### ChartLlama: A Multimodal LLM for Chart Understanding and Generation



#### Anonymous ACL submission

Figure 1: **Capability demonstration of ChartLlama.** An instruction-tuning dataset is created based on our proposed data generation pipeline. We train ChartLlama on this dataset and achieve the abilities shown in the figure.

#### Abstract

Multi-modal large language models have demonstrated impressive performances on most 003 vision-language tasks. However, the model generally lacks the understanding capabilities for specific domain data, particularly when it comes to interpreting chart figures. This is mainly due to the lack of relevant multi-modal instruction tuning datasets. In this article, we create a high-quality instruction-tuning dataset leveraging GPT-4. We develop a multi-step data generation process in which different steps are responsible for generating tabular data, creating chart figures, and designing instruction tuning data separately. Our method's flexibility enables us to generate diverse, high-quality instruction-tuning data consistently and efficiently while maintaining a low resource expenditure. Additionally, it allows us to incorporate a wider variety of chart and task types not yet featured in existing datasets. Next, we introduce ChartLlama, a multi-modal large language model that we've trained using our 022 created dataset. ChartLlama outperforms all prior methods in ChartQA, Chart-to-text, and Chart-extraction evaluation benchmarks. Additionally, ChartLlama significantly improves upon the baseline in our specially compiled chart dataset, which includes new chart and task types. The results of ChartLlama confirm the value and huge potential of our proposed data generation method in enhancing chart comprehension. 024

025

027

031

032

033

034

036

041

043

044

### 1 Introduction

In the past year, the field of artificial intelligence has undergone remarkable advancements. A key highlight is the emergence of large language models (LLMs) like GPT-4 (OpenAI, 2023). These models (Ouyang et al., 2022; Zeng et al., 2022; Team, 2023; Baichuan, 2023; Touvron et al., 2023a,b) have demonstrated a remarkable capability to comprehend and generate intricate textual data, opening doors to myriads of applications in both academia and industry. Taking this progress a step further, the introduction of GPT-4V (Yang

880

096

et al., 2023) marked another milestone. It endows LLMs with the ability to interpret visual information, essentially providing them with a vision. As a result, they can now extract and analyze data from images, marking a significant evolution in the capacities of these models.

However, despite the achievements and potentials of models like GPT-4V, the details behind GPT-4V's architecture remain a mystery. This opacity has given rise to questions within the academic world about the best practices for designing multi-modal LLMs. Notably, pioneering research initiatives, like LLaVA (Liu et al., 2023c,b) and MiniGPT (Zhu et al., 2023; Chen et al., 2023), provide insightful directions in this regard. Their findings suggest that by incorporating visual encoders into existing LLMs and then fine-tuning them using multi-modal instruction-tuning datasets, LLMs can be effectively transformed into multi-modal LLMs. It's noteworthy that these multi-modal datasets are typically derived from established benchmarks, presenting a cost-effective method for accumulating data required for instruction tuning.

Datasets grounded on established benchmarks, such as COCO (Lin et al., 2014), have significantly enhanced the abilities of multi-modal LLMs to interpret everyday photographs adeptly. However, when confronted with specialized visual representations, such as charts, they reveal a noticeable limitation (Yang et al., 2023; Liu et al., 2023a). Charts are important visual instruments that translate complex data sets into digestible visual narratives, playing a crucial role in facilitating understanding, shaping insights, and efficiently conveying information. Their pervasive presence, from academic publications to corporate presentations, underscores the essentiality of enhancing the capability of multimodal LLMs in interpreting charts. Indeed, gathering data specifically to refine instructions for understanding charts presents several challenges. These typically stem from two areas: understanding and generation. An effective chart understanding model should be capable of extracting and summarizing data from various types of charts and making predictions based on this information.

However, most existing datasets (Masry et al., 2022; Kantharaj et al., 2022; Methani et al., 2020;
Masry et al., 2023) only provide support for simple question-answering or captioning, primarily due to the absence of detailed chart information and annotations that provide a high-level understanding of raw data. The high dependency on manually an-

notated charts gathered by web crawlers negatively affects the quality of these datasets. Thus, the previous annotating methods could only result in chart datasets with lower quality and less comprehensive annotations. Compared with chart understanding, generating chart figures is a more challenging task for the model because existing deep-learningbased generation methods (Ramesh et al., 2021; Rombach et al., 2021) struggle to accurately create images based on instructions. Using Python code to generate charts seems promising which needs the corresponding annotations to supervise models. Most charts obtained from the web are devoid of detailed annotations, making it challenging to annotate the generation code. The absence of code annotations makes it challenging to supervise models in code generation. These issues combined impede the model's ability to understand charts and learn generation jointly.

097

098

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

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

To address this, we introduce an adaptive and innovative data collection approach exclusively tailored to chart understanding and generation. At the heart of our methodology is the strategic employment of GPT-4's robust linguistic and coding capabilities, which facilitate the creation of rich multi-modal datasets. This innovative integration not only optimizes data accuracy but also ensures its wide-ranging diversity. Specifically, our method comprises three main phases:

1) **Chart Data Generation**. Our strategy for data collection stands out for its flexibility. Rather than limiting data collection to conventional data sources such as the web or existing datasets, we harness the power of GPT-4 to produce synthesized data. By providing specific characteristics such as topics, distributions, and trends, we guide GPT-4 to produce data that is both diverse and precise.

2) **Chart Figure Generation**. Subsequently, GPT-4's commendable coding skills are utilized to script chart plots using the open-sourced library, like Matplotlib, given the data and function documentation. The result is a collection of meticulously rendered charts that span various forms, each accurately representing its underlying data.

3) **Instruction data generation**. Beyond chart rendering, GPT-4 is further employed to interpret and narrate chart content, ensuring a holistic understanding. It is prompted to construct relevant question-answer pairs correlating with the charts. This results in a comprehensive instruction-tuning corpus, amalgamating the narrative texts, questionanswer pairs, and source or modified codes of the charts.

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

165

166

167

168

169

171

172

173

174

175

177

178

179

181

182

184

185

A standout feature of our methodology is its flexibility, which diminishes the potential for bias while simultaneously offering scalability. Building on this robust methodology, we've crafted a benchmark dataset, which is made available for public access. This dataset stands out, not only for its superior quality but also its unparalleled diversity. To showcase the superiority of our benchmark, we introduced a multi-modal Large Language Model (LLM) named ChartLlama trained with our established benchmarks. Our extensive experiments evaluated on multiple existing benchmark datasets show that our model outperforms previous methods with remarkable advantages and considerably less training data. Additionally, ChartLlama is equipped with several unique capabilities, including the ability to support a wider range of chart types, infer across multiple charts, undertake chart de-rendering tasks, and even edit chart figures.

Our main contributions are summarized as follows:

• We introduce a novel multi-modal data collection approach specifically designed for chart understanding and generation. The proposed data collection method boasts superior flexibility and scalability, enabling easy migration to different types of charts and various tasks.

• Through our innovative data collection approach, we create a benchmark dataset that stands out in terms of both quality and diversity. We make this dataset publicly available to catalyze further advancements in the field.

• We develop ChartLlama, a multi-modal LLM that not only surpasses existing models on various existing benchmarks but also possesses a diverse range of unique chart understanding and generation capabilities.

### 2 Related work

### 2.1 Large Language Model

The series of LLM models, such as GPT-189 3.5 (Ouyang et al., 2022) and GPT-4 (OpenAI, 190 2023), have demonstrated remarkable reasoning 191 and conversational capabilities, which have gar-192 193 nered widespread attention in the academic community. Following closely, a number of open-194 source LLM (Baichuan, 2023; Touvron et al., 195 2023a,b; Zeng et al., 2022; Bai et al., 2023a) models emerged, among which Llama (Touvron et al., 197

2023a) and Llama 2 (Touvron et al., 2023b) are notable representatives. With extensive pre-training on large-scale datasets and carefully designed instruction datasets, these models have also showcased similar understanding and conversational abilities. 198

199

200

201

202

203

204

205

206

207

208

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

### 2.2 Multi-modal Large Language Model

Concurrently, the academic community has witnessed a surge of development in multi-modal LLMs (Li et al., 2023a; Ye et al., 2023; Li et al., 2023c,b; Zhang et al., 2023b; Hu et al., 2023; Zhao et al., 2023; Bai et al., 2023b; Zhang et al., 2023a; Liu et al., 2023b,c; Chen et al., 2023; Zhu et al., 2023) built upon existing open-source models. With the exploration of training strategies and an increase in dataset scale, the performance of these new models has steadily improved, reaching comparable levels to GPT-4V in specific evaluation metrics. Notably, LLaVA-1.5 (Liu et al., 2023b), an iterative version of LLaVA, has gained popularity as a baseline due to its user-friendly training framework, superior performance, and data efficiency. Our work is also based on LLaVA-1.5.

### 2.3 Chart Understanding

In evaluations such as the report of GPT-4V (Yang et al., 2023) and HallusionBench (Liu et al., 2023a), it is evident that current multi-modal LLMs still struggle with complex chart-related problems. There are already some datasets (Methani et al., 2020; Kantharaj et al., 2022; Masry et al., 2022) available for evaluating models' chart understanding capabilities, mainly divided into two categories. One category measures through simple questionand-answer tasks, such as ChartQA (Masry et al., 2022), which has high-quality questions and answers annotated by humans. The other category converts charts into textual descriptions, with Chartto-text (Kantharaj et al., 2022) being a representative work in this field. The charts and annotations in these datasets are derived from the real world, ensuring higher quality, and encouraging models to delve deeper into the trends and meanings behind the charts. Previous works focusing on chart understanding tasks can be divided into two main kinds of approaches. One kind of approach is using a single model to understand the charts and answer questions in natural language, for example, (Masry et al., 2023; Liu et al., 2022b). The other kind of approach, such as (Liu et al., 2022a; Xia et al., 2023), is to first utilize the model to convert the

Stage 1: C	Chart Data (	Generation	Stage 2: Chart Figure Generation	Stage 3: Instruction Data Generation
Input Theme: Global av Trend: Rapid incre  Output	erage temperature, Da ease, Slow increase,	ily traffic, (5)	Input In context examples: Raw data: tabular data from Stage 1.  Output	Input The descriptions: The chart presents the variation in The raw data: Year, Amazon, Siberian'n 2010, 500 Instruction tuning data
Detailed descrip variation in forest co Rainforest and the Si fluctuations and sude Raw Data:	vtions about data: over over time, specifi- iberian Taigashow den drops in forest cor	the chart presents the cally for the Amazon cases the irregular verage for	Detailed descriptions anount charts: the piot has labels for xand y axis as "Year' and 'Area" (Square Kilometers), respectively, and the title of the plot is 'Comparison of Amazon Rainforest and Siberian Taiga Area'. A legend is placed at the upper right corner Generated figures:	<ul> <li>Q1: What is named to Fop songs in the 2000s decoding to the chart? A1: 50</li> <li>Q2: From the chart, can we infer any potential reasons for the more significant reduction in forest coverage? A2: It could</li> <li>Q3: Extract the raw data from the given chart. A3:</li> <li>Q4: Redraw the given chart figure. A4:</li> </ul>
Year	Amazon	Siberian	Residence of listicity business and provide the free of the second secon	Q5: Draw a funnel chart based on given raw data. A5:
2010	500	200	, interview of the second seco	V6: Remove the grids in the given chart figure. A6:
2011	600	300 		Abilities: Q&A, Chart Descriptions,

Figure 2: **Pipeline of our data generation method**. The innovative data generation process we proposed consists of three important steps relying on GPT-4. The dataset generated using this process exhibits significant advantages compared to previous datasets in terms of data diversity, quality, the number of chart types, and the variety of tasks. ChartLlama, which is trained on this dataset, has the ability to perform various tasks based on the design of the instruction-tuning data.

charts into structured data and then analyze and answer questions based on the structured data using existing large models. In our work, we primarily explore the former kind, aiming to leverage a single model to complete the entire process of chart understanding.

### 3 Method

248

249

251

254

256

260

261

262

263

264

265

267

269

271

272

273

275 276

277

278

In this section, we detail our unique approach to chart understanding and generation. Our method involves three interconnected steps: data collection, chart figure generation, and instruction data generation. We illustrate this process in Fig. 2. These steps are detailed in the following subsections.

### 3.1 Chart Data Generation

Our primary goal in chart data collection is to collect diverse and high-quality data. We employ two main strategies for this purpose: 1) Data Generation from Scratch Using GPT-4: To collect a diverse and high-quality dataset, we initially generate tabular data from scratch using GPT-4. We instruct GPT-4 to create data tables based on specific themes, distributions, and other characteristics like the size of the dataset in terms of rows and columns. This process ensured the creation of data with known and controlled characteristics, which can be essential for generating reliable instructionanswer pairs. Moreover, by managing these characteristics, we can intentionally minimize bias, leading to a more balanced dataset. 2) Synthesizing Data from Existing Chart Datasets. Our second strategy is to synthesize data by referencing existing chart datasets. These datasets already encompass a range of topics and characteristics, providing a solid base for data generation. By prompting GPT-4 with these datasets, we guide it to generate reasonable data that complements its existing knowledge base. This method added variety to our dataset and improved its overall quality. 281

283

284

285

287

288

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

Generating diverse data at scale using the LLM is not an easy task. When the prompt is designed improperly, the model tends to generate repetitive and meaningless data that deviates from the distribution of real-world data and thus lacks valuable insights that could be important for designing meaningful question-and-answer tasks. If we simply provide a set of data and require the model to imitate without any additional guidance, the model will probably just repeat the reference data. Therefore, in this step, it is necessary to provide the model with additional information, such as the topic and distribution, to ensure that it can be properly guided to generate meaningful data. We will now explain these pieces of information in detail.

**Chart theme:** We first generate hundreds of possible themes, which are all short phrases. When we generate data, we randomly select one from all those themes, which makes the data meaningful and diverse. This also makes it much more easy to generate questions and responses for instruction tuning.

**Data trends:** Another important characteristic of the data is the trends. We first generate several typical trend descriptions, like steadily increasing and suddenly dropping, then randomly select a few trends and require the model to generate data following them. If lacking such characteristics, the model will tend to generate several sets of data with meaningless distributions.

Column and row lengths: The lengths of columns
and lengths are also necessary for data generation.
Without specific constraints, LLMs tend to generate
excessively long or even repetitive data, which is
difficult to present in a meaningful way through
charts.

**Chart types:** Charts of different types usually share different characteristics. For example, the sum of the values in pie charts should be 100%. If not specify the type of chart, we might end up generating data that doesn't comply with the corresponding chart standards.

#### 3.2 Chart Figure Generation

322

323

325

326

327

The next step is to transform our dataset into visual 329 charts using GPT-4's coding capabilities. We used 330 popular chart plotting libraries, such as Matplotlib, 331 as our primary tools. When prompting GPT-4, we 332 provide the collected data, relevant function documentation, and in-context examples. We also give 334 detailed instructions on diversifying aspects like color schemes and line types to enhance the visual appeal of the charts. To increase the diversity and success rate of our chart generation, we ran-338 domly sample successfully generated codes as incontext examples in the prompts. Compared with previous automated chart generation efforts that 341 relied on templates, our approach offers greater va-342 riety and better visual appeal. It also enables us to 343 generalize across different chart types effectively. The result was a collection of meticulously crafted charts, each accurately representing its data and visually appealing, showcasing the effectiveness of our method. The necessary input for the prompts in this stage is listed below.

Chart data: This is the most essential input for the task. The chart data is the information that will be visualized in the chart. Without it, no meaningful chart can be made.

**Related function documentation:** This is an important reference for generating the Python code. It provides information about the available functions and features that can be used to create the chart. With the documentation, the model could even create charts in new styles that are not in the in-context examples.

In context example: These in-context examples
are sampled from pre-selected high-quality code.
This helps to facilitate the construction of the
Python code. When there is new generated code in
high quality, we can save and sample it, which is
used as in-context examples later.

Other requirements: To ensure that the final gen-367 erated code is suitable for batch processing and 368 execution, we also need to include several require-369 ments in the prompt. For example, the data is re-370 quired to be listed in the code to make the gener-371 ated code self-contained and executable without 372 the need for external files. We also set the require-373 ments for the title, axis labels, legend, and text 374 annotations. They provide context about what the 375 chart represents and make it easier to understand 376 the data. Without them, the chart can be confusing 377 and difficult to interpret. 378

379

381

382

383

384

385

386

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

#### 3.3 Instruction data generation

After completing the first two stages, we gathered comprehensive information about each chart, including precise tabular data, various characteristics from various perspectives, and the chart plotting code. Leveraging this rich information, we move on to generating a wide range of instruction-answer data with the assistance of GPT-4, significantly enhancing the capabilities of models trained on this dataset. In addition to fundamental chart understanding functionalities such as Q&A and summarization, our approach allows us to construct instructions and answers for more complex tasks, such as accurate data extraction, detailed chart descriptions, chart code generation, and even chart editing. Compared to previous pipelines for instruction data generation that often rely on human annotation, our methods yield significant time savings while enhancing diversity and quality in the resulting dataset.

Here are more details about the data that needs to be filled into the prompt.

**Chart descriptions and raw data:** providing these descriptions helps the model understand the context better. The first description helps the model to understand the nature of the data, and the second description assists in understanding the visual representation of the data. The raw data feeds the model with the actual values to base its responses on. All the descriptions and raw data are generated in the first and second stages.

**Characteristics to be asked about:** This requirement ensures that the model asks diverse and relevant questions about the chart. It prompts the model to explore different features of the data and its representation.

Mathad	Chartqa			Chart-to-text		Chart extraction (human)		Chart extraction (augmented)	
Wethod	Human	Augmented	Average	Pew	Statista	Precision	F1	precision	F1
Pix2struct (Lee et al., 2023)	30.50	81.60	56.00	10.30	38.00	-	-	-	-
Matcha (Liu et al., 2022b)	38.20	90.20	64.20	12.20	39.40	-	-	_	-
DePlot (Liu et al., 2022a)	-	-	-	-	-	81.32	81.15	93.42	93.29
Unichart (Masry et al., 2023)	43.92	88.56	66.24	12.48	38.21	61.51	35.20	79.59	70.21
Baseline* (Liu et al., 2023b)	37.68	72.96	55.32	7.16	24.65	53.48	48.39	55.17	49.50
ChartLlama	48.96	90.36	69.66	14.23	40.71	84.92	84.89	94.94	94.78

Table 1: **Results on traditional tasks.** We compare our work with the previous three open-source models and also compare it with Baseline\* trained on the training split of respective benchmarks.

Method	Detailed Description	Chart-to-chart		Text-to-chart		Chart-editing		Chart-to-text	
Wethod	Detailed Description	GPT Score Success Rate (%)		GPT Score	Success Rate (%)	GPT Score	Success Rate (%)	Pew	Statista
LLaVA-1.5 (Liu et al., 2023b)	67.2	64.8	46	62.2	77	51.6	38	65.8	73.4
ChartLlama	74.2	74.4	73	81.6	81	75.6	71	81.0	92.6

Table 2: **Results on new tasks.** We primarily compared our work with the baseline model LLaVA-1.5. For the proposed new task, we used GPT for evaluation and validated the effectiveness of our proposed dataset. Evaluation of Chart-to-text using ChatGPT is also listed.

Chart type	Unichart	Baseline*	ChartLlama
Funnel	18.30	49.32	70.59
Gantt	9.80	40.17	56.64
Heatmap	25.43	38.18	53.18
Scatter	26.32	37.91	54.97
Box	16.67	28.33	37.33
Candlestick	15.79	25.69	46.20

Table 3: **Performances of Q&A on more categories of chart.** Baseline\* means a modified version of LLaVA-1.5, which is further trained on the ChartQA dataset. We evaluate the performance of Baseline\* and the previous state-of-the-art model Unichart on these new chart types.

#### 4 Experiment

415

416

417

418

419

420

421

422

423

424

#### 4.1 Implementation details

**Implementation details.** We train ChartLlama based on LLaVA-1.5 which provides fundamental abilities crucial for chart understanding and generation, including the OCR functionality. The projection layer and LLM are trained on our proposed dataset. Details of the model architecture and training hyper-parameters can be referred to in our appendix.

Dataset statistics. We show the statistics of 425 our generated dataset in the Appendix. In our 426 instruction-tuning data, Q&A dominates while the 427 other tasks correspond to similar proportions of 428 data. This is mainly because a single chart could be 429 utilized to construct multiple Q&A data. Previous 430 datasets usually gather only three types of charts: 431 432 bar charts, line charts, and pie charts. Unlike them, we support a wide range of chart types. This is 433 mainly due to the strong flexibility of our data con-434 struction method. It's worth noting that we can 435 continue to expand on more data and chart types. 436

	Human	ChatGPT	p-value	win-rate
LLaVA-1.5	66.78	62.32	1.32e-5	68%
ChartLlama	78.01	77.36	3.98e-7	32%

Table 4: **More evaluation metrics**. The evaluation scores of our evaluation metrics and human evaluation metrics are highly aligned and verify the effectiveness of our evaluation metrics.

#### 4.2 Evaluation Benchmark and Metrics

We evaluate possible models on seven tasks, including both the traditional tasks and novel tasks which verifies that our data generation pipeline has good scalability towards various tasks and chart types. 437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

**Traditional Tasks.** Three traditional tasks are evaluated, namely ChartQA, Chart-to-text, and Chartextraction. We use Relaxed Accuracy (Masry et al., 2022), GPT-4 as metrics for evaluation, and Precision and F1 scores as metrics, respectively. The details are listed in the Appendix.

**New tasks.** In addition to traditional tasks, we have devised four additional innovative tasks, three of which are targeted at chart generation to verify the scalability of novel tasks.

1) Detailed description. This task necessitates a comprehensive description of the given chart figure in a detailed manner, rather than summarizing it. The evaluation metric for detailed description is similar to the evaluation metric in Chart-to-text using GPT-4. We include detailed descriptions of the data and chart figures as conditions for GPT-4 to assist evaluation.

2) Chart-to-chart. This task aims to reconstruct the given chart figure. We design comprehensive evaluation metrics for code generation and utilize GPT-4 to measure the quality of the code. For the

464 chart-to-chart task, we evaluated the precision of
465 data, axes, colors, chart types, and titles, rating
466 from 0 to 5. Then we average them as the score for
467 each sample. Finally, we normalize it to a range of
468 0 to 100 for easier analysis and report the average
469 score across the entire test set.

3) Text-to-chart. The task aims at generating chart 470 figures according to instructions and tabular data. 471 We provide the input instructions and the gener-472 ated code as conditions for evaluation criteria. The 473 evaluation focuses mainly on visual similarity, com-474 pleteness, accuracy, and aesthetics. Each standard 475 is equally rated from 1 to 5 points. After averaging 476 and normalization, we get the final score. 477

4) Chart-editing. The input condition for this task is 478 a chart figure and an instruction describing how to 479 edit the chart. It is expected to create a new figure 480 that has been modified according to instructions 481 based on the given chart figure. The evaluation 482 method for chart-editing uses a similar process to 483 previous chart generation-related tasks. The in-484 put conditions include the code of the chart to be 485 modified, instructions, and the generated code of 486 the model. The data accuracy, completeness, aes-487 488 thetics, and instruction following performance are scored on a scale from 0 to 5. After averaging and 489 normalization, the final result is obtained. 490

The alignment between our Evaluation Metrics and Human Evaluation. The results are listed in Table 4. We randomly select 100 samples for evaluation. Three volunteers were selected to evaluate the performance using the same criteria that were applied to GPT-4. Furthermore, we calculate the correlation between the scores given by humans and those given by ChatGPT to demonstrate that our metrics, which use ChatGPT as the evaluator, are reasonable. We also show the win-rate evaluation.

#### 4.3 Results

491

492

493

494

495

496

497

498

499

500

503

504

506

510

511

512

513

514

We first compare our methods with existing chart understanding models, such as Pix2Struct (Lee et al., 2023), Matcha (Liu et al., 2022b), unichart (Masry et al., 2023). Then we further construct Baseline\* using the same model architecture (Liu et al., 2023b) as ours, but is trained on the training split of each dataset separately. On traditional tasks, we have also tried to compare with existing multimodal large language models such as InternLM-XComposer (Zhang et al., 2023a), MiniGPT-v2 (Chen et al., 2023b). However, we found the limitation of their instruction-following ability makes it hard to be evaluated by existing metrics.

ChartQA. ChartLlama achieves the best performance on both human and augmented splits of ChartQA (Masry et al., 2022) as listed in Table 1. Our ChartLlama also succeeds LLaVA-1.5 trained on the Unichart dataset, which is shown in the Appendix in detail. Previous methods typically involved pretraining on larger datasets and then finetuning on the training split of the same datasets to achieve better results, while ChartLlama does it in a zero-shot way after training on our dataset. Notably, although previous methods are trained on the ChartQA's training split, our method achieves significant advantages using much less data as shown in Table 7. Besides, we also evaluate our model on charts of novel types as shown in Table 3. Our model gains significant improvement towards Unichart and the Baseline\*. This shows the superiority of ChartLlama in the ability to understand charts in novel charts.

**Chart-to-text.** As shown in Table 1 and Table 2, our method consistently outperforms the previous state-of-the-art approaches under different evaluation metrics and splits in Chart-to-text (Kantharaj et al., 2022). The improvement in our performance primarily stems from the model's ability to handle long texts. Previous works often encountered meaningless repetitions at the end of sentences when dealing with relatively longer texts.

**Chart extraction.** Our model performed the best in this task on ChartQA (Masry et al., 2022) as listed in Table 1. ChartLlama has been trained on a variety of instruction-tuning data, which greatly improved its ability to understand chart figures. This is the reason why it can significantly outperform LLaVA-1.5 in terms of performance.

**Detailed description.** ChartLlama gains significant performance improvement over LLaVA-1.5 which is shown in Table 2. The detailed description task requires the model's ability to understand image details, which can be significantly improved during the training for tasks related to chart figures. **Chart generation and modification.** In Table 2, we compare our method with the original LLaVA-1.5, and we can see that our model gains consistent improvement over three tasks. LLaVA-1.5, which is the base model of ChartLlama, processes strong abilities to follow instructions and generate Python code, and thus also gains reasonable performances on chart generation and modification tasks.



Figure 3: **Qualitative comparison for Chart-to-chart and Chart editing tasks.** We present the output results of LLaVA-1.5 and ChartLLaMA for the same chart given different instructions. The instruction in the first row requires the model to output the original chart, performing the chart-to-chart task. The instruction in the second row requires the model to output a horizontal bar chart, performing the chart editing task.



Figure 4: **Qualitative comparison for Text-to-chart task.** We have presented the generated images by ChartLLaMA and LLaVA-1.5 given the tabular data and the specified requirements.

#### 4.4 Qualitative results

571

574

576

578

579

Figure 3 visualizes the chart-to-chart and chartediting results. ChartLlama plots with the correct color and chart type, while LLaVA-1.5 cannot guarantee correctness. Figure 4 shows the text-to-chart results of ChartLlama and LLaVA-1.5. In the first example, ChartLlama successfully generates a funnel chart following the instructions and plots correct values. But LLaVA-1.5 even cannot draw funnel charts. In the second example, it is obvious that the result of ChartLlama contains more details and adds data values for human convenience. Both two examples show the strong ability of chartgenerating and editing abilities of ChartLlama.

#### 5 Conclusion

582 In this paper, we propose a flexible and robust ap-583 proach for synthesizing instruction-tuning data for chart data. Then we train ChartLlama on the proposed dataset. The data generation flow we propose greatly reduces the difficulty of generating chartrelated data and improves the controllability and diversity of the generated data. Experiments conducted on both traditional tasks and our newly constructed tasks validate the outstanding performance of ChartLlama. Thanks to the diverse instructiontuning data in our dataset, ChartLlama possesses various capabilities that were absent in previous models. Moreover, its ability to comprehend both instructions and figures can easily extend to new categories of chart figures or tasks. We believe that our data generation process can make significant contributions to multimodal LLM in tasks related to chart understanding. Furthermore, it will facilitate the application of similar data generation processes in other domains.

584

585

586

587

588

589

591

593

594

595

597

598

599

600

698

699

700

701

652

653

654

655

656

657

658

659

660

#### 6 Limitations

602

611

614

615

616

617

618

619

628

630

631

633

635

647

648

651

The current version of ChartLlama's vision encoder lacks the ability to handle multilingual OCR tasks, restricting the model's utility for charts containing non-English text. To overcome this limitation, we are contemplating the creation of a novel vision encoder that boasts proficiency in multilingual OCR tasks.

#### Ethics Statement

This section provides a comprehensive reflection on the broader impacts and ethical considerations associated with our multi-modal language model, ChartLlama. ChartLlama, designed to improve machine understanding of chart figures by using a high-quality instruction-tuning dataset, prompts significant ethical contemplation, especially in terms of data privacy, data representation, accuracy, accessibility, and potential misuse.

> **Data Privacy**: ChartLlama's design is founded on the principle of respecting user data privacy. The model interacts with data in a manner that does not require access to sensitive information, thereby ensuring user data confidentiality. We are committed to continually improving these privacy measures to protect user data.

**Data Representation and Accuracy**: The data used in generating the instruction-tuning dataset is carefully selected to ensure diversity and eliminate bias. Moreover, we recognize the vital importance of accuracy in our model's interpretation of chart figures. We are dedicated to continually refining our model to enhance its interpretative accuracy and reliability.

Accessibility: Our model aims to make chart figures more accessible and understandable to users who may find them challenging to interpret, including individuals with visual impairments and those not familiar with data visualization techniques. We acknowledge the current limitations of our model in this regard and are committed to improving its capabilities to make it more inclusive.

**Sustainability**: Aware of the environmental impact associated with training large AI models, our proposed data generation method aims to generate high-quality tuning data efficiently while maintaining low resource expenditure. This approach is part of our commitment to minimizing the environmental footprint of our research and development activities.

Potential Misuse: An important ethical concern

is the potential misuse of ChartLlama for malicious purposes, such as misrepresenting data or manipulating chart figures. We strictly oppose the use of our model for any unethical practices. To combat this, we advocate for:

- Monitoring and Detection: Implementing tools to detect the misuse of ChartLlama, particularly in data misrepresentation.
- Ethical Guidelines and Governance: Establishing and enforcing ethical guidelines for the use of ChartLlama, ensuring researchers and developers are aware of the ethical implications of their work.
- Collaboration with Stakeholders: Working with data scientists, researchers, and users to align ChartLlama's use with ethical standards and best practices in data visualization and interpretation.

Through this ethics statement, we aim to underscore our dedication to responsible innovation, emphasizing the importance of safeguarding against potential misuse of technology while promoting the beneficial potentials of AI in chart comprehension and data visualization.

#### References

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023a. Qwen technical report. arXiv preprint arXiv:2309.16609.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Baichuan. 2023. Baichuan 2: Open large-scale language models. arXiv preprint arXiv:2309.10305.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechu Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed

808

809

810

811

as a unified interface for vision-language multi-task learning. arXiv preprint arXiv:2310.09478. Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In International Conference on Learning Representations. Wenbo Hu, Yifan Xu, Yi Li, Weiyue Li, Zeyuan Chen, and Zhuowen Tu. 2023. Bliva: A simple multimodal llm for better handling of text-rich visual questions. 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. Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In International Conference on Machine Learning, pages 18893-18912. PMLR. Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726. Juncheng Li, Kaihang Pan, Zhiqi Ge, Minghe Gao, Hanwang Zhang, Wei Ji, Wenqiao Zhang, Tat-Seng Chua, Siliang Tang, and Yueting Zhuang. 2023b. Fine-tuning multimodal llms to follow zeroshot demonstrative instructions. arXiv preprint arXiv:2308.04152. Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023c. BLIP-2: bootstrapping language-image pretraining with frozen image encoders and large language models. In ICML. Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In Computer Vision-ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V13, pages 740-755. Springer. Fangyu Liu, Julian Martin Eisenschlos, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Wenhu Chen, Nigel Collier, and Yasemin Altun. 2022a. Deplot: One-shot visual language reasoning by plot-to-table translation. arXiv preprint arXiv:2212.10505. Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenschlos.

Elhoseiny. 2023. Minigpt-v2: large language model

702

703

706

711

712

714

715

716

717

718

719

721

722

723

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

- 2022b. Matcha: Enhancing visual language pretraining with math reasoning and chart derendering. *arXiv preprint arXiv:2212.09662*.
- Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2023a. Hallusionbench: You see what you think? or you think what you see? an image-context reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models. *arXiv preprint arXiv:2310.14566*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023b. Improved baselines with visual instruction tuning.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning.
- Yadong Lu, Chunyuan Li, Haotian Liu, Jianwei Yang, Jianfeng Gao, and Yelong Shen. 2023. An empirical study of scaling instruct-tuned large multimodal models.
- 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*.
- Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. ChartQA: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263– 2279, Dublin, Ireland. Association for Computational Linguistics.
- Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In *The IEEE Winter Conference on Applications of Computer Vision (WACV)*.
- R OpenAI. 2023. Gpt-4 technical report. arxiv 2303.08774. *View in Article*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.
- 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.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, M Chen, R Child, V Misra, P Mishkin, G Krueger, S Agarwal, et al. 2021. Dall· e: Creating images from text. *OpenAI blog. https://openai. com/blog/dall-e.*

812

els.

Robin Rombach, Andreas Blattmann, Dominik Lorenz,

resolution image synthesis with latent diffusion mod-

InternLM Team. 2023. Internlm: A multilingual lan-

ties. https://github.com/InternLM/InternLM.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier

Martinet, Marie-Anne Lachaux, Timothée Lacroix,

Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal

Azhar, Aurelien Rodriguez, Armand Joulin, Edouard

Grave, and Guillaume Lample. 2023a. Llama: Open

Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-

bert, Amjad Almahairi, Yasmine Babaei, Nikolay

Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti

Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton

Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,

Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,

Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan

Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,

Isabel Kloumann, Artem Korenev, Punit Singh Koura,

Marie-Anne Lachaux, Thibaut Lavril, Jenva Lee, Di-

ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-

bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-

stein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu,

Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Ro-

driguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and

Renqiu Xia, Bo Zhang, Haoyang Peng, Ning Liao, Peng Ye, Botian Shi, Junchi Yan, and Yu Qiao.

Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng

Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming

Yan, Yiyang Zhou, Junyang Wang, Anwen Hu,

Pengcheng Shi, Yaya Shi, Chaoya Jiang, Chenliang

Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. 2023. mplug-owl: Modularization empowers large language models with

Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. arXiv preprint

11

Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan

Wang. 2023. The dawn of lmms: Preliminary explorations with gpt-4v (ision). arXiv preprint

2023. Structchart: Perception, structuring, reasoning for visual chart understanding. arXiv preprint

fine-tuned chat models.

arXiv:2309.11268.

arXiv:2309.17421.

multimodality.

arXiv:2210.02414.

and efficient foundation language models.

guage model with progressively enhanced capabili-

High-

Patrick Esser, and Björn Ommer. 2021.

Pan Zhang, Xiaoyi Dong, Bin Wang, Yuhang Cao, Chao

Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding,

Songyang Zhang, Haodong Duan, Wenwei Zhang,

Hang Yan, Xinyue Zhang, Wei Li, Jingwen Li, Kai

Chen, Conghui He, Xingcheng Zhang, Yu Qiao,

xcomposer: A vision-language large model for ad-

vanced text-image comprehension and composition.

Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and

Yu Qiao. 2023b. Llama-adapter: Efficient fine-tuning

of language models with zero-init attention. arXiv

Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian

Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng

Wang, Wenjuan Han, and Baobao Chang. 2023.

Mmicl: Empowering vision-language model with

multi-modal in-context learning. arXiv preprint

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and

Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing

vision-language understanding with advanced large

language models. arXiv preprint arXiv:2304.10592.

Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu,

Dahua Lin, and Jiaqi Wang. 2023a.

preprint arXiv:2303.16199.

arXiv:2309.07915.

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

890

891

892

Internlm-

- 820 821

- 825

- 832
- 833

- 835

839

- 836

- 834

- 838
- 842

845

853

854

857

858

860

Method		ChartQA			ChartQA on special charts				
Method	Human	Augmented	Average	Funnel	Gantt	Heatmap	Scatter	Box	Candlestick
Unichart (Masry et al., 2023)	43.92	88.56	66.24	18.30	9.80	25.43	26.32	16.67	15.79
InternLM-XComposer-VL (Zhang et al., 2023a)	8.48	7.36	7.92	12.42	6.36	16.18	18.13	15.33	16.96
Mini-GPT-v2 (Chen et al., 2023)	15.60	8.40	12.00	26.7	15.03	28.32	28.65	21.33	17.54
Qwen-VL (Bai et al., 2023b)	37.60	63.76	50.68	6.54	9.83	13.29	7.02	8.00	1.75
mPLUG-Owl2 (Ye et al., 2023)	21.20	22.0	21.60	23.53	27.75	19.08	16.37	15.33	19.30
Baseline* (Liu et al., 2023b)	37.68	72.96	55.32	49.32	40.17	38.18	37.91	28.33	25.69
ChartLlama	48.96	90.36	69.66	70.59	56.64	53.18	54.97	37.33	46.20

Table 5: Results on traditional tasks. We compare our work with the previous three open-source models and also compare it with Baseline\* trained on the training split of respective benchmarks.



Ouestion: What's the average of all the values in the green bars (round to

decrease in popularity from the Country genre to the Classical genre?

LLaVA-1.5: 50% ChartLLaMA: 75%

Figure 5: Visualization on the ChartQA task. Here are two examples of the predictions of Unichart, LLaVA-1.5, and ChartLlama. Our proposed ChartLlama could follow the long instructions and do calculations to get the correct results.

#### Appendix

#### Model architecture Α

To elucidate our training strategies, we provide some clarification about the modifications in LLaVA-1.5 (Liu et al., 2023b), and introduce its essential model architectures.

Vision encoder: LLaVA-1.5 incorporates CLIP's vision encoder (Radford et al., 2021). The primary distinction is that LLaVA-1.5 employs ViT-L/14@336px, while LLaVA uses ViT-L/14@224px. Another notable alteration concerns the image processor. Eschewing traditional center cropping, LLaVA-1.5 adopts padding as an image pre-processing technique, ensuring that all information in the provided image can be apprehended.

**Projection layer:** In LLaVA-1.5, the initial single linear layer is substituted with a two-layer MLP, resulting in improved performance.

Lora Layer: Based on experiments in (Lu et al., 911 2023; Liu et al., 2023b), implementing Lora (Hu 912 et al., 2022) layers is sufficient to achieve perfor-913 mance comparable to full fine-tuning strategies. 914

Prompt Design	Successful Rate
Original	85%
w/o In context example	43%
<i>w/o</i> Documentation	65%
w/o Both	28%

Table 6: Ablations on Prompt of Stage Two. The first row shows the successful rate of our proposed data generation method in the second stage. Then we evaluate the generated results when removing the in-context examples, the documentation, and both of them, respectively.

Datasets	#Chart type	#Chart figure	#Instruction tuning data	#Task type
ChartQA (Masry et al., 2022)	3	21.9K	32.7K	1
PlotQA (Methani et al., 2020)	3	224K	28M	1
Chart-to-text (Kantharaj et al., 2022)	6	44K	44K	1
Unichart (Masry et al., 2023)	3	627K	7M	3
StructChart (Xia et al., 2023)	3	9K	9K	1
ChartLlama	10	11K	160K	7

Table 7: Dataset statistics. Thanks to the flexibility of our data construction method, our proposed dataset supports a wider range of chart types and tasks. We can generate more diverse instruction-tuning data based on specific requirements.

For the original LLaVA (Liu et al., 2023c), Lora layers with a Lora rank of 64 suffice, whereas for LLaVA-1.5 (Liu et al., 2023b), the Lora rank needs to exceed 128.

#### **Evaluation Metrics for Traditional** B Tasks

1) For ChartQA (Masry et al., 2022), we evaluate relaxed accuracy on human and augmentation split, respectively.

2) Chart-to-text contains two separate datasets for training and evaluation. BLEU-4 and GPT-4 serve as metrics for evaluation. BLEU-4 is widely used in many NLP tasks. However, the Chart-to-text datasets contain too few ground-truth references, which harms the evaluation metrics. To facilitate more reasonable evaluations, we propose a new evaluation metric based on GPT-4, referring to the GPTScore (Fu et al., 2023). We designed scoring criteria that require the ground-truth reference and

907 908

909

910

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932



Figure 6: **Visualization of Chart-extraction.** We find that ChartLlama is especially good at long text processing. While the previous SOTA, Unichart, will generate meaningless redundant words when the output is too long.



Figure 7: **Visualization of Chart-to-text.** We select one image from the Pew Dataset and show the results of Unichart, LLaVA-1.5, and ChartLlama. We find that Unichart easily falls into repeated words again and LLaVA-1.5 suffers from hallucination.

raw data as input conditions. Details can be found in the appendix.

934 935

937

939

941

947

948

949

952

953

956

960

3) Chart extraction aims to extract the tabular data from the given chart figure. We follow the evaluation framework of DePlot (Liu et al., 2022a) and report the Precision and F1 scores on the challenging ChartQA dataset, which also provides the tabular data for each chart figure.

#### C Dataset Scale and Statistics

Human Evaluation of Data Quality. We sample 100 instances from our proposed dataset to assess the quality and then establish four criteria for volunteers to rate the quality of the provided training samples. These criteria gauge the aesthetics of the generated figure, the presence of chart occlusion, the accuracy of the given answer, and the quality of the image design. Each scoring criterion carries a maximum of one point. The average score for the provided images is 3.7/4.0. We also evaluate UniChart using the same criteria, and its average score is 3.1/4.0. This indicates the superior quality of our generated chart dataset.

The dataset's statistics are shown in Table 7 and Fig. 16. The model's training process is broken down into two critical stages: pretraining and finetuning. The primary objective of pretraining is to effectively initialize the vision projector while finetuning steers the Language Learning Model (LLM) to adhere to the provided instructions.

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

In the pretraining phase, LLaVA-1.5 utilizes approximately 558k image-caption pairs to train the projection layer. It is anticipated that the vision features will align with the language features to a certain extent. This dataset originates from a subset of around 558K image-text pairs from LAION-CC-SBU, each paired with a BLIP caption.

The fine-tuning phase involves further training of the model on 665k instruction-following data pairs. LLaVA-1.5 manifests an array of capabilities during this stage. The instruction-following data pairs are meticulously generated to encompass the required abilities. To enhance the model's capacities in varied contexts, additional academictask-focused Visual Question-Answering (VQA) datasets for VQA, Optical Character Recognition (OCR), and region-level perception are incorporated. The final compilation includes several datasets: OpenKnowledge VQA (OKVQA, A-OKVQA), Region-level VQA (Visual Genome, RefCOCO), and OCR (OCRVQA, TextCaps). A-OKVQA is transformed into multiple-choice questions, employing a specific response formatting prompt: answer by directly specifying the option's letter from the provided choices.



Figure 8: **The prompt template for Stage One.** This template is used for Chart Data Generation. Utilizing this template could guide GPT-4 to generate diverse raw tabular data and detailed descriptions of the content.

# SchatGPT for Chart Figure Plotting

You are a specialist in two aspects, drawing charts with matplotlib or plotly, and providing detailed descriptions about the chart. You receive the data in the format of csv table. In addition, you are provided with an example of Python code drawing a chart for reference. You also receive some parameters that could be used to increase the diversity. You need to generate Python code to plot the given data as a chart figure and providing detailed description about the figure. Additional requirements:

The chart should have the title, labels on x-axis and y-axis. The chart should have legend. You can annotate data values above the point on the chart figure. Do not use show function to show the figure. The csv data should be listed in the code.

The output contains two parts. The first part is the generated Python code wrapped in <code start> and <code end>. Next is the detailed description about the chart wrapped in <description start> and <description end>. The code should be able to be executed without external files.

The given data: [Chart data].

989

991

995

997

1002

The given code example: [Code example].

As for additional parameters, you could consider: [Related function documentation].

Figure 9: The prompt template for Chart Figure Plotting. Following such instructions, GPT-4 could generate codes that could draw chart figures using Python packages.

### **D** Generation Prompt for ChartLlama

As listed in Figure 8, Figure 9, and Figure 10, we have provided standard prompts for data generation in three stages. The text in black color in the figure denotes the fixed prompt template, while the text in red color brackets requires filling in, which serves to enhance the diversity and controllability of the generated results. The detailed meanings of the different variables have already been discussed in the main text, thus we will not elaborate further.

### E Ablation Study on the Conditions of Generation Prompt for ChartLlama

In order to verify the impact of our proposed generation process on the results, we designed an ablation experiment on the prompt for the second step, diagram construction, which is shown in Table 6. 1003 Specifically, we removed the in-context examples 1004 and the description of the function, then retested the 1005 probability of successful generation. The results 1006 show that combining both in-context examples and 1007 documentation could significantly improve the suc-1008 cessful rate of plotting figures. Also, we observe that the diversity could also improve a lot, which is 1010 hard to quantify. 1011

### F Filtering Mechanism

The data generation process may produce some erroneous samples, but filtering and correcting these1013roneous samples, but filtering and correcting these1014samples can be challenging because the samples1015contain figures that cannot be processed by GPT-10164. We only performed basic error correction, in-1017

# S ChatGPT for Instruction Data Generation

You are an AI visual assistant that can analyze chart figures. You receive two detailed descriptions and raw data about the same chart. The first description is the information about the raw data in the chart. The second description is about the chart figure based on Python code. In addition, raw data values within the chart is given. Answer all questions as you are seeing the chart figure. Design a question-answer pair between you and a person asking about this chart figure. The answers should be a single word or phrase, and in a tone that a visual AI assistant is seeing the chart figure and answering the question.

Ask diverse questions and give corresponding answers. Include questions asking about [Characteristics] and so on. Only include questions that have definite answers:(1) one can see in the chart figure that the question asks about and can answer confidently;(2) one can determine confidently from the chart figure that it is not in the chart figure. Do not ask any question that cannot be answered confidently. The answers should be a single word or phrase. Here are some examples and remember to follow their format: [In context examples]. The first description: [Description about chart data]. The second description: [Description about chart figure].

The raw data: [Raw data].

1018

1019 1020

1021

1022

1023

1024

1026

1027

1028

1029

1030

1031

1032

1034

1035

1036

1037

1039

1040

1042

Figure 10: **The prompt template for Instruction Data Generation.** This step is targeted at generating various training data. To guarantee the quality and diversity of the generated samples, we need to give enough information on the chart figure and in-context examples.

cluding checking the data generation format and verifying the correct execution of the code. The data generation format check involves confirming whether the model has separated different data results with different markers according to our requirements. The check for correct code execution involves running the generated plotting script. If this script fails to run, we no longer use the training sample corresponding to that plot. Such basic data screening is sufficient to ensure the quality of the generated dataset. We are also considering incorporating more effective automatic screening mechanisms to avoid contamination of the dataset by poor-quality samples.

#### G Evaluation Prompt for ChartLlama

We have prepared five evaluation prompts in total, each tailored for a specific task: chart-to-text in Figure 15, detailed description in Figure 11, chartto-chart in Figure 12, text-to-chart in Figure 13, and chart-editing in Figure 14. We have designed distinctive scoring criteria for different tasks and provided reference information based on the additional annotations in the dataset. Ultimately, we employed GPT-4 for scoring purposes.

#### H Comparison with Multi-modal LLMs

1043Comparison with LLaVA-1.5 Trained on Al-1044ternate Datasets. Of all the pre-existing chart-1045comprehending datasets, UniChart is the largest1046and encompasses the most task types. Thus,1047we compare our ChartLlama with a novel model

Dataset	Cł	nartQA	Chart Extraction			
Dataset	human	augmented	human	augmented		
UniChart	29.36	40.0	26.54	60.40		
Ours	48.96	48.96 90.36		94.78		

Table 8: Comparison with LLaVA-1.5 trained on thedataset proposed in Unichart.

trained on UniChart using the LLaVA-1.5 framework. The results pertinent to ChartQA-related tasks are presented in the table above. The evaluation metrics of ChartQA and Chart Extraction are Relaxed Accuracy and F1 score, respectively. As expected, the performance of LLaVA-1.5 trained on UniChart is subpar, because each sample only includes simple question-answer pairs or succinct captions, which fail to provide sufficient information for the model to interpret the provided chart figures.

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

Traditional Tasks: The Table 5 includes a comparison of existing state-of-the-art (SOTA) models, 1060 illustrating their respective performances. Interest-1061 ingly, some models (Zhang et al., 2023a) show un-1062 expectedly low performance. This outcome is not 1063 a consequence of our experimental configuration. 1064 Rather, it derives from the fact that these models 1065 have not been trained on corresponding instruction-1066 following tasks, which results in outputs that are 1067 incompatible with the evaluation framework. We 1068 argue that training these models specifically on 1069 instruction-following tasks using specific datasets 1070 would likely yield improved performance. Another notable observation is the performance gap 1072

S ChatGPT for Evaluation of Detailed Description
You are an expert tasked with evaluating the descriptions generated by a model. I will provide you with the ground-truth description and raw data, as shown below: [Ground truth description]
IRaw datal
The description generated by the model is as follows:
[Predicted description]
Please refer to the above content and score the model on the given criteria. The criteria are as follows: **0 points:**
- The model's description doesn't refer to the chart at all, or is completely irrelevant.
- It doesn't show any understanding of the chart figures or raw data. **1 point:**
- The model's description refers to the chart but the details are largely incorrect.
- There is minimal understanding of the chart figures and raw data.
**2 points:**
- The model's description refers to the chart and some details are correct, but key elements are missing or incorrect.
- There is a basic understanding of the chart figures and raw data, but significant errors are present. **3 points:**
- The model's description accurately refers to the chart and most details are correct.
- The model shows a good understanding of the chart figures and raw data, but there are some minor errors or omissions. **4 points:**
- The model's description accurately refers to the chart and all details are correct.
- The model shows a very good understanding of the chart figures and raw data.
- There might be minor improvements possible in the description's clarity or completeness. **5 points:**
- The model's description accurately and comprehensively describes the chart.
- The model shows an excellent understanding of the chart figures and raw data, with no errors or omissions.
- The description is clear, detailed, and precise. It could be used as a standalone explanation of the chart.
First return a single value (from 0 to 5) in the first line, then reply with your reason in the second line.

Figure 11: **The prompt template used for evaluation on the Detailed-description task.** The input conditions are the ground-truth description, raw data, and predicted description. GPT-4 will follow the given criteria and generate the final score and reasons.

of Qwen-VL between the ChartQA test splits and 1073 1074 the ChartQA on our specially generated charts. Despite being trained on ChartQA, Qwen-VL under-1075 performs on the specially generated charts, under-1076 scoring the effectiveness and need for our proposed benchmark. However, the lack of general training 1079 scripts provided by many models poses a challenge to our fine-tuning efforts. Nonetheless, our hypoth-1080 esis finds support in the model LLaVA-1.5. Ini-1081 tially, LLaVA-1.5 performed poorly on the dataset but showed significant improvement when trained 1083 on the designated dataset. 1084

Novel Tasks: We also conducted tests on the 1085 newly proposed tasks. However, most of the given 1086 dataset cannot generate executable Python code ex-1087 cept LLaVA-1.5 (Liu et al., 2023b). We speculate 1088 1089 that this is because these large multimodal models have been overtrained on visual language datasets, resulting in the loss of their code generation ca-1091 pabilities in Language Learning Models (LLMs); while LLaVA-1.5 adopted a series of optimization 1093

measures during its training process. For instance, 1094 compared to other large multimodal language models, LLaVA-1.5 has a shorter training time, fewer 1096 training parameters, a more moderate dataset scale, 1097 and incorporates pure text data during training to 1098 maintain the basic capabilities of LLMs. This ex-1099 periment also suggests that if we expect the model 1100 to have a certain level of generalization ability, we 1101 should avoid making excessive adjustments to the 1102 LLMs. This is also why our ChartLlama model 1103 chose to train with fewer parameters. 1104

### I More Qualitative Results

**ChartQA.** As shown in Figure 5, we compare 1106 our ChartLLaMA with Unichart and LLaVA-1.5. 1107 The given examples are both related to longer ques-1108 tions and calculations, which is hard for Unichart. 1109 What's more, without the language understanding 1110 ability, Unichart even cannot follow complex in-1111 structions. In Example 2, the answer of Unichart 1112 is even not a percentage. Although LLaVA-1.5 has 1113 the ability of OCR and instruction-tuning, it cannot 1114

- identify which part of the image is related to thequestion because it has not been trained on chartfigures. Thus, it fails in both examples, either.
- Chart Extraction. As depicted in Figure 6, 1118 ChartLLaMA also possesses the capability to con-1119 vert charts into structured data. Both the output 1120 results of Unichart and ChartLLaMA are a string 1121 of characters and we visualize it as tables for con-1122 venience. The first mistake of Unichart is reversing 1123 the order of years. Another mistake in Unichart 1124 is the persistent output of repetitive and meaning-1125 less characters at the end. Meanwhile, our pro-1126 posed model, ChartLLaMA, benefits from strong 1127 language comprehension and output capabilities, 1128 which prevent the occurrence of such errors. 1129
- Chart Description. In Figure 7, we visualize 1130 the results of Unichart, LLaVA-1.5, and ChartL-1131 LaMA on the Chart-to-text task. The results from 1132 Unichart contain incorrect values and meaningless 1133 repetitions when generating long texts. LLaVA-1134 1.5 performs better for long output sequences but 1135 suffers from wrong OCR recognition results and 1136 hallucinations. Our proposed ChartLLaMA per-1137 forms best among these three models. 1138

# S ChatGPT for Evaluation of Chart-to-chart

You are an expert tasked with evaluating the Python code generated by a model. I will provide you with the groundtruth code used for generating the chart figure, and the predicted python code. You need to evaluate the predicted Python code and score it from 0 points to 5 points. **Here is the criteria:** 

\*\*0 points:\*\*

- The model's generated Python code either does not produce a chart at all, or the chart is entirely unrelated to the original.

- It doesn't show any understanding of the chart figure or the Python code used to produce it.

\*\*1 point:\*\*

- The model's generated Python code refers to the original chart but the details are largely incorrect.

- There is minimal understanding of the chart figure and the Python code used to produce it.

\*\*2 points:\*\*

- The model's generated Python code refers to the original chart and some details are correct, but key elements are missing or incorrect.

- There is a basic understanding of the chart figure and the Python code used to produce it, but significant errors are present.

\*\*3 points:\*\*

The model's generated Python code accurately refers to the original chart and most details are correct.
The model shows a good understanding of the chart figure and the Python code used to produce it, but there are some minor errors or omissions.

\*\*4 points:\*\*

- The model's generated Python code accurately refers to the original chart and all details are correct.

- The model shows a very good understanding of the chart figure and the Python code used to produce it.
- There might be minor improvements possible in the code's clarity or completeness.

\*\*5 points:\*\*

- The model's generated Python code accurately and comprehensively reproduces the original chart.

- The model shows an excellent understanding of the chart figure and the Python code used to produce it, with no errors or omissions.

- The generated Python code is clear, detailed, and precise. It could be used as a standalone code to draw the chart.

The ground-truth Python code: [The Ground-truth code]

The predicted Python code: [The predicted code]

First return a single value (from 0 to 5) in the first line, then reply with your reason in the second line.

Figure 12: The prompt template used for evaluation on the Chart-to-chart task. The input conditions are ground-truth code and predicted code. Following the given criteria, GPT-4 generates the score and corresponding reason.

# S ChatGPT for Evaluation of Text-to-chart

You are an expert tasked with evaluating the Python code generated by an LLM. The LLM could generate Python code for chart figure based on input raw csv data and instructions. I will provide you with the input raw data and the instruction. Also, I will give you a reference code, and the predicted python code by the LLM. You need to evaluate the predicted Python code and score it from 0 points to 5 points. **Here is the criteria:** 

1. **\*\*Correctness\*\*:** This metric evaluates whether the generated code accurately fulfills the given instructions. The score could be binary (1 for correct, 0 for incorrect) or based on a proportion of test cases passed.

Scoring Standard:

- Score 5: The code fulfills all tasks perfectly.
- Score 3: The code fulfills some tasks but has minor errors.
- Score 0: The code does not fulfill the tasks or is entirely incorrect.

2. **\*\*Readability\*\*:** This metric assesses whether the code is easy to read and understand, which includes appropriate use of comments, variable names, and code structure.

Scoring Standard:

- Score 5: The code is very readable with good structure, comments, and variable names.
- Score 3: The code is somewhat readable but could be improved.
- Score 0: The code is not readable or poorly structured.

3. **\*\*Visualization Aesthetics and Detailing**\*\*: This metric evaluates the level of detailing in the generated figures and the aesthetics of the visualization. It assesses how well the code incorporates elements like color, labels, annotations, and other features to improve the look and interpretability of the graphs.

Scoring Standard:

- Score 5: The code consistently generates figures with excellent detailing and aesthetics. Graphs have appropriate and diverse color schemes, clear labels, and annotations, making them easy to interpret and visually appealing.

- Score 3: The code generates figures with adequate detailing and aesthetics. Some elements like color, labels, or annotations could be improved for better interpretability and visual appeal.

- Score 0: The code does not generate figures, or the figures generated lack any form of detailing or aesthetics, making them uninterpretable and visually unappealing.

The raw data and instruction: [Raw tabular data and Instructions]

The reference Python code: [Reference code]

The predicted Python code: [Predicted code]

The output should first give the average score based on three criteria, then output scores for each criteria. The output should follow this format: Average: number Correctness: number Readability: number Visualization Aesthetics and Detailing: number

Figure 13: **The prompt template used for evaluation on the Text-to-chart task.** The input conditions are the raw tabular data and instructions, the reference code, and the predicted code. Finally, GPT-4 will return with the average score and scores for each criterion.

# S ChatGPT for Evaluation of Chart-editing

You are an expert evaluator tasked with assessing the performance of a model on a chart-editing task. You will be provided with the original code of the chart, the instructions given to the model, and the code generated by the model.

The original code: [Original code] Instructions: [Instructions] The generated code: [Generated code]

Using the criteria below, please score the model's performance:

#### \*\*Data Accuracy\*\*

0 points: The model makes no modifications to the chart based on the instructions.

- 1 point: The model makes some modifications to the chart, but they are largely incorrect based on the instructions.
- 2 points: The model makes modifications to the chart and some are correct, but key elements are missing or incorrect.
- 3 points: The model makes accurate modifications to the chart and most are correct based on the instructions.
- 4 points: The model makes accurate modifications to the chart and all are correct based on the instructions.
- 5 points: The model makes accurate and comprehensive modifications to the chart based on the instructions.

#### \*\*Completeness\*\*

- 0 points: The generated code is incomplete and shows no understanding of the original chart or instructions.
- 1 point: The generated code is partially complete but shows minimal understanding of the original chart and instructions.
- 2 points: The generated code is mostly complete but lacks some key elements or shows significant errors.
- 3 points: The generated code is largely complete with only minor elements missing or incorrect.
- 4 points: The generated code is almost entirely complete with only minor improvements possible.
- 5 points: The generated code is completely detailed and precise.

#### \*\*Aesthetics\*\*

- 0 points: The model fails to maintain or improve the aesthetics of the original chart.
- 1 point: The model makes some aesthetic modifications, but they are largely incorrect or inappropriate.
- 2 points: The model makes aesthetic modifications and some are correct, but key elements are missing or incorrect.
- 3 points: The model maintains or improves the aesthetics of the chart with some minor errors or omissions.
- 4 points: The model significantly enhances the aesthetics of the chart with only minor improvements possible.
- 5 points: The model excellently enhances the aesthetics of the chart with no improvements needed.

\*\*Instruction Following Performance\*\*

- 0 points: The model fails to follow the instructions at all.
- 1 point: The model follows some parts of the instructions but misses out on major aspects.
- 2 points: The model follows the instructions to a basic extent but misses out on or incorrectly interprets key elements.
- 3 points: The model largely follows the instructions with only minor elements missed or incorrectly interpreted.
- 4 points: The model follows the instructions almost entirely with only minor improvements possible.
- 5 points: The model follows the instructions excellently with no elements missed or incorrectly interpreted.

The output should first give the average score based on three criteria, then output scores for each criteria

Figure 14: **The prompt template used for evaluation on the Chart-editing task.** Input conditions include the original code corresponding to the given chart figure, the instructions that describe how to edit the figure, and the generated code. The output will contain the final average score and scores for each criterion.

ChatGPT for Evaluation of Chart-to-text
You are an expert evaluator assessing the performance of multi-modal LLM in generating detailed descriptions based on chart figure. Your task is to evaluate the description generated by multi-modal LLM. You receive the ground truth description and raw data for reference. Here is the ground truth description for reference: [Ground-truth description]
Here is the raw data: [Raw data]
Now, here is the description generated by the multi-modal LLM: [Predicted description]
Using the reference and generated descriptions above, please rate the performance of multi-modal LLM on a scale of 0 to 5 based on the following criteria:
<ul> <li>**0 Points:**</li> <li>The generated description doesn't reference the chart data at all or is completely irrelevant.</li> <li>Multi-modal LLM doesn't show any understanding of the chart figures or raw data.</li> <li>**1 Point:**</li> </ul>
<ul> <li>The generated description refers to the chart, but most details are incorrect.</li> <li>Multi-modal LLM displays minimal understanding of the chart figures and raw data.</li> <li>*2 Points:**</li> </ul>
- The generated description refers to the chart, and some details are correct, but important elements are missing or incorrect.
- Multi-modal LLM shows basic understanding of the chart figures and raw data, but significant errors are present.
- The generated description accurately refers to the chart, and most details are correct.
<ul> <li>Multi-modal LLM shows a good understanding of the chart figures and raw data, but there are some minor errors or omissions.</li> <li>**4 Points:**</li> </ul>
- The generated description accurately refers to the chart, and all details are correct.
- Multi-modal LLM shows a strong understanding of the chart figures and raw data.
<ul> <li>There could be minor improvements in the clarity or completeness of the description.</li> <li>**5 Points:**</li> </ul>
<ul> <li>The generated description accurately and comprehensively describes the chart.</li> <li>Multi-modal LLM shows an excellent understanding of the chart figures and raw data, with no errors or omissions.</li> <li>The description is clear, detailed, and precise. It could be used as a standalone explanation of the chart.</li> </ul>
Please provide your score in the first line and explain your rating in the second line.

Figure 15: **The prompt template used for evaluation on the Chart-to-text task.** The ground-truth description, raw data, and predicted description are input conditions. This evaluation prompt requires GPT-4 to give the final score and explanations of the given score.



Figure 16: **Distributions of different types of data in our dataset**. The top and bottom pie charts show the distribution of task types and chart types, respectively. (The illustration is generated by our proposed ChartL-lama.)