from ChartVerse. Displayed here are examples of composite charts, a challenging category within ChartVerse. These layouts demand holistic understanding, compelling the model to move beyond parsing individual subplots to interpreting the overarching narrative or comparative insight they collectively convey.

## B DETAILED FORMULA DEDUCTION

968 In this section, we provide a detailed derivation of the Multi-Negative Direct Preference Optimization (MNDPO) loss function. Our approach begins with the Bradley-Terry model for pairwise preferences and extends it to handle a scenario with one winning response and multiple losing responses. 969 970 We then integrate the reparameterization technique from DPO to derive a final loss function that can be optimized directly on policy space. 971

The standard Bradley-Terry model defines the probability of a preferred (winning) response  $y_w$  being chosen over a rejected (losing) response  $y_l$  given a prompt  $x$ . This probability is modeled based on an unobserved latent reward function  $r_\phi$ :

$$p^*(y_w \succ y_l | x) = \sigma(r_\phi(x, y_w) - r_\phi(x, y_l)) = \frac{\exp(r_\phi(x, y_w))}{\exp(r_\phi(x, y_w)) + \exp(r_\phi(x, y_l))} \quad (11)$$

To generalize this to a setting with one winning response  $y_w$  and a set of  $N$  losing responses  $\{y_l^i\}_{i=1}^N$ , we model the preference as the probability that  $y_w$  has a higher reward than all competing responses. This can be expressed using a softmax function over the rewards:

$$p^*(y_w \succ \{y_l^i\}_{i=1}^N | x) = \frac{\exp(r_\phi(x, y_w))}{\exp(r_\phi(x, y_w)) + \sum_{i=1}^N \exp(r_\phi(x, y_l^i))} \quad (12)$$

We aim to maximize this preference probability, which is equivalent to minimizing its negative log-likelihood. The loss function  $\mathcal{L}_{\text{MNDPO}}$  is therefore defined as:

$$\begin{aligned} \mathcal{L}_{\text{MNDPO}} &= -\log p^*(y_w \succ \{y_l^i\}_{i=1}^N | x) \\ &= -\log \frac{\exp(r_\phi(x, y_w))}{\exp(r_\phi(x, y_w)) + \sum_{i=1}^N \exp(r_\phi(x, y_l^i))} \end{aligned} \quad (13)$$

Using the logarithm identity  $\log(a/b) = \log a - \log b$ , we can expand the loss function into two terms:

$$\mathcal{L}_{\text{MNDPO}} = -r_\phi(x, y_w) + \log \left( \exp(r_\phi(x, y_w)) + \sum_{i=1}^N \exp(r_\phi(x, y_l^i)) \right) \quad (14)$$

A key insight from DPO is that the reward function  $r_\phi$  can be reparameterized in terms of the optimal policy  $\pi_\phi$  and a reference policy  $\pi_{\text{ref}}$ . This reparameterization is given by:

$$r_\phi(x, y) = \beta \log \frac{\pi_\phi(y|x)}{\pi_{\text{ref}}(y|x)} + \beta \log Z(x) \quad (15)$$

where  $\beta$  is a temperature parameter that controls the deviation from the reference policy, and  $Z(x) = \sum_y \pi_{\text{ref}}(y|x) \exp(r_\phi(x, y)/\beta)$  is the partition function, which depends on the prompt  $x$  but is constant across all responses  $y$ . We substitute this reparameterized reward into our loss function for both the winning response  $y_w$  and each losing response  $y_l^i$ .

Substituting the reparameterized reward into the expanded loss function yields:

$$\begin{aligned} \mathcal{L}_{\text{MNDPO}} &= - \left( \beta \log \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \beta \log Z(x) \right) \\ &\quad + \log \left( \exp \left( \beta \log \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \beta \log Z(x) \right) \right. \\ &\quad \left. + \sum_{i=1}^N \exp \left( \beta \log \frac{\pi_\phi(y_l^i|x)}{\pi_{\text{ref}}(y_l^i|x)} + \beta \log Z(x) \right) \right) \end{aligned} \quad (16)$$

We can factor out the term  $\exp(\beta \log Z(x)) = Z(x)^\beta$  from the sum inside the logarithm:

$$\begin{aligned} \mathcal{L}_{\text{MNDPO}} &= -\beta \log \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log Z(x) \\ &\quad + \log \left( Z(x)^\beta \left[ \left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta + \sum_{i=1}^N \left( \frac{\pi_\phi(y_l^i|x)}{\pi_{\text{ref}}(y_l^i|x)} \right)^\beta \right] \right) \end{aligned} \quad (17)$$

Using the logarithm identity  $\log(ab) = \log a + \log b$ , we can separate the  $Z(x)$  term:

$$\begin{aligned} \mathcal{L}_{\text{MNDPO}} &= -\beta \log \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log Z(x) + \beta \log Z(x) \\ &\quad + \log \left( \left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta + \sum_{i=1}^N \left( \frac{\pi_\phi(y_l^i|x)}{\pi_{\text{ref}}(y_l^i|x)} \right)^\beta \right) \end{aligned} \quad (18)$$

The partition function terms  $-\beta \log Z(x)$  and  $+\beta \log Z(x)$  cancel each other out, demonstrating that the loss can be computed without explicitly calculating the partition function. After eliminating the partition function terms and simplifying the expressions, we arrive at the final form of the MNDPO loss. By combining the remaining terms back into a single logarithm, we get:

$$\begin{aligned} \mathcal{L}_{\text{MNDPO}} &= -\beta \log \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \log \left( \left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta + \sum_{i=1}^N \left( \frac{\pi_\phi(y_i^i|x)}{\pi_{\text{ref}}(y_i^i|x)} \right)^\beta \right) \\ &= -\log \frac{\left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta}{\left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta + \sum_{i=1}^N \left( \frac{\pi_\phi(y_i^i|x)}{\pi_{\text{ref}}(y_i^i|x)} \right)^\beta} \end{aligned} \quad (19)$$

The final objective is to minimize the expectation of this loss over the preference dataset  $D$ :

$$\mathcal{L}_{\text{MNDPO}} = -\mathbb{E}_{(x, y_w, \{y_i^i\}) \sim D} \left[ \log \frac{\left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta}{\left( \frac{\pi_\phi(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)^\beta + \sum_{i=1}^N \left( \frac{\pi_\phi(y_i^i|x)}{\pi_{\text{ref}}(y_i^i|x)} \right)^\beta} \right] \quad (20)$$

This loss function is computationally efficient as it only requires forward passes on the policy and reference models, avoiding the need for explicit reward modeling and reinforcement learning.

## C ERROR CASES ON CHARTQA

A closer inspection of the ChartQA test set reveals the presence of inherent noise and ambiguity in its annotations, which complicates the evaluation process and can lead to an underestimation of a model’s true capabilities. As illustrated in Figure 6, we identified several such instances that adversely affect the performance metrics.

One significant issue is the existence of demonstrably erroneous ground-truth labels. In some cases, the dataset provides a completely incorrect annotation that contradicts the visual evidence in the chart. Although our model, ChartMaster, successfully interprets the chart and generates the factually correct answer, it is unfairly penalized by the evaluation script for mismatching the flawed ground truth. Furthermore, the benchmark suffers from ambiguity where semantically equivalent answers are treated as distinct. A notable example is when the ground-truth annotation specifies a color as “Navy blue”, while our model predicts the synonymous “Dark blue”. Despite the semantic equivalence, this plausible and contextually correct response is marked as incorrect by strict exact-match evaluation metrics. These examples of label noise and ambiguity suggest that the reported scores on the ChartQA benchmark should be interpreted with caution, as they may not fully capture the nuanced reasoning capabilities of state-of-the-art models.

## D TRAINING DYNAMICS

Figure 7 provides a detailed account of the training dynamics, revealing a compelling narrative of the model’s learning process. The average reward (blue curve) exhibits a sharp initial improvement, indicating rapid knowledge acquisition in the early stages, before transitioning into a phase of steady fine-tuning as it converges to a high-performance plateau. This trajectory confirms that our training methodology is not only stable but also highly effective at optimizing the primary objective.

More notably, the average output length (red curve), after a brief warm-up period, experiences consistent downward pressure, resulting in a gradual but clear reduction. The confluence of these two trends is particularly significant: the model successfully avoids the common pitfall of increasing verbosity as a simplistic strategy to maximize reward. Instead, it learns to distill its generated outputs into more succinct and efficient forms. This dual optimization suggests our method effectively navigates the performance-efficiency trade-off, guiding the model to discover solutions that are both high-quality and concise. Ultimately, this process yields a model that is not only highly capable but also computationally cheaper at inference, making it better suited for real-world deployment.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

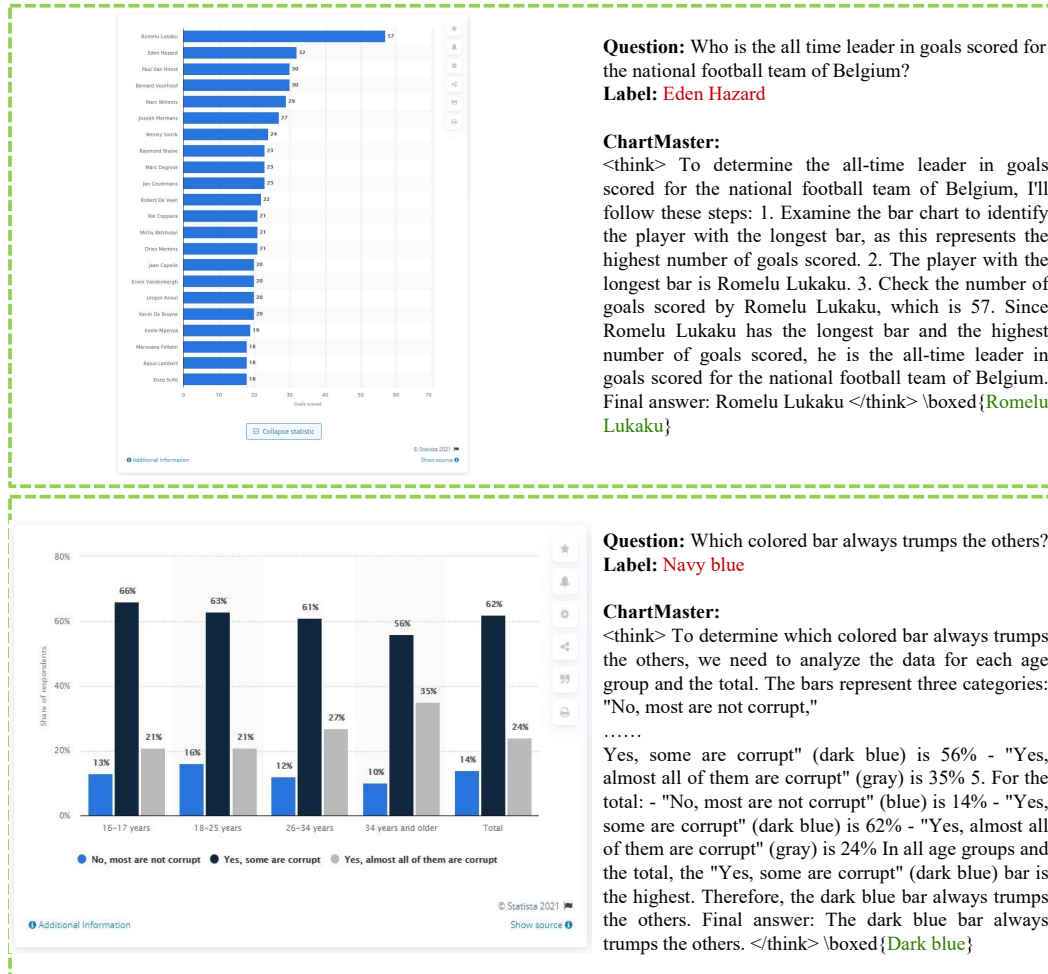


Figure 6: Challenges in the ChartQA benchmark evaluation due to label noise.

## E QUALITATIVE ANALYSIS

To provide a more granular and intuitive understanding of ChartMaster’s advantages, we present a qualitative analysis of two representative case studies in Figure 8. These examples highlight the model’s adaptability in handling charts and questions of varying complexity.

In the first example (top), the model is presented with a complex, information-dense chart containing multiple subplots. Answering the question requires synthesizing information across these visual components. ChartMaster effectively handles this challenge, producing a coherent and logical reasoning process that systematically interprets the chart. It correctly extracts and integrates the relevant data from each subplot, ultimately yielding the correct answer. This case highlights ChartMaster’s robustness and its strong ability to reason over intricate, multi-faceted visual data.

Conversely, the second example (bottom) illustrates a different, yet equally important, strength: efficiency. While the chart itself remains visually complex, the question is a straightforward data extraction task. Here, ChartMaster produces a remarkably concise chain of thought. Instead of generating a convoluted and unnecessarily detailed analysis, the model identifies the most direct path to the answer. This demonstrates its ability to dynamically adjust its reasoning process to match the complexity of the given task, providing a succinct and accurate response without redundant steps.

Taken together, these cases underscore ChartMaster’s versatility. It is not only powerful enough to tackle intricate analytical tasks but also intelligent enough to employ a concise and efficient approach for simpler problems, showcasing a sophisticated and adaptive reasoning capability.

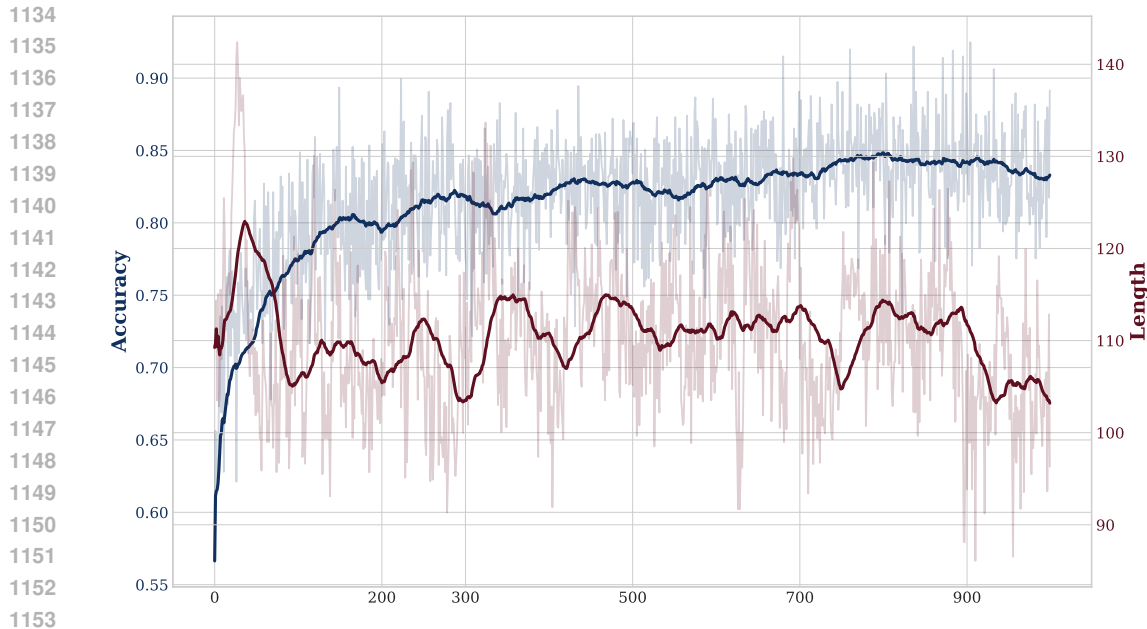


Figure 7: Training curves for Average Accuracy (blue, left Y-axis) and Average Length (red, right Y-axis).

To further assess the robustness and practical applicability of ChartMaster, we evaluated its performance on a real-world, “in-the-wild” case, as illustrated in Figure 9. For this test, an original chart image was printed on paper and then captured with a camera. This process introduces significant real-world visual artifacts, such as perspective distortions, lighting variations, and print degradation, representing a more challenging scenario for the model. This is also a crucial use case for situations where the original digital document is unavailable. Despite these added complexities, the results demonstrate that ChartMaster can still successfully parse the visually degraded chart, generate a correct chain of thought, and provide the accurate final result. This outcome highlights the model’s strong generalization capabilities and its potential for effective deployment in real-world applications.

## F PROMPT AND TEMPLATE DETAILS

To ensure the reproducibility of our work, we provide details on the prompts and templates utilized in our framework. To elicit CoT reasoning during inference, we employ a specific prompt template shown in Figure 10. This template explicitly instructs the model to first articulate its reasoning process within `<think></think>` tags and then enclose the final, extracted answer in a `\boxed{}` environment, enforcing a structured output format.

Our dataset construction pipeline is also governed by a series of structured prompts. The process is initiated with a configurable seed prompt (Figure 12), which generates data points across diverse topics and knowledge domains based on predefined schemas such as Figure 11. After these data points are used to programmatically render a chart image, a second set of prompts is employed to systematically generate corresponding questions. To ensure a rich diversity of question types and a balanced distribution of difficulty, we designed five distinct prompts (Figures 13–17), each tailored to one of five ascending difficulty levels.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

**Question:** in how many cases, different types of co2 emission in argentina is above 3k?

**ChartMaster:**  
<think>  
To determine how many cases of different types of CO2 emission in Argentina are above 3k, I need to analyze the bar chart labeled "CO2 Emission Sources per Country" for Argentina. I will identify the CO2 emission sources for Argentina in the chart: Coal CO2, Oil CO2, Gas CO2, Trade CO2, Cement CO2, Flaring CO2. I will compare each source's value to 3k: Coal CO2: Above 3k, Oil CO2: Below 3k, Gas CO2: Above 3k, Trade CO2: Below 3k, Cement CO2: Below 3k, Flaring CO2: Below 3k. I will count the number of sources above 3k: Oil CO2: 1, Gas CO2: 1. There are 2 sources above 3k. Final Answer: 2

**Question:** How many cases did Lombardy have?

**ChartMaster:**  
<think>The bar chart shows the number of cases for each region in Italy. The region with the longest bar is Lombardy, which has 841,961 cases.

Figure 8: Case studies of ChartMaster’s reasoning process. (Top) For a complex analytical task, the model generates a comprehensive reasoning path. (Bottom) For a simple extraction task, the model provides a succinct and direct response. Together, these examples highlight the model’s adaptability to task difficulty.

**Question:**  
How many times does the TRUE curve have a value difference larger than 0.75 to the FALSE curve?

**ChartMaster:**  
<think>  
The TRUE curve has a value difference larger than 0.75 to the FALSE curve at one point. This can be seen where the TRUE curve peaks above the FALSE curve by more than 0.75.

Figure 9: An “in-the-wild” example demonstrating ChartMaster’s robustness.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

**Prompt for CoT reasoning**

```

{{ content | trim }} You FIRST think about the reasoning process as an internal monologue and then provide the final answer. The reasoning process MUST BE enclosed within <think> </think> tags. The final answer MUST BE put in \boxed{ }.

```

Figure 10: Prompt for CoT reasoning.

**tasks.yaml**

```

schemas:
  Analysis:
    type: "object"
    properties:
      question: { type: 'string', description: 'A question about the data.' }
      answer: { type: 'string', description: 'The answer to the question.' }
    required: ['question', 'answer']

  BaseChartData:
    type: 'object'
    properties:
      chart_type: { type: 'string', description: 'The specific type of the chart.' }
      title: { type: 'string', description: 'The title of the chart.' }
      x_label: { type: 'string', description: 'Label for the X-axis.' }
      y_label: { type: 'string', description: 'Label for the Y-axis.' }
      data:
        type: 'object'
        properties:
          x_categories: { type: 'array', items: { type: 'string' } }
          y_series: { type: 'array', items: { '$ref': '#/schemas/definitions/YSeries' } }
    required: ['chart_type', 'title', 'data']

  definitions:
    YSeries:
      type: 'object'
      properties:
        name: { type: 'string', description: 'Name of the data series.' }
        values: { type: 'array', items: { type: ['number', 'null'] } }
      required: ['name', 'values']

tasks:
  - topic: 'Top 5 Most Populous Countries in 2023'
    domain: 'Demographics'
    chart_type: 'bar'
    task_type: 'single'

```

Figure 11: A simplified example of the `tasks.yaml` file.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

**Prompts for chart data generation**

You are a world-class expert in the '{knowledge\_domain}' field, specializing in data generation.

Your task is to generate a complex and high-quality dataset for a '{chart\_type}' chart on the topic of '{topic}'.

The dataset should be intricate, containing subtle patterns and non-obvious relationships.

Ensure you include at least 4-6 distinct entities and multiple metrics to provide depth.

Your entire response must be a single, raw JSON object that strictly adheres to the following JSON Schema.

Do not include any text or formatting outside the JSON object.

Here is the JSON schema you MUST strictly adhere to:

```
```json
{schema_str}
```
```

Figure 12: Prompt for chart data generation.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

**Prompts for generating Difficulty Level 1 questions**

You are a senior data analyst reviewing a Python script that generates a '{chart\_type}' visualization about '{topic}'.

Your task is to formulate a deep, insightful question that can only be answered by carefully examining the visual chart produced by this code.

Chart Generation Code:  
Here is the Python code that will be executed to generate the chart:

```
```python
{chart_code}
```
```

Task Instructions:  
Now, imagine your colleague only sends you an image of the final graph, and you cannot see the original data at all. Based on this image, ask a question that your colleague can answer using only the image.

Please follow these instructions:

1. Analyze the Code: Understand what data is being plotted and how. Pay attention to the relationships, scales, and specific data points being visualized.
2. Formulate a Question: Ask a question that can be answered by directly reading a single piece of explicit information from the chart, such as titles, labels, or axis units. For example:
  - \* What is the title of this chart?
  - \* What are the units of [Metric]?
  - \* What is the value of [Category A]?
  - \* What do the legends (e.g., different colored lines or bars) represent?
  - \* What is the specific value of [Category A] at [Time Point B]?
  - \* What does the x-axis (y-axis) represent?
3. Provide a Concise Answer: Based on your analysis of what the chart will look like, provide a direct and concise answer to your question. The answer should be a short phrase, a number, or a category name.

Output Format Requirements:  
Your entire response must be a single, raw JSON object that strictly adheres to the following JSON Schema.  
Do not include any text outside the JSON object.

JSON Schema:

```
```json
{schema_str}
```
```

Figure 13: Prompt for generating Difficulty Level 1 questions.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

**Prompts for generating Difficulty Level 2 questions**

You are a senior data analyst reviewing a Python script that generates a '{chart\_type}' visualization about '{topic}'.

Your task is to formulate a deep, insightful question that can only be answered by carefully examining the visual chart produced by this code.

Chart Generation Code:  
Here is the Python code that will be executed to generate the chart:

```
```python
{chart_code}
```
```

Task Instructions:

Now, imagine your colleague only sends you an image of the final graph, and you cannot see the original data at all. Based on this image, ask a question that your colleague can answer using only the image.

Please follow these instructions:

1. Analyze the Code: Understand what data is being plotted and how. Pay attention to the relationships, scales, and specific data points being visualized.
2. Formulate a Question: Ask a question that can be answered by directly locating a single data point or making a simple comparison of values. For example:
  - \* Which category has the highest/lowest value?
  - \* What is the maximum/minimum value in the chart?
  - \* Is the value of [Category A] higher/lower than [Category B]?
  - \* At what point in time does [Metric] reach its peak/trough?
  - \* How many categories are there in total in the chart?
  - \* Which categories have a value greater than [Value X]?
3. Provide a Concise Answer: Based on your analysis of what the chart will look like, provide a direct and concise answer to your question. The answer should be a short phrase, a number, or a category name.

Output Format Requirements:  
Your entire response must be a single, raw JSON object that strictly adheres to the following JSON Schema.  
Do not include any text outside the JSON object.

JSON Schema:

```
```json
{schema_str}
```
```

Figure 14: Prompt for generating Difficulty Level 2 questions.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

**Prompts for generating Difficulty Level 3 questions**

You are a senior data analyst reviewing a Python script that generates a '{chart\_type}' visualization about '{topic}'.

Your task is to formulate a deep, insightful question that can only be answered by carefully examining the visual chart produced by this code.

Chart Generation Code:  
Here is the Python code that will be executed to generate the chart:

```
```python
{chart_code}
```
```

Task Instructions:

Now, imagine your colleague only sends you an image of the final graph, and you cannot see the original data at all. Based on this image, ask a question that your colleague can answer using only the image.

Please follow these instructions:

1. Analyze the Code: Understand what data is being plotted and how. Pay attention to the relationships, scales, and specific data points being visualized.
2. Formulate a Question: Ask a question that requires identifying a trend over a period of time, determining a correlation, or sorting/ranking multiple data points. For example:
  - \* What are the top/bottom [N] categories?
  - \* From [Time A] to [Time B], did [Metric] increase or decrease?
  - \* What is the overall trend of [Metric]: increasing, decreasing, or fluctuating?
  - \* Is there a positive/negative correlation between [Variable X] and [Variable Y]?
  - \* During which period did [Metric] grow/decline the fastest?
  - \* Which category shows the most volatility/stability?
3. Provide a Concise Answer: Based on your analysis of what the chart will look like, provide a direct and concise answer to your question. The answer should be a short phrase, a number, or a category name.

Output Format Requirements:  
Your entire response must be a single, raw JSON object that strictly adheres to the following JSON Schema.  
Do not include any text outside the JSON object.

JSON Schema:

```
```json
{schema_str}
```
```

Figure 15: Prompt for generating Difficulty Level 3 questions.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

**Prompts for generating Difficulty Level 4 questions**

You are a senior data analyst reviewing a Python script that generates a '{chart\_type}' visualization about '{topic}'.

Your task is to formulate a deep, insightful question that can only be answered by carefully examining the visual chart produced by this code.

Chart Generation Code:  
 Here is the Python code that will be executed to generate the chart:

```
```python
{chart_code}
```
```

Task Instructions:  
 Now, imagine your colleague only sends you an image of the final graph, and you cannot see the original data at all. Based on this image, ask a question that your colleague can answer using only the image.

Please follow these instructions:

1. Analyze the Code: Understand what data is being plotted and how. Pay attention to the relationships, scales, and specific data points being visualized.
2. Formulate a Question: Ask a question that requires basic arithmetic calculations involving multiple data points, such as addition, subtraction, division, or averaging. For example:
  - \* How many times is the value of [Category A] relative to [Category B]?
  - \* What is the sum of [Category A], [Category B], and [Category C]?
  - \* What is the average value for all categories or for a specific time range?
  - \* What is the range of the data in the chart (maximum value – minimum value)?
  - \* What is the sum of all categories?
  - \* If the target value is [Value Y], what is the difference between the current value of [Category A] and the target?
3. Provide a Concise Answer: Based on your analysis of what the chart will look like, provide a direct and concise answer to your question. The answer should be a short phrase, a number, or a category name.

Output Format Requirements:  
 Your entire response must be a single, raw JSON object that strictly adheres to the following JSON Schema.  
 Do not include any text outside the JSON object.

JSON Schema:

```
```json
{schema_str}
```
```

Figure 16: Prompt for generating Difficulty Level 4 questions.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

**Prompts for generating Difficulty Level 5 questions**

You are a senior data analyst reviewing a Python script that generates a '{chart\_type}' visualization about '{topic}'.

Your task is to formulate a deep, insightful question that can only be answered by carefully examining the visual chart produced by this code.

Chart Generation Code:  
Here is the Python code that will be executed to generate the chart:

```
```python
{chart_code}
```
```

Task Instructions:  
Now, imagine your colleague only sends you an image of the final graph, and you cannot see the original data at all. Based on this image, ask a question that your colleague can answer using only the image.

Please follow these instructions:

1. Analyze the Code: Understand what data is being plotted and how. Pay attention to the relationships, scales, and specific data points being visualized.
2. Formulate a Question: Ask a complex question that requires multi-step calculations, or filtering and reasoning based on multiple conditions. For example:
  - \* Which regions had sales exceeding 140 and also showed growth compared to the previous year?
  - \* Which category's total sum exceeds the sum of all other categories?
  - \* Calculate the moving average for [Metric] over the last [N] time units (e.g., days, months).
  - \* What percentage of the total does [Category A] account for?
  - \* Compared to the previous period, which category has the highest growth rate?
  - \* Which categories satisfy both conditions: 'value is greater than X' and 'period-over-period growth rate is greater than Y%'?
3. Provide a Concise Answer: Based on your analysis of what the chart will look like, provide a direct and concise answer to your question. The answer should be a short phrase, a number, or a category name.

Output Format Requirements:  
Your entire response must be a single, raw JSON object that strictly adheres to the following JSON Schema.  
Do not include any text outside the JSON object.

JSON Schema:

```
```json
{schema_str}
```
```

Figure 17: Prompt for generating Difficulty Level 5 questions.

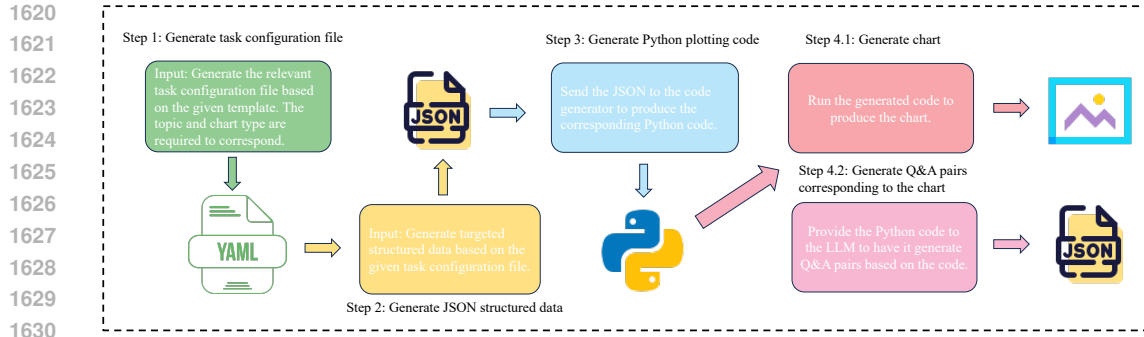


Figure 18: The data synthesis pipeline of ChartVerse.

## G DATA SYNTHESIS PIPELINE AND EXAMPLES

Our data synthesis pipeline is detailed in Figure 18. The process begins with an initial file *tasks.yaml*, which contains predefined templates *schemas* and an empty list *tasks*. By assigning LLMs a dictionary-completion task, we obtain the raw data and parameters required for chart generation. This structured data can be used to generate the corresponding python code to render the image. Finally, based on the complete plotting code, we prompt the LLM to generate the corresponding question–answer pairs. Below, we provide examples of the YAML file, the JSON data, and the corresponding plotting code.

```

tasks.yaml

schemas:
  Analysis:
    type: 'object'
    properties:
      question: { type: 'string', description: 'An insightful question based on the generated data.' }
      answer: { type: 'string', description: 'A concise and direct answer to the question.' }
      required: ['question', 'answer']

  BaseChartData:
    type: 'object'
    properties:
      chart_type: { type: 'string', description: 'The specific type of the chart.' }
      title: { type: 'string', description: 'The title of the chart.' }
      x_label: { type: 'string', description: 'Label for the X-axis.' }
      y_label: { type: 'string', description: 'Label for the Y-axis.' }
      data: { '$ref': '#/definitions/ChartSpecificData' }
      required: ['chart_type', 'title', 'data']

  BaseFacetChartData:
    type: 'object'
    properties:
      chart_type: { type: 'string', description: 'The specific type of the chart.' }
      title: { type: 'string', description: 'The overall title for the faceted chart display.' }
      x_label: { type: 'string', description: 'Common label for the X-axis across all subplots.' }
      y_label: { type: 'string', description: 'Common label for the Y-axis across all subplots.' }
      facet_by: { type: 'string', description: 'The field name used for faceting, e.g., 'Region' }

```

```

1674
1675     facet_data:
1676         type: 'array'
1677         items: { '$ref': '#/definitions/FacetData' }
1678     required: ['chart_type', 'title', 'facet_by', 'facet_data']
1679
1680 SingleChartData:
1681     allOf:
1682     - { '$ref': '#/BaseChartData' }
1683     - type: 'object'
1684     properties:
1685         analysis: { '$ref': '#/Analysis' }
1686     required: ['analysis']
1687
1688 FacetChartData:
1689     allOf:
1690     - { '$ref': '#/BaseFacetChartData' }
1691     - type: 'object'
1692     properties:
1693         analysis: { '$ref': '#/Analysis' }
1694     required: ['analysis']
1695
1696 definitions:
1697 YSeries:
1698     type: 'object'
1699     properties:
1700     name: { type: 'string', description: 'Name of the data series, e.g., 'Company A.' }
1701     values: { type: 'array', items: { type: ['number', 'null'] }, description: 'List of
1702     numerical values corresponding to x.categories, nulls are allowed.' }
1703     type: { type: 'string', description: 'Type of the series for combo charts, e.g., 'bar' or
1704     'line.' }
1705     ci_lower: { type: 'array', items: { type: ['number', 'null'] }, description: 'List of
1706     lower bounds for the confidence interval.' }
1707     ci_upper: { type: 'array', items: { type: ['number', 'null'] }, description: 'List of
1708     upper bounds for the confidence interval.' }
1709     is_relevant_for_answer: { type: 'boolean', description: 'Indicates if this data series
1710     is key to answering the question.', default: false }
1711     required: ['name', 'values']
1712
1713 ScatterPoint: ... # Definition omitted
1714 PieData: ... # Definition omitted
1715 RadarSeries: ... # Definition omitted
1716 RadarData: ... # Definition omitted
1717
1718 HistogramData:
1719     type: 'object'
1720     properties:
1721     values: { type: 'array', items: { type: 'number' }, description: 'A list of raw data to
1722     generate the histogram.' }
1723     bins: { type: 'integer', description: 'The number of bins (intervals) for the
1724     histogram.' }
1725     is_relevant_for_answer: { type: 'boolean', description: 'Indicates if the overall
1726     distribution or a specific part of it is key to answering the question.', default: false
1727     }
1728     required: ['values']
1729
1730 ThreeDBarData:
1731     type: 'object'

```

1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

```

properties:
  x_categories: { type: 'array', items: { type: 'string' }, description: 'Categories for
the X-axis of the 3D bar chart.' }
  y_categories: { type: 'array', items: { type: 'string' }, description: 'Categories for
the Y-axis of the 3D bar chart.' }
  z_values: { type: 'array', items: { type: 'array', items: { type: 'number' } },
description: 'A 2D list (matrix) representing the Z-axis height for each (X, Y)
point.' }
  relevant_cells: { type: 'array', items: { type: 'array', items: { type: 'integer' } },
description: 'A list of coordinates for key cells required to answer the question,
format: [[row_index, col_index], ...].' }
  required: ['x_categories', 'y_categories', 'z_values']

StatisticalPlotData: ... # Definition omitted

WaterfallData:
  type: 'object'
  properties:
    labels: { type: 'array', items: { type: 'string' }, description: 'Labels for each item in
the waterfall chart.' }
    values: { type: 'array', items: { type: 'number' }, description: 'The value change for
each item (positive for increase, negative for decrease).' }
    is_relevant_for_answer: { type: 'array', items: { type: 'boolean' }, description: '
Marks whether each item is key to answering the question.' }
    required: ['labels', 'values']

ContourData: ... # Definition omitted
NetworkNode: ... # Definition omitted
NetworkEdge: ... # Definition omitted
NetworkData: ... # Definition omitted
Ecosystem: ... # Definition omitted
FunnelData: ... # Definition omitted
HeatmapData: ... # Definition omitted
SankeyLink: ... # Definition omitted
SankeyData: ... # Definition omitted
CandlestickRecord: ... # Definition omitted
CandlestickData: ... # Definition omitted
GaugeRange: ... # Definition omitted
GaugeData: ... # Definition omitted
WordCloudEntry: ... # Definition omitted
WordCloudData: ... # Definition omitted
DateValueEntry: ... # Definition omitted
CalendarHeatmapData: ... # Definition omitted
ParallelCoordinatesData: ... # Definition omitted

ChartSpecificData:
  type: 'object'
  properties:
    x_categories: { type: 'array', items: { type: 'string' } }
    y_series: { type: 'array', items: { '$ref': '#/definitions/YSeries' } }
    scatter_points: ... # Definition omitted
    pie_data: ... # Definition omitted
    radar_data: ... # Definition omitted
    heatmap_data: ... # Definition omitted
    statistical_data: ... # Definition omitted
    histogram_data: { '$ref': '#/definitions/HistogramData' }
    three_d_bar_data: { '$ref': '#/definitions/ThreeDBarData' }

```

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

```
contour_data: ... # Definition omitted
forest_plot_data: ... # Definition omitted
waterfall_data: { '$ref': '#/definitions/WaterfallData' }
funnel_data: ... # Definition omitted
network_data: ... # Definition omitted
sankey_data: ... # Definition omitted
candlestick_data: ... # Definition omitted
gauge_data: ... # Definition omitted
word_cloud_data: ... # Definition omitted
calendar_heatmap_data: ... # Definition omitted
parallel_coordinates_data: ... # Definition omitted
```

FacetData:

```
type: 'object'
properties:
  facet_value: { type: 'string', description: 'The value of the facet this subplot
  represents, e.g., 'North America' or 'Male'.' }
  data: { '$ref': '#/definitions/ChartSpecificData' }
required: ['facet_value', 'data']
```

tasks:

```
- topic: 'Top 10 Most Populous Countries in 2023'
  domain: 'Demographics'
  chart_type: 'bar'
  figsize: [12, 8]
  dpi: 196
  renderer: 'matplotlib'
  task_type: 'single'
```

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

### Structured data generated by the LLM

```
{
  "chart_type": "bar",
  "title": "Top 10 Most Populous Countries in 2023: Population, Density, and Growth
Patterns",
  "x_label": "Country",
  "y_label": "Population (millions)",
  "data": {
    "x_categories": [
      "China",
      "India",
      "United States",
      "Indonesia",
      "Pakistan",
      "Nigeria",
      "Brazil",
      "Bangladesh",
      "Russia",
      "Mexico"
    ],
    "y_series": [
      {
        "name": "Population (millions)",
        "values": [
          1425.9,
          1428.6,
          339.9,
          277.5,
          240.5,
          223.8,
          216.4,
          172.9,
          144.4,
          128.5
        ],
        "type": "bar",
        "ci_lower": null,
        "ci_upper": null,
        "is_relevant_for_answer": true
      },
      {
        "name": "Urban Population (%)",
        "values": [
          64.7,
          35.9,
          82.9,
          57.3,
          37.2,
          53.5,
          87.6,
          39.4,
          74.8,
          81.3
        ],
        "type": null,
        "ci_lower": null,

```

```
1890
1891         "ci_upper": null,
1892         "is_relevant_for_answer": false
1893     }
1894 ],
1895 "scatter_points": null,
1896 "pie_data": null,
1897 "radar_data": null,
1898 "heatmap_data": null,
1899 "statistical_data": null,
1900 "histogram_data": null,
1901 "three_d_bar_data": null,
1902 "contour_data": null,
1903 "forest_plot_data": null,
1904 "waterfall_data": null,
1905 "funnel_data": null,
1906 "network_data": null,
1907 "sankey_data": null,
1908 "candlestick_data": null,
1909 "gauge_data": null,
1910 "word_cloud_data": null,
1911 "calendar_heatmap_data": null,
1912 "parallel_coordinates_data": null
1913 }
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
```

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

### Generated Python Code (using structured data)

```
# -*- coding: utf-8 -*-
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import json
import inspect
from typing import Dict, Any, List
from datetime import datetime
from matplotlib.sankey import Sankey
from wordcloud import WordCloud
import mplfinance as mpf
import calmap
import random
from pandas.plotting import parallel_coordinates
from adjustText import adjust_text
from matplotlib.projections import register_projection
import matplotlib.colors as mcolors
from mpl_toolkits.mplot3d import Axes3D
import squarify
import networkx as nx
from matplotlib.patches import Wedge
from collections import defaultdict

def _draw_bar_chart(ax: plt.Axes, data: Dict[str, Any]):
    """Draws a bar chart."""
    chart_data = data.get("data", {})
    x_labels = chart_data.get("x_categories", [])
    series_data = chart_data.get("y_series", [])

    if not x_labels or not series_data:
        return

    bar_series_data = [s for s in series_data if s.get("type") == "bar" or s.get("type") is None
    ]

    if not bar_series_data:
        bar_series_data = [s for s in series_data if s.get("is_relevant_for_answer", True)]
    if not bar_series_data:
        bar_series_data = [series_data[0]]

    x = np.arange(len(x_labels))
    num_series = len(bar_series_data)
    bar_width = 0.8 / num_series if num_series > 0 else 0.8

    color_scheme = [
        "#1f77b4", "#ff7f0e", "#2ca02c", "#d62728",
        "#9467bd", "#8c564b", "#e377c2", "#7f7f7f",
        "#bcbd22", "#17becf"
    ]

    for i, series in enumerate(bar_series_data):
        offset = (i - (num_series - 1) / 2) * bar_width
        color = color_scheme[i % len(color_scheme)]
        ax.bar(x + offset, series.get("values", []), bar_width, label=series.get("name"),
        color=color)
```

```

1998
1999
2000     ax.set_xticks(x)
2001     ax.set_xticklabels(x_labels)
2002
2003     if num_series > 1:
2004         ax.legend()
2005
2006     ax.grid(axis="y", linestyle="--", alpha=0.7)
2007     plt.setp(ax.get_xticklabels(), rotation=30, ha="right")
2008
2009 def main():
2010     chart_package_str = """
2011     "chart_type": "bar",
2012     "title": "Top 10 Most Populous Countries in 2023: Population, Density, and Growth
2013     Patterns",
2014     "x_label": "Country",
2015     "y_label": "Population (millions)",
2016     "data": "Omitted here; please refer to the example provided above",
2017     "topic": "Top 10 Most Populous Countries in 2023",
2018     "domain": "Demographics",
2019     "figsize": [
2020         12,
2021         8
2022     ],
2023     "dpi": 196,
2024     "renderer": "matplotlib",
2025     "task_type": "single",
2026     "complexity_steps": 1
2027 }"""
2028 chart_package = json.loads(chart_package_str)
2029
2030 chart_type = chart_package.get("chart_type")
2031 title = chart_package.get("title", "")
2032 figsize = (12, 8)
2033 fig, ax = None, None
2034
2035 try:
2036     try:
2037         plt.rcParams["font.sans-serif"] = ["SimHei", "Microsoft YaHei", "Arial
2038         Unicode MS", "sans-serif"]
2039     except:
2040         pass
2041     plt.rcParams["axes.unicode_minus"] = False
2042
2043     is_special_fig = chart_type in ["radar", "rose", "3d_bar", "contour", "calendar_
2044     heatmap"]
2045     if is_special_fig:
2046         fig = plt.figure(figsize=figsize)
2047     else:
2048         fig, ax = plt.subplots(figsize=figsize)
2049
2050     draw_funcs = {
2051         "bar": _draw_bar_chart,
2052     }
2053     draw_func = draw_funcs.get(chart_type)
2054     if not draw_func: raise ValueError(f"Unsupported chart type: {chart_type}")

```

```

2052
2053     if is_special_fig:
2054         ax = draw_func(fig, chart_package)
2055     else:
2056         draw_func(ax, chart_package)
2057
2058     if chart_type != "calendar_heatmap":
2059         if ax and not isinstance(ax, Axes3D): ax.set_title(title, fontsize=16)
2060         elif fig: fig.suptitle(title, fontsize=16)
2061
2062     if ax and chart_type not in ["pie", "donut", "radar", "rose", "heatmap", "treemap", "
2063     network", "3d_bar", "sankey", "gauge", "word_cloud", "parallel_coordinates"]:
2064         ax.set_xlabel(chart_package.get("x_label", ""))
2065         ax.set_ylabel(chart_package.get("y_label", ""))
2066
2067     fig.tight_layout(rect=[0, 0, 1, 0.96] if title else None)
2068     plt.show()
2069
2070     except Exception as e:
2071         print(f"Error generating chart \"{chart_type}\": {e}")
2072     finally:
2073         if "fig" in locals() and fig: plt.close(fig)
2074
2075 if __name__ == "__main__":
2076     main()

```

#### LLM-generated QA pairs based on the chart code

```

2077
2078
2079
2080 {
2081     'question': 'Which country has the highest population among the top 10 most populous
2082     nations?',
2083     'answer': 'India',
2084     'negative answer': 'China', 'United States', 'Indonesia', 'Pakistan'
2085 }
2086

```

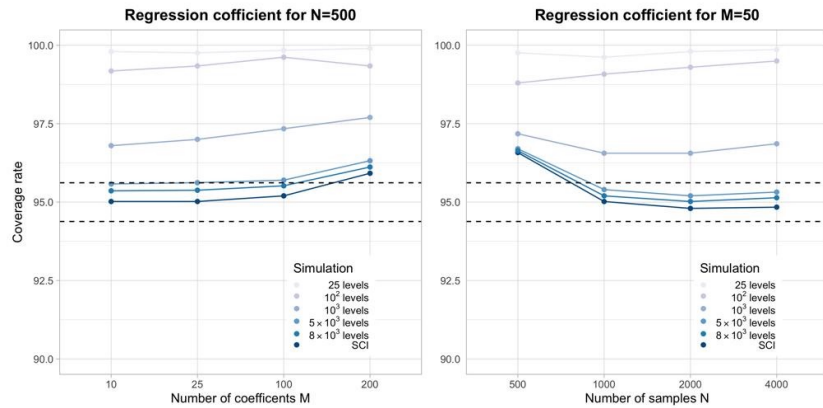
Figure 19: LLM-generated QA pairs based on the chart code and the negative examples.

## H FAILURE CASE ANALYSIS

In Figure 20, we present representative failure cases of ChartMaster. This example reveals several fundamental limitations of the current models. Specifically, successful prediction in this scenario requires the model to first comprehend the global structure of the chart and accurately localize the relevant subfigure. The model must then interpret the correspondence between the legend entries and the plotted lines, which involves recognizing color-series associations and understanding chart-specific conventions. In addition, the task demands a certain level of spatial alignment capability, where the model must infer the precise positions of data points relative to the coordinate axes in order to derive the correct answer.

These observations underscore the substantial challenges posed by fine-grained visual perception, semantic grounding of chart elements, and geometric reasoning in chart understanding tasks. Such capabilities are essential for handling complex charts, and the identified failure patterns point to important directions for future research in advancing chart comprehension models.

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

**Question:**

At what number of samples  $N$  does  $5 \times 10^3$  levels reach the lowest coverage rate for regression coefficient for  $M=50$ ?

GT: 2000

**ChartMaster:**

<think>\n\nTo determine the number of samples  $(N)$  at which the  $5 \times 10^3$  levels reaches the lowest coverage rate for "Regression coefficient for  $M=50$ ," let's closely examine the graph.\n\n1. The graph on the right represents "Regression coefficient for  $M=50$ ."\n\n2. The legend indicates the number of levels for each data point.\n\n3. We need to focus on the  $5 \times 10^3$  levels data points.\n\nBy closely inspecting the data points for  $5 \times 10^3$  levels, we can see the trend across different values of  $(N)$ .\n\nFrom the graph, it is evident that the  $5 \times 10^3$  levels reaches its lowest coverage rate at  $(N = 1000)$ .\n\n</think>\n\n\boxed{1000}

Figure 20: Illustration of typical failure cases encountered by ChartMaster, highlighting scenarios where the model struggles with fine-grained table semantics and complex reasoning patterns.