SIMPLOT: Enhancing Chart Question Answering by Distilling Essentials

Wonjoong Kim^{1*}, Sangwu Park^{2*}, Yeonjun In¹, Seokwon Han¹, Chanyoung Park^{1†}

¹KAIST ²Chung-Ang University

{wjkim, yeonjun.in, dulwich10, cy.park}@kaist.ac.kr sangwu99@cau.ac.kr

Abstract

Recently, interpreting complex charts with logical reasoning has emerged as challenges due to the development of vision-language models. A prior state-of-the-art (SOTA) model has presented an end-to-end method that leverages the vision-language model to convert charts into table format utilizing Large Language Model (LLM) for reasoning. However, unlike natural images, charts contain a mix of essential and irrelevant information required for chart reasoning, and we discover that this characteristic can lower the performance of chartto-table extraction. In this paper, we introduce SIMPLOT, a method designed to extract only the elements necessary for chart reasoning. The proposed method involves two steps: 1) training to mimic a simple plot that contains only the essential information from a complex chart for table extraction, followed by 2) performing reasoning based on the table. Our model enables accurate chart reasoning without the need for additional annotations or datasets, and its effectiveness is demonstrated through various experiments. Furthermore, we propose a novel prompt mimicking how human interpret charts for more accurate reasoning. Our source code is available at https://github.com/sangwu99/Simplot.

1 Introduction

The rapid advancements in vision-language models have accelerated a wave of research into models capable of handling data that integrates both images and text (Zhang et al., 2021; Zhou et al., 2020; Kim et al., 2023, 2024), thus undertaking a variety of tasks. Among these tasks, the interest in models capable of advanced reasoning has significantly increased. This increasing field has seen considerable success in addressing visual question answering (VQA) tasks targeted at natural images,

marking a trend towards models that can engage in complex reasoning based on images (Antol et al., 2015; Shao et al., 2023; Gardères et al., 2020). Despite these advancements, the domain of mathematical multimodal reasoning such as interpreting charts remains relatively unexplored. Mathematical reasoning in question answering poses unique challenges, as models proficient in natural images struggle with specific types, such as charts. Charts, with their unique formats and the need for logical interpretation, necessitate a different approach to learning compared to conventional VQA models targeting natural images.

Prior chart reasoning methods mainly rely on heuristic rule-based systems, and thus they are not only limited to pre-defined chart formats but also struggle with novel chart types without additional rule formulation (Luo et al., 2021). Moreover, the performance of models utilizing OCR or key-point detection modules are highly dependent on the performance of these modules, and also face significant annotation costs and are typically unable to perform end-to-end reasoning (Methani et al., 2020; Poco and Heer, 2017). In response to these limitations, recent approaches have adopted visionlanguage models trained in an end-to-end manner without heuristic rules (Cheng et al., 2023; Liu et al., 2022b), which, however, require fine-tuning for each specific downstream task, limiting their

A novel approach, called Deplot (Liu et al., 2022a), that combines vision-language models with Large Language Models (LLMs) has emerged as an approach to addressing these issues. Specifically, Deplot first transforms charts into tables (i.e., chart-to-table extraction), and then employs the extracted tables along with the LLMs for reasoning. This approach not only aims to solve the inherent problems of previous methodologies, but also leverages the capability of the LLMs to enhance performance in chart question answering. Converting charts into

^{*}These authors contributed equally.

[†]Corresponding author.

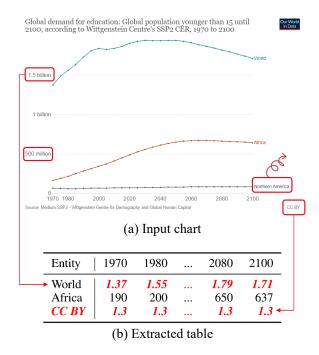


Figure 1: An example of chart-to-table extraction (from (a) to (b)) by Deplot (Liu et al., 2022a) on ChartQA.

tables before reasoning offers several advantages, including improved interpretability and the ability to achieve high performance in table reasoning, thus facilitates more accurate and precise reasoning compared to a direct image-based QA.

Despite its effectiveness, Deplot suffers from the following two limitations that still need to be addressed. First, we discovered that Deplot struggles to fully utilize the textual context within a chart for chart-to-table extraction (See Fig. 1). For example, the counting unit (e.g., "million" or "billion") is overlooked when converting a chart (Fig. 1(a)) into a table (Fig. 1(b)). Moreover, irrelevant information (e.g., "CC BY") is extracted as an entity in the table, while relevant information (e.g., "Northern America") is overlooked. However, as real-world charts often contain information that might not be helpful for chart reasoning (e.g., source credit and "CC BY"), it is crucial for the model to be able to differentiate between relevant and irrelevant information. Owing to these challenges, Deplot is prone to eventually extracting inaccurate values within a table, failing in effectively interpreting the information of the chart.

The second limitation of Deplot is that it solely relies on tables to solve chart reasoning tasks, while overlooking the visual information associated with the tables. For this reason, Deplot fails to answer questions regarding the information that cannot be obtained from the extracted table itself (e.g., "What is the value of the third bar from the top?", "What year does the orange line represent?"), indicating a significant shortfall in the model's capability to interpret visual data.

In this paper, we propose a simple yet effective method, named SIMPLOT, to address the aforementioned limitations of Deplot. To handle the first limitation, we develop two methods aimed at providing explicit supervision to the model. The first method, named row-col rendering, involves explicitly inserting information about the rows and columns that an extracted table should contain, allowing the model to perceive more accurate textual context. The second method involves conveying only the essential information in the chart to the model, which is built on our observation that converting an original chart containing irrelevant information into a simple chart containing only the essential information significantly improves the performance of table extraction.

To address the second limitation, we present a novel prompt named *Human-oriented chart instruction* while leveraging a Large Multimodal Model (LMM) to utilize the visual attributes of a chart. Since more advanced reasoning is required when interpreting charts, LMM requires a prompt that is specifically designed for the chart reasoning task, differentiated from natural images. Hence, we provide instructions to the model in a way that is similar to how humans interpret charts for precise reasoning.

Our contributions are summarized as follows:

- 1. We provide guidelines for utilizing textual information within charts for reasoning that was previously overlooked.
- SIMPLOT extracts only essential information from complex charts, preventing irrelevant information from entering the model, resulting in detailed reasoning.
- 3. We present a prompt specifically designed for chart reasoning. Through this prompt, LMMs perform more accurate reasoning by mimicking how humans interpret charts.
- 4. Extensive experiments validate that SIM-PLOT successfully addresses the limitations of Deplot. A further appeal of SIMPLOT is that it is model-agnostic, i.e., it can be applied to any existing model that involves chart-to-table extraction for chart reasoning beyond Deplot.

2 Proposed Method: SIMPLOT

In this section, we describe our proposed method, SIMPLOT. Fig 2 presents the overall framework of SIMPLOT, where SIMPLOT eventually performs reasoning with an LMM given a table extracted from a chart (Inference stage). In this work, we focus on chart-to-table extraction (Training stage) as a high-quality table enables precise reasoning with a powerful LMM. It is important to note that SIMPLOT can be combined with diverse LMM variants, and various techniques such as Chain-of-Thought (CoT) (Wei et al., 2022) can further improve the accuracy of the chart reasoning, which we leave as future work.

2.1 Preprocessing Stage

Before the training stage, we conduct two preprocessing steps, i.e., 1) Simple Chart Generation, and 2) Row-Column Rendering. Note that as these steps can be readily done offline, they do not increase the training time.

1) Simple Chart Generation. We generate a simple chart by excluding irrelevant information from the original chart, keeping only essential elements required for reasoning. Specifically, since each of the original charts in the dataset is annotated with a table in the CSV format, we use a Python library (i.e., Matplotlib) to plot a chart based on the table. The chart generated in this manner is considered as a simplified version of the original chart, i.e., simple chart. This process requires no separate training or additional costs, and it can be executed offline by running a simple code snippet.

2) Row-Column Rendering. Inspired by rendering questions over images used for QA (Lee et al., 2023), we aim to improve the table extraction accuracy by rendering information about the rows and columns that should be included in the table onto the image, enabling the model to utilize this information to extract more accurate tables when converting charts to tables. Note that in the training stage, since we are given the ground-truth charttable pairs, we simply render rows and columns of each table onto its paired image containing the chart. However, it would not be feasible in the inference stage, since we would be only given the charts without tables. Instead, in the inference stage, we utilize a LMM to extract rows and columns from the chart and render them onto the image. Despite the relative low performance of LMM such as GPT-4 (Achiam et al., 2023) in reasoning about

charts, we found they can accurately extract row and column information from charts, since such information is given as text and structured in a relatively simple manner. Detailed description of the process is presented in Appendix B.

2.2 Training Stage: Chart-to-Table Extraction

Our proposed chart-to-table extraction approach consists of two phases: **Phase 1**) Training a teacher encoder and a table decoder by performing the chart-to-table extraction given simple charts containing only the essential information for reasoning, rather than complex original charts; **Phase 2**) Training a student encoder via contrastive learning by extracting the table given the original chart, while being distilled the knowledge to embed original charts to the embedding space of simple charts.

2.2.1 Phase 1: Learning with Essential Part from Simple Chart

Our model learns the process of extracting tables from previously generated simple charts. Specifically, we use Deplot (Liu et al., 2022a) as our backbone model, which is our baseline model that consists of an image encoder and a text decoder, and fine-tune it on the generated simple chart-table pairs. This facilitates the image encoder to obtain representations containing only essential information within the chart, while the text decoder converts this representation into a table format. In this paper, we name the encoder and decoder as *chart encoder* E_{chart} and *table decoder* D_{table} , respectively.

To corroborate our model design of focusing on the essential parts when learning from charts, in Fig. 3, we compare the chart-to-table extraction performance when using either one of the two chart types (i.e., original chart (in green), and simple chart containing essential information only (in red)) for both training and inference. We observe that using simple charts for chart-to-table extraction greatly outperforms the case when original charts are used.

Note that the chart encoder trained here serves as the teacher encoder $E_{chart}^{teacher}$, providing guidelines for the representation that the student encoder $E_{chart}^{student}$ should learn in the subsequent learning process to be described in Sec. 2.2.2. Following Deplot, we adopt ViT (Dosovitskiy et al., 2020) as the chart encoders. Note that the trained table decoder D_{table} and a fully connected layer (FC) remain frozen in the next stage.

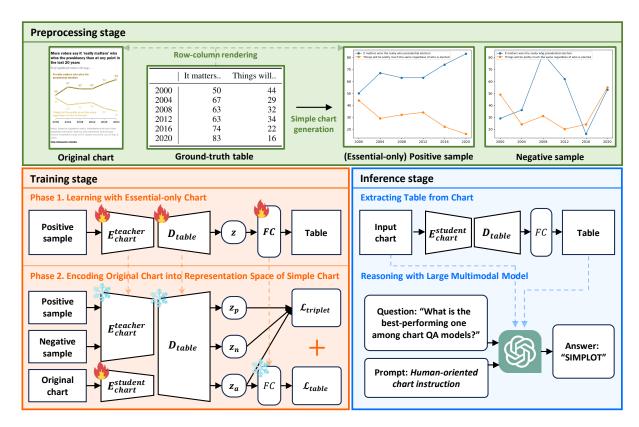


Figure 2: Overall framework of SIMPLOT. Upper box presents the preprocessing stage, which involves generating a simple positive sample containing only essential information from the original chart, as well as a negative sample, along with row and column rendering. Lower left box illustrates the training stage including two phases. In Phase 1 of the training stage, a teacher encoder and a table decoder are trained using a simple chart, and in Phase 2, a student encoder is trained with the original chart, while being distilled the knowledge of the teacher encoder on how to generate a table from a simple chart. Lower right box illustrates the inference stage, where an LMM receives the original chart and the extracted table along with prompts for reasoning.

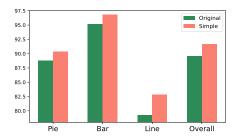


Figure 3: Accuracy of chart-to-table extraction using the original (in green) and simple charts (in red).

2.2.2 Phase 2: Encoding Original Chart into Representation Space of Simple Chart

As shown in Fig. 2, after training the teacher encoder $E_{chart}^{teacher}$ with simple charts as inputs, our goal is to train a student encoder $E_{chart}^{student}$ in a way that the representation of an original chart image obtained from $E_{chart}^{student}$ closely matches the representation of its corresponding simplified chart obtained from $E_{chart}^{teacher}$. This is to enable $E_{chart}^{student}$ to encode any original chart into the representa-

tion space of simplified chart images, thereby improving the accuracy of table generation. Then, given the chart representations, we generate table representations z using the frozen table decoder D_{table} that is trained in Phase 1. More precisely, given an original chart A and the simplified chart P as a positive sample, we generate a negative sample N by randomly shuffling values in the simple chart P. Then, we define a triplet loss using the representations of the original chart z_a , its corresponding positive sample z_p , and negative sample z_n obtained from the sequentially passing A, P and N into chart encoders and the table decoder D_{table} as follows:

$$\mathcal{L}_{triplet}(A, P, N) = \max\{d(z_a, z_p) - d(z_a, z_n) + m, 0\},$$

$$z_a = D_{table}(E_{chart}^{student}(A)),$$
where
$$z_p = D_{table}(E_{chart}^{teacher}(P)),$$

$$z_n = D_{table}(E_{chart}^{teacher}(N)).$$
(1)

where d denotes the distance defined as $d(z_i, z_j) = \|z_i - z_j\|_2$, and m denotes the margin which induces the anchor z_a to be closer to the positive

sample z_p and farther from the negative sample z_n . Through this process, the student chart encoder $E_{chart}^{student}$ is trained to extract only the essential information from original charts, even if they contain information that is irrelevant to chart reasoning. By optimizing Eq. 1, $E_{chart}^{student}$ focuses on subtle differences in the representations of positive and negative samples that would be highly similar in the representation space, which is expected to enhance the precise value mapping for table generation.

The representation z_a obtained after passing the original chart A through the student encoder $E_{chart}^{student}$ and the table decoder D_{table} is then fed into the frozen fully connected layer (FC) trained in Phase 1 to generate a table.

Furthermore, since our goal is not just to train $E_{chart}^{student}$ to mimic the representation of simple charts, but to eventually generate tables effectively, we also introduce a cross-entropy loss for table generation as follows:

$$T = [\hat{y}_1, \dots, \hat{y}_n] = FC(z_a), \tag{2}$$

$$\mathcal{L}_{table} = -\frac{1}{N} \sum_{i=1}^{n} \sum_{c=1}^{C} y_{i,c} \log \left(\frac{\exp(\hat{y}_{i,c})}{\sum_{j=1}^{C} \exp(\hat{y}_{i,j})} \right), \quad (3)$$

where T a linearized textual sequence of the generated table, n and C denote the length of the linearized textual sequence of the table T, and the number of classes, respectively, and y_i and \hat{y}_i denote the i-th ground-truth and the predicted token each belonging to one of the C classes, i.e., $y_{i,c} = 0$ if y_i belongs to class c and 0 otherwise. Note that, as in Deplot, lines of T are separated by '<0x0A>' (line break), and cells are distinguished by 'l'.

The final loss function for chart-to-table extraction is defined as follows:

$$\mathcal{L}_{final} = \lambda \mathcal{L}_{triplet} + (1 - \lambda) \mathcal{L}_{table}. \tag{4}$$

where λ is a hyperparameter balancing the two losses, i.e., $\mathcal{L}_{triplet}$ and \mathcal{L}_{table} . Note that considering the scale of the two losses, we set λ to 0.1 to balance between them. For detailed analysis of hyperparameter λ , please refer to Appendix G.

In summary, in Phase 2, we align the representation space of $E_{chart}^{student}$, which is trained with original charts, and $E_{chart}^{teacher}$, which is trained with simple charts, thereby allowing $E_{chart}^{student}$ to generate chart representations that mainly contain essential information for chart reasoning. Based on the representations obtained from $E_{chart}^{student}$, we use the frozen decoder D_{table} to generate tables that will be passed on to the LMM in the inference stage.

2.3 Inference Stage: Reasoning with Extracted Table

In the inference stage, we use a Large Multimodal Model (LMM) to perform various reasoning tasks given the tables obtained from the training stage as input. However, relying solely on tables without the associated visual information to solve chart reasoning tasks makes it impossible for the model to answer questions regarding visual information in charts (e.g., "What is the value of the third bar from the top?", "What year does the orange line represent?"). To address this issue, we additionally provide the original chart as another input to the LMM to enable the model to provide answers to a wider range of questions. By doing so, we expect the model to not only answer to questions regarding precise chart values by referring to the tables, but also to answer about visual attributes not contained in the table by referring to the original chart. Please refer to Appendix F for experimental results on the effects of using the table and original chart together.

Human-oriented chart instruction. Furthermore, we present a novel prompt that is specifically designed for chart reasoning, named Human-oriented chart instruction. Although prompt engineering is a widely studied topic for processing natural images (Kim et al., 2023; Li et al., 2022), it is crucial to develop specialized prompts specifically tailored for charts, due to the inherent nature of chart reasoning tasks, i.e., they require more advanced advanced thinking steps for interpretation. Our proposed prompt provides instructions that mimic how humans interpret charts, enabling more accurate reasoning within the LMM. To our best knowledge, this is the first work to study a prompt specifically targeting charts, and its effectiveness will be demonstrated in Section 3.1.2. We provide a snippet of our proposed prompt in Fig. 4. For detailed description of Human-oriented chart instruction, please refer to Appendix C.

3 Experiments

Dataset. In this paper, we evaluate SIMPLOT using widely used dataset for chart reasoning, ChartQA (Masry et al., 2022) comprising three types of charts (i.e., pie, bar, and line), and PlotQA (Methani et al., 2020) including three types of charts (i.e., dot line, line, bar). ChartQA also includes human-authored and LLM-augmented QA pairs, while PlotQA consists of QA pairs generated from templates created by manually analyzed sam-

Instruction for bar chart: 1. Firstly, bars of the same color represent the same column. Therefore, distinguishing colors and identifying corresponding columns is crucial (usually displayed around the main chart in the form of a legend). 2. Next, determine the location of rows. For vertical bar charts, rows are typically annotated at the bottom of the main chart, while for horizontal bar charts, they are annotated on the left or right side of the main chart. 3. Then, combine the colors of the nearest bars with annotated rows to determine which row and column the bars correspond to in the table. 4. ... Instruction for line chart: 1. In the case of a line chart, ...

Figure 4: A snippet of the proposed prompt.

ple questions collected through crowd-sourcing. Detailed explanation for dataset is presented in Appendix D.1.

Baselines. In this paper, we use various models as baselines that can perform chart reasoning. The baselines we employ are divided into three categories including Vision-Language pre-trained models (VLP), fully supervised models for chart reasoning, and models which extract table and perform reasoning using LLMs. Detailed descriptions of the models are provided in the Appendix D.2.

Evaluation protocol. For question answering based on charts, we report Relaxed Accuracy (RA) of 2,500 questions in the test set following previous works (Lee et al., 2023; Kantharaj et al., 2022b; Masry et al., 2022; Kantharaj et al., 2022a; Liu et al., 2022b).

Furthermore, we propose a proper metric to evaluate the chart extraction performance, employing a slightly modified metric derived from Relative Mapping Similarity (RMS), presented by Deplot (Liu et al., 2022a). We discover that the vanilla RMS does not accurately measure the chart-totable performance in certain cases, and thus we introduce a metric named Relative Distance (RD). For the differences between RMS and RD, please refer to the Appendix D.3.

3.1 Experimental Results

3.1.1 Results on Chart-to-Table Extraction

In Table 1, we observe that SIMPLOT exhibits superior chart-to-table extraction performance com-

Table 1: Chart-to-table extraction performance (RD_{F1}) on the ChartQA dataset over various chart types.

Models		Overall		
Models	Pie	Bar	Line	Overan
UniChart	84.86	92.58	85.16	88.03
Deplot	88.82	96.37	82.25	90.95
SIMPLOT	91.41	96.87	<u>84.74</u>	92.32

Table 2: Chart question answering performance (RA) on the ChartQA dataset.

	Models	Data type					
	Models	Human	Augmented	Overall			
	TaPas	28.72	53.84	41.28			
	V-TaPas	29.60	61.44	45.52			
lels	T5	25.12	56.96	41.04			
10 d	VL-T5	26.24	56.88	41.56			
VLP models	PaLI	30.40	64.90	47.65			
V	Mini-GPT	8.40	15.60	12.00			
	LLaVa	37.68	72.96	55.32			
	GPT-4V	56.48	63.04	59.76			
	ChartQA	40.08	63.60	51.84			
eq	ChartT5	31.80	74.40	53.10			
vis	Pix2Struct	30.50	81.60	56.05			
Supervised	MatCha	38.20	90.20	64.20			
\mathbf{Su}	Unichart	43.92	88.56	66.24			
	ChartLlama	48.96	90.36	69.66			
e	Deplot	62.71	78.63	70.67			
Table	Unichart ¹	<u>67.04</u>	69.92	68.48			
Ξ	SIMPLOT	78.07	88.42	83.24			

pared to existing methods including Deplot (Liu et al., 2022a) and UniChart (Masry et al., 2023). We argue that the superiority of SIMPLOT in terms of chart-to-table extraction eventually leads to more precise reasoning as will be shown in Section 3.1.2. While the difference in chart-to-table performance may seem small, we argue that even minor errors such as slightly incorrect values or column and row names during the chart-to-table extraction may entail completely incorrect answers, as shown in Table 2.

3.1.2 Results on Chart Reasoning

The performance of question answering on the ChartQA dataset is summarized in Table 2. 1) In general, we observe that methods designed with a focus on chart reasoning (i.e., 'Supervised' and 'Table') outperform VLP in QA tasks. This demonstrates that despite the superior performance of vision-language models in performing a wide range

¹The authors also provide pre-trained model for chart to table extraction. We use this model to extract the table and conduct reasoning in the same setting of Deplot and SIMPLOT.

of tasks in the traditional natural image domain, they face difficulties in interpreting charts due to the significantly different characteristics of natural images and charts. 2) Furthermore, noticeable performance difference is observed among models aimed at chart reasoning. Methods using tables such as Deplot, Unichart, and SIMPLOT, generally outperform models that conduct reasoning based solely on images (i.e., 'Supervised'). This suggests that going through tables allows for more detailed reasoning, aligning well with our motivation that emphasizes the importance of effectively extracting tables. Additionally, models utilizing extracted tables and LMM can efficiently handle various tasks such as text summarization without any additional data and training. 3) Among the table-based reasoning models, SIMPLOT performs the best by utilizing textual information, which is essential for obtaining the most accurate table, and by employing prompts mimicking the human reasoning process. Additional experimental results on PlotQA dataset and effectiveness of our method are presented in Appendix E and 3.2, respectively.

Discussion. The overall performance of SIM-PLOT is consistently higher than other baselines; however, on the augmented set, it sometimes exhibits slightly lower performance compared to some models. This can be attributed to the characteristic of the augmented set, which is obtained by fine-tuning the T5 model on the SQuAD QA dataset (Rajpurkar et al., 2016) to generate questions and answers. Consequently, the questions in the augmented set are mostly simpler and relatively easier to answer compared to those in the human set. Models that are fully fine-tuned on such questions tend to show notably higher performance on augmented data. For instance, MatCha (Liu et al., 2022b) shows a performance of 90.2 on the augmented set, but its performance drops significantly to 38.2 on the human set which is composed of more complex QA pairs, resulting in an overall performance of 64.2. In contrast, SIMPLOT's high performance even on the human set without additional fine-tuning suggests that proposed method is effective for more complex reasoning tasks.

3.2 Further Analysis

Ablation Study. We conduct ablation studies to determine the impact of each module of SIM-PLOT on the performance. Upper part of the Table 3 shows the effects of the two modules we introduced

Table 3: Ablation studies on each component of SIM-PLOT for table extraction and QA tasks.

Row-col rendering	Simple chart	Prompt	RD_{F1}	RA
×	Х	-	90.95	-
✓	×	-	91.40	-
X	✓	-	91.86	-
✓	✓	-	92.32	-
-	-	Х	-	79.79
-	-	✓	-	83.24

for chart-to-table extraction: 1) row-column rendering and 2) distillation from a simple chart. These results demonstrate that each module enhances the performance of table extraction, thereby proving the effectiveness of our design.

Additionally, to verify the effect of the proposed prompt, we also present the QA performance with and without the *Human-oriented chart instruction* (lower part of Table 3). We observe that our prompt that mimics how humans think enhances the performance of the QA task, which confirms the importance of using appropriate prompts tailored to the task when employing LMMs. It is notable that providing more diverse in-context examples could lead to even greater performance gains.

SIMPLOT is Model-agnostic. While SIMPLOT is originally designed to address the limitations of Deplot, it is model-agnostic, i.e., SIMPLOT can be applied to any existing table-based reasoning model to enhance its performance. Here, we applied SIMPLOT to one of our baselines, Unichart (Masry et al., 2023), which is pre-trained to perform chart-to-table extraction task. Specifically, Fig. 5 (a) demonstrates significant performance enhancements with SIMPLOT applied to Deplot and Unichart, respectively. Consequently, Fig. 5 (b) show a substantial increase in question answering performance. These results verifies the generality and practicality of SIMPLOT.

Broad Applicability of the Proposed Prompt. To verify the effectiveness of our proposed prompt, we apply it to Deplot with images (upper part of Table 5). We observe that applying the proposed prompt to the 'Deplot+img' indeed greatly enhances its performance.

Using Images as Input to Unichart/Deplot for Inference. Here, we further demonstrate that using both tables and images is essential for chart reasoning, as emphasized earlier. Table 4 shows that the performance of the vanilla Deplot and the vanilla Unichart, both of which only use tables

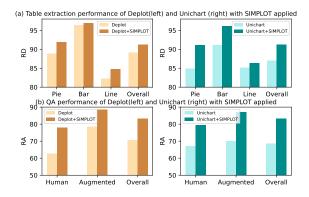


Figure 5: Improvements in performance observed when SIMPLOT is applied across diverse models.

Table 4: Chart question answering performance (RA) with/without image on ChartQA dataset.

Models	Human	Augmented	Overall
Unichart Unichart + img. Unichart + SIMPLOT w/o prompt Unichart + SIMPLOT	67.04	69.92	68.48
	75.04	88.82	81.93
	<u>76.56</u>	<u>88.64</u>	<u>82.60</u>
	79.56	87.18	83.37
Deplot Deplot + img. Deplot + SIMPLOT w/o prompt Deplot + SIMPLOT	62.71	78.63	70.67
	72.39	85.01	78.70
	73.91	<u>85.67</u>	<u>79.79</u>
	76.70	88.42	82.56

for reasoning, is improved when images are also utilized during inference. This indicates a clear drawback of relying solely on tables and proves the necessity of using images with specifically designed prompts. Moreover, even without the prompt, SIMPLOT outperforms the baselines with images used during inference, thanks to the superior table extraction capability of SIMPLOT.

Performance on more Challenging Questions.

It is worth noting that 'Deplot+img+prompt' now performs competitively to SIMPLOT. However, we argue that this is mainly due to the simplicity of the ChartQA dataset, which makes the dataset insufficient to evaluate the table extraction capability of SIMPLOT. For example, ChartQA contains questions that can be answered by referencing just one row of the extracted table.

To evaluate the models with more challenging questions, we randomly sample 100 test images and then provide GPT-4 with the following instruction: "Create a challenging question-answer pair that requires referencing at least two rows and two columns to solve.". Through this process, we generated 100 question-answer pairs and eventually obtained 85 pairs after manually filtering out inaccurate pairs, and then evaluate 'Deplot+img+prompt' and SIMPLOT on these question-answer pairs (lower part of Table 5). We observe

Table 5: Effectiveness of proposed prompt on ChartQA.

	Models	Human	Augmented	Overall
<u></u>	Deplot + img.	72.39	85.01	78.70
Easy	Deplot + img. + prompt	77.75	88.30	83.03
1	SIMPLOT	78.07	88.42	83.24
rd	Deplot + img. + prompt	-	-	49.41
Har	SIMPLOT	-	-	65.88

that the performance gap greatly increases as the questions get more challenging, which demonstrates the importance of extracting accurate tables when encountering challenging questions. For the specific example of the challenging question and how the extracted tables from SIMPLOT and Deplot lead to differences in complex QA performance, please refer to Appendix I.

Case Studies. We provide case studies showing how SIMPLOT address the concerns we have raised, compared to our main baseline, Deplot (Liu et al., 2022a).

Firstly, we verify whether SIMPLOT correctly captures textual information within charts. Fig. 7 illustrates that the table extracted by Deplot fails to consider the counting unit, i.e., "million/billion," present in the chart, while SIMPLOT correctly captures the textual information.

We also confirm that rendering columns and rows of a table onto the image is beneficial for extracting an accurate table. In Fig. 8, Deplot (Liu et al., 2022a) omits a relevant row (i.e., "Montenegro"), redundantly extracts several rows, and even generates entirely irrelevant rows such as "CC BY," which would result in a poor performance in reasoning tasks. In contrast, SIMPLOT accurately generates a table by extracting only the relevant rows, ensuring accurate reasoning.

Furthermore, we demonstrate the effectiveness of *Human-oriented chart instruction*. Fig. 6 illustrates the difference between interpreting the same chart with/without using the proposed prompt. While the model fails to produce an accurate answer without the prompt, the model generates a correct answer in a process similar to how humans interpret charts when using the prompt. This highlights the effectiveness of our prompt designed to mimic human reasoning. Please refer to Appendix H for further case studies and Appendix J for error analysis.

4 Related Work

Early chart processing methods relied on rulebased approaches, which lacked flexibility for var02 04 06 08 10 12 CENTER 2012 News Consumption Survey

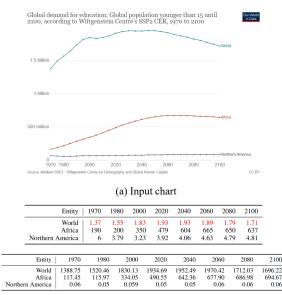




Explanation: The percentage of people who regularly watch CNN in the year '02 is found in the chart under the CNN column at the point where the year '02 is marked. The chart shows that the value is 21%.

Explanation: To answer the question, we must refer to the line chart. The x-axis represents the years, which correspond to the rows in the table. Each line represents a different news network, which corresponds to the columns in the table. The value for CNN in a specific year can be found by locating the point on the CNN line (red) that aligns with the year on the x-axis and then reading the value from the y-axis The chart shows the percentage of people who regularly watch CNN, Fox News, and MSNBC over various years. The red line represents CNN. In 2002, the red line is at the 25% mark on the y-axis.

Figure 6: Case study. Interpreting charts similar to how humans do using *Human-oriented chart instruction*.

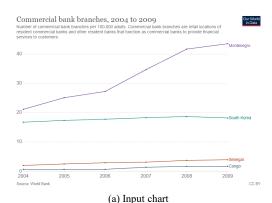


(b) Extracted table from Deplot (upper) & SIMPLOT (lower)

Figure 7: SIMPLOT captures contextual information (billion/million) and extract precise values in the chart.

ious chart types (Balaji et al., 2018; Choi et al., 2019). OCR and key-point detection methods are utilized to address these limitations, but they depend on data annotation and module performance, making them time-consuming (Poco and Heer, 2017; Zhou et al., 2023; Xue et al., 2023).

Recent papers train end-to-end vision-language models to improve chart reasoning without heuristic rules, but they need fine-tuning for specific tasks (Meng et al., 2024; Huang et al., 2023; Masry et al., 2023; Han et al., 2023), so recent approaches transform charts into tables and utilized LLM for question answering while enhancing interpretability and allowing the model apply to various downstream task without additional training (Liu et al., 2022a; Xia et al., 2023). For complete related works, please refer to the Appendix A.



Entity	Com	mercial	bank bi	anches,	2004 to	2009
South Korea						16.86
Congo						4.16
Senegal						3.12
Congo						3.7
Senegal						4.32
Como	1 .					4.0
CC BY	+1	missing	row			4.19
Entity	2004	2005	2006	2007	2008	2009
Montenegro	21.4	25.5	27.0	34.5	42.1	43.2
South Korea	16.8	17.0	17.9	18.3	18.7	18.2
Senegal	2.0	2.1	2.2	2.3	2.7	2.8
Congo	0.3	0.3	0.3	0.3	0.3	0.3

(b) Extracted table from Deplot (upper) & SIMPLOT (lower)

Figure 8: Case study. Extract precise rows without unnecessary information in the chart.

Conclusion

We propose SIMPLOT, which extracts only the essential information required for chart reasoning. Furthermore, by leveraging textual information inside the charts, SIMPLOT enables accurate chart reasoning. Finally, we propose a novel prompt specifically designed for chart reasoning. It provides instruction to LMM imitating how humans interpret chart, enabling more precise reasoning. Through extensive experiments, we demonstrate that SIMPLOT effectively handle concerns raised by existing works while improving performance, and also can be applied to other existing table-based reasoning model.

Limitations

Despite the outstanding performance of SIM-PLOT in chart-to-table extraction, there is room for improvement. During the process of rendering columns and rows onto the images, inaccuracies in extracting columns and rows can lead to failures in effectively converting charts into tables. Particularly, it can be challenging to extract appropriate columns and rows when the texts in the chart are complex. We provide a detailed error analysis related to this issue in the Appendix J. To address these problems, future work could involve verifying and ensuring that the extracted rows and columns maintain a coherent context, aiming for more accurate rendering and extraction.

Ethics Statement

SIMPLOT has not been trained with private or sensitive data. This significantly reduces the risk of generating harmful or misleading content. ChartQA and PlotQA dataset, which we utilize to collect chart-table pairs, is publicly accessible for research purposes. To make sure the transparency and reproducibility of our experiments, we provide detailed information on hyperparameter configurations in our paper and publicly share our source code. This careful approach mitigates the potential for unethical outcomes associated with data usage. While our models demonstrate state-of-the-art performance on the both dataset, we acknowledge the possibility of their misuse. There exists a risk that our models could be exploited to mislead the public about the content and implications of charts. Despite our models' high performance, we cannot guarantee their outputs will always be accurate, emphasizing the need for critical interpretation and verification of results.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Abhijit Balaji, Thuvaarakkesh Ramanathan, and Venkateshwarlu Sonathi. 2018. Chart-text: A fully automated chart image descriptor. *arXiv preprint arXiv:1812.10636*.
- Ali Furkan Biten, Ruben Tito, Andres Mafla, Lluis Gomez, Marçal Rusinol, Minesh Mathew, CV Jawahar, Ernest Valveny, and Dimosthenis Karatzas. 2019. Icdar 2019 competition on scene text visual question answering. In 2019 International Conference on Document Analysis and Recognition (ICDAR), pages 1563–1570. IEEE.
- Hermann Bujard, Reiner Gentz, Michael Lanzer, Dietrich Stueber, Michael Mueller, Ibrahim Ibrahimi, Marie-Therese Haeuptle, and Bernhard Dobberstein. 1987. [26] a t5 promoter-based transcription-translation system for the analysis of proteins in vitro and in vivo. In *Methods in enzymology*, volume 155, pages 416–433. Elsevier.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. 2022a. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*.
- Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. 2022b. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer.
- Zhi-Qi Cheng, Qi Dai, Siyao Li, Jingdong Sun, Teruko Mitamura, and Alexander G Hauptmann. 2023. Chartreader: A unified framework for chart derendering and comprehension without heuristic rules. *arXiv* preprint arXiv:2304.02173.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See https://vicuna. lmsys. org (accessed 14 April 2023).
- Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. 2021. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*, pages 1931–1942. PMLR.
- Jinho Choi, Sanghun Jung, Deok Gun Park, Jaegul Choo, and Niklas Elmqvist. 2019. Visualizing for

- the non-visual: Enabling the visually impaired to use visualization. In *Computer Graphics Forum*, volume 38, pages 249–260. Wiley Online Library.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- François Gardères, Maryam Ziaeefard, Baptiste Abeloos, and Freddy Lecue. 2020. Conceptbert: Concept-aware representation for visual question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 489–498.
- 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.
- Johan Holmgren, Paul Davidsson, Jan A Persson, and Linda Ramstedt. 2012. Tapas: A multi-agent-based model for simulation of transport chains. *Simulation Modelling Practice and Theory*, 23:1–18.
- Kung-Hsiang Huang, Mingyang Zhou, Hou Pong Chan, Yi R Fung, Zhenhailong Wang, Lingyu Zhang, Shih-Fu Chang, and Heng Ji. 2023. Do lvlms understand charts? analyzing and correcting factual errors in chart captioning. *arXiv* preprint arXiv:2312.10160.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR.
- Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656.
- Shankar Kantharaj, Xuan Long Do, Rixie Tiffany Ko Leong, Jia Qing Tan, Enamul Hoque, and Shafiq Joty. 2022a. Opencqa: Open-ended question answering with charts. *arXiv preprint arXiv:2210.06628*.
- Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. 2022b. Chart-to-text: A large-scale benchmark for chart summarization. *arXiv preprint arXiv:2203.06486*.
- Kibum Kim, Kanghoon Yoon, Yeonjun In, Jinyoung Moon, Donghyun Kim, and Chanyoung Park. 2024. Adaptive self-training framework for finegrained scene graph generation. *arXiv preprint arXiv:2401.09786*.

- Kibum Kim, Kanghoon Yoon, Jaehyeong Jeon, Yeonjun In, Jinyoung Moon, Donghyun Kim, and Chanyoung Park. 2023. Llm4sgg: Large language model for weakly supervised scene graph generation. *arXiv e-prints*, pages arXiv–2310.
- 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.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR.
- Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 121–137. 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. 2022b. Matcha: Enhancing visual language pretraining with math reasoning and chart derendering. *arXiv* preprint arXiv:2212.09662.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. Advances in neural information processing systems, 32.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023a. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. *arXiv e-prints*, pages arXiv–2310.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023b. Chameleon: Plug-and-play compositional reasoning with large language models. arXiv preprint arXiv:2304.09842.
- Junyu Luo, Zekun Li, Jinpeng Wang, and Chin-Yew Lin. 2021. Chartocr: Data extraction from charts

- images via a deep hybrid framework. In *Proceedings* of the IEEE/CVF winter conference on applications of computer vision, pages 1917–1925.
- 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 Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.
- 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.
- Jorge Poco and Jeffrey Heer. 2017. Reverse-engineering visualizations: Recovering visual encodings from chart images. In *Computer graphics forum*, volume 36, pages 353–363. Wiley Online Library.
- 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.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. 2023. Prompting large language models with answer heuristics for knowledge-based visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14974–14983.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. Vl-bert: Pre-training of generic visual-linguistic representations. *arXiv* preprint arXiv:1908.08530.
- Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv* preprint arXiv:1908.07490.
- 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.

- Renqiu Xia, Bo Zhang, Haoyang Peng, Ning Liao, Peng Ye, Botian Shi, Junchi Yan, and Yu Qiao. 2023. Structchart: Perception, structuring, reasoning for visual chart understanding. *arXiv preprint arXiv:2309.11268*.
- Wenyuan Xue, Dapeng Chen, Baosheng Yu, Yifei Chen, Sai Zhou, and Wei Peng. 2023. Chartdetr: A multishape detection network for visual chart recognition. *arXiv preprint arXiv:2308.07743*.
- Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5579–5588.
- Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13041–13049.
- Mingyang Zhou, Yi R Fung, Long Chen, Christopher Thomas, Heng Ji, and Shih-Fu Chang. 2023. Enhanced chart understanding in vision and language task via cross-modal pre-training on plot table pairs. arXiv preprint arXiv:2305.18641.
- 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*.

A Complete Related Work

Pre-trained vision-language models. Vision-language pre-training seeks to enhance the performance of downstream vision and language tasks such as visual question answering (VQA) (Tan and Bansal, 2019; Alayrac et al., 2022), image-text retrieval (Chen et al., 2020; Jia et al., 2021), and image captioning (Li et al., 2020, 2022).

Several models learn cross-modal representation from visual features and language tokens (Chen et al., 2020; Su et al., 2019). Specifically, Lu et al. (2019) employs two single-modal networks, alongside a cross-modal transformer layer integrating information from input sentences and images.

Radford et al. (2021) and Jia et al. (2021) facilitated zero-shot adaptation utilizing contrastive learning for various tasks on multimodal data. However, their drawback resides in their limitation to generate language, diminishing their suitability for open-ended tasks such as VQA.

Others propose a different approach to develop an enhanced Large Multimodal Model (LMM). For example, Liu et al. (2023) employs a visual encoder and a language decoder from other studies (Radford et al., 2021; Chiang et al., 2023). It is capable of performing various QA tasks including conversation, detail description, complex reasoning and focuses on chat capability.

Chart reasoning models. Early methods for processing charts predominantly utilized rule-based approaches (Balaji et al., 2018; Choi et al., 2019). However, they were specifically tailored to predefined chart formats, limiting their applicability only to certain types of charts. Moreover, incorporating new chart designs require additional rules, leading to the impracticality of immediate adaptation to new chart formats.

Due to above limitations, researchers utilize modules such as OCR and key-point detection (Poco and Heer, 2017; Zhou et al., 2023; Xue et al., 2023). However, these methods can be time-consuming and require dataset annotation for labeling, and rely on the performance of modules. For example, Luo et al. (2021) suffers from a drawback as it solely predicts the raw data values without establishing connections to their respective axes or legends.

In an effort to address these challenges, researchers develop vision-language models (Meng et al., 2024; Huang et al., 2023) without heuristic rules. Some approaches comprehend the chart and

respond to questions in natural language (Masry et al., 2023; Han et al., 2023). However, such models require fine-tuning for each downstream task, constraining their adaptability for diverse tasks. Meanwhile, studies such as Methani et al. (2020); Lu et al. (2023a) have created datasets to enable more realistic mathematical reasoning.

To tackle these challenges, recent studies Liu et al. (2022a); Xia et al. (2023) introduce a two-step approach, chart-to-table extraction and reasoning with LLM. Recent advancements like Liu et al. (2022a); Xia et al. (2023) have introduced a two-step approach to tackle challenges in chart analysis. Firstly, charts are transformed into tables, improving interpretability and aiding in identifying inaccuracies. Lu et al. (2023b) utilizes Program of Thought (PoT), achieving significant performance gains. This approach enables more accurate and interpretable reasoning compared to traditional methods, indicating that table extraction is crucial.

B Row-column Rendering

Fig. 9 presents a detailed example of our row-column rendering process. For a simple chart and an original chart used in the training stage, we extract rows and columns from their associated table and directly render them onto the charts since the ground-truth table is provided in this stage. However, for the original chart used in the inference stage, since we cannot access the ground-truth table, we use an LMM to extract rows and columns from the input chart and rendered onto the image.

C Human-oriented Prompt for Chart Reasoning

We present instructions in the prompt to mimic the way humans commonly perceive charts. We name this approach *Human-oriented chart instruction*, which is a simple yet effective methodology as shown in the QA results. The detailed prompt is shown in Fig. 10.

In the *Human-oriented chart instruction*, universal instructions applicable to all chart types are provided first. They include components contained in every chart such as chart title, y-axis, x-axis, legend, and data labels, facilitating the LMM to comprehend the chart. Furthermore, we guide the alignment between a chart and the generated table by instructing the system within the prompt to effectively comprehend the positions of columns and rows as well as the outline of the chart, referring to

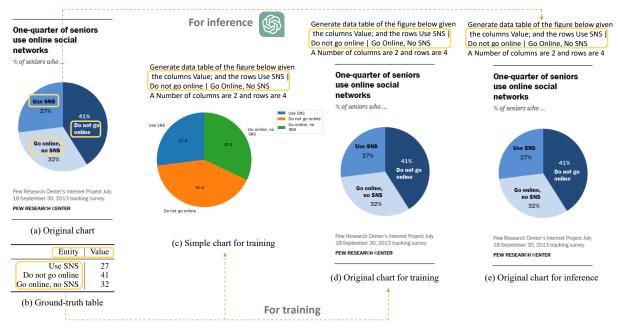


Figure 9: An example of row-column rendering for an input image in training and inference stage.

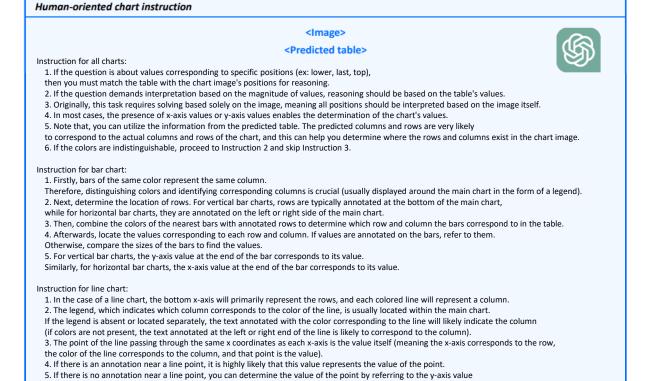


Figure 10: Overall prompt of *Human-oriented chart instruction*.

6. In a line chart, it is crucial to understand the flow of the line. Lines can show trends of decreasing, increasing, or remaining constant, and when multiple lines intersect, it is important to identify which line corresponds to which column based on their colors.

2. Each section has one color, and the row it corresponds to is likely indicated by text either inside the section or close to it (if not nearby, it can be identified through the legend or connected to the corresponding text by lines or markers).

corresponding to the y coordinate of the point.

1. In a pie chart, it is very important to determine which color corresponds to which row.

3. In the case of a pie chart, the values are usually annotated on each section of the pie chart.

Instruction for pie chart:

the table generated by SIMPLOT.

Subsequently, instructions are described for each chart type. Here, we describe the instructions designed for bar charts, but those for line charts and pie charts are designed in a similar manner.

Bar Chart: Instruction 1. When humans recognize a bar chart, the initial observation would be columns distinguished by their colors. Following this observation, the legend within the image is consulted to identify which color denotes which column. If the colors are indistinguishable, proceed to Instruction 2 and skip Instruction 3.

Bar Chart: Instruction 2. Having identified the columns, humans would proceed to observe which row each bar corresponds to. Specifically, for vertical bar graphs, rows are usually annotated below the main chart, whereas for horizontal bar graphs, annotations on the lateral sides determine the rows.

Bar Chart: Instruction 3. Humans then would match the color of the bar with the annotated row, thereby aligning each bar with its corresponding column and row.

Bar Chart: Instruction 4 & 5. Lastly, humans would figure out the value corresponding to each row and column by referencing the value annotated on the bar, the relative sizes between bars, and the x-axis and y-axis. Notably, it is understood intuitively that the value at the point where the bar terminates on the x-axis or y-axis represents the value of that bar.

We argue that developing a language model that captures the sequential and unconscious reasoning process involved in human chart recognition through *Human-oriented chart instruction* not only mimics human cognition effectively but also ensures the alignment of predicted tables with chart images. This leads to substantial performance enhancements at a lower cost.

D Details regarding Experiments

D.1 Dataset Description

In this paper, we mainly evaluate our proposed method with a widely-used dataset, ChartQA (Masry et al., 2022). It consists of real-world charts containing complex information, which aligns well with our motivation compared to other synthetic datasets such as PlotQA (Methani et al., 2020) and DVQA (Kafle et al., 2018), both of which only include simple charts. Furthermore, ChartQA is comprised of three chart types (i.e., pie, bar, and

line), whereas DVQA and PlotQA comprise only one or two types of chart. Therefore, ChartQA is a more challenging dataset and is desirable for evaluating performance on real-world chart reasoning tasks.

For these reasons, we mainly focus on ChartQA dataset in this paper since other datasets are composed only of very simple charts, which are not suitable for our model that is designed to handle real-world charts containing noisy information unnecessary for reasoning. Nevertheless, to demonstrate that SIMPLOT can be applied to any dataset, we conduct additional experiments on the PlotQA dataset in the Appendix E.

ChartQA also includes human-authored and LLM-augmented QA pairs for reasoning. Please note that we do not use QA pairs for training, but train the model to generate tables only from the charts in the training set, then use QA pairs from the test set for measuring reasoning performance.

Detailed statistics of ChartQA dataset are presented in Table 6, and statistics of PlotQA dataset are presented in Table 7. To reduce the training cost and balance the dataset with ChartQA, we stratified and sampled 10% of the images from the PlotQA dataset based on type, which were then used for training and inference. Additionally, for QA pairs, we used one pair per image.

Table 6: Statistics of ChartQA dataset.

type	Pie	Bar	Line	QA pair
Train set	541	15,581	2,195	-
Validation set	48	837	171	-
Test set	78	1,230	211	2,500

Table 7: Statistics of PlotQA dataset.

type	Dot line	Line	Bar	QA pair
Train set	26,010	25,897	105,163	-
Validation set	5,571	5,547	22,541	-
Test set	5,574	5,549	22,534	4,342,514

D.2 Compared Methods

We use a wide range of models capable of performing chart reasoning as baselines, categorizing them into three distinct categories.

The first category includes Vision-Language Pretrained models (VLP) as follows:

• TaPas (Holmgren et al., 2012)

- V-TaPas (Masry et al., 2022)
- T5 (Bujard et al., 1987)
- VL-T5 (Cho et al., 2021)
- PaLI (Chen et al., 2022b)
- Mini-GPT (Zhu et al., 2023)
- LLaVa (Liu et al., 2023)
- GPT-4V (Achiam et al., 2023)

The second category consists of supervised models, including followings:

- ChartQA (Masry et al., 2022)
- ChartT5 (Zhou et al., 2023)
- Pix2Struct (Lee et al., 2023)
- MatCha (Liu et al., 2022b)
- Unichart (Masry et al., 2023)
- ChartLlama (Han et al., 2023)

The third category consists of table-based reasoning models, including Deplot (Liu et al., 2022a) and Unichart (Masry et al., 2023). Deplot converts charts to tables similar to SIMPLOT and employs prompt-based methodologies such as Chain-of-Thought (CoT) (Wei et al., 2022), Program-of-Thought (PoT) (Chen et al., 2022a) alongside using an LLM for question reasoning. Although Unichart is trained to perform QA directly from charts, it is also pre-trained to generate tables. In this paper, we use this variant of Unichart to extract tables and apply it in the same setting as Deplot.

D.3 Details regarding Proposed Metric

As illustrated in Fig. 11, columns of extracted table from SIMPLOT (e.g., "Percentage (%)" in the figure) often carry the same meaning as the ground truth (e.g., "Value" in the figure) but have different Levenshtein distances. Additionally, as shown in Fig.12, in cases where columns are not explicitly stated in the image, an SIMPLOT utilizing LMM without further fine-tuning generates arbitrary columns. We confidently assert that there are unequivocally no problems associated with chart analysis when employing an LMM in both scenarios. However, given that RMS is composed of the product of Levenshtein distance (from columns and rows) and relative distance (from values), this leads

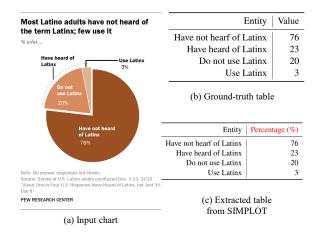


Figure 11: An example of an extracted table that poses no issues for reasoning despite having a different column name from that in the ground-truth table.

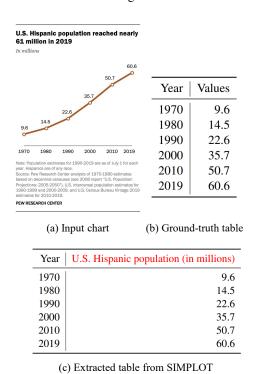


Figure 12: An example where reasoning remains unaffected, yet a column extracted by SIMPLOT differs significantly due to the lack of clear column description in the image.

to an apparent degradation in performance of chart-to-table task even when there are no issues with QA. Consequently, we construct a Levenshtein distance matrix by concatenating column and row predictions with the ground truth, followed by applying minimal cost matching to consider only the relative distance of each pair. We adopt this Relative Distance (RD) as our main metric. This approach allows us to accurately assess performance in value mapping.

We follow Deplot (Liu et al., 2022a) to define RD as in the following equation, with most formulas extracted from this paper. Following Liu et al. (2022a), we use basic concept of Related Mapping Similarity (RMS), regarding tables as unordered sets of mappings from row and column headers (r, c) to a single value v. $p_i = (p_i^r, p_i^c, p_i^v)$ and $t_j=(t_j^r,t_j^c,t_j^v)$ indicates each entity for the predicted table $\mathcal{P}=\{p_i\}_{1\leq i\leq N}$ and the ground truth table $\mathcal{T} = \{t_i\}_{1 \leq i \leq M}$, respectively. Utilizing Normalized Levenshtein Distance (NL_{τ}) (Biten et al., 2019), we denote the distance between two keys p_i and t_j as $NL_{\tau}(p^r||p^c, t^r||t^c)$ where || indicates concatenation of string. The distance between values of table is computed with relative distance $D_{\theta}(p^{v}, t^{v}) = min(1, ||p^{v} - t^{v}||/||t^{v}||)$ and distances larger than θ are set to the maximum of 1. Incorporating these two distances, we obtain the similarity between two entities in a mapping $D_{\tau,\theta}(p,t)$ as $(1 - NL_{\tau}(p^r || p^c, t^r || t^c))(1 - D_{\theta}(p^v, t^v))$.

For RMS, Liu et al. (2022a) computes the pairwise similarity in \mathcal{P} and \mathcal{T} using the cost function $(1-NL_{\tau}(p^r||p^c,t^r||t^c))$, obtaining a similarity matrix in shape of $N\times M$ and identifying the minimal cost matching $\mathbf{X}\in\mathbb{R}^{N\times M}$ between the keys (in the binary matrix form).

While precision and recall is obtained from two mappings in Deplot (Liu et al., 2022a), we only use distance of the numeric entries to compute RD:

$$RD_{precision} = 1 - \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} \mathbf{X}_{ij}(D_{\theta}(p^{v}, t^{v}))}{N}, \quad (5)$$

$$RD_{recall} = 1 - \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} \mathbf{X}_{ij} (D_{\theta}(p^{v}, t^{v}))}{M}.$$
 (6)

Similar to Deplot, we compute RD_{F1} as the harmonic mean of the precision and recall.

E Additional Experiments with PlotQA

To demonstrate that our method can be applied to a variety of data beyond the ChartQA dataset, we present additional experimental results on the PlotQA dataset in Table 8. Our proposed method exhibits superior performance across various chart types, indicating that it can be applied to any chart type. However, unlike ChartQA, which consists of real-world charts, PlotQA is composed only of very simple charts, so the performance differences between SIMPLOT and other methods are not as significant as in the result on the ChartQA dataset. We

emphasize once again that SIMPLOT is designed to handle real-world charts containing unnecessary and noisy information for chart question answering.

Table 8: Chart question answering performance (RA) on the PlotQA dataset.

Models	Dot line	Line	Bar	Overall
GPT-4V	50.53	58.84	53.85	54.11
Unichart	58.78	53.26	60.10	58.74
Deplot	66.66	55.59	61.73	61.53
SIMPLOT	60.93	65.57	73.84	70.32

F Effectiveness of Utilizing Image with Extracted Table

Using tables for chart reasoning allows for performing various downstream tasks (such as chart summarization, chartQA, etc.). However, there is an inherent limitation of using tables for chart reasoning, i.e., incorrect extraction of the table leads to failure in all subsequent downstream tasks. Taking this limitation into account, unlike Deplot that solely relies on tables for reasoning, SIMPLOT simultaneously considers both the extracted table and the image itself to compensate for the incorrectly extracted table.

Table 9: The number of times QA succeeded or failed when both methods failed to extract the table properly.

SIMPLOT	Success	Fail
Success	-	7
Fail	48	-

To further verify the effectiveness of using an image alongside the generated table, we conduct additional experiments. Table 9 proves that using images along with tables increases the success rate of QA compared to using only tables. We assume both Deplot and SIMPLOT failed to accurately extract the table in cases where the Relative Distance (RD) was below 0.7 (i.e., when the SIMPLOT does not benefit from utilizing simple charts) and the difference between the two RDs was within 0.1, and there are 155 such cases. Among these cases, even though we provide images to Deplot and removed the prompt from SIMPLOT for a fair comparison, while Deplot failed in QA, SIMPLOT succeeded in 48 cases, which is more than 30% additional correct answers. Conversely, when SIMPLOT failed in QA, Deplot succeeded in only 7 cases, which is

less than 5%. This demonstrates that utilizing both plot and table can effectively address the limitation of solely relying on tables for reasoning.

G Hyperparameter Analysis

We conduct a hyperparameter analysis by setting various values for the balancing parameter λ and analyzing the result changes in RD. As illustrated in Fig. 13, the changes in RD values due to variations in λ are minimal. This indicated that our model is robust and performs consistently well regardless of the λ value.

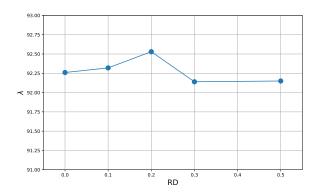


Figure 13: Table extraction performance (RD_{F1}) with varying λ .

H More Case Studies

Here, we present more case studies.

Fig. 14 and 15 show that in both pie charts and bar charts, Deplot (Liu et al., 2022a) fails to capture relevant rows, generates inaccurate rows, thus failing to create values. On the other hand, SIMPLOT consistently produces accurate tables. This directly proves the generality of SIMPLOT across various chart types.

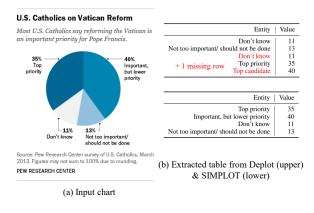


Figure 14: Case study. Extract all necessary rows in the pie chart.

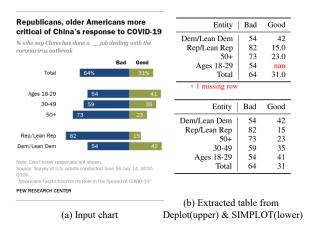


Figure 15: Case study. Extract all necessary rows and precise value in the bar chart.

Furthermore, Fig. 17 demonstrates that our *Human-oriented chart instruction* works well by imitating the way of human interpretation for charts. Without the proposed prompt, an LLM fails to detect the exact location of the line component "Favor", leading to a wrong answer. However, utilizing *Human-oriented chart instruction*, an LMM can precisely locate the necessary component, resulting in a correct answer.

I Evaluations on more Challenging Questions

While conducting case studies, we observed that Deplot and Simplot often exhibited similar QA performance, despite Deplot failing to accurately extract tables that SIMPLOT successfully extract. We attribute this mainly to the simplicity of the questions in the ChartQA dataset. For example, in Fig. 15, although Deplot fails to accurately extract the table, if the question happens to be about 'Dem/Lean Dem', which is not related to the inaccurate part of the extracted table, Deplot can still generate the correct answer. This proves that the performance of our method can be further differentiated with more complex questions.

Hence, we generated more challenging questions, and evaluated models in Section 3.2. Fig. 16 shows an example of challenging question, where Deplot fails to answer correctly due to its poorly extracted table, while SIMPLOT generates a precise answer based on its accurately extracted table. The question is "What is the sum of the birth rate in China in 1955 and the death rate in China in 1965?", which requires more than one row and column to answer. In this example, we verify that

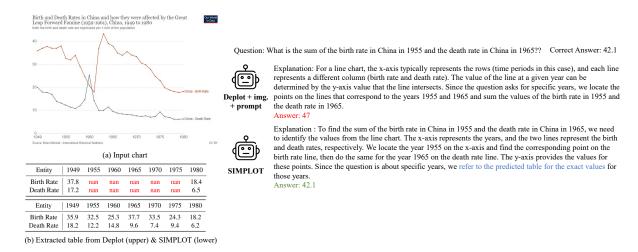


Figure 16: Case study. Effectiveness of precise extracted table for challenging question.

extracting accurate tables is crucial for solving challenging questions, and that the ChartQA dataset is insufficient to precisely evaluate the reasoning models.

J Error Analysis

Here we present examples of SIMPLOT failures, their reasons, and possible solutions. In Fig. 18 (c), it can be observed that the extracted table from the proposed method contains inaccurate rows. This is attributed to a line break between the items "Exposure to false or incorrect information" and "Losing the ability to Communicate face-to-face" in the input chart, leading to the misinterpretation of each item as two separate entities during row extraction. To address such issues, one possible solution would be to refine the rows and columns by considering their contextual information before rendering.

In addition, we find some examples where the LMM lacks specific abilities, such as color detection. In Fig. 19, there is a question asking about the color of "Jamaica" with the correct answer being orange, but the model recognizes it as red and answers incorrectly. One possible solution would be to broaden the range of similar colors for the model to answer together.

K Responsible NLP Research Checklist

Hardware.

We train SIMPLOT using a NVIDIA RTX A6000. The training time for phase 1 (Section 2.2.1) and phase 2 (Section 2.2.2) are about 2 hours and 4 hours per epoch, respectively. We train phase 1 for 7 epochs and phase 2 for 9 epochs.

Parameters.

SIMPLOT has 374M parameters, and the exact number of parameters in GPT-4 has not been publicly reported. For LMM prompting, we use temperature of 0.1.

(a) Input chart

Question: Which group favor or oppose has the highest value in the year Sept 2008?

Correct Answer: Favor



Explanation: According to the table, in Sept 2008, the value for 'Favor' is 0 and the value for 'Oppose' is 28. Therefore, the 'Oppose' group has the highest value in the year Sept 2008.

Answer: Oppose



Explanation: For a line chart, the x-axis typically represents the rows (time periods in this case), and each line represents a different category or column. The points on the lines correspond to the values for each category at a given time. To find the highest value for a specific time period, you need to look at the points on the chart that align with that time period on the x-axis and compare the y-axis values of those points. Based on the line chart, in Sept 2008, the point on the 'Favor' line is higher than the point on the 'Oppose' line. The 'Favor' line is at approximately 67, and the 'Oppose' line is at 28, as indicated by the annotations on the chart

Answer: Favor

Figure 17: Case study. Interpreting charts similar to how humans do using *Human-oriented chart instruction*.

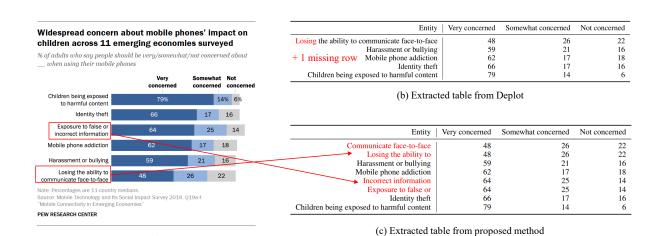


Figure 18: Error Analysis. Failure in row / column rendering.

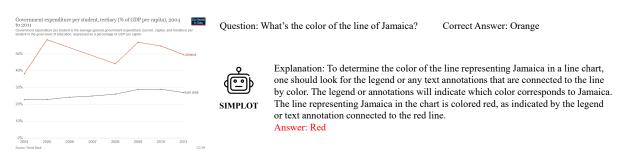


Figure 19: Error Analysis. Failure in detecting a color in the chart.