

# DCQA: DOCUMENT-LEVEL CHART QUESTION ANSWERING TOWARDS COMPLEX REASONING AND COMMON-SENSE UNDERSTANDING

Anran Wu<sup>\*</sup>, Luwei Xiao<sup>\*</sup>, Xingjiao Wu<sup>†</sup>, Shuwen Yang<sup>\*</sup>, Junjie Xu<sup>\*</sup>, Zisong Zhuang<sup>\*</sup>, Nian Xie<sup>‡</sup>, Cheng Jin<sup>†</sup>, Liang He<sup>\*</sup>

<sup>\*</sup> East China Normal University, Shanghai, China

<sup>†</sup> Fudan University, Shanghai, China

<sup>‡</sup> Huawei Noah’s Ark Lab, Shenzhen, China

## ABSTRACT

Visually-situated languages such as charts and plots are omnipresent in real-world documents. These graphical depictions are human-readable and are often analyzed in visually-rich documents to address a variety of questions that necessitate complex reasoning and common-sense responses. Despite the growing number of datasets that aim to answer questions over charts, most only address this task in isolation, without considering the broader context of document-level question answering. Moreover, such datasets lack adequate common-sense reasoning information in their questions. In this work, we introduce a novel task named document-level chart question answering (DCQA). The goal of this task is to conduct document-level question answering, extracting charts or plots in the document via document layout analysis (DLA) first and subsequently performing chart question answering (CQA). The newly developed benchmark dataset comprises 50,010 synthetic documents integrating charts in a wide range of styles (6 styles in contrast to 3 for PlotQA and ChartQA) and includes 699,051 questions that demand a high degree of reasoning ability and common-sense understanding. Besides, we present the development of a potent question-answer generation engine that employs table data, a rich color set, and basic question templates to produce a vast array of reasoning question-answer pairs automatically. Based on DCQA, we devise an OCR-free transformer for document-level chart-oriented understanding, capable of DLA and answering complex reasoning and common-sense questions over charts in an OCR-free manner. Our DCQA dataset is expected to foster research on understanding visualizations in documents, especially for scenarios that require complex reasoning for charts in the visually-rich document. We implement and evaluate a set of baselines, and our proposed method achieves comparable results.

**Index Terms**— Document Layout Analysis, Chart Question Answering, Common-sense Reasoning.

Anran Wu and Luwei Xiao contributed equally to this work. Corresponding author: Xingjiao Wu (e-mail: xjwu\_cs@fudan.edu.cn).

## 1. INTRODUCTION

The emergence of visual language as a novel communicative tool, characterized by a tightly integrated interplay of visual and textual components, can be attributed to a confluence of factors, notably globalization, the growing intricacy of commerce and technology, and the convergence of lexicons from diverse fields that were once disparate [7]. The prevalence of visually-situated language in various document types, such as academic, business, medical, and others, is markedly high [8–10]. Gaining a comprehensive understanding of these graphical representations, such as charts and plots, is essential in extracting valuable and pragmatic insights from data [11]. To conduct data analysis, individuals frequently pose intricate queries that require common-sense and arithmetic or logical operations pertaining to graphical representations. Answering such inquiries demands a substantial level of cognitive and reasoning exertion, as individuals are required to be aware of common sense and integrate numerous logical operations, including but not limited to retrieving entities, comparing trends, calculating averages, finding extremum, etc. Typically, the chart question answering (CQA) system [12] aims to generate the desired answer by taking a chart-question pair as input, constituting a fundamental function within the domain of intelligent document understanding (IDU) [13].

Despite the CQA task has drawn ever-growing attention from visual question answering communities in recent years, existing datasets has encountered certain obstacles: (i) Notably, while charts constitute crucial components of documents, the majority of current datasets treat the CQA task solely at the question-answering level, without taking into account its significance as a document-level task. (ii) Questions generally prioritize reasoning or visual features, potentially losing sight of common sense information that individuals typically consider when posing questions, which is a misalignment with the typical questioning habits of individuals. (iii) The quantity of chart types, as exemplified by PlotQA and ChartQA datasets, is comparatively restricted (only three). Such a limited representation fails to capture the broad range

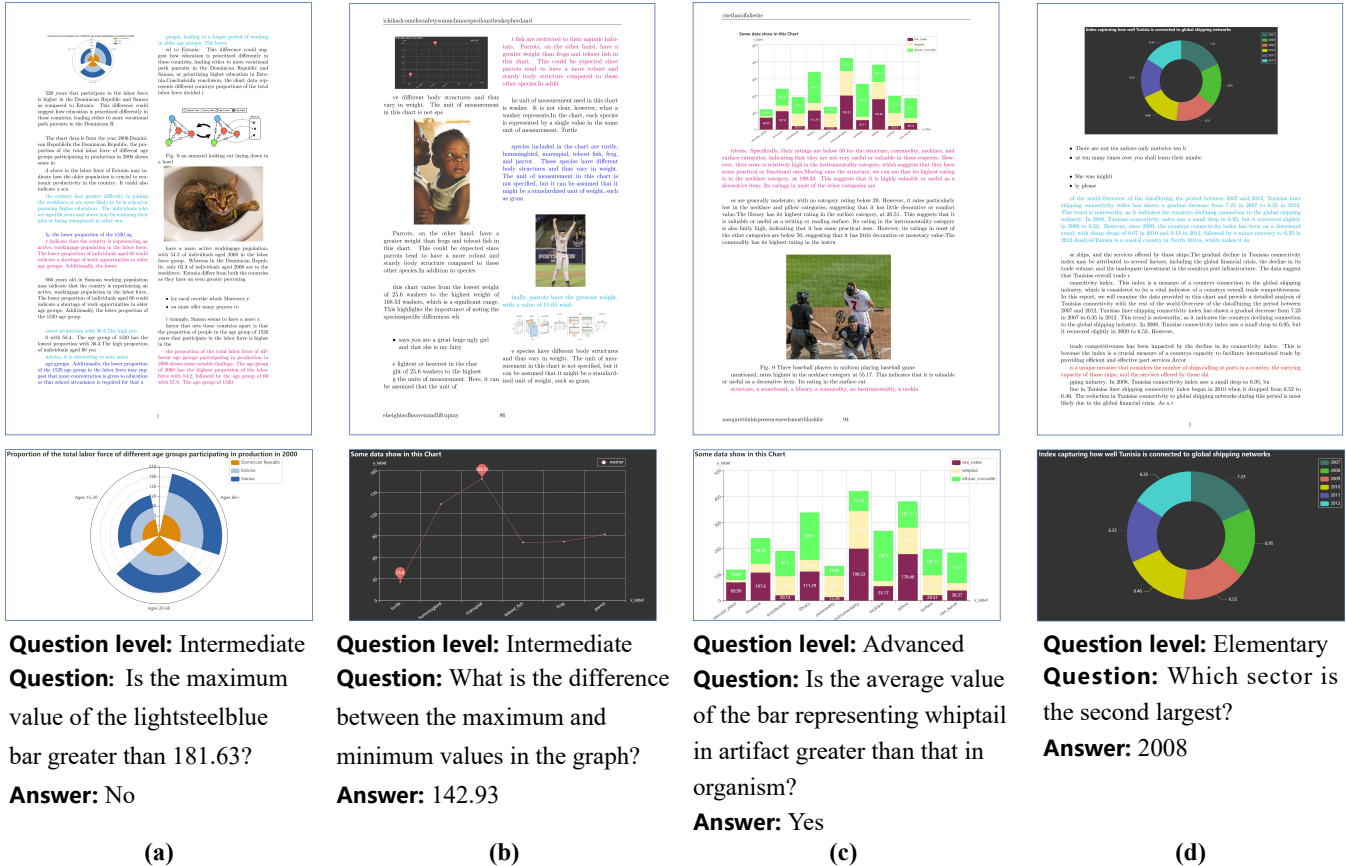


Fig. 1: Examples from DCQA.

of chart styles that are present in real-world documents.

Furthermore, in real-world settings, users typically first identify the location of charts in documents before querying them. However, directly analyzing document layout poses challenges, as charts lack explicit annotations. Manually annotating datasets for document layout analysis related to charts is remarkably laborious and time-consuming. This motivates an automated system to generate annotations associated with charts, eliminating the need for costly labeled data collection. Such a system would locate charts in arbitrary documents without annotations and produce bounding boxes highlighting them, providing chart-specific layout information without human intervention. Additionally, certain baseline models rely on obtaining high-quality optical character recognition (OCR) outcomes to extract the data table structure from the chart image. Therefore, current models' efficacy generally relies upon the accuracy of OCR results. Nevertheless, incorporating an OCR-dependent approach for CQA system poses significant challenges. For one thing, commercially available OCR techniques often exhibit limited adaptability in addressing diverse languages or changes in the domain, which are commonly encountered in the context of charts. Such limitations may impede the generalization abil-

ity of these methods. For another, the occurrence of errors during the OCR process is unavoidable, and such erroneous outcomes have the potential to propagate to the CQA system, thereby adversely affecting subsequent processes [14].

To alleviate above issues, we go beyond the traditional dataset by presenting a large-scale document-level chart question answering dataset (DCQA). DCQA comprises 50,010 synthetic documents and 699,051 question-answer pairs generated using our customized semantic-rich question-answer generation engine. The dataset includes questions that focus on vision, complex reasoning, and common-sense knowledge. Common-sense knowledge reasoning primarily involves evaluating the ability of CQA models to distinguish between legend labels and entity names belonging to specific parent classes, and subsequently performing reasoning operations based on this discriminative ability. Each document in the DCQA includes a chart, unrelated images and a descriptive caption related to the chart. The language used in the DCQA dataset is English. The chart types exhibit a diverse range of styles and can be broadly categorized into six major types, namely Bar chart, Line plot, Pie chart, Scatter plot, Box plot, and Mixed chart, each of which is further divided into subtypes, yielding a total of 30 chart subtypes. Figure 1



displays some examples from DCQA. More examples are provided in Appendix E.

Drawing upon the DCQA dataset, we further devise a transformer-based OCR-free architecture to perform document layout analysis and chart question answering. Initially, we exploit swin transformer [15] as the vision backbone to extract visual features of the input document. Next, the extracted features are fed into the detection component to perform document layout analysis [16]. Upon successfully identifying the chart image, we extract the relevant visual content from the chart, which is then utilized as input to the textual decoder for answer prediction. This novel OCR-free architecture provides a plug-and-play solution for performing chart question answering directly from the document.

In a nutshell, our contributions are as follows:

- We present a comprehensive and extensive document-level chart question answering dataset, DCQA, which features a wide range of chart styles and includes question-answer pairs that incorporate complex reasoning and common-sense knowledge. The dataset’s scale and diversity make it a valuable resource for researchers interested in developing and evaluating chart question answering models.
- We conceptualize chart question answering as a document-level task and propose a transformer-based OCR-free model to effectively address this task.
- We perform comprehensive experiments and thorough analyses on DCQA, verifying the efficacy of our model.

## 2. RELATED WORK

### 2.1. Chart Question Answering

Given the critical role of the CQA subtask in Document AI and the widespread usage of OCR in CQA baseline systems for utilizing layout information, this paper proposes categorizing existing frameworks into two categories: OCR-free and OCR-based. For OCR-free methods, a series of previous studies [1, 17, 18] utilized a combination of Convolutional Neural Networks (CNN) and Long Short-term Memory (LSTM) units to respectively model the visual and textual representations. These modalities were then inputted into a feed-forward network or a Relation Network [19] to predict the answer. Regarding OCR-based approaches, early works [2–4, 20] applied visual encoders to represent the visual content extracted from the chart image. These methods predominantly utilized OCR-based dynamic encoding techniques to model the question representation by incorporating positional information of the textual elements present in the chart. Subsequently, an attention mechanism was employed to interactively learn the features from the chart and the

question, which were then exploited by a softmax classification layer. Recently, a few studies [5, 6] proposed to use OCR to extract visual elements (e.g., legends, x-axis-label, y-axis-label) from the plot and utilize them to reconstruct the underlying data table. They then utilized table QA algorithms in conjunction with the visual backbone to address this task. Although the OCR-based approaches have demonstrated commendable performance, they are still vulnerable to the impact of OCR noise and entail significant computational complexity.

### 2.2. Chart Question Answering Datasets

To date, only a limited number of datasets have been explicitly designed for chart question answering. These datasets include FigureQA [1], DVQA [2], LEAF-QA [3], LEAFQA++ [4], PlotQA [5] and ChartQA [6]. Despite consisting of a diverse set of synthetic charts, FigureQA suffers from a lack of specificity in terms of chart element labeling, utilizing only generic titles and color names. Furthermore, the questions are limited to a few template-based formats with binary “yes/no” answers. DVQA is limited to a single chart type, namely the Bar chart, and suffers from inadequate semantic relations between the textual elements. (e.g., bar and legend labels are randomly selected words) as well as restricted Y-axis value ranges. Numeric answers are primarily integers in both the train and test sets and share the same values. As with FigureQA, all bar plots in DVQA are artificially generated, and the questions are based on a small number of templates. Both LEAF-QA and its extended version, LEAFQA++, are not publicly available. Besides, they share a significant limitation: the absence of regression question-answering pairs. This is evident from the question templates described in their reference and the discrete answer set employed. Although PlotQA is currently the largest publicly available dataset for CQA, it is limited by imbalanced question distribution, as it is heavily weighted towards data-related questions and lacks an appropriate proportion of queries pertaining to the visual characteristics of chart elements, including color and shape. Regarding ChartQA, it is the pioneer dataset to compile real-world charts with a blend of human-created QA pairs and machine-generated QA pairs. However, despite its innovative contribution, ChartQA is characterized by a limited size and encompasses only three distinct plot types. Furthermore, the paucity of question-answer pairs (two) per chart undermines the potential for a comprehensive understanding of the underlying information conveyed by these visualizations. This work presents a novel and intricate CQA dataset, which diverges from prior datasets in several respects. Firstly, DCQA is introduced, which reformulates the CQA task by integrating document layout analysis and chart question answering. Secondly, in addition to visual and complex reasoning questions, DCQA incorporates common sense-aware questions. Last but not least, DCQA covers a broad range of chart types

**Table 1:** A comprehensive comparison between existing datasets and our proposed DCQA.

Datasets	Chart				Question				Available
	Document-level	Types	Num.	Source	Com. sense	Types	Templates	Num.	
FigureQA [1]	✗	4	180,000	Colour set	✗	Template based	15	2,388,698	✓
DVQA [2]	✗	1	300,000	Brown corpus	✗	Template based	26	3,487,194	✓
LEAF-QA [3]	✗	5	245,903	Real-world	✗	Template based	35	1,906,486	✗
LEAF-QA++ [4]	✗	5	245,903	Real-world	✗	Template based	75	2,589,355	✗
PlotQA-D1 [5]	✗	3	224,377	Real-world	✗	Template based	74	8,190,674	✓
PlotQA-D2 [5]	✗	3	224,377	Real-world	✗	Template based	74	28,952,641	✓
ChartQA-H [6]	✗	3	4,804	Real-world	✗	Human authored	None	9,608	✓
DCQA	✓	6	50,010	Real-world & WordNet	✓	Template based	324	699,051	✓

**Fig. 2:** The illustration of chart types.

(6 main types and 30 sub-types) with a wealth of color types (595) and different backgrounds (white and black). A comprehensive comparison between existing datasets and DCQA is presented in Table 1.

### 3. DCQA DATASET

In this section, we describe the construction of DCQA from the following four aspects: (i) Data collection, (ii) Chart generation, (iii) Automatic QA pair generation engine, and (iv) Document generation. The general workflow of the DCQA generation process is shown in Figure 3. A detailed generation procedure is provided in Appendix A.

#### 3.1. Data Collection

Given the variability of chart styles in real-world scenarios, integrating real-world sources and randomly generated data for producing charts can augment the models robustness and

adaptability to various chart formats encountered in practical scenarios. Drawing upon this observation, charts included in our dataset was derived by two means: utilizing real-world sources and randomly generated data. The detailed process of data collection is shown in Appendix A.

#### 3.2. Chart Generation

We exploit Pyecharts<sup>1</sup>, a Python visualization tool library based on the Echarts [21] charting library, to generate our charts. Our DCQA contains six different chart styles: bar chart, line chart, scatter plot, box plot, pie chart, and combination chart (line and bar). These chart styles can be further divided into 30 sub-types (As shown in Figure 2). The color of chart elements is randomly picked up from a color set JohnDecember, which covers a wide range of colors (595). Besides, the chart presents two distinct styles of background, namely dark and light, of which the former was not previously observed in any of the CQA datasets. Detailed chart information is provided in Appendix B.

#### 3.3. Automatic QA Pair Generation Engine

##### 3.3.1. Question Collection via Crowd-sourcing.

Since the generated charts are from disparate data sources and encompass a wide range of topics, engaging a cadre of individuals with diverse backgrounds, experiences, and expertise is necessary to craft questions about the corresponding charts. We have meticulously curated a corpus of 573 charts spanning six categories, comprising data extracted from real-world and randomly generated sources, which serve as paradigmatic instances from which questions can be formulated. We commission a cohort of post-graduate students affiliated with our academic institution, and employees from Huawei company, to generate ten distinct questions for each of the selected charts, with an emphasis on reasoning and common sense awareness. We have obtained 5730 questions.

<sup>1</sup><https://github.com/pyecharts/pyecharts>

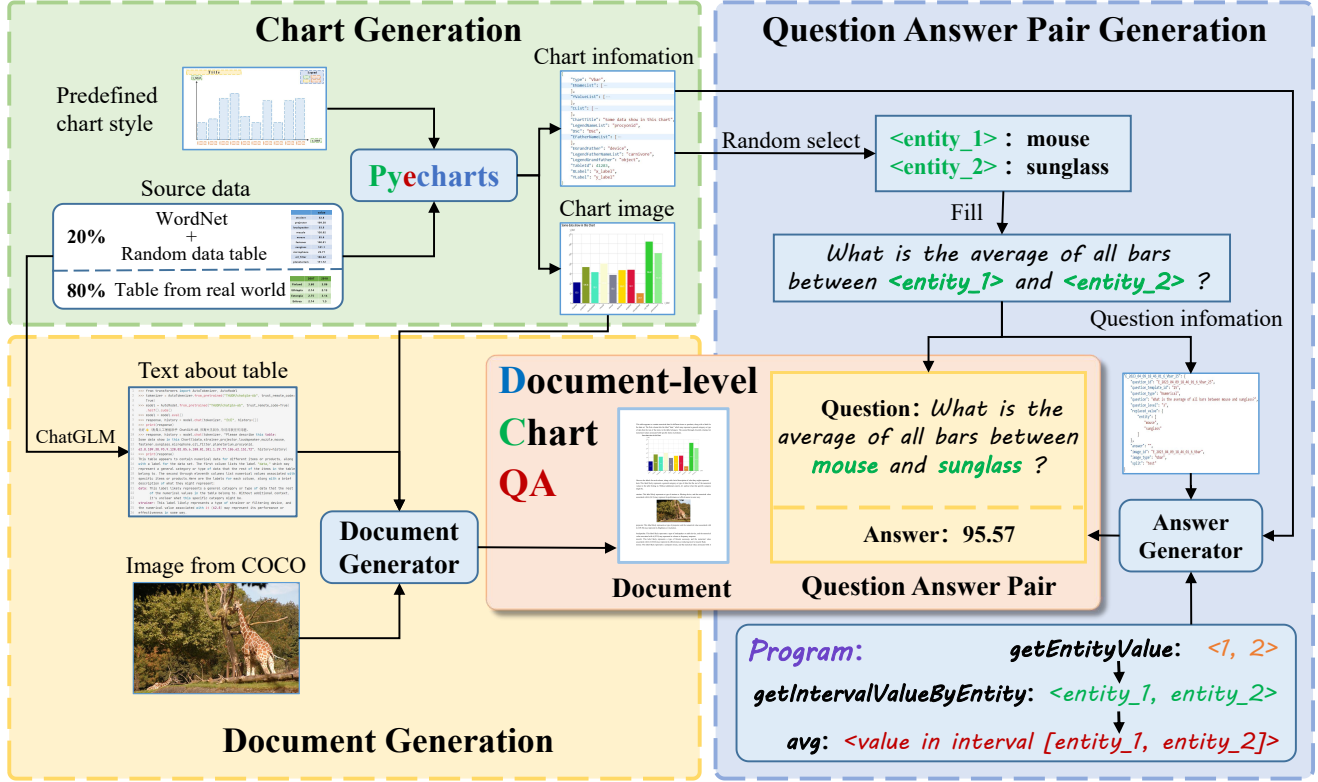


Fig. 3: The generation workflow of DCQA.

### 3.3.2. Question Generation.

After an exhaustive process of meticulous meetings and in-depth discussions, we have successfully distilled a total of 324 question templates from the original pool of 5730 questions. Out of these templates, 204 are specifically tailored for visual and numeric reasoning, while 120 templates are dedicated to common sense reasoning. Code will be publicly available at github . Table 2 displays the statistics of the dataset. Details in Appendix C.2.

**Visual and numeric reasoning:** These kinds of questions necessitate combining visual elements understanding and numerical reasoning techniques (e.g., sum, multiple, average, etc.). Integrating visual and numerical reasoning inquiries can facilitate CQA systems’ comprehension of chart content, as it encourages the concurrent utilization of their visual and analytic faculties, thereby enabling them to engage in a more profound exploration of the underlying message conveyed by the data and achieve a more precise interpretation of chart figures. Examples of this type of question are presented in Figure 1 (a) and (b).

**Common sense reasoning:** Questions of this type demand combining common sense knowledge and numerical operation. Common sense is able to serve as a facilitator for CQA systems to gain a more profound understanding of the real-life background and context reflected by the data,

thus accurately inferring the meaning behind the data. Meanwhile, numerical reasoning skills can allow readers to fathom the underlying interconnections and relationships of the data and infer potential outcomes and trends. The combination of both proficiencies can profoundly equip CQA models with diverse conceptualizations of chart content and enable them to increase the usefulness of data in comprehensively examining and scrutinizing data, identifying patterns and trends, and making predictions and decisions. Examples of this type of question are presented in Figure 1 (c).

**Construction of the hierarchical entity database:** Common sense reasoning is a crucial aspect of DCQA, which primarily involves evaluating a CQA model’s ability to discriminate between legend labels and entity names that belong to a specific parent class and then perform reasoning operations based on this discriminatory capacity. Therefore, a hierarchical entity database with a tree-structured architecture and a well-defined set of parent-child relationships is necessary to serve as a source for both entity names and legend labels. The construction of the hierarchical entity database is expounded upon in Appendix C.1.

**Categorization of question difficulty levels:** Additionally, the entire set of question templates has been classified into five distinct levels delineated by their respective levels of complexity, namely, beginner, elementary, intermediate, ad-

**Table 2:** Detailed statistics of DCQA.

Dataset split	Chart style						Question type		Answer type		
	Bar	Line	Pie	Scatter	Box	Comb.	Reasoning	Com. sense	Yes/No	Elements	Open vocab.
Training	14,674	9,338	5,336	5,336	4,002	1,334	496,241	59,082	171,447	51,215	332,661
Validation	1,826	1,162	664	664	498	166	61,752	9,058	22,638	6,513	41,659
Test	1,837	1,169	668	668	501	167	62,124	10,794	23,261	6,688	42,969
Total	18,337	11,669	6,668	6,668	5,001	1,667	620,117	78,934	217,346	64,416	417,289

**Table 3:** The distribution of different question levels.

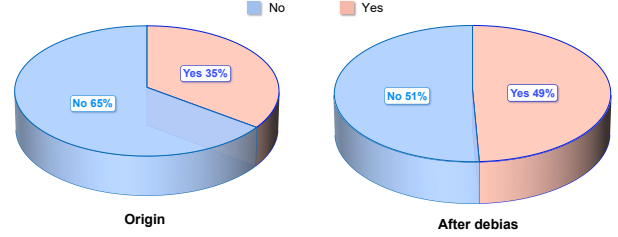
Split	Question levels				
	Beginner	Elementary	Intermediate	Advanced	Expert
Train	49,358	157,679	239,720	96,637	11,929
Val.	6,142	19,630	29,838	13,244	1,956
Test	6,179	19,746	30,053	14,071	2,869
Total	61,679	197,055	299,611	123,952	16,754

vanced, and expert.

The statistic of question levels is displayed in Table 3. The difficulty levels of the question templates are manually annotated based on the following criteria: (1) Beginner includes questions about the overall nature of a chart image, such as whether it is horizontal or vertical or the number of columns it contains, as well as retrieval for the value of a specific chart element. (2) Elementary primarily involves questions carrying out some form of operation on all chart elements within a chart, such as determining the maximum, minimum, median, or mean value. (3) Intermediate refers to questions that involve applying specific operations to chart elements that satisfy predetermined criteria, including but not limited to identifying the maximum, minimum, median, mean, sum, or difference of chart elements based on their color, legend, or numerical value. (4) Advanced questions demand performing composite operations on chart elements that meet a specific property (building upon the operations mentioned for intermediate questions), such as finding the sum or difference of two maximum values after they have been identified. (5) Building upon advanced questions, questions that involve common sense will be categorized as expert-level.

### 3.3.3. Answer Generation.

In this paper, answers are generated through an automated process based on a customized set of procedures. Specifically, a solution step is designed for each question template, with each step representing an atomic operation that is implemented using specific functions to achieve its intended functionality. When solving specific questions to generate answers, the designed solution steps are followed by invoking corresponding functions, resulting in answers for the respec-

**Fig. 4:** Comparison of Yes and No answer proportions for biased Yes/No questions before and after debiasing.

tive questions. The process of obtaining answers and the advantages of automatic answer generation are elaborated upon in Appendix A.

### 3.3.4. Debiasing of Question-answer Pairs.

Upon completing the question-answer pair generation process, extra analysis is conducted on the distribution of every answer type. It is noted that the highest proportion of the answer type in the dataset is the binary classification yes or no. However, the ratio between yes and no is severely imbalanced, with a vastly larger number of no responses compared to yes. As a result, a post-processing adjustment is necessary to address this language bias and prevent the model from exploiting the answer distribution pattern to output the answer without paying attention to the visual content.

The debiasing procedure can be concisely described as: Firstly, filter out all question templates with Yes/No answers. Secondly, count the number of Yes/No answers for each Yes/No question template in the dataset. For each Yes/No question template, determine whether the number of Yes answers exceeds the number of No answers or vice versa. For the question with a larger proportion of answers, adjust the values of the replaceable modules in the question to change the answer to the less frequent one, and iterate this process until the number of Yes and No answers for the template are equal or differ by only 1. Most of the Yes/No question templates can be balanced by modifying the answer using the above approach. However, a few questions cannot be balanced this way, and their answers are not significantly

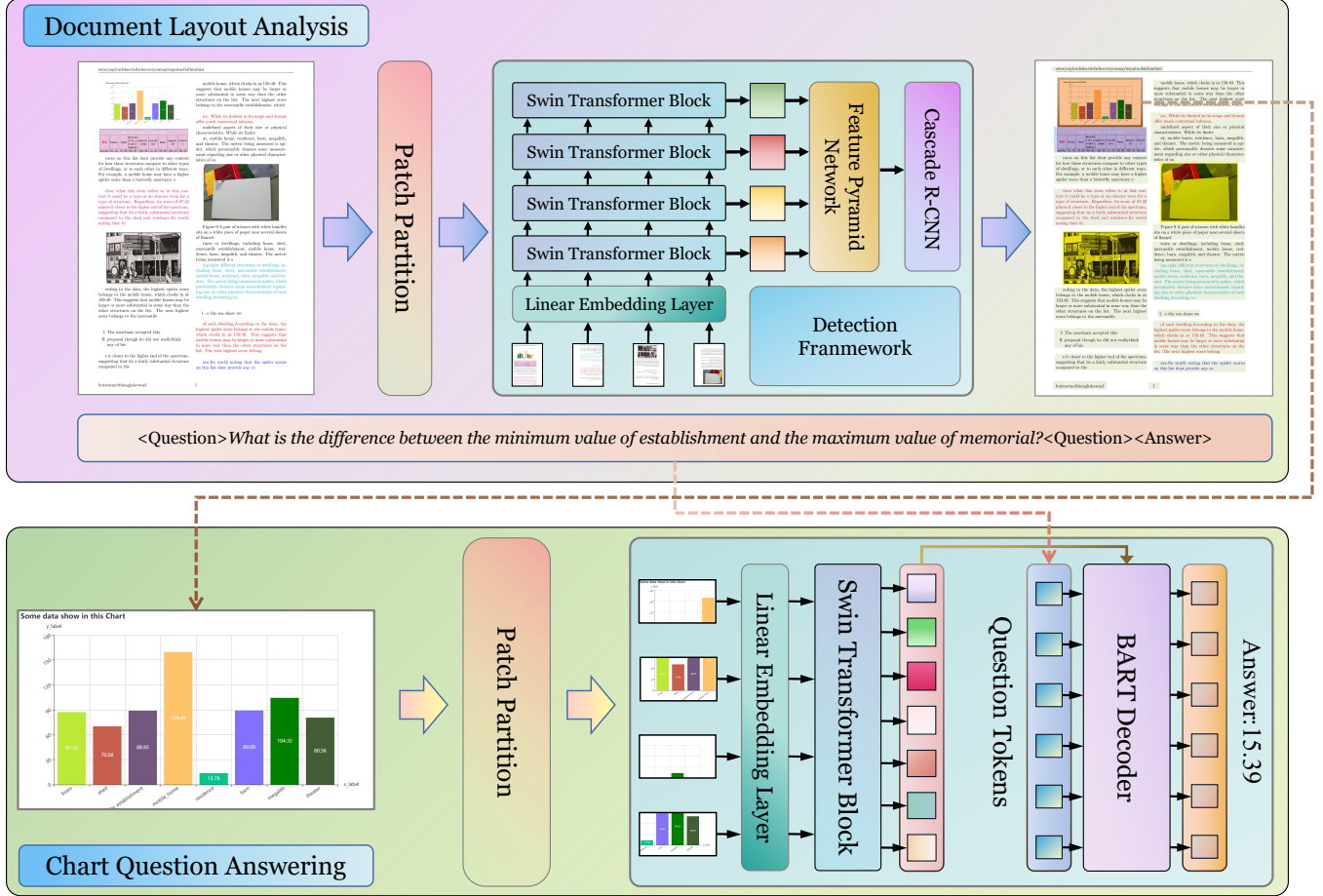


Fig. 5: The overall architecture of the proposed TOT-Doctor.

imbalanced. Therefore, we do not handle such question templates in this paper for now. Before debiasing, the dataset's overall proportion of yes and no answers was 35.16% and 64.84%, respectively. After debiasing, the overall proportion of yes and no answers in the dataset became 49.26% and 50.74%, respectively. The effectiveness of the debiasing is compared in Figure 4, where the noticeable changes in the answer distribution of "Yes" and "No" before and after debiasing demonstrate that the debiasing method employed in this study effectively addresses the answer imbalance in binary questions.

### 3.4. Document Generation

In accordance with the methodology outlined in [22], we utilize LaTeX to generate synthetic documents that include diverse multimedia elements. In addition to the chart image, the generated document also incorporates other visual content<sup>2</sup> and textual content produced by ChatGLM [23, 24]. Notably, integrating these elements augments the informational value

of the synthetic document beyond the mere inclusion of the chart image, which is more consistent with real-world documents. Detailed document information is provided in Appendix D.

## 4. METHOD

The DCQA dataset contains samples that more closely resemble real-world scenarios, considering the omnipresence of charts in documents. Our proposed approach for effectively comprehending chart-oriented documents is through the application of a novel OCR-free multi-modal CQA model, TOT-Doctor, **Two-stage OCR-free transformer for document-level chart-oriented understanding**. An overall architecture of the TOT-Doctor is displayed in Figure 5.

### 4.1. Document Layout Analysis

With the input document image denoted as  $I \in \mathbb{R}^{3 \times H \times W}$ , a patch partition module is utilized to divide it into non-overlapping patches. Subsequently, these patches are passed

<sup>2</sup><https://cocodataset.org>

**Table 4:** Hyperparameters of the baselines and the proposed method that were used on DCQA. The term "CE" refers to Cross Entropy.

Model	LayoutLMv2	LayoutLMv3	Pix2struct	Matcha	Ours
Training Loss	CE	CE	CE	CE	CE
Batch Size	12	12	16	12	16
Learning Rate	$10^{-5}$	$2 \times 10^{-5}$	$2 \times 10^{-5}$	$2 \times 10^{-6}$	$4 \times 10^{-5}$
Optimizer	AdamW	AdamW	AdamW	AdamW	AdamW
Scheduler	multi step	multi step	cosine schedule	cosine schedule	cosine schedule

**Table 5:** Performance of the document layout analysis stage of ToT-Doctor on the DCQA test dataset (%). PH stands for page header, PF stands for page footer, and PG represents page number.

	Chart imgae	Text	Picture	Caption	List	PH	PF	PG	Total
AP	99.901	92.332	99.667	92.968	95.760	93.478	90.775	81.580	<b>93.971</b>

**Table 6:** Performance of our TOT-Doctor and other baselines on DCQA. Best results are in bold. L2\_B indicates LayoutLMv2-base, L3\_B means LayoutLMv3-base.

Method	Dev	Test
L2_B [25]	4.844%	4.816%
L3_B [26]	15.768%	15.760%
Pix2Struct-base [9]	19.487%	19.294%
MATCHA-base [27]	19.249%	19.132%
TOT-Doctor	<b>40.658%</b>	<b>40.090%</b>

into a detection framework comprising several Swin Transformer-base blocks, a Feature Pyramid Network (FPN) [28], and a cascade R-CNN [29]. Initially, a linear embedding layer is employed to map the raw-valued RGB features of patches into a fixed-dimensional representation and then fed into Swin Transformer blocks. Since the scale of feature maps generated by different Swin Transformer blocks is distinct, they are fed into the multi-scale FPN. Next, the output from FPN is taken as the input for a Cascade R-CNN. Utilizing the prediction results of document objects, we extract the chart area and combine it with questions for chart understanding.

## 4.2. Chart Question Answering

After we obtain the chart image  $I \in \mathbb{R}^{3 \times H_0 \times W_0}$  from the document via document layout analysis, at this stage, we utilize chart-question pairs to conduct chart question answering. Concretely, we model the visual features by dividing the chart image into non-overlapping patches via the patch partition module and passing through them into Swin Transformer-

base. Following the teacher-forcing scheme [30], the chart visual features  $V \in \mathbb{R}^{m \times d}$  obtained from the last Swin Transformer block is denoted as  $K$  for BART decoder [31]. In the training stage, the question representation  $S \in \mathbb{R}^{n \times d}$  and the ground-truth  $A \in \mathbb{R}^{k \times d}$  are concatenated as  $Q$  for BART decoder. During the test phase, following the prompt scheme by GPT-3 [32], we add special tokens (as shown in Figure 5) to assist the model in generating the sequence. Subsequently, we feed them jointly into the BART decoder to generate the predicted sequence  $(y_i)_{i=1}^z$ ,  $y \in \mathbb{R}^g$  is the one-hot vector refers to the  $i$ -th generated token,  $g$  denotes the size of the vocabulary,  $z$  is a hyperparameter and indicates the number of tokens generated in the sequence.

## 5. EXPERIMENTS

In this section, we provide a comprehensive analysis of the experimental results to establish the validity of the recently developed DCQA dataset and verify the excellent efficacy of our proposed TOT-Doctor model through a comparative evaluation against other baselines.

### 5.1. Baselines

We compare the TOT-Doctor with the classical approaches include LayoutLMv2 [25], LayoutLMv3 [26], Pix2Struct [9], and MATCHA [27].

- LayoutLMv2 [25] is a Transformer-based encoder that models multi-modal information, including text, layout, and image, by incorporating a spatial-aware self-attention mechanism along with two innovative pre-training strategies.



**Table 7:** Performance of our TOT-Doctor and other baselines on different question levels. Best results are in bold.

Question levels	LayoutLMv2		LayoutLMv3		Pix2Struct		MATCHA		TOT-Doctor	
	Dev	Test	Dev	Test	Dev	Test	Dev	Test	Dev	Test
Beginner	0.928%	0.955%	11.071%	10.989%	22.485%	22.706%	19.065%	18.887%	<b>44.302%</b>	<b>42.337%</b>
Elementary	3.316%	3.150%	14.544%	14.671%	18.945%	18.373%	17.835%	17.371%	<b>42.792%</b>	<b>41.948%</b>
Intermediate	8.312%	8.312%	21.084%	21.033%	24.073%	23.771%	25.028%	25.159%	<b>47.543%</b>	<b>47.153%</b>
Advanced	1.608%	2.132%	8.744%	9.580%	9.906%	11.321%	9.665%	9.665%	<b>24.268%</b>	<b>26.778%</b>
Expert	1.483%	1.150%	9.254%	8.609%	10.429%	10.491%	10.429%	10.422%	<b>13.753%</b>	<b>13.768%</b>

**Table 8:** Performance of our TOT-Doctor and other baselines on different Answer types. Best results are in bold.

Answer types	LayoutLMv2		LayoutLMv3		Pix2Struct		MATCHA		TOT-Doctor	
	Dev	Test	Dev	Test	Dev	Test	Dev	Test	Dev	Test
Yes/No	15.152%	15.098%	48.944%	48.996%	49.249%	49.142%	50.362%	50.209%	<b>78.103%</b>	<b>77.529%</b>
Numerical	0%	0%	0.154%	0.142%	5.960%	5.879%	5.202%	5.153%	<b>24.650%</b>	<b>24.006%</b>
String	0%	0%	0.322%	0.508%	2.564%	1.675%	0.952%	0.867%	<b>12.897%</b>	<b>13.218%</b>

- LayoutLMv3 [26] is a multi-modal Transformer framework without the vision backbone that leverages reconstructive objectives for cross-modal alignment learning, showcasing notable generality in the context of document vision tasks.
- Pix2Struct [9] is an image-to-text model tailored for visual language understanding, which is pre-trained on visually-rich screenshots of web pages with *screenshot parsing* objective.
- MATCHA [27] is a Pix2Struct-based model pre-trained for chart underlying structure understanding and mathematical reasoning.

## 5.2. Evaluation Metrics

Our study employs accuracy the primary evaluation metric, wherein the assessment of the predicted answers correctness depends on the nature of the answer type. In the case of textual answers, such as binary responses, entities, and integers, the evaluation criterion mandates that the predicted answer should match the ground truth exactly. For numerical answers in the form of floating-point values, it is not always feasible to expect that the predicted answer will precisely match the correct answer. Therefore, we consider the answer correct if it falls within 5% of the expected value.

## 5.3. Experiment Setup

This section primarily supplements experimental configurations, including parameters such as batch size and learning rate, and outlines the preprocessing steps undertaken to ensure a fair comparison for the extractive model.

Table 4 shows the detailed experimental setup, in which CE refers to Cross Entropy. All models were verified every 5000 iterations during training. In our implementation of LayoutLMv2 and v3, we employ a multi-step learning rate schedule. More specifically, we gradually decrease the learning rate by a factor of 2 after each epoch of training. For other models, we use the cosine scheduler to adjust the learning rate, where the number of warm-up steps is set to 1000. The LayoutLM series employs a model-based extraction approach, which requires the system to select answers from the optical character recognition (OCR) results. To enable the model to perform tasks such as binary classification (e.g., Yes/No), we have implemented a simple yet effective solution: we add two special characters, "Yes" and "No", to the OCR results.

## 5.4. Model Parameter Quantity

TOT-Doctor consists of two main components, namely the document layout analysis and the chart question answering. The model parameters of TOT-Doctor are calculated and found to be as follows: the encoder of the Swin Transformer used in the document layout analysis has 74M parameters, while the detector component has 48M parameters. The encoder of the Swin Transformer used in the chart question answering phase has 74M parameters, and the BART has 127M parameters.

## 5.5. Main Results

### 5.5.1. Results of Document Layout Analysis

Based on the fine-grained document element annotations present within the DCQA dataset introduced in this paper, encompassing various elements such as chart image, picture,

textual content, list, caption, header, footer, and page number, the document layout analysis model of the proposed ToT-Doctor framework was trained. The conclusive results of the document layout analysis testing on the test set are presented in Table 5. Notably, the detection accuracy for chart image reached an impressive 99.901%, exhibiting a near-complete precision in accurately identifying their respective spatial position. This achievement serves as a robust foundational prerequisite for facilitating subsequent stage of chart question answering task.

### 5.5.2. Results of Question Answering

The experiment results of our proposed TOT-Doctor and other baselines are displayed in Table 6. Due to the incapacity of the baseline to perform document layout analysis, we add our document layout analysis framework to them before conduct chart question answering.

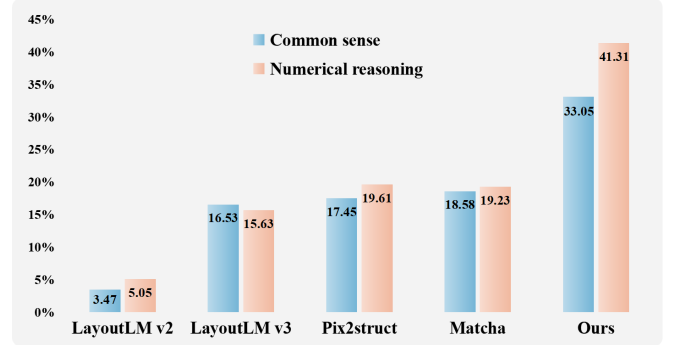
Based on the data listed in Table 6, it can be seen that TOT-Doctor consistently surpasses its counterparts concerning the accuracy in both the validation and test sets, corroborating the efficacy of TOT-Doctor. It is noteworthy that despite LayoutLMv2 and LayoutLMv3 being reliant on OCR for obtaining the answers, their performance continues to lag behind our OCR-free TOT-Doctor. This observation proves the robustness and OCR error mitigation capabilities of the TOT-Doctor proposed in this study. Furthermore, in comparison to the latest state-of-the-art (SOTA) pre-trained visual language understanding model Pix2Struct, TOT-Doctor demonstrates superior performance, achieving a significant improvement of approximately 21.171% on the development dataset, and a respectable enhancement of 20.796% on the test dataset, respectively. The outstanding performance of our TOT-Doctor model underscores the significance of integrating vision and language features in an OCR-free manner to address the questions posed in DCQA effectively.

## 5.6. Ablation Study

### 5.6.1. Evaluation on Different Question Levels.

The DCQA dataset incorporates five distinct question levels. In order to better discern the effectiveness of our proposed TOT-Doctor and other baselines, we conduct a performance analysis on each question level. The results of our analysis are presented in Table 7 for reference. It is evident from the results that TOT-Doctor consistently outperforms other baselines in all five levels of questions. Notably, TOT-Doctor exhibits superior performance on intermediate-level questions, demonstrating its efficacy in directly applying specific operations such as identifying maximum, minimum, median, mean, and sum to chart elements or analyzing the differences of chart elements based on their color, legend, or numerical value. However, when encountering the sub-difficult advanced-level questions involving composite oper-

ations and the most challenging expert-level questions that necessitate commonsense understanding, the performance of all baselines, including TOT-Doctor, considerably decreases compared to the other three more tractable question levels. This implies that the ability for complex reasoning and commonsense understanding still requires further improvement for analyzing complex documents.



**Fig. 6:** Performance comparison between common sense and numerical reasoning questions on DCQA test set.

### 5.6.2. Evaluation on Different Question Types.

We conduct additional experiments to evaluate the performance of the TOT-Doctor and baseline models on different question types. As discussed before, DCQA comprises two primary question types: visual and numeric reasoning and commonsense reasoning. Figure 6 presents the results of all baselines for each question type on the test set. Our proposed TOT-Doctor outperforms other baselines significantly, particularly in numerical reasoning questions. To top it all off, TOT-Doctor demonstrates more proficiency in numerical reasoning compared to commonsense understanding.

### 5.6.3. Evaluation on Different Answer Types.

To further investigate our TOT-Doctor model, we assess the ability of the TOT-Doctor to generate different answer types. From Table 8, we discover that except for LayoutLMv2, all other baselines perform better in answering Yes/No questions. We observe that LayoutLMv2 and LayoutLMv3 exhibit frustrating performance in generating numerical and string answers. We speculate that this is mainly because LayoutLMv2 and LayoutLMv3 are extractive models, which means that they cannot generate answers that have not appeared in the document. This precisely explains their poor performance on the DCQA dataset, where the answers are largely obtained through data reasoning and involve numerical values. It is noted that EasyOCR<sup>3</sup> is utilized as the OCR system for these two baselines, which deviates from the OCR

<sup>3</sup><https://github.com/JaidedAI/EasyOCR>

employed in their original versions. Based on this observation, we posit that the accuracy of the OCR system may have contributed to the subpar performance of the models. Moreover, TOT-Doctor is well versed in generating answers in terms of numerical or string, which verifies the robustness of TOT-Doctor. However, overall, all baselines achieve abysmal accuracy in generating numerical and string answers, highlighting the significant challenge posed by document-level chart understanding, which calls for further research efforts.

## 6. CONCLUSION

This paper introduces a newly developed and demanding dataset for document-level chart question answering called DCQA. The dataset comprises a sizeable body of documents containing a wide variety of chart styles. The questions posed within this dataset require complex reasoning skills and a strong capacity for common-sense understanding with regard to the information contained within the charts. We developed a specialized question-answer generation engine capable of automatically producing all question-answer pairs. Additionally, to tackle the highly challenging DCQA dataset, this paper proposes a novel OCR-free transformer-based framework, TOT-Doctor, designed explicitly for document-level chart question answering without relying on optical character recognition techniques. Experimental results demonstrate the veracity of our proposed TOT-Doctor in addressing complex document-level chart question answering tasks. However, the DCQA dataset has some limitations, such as the single background style of the charts, relatively simple document layouts, and the inability to exhaustively cover all possible questions generated by the templates. In future work, we plan to address these limitations by enriching the chart backgrounds with various colors and textures, considering more complex document layouts, and thereby enhancing the scalability of our method.

## 7. REFERENCES

- [1] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio, “Figureqa: An annotated figure dataset for visual reasoning,” *arXiv preprint arXiv:1710.07300*, 2017.
- [2] Kushal Kaffle, Brian Price, Scott Cohen, and Christopher Kanan, “Dvqa: Understanding data visualizations via question answering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5648–5656.
- [3] Ritwick Chaudhry, Sumit Shekhar, Utkarsh Gupta, Pranav Maneriker, Prann Bansal, and Ajay Joshi, “Leaf-qa: Locate, encode & attend for figure question answering,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 3512–3521.
- [4] Hrituraj Singh and Sumit Shekhar, “Stl-cqa: Structure-based transformers with localization and encoding for chart question answering,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 3275–3284.
- [5] Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar, “Plotqa: Reasoning over scientific plots,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 1527–1536.
- [6] Ahmed Masry, Do Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque, “Chartqa: A benchmark for question answering about charts with visual and logical reasoning,” in *Findings of the Association for Computational Linguistics: ACL 2022*, 2022, pp. 2263–2279.
- [7] Robert E Horn, “Visual language,” *MacroVu Inc. Washington*, 1998.
- [8] Dae Hyun Kim, Enamul Hoque, and Maneesh Agrawala, “Answering questions about charts and generating visual explanations,” in *Proceedings of the 2020 CHI conference on human factors in computing systems*, 2020, pp. 1–13.
- [9] Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova, “Pix2struct: Screenshot parsing as pretraining for visual language understanding,” *arXiv preprint arXiv:2210.03347*, 2022.
- [10] Weihong Ma, Hesuo Zhang, Shuang Yan, Guangshun Yao, Yichao Huang, Hui Li, Yaqiang Wu, and Lianwen Jin, “Towards an efficient framework for data extraction from chart images,” in *Document Analysis and Recognition-ICDAR 2021: 16th International Conference, Lausanne, Switzerland, September 5–10, 2021, Proceedings, Part I*. Springer, 2021, pp. 583–597.
- [11] Kenny Davila, Srirangaraj Setlur, David Doermann, Bhargava Urala Kota, and Venu Govindaraju, “Chart mining: A survey of methods for automated chart analysis,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 11, pp. 3799–3819, 2020.
- [12] Matan Levy, Rami Ben-Ari, and Dani Lischinski, “Classification-regression for chart comprehension,” in *Computer Vision-ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI*. Springer, 2022, pp. 469–484.
- [13] Junlong Li, Yiheng Xu, Tengchao Lv, Lei Cui, Cha Zhang, and Furu Wei, “Dit: Self-supervised pre-training for document image transformer,” in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 3530–3539.
- [14] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park, “Ocr-free document understanding transformer,” in *Computer Vision-ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVIII*. Springer, 2022, pp. 498–517.
- [15] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 10012–10022.

- [16] Galal M Binmakhashen and Sabri A Mahmoud, “Document layout analysis: a comprehensive survey,” *ACM Computing Surveys (CSUR)*, vol. 52, no. 6, pp. 1–36, 2019.
- [17] Revanth Reddy, Rahul Ramesh, Ameet Deshpande, and Mitesh M Khapra, “FigureNet: A deep learning model for question-answering on scientific plots,” in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.
- [18] Jialong Zou, Guoli Wu, Taofeng Xue, and Qingfeng Wu, “An affinity-driven relation network for figure question answering,” in *2020 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2020, pp. 1–6.
- [19] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillcrap, “A simple neural network module for relational reasoning,” *Advances in neural information processing systems*, vol. 30, 2017.
- [20] Kushal Kafle, Robik Shrestha, Scott Cohen, Brian Price, and Christopher Kanan, “Answering questions about data visualizations using efficient bimodal fusion,” in *Proceedings of the IEEE/CVF Winter conference on applications of computer vision*, 2020, pp. 1498–1507.
- [21] Deqing Li, Honghui Mei, Yi Shen, Shuang Su, Wenli Zhang, Junting Wang, Ming Zu, and Wei Chen, “Echarts: a declarative framework for rapid construction of web-based visualization,” *Visual Informatics*, vol. 2, no. 2, pp. 136–146, 2018.
- [22] Xingjiao Wu, Yingbin Zheng, Tianlong Ma, Hao Ye, and Liang He, “Document image layout analysis via explicit edge embedding network,” *Information Sciences*, vol. 577, pp. 436–448, 2021.
- [23] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang, “Glm: General language model pretraining with autoregressive blank infilling,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022, pp. 320–335.
- [24] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al., “Glm-130b: An open bilingual pre-trained model,” *arXiv preprint arXiv:2210.02414*, 2022.
- [25] Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, et al., “Layoutlmv2: Multi-modal pre-training for visually-rich document understanding,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2021, pp. 2579–2591.
- [26] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei, “Layoutlmv3: Pre-training for document ai with unified text and image masking,” in *Proceedings of the 30th ACM International Conference on Multimedia*, 2022, pp. 4083–4091.
- [27] Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenschlos, “Matcha: Enhancing visual language pretraining with math reasoning and chart derendering,” *arXiv preprint arXiv:2212.09662*, 2022.
- [28] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, “Feature pyramid networks for object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [29] Zhaowei Cai and Nuno Vasconcelos, “Cascade r-cnn: Delving into high quality object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6154–6162.
- [30] Ronald J Williams and David Zipser, “A learning algorithm for continually running fully recurrent neural networks,” *Neural computation*, vol. 1, no. 2, pp. 270–280, 1989.
- [31] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer, “Multilingual denoising pre-training for neural machine translation,” *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 726–742, 2020.
- [32] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al., “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

This supplementary material will provide more description, experiment and the details of DCQA dataset, which are organized as outlined below:

- **Section A: Dataset Generate Workflow.** This section provides further details of the DCQA generation workflow, supplemented by specific examples to facilitate a deeper comprehension.
- **Section B: Chart Information.** This section presents the abundant chart information encompassed within the DCQA dataset, thereby furnishing researchers with a wealth of potential resources for future investigations.
- **Section C: Question Generation.** This section delineates further intricacies in question generation within the DCQA dataset, alongside endeavors undertaken to enhance the robustness of question-answer pairs therein.
- **Section D: Document Information.** This section elucidates the specifics of fine-grained annotation information within the DCQA dataset documents, which constitute a fundamental underpinning for the chart-oriented document layout analysis task.
- **Section E: Document Example.** This section showcases examples of fine-grained categorization for chart types within the DCQA dataset, encompassing a total of 6 primary categories and 30 subcategories of charts.

## A. DATASET GENERATE WORKFLOW

The DCQA dataset generation process, as depicted in Figure A1, encompasses four main components: chart generation, question generation, answer generation, and document generation. In addition, this section also encompasses the specifics of data collection.

**Data Collection.** For the real-world data sources, we crawl data from the following on-line sources: (i) World Bank Open Data (<https://data.worldbank.org>). (ii) United Nations Data (<https://data.un.org/>). (iii) International Monetary Fund Data (<https://data.imf.org/>). For the randomly generated data, we adopt part of the WordNet<sup>4</sup> as the chart element (e.g., legend names, axes ticks) set and randomly select and assign these labels to the entities. The values of the entities are floating-point numbers rounded to two decimal places, randomly selected from a uniform distribution ranging from 1 to 200.

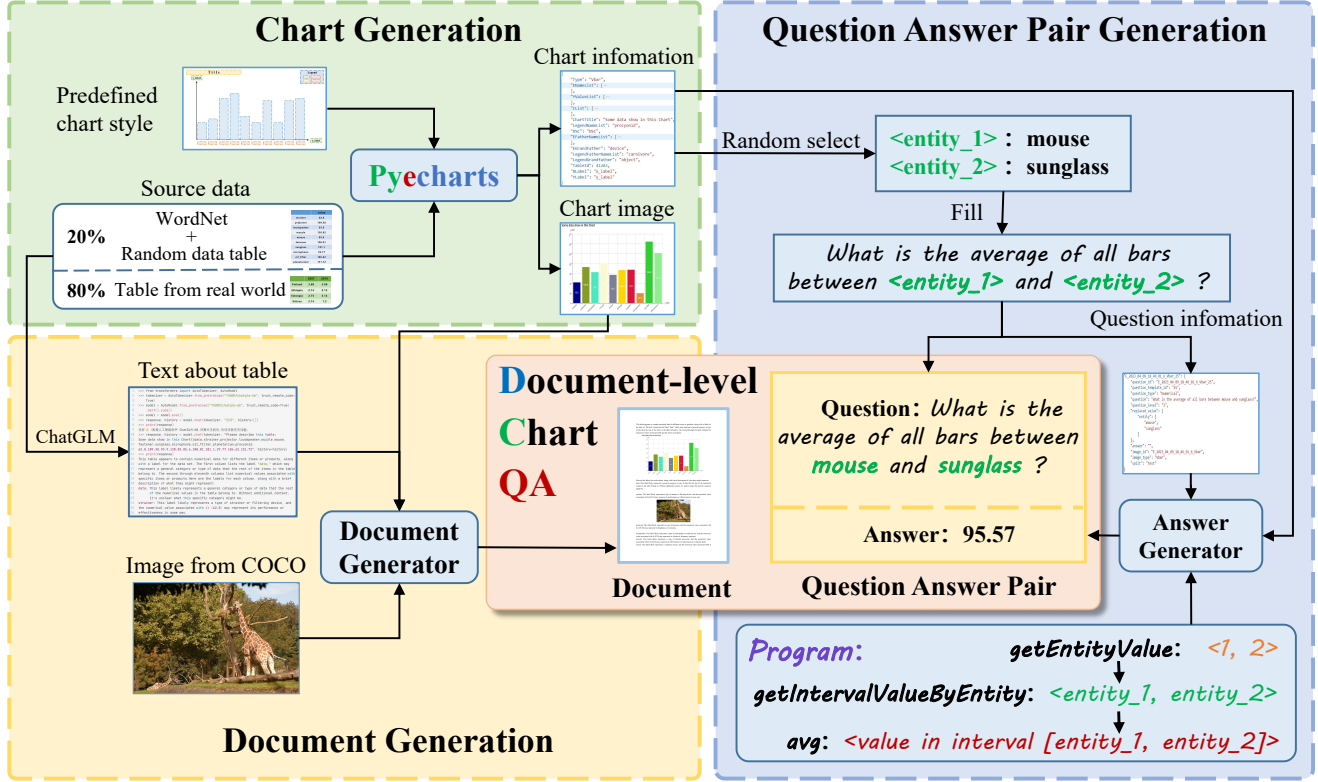
**Chart Generation.** The initial step of the process involves chart generation, whereby predefined chart styles are used for specific chart types. These styles are pre-configured using Pyecharts and the corresponding source data tables are prepared for each chart. The source data table has a 20% chance of being randomly generated and an 80% chance of being sourced from real-world tabular data. The chart in Figure A1 are generated using randomly generated tabular data. The entity names of the table are extracted from the hierarchical entity library adopted in this paper. The specific extraction method is described in detail in Appendix C.1. Upon the availability of the source data table and the predefined chart styles, the Pyecharts tool is employed to generate the chart and retrieve specific chart information, as demonstrated in Figure A2, which presents the intricate details of the chart.

<sup>4</sup><https://wordnet.princeton.edu>

**Question Generation.** The question generation involves the preparation of question templates and chart information. As previously discussed, we distill and abstract the 5730 questions into 324 templates. Concretely, the template structure adheres to the following standard: *"How many bars are greater than < value > in the < entity >?"*. The placeholders represented by *< entity >* are interchanged with the various graphical elements such as legends, X/Y ticks, the color of the legends/ticks and etc. The selection of values designated by *< value >* is achieved via randomization that relates to the extent of range within the X/Y axis of the present chart. Based on the type of replaceable modules in the question template, chart-specific questions are generated by randomly selecting elements from the range of selectable replaceable module values in the chart information and filling them in the replaceable modules of the question template. The corresponding question information, as depicted in Figure A3, can then be obtained. For instance, consider the template *"What is the average of all bars between < entity\_1 > and < entity\_2 >?"* depicted in Figure A1. The two replaceable modules, *< entity\_1 >* and *< entity\_2 >*, require substitution with specific entity names. To accomplish this, we randomly select two entity names from the entity name list provided in the chart information, namely "mouse" and "sunglass", respectively. Subsequently, we replace the placeholders with the chosen entities, resulting in the question *"What is the average of all bars between mouse and sunglass?"*. The average length of the questions is 11.46 words.

**Answer Generation.** The process of obtaining answers requires the provision of input parameters for solution derivation, which encompass the following parameters: (1) Image ID of the chart; (2) Question ID; and (3) Specific values to be filled into the replaceable modules within the question. The answer generation is achieved by designing solution steps and atomic operations. With this customized set of procedures, all answers for questions are created automatically without any human annotation, alleviating the probable errors and yielding cost-saving benefits. Moreover, we rigorously test the procedure for answer generation, with checks performed on each function and category of question, guaranteeing the correctness of the generated answer. The process of answer generation necessitates specific solution steps based on the question template, atomic operation functions that correspond to each step, filled-in replaceable values from the question information, and chart information. The atomic operation functions are invoked sequentially according to the question-solving steps, with the necessary parameters passed to the functions. These parameters are derived from the chart information and the replaced values of the substitutable modules for the question, resulting in the calculation of the answer. As an illustration in Figure A1, consider the question *"What is the average value of all bars between mouse and sunglass?"*, according to the solving procedure, the first step involves calling the `getEntityValue` function with arguments 1 and 2, which indicate the retrieval of the replacement values for *< entity\_1 >* and *< entity\_2 >*, respectively corresponding to "mouse" and "sunglass". Subsequently, the `getIntervalValueByEntity` function is executed with "mouse" and "sunglass" as its parameters to retrieve the values of all bars located between these entities from the chart information, resulting in a list of values [85.6, 100.01, 101.1]. Ultimately, the function `avg` is invoked with the parameter [85.6, 100.01, 101.1] to compute the mean of all values in the list, yielding the answer to the query, which is 95.57.

**Document Generation.** The document generation process entails



**Fig. A1:** The generation workflow of DCQA. A larger version of Figure 2 is available in the body of the paper for a more detailed view of the process.

the preparation of chart images, corresponding source data tables, and images from the COCO dataset. Initially, ChatGLM is employed to generate text content related to the table, which is then utilized as the primary textual information for the document. During the document generation process, based on LaTeX, a section of the document’s y-axis is randomly selected as the next range to be filled. Then, either text, chart images, or images from the COCO dataset are randomly chosen and inserted into the selected section. The aforementioned process is repeated until the entirety of the y-axis range in the document is filled. For documents with two or three columns, the above process is executed sequentially for each column to complete the document generation. As the document is progressively populated, the positional coordinates of each element within the document are systematically recorded and stored as annotation information within an XML document. This meticulous process allows us to obtain documents enriched with granular positional annotations. Further details regarding document information can be found in Appendix D.

## B. CHART INFORMATION

Chart information is stored in dictionary, where each chart’s unique ID serves as the key. This allows for easy retrieval of chart information using its ID. An illustration of chart information is presented in Figure A2.

Figure A2 shows the chart information of a grouped horizontal

bar chart. The chart ID of this chart, L\_2023.04.08.18\_48.12.8.GHbar is composed of the Machine ID\_ Generation time(Year)\_Generation time(Month)\_Generation time(Day)\_Generation time(Hour)\_Generation time(Minute)\_Generation time(Second)\_Random(0-9)\_Chart type used to generate the chart. The remaining content constitutes the chart information, including the chart type, title, entity name list (i.e., axis label), legend label list, data, element color name and color value, entity’s parent class list, entity’s grandparent class, legend label’s parent class list, legend label’s grandparent class, x-axis title, y-axis title, the table ID used to generate the chart, and some additional information (DSC). The additional information is provided for special chart types and is used to store some information not typically included in conventional charts, such as the value of marker lines in Marker Single Line and the range of each highlighted interval in Interval Highlight Single Line.

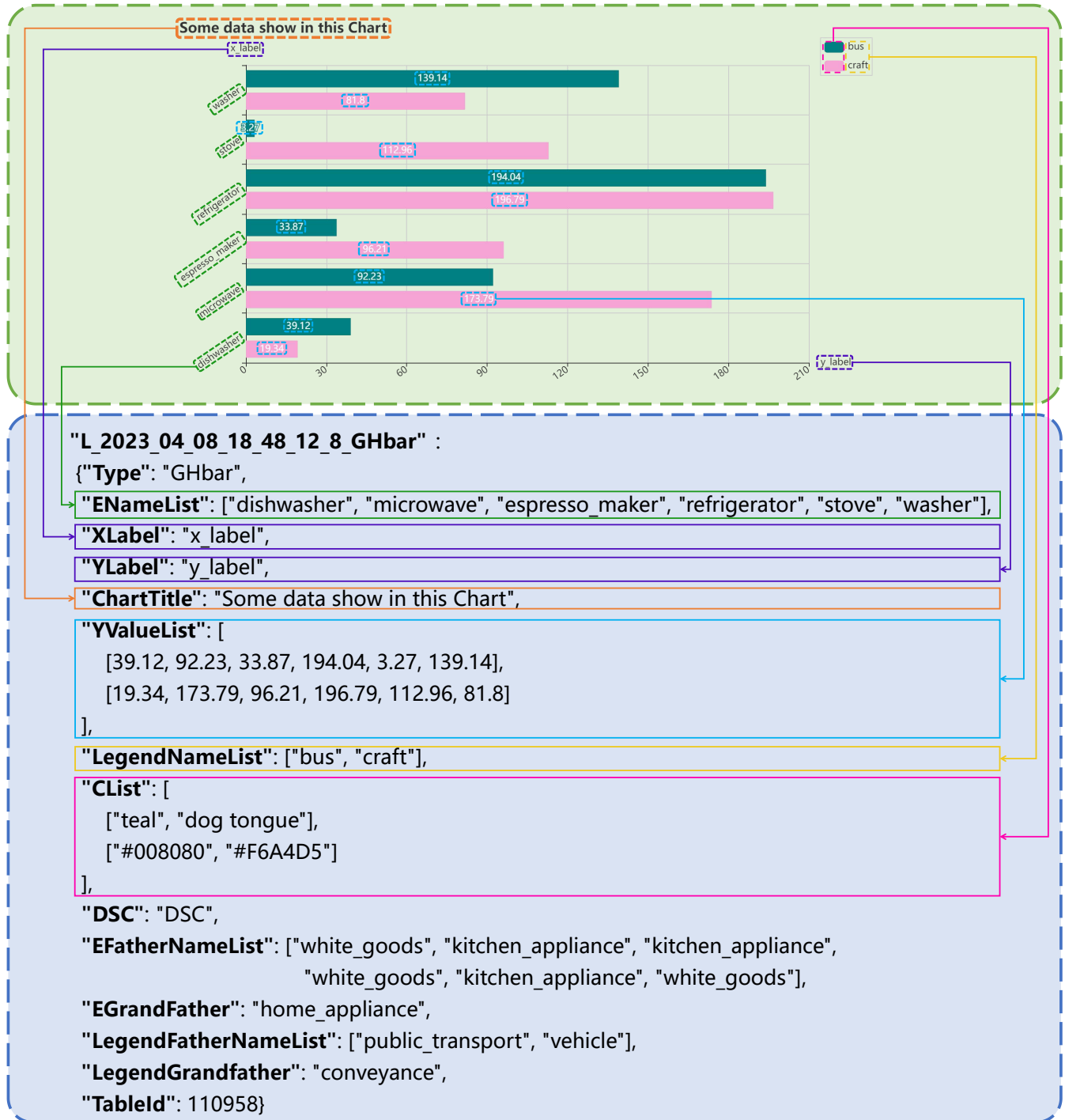
## C. QUESTION GENERATION

This section aims to provide a comprehensive introduction to two crucial components in the field of question generation, namely the hierarchical entity database and question information.

### C.1. Hierarchical Entity Database

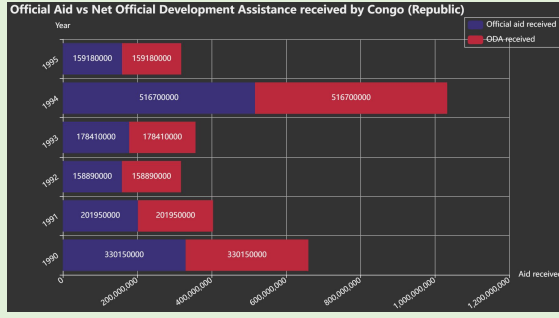
The hierarchical entity database presented in this paper is a tree-like structure constructed from the labels found in ImageNet and





**Fig. A2:** An example of chart information. The information includes the chart's title, entity names, legend labels, values, color names, and corresponding color values, etc.

**Chart : "I\_2023\_04\_06\_00\_32\_33\_1\_SHbar"**



**Question:** Is the maximum value of ODA received less than 194705746.02?

#### Question information:

```
"I_2023_04_06_00_32_33_1_SHbar_41" : {
  "question": "Is the maximum value of ODA received
less than 194705746.02?",
  "answer": "No",
  "image_id": "I_2023_04_06_00_32_33_1_SHbar",
  "question_id": "I_2023_04_06_00_32_33_1_SHbar_41",
  "image_type": "SHbar",
  "question_type": "Numerical",
  "question_template_id": "41",
  "question_level": "3",
  "replaced_value": {
    "legend": ["ODA received"],
    "value": [194705746.02] },
  "split": "test"}
```

**Fig. A3:** An example of question information. The information is composed of the question, answer, question ID, corresponding chart ID, template ID used for the question, difficulty level, and replaceable values, among other details.

WordNet, where the nodes with parent-child relationships represent hierarchical dependencies, i.e., a parent node includes child nodes and the parent node is the parent class of the child node. Specifically, the construction of a hierarchical entity database are as follows: (1) We opt for labels from ImageNet as fundamental entities, the hypernym query function in WordNet is employed to collect the hypernyms (i.e., parent nodes and ancestor nodes) of all the fundamental entities. (2) Entities with only a single child node among their hypernyms are excluded, and their parent and child nodes are connected ( $A \rightarrow B \rightarrow C \Rightarrow A \rightarrow C$ ) to eliminate redundant entities. (3) Finally, we obtain a directed acyclic graph (DAG), as some nodes may have multiple parent nodes. In this case, we retain only one parent node for these nodes and remove edges to other parent nodes. We save this entity repository as a tree structure in json file, which comprises noun words from WordNet that are utilized in the ImageNet dataset and their related words. All of the entities for generated data are randomly picked up from this hierarchical entity database. An example of the hierarchical entity database is provided in Appendix C.1. When generating charts, multiple parent entities under the same grandparent class are selected, and multiple child entities are then extracted from each parent entity as entity names or legend labels. The entity names/legend labels and their corresponding parent and grandparent entities are recorded. When generating questions from randomly generated data, the replaceable modules in the templates are randomly selected from the parent entities for filling. For example, in the template "What is the sum of < legendsort > in < entitysort >?", if encountering < entitysort >, a parent entity name is randomly selected from the collection of parent entities of the entity name to fill in < entitysort >. If encountering < legendsort >, a parent entity name is randomly selected from the collection of parent entities of the legend label to fill in < legendsort >.

Figure A4 illustrates a part of the hierarchical entity database in this paper, where "organism" has three subclasses: "animal," "vascular\_plant," and "fungus," and each of these subclasses has its own subclasses. To extract the desired entities from the hierarchical en-

tity database, a grandfather class is first selected, which in this case is "vertebrate" as shown in Figure A4. Then, nodes are selected from the grandfather class's child nodes to serve as parent classes, and finally, the subclasses of the selected parent classes are identified as the required entities. The entities and their corresponding parent and grandfather classes are recorded in the chart information, as depicted in Figure A2. Entities extracted from the hierarchical entity database can serve as entity names or legend labels in charts. They are extensively utilized in the chart generation, question formulation, and answer creation for charts generated from random data sources, and are present throughout the entirety of the dataset generation process.

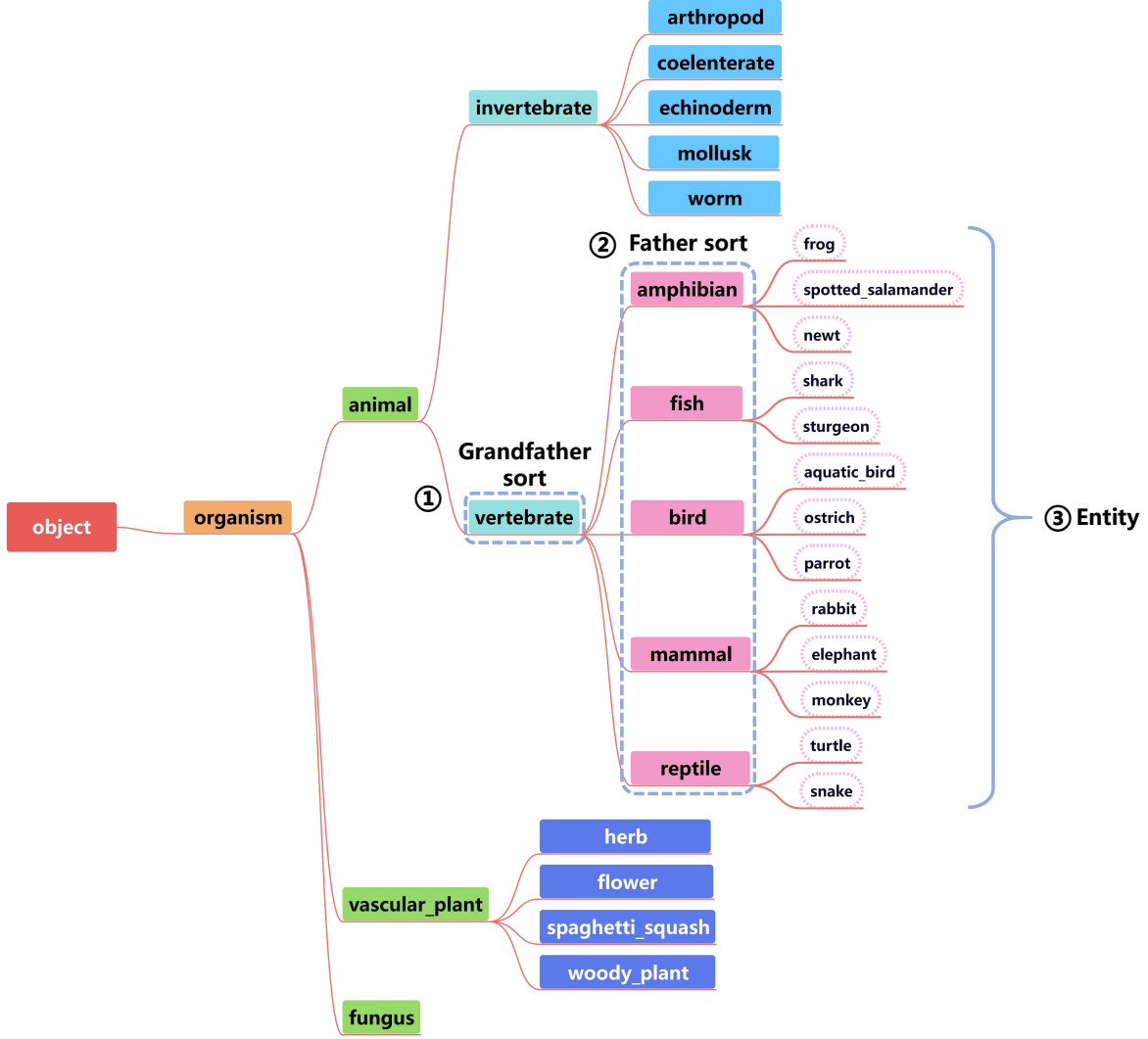
## C.2. Question Information

In a manner similar to the chart information, each piece of question information is systematically structured as a dictionary value, with a unique ID for each question serving as the key. This key-based organization enables the retrieval of information associated with any specific question by accessing its corresponding ID. Figure A3 illustrates an instance of question information along with its corresponding chart.

Figure A3 displays an illustrative example of a stacked horizontal bar chart pertaining to an extremum question, with a question ID of "I\_2023\_04\_06\_00\_32\_33\_1\_S-Hbar\_41" generated by concatenating the chart ID and the question template ID. The dictionary value comprises comprehensive details about the question, including the question statement, answer, corresponding chart ID, chart type, question ID, question type, question template ID, question difficulty, fill-in values for replaceable modules in the question, and the dataset partition. These extensive details cater to the diverse requirements of researchers.

## D. DOCUMENT INFORMATION

The information within a document is stored in XML format, where each element's details are organized using tags. The elements within



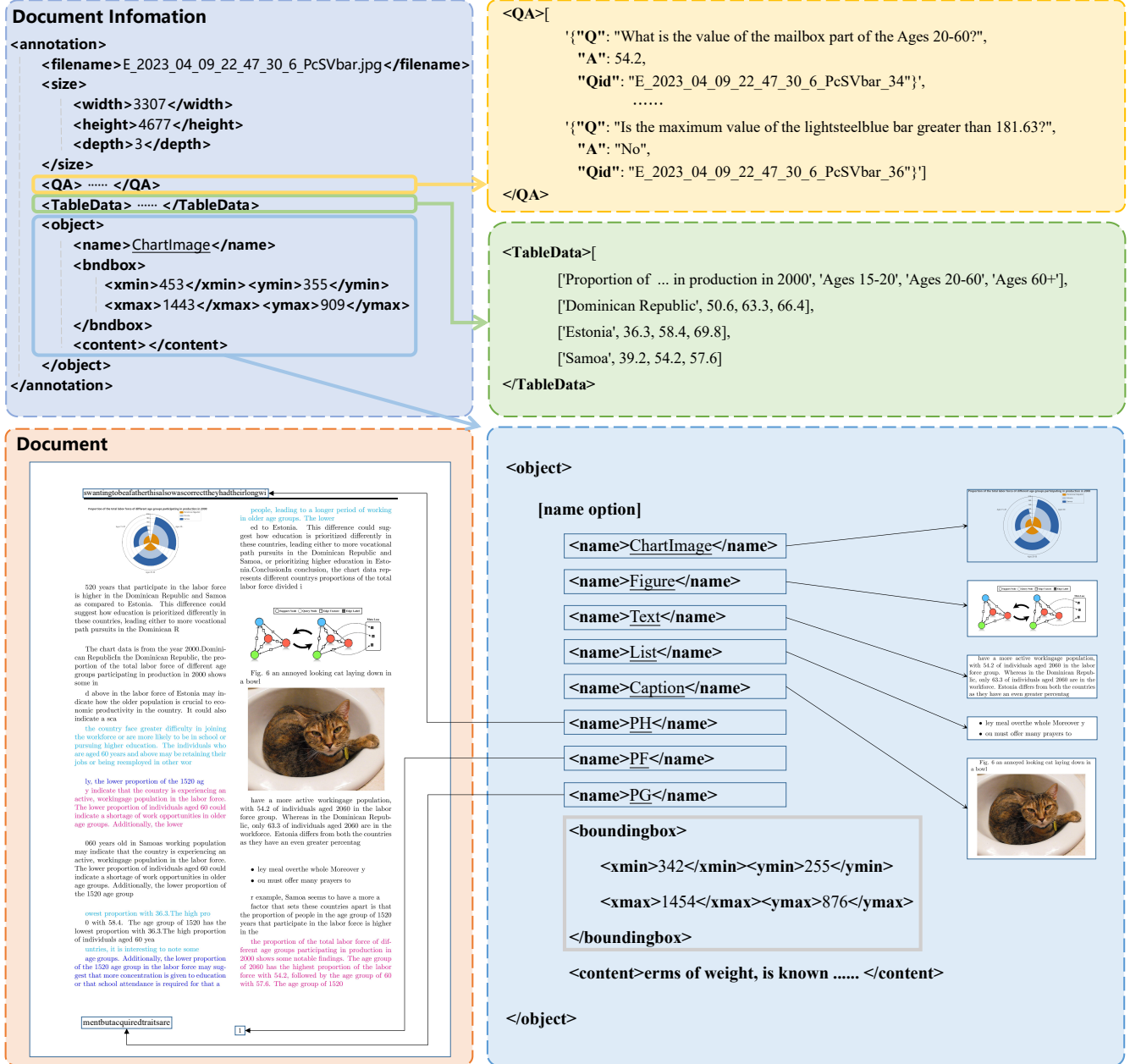
**Fig. A4:** An example of our hierarchical entity library. This example is a part of the hierarchical entity library used in this paper. Vertebrates are selected as the grandparent node, and the children nodes under vertebrates are chosen as parent nodes. The resulting children nodes are used as entity names for the chart.

the document can be categorized into various types, such as chart images, regular images from the COCO dataset, text information pertaining to tables, textual content arranged as ordered or unordered lists, image descriptions, headers, footers, and page numbers. Each document element is encapsulated within an `< object >` tag and comprises `< name >`, `< boundingbox >`, and `< content >` components. The `< name >` tag specifies the element's classification, while the `< boundingbox >` indicates its position within the document. The position is determined by the pixel coordinates of the element's top-left and bottom-right corners, with the top-left corner of the document serving as the reference point. The `< content >` tag is utilized to store textual information associated with text elements, remaining empty for non-textual elements. In addition to the document elements, each document also retains the tabular data and all corresponding question-answer pairs related to the associated chart, as depicted in Figure A5.

## E. DOCUMENT EXAMPLES

The DCQA dataset encompasses a broad spectrum of chart types, encompassing 6 primary categories: bar chart, line chart, pie chart, scatter plot, box plot, and combination chart. Each category is further subdivided into a total of 30 subtypes. In this section, we will present a comprehensive enumeration of each subtype, accompanied by illustrative examples depicted in Figures A6 to A10.

**Bar Chart.** Within the DCQA dataset, the bar chart category consists of 11 distinct subtypes: namely Vertical Bar, Horizontal Bar, Stack Vertical Bar, Stack Horizontal Bar, Group Vertical Bar, Group Horizontal Bar, Polar Coordinates Vertical Bar, Polar Coordinates Horizontal Bar, Polar Coordinates Stack Vertical Bar, Polar Coordinates Stack Horizontal Bar, and Waterfall Bar. Each subtype of the bar chart is exemplified in Figure A6, accompanied by a corre-



**Fig. A5:** An example of document information. The document image is situated in the lower-left corner, while a summarized version of the document information is located in the upper-left corner. For each highlighted section, detailed information is presented within a corresponding dashed box, labeled with the same color as the corresponding arrow.

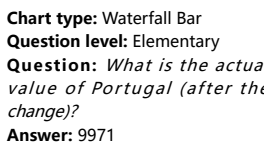
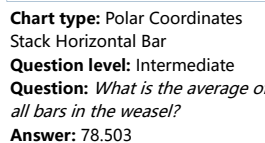
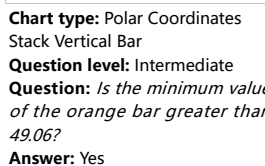
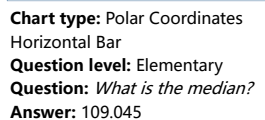
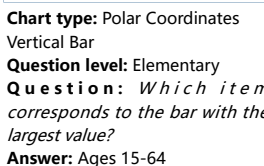
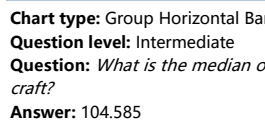
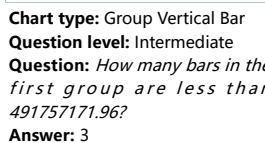
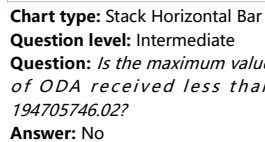
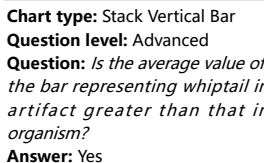
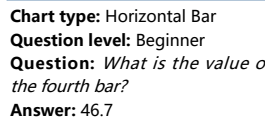
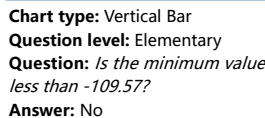
sponding question-answer pair and an assigned difficulty level. The objective of this presentation is to improve understanding and perception of the dataset's style and level of difficulty.

**Line Chart.** The line chart category primarily encompasses 7 subtypes, namely Single Line, Smooth Single Line, MultiLine, Marker Single Line, Best Value Single Line, Best Value MultiLine, and Interval Highlight Single Line. Each subtype of the line chart is exemplified by an example depicted in Figure A7.

**Pie Chart.** Pie chart comprise 4 main subtypes, specifically Simple Pie, Ring Pie, Rose Pie, and Nesting Pie. Examples of each pie chart subtype is provided in Figure A8.

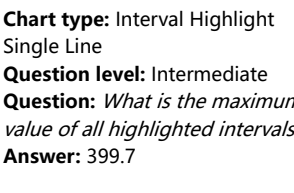
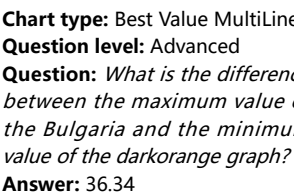
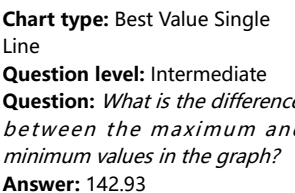
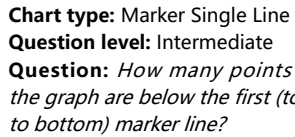
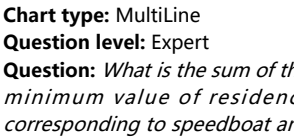
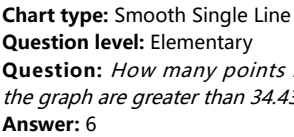
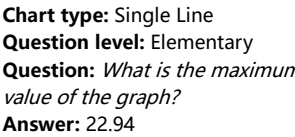
**Scatter Plot.** Scatter plot comprise 4 main subtypes, namely Simple Scatter, Multi Scatter, Bubble Scatter, and Check Bubble Scatter. An example of each subtype is illustrated in Figure A9.

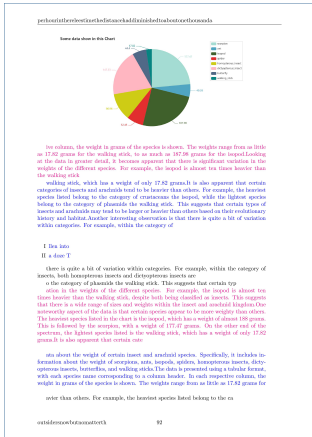
**Box Plot and Combination Chart.** Box plot encompass 3 distinct subtypes, namely Vertical Boxplot, Horizontal Boxplot, and Multi Boxplot. In contrast, combination chart comprise a single subtype, specifically Line Bar. An illustrative example for each subtype of the box plot and combination chart is presented in Figure A10.



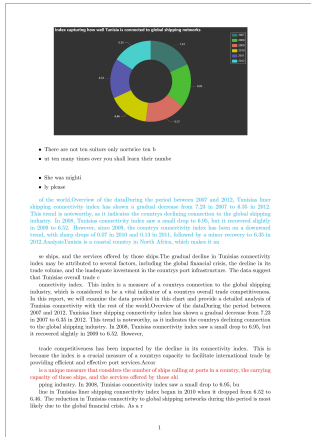
20



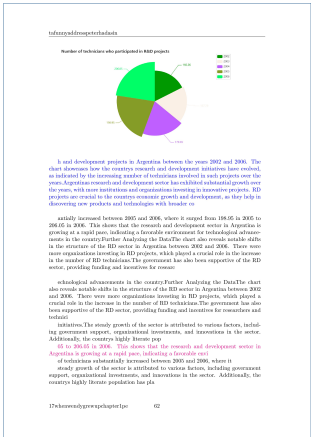




**Chart type:** Simple Pie  
**Question level:** Elementary  
**Question:** Which sector is the second largest?  
**Answer:** scorpion



**Chart type:** Ring Pie  
**Question level:** Intermediate  
**Question:** How many times is 2009 compared to 2011?  
**Answer:** 1.03

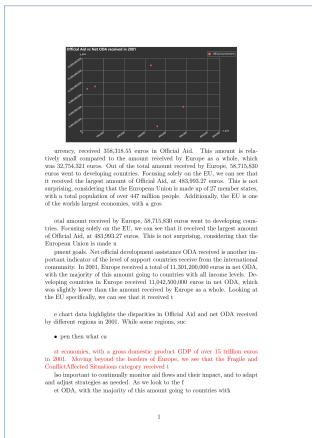


**Chart type:** Rose Pie  
**Question level:** Intermediate  
**Question:** What is 2005's ranking?  
**Answer:** 2

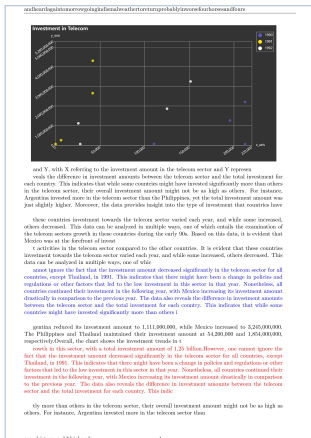


**Chart type:** Nesting Pie  
**Question level:** Expert  
**Question:** How many times is the sum of snake in the outer circle greater than the sum of arthropod in the inner circle?  
**Answer:** 0.45

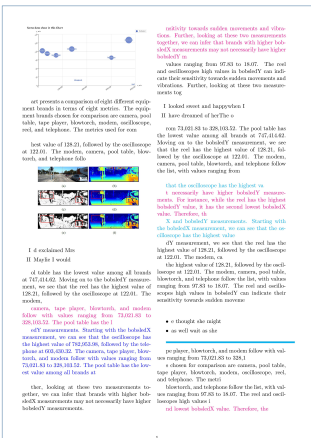
**Fig. A8:** Examples of all subtypes of pie chart in the generated document. Each chart is accompanied by a question and corresponding answer.



**Chart type:** Simple Scatter  
**Question level:** Elementary  
**Question:** Are more than 30% of the scatter points concentrated in the upper half of the chart?  
**Answer:** Yes



**Chart type:** Multi Scatter  
**Question level:** Beginner  
**Question:** How many colors are the points in the scatter chart?  
**Answer:** 3

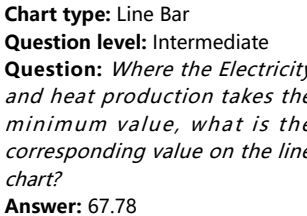
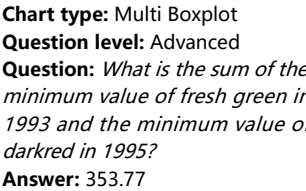
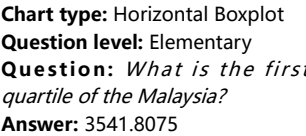
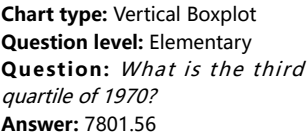


**Chart type:** Bubble Scatter  
**Question level:** Intermediate  
**Question:** Does the scatter that takes the maximum value belong to photographic equipment?  
**Answer:** No



**Chart type:** Check Bubble Scatter  
**Question level:** Intermediate  
**Question:** At which time of the working day does the sum of the number of check-ins is the largest?  
**Answer:** 9pm

**Fig. A9:** Examples of all subtypes of scatter plot in the generated document. Each chart is accompanied by a question and corresponding answer.



23