

ChartCoder: Advancing Multimodal Large Language Model for Chart-to-Code Generation

Anonymous ACL submission

Abstract

Multimodal Large Language Models (MLLMs) have demonstrated remarkable capabilities in chart understanding tasks. However, interpreting charts with textual descriptions often leads to information loss, as it fails to fully capture the dense information embedded in charts. In contrast, parsing charts into code provides loss-less representations that can effectively contain all critical details. Although existing open-source MLLMs have achieved success in chart understanding tasks, they still face two major challenges when applied to chart-to-code tasks.: (1) Low executability and poor restoration of chart details in the generated code and (2) Lack of large-scale and diverse training data. To address these challenges, we propose **ChartCoder**, the first dedicated chart-to-code MLLM, which leverages Code LLMs as the language backbone to enhance the executability of the generated code. Furthermore, we introduce **Chart2Code-160k**, the first large-scale and diverse dataset for chart-to-code generation, and propose the **Snippet-of-Thought (SoT)** method, which transforms direct chart-to-code generation data into step-by-step generation. Experiments demonstrate that ChartCoder, with only 7B parameters, surpasses existing open-source MLLMs on chart-to-code benchmarks, achieving superior chart restoration and code executability.

1 Introduction

Recently, Multimodal Large Language Models (MLLMs) have demonstrated remarkable capabilities in addressing a wide range of visual tasks, such as captioning and question answering (Zhang et al., 2024c; Wang et al., 2024b; Chen et al., 2023; Zeng et al., 2024). However, current models still face significant challenges in understanding and analyzing the dense visual information present in complex and informative images. As a significant form of information-intensive images, charts contain com-

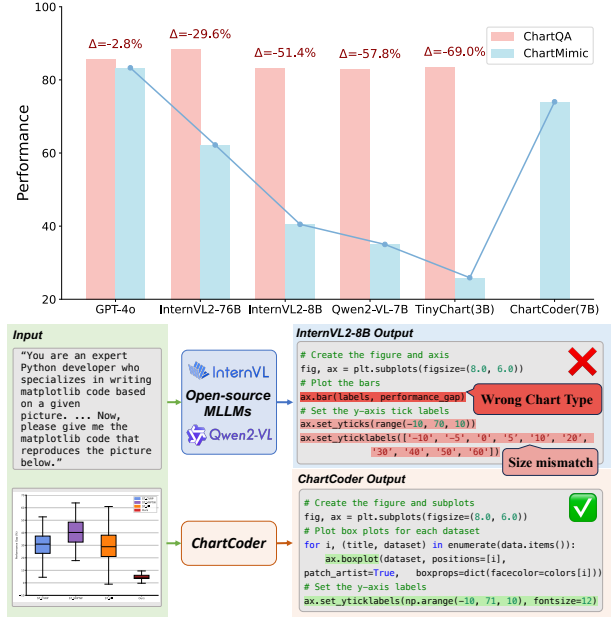


Figure 1: Comparison of existing MLLMs performance on ChartQA and ChartMimic benchmarks. In the chart-to-code task, open-source MLLMs struggle with mismatches in chart types and sizes, whereas ChartCoder generates accurate code.

plex information such as data and structures, playing a pivotal role in effectively presenting details. The automation of chart comprehension and summarization has garnered significant attention from the research community. To advance chart understanding tasks, current studies leverage existing MLLMs and perform supervised fine-tuning (SFT) on various large-scale chart understanding datasets, such as chart question answering (Methani et al., 2020) and chart-to-text generation (Kantharaj et al., 2022), achieving state-of-the-art performance on existing chart understanding benchmarks.

However, existing works generally treat charts as natural images and fine-tune models by generating natural language descriptions (Zhang et al., 2024b; Han et al., 2023; Meng et al., 2024). This inevitably overlooks the dense information embedded within the charts, resulting in inefficiencies in analysis and

comprehension. On the other hand, parsing a chart into code offers a lossless representation, providing a more efficient and comprehensive approach to understanding the chart by accurately capturing and summarizing all its information. Recent works (Shi et al., 2024; Wu et al., 2024; Xia et al., 2024) have proposed various chart-to-code benchmarks, aiming to evaluate the chart reasoning abilities through code. However, current open-source MLLMs are not well-aligned with code generation tasks (Zhang et al., 2024a), resulting in poor performance in parsing charts into corresponding code and limited execution rate of the generated code. As shown in Figure 1, the InternVL2-8B suffers from chart type errors and coordinate size mismatches when converting boxplots to code.

To overcome the above challenges in chart-to-code generation, we first conduct an exploratory attempt by leveraging Code LLMs as the language backbone of the MLLMs and propose **ChartCoder**, the first dedicated chart-to-code MLLM, which incorporates a two-stage training paradigm that contains chart-to-text alignment and chart-to-code instruction tuning. However, compared to chart-to-text, the available chart-to-code dataset is significantly smaller in scale, making it insufficient to support effectively supervised fine-tuning. Therefore, to address the scarcity of data for the chart-to-code domain and train our proposed ChartCoder, we propose the first large-scale diverse and high-quality chart-to-code dataset named **Chart2Code-160k** along with the model, which contains 160k diverse chart-code pairs with 27 chart types. To enhance the model’s capacity to capture critical information, such as chart types and data values, and strengthen its reasoning ability, we propose the **Snippet-of-Thought (SoT)** method, which emphasizes critical information and optimizes the chart-to-code reasoning process. Specifically, we sample 50k chart-code pairs from the Chart2Code-160k, then utilize Chain-of-Thought (CoT) (Wei et al., 2022) and Program-of-Thought (PoT) (Chen et al., 2022) methods to extend direct generation to step-by-step generation, which aims to emphasize critical information in each step. Experimental results show that by utilizing our proposed Chart2Code-160k with the SoT method, ChartCoder, which, with only 7B parameters, outperforms all open-source MLLMs across various chart-to-code benchmarks. As shown in Figure 1, ChartCoder demonstrates a significantly higher ability to generate correct and executable code.

In summary, the main contributions of this work are as follows:

- We propose **ChartCoder**, the first chart-to-code MLLM, which leverages Code LLMs as language backbones. With only 7B parameters, ChartCoder outperforms existing open-source MLLMs on chart-to-code benchmarks.
- We introduce **Chart2Code-160k**, the first large-scale and diverse chart-to-code dataset, consisting of 160k chart-code pairs across 27 chart types.
- We propose **Snippet-of-Thought (SoT)**, transforming direct generation to step-by-step generation to emphasize critical information and enhance reasoning capabilities.

2 Related Works

2.1 Chart Understanding

Chart understanding is a crucial area of research that encompasses both low-level tasks, such as extracting data from visual charts, and high-level tasks, including QA (Masry et al., 2022), summarization (Kantharaj et al., 2022), and re-rendering chart content (Han et al., 2023). Previous approaches (Singh et al., 2019; Methani et al., 2020) have typically relied on pipeline-based methods. However, these pipeline approaches often struggle with error accumulation across different stages, which limits their overall effectiveness and flexibility. Recent works have led to the development of end-to-end MLLMs (Liu et al., 2023b,c; Han et al., 2023; Ye et al., 2023) specifically designed for chart-related tasks. Trained on extensive chart-specific datasets, these chart-domain MLLMs (Xia et al., 2024; Zhang et al., 2024b; Masry et al., 2024; Yang et al., 2024b) have achieved superior performance across various chart-related tasks. However, existing studies typically describe charts in natural language, which inevitably overlooks the dense information embedded within them, leading to inefficiencies in analysis and understanding. In contrast, code serves as a lossless representation of charts, offering a more effective and expressive approach to capturing chart information, thereby facilitating the solution of various chart-related tasks.

2.2 MLLMs For Code

Multimodal code generation has recently garnered much more attention. Several works, such as MM-Code (Li et al., 2024b) and HumanEval-V (Zhang

et al., 2024a), have been developed to evaluate the capability of MLLMs in solving code problems that incorporate visual elements. Design2Code (Si et al., 2024) and Web2Code (Yun et al., 2024) evaluate the performance of MLLMs by focusing on code generation for HTML web page creation. Among the emerging tasks in this domain, chart-to-code generation has attracted significant interest as the visual elements of charts are more complex. This task challenges MLLMs to generate code that accurately reproduces a given chart or visual representation. Recent works like ChartMimic (Shi et al., 2024) evaluate the reasoning ability of MLLMs in this context. Similarly, Plot2Code (Wu et al., 2024) and ChartX (Xia et al., 2024) also evaluate MLLMs code generation ability, especially for text and data reproducibility. To date, no dedicated research has focused on solving the chart-to-code generation problem. Our work is the first to attempt to address this challenge.

3 Chart2Code-160k Dataset

3.1 Direct Chart-to-code Generation

Despite the availability of many chart reasoning instruction-tuning datasets, there is a notable lack of datasets specifically for chart-to-code tasks. Compared to chart reasoning data, chart-to-code data have the following distinct characteristics: (1) *One-to-One Mapping*: Unlike chart reasoning datasets, which could derive multiple question-answer pairs from a single chart, chart-to-code datasets require a one-to-one correspondence, demanding a large number of chart images for training. (2) *Diversity Reflect on Charts*: Unlike the diversity of chart reasoning data, which can be reflected in instructions, the diversity of chart-to-code data primarily lies in the variety of chart types and structures. (3) *Syntax Constraints*: Unlike the flexible natural language answers in chart reasoning tasks, the output code must strictly adhere to programming syntax to ensure executability.

Therefore, collecting a large number of chart-code pairs that meet the above requirements is challenging. Recent studies have demonstrated the feasibility of generating code with LLMs (Xu et al., 2023; Zhang et al., 2024c). Leveraging the one-to-one mapping property of chart-to-code data, we generate code first and execute it to produce the corresponding charts. In this way, we construct the first large-scale and diverse chart-to-code dataset, named **Chart2Code-160k**.

Dataset	Train/Eval	Chart Type	Number
ChartX	Eval	18	6k
Plot2Code	Eval	6	132
ChartMimic	Eval	22	2.4k
ChartLlama	Train	10	7.8k
Chart2Code-160k	Train	27	160k

Table 1: Comparisons of existing chart-to-code datasets.

Specifically, we generate chart-to-code data through the following steps: First, we prompt the LLM to generate keywords within a specific domain and guide it to generate simulated data related to these keywords. Then, to ensure the diversity of chart types, we identify 27 commonly used chart types and manually write seed codes for each as in-context demonstrations. We further provide available functions such as `plt.text()` and parameters such as `hatch='/'` to encourage the generation of more diverse functions and parameters (Xu et al., 2023), resulting in the chart structures more diversely. To enhance the generality of generated code, LLMs are encouraged to use standard libraries such as Matplotlib and Seaborn. Additionally, we explicitly define all parameters within the code itself, eliminating the need for external files such as CSVs. This ensures that the code can be executed directly and accurately to represent the chart details. The final step involved executing the generated code to produce the corresponding chart. We utilize the above process to yield 200k code snippets for charts. After executing the code and filtering out problematic charts, such as those with excessive pixels or ticks, we construct a high-quality dataset of 160k diverse chart-to-code pairs. These pairs are formatted as multimodal instruction-tuning samples in the unified structure of `<chart image, input instruction, output code>`.

3.2 Step-by-step Chart-to-code Generation

Although the dataset described above includes various chart types and structures, most of the generated code follows a similar template, with only certain details (such as colors and values) providing the essential distinguishing information. This can cause chart-to-code generation models to overlook these critical details and thus produce incomplete or incorrect results. To address the above challenge and further improve the reasoning ability of MLLMs, we propose the Snippet-of-Thought (SoT) method to expand direct chart-to-code generation into step-by-step generation formats, which have

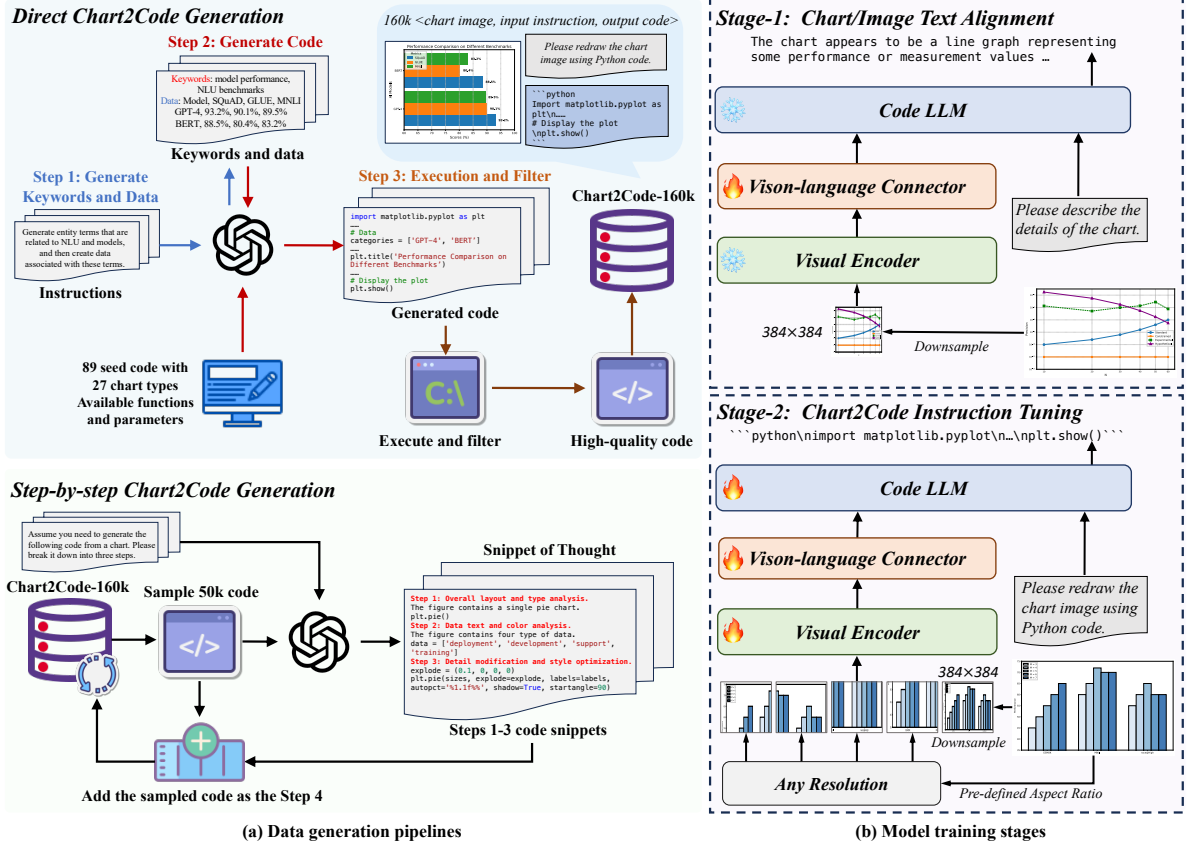


Figure 2: Illustration of Chat2Code dataset generation process and the ChartCoder training process. The dataset generation process is divided into two stages: direct generation and step-by-step generation. In the step-by-step generation, the code processed by the Snippet-of-Thought method is sampled from the Chart2Code-160k generated in the direct generation process. The training process of the ChartCoder also consists of two stages: alignment and instruction tuning.

demonstrated effectiveness in text-to-code generation tasks (Zheng et al., 2023).

Specifically, we adopt the SoT to imitate the human reasoning process, deriving the final code step by step. This process is divided into four steps: Step 1: Generate the chart type and layout, such as `plt.bar()` and `plt.subplot()`. Step 2: Generate the data and corresponding colors used in the chart, such as `data=[10, 15]` and `colors=['#FF0000', '#00FF00']`. Step 3: Generate critical details of the chart, such as `hatch='/'` and `loc='upper left'`. Step 4: Generate the complete and final code. Different from CoT and PoT, we incorporate textual explanations and code snippets for each step to emphasize key information enhance the reasoning process and produce comprehensive outputs.

However, directly instructing the LLM to generate step-by-step code may lead to hallucinations, causing inconsistencies between intermediate code snippets and the final executable code. To maintain consistency among code snippets, we reformulated the step-by-step code data generation into a two-

step process involving code generation and decomposition. We sample 50k chart-code pairs from the previously generated 160k data pairs and encourage the LLM to decompose the original code into the required textual explanations and code snippets of Steps 1–3, then concatenate the complete code in Step 4. To further mitigate hallucinations, such as undefined values or parameters in Steps 1 and 2, we used placeholder or default parameters during code decomposition to ensure the construction of consistent and reliable step-by-step code.

3.3 Dataset Analysis

Chart2Code-160k dataset provides three key advantages: (1) *The First Large-Scale Dataset*: It contains 160k data pairs for instruction tuning, significantly surpassing the size of previous datasets. (2) *Diverse Chart Structures and Types*: It includes 27 different chart types, with diverse structures enabled by a wide variety of functions and parameters in the code. (3) *Syntactically Correct and Executable Code*: All corresponding code is syntactically correct and executable, with explicitly

defined parameters that ensure precise alignment between chart structures and code representations. The comparisons of Chart2Code-160k with relevant chart-to-code datasets are listed in Table 1. Although numerous chart-related datasets have been proposed (Liu et al., 2023a; Han et al., 2023), none have sufficiently addressed the challenge of chart-to-code tasks. Our proposed Chart2Code-160k fills this gap, equipping the model with advanced capabilities for downstream chart tasks.

4 ChartCoder Model

After constructing the Chart2Code-160k, we aimed to leverage the data to enhance the capacities of MLLMs to generate code from charts. Unlike previous methods that rely on general LLMs with a low proportion of code in their training corpus, which limits their coding capabilities, we pioneer the use of Code LLMs to enhance the coding abilities of MLLMs from scratch.

4.1 Model Architecture

Following the standard architecture of MLLMs, ChartCoder consists of three modules: a pre-trained vision encoder (SigLIP-384 (Zhai et al., 2023a)), a vision-language connector (two-layer MLP) and a Code LLM backbone (DeepSeek Coder 6.7B (Guo et al., 2024)). The vision encoder extracts the input image into visual features, and the connector projects it into the word embedding space. LLM backbone then combines visual and textual features to generate responses.

Previous works emphasize the importance of high-resolution input for chart understanding (Liu et al., 2024), as details like textual words may lost in low-resolution images. However, vision transformers (ViTs) like CLIP (Radford et al., 2021) and SigLIP (Zhai et al., 2023b) are constrained to resolutions of 224^2 and 384^2 respectively, which limits their capacities to encode chart images with sufficient detail. To address this, we utilize the Any Resolution strategy (Liu et al., 2024) to resize and patchify chart images to ensure ChartCoder processes high-resolution chart images effectively. Specifically, the input chart image is first resized to a pre-defined optimal aspect ratio, whose height and width are integer multiples of the image resolution. The resized image is then divided into patches of standard resolution and concatenated with a directly downsampled version of the image. This approach preserves both general and

detailed information without requiring the original high-resolution image to be resized into a standard square, thereby avoiding the loss of fine details. Details are shown in Figure 2.

4.2 Model Training

Since we propose to use Code LLMs as the language backbone to enhance the code abilities of MLLMs, existing models do not meet our requirements as their backbones are general LLMs. Thus, to align charts with text and perform supervised fine-tuning for chart-to-code tasks, we adopt the following two-stage training process.

Chart-to-text Alignment. The alignment process aims to endow the model with chart structure perception capability. In this stage, we freeze the language and vision encoder models and pre-train the vision-language connector (Liu et al., 2023c). We collect and filter public chart corpora for alignment, which contains multiple tasks like chart caption and chart-to-table. Specifically, we use the following corpora: (1) UniChart (Masry et al., 2023), (2) Chart-to-Text (Kantharaj et al., 2022), (3) SciCap (Hsu et al., 2021), and (4) SciCap+ (Yang et al., 2024a). Additionally, we incorporate the LLaVA pre-training dataset (Liu et al., 2023c) and our proposed Chart2Code-160k to achieve a more balanced coverage of concepts.

Chart-to-code Instruction-tuning. The second stage focuses on enhancing the model’s capabilities in chart-to-code tasks. In this stage, all three modules are jointly fine-tuned with our proposed Chart2Code-160k, and additional code-related data, such as ChartQA PoT (Zhang et al., 2024b) and ChartLlama chart-to-chart (Han et al., 2023).

5 Experiments

5.1 Baselines and Benchmarks

We compare ChartCoder with existing models in three setups (1) General-domain open-source MLLMs including InternVL2(4B, 8B, 26B, 76B) (Chen et al., 2024), Qwen2-VL(7B, 72B) (Wang et al., 2024a), DeepSeek-VL-7B (Lu et al., 2024), LLaVA-Next(7B) (Li et al., 2024a) and MiniCPM-Llama3-V2.5 (Yao et al., 2024). (2) Proprietary models include GeminiProVision (Team et al., 2023), Claude-3-opus (Anthropic, 2024), GPT-4V (OpenAI, 2023), and GPT-4o (OpenAI, 2024). (3) Chart-domain MLLMs including ChartLlama (Han et al., 2023), ChartAssistant (Meng et al., 2024), Tinychart (Zhang et al., 2024b) and ChartVLM

Model	Params	ChartMimic			Plot2Code			ChartX
		Exec.Rate	Low-Level	High-Level	Pass Rate	Text-Match	Rating	GPT-score
Full score	-	100	100	100	100	100	10	5
<i>Proprietary</i>								
GeminiProVision	-	68.2	53.8	53.3	68.2	53.6	3.69	-
Claude-3-opus	-	83.3	60.5	60.1	84.1	57.5	3.80	-
GPT-4V	-	91.2	76.4	78.9	84.1	57.7	5.58	2.63
GPT-4o	-	93.2	79.0	83.5	88.6	56.3	5.71	-
<i>Open-Source General-Domain</i>								
DeepSeek-VL-7B	7.3B	41.3	19.0	17.6	64.4	32.6	2.26	-
LLaVA-Next-Mistral-7B	7.6B	59.7	20.7	21.3	72.0	38.7	2.87	-
Qwen2-VL-7B	7.0B	67.0	32.9	35.0	68.2	33.8	3.10	1.50
InternVL2-4B	4.2B	66.2	33.8	38.4	66.3	33.4	2.52	1.57
InternVL2-8B	8.1B	61.8	34.4	38.9	77.3	37.1	2.78	1.63
MiniCPM-Llama3-V2.5	8.4B	80.3	36.6	42.1	76.3	37.3	2.61	1.66
InternVL2-26B	26.0B	69.3	41.4	47.4	81.3	43.1	3.42	1.70
Qwen2-VL-72B	72.0B	73.3	54.4	50.9	72.0	53.4	4.26	1.69
InternVL2-Llama3-76B	76.0B	83.2	54.8	62.2	85.6	46.6	3.89	1.74
<i>Open-Source Chart-Domain</i>								
ChartLlama	13B	57.5	24.8	28.1	58.4	40.3	2.32	0.94
ChartAssistant	13B	-	-	-	-	-	-	0.82
TinyChart	3B	42.5	26.3	25.9	43.2	44.6	2.19	1.89
ChartVLM-L	14.3B	19.5	15.8	13.9	-	-	-	1.58
ChartCoder (Ours)	7.0B	91.4	77.4	74.0	87.9	54.5	4.50	2.09

Table 2: Evaluation results of various baseline models. Unless otherwise specified, we directly use the results in the relevant benchmarks. We evaluate models that are not reported in those benchmarks. The best performances of open-source MLLMs are indicated in **bold**.

Model	Chart Types	Layout	Text Content	Data	Style	Clarity
Full score	20	10	20	20	20	10
GPT-4o	18.96	9.59	17.16	15.68	14.66	8.84
InternVL2-Llama3-76B	13.06	8.44	12.59	10.51	8.74	7.87
Qwen2-VL-72B	10.45	7.83	9.92	8.14	7.10	7.47
InternVL2-8B	7.20	6.82	8.81	5.74	5.42	6.64
TinyChart	4.16	5.06	5.22	2.74	3.21	5.58
ChartVLM-L	0.97	3.53	2.48	0.81	0.90	5.25
ChartCoder (Ours)	16.83	9.13	14.77	12.41	12.68	8.29

Table 3: Detailed results of high-level scores on ChartMimic Direct Mimic task. All the subscores of ChartCoder are close to GPT-4o.

(Xia et al., 2024). All the methods are evaluated on the benchmarks ChartMimic (Shi et al., 2024), Plot2Code (Wu et al., 2024) and ChartX (Xia et al., 2024). We revise the Rating calculation in Plot2Code. The original evaluation only considers charts corresponding to executable code, which leads to higher ratings for only generating simple charts. We calculate all the results, which better reflect the impact of complex charts. For all methods, the zero-shot setting was adopted during the evaluation. Details about these benchmarks are shown in the Appendix A.2.

5.2 Main Results

As indicated in Table 2 ChartCoder achieves the best performance among open-source MLLMs

in all the chart-to-code tasks and even better than some proprietary models. Notably, on the most challenging ChartMimic task, ChartCoder surpasses leading small-scale general-domain MLLMs (<20B) such as MiniCPM-Llama3-V2.5 and InternVL2-8B with average scores of **26.7** and **34.6** respectively. The improvement achieved by ChartCoder highlights the effectiveness of our proposed Code LLM as the language backbone, combined with the Chart2Code-160k dataset, in enabling MLLMs to excel in chart understanding and code generation tasks. In addition, ChartCoder also performs better than existing state-of-the-art large-scale MLLMs such as InternVL2-Llama3-76B and chart-domain MLLMs such as TinyChart.

We further illustrate the detailed high-level and

Model	Text	Layout	Type	Color
Full score	100	100	100	100
GPT-4o	81.5	89.8	77.3	67.2
InternVL2-Llama3-76B	54.1	74.5	49.2	41.5
Qwen2-VL-72B	43.2	80.5	54.6	39.4
InternVL2-8B	31.5	51.1	28.6	26.2
TinyChart	9.8	48.2	32.9	14.2
ChartVLM-L	7.7	33.7	17.6	5.2
ChartCoder (Ours)	67.2	95.0	78.5	69.0

Table 4: Detailed results of low-level scores on ChartMimic Direct Mimic task. Three out of four subscores of ChartCoder are even higher than GPT-4o.

low-level scores for the ChartMimic benchmark. The high-level score utilizes GPT-4o to evaluate the detailed similarity between the ground truth and generated chart images in six aspects: chart types, layout, text content, data, style, and clarity. The low-level score is calculated based on a comparison between the ground truth and the generated code, focusing on the code similarities in four aspects: text, layout, type, and color.

Table 3 denotes the high-level results. ChartCoder is the model most comparable to GPT-4o, as the evaluations were conducted by GPT-4o itself, suggesting the actual performance gap may not be as pronounced as it appears. Notably, ChartCoder shows the largest gap with GPT-4o in the ‘data’ category, which highlights the complexity of extracting numerical values from charts, aligning with conclusions from existing chart understanding benchmarks: current MLLMs struggle to directly and accurately extract complete data from complex charts (Wang et al., 2024c; Zhang et al., 2024b).

Table 4 shows the low-level results. ChartCoder even slightly outperforms GPT-4o in layout, type and color, highlighting the diversity of our proposed Chart2Code-160k dataset. However, the text score of the ChartCoder is lower than GPT-4o, which is similar to the results of high-level scores. We believe this is due to the lack of specialized chart OCR-oriented training for our model. Nevertheless, our text accuracy still surpasses that of open-source models, indicating the effectiveness of our proposed ChartCoder model and Chart2Code-160k dataset. We further present some case studies on ChartMimic and compare ChartCoder with GPT-4o, Qwen2-VL-72B, and InternVL-8B. The results are shown in Figure 4, the outputs of ChartCoder are much more similar to the ground truth chart than open-source MLLMs.

Methods	ChartMimic		
	Exec.Rate	Low-Level	High-Level
ChartCoder	91.4	77.4	74.0
<i>Code LLM → General LLM</i>			
DeepSeek LLM	80.6	61.4	63.4
△	-10.8	-16.0	-10.6
<i>Different Visual Encoders</i>			
CLIP-336	91.6	77.3	70.3
△	+0.2	-0.1	-3.7
<i>Without Step-by-step Generation</i>			
w/o SoT	89.2	70.1	65.4
△	-2.2	-7.3	-8.6
<i>Open-source MLLM Finetuned on Chart2Code-160k</i>			
Qwen2-VL-7B	67.0	32.9	35.0
Finetuned Model	83.6	73.4	68.2
△	+16.7	+40.5	+33.2

Table 5: The ablation studies on model architecture and data. The results show that the effectiveness of our proposed code LLM backbone and dataset.

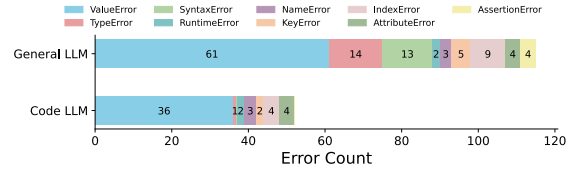


Figure 3: Comparison of error types on ChartMimic direct generation tasks with code and general LLMs as language backbone, respectively.

5.3 Ablation Study

We perform extensive ablation experiments to validate the effectiveness of our proposed model and dataset. We divide the ablation study into three parts, and the results are shown in Table 5 (1) *Code or general LLMs*. To investigate whether employing Code LLMs as language backbone provides specific advantages in chart-to-code tasks and identify the nature of these potential benefits, we replace the Code LLM, DeepSeek Coder 6.7B, with general LLM, DeepSeek LLM 7B (Bi et al., 2024), maintaining the same two-stage training procedures. The result shows that compared with general LLM, utilizing code LLM as the language backbone could significantly improve the execution rate, as well as the low-level and high-level scores. We further analyze the types of errors in the code that failed to execute and find that utilizing code LLMs significantly reduces syntax errors like missing closing quotation marks and type errors like incorrect argument type. (2) *Resolution of vision encoders*. Previous studies have indicated

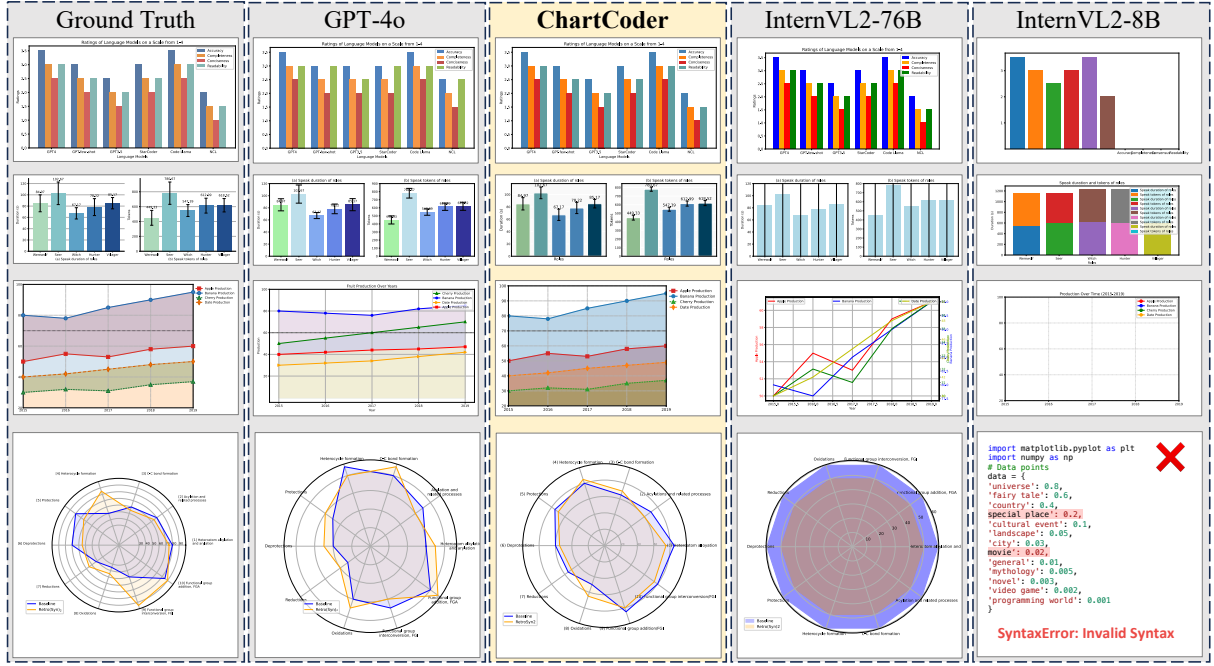


Figure 4: Generated charts of different model outputs after code execution. Our proposed ChartCoder performs significantly better than InternVL2-8B of a similar model scale.

that performance on chart understanding tasks is resolution-dependent, with lower resolutions negatively impacting model performance (Liu et al., 2024). To verify whether resolution affects chart-to-code tasks, we replace SigLIP-384 with CLIP-336 (Radford et al., 2021) and maintain the other setting. The result shows that the resolution of the vision encoder generally does not affect the output code execution rate but slightly influences the high-level chart similarity. Through our analysis, we find that, similar to the challenges in chart understanding, this issue is caused by the negative impact of low resolution on the recognition of text and special symbols. However, as we utilize the Any Resolution strategy, this impact has been reduced significantly. (3) *Dataset effectiveness*. We design two scenarios to illustrate our proposed Chart2Code-160k dataset. Firstly, to evaluate our proposed SoT method to emphasize the critical information in the chart, we remove the 50k step-by-step generation data and train the model using only the direct generation data. The result shows it influences the low-level and high-level scores notably, especially in text content and data, which shows the role of emphasising critical information. Secondly, we select Qwen2-VL-7B as the baseline of open-source MLLM and directly fine-tune it on our proposed Chart2Code-160k datasets. The result illustrates that after fine-tuning, the performance improves significantly on all the metrics, demonstrating the

Benchmark	Image	Image+Code	Δ
MMC T/F	0.76	0.81	6.5%

Table 6: Comparison of the impact of using ChartCoder-generated code as auxiliary information on the MMC True/False benchmark.

effectiveness of Chart2Code-160k.

5.4 Analysis

We further evaluate the role of code in the chart understanding task. We use MiniCPM-Llama3-V2.5 to evaluate two input forms, Image only and Image with Code, on the MMC True/False benchmark (Liu et al., 2023a). The result in Table 6 shows that using code helps the model better understand chart details, especially the chart types and the data they contain. A case study is shown in Figure 6.

5.5 Conclusion

This work aims to tackle the challenge of chart-to-code tasks with MLLMs. First, we propose the ChartCoder, which utilizes Code LLM as the language backbone dedicated to chart-to-code tasks. Second, to solve the scarcity of chart-to-code data, we present the first large-scale and diverse chart-to-code dataset, Chart2Code-160k. Finally, to emphasize the key information, we propose the Snippet-of-Thought (SoT) method to generate step-by-step data. Experiments show that ChartCoder outperforms existing open-source MLLMs.

Limitation

Our study is comprehensive but has certain limitations that we aim to address in future research. Due to constraints in computational resources, we only trained ChartCoder with 7B parameters, which has demonstrated sufficiently good results for now. A larger model could potentially achieve even better performance. Future work may focus on exploring more complex and diverse charts and codes while also experimenting with other image types, such as HTML, to develop a comprehensive multi-modal code large language model.

Ethical Statement

Our research employs publicly available models and datasets with proper citations. This approach minimizes the risk of generating toxic content, leveraging the widely used and non-toxic nature of our datasets and prompts.

References

- Anthropic. 2024. [Introducing the next generation of claude](#).
- DeepSeek-AI Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qishi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wen-Hui Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng Liang, Fangyun Lin, Aixin Liu, Bo Liu (Benjamin Liu), Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghao Lu, Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren, Zehui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Jun-Mei Song, Xuecheng Su, Jingxiang Sun, Yaofeng Sun, Min Tang, Bing-Li Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji Wang, Tong Wu, Yu Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yi Xiong, Hanwei Xu, Ronald X Xu, Yanhong Xu, Dejian Yang, Hu mei You, Shuiping Yu, Xin yuan Yu, Bo Zhang, Haowei Zhang, Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghu Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shunfeng Zhou, Qihao Zhu, and Yuheng Zou. 2024. [Deepseek llm: Scaling open-source language models with longtermism](#). *ArXiv*, abs/2401.02954.
- Lin Chen, Jinsong Li, Xiao wen Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. [Sharegpt4v: Improving large multi-modal models with better captions](#). In *European Conference on Computer Vision*.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022. Program of thoughts

prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*.

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

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.

Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. 2023. Chartllama: A multimodal llm for chart understanding and generation. *arXiv preprint arXiv:2311.16483*.

Ting-Yao Hsu, C Lee Giles, and Ting-Hao’Kenneth’ Huang. 2021. Scicap: Generating captions for scientific figures. *arXiv preprint arXiv:2110.11624*.

Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. 2022. Chart-to-text: A large-scale benchmark for chart summarization. *arXiv preprint arXiv:2203.06486*.

Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024a. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*.

Kaixin Li, Yuchen Tian, Qisheng Hu, Ziyang Luo, Zhiyong Huang, and Jing Ma. 2024b. Mmcode: Benchmarking multimodal large language models for code generation with visually rich programming problems. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 736–783.

Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. 2023a. Mmc: Advancing multimodal chart understanding with large-scale instruction tuning. *arXiv preprint arXiv:2311.10774*.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023b. Improved baselines with visual instruction tuning.

Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. [Llava-next: Improved reasoning, ocr, and world knowledge](#).

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning.

645	Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	699
646	Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhu-	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	700
647	oshu Li, Hao Yang, et al. 2024. Deepseek-vl: towards	Schalkwyk, Andrew M Dai, Anja Hauth, Katie	701
648	real-world vision-language understanding. <i>arXiv</i>	Millican, et al. 2023. Gemini: a family of	702
649	<i>preprint arXiv:2403.05525</i> .	highly capable multimodal models. <i>arXiv preprint</i>	703
		<i>arXiv:2312.11805</i> .	704
650	Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do,	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	705
651	Enamul Hoque, and Shafiq Joty. 2023. Unichart:	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin	706
652	A universal vision-language pretrained model for	Wang, Wenbin Ge, et al. 2024a. Qwen2-vl: Enhanc-	707
653	chart comprehension and reasoning. <i>arXiv preprint</i>	ing vision-language model’s perception of the world	708
654	<i>arXiv:2305.14761</i> .	at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	709
655	Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty,	Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu,	710
656	and Enamul Hoque. 2022. Chartqa: A benchmark	Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie	711
657	for question answering about charts with visual and	Zhou, Yu Qiao, et al. 2024b. Visionllm: Large	712
658	logical reasoning. <i>arXiv preprint arXiv:2203.10244</i> .	language model is also an open-ended decoder for	713
		vision-centric tasks. <i>Advances in Neural Information</i>	714
659	Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan	<i>Processing Systems</i> , 36.	715
660	Parvez, Enamul Hoque, and Shafiq Joty. 2024.	Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen,	716
661	Chartinstruct: Instruction tuning for chart com-	Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu,	717
662	prehension and reasoning. <i>arXiv preprint</i>	Haotian Liu, Sadhika Malladi, et al. 2024c. Charxiv:	718
663	<i>arXiv:2403.09028</i> .	Charting gaps in realistic chart understanding in mul-	719
		timodal llms. <i>arXiv preprint arXiv:2406.18521</i> .	720
664	Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	721
665	Kaipeng Zhang, Yu Qiao, and Ping Luo. 2024. Char-	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	722
666	tassissant: A universal chart multimodal language	et al. 2022. Chain-of-thought prompting elicits rea-	723
667	model via chart-to-table pre-training and multitask	soning in large language models. <i>Advances in neural</i>	724
668	instruction tuning. <i>arXiv preprint arXiv:2401.02384</i> .	<i>information processing systems</i> , 35:24824–24837.	725
669	Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and	Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang,	726
670	Pratyush Kumar. 2020. Plotqa: Reasoning over sci-	Zhixuan Liang, Zeyu Lu, Ying Shan, and Ping Luo.	727
671	entific plots. In <i>Proceedings of the IEEE/CVF Win-</i>	2024. Plot2code: A comprehensive benchmark for	728
672	<i>ter Conference on Applications of Computer Vision</i> ,	evaluating multi-modal large language models in	729
673	pages 1527–1536.	code generation from scientific plots. <i>arXiv preprint</i>	730
		<i>arXiv:2405.07990</i> .	731
674	OpenAI. 2023. Gpt-4v(ision) system card .	Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao	732
675	OpenAI. 2024. Gpt-4o . Accessed: 2024-05-13.	Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min	733
676	Qwen Team. 2024. Qwen2.5: A party of foundation	Dou, Botian Shi, Junchi Yan, et al. 2024. Chartx	734
677	models .	& chartvlm: A versatile benchmark and founda-	735
678	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	tion model for complicated chart reasoning. <i>arXiv</i>	736
679	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	<i>preprint arXiv:2402.12185</i> .	737
680	try, Amanda Askell, Pamela Mishkin, Jack Clark,	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,	738
681	et al. 2021. Learning transferable visual models from	Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin	739
682	natural language supervision. In <i>International confer-</i>	Jiang. 2023. Wizardlm: Empowering large lan-	740
683	<i>ence on machine learning</i> , pages 8748–8763. PMLR.	guage models to follow complex instructions. <i>arXiv</i>	741
684	Chufan Shi, Cheng Yang, Yaxin Liu, Bo Shui, Junjie	<i>preprint arXiv:2304.12244</i> .	742
685	Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li,	Zhishen Yang, Raj Dabre, Hideki Tanaka, and Naoaki	743
686	Yuxiang Zhang, et al. 2024. Chartmimic: Evaluating	Okazaki. 2024a. Scicap+: A knowledge augmented	744
687	Imm’s cross-modal reasoning capability via chart-to-	dataset to study the challenges of scientific figure	745
688	code generation. <i>arXiv preprint arXiv:2406.09961</i> .	captioning. <i>Journal of Natural Language Processing</i> ,	746
689	Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruiibo	31(3):1140–1165.	747
690	Liu, and Diyi Yang. 2024. Design2code: How far are	Zhiyu Yang, Zihan Zhou, Shuo Wang, Xin Cong,	748
691	we from automating front-end engineering? <i>arXiv</i>	Xu Han, Yukun Yan, Zhenghao Liu, Zhixing Tan,	749
692	<i>preprint arXiv:2403.03163</i> .	Pengyuan Liu, Dong Yu, et al. 2024b. Matplota-	750
693	Amanpreet Singh, Vivek Natarajan, Meet Shah,	gent: Method and evaluation for llm-based agen-	751
694	Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,	tic scientific data visualization. <i>arXiv preprint</i>	752
695	and Marcus Rohrbach. 2019. Towards vqa models	<i>arXiv:2402.11453</i> .	753
696	that can read. In <i>Proceedings of the IEEE/CVF con-</i>		
697	<i>ference on computer vision and pattern recognition</i> ,		
698	pages 8317–8326.		

- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yi Zhou, Junyan Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qiang Qi, Ji Zhang, and Feiyan Huang. 2023. [mplug-owl: Modularization empowers large language models with multimodality](#). *ArXiv*, abs/2304.14178.
- Sukmin Yun, Haokun Lin, Rusiru Thushara, Mohammad Qazim Bhat, Yongxin Wang, Zutao Jiang, Mingkai Deng, Jinhong Wang, Tianhua Tao, Junbo Li, et al. 2024. Web2code: A large-scale webpage-to-code dataset and evaluation framework for multi-modal llms. *arXiv preprint arXiv:2406.20098*.
- Fanhu Zeng, Fei Zhu, Haiyang Guo, Xu-Yao Zhang, and Cheng-Lin Liu. 2024. Modalprompt: Dual-modality guided prompt for continual learning of large multi-modal models. *arXiv preprint arXiv:2410.05849*.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023a. [Sigmoid loss for language image pre-training](#). *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11941–11952.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023b. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986.
- Fengji Zhang, Linqun Wu, Huiyu Bai, Guancheng Lin, Xiao Li, Xiao Yu, Yue Wang, Bei Chen, and Jacky Keung. 2024a. Humaneval-v: Evaluating visual understanding and reasoning abilities of large multi-modal models through coding tasks. *arXiv preprint arXiv:2410.12381*.
- Liang Zhang, Anwen Hu, Haiyang Xu, Ming Yan, Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang. 2024b. Tinychart: Efficient chart understanding with visual token merging and program-of-thoughts learning. *arXiv preprint arXiv:2404.16635*.
- Wenqi Zhang, Zhenglin Cheng, Yuanyu He, Mengna Wang, Yongliang Shen, Zeqi Tan, Guiyang Hou, Mingqian He, Yanna Ma, Weiming Lu, et al. 2024c. Multimodal self-instruct: Synthetic abstract image and visual reasoning instruction using language model. *arXiv preprint arXiv:2407.07053*.
- Wenqing Zheng, S P Sharan, Ajay Jaiswal, Kevin Wang, Yihan Xi, Dejie Xu, and Zhangyang Wang. 2023. [Outline, then details: Syntactically guided coarse-to-fine code generation](#). In *International Conference on Machine Learning*.

A Appendix

A.1 Implementation Details

ChartCoder is initialized with SigLIP-384 (Radford et al., 2021) as the vision encoder and DeepSeek Coder 6.7B (Guo et al., 2024) as the large language model. The whole training process is divided into alignment and instruction tuning. During the alignment stage, we only train the vision-language connector with the chart-to-text alignment data. The learning rate is set to 1e-3. In the instruction tuning stage, we train the entire model for 1 epoch with a batchsize of 128. The learning rate of SigLIP and other modules are 5e-6 and 1e-5 respectively, with a warmup at the beginning of 3%, then decays to 0 at the end of training. The alignment and instruction tuning processes cost 12 and 5 hours on 32 Tesla A100 GPUs with 80 GB VRAMs.

A.2 Benchmark Details

ChartMimic (Shi et al., 2024) focuses on evaluating the ability of MLLMs to redraw charts from ArXiv papers, emphasizing the preservation of the original style and appearance. It consists of two subsets: testmini and test. Following the settings in the original paper, we adopt the Direct Mimic task on the testmini subset as the default evaluation standard, reporting execution success rates alongside low-level and high-level scores.

Plot2Code (Wu et al., 2024) aims to evaluate models’ abilities to generate code corresponding to charts from the available Matplotlib galleries, with a focus on textual similarity. We evaluate models on its Direct Asking task using three metrics: Pass Rate, Text-Match, and Rating.

ChartX (Xia et al., 2024) contains various tasks with synthesis chart images, including Question Answering, Summarization, Description and Redrawing. We choose the Redrawing task and report the GPT score as the metrics in ChartX.

A.3 More Ablation Studies

We also perform more ablation studies on the language backbone and further choose Qwen2.5-7B and Qwen2.5 Coder-7B (Qwen Team, 2024) for comparison. The results also show that using Code LLM as the language backbone is better than using general LLM. However, we find that using the Qwen2.5 Coder as the backbone does not perform as well as using the DeepSeek Coder. This observation seems counterintuitive, as the official evaluation suggests that the performance of the Qwen2.5

Methods	ChartMimic		
	Exec.Rate	Low-Level	High-Level
ChartCoder	91.4	77.4	74.0
<i>Replace Language Backbone</i>			
Qwen2.5	88.1	73.4	67.9
Δ	-3.3	-4.0	-6.1
Qwen2.5 Coder	90.3	76.8	69.7
Δ	-1.1	-0.6	-4.3

Table 7: The ablation studies on model architecture and data. The results show that the effectiveness of our proposed model architecture and dataset.

Type	pie	line	bar	bar_num
Percent	8.0%	9.7%	8.3%	3.3%
Type	3d	area	box	bubble
Percent	5.6%	3.9%	4.4%	2.8%
Type	candlestick	funnel	heatmap	multi-axes
Percent	2.8%	2.7%	3.9%	3.8%
Type	rader	ring	pie	rose
Percent	3.8%	2.7%	2.8%	3.9%
Type	treemap	violin	scatter	quiver
Percent	3.9%	3.9%	3.8%	5.2%
Type	inset	histogram	graph	error bar
Percent	1.2%	1.2%	1.2%	1.6%
Type	error point	density	Combination	Total
Percent	1.6%	1.2%	2.8%	100%

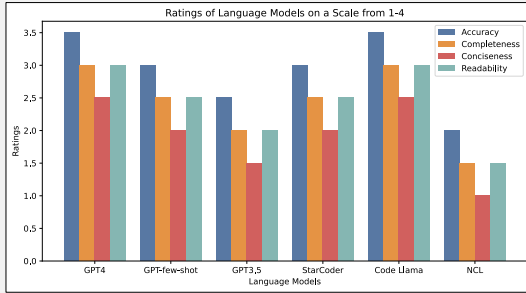
Table 8: Type distributions of the Chart2Code-160k instruction-tuning dataset.

Coder is better than the DeepSeek Coder. We analyze experimental results and find that the code generated by Qwen2.5 is more standardized. For instance, the DeepSeek Coder backbone tends to use `ax[0]`, `ax[1]`, while the Qwen2.5 Coder backbone prefers a more standardized approach, such as using `for i in range(2): ax[i]`. However, in some complex scenarios, using a `for` loop may lead to errors, such as `ax[0]` and `ax[1]` do not have same number of bars.

A.4 Chart2Code-160k Analysis

We count the proportion of different charts in the Chart2Code-160k dataset in Table 8.

Instruction



You are an expert Python developer who specializes in writing matplotlib code based on a given picture. I found a very nice picture in a STEM paper, but there is no corresponding source code available. I need your help to generate the Python code that can reproduce the picture based on the picture I provide. Note that it is necessary to use `figsize=(9.0, 5.0)` to set the image size to match the original size. Now, please give me the matplotlib code that reproduces the picture below.

GPT-4o Output

```
import matplotlib.pyplot as plt
import numpy as np

# Data
models = ['GPT4', 'GPT-few-shot', 'GPT3.5', 'StarCoder', 'Code Llama', 'NCL']
accuracy = [3.5, 3.0, 2.5, 3.0, 3.5, 2.5]
completeness = [3.0, 2.5, 2.0, 2.5, 3.0, 2.0]
conciseness = [2.5, 2.0, 1.5, 2.0, 2.5, 1.0]
readability = [3.0, 3.0, 2.5, 3.0, 3.0, 1.5]
x = np.arange(len(models)) # the 'label' locations
width = 0.2 # the width of the bars

# Plot
fig, ax = plt.subplots(figsize=(9.0, 5.0))
rects1 = ax.bar(x - 1.5*width, accuracy, width, label='Accuracy', color='#4682b4')
rects2 = ax.bar(x - 0.5*width, completeness, width, label='Completeness', color='#ff7f0e')
rects3 = ax.bar(x + 0.5*width, conciseness, width, label='Conciseness', color='#c0504d')
rects4 = ax.bar(x + 1.5*width, readability, width, label='Readability', color='#9467bd')
# Add some text for labels, title and custom x-axis tick labels, etc.
ax.set_xlabel('Language Models')
ax.set_ylabel('Ratings')
ax.set_title('Ratings of Language Models on a Scale from 1-4')
ax.set_xticks(x)
ax.set_xticklabels(models)
ax.legend()
fig.tight_layout()
```

ChartCoder Output

```
import matplotlib.pyplot as plt
import numpy as np

# Data
models = ['GPT4', 'GPT-few-shot', 'GPT3.5', 'StarCoder', 'Code Llama', 'NCL']
accuracy = [3.5, 3.0, 2.5, 3.0, 3.5, 2.0]
completeness = [3.0, 2.5, 2.0, 2.5, 3.0, 1.5]
conciseness = [2.5, 2.0, 1.5, 2.0, 2.5, 1.0]
readability = [3.0, 2.5, 2.0, 2.5, 3.0, 1.5]
barWidth = 0.2

r1 = np.arange(len(models))
r2 = [x + barWidth for x in r1]
r3 = [x + barWidth for x in r2]
r4 = [x + barWidth for x in r3]
# Create the figure and the axes
fig, ax = plt.subplots(figsize=(10.0, 5.0))
# Plot the bars
ax.bar(r1, accuracy, color='#4682b4', width=barWidth, label='Accuracy')
ax.bar(r2, completeness, color='#ff7f0e', width=barWidth, label='Completeness')
ax.bar(r3, conciseness, color='#c0504d', width=barWidth, label='Conciseness')
ax.bar(r4, readability, color='#9467bd', width=barWidth, label='Readability')
# Add xticks on the middle of the group bars
ax.set_xlabel('Language Models', fontweight='medium')
ax.set_ylabel('Ratings', fontweight='medium')
ax.set_title('Ratings of Language Models on a Scale from 1-4', fontweight='medium')
ax.set_xticks([r + barWidth for r in range(len(models))])
ax.set_xticklabels(models)
# Create legend & show graphic
ax.legend(loc='upper right', fontsize='large')
plt.tight_layout()
```

InternVL2-Llama3-76B Output

```
import matplotlib.pyplot as plt
import numpy as np

# Data
models = ['GPT4', 'GPT-few-shot', 'GPT3.5', 'StarCoder', 'Code Llama', 'NCL']
accuracy = [3.5, 3.0, 2.5, 3.0, 3.5, 2.0]
completeness = [3.0, 2.5, 2.0, 2.5, 3.0, 1.5]
conciseness = [2.5, 2.0, 1.5, 2.0, 2.5, 1.0]
readability = [3.0, 2.5, 2.0, 2.5, 3.0, 1.5]

# Plotting
fig, ax = plt.subplots(figsize=(9.0, 5.0))
x = np.arange(len(models))
width = 0.2

ax.bar(x - 1.5*width, accuracy, width, label='Accuracy', color='blue')
ax.bar(x - 0.5*width, completeness, width, label='Completeness', color='orange')
ax.bar(x + 0.5*width, conciseness, width, label='Conciseness', color='red')
ax.bar(x + 1.5*width, readability, width, label='Readability', color='green')

ax.set_ylabel('Ratings')
ax.set_title('Ratings of Language Models on a Scale from 1-4')
ax.set_xticks(x)
ax.set_xticklabels(models)
ax.legend()
```

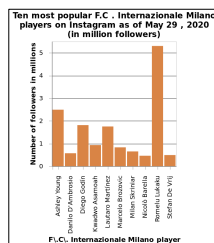
InternVL2-B Output

```
import matplotlib.pyplot as plt

# Data for the bar chart
ratings = {
    'GPT4': {'Accuracy': 3.5, 'Completeness': 3.0, 'Consensus': 2.5, 'Readability': 3.0},
    'GPT-few-shot': {'Accuracy': 3.0, 'Completeness': 2.5, 'Consensus': 2.0, 'Readability': 2.5},
    'GPT3.5': {'Accuracy': 2.5, 'Completeness': 2.0, 'Consensus': 1.5, 'Readability': 2.0},
    'StarCoder': {'Accuracy': 3.0, 'Completeness': 2.5, 'Consensus': 2.0, 'Readability': 2.5},
    'Code Llama': {'Accuracy': 3.5, 'Completeness': 3.0, 'Consensus': 2.5, 'Readability': 3.0},
    'NCL': {'Accuracy': 2.0, 'Completeness': 1.5, 'Consensus': 1.0, 'Readability': 1.5}
}

# Set up the figure and axes
fig, ax = plt.subplots(figsize=(9.0, 5.0))
# Plot the bars
for model, ratings in ratings.items():
    accuracy = ratings['Accuracy']
    completeness = ratings['Completeness']
    consensus = ratings['Consensus']
    readability = ratings['Readability']
    ax.bar(model, [accuracy, completeness, consensus, readability])
# Set the x and y tick labels
ax.set_xlabel('Model')
ax.set_ylabel('Rating')
ax.set_xticks(range(1, 4))
# Set the x tick labels
ax.set_xticklabels(['Accuracy', 'Completeness', 'Consensus', 'Readability'])
# Set the y tick labels
ax.set_yticklabels(['1', '2', '3', '4'])
# Set the title
ax.set_title('Ratings of Language Models on a Scale from 1-4')
```

Figure 5: A example of comparing the code corresponding to the bar chart generated by different models.



ChartCoder Output

```
< Code >
'''python/nimport matplotlib.pyplot as
plt/nimport os/n# Data/nplayers = ['Ashley
Young', 'Danilo Ambrosio', 'Diego Godinho',
'Kwadwo Asamoah', 'Laurento Brozovic',
'Marcelo Brozovic', 'Milan Skriniar',
'Brozovic', 'Nicolo Barella', 'Romelu Lukaku',
'Stefan De Vrij']/nfollowers = [2.5, 0.7, 1.8,
0.9, 1.9, 0.8, 0.6, 0.4, 5.4, 0.5] .....'''
```

<Original Instruction > :

```
{
  "title": "Ten most popular F.C. Internazionale Milano players on Instagram as of May 29, 2020 (in million followers)",
  "players": [
    { "name": "Lukaku", "followers": 5.3 },
    { "name": "Young", "followers": 2.49 },
    { "name": "Godin", "followers": 1.81 },
    { "name": "Asamoah", "followers": 0.93 },
    { "name": "Brozovic", "followers": 1.75 },
    { "name": "Skriniar", "followers": 0.65 },
    { "name": "D'Ambrosio", "followers": 0.83 },
    { "name": "De Vrij", "followers": 0.49 },
    { "name": "Barella", "followers": 0.46 }
  ]
}
```

< Label > :

False

MiniCPM-Llama3-V2.5 Output

< Original Response > :

True

< Instruction with Code > :

You can use the above code as the reference to solve the problem, but focus on extracting information from the image to answer the following problem.

MiniCPM-Llama3-V2.5 Output

< Response with Code > :

The provided JSON data is not related to the original image or the bar chart displayed in the code. It represents a different set of data concerning the per capita real GDP of Massachusetts from 2000 to 2019, measured in chained 2012 U.S. dollars. The data in the JSON file does not correspond to the content of the bar chart, which displays Instagram followers for soccer players from Internazionale Milano as of May 29, 2020. Therefore, the statement that this is a correct JSON file transferred from the image is false.

Figure 6: A case study comparing the outputs of utilizing image and image+code as the inputs