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 lossless representations that can effectively contain all critical details. Although existing opensource 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 restora-013 tion 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 017 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-021 scale and diverse dataset for chart-to-code generation, and propose the Snippet-of-Thought (SoT) method, which transforms direct chartto-code generation data into step-by-step generation. Experiments demonstrate that Chart-Coder, with only 7B parameters, surpasses existing open-source MLLMs on chart-to-code benchmarks, achieving superior chart restoration and code excitability.

1 Introduction

037

041

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). 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 complex information such as data



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

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 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.

061

062

065

072

090

091

100

101

103

104

105

106 107

108

109

110

111

To overcome the above challenges in chart-tocode 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-totext, 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) 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-tocode MLLM, which leverages Code LLMs as
 language backbones. With only 7B parameters, ChartCoder outperforms existing opensource MLLMs on chart-to-code benchmarks.

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

- 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 and high-level tasks. 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) specifically designed for chartrelated tasks. Trained on extensive chart-specific datasets, these chart-domain MLLMs (Xia et al., 2024; Zhang 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-159 to-code generation has attracted significant interest 160 as the visual elements of charts are more complex. 161 This task challenges MLLMs to generate code that 162 accurately reproduces a given chart or visual representation. Recent works like ChartMimic (Shi et al., 164 2024) evaluate the reasoning ability of MLLMs in 165 this context. Similarly, Plot2Code (Wu et al., 2024) 166 and ChartX (Xia et al., 2024) also evaluate MLLMs code generation ability, especially for text and data 168 reproducibility. To the best of our knowledge, no dedicated research has focused on solving the chart-170 to-code generation problem. Our work is the first 171 to attempt to address this challenge. 172

3 Chart2Code-160k Dataset

173

174

175

176

178

179

181

183

184

187

188

190

192

193

196

197

198

204

205

208

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 questionanswer 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-tocode 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 chartcode 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 oneto-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**.

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

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.

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

chart types and manually write 79 template codes for each as in-context demonstrations. These template codes contain almost all common chart formats. We further provide available functions such as plt.text() and parameters such as hatch='/' to encourage the generation of more diverse functions and parameters, 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, instruction, 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 has demonstrated effectiveness in text-to-code generation tasks (Zheng et al., 2023; Luo et al., 2024).

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 lay-



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.

out, 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.

253

259

260

261

263

270

271

272

274

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 twostep 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. 275

276

277

278

279

280

281

283

284

285

286

289

290

292

293

294

296

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. To ensure data quality, we randomly sample around 1k instances and evaluate the quality of the chart

Dataset	Source	Chart Quality	
		Mean μ	$\mathrm{SD}\sigma$
Chart2Code-160k sample	Generated	77.32	16.34
ChartMimic testmini	Real-world	78.96	13.66

Table 2: Quantitative evaluation of the chart quality, comparing with real-world charts. SD is the abbreviation version for standard deviation.

310

312

314

316

319

321

322

326

327

328

332

334

338

297

images manually during the dataset construction period. Given the strong generation capabilities of LLM, we reckon the generated charts are suitable for training purposes. Furthermore, to quantitatively evaluate the chart quality, we also sample 8k data pairs (5% of the total) from Chart2Code-160k and utilize gpt-4o-2024-08-06 to evaluate them on four criteria: *Aesthetics, Readability, Reproducibility*, and *Data Presentation Simplicity*. The results in Table 2 show that the overall scores are broadly the same as real-world charts in Chart-Mimic. The detailed prompt is in the Figure 8. Our proposed Chart2Code-160k fills the gap between chart and code, 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, 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 pretrained 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; Guo et al., 2025), 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 224² and 384² 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. 339

340

341

343

344

345

346

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

385

387

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., 2024). 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

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
			Propr	ietary				
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-40	-	93.2	79.0	83.5	88.6	56.3	5.71	-
		6	Open-Source G	eneral-Domain	n			
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
	Open-Source Chart-Domain							
ChartLlama	13B	57.5	24.8	28.1	58.4	40.3	2.32	0.94
ChartAssisstant	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 3: 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-40	18.96	9.59	17.16	15.68	14.66	8.84
InternVL2-Llama3-76B Qwen2-VL-72B InternVL2-8B	13.06 10.45 7.20	8.44 7.83 6.82	12.59 9.92 8.81	10.51 8.14 5.74	8.74 7.10 5.42	7.87 7.47 6.64
TinyChart ChartVLM-L ChartCoder (Ours)	4.16 0.97 16.83	5.06 3.53 9.13	5.22 2.48 14.77	2.74 0.81 12.41	3.21 0.90 12.68	5.58 5.25 8.29

Table 4: Detailed results of high-level scores on ChartMimic Direct Mimic task. All the subscores of ChartCoder are close to GPT-40.

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-40 (OpenAI, 2024). (3) Chart-domain MLLMs including ChartLlama (Han et al., 2023), ChartAssisstant (Meng et al., 2024), Tinychart (Zhang et al., 2024b) and ChartVLM (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,

389

395

400

401

402

403

404

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. 405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

5.2 Main Results

As indicated in Table 3 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

Model	Text	Layout	Туре	Color
Full score	100	100	100	100
GPT-40	81.5	89.8	77.3	67.2
InternVL2-Llama3-76B Qwen2-VL-72B InternVL2-8B	54.1 43.2 31.5	74.5 80.5 51.1	49.2 54.6 28.6	41.5 39.4 26.2
TinyChart ChartVLM-L ChartCoder (Ours)	9.8 7.7 67.2	48.2 33.7 95.0	32.9 17.6 78.5	14.2 5.2 69.0

Table 5: Detailed results of low-level scores on Chart-Mimic Direct Mimic task. Three out of four subscores of ChartCoder are even higher than GPT-40.

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 largescale MLLMs such as InternVL2-Llama3-76B and chart-domain MLLMs such as TinyChart.

421

422

423 424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454 455

456

457

458

459

We further illustrate the detailed high-level and low-level scores for the ChartMimic benchmark. The high-level score utilizes GPT-40 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 4 denotes the high-level results. Chart-Coder is the model most comparable to GPT-40, as the evaluations were conducted by GPT-40 itself, suggesting the actual performance gap may not be as pronounced as it appears. Notably, Chart-Coder shows the largest gap with GPT-40 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 5 shows the low-level results. ChartCoder even slightly outperforms GPT-40 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-40, 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

		ChartMimic			
Methods	Exec.Rate	Low-Level	High-Level		
ChartCoder	91.4	77.4	74.0		
Code	$e LLM \rightarrow G$	eneral LLM			
DeepSeek LLM	80.6	61.4	63.4		
\bigtriangleup	-10.8	-16.0	-10.6		
Dif	Different Visual Encoders				
CLIP-336	91.6	77.3	70.3		
\bigtriangleup	+0.2	-0.1	-3.7		
Withou	ıt Step-by-st	ep Generatio	n		
w/o SoT	89.2	70.1	65.4		
\bigtriangleup	-2.2	-7.3	-8.6		
Open-source MLLM Finetund on Chart2Code-160k					
Qwen2-VL-7B	67.0	32.9	35.0		
Finetuned Model	83.6	73.4	68.2		
\bigtriangleup	+16.7	+40.5	+33.2		

Table 6: The ablation studies on model architecture and data. The results show that the effectiveness of our proposed code LLM backbone and dataset.

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 existing MLLMs. The results are shown in Figure 3, the outputs of ChartCoder are much more similar to the ground truth chart than open-source models. 460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

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 6. (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 that performance on chart understanding tasks is resolution-dependent, with lower resolu-



Figure 3: Generated charts of different model outputs after code execution. Our proposed ChartCoder performs significantly better than InternVL2-8B of a similar model scale.

Model	Image	Image+Code	\triangle
MiniCPM-Llama3-V2.5	0.76	0.81	6.5%
InternVL2-8B	0.79	0.82	3.8%

Table 7: Comparison of the impact of using code as auxiliary contexts on the MMC True/False task.

491

492

493

494

495

496

497

499

502

504

510

511

512

513

514

tions negatively impacting model performance (Liu et al., 2024). To verify whether resolution affects chart-to-code tasks, we replace SigLIP-384 with CLIP-336 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 effectiveness of Chart2Code-160k.

515

516

517

518

520

521

522

523

524

525

526

528

529

530

531

532

534

535

536

537

538

540

541

542

5.4 Analysis

We further evaluate the role of code in the chart understanding task. We use two MLLMs 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 7 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 5. Also, we utilize LLM to evaluate the readability of ChartCoder output code, and details are in the Appendix A.4.

6 Conclusion

This work aims to tackle the challenge of chartto-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-tocode dataset, Chart2Code-160k. Finally, to emphasize the key information, we propose the Snippetof-Thought (SoT) method to generate step-by-step data. Experiments show that ChartCoder outperforms existing open-source MLLMs.

543 Limitation

Our study is comprehensive but has certain limitations that we aim to address in future research. Due 545 to constraints in computational resources, we only 546 trained ChartCoder with 7B parameters, which has 547 demonstrated sufficiently good results for now. A 549 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 552 HTML, to develop a comprehensive multi-modal 553 code large language model. 554

555 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

556

560

564

566

567

570

571

572

573

574

575

577

578

580

581

582

583

584

585

587

591

594

595

- 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, Qiushi 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, Yu 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.
 - 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*. 596

597

599

600

601

602

603

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

- Zonghao Guo, Ruyi Xu, Yuan Yao, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, and Gao Huang. 2025. Llava-uhd: an lmm perceiving any aspect ratio and high-resolution images. In *European Conference on Computer Vision*, pages 390–406. Springer.
- 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. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. 2024. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*.

- 651

- 672
- 673
- 674

676

- 677 678 679

686

697

702

- Xianzhen Luo, Qingfu Zhu, Zhiming Zhang, Libo Qin, Xuanyu Zhang, Qing Yang, Dongliang Xu, and Wanxiang Che. 2024. Python is not always the best choice: Embracing multilingual program of thoughts. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 7185–7212.
- Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq Joty. 2023. Unichart: A universal vision-language pretrained model for chart comprehension and reasoning. arXiv preprint arXiv:2305.14761.
- Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. 2024. Chartassisstant: A universal chart multimodal language model via chart-to-table pre-training and multitask instruction tuning. arXiv preprint arXiv:2401.02384.
- Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 1527-1536.
- OpenAI. 2023. Gpt-4v(ision) system card.
- OpenAI. 2024. Gpt-4o. Accessed: 2024-05-13.
 - Qwen Team. 2024. Qwen2.5: A party of foundation models.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748–8763. PMLR.
 - Chufan Shi, Cheng Yang, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, et al. 2024. Chartmimic: Evaluating lmm's cross-modal reasoning capability via chart-tocode generation. arXiv preprint arXiv:2406.09961.
 - Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2code: How far are we from automating front-end engineering? arXiv preprint arXiv:2403.03163.
 - Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vga models that can read. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 8317-8326.
 - Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024a. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191.

703

704

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

750

751

752

753

754

755

- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. 2024b. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. Advances in Neural Information Processing Systems, 36.
- Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. 2024c. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. arXiv preprint arXiv:2406.18521.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.
- Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and Ping Luo. 2024. Plot2code: A comprehensive benchmark for evaluating multi-modal large language models in code generation from scientific plots. arXiv preprint arXiv:2405.07990.
- Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min Dou, Botian Shi, Junchi Yan, et al. 2024. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. arXiv preprint arXiv:2402.12185.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. arXiv preprint arXiv:2304.12244.
- Zhishen Yang, Raj Dabre, Hideki Tanaka, and Naoaki Okazaki. 2024. Scicap+: A knowledge augmented dataset to study the challenges of scientific figure captioning. Journal of Natural Language Processing, 31(3):1140-1165.
- 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.
- 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 webpageto-code dataset and evaluation framework for multimodal llms. arXiv preprint arXiv:2406.20098.

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.

757

758

760

761

762

763

767

770

771

772 773

774

775

776

777

778 779

781

784

785

786

- 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, Linquan 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 multimodal 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, Dejia Xu, and Zhangyang Wang. 2023.
 Outline, then details: Syntactically guided coarse-tofine code generation. In *International Conference on Machine Learning*.

A Appendix

789

790

791

793

794

798

802

811

813

814

815

817

818

820

821

823

825

827

833

834

837

A.1 Implementation Details

In the data generation stage, we utilize gpt-4o-2024-08-06 as the LLM for both direct and step-by-step generation processes.

In the training stage, 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

Methods		ChartMimic	;	
wiethous	Exec.Rate	Low-Level	High-Level	
ChartCoder	91.4	77.4	74.0	
Rep	Replace Language Backbone			
Qwen2.5	88.1	73.4	67.9	
\bigtriangleup	-3.3	-4.0	-6.1	
Qwen2.5 Coder	90.3	76.8	69.7	
\bigtriangleup	-1.1	-0.6	-4.3	

Table 8: The ablation studies on model architecture and data. The results show that the effectiveness of our proposed model architecture and dataset.

Dataset	Source	Chart Quality	
		Mean μ	$\mathrm{SD}\sigma$
Qwen2-VL-7B Output	Model generated	82.48	6.81
ChartCoder Output	Model generated	85.22	6.78
ChartMimic Source	Human written	87.66	4.30

Table 9: Performance Comparison of model outputs and human-written sources. SD is the abbreviation version for standard deviation.

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 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 Output Code Analysis

To evaluate the output code readability, we conduct an ablation experiment, utilizing gpt-4o-2024-08-06 to evaluate the output code readability. We evaluate four aspects of the generated code, including *Naming Conventions, Code Structure, Comments,* and *Logical Clarity,* with a total score of 100. We choose the generated code from the ChartMimic task (ChartCoder output) and the ground truth code (human-annotated) in the ChartMimic dataset. The results are as shown in Table 9. We also evaluate the error types on ChartMimic direct generation tasks with code and general LLMs as the language backbone. The results are shown in Figure 4.

860

861

862

863

864

865



Figure 4: Comparison of error types on ChartMimic direct generation tasks with code and general LLMs as language backbone, respectively.

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

Table 10: Type distributions of the Chart2Code-160k instruction-tuning dataset.

868 869 870 871 872

867

A.5 Chart2Code-160k Analysis

We count the proportion of different charts in the Chart2Code-160k dataset in Table 10. Also, we utilize gpt-4o-2024-08-06 to evaluate the quality of the charts in the Chart2Code-160k and compare them with the real-world chart. The prompts are shown in Figure 8.





0.9, 1.9, 0.8, 0.6, 0.4, 5.4, 0.5]`



Figure 6: A example of comparing the code corresponding to the bar chart generated by different models.

Prompt for Code Readability Evaluation

Please score the code's readability based on the following four aspects. Each aspect is worth 25 points, for a total of 100 points. Naming Conventions (25 points) Score: [X]/25 Explanation: [Provide a brief explanation of how well the variable, function, and class names convey their purpose and whether the naming style is consistent across the codebase.] Code Structure (25 points) Score: [X]/25 Explanation: [Explain whether functions are concise, whether the code uses indentation and blank lines appropriately, and whether the code is modularized effectively.] Comments (25 points) Score: [X]/25 Explanation: [Discuss the clarity and appropriateness of the comments, and whether functions/methods have proper documentation comments explaining inputs, outputs, and functionality.] Logical Clarity (25 points) Score: [X]/25 Explanation: [Evaluate the intuitiveness of the code, whether it's easy to understand, and whether the control flow is simple and avoids unnecessary complexity.] Total Score: [X]/100 Summary: [Provide a brief overall assessment of the code's readability, pointing out strengths and potential areas for improvement.]

Figure 7: Prompt for dataset quality evaluation.

Prompt for Code Readability Evaluation

You are a professional chart analyser. Please evaluate the image based on the following four criteria: aesthetics, readability, reproducibility, and data presentation simplicity. Provide a score for each criterion and include an overall score along with a brief evaluation. Scoring Criteria and Requirements: Aesthetics (25 points) Requirements:
The chart design should be simple and clear, avoiding complex decorations, and should effectively communicate information
Colors should be harmonious and have high contrast, making it easy to differentiate between different data groups.
Legends and labels should be clear, with appropriately sized fonts, avoiding visual clutter. Scoring: [X]/25
Readability (30 points) Requirements:
The chart should have clear titles, axis labels, and legends, enabling quick communication of the main message.
Data curves or point annotations should avoid being overly dense or overlapping, maintaining good readability.
The overall layout should follow a logical structure without any confusing elements. Scoring: [X]/30
Reproducibility (30 points) Requirements:
The chart design should be easy to replicate using common tools.
point deductions.
recreate the chart from scratch.
or complex elements.
Data Presentation Simplicity (15 points)
Data presentation should be concise, avoiding redundant information.
Data should be displayed in an intuitive way, without excessive curves, points, or annotations. High-scoring charts should focus on the main data being presented, avoiding decorative or unrelated information
Scoring: [X]/15

Figure 8: Prompt for dataset quality evaluation.