

# ChartVerse: Scaling Chart Reasoning via Reliable Programmatic Synthesis from Scratch

Zheng Liu<sup>1,2,\*</sup>, Honglin Lin<sup>1,3,\*</sup>, Xiaoyang Wang<sup>1</sup>, Xin Gao<sup>1</sup>, Yu Li<sup>1</sup>, Mengzhang Cai<sup>1</sup>,  
Yun Zhu<sup>1</sup>, Zhanping Zhong<sup>1</sup>, Qizhi Pei<sup>1</sup>, Zhuoshi Pan<sup>1</sup>, Xiaoran Shang<sup>1</sup>,  
Conghui He<sup>1</sup>, Bin Cui<sup>2,5</sup>, Wentao Zhang<sup>2,4,5,†</sup>, Lijun Wu<sup>1,†</sup>

<sup>1</sup>Shanghai AI Laboratory, <sup>2</sup>Peking University,

<sup>3</sup>Shanghai Jiao Tong University, <sup>4</sup>Zhongguancun Academy,

<sup>5</sup>Beijing Key Laboratory of Software and Hardware Cooperative Artificial Intelligence Systems

## Abstract

Chart reasoning is a critical capability for Vision Language Models (VLMs). However, the development of open-source models is severely hindered by the lack of high-quality training data. Existing datasets suffer from a dual challenge: synthetic charts are often simplistic and repetitive, while the associated QA pairs are prone to hallucinations and lack the reasoning depth required for complex tasks. To bridge this gap, we propose **ChartVerse**, a scalable framework designed to synthesize complex charts and reliable reasoning data from scratch. (1) To address the bottleneck of simple patterns, we first introduce *Rollout Posterior Entropy (RPE)*, a novel metric that quantifies chart complexity. Guided by RPE, we develop *complexity-aware chart coder* to autonomously synthesize diverse, high-complexity charts via executable programs. (2) To guarantee reasoning rigor, we develop *truth-anchored inverse QA synthesis*. Diverging from standard generation, we adopt an answer-first paradigm: we extract deterministic answers directly from the source code, generate questions conditional on these anchors, and enforce strict consistency verification. To further elevate difficulty and reasoning depth, we filter samples based on model fail-rate and distill high-quality Chain-of-Thought (CoT) reasoning. We curate ChartVerse-SFT-600K and ChartVerse-RL-40K using Qwen3-VL-30B-A3B-Thinking as the teacher. Experimental results demonstrate that ChartVerse-8B achieves state-of-the-art performance, notably surpassing its teacher and rivaling the stronger Qwen3-VL-32B-Thinking. We release our code, model weights, and datasets in <https://chartverse.github.io>.

## 1 Introduction

Recent Vision Language Models (VLMs) (Bai et al., 2025; Zhang et al., 2025; OpenAI, 2025;

Google DeepMind, 2025) exhibit strong reasoning capabilities, yet their performance on chart reasoning remains far from satisfactory. A central limitation lies in the lack of high-quality data that simultaneously achieves scale, diversity, and complexity. For example, manually curated datasets (Masry et al., 2022; Methani et al., 2020) provide realistic chart distributions but are severely limited by annotation cost. Synthetic approaches—ranging from template-based rendering (Kafle et al., 2018; Han et al., 2023) to recent code-driven pipelines (He et al., 2025; Zhao et al., 2025; Kondic et al., 2025)—offer scalability but typically rely on fixed templates or seed charts, resulting in limited visual diversity and poor coverage of real-world, long-tail chart patterns. Although proprietary systems (Yang et al., 2025a) can produce visually rich charts, their dependence on closed commercial APIs prevents large-scale, open-source adoption.

Beyond visual diversity, the quality of reasoning supervision presents a second major challenge. Existing chart Question-Answering (QA) datasets are constrained in both difficulty and reliability. Visually simplistic charts inherently limit question complexity, making it difficult to construct multi-step reasoning problems. Moreover, answer correctness is also hard to guarantee at scale: manual verification (Masry et al., 2022) is infeasible, while LLM-generated QA pairs (Xu et al., 2025a; Yang et al., 2025a; Chen et al., 2025) frequently suffer from hallucinations and numerical errors. As a result, current pipelines struggle to provide large-scale, high-fidelity chart reasoning data with rigorous and error-free supervision for model training.

To address the above challenges, we propose **ChartVerse**, a scalable, code-driven framework for chart reasoning data synthesis. ChartVerse is designed to jointly improve visual complexity, distributional diversity, and reasoning reliability, enabling the construction of large-scale, challenging chart QA data without relying on fixed templates

\* Equal contribution.

† Corresponding authors.

Table 1: Comparison of ChartVerse-SFT-600K with existing chart datasets.

Datasets	Chart Properties			Q&A Properties				Complexity & Diversity		
	Chart Count	Data Source	Textual Data	QA Count	Total Tokens	Answer Accuracy	Reasoning Data	Rollout Posterior Entropy	Color Entropy	Semantic Emb Spread
ChartQA (Masry et al., 2022)	18K	Real-world, Synthetic	Table	28K	100K	✗	✗	0.26	2.03	0.37
PlotQA (Methani et al., 2020)	156K	Real-world, Synthetic	Table	20M	117M	✗	✗	0.21	0.83	0.30
FigureQA (Kahou et al., 2018)	100K	Synthetic	-	1.3M	2.6M	✗	✗	0.25	1.04	0.24
ReachQA (He et al., 2025)	3K	Synthetic	Code	20K	1.2M	✗	✓	0.29	1.91	0.47
ECD (Yang et al., 2025b)	10K	Synthetic	Code	321K	18.6M	✗	✓	0.31	2.13	0.40
CoSyn (Yang et al., 2025a)	116K	Synthetic	Code	1.1M	5.4M	✗	✗	0.35	1.73	0.54
START (Liu et al., 2025e)	400K	Synthetic	Code	400K	143M	✗	✗	0.33	1.63	0.27
ChartGen (Kondic et al., 2025)	216K	Synthetic	Code	-	-	✗	✗	0.30	1.31	0.39
ChartCoder (Zhao et al., 2025)	163K	Synthetic	Code	-	-	✗	✗	0.33	1.48	0.38
ChartVerse-SFT-600K(Ours)	412K	Synthetic	Code	603K	3.9B	✓	✓	0.44	3.17	0.51

or manual verification.

We first introduce **Rollout Posterior Entropy (RPE)**, a metric for estimating intrinsic chart difficulty. Given a chart image, we repeatedly prompt an existing VLM to parse the image into executable chart code and measure the consistency across multiple rollouts. Charts that admit diverse, inconsistent code reconstructions exhibit higher posterior entropy, indicating greater visual and semantic ambiguity. RPE thus provides an objective, model-agnostic signal for identifying challenging charts and guiding data selection.

Using RPE, we first filter existing chart corpora to construct a high-difficulty seed dataset, which is then used to train a **complexity-aware chart coder**. Unlike prior template-based or seed-conditioned approaches, our chart coder generates chart code from scratch via high-temperature sampling, enabling broad exploration of the long-tail chart distribution. Through an iterative loop of generation, difficulty filtering, and retraining, the chart coder progressively improves both code quality and visual complexity, producing diverse, realistic charts.

Further, to ensure both correctness and reasoning difficulty for QA, ChartVerse adopts a **truth-anchored inverse QA synthesis** pipeline. Ground-truth answers are first computed directly from chart code, guaranteeing numerical and semantic accuracy. Questions are then reverse-synthesized conditioned on these answers. To control difficulty, we evaluate the resulting QA pairs using a strong VLM and retain samples with non-trivial fail-rate, ensuring that the final dataset emphasizes genuinely challenging reasoning cases.

Leveraging Qwen3-VL-30B-A3B-Thinking as the teacher model, we finally distill high-quality chain-of-thought (CoT) supervision and curate two datasets: ChartVerse-SFT-600K and ChartVerse-RL-40K. Initializing from the Qwen3-VL Instruct

series (Bai et al., 2025), we train ChartVerse models at 2B, 4B, and 8B scales. Despite their smaller sizes, these models achieve state-of-the-art (SOTA) performance on chart reasoning benchmarks. Notably, ChartVerse-4B surpasses Qwen3-VL-8B-Thinking, while ChartVerse-8B even outperforms its 30B teacher, approaching the performance of Qwen3-VL-32B-Thinking.

## 2 Related Work

**Chart VLMs.** Leading proprietary models (OpenAI, 2025; Google DeepMind, 2025; Anthropic, 2025; Bai et al., 2025) have demonstrated exceptional reasoning capabilities in complex chart understanding. In the open-source domain, early methods like ChartLlama (Han et al., 2023), ChartGemma (Masry et al., 2024b), and EvoChart (Huang et al., 2025) primarily focused on visual instruction tuning to handle fundamental structural extraction tasks like OCR and element counting. Recent works have transitioned towards reasoning-centric designs. TinyChart (Zhang et al., 2024a) employs template-based Program-of-Thought synthesis. ChartReasoner (Jia et al., 2025) and Chart-R1 (Chen et al., 2025) integrate reinforcement learning with supervised finetuning, directly incentivizing the model to generate robust CoT paths for complex problem-solving.

**Chart Image Synthesis.** Conventional approaches (Kafle et al., 2018; Han et al., 2023) rely on rule-based rendering engines to visualize data tables. These methods often yield charts with rigid layouts and limited diversity. Recent works (Zhao et al., 2025; He et al., 2025; Xu et al., 2025a; Jiang and Luo, 2025; Kondic et al., 2025; Yang et al., 2025b) leverage LLMs to generate executable code or augment fixed visual seeds. Despite offering better flexibility, these methods remain constrained by finite and simple templates. The charts exhibit

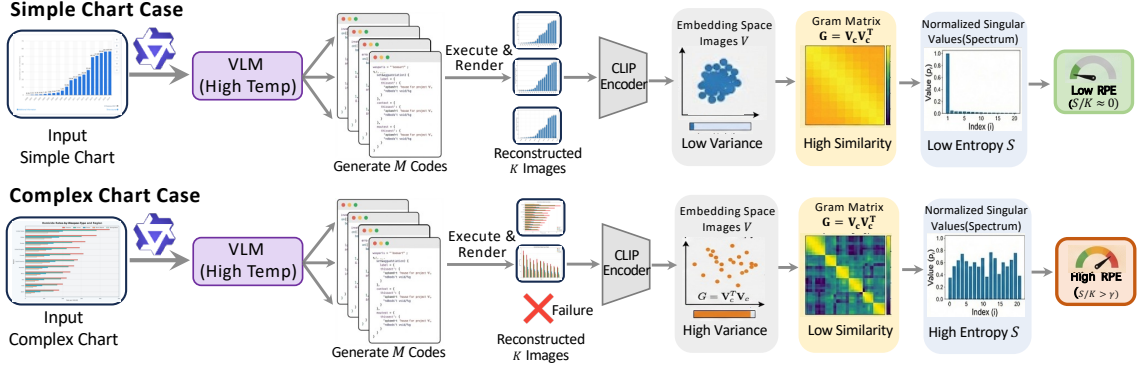


Figure 1: Illustration of Rollout Posterior Entropy. We quantify complexity via generative stability: simple charts yield consistent reconstructions (low RPE), while complex charts result in divergent outcomes (high RPE).

trivial topologies that are insufficient to elicit the complex reasoning capabilities of VLMs. Methods like Cosyn (Yang et al., 2025a) employ proprietary model Claude-3.5-Sonnet for end-to-end synthesis to avoid template restrictions. However, the dependence on commercial APIs introduces prohibitive costs and scalability bottlenecks.

**Chart QA Construction.** High-quality QA pairs is pivotal for chart understanding. While human annotation (Masry et al., 2022; Methani et al., 2020) guarantees accuracy, its unscalability has necessitated a shift towards automated generation. Recent approaches (Yang et al., 2025a; Chen et al., 2025; Masry et al., 2024a; Liu et al., 2024; Meng et al., 2024; Liu et al., 2025b,c; Liu et al.; Lin et al., 2026a; Liu et al., 2025d; Niu et al., 2025) typically prompt leading VLMs to synthesize QA pairs from images, with some works (He et al., 2025; Lin et al., 2026b) leveraging auxiliary code. Despite their popularity, these methods are constrained by the stochastic nature of LLMs, which leads to unavoidable hallucinations and inaccuracies. Lacking ground truth for verification, these methods predominantly yield uncalibrated, simplistic questions and lack multi-step CoT, accurate answers.

### 3 Methodology

In this section, we introduce our ChartVerse framework for scalable chart reasoning data synthesis. We first introduce a difficulty metric to quantify intrinsic chart complexity, which guides data selection and generation. Based on this metric, we develop a complexity-aware chart coder and a truth-anchored inverse QA synthesis pipeline to produce challenging and reliable training data. Figure 2 shows the complete process.

#### 3.1 Quantifying Intrinsic Chart Complexity

A prerequisite for training a robust chart generator is the identification and filtration of high-complexity samples. However, distinguishing complexity via direct pixel analysis is intractable, as visual density does not correlate with structural difficulty. When humans interpret a complex chart, they often form diverse understandings and struggle to capture every detail. This implies that the decoding space for complex visual data is variable, making consistent reproduction difficult.

Guided by this insight, we propose Rollout Posterior Entropy (RPE) to measure the variability. Specifically, we employ a VLM to generate executable code for charts as a proxy for structural interpretation. As illustrated in Figure 1, simple charts yield consistent code that renders into identical images, whereas complex charts induce high instability, resulting in divergent reconstructions.

**Modeling Structural Variability.** Given a chart, we utilize Qwen3-VL-2B-Thinking to generate plotting code 8 times with the temperature set to 1.0. We then execute these codes to render reconstructed images. Let  $K$  denote the number of successfully executed samples. We utilize CLIP (Radford et al., 2021) to extract  $1 \times d$  dimensional features for each image and stack them into an image matrix  $V \in \mathbf{R}^{K \times d}$ .

To quantify the variability, we first center  $V$  to remove translational shifts and compute the Gram matrix  $G$  to model pairwise similarities:

$$V_c = \left( I - \frac{1}{K} \mathbf{1}\mathbf{1}^T \right) V, \quad G = V_c V_c^T. \quad (1)$$

Next, we calculate the singular values  $\sigma_i$  of  $G$ . We normalize these values and compute the spec-

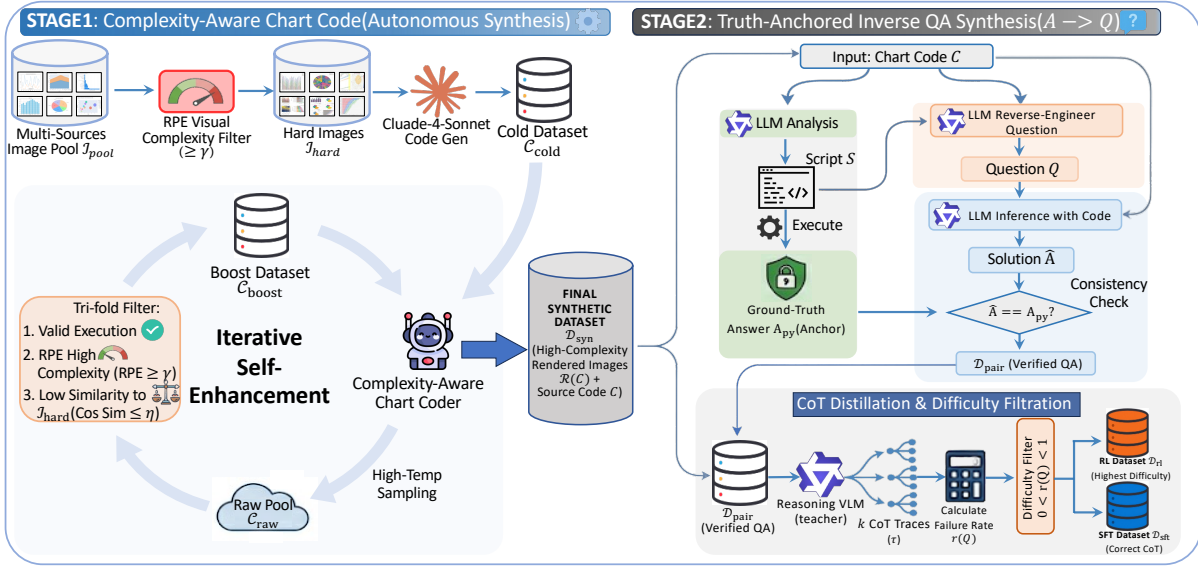


Figure 2: Overview of ChartVerse. Stage 1 trains a complexity-aware chart coder via RPE-guided iterative self-enhancement to synthesize diverse, high-complexity chart codes. Stage 2 generates verifiable QA pairs using truth-anchored inverse synthesis and CoT distillation, followed by failure-rate filtering to guarantee difficulty.

tral entropy to measure the degree of variation.

$$p_i = \sigma_i / \sum \sigma_j, \quad S = - \sum_{i=1}^K p_i \log p_i. \quad (2)$$

Since execution failures inherently signal difficulty, we normalize  $S$  by the valid count  $K$  to obtain final RPE:

$$RPE = \frac{S}{K}. \quad (3)$$

Higher  $RPE$  denotes greater reproduction difficulty, reflecting intrinsic chart complexity.

**Diagnosing Existing Datasets.** Leveraging RPE, we evaluate the complexity of mainstream datasets. As shown in Table 1, nearly all existing datasets exhibit low RPE scores, indicating that such charts are trivial. To drive further progress, it is critical to synthesize data that transcends simple visual variation, offering the complexity required to challenge the consistency of modern VLMs.

### 3.2 Complexity-Aware Chart Coder

RPE suggests that executable code is a high-fidelity proxy of chart structure, we therefore synthesize complex charts via code and propose a complexity-aware chart coder that can autonomously generate non-standard, high-complexity charts from scratch.

**Difficulty-Filtered Cold Start.** We first construct a high-quality cold-start seed code  $C_{cold}$ . We

aggregate chart images from existing datasets in Table 1 to form an image pool  $I_{pool}$  and apply RPE to quantify visual complexity. Based on the statistics in Table 1, we set an RPE threshold of 0.4 to retain only high RPE samples:

$$I_{hard} = \{x \in I_{pool} \mid RPE(x) \geq 0.4\}, \quad (4)$$

Since real-world charts in  $I_{hard}$  lack ground-truth source code, we employ Claude-4-Sonnet to infer candidate code  $C(x)$  for each image  $x \in I_{hard}$ . We discard samples that trigger execution errors, obtaining the cold-start training set:

$$C_{cold} = \{C(x) \mid x \in I_{hard}\}. \quad (5)$$

Ultimately, we curate a cold-start set  $C_{cold}$  comprising 60k high-quality samples, which serves as the foundational corpus for our chart coder.

Unlike existing methods that rely on code as rigid templates, we train a strong coder LLM, Qwen2.5-Coder-7B (Hui et al., 2024), to function as a flexible generative model. We structure each training instance by pairing a concise system instruction  $\mathcal{T}$  (see Prompt 12) as the model input with the chart code sequence  $C = (c_1, c_2, \dots, c_L) \in C_{cold}$  as the target answer. The model is optimized using the standard cross-entropy loss to internalize the structural and logic patterns of executable chart codes. During inference, by applying a high sampling temperature, the coder leverages its aligned internal knowledge and learned patterns to generate diverse chart codes. This approach enables

the model to move beyond simple replication and synthesize novel, complex chart code.

**Iterative Self-Enhancement Loop.** Despite its quality,  $\mathcal{C}_{cold}$  is typically scarce, and the coder trained on it exhibited limited structural complexity, making large-scale high-difficulty synthesis non-trivial. To address this, we introduce a self-enhancement mechanism that iteratively expands data and strengthens the coder capability.

Specifically, we first employ the cold-start coder to perform large-scale sampling at a temperature of 1.0, generating a raw candidate pool  $\mathcal{C}_{raw}$  of 2M code samples. To ensure quality and diversity, we discard codes that trigger execution errors, then render the charts and filter out low-complexity samples using RPE. To ensure diversity, we exclude charts that exhibit high visual similarity to  $\mathcal{I}_{hard}$  based on image embeddings from CLIP. We apply an RPE threshold  $\geq 0.4$  and a similarity limit  $\leq 0.65$  to retain only complex, non-redundant samples:

$$\mathcal{C}_{boost} = \{C \in \mathcal{C}_{raw} \mid RPE(\mathcal{R}(C)) \geq 0.4 \\ \wedge \text{Sim}(\mathcal{R}(C), \mathcal{I}_{hard}) \leq 0.65\}, \quad (6)$$

where  $\mathcal{R}(C)$  denotes the rendered image, and  $\text{Sim}(x, \mathcal{I}_{hard}) = \max_{x' \in \mathcal{I}_{hard}} \cos(e(x), e(x'))$  calculates the maximum cosine similarity between the CLIP embedding  $e(x)$  of a candidate and  $\mathcal{I}_{hard}$ .

Ultimately, we curate a boosted set  $\mathcal{C}_{boost}$  of 200k samples and retrain the coder on  $\mathcal{C}_{cold} \cup \mathcal{C}_{boost}$  to obtain a stronger model. We repeat this procedure for two iterations to ensure the model’s generative stability and structural diversity.

**Synthetic Chart Dataset.** We utilize the final complexity-aware coder to perform large-scale sampling at a temperature of 1.0, generating 1M raw chart code candidates  $\mathcal{C}_{syn}$ . We again apply RPE filtering to retain only high-complexity samples. This process yields a final dataset  $\mathcal{D}_{syn}$  comprising 700k pairs of rendered images and their corresponding executable codes:

$$\mathcal{D}_{syn} = \{(\mathcal{R}(C), C) \mid C \in \mathcal{C}_{syn} \\ \wedge RPE(\mathcal{R}(C)) \geq 0.4\}. \quad (7)$$

This extensive corpus of complex charts serves as the foundation for our QA generation.

### 3.3 Truth-Anchored Inverse QA Synthesis

While the complexity of  $\mathcal{D}_{syn}$  provides a foundation for constructing difficult questions, it introduces

significant challenges in accurate answer extraction. To generate QA pairs that preserve both *accuracy* and *difficulty*, we propose to invert the generation flow to enforce "Answer as the Source of Truth".

**Inverse Logic Construction.** Existing QA generation approaches, which first synthesize a question and then employ a VLM to derive the answer ( $Q \rightarrow A$ ), suffer from stochastic hallucinations and cannot guarantee correctness. To address this, we propose an Inverse Synthesis paradigm ( $A \rightarrow Q$ ).

An advantage of  $\mathcal{D}_{syn}$  is that the images are rendered from code, ensuring that codes encapsulate the absolute ground truth of all visual data. Leveraging this, we synthesize the answer first. As shown in Figure 2, to ensure absolute precision, we employ Qwen3-30B-A3B-Thinking  $\pi_{LLM}$  to analyze the chart code  $C$  and generate a Python script  $S$  that performs meaningful operations on the data in chart code, yielding a precise numerical or categorical result. By executing this script in a deterministic environment  $\mathcal{E}$ , we obtain the ground-truth answer  $A_{py}$ , completely bypassing numerical errors from direct VLM answering:

$$S \sim \pi_{LLM}(C), \quad A_{py} = \mathcal{E}(S). \quad (8)$$

We then instruct the LLM again to reverse-engineer a question  $Q$  that strictly aligns with chart code  $C$  and Script  $S$ :

$$Q \sim \pi_{LLM}(C, S). \quad (9)$$

To guarantee that the synthesized question is unambiguously solvable, we perform a consistency check. We feed the question  $Q$  and chart code  $C$  back into the LLM to produce a solution  $\hat{A}$ ,

$$\hat{A} \sim \pi_{LLM}(C, Q). \quad (10)$$

We strictly retain samples where the inferred solution matches the programmatic ground truth:

$$\mathcal{D}_{pair} = \{(\mathcal{R}(C), Q, A_{py}) \mid \hat{A} = A_{py}\}. \quad (11)$$

Through this rigorous inverse pipeline, we ensure that every sample is logically sound and factually accurate. It is worth noting that this logical construction phase operates exclusively within the textual domain. By relying solely on the code-text modality, we achieve high-efficiency synthesis while eliminating potential visual encoding errors.

Table 2: Comparison of different models. The shaded rows denote our ChartVerse models.

Model	ChartQA Pro	CharXiv (RQ)	CharXiv (DQ)	Chart Museum	ChartX	Evo Chart	ChartBench (GPT-acc)	Avg
ECD-7B (Yang et al., 2025b)	45.1	41.7	74.9	24.5	59.3	56.0	48.6	50.0
START-7B (Liu et al., 2025e)	43.5	46.3	76.8	29.7	57.3	63.8	50.0	52.5
Chart-R1-7B (Chen et al., 2025)	45.6	47.7	70.0	33.4	60.9	67.1	50.5	53.6
ChartVerse-2B	48.2	46.9	71.2	37.5	60.5	66.8	49.1	54.3
InternVL3.5-38B (Wang et al., 2025)	47.8	45.5	86.5	34.2	55.5	65.0	49.0	54.8
InternVL3.5-241B-A28B	50.7	47.5	88.1	36.1	59.1	67.4	48.6	56.8
Qwen3-VL-8B-Thinking (Bai et al., 2025)	53.9	53.0	85.9	44.3	59.6	74.1	49.1	60.0
ChartVerse-4B	55.2	56.2	84.1	45.9	63.7	75.0	52.9	61.9
Qwen3-VL-30B-A3B-Thinking	55.8	56.6	86.9	49.2	62.3	77.2	52.4	62.9
ChartVerse-8B	56.2	60.8	88.0	49.2	63.9	76.2	54.2	64.1
Qwen3-VL-32B-Thinking	58.8	65.2	90.2	55.9	64.1	80.8	54.3	67.0
Qwen3-VL-235B-A30B-Thinking	60.0	66.1	90.5	60.0	64.5	79.9	54.5	67.9

**CoT Distillation & Difficulty Filtration.** Correctness alone does not imply training value. We further estimate visual reasoning difficulty via repeated VLM rollouts and keep only samples that are neither trivial nor impossible. We employ Qwen3-VL-30B-A3B-Thinking  $\pi_{VLM}$  to generate 3 CoT reasoning traces  $\{\tau_1, \tau_2, \tau_3\}$ .

As shown in Figure 2, we extract the predicted answer  $\hat{a}_j$  from each trace  $\tau_j$  and calculate the failure rate  $r(Q)$  against the ground truth:

$$r(Q) = 1 - \frac{1}{3} \sum_{j=1}^3 \text{Match}(\hat{a}_j, A_{py}). \quad (12)$$

We construct  $\mathcal{D}_{\text{hard}}$  by retaining samples where  $r(Q)$  is neither 0 nor 1. This ensures that samples are neither impossible to solve nor trivially simple.

$$\mathcal{D}_{\text{hard}} = \{(\mathcal{R}(C), Q, A_{py}, \tau) \mid 0 < r(Q) < 1\}. \quad (13)$$

We prioritize the most challenging samples with the highest  $r(Q)$  in  $\mathcal{D}_{\text{hard}}$  for the Reinforcement Learning dataset  $\mathcal{D}_{\text{rl}}$ , while utilizing the remaining samples enriched with correct CoT traces  $\tau$  for the Supervised Finetuning dataset  $\mathcal{D}_{\text{sft}}$ . After the processing, we finally obtain the ChartVerse-SFT-600K and ChartVerse-RL-40K datasets. Details and prompts on synthesis pipelines and coder training are available in Appendix C.

## 4 Experiments

### 4.1 Experimental Setup

**Training Details.** All ChartVerse models are initialized from the Qwen3-VL-Instruct (Bai et al., 2025) series. We first perform Supervised Fine-Tuning (SFT) on the ChartVerse-SFT-600K dataset using LLaMA-Factory (Zheng et al., 2024). Subsequently, Reinforcement Learning (RL) is

conducted on the ChartVerse-RL-40K dataset with VeRL (Sheng et al., 2025), adopting the GSPO (Zheng et al., 2025) algorithm. Comprehensive training configurations and hyperparameters are provided in Appendix A.

**Baselines.** We compare our tuned ChartVerse models against two categories of models: (1) General open-source VLMs: leading models including Qwen3-VL-Thinking series (8B, 30B, 32B, 235B) (Bai et al., 2025), InternVL3.5 series (38B, 241B) (Wang et al., 2025). (2) Specialized chart-domain VLMs: ECD (Yang et al., 2025b), START (Liu et al., 2025e), and Chart-R1 (Chen et al., 2025).

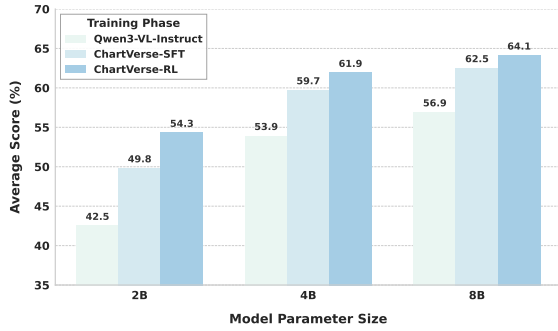
**Benchmarks.** We evaluate performance across 6 benchmarks necessitating complex chart understanding and reasoning: ChartQA-Pro (Masry et al., 2025), CharXiv (Wang et al., 2024), ChartMuseum (Tang et al., 2025), ChartX (Xia et al., 2025), ChartBench (Xu et al., 2024), and EvoChart (Huang et al., 2025).

### 4.2 Main Results

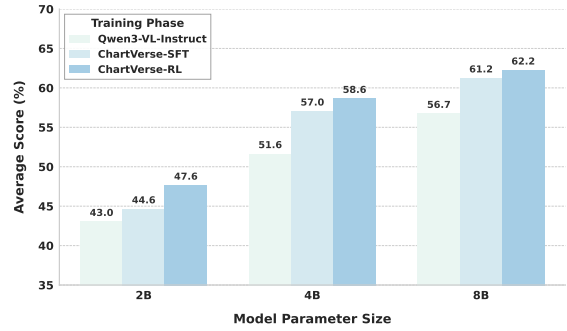
**Overall Results.** Table 2 compares ChartVerse models with above baselines. Overall, ChartVerse consistently delivers strong performance across all model scales, demonstrating the effectiveness of our difficulty-aware data synthesis framework. We highlight three key observations:

**(1) Competitive performance at small scale.** ChartVerse-2B achieves an average score of 54.3, exceeding all chart-specific baselines, including ECD-7B, START-7B, and Chart-R1-7B. This shows that complexity-controlled chart data can substantially offset model size limitations.

**(2) Data quality over model scale.** ChartVerse-4B attains an average score of 61.9, outperform-

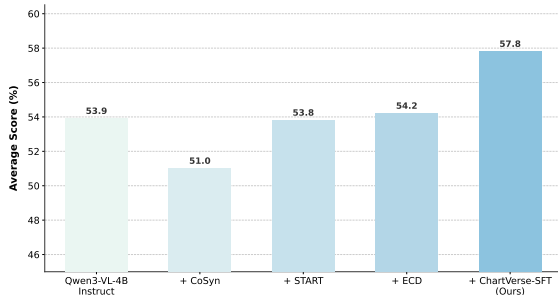


(a) Performance on chart-related benchmarks

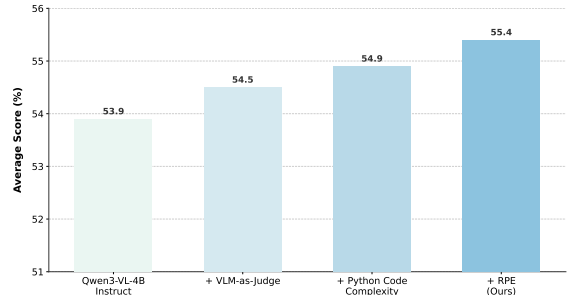


(b) Performance on STEM-related benchmarks

Figure 3: Performance evolution of ChartVerse across different training phases on chart and STEM benchmarks.



(a) Comparison of different datasets



(b) Ablation study on RPE strategy.

Figure 4: Ablation study results on different datasets and the proposed RPE strategy.

ing Qwen3-VL-8B-Thinking (60.0) despite using only half the parameters. This gap highlights the dominant role of data quality in improving chart reasoning performance.

**(3) Beyond the teacher model.** ChartVerse-8B further improves the average score to 64.1, surpassing its teacher Qwen3-VL-30B-A3B-Thinking (62.9) and approaching Qwen3-VL-32B-Thinking (67.0). This result indicates that ChartVerse enables student models to exceed the distillation ceiling.

Collectively, these observations validate the advantages of ChartVerse’s rigorously controlled chart and QA synthesis pipeline, showing that scalable generation of training data translates into consistent gains across model scales.

### 4.3 Training Stages and Generalization

**SFT and RL are Both Effective.** Figure 3a reports results on chart-related benchmarks under different training stages. Across all model sizes, ChartVerse-SFT yields substantial improvements over the Qwen3-VL-Instruct baseline, and ChartVerse-RL further provides consistent gains. For instance, ChartVerse-2B improves from 42.5 to 49.8 after SFT and reaches 54.3 after RL, while ChartVerse-8B increases from 56.9 to 62.5 and

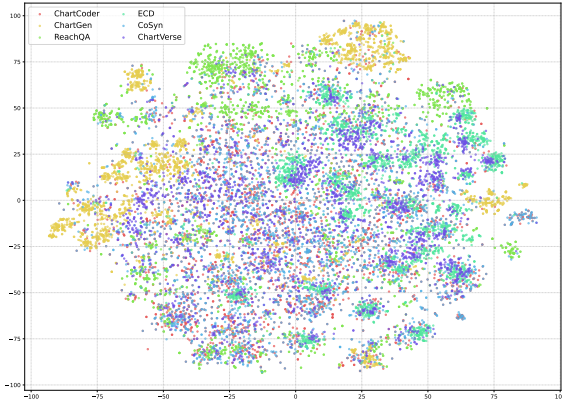
64.1, respectively. These results indicate that SFT establishes strong chart reasoning foundations, and RL further enhances performance by focusing on more challenging samples.

**Strong Generalization to STEM tasks.** We further evaluate ChartVerse on STEM-related benchmarks, including MathVista (Lu et al., 2023), DynaMath (Zou et al., 2024), MathVerse (Zhang et al., 2024b), LogicVista (Xiao et al., 2024), and VisuLogic (Xu et al., 2025b). As shown in Figure 3b, ChartVerse-trained models outperform the Instruct baseline across all scales. Notably, ChartVerse-8B improves from 56.7 to 61.2 after SFT and reaches 62.2 with RL, demonstrating that the reasoning skills learned from ChartVerse data transfer effectively to out-of-domain STEM reasoning tasks.

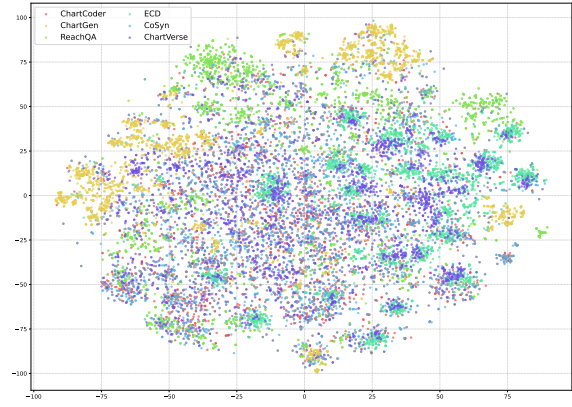
## 5 Ablations and Analysis

### 5.1 Fair Comparison with Existing Datasets

To control for data scale, we compare different training datasets under a fixed budget of 100K samples. As summarized by the Avg scores in Table 4a, most existing datasets provide little benefit to Qwen3-VL-4B-Instruct: CoSyn even degrades performance (53.9  $\rightarrow$  51.0), while START and ECD yield only marginal gains (53.8 and 54.2,



(a) Visualization of chart image embeddings



(b) Visualization of chart code embeddings

Figure 5: t-SNE visualizations of feature distributions for ChartVerse-SFT and existing datasets.

Table 3: Performance comparison of different QA synthesis strategies. All models are trained on Qwen3-VL-4B-Instruct with 100k synthesized samples.

Training Data (100k)	ChartQA Pro	CharXiv (RQ)	CharXiv (DQ)	Chart Museum	ChartX	Evo Chart	ChartBench (GPT-acc)	Avg
Qwen3-VL-4B-Instruct	53.7	39.7	76.2	37.2	57.2	68.2	45.1	53.9
+ Image-Space Gen.	53.1	46.6	81.3	38.4	58.4	70.4	47.3	56.5
+ Code-Space Gen.	53.9	47.0	81.2	38.6	58.9	70.3	47.8	56.8
+ Truth-Anchored Gen. w/o Failrate	54.0	47.4	81.7	39.8	59.7	70.7	48.2	57.4
+ Truth-Anchored Gen. w/ Failrate (Ours)	53.4	48.8	82.6	40.3	60.1	70.9	48.7	57.8

respectively), barely surpassing the baseline. In contrast, ChartVerse-SFT leads to a substantial improvement, boosting the average score to 57.8 (+3.9). This clear gap under identical data budgets indicates that the performance gains stem from data quality and task structure rather than scale, highlighting ChartVerse-SFT as significantly more effective for enhancing strong modern VLMs.

## 5.2 Effectiveness of Rollout Posterior Entropy

To verify the effectiveness of RPE, we compare RPE with two representative data selection strategies: *VLM-as-Judge*, which selects samples with high model-assessed visual complexity, and *Python Code Complexity*, which favors charts with dense Matplotlib structures. All strategies select 100K samples from the same data pool to ensure a fair comparison.

As shown in Figure 4b, RPE identifies substantially harder samples, as reflected by the highest failure rate of Qwen3-VL-4B-Instruct (27.6%), compared to 21.1% for VLM-as-Judge and 23.5% for Python Code Complexity. This suggests that RPE goes beyond surface-level visual or structural heuristics and captures samples that are intrinsically more challenging for the model. Crucially, this increased difficulty translates into better

downstream performance. Fine-tuning on RPE-selected data achieves the highest average score (55.4), outperforming both competing strategies. In contrast, methods based on perceived complexity yield smaller gains despite selecting easier samples. These results indicate that RPE is more effective at selecting *learning-relevant* hard examples, leading to stronger performance improvements under a fixed data budget.

## 5.3 Diversity and Distribution Analysis

To analyze the quality and diversity of ChartVerse-SFT, we compare it with representative chart datasets from multiple perspectives, including data scale, supervision signals, and distributional diversity. For feature-level analysis, we extract embeddings from chart code and images using Qwen3-8B-Embedding and CLIP, and visualize their distributions with t-SNE.

As illustrated in Figure 5, ChartVerse-SFT covers a substantially broader feature space than prior datasets, largely subsuming the distributions of existing chart corpora. This indicates that our dataset is not concentrated around a narrow set of visual or structural patterns, but instead spans diverse chart styles and semantics.

We further quantify diversity using complemen-

tary metrics. Color Distribution Entropy measures the richness of visual appearance, while Semantic Embedding Spread captures structural and semantic variation across samples. As reported in Table 1, ChartVerse-SFT consistently attains the highest values on both metrics, exceeding all existing datasets by a clear margin. Together with its higher rollout posterior entropy, these results suggest that ChartVerse-SFT is not only larger in scale, but also more diverse and information-rich, offering broader coverage of chart layouts, visual styles, and semantic structures.

#### 5.4 Breaking the Distillation Ceiling

A central goal of ChartVerse is to enable the student model to outperform its teacher, which we achieve through truth-anchored inverse QA synthesis. We compare this method with two common QA construction strategies—direct image-based generation and direct code-based generation—and include an ablation that removes fail-rate based selection.

As shown in Table 3, all QA synthesis strategies improve upon the baseline, indicating that synthetic supervision is generally beneficial. However, the magnitude of improvement varies substantially across methods. Direct image-space generation yields only moderate gains (Avg 56.5), while code-space generation performs slightly better (56.8), suggesting that structural cues from code are helpful but still limited.

In contrast, truth-anchored generation leads to a clear performance jump (57.4), demonstrating the advantage of anchoring QA synthesis to verifiable ground truth rather than surface-level patterns. Importantly, incorporating failure-rate-based selection further boosts the average score to 57.8, achieving the best performance. This additional gain highlights the role of hard-sample mining in maximizing the effectiveness of synthesized data, and confirms that our strategy produces not only more accurate but also more effective supervision.

## 6 Conclusion

We introduce ChartVerse, a scalable framework addressing the scarcity of complex data for chart reasoning. We introduce Rollout Posterior Entropy to quantify chart complexity, guiding Complexity-Aware Chart Coder to synthesize diverse, non-trivial charts. We propose Truth-Anchored Inverse QA Synthesis, ensuring accuracy via code-derived ground truths, reverse-synthesized questions, and

consistency checks. Trained on our ChartVerse-SFT-600k and ChartVerse-RL-40k, ChartVerse-8B surpasses its teacher model Qwen3-VL-30B-A3B-Thinking and approaches Qwen3-VL-32B-Thinking, demonstrating the superior quality of our data.

## Limitations

**Computational Overhead of RPE.** While Rollout Posterior Entropy effectively quantifies complexity, its reliance on large-scale sampling entails computational costs. Despite mitigation via vLLM, developing lightweight approximation methods to reduce resource demands remains a valuable direction for future optimization.

**Scope of Answer Types.** Our Truth-Anchored Inverse QA prioritizes rigorous correctness through code execution, resulting in datasets focused primarily on precise numerical and logical reasoning. Future work aims to extend this pipeline to include qualitative visual interpretations while preserving the same standards of verification.

**Teacher Model Capacity.** Our current synthesis employs Qwen3-VL-30B-A3B-Thinking. While effective, a capability gap remains compared to frontier models. Leveraging significantly larger models, such as Qwen3-VL-235B-A30B-Thinking, could further enhance the reasoning depth and complexity of the generated data.

## Acknowledgments

This work is supported by Shanghai Artificial Intelligence Laboratory, Fundamental and Interdisciplinary Disciplines Breakthrough Plan of the Ministry of Education of China (JYB2025XDXM113), National Natural Science Foundation of China (92470121, 62402016), National Key R&D Program of China (2024YFA1014003), Zhongguncun Academy (C20250204, C20250602), Beijing Major Science and Technology Project (Z251100008125043, Z251100008425023), and High-performance Computing Platform of Peking University.

## References

- Anthropic. 2025. [Claude 3.5 technical report](#). Technical report, Anthropic.
- Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhi-fang Guo, Qidong Huang, Jie Huang, Fei Huang,

- Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025. [Qwen3-vl technical report](#). *Preprint*, arXiv:2511.21631.
- Mengzhang Cai, Xin Gao, Yu Li, Honglin Lin, Zheng Liu, Zhuoshi Pan, Qizhi Pei, Xiaoran Shang, Mengyuan Sun, Zinan Tang, Xiaoyang Wang, Zhanping Zhong, Yun Zhu, Dahua Lin, Conghui He, and Lijun Wu. 2025. [Opendataarena: A fair and open arena for benchmarking post-training dataset value](#). *Preprint*, arXiv:2512.14051.
- Lei Chen, Xuanle Zhao, Zhixiong Zeng, Jing Huang, Yufeng Zhong, and Lin Ma. 2025. [Chart-r1: Chain-of-thought supervision and reinforcement for advanced chart reasoner](#). *Preprint*, arXiv:2507.15509.
- Haodong Duan, Xinyu Fang, Junming Yang, Xiangyu Zhao, Yuxuan Qiao, Mo Li, Amit Agarwal, Zhe Chen, Lin Chen, Yuan Liu, Yubo Ma, Hailong Sun, Yifan Zhang, Shiyin Lu, Tack Hwa Wong, Weiyun Wang, Peiheng Zhou, Xiaozhe Li, Chaoyou Fu, and 13 others. 2025. [Vlmevalkit: An open-source toolkit for evaluating large multi-modality models](#). *Preprint*, arXiv:2407.11691.
- Google DeepMind. 2025. [Gemini 3 technical report](#). Technical report, Google.
- Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. 2023. [Chartllama: A multimodal llm for chart understanding and generation](#). *Preprint*, arXiv:2311.16483.
- Wei He, Zhiheng Xi, Wanxu Zhao, Xiaoran Fan, Yiwen Ding, Zifei Shan, Tao Gui, Qi Zhang, and Xuanjing Huang. 2025. [Distill visual chart reasoning ability from llms to mllms](#). *Preprint*, arXiv:2410.18798.
- Muye Huang, Han Lai, Xinyu Zhang, Wenjun Wu, Jie Ma, Lingling Zhang, and Jun Liu. 2025. [Evochart: A benchmark and a self-training approach towards real-world chart understanding](#). *Preprint*, arXiv:2409.01577.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, Kai Dang, Yang Fan, Yichang Zhang, An Yang, Rui Men, Fei Huang, Bo Zheng, Yibo Miao, Shanghaoran Quan, and 5 others. 2024. [Qwen2.5-coder technical report](#). *Preprint*, arXiv:2409.12186.
- Caijun Jia, Nan Xu, Jingxuan Wei, Qingli Wang, Lei Wang, Bihui Yu, and Junnan Zhu. 2025. [Chartreasoner: Code-driven modality bridging for long-chain reasoning in chart question answering](#). *Preprint*, arXiv:2506.10116.
- Gongyao Jiang and Qiong Luo. 2025. [Chart-coca: Self-improving chart understanding of vision lms via code-driven synthesis and candidate-conditioned answering](#). *Preprint*, arXiv:2508.11975.
- Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. [Dvqa: Understanding data visualizations via question answering](#). *Preprint*, arXiv:1801.08163.
- Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Akos Kadar, Adam Trischler, and Yoshua Bengio. 2018. [Figureqa: An annotated figure dataset for visual reasoning](#). *Preprint*, arXiv:1710.07300.
- Jovana Kondic, Pengyuan Li, Dhiraj Joshi, Zexue He, Shafiq Abedin, Jennifer Sun, Ben Wiesel, Eli Schwartz, Ahmed Nassar, Bo Wu, Assaf Arbelle, Aude Oliva, Dan Gutfreund, Leonid Karlinsky, and Rogerio Feris. 2025. [Chartgen: Scaling chart understanding via code-guided synthetic chart generation](#). *Preprint*, arXiv:2507.19492.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). *Preprint*, arXiv:2309.06180.
- Honglin Lin, Zheng Liu, Yun Zhu, Chonghan Qin, Juekai Lin, Xiaoran Shang, Conghui He, Wentao Zhang, and Lijun Wu. 2026a. [Mmfinereason: Closing the multimodal reasoning gap via open data-centric methods](#). *arXiv preprint arXiv:2601.21821*.
- Honglin Lin, Chonghan Qin, Zheng Liu, Qizhi Pei, Yu Li, Zhanping Zhong, Xin Gao, Yanfeng Wang, Conghui He, and Lijun Wu. 2026b. [Scientific image synthesis: Benchmarking, methodologies, and downstream utility](#). *arXiv preprint arXiv:2601.17027*.
- Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. 2024. [Mmc: Advancing multimodal chart understanding with large-scale instruction tuning](#). *Preprint*, arXiv:2311.10774.
- Shudong Liu, Hongwei Liu, Junnan Liu, Linchen Xiao, Songyang Gao, Chengqi Lyu, Yuzhe Gu, Wenwei Zhang, Derek F. Wong, Songyang Zhang, and Kai Chen. 2025a. [Compassverifier: A unified and robust verifier for llms evaluation and outcome reward](#). *Preprint*, arXiv:2508.03686.
- Zheng Liu, Hao Liang, Bozhou Li, Wentao Xiong, Chong Chen, Conghui He, Wentao Zhang, and Bin Cui. 2025b. [Synthvlm: Towards high-quality and efficient synthesis of image-caption datasets for vision-language models](#). In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 12807–12814.
- Zheng Liu, Mengjie Liu, Jingzhou Chen, Jingwei Xu, Bin Cui, Conghui He, and Wentao Zhang. [Flare: Fully integration of vision-language representations for deep cross-modal understanding](#). In *The Fourteenth International Conference on Learning Representations*.

- Zheng Liu, Mengjie Liu, Jingzhou Chen, Jingwei Xu, Bin Cui, Conghui He, and Wentao Zhang. 2025c. Fusion: Fully integration of vision-language representations for deep cross-modal understanding. *arXiv preprint arXiv:2504.09925*.
- Zheng Liu, Mengjie Liu, Siwei Wen, Mengzhang Cai, Bin Cui, Conghui He, and Wentao Zhang. 2025d. From uniform to heterogeneous: Tailoring policy optimization to every token's nature. *arXiv preprint arXiv:2509.16591*.
- Zhuoming Liu, Xiaofeng Gao, Feiyang Niu, Qiaozhi Gao, Liu Liu, and Robinson Pirmuthu. 2025e. [Start: Spatial and textual learning for chart understanding](#). *Preprint*, arXiv:2512.07186.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kaiwei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- Ahmed Masry, Mohammed Saidul Islam, Mahir Ahmed, Aayush Bajaj, Firoz Kabir, Aaryaman Kartha, Md Tahmid Rahman Laskar, Mizanur Rahman, Shadikur Rahman, Mehrad Shahmohammadi, Megh Thakkar, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. 2025. [Chartqapro: A more diverse and challenging benchmark for chart question answering](#). *Preprint*, arXiv:2504.05506.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. [Chartqa: A benchmark for question answering about charts with visual and logical reasoning](#). *Preprint*, arXiv:2203.10244.
- Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. 2024a. [Chartinstruct: Instruction tuning for chart comprehension and reasoning](#). *Preprint*, arXiv:2403.09028.
- Ahmed Masry, Megh Thakkar, Aayush Bajaj, Aaryaman Kartha, Enamul Hoque, and Shafiq Joty. 2024b. [Chartgemma: Visual instruction-tuning for chart reasoning in the wild](#). *Preprint*, arXiv:2407.04172.
- Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. 2024. [Chartassisstant: A universal chart multimodal language model via chart-to-table pre-training and multitask instruction tuning](#). *Preprint*, arXiv:2401.02384.
- Nitish Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. 2020. [Plotqa: Reasoning over scientific plots](#). *Preprint*, arXiv:1909.00997.
- Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, and 1 others. 2025. Mineru2. 5: A decoupled vision-language model for efficient high-resolution document parsing. *arXiv preprint arXiv:2509.22186*.
- OpenAI. 2025. [Gpt-5 system card](#). Technical report, OpenAI.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). *Preprint*, arXiv:2103.00020.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. [Hybridflow: A flexible and efficient rlhf framework](#). In *Proceedings of the Twentieth European Conference on Computer Systems*, EuroSys '25, page 1279–1297. ACM.
- Liyan Tang, Grace Kim, Xinyu Zhao, Thom Lake, Wenxuan Ding, Fangcong Yin, Prasann Singhal, Manya Wadhwa, Zeyu Leo Liu, Zayne Sprague, Ramya Namuduri, Bodun Hu, Juan Diego Rodriguez, Puyuan Peng, and Greg Durrett. 2025. [Chartmuseum: Testing visual reasoning capabilities of large vision-language models](#). *Preprint*, arXiv:2505.13444.
- Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, Zhaokai Wang, Zhe Chen, Hongjie Zhang, Ganlin Yang, Haomin Wang, Qi Wei, Jinhui Yin, Wenhao Li, Erfei Cui, and 56 others. 2025. [Internvl3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency](#). *Preprint*, arXiv:2508.18265.
- Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. 2024. [Charxiv: Charting gaps in realistic chart understanding in multimodal llms](#). *Preprint*, arXiv:2406.18521.
- Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Peng Ye, Min Dou, Botian Shi, Junchi Yan, and Yu Qiao. 2025. [Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning](#). *Preprint*, arXiv:2402.12185.
- Yijia Xiao, Edward Sun, Tianyu Liu, and Wei Wang. 2024. [Logicvista: Multimodal llm logical reasoning benchmark in visual contexts](#). *arXiv preprint arXiv:2407.04973*.
- Duo Xu, Hao Cheng, Xin Lin, Zhen Xie, and Hao Wang. 2025a. [Chartm<sup>3</sup>: A multi-stage code-driven pipeline for constructing multi-dimensional and multi-step visual reasoning data in chart comprehension](#). *Preprint*, arXiv:2511.02415.
- Weiye Xu, Jiahao Wang, Weiyun Wang, Zhe Chen, Wengang Zhou, Aijun Yang, Lewei Lu, Houqiang Li, Xiaohua Wang, Xizhou Zhu, and 1 others. 2025b. [Visulogic: A benchmark for evaluating visual reasoning in multi-modal large language models](#). *arXiv preprint arXiv:2504.15279*.
- Zhengzhuo Xu, Sinan Du, Yiyan Qi, Chengjin Xu, Chun Yuan, and Jian Guo. 2024. [Chartbench: A benchmark for complex visual reasoning in charts](#). *Preprint*, arXiv:2312.15915.

- Yue Yang, Ajay Patel, Matt Deitke, Tanmay Gupta, Luca Weihs, Andrew Head, Mark Yatskar, Chris Callison-Burch, Ranjay Krishna, Aniruddha Kembhavi, and Christopher Clark. 2025a. [Scaling text-rich image understanding via code-guided synthetic multimodal data generation](#). *Preprint*, arXiv:2502.14846.
- Yuwei Yang, Zeyu Zhang, Yunzhong Hou, Zhuowan Li, Gaowen Liu, Ali Payani, Yuan-Sen Ting, and Liang Zheng. 2025b. [Effective training data synthesis for improving mllm chart understanding](#). *Preprint*, arXiv:2508.06492.
- Kaichen Zhang, Keming Wu, Zuhao Yang, Kairui Hu, Bin Wang, Ziwei Liu, Xingxuan Li, and Lidong Bing. 2025. Openmmreasoner: Pushing the frontiers for multimodal reasoning with an open and general recipe. *arXiv preprint arXiv:2511.16334*.
- Liang Zhang, Anwen Hu, Haiyang Xu, Ming Yan, Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang. 2024a. [Tinchart: Efficient chart understanding with visual token merging and program-of-thoughts learning](#). *Preprint*, arXiv:2404.16635.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, and 1 others. 2024b. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pages 169–186. Springer.
- Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2025. [Chartcoder: Advancing multimodal large language model for chart-to-code generation](#). *Preprint*, arXiv:2501.06598.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, Jingren Zhou, and Junyang Lin. 2025. [Group sequence policy optimization](#). *Preprint*, arXiv:2507.18071.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyuan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. [Llamafactory: Unified efficient fine-tuning of 100+ language models](#). *Preprint*, arXiv:2403.13372.
- Yun Zhu, Xiaoran Shang, Qifan Yu, Kaihang Pan, Boci Peng, Xiaotang Wang, Xiaoyang Wang, Yu Li, Honglin Lin, Zheng Liu, and 1 others. 2026. From factors to methods: A comprehensive survey on visual instruction tuning data selection. *Authorea Preprints*.
- Chengke Zou, Xingang Guo, Rui Yang, Junyu Zhang, Bin Hu, and Huan Zhang. 2024. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models. *arXiv preprint arXiv:2411.00836*.

# Appendix

## A More Training Details

We implement our training pipeline using the LLaMA-Factory and Verl frameworks, conducted on a cluster of 4 nodes equipped with 8 A100 GPUs each (32 GPUs in total).

**Supervised Fine-Tuning.** We perform full-parameter fine-tuning on the Qwen3-VL-Instruct models. The visual encoder is configured to handle dynamic resolutions ranging from 802, 816 to 3, 211, 264 pixels. To support extensive chart reasoning, we set a large context cutoff length of 22,000 tokens. To enhance training efficiency, we employ sequence packing and configure the training with a global batch size of 128. The model is trained for 1 epoch using the cosine learning rate scheduler with a peak learning rate of  $1.0 \times 10^{-5}$  and a warmup ratio of 0.1. We utilize DeepSpeed ZeRO-3 optimization and BF16 to maximize memory efficiency during the full-parameter update. The entire fine-tuning process for the 8B model takes approximately 1.5 days.

**Reinforcement Learning.** Following SFT, we apply Group Relative Policy Optimization (GRPO) to further align the model’s reasoning capabilities. We employ vLLM as the rollout engine, sampling  $G = 16$  distinct responses for each prompt with a temperature of 1.0 and no top- $k$  restriction. The context window is expanded to support long-chain reasoning, with a maximum prompt length of 16,384 tokens and a maximum response length of 16,384 tokens. The actor model is optimized with a learning rate of  $1.0 \times 10^{-6}$ , utilizing a global batch size of 128. To ensure stability with such long contexts, we disable the standard KL-divergence penalty ( $\beta_{KL} = 0$ ) and instead rely on a conservative clipping strategy (ratio  $\epsilon \in [3 \times 10^{-4}, 4 \times 10^{-4}]$ ). The RL process runs for approximately 300 steps. Training the 8B model in this stage requires approximately 4 days.

## B Evaluation Details

We conduct a comprehensive evaluation using the VLMEvalKit (Duan et al., 2025). Since specific domain-specialized datasets such as ChartX, EvoChart, and ChartBench are not natively integrated, we preprocessed and converted these benchmarks into the standard VLMEvalKit interface

format. For evaluation metrics, we replace traditional exact string matching with the compass-verifier (Liu et al., 2025a), which employs an LLM-as-a-Judge to accurately assess response correctness. Regarding rollout settings, following the guidelines from OpenDataArena (Cai et al., 2025; Zhu et al., 2026) and the official Qwen documentation, we set the temperature to 0.6, top- $p$  to 0.95, and top- $k$  to 20. Additionally, we apply a repetition penalty of 1.05 and limit the maximum response length to 32768 tokens. Furthermore, we exclude all system prompts during inference.

## C Data Synthesis Details

### C.1 Complexity-Aware Chart Coder

We first illustrate the distribution of our cold start dataset  $\mathcal{C}_{\text{cold}}$ , in Table 4. This diverse collection serves as the foundational corpus for initializing our coders. Furthermore, to generate high-quality code representations for complex charts, we utilize Claude-4-Sonnet. The specific prompt designed for this generation process is presented in Table 11.

Table 4: Data distribution of the cold start dataset  $\mathcal{C}_{\text{cold}}$ .

Source	Count	Percentage
CoSyn	22088	35.67%
ChartGen	15337	24.77%
ChartCoder	10214	16.50%
ChartQA	5602	9.05%
FigureQA	4214	6.81%
PlotQA	2150	3.47%
ECD	1771	2.86%
ReachQA	543	0.88%
<b>Total</b>	<b>61919</b>	<b>100.00%</b>

**Training Configuration.** We train the Complexity-Aware Chart Coder using Qwen2.5-Coder-7B-Instruct as the backbone model. During both training and inference, we employ the identical system prompt  $\mathcal{T}$ , as presented in Table 12. The training pipeline is implemented on the LLaMA-Factory framework. We apply the same training configuration across every iterative training stage to ensure consistency. We perform full-parameter fine-tuning with a global batch size of 16 and a context cutoff length of 4,096 tokens, utilizing sequence packing to maximize computational efficiency. The optimization process runs for 5 epochs using a

Table 5: Performance evolution of ChartVerse across different training phases on chart and STEM benchmarks.

Model	Chart-Related								STEM-Related					
	ChartQA-Pro	CharXiv(DQ)	CharXiv(RQ)	ChartMuseum	ChartX	EvoChart	ChartBench	Avg	MathVista	DynaMath	MathVerse	LogicVista	VisuLogic	Avg
Qwen3-VL-Instruct-2B	42.1	26.8	62.3	23.9	49.8	53.6	38.8	42.5	61.3	52.1	54.2	35.8	11.5	43.0
+ SFT	44.4	40.8	69.1	30.0	56.9	61.0	46.4	49.8	58.7	47.6	55.8	43.2	17.5	44.6
+ RL	48.2	46.9	71.2	37.5	60.5	66.8	49.1	54.3	60.4	49.2	56.8	48.5	23.1	47.6
Qwen3-VL-Instruct-4B	53.7	39.7	76.2	37.2	57.2	68.2	45.1	53.9	73.7	46.8	65.3	53.2	19.0	51.6
+ SFT	52.9	52.8	84.5	42.5	61.9	73.3	50.3	59.7	70.9	60.7	71.1	57.3	25.0	57.0
+ RL	55.2	56.2	84.1	45.9	63.7	75.0	52.9	61.9	72.8	62.0	71.8	59.3	27.0	58.6
Qwen3-VL-Instruct-8B	54.4	46.4	83.0	39.6	58.2	70.2	46.6	56.9	77.2	62.1	67.7	55.3	22.5	56.7
+ SFT	55.5	56.2	88.3	47.5	61.0	76.7	52.2	62.5	75.0	67.4	75.6	61.3	26.5	61.2
+ RL	56.2	60.8	88.0	49.2	63.9	76.2	54.2	64.1	75.6	69.0	76.5	62.6	27.1	62.2

Table 6: Comparison of Different Datasets. All models are trained on Qwen3-VL-4B-Instruct using 100K samples.

Training Data (100K)	ChartQA Pro	CharXiv (RQ)	CharXiv (DQ)	Chart Museum	ChartX	Evo Chart	ChartBench (GPT-acc)	Avg
Qwen3-VL-4B-Instruct	53.7	39.7	76.2	37.2	57.2	68.2	45.1	53.9
+ CoSyn	50.3	35.1	74.4	33.5	55.8	62.4	45.8	51.0
+ START	53.5	39.2	77.2	37.9	56.4	68.0	44.5	53.8
+ ECD	54.4	38.3	77.5	37.7	58.3	68.8	44.7	54.2
+ ChartVerse-SFT (Ours)	53.4	48.8	82.6	40.3	60.1	70.9	48.7	57.8

Table 7: Comparison of data selection strategies. The Fail Rate (left) indicates sample difficulty, while the right section shows the results of Qwen3-VL-4B-Instruct trained on 100k selected data.

Strategy	Difficulty	Downstream Performance (%)							
	Fail Rate (N=500)	ChartQA Pro	CharXiv (RQ)	CharXiv (DQ)	Chart Museum	ChartX	Evo Chart	ChartBench (GPT-acc)	Avg
Qwen3-VL-4B-Instruct	-	53.7	39.7	76.2	37.2	57.2	68.2	45.1	53.9
+ VLM-as-Judge	21.1%	53.4	40.3	75.7	38.8	58.6	69.2	45.4	54.5
+ Python Code Complexity	23.5%	53.7	40.5	77.7	39.0	58.5	70.3	44.7	54.9
+ RPE (Ours)	27.6%	53.3	41.0	78.6	39.2	58.5	70.8	46.3	55.4

cosine learning rate scheduler, configured with a peak learning rate of  $2.0 \times 10^{-5}$  and a warmup ratio of 0.05. We employ DeepSpeed ZeRO-3 optimization with BF16 precision.

**Sampling Configuration.** During the large-scale sampling phase, we set the sampling temperature to 1.0, top- $p$  to 0.95, and top- $k$  to 20, ensuring diverse and robust generation. Regarding the RPE, we processed a total of approximately 4 million samples. Leveraging vLLM, our RPE computation is highly efficient, completed in 3 days on 64 A100 GPUs.

**RPE Data Generation.** Regarding the RPE, we processed a total of approximately 4 million samples. This process involved approximately 32 million inference calls to the Qwen3-VL-2B-Thinking model. Leveraging vLLM (Kwon et al., 2023), our RPE computation is highly efficient, completed in 3 days on 64 A100 GPUs

**Data Decontamination.** To ensure zero overlap between training and evaluation data, we perform rigorous decontamination on the initial image pool

$\mathcal{I}_{pool}$ , ChartVerse-SFT-600K, and ChartVerse-RL-40K. Specifically, we discard any training sample whose chart image exhibits a CLIP cosine similarity exceeding 0.65 with any benchmark image. Notably, because our training charts are synthesized from scratch via code generation rather than sourced from existing datasets, the number of filtered images is negligibly small.

## C.2 Truth-Anchored Inverse QA Synthesis

**Inverse Logic Construction.** Our synthesis pipeline relies on a sequence of carefully designed prompts to ensure logical consistency and factual accuracy. We present the specific prompts utilized for Python script generation, reverse question synthesis, and consistency checking in Table 13, Table 14, and Table 15, respectively.

Specifically, in the Python script generation phase (Table 13), we enforce a strict constraint: the generated code must execute to yield a precise numerical value or a definitive categorical label, ensuring the answer is deterministic. Following this, during the reverse question synthesis (Table 14), we mandate that the synthesized text must be a

Table 8: Comparison of ChartVerse performance across Instruct and Thinking model bases.

Model Base	ChartQA Pro	CharXiv (RQ)	CharXiv (DQ)	Chart Museum	ChartX	Evo Chart	Chart Bench	Avg
Qwen3-VL-4B-Instruct	55.2	56.2	84.1	45.9	63.7	75.0	52.9	61.9
Qwen3-VL-4B-Thinking	55.8	57.6	85.0	47.7	64.1	76.1	54.5	63.0
Qwen3-VL-8B-Instruct	56.2	60.8	88.0	49.2	63.9	76.2	54.2	64.1
Qwen3-VL-8B-Thinking	57.3	61.6	88.2	51.4	65.2	77.0	55.1	65.1

Table 9: Ablation of teacher model scale. All results use Qwen3-VL-8B-Instruct as the student model.

Training Data	Teacher Model	ChartQA Pro	CharXiv (RQ)	CharXiv (DQ)	Chart Museum	ChartX	Evo Chart	Chart Bench	Avg
ChartVerse-SFT-600K	30B-A3B-Thinking	55.5	56.2	88.3	47.5	61.0	76.7	52.2	62.5
ChartVerse-SFT-100K	235B-A30B-Thinking	56.9	59.2	88.7	49.9	63.4	76.5	54.0	64.1

valid interrogative sentence, and the answer to this question must correspond exactly to the execution result of the preceding Python script.

In our specific implementation, we synthesize two distinct Python scripts for each chart code, followed by separate reverse question synthesis and consistency checks. This process involved a total of 4 million calls to Qwen-30B-A3B-Thinking, taking approximately 4 days on 128 A100 GPUs.

**CoT Distillation & Difficulty Filtration.** We display the prompt used for distilling Chain-of-Thought (CoT) reasoning traces in Table 16.

To ensure high-quality and non-redundant reasoning traces, we apply a multi-stage filtering procedure to the distilled outputs:

- **Template and Length Validation:** We first impose strict structural validation to ensure the usability of the distilled output. Specifically, we filter out any reasoning trace that fails to adhere to the mandated `<think>...</think>` and `<answer>...</answer>` output template. Furthermore, to prevent the retention of superficial or trivial rationales, we enforce a minimum length constraint, discarding traces that are shorter than 100 words.
- **N-gram De-duplication:** We detect and remove templated or overly repetitive CoTs using an n-gram overlap criterion. Concretely, we flag CoTs that contain any 50-gram that repeats at least 3 times. Flagged traces are discarded.

In this step, we processed approximately 1 million samples. With each sample requiring three

distinct inference calls, this amounted to a total of 3 million calls to Qwen3-VL-30B-A3B-Thinking, consuming approximately 5 days on 128 A100 GPUs.

Throughout all processes in this stage, we set the LLM inference parameters to a temperature of 0.6, top- $p$  of 0.95, top- $k$  of 20, and a maximum output length of 32768.

## D Additional Results

### D.1 Detailed Performance Analysis of Training Phases

We provide the specific performance corresponding to Figure 3 in Table 5. The results across the 2B, 4B, and 8B scales demonstrate a consistent upward trajectory. ChartVerse-8B achieves a significant leap on CharXiv (RQ), rising from a baseline of 46.4 to 60.8. Similarly, on ChartMuseum, it reaches a peak performance of 49.2. The reasoning skills acquired from chart-specific data also generalize effectively to out-of-domain STEM tasks. ChartVerse-8B improves its performance on MathVerse from 67.5 to 76.5 and on DynaMath from 62.1 to 69.0. These gains confirm that the logical reasoning capabilities learned through our framework transfer robustly to complex scientific contexts.

### D.2 Comparison with Existing Datasets

Complementing the visual comparison in Figure 4a, Table 6 evaluates ChartVerse-SFT against competing datasets under a fixed 100K sample budget. While datasets like CoSyn and START degrade performance on specific benchmarks, ChartVerse-SFT consistently drives improvements. Specifically, our dataset boosts ChartMuseum to 40.3 and signifi-

Table 10: Comparison of different models for RPE calculation. Results are computed on a random subset of 10K charts. **Range** denotes the difference between the 75th and 25th percentiles.

Model	Mean RPE	Std RPE	25th Perc.	75th Perc.	Effective Range	Execution Success Rate
Qwen3-VL-2B-Instruct	0.48	0.16	0.38	0.55	0.17	63.2%
Qwen3-VL-2B-Thinking	0.41	0.18	0.32	0.51	0.19	88.6%
Qwen3-VL-4B-Thinking	0.35	0.10	0.30	0.42	0.12	93.3%
Qwen3-VL-8B-Thinking	0.32	0.09	0.27	0.37	0.10	96.1%

cantly lifts CharXiv (RQ) performance from 39.7 to 48.8, confirming that our gains stem from superior data quality rather than scale.

### D.3 Ablation on RPE Strategy

The effectiveness of the Rollout Posterior Entropy strategy, as shown in Figure 4b, is further quantified in Table 7. RPE successfully identifies intrinsically harder samples, evidenced by a model failure rate of 27.6, compared to only 21.1 for VLM-as-Judge. This focus on difficulty translates directly to downstream accuracy: on the CharXiv (DQ) benchmark, the RPE-trained model achieves 78.6, outperforming both VLM-as-Judge 75.7 and Python Code Complexity 77.7.

### D.4 Impact of the Model Base

We investigate the compatibility of ChartVerse data with different model architectures by comparing Instruct and Thinking base models. Our primary experiments employ the Instruct series to rigorously validate that the ChartVerse dataset can inject complex reasoning capabilities from scratch, without relying on built-in extended thinking pre-training. To evaluate generalizability, we apply the full training recipe (ChartVerse-SFT-600K + ChartVerse-RL-40K) to both Instruct and Thinking models at the 4B and 8B scales.

As shown in Table 8, the Thinking bases yield additional improvements over the Instruct bases, confirming that our synthesized data is highly compatible with and beneficial for thinking-capable models. However, the performance gap between the two bases is relatively narrow. We attribute this to the extensive scale of ChartVerse-SFT-600K, which substantially reshapes the Instruct model’s reasoning behavior, endowing it with analytical capabilities comparable to the native Thinking variants.

### D.5 Impact of the Teacher Model Scale

Our primary data synthesis pipeline utilizes Qwen3-VL-30B-A3B-Thinking as the teacher to balance data quality and computational feasibility. By leveraging truth-anchored synthesis and failure-rate filtering, our framework successfully enables the student model to surpass the teacher despite the teacher’s relatively constrained capacity.

To explore the effect of scaling the teacher model, we conduct an ablation study using a significantly larger teacher, Qwen3-VL-235B-A30B-Thinking, to synthesize a subset of 100K high-quality SFT samples. As presented in Table 9, training the 8B Instruct model on this 100K dataset achieves an average score of 64.1, which already outperforms the model trained on the full 600K dataset generated by the 30B teacher. This confirms that a stronger teacher directly enhances the synthesized data quality, and validates that our ChartVerse framework scales effectively with teacher capacity, yielding superior student performance even with a reduced sample size.

## E Model Selection for RPE Computation

The Relative Prediction Entropy (RPE) serves as our core metric for assessing chart complexity. We select **Qwen3-VL-2B-Thinking** as the base model for RPE calculation with 8 inference passes ( $K = 8$ ), based on the following three considerations:

**Computational Efficiency.** RPE requires multiple forward passes per image. Given our initial pool of over 3M images, the 2B-scale model offers the most practical trade-off, enabling large-scale entropy estimation within a reasonable computational budget.

**Complexity Discrimination.** Smaller models exhibit higher sensitivity to structural complexity. As shown in Table 10, larger models (4B, 8B) reconstruct complex charts more consistently, which compresses the RPE distribution (standard devi-

ation decreases from 0.18 to 0.09). The 2B model provides the widest RPE range and the highest standard deviation, offering the sharpest discriminative power to separate simple and complex charts.

**Execution Reliability.** While Instruct and Thinking variants at the 2B scale show similar entropy spreads, the Thinking variant achieves a significantly higher code execution success rate (88.6% vs. 63.2%). Since RPE relies on successfully rendered images to compute valid entropy, the structural reasoning of the Thinking model minimizes noise from failed executions, ensuring more stable and trustworthy complexity scores.

## F ChartVerse-SFT Analysis

We first demonstrate the visual complexity of our synthesized training corpus, ChartVerse-SFT. As illustrated in Figure 6, the dataset exhibits exceptional diversity by covering a wide spectrum of visualization types. These range from standard statistical graphs such as violin plots and radar charts to high-dimensional 3D representations and hierarchical structures like treemaps. Crucially, a significant portion of the data features intricate multi-subplot layouts and mixed-type dashboards. This structural complexity compels the model to handle diverse visual encoding schemas and perform fine-grained reasoning across spatially distinct panels, significantly surpassing the difficulty of traditional chart datasets.

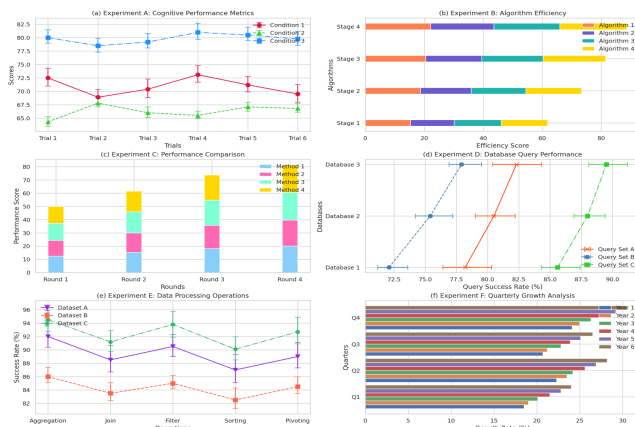
Furthermore, we showcase the intricate QA pairs generated within our dataset through specific case studies shown in Figures 7, 8, and 9. Leveraging the high information density of the synthesized charts, our pipeline is capable of formulating highly challenging queries. As evident in these examples, the questions transcend simple data retrieval by necessitating rigorous multi-step reasoning and cross-subplot integration.

For instance, the example in Figure 7 presents a multi-condition verification task where the model must sequentially validate three separate metrics including average scores and total efficiency across distinct experiment sub-panels. Similarly, Figure 8 demonstrates a demand for complex derived calculations. Here the model must visually extract statistical markers such as quartiles and medians to determine the region with the greatest proportional variation. This design ensures that the model must develop a holistic understanding of the visual context rather than relying on local pattern matching.

Finally, Figure 9 requires the model to aggregate information from multiple spatially distinct plots. To identify the element with the optimal trade-off, the model must synthesize data regarding power levels, temperature stability, and complexity scores from three separate charts.



Figure 6: Overview of the ChartVerse-SFT dataset. The samples demonstrate high diversity in chart types including 3D, hierarchical, and statistical plots, as well as structural complexity featuring various multi-subplot layouts.



### Code Solution

```

1 # Cognitive performance data (Experiment A) for three conditions over 6 trials
2 cognitive_condition_a_scores = [72.5, 68.9, 78.4, 73.1, 71.2, 69.3]
3 cognitive_condition_b_scores = [64.3, 67.8, 66.0, 65.5, 67.1, 66.8]
4 cognitive_condition_c_scores = [88.0, 78.5, 79.2, 81.0, 88.5, 79.8]
5
6 # Calculate average scores per cognitive condition
7 condition_a_average = sum(cognitive_condition_a_scores) / len(cognitive_condition_a_scores)
8 condition_b_average = sum(cognitive_condition_b_scores) / len(cognitive_condition_b_scores)
9 condition_c_average = sum(cognitive_condition_c_scores) / len(cognitive_condition_c_scores)
10
11 # Determine highest cognitive performance average
12 best_cognitive_average = max(condition_a_average, condition_b_average, condition_c_average)
13
14 # Algorithm efficiency data (Experiment B) - 4 algorithms across 4 stages
15 algorithm_efficiency_by_stage = [
16     [15.3, 18.7, 20.5, 22.1], # Algorithm 1
17     [14.8, 17.2, 19.0, 21.5], # Algorithm 2
18     [16.0, 18.5, 20.8, 22.3], # Algorithm 3
19     [15.7, 18.9, 21.2, 22.8] # Algorithm 4
20 ]
21
22 # Calculate total efficiency for final stage (Stage 4)
23 stage4_efficiency = sum(stage[3] for stage in algorithm_efficiency_by_stage)
24
25 # Quarterly growth data (Experiment F) - 6 years, 4 quarters each
26 quarterly_growth_rates = [
27     [18.5, 22.3, 20.7, 24.1], # Year 1
28     [19.0, 23.5, 21.2, 25.8], # Year 2
29     [20.1, 24.2, 22.8, 26.3], # Year 3
30     [21.5, 25.6, 23.9, 27.8], # Year 4
31     [22.8, 26.9, 25.1, 29.2], # Year 5
32     [24.0, 28.2, 26.5, 30.5] # Year 6 (latest)
33 ]
34
35 # Calculate average growth rate for latest year (Year 6)
36 year6_average_growth = sum(quarterly_growth_rates[5]) / len(quarterly_growth_rates[5])
37
38 # Business thresholds for validation
39 cognitive_threshold = 79
40 stage4_threshold = 80
41 year6_threshold = 20
42
43 # Validate all metrics against thresholds
44 is_cognitive_exceeded = best_cognitive_average > cognitive_threshold
45 is_stage4_exceeded = stage4_efficiency > stage4_threshold
46 is_year6_exceeded = year6_average_growth > year6_threshold
47
48 # Final result: "Yes" if all thresholds met, else "No"
49 result = "Yes" if (is_cognitive_exceeded and is_stage4_exceeded and is_year6_exceeded) else "No"
50 print(result)

```

**Q:** Does the top cognitive condition's average score exceed 79, the final algorithm stage's total efficiency surpass 80, and the latest year's average growth rate exceed 20?

**A:** Yes

### CoT Solution

I need to solve this problem which has three parts to check:

1. Does the top cognitive condition's average score exceed 79?
2. Does the final algorithm stage's total efficiency surpass 80?
3. Does the latest year's average growth rate exceed 20?

I'll analyze each part using the provided graphs.

First, let's look at part 1: "Does the top cognitive condition's average score exceed 79?"

This refers to graph (a) "Experiment A: Cognitive Performance Metrics". This graph shows scores for three conditions (1, 2, and 3) across 6 trials.

I need to find the top cognitive condition. Looking at the graph:

- Condition 1 (red line with circles) - seems to have the highest scores overall
- Condition 2 (green line with triangles)
- Condition 3 (blue line with squares)

To find the average score for the top condition, I need to calculate the average of Condition 1's scores across all 6 trials.

From the graph (a):

- Trial 1: ~73
- Trial 2: ~69
- Trial 3: ~71
- Trial 4: ~74
- Trial 5: ~72
- Trial 6: ~70

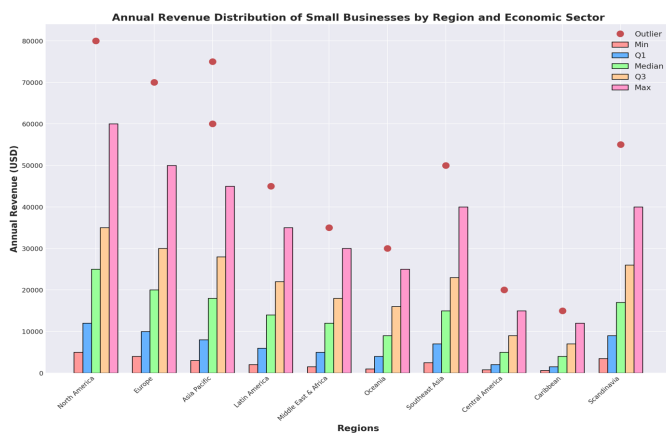
Let me calculate the average:

$$(73 + 69 + 71 + 74 + 72 + 70) / 6 = 429 / 6 = 71.5$$

So the average score for Condition 1 is approximately 71.5, which is less than 79.

Wait, I need to double-check this. The graph shows that Condition 3 (blue line) actually has the highest scores. Let me look more carefully .....

Figure 7: ChartVerse-SFT QA Example-1



**Q:** Which region demonstrates the greatest proportional variation in annual revenue compared to its typical revenue level?

**A:** Central America

```

Code Solution

1: regions = [
2:     "North America",
3:     "Europe",
4:     "Asia Pacific",
5:     "Latin America",
6:     "Middle East & Africa",
7:     "Oceania",
8:     "Southeast Asia",
9:     "Central America",
10:    "Caribbean",
11:    "Scandinavia"
12: ]
13:
14: economic_sectors = [
15:     "Healthcare",
16:     "Technology",
17:     "Finance",
18:     "Manufacturing",
19:     "Retail",
20:     "Education",
21:     "Energy",
22:     "Agriculture",
23:     "Hospitality",
24:     "Telecommunications"
25: ]
26:
27: min_vals = [5000, 4000, 3000, 2000, 1500, 1000, 2500, 800, 600, 3500]
28: q1_vals = [12000, 10000, 8000, 6000, 4000, 5000, 4000, 7000, 2000, 1500, 9000]
29: median_vals = [20000, 20000, 18000, 14000, 12000, 9000, 12000, 5000, 3000, 4000, 17000]
30: q3_vals = [32000, 30000, 28000, 22000, 18000, 16000, 23000, 9000, 7000, 26000]
31: max_vals = [60000, 50000, 45000, 35000, 30000, 25000, 40000, 15000, 12000, 40000]
32:
33: outliers = [
34:     [80000],
35:     [70000],
36:     [60000, 75000],
37:     [45000],
38:     [35000],
39:     [20000],
40:     [50000],
41:     [10000],
42:     [15000],
43:     [55000]
44: ]
45:
46: # Calculate interquartile range for each region
47: region_interquartile_ranges = [q3_vals[i] - q1_vals[i] for i in range(len(regions))]
48:
49: # Calculate normalized revenue dispersion (IQR relative to median)
50: region_revenue_dispersion = [
51:     region_interquartile_ranges[i] / median_vals[i]
52:     for i in range(len(regions))
53: ]
54:
55: # Identify region with highest normalized dispersion
56: best_region_index = region_revenue_dispersion.index(max(region_revenue_dispersion))
57: best_region = regions[best_region_index]
58:
59: print(best_region)

```

**CoT Solution**

Okay, let's tackle this problem step by step. The question is asking which region has the greatest proportional variation in annual revenue compared to its typical revenue level. First, I need to understand what "proportional variation" means here. I think it refers to the ratio of the spread (like the interquartile range or the range) to the typical revenue, which is probably the median. So, higher proportional variation would mean that the spread is larger relative to the median.

Looking at the chart, it's a box plot with some outliers. The key shows that the colors correspond to Min (light red), Q1 (blue), and Outliers. But wait, the chart actually has bars for each of these, so each region has multiple bars: Min (light red), Q1 (blue), Median (green), Q3 (orange), Max (pink), and the red dots are outliers. So for each region, the data is represented by these five values plus outliers.

But to find the proportional variation, I need to calculate the range (Max - Min) or maybe the interquartile range (Q3 - Q1) divided by the median. The problem says "proportional variation compared to its typical revenue level". The typical revenue level is likely the median, so the variation would be (range / median) or (IQR / median). But the question says "greatest proportional variation", so we need to see which region has the highest ratio of the spread to the median.

Let me check the chart for each region. Let's list out the regions and their values. The regions are: North America, Europe, Asia Pacific, Latin America, Middle East & Africa, Oceania, Southeast Asia, Central America, Caribbean, Scandinavia.

For each region, I need to find:

1. The median (green bar)
2. The spread (maybe Max - Min, since that's the total range, or Q3 - Q1 as the IQR. But since the question says "proportional variation", which could be the coefficient of variation, but in the context of box plots, the range (Max - Min) divided by median might be a good measure here. Alternatively, the IQR over median. Let's see.)

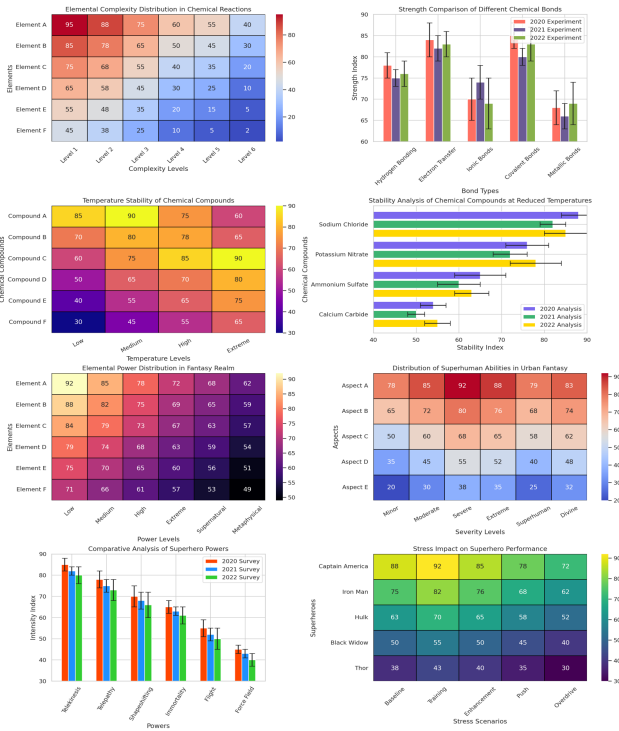
But the problem says "proportional variation in annual revenue compared to its typical revenue level". So the typical level is the median. The variation could be the range (Max - Min) or the IQR (Q3 - Q1). Let's check both.

First, let's get the values from the chart. I need to read the approximate values for each region.

Starting with North America:

- Min (red bar): ~4,000 - Q1 (blue): ~12,000 - Median (green): ~25,000 - Q3 (orange): ~35,000 - Max (pink): ~60,000
- Outliers: 80,000 and 70,000 (but those are outliers, so Max is 60,000) .....

Figure 8: ChartVerse-SFT QA Example-2



```

Code Solution

1 # Define datasets from visualization code
2 element_complexity_data = [
3     [95, 88, 75, 60, 55, 40],
4     [85, 78, 65, 50, 45, 30],
5     [75, 68, 55, 40, 35, 20],
6     [65, 58, 45, 30, 25, 10],
7     [55, 48, 35, 20, 15, 5],
8     [45, 38, 25, 10, 5, 2]
9 ]
10
11 element_power_data = [
12     [92, 85, 78, 72, 68, 62],
13     [88, 82, 75, 69, 65, 59],
14     [84, 79, 73, 67, 63, 57],
15     [79, 74, 68, 63, 59, 54],
16     [75, 70, 65, 60, 56, 51],
17     [71, 66, 61, 57, 53, 49]
18 ]
19
20 compound_stability_data = [
21     [85, 90, 75, 60],
22     [78, 88, 76, 63],
23     [68, 75, 65, 50],
24     [58, 65, 70, 80],
25     [50, 55, 65, 75],
26     [38, 45, 55, 65]
27 ]
28
29 # Extract required metrics
30 element_complexity_level = [row[0] for row in element_complexity_data]
31 element_power_low = [row[0] for row in element_power_data]
32 element_stability_extreme = [row[3] for row in compound_stability_data]
33
34 # Calculate efficiency metric: (Power * Stability) + Complexity
35 efficiency_scores = [
36     (power * stability) / complexity
37     for power, stability, complexity in zip(
38         element_power_low, element_stability_extreme, element_complexity_level
39     )
40 ]
41
42 # Identify element with highest efficiency
43 element_names = ['Element A', 'Element B', 'Element C', 'Element D', 'Element E', 'Element F']
44 top_element = element_names[efficiency_scores.index(max(efficiency_scores))]
45
46 # Output final result
47 print(top_element)
48
49

```

**Q:** Which element achieves the greatest balance of high power and extreme-temperature stability relative to its inherent complexity?

**A:** Element F

**CoT Solution**

Okay, let's tackle this problem step by step. The question is asking which element has the greatest balance of high power, extreme-temperature stability, and low inherent complexity. So, I need to look at three different charts for each element: Elemental Power Distribution (which gives power levels), Temperature Stability of Chemical Compounds (for stability at extreme temperatures), and Elemental Complexity Distribution (for complexity levels).

First, I need to identify which elements are present in all three charts. The elements mentioned in the first chart (Elemental Complexity) are A, B, C, D, E, F. The same elements are in the Elemental Power Distribution chart (A-F) and the Temperature Stability chart also has compounds, but wait, the Temperature Stability chart is for compounds (A-F), not elements. Hmm, maybe the problem is using "element" as a general term here, but the Temperature Stability chart is about compounds. Wait, the problem says "element", but the Temperature Stability chart is labeled "Chemical Compounds". Maybe there's a mix-up, but looking at the problem statement again: "which element achieves the greatest balance of high power and extreme-temperature stability relative to its inherent complexity?"

Wait, the Elemental Power Distribution chart is about "Elements" (A-F), the Temperature Stability chart is about "Chemical Compounds" (A-F), but the problem is asking about elements. Maybe the compounds here are related to the elements, but perhaps the problem is using "element" in a non-technical way, or maybe it's a mistake. But since all three charts have A-F in the elements/compounds, I'll assume that the elements A-F in the power and complexity charts correspond to the compounds A-F in the stability chart.

So, for each element (A-F), I need to get:

1. Power at "Extreme" level (from Elemental Power Distribution: "Power Levels" with "Extreme" as the fourth column? Let me check the chart. The Elemental Power Distribution has columns: Low, Medium, High, Extreme, Supernatural, Metaphysical. So "Extreme" is the 4th column. The values for each element at Extreme level:

- Element A: 72 (4th column: 92,85,78,72,68,62 → 72)
- Element B: 69 (88,82,75,69,65,59 → 69)
- Element C: 67 (84,79,73,67,63,57 → 67)
- Element D: 63 (79,74,68,63,59,54 → 63)
- Element E: 60 (75,70,65,60,56,51 → 60)
- Element F: 57 (71,66,61,57,53,49 → 57) .....

Figure 9: ChartVerse-SFT QA Example-3

Table 11: Prompt for Claude-4-Sonnet to generate code

```
Prompt for Claude-4-Sonnet to generate code

## Role
You are a matplotlib expert. Given a chart image, generate Python code that accurately reproduces this visualization.
## Inputs
Chart Image: The visualization to reproduce
## Requirements

1. Code Structure:

- Import necessary libraries
- Include clear comments explaining the code only (no words like "reproduce", "match", "original", etc.)
- Ensure code is fully executable
- DO NOT use read_csv or any file reading operations
- Define all data directly in the code as variables/lists/dictionaries
- DO NOT use plt.show(), save the figure as 'image.png' instead



2. Visual Elements to Match:

- Chart type and all data points
- Colors, styling, and proportions
- Title, axis labels, legend, and grid



3. Text Positioning:

- Ensure all text labels are clearly visible and DO NOT overlap
- Include loc='best' parameter when setting legend to avoid text conflicts
- Position each label carefully in its appropriate location
- Adjust text placement as needed to maintain readability
- Use plt.tight_layout() to optimize overall spacing



4. Output Format:

- First briefly analyze the chart type and data structure.
- Then provide complete matplotlib code:


```
```python
[Code]
```
```



## Key Points


- IMPORTANT: ensure data is reasonable(avoid nan)
- Text labels must not overlap - adjust positions to keep them readable
- Match visual style as closely as possible

```

Table 12: System Prompt for Complexity-Aware Chart Coder

```
System Prompt for Complexity-Aware Chart Coder

## Role
You are a Python visualization expert.
## Goal
Generate a random Python visualization code focusing on charts, tables, or diagrams.
## Requirements

1. Scope & Tools:
    • Choose any visualization type (chart, table, flowchart, diagram, etc.).
    • Create comprehensive sample data.
    • Use standard Python visualization libraries (e.g., matplotlib, graphviz).

2. Visual Design:
    • Ensure the visualization is visually appealing with proper labels, titles, and color schemes.
    • Include sufficient visual elements to demonstrate complexity.

3. Layout & Clarity (Critical):
    • Carefully design the layout to avoid any overlapping text or elements.
    • Adjust figure size, margins, and spacing parameters for optimal clarity.

4. Output Format:
    • Only output the Python visualization code wrapped in:
      ```python
      [Code]
      ```
```

Table 13: Prompt for Inverse Logic Construction - Step 1: Logic Architecture & Code Generation

### Step 1: Logic Architecture & Code Generation

#### ## Role

You are an expert Data Scientist specializing in **Algorithmic Synthesis**. Your task is to look at raw data arrays from visualization code and engineer sophisticated, multi-step quantitative operations without needing instructions.

#### ## Objective

You will be provided with a **Visualization Code Snippet**. Your workflow is:

1. **Decode**: Extract the raw data signals from the visualization code.
2. **Engineer**: Design a **multi-stage analysis** logic.
3. **Implement**: Write a standalone Python script that executes this logic to produce a single significant **result (a specific number or a text label)**.

#### ## Design Constraints (Difficulty Control)

- **Target Level**: The logic must be **accessible to a High School student**.
- **Linearity**: Focus on clear, linear logic, but avoid trivial operations.

#### ## Logic Architecture (The "Triangle of Complexity") You must invent a calculation path that adheres to these rules:

1. **Cross-Variable Correlation**: You must use interactions between at least three different data lists.
2. **Multi-Subplot Integration (CRITICAL)**: If the input code contains multiple subplots, your calculation logic **MUST synthesize data across at least 2 subplots**.
3. **Structural Dependency**: Do not simply process the dataset uniformly.
4. **Derived Aggregation**: The final output must be a single synthesized value (numerical scalar OR categorical label).

#### ## Output Requirements

1. **Data Context**: A brief professional summary of what the data represents.
2. **Logic Blueprint**: A step-by-step text description of the algorithm.
3. **The Python Script (The Solution)**:
  - **Self-Contained**: Define all data lists explicitly inside the script.
  - **Narrative Variable Naming**: Variable names **MUST** serve as documentation.
  - **Atomic Output**: The script must end with exactly **ONE** `print()` statement showing **ONLY** the final result.
  - Wrap the code in: `<answer>[Your Code] </answer>`

#### ## Input Data

##### **Input Code for Analysis:**

```
{chart code}
```

Table 14: Prompt for Inverse Logic Construction - Step 2: Semantic Reverse-Engineering

### Step 2: Semantic Reverse-Engineering (Question Synthesis)

## Role

You are an expert in **Reverse-Engineering** complex code logic into natural, high-level business or scientific inquiries.

## Objective

Your task is to formulate the **Analytical Question** that would prompt a human to write exactly the provided Solution Script based on the Original Visualization Code.

## Instructions

#### Phase 1: Logic Analysis

Analyze the "Solution Script" to understand *what* is being calculated. Look at the variable names to determine the **INTENT**.

#### Phase 2: Question Formulation

Write the question that requires this exact logic to solve.

#### Constraint A: Explicit Target (Crucial)

- The question must ask for a specific number OR a specific entity.
- **Avoid Vague Phrasing:** Ask for the specific metric calculated.

#### Constraint B: Strict Logic Mapping

- **The Code IS The Answer:** Ensure the provided Python script is the **exact, step-by-step solution** to your question.

#### Constraint C: Semantic Abstraction (Most Important)

- **PROHIBITION:** Do NOT describe mechanical steps (e.g., "Add list A and B").
- **REQUIREMENT:** Describe the **INTENT** and **DOMAIN MEANING** (e.g., "Calculate the profit margin").

## Output Format Structure your response strictly as follows:

1. **Logic interpretation:** Briefly explain the semantic meaning of the Python script.
2. **The Analytical Question:** The abstract, scenario-based question wrapped in <question>... </question>.

## Input Data

#### 1. Original Visualization Code:

{chart code}

#### 2. The Solution Script (The Answer Logic):

{generated\_python\_code}

Table 15: Prompt for Inverse Logic Construction - Step 3: LLM Inference

### Step 3: LLM Inference Prompt

#### ## Role

You are an expert in **Chart Code Comprehension** and **Data Reasoning**. Your task is to answer a specific question by interpreting the raw data structures defined in a Python code snippet.

#### ## Core Principles

1. **Code as Data Source:** Treat the Python code as a structured document containing the ground truth data. You do not need to "run" the code, but you must "read" and understand the variables (lists, dictionaries, values).
2. **Semantic Mapping:** Map the terms in the question to the corresponding variables and indices in the code.
3. **Logical Derivation:** Based on the extracted data, perform the necessary logical reasoning to answer the question.

#### ## Solution Framework

##### Phase 1: Code Structure Analysis

- Scan the code to identify key variables and data lists.
- Understand the relationship between different lists (e.g., `x_axis` usually corresponds to `y_axis` by index).

##### Phase 2: Information Extraction & Reasoning

- **Locate Data:** Pinpoint the specific data points (values, labels, axis limits) in the code required by the question.
- **Step-by-Step Derivation:** Execute a clear, logical reasoning process. First, explicitly quote the values found in the code. Then, perform the necessary logical comparisons or arithmetic calculations step-by-step to derive the answer.
- **Note:** Rely **ONLY** on the information explicitly present in the code.

##### Phase 3: Answer Formulation

- Formulate a clear, concise answer based on your findings.
- Ensure the answer directly addresses the inquiry.

#### ## Strict Formatting Requirements

1. The final result goes **ONLY** inside `<answer>...</answer>` tags.
2. The last line of your response must be **EXACTLY**: "Therefore, the final answer is `<answer>ANSWER</answer>`."

#### ## Input Data

##### 1. Chart Code:

{chart code}

##### 2. Question:

{generated\_question}

Table 16: Prompt for Vision-Language Model (VLM) CoT Distillation

CoT Distillation Prompt
<p><b>## Role</b>            You are an expert in science and visual reasoning with advanced capabilities in multimodal analysis. Your response will be used as a high-quality example to train a new AI model. Solve the problem efficiently and clearly by integrating <b>ALL</b> information from multimodal inputs.</p> <p><b>## Core Principles</b></p> <ol style="list-style-type: none"> <li>1. <b>Equal Weight to All Inputs:</b> Information from images is AS IMPORTANT as text. Never ignore visual elements.</li> <li>2. <b>Systematic Analysis:</b> Follow a rigorous, reproducible approach.</li> <li>3. <b>Precision and Accuracy:</b> Double-check all calculations.</li> <li>4. <b>Adaptive Reasoning:</b> Choose the most appropriate method based on context.</li> </ol> <p><b>## Solution Framework</b></p> <p><b>Phase 1: Comprehensive Information Extraction</b></p> <ul style="list-style-type: none"> <li>• Carefully analyze ALL text content for requirements and constraints.</li> <li>• Thoroughly examine ALL visual elements, extracting every piece of relevant information.</li> <li>• Explicitly connect visual and textual information.</li> </ul> <p><b>Phase 2: Strategic Problem Setup</b></p> <ul style="list-style-type: none"> <li>• Compile information and clearly state what needs to be found.</li> <li>• Identify relevant principles and state assumptions.</li> </ul> <p><b>Phase 3: Rigorous Solution Execution</b></p> <ul style="list-style-type: none"> <li>• Present solution with complete logical flow and proper notation.</li> <li>• Reference specific parts of visual inputs to support reasoning.</li> <li>• Maintain unit consistency and precision.</li> </ul> <p><b>Phase 4: Solution Validation</b></p> <ul style="list-style-type: none"> <li>• Verify the answer makes scientific/logical sense.</li> <li>• Ensure dimensional analysis is correct.</li> </ul> <p><b>## Strict Formatting Requirements</b></p> <ol style="list-style-type: none"> <li>1. The final result goes ONLY inside &lt;answer&gt;...&lt;/answer&gt; tags.</li> <li>2. Include units inside the tags when required (e.g., &lt;answer&gt;5.2 m/s&lt;/answer&gt;).</li> <li>3. The last line must be EXACTLY: "Therefore, the final answer is &lt;answer&gt;ANSWER&lt;/answer&gt;."</li> </ol> <p><b>## Input Data</b>  <b>Problem:</b>            {question}  <i>(Visual Inputs are provided via the vLLM interface)</i></p>