

PLOT CRAFT: PUSHING THE LIMITS OF LLMs FOR COMPLEX AND INTERACTIVE DATA VISUALIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent Large Language Models (LLMs) have demonstrated remarkable proficiency in code generation. However, their ability to create complex visualizations for scaled and structured data remains largely unevaluated and underdeveloped. To address this gap, we introduce **PlotCraft**, a new benchmark featuring 1k challenging visualization tasks that cover a wide range of topics, such as finance, scientific research, and sociology. The benchmark is structured around seven high-level visualization tasks and encompasses 48 distinct chart types. Crucially, it is the first to systematically evaluate both single-turn generation and multi-turn refinement across a diverse spectrum of task complexities. Our comprehensive evaluation of 23 leading LLMs on PlotCraft reveals obvious performance deficiencies in handling sophisticated visualization tasks. To bridge this performance gap, we develop **SynthVis-30K**, a large-scale, high-quality dataset of complex visualization code synthesized via a collaborative agent framework. Building upon this dataset, we develop **PlotCrafter**, a novel code generation model that achieves strong capabilities in complex data visualization with a remarkably small size. Across VisEval, PandasPlotBench, and our proposed PlotCraft, PlotCrafter shows performance comparable to that of leading proprietary approaches. Especially, on hard task, Our model achieves over 50% performance improvement. We will release the benchmark, dataset, and code at <https://anonymous.4open.science/r/PlotCraft-E320>.

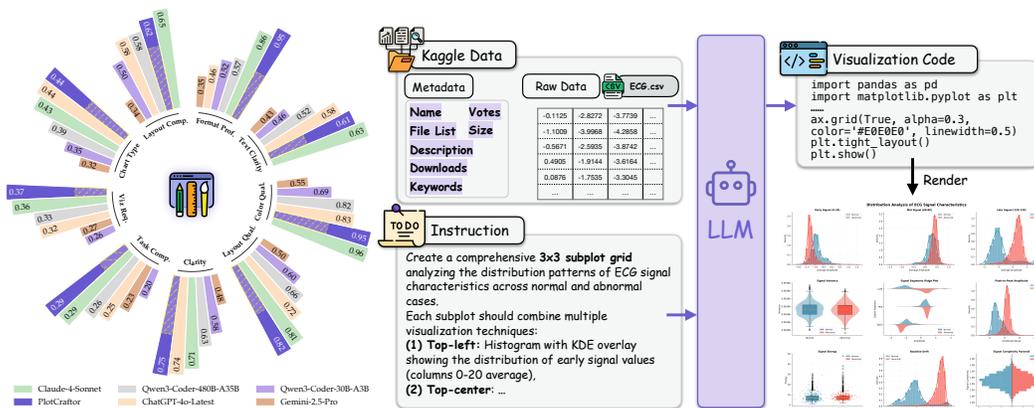


Figure 1: An overview of the PlotCraft benchmark and the performance of several leading LLMs and our model, PlotCrafter. (Left) A polar bar chart compares the performance of PlotCrafter against five leading baseline models across all of our proposed sub-metrics. The purple area explicitly highlights the performance gains of PlotCrafter relative to its base model, Qwen3-Coder-30B-A3B. (Right) An example task from PlotCraft, which requires an LLM to process raw Kaggle data and a complex, human-written instruction to generate visualization code, which is then rendered into the final chart. PlotCraft benchmark comprises 1k high-quality evaluation instances.

1 INTRODUCTION

Recent advancements in AI research have demonstrated the powerful code generation capabilities of LLMs (OpenAI, 2023; 2025; Anthropic, 2023; Team, 2024; Rozière et al., 2023; Hui et al., 2024; MistralAI, 2024; Team et al., 2025b), which have solved coding challenges in domains such as software engineering (Jimenez et al., 2023; Zhang et al., 2025; Pan et al., 2025b), code completion Ding et al. (2023); Yang et al. (2024a); Gong et al. (2024), and algorithmic problem-solving (Chen et al., 2021a; Zhuo et al., 2025; Jain et al., 2024). However, in the domain of data visualization, the capabilities of LLM have yet to be fully explored. We observe that while LLMs excel at creating simple, single-panel charts, they struggle to generate plots with multiple subplots and intricate composite layouts from large, complex structured data. As illustrated in **Figure 2**, when faced with these complex plotting tasks, LLMs often generate code that results in chaotic layouts, overlapping elements, and obscured axis labels, or they completely ignore instructions about certain areas of plots.

Prior works on chart generation have primarily focused on relatively simple, text to visualization (Text2Vis) tasks (e.g., VisEval (Chen et al., 2024)/PandasPlotBench (Galimzyanov et al., 2025)) or on understanding and reproducing existing data charts (Chart2Code) (e.g., Plot2Code (Wu et al., 2024), ChartMimic (Yang et al., 2025b)). However, these efforts do not adequately assess the ability of LLMs to synthesize complex visualization code from structured raw data. To systematically evaluate these capabilities of LLMs, we introduce **PlotCraft**, a large-scale benchmark comprising 982 instances. PlotCraft is characterized by its use of instructional prompts of varying difficulties, a diverse range of chart types, and multi-level evaluation metrics. Specifically, we design and collect tasks of varying difficulty based on the number of subplots to be generated, chart complexity, data volume, and the level of detail in the instructions. Furthermore, we introduce an additional multi-turn refinement task, which allows the model to refine the chart over multiple conversational turns, thereby assessing its ability to debug code and perform iterative optimization based on user feedback. **Figure 1** provides an overview of our benchmark.

We evaluate 23 prominent LLMs on PlotCraft benchmark, including 6 proprietary models and 17 open-weight models. We observe that most models perform well on simple chart generation tasks, but fail on medium and hard tasks, which require processing larger data and generating composite plots with multiple, properly arranged chart types. Results show that existing LLMs have limited capabilities in complex scientific visualization code generation tasks. To further validate our findings, we score a set of results with human evaluation and a correlation analysis (**Section 5.3**) demonstrates a high correlation between our multi-level metrics and human evaluation.

To address the limited visualization code generation capabilities, we construct a new dataset **SynthVis-30K**, and a novel model **PlotCrafter**. Specifically, we collect 30k multi-modal chart data instances covering 31 topics, 48 chart types, and 7 tasks. Each data instance contains a human instruction, some data files for visualization (CSV/XLSX), the resulting chart image, and the corresponding source code. A detailed breakdown of the dataset’s coverage and the curation pipeline is presented in Figure 3. Based on SynthVis-30K, we construct PlotCrafter that achieves strong capabilities in complex data visualization with a minimal model size.

Experimental results demonstrate that PlotCrafter improves performance on PlotCraft by 25%, achieving performance comparable to leading proprietary LLMs. Concurrently, we adapt other benchmarks to our task settings, including VisEval (Chen et al., 2024), PandasPlotBench (Galimzyanov et al., 2025). PlotCrafter demonstrate a 7% higher performance compared to the much larger Qwen3-Coder-480B. These results validate PlotCrafter’s powerful abilities in different text-to-visualization code generation tasks.

2 RELATED WORKS

2.1 CODE GENERATION

Recent advances in Large Language Models (LLMs), including general-purpose models (e.g., GPT (OpenAI, 2023), Claude (Anthropic, 2023), Gemini (Team, 2024)) and specialized code models (e.g., Qwen-Coder (Hui et al., 2024), DeepSeek-Coder (Guo et al., 2024), Codestral (MistralAI, 2024)), have demonstrated powerful coding capabilities (DeepSeek-AI & etc., 2024; Rozière et al., 2023; Team et al., 2025b;a). While conventional tasks like algorithmic problem-solving (Chen et al.,

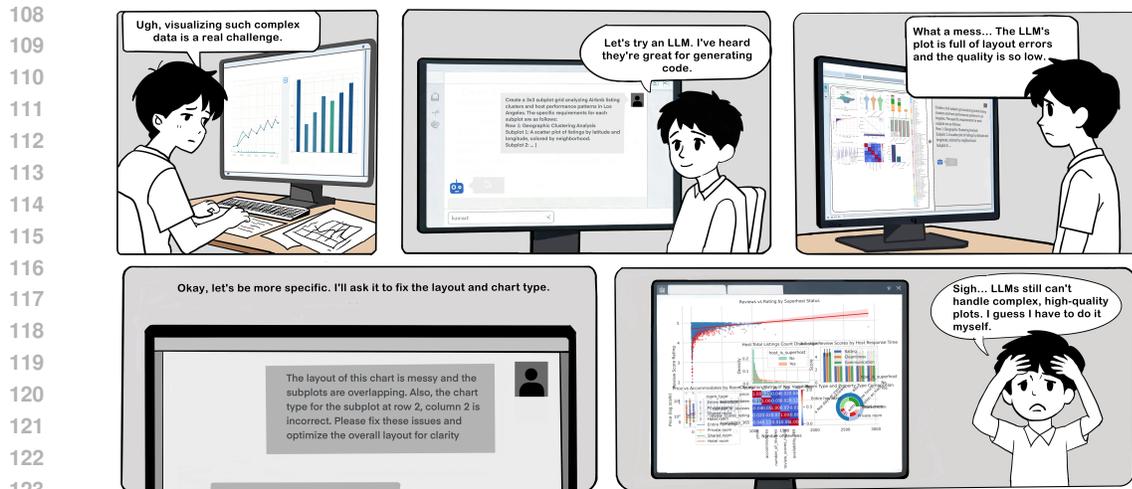


Figure 2: A real-world example illustrating the limitations of LLMs on complex visualization tasks. When presented with a sophisticated request, the model generates a low-quality output and struggles to make effective improvements during the subsequent refinement process.

2021a; Zhuo et al., 2025) and software engineering (Jimenez et al., 2023; Zhang et al., 2025) are evaluated on functional correctness (Chen et al., 2021b), this paper focuses on data visualization, a domain where generated code must produce a visually accurate output, a requirement shared by front-end design (Xu et al., 2025; Lu et al., 2025) and SVG generation (Xing et al., 2025).

2.2 DATA VISUALIZATION

Prior LLM-based data visualization research spans three areas: The first, chart understanding (e.g., QA and captioning), focuses on interpreting visual information from plots, such as for question answering or summary generation (Li et al., 2024; Zeng et al., 2024; Rahman et al., 2023; Kantharaj et al., 2022; Jia et al., 2025). The second, Chart-to-Code (Chart2Code) generation, involves reverse-engineering a visualization by generating the code required to replicate it (Wu et al., 2024; Yang et al., 2025b; Zhao et al., 2025), and the third, Text-to-Visualization (Text2Vis), concerns generating visualization specifications or code from natural language descriptions (Luo et al., 2025; Galimzyanov et al., 2025; Ni et al., 2025). However, these works predominantly focus on simple, single-panel plots, offering limited differentiation for advanced LLMs’ ability to handle complex layouts and high information density. We introduce PlotCraft to address this critical evaluation gap.

3 THE PLOT CRAFT BENCHMARK

3.1 TASK DEFINITION

We define the data visualization task as a conditional code generation problem. Formally, given a natural language instruction I , metadata M , and a raw dataset D , a Large Language Model (\mathcal{F}) must generate an executable code snippet C . This code must use D to render a visualization that satisfies all requirements outlined in I . This task is formulated as $C = \mathcal{F}(I, M, D)$. The PlotCraft benchmark evaluates models on two distinct variants of this core task.

The two variants are: **(1) Single-Turn Generation**, where the model must generate the complete visualization code in a single step from the initial instruction, thereby measuring its ability to execute complex requests from scratch; and **(2) Multi-Turn Refinement**, which assesses a model’s ability to debug and iteratively improve existing code based on a conversation history and a new modification request, mirroring a realistic user interaction workflow.

3.2 BENCHMARK COVERAGE ANALYSIS

The PlotCraft benchmark comprises 982 evaluation instances, evenly divided into 491 for single-turn generation and 491 for multi-turn refinement. Both single-turn and multi-turn tasks are based on the

Table 1: A comparison of PlotCraft with existing benchmarks. A ✓ indicates the presence of a feature, while a ✗ indicates its absence.

Benchmarks	# of Test Instances	Composite Types	Multiple Subplots	Multi-Turn	Evaluation Metric
<i>Chart Understanding Benchmarks</i>					
ChartQA (Masry et al., 2022)	10K	✗	✗	✗	Accuracy
ChartSumm (Rahman et al., 2023)	84K	✗	✗	✗	Match-based
CharArXiv (Wang et al., 2024)	93K	✗	✓	✗	MLLM Score
ChartX (Xia et al., 2025)	1152	✓	✗	✗	Multi-Level
<i>Chart to Code Benchmarks</i>					
Plot2Code (Wu et al., 2024)	132	✗	✓	✗	MLLM Score
Design2Code (Si et al., 2024)	484	✗	✗	✗	Multi-Level
ChartMimic (Yang et al., 2025b)	4,800	✓	✓	✗	Multi-Level
<i>Text to Visualization Benchmarks</i>					
MatPlotBench Yang et al. (2024b)	100	✗	✗	✗	MLLM Score
VisEval Chen et al. (2024)	2300	✗	✗	✗	Multi-Level
PandasPlotBench Galimzyanov et al. (2025)	150	✗	✗	✗	Multi-Level
PlotCraft (Ours)	982	✓	✓	✓	Multi-Level

same underlying visualization goal. We describe the benchmark’s coverage across three dimensions: Chart Types, Thematic Topics, and Task Complexity (the same as shown in Figure 3).

Chart Types Tasks in PlotCraft are categorized according to a high-level Visualization Intent Taxonomy, encompassing seven primary classes: Correlation, Deviation, Ranking, Distribution, Composition, Change, and Groups. These categories contain 48 distinct plot types, including scatter plots, bubble plots, pie charts, dendrograms, violin plots, and so on. As the majority of tasks require the combination of multiple plot types within a single figure, a single, exclusive plot type label is often not applicable. We provide illustrative examples for each category in Appendix A.1.

Thematic Coverage PlotCraft spans a diverse range of thematic topics to ensure broad applicability. The benchmark is organized into eight high-level domains: Finance (Business & Finance), Industry (Industry & E-commerce), Health (Health & Medicine), Research (Science & Research), Society (Government & Society), Media (Culture & Media), Tech. (Technology & Computing), and Env. (Environment & Geospatial). These domains are further divided into 31 fine-grained sub-topics, a detailed breakdown of which is provided in Appendix A.2.

Task Complexity To systematically evaluate model capabilities, tasks are stratified into three levels of complexity: (1) Easy: The task requires generating a single, standard chart (e.g., a bar chart or a line plot) to visualize the data. (2) Medium: The task involves creating a composite visualization, which could be either a single chart combining multiple plot types or a subplot grid where each individual subplot is of a simple nature. (3) Hard: The task demands the creation of a complex subplot grid where each individual subplot is itself a composite chart, requiring advanced spatial and logical reasoning. We provide illustrative examples for different complexity in Appendix A.3.

3.3 DATA CURATION PROCESS

The construction of PlotCraft adheres to four guiding principles: (1) Grounded in Real Data: All tasks are based on real-world datasets to avoid artifacts of synthetic data. (2) Built from Scratch: All tasks and reference solutions are newly created to prevent data leakage from existing sources. (3) Zero-Reference Generation: Tasks provide no sample images or code, requiring models to generate visuals from abstract instructions. (4) Compositional Complexity: The benchmark spans a wide spectrum of complexities, including tasks with intricate layouts and multiple chart types within a single figure. Adhering to these principles, we curate the benchmark through a four-step pipeline, which we overview here. Further details are provided in Appendix B.

Data Filtering We source open-source datasets from Kaggle. The use of these datasets in our work is for scientific research purposes only; other uses are subject to their original licenses. Using metrics such as vote counts, download counts, and usability ratings, we performed an initial screening of 7,162 datasets, comprising over 95,000 files and 25.6 billion data rows. From this pool, we conduct a second filtering stage to select 140 core datasets for the benchmark. This selection was curated to ensure diverse topics, a wide distribution of data volume and complexity, and minimal overlap with datasets used in other visualization tasks to mitigate data leakage. The final collection of benchmark datasets contains 1,874 raw data files and approximately 462 million data rows.

Task and Instruction Writing Using the selected benchmark datasets, we authored 491 unique visualization tasks. The design of these tasks was guided by a combination of seven high-level visualization intents and 48 distinct chart types. We ensure a roughly balanced distribution of tasks across three complexity levels: simple (159), medium (163), and hard (169). The accompanying instructions are written to be abstract, specifying only the visualization requirements without providing any guidance on code implementation.

Reference Code Writing For each task, a reference solution is implemented by a team of five senior Python developers using Matplotlib and its associated libraries. It is important to note that this code serves as a valid reference implementation that achieves a high-quality result, rather than the single, optimal solution. The details of the sandboxed execution environment used for our evaluation are provided in Appendix A.4.

Multi-turn Conversation Synthesis To simulate realistic refinement scenarios, we first generate an initial code draft for each task using a less capable model. Human experts then review these drafts and intentionally introduced common, realistic errors—such as incorrect chart types or overlapping visual components—to create a faulty, low-quality implementation. Examples of such faulty code are provided in Appendix B.4. This faulty code, along with its rendered image and the original instruction, was presented to human annotators who then wrote a natural language modification request to correct the errors. Each multi-turn task instance is thus composed of the original instruction, data metadata, the faulty code, and this human-authored refinement prompt. This process yield 491 multi-turn tasks, which, combined with the 491 single-turn tasks, constitute the complete PlotCraft benchmark of 982 instances.

3.4 EVALUATION METRICS

Our evaluation pipeline first assesses functional correctness. All generated code must execute successfully in our sandboxed environment and produce a valid, non-empty image to be eligible for further analysis. Any code that fails this initial step receives a score of zero. For all valid outputs, we follow established practices and employ Gemini-2.5-Pro as an automated visual judge. [We acknowledge that relying on an automated judge, especially a proprietary model, introduces inherent limitations and potential biases in visual and aesthetic judgment \(Xie et al., 2025; Pan et al., 2025a; Shi et al., 2025\).](#) This judge scores the generated charts based on a two-dimensional framework: **Task Compliance** and **Chart Quality**. To ensure a granular assessment, Task Compliance is decomposed into four distinct sub-metrics evaluated on a binary scale (0=fail, 1=pass): Layout Compliance, Chart Type Compliance, Requirement Fulfillment, and Complete Task Fulfillment. Similarly, Chart Quality is assessed on a 3-point scale (0, 1, or 2) across five criteria: Clarity (the absence of overlapping elements), Layout Quality, Color Quality, Text Readability, and Professional Formatting. The detailed scoring rubrics for all metrics, along with the specific prompts used to query our automated judge and specific judgement cases, are available in Appendix C.

4 THE PLOTRAFTOR

To address the performance gap observed in our benchmark evaluations, we develop PlotCrafter, a model specifically fine-tuned for complex data visualization. The development process consists of two key stages: the creation of a large-scale, high-quality synthetic dataset, which we call SynthVis-30K, and the subsequent Supervised Fine-Tuning (SFT) of a base model on this data.

4.1 SYNTHVIS-30K DATASET

The SynthVis-30K dataset was created using a multi-agent framework, illustrated in Figure 3, which comprises two main stages: **Task Generation** and **Code Generation**. Details are in Appendix E.

Task Generation. We first sourced a collection of Kaggle datasets ([with CC BY 4.0 or Apache 2.0 licenses](#)), ensuring no overlap with those used in the PlotCraft benchmark to prevent data contamination. A Data Analyzer agent processes each dataset to extract structured metadata, including data types, column names, and sample rows. This formatted data then enters a Task Cycle, an iterative loop between a Task Generator and a Task Judge. The Task Generator, prompted with few-shot examples from PlotCraft, proposes 3-6 visualization tasks of varying difficulty for a given dataset. The

Table 2: **Quantitative results on PlotCraft for 16 primary LLMs across two settings: Single-Turn Generation and Multi-Turn Refinement.** Task-Comp. and Quality denote the total scores for the Task Compliance (out of a maximum of 4) and Chart Quality (out of a maximum of 10) sub-metrics, respectively. AVG score is the average of these scores across both single-turn and multi-turn evaluations. Within each model category, the best score is **bolded** and the second-best is underlined.

Model	Single-Turn Generation			Multi-Turn Refinement			AVG score
	Pass Rate (%)	Task-Comp.	Quality	Pass Rate (%)	Task-Comp.	Quality	
<i>Closed-source LLMs</i>							
Claude-4.1-Opus (Anthropic, 2023)	76.20	1.93	4.20	81.44	2.05	5.22	6.70
Claude-4-Sonnet (Anthropic, 2023)	68.84	1.73	<u>3.99</u>	<u>78.41</u>	<u>1.88</u>	<u>4.80</u>	<u>6.20</u>
Gemini-2.5-Pro (Team, 2024)	41.34	1.15	2.31	58.86	1.51	3.80	4.39
ChatGPT-4o-Latest (OpenAI, 2023)	63.54	1.60	3.33	69.25	1.51	4.23	5.33
GPT-5 (OpenAI, 2025)	<u>69.86</u>	<u>1.76</u>	<u>2.87</u>	74.13	1.80	4.33	5.38
Grok-4 (xAI, 2025)	64.52	1.63	3.66	70.24	1.61	4.42	5.65
<i>Open-source LLMs</i>							
Kimi-K2 (Team et al., 2025b)	60.13	1.52	<u>3.36</u>	61.03	1.49	4.05	5.21
DeepSeek-V3.1 (DeepSeek-AI & etc., 2024)	55.71	1.49	3.20	56.86	1.50	3.99	5.09
GLM-4.5 (Team et al., 2025a)	43.38	1.25	2.48	64.97	1.54	3.91	4.59
GPT-oss-120B (Agarwal et al., 2025)	48.23	1.32	2.68	29.69	0.77	1.87	3.32
GPT-oss-20B (Agarwal et al., 2025)	44.68	1.21	2.43	34.02	0.89	2.14	3.34
Seed-Coder-8B (Seed et al., 2025)	32.38	0.87	1.74	57.03	1.26	3.22	3.55
VisCoder-7B (Ni et al., 2025)	25.46	0.67	1.50	51.73	0.99	2.96	3.06
Qwen3-Coder-480B-A35B (Hui et al., 2024)	<u>61.30</u>	<u>1.56</u>	3.21	<u>75.97</u>	<u>1.75</u>	<u>4.46</u>	<u>5.49</u>
Qwen3-Coder-30B-A3B (Hui et al., 2024)	52.55	1.32	2.85	73.12	1.55	4.18	4.95
PlotCrafter-30B-A3B (Ours)	64.36	1.73	4.09	77.11	1.76	4.74	6.16

Table 3: Quantitative results for 16 LLMs across two benchmarks, VisEval and PandasPlotBench.

Model	VisEval					PandasPlotBench		
	Invalid Rate (%)	Illegal Rate (%)	Pass Rate (%)	Readability	Quality	Pass Rate (%)	Vis	Task
<i>Closed-source LLMs</i>								
Claude-4.1-Opus (Anthropic, 2023)	0.87	5.42	96.14	4.39	3.91	100	78	95
Claude-4-Sonnet (Anthropic, 2023)	1.82	6.78	<u>94.51</u>	<u>4.27</u>	<u>3.85</u>	98.3	75	92
Gemini-2.5-Pro (Team, 2024)	6.79	15.83	83.71	3.76	3.09	85.1	66	77
ChatGPT-4o-Latest (OpenAI, 2023)	3.39	14.54	84.47	4.01	3.34	93.7	72	89
GPT-5 (OpenAI, 2025)	5.01	12.98	87.71	4.22	3.52	82.9	66	79
Grok-4 (xAI, 2025)	4.02	8.53	94.41	4.30	3.69	94.6	<u>75</u>	88
<i>Open-source LLMs</i>								
Kimi-K2 (Team et al., 2025b)	2.45	8.91	<u>93.62</u>	<u>4.24</u>	<u>3.78</u>	95.4	76	91
DeepSeek-V3.1 (DeepSeek-AI & etc., 2024)	5.94	9.82	92.31	4.17	3.59	93.7	74	86
GLM-4.5 (Team et al., 2025a)	5.03	13.07	88.54	4.19	3.48	47.4	39	44
GPT-oss-120B (Agarwal et al., 2025)	3.91	14.88	90.13	4.03	3.35	92.6	71	88
GPT-oss-20B (Agarwal et al., 2025)	5.23	15.12	86.34	3.95	3.28	85.7	69	87
Seed-Coder-8B (Seed et al., 2025)	8.15	18.25	75.40	3.45	2.81	54.3	65	80
VisCoder-7B (Ni et al., 2025)	6.45	16.05	82.88	3.81	3.12	87.7	66	78
Qwen3-Coder-480B-A35B (Hui et al., 2024)	2.89	11.53	91.75	4.18	3.55	97.1	73	87
Qwen3-Coder-30B-A3B (Hui et al., 2024)	3.55	13.90	89.67	4.10	3.41	94.3	71	86
PlotCrafter-30B-A3B (ours)	2.61	9.76	93.81	4.26	3.80	96.0	74	91

5 EXPERIMENTS

5.1 EXPERIMENTS SETUP

Baseline Setup We evaluate a total of 24 models, including our proposed **PlotCrafter** and 23 widely-used proprietary and open-source Large Language Models. For proprietary models, we select six representative systems: Anthropic’s flagship models, Claude-4.1-Opus and Claude-4-Sonnet (Anthropic, 2023); OpenAI’s most advanced models, GPT-5 and GPT-4o (latest version) (OpenAI, 2025); Google’s Gemini 2.5 Pro (Team, 2024); and xAI’s Grok-4 (xAI, 2025). For open-weight models, we choose nine competitive models with parameter sizes ranging from 7B to 1T: Kimi-K2 (Team et al., 2025b) (1T, A32B), Deepseek-V3.1 (DeepSeek-AI & etc., 2024) (671B, A37B), GLM-4.5 (Team et al., 2025a) (355B, A32B), GPT-oss (Agarwal et al., 2025) (120B, A5.1B; 20B, A3.6B), Qwen3-Coder (Hui et al., 2024) (480B, A22B; 30B, A3B), Qwen2.5-Coder (Hui et al., 2024), DeepSeek-Coder (Guo et al., 2024), and Qwen3 (Yang et al., 2025a). Additionally, we include Seed-Coder-8B (Seed et al., 2025) and Viscoder-7B (Ni et al., 2025) as baseline models of a comparable size to our PlotCrafter (30B, A3B) for a more direct comparison.

Benchmark Setup To ensure a comprehensive evaluation, in addition to our primary PlotCraft benchmark, we assessed model performance on two other established text2vis benchmarks. The first, VisEval (Chen et al., 2024), is a large-scale benchmark composed of high-quality visualization

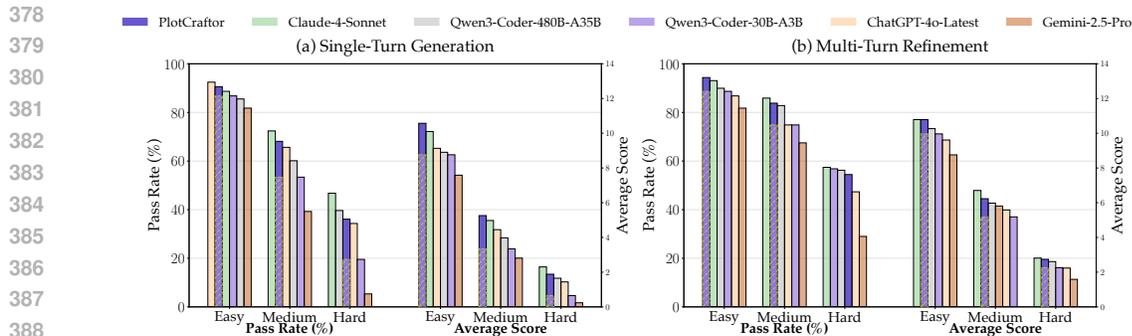


Figure 4: Performance comparison of PlotCrafter and five leading LLMs on tasks of varying difficulty. The figure is split into two subplots: (a) Single-Turn Generation and (b) Multi-Turn Refinement. Within each subplot, we report the Pass Rate (%) and Average Score for tasks categorized as Easy, Medium, and Hard. The yellow hatched area within the PlotCrafter bars indicates the score contribution from its base model, Qwen3-Coder-30B-A3B.

tasks curated from nvBench (Luo et al., 2025). It employs a multi-faceted evaluation protocol that includes execution pass rate, GPT-4 based scoring, and SVG based layout checks. The second, PandasPlotBench (Galimzyanov et al., 2025), is a benchmark containing 175 test cases specifically designed to evaluate code generation for Matplotlib and its related libraries.

Evaluation Details All experiments were conducted by sending requests in the standard OpenAI API format, with the chat structure following the ChatML format (OpenAI, 2022). We employed a greedy decoding strategy and set the maximum output length to 131,072 tokens; if a model did not support this length, we used its specified maximum limit. The final reported results are the average of three experimental runs. For proprietary models, we queried their official APIs directly. For open-weight models, we utilized the vLLM framework for serving. It is important to note a potential limitation regarding the GPT-oss series of models: as they are trained on the harmony format (Agarwal et al., 2025), which natively incorporates multi-turn reasoning and tool invocation, our standardized ChatML-based setup may not fully leverage their intended capabilities. Additional evaluation details and specific prompts are in Appendix F.

5.2 MAIN RESULTS

This section presents the results of 16 leading LLMs on the PlotCraft, VisEval, and PandasPlotBench benchmarks. Performance on PlotCraft is detailed in Table 2, with complete 24 models results available in Appendix Table 7. Results for VisEval and PandasPlotBench are shown in Table 3. For a granular analysis, Figures 1 and 4 break down performance by sub-metrics and task difficulty, respectively. The key findings are as follows.

Claude-4.1-Opus leads proprietary models, while PlotCrafter is the top-performing open-weight model, achieving results comparable to Claude-4-Sonnet. Among proprietary models, Claude-4.1-Opus demonstrates a significant advantage across all tasks. It achieved an average score of 6.70 on PlotCraft and a perfect 100% pass rate on PandasPlotBench, substantially outperforming other models.

Among open-weight models, our PlotCrafter achieves SOTA performance on PlotCraft. It surpasses its base model, Qwen3-Coder-30B-A3B, by a margin of 25% and nearly matches the performance of the top-tier Claude 4 Sonnet, with a score difference of less than 1%. PlotCrafter also secures leading results among open models on the VisEval and PandasPlotBench benchmarks. An analysis of performance by task difficulty, focusing on 6 key models (4 widely-used coding

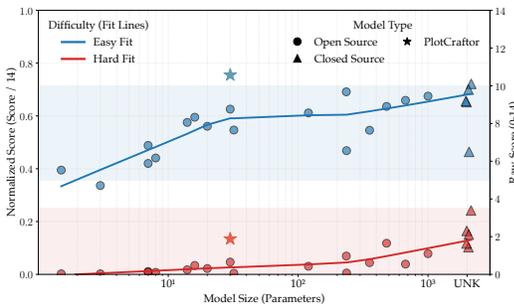


Figure 5: Performance scaling on Easy vs. Hard.

Table 4: Cohen’s Kappa scores for agreement between our Gemini-2.5-Pro judge and human evaluations, categorized by Compliance and Quality metrics.

Compliance Metrics				Quality Metrics				
Layout	Type	Visual	Task	Clarity	Layout	Color	Text	Format
0.90	0.78	0.72	0.80	0.73	0.71	0.69	0.65	0.73

models, our baseline, and PlotCrafter), reveals that while GPT-4o excels on simple, single-turn tasks, Claude 4.1 Opus shows superior performance on more complex and demanding tasks. A more detailed discussion and comparison between different models is available in Appendix G.

Strong performance on simple benchmarks does not guarantee success on complex visualization tasks. We observe a significant performance disparity between models on different benchmarks. Models such as Kimi-K2 and GPT-4o, which perform exceptionally well on the simpler tasks found in VisEval and PandasPlotBench, experience a notable performance degradation on the complex, multi-faceted tasks within PlotCraft. In contrast, the Claude series and our PlotCrafter demonstrate consistently high performance across all three benchmarks, indicating their robust capability in handling both simple and sophisticated visualization challenges.

Task difficulty significantly impacts the effectiveness of model scaling and fine-tuning. As shown in Figure 5, the benefits of model scaling are highly dependent on task difficulty. For models under 100B parameters, performance on Easy tasks scales rapidly with size, while performance on Hard tasks remains flat, only improving for models beyond the 100B threshold. This disparity is mirrored in supervised fine-tuning (SFT): smaller models can be fine-tuned to near-proprietary levels on Easy tasks, yet SFT provides minimal benefit for Hard tasks. This indicates that solving complex visualization challenges may rely more on the emergent reasoning abilities that come with scale than on task-specific fine-tuning.

5.3 CORRELATION WITH HUMAN EVALUATION

To validate the reliability of our proposed multi-dimensional metrics, we measured their correlation with human expert judgments. We conducted a scoring experiment on 500 charts randomly sampled from the output of PlotCrafter. For each chart, three human annotators provided scores for all metrics, with the final ground-truth score determined by a majority vote. We then calculated the Cohen’s Kappa coefficient to assess the inter-rater agreement between these human-derived scores and the automated scores provided by our Gemini-2.5-Pro based judge. As presented in Table 4, the average Kappa score across all evaluation metrics reached a substantial level of agreement, demonstrating the reliability of our automated evaluation framework. A comparative analysis of other models as evaluators is detailed in Appendix H, where we found that models such as Claude and GPT-4o were significantly less effective at replicating human scoring patterns.

5.4 ERROR ANALYSIS

To understand the primary failure modes on the PlotCraft benchmark, we analyzed common errors and categorized them into three main types. (1) **Code-Level Errors**, includes functionally flawed code that causes runtime exceptions or critical rendering failures, such as using non-existent object attributes or creating conflicts between incompatible layout managers (e.g., `constrained.layout` with `fig.legend()`). (2) **Task Compliance Failures**, occurs when the visual output disregards explicit instructions, such as generating an incorrect subplot layout or using the wrong chart types. (3) **Deficiencies in Chart Quality**, refers to aesthetic issues that degrade the visualization’s readability and professional appearance, including element overlap, sub-optimal color choices, and incorrect text. A more detailed, case-by-case analysis of these errors is provided in Appendix I.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

6 CONCLUSION AND FUTURE DIRECTION

In this work, we addressed the significant gap in the ability of Large Language Models to handle complex data visualization. We introduced PlotCraft, a novel and challenging benchmark designed to evaluate both single-turn generation and multi-turn refinement capabilities. Our comprehensive evaluation over 23 leading models on PlotCraft demonstrated that current LLMs struggle with sophisticated visualization tasks. To bridge this performance gap, we developed the SynthVis-30K dataset and used it to train PlotCrafter, a novel and efficient model. Experimental results show that PlotCrafter achieves state-of-the-art performance among open-weight models, delivering results comparable to leading proprietary LLMs. Meanwhile, we observed that there are still limitations in data generation and evaluation. The synthesis of high-quality data relies on costly multi-agent systems, and MLLM-based judges may fail to detect subtle visual errors, like minor overlaps or color style. This might be an area for future research. [Furthermore, while PlotCraft focuses on detailed instructions to rigorously test capability limits, exploring model performance under abstract or open-ended goals remains a promising direction for future work.](#)

7 ETHICS STATEMENT

The data used in the PlotCraft benchmark is sourced exclusively from public repositories that are governed by licenses permitting their use in software and research. Our contributions fully adhere to the terms of these licenses. We did not use any data beyond what is publicly available and downloadable via the standard Kaggle API. Our work did not involve the participation of any human subjects; we did not use crowdsourcing or recruit any external human workers for any part of the PlotCraft benchmark’s creation. All work, including the environment configuration, data curation, synthetic data generation, and the writing of this paper, was conducted entirely by the author team.

8 REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our work, we provide a complete codebase with detailed instructions for replicating the PlotCraft benchmark results and the PlotCrafter training process. The evaluation framework for PlotCraft is described in Section 3.4, with the full scoring rubrics and judge prompts available in Appendix C. The training methodology and hyperparameters for PlotCrafter are detailed in Section 4.2, with additional specifics on the SynthVis-30K data synthesis process provided in Appendix E. The complete specifications of our sandboxed execution environment are detailed in Appendix A.4. To further facilitate community engagement and standardized evaluation, we plan to release a PyPI package and host a public leaderboard for the benchmark.

9 LLM USAGE

The use of Large Language Models (LLMs) in this work was limited. A multimodal large model was employed to generate the layout for some of the images shown within the monitor displays in Figure 2. Additionally, an LLM was used sparingly for writing assistance.

REFERENCES

- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*, 2025.
- Anthropic. Introducing Claude, 2023. URL <https://www.anthropic.com/index/introducing-claude>.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021a. URL <https://arxiv.org/abs/2107.03374>.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021b.
- Nan Chen, Yuge Zhang, Jiahang Xu, Kan Ren, and Yuqing Yang. Viseval: A benchmark for data visualization in the era of large language models, 2024. URL <https://arxiv.org/abs/2407.00981>.
- DeepSeek-AI and etc. Deepseek-v3 technical report, 2024. URL <https://arxiv.org/abs/2412.19437>.

- 594 Yangruibo Ding, Zijian Wang, Wasi Uddin Ahmad, Hantian Ding, Ming Tan, Nihal Jain, Murali Kr-
595 ishna Ramanathan, Ramesh Nallapati, Parminder Bhatia, Dan Roth, and Bing Xiang. Cross-
596 codeeval: A diverse and multilingual benchmark for cross-file code completion. In Alice Oh,
597 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Ad-
598 vances in Neural Information Processing Systems 36: Annual Conference on Neural Information
599 Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- 600 Timur Galimzyanov, Sergey Titov, Yaroslav Golubev, and Egor Bogomolov. Drawing pandas: A
601 benchmark for llms in generating plotting code, 2025. URL [https://arxiv.org/abs/
602 2412.02764](https://arxiv.org/abs/2412.02764).
- 603 Linyuan Gong, Sida Wang, Mostafa Elhoushi, and Alvin Cheung. Evaluation of llms on syntax-
604 aware code fill-in-the-middle tasks. In *Forty-first International Conference on Machine Learning,
605 ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- 607 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao
608 Bi, Y Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming-
609 the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024. URL [https://arxiv.
610 org/abs/2401.14196](https://arxiv.org/abs/2401.14196).
- 611 Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,
612 Bowen Yu, Kai Dang, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*,
613 2024.
- 615 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
616 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
617 evaluation of large language models for code, 2024. URL [https://arxiv.org/abs/
618 2403.07974](https://arxiv.org/abs/2403.07974).
- 619 Caijun Jia, Nan Xu, Jingxuan Wei, Qingli Wang, Lei Wang, Bihui Yu, and Junnan Zhu. Chartrea-
620 soner: Code-driven modality bridging for long-chain reasoning in chart question answering, 2025.
621 URL <https://arxiv.org/abs/2506.10116>.
- 622 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
623 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint
624 arXiv:2310.06770*, 2023.
- 625 Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, et al. Chart-to-text: A large-scale bench-
626 mark for chart summarization. *arXiv preprint arXiv:2203.06486*, 2022.
- 627 Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim,
628 Sungyoung Ji, Byungju Lee, Xifeng Yan, et al. Mmsci: A multimodal multi-discipline dataset for
629 phd-level scientific comprehension. *arXiv preprint arXiv:2407.04903*, 2024.
- 630 Zimu Lu, Yunqiao Yang, Houxing Ren, Haotian Hou, Han Xiao, Ke Wang, Weikang Shi, Aojun
631 Zhou, Mingjie Zhan, and Hongsheng Li. Webgen-bench: Evaluating llms on generating interac-
632 tive and functional websites from scratch, 2025. URL [https://arxiv.org/abs/2505.
633 03733](https://arxiv.org/abs/2505.03733).
- 634 Tianqi Luo, Chuhan Huang, Leixian Shen, Boyan Li, Shuyu Shen, Wei Zeng, Nan Tang, and Yuyu
635 Luo. nvbench 2.0: Resolving ambiguity in text-to-visualization through stepwise reasoning, 2025.
636 URL <https://arxiv.org/abs/2503.12880>.
- 637 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
638 mark for question answering about charts with visual and logical reasoning. *arXiv preprint
639 arXiv:2203.10244*, 2022.
- 640 MistralAI. Codestral. <https://mistral.ai/news/codestral>, 2024. 2024.05.29.
- 641 Yuansheng Ni, Ping Nie, Kai Zou, Xiang Yue, and Wenhui Chen. Viscoder: Fine-tuning llms for
642 executable python visualization code generation, 2025. URL [https://arxiv.org/abs/
643 2506.03930](https://arxiv.org/abs/2506.03930).

- 648 OpenAI. ChatML, 2022. URL <https://github.com/openai/openai-python/blob/e389823ba013a24b4c32ce38fa0bd87e6bccae94/chatml.md>.
- 649
- 650
- 651 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. URL <https://arxiv.org/abs/2303.08774>.
- 652
- 653 OpenAI. Gpt-5 system card. Technical report, 2025. URL <https://cdn.openai.com/gpt-5-system-card.pdf>.
- 654
- 655
- 656 Bo Pan, Yixiao Fu, Ke Wang, Junyu Lu, Lunke Pan, Ziyang Qian, Yuhan Chen, Guoliang Wang,
657 Yitao Zhou, Li Zheng, et al. Vis-shepherd: Constructing critic for llm-based data visualization
658 generation. *arXiv preprint arXiv:2506.13326*, 2025a.
- 659 Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe
660 Zhang. Training software engineering agents and verifiers with swe-gym, 2025b. URL <https://arxiv.org/abs/2412.21139>.
- 661
- 662 Raian Rahman, Rizvi Hasan, Abdullah Al Farhad, et al. Chartsumm: A comprehensive
663 benchmark for automatic chart summarization of long and short summaries. *arXiv preprint*
664 *arXiv:2304.13620*, 2023.
- 665
- 666 Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi
667 Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code Llama: Open foundation models for code.
668 *arXiv preprint arXiv:2308.12950*, 2023. URL <https://arxiv.org/abs/2308.12950>.
- 669 ByteDance Seed, Yuyu Zhang, Jing Su, Yifan Sun, Chenguang Xi, Xia Xiao, Shen Zheng, Anxiang
670 Zhang, Kaibo Liu, Daoguang Zan, et al. Seed-coder: Let the code model curate data for itself.
671 *arXiv preprint arXiv:2506.03524*, 2025.
- 672
- 673 Shengze Shi, Tao Ren, Guoliang Zhu, Guandong Feng, and Jun Hu. Closing the feedback loop in
674 text2vis: Refining visualization with vision-language models. In *Proceedings of the 33rd ACM*
675 *International Conference on Multimedia*, pp. 9053–9061, 2025.
- 676
- 677 Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan
678 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model par-
allelism, 2020. URL <https://arxiv.org/abs/1909.08053>.
- 679
- 680 Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. Design2code: How far
681 are we from automating front-end engineering? *arXiv preprint arXiv:2403.03163*, 2024.
- 682
- 683 5 Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang
684 Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang, Lucen Zhong, Mingdao Liu, Rui Lu,
685 Shulin Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei, Yean Cheng, Yifan An, Yilin Niu,
686 Yuanhao Wen, Yushi Bai, Zhengxiao Du, Zihan Wang, Zilin Zhu, Bohan Zhang, Bosi Wen,
687 Bowen Wu, Bowen Xu, Can Huang, Casey Zhao, Changpeng Cai, Chao Yu, Chen Li, Chendi
688 Ge, Chenghua Huang, Chenhui Zhang, Chenxi Xu, Chenzheng Zhu, Chuang Li, Congfeng Yin,
689 Daoyan Lin, Dayong Yang, Dazhi Jiang, Ding Ai, Erle Zhu, Fei Wang, Gengzheng Pan, Guo
690 Wang, Hailong Sun, Haitao Li, Haiyang Li, Haiyi Hu, Hanyu Zhang, Hao Peng, Hao Tai, Haoke
691 Zhang, Haoran Wang, Haoyu Yang, He Liu, He Zhao, Hongwei Liu, Hongxi Yan, Huan Liu, Hui-
692 long Chen, Ji Li, Jiajing Zhao, Jiamin Ren, Jian Jiao, Jiani Zhao, Jianyang Yan, Jiaqi Wang, Jiayi
693 Gui, Jiayue Zhao, Jie Liu, Jijie Li, Jing Li, Jing Lu, Jingsen Wang, Jingwei Yuan, Jingxuan Li,
694 Jingzhao Du, Jinhua Du, Jinxin Liu, Junkai Zhi, Junli Gao, Ke Wang, Lekang Yang, Liang Xu, Lin
695 Fan, Lindong Wu, Lintao Ding, Lu Wang, Man Zhang, Minghao Li, Minghuan Xu, Mingming
696 Zhao, Mingshu Zhai, Pengfan Du, Qian Dong, Shangde Lei, Shangqing Tu, Shangtong Yang,
697 Shaoyou Lu, Shijie Li, Shuang Li, Shuang-Li, Shuxun Yang, Siboyi, Tianshu Yu, Wei Tian,
698 Weihang Wang, Wenbo Yu, Weng Lam Tam, Wenjie Liang, Wentao Liu, Xiao Wang, Xiaohan Jia,
699 Xiaotao Gu, Xiaoying Ling, Xin Wang, Xing Fan, Xingru Pan, Xinyuan Zhang, Xinze Zhang,
700 Xiuqing Fu, Xunkai Zhang, Yabo Xu, Yandong Wu, Yida Lu, Yidong Wang, Yilin Zhou, Yiming
701 Pan, Ying Zhang, Yingli Wang, Yingru Li, Yinpei Su, Yipeng Geng, Yitong Zhu, Yongkun Yang,
Yuhang Li, Yuhao Wu, Yujiang Li, Yunan Liu, Yunqing Wang, Yuntao Li, Yuxuan Zhang, Zezhen
Liu, Zhen Yang, Zhengda Zhou, Zhongpei Qiao, Zhuoer Feng, Zhuorui Liu, Zichen Zhang, Zi-
han Wang, Zijun Yao, Zikang Wang, Ziqiang Liu, Ziwei Chai, Zixuan Li, Zuodong Zhao, Wen-
guang Chen, Jidong Zhai, Bin Xu, Minlie Huang, Hongning Wang, Juanzi Li, Yuxiao Dong,

- 702 and Jie Tang. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models, 2025a. URL
703 <https://arxiv.org/abs/2508.06471>.
704
- 705 Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of con-
706 text, 2024. URL <https://arxiv.org/abs/2403.05530>.
707
- 708 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijie Chen,
709 Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong,
710 Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao,
711 Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang
712 Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu,
713 Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin,
714 Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao
715 Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin
716 Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu,
717 Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe
718 Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo
719 Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi,
720 Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng
721 Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaying Wang,
722 Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang,
723 Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu,
724 Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing
725 Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie
726 Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao,
727 Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang
728 Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang,
729 Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng
730 Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou,
731 Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence,
732 2025b. URL <https://arxiv.org/abs/2507.20534>.
733
- 734 Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi
735 Wu, Haotian Liu, Sadhika Malladi, et al. Charxiv: Charting gaps in realistic chart understanding
736 in multimodal llms. *arXiv preprint arXiv:2406.18521*, 2024.
737
- 738 Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and
739 Ping Luo. Plot2code: A comprehensive benchmark for evaluating multi-modal large language
740 models in code generation from scientific plots. *arXiv preprint arXiv:2405.07990*, 2024.
741
- 742 xAI. Grok 4. Technical report, 2025. URL <https://x.ai/news/grok-4>.
743
- 744 Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Peng
745 Ye, Min Dou, Botian Shi, Junchi Yan, and Yu Qiao. Chartx chartvlm: A versatile benchmark and
746 foundation model for complicated chart reasoning, 2025. URL <https://arxiv.org/abs/2402.12185>.
747
- 748 Yupeng Xie, Zhiyang Zhang, Yifan Wu, Sirong Lu, Jiayi Zhang, Zhaoyang Yu, Jinlin Wang, Sirui
749 Hong, Bang Liu, Chenglin Wu, et al. Visjudge-bench: Aesthetics and quality assessment of
750 visualizations. *arXiv preprint arXiv:2510.22373*, 2025.
751
- 752 Ximing Xing, Juncheng Hu, Guotao Liang, Jing Zhang, Dong Xu, and Qian Yu. Empowering llms to
753 understand and generate complex vector graphics, 2025. URL <https://arxiv.org/abs/2412.11102>.
754
- 755 Kai Xu, Yiwei Mao, XinYi Guan, and ZiLong Feng. Web-bench: A llm code benchmark based on
756 web standards and frameworks, 2025. URL <https://arxiv.org/abs/2505.07473>.
757
- 758 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
759 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
760 *arXiv:2505.09388*, 2025a.

- 756 Cheng Yang, Chufan Shi, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu,
757 Siheng Li, Yuxiang Zhang, Gongye Liu, Xiaomei Nie, Deng Cai, and Yujiu Yang. Chartmimic:
758 Evaluating llm’s cross-modal reasoning capability via chart-to-code generation, 2025b. URL
759 <https://arxiv.org/abs/2406.09961>.
- 760 Jian Yang, Jiajun Zhang, Jiayi Yang, Ke Jin, Lei Zhang, Qiyao Peng, Ken Deng, Yibo Miao, Tianyu
761 Liu, Zeyu Cui, et al. Execrepobench: Multi-level executable code completion evaluation. *arXiv*
762 *preprint arXiv:2412.11990*, 2024a.
- 764 Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong, Xu Han, Yukun Yan, Zhenghao Liu, Zhixing
765 Tan, Pengyuan Liu, Dong Yu, et al. Matplotagent: Method and evaluation for llm-based agentic
766 scientific data visualization. *arXiv preprint arXiv:2402.11453*, 2024b.
- 767 Xingchen Zeng, Haichuan Lin, Yilin Ye, and Wei Zeng. Advancing multimodal large language
768 models in chart question answering with visualization-referenced instruction tuning. *IEEE Trans-*
769 *actions on Visualization and Computer Graphics*, 2024.
- 771 Lei Zhang, Jiayi Yang, Min Yang, Jian Yang, Mouxiang Chen, Jiajun Zhang, Zeyu Cui, Binyuan
772 Hui, and Junyang Lin. Swe-flow: Synthesizing software engineering data in a test-driven manner.
773 *arXiv preprint arXiv:2506.09003*, 2025.
- 774 Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Zhiyuan Liu, and Maosong Sun. Chart-
775 coder: Advancing multimodal large language model for chart-to-code generation, 2025. URL
776 <https://arxiv.org/abs/2501.06598>.
- 777 Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widayarsi, Imam
778 Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen Gong, Thong
779 Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan
780 Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang,
781 Binyuan Hui, Niklas Muennighoff, David Lo, Daniel Fried, Xiaoning Du, Harm de Vries, and
782 Leandro Von Werra. Bigcodebench: Benchmarking code generation with diverse function calls
783 and complex instructions, 2025. URL <https://arxiv.org/abs/2406.15877>.
- 784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810	APPENDIX	
811		
812		
813	A Benchmark Coverage Details	17
814	A.1 Chart Types	17
815	A.2 Thematic Coverage	19
816	A.3 Task Complexity	25
817	A.4 Sandboxed Environment	25
818		
819		
820		
821	B Benchmark Data Curation Details	26
822	B.1 Design Principles	26
823	B.2 Data Filtering	27
824	B.3 Task and Instruction Writing	27
825	B.4 Multi-turn Conversation	28
826	B.5 Quality Control	28
827		
828		
829		
830	C Evaluation Metrics	34
831	C.1 LLM Judge Cases	39
832		
833		
834	D Additional Results	44
835		
836	E SynthVis-30K Details	44
837		
838	F Evaluation Details	45
839		
840		
841	G Discussion	46
842	G.1 Model Performance Comparison	46
843	G.2 Scaling Comparison	51
844		
845		
846	H Correlation with Human Evaluation Details	56
847		
848	I Error Analysis	56
849		
850		
851		
852		
853		
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		

A BENCHMARK COVERAGE DETAILS

This section provides a detailed breakdown of the PlotCraft benchmark’s coverage. Table 5 presents the frequency distribution of 48 distinct chart types as they are mentioned within the instructions for each visualization task in the PlotCraft dataset. This distribution highlights the diversity of chart creation tasks included in the benchmark.

Table 5: This table illustrates the frequency distribution of 48 chart types mentioned within the instructions of visualization tasks from the PlotCraft dataset. To prevent data skew from repetition, each distinct task was analyzed and counted a single time, given that the same set of tasks is utilized for both single-turn and multi-turn interactions.

Type	Count	Type	Count	Type	Count	Type	Count
Line Chart	279	Scatter Plot	225	Time Series	179	Area Chart	164
Correlation Heatmap	145	Violin Plot	132	Box Plot	130	Stacked Area Chart	127
Error bar	96	Bubble Plot	91	Time Series Plot	88	Density Plot	85
Stacked Bar Chart	79	Network Graph	75	Parallel Coordinates	71	Histograms	70
Pie Chart	68	Radar Chart	67	grouped charts	67	Counts Plot	57
Dendrogram	56	Slope Chart	56	Seasonal Plot	54	Treemap	54
Cluster Plot	44	Population Pyramid	24	Pair Plot	17	Lollipop Chart	15
Dumbbell Plot	14	Autocorrelation	10	Cross-Correlation Plot	9	Waffle Chart	9
2D Histogram	7	Diverging Lollipop Chart	7	Pairwise Plot	7	Dot Plot	6
Jitter Plot	6	Partial Autocorrelation Plot	6	Categorical Plots	4	Joy Plot	4
Diverging Bars Chart	3	Ridgeline Plot	2	Stripplot	2	Calendar Heat Map	1
Correlogram	1	Distributed Dot Plot	1	Errorpoint Chart	1	Ordered Bar Chart	1

A.1 CHART TYPES

Distribution Tasks in this category require the visualization of data spread and grouping. Figure 6 presents a representative example of a class distribution task. The associated instruction demands the generation of a complex, multi-panel figure: “Create a comprehensive 3x3 subplot grid analyzing the distribution patterns of ECG signal characteristics across normal and abnormal cases.”

Correlation This category focuses on tasks that require visualizing the relationship and interdependence between two or more variables. Figure 7 illustrates a representative example, for which the model must generate a complex multi-panel visualization. The instruction for this task is: “Create a comprehensive 3x3 subplot grid analyzing the relationship between economic indicators and Olympic performance across different regions.”

Groups This category involves tasks that require visualizing clusters, hierarchies, and other logical groupings within a dataset. Figure 8 shows a complex example where the model is instructed to represent intricate relationships with the prompt: “Create a comprehensive 3x3 subplot grid analyzing the clustering patterns and hierarchical relationships within the Brawl Stars competitive ecosystem.”

Change Tasks in the Change category focus on visualizing trends, time-series data, and the evolution of metrics over a period. A representative example is presented in Figure 9, where the instruction requires a multi-faceted temporal analysis: “Create a comprehensive 3x2 subplot grid analyzing the temporal evolution of NIRF rankings across different educational categories from 2016-2021.”

Composition This category requires visualizing the composition of a whole, illustrating the proportions and breakdown of its constituent parts. Figure 10 provides an example of a multi-dimensional composition task, with the instruction: “Create a composite visualization showing Netflix content composition across different dimensions, designed as a 2x2 subplot grid.”

Ranking Ranking tasks involve comparing and ordering items based on one or more quantitative metrics. Figure 11 showcases an example where the model is prompted to perform a comparative ranking with the instruction: “Create a comprehensive ranking visualization that compares Udemy course performance across different metrics. Design a 2x2 subplot grid.”

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

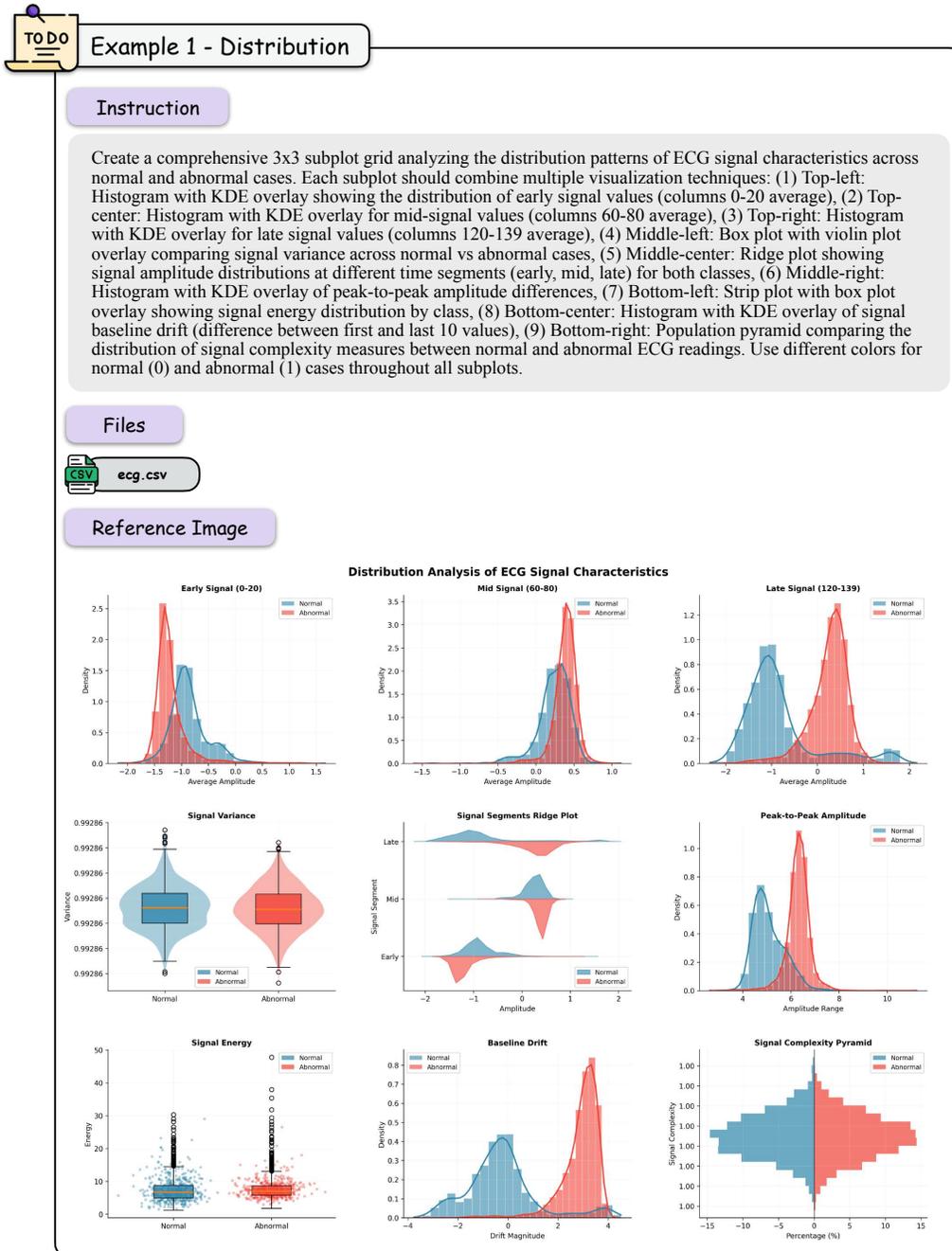
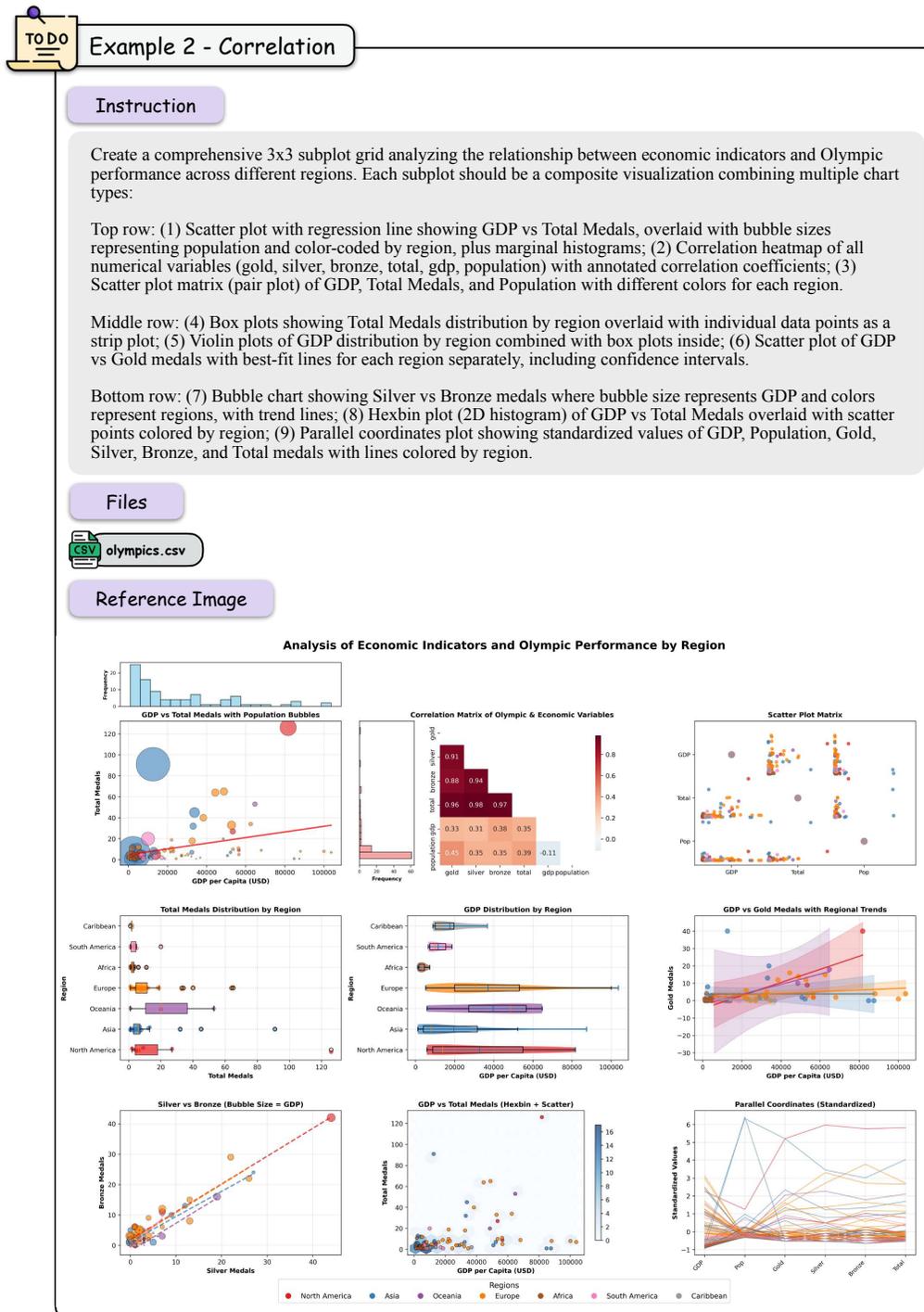


Figure 6: An example of a complex distribution visualization task from the PlotCraft benchmark. The model is required to generate a 3x3 grid of plots to compare ECG signal distributions, showcasing the benchmark’s focus on sophisticated, multi-panel figures.

Deviation The Deviation category includes tasks that visualize the difference or variance of data points against a fixed baseline or between different groups. An example is shown in Figure 12, where the task is to analyze prediction errors: “Create a 3x3 grid of composite visualizations analyzing the deviation patterns in Titanic passenger survival predictions.”

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025



1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

TODO

Example 3 - Groups

Instruction

Create a comprehensive 3x3 subplot grid analyzing the clustering patterns and hierarchical relationships within the Brawl Stars competitive ecosystem. Each subplot must be a composite visualization combining multiple chart types:

Row 1: Player Performance Clustering Analysis

- Subplot 1: Scatter plot with KDE contours showing the relationship between total trophies and highest trophies, with points colored by experience level and sized by 3vs3 victories
- Subplot 2: Violin plot overlaid with box plots comparing trophy distributions across different player name color categories, with individual data points as a swarm plot
- Subplot 3: Parallel coordinates plot showing the relationships between trophies, experience level, solo victories, duo victories, and 3vs3 victories for top players

Row 2: Club Ecosystem and Member Dynamics

- Subplot 4: Network-style scatter plot showing club trophies vs member count, with bubble sizes representing required trophies and colors indicating club type, overlaid with trend lines for each club type
- Subplot 5: Stacked bar chart showing club member role distributions (President, Vice President, Senior, Member) combined with a line plot overlay showing average member trophies per club
- Subplot 6: Hierarchical cluster heatmap showing the correlation matrix between club characteristics (trophies, member count, required trophies) and member performance metrics

Row 3: Brawler Performance Segmentation

- Subplot 7: Grouped violin plots comparing trophy distributions for different brawler power levels (9, 10, 11), overlaid with median lines and quartile markers
- Subplot 8: Multi-dimensional scatter plot matrix showing relationships between brawler rank, trophies, and highest trophies for selected high-performing brawlers, with regression lines
- Subplot 9: Dendrogram-style cluster analysis combined with a heatmap showing player performance patterns across different brawler categories (tank, support, damage dealer, etc.)

Each subplot should reveal distinct clustering patterns, hierarchical relationships, or natural groupings within the competitive Brawl Stars community structure.

Files

global_club_info.csv

global_player_rankings.csv

global_player_info.csv

global_club_rankings.csv

Reference Image

Brawl Stars: Club Performance and Player Clustering Analysis

Figure 8: An example of a complex Groups visualization task. The model must generate a 3x3 grid to display clustering and hierarchical structures, testing its ability to render complex relational plots like dendrograms or network graphs.

20

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

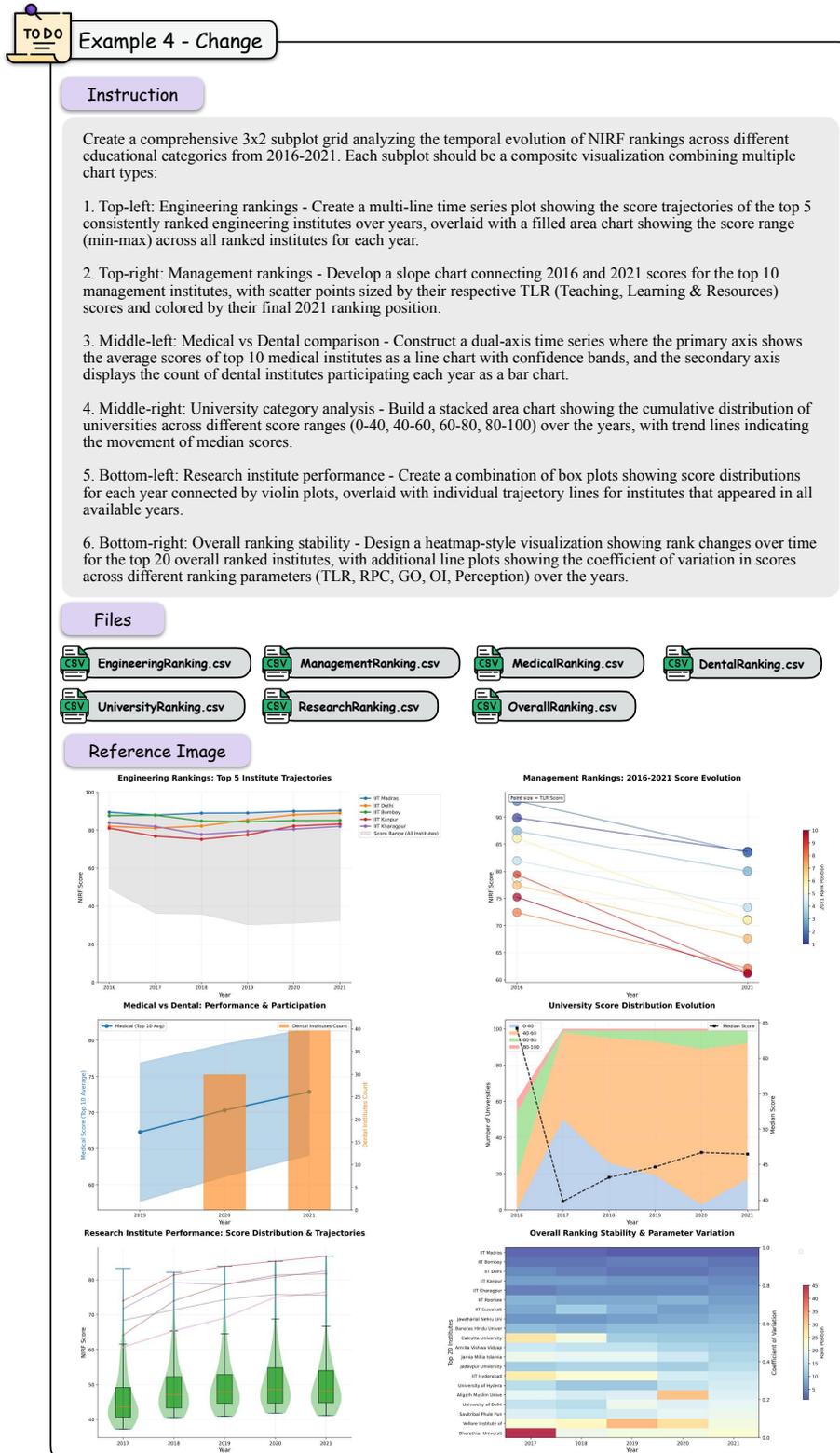


Figure 9: An example of a Change visualization task from the benchmark. This prompt requires a 3x2 subplot grid to track the temporal evolution of rankings, evaluating the model’s ability to create detailed time-series visualizations.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

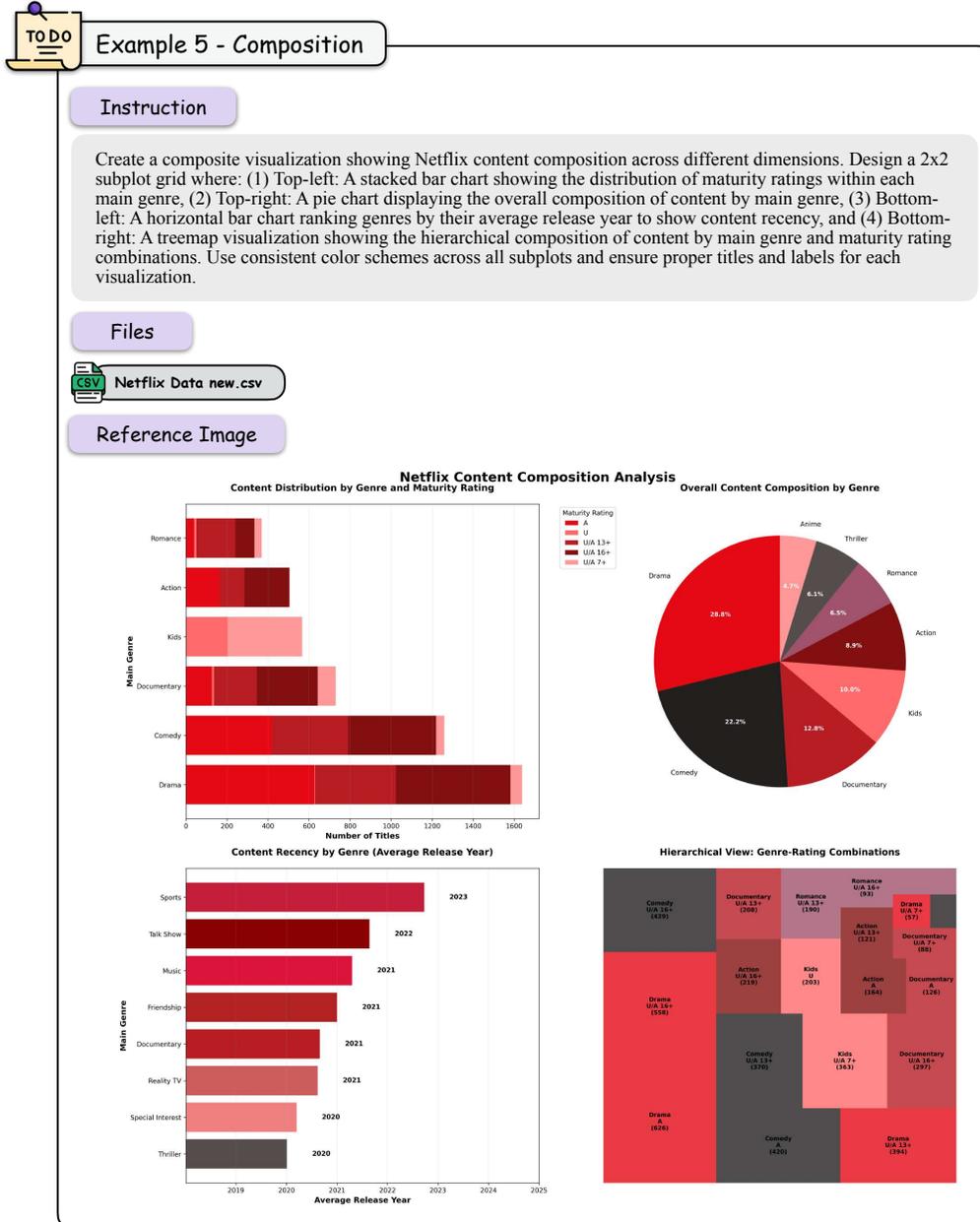
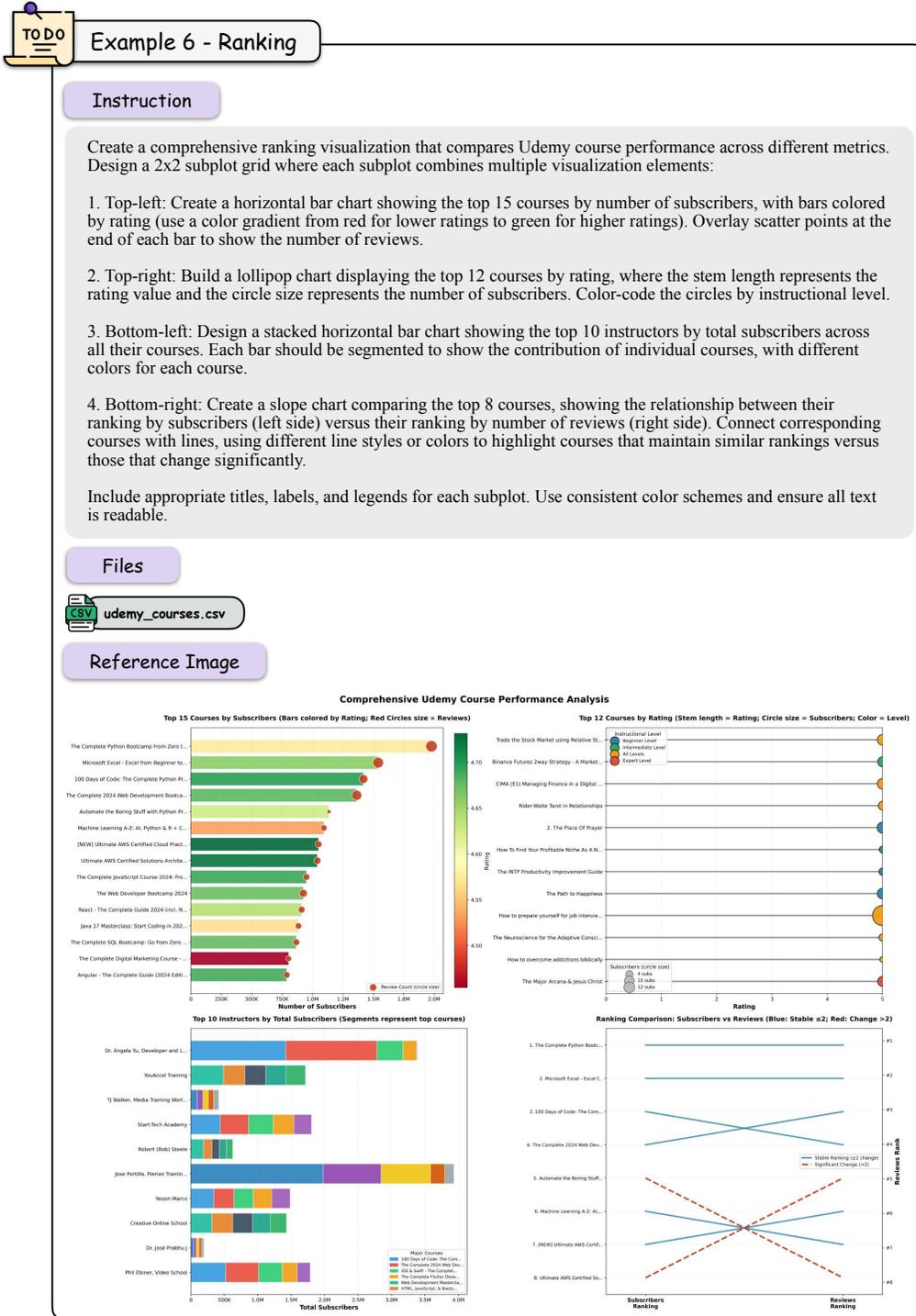


Figure 10: A Composition task example requiring the model to break down a dataset into its constituent parts. The instruction to use a 2x2 grid tests the ability to create comparative part-to-whole visualizations across different categories.

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



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

To Do

Example 7 - Deviation

Instruction

Create a 3x3 grid of composite visualizations analyzing the deviation patterns in Titanic passenger survival predictions. Each subplot should combine multiple chart types to show different aspects of deviation from expected survival patterns. Top row: (1) Diverging bar chart with error bars showing PassengerId ranges vs survival rate deviations from 50% baseline, overlaid with a line plot of cumulative deviation trends, (2) Dumbbell plot comparing actual vs expected survival rates across passenger ID quartiles with scatter points showing individual passenger deviations, (3) Area chart showing deviation magnitude distribution with overlaid violin plot of deviation densities. Middle row: (4) Radar chart displaying deviation metrics across multiple dimensions (early vs late passengers, ID clustering patterns) with overlaid polar scatter plot, (5) Slope chart connecting baseline expectations to actual outcomes across passenger segments with error band visualization, (6) Stacked area chart showing cumulative positive/negative deviations with overlaid trend line and confidence intervals. Bottom row: (7) Diverging lollipop chart of passenger ID bins vs deviation scores with histogram overlay, (8) Combined errorbar and stripplot showing deviation spread patterns across passenger groups, (9) Multi-layered deviation heatmap with contour lines and marginal distribution plots showing the relationship between PassengerId position and survival prediction accuracy.

Files

📄

All 1.csv

Reference Image

Figure 12: A Deviation task from the benchmark, requiring the analysis of prediction errors. The instruction to create a 3x3 grid tests the model’s ability to visualize variance and differences, such as in divergence bars or error plots.

ence & Research), **Society** (Government & Society), **Media** (Culture & Media), **Tech.** (Technology & Computing), and **Env.** (Environment & Geospatial). These domains are further subdivided into 31 fine-grained sub-topics. A comprehensive list of all included sub-topics is presented in Table 6.

Table 6: This table illustrates the coverage of 31 fine-grained thematic topics within the PlotCraft dataset.

Type	Count	Type	Count	Type	Count	Type	Count
Manufacturing, Logistics	18	Supply Chain, Retail	17	Sales, Customer Behavior	22	Energy	11
Medical Imaging	25	Clinical Trials	16	Public Health	20	Genomics for Medicine	11
Digital Health Records	19	Physics & Astronomy	14	Biology (non-medical)	15	Chemistry	13
Material Science	12	Social Sciences Research	15	Civics & Elections	10	Demographics	12
Urban Planning	14	Transportation	16	Education	18	Arts & Literature	11
Movies & TV	13	Music, Sports	12	News & Journalism	19	Computer Science & AI	44
Software & Code	28	Internet & Social Media	21	Cybersecurity	14	Climate & Weather	12
Satellite Imagery	7	Ecology & Agriculture	12	Geographic Information	-	-	-

A.3 TASK COMPLEXITY

Tasks within the PlotCraft benchmark are stratified into three distinct levels of complexity:

- **Easy:** Tasks require the generation of a single, standard chart (e.g., one bar chart, one line chart, or one scatter plot) to visualize the data.
- **Medium:** Tasks involve creating a composite visualization. This can be either: (a) a single figure that integrates multiple plot types or variables (e.g., a bar chart with a line plot overlay), or (b) a multi-panel grid (e.g., 2x1, 2x2) composed of simple, individual plots.
- **Hard:** Tasks demand the creation of a complex, multi-panel grid (e.g., 2x2, 3x3) wherein each individual subplot is itself a composite chart. For example, a 2x2 grid where each of the four plots contains both a histogram and a Kernel Density Estimation (KDE) curve.

To illustrate these levels, Figure 13 displays the instructions and raw data for three tasks of varying difficulty, all derived from the `firethorn-10-world-co2-emission-analysis` dataset. Figure 14 then shows the corresponding ground-truth reference images for each of these three instructions.

A.4 SANDBOXED ENVIRONMENT

To ensure the safe, consistent, and reproducible execution of all model-generated code, we constructed a dedicated sandboxed environment containerized using Docker. This approach provides an isolated and standardized platform, guaranteeing that all models are evaluated under identical conditions.

The environment is based on Python 3.13 and comes pre-installed with a comprehensive suite of libraries essential for data analysis and visualization. The foundational libraries include Pandas for data manipulation and NumPy for numerical operations. For visualization, the environment is equipped with Matplotlib as the primary plotting library, complemented by a wide array of higher-level and specialized libraries to support diverse charting requirements. These include:

- **Seaborn:** For high-level statistical graphics.
- **Plotly:** For interactive charts.
- **Squarify:** For creating treemaps.
- **scikit-learn:** For generating machine learning-related plots like confusion matrices.
- **statsmodels:** For advanced statistical visualizations.

All evaluations are conducted on a server equipped with 128 CPU cores and 1024 GB of RAM. Each code execution is performed within its container without network access and is subject to a strict execution timeout of 120 seconds to prevent runaway processes. This fully-specified, sandboxed setup eliminates variability from system configurations and provides a fair and secure assessment of each model’s code generation capabilities.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386

The screenshot shows a 'TODO' list item titled 'Difficulty Comparison - Part 1'. It contains three task levels: 'Easy Instruction', 'Medium', and 'Hard Instruction'. At the bottom, there are two file icons labeled 'World CO2 Emission Data.csv' and 'World CO2 Emission MetaData.csv'.

Dataset Name
aniruddha1995/firethorn-10-world-co2-emission-analysis

Easy Instruction
Create a line chart showing the evolution of CO2 emissions per capita over time from 1990 to 2020 for all four countries in the dataset. Extract and clean the CO2 emissions per capita data from the appropriate series, handle missing values by interpolation where reasonable, and display each country as a separate colored line with proper labels, legend, and title.

Medium
Create a composite visualization showing the temporal evolution of CO2 emissions across different countries. Design a subplot layout with two complementary charts: (1) a line chart displaying CO2 emissions trends over time (1990-2020) for the top 3 countries by total emissions, and (2) a stacked area chart showing the cumulative contribution of these same countries to global CO2 emissions over the same time period. Both charts should highlight how emission patterns have changed over the three-decade span and reveal which countries have been the dominant contributors to global CO2 emissions.

Hard Instruction
Create a comprehensive 3x3 subplot grid analyzing the temporal evolution of CO2 emissions and environmental indicators across different countries from 1990-2020. Each subplot must be a composite visualization combining multiple chart types:
Top row: (1) Line chart with confidence bands showing CO2 emissions per capita trends overlaid with scatter points for key milestone years, (2) Stacked area chart showing composition of different emission sources with trend lines for total emissions, (3) Dual-axis plot combining bar chart of annual emission changes with line plot of cumulative emissions.
Middle row: (4) Time series decomposition plot showing seasonal patterns in methane emissions with moving averages and trend components, (5) Slope chart connecting 1990 and 2020 values with intermediate milestone markers overlaid on a background heatmap of emission intensity, (6) Multi-line time series with error bands comparing agricultural vs energy sector emissions with filled areas showing the gap between them.
Bottom row: (7) Calendar heatmap of monthly emission data overlaid with box plots showing quarterly distributions, (8) Autocorrelation plot combined with partial autocorrelation analysis showing emission pattern dependencies over time, (9) Cross-correlation matrix heatmap between different emission types with time-lagged correlation coefficients displayed as a network graph overlay.
Each subplot must include country-specific trend analysis, statistical significance indicators, and highlight periods of major environmental policy changes or economic events that influenced emission patterns.

Files
World CO2 Emission Data.csv World CO2 Emission MetaData.csv

1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Figure 13: Example tasks of varying complexity levels from the PlotCraft benchmark. For a single dataset (firethorn-10-world-co2-emission-analysis), this figure shows the raw data and the natural language instructions for an Easy, Medium, and Hard task.

B BENCHMARK DATA CURATION DETAILS

B.1 DESIGN PRINCIPLES

The PlotCraft benchmark was designed to address critical gaps in existing evaluation methodologies, guided by four core principles for a more realistic and comprehensive assessment of LLM visualization capabilities.

- **Grounded in Real Data:** To ensure practical relevance and ecological validity, all tasks are constructed using authentic, real-world datasets. This principle allows the benchmark to circumvent the limitations and potential artifacts of synthetic data, which often lacks the noise, complexity, and inherent quirks found in genuine data sources. By grounding tasks in reality, PlotCraft provides a more robust evaluation of a model's ability to handle practical data visualization scenarios.

- 1404 • **Built from Scratch:** To mitigate data contamination and prevent benchmark leakage, all
1405 tasks and corresponding code in PlotCraft are built from scratch. We utilize open-source
1406 datasets from Kaggle to generate novel visualization challenges, ensuring no overlap with
1407 existing benchmarks or code repositories. This approach guarantees a more accurate as-
1408 sessment of a model’s true generalization capabilities.
- 1409 • **Zero-Reference Generation:** All tasks are initiated from text-only instructions, with no
1410 reference images or code provided. This setup compels the model to generate a visualiza-
1411 tion from an abstract concept, mirroring the creative and interpretive workflow of a data
1412 analyst, rather than the simpler task of replicating a known visual pattern.
- 1413 • **Compositional Complexity:** PlotCraft’s tasks feature a wide spectrum of complexity, re-
1414 quiring the generation of multi-panel plots with intricate layouts (e.g., 2×2, 3×3 grids),
1415 shared axes, complex legends, and the combination of multiple chart types within a single
1416 figure. This focus on composition directly probes the spatial and logical reasoning skills
1417 that are undertested by current benchmarks.

1419 B.2 DATA FILTERING

1420 Our data curation process began with sourcing datasets from Kaggle. This process was structured
1421 into two distinct phases to ensure the final benchmark was of high quality, diverse, and robust.

1422 The first phase consisted of a large-scale, automated pre-screening of an initial pool of 7,162 can-
1423 didate datasets, which collectively comprised over 95,000 files and 25.6 billion data rows. To filter
1424 for quality and relevance, we applied quantitative thresholds for community engagement metrics
1425 (vote counts ≥ 20 and download frequency ≥ 100) and official Kaggle usability ratings. Datasets
1426 that did not meet these minimum criteria for community validation and documentation quality were
1427 programmatically excluded.

1428 In the second phase, the resulting pool of high-quality datasets underwent a rigorous manual curation
1429 to select the final 140 core datasets for the benchmark. This selection was guided by three key
1430 principles. First, to ensure thematic diversity, we employed a stratified approach based on domain
1431 tags, guaranteeing broad coverage across topics like finance, healthcare, and technology. Second, we
1432 deliberately selected for a wide distribution of data volume and complexity, ensuring the inclusion
1433 of everything from small, clean tables to large, multi-file relational datasets. Finally, to mitigate
1434 data leakage, we conducted a thorough review to exclude datasets known to be prevalent in major
1435 pre-training corpora or other common visualization benchmarks.

1436 The final curated collection for PlotCraft consists of 1,874 raw data files, totaling approximately 462
1437 million data rows, providing a robust and novel foundation for evaluating visualization models.

1440 B.3 TASK AND INSTRUCTION WRITING

1441 The creation of our benchmark’s 491 unique visualization tasks was a meticulous, two-phase process
1442 designed to ensure depth, diversity, and a rigorous test of model capabilities.

1443 The first phase, Comprehensive Task Generation, was conducted by a team of three data visualiza-
1444 tion experts. Adopting a systematic approach, the experts were tasked with authoring prompts for
1445 each of the 140 curated datasets. For each dataset, they targeted seven distinct, high-level visual-
1446 ization intents (Correlation, Deviation, Ranking, etc.) and aimed to create a variant for each of the
1447 three complexity levels (Easy, Medium, and Hard). To guide this process, the experts utilized a rich
1448 taxonomy of over 50 distinct chart types. This initial generation phase resulted in a large candidate
1449 pool of nearly 2,940 tasks (140 datasets \times 7 intents \times 3 complexity levels), providing comprehensive
1450 coverage of datasets, analytical goals, and difficulty.

1451 The second phase, Collaborative Curation and Refinement, involved multiple rounds of review
1452 where the expert team discussed and filtered the large initial pool down to the final 491 tasks. Tasks
1453 were selected based on several criteria, including the clarity of the prompt, the feasibility of the
1454 visualization, and its analytical value. A key goal of this curation was to select the highest-quality
1455 examples while maintaining the balanced distribution across the three complexity levels established
1456 during generation. The final set reflects this, with 159 Easy, 163 Medium, and 169 Hard tasks.

Crucially, all final instructions were refined to be abstract and goal-oriented. They articulate the analytical requirements of the visualization, such as the desired layout, markings, and data transformations—without providing any guidance on code implementation. This approach compels the model to reason about the task from first principles, mirroring a more realistic human workflow.

The seven high-level visualization intents, along with the extensive range of chart types considered for each, are detailed below:

- **Correlation:** Scatter Plot, Bubble Plot, Scatter with Best Fit Line, Jitter Plot, Counts Plot, Scatter with Marginal Plots, Correlogram, Pairwise Plot, Network Graph, and Cluster Plot.
- **Deviation:** Diverging Bars Chart, Diverging Lollipop Chart, Slope Chart, Dumbbell Plot, Area Chart, Radar Chart, and Errorbar / Errorpoint Chart.
- **Ranking:** Ordered Bar Chart, Lollipop Chart, Dot Plot, Slope Chart, Dumbbell Plot, Stacked Bar Chart, and Diverging Lollipop Chart.
- **Distribution:** Histogram, Density Plot, Box Plot, Violin Plot, Joy Plot / Ridgeline Plot, Population Pyramid, Jitter Plot / Stripplot, and Categorical Plots.
- **Composition:** Stacked Bar Chart, Stacked Area Chart, Pie Chart, Treemap, and Waffle Chart.
- **Change:** Line Chart / Time Series Plot, Area Chart, Time Series with Error Bands, Calendar Heatmap, Seasonal Plot, Slope Chart, Dumbbell Plot, and Time Series Decomposition Plot.
- **Groups:** Dendrogram, Cluster Plot, Network Graph, Radar Chart, Treemap, Parallel Coordinates Plot, and Grouped Charts.

B.4 MULTI-TURN CONVERSATION

The following figures 15, 16, 17, and 18 provide several examples of the setup for our multi-turn refinement tasks. Each example consists of two key components: (1) the rendered image produced by an initial, intentionally flawed code implementation, and (2) the corresponding human-authored natural language request for its modification.

B.5 QUALITY CONTROL

To guarantee the robustness of PlotCraft, we implemented a multi-stage quality control protocol involving cross-validation and consensus-based adjudication.

Multi-turn Refinement Verification Ensuring the correctness of the multi-turn tasks required verifying the logical consistency between the faulty image, the error description, and the refinement instruction. We employed a *triple-verification strategy* for this phase:

1. **Triplet Validation:** Each generated triplet (Original Instruction, Faulty Image, Refinement Request) was reviewed by three independent annotators.
2. **Factuality & Consistency Check:** Annotators verified two specific conditions: First, *Error Existence*—confirming that the errors described in the refinement request (e.g., “the text overlaps”) were visually present in the rendered image. Second, *Instruction Consistency*—ensuring the refinement request did not contradict the original visualization goal.
3. **Majority Vote & Consensus:** A task was only accepted if it received a unanimous “Pass” vote. In cases of partial agreement (2/3), the triplet underwent a consensus discussion phase to determine if the refinement prompt could be polished. Triplets failing to reach consensus were discarded.

This strict verification loop ensures that all multi-turn interactions in PlotCraft reflect grounded, logical, and realistic user behaviors.

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

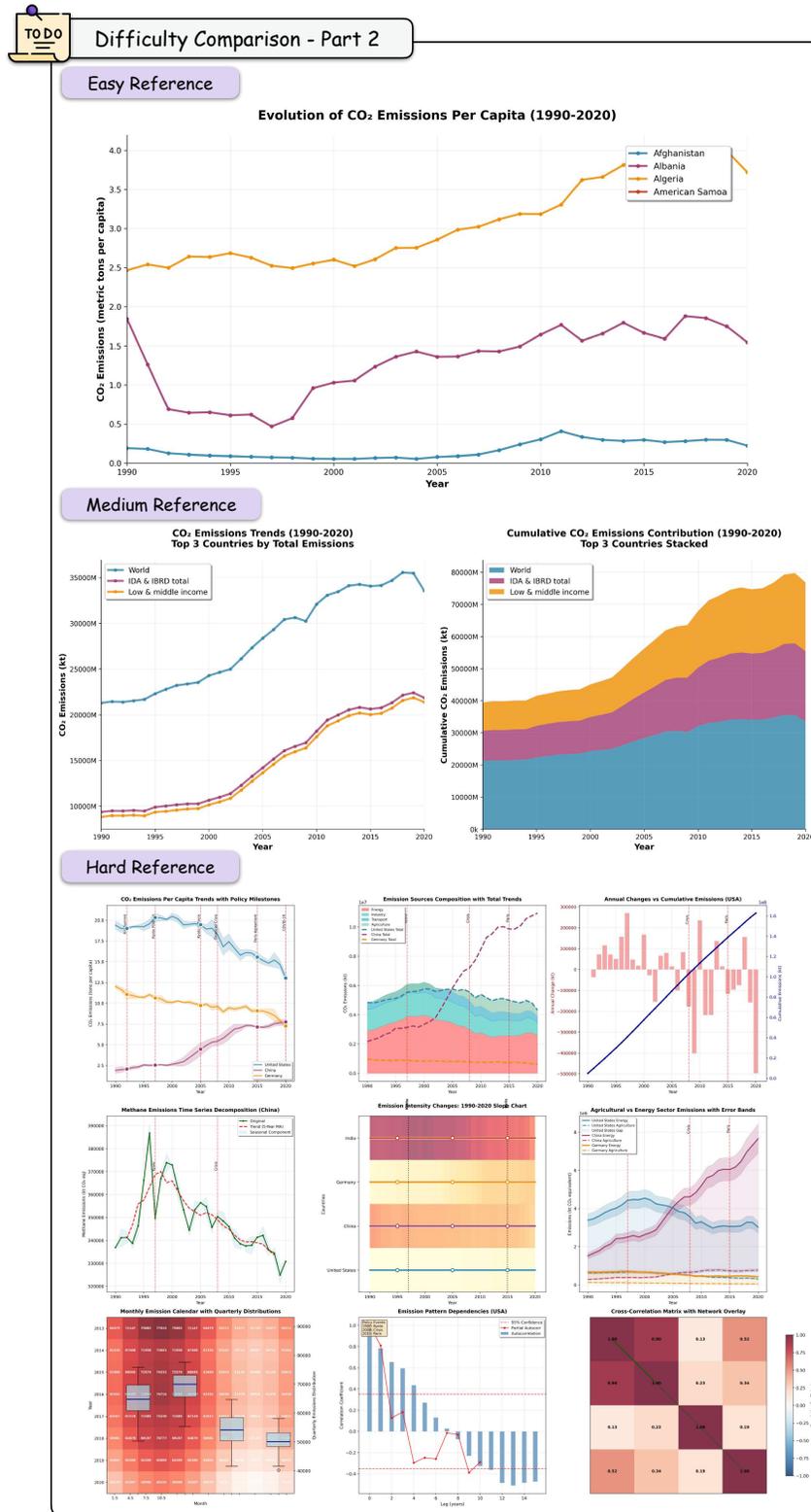


Figure 14: Reference images corresponding to the three task instructions shown in Figure 13. These images represent the ground-truth visualizations for the Easy, Medium, and Hard tasks, respectively, illustrating the expected output for each complexity level.

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

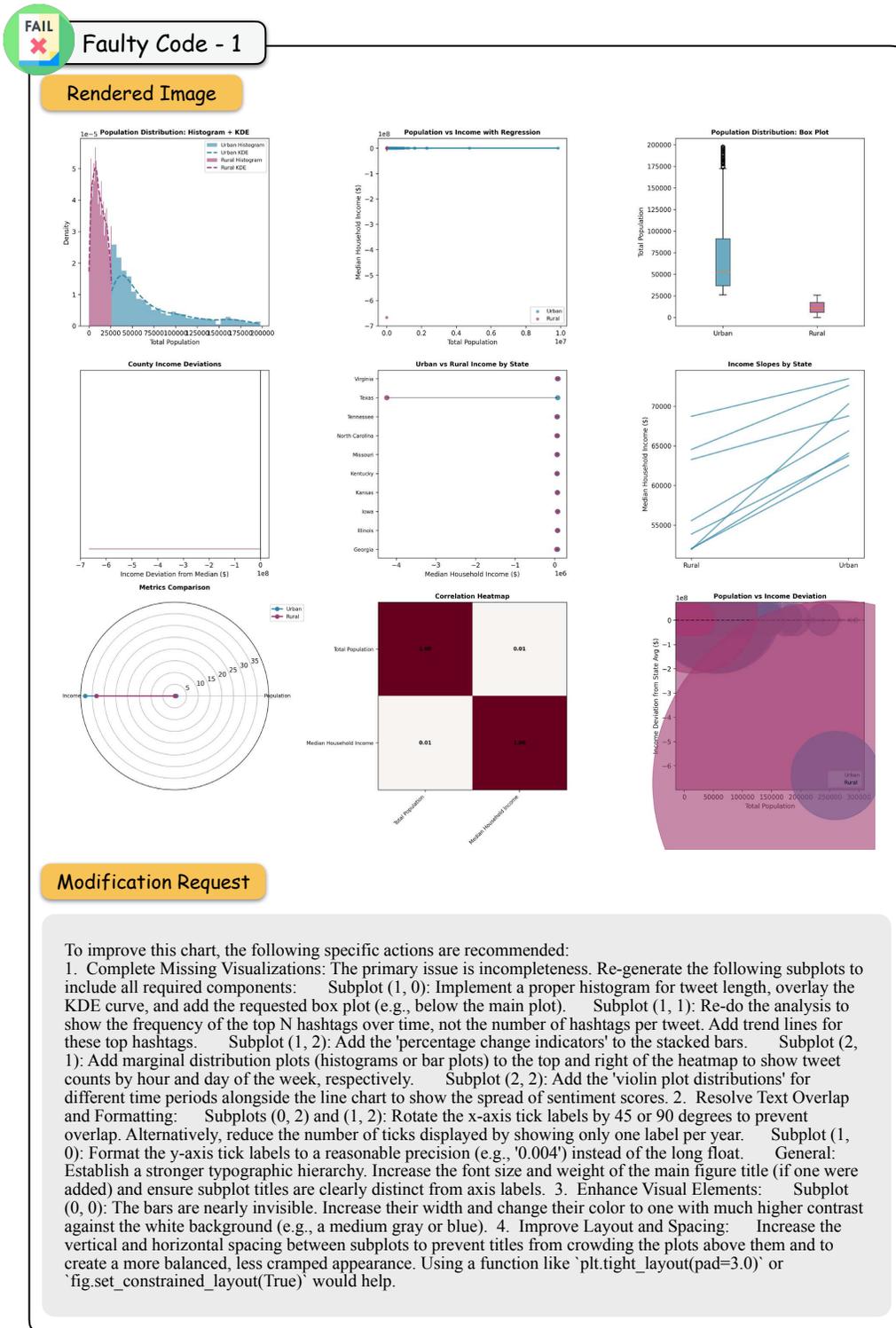


Figure 15: A multi-turn task example with multiple errors. Subplot (1, 0) fails to implement the required histogram, KDE curve, and boxplot. Furthermore, subplot (2, 2) exhibits severe rendering artifacts, including element overlap and content extending beyond the subplot boundaries.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

FAIL

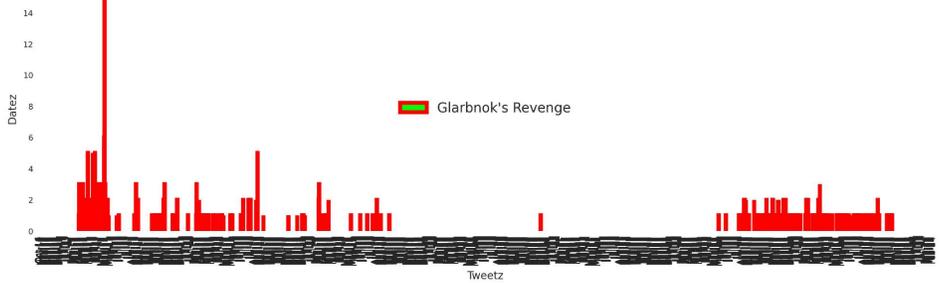
Faulty Code - 2

Rendered Image

Uniform Migration Patterns

Banana Frequency Over Time





Modification Request

The generated chart is fundamentally flawed and cannot be salvaged. It must be recreated from scratch to meet the user's requirements. 1. Correct Chart Type: Replace the bar chart and pie chart with a single line chart as requested. 2. Correct Data and Labels: The chart must plot the frequency of Khaadi tweets over time. The x-axis should be labeled 'Date' and the y-axis should be labeled 'Number of Tweets' or 'Tweet Frequency'. 3. Add a Relevant Title: The title should be descriptive and relevant, for example, 'Daily Tweet Frequency for Khaadi'. 4. Ensure Text Readability: The x-axis date labels must be formatted and spaced properly to be legible. If there are too many dates, consider rotating the labels or plotting data at a weekly interval. 5. Professional Aesthetics: Remove all nonsensical elements (the pie chart, the incorrect legend). Use a simple, professional color palette, such as a single shade of blue for the line. Ensure the title is bold and larger than the axis labels to establish a clear visual hierarchy. The layout should be balanced, utilizing the available space effectively without excessive white space.

Figure 16: An example of incorrect chart type and poor aesthetics. The generated code produces pie charts, which violates the instruction, and utilizes a color palette with excessively high contrast, diminishing the visual quality.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

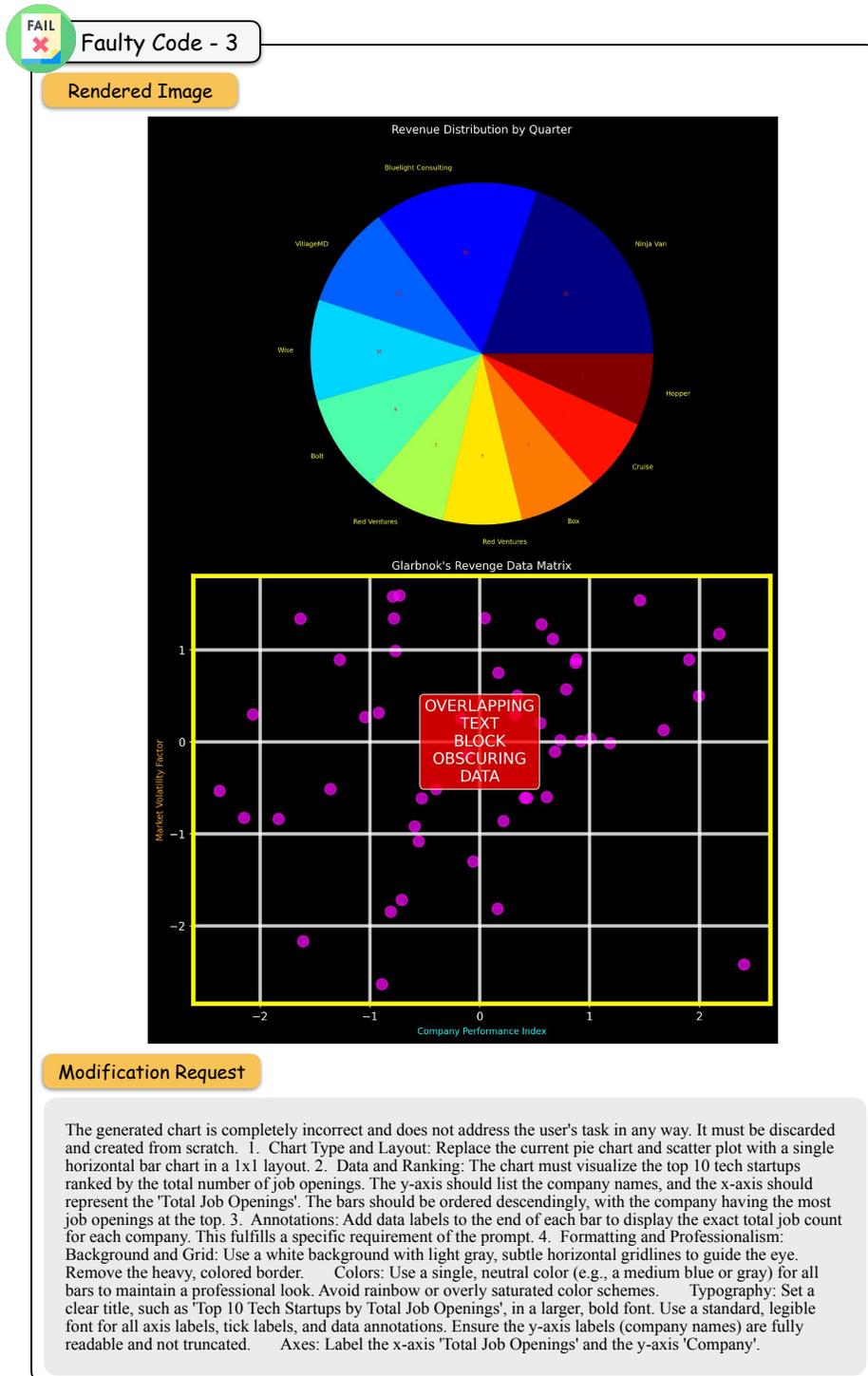


Figure 17: An example of incorrect layout and suboptimal color choice. The code fails to adhere to the specified 1x1 grid layout, and the visualization suffers from an oversaturated color scheme that compromises professional appearance.

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

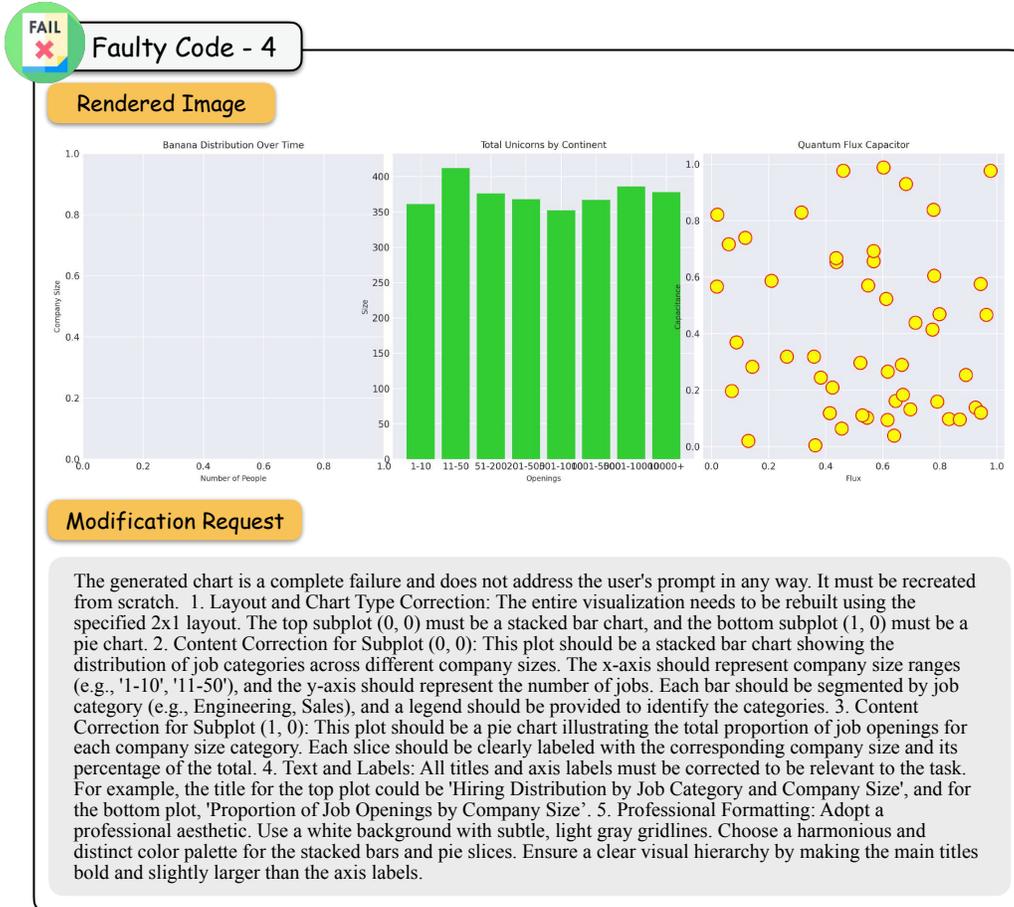


Figure 18: An example where chart types are correct but formatting is poor. The visualization suffers from overlapping x-axis tick labels, which impairs readability, and uses a default gray background that lacks professional polish.

C EVALUATION METRICS

This section details the comprehensive prompt in Figure C, C, C, C and scoring rubric used to evaluate the generated visualizations. Our evaluation is structured around two primary dimensions: **Task Compliance** and **Chart Quality**, each with a set of granular sub-metrics as defined below.

1. TASK COMPLIANCE (BINARY SCORING: 0/1)

For each criterion, a binary score is assigned: 1 for compliance (requirement met) or 0 for non-compliance (requirement not met).

1. **Layout Compliance:**

- *Question:* Does the chart follow the required layout specification (e.g., 1x1, 2x2)?
- *Score:* 1 if the layout matches the requirement exactly; 0 otherwise.

2. **Chart Type Compliance:**

- *Question:* Does the chart use the correct specified chart type(s)?
- *Score:* 1 if all chart types match the requirement; 0 otherwise.

3. **Visualization Requirement Fulfillment:**

- *Question:* Does the chart fulfill the core visualization goal (e.g., showing a relationship, displaying a trend)?
- *Score:* 1 if the core visualization requirement is met; 0 otherwise.

4. **Complete Task Fulfillment:**

- *Question:* Are all specified requirements from the task instruction completed?
- *Score:* 1 if all requirements are fulfilled; 0 if any requirement is missing.

2. CHART QUALITY (3-LEVEL SCORING: 0/1/2)

For each criterion, a score from 0 to 2 is assigned based on the following rubric.

1. **Clarity (No Overlap):**

Score 2 No overlap exists between any elements; all components are clearly separated.

Score 1 Minor overlap exists but does not significantly impact readability.

Score 0 Significant overlap exists, severely affecting the chart's readability.

2. **Layout Quality:**

Score 2 Excellent layout with well-proportioned elements and optimal spacing.

Score 1 Good layout with acceptable element sizes and spacing, but with minor imperfections.

Score 0 Poor layout with inappropriately sized, cramped, or poorly spaced elements.

3. **Color Quality:**

Score 2 Excellent, harmonious color scheme with appropriate contrast and visual appeal.

Score 1 Good color scheme with acceptable harmony and contrast, but with minor issues.

Score 0 Poor color scheme with clashing, harsh, or overly dull colors.

4. **Text Readability:**

Score 2 All text content is correct, appropriately sized, and clearly legible.

Score 1 Text contains minor issues (e.g., small font size, minor typos) that do not significantly impair understanding.

Score 0 Text has significant correctness or legibility issues (e.g., wrong labels, unreadable font).

5. **Professional Formatting:**

Score 2 (Publication-Ready) The chart is highly professional and polished, adhering to formal publication standards (e.g., white background, subtle gridlines, clear typographic hierarchy).

Score 1 (Needs Revision) The chart is functional but lacks professional refinement. It may use a plain default style, have a weak visual hierarchy, or heavy borders.

Score 0 (Unacceptable) The chart uses a non-standard, themed style inappropriate for formal contexts (e.g., non-white background, high-contrast gridlines).

Prompt C.1: Evaluation Prompt - Part 1

****Role and Goal:****

You are a meticulous data visualization quality inspector. Your task is to evaluate a generated chart based on a user's task prompt and the chart's visual properties. The evaluation is divided into two main categories: Task Compliance and Chart Quality.

****Task Prompt Format:****

The user's request follows this format:

```

# Task

Category: {task\_category}

Instruction: {task\_instruction}

```

****Evaluation Categories:****

1. Task Compliance (Binary Scoring: 0/1)

For each criterion, provide a binary score: '1' for compliance (requirement met) or '0' for non-compliance (requirement not met).

1. ****Layout Compliance:****

- ****Question****: Does the chart follow the required layout specification (e.g., 1x1, 2x2, 2x3)?
- ****Score****: '1' if the layout matches the requirement exactly. '0' if the layout differs from what was specified.

2. ****Chart Type Compliance:****

- ****Question****: Does the chart use the correct chart type as specified (e.g., bar chart, line chart, scatter plot)?
- ****Score****: '1' if the chart type matches the requirement. '0' if a different chart type was used.

3. ****Visualization Requirement Fulfillment:****

- ****Question****: Does the chart fulfill the specific visualization requirement (e.g., showing relationship between two variables, displaying trends over time)?
- ****Score****: '1' if the core visualization requirement is met. '0' if the requirement is not addressed.

4. ****Complete Task Fulfillment:****

- ****Question****: Are all specified requirements from the task instruction completed?
- ****Score****: '1' if all requirements are fulfilled. '0' if any requirement is missing or incomplete.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Prompt C.2: Evaluation Prompt - Part 2

2. Chart Quality (3-Level Scoring: 0/1/2)

For each criterion, provide a score from 0-2 with detailed explanations for each level.

- Clarity (No Overlap):**
 - Score 2:** No overlap exists between any elements; all subplots, titles, axis labels, tick marks, legends, and text boxes are clearly separated.
 - Score 1:** Minor overlap between text boxes or legends with plot content (data points, lines, bars) or border lines, but doesn't significantly impact readability.
 - Score 0:** Significant overlap between subplots, titles, axis labels, tick marks, or other text elements that severely affects readability.
- Layout Quality:**
 - Score 2:** Excellent layout with well-proportioned elements, optimal spacing, balanced white space distribution, and outstanding overall visual appeal.
 - Score 1:** Good layout with reasonable element sizes and spacing, acceptable visual balance with minor imperfections.
 - Score 0:** Poor layout with inappropriately sized elements, cramped or excessive spacing, unbalanced composition, or unappealing visual presentation.
- Color Quality:**
 - Score 2:** Excellent color scheme with harmonious palette, appropriate contrast, visually appealing combinations, and effective use of distinct colors for differentiation.
 - Score 1:** Good color scheme with acceptable harmony and contrast, minor issues that don't significantly impact aesthetics.
 - Score 0:** Poor color scheme with clashing colors, excessive harsh contrasts, overly dull/muted colors, or ineffective use of similar colors that lack distinction.
- Text Clarity:**
 - Score 2:** All text content is correct and appropriate, including accurate axis labels, proper titles, correct tick mark text, clear legends, and accurate text box content.
 - Score 1:** Most text content is correct with minor issues that don't significantly impair understanding or convey wrong information.
 - Score 0:** Text content has significant correctness issues including incorrect axis labels, inappropriate titles, wrong tick marks, unclear legends, or inaccurate text box content.

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

Prompt C.3: Evaluation Prompt - Part 3

5. **Formatting and Professional Standards:**
- **Score 2:** (Excellent / Publication-Ready): The chart's formatting is highly professional, clean, and adheres strictly to formal publication standards. It appears polished and intentionally designed, not like a default software output. Key characteristics include:
 - A white background is used.
 - Gridlines, if present, are subtle (thin, light gray) and do not distract from the data.
 - A clear typographic hierarchy is established, with the main title being visually distinct (e.g., bolded and/or larger) from axis labels and other text.
 - All non-data elements (axes, ticks) are appropriately weighted and do not appear heavy or clumsy.
 - **Score 1:** (Good / Needs Revision): The chart is functional but lacks professional refinement and contains minor stylistic issues. It is clear but would require formatting adjustments before formal publication. Key characteristics include:
 - The background is white, but the overall aesthetic is plain or unpolished.
 - The title lacks emphasis (e.g., is not bolded), making the visual hierarchy weak.
 - It may have slightly heavy chart borders (spines) or default styling that feels more like a "first draft" than a final product.
 - **Score 0:** (Poor / Unacceptable for Formal Use): The chart uses a non-standard, themed style that is inappropriate for academic or professional contexts. It is immediately identifiable as a default output from an analysis tool. Key characteristics include:
 - It features a non-white background (e.g., gray, beige).
 - It employs high-contrast, distracting elements like white gridlines on a colored background.
 - The overall visual style is cluttered or stylized in a way that detracts from a formal, serious tone.

Prompt C.4: Evaluation Prompt - Part 4

```

1998
1999
2000
2001 **Output Format:**
2002 Your response **MUST** be a single, valid JSON object, without
2003 any additional text before or after it. Use the exact
2004 structure below:
2005
2006 ```json
2007 {
2008   "task_compliance": {
2009     "layout_compliance": {
2010       "score": <0_or_1>,
2011       "reason": "<State whether the layout matches requirements
2012         and specify deviations if score is 0.>"
2013     },
2014     "chart_type_compliance": {
2015       "score": <0_or_1>,
2016       "reason": "<State whether the chart type matches
2017         requirements and specify what was expected vs. actual
2018         if score is 0.>"
2019     },
2020     "visualization_requirement_fulfillment": {
2021       "score": <0_or_1>,
2022       "reason": "<State whether the core visualization
2023         requirement is met and specify what's missing if score
2024         is 0.>"
2025     },
2026     "complete_task_fulfillment": {
2027       "score": <0_or_1>,
2028       "reason": "<State whether all requirements are completed
2029         and list missing items if score is 0.>"
2030     }
2031   },
2032   "chart_quality": {
2033     "clarity_no_overlap": {
2034       "score": <0_1_or_2>,
2035       "reason": "<Describe the overlap situation and justify the
2036         score level.>"
2037     },
2038     "layout_quality": {
2039       "score": <0_1_or_2>,
2040       "reason": "<Evaluate element sizing, spacing, and overall
2041         visual balance.>"
2042     },
2043     "color_quality": {
2044       "score": <0_1_or_2>,
2045       "reason": "<Assess color harmony, contrast, and aesthetic
2046         appeal.>"
2047     },
2048     "text_clarity": {
2049       "score": <0_1_or_2>,
2050       "reason": "<Evaluate text readability, correctness, and
2051         positioning.>"
2052     },
2053     "formatting_and_professional_standards": {
2054       "score": <0_1_or_2>,
2055       "reason": "Evaluate formatting and professional standards
2056         .>"
2057     }
2058   }
2059 }
2060 ```

```

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

C.1 LLM JUDGE CASES

Case - 1 This case, presented in Figure 19, illustrates a complete evaluation performed by our Gemini-2.5-Pro judge, including the original instruction, the model-generated image, and the resulting scores. The example is notable because the generated visualization is of high aesthetic quality but exhibits significant Task Compliance failures. This highlights the importance of decoupling the evaluation of quality from correctness.

Specifically, several subplots do not adhere to the prompt’s requirements:

- Subplot (1,0) is incorrectly implemented as a scatter plot instead of the required dumbbell plot.
- Subplot (1,1) is rendered as a line chart with error bands rather than the specified area chart.
- Subplots (0,0) and (0,1) fail to meet the "diverging" chart criteria, as they employ a monochromatic color scheme that does not differentiate between positive and negative values.

Case - 2 This case, presented in Figure 20, illustrates a generated visualization with a cascade of failures across multiple evaluation criteria, rendering it both incorrect and uninterpretable. The primary issues are categorized below.

- **Chart Type Non-Compliance:** The model fails to implement the specified chart types correctly. In subplot (1, 0), the bubbles are of uniform size and do not encode the transaction amount as required for a proper bubble chart. In subplot (2, 0), the treemap omits the required embedded bar charts, merely subdividing the areas instead.
- **Severe Element Overlap:** The visualization suffers from pervasive element occlusion that severely impacts readability. This includes illegible, overlapping x-axis tick labels in subplots (1, 0) and (2, 2); a legend in subplot (2, 1) that obstructs data points; and a dendrogram in subplot (2, 2) that clashes with its corresponding heatmap.
- **Inconsistent Color Scheme:** The use of color is critically flawed and misleading. In the first row, 'Fraud' is represented by red, but this is reversed in the parallel coordinates plot (1, 1), where 'Fraud' is incorrectly labeled as light blue. This semantic inconsistency makes the chart actively deceptive.
- **Poor Text Readability:** The chart has significant text-related issues. Multiple axis labels are illegible due to overlap, text within the treemap (2, 0) is too small to read, and some axes use raw, unformatted variable names (e.g., 'INIT_BALANCE'), which is unprofessional.

Case - 3 This case, presented in Figure 21, showcases a generated visualization that is highly successful in terms of both task compliance and overall chart quality. Its sole deficiency lies in the Professional Formatting sub-metric. The chart utilizes a non-standard, themed style that is inappropriate for formal publication; it features a gray background with high-contrast white gridlines, which is characteristic of a default software output (e.g., from Seaborn) rather than a polished, professional graphic. For publication-ready figures, a clean white background with subtle, non-distracting gridlines is the expected standard. This example underscores the importance of evaluating not just correctness, but also the fine-grained stylistic details that separate a functional plot from a professional one.

Case - 4 This case, presented in Figure 22, highlights how an otherwise high-quality visualization can be penalized for subtle but important flaws in its layout and clarity. While the chart successfully fulfills the core task requirements, its final score is reduced due to two specific issues. First, the layout is suboptimal, with excessive vertical white space between the main figure title and the subplot grid, creating a disjointed appearance. Second, the chart suffers from a minor clarity issue, as several x-axis tick labels exhibit slight overlap, which can impede readability. This example demonstrates the importance of meticulous polishing, a nuance that even capable models can overlook.

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

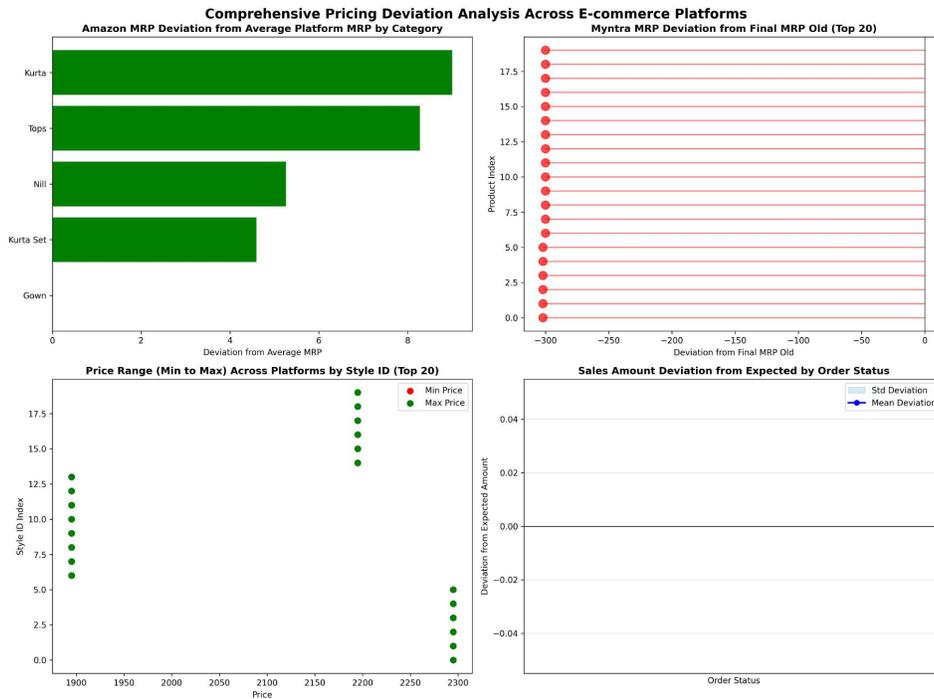


Judgement - 1

Instruction

Create a comprehensive pricing deviation analysis across multiple e-commerce platforms by combining data from three key files. Design a 2x2 subplot grid where each subplot shows the deviation of platform-specific MRPs from a baseline reference price. In the top-left, create a diverging bar chart showing how Amazon MRP deviates from the average MRP across all platforms for each product category. In the top-right, display a diverging lollipop chart comparing Myntra MRP deviations from the Final MRP Old baseline. In the bottom-left, construct a dumbbell plot showing the price range between the lowest platform MRP and highest platform MRP for each Style ID, highlighting the pricing spread. In the bottom-right, create an area chart with error bands showing the deviation of actual sale amounts from expected amounts (based on MRP) over different order statuses, using the Amazon sales data. Each subplot should use different colors to distinguish positive and negative deviations, and include reference lines at zero deviation.

Generation



Judgement

Layout Compliance	1	Clarity	2
Chart Type Compliance	0	Layout Quality	1
Requirement Fulfillment	0	Color Quality	2
Complete Task Fulfillment	0	Text Readability	2
		Professional Formatting	1

Figure 19: A complete evaluation case from our Gemini-2.5-Pro judge, displaying the original instruction, a generated image, and the final scores. This example highlights a common failure mode where a model produces a high-quality, aesthetically pleasing chart that nonetheless fails to comply with several key requirements of the prompt, particularly regarding the use of correct chart types.

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213



Judgement - 2

Instruction

Create a comprehensive 3x3 subplot grid analyzing money laundering patterns and account clustering in the AML dataset. Each subplot should combine multiple visualization techniques:

Row 1: Account Analysis

- Subplot 1: Combine a scatter plot of account initial balance vs fraud status with overlaid box plots showing balance distribution by fraud category
- Subplot 2: Create a network graph showing account relationships through transactions, with nodes colored by fraud status and sized by transaction volume
- Subplot 3: Overlay a histogram and KDE curve showing the distribution of account balances, with separate curves for fraudulent and non-fraudulent accounts

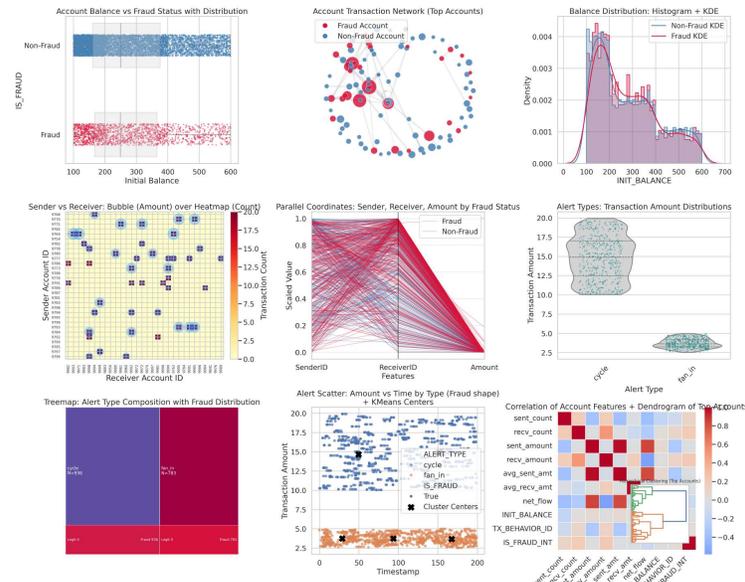
Row 2: Transaction Flow Analysis

- Subplot 4: Combine a bubble chart (sender vs receiver accounts) where bubble size represents transaction amount, overlaid with a heatmap showing transaction density
- Subplot 5: Create a parallel coordinates plot showing the relationship between sender account, receiver account, transaction amount, and fraud status
- Subplot 6: Overlay violin plots and strip plots showing transaction amount distributions across different alert types

Row 3: Alert Pattern Investigation

- Subplot 7: Combine a treemap showing alert type composition with embedded bar charts showing fraud distribution within each alert type
- Subplot 8: Create a cluster analysis plot using transaction amounts and timestamps, with points colored by alert type and shaped by fraud status
- Subplot 9: Overlay a correlation heatmap of numerical variables with a dendrogram showing hierarchical clustering of accounts based on transaction patterns

Generation



Judgement

Layout Compliance	1	Clarity	0
Chart Type Compliance	0	Layout Quality	0
Requirement Fulfillment	1	Color Quality	1
Complete Task Fulfillment	0	Text Readability	1
		Professional Formatting	0

Figure 20: A case study of a generated visualization exhibiting a cascade of failures. The output suffers from non-compliance with chart type requirements, severe element overlap, a critically inconsistent and misleading color scheme, and poor text formatting, rendering it both incorrect and uninterpretable.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

Judgement - 3

Instruction

Create a histogram showing the distribution of post scores in the Reddit r/datascience dataset. Use appropriate binning to reveal the underlying pattern of how posts are scored, and include proper axis labels and a title that describes what the visualization shows.

Generation

Distribution of Post Scores in r/datascience Subreddit
Showing Strong Right-Skewed Pattern with Most Posts Receiving Low Scores

Judgement

Layout Compliance	1	Clarity	2
Chart Type Compliance	1	Layout Quality	2
Requirement Fulfillment	1	Color Quality	2
Complete Task Fulfillment	1	Text Readability	2
		Professional Formatting	0

Figure 21: An example of a high-quality visualization that is penalized for a lack of professional formatting. While the chart correctly adheres to all task requirements, its use of a default gray background and high-contrast gridlines prevents it from meeting publication-ready standards.

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321



Judgement - 4

Instruction

Create a comprehensive 3x3 subplot grid analyzing university performance clusters and regional patterns across Asian countries. Each subplot should be a composite visualization: (1) Top-left: Scatter plot with Academic Reputation vs Employer Reputation, overlaid with country-based color coding and bubble sizes representing Overall Score, plus marginal histograms showing distribution of each metric; (2) Top-center: Stacked bar chart showing count of universities by country, overlaid with a line plot showing average Overall Score per country; (3) Top-right: Box plot showing Citations per Paper distribution by country, overlaid with violin plots to show density distributions; (4) Middle-left: Radar chart comparing average performance metrics (Academic Reputation, Employer Reputation, International Students, International Faculty, Faculty Student Ratio) for top 5 countries by university count; (5) Middle-center: Heatmap showing correlation matrix of all numerical performance metrics, overlaid with hierarchical clustering dendrogram; (6) Middle-right: Parallel coordinates plot showing performance profiles of top 20 universities across key metrics (Academic Reputation, Employer Reputation, Citations per Paper, International Students), with lines colored by country; (7) Bottom-left: Treemap showing university count by country and city hierarchy, with cell sizes representing total universities and colors representing average Overall Score; (8) Bottom-center: Network graph showing country relationships based on similarity in performance metrics, with node sizes representing university count and edge weights representing similarity scores; (9) Bottom-right: Cluster scatter plot using PCA on all performance metrics, with points colored by country and shaped by ranking tiers (1-20, 21-50, 51-100), overlaid with cluster boundaries and centroids.

Generation



Judgement

Layout Compliance	1	Clarity	1
Chart Type Compliance	1	Layout Quality	1
Requirement Fulfillment	1	Color Quality	2
Complete Task Fulfillment	1	Text Readability	1
		Professional Formatting	1

Figure 22: An example of a high-quality visualization penalized for subtle layout and clarity issues. The excessive white space between the title and the plots, along with slightly overlapping x-axis labels, detracts from its overall professional quality.

Model	Single-Turn Generation			Multi-Turn Refinement			AVG score
	Pass Rate (%)	Task-Comp.	Quality	Pass Rate (%)	Task-Comp.	Quality	
<i>Closed-source LLMs</i>							
Claude-4.1-Opus (Anthropic, 2023)	76.20	1.93	4.20	81.44	2.05	5.22	6.70
Claude-4-Sonnet (Anthropic, 2023)	68.84	1.73	3.99	78.41	1.88	4.80	6.20
Gemini-2.5-Pro (Team, 2024)	41.34	1.15	2.31	58.86	1.51	3.80	4.39
ChatGPT-4o-Latest (OpenAI, 2023)	63.54	1.60	3.33	69.25	1.51	4.23	5.33
GPT-5 (OpenAI, 2025)	69.86	1.76	2.87	74.13	1.80	4.33	5.38
Grok-4 (xAI, 2025)	64.52	1.63	3.66	70.24	1.61	4.42	5.65
<i>Open-source LLMs</i>							
Kimi-K2 (Team et al., 2025b)	60.13	1.52	3.36	61.03	1.49	4.05	5.21
DeepSeek-V3.1 (DeepSeek-AI & etc., 2024)	55.71	1.49	3.20	56.86	1.50	3.99	5.09
DeepSeek-Coder-V2 (Guo et al., 2024)	47.45	1.23	2.68	68.23	1.47	4.12	4.76
DeepSeek-Coder-V2-Lite (Guo et al., 2024)	32.79	0.77	1.90	47.45	0.97	2.73	3.19
GLM-4.5 (Team et al., 2025a)	43.38	1.25	2.48	64.97	1.54	3.91	4.59
GPT-oss-120B (Agarwal et al., 2025)	48.23	1.32	2.68	29.69	0.77	1.87	3.32
GPT-oss-20B (Agarwal et al., 2025)	44.68	1.21	2.43	34.02	0.89	2.14	3.34
Seed-Coder-8B (Seed et al., 2025)	32.38	0.87	1.74	57.03	1.26	3.22	3.55
VisCoder-7B (Ni et al., 2025)	25.46	0.67	1.50	51.73	0.99	2.96	3.06
Qwen2.5-Coder-1.5B (Hui et al., 2024)	22.81	0.55	1.38	36.86	0.66	1.71	2.15
Qwen2.5-Coder-3B (Hui et al., 2024)	17.92	0.50	1.17	36.46	0.77	2.00	2.22
Qwen2.5-Coder-7B (Hui et al., 2024)	29.94	0.79	1.79	46.64	0.98	2.65	3.10
Qwen2.5-Coder-14B (Hui et al., 2024)	38.29	1.09	2.29	51.53	1.17	3.04	3.80
Qwen2.5-Coder-32B (Hui et al., 2024)	37.68	1.05	2.21	48.47	1.20	3.11	3.79
Qwen3-235B-A22B-2507 (Yang et al., 2025a)	56.01	1.53	3.31	71.69	1.70	4.49	5.51
Qwen3-Coder-480B-A35B (Hui et al., 2024)	61.30	1.56	3.21	75.97	1.75	4.46	5.49
Qwen3-Coder-30B-A3B (Hui et al., 2024)	52.55	1.32	2.85	73.12	1.55	4.18	4.95
PlotCrafter-30B-A3B (Ours)	64.36	1.73	4.09	77.11	1.76	4.74	6.16

Table 7: The complete quantitative results on PlotCraft for 24 LLMs across two settings: Single-Turn Generation and Multi-Turn Refinement.

D ADDITIONAL RESULTS

Table 7 presents the complete quantitative results for all 24 evaluated LLMs on the PlotCraft benchmark. The data reinforces the findings from our main analysis and further highlights the strong performance of our model.

Among all open-source models, PlotCrafter demonstrates a clear advantage, achieving an average score of 6.16. This result significantly surpasses other leading open-weight models, including the much larger Qwen3-235B (5.51) and Qwen3-Coder-480B (5.49). More importantly, PlotCrafter’s performance closes the gap with top-tier proprietary systems, achieving a score nearly identical to that of Claude-4-Sonnet (6.20). These comprehensive results validate that PlotCrafter provides SOTA capabilities within the open-source community for complex data visualization tasks.

E SYNTHVIS-30K DETAILS

This section provides further implementation details for the multi-agent framework used to create the SynthVis-30K dataset.

Task Generation Agents and Criteria. The roles of both the Task Generator and the Task Judge were fulfilled by an ensemble of Large Language Models, primarily consisting of Claude-4-Sonnet and Qwen3-Coder-480B-A22B. The iterative refinement cycle for a given task was designed to be rigorous; a task was only finalized and accepted when the Task Judge agent could no longer identify any logical inconsistencies or feasibility issues with the proposed instructions.

Code Generation Agents and Termination Criteria. The Code Generator agent was also comprised of an ensemble of Claude-4-Sonnet and Qwen3-Coder-480B-A22B. The crucial role of the Visual Judge, responsible for assessing the quality of the rendered images, was performed by Gemini-2.5-Pro. The code generation process was constrained to a maximum of 10 refinement iterations. The cycle was considered successful and terminated early if the generated visualization achieved a perfect Task Compliance score (4 out of 4) and a Chart Quality score of at least 6 (out of 10), with the additional constraint that each of the five quality sub-metrics must score at least 1. If these criteria were not met within the 10-iteration limit, the entire task-code pair was discarded to maintain the high-quality standard of the final dataset.

2376 **SFT Trajectory Synthesis Details.** The Chain-of-Thought (CoT) rationales for the single-turn
2377 training instances were generated using Claude-4-Sonnet. The final SynthVis-30K dataset is com-
2378 posed of 30,000 instances with a 2:1 split between single-turn and multi-turn trajectories, resulting
2379 in 20,000 single-turn generation instances and 10,000 multi-turn refinement instances.

2380

2381

2382 F EVALUATION DETAILS

2383

2384

2385

2386

2387

2388

2389

2390

2391

2392

2393

2394

Prompt F.1: Generation Prompt

2395

2396

2397

2398

2399

2400

2401

2402

2403

2404

2405

2406

2407

2408

2409

2410

2411

2412

2413

2414

2415

2416

2417

2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2429

```
You are an expert Python data visualization developer
specializing in matplotlib and seaborn. Your task is to
generate high-quality, executable Python code for data
visualization based on the given task description and
dataset information.

Output format:
- Provide only the Python code wrapped in ```python and ```
  markers
- Ensure the code can run independently
```

2430 G DISCUSSION

2431
2432
2433
2434
2435
2436
2437
2438

2439 G.1 MODEL PERFORMANCE COMPARISON

2440
2441
2442
2443
2444

2445 **Comparison 1** For the task detailed in Instruction G.1, Figure 23 compares the outputs of several
2446 leading models. Our model, PlotCrafter, successfully generates a visualization that meets all speci-
2447 fied requirements. In contrast, the base model, Qwen3-Coder-30B-A3B, exhibits both chart type and
2448 factual errors and uses an unprofessional default gray background. GPT-5 clutters the visualization
2449 with excessive legends, gridlines, and text, leading to poor clarity and element overlap. Similarly,
2450 GLM-4.5 and Cluade-4.1-Opus produce charts with unpolished gray backgrounds and poor color
2451 choices, with the latter also failing on chart type compliance. Gemini-2.5-Pro fails to generate any
2452 output, resulting in a blank image.

2453
2454
2455
2456
2457
2458
2459
2460
2461
24622463 **Instruction G.1: Comparison - 1**2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

Create a comprehensive temporal analysis of Kaggle tweet engagement patterns from 2010–2021. Design a 2x2 subplot grid where each subplot combines multiple visualization elements: (1) Top-left: A line chart showing yearly tweet volume trends overlaid with a bar chart displaying average engagement metrics (likes + retweets) per year, (2) Top-right: An area chart depicting the cumulative distribution of tweet languages over time with stacked areas for the top 5 languages, (3) Bottom-left: A dual-axis plot combining a line chart of monthly tweet frequency patterns overlaid with a scatter plot showing seasonal engagement spikes, and (4) Bottom-right: A time series decomposition showing trend, seasonal, and residual components of daily tweet activity across the entire dataset period. Each subplot should include appropriate legends, annotations for significant events or patterns, and use consistent color schemes to highlight the evolution of Kaggle’s social media presence and community engagement over the decade.

2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537

Instruction G.2: Comparison - 2

Create a comprehensive 3x2 subplot grid analyzing profit margin deviations and pricing strategies across multiple e-commerce platforms. Each subplot should be a composite visualization combining multiple chart types:

Top row (3 subplots):

1. Left: Create a diverging bar chart showing the deviation of each platform's MRP from the average MRP, overlaid with error bars representing the standard deviation of pricing across different style categories
2. Center: Design a dumbbell plot comparing TP1 vs TP2 costs for different product categories, with a secondary y-axis line plot showing the profit margin percentage deviation from the overall average margin
3. Right: Build a slope chart showing MRP changes from "MRP Old" to "Final MRP Old" for top 10 style IDs, combined with scatter points indicating the magnitude of price adjustment

Bottom row (2 subplots):

4. Left: Construct a diverging lollipop chart displaying how each platform's MRP deviates from the baseline Amazon MRP, with horizontal reference lines showing $\pm 10\%$ and $\pm 20\%$ deviation thresholds
5. Right: Generate a radar chart comparing normalized pricing metrics (TP1, TP2, various platform MRPs) for the top 5 most frequent style categories, overlaid with area fill showing the deviation range from the category median

Use a consistent color scheme where positive deviations are shown in green tones and negative deviations in red tones. Include proper titles, legends, and annotations highlighting the most significant deviations. The visualization should reveal pricing inconsistencies, profit margin variations, and strategic pricing patterns across different e-commerce platforms.

Instruction G.3: Comparison - 3

Create a composite visualization showing the composition of cosmetic products by chemical content. Design a subplot with two complementary charts: (1) a stacked bar chart displaying the top 8 companies by total product count, with each bar segment colored by primary category to show the product type distribution within each company, and (2) a pie chart showing the overall market share of these top 8 companies based on their total number of reported products. Include proper legends, titles, and ensure the color schemes are consistent between both charts.

2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591

Instruction G.4: Comparison - 4

Create a comprehensive 3x3 subplot grid analyzing the temporal evolution of London Underground station usage from 2007–2017. Each subplot should be a composite visualization combining multiple chart types:

Top row (2007–2009): For each year, create a scatter plot showing the relationship between weekday entries and annual usage (in millions), with bubble sizes representing weekend activity levels (Saturday + Sunday entries), overlaid with a regression line and confidence intervals.

Middle row (2010–2012): For each year, create a dual-axis plot combining a histogram of annual usage distribution with a KDE curve overlay, and add vertical lines marking the 25th, 50th, and 75th percentiles of usage.

Bottom row (2013–2015): For each year, create a combination plot showing both a box plot of weekday vs weekend entry ratios by borough (grouped by top 10 boroughs by station count) and overlay violin plots to show the distribution density.

Each subplot should include year-specific titles, appropriate color schemes that evolve across the timeline, and statistical annotations (correlation coefficients for scatter plots, percentile values for histograms, and median values for box plots). The overall visualization should reveal how station usage patterns, distributions, and borough-level variations evolved during this decade.

Instruction G.5: Comparison - 5

Create a composite visualization showing the temporal evolution of CO2 emissions across different countries. Design a subplot layout with two complementary charts: (1) a line chart displaying CO2 emissions trends over time (1990–2020) for the top 3 countries by total emissions, and (2) a stacked area chart showing the cumulative contribution of these same countries to global CO2 emissions over the same time period. Both charts should highlight how emission patterns have changed over the three-decade span and reveal which countries have been the dominant contributors to global CO2 emissions.

Comparison 2 The task in Instruction G.1 requires a sophisticated 3x2 grid with composite charts (3 subplots for top row and 2 for bottom row). As shown in Figure 24, PlotCrafter correctly renders this complex layout. The other models struggle significantly: Qwen3-Coder-30B-A3B implements an incorrect layout and suffers from extensive element overlap. GPT-5 produces plots with distorted aspect ratios, awkward typography, and significant text overlap. GPT-oss-120B also fails to generate the correct 3x2 layout. While Cluade-4.1-Opus generates a mostly correct visualization, it misplaces the legend, affecting the overall composition. Gemini-2.5-Pro again produces a blank image.

2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645

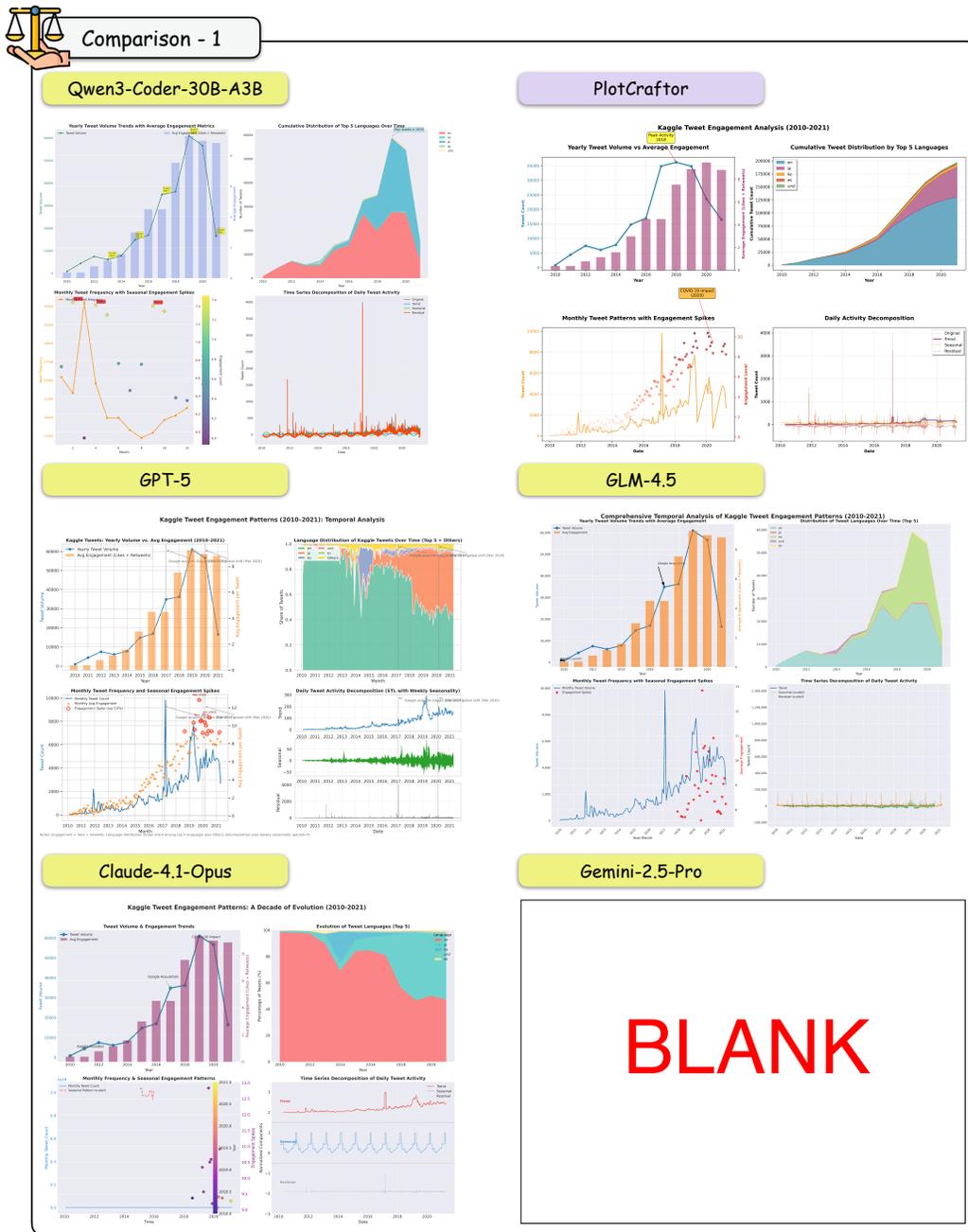


Figure 23: Qualitative comparison for the task in Instruction G.1. PlotCrafter produces a correct and complete visualization, while other models exhibit a range of failures, including incorrect chart types (Qwen3-Coder, Claude-4.1-Opus), excessive clutter (GPT-5), poor formatting (GLM-4.5), and a blank output (Gemini-2.5-Pro).

Comparison 3 The task in Instruction G.1, which requires a simpler composite chart, demonstrates that many high-performing models can handle less complex requests. As seen in Figure 25, PlotCrafter, Kimi-K2, Claude-4-Sonnet, and Qwen3-Coder-480B-A22B all produce satisfactory results. This comparison highlights the specific failure modes of other models on what should be a manageable task: Qwen3-Coder-30B-A3B uses an incorrect chart type, while GPT-5 fails to produce any visualization.

2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699

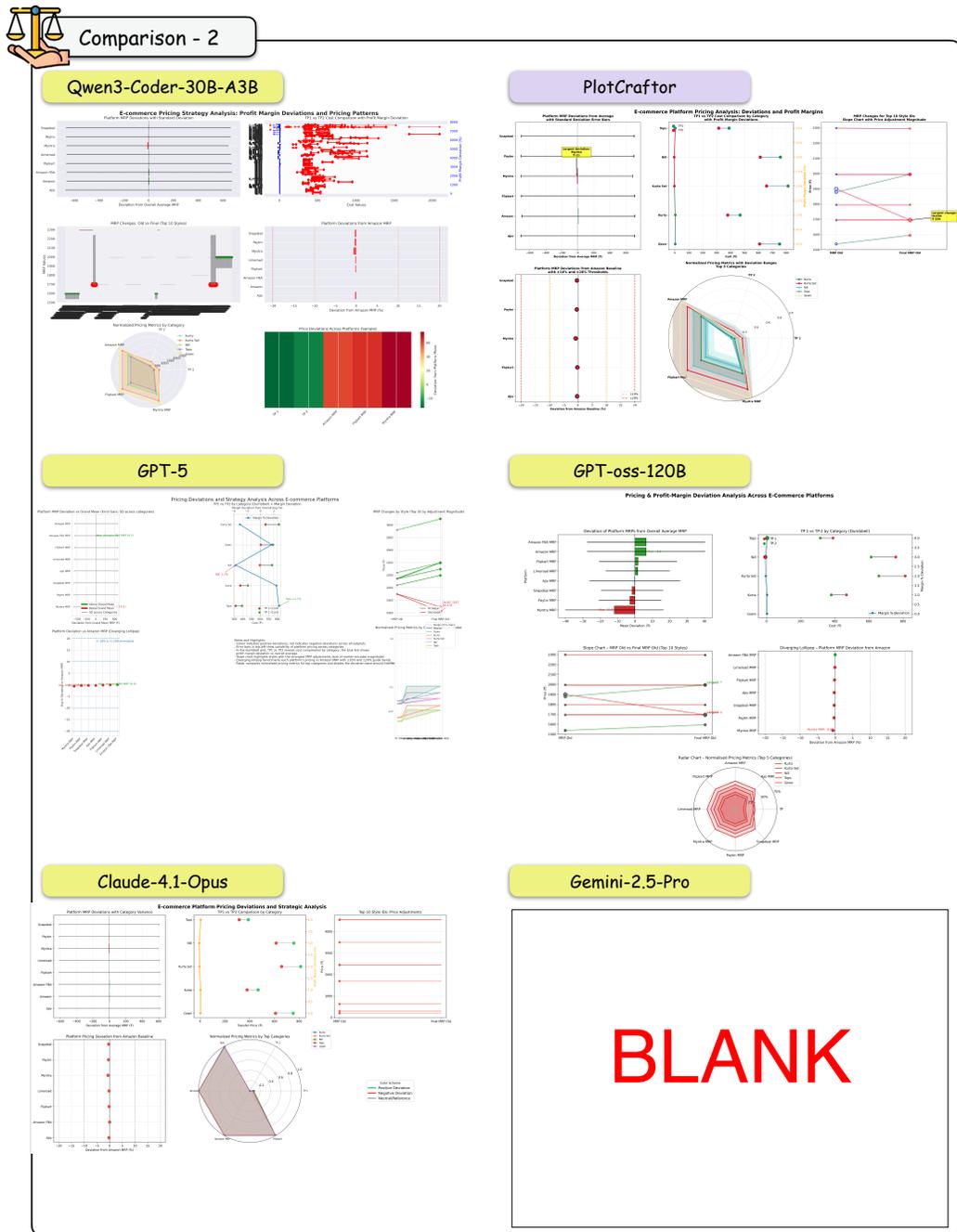


Figure 24: Qualitative comparison for the task in Instruction G.1. PlotCrafter correctly generates the complex 3x2 composite grid, whereas other models fail on layout generation (Qwen3-Coder, GPT-oss-120B), produce distorted outputs (GPT-5), have minor compositional flaws (Claude-4.1-Opus), or fail completely (Gemini-2.5-Pro).

Comparison 4 For the demanding 3x3 temporal grid specified in Instruction G.1, Figure 26 shows that PlotCrafter is the only model to produce a fully correct visualization. The base model, Qwen3-Coder-30B-A3B, fails to generate any output. The other models, including GPT-5, GLM-4.5, Claude-4-Sonnet, and Gemini-2.5-Pro, all manage to generate visualizations but fail to adhere to the prompt, exhibiting various chart type errors across the subplots.

2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753

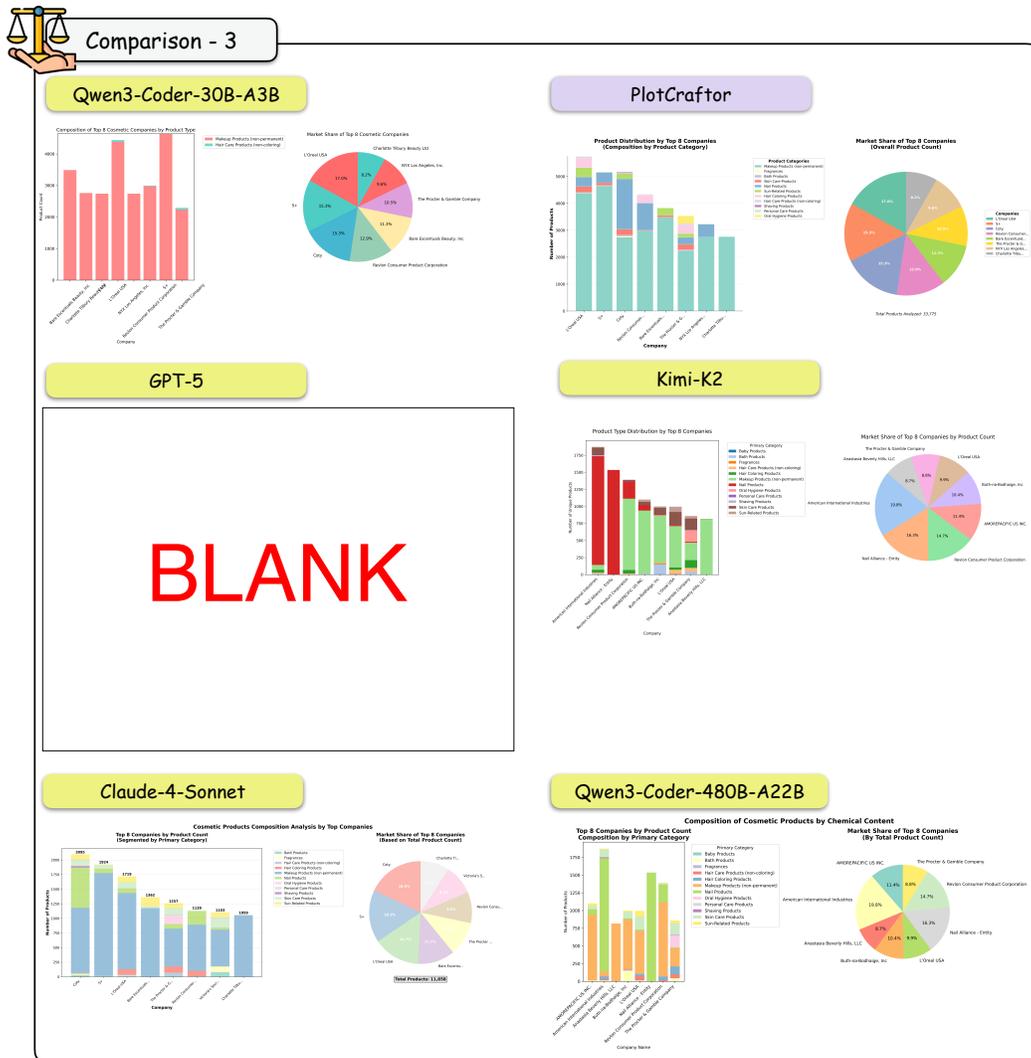


Figure 25: Qualitative comparison for the simpler task in Instruction G.1. While PlotCraftor and several other capable models generate correct visualizations, this figure highlights key failures, including an incorrect chart type from Qwen3-Coder-30B-A3B and a blank output from GPT-5.

Comparison 5 The task in Instruction G.1 requires a composite plot with line and stacked area charts. Figure 27 illustrates that PlotCraftor successfully generates the required visualization with professional formatting. In contrast, Qwen3-Coder-30B-A3B generates an incorrect chart type, failing to meet the core requirement. GPT-5 produces a visualization that, while functionally similar, suffers from a cramped and poorly organized layout. Other powerful models like ChatGPT-4o, Claude-4.1-Opus, and Gemini-2.5-Pro also attempted the task with varying degrees of success and failure as depicted.

G.2 SCALING COMPARISON

Figure 28 and Figure 29 present scatter plots of average model scores as a function of model size on PlotCraft’s Easy and Hard tasks, respectively. These plots visually confirm that the benefits of model scaling are highly dependent on task difficulty. For models under 100B parameters, performance on Easy tasks scales rapidly with size, while performance on Hard tasks remains flat, only improving for models beyond the 100B threshold. This disparity is mirrored in supervised fine-tuning (SFT): smaller models can be fine-tuned to near-proprietary levels on Easy tasks, yet SFT provides minimal

2754 benefit for Hard tasks. This indicates that solving complex visualization challenges may rely more
2755 on the emergent reasoning abilities that come with scale than on task-specific fine-tuning.
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807

2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861

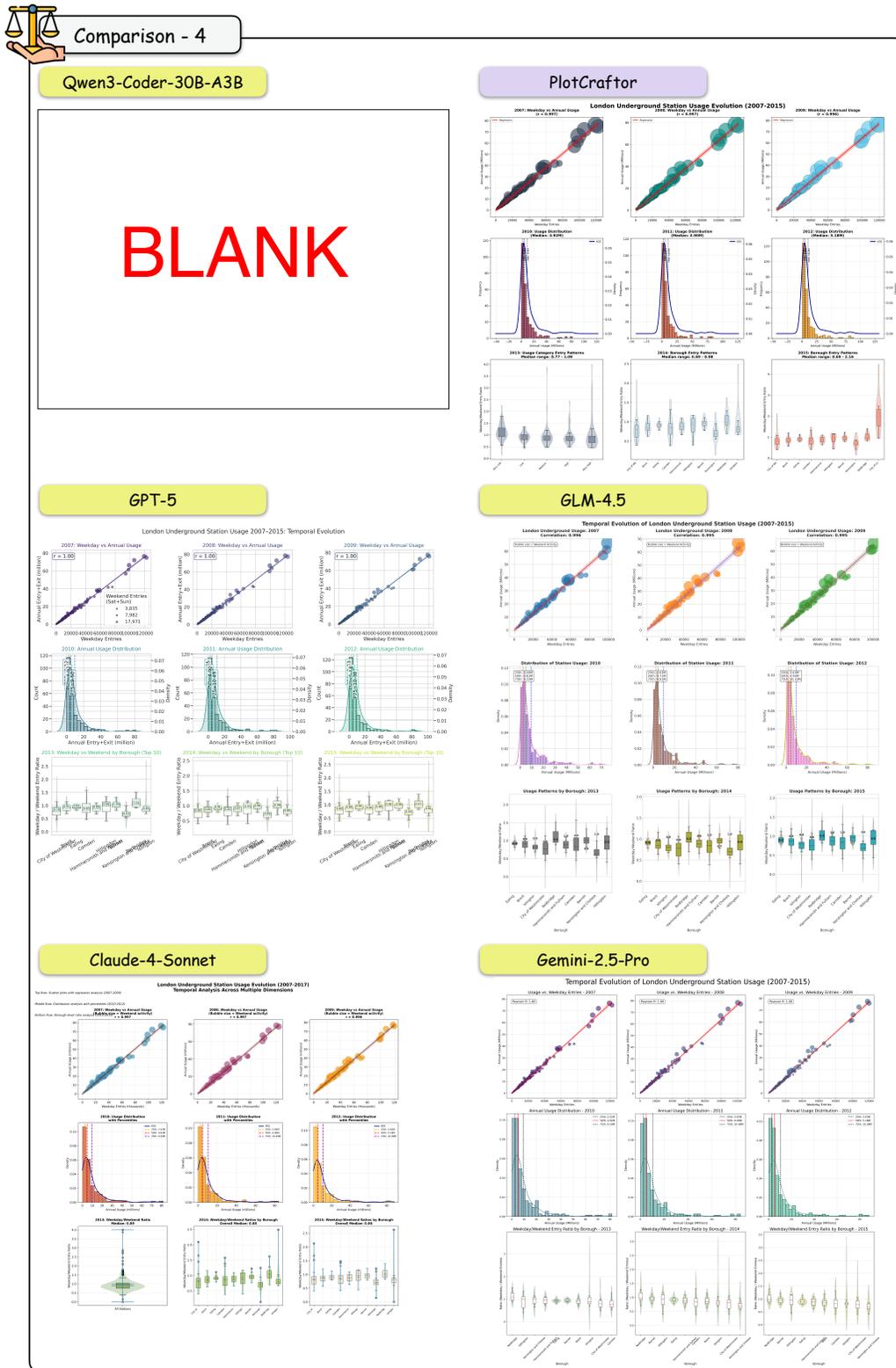


Figure 26: Qualitative comparison for the highly complex 3x3 grid in Instruction G.1. PlotCrafter is the only model to successfully fulfill all requirements. Other models either failed completely (Qwen3-Coder-30B-A3B) or produced visualizations with significant chart type errors (GPT-5, GLM-4.5, Claude-4-Sonnet, Gemini-2.5-Pro).

2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915

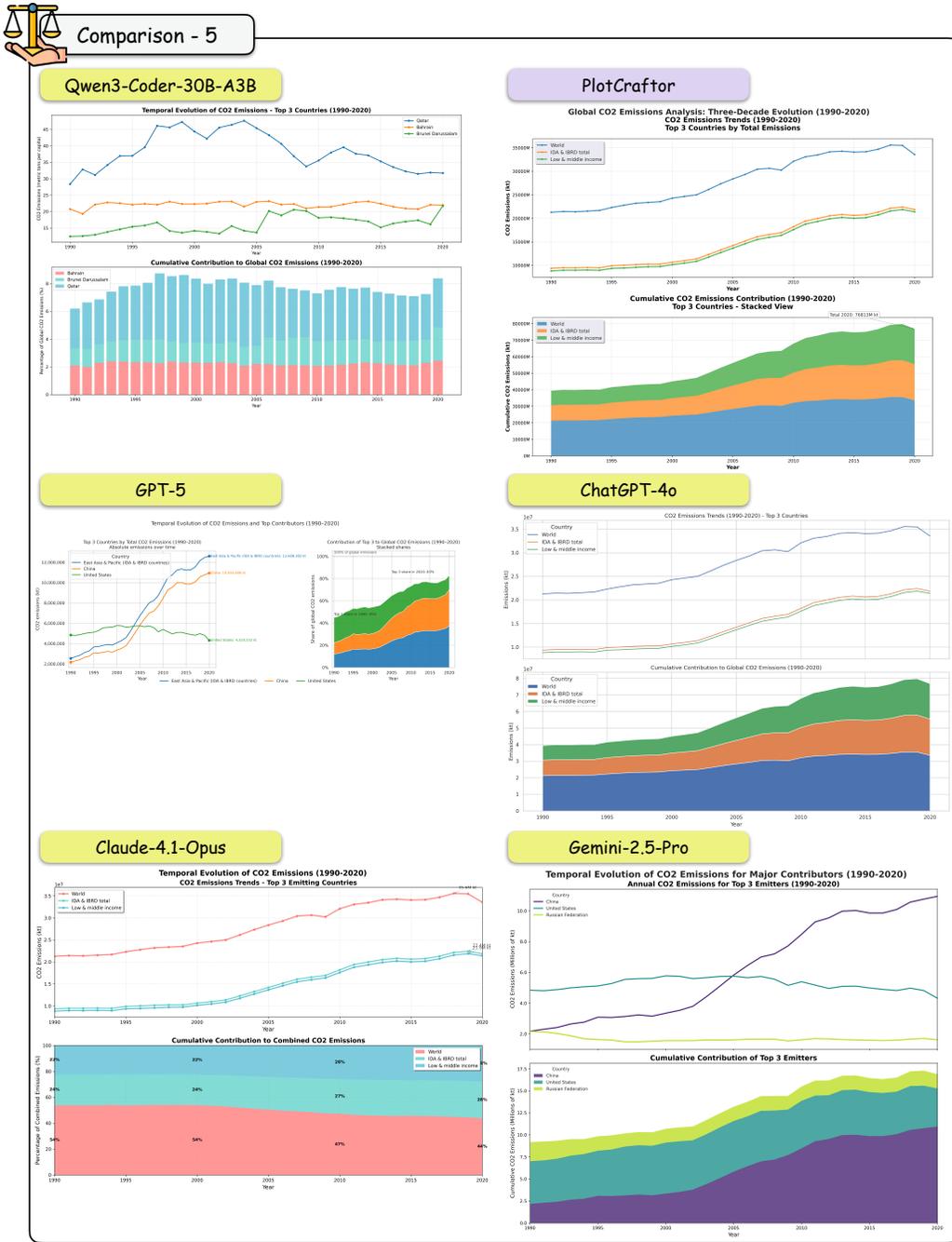


Figure 27: Qualitative comparison for the composite time-series task in Instruction G.1. PlotCrafter’s correct and well-formatted output is contrasted with other models that produced incorrect chart types (Qwen3-Coder-30B-A3B) or suffered from poor, cramped layouts (GPT-5).

2970
2971
2972
2973
2974

Compliance Metrics				Quality Metrics				
Layout	Type	Visual	Task	Clarity	Layout	Color	Text	Format
0.61	0.58	0.41	0.50	0.39	0.51	0.63	0.58	0.70

2975
2976
2977

Table 8: Cohen’s Kappa scores for agreement between our Claude-4-Sonnet judge and human evaluations, categorized by Compliance and Quality metrics.

2978
2979
2980
2981

Compliance Metrics				Quality Metrics				
Layout	Type	Visual	Task	Clarity	Layout	Color	Text	Format
0.64	0.62	0.53	0.55	0.36	0.52	0.54	0.61	0.67

2982
2983
2984

Table 9: Cohen’s Kappa scores for agreement between our ChatGPT-4o judge and human evaluations, categorized by Compliance and Quality metrics.

2985
2986
2987

H CORRELATION WITH HUMAN EVALUATION DETAILS

2988
2989
2990
2991

In our study, we also evaluated the reliability of other prominent models, Claude-4-Sonnet and ChatGPT-4o, as potential automated judges. The agreement between these models and our human evaluations, as measured by Cohen’s Kappa scores, is presented in Table 9 and Table app-tab:gpt-judge, respectively.

2992
2993
2994
2995
2996
2997
2998
2999

The results reveal a consistent and critical weakness in both models’ visual analysis capabilities. While they achieve moderate to substantial agreement on most compliance and formatting metrics, their reliability drops sharply on the Clarity metric, which is specifically designed to assess element overlap. With Kappa scores of only 0.39 for Claude-4-Sonnet and 0.36 for ChatGPT-4o on this metric, it is evident that both models struggle to effectively identify and penalize overlapping elements. This deficiency often leads them to assign inaccurately high quality scores to charts with significant readability issues, limiting their viability as standalone judges for complex data visualizations.

3000
3001

I ERROR ANALYSIS

3002
3003
3004

Error Case 1: Severe Element Overlap

3005
3006
3007
3008
3009
3010
3011

This case demonstrates a critical failure in spatial reasoning, resulting in severe overlap between chart elements and text. As shown in Figure 30, the generated visualization suffers from widespread element occlusion that renders it largely unreadable. Specific failures include: overlapping pie chart labels in subplot (0,1); annotations clashing with the line plot in subplot (0,0); indecipherable stacked text in subplot (0,2); a legend obscuring the plot in subplot (1,0); and a complete overlap of two distinct charts in subplot (1,2). These errors indicate a fundamental inability of the model to manage spatial allocation within a complex multi-plot layout.

3012
3013

Error Case 2: Conflicting Layout Managers

3014
3015
3016
3017

This case illustrates a technical failure where incompatible layout management commands lead to a complete collapse of the plotting canvas. The model-generated code produces a blank image because of a conflict between Matplotlib’s `constrained_layout` engine and the subsequent addition of figure-level elements, particularly `fig.legend()`.

3018
3019

The core of the issue lies in the subplot initialization:

3020
3021
3022
3023

```
fig, (ax_top, ax_bottom) = plt.subplots(
2, 1, figsize=(12, 16), sharex=True, \texttt{constrained\_layout=}
True),
gridspec_kw=dict(height_ratios=[1, 1])
)
```

3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077

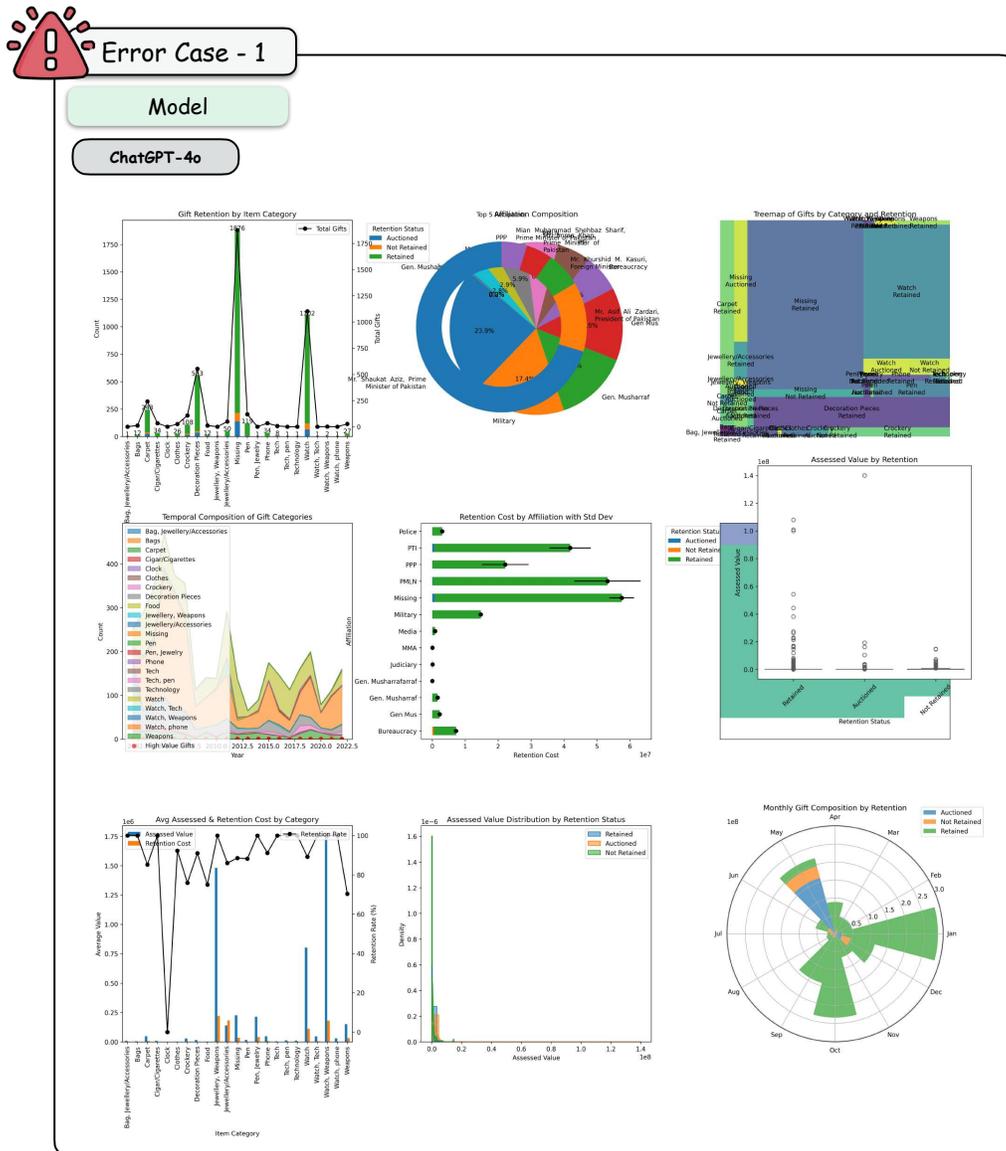


Figure 30: An error case generated by **ChatGPT-4o**, characterized by severe and pervasive element overlap. The visualization fails due to multiple instances of text, labels, legends, and entire subplots occluding one another, making the chart uninterpretable and highlighting a deficiency in layout management.

Here, `constrained_layout=True` activates a sophisticated automatic layout engine designed to prevent the overlap of axes labels and titles by adjusting subplot positions. However, this engine has known limitations when interacting with elements added directly to the figure canvas, such as:

```
fig.suptitle(...)
fig.text(...)
fig.legend(...)
fig.text(...)
```

The conflict arises because `constrained_layout` is primarily designed to manage subplots (Axes) and their immediate decorations. It cannot properly account for the space consumed by figure-level objects like `fig.legend()`, which are placed in the figure's coordinate system independently of the subplot grid.

3078 The failure proceeds as follows:
3079

- 3080 1. **Engine Activation:** `constrained_layout=True` instructs Matplotlib to manage all
3081 subplot layouts automatically.
- 3082 2. **Content Preparation:** The plotting functions successfully prepare the bar charts and text
3083 within the `ax_top` and `ax_bottom` axes objects in memory.
- 3084 3. **Conflict Introduction:** The code then adds a `fig.legend()` directly to the figure. The
3085 layout engine does not natively know how to reserve space for this object while also opti-
3086 mizing the subplot positions.
- 3087 4. **Layout Calculation Failure:** When `plt.show()` is called, the rendering backend exe-
3088 cutes the `constrained_layout` algorithm. It attempts to find a solution that accommo-
3089 dates the subplots, their labels, and the "external" figure-level elements. Unable to converge
3090 on a stable solution, the algorithm fails.
- 3091 5. **Canvas Collapse:** A common outcome of this failure is that the layout engine allocates
3092 zero (or a near-zero) height and width to the subplots in a misguided attempt to make space
3093 for the other elements. Consequently, the primary drawing areas vanish.
3094

3095 The final rendered output is a blank figure canvas. While the figure title or legend might be present,
3096 the core subplots are invisible because their dimensions have been reduced to zero, as shown in
3097 Figure 31.

3098
3099 **Error Case 3: Critical Rendering and Layout Failures** This case, generated by Gemini-2.5-
3100 Pro, demonstrates a combination of critical rendering failures and severe layout issues that render
3101 the chart unusable, as shown in Figure 32.

3102 The most significant issue is a rendering failure in the main plot, which displays as a solid gray
3103 area instead of the intended visualization, indicating a fundamental error in the code's data-to-visual
3104 mapping. In addition to this critical failure, the chart suffers from severe clarity problems. All x-
3105 axis tick labels are collapsed and stacked on top of one another, making them completely illegible.
3106 Furthermore, the legend in the supplementary plot overlaps with the chart's content, obscuring key
3107 information and making that part of the visualization difficult to interpret.

3108
3109
3110 **Reference Visualization for Error Case 3** Figure 33 displays the ground truth for the task previ-
3111 ously shown in Figure 32. Unlike the model's output, which suffered from severe overlapping, the
3112 reference code achieves a clean and readable layout by strategically selecting top cities based on the
3113 data description. This demonstrates that the layout issues were due to the model's lack of reasoning
3114 capability rather than task infeasibility.

3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131

3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185



Figure 31: A failure case generated by **GPT-5**, resulting from a conflict between Matplotlib’s layout managers. The use of `constrained_layout=True` in conjunction with figure-level elements like `fig.legend()` causes the layout engine to fail, collapsing the subplot dimensions to zero and producing a blank image.

3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239

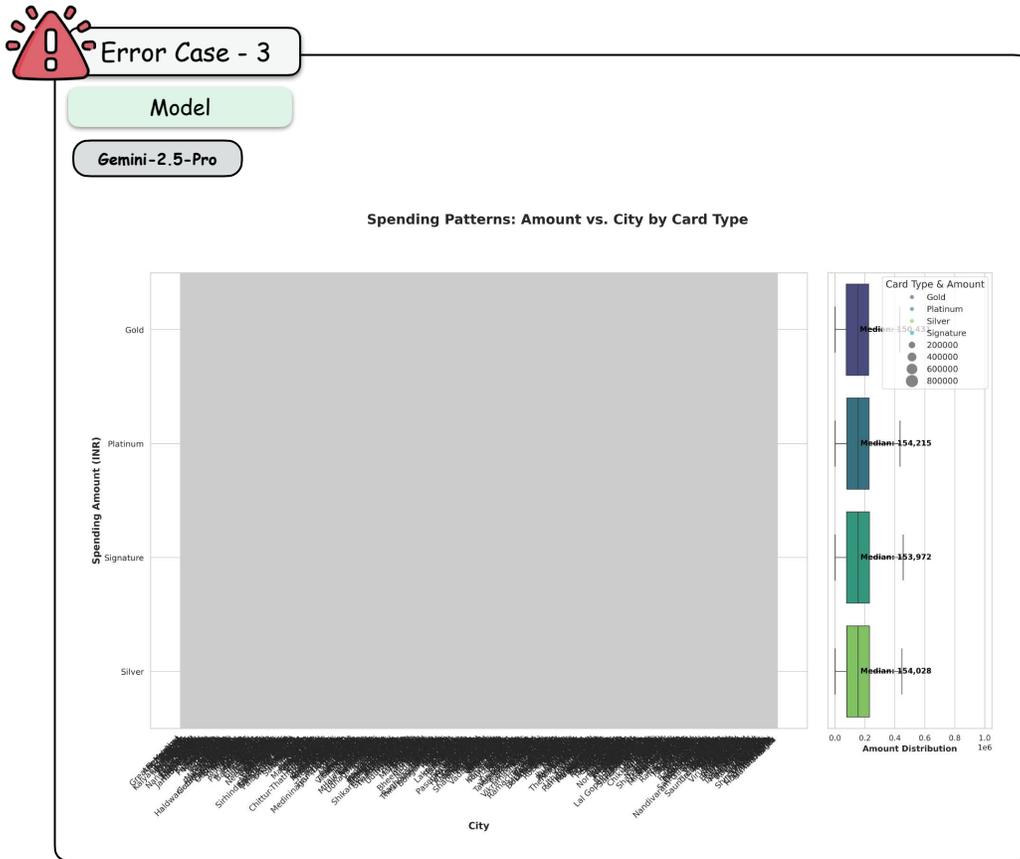


Figure 32: An error case generated by **Gemini-2.5-Pro** exhibiting both a critical rendering failure and severe layout issues. The main plot is incorrectly rendered as a solid gray block, while overlapping x-axis labels and legends make other parts of the visualization unreadable.

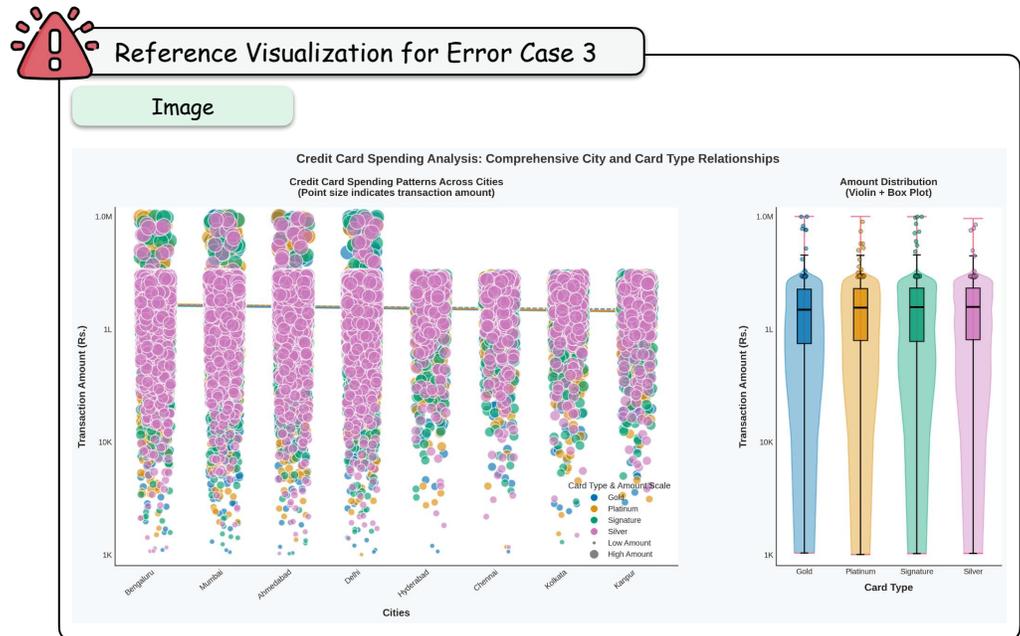


Figure 33: Human-written reference visualization for Error Case 3.

3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293

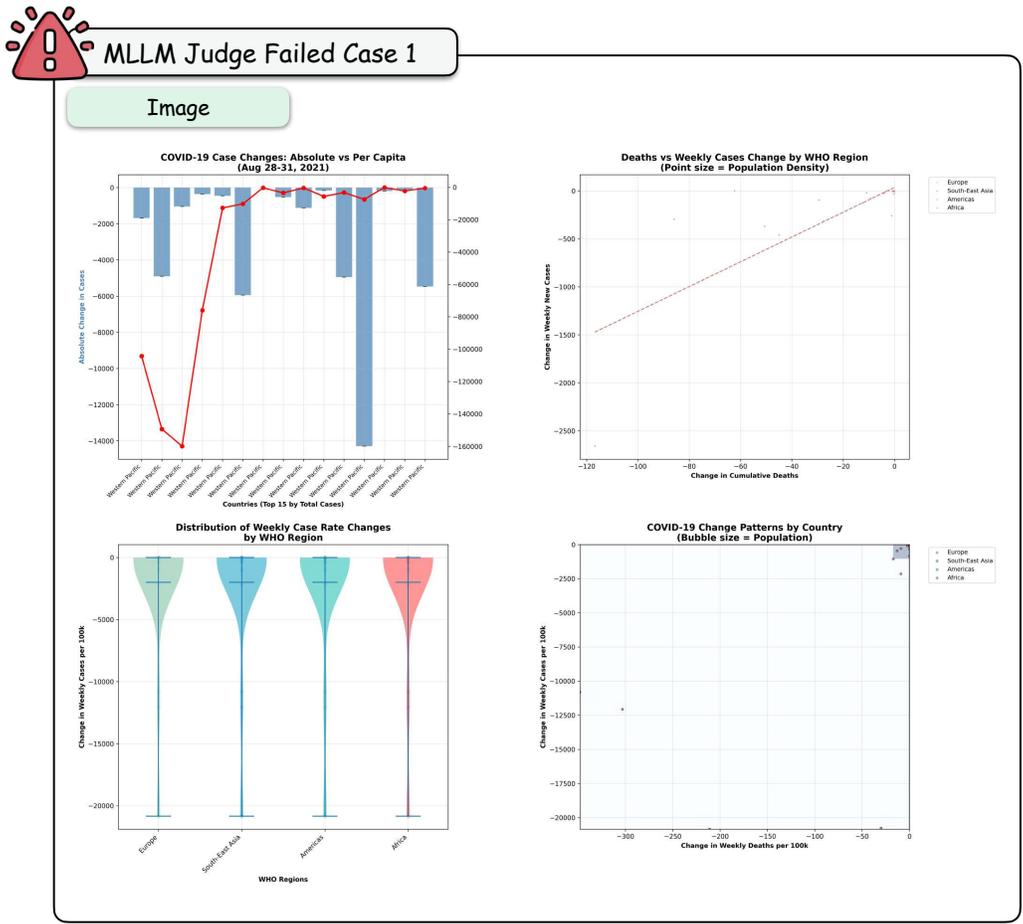


Figure 34: MLLM visual judge failed case 1

Table 10: **Quantitative results on PlotCraft across three difficulty settings: Hard, Medium, and Easy.** The presented scores are the average results over 5 independent evals. Task-Comp. and Quality denote the total scores for the Task Compliance (out of a maximum of 4) and Chart Quality (out of a maximum of 10) sub-metrics, respectively. The metrics are presented separately for Single-Turn Generation and Multi-Turn Refinement.

Model	Single-Turn Generation			Multi-Turn Refinement		
	Pass Rate (%)	Task-Comp.	Quality	Pass Rate (%)	Task-Comp.	Quality
HARD						
<i>Closed-source LLMs</i>						
Claude-4.1-Opus (Anthropic, 2023)	64.85	0.79 ± 0.02	2.59 ± 0.01	68.26	0.89 ± 0.01	2.73 ± 0.02
Claude-4-Sonnet (Anthropic, 2023)	46.75	0.59 ± 0.01	1.71 ± 0.03	57.40	0.65 ± 0.04	2.17 ± 0.02
Gemini-2.5-Pro (Team, 2024)	5.33	0.07 ± 0.09	0.17 ± 0.08	28.99	0.29 ± 0.02	1.30 ± 0.11
ChatGPT-4o-Latest (OpenAI, 2023)	34.32	0.31 ± 0.01	1.12 ± 0.03	47.34	0.37 ± 0.05	1.86 ± 0.02
GPT-5 (OpenAI, 2025)	44.97	0.98 ± 0.01	1.12 ± 0.01	53.85	0.98 ± 0.02	2.03 ± 0.01
<i>Open-source LLMs</i>						
Kimi-K2 (Team et al., 2025b)	29.70	0.28 ± 0.01	0.82 ± 0.06	32.93	0.40 ± 0.01	1.37 ± 0.01
DeepSeek-V3.1 (DeepSeek-AI & etc., 2024)	11.83	0.14 ± 0.02	0.40 ± 0.00	10.06	0.11 ± 0.03	0.43 ± 0.05
GLM-4.5 (Team et al., 2025a)	15.98	0.18 ± 0.02	0.44 ± 0.04	43.79	0.51 ± 0.01	1.64 ± 0.01
GPT-oss-120B (Agarwal et al., 2025)	12.12	0.12 ± 0.01	0.32 ± 0.02	13.17	0.13 ± 0.03	0.42 ± 0.04
GPT-oss-20B (Agarwal et al., 2025)	8.48	0.10 ± 0.01	0.21 ± 0.02	12.57	0.11 ± 0.01	0.31 ± 0.03
Seed-Coder-8B (Seed et al., 2025)	3.55	0.02 ± 0.00	0.08 ± 0.02	37.28	0.30 ± 0.01	1.22 ± 0.01
VisCoder-7B (Ni et al., 2025)	5.33	0.02 ± 0.00	0.13 ± 0.00	50.30	0.38 ± 0.01	2.04 ± 0.01
Qwen3-Coder-480B-A35B (Hui et al., 2024)	39.64	0.48 ± 0.02	1.17 ± 0.06	56.21	0.63 ± 0.01	1.98 ± 0.07
Qwen3-Coder-30B-A3B (Hui et al., 2024)	19.53	0.18 ± 0.01	0.47 ± 0.02	56.80	0.50 ± 0.01	1.76 ± 0.06
PlotCrafter-30B-A3B (Ours)	36.09	0.43 ± 0.01	1.45 ± 0.01	54.44	0.63 ± 0.01	2.11 ± 0.01
MEDIUM						
<i>Closed-source LLMs</i>						
Claude-4.1-Opus (Anthropic, 2023)	72.33	1.97 ± 0.02	4.01 ± 0.02	82.61	2.06 ± 0.02	5.20 ± 0.02
Claude-4-Sonnet (Anthropic, 2023)	72.39	1.58 ± 0.02	3.39 ± 0.02	85.89	1.96 ± 0.02	4.74 ± 0.02
Gemini-2.5-Pro (Team, 2024)	39.26	1.04 ± 0.01	1.77 ± 0.01	67.48	1.74 ± 0.02	4.07 ± 0.02
ChatGPT-4o-Latest (OpenAI, 2023)	65.64	1.48 ± 0.01	2.96 ± 0.02	74.85	1.48 ± 0.02	4.10 ± 0.02
GPT-5 (OpenAI, 2025)	77.30	2.15 ± 0.02	3.29 ± 0.02	80.37	2.04 ± 0.02	4.76 ± 0.02
<i>Open-source LLMs</i>						
Kimi-K2 (Team et al., 2025b)	62.89	1.38 ± 0.01	2.99 ± 0.02	68.94	1.47 ± 0.02	4.05 ± 0.02
DeepSeek-V3.1 (DeepSeek-AI & etc., 2024)	53.99	1.15 ± 0.01	2.54 ± 0.01	42.33	1.01 ± 0.01	2.50 ± 0.01
GLM-4.5 (Team et al., 2025a)	44.17	1.20 ± 0.01	1.96 ± 0.01	68.10	1.50 ± 0.02	3.66 ± 0.02
GPT-oss-120B (Agarwal et al., 2025)	51.57	1.14 ± 0.01	2.13 ± 0.01	32.30	0.76 ± 0.01	1.69 ± 0.01
GPT-oss-20B (Agarwal et al., 2025)	45.91	1.08 ± 0.01	1.94 ± 0.01	36.65	0.93 ± 0.01	1.99 ± 0.01
Seed-Coder-8B (Seed et al., 2025)	28.83	0.60 ± 0.01	1.14 ± 0.01	50.31	0.99 ± 0.01	2.50 ± 0.01
VisCoder-7B (Ni et al., 2025)	11.04	0.17 ± 0.00	0.47 ± 0.00	39.88	0.66 ± 0.01	1.92 ± 0.01
Qwen3-Coder-480B-A35B (Hui et al., 2024)	60.12	1.35 ± 0.01	2.63 ± 0.01	82.82	1.66 ± 0.02	4.31 ± 0.02
Qwen3-Coder-30B-A3B (Hui et al., 2024)	53.37	1.07 ± 0.01	2.27 ± 0.01	74.85	1.48 ± 0.02	3.71 ± 0.02
PlotCrafter-30B-A3B (Ours)	68.10	1.72 ± 0.02	3.54 ± 0.02	83.80	1.79 ± 0.02	4.44 ± 0.02
EASY						
<i>Closed-source LLMs</i>						
Claude-4.1-Opus (Anthropic, 2023)	92.26	3.08 ± 0.03	6.10 ± 0.03	94.27	3.29 ± 0.03	7.89 ± 0.03
Claude-4-Sonnet (Anthropic, 2023)	88.68	3.09 ± 0.03	7.01 ± 0.03	93.08	3.11 ± 0.03	7.67 ± 0.03
Gemini-2.5-Pro (Team, 2024)	81.76	2.43 ± 0.02	5.16 ± 0.02	81.76	2.58 ± 0.02	6.18 ± 0.03
ChatGPT-4o-Latest (OpenAI, 2023)	92.45	3.08 ± 0.03	6.06 ± 0.03	86.79	2.73 ± 0.02	6.88 ± 0.03
GPT-5 (OpenAI, 2025)	88.68	2.20 ± 0.02	4.30 ± 0.02	89.31	2.43 ± 0.02	6.34 ± 0.03
<i>Open-source LLMs</i>						
Kimi-K2 (Team et al., 2025b)	89.68	2.99 ± 0.03	6.46 ± 0.03	82.80	2.66 ± 0.02	6.92 ± 0.03
DeepSeek-V3.1 (DeepSeek-AI & etc., 2024)	88.68	2.97 ± 0.03	6.25 ± 0.03	59.75	1.96 ± 0.02	4.69 ± 0.02
GLM-4.5 (Team et al., 2025a)	71.70	2.44 ± 0.02	5.20 ± 0.02	84.28	2.67 ± 0.02	6.60 ± 0.03
GPT-oss-120B (Agarwal et al., 2025)	83.23	2.79 ± 0.02	5.77 ± 0.02	44.59	1.46 ± 0.01	3.59 ± 0.01
GPT-oss-20B (Agarwal et al., 2025)	81.94	2.55 ± 0.02	5.30 ± 0.02	54.14	1.68 ± 0.02	4.25 ± 0.02
Seed-Coder-8B (Seed et al., 2025)	66.67	2.04 ± 0.02	4.13 ± 0.02	84.91	2.57 ± 0.02	6.09 ± 0.03
VisCoder-7B (Ni et al., 2025)	61.64	1.86 ± 0.02	4.02 ± 0.02	65.41	2.00 ± 0.02	5.02 ± 0.02
Qwen3-Coder-480B-A35B (Hui et al., 2024)	85.53	2.92 ± 0.03	5.98 ± 0.03	89.94	3.03 ± 0.03	7.24 ± 0.03
Qwen3-Coder-30B-A3B (Hui et al., 2024)	86.79	2.77 ± 0.02	5.99 ± 0.03	88.68	2.74 ± 0.02	7.23 ± 0.03
PlotCrafter-30B-A3B (Ours)	90.57	3.13 ± 0.03	7.45 ± 0.03	94.34	2.94 ± 0.03	7.85 ± 0.03