

**CHART QUESTION ANSWERING WITH VISUAL AND LOGICAL
REASONING**

AHMED MASRY

A THESIS
SUBMITTED TO THE FACULTY OF GRADUATE STUDIES
IN PARTIAL FULFILMENT OF THE REQUIREMENTS
FOR THE DEGREE OF MASTER OF SCIENCE

GRADUATE PROGRAM IN
ELECTRICAL ENGINEERING AND COMPUTER SCIENCE
YORK UNIVERSITY
TORONTO, ONTARIO
JUNE 2022
© AHMED MASRY, 2022

Abstract

Charts are very popular for analyzing data. When exploring charts, people often ask complex reasoning questions that involve several logical and arithmetic operations. They also commonly refer to visual features of a chart in their questions. However, most existing datasets do not focus on such complex reasoning questions as their questions are template-based and answers come from a fixed-vocabulary. In this thesis work, we present a large-scale benchmark covering 9.6K human-written questions and 23.1K questions generated from human-written chart summaries. To address the unique challenges in our benchmark involving visual and logical reasoning, we present transformer-based models that combine visual features and the data table of the chart. Moreover, we propose chart-specific pretraining tasks that improve the visual and logical reasoning skills of our models. While our models achieve the state-of-the-art results on the previous datasets and our benchmark, the evaluation also reveals several challenges in answering complex reasoning questions.

Dedication

I would like to dedicate this thesis to my parents, **Gamal Masry** and **Salwa Masry**, who instilled a love of science, learning, and knowledge in me since a very young age. I am very thankful to them for supporting me with all possible means to excel in my education journey and achieve my goals. I hope this achievement will fulfill the dreams they envisioned for me. May Allah bless them with good health and endless happiness.

Acknowledgements

First of all, I would like to express my sincerest gratitude to my thesis supervisor, **Prof. Enamul Hoque Prince**, for his endless support from the very first day and throughout my masters degree. His valuable feedback and suggestions helped me much to improve the quality of this thesis work and publish parts of it in several prestigious venues. I am very thankful to him for teaching and guiding me to conduct high-quality research in ML/NLP. It was a great honor for me to work with him for two years.

I would like to thank **Prof. Aijun An** for being a committee member of my thesis and spending her valuable time on it.

I also want to thank **Prof. Shafiq Joty**, who is one of the co-authors of our ACL papers, for his great feedback and comments. I would like to also thank **Do Xuan Long and Jia Qing Tan** for their collaboration in my research.

I would like to thank **Prof. Baris Akgun** who is an assistant professor at Koc University in Turkey where I completed my BSc in Computer Engineering degree. I started my very first research experience with him during my time at Koc University, and

learned much from him about AI/ML. I want to also thank him for supporting me during the York MSc. Computer Science program application process and recommending me.

Finally, I would like to thank Compute Canada for providing the computational resources which were needed to run the experiments in my thesis.

Preface

This thesis work was mainly done by **Ahmed Masry** under the supervision of **Prof. Enamul Hoque Prince**. Parts of this thesis work has been accepted for publication at prestigious venues as described below. In the publications where I am listed as the first author [1, 2], I was responsible for proposing and implementing the methodology, setting up the experimental setup, running the experiments, and analyzing the results. My co-authors (**Enamul Hoque** and **Shafiq Joty**) assisted in editing the manuscripts and providing valuable feedback to improve the methodology and experimental setup. Moreover, **Do Xuan Long** and **Jai Qing Tan** assisted in the data collection and cleaning process as well as running some of the experiments. In the other papers [3,4] where I am listed as a co-author, I assisted in the dataset collection and annotation process, running some of the experiments, and editing the manuscripts.

1. **Ahmed Masry**, Do Xuan Long, Jia Qing Tan, Shafiq Joty, Enamul Hoque, “ChartQA: A Benchmark for Question Answering about Charts with Visual and Logical Reasoning”, accepted in Findings of the 59th Annual Meeting of the Association for Computational

Linguistics (ACL'22 Findings) 2022.

2. **Ahmed Masry**, Enamul Hoque Prince, “Integrating Image Data Extraction and Table Parsing Methods for Chart Question answering,” Chart Question Answering Workshop, in conjunction with the Conference on Computer Vision and Pattern Recognition (**CVPR**), 2021 [**Best Paper Award**]
3. Shankar Kantharaj* ¹, Rixie Tiffany Ko Leong*, Xiang Lin*, **Ahmed Masry***, Megh Thakkar*, Enamul Hoque, Shafiq Joty, “Chart-to-Text: A Large-Scale Benchmark for Chart Summarization”, accepted in Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL'22) 2022.
4. Enamul Hoque, Parsa Kavehzadeh, **Ahmed Masry**, “Chart Question Answering: State of the Art and Future Directions”, In EuroVis 2022.

¹* Equal contribution. Listing order is based on the alphabetical ordering of author surnames

Table of Contents

Abstract	ii
Dedication	iii
Acknowledgements	iv
Preface	vi
Table of Contents	viii
List of Tables	xii
List of Figures	xiv
1 Introduction	1
1.1 Motivation	1
1.2 Our Approach	3
1.3 Contributions	6

1.4	Outline	7
2	Literature Review	8
2.1	Question Answering	8
2.2	Chart Question Answering	12
2.2.1	Datasets	12
2.2.2	Models	14
2.3	Chart Data Extraction	17
2.4	Transformer-based Models	19
2.4.1	Transformer	19
2.4.2	BERT	22
2.4.3	TaPas	25
2.4.4	T5	27
2.5	Discussion	28
3	Chart Question Answering Benchmark	30
3.1	ChartQA Dataset	30
3.1.1	Data Collection and Preparation	30
3.1.2	Dataset Annotation	31
3.1.3	Dataset Analysis	36
3.2	ChartQA Methodology	39

3.2.1	Problem Formulation	39
3.2.2	Data Extraction	40
3.2.3	Models	42
3.3	ChartQA Evaluation	47
3.3.1	Datasets	47
3.3.2	Baselines	48
3.3.3	Evaluation Metrics	49
3.3.4	Results	50
3.3.5	Ablation Studies	54
3.3.6	Qualitative Analysis	56
3.4	Discussion	57
4	Transfer Learning for Chart Question Answering	59
4.1	Motivation	59
4.2	Problem Definition & Input Representation	60
4.3	Chart-T5 Model	61
4.4	Pretraining	65
4.4.1	Tasks	65
4.4.2	Pretraining Dataset	67
4.5	Evaluation	68

4.5.1	Datasets, Baselines & Metrics	68
4.5.2	Results	69
4.5.3	Qualitative Analysis	70
4.6	Discussion	73
5	Conclusions and Future Work	75
5.1	Conclusions	75
5.2	Future Work	77
	Bibliography	79
A	Appendix	89
A.1	Pretraining Tasks Templates	89

List of Tables

2.1	Comparison between existing datasets and our new ChartQA benchmark	12
3.1	Sample question answer pairs generated from human-written summaries in Statista.	35
3.2	Our dataset statistics for each split.	35
3.3	Number of charts from each source. Statista-H and Statista-M refer to the datasets with human-written and machine generated questions respectively from Statista	36
3.4	ChartQA benchmark statistics.	37
3.5	Distribution of questions types of among 300 randomly chosen human written questions (blue-colored tokens make visual references to the chart).	37
3.6	Usage of visual references in visual questions among 300 randomly chosen questions	38
3.7	Accuracies of our data extraction algorithm on the test sets of DVQA, PlotQA, and ChartQA. Since the gold data table is not available in FigureQA, we report the results on the Validation2 set.	51

3.8	Evaluation results for different models. For DVQA, we have reported the results with and without using Oracle for OCR. We do not evaluate on FigureQA test sets with the gold data table setup since they do not have ground data tables.	52
3.9	Accuracy of the different models on our benchmark. VisionTaPas [†] does not support difference and ratio operations. VisionTaPas [*] and VL-T5 [*] are trained on PlotQA and evaluated directly on ChartQA.	55
3.10	Results for VisionTaPas and VL-T5 on the ChartQA test set by chart type.	55
3.11	Accuracies of VisionTaPas and VL-T5 on the ChartQA-H test set by question type on 200 random samples.	55
4.1	Flattened visual data table.	64
4.2	Input and Output formats for our pretraining tasks. The input consists of the task prefix and query which we later concatenate with the chart visual data table.	65
4.3	A comparison between our best performing model from Chapter 3, VL-T5 Pretrained and the Chart-T5 model. Pretrained refers to the models that are first pretrained on PlotQA and then finetuned on our benchmark, ChartQA.	69

List of Figures

1.1	Sample questions in our benchmark.	2
2.1	The attention mechanism used the in the transformer model. This picture is taken from the original "Attention is all you need" paper [68].	20
2.2	The transformer model architecture. This picture is taken from the original "Attention is all you need" paper [68].	21
2.3	BERT [12] model architecture and two phases: pretraining and finetuning.	23
2.4	BERT Input Embeddings: Token, Segment, and Positional.	23
2.5	TaPas [20] model architecture. The output of the [CLS] token is passed to a classification layer that selects the operator. Also, the outputs from the table cell tokens are passed to a cell selection layer that selects the relevant cells. The predicted operator is then applied on the selected cells.	25
2.6	T5 [54] text-to-text multi-task finetuning framework for different NLP tasks.	28

3.1	The user interface for the annotation task	33
3.2	Distribution of topics in the datasets.	38
3.3	Our approach for question answering over charts. If not provided, the underlying data table is first extracted from the chart image using ChartOCR. We then pass the extracted data table in addition to the question and the image features to the ChartQA model where the ChartQA model represents one of the following: TaPas, VisionTaPas, T5, and VL-T5.	39
3.4	Data Extraction Process	41
3.5	Data extraction examples from OWID and Pew.	41
3.6	T5 and VL-T5 neural models for ChartQA. Data tables are first flattened and fed into the model along with the question (and visual features in VL-T5).	42
3.7	TaPas and VisionTaPas models. TaPas adds positional embeddings to the tokens to encode the tabular structure of the data table. VisionTaPas uses a cross-modality encoder to combine visual features from ViT and outputs from TaPas encoders.	44
3.8	Example of errors from VisionTaPas	56
3.9	Sample outputs of our model VisionTaPas on our new ChartQA test set. Answers in green are correct and answers in red are incorrect.	58

4.1	A bar chart image with its corresponding data table and visual data table. The data table consists of three main cell types: data values (green), column headers (yellow), and row names (blue). Each one of these cells is then extended to include the visual features so that we can construct the visual data table. We obtain the visual features of each cell by cropping its relevant element (e.g., bar) from the chart image. Each cell in Figure 4.1c contains the cropped image and the underlying data value of each chart element (separated by a vertical black line). They also contain the cell type (e.g., x-axis label, bar, legend mark, .etc)	62
4.3	Our approach to obtain the Visual Token of the chart elements. The chart element (e.g., bar) is first cropped from the chart image. Then, we pass it through a CNN network to extract the visual features from it. After that, we add the visual features and bounding box vector to obtain the final Visual Token that encodes both the visual features and positional information (bounding box)	64
4.4	Sample outputs of the Chart-T5 and VL-T5 (Pretrained) models on our new ChartQA test set. Answers in green are correct and answers in red are incorrect.	71
4.5	Some charts with many data points from the ChartQA-M dataset.	72
4.6	Example of errors from the Chart-T5 model.	74

1 Introduction

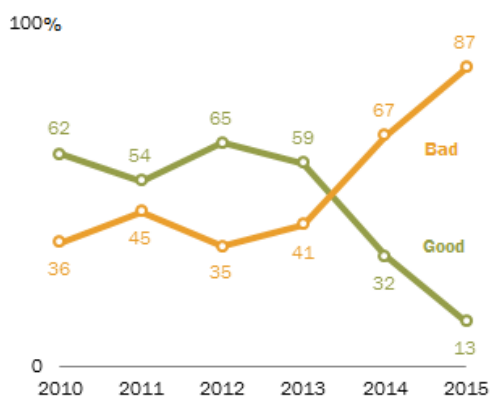
1.1 Motivation

Data visualizations such as bar charts and line charts have become popular in analyzing data and making informed decisions. To analyze data, often people ask complex reasoning questions about charts involving arithmetic and logical operations [31]. Answering such questions requires a significant amount of perceptual and cognitive efforts as people need to combine multiple operations such as retrieving values, comparing values, finding maximum, calculating sums and differences of values. For example, the question Q1 in Figure 1.1 requires the user to compute the differences between the two lines for each year and find the year with the highest difference.

The goal of a Chart Question Answering (ChartQA) system is to help users by taking a chart and a natural language question as input and predicting the answer. This task differs from other QA tasks such as QA on texts [55] and tables [50] because the input for ChartQA is a visual representation of data that can draw a reader's attention to various prominent features such as trends and outliers [31, 32]. Also, people tend to ask questions

Rapid Decline in Brazilians' Assessment of Economy

Current economic situation in Brazil is ...



Q1: Which year has the most divergent opinions about Brazil's economy?

Answer: 2015

Q2: What is the peak value of the orange line?

Answer: 87

Figure 1.1: Sample questions in our benchmark.

by referring to visual attributes of marks. For example, in Figure 1.1, Q2 refers to the color of a mark ('line') and its attribute ('peak') in the chart.

While the task of ChartQA has received growing attentions in recent years, existing datasets have several major limitations: (i) the questions are generated automatically using pre-defined templates [28, 26, 6, 64] which lack naturalness, (ii) the charts are created automatically using a programming tool like Matplotlib [64] which do not reflect the diverse styles of many real-world charts, and finally, (iii) in most datasets, the answer comes from a small fixed sized vocabulary (*e.g.*, chart axis labels, 'yes', 'no'), ignoring many complex reasoning questions where the answer is derived through various mathematical operations such as aggregation and comparison. Since most datasets only support *fixed vocabulary* questions, existing models usually treat the task as a classification problem and rely on dynamic encoding techniques with the questions and answers encoded in

terms of spatial positions of chart elements (*e.g.*, *x-axis-label-1*). Such approaches do not work when the OCR model generates errors or when the question refers to chart elements using synonyms (*e.g.*, US vs. United States). PlotQA [44] attempts to support *open vocabulary* questions by applying a TableQA model [50] but it does not consider any visual features of a chart which are critical for answering visual reasoning questions. To sum up, most existing datasets are synthetic, lack naturalness, and do not reflect the plethora of the questions or the charts’ visual styles in the ChartQA task. Consequently, the previously proposed approaches are tweaked to handle limited aspects of the task (*e.g.*, *fixed vocabulary* questions) and do not generalize well. Moreover, they are mostly adapted from relevant domains (*e.g.*, Visual Question Answering and TableQA), so they do not fully utilize the underlying structure, visual properties, and data of the chart images.

1.2 Our Approach

This thesis work aims to address the limitations of both the existing datasets and approaches. Our first goal is to build a ChartQA dataset with human-authored natural questions and real-world charts that capture a large span of the charts’ visual styles on the Web. Our second goal is to propose a novel approach that utilizes both the visual features and the underlying data values of the chart images to comprehend the chart images’ underlying structure and reason about them.

With respect to the first goal, we have developed a large-scale benchmark covering 9,608 human-written questions focusing on logical and visual reasoning questions. Since human annotations are costly, we also generated another 23,111 questions automatically from human-written chart summaries using a T5 model [54] and manually validated a subset of it for quality assurance. In this way, we collect a large number of questions automatically while maintaining rich variations in language as they were generated from human-written summaries. Our benchmark consists of 20,882 charts which are curated from four different online sources to ensure variety in visual styles and topics.

Our benchmark introduces several challenges which existing approaches may fail to address. First, most of our answers are Open-Vocabulary which may require aggregating the underlying data values, so classification-based models' performance is expected to fall dramatically. Second, our human-authored questions use synonyms and usually address the charts' textual elements in different ways (e.g., GDP vs. Gross Domestic Product), which dynamic encoding-based approaches can not identify. Third and most importantly, our questions involve complex reasoning and visual references to charts, which TableQA-based models can not handle since they do not consider the visual features of the chart images. To address these challenges, we propose novel approaches that consider both the chart images' visual features and underlying data values. Furthermore, we are interested in exploring various task setups by utilizing different representations of charts in the input to our models. Is the underlying data table of the chart sufficient to answer the given

question? How can we effectively combine the charts' visual features and data table in the input to our models? What is the optimal chart input representation that best conveys the underlying structure and properties of the charts to our models?

To answer these questions, we develop a pipeline approach that first extracts the underlying data table from the chart image by adapting the ChartOCR model [40] as well as the visual features from the chart image using neural models. Then, we adapt two transformer-based QA models where we utilize both the extracted data table and visual features of the charts in a unified way. Moreover, we extend this approach by proposing a novel input representation of charts that reflect the chart images' underlying structure by associating the visual features of the charts elements with their corresponding data values in a coherent way. Transfer learning led to several breakthroughs in the natural language processing field [12, 52, 53, 5] as well as the vision domain [13, 39]. Pretraining such models on large corpora of data using self-supervised objectives can teach them general features and knowledge that can be adapted by finetuning the models on the downstream tasks. This inspired us to explore transfer learning for the ChartQA task by proposing a set of chart-specific pretraining tasks that can inject several numerical and visual reasoning skills into our models which are critical in the ChartQA task. Finally, we show the effectiveness of our models by comparing them with the previous models on the previous datasets and our newly created benchmark.

1.3 Contributions

In this thesis work, our contributions can be summarized as follows:

- We introduce a large-scale ChartQA dataset with real-world charts and human-authored question-answer pairs. Moreover, we augment our dataset with machine-generated questions from human-written chart summaries using the T5 model [54].
- We propose an approach that utilize both the visual features and automatically extracted data from charts in the inputs to transformer-based QA models that provide state-of-the-art results.
- We extend the state-of-the-art charts data extraction model, ChartOCR [40], to output the fully-structured data table from the charts images which we utilize in the input to our QA models.
- We further extend and improve our ChartQA approach by combining the visual features and the underlying data table of the chart images into a unified format in the input representation to our model. Moreover, we define a set of pretraining tasks that inject numerical and visual reasoning skills into our model which in turn improves its performance.
- We perform an extensive analysis and evaluation of the performance of our models.

1.4 Outline

The rest of the thesis is structured as follows. Chapter 2 gives a literature review of the chart question answering task, the existing ChartQA datasets, ChartQA models, and Chart Data Extraction. In the same chapter, we discuss several transformer-based models (e.g., BERT, TaPas, T5) which we utilize in our approach. Chapter 3 describes our proposed approach to the ChartQA task and our newly created ChartQA dataset collection process. Moreover, we also present the evaluation metrics, the baselines, the datasets we use for evaluation, and our results analysis. In Chapter 4, we discuss our transfer learning follow-up work that extends and improves our approach to the ChartQA task as well as our pretraining tasks. Finally, in Chapter 5, we discuss our conclusions and propose directions for future work.

2 Literature Review

This chapter begins with a brief introduction to the question answering task over the different data formats (texts, images, tables, charts, .etc). Then, we conduct a literature review on the chart question answering task by discussing the previous datasets and models limitations. In addition, we discuss the chart data extraction approaches and the possible areas for improvements. Finally, we conduct a literature review on the transformer-based models that we utilize to build our ChartQA approach.

2.1 Question Answering

Question answering is a natural language processing task in which models provide relevant answers to questions about a given context. Unlike traditional information retrieval, where the systems return relevant and complete documents to the given queries, question answering requires the systems to output concise and coherent answers after comprehending the given context. Due to the rapidly growing amount of data on the Web, the task of question answering over text has received growing interest in the past few years.

Many datasets and benchmarks have been introduced to investigate the key challenges of the task such as SQuAD[55], CNN/Daily Mail [19], WikiQA[72], and TriviaQA [17]. While researchers have proposed several deep learning models architectures to tackle this task, transformer-based models [54, 12, 70, 36, 15] are considered the best performing approaches so far.

Apart from text, question answering has been applied to several data formats (e.g., data tables, images, videos, etc.). Table Question Answering (TableQA) is a natural language processing task in which a model produces an answer to a given question about a given structured data table. Due to its growing interest in the past few years, several data-sets have been proposed (e.g., WikiTQ [49], WikiSQL [77], SQA [22], TabFact[14]) that reveals the challenges of the task. Earlier attempts [49, 77, 11] treated the TableQA task as a semantic parsing task in which models transform the given question to a logical query and executes it over a given data table . However, obtaining such logical forms annotations can be costly and error-prone. In contrast, other recent approaches employ weakly-supervised models that only rely on the final answer as supervision and do not require logical forms to produce the answers [46, 20, 74, 75].

Visual Question Answering (VQA) is a multidisciplinary task that combines natural language processing and computer vision techniques to answer questions about images. An example question is “what is the color of the man’s shirt?” about an image in which a man is wearing a shirt with a specific color. Such questions require understanding

the visual content (colors, objects, positions) and reason about it. Several datasets have been introduced to explore the various challenges in this domain including, but not limited to, COCO-QA [38], DAQUAR [41], and VQA [1]. Moreover, several neural models have been proposed to address the challenges of this multi-modal task. VQA models [65, 39, 7, 37, 76] usually consist of encoder modules to encode the textual input (question) and visual input (image) before aggregating the extracted features to select the final answer from a finite set of answers.

Chart Question Answering (ChartQA) is quite different from the above-mentioned QA tasks. The goal of chart question answering is to answer a natural-language question about a chart image to facilitate data analysis. Unlike TableQA, where the data is presented in tabular format, charts present the data in visual representations that highlight notable features (e.g., trends, outliers). Moreover, the questions in the ChartQA task usually make visual references to the data elements and their visual properties in the charts. Consequently, many questions may not be answerable from the underlying data table only (e.g., “What is the value of the blue bar?”). While the task of ChartQA may be viewed as a subarea of the VQA task, there are several core differences between them. First, the vocabulary used in VQA questions refers to immutable semantic concepts across the different images [26]. For example, in the above-mentioned question, “what is the color of the man’s shirt?”, the meaning and visual features of the word “shirt” are usually consistent across the different images that contain shirts. However, in ChartQA, the

questions' vocabulary usually refers to the chart's textual elements (*e.g.*, x axis labels, legend labels) with arbitrary texts. The same texts may refer to different visual elements (*e.g.*, bars, line points, pie segments) and different data values in different charts.

Moreover, charts have much richer semantic content that the models need to reason about than the regular images used in the VQA datasets [26]. For example, the question “what is the color of the man's shirt?” can be easily answered by directly checking the color tones and features of the man standing in the image. In contrast, the question “What is the difference between Google's and Apple's annual revenue?” requires the models to perform a multi-stage reasoning process. First, the model needs to locate the labels of “Google” and “Apple” in the chart. Then, it should identify their corresponding value encoding marks (*e.g.*, bars) and map their visual channels (*e.g.*, length) to the underlying data values. Finally, the model takes the difference between the two data values. Consequently, ChartQA models should be able to extract the underlying data from the chart images, comprehend the charts' underlying structure and visual properties, and have the capability to reason about those extracted data and visual features of the charts.

Datasets	Question Types	Answer Types	Real-world Data	Real-world Charts	#Charts/ #QA pairs
FigureQA [28]	Template-based	Fixed	✗	✗	180K/2.3M
DVQA [26]	Template-based	Fixed	✗	✗	300K/3.4M
LEAF-QA [6]	Template-based	Fixed	✓	✗	240K/2M
LEAFQA++ [64]	Template-based	Fixed	✓	✗	244K/2.5M
PlotQA [44]	Template-based	Open	✓	✗	224K/28M
ChartQA-H (ours)	Human-authored	Open	✓	✓	4.8K/9.6K
ChartQA-M (ours)	Machine generated	Open	✓	✓	17.1K/23.1K

Table 2.1: Comparison between existing datasets and our new ChartQA benchmark

2.2 Chart Question Answering

2.2.1 Datasets

An overall comparison between the existing dataset is shown in Table 2.1. Kahou et al. [28] first introduced the FigureQA dataset, where the questions were generated using 15 templates (yes/no questions). Kafle et al. [26] then designed the DVQA dataset, consisting of 25 templates covering different questions types: structural, data retrieval, and reasoning (and had a fixed number of answers, 1576). However, the charts’ data were randomly generated in both datasets and the charts were synthetically plotted using a software (Matplotlib). While FigureQA uses the X11 named color set ² for the charts’ textual labels (*e.g.*, *legend-label*), DVQA uses the most frequent 1000 words in the Brown Corpus. Moreover, FigureQA has 2.3M QA pairs and 180K charts covering 3 different chart types: bar, line, and pie charts. In contrast, DVQA consists of 3.4M QA pairs and 300K bar charts.

²<https://cgkit.freedesktop.org/xorg/app/rgb/tree/rgb.txt>

LeafQA [6] was the first dataset to utilize real-world data to plot their chart images. Similar to the previous datasets, they automatically generated their QA pairs using 35 pre-defined templates and 75 fixed answers. The dataset consists of 2M QA pairs and covers more chart types (*e.g.*, scatter-plots, box-plots). LeafQA++ [64] is an extended version that introduces more templates, 75, and more style variations in the plotted chart images (*e.g.*, grid lines). It also balances the distribution of the answers so that common answers do not dominate the dataset and lead to over-fitting of the models. A major limitation of all these datasets is that the answers come from a fixed set of vocabulary (Fixed-Vocab). To address this limitation, Methani et al. [44] constructed the PlotQA dataset, which contains template-based questions that require applying mathematical operations on the aggregated data values of the chart image (*e.g.*, what is the difference between the highest and lowest revenue?).

The questions in all these datasets are automatically generated using hand-crafted templates, and the charts are synthetically plotted using the same software (*e.g.*, Matplotlib). Consequently, the questions do not reflect the language variations in the human-authored questions, and the charts are limited in terms of visual styles and variations compared to the real-world charts found on the Web. Moreover, they do not focus on the visual reasoning questions, which humans usually ask to draw useful insights and conclusions from the charts.

Kim et al. [30] also conducted a formative study to analyze how people ask questions

about charts where they collected a small-scale set of 629 human-authored QA pairs and 52 chart images. Given a chart image, they asked the workers to ask a question and provide both the answer and the explanation. According to their findings, humans tend to refer to the visual marks of the charts in their questions as well as their explanations. This highlights the importance of the visual reasoning questions which the existing datasets usually lack. To the best of our knowledge, there is no large-scale ChartQA dataset with human-authored questions, focusing on visual reasoning, and real-world charts which motivated us to build our dataset.

2.2.2 Models

There are two main approaches for ChartQA. The first approach utilizes classification-based visual question answering (VisualQA) models [28, 26, 6, 27, 64] that consist of three main modules: image encoder, question encoder, and attention or aggregation mechanism followed by a final classification layer. Kahou et al. [28] utilized relation networks [58] to reason about the relations between the chart elements. They extract the features of the chart elements using a CNN, and embed the question using an LSTM before feeding them to the relation network. Kafle et al. [26] and Chaudhry et al. [6] also adapted the Stack Attention Networks [73] by employing dynamic encoding techniques where the question and answers are encoded in terms of the spatial information of the chart textual labels (*e.g.* *x-axis-label-1*). Kafle et al. [27] proposed the PReFIL model

that fuses the image and question low-level and high-level features in parallel to learn bi-modal embeddings, which helps the model perform multi-stage reasoning processes in a semi-recurrent fashion. The embeddings are then aggregated and projected into a final classification layer. STL-CQA [64] was the first work to apply transformer-based models on the ChartQA task adapted from several vision-language transformer-based approaches [65, 13, 39]. STL-CQA model consists of three transformer-encoder modules: question, chart image, and reasoning encoders.

Despite the promising results and progress in this VisualQA-based approach, there are three main drawbacks. First, it does not utilize the underlying data table of the chart, which can be pivotal when answering the data dependent reasoning questions (*e.g.*, *What is the average of the annual expenses?*). The second drawback is that it can only handle the Fixed-Vocab questions which have a limited number of answers since it relies on classification-based models. Finally, they mostly utilize dynamic encoding techniques which are prone to the OCR noise and does not work if the question uses synonyms to address the chart elements (*e.g.*, GDP vs. Gross Domestic Product).

The second approach applies table question answering (TableQA) models by utilizing the underlying data table of the chart image [44, 42, 30]. Kim et al. [30] and Methani et al. [44] employ the SEMPRE model which needs to translate the question into a logical query as an intermediate step in retrieving the answer. Generating such logical forms can be difficult (*e.g.* costs for collecting logical forms and label bias problem) and answers

can be limited (e.g. cannot generate 'yes'/'no' answers). Consequently, Masry and Hoque [42] employed TaPas [20] which is trained in a weakly-supervised fashion using the final answer as the only supervision. TaPas selects the relevant cells from the data table and applies an aggregation operation (SUM, COUNT, AVERAGE, NONE) on the selected cells. While Kim et al. [30] and Masry and Hoque [42] assume the gold data table of the chart image is given, Methani et al. [44] proposed a pipeline approach to extract the underlying data table from the chart image using computer vision techniques. Although this TableQA-based approach can successfully handle the Open-Vocab questions that requires applying aggregation/logical operations, it has two main drawbacks. First, it does not consider the visual features of the chart images (*e.g.*, color, positions) when answering the questions. Consequently, it fails when the questions make complex visual references to the chart components (*e.g.*, grey bar, leftmost point). Second, this approach requires access to the underlying data table of the chart image which is not accessible in most real-world charts on the Web. Despite the previous attempts to extract the data from the charts [40, 24, 44], the performance is still quite limited.

To sum up, current ChartQA approaches utilizes either the chart image features or the underlying data table in their inputs, but not both. Such input representations drop important features of the charts. In our work, we combine both input representations (data table and visual features) to guide our model to address the complex charts' questions, especially the visual reasoning questions.

2.3 Chart Data Extraction

Early attempts were extracting data from chart images, usually in a semi-automatic manner, using a set of hand-crafted rules [59, 24]. Savva et al. [59] previously presented the ReVision system to re-plot poorly-designed charts. Their system consists of three stages: charts classification, data and marks extraction, and finally re-plotting the chart data in higher quality format using a visualization tool. Although their system can classify ten types of charts, it can extract the data only from bar and pie charts. Moreover, their system can not parse legends, so they overlook the important legend labeling information in complex bar charts (e.g. grouped or stacked). They assume that the text information in the images are given and their data extraction system's accuracy is also quite low. In addition to that, they have constructed their own corpus of over 2500 chart images and made it publicly available³. Jung et al. [23] presented the ChartSense tool which determine the chart type and extract the raw data in a semi-automatic manner by relying on the user input and assistance to specify and detect different components in the chart image (e.g. axis labels and positions). Moreover, they have augmented the Revision corpus with additional chart images collected from the Web and used it in their experiments. Although their classification model works for 10 types of charts, their data extraction module is only limited to 6 types. FigureSeer [63] is an end-to-end framework that can detect the figures from research papers, classifies them, and extract the underlying data from line charts.

³<https://old.datahub.io/dataset/vis-revision-corpus>

Their data extraction process works in the following manner. They first detect and parse both the axes and the legend area by locating the text labels using OCR modules. Then, they formulate lines detection as an optimal-path finding problem and they differentiate between the different lines using similarity features (e.g. colors). Yet, their data extraction algorithms performed poorly.

Several web tools have also been designed to extract the data from the chart images semi-automatically [57, 43]. WebPlotDigitizer [57] is a semi-automatic open-source tool that relies on the user assistance. The user needs to locate several chart components before running the extraction algorithm (e.g. y-axis start and end point, bars top points, line points, etc). iVoLVER [43] is another data extraction tool that requires extensive user interaction to acquire several data elements from chart images (e.g. text, colors, quantities).

Poco and Heer [51] then designed the first fully-automated pipeline that extracts the visual encodings from the chart images. Their pipeline consists of the following phases: text detection and recognition, text role classification, chart type detection, and encodings induction. However, their approach is limited to extracting the visual encoding and does not extract the underlying data values. Following up on this work, Choi et al. [10] then extend their pipeline to extract the raw data from the chart images by detecting the visual marks (e.g., bars) and estimating their data values. Yet, their approach's performance is still limited and supports only three types of charts: line, simple vertical bar and pie

charts.

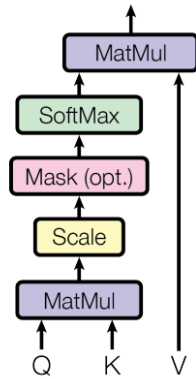
Luo et al. [40] then proposed a hybrid approach, ChartOCR, that extracts the raw data values using a combination of deep-learning and rule-based methods with high accuracy. They utilize key-point detection networks [33] to locate the chart key components (e.g., bars, line points) and estimate their values using the y-axis labels. Still, they only extract the raw data values of the chart encoding marks (e.g., bar, line) without associating them with their textual labels (e.g., *legend-label*, *axis-label*). In this thesis work, we extend their approach and utilize the state-of-the-art OCR model, CRAFT [3], to output the fully-structured underlying data table of the chart images.

2.4 Transformer-based Models

2.4.1 Transformer

Traditional Sequence-to-Sequence architectures mostly relied on RNN and LSTM to construct their encoders and decoders. LSTM and RNN are special types of neural networks in which the output from the previous step is passed in the input to the model along with the current step input. Hence, RNN and LSTM were suitable at handling sequences (e.g., texts). However, they suffer from two major limitations due to their recurrent nature. First, they are computationally expensive since the recursive process does not allow parallelism (the input to the model at each step is dependent on the

Scaled Dot-Product Attention



Multi-Head Attention

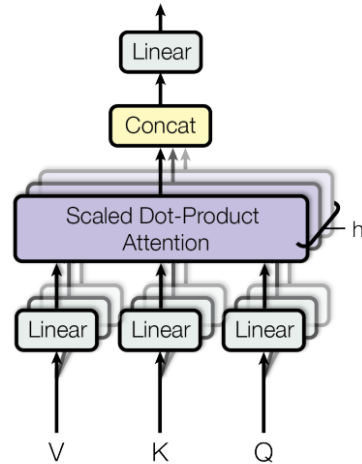


Figure 2.1: The attention mechanism used the in the transformer model. This picture is taken from the original "Attention is all you need" paper [68].

model's output at the previous step). Second, they are not efficient in handling long-term dependencies. Real-world texts are often long. Consequently, as the LSTM or RNN is going over the text token by token, it keeps dropping or forgetting the initial context and tokens and in the text.

To overcome these limitations, Vaswani et al. [68] introduced the Transformer model which is an encoder-decoder architecture that relies on the attention mechanism instead of recursion. The overall model architecture is shown in Figure 2.2.

Scaled Dot-Product Attention is the mainly used attention mechanism in the transformers model. As shown in Figure 2.1, the input is composed of three main components: queries (Q), keys (K), and values (V). First, the queries and keys are multiplied together. Second, the output is scaled by the dimension of the keys, d_k , and passes through a softmax layer.

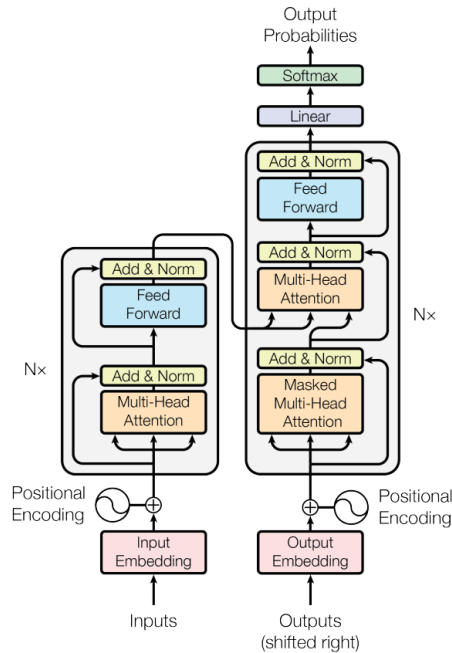


Figure 2.2: The transformer model architecture. This picture is taken from the original "Attention is all you need" paper [68].

Finally, the output from this softmax is multiplied with the values, V , to compute the final output as follows:

$$Attention(Q, K, V) = SoftMax\left(\frac{Q^T K}{\sqrt{d_k}}\right)V$$

Multi Head Attention projects the input queries, keys, and values to different linear projections which are utilized by different attention layers in parallel. The outputs from these attention layers are then concatenated and linearly projected to produce the final output as shown in Figure 2.1.

Transformer Encoder consists of several identical layers which are stacked on top of

each other. Each one of these layers contain two sub-layers: multi head attention and feed-forward layers which are both followed by a residual connection and layer normalization [2].

Transformer Decoder has a similar architecture to the encoder. It also consists of several stacked identical layers. However, the decoder layers contains an additional multi head attention layer in which the queries comes from the previous decoder layer and the keys and values come from the encoder output. This can allow the model to attend to the input sequence tokens while generating the output sequence tokens. Moreover, the attention layer output are also masked to prevent the decoder from attending to the next unseen tokens in the training process.

To account for the positional information of the input tokens, the transformer model uses a special type of embedding called positional encoding which are produced using a sinusoidal function. The resulting embeddings are then added to the original text embeddings and fed to the model as shown in Figure 2.2.

2.4.2 BERT

BERT (Bidirectional Encoder Representations from Transformers) [12] is a bidirectional language model that is based on the transformer [68] encoder. BERT's key advantage is that it's pretrained on large-scale unlabelled text data to learn bidirectional representations of the language words within their context. The overall model architecture is shown in

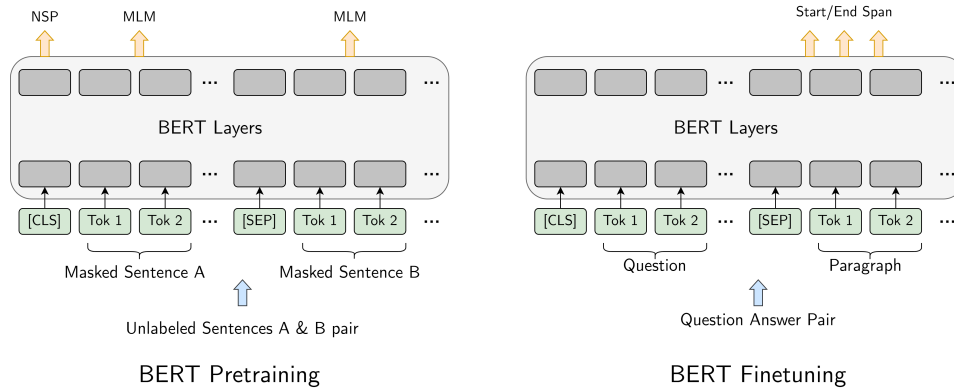


Figure 2.3: BERT [12] model architecture and two phases: pretraining and finetuning.

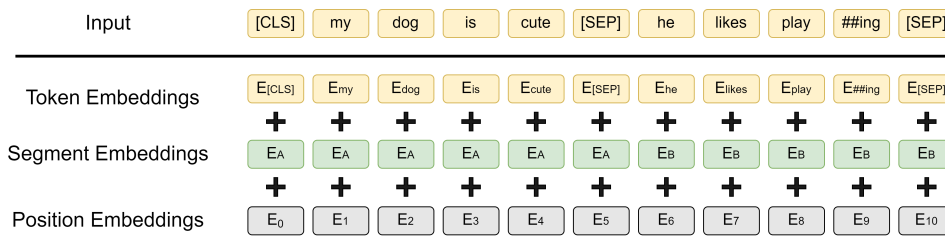


Figure 2.4: BERT Input Embeddings: Token, Segment, and Positional.

Figure 2.3. The input consists of two sentences which are separated by a special token [SEP]. Moreover, the first token in the input, [CLS], is also a special type of token that represents the overall context. In addition to the transformer model embeddings (token and positional), BERT [12] also adds another type of embedding called Segment Embedding which differentiate between the two input sentences as shown in Figure 2.4. The model has two phases: pretraining and finetuning.

Pretraining: Traditional language models are usually trained on the next word prediction task. Given a previous sequence of words (“Canada is located north of the ...”), the model is trained to predict the next word (“US”). To mimic this goal, BERT [12]

introduced the Masked Language Modeling task, **MLM**, in which the model predicts the masked input tokens. In Masked Language Modeling, 15% of the input tokens are masked randomly by replacing their token with a special token called [MASK]. The model is then pretrained to predict the original token of the masked one by feeding its embedding from the output to a classification layer which acts as the language modeling head. Moreover, in order to teach the model the relations between the two input sentences, BERT [12] uses another pretraining task called Next Sentence Prediction (NSP). Given two sentences A and B from some context, the model is pretrained to predict whether B is the following sentence to A or not. To achieve this goal, the output of the [CLS] token, which represents the overall context of the input, is fed to a final binary classification layer.

Finetuning: After pretraining BERT with the two above-mentioned tasks, the model can be easily adapted and finetuned on a variety of NLP downstream tasks by adapting its input and output. For example, for the extractive question answering task, the input to the model consists of the question and the paragraph, separated by the [SEP] token. In the output, we can train the model to predict the start/end tokens for the answer span in the paragraph. We can achieve this by replacing the language modeling head with another classification layer that predicts whether the token is start, end, or not.

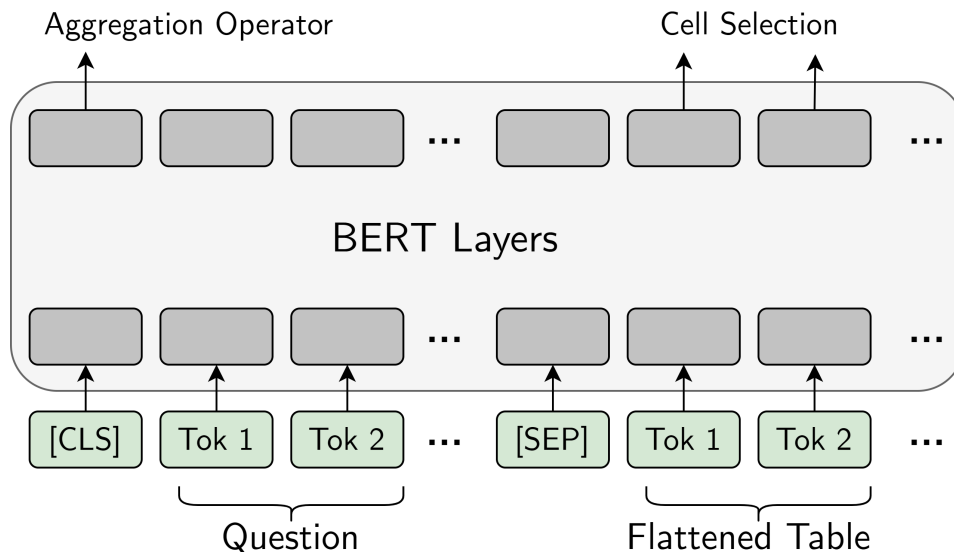


Figure 2.5: TaPas [20] model architecture. The output of the [CLS] token is passed to a classification layer that selects the operator. Also, the outputs from the table cell tokens are passed to a cell selection layer that selects the relevant cells. The predicted operator is then applied on the selected cells.

2.4.3 TaPas

Previous table question answering (TableQA) models [50] relied on parsing methods which need to translate the question into a logical query as an intermediate step in retrieving the answer from the data tables. Generating such logical forms can be difficult (e.g. costs for collecting logical forms and label bias problem) and answers can be limited (e.g. cannot generate 'yes'/'no' answers). To address this issue, TaPas [20] alleviates the need for logical forms generation by predicting the selection of relevant table cells and then applying an aggregation operator to such selection.

Tapas [20] is one of the state-of-the-art models in the Table QA literature. It is based

on the BERT architecture [12] and uses additional embeddings to encode the tabular data structure such as Position ID, Segment ID, Column/Row ID and Rank ID [20]. As shown in Figure 2.5, the model takes the flattened data table and the question as input. Similar to BERT, TAPAS adds the [CLS] token at the beginning and uses the [SEP] token to separate between the question and the table cells. The output embeddings of the table cells tokens are fed into a classification layer to select the relevant table cells to the question. Moreover, the output of the first token, [CLS] is fed into another classification layer which predicts the required mathematical operation (e.g. SUM, AVERAGE, and COUNT). The operator is then applied on the selected cells and the final answer is computed.

Since the TableQA datasets are usually small, pretraining TaPas [20] is essential to teach the model the correlations between the text and table content, the relations between the table cells, and the table structure. Inspired by BERT [12], TaPas [20] is also pretrained on a large-scale table-text pairs dataset collected from online sources (Wikipedia). The main pretraining objective of TaPas [20] is masked language modeling. Unlike BERT [12] where 15% of the input words are masked, TaPas [20] is pretrained to predict the masked table cells. The authors also tried different pretraining objectives such as Next Sentence Prediction in which the model predicts whether a given text is relevant to the table or not. However, it did not have any noticeable impact on the performance of the model.

2.4.4 T5

T5 [54] is an encoder-decoder model which unifies the NLP tasks as text-to-text generation using the same architecture and loss function. The architecture is almost identical to the original Transformer model [68], but it differs in three components: the positional embedding scheme, the layer normalization structure, and the order of the layer normalization and residual connection. Unlike the Transformer [68] which adds a positional embedding vector to each token in the input to indicate its position in the input sequence, T5 [54] uses relative position embeddings [62] to encode the different positional relations between the input elements in the keys and queries of the attention layers. Furthermore, T5 removes the bias from the layer normalization. It also places the layer normalization before the residual connection in the encoder and decoder layers.

Pretraining: T5 has been pre-trained on the Colossal Clean Crawled Corpus (C4) with a self-supervised denoising objective that mimics the Masked Language Modeling (MLM) task. Since the output of the T5 model is text, the T5 uses a different mechanism from BERT [12]. First, 15% of the input tokens are masked by replacing them with special sentinel tokens (e.g., Canada is located <X> of the <Y>) that identify each masked token. Then, the model is trained to predict the text that corresponds to each masked token (e.g., ”<X> north <Y> US”).

Finetuning: T5 is then finetuned in a multi-task setup on eight different downstream

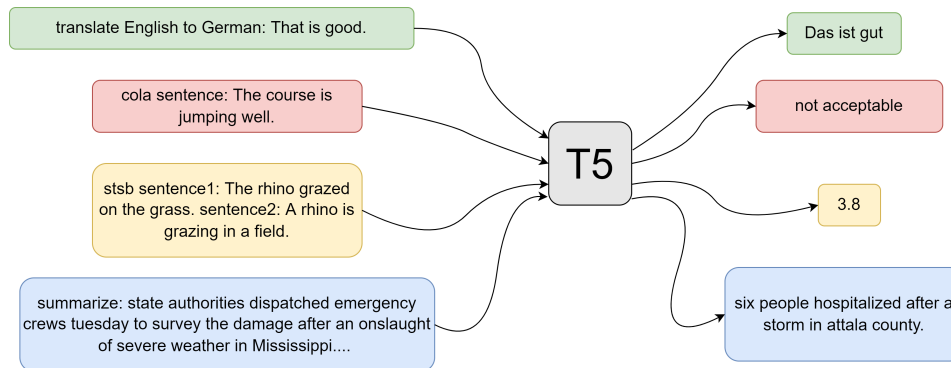


Figure 2.6: T5 [54] text-to-text multi-task finetuning framework for different NLP tasks.

tasks: Sentence acceptability judgment, Sentiment analysis, Paraphrasing/sentence similarity, Natural language inference, Coreference resolution, Sentence completion, Word sense disambiguation, and Question answering. The overall finetuning framework is shown in Figure 2.6. To distinguish each different tasks in the input, a prefix (e.g., "summarize:") is added at the beginning of the input. This enables the T5 model to unify all the NLP tasks and use the same architecture, parameters, and loss function for all the tasks.

2.5 Discussion

In this chapter, we discussed the chart question answering task and how it is different from other question answering tasks (e.g., TableQA and VQA). Moreover, we analyzed the previous ChartQA datasets limitations and showed that they mostly use synthetic chart images and template-based questions which lack naturalness. We also showed that the previous models either utilize the chart image (e.g., VQA models) or the data table of

the chart (e.g., TableQA), but not both, which limits their visual and logical reasoning capabilities. In addition, we discussed the previous chart data extraction approaches and how we extend them in our approach. Finally, we presented the transformer-based models which our ChartQA approach rely on and extend, particularly TaPas [20] and T5 [54].

3 Chart Question Answering Benchmark

In this chapter, we present two contributions. The first contribution is constructing a ChartQA dataset where the questions are human-authored and the charts are curated from different online sources. The second contribution is designing an approach to tackle the ChartQA task by combining both the visual features of the chart image in addition to the underlying data values.

3.1 ChartQA Dataset

3.1.1 Data Collection and Preparation

To ensure that our benchmark covers various topics and charts with a diverse range of styles, we crawled charts from four different sources: (i) Statista ([statista.com](https://www.statista.com)) is an online platform that presents charts covering a variety of topics including economy, politics, and industry. (ii) The Pew research ([pewresearch.org](https://www.pewresearch.org)) publishes report about social and economic issues, demographic trends and public opinion with a wide variety of charts. (iii) Our World In Data or OWID (ourworldindata.org) is another platform that contains

thousands of charts about different global issues such as economy, finance, and society. (iv) Organisation for Economic Co-operation and Development or OECD (oecd.org) is a global organization which shares reports and data analysis for policymaking about several global issues.

For the Pew dataset, we only crawled chart images since the underlying data tables are not available. For the other three, we extracted the underlying data tables, metadata (*e.g.*, title, chart type), SVG file and associate text description. Finally, we extracted the bounding boxes information of the different chart elements (*e.g.*, x-axis labels) from the SVG files to train our data extraction models. Moreover, we manually annotate 825 charts from Pew to get the required annotations for our models.

3.1.2 Dataset Annotation

We have two main annotations procedures: (i) collect human-authored QA pairs using Amazon Mechanical Turk (AMT) and (ii) generate QA pairs from the Statista human-written summaries.

- **Human-authored QA annotation** To create human-authored QA pairs, we designed an AMT task in which we asked the crowdworkers to focus on two types of questions for each chart image: compositional and visual questions. Compositional questions contain at least two mathematical/logical operations like *sum*, *difference* and *average*, while visual

questions refer to the visual attributes such as *color*, *height*, and *length* of graphical marks (*e.g.*, *bars*) in the chart. We focus on these two types of questions because people tend to ask them commonly [31, 21] and previous datasets mostly do not focus on such complex visual and logical reasoning questions. To ensure quality, we selected workers with an acceptance rate of 95% and total accomplished HITs of 5000. Moreover, we further filtered the workers by giving them a pre-test to select the best qualified workers for this task. The data collection interface is shown in Figure 3.1. While presenting the chart, we ensure that the data labels of chart elements are visible to workers so that they can accurately perform the necessary arithmetic and logical operations to provide and answer the questions successfully.

For each chart, the workers provide two questions with the answers. The same questions are then answered by another annotator. If both workers' answers exactly match, we consider the answer to be correct. Otherwise, we manually check the answers to select the final correct answer. Overall, the agreement between the crowd workers based on exact matches was 61.04%.

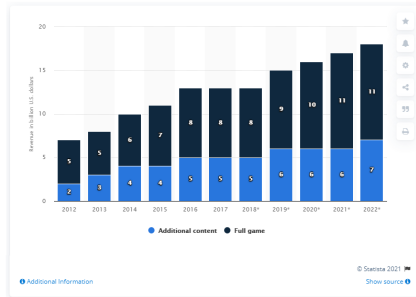
• **Dataset Augmentation** Prior work on QA has performed data augmentation by either creating template-based or machine generated questions, *e.g.*, for visual QA [25] and textual QA [35]. Template-based questions generally lack rich linguistic variations. On the other hand, large-scale language models like T5 [54] which are trained on very large

Instructions

[View instructions](#)

First Section. Answer the given questions!

Chart Title: PC and console games revenue worldwide from 2012 to 2022, by type (in billion U.S. dollars)



Question
In which year the highest revenue was generated by selling PC and console games ??

Please write the answer

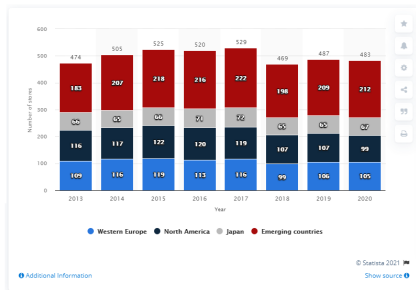
Answer

Question
In which year the difference between selling of additional content and full game was maximum??

Please write the answer

Answer

Chart Title: Number of directly operated Gucci stores worldwide from 2013 to 2020, by region



Question
Which region had maximum number of Gucci stores in 2016??

Please write the answer

Answer

Question
What is the difference between maximum number of Gucci stores operated in Emerging Countries over the years and minimum number of stores operated in Western Europe over the years??

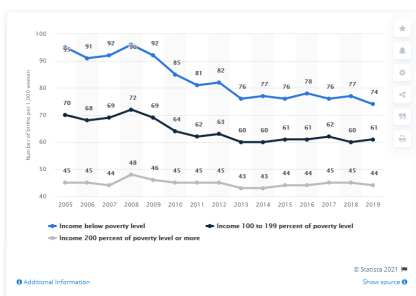
Please write the answer

Answer

Second Section. For each chart, ask one visual question and one compositional question as explained in the instructions and provide the answers.

While there are some Q/A examples, they are merely for the purpose of tutorial. When you create your own questions, please try to be creative by forming questions that are preferably different from the given Q/A examples in terms of operations/wordings

Chart Title: Birth rate in the United States from 2005 to 2019, by poverty status



Please ask a visual question that refers to the visual attributes of graphical marks in the chart

Question

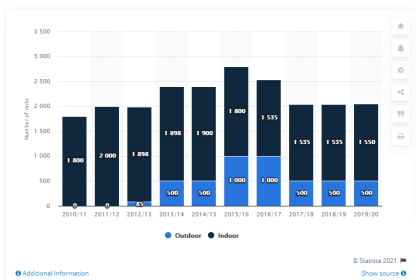
Please write the answer to the question

Please ask a compositional question that requires at least two mathematical/logical operations

Question

Please write the answer to the question

Chart Title: Total number of ice hockey rinks in the United States from 2010/11 to 2019/20



Please ask a visual question that refers to the visual attributes of graphical marks in the chart

Question

Please write the answer to the question

Please ask a compositional question that requires at least two mathematical/logical operations

Question

Please write the answer to the question

Figure 3.1: The user interface for the annotation task

data from various web sources can learn general linguistic properties and variations [4]. Therefore, we opt for the latter.

Specifically, we fine-tune a pre-trained T5 model on the SQuAD QA dataset [55] and apply to the human-written chart summaries that come with the charts from Statista to automatically generate questions that are human-like with sufficient lexical and syntactic variations. The process involves training and applying two T5 models: one for *answer extraction* and the other for answer-aware *question generation*. For answer extraction, the T5 model is trained to generate possible answers separated by [SEP] token given the textual summary as input (*i.e.*, trained on SQuAD’s *passage* → *answer* pairs). For question generation, the proposed answer is first concatenated with the summary in the format: Answer: Answer Context : Chart Summary. Then, the T5 model is trained to generate a question from the given question using the chart summary. This model is trained on SQuAD’s (*passage, answer*) → *question* pairs. Since the summaries are human-written, the generated questions are similar to the human-authored questions as shown in Table 3.1.

However, the T5 question generation model may still generate invalid questions because of the mismatch in training and test domains. We notice that some questions are either incomplete or not answerable from the chart (*e.g.*, ‘What province includes Cape Town?’ is not answerable because it requires knowledge outside of the chart). To filter out such invalid questions, we developed a simple heuristic where we filter out the question

Question Type	Human-written Summary	Generated Question	Answer
Compositional	Cancer was the leading cause of death among state prisoners in the United States, which killed 1,137 state prisoners in 2018. Heart disease was the second leading cause of death in that year, accounting for 1,052 deaths.	What was the second leading cause of death among state prisoners in 2018?	Heart disease
Compositional	This statistic shows the number of tourist arrivals at accommodation establishments in Latvia from 2006 to 2019. Since 2009 there has been an increasing trend in arrivals.	Since what year has there been an increasing trend in arrivals?	2009
Data Retrieval	The statistic shows the youth unemployment rate in the Gambia from 1999 to 2019. According to the source, the data are ILO estimates. In 2019, the estimated youth unemployment rate in the Gambia was at 12.44 percent.	What was the youth unemployment rate in the Gambia in 2019?	12.44 percent
Data Retrieval	This statistic shows the total population of Portugal from 2016 to 2020, with projections up until 2026. In 2020, the total population of Portugal was at approximately 10.29 million inhabitants.	In what year did Portugal's population reach 10.29 million?	2020

Table 3.1: Sample question answer pairs generated from human-written summaries in Statista.

Split	ChartQA-H		ChartQA-M	
	Charts	Questions	Charts	Questions
Training	3,699	7,398	15,474	20,901
Validation	480	960	680	960
Test	625	1,250	987	1,250
Total	4,804	9,608	17,141	23,111

Table 3.2: Our dataset statistics for each split.

if the answer cannot be found in the chart data table. This heuristic was inspired by the fact that most answers to the generated questions were values/labels of chart elements. After applying the heuristic, we manually analyzed 1,250 QA pairs and found that 86.64% of them were complete and answerable given the chart. Moreover, for the sake of fair evaluation, we manually cleaned the test set of the machine generated dataset by removing invalid questions.

- **Data split** We randomly split both of the human-written (ChartQA-H) and machine generated (ChartQA-M) QA pairs into train, validation, and test sets as shown in Table 3.2.

Type	Statista-H	Pew	OWID	OECD	Statista-M
Bar	1,696	783	507	128	15,223
Line	401	249	279	103	1,768
Pie	387	271	0	0	150
Total	2,484	1,303	786	231	17,141

Table 3.3: Number of charts from each source. Statista-H and Statista-M refer to the datasets with human-written and machine generated questions respectively from Statista

3.1.3 Dataset Analysis

Our dataset has three commonly used chart types: bar, line, and pie charts (Table 3.3). Bar is the most common type of chart across all datasets as they are quite prevalent in real-world sources. We further categorize the bar and line charts into simple vs complex where data tables of simple charts have only two columns where complex charts involve multiple columns (*e.g.*, stacked or grouped bars and multi-line charts). Among bar charts, 79.4% were simple and 29.6% were complex. For line charts, 61.0% were simple and 39.0% were complex.

We have also analyzed the basic linguistic statistics about our benchmark (Table 3.4). Unlike previous datasets, our benchmark has more unique tokens on both types of QA pairs and on both questions and answers – 6,150 and 4,319 unique tokens in questions and answers respectively in ChartQA-H whereas 12,379 and 11,979 unique tokens in questions and answers respectively in ChartQA-M. We also observe that questions cover a variety of syntactic structure and sometimes exhibit informal languages and typos. Overall, this suggests the richness of language variations which may introduce more challenges to

Type	ChartQA-H	ChartQA-M
Avg. Character per question	60.53	67.82
Avg. Character per answer	5.31	5.0
Avg. Token per question	12.32	13.18
Avg. Token per answer	1.31	1.08
Unique tokens in questions	6,150	12,379
Unique tokens in answers	4,319	11,979
Numeric answers	6,583	19,622
Non-numeric answers	3,025	3,489

Table 3.4: ChartQA benchmark statistics.

Type	Example	%
<i>Data retrieval</i>	What’s the percentage of men who thinks Valentine’s Day is overrated?	13.0
<i>Visual</i>	What is the value of the rightmost light blue bar ?	10.7
<i>Compositional</i>	How many years does the poverty percentage rose above 11%?	43.0
<i>Both visual and compositional</i>	Between the second and the third age groups from the left , which opinion deviates the most?	33.3

Table 3.5: Distribution of questions types of among 300 randomly chosen human written questions (blue-colored tokens make visual references to the chart).

the task. Finally, the topic distribution in our data is quite diverse as it is constructed from four different sources. Politics is a common topic among all sources but particularly in the Pew dataset where nearly half of charts are about U.S. Politics & Policy (45.4 %). The most frequent topic from OECD and OWID is Society (34.0 % and 26.0 % respectively). Other common topics include economy, health, and society.

To analyze the nature of questions, we randomly selected 300 QA pairs from our benchmark and categorized them into four types (Table 3.5). We see that the vast majority of questions (76.33% in total) are either compositional or both visual and compositional, which reflects the real-world scenarios where people ask complex reasoning questions.

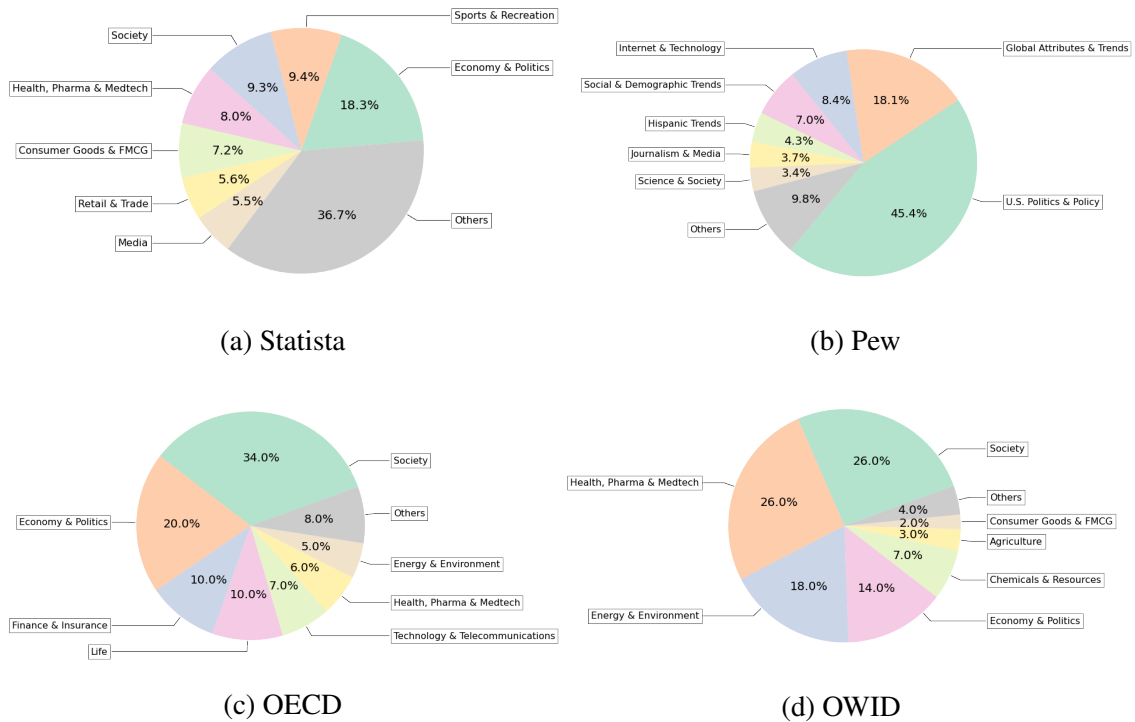


Figure 3.2: Distribution of topics in the datasets.

Type	Examples	Percentage
Color	green line, red bar	44.70%
Length	tallest bar	40.15%
Size	largest pie slice	11.36%
Position	rightmost, topmost	8.33%
Counting marks	how many green bars	3.03%
Unit of a mark	bar unit	0.76%

Table 3.6: Usage of visual references in visual questions among 300 randomly chosen questions

We also find that people make visual references to a variety of visual attributes of marks (see Table 3.6), most commonly to *color* (e.g., ‘orange line’) and *length* (e.g., ‘tallest bar’) followed by *size* (e.g., ‘largest slice’) and *position* (e.g., ‘leftmost bar’).

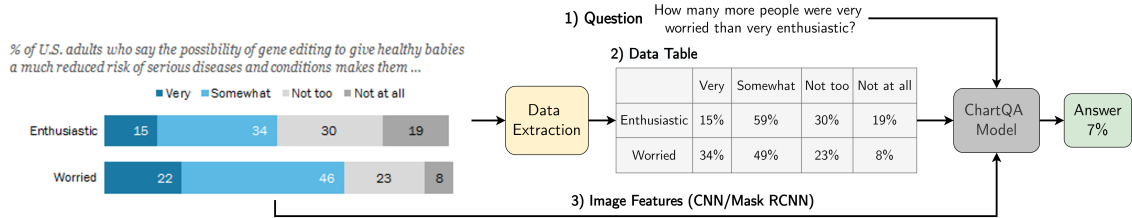


Figure 3.3: Our approach for question answering over charts. If not provided, the underlying data table is first extracted from the chart image using ChartOCR. We then pass the extracted data table in addition to the question and the image features to the ChartQA model where the ChartQA model represents one of the following: TaPas, VisionTaPas, T5, and VL-T5.

3.2 ChartQA Methodology

3.2.1 Problem Formulation

The overall process of our ChartQA system is shown in Figure 3.3. We consider two problem settings for ChartQA. The first setting assumes that the underlying data table of the chart image is available. Formally, we are given a dataset with N examples $\mathcal{D} = \{c_i, t_i, q_i, a_i\}_{i=1}^N$, where c_i represents a chart image, t_i represents the underlying data table, q_i represents a question over c_i , and a_i represents the answer to the question. The ChartQA models learn to predict the answer a_i given c_i , t_i and q_i .

The gold data tables are not generally accessible in most real-world scenarios. Thus we consider the second setup where the underlying data table t_i for chart image c_i is extracted by adapting a state-of-the-art ChartOCR [40].

3.2.2 Data Extraction

We extend ChartOCR [40] which relies on both deep-learning models and rule-based techniques to parse the chart image into the underlying data table. The chart image is parsed in three main stages. In the first stage, key-point detection networks, adapted from [33], locates the chart visual marks (*e.g.*, bars, plot area, line points). Ideally, the network locates the top-left point and bottom-right points for the rectangular objects (*e.g.*, bar, plot area). In line charts, the detection network locates the coordinates of the points connecting the line segments. In pie charts, the network locates the intersection points between the pie segments along the pie perimeter. We extend their detection networks to also locate the chart textual elements (*e.g.* *x-axis-label*, *legend-label*) as shown in Figure 3.4a and utilize the CRAFT model [3] to read their underlying texts. In the second stage, the chart scale is estimated using the *y-axis-labels* value for line and bar charts, Figure 3.4b. For pie charts, the value of each segment is estimated by calculating the angle between its borderlines. Finally, the model aggregates the extracted data values (using color and proximity heuristics) to output the final *raw data values*. We extend their approach to extract the *fully-structured* data table with the textual labels (*e.g.* column headers). As shown in Figure 3.4, we associate the estimated bars data values (*e.g.*, ‘17.13’, ‘40.14’) with their closest *x-axis-label* (‘Snapchat’). Moreover, if the chart has more than one data series (dark bars or blue bars values), each data series is matched with its *legend-label*

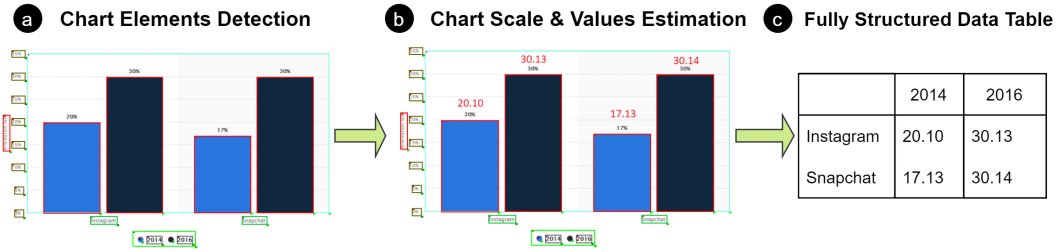


Figure 3.4: Data Extraction Process

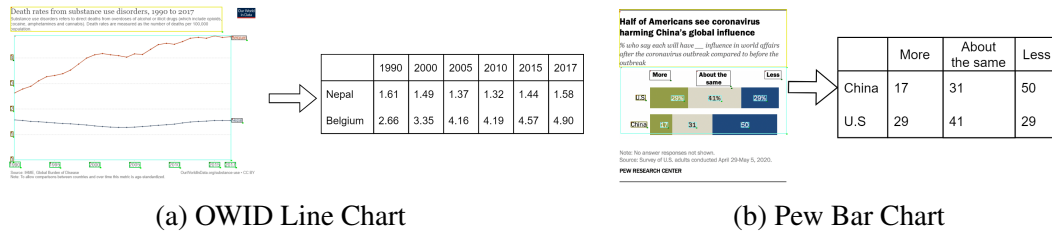


Figure 3.5: Data extraction examples from OWID and Pew.

(*e.g.*, '2016', '2014') based on the color of the *legend mark* and data-encoding marks (*e.g.*, bars). If we cannot match data values with legends by colors (*e.g.*, when all legend marks have the same color or there are no legend marks), we use other criteria that associate data-encoding marks with legend marks (*e.g.*, proximity, alignment). For example, in Figure 3.5b, 'More' is matched with '17' and '29' since they are vertically aligned. Similarly, for line charts if there is no explicit legend mark for a line series we associate the legend labels with the points of their closest lines as shown in Figure 3.5a.

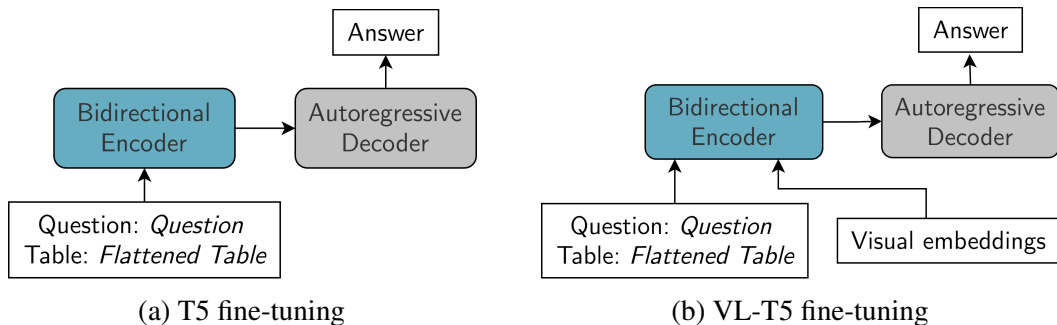


Figure 3.6: T5 and VL-T5 neural models for ChartQA. Data tables are first flattened and fed into the model along with the question (and visual features in VL-T5).

3.2.3 Models

Our approach to ChartQA builds on two of the state-of-the-art TableQA models: T5 [54, 47] and TAPAS [20]. The input to these models consists of the question q_i and the data table t_i . Different from TableQA, ChartQA often involves extracting visual information from chart images. For this, we also experiment with the visual counterparts of the TableQA models that also take the chart image features into account. While T5 has a visual variant, VL-T5 [8], TAPAS does not. In this work, we extend Tapas to consider the image features and call it VisionTAPAS.

- T5 [54] is an encoder-decoder model which unifies the NLP tasks as text-to-text generation using the same architecture and loss function. It has been pre-trained on massive amount of unlabelled data with a self-supervised denoising objective. To fine-tune T5 on our ChartQA task, we flatten the data table and feed it along with the question as: "Question: *Question tokens* Table: *Flattened table tokens*", and the model is

trained to generate the answer directly.

- VL-T5 [8] is an extension of T5 that unifies the Vision-Language (VL) tasks as text generation conditioned on multimodal inputs. The input consists of both textual tokens and visual features of the objects extracted from the image using Faster R-CNN [56]. The model is pre-trained on multiple multimodal tasks such as language modeling, visual QA, and visual grounding. We utilize VL-T5 for our ChartQA task in the following manner. For the textual input, we do the same as T5 where we flatten the data table of the chart image and concatenate it with the question text. For the visual input, we extract the visual features of different marks in the chart image (*e.g.*, bars, lines) using Mask R-CNN [18] with Resnet-101 as its backbone. We train the model to detect the following 15 objects: *'Legend'*, *'yAxisTitle'*, *'ChartTitle'*, *'xAxisTitle'*, *'LegendPreview'*, *'PlotArea'*, *'yAxisLabel'*, *'xAxisLabel'*, *'LegendLabel'*, *'PieLabel'*, *'bar'*, *'pie'*, *'line'*, *'pieSlice'*, and *'dotLine'*. For the bounding boxes annotations, we use the available bboxes. For the masks, we generate them easily using the bounding boxes for all the rectangular objects. For *'pie'* and *'pieSlice'*, we follow a similar approach to [64] where we generate the masks by projecting the radius along the pie perimeter from the starting to the ending points of each slice. We use the detectron2 library [69] and initialize the model with pre-trained weights on the COCO dataset [38]. Unlike the original VL-T5 where a fixed number of objects is provided (36), the number of elements varies from one chart to another. To account for this, we pad the extracted visual features with zeros to have a fixed length of 36.

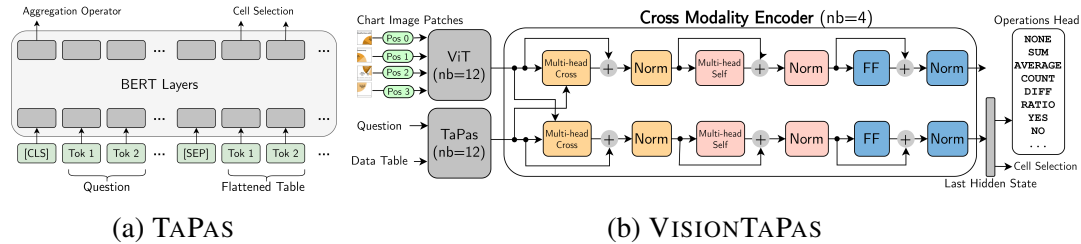


Figure 3.7: TaPas and VisionTaPas models. TaPas adds positional embeddings to the tokens to encode the tabular structure of the data table. VisionTaPas uses a cross-modality encoder to combine visual features from ViT and outputs from TaPas encoders.

- TAPAS [20] extends a BERT [12] architecture with additional positional embeddings for rows and columns to encode a table. As shown in fig:tapas, the input to the model has the following format: $[CLS] \textit{Question tokens} [SEP] \textit{Flattened table tokens}$. The tokens are encoded with the table-specific positional embeddings in addition to BERT’s segment and positional embeddings. The model has two output heads: aggregation operation head and cell selection head. The aggregation operation head predicts an operation (e.g., COUNT, SUM, AVERAGE, NONE) which is then applied to the cell values selected by the cell selection head. Depending on the operation type, the selected cells can constitute the final answer or the input used to infer the final answer.

TaPas is first pre-trained on masked language modeling objective using table-text pairs crawled from Wikipedia where table cells are randomly masked and the model is trained to predict them. It is then fine-tuned in a weakly-supervised manner (using answers as the only supervision) with end-to-end differentiable objectives.

- **VisionTaPas** is our extension of TaPas for QA over charts. It consists of three main

components: a vision transformer encoder for encoding the chart image, a TaPas encoder for encoding the question and data table and a cross-modal encoder (VisionTapas).

Vision Transformer or ViT [13] utilizes the transformer encoder architecture [67] in vision tasks. Given a 2D chart image, the image is divided into a sequence of 2D patches $\{\mathbf{p}_1, \dots, \mathbf{p}_n\}$. Each patch is then flattened and linearly projected into a d -dimensional embedding vector. To incorporate the positional information of the patches, 1D learnable positional embeddings are added to the image features. An L -layer ViT encoder produces a sequence of embeddings $\mathbf{H} = \{\mathbf{h}_{\text{cls}}^L, \mathbf{h}_1^L, \dots, \mathbf{h}_n^L\}$ representing the special [CLS] token and the image patches. We initialize the ViT module with the pre-trained weights from [13].

The **TaPas** encoder is utilized in the same manner as described above to encode the tokens in the question and the data table. For an input token sequence $\{w_{\text{cls}}, w_1, \dots, w_m\}$, an L -layer TaPas generates the corresponding encodings $\mathbf{Z} = \{z_{\text{cls}}^L, z_1^L, \dots, z_m^L\}$. This module is initialized with the TaPas weights [20] pre-trained on the WikiTQ dataset [50].

The **Cross-modality Encoder** takes the output of ViT and TaPas encoders (\mathbf{H} and \mathbf{Z}) and compute multimodal encodings. It has four blocks, each containing a visual branch and a textual-tabular branch. The input first passes through the multi-headed cross attention layers in parallel, where in the visual branch the query vectors are the visual features, and the key and context vectors are the textual-tabular features and vice versa in the textual-tabular branch. The cross-attended features are then passed through a

self-attention layer followed by a fully connected layer. Similar to the transformer model, each layer applies layer normalization [2] and is wrapped with a residual connection. Finally, we append the aggregation operation and the cell selection heads of TaPas to the final layer at the textual-tabular branch.

Extension to Other Operations Many questions in our ChartQA dataset require performing a subtraction or ratio operation, which the original TaPas model does not support. We thus extend the operation head to add those two operations (Figure 3.7b). However, instead of training them in a weakly-supervised manner based on the final answer (as done in TaPas), we find it more effective when provided with more direct but potentially noisy supervision on the cells to consider. We rely on some heuristics to generate such supervision in our training data. For example, given a question “What’s the difference between A and B?”, an answer 5, and data values “3, 6, 8”, we look for two values between which the difference is 5 (i.e. 8 and 3). While this may yield noisy supervision, similar approaches have been successfully exploited to inject reasoning capability in neural models [16, 60]; on a random sample of 100 such questions, a manual checking shows 24% noise with our heuristics. To handle the fixed vocabulary answers (*e.g.*, ‘Yes’, ‘No’), we further extend the operation head to include those classes.

3.3 ChartQA Evaluation

In this section, we describe the datasets we use in our experiments to evaluate our models as well as the baselines we compare our results against. Then, we describe the evaluation metrics we use to evaluate both our data extraction model and the different Chart Question Answering models. Finally, we analyze our results to identify the key strengths of our approach and the directions for future improvements.

3.3.1 Datasets

We evaluate our models on three datasets from previous work namely, FigureQA [28], PlotQA [44] and DVQA [26], as well as our newly created ChartQA dataset.⁴

- **FigureQA** [28] is a synthetic ChartQA dataset consisting of 180K synthetically plotted charts and 2.3M QA pairs automatically generated from 15 templates (yes/no questions). Their data values were randomly generated and they utilized the X11 color set for their charts textual labels (*e.g.*, 'legend-label', 'x-axis-label').

- **DVQA** [26] is a synthetic ChartQA dataset consisting of 300K synthetically plotted charts and 3.4M template-based QA pairs. The dataset covers different question types: structural, data retrieval, and reasoning questions. However, it also has a limited number of answers, 1576, and one chart type: bar charts. The data values were randomly generated and the chart labels were selected from the most frequent 1000 vocabulary in the Brown

⁴Two other datasets (LeafQA, LeafQA++) are not publicly available

Corpus.

- **PlotQA** [44] is the only synthetic ChartQA synthetic dataset that supports open-vocabulary questions which require aggregating over the underlying data values of the chart image. It consists of 28M QA pairs which are synthetically generated using 74 templates, and 224K synthetically plotted chart images.
- **ChartQA** is our newly created dataset covering 9,608 human-written questions focusing on logical and visual reasoning questions, and 23,111 questions automatically-generated from human-written chart summaries using the T5 model [54]. Moreover, our chart images were crawled from four different online sources (Statista, OWID, OECD, and Pew) to ensure diversity in topics and visual styles.

3.3.2 Baselines

We compare our benchmarking models with the two following baselines⁵:

- **PREFIL** [27] is a classification approach that fuses the question and image low-level and high-level features in parallel. The features are then aggregated and projected into a final classification layer.
- **PLOTQA*** is our reimplementation of PlotQA [44]. It parses the chart image to extract the underlying data table and then employs a TableQA model from [50]. However, since their data extraction approach is specific to their synthetic dataset that does not generalize

⁵Two other baselines (STL-CQA, LEAF-NET) are not publicly available

well to real-world charts, we use data tables extracted according to our method (Section 3.2.1) to evaluate their approach.

3.3.3 Evaluation Metrics

3.3.3.1 Chart Question Answering

Following [44], we use a relaxed accuracy measure for the numeric answers to allow a minor inaccuracy that may result from the automatic data extraction process. We consider an answer to be correct if it is within 5% of the gold answer. For non-numeric answers, we still need an exact match to consider an answer to be correct.

3.3.3.2 Chart Data Extraction

Our evaluation metric is adapted from ChartOCR [40]. The distance between any two data values is estimated as follows:

$$D(gt, pr) = \min(1, \left| \frac{gt - pr}{gt} \right|)$$

where gt is the ground truth value and pr is the predicted value. For each chart, the cost matrix C , where $C_{n,m} = D(gt_n, pr_m)$ is computed and the total minimum cost is

calculated by solving the following linear sum assignment problem

$$Cost = \sum_{i=1}^K \sum_{j=1}^K C_{i,j} X_{i,j}$$

Where $K = \max(N, M)$ and X is a binary assignment matrix. The final overall score is then estimated as follows:

$$Overall\ Score = \frac{1}{L} \sum_{i=1}^L 1 - \frac{cost}{K_i}$$

where L is the total number of charts.

3.3.4 Results

3.3.4.1 Data Extraction

Our evaluation results are shown in Table 3.7. We have noticed that the accuracy is specifically lower on line and dot line charts in FigureQA and PlotQA. In DVQA, the extracted tables from logarithmic-scale charts were quite noisy since ChartOCR does not support them. Moreover, PlotQA has many charts with very large values (usually written in E notation). Hence, errors in such figures have higher impact on the overall accuracy. Overall, the accuracy on PlotQA and ChartQA are generally lower since they have more complex charts (PlotQA has numerous charts with very large values (*e.g.*, $1e^6$))

and ChartQA has real-world challenging charts). A major limitation of evaluation metrics for the chart data extraction is that they do not take the extracted textual tokens into consideration (which are much more noisy in real-world figures). Hence, better metrics are still needed in the future.

Dataset	Accuracy
FigureQA	95.05%
DVQA	89.98%
PlotQA	80.88%
ChartQA	83.85%

Table 3.7: Accuracies of our data extraction algorithm on the test sets of DVQA, PlotQA, and ChartQA. Since the gold data table is not available in FigureQA, we report the results on the Validation2 set.

3.3.4.2 ChartQA

Previous Datasets When the gold data table is provided, VisionTaPas and VL-T5 achieve near perfect results, however, the performance slightly decreases when it is not provided (tab:evaluation-table). Still, VisionTaPas and VL-T5 achieve state-of-the-art results on DVQA (fully-automated setup) and PlotQA V1 datasets, respectively. For example, VisionTaPas achieves 94.54% accuracy in the DVQA test set (14.5% margin over PReFIL). Moreover, our approach proved to be more robust to OCR noise. Unlike PReFIL whose performance significantly dropped by 16.49% when using OCR outputs instead of ORACLE, VisionTaPas only witnessed a marginal decrease in performance (0.92%). Similarly, in the PlotQA dataset, both models have outperformed the PlotQA

Models	FigureQA			DVQA (ORACLE / OCR)			PlotQA		ChartQA	
	Val1	Val2	Test1	Test2	Test-Familiar	Test-Novel	Test V1	Test V2	Val	Test
Gold Data Table Provided										
TaPas	98.10%	98.09%	-	-	53.40%	53.40%	21.56%	19.55%	49.16%	51.80%
VisionTaPas	97.59%	97.96%	-	-	99.36%	99.37%	80.18%	58.29%	59.32%	61.84%
T5	95.75%	95.75%	-	-	94.33%	81.42%	93.24%	85.99%	59.11%	59.80%
VL-T5	96.45%	96.43%	-	-	98.90%	80.18%	96.38%	84.70%	58.80%	59.12%
Gold Data Table Not Provided										
TaPas	90.32%	90.43%	89.52%	89.57%	50.28% / 48.82%	50.24% / 48.68%	15.09%	12.90%	39.68%	41.28%
VisionTaPas	91.46%	91.45%	90.68%	90.64%	95.38% / 94.43%	95.46% / 94.54%	65.30%	42.50%	42.60%	45.52%
T5	87.97%	87.83%	87.56%	87.57%	90.20% / 89.01%	77.97% / 76.89%	72.62%	56.22%	40.15%	41.04%
VL-T5	88.60%	88.49%	88.20%	88.18%	94.80% / 93.75%	77.04% / 76.14%	75.90%	56.02%	38.43%	41.56%
PReFIL	94.84%	93.26%	94.88%	93.16%	96.37% / 80.88%	96.53% / 80.04%	-	-	4.53%	4.8%
PlotQA*	-	-	-	-	----- / 57.99%	----- / 59.54%	53.96% ²	22.52%	36.15%	38.00%
STL-CQA	-	-	-	-	97.35% / -----	97.51% / -----	-	-	-	-

Table 3.8: Evaluation results for different models. For DVQA, we have reported the results with and without using Oracle for OCR. We do not evaluate on FigureQA test sets with the gold data table setup since they do not have ground data tables.

model by wide margins. Another observation is that the improvement of VL-T5 over T5 is limited only to the PlotQA V1 dataset likely due to the lack of visual reasoning questions. In fact, the performance of both models is quite similar on PlotQA V2 test set where the majority of the questions are not visual. Finally, while the TaPas model achieves the best results on FigureQA (Gold Table setup), it does not perform very well on DVQA and PlotQA. This is likely because most questions in FigureQA are answerable from the data table alone. In PlotQA, however, questions are not always answerable from the data table alone and may involve the difference and ratio operations which are not supported by TaPas. This highlights the importance of the extensions we have made in the VisionTaPas model.

ChartQA Dataset We observe that VisionTaPas achieves state-of-the-art performance on both problem scenarios. PReFIL performs poorly (4.8%) as it is a classification model

²The result was reported by Levy et al. [34].

which does not work well for the open-vocabulary questions in our dataset. We also notice VL-T5 does not necessarily improve over T5, likely because many visual questions in our new dataset involve multiple references to chart elements and VL-T5 cannot effectively capture such references. Overall, the accuracies of different models are generally lower in our dataset compared to previous datasets, suggesting the challenges introduced with the human-written visual and logical reasoning questions. Finally, the performance of our models decreases when the gold data table was not given. This highlights the increasing challenge of automatic data extraction from real-world charts with diversity in styles.

We also evaluate the *transferability* of the models and the datasets, where we first pretrain the two top performing models (VisionTaPas and VL-T5) on the PlotQA dataset and then fine-tune them on ChartQA. From tab:pretrained, we notice that the accuracy increased from 41.56% to 51.84% for VL-T5 while the improvement for VisionTaPas was marginal (1.56%). One possible explanation is that VisionTaPas does not support nested arithmetic operations which are prevalent in ChartQA, so pretraining does not have a substantial effect. In contrast, we observe that the performance gain for VL-T5 were mainly for the compositional questions that do not require nested operations. Overall, this suggests that large datasets like PlotQA can be useful for pretraining the model even if the questions are generated from a small number of templates.

We also performed another experiment in which we train the VL-T5 and VisionTaPas on the PlotQA dataset and evaluate directly on the ChartQA dataset without any fine-

tuning. As shown in Table 3.9, the performance of the models decreased by wide margins when they are trained on the PlotQA dataset instead of the target dataset (e.g., 45.52% to 31.96% for VisionTaPas). This supports our hypothesis that our newly created dataset, ChartQA, introduces more challenging visual and compositional questions and more lexical variations which the previous datasets lack.

3.3.5 Ablation Studies

In the context of deep learning, ablation studies usually refer to removing certain components of the neural network in order to assess the importance of these components in the model [45]. In this section, we follow the TaPas [20] model’s ablation studies where the authors removed the aggregation operators and evaluated the model without them. In particular, we remove the supervision for ‘difference’ and ‘ratio’ operations from the VisionTaPas model to show the importance of the extensions we made in VisionTaPas on the ChartQA Task. The overall accuracy dropped by 1.80% and the accuracy on ChartQA-H (which have many such questions) dropped by 4.76% which suggests the usefulness of these operations (Table 3.9).

We further analyze the performance by chart types (Table 3.10) and question types (Table 3.11). VisionTapas and VL-T5 perform better on bar charts while the performance decreases for other charts mainly due to higher data extraction errors, especially for pie charts which are less common in our dataset. To analyze question types, we randomly

Model	ChartQA-H	ChartQA-M	Overall
TaPas	28.72%	53.84%	41.28%
VisionTaPas	29.60%	61.44%	45.52%
VisionTaPas [†]	24.84%	61.60%	43.72%
T5	25.12%	56.96%	41.04%
VL-T5	26.24%	56.88%	41.56%
VisionTaPas [*]	25.12%	38.80%	31.96%
VL-T5 [*]	22.08%	19.84%	20.96%
VisionTaPas Pretrained	32.56%	61.60%	47.08%
VL-T5 Pretrained	40.08%	63.60%	51.84%

Table 3.9: Accuracy of the different models on our benchmark. VisionTaPas[†] does not support difference and ratio operations. VisionTaPas^{*} and VL-T5^{*} are trained on PlotQA and evaluated directly on ChartQA.

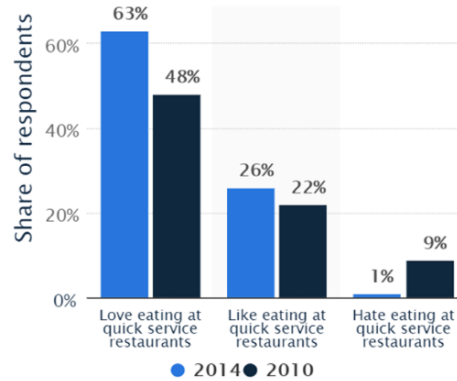
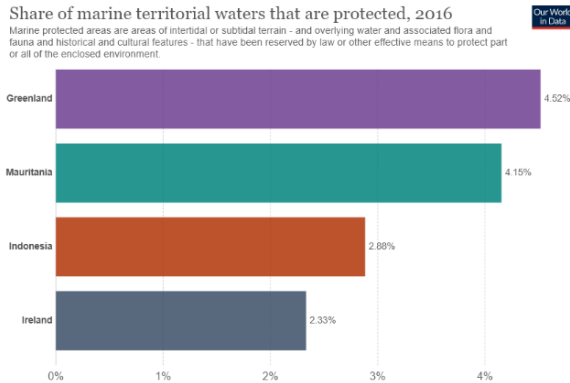
sampled 200 human-written questions and classified them into four main categories. As expected, the performance is much higher on the data retrieval questions that do not require mathematical reasoning while the performance is lower for visual questions which refers to chart elements.

Model	Bar	Line	Pie	Overall
VisionTaPas	49.80%	38.20%	24.41%	45.52%
VL-T5	45.82%	35.40%	25.00%	41.56%

Table 3.10: Results for VisionTaPas and VL-T5 on the ChartQA test set by chart type.

Model	Data Retrieval	Visual Compositional	Compositional	Visual	Overall
VisionTaPas	60.00%	29.78%	34.88%	16.21%	34.00%
VL-T5	50.00%	19.14%	24.41%	21.62%	26.50%

Table 3.11: Accuracies of VisionTaPas and VL-T5 on the ChartQA-H test set by question type on 200 random samples.



Q1: what is the difference between the sum shares of indonesia and ireland and share of mauritania?

A: 1.06 **Output:** 0.547

Q2: what is the least difference between light blue bar and dark blue bar?

A: 4 **Output:** 39.32

Figure 3.8: Example of errors from VisionTaPas

3.3.6 Qualitative Analysis

We have manually analyzed model predictions to investigate the key challenges existing models face (see sample predictions in Figure 3.9).

Logical Inference with Nested Operations While VisionTaPas and VL-T5 handle various mathematical/logical operations, still they cannot effectively handle nested operations. For example, *Q1* in results requires the model to add two numbers and then subtract from another number, but our model only outputs the difference between two numbers. In future, we will extend the VisionTaPas model (by possibly training it in a sequential fashion [9]) to address the issue.

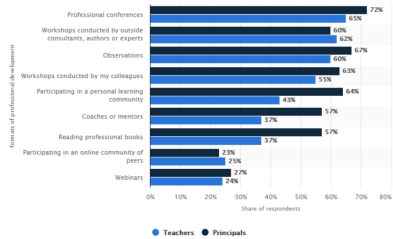
Input Representation Complex visual compositional questions may require a multi-stage

reasoning process (*e.g.*, $Q2$ in results). Currently, our models take the data table and the visual features of the chart separately and then combine them. Such representation does not fully capture the chart structure. In future, we will develop better representations including semantic graph representations [66] that can exploit the relations among the question, chart objects, and data values.

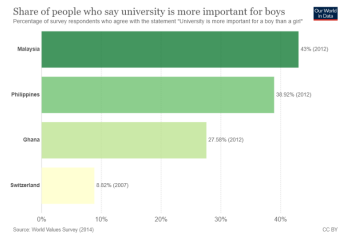
Computer Vision Challenges Table 3.8 indicates that performance of our models decrease when the gold table is not given, suggesting the need for more accurate data extraction. Current approaches for automatic data extraction are modular and combine deep learning and rule-based methods which are error-prone. An end-to-end deep learning approach could help improve the performance and generalize well to different chart styles.

3.4 Discussion

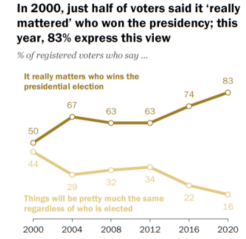
In this chapter, we presented our chart question answering benchmark. First, we discussed our newly created ChartQA dataset by showing the dataset preparation, collection, and annotation processes. We also analyzed the key statistics of our dataset by showing the diversity of our real-world charts' visual styles and the human-authored questions' types which distinguish our dataset from the previous ones and make it more challenging. Second, we presented our ChartQA approach and a set of transformer-based models. Our proposed model VisionTaPas achieves the state-of-the-art performance on the previous datasets as well as on our newly created benchmark, ChartQA. Finally, we analyzed the



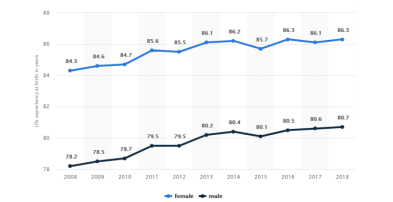
Q1: What is the most effective teaching format for Principals?
A: Professional conferences
Output: Professional conferences



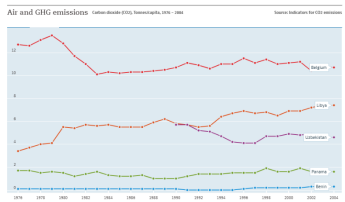
Q2: What is the average share of people in Philippines and Ghana who think University is more important for boys?
A: 33.25
Output: 33.27



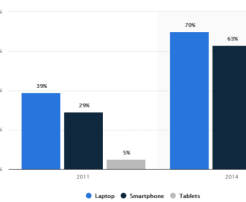
Q3: What's the peak value of dark brown graph?
A: 83
Output: 83



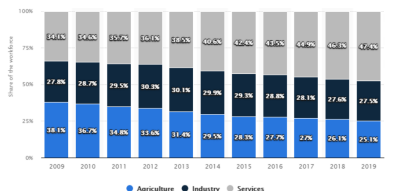
Q4: What is the difference between the highest and lowest life expectancy at birth from 2008 to 2018 for female?
A: 2
Output: 2.01



Q5: Which country recorded the highest Air and GHG emissions over the years?
A: Belgium
Output: Belgium



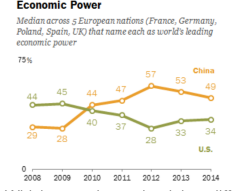
Q6: Which digital device has most explosive increase in ownership from 2011 to 2014?
A: Tablets
Output: Laptop



Q7: In which economic sector the workforce distribution was steadily increasing over the years?
A: Services
Output: Industry



Q8: Which year saw the sharpest drop in Estimated revenue?
A: 2009
Output: 2011



Q9: Which year shows the tiniest difference in values between China and US being seen as leading economic power across all the years?
A: 2010
Output: 2012

Figure 3.9: Sample outputs of our model VisionTaPas on our new ChartQA test set. Answers in green are correct and answers in red are incorrect.

results of our models by highlighting their limitations and showing the possible areas for improvements.

4 Transfer Learning for Chart Question Answering

4.1 Motivation

As we previously discussed, our benchmark, ChartQA, has complex questions that require the model to reason over the visual properties of the chart and apply mathematical/logical operations on the underlying data values. Our previous models (*e.g.*, T5 and VL-T5) do not have an effective understanding of the structural and visual properties of the chart since they were pretrained on text data or vision language tasks that mainly use natural images which are different from charts. Our previous models also lack mathematical/logical reasoning abilities which are critical in the ChartQA task. In the previous chapter, we have showed that pretraining our models on large datasets (*e.g.*, PlotQA) can be useful to inject several skills into our models (*e.g.*, numerical reasoning skills). Finetuning these pretrained models on the ChartQA dataset lead to significant improvement in the performance. In this chapter, we extend this transfer learning approach by defining more complex pretraining tasks (*e.g.*, visual reasoning) that can address the challenges in our dataset. These pretraining tasks have three main goals: (*i*) teach the model to retrieve the

visual properties of the chart elements such as colors (e.g., blue bar) and positions (e.g., leftmost bar), (ii) perform mathematical/logical operations (e.g., sum, difference) which can be combined in compositional questions, (iii) reason over the visual properties of the chart elements.

Another major limitation in our previous approach is that the input representation to our models does not relate between the visual features and the data table values. For example, in the VL-T5 and VisionTaPas models, we feed the data table and the visual features separately without relating them. In this chapter, we define a new input representation by fusing both the visual features and data table values in a unified format. Our input representation associates each data value in the underlying data table cells with the visual features of its corresponding element (e.g., bar, line point, x-axis-label) in the chart image. Consequently, our transfer learning approach outperforms our previously described approach (Chapter 3) by wide margins and achieves the state-of-the-art results on our newly created dataset, ChartQA.

4.2 Problem Definition & Input Representation

In our task setup, we are given a dataset with N examples $\mathcal{D} = \{c_i, t_i, q_i, a_i\}_{i=1}^N$, where c_i represents a chart image, t_i represents the underlying data table, q_i represents a question over c_i , and a_i represents the answer to the question. Our approach first fuses c_i and t_i to construct our unified input representation vt_i which represents the visual data table. Each

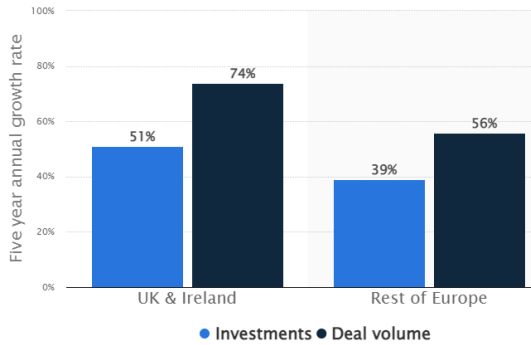
cell in the visual data table, vt_i , contains the data value of a chart element (e.g., bar) in addition to its visual features. Our ChartQA model then learns to predict the answer a_i given vt_i and q_i .

As shown in Figure 4.1, the data table, t_i , consists of three main cell types: data values, column headers, and row names. While the data values usually corresponds to the bars, line points, or pie segments values in the chart image, c_i , the column headers and row names correspond to the x-axis-labels and legend labels respectively. Consequently, in order to construct the visual data table, vt_i , we crop the chart element (e.g., bar) that corresponds to each data table cell from the chart image and combine it with the cell value. Since the legend marks colors are essential in connecting the legend labels to their corresponding data elements (e.g., bars), we also fuse them in the visual data table, vt_i , as shown in Figure 4.1 by connecting them to their corresponding legend labels. To sum up, our new unified format associate the visual features with their corresponding cells in the underlying data table which can facilitate the visual reasoning process for our model.

4.3 Chart-T5 Model

Our new model, Chart-T5, also builds on one of the state-of-the-art Vision-Language (VL) models, VL-T5 [8]. Unlike VL-T5 [8] where the chart representation in the input consists of two separate components: the visual features and the data table, Chart-T5 utilizes our newly constructed representation, the visual data table (vt_i). The input to

Five-year compound annual growth rate on investments and deals in the fintech industry in Europe as of 2014



(a) Bar Chart from the Statista dataset

	UK & Ireland	Rest of Europe
Investments	51%	39%
Deal Volume	74%	56%

(b) Data table of the chart

		X-axis-label UK & Ireland UK & Ireland	X-axis-label Rest of Europe Rest Of Europe
Legend Label Investments Investments	Legend Mark ● <legend-mark>	Bar ■ 51%	Bar ■ 39%
Legend Label Deal volume Deal volume	Legend Mark ● <legend-mark>	Bar ■ 74%	Bar ■ 56%

(c) Visual data table of the chart

Figure 4.1: A bar chart image with its corresponding data table and visual data table. The data table consists of three main cell types: data values (green), column headers (yellow), and row names (blue). Each one of these cells is then extended to include the visual features so that we can construct the visual data table. We obtain the visual features of each cell by cropping its relevant element (e.g., bar) from the chart image. Each cell in Figure 4.1c contains the cropped image and the underlying data value of each chart element (separated by a vertical black line). They also contain the cell type (e.g., x-axis label, bar, legend mark, .etc)

our model consists of a one-dimensional sequence of tokens. Therefore, we flatten our two-dimensional visual data table into a series of sequences: the table header sequence that starts with the [SEQX] token and table row sequences that start with the [SEQL]

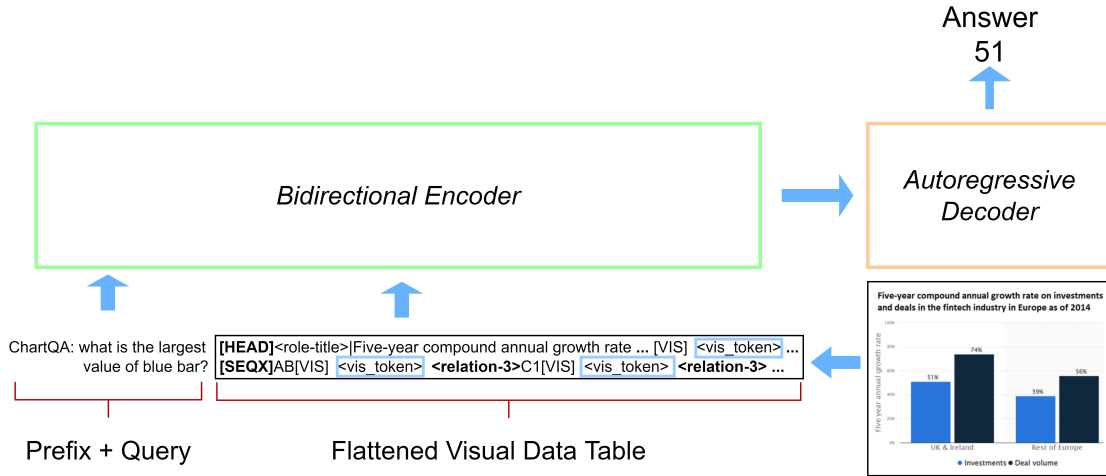


Figure 4.2: Chart-T5 model is a seq2seq model consisting of bidirectional encoder and autoregressive decoder. The input to the encoder consists of the flattened visual data table in addition to the prefix and question tokens, and the model is finetuned to generate the answer.

tokens as shown in Table 4.1. Moreover, we add another head sequence (starts with [HEAD]) that contains the chart main elements (plot area, title, legend area, x-axis title, y-axis title). We separate between the tokens of each cell using special types of relations (e.g., relation-1, relation-2) depending on the type of the connected cells. Each cell consists of two types of tokens: text tokens and a visual token which are concatenated as follows: "Text Tokens [VIS] Visual Token" where [VIS] is a special token that separates between the text and visual tokens. Similar to T5 [54] and VL-T5 [8], we obtain the text tokens by tokenizing the cell text (e.g., data value) into a series of tokens (subwords). To obtain the visual token of the cell, we crop its corresponding element from the chart image (e.g., bar) and pass through a CNN network to produce the visual features (see Figure 4.3). Then, the visual features vector is added to the bounding box vector to produce the visual

```

[HEAD] <type-title>|Title[VIS]<vis.token><relation-1><type-plot-area>[VIS]<vis.token>...
[SEQX] UK & Ireland[VIS]<vis.token><relation-2>Rest of Europe[VIS]<vis.token> ...
[SEQL] Investments[VIS]<vis.token><relation-3>51%[VIS]<vis.token><relation-3>...
[SEQL] Deal Volume[VIS]<vis.token><relation-3>74%[VIS]<vis.token><relation-3>...

```

Table 4.1: Flattened visual data table.

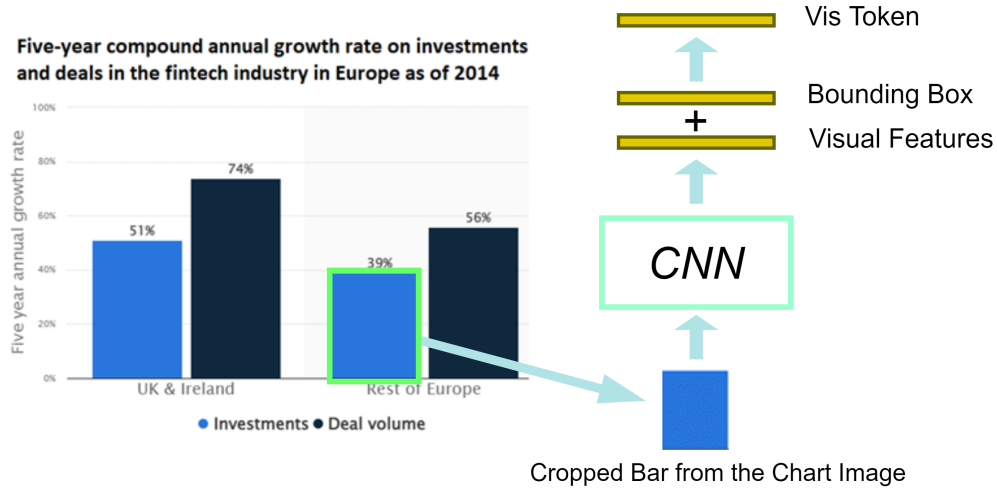


Figure 4.3: Our approach to obtain the Visual Token of the chart elements. The chart element (e.g., bar) is first cropped from the chart image. Then, we pass it through a CNN network to extract the visual features from it. After that, we add the visual features and bounding box vector to obtain the final Visual Token that encodes both the visual features and positional information (bounding box)

token. Finally, as shown in Figure 4.2, we concatenate all of the flattened visual data table sequences and feed it along with the task prefix and the question as: “ChartQA: *Question tokens - Flattened Visual table tokens*”, to our model and train it to generate the answer directly.

Pretraining Tasks	Chart Image	Task Prefix + Query	Output
Visual Attribute Retrieval		color: 51% bar position: rightmost bar position: second bar from the left	Blue 56% 74%
Numerical Reasoning		reasoning: what is the difference between the investments in Europe and the UK?	12%
Visual Reasoning		visual reasoning: what is the sum of the dark blue bars?	130

Table 4.2: Input and Output formats for our pretraining tasks. The input consists of the task prefix and query which we later concatenate with the chart visual data table.

4.4 Pretraining

4.4.1 Tasks

In order to teach our model the structural and visual properties of the chart as well as the numerical reasoning skills, we pretrain our model under a multi-task setup using three pretraining tasks: (i) Visual Attribute Retrieval (ii) Numerical Reasoning (iii) Visual Reasoning. We explain these tasks in detail below.

• Visual Attribute Retrieval

Visual Attribute Retrieval is a pretraining task that teaches the model to extract and utilize the visual attributes (e.g., color and position) of the chart different elements. Extracting and reasoning over the colors and positions of the chart elements can be of great importance for our model when answering the visual questions that refers to the visual attributes of the chart elements (e.g., what is the value of the red bar?, what is the value of the leftmost bar?). To construct the queries for this pretraining task, we manually analyzed the human-authored questions in our dataset to identify the most common visual

attributes references. We have noticed that people mostly refer to the chart elements using their colors or positions (e.g. blue bar, leftmost bar). Consequently, we designed **18** templates (see Table 4.2 and the Appendix A.1 for all the templates) which we utilize to generate the queries for the chart images.

- **Numerical Reasoning**

The ChartQA task involves questions with various mathematical/logical operations that need to be applied on the chart data values to obtain the answer. Moreover, many questions in our dataset are compositional and may require applying more than one operation which our previous models, VL-T5 and VisionTaPas, fail to address (e.g., what is the difference between the total revenue in 2017 and 2018 of all countries?). Hence, it is essential to infuse the numerical reasoning skills into our model. To this extent, we design a template-based numerical reasoning task where the model is trained to execute/perform the most common mathematical operations over the charts data values. We manually analyzed our dataset to find the most common operations (e.g., sum, average, difference, etc.) and constructed **27** templates that we utilize to generate the task questions (see Table 4.2 and the Appendix A.1 for all the templates).

- **Visual Reasoning**

The ChartQA task also requires the models to apply mathematical operations on the data values of the chart elements which are referred to by their visual attributes in the questions (e.g., what is the difference between the leftmost and rightmost bars?).

Consequently, we introduce the visual reasoning pretraining task which can be viewed as the combination of the two above-mentioned tasks. After analyzing the human-authored questions, we created **44** visual reasoning templates (see Table 4.2 for examples and the Appendix A.1 for all the templates).

4.4.2 Pretraining Dataset

We aggregate our pretraining data from five different chart datasets: OWID, OECD, Statista, PlotQA, and Pew to ensure variety in the different visual styles of the chart images. We first carefully filter these datasets to remove any duplicate or badly designed charts for which data extraction is not possible due to poor resolution or unconventional styles. Then, we utilize the available underlying data tables and annotations of the chart images to generate the pretraining tasks queries from our designed templates. Moreover, we have also added from **1,392,471** QA pairs from the PlotQA V1 dataset questions to our pretraining corpus since they contain a lot of mathematical/logical questions. Overall, our pretraining dataset contains **331,696** chart images and **3,487,893** queries. We carefully split the pretraining data into training and validation sets such that the test images in our downstream task, ChartQA, are not seen by the model during the pretraining process.

4.5 Evaluation

4.5.1 Datasets, Baselines & Metrics

Looking at the results of our previous models (VL-T5 and VisionTaPas) on the FigureQA [28], DVQA [26], and PlotQA [44] datasets (see Table 3.8), we notice that our models give near perfect results given the underlying gold data table. However, the results on our new dataset, ChartQA, is still limited due to several open challenges such as the visual reasoning and compositional questions. Consequently, we mainly focus our evaluation of the Chart-T5 model on the ChartQA dataset. We first pretrain our Chart-T5 model on our pretraining dataset for 10 epochs and then we finetune it on the ChartQA dataset for 30 epochs. In both experiments, we use a learning rate of 0.0001 and a batch size of 96. Moreover, during pretraining we randomly shuffle the ordering of the [SEQL] sequences in our visual data table (see Table 4.1) to make our model insensitive to their ordering. Since we are using the ground truth annotations to construct the visual data tables of the chart images, we compare the Chart-T5 model results with the gold data table setup to ensure a fair comparison. Specifically, we compare our Chart-T5 model with our previous best performing model, VL-T5 (pretrained on PlotQA). Following our previous approach and PlotQA [44], we use the relaxed accuracy metric to evaluate our models. For the numeric answers, we consider the answer to be correct if it is within 5% of the gold answer. For the non-numeric answers, an exact match is required to consider

the answer correct.

Model	ChartQA-H	ChartQA-M	Overall
VL-T5 Pretrained	49.92%	92.40%	71.16%
Chart-T5	60.88%	86.96%	73.92%

Table 4.3: A comparison between our best performing model from Chapter 3, VL-T5 **Pretrained** and the Chart-T5 model. **Pretrained** refers to the models that are first pretrained on PlotQA and then finetuned on our benchmark, ChartQA.

4.5.2 Results

Looking at Table 4.3, we can observe that Chart-T5 outperforms our previous best performing model, VL-T5 (pretrained on PlotQA), and sets the new state-of-the-art results (73.92%). We can notice a significant improvement in performance on the human-authored questions, ChartQA-H, (49.92% to 60.88%) mainly due to the numerical and visual reasoning skills injected into our model, Chart-T5, through our pretraining tasks and our new input representation, the visual data table. Moreover, we manually analyzed both our models, Chart-T5 and VL-T5, predictions on the human-authored test set, and noticed that the performance gain for Chart-T5 was mainly for the visual reasoning questions. As shown in Figure 4.4, Chart-T5 is superior on the visual reasoning questions. Finally, Chart-T5 performance slightly decreases on the machine-generated questions, ChartQA-M. Some of the chart images in the ChartQA-M split are very large; they present so many data points (see Figure 4.5 for examples). Consequently, the input representation, visual data table, becomes very long and exceeds the input limit to our Chart-T5 model,

which is 512 tokens. On the other hand, the regular data table used by VL-T5 are often shorter and have less tokens since they only encode a fixed number of visual tokens, 36, regardless of the number of chart elements. Moreover, ChartQA-M mostly consists of data retrieval questions and do not contain visual references to the chart elements. Hence, our new input representation and pretraining tasks have minimal effect on them.

4.5.3 Qualitative Analysis

We have manually analyzed the Chart-T5 model predictions on the ChartQA test set to identify the key challenges and limitations.

- **Numerical Hallucinations** Although our numerical reasoning tasks help the Chart-T5 model to understand and perform the common mathematical/logical operations, it tends to hallucinate and deviate with the answers that may contain many digits. For example, as shown in Figure 4.6a, the model’s predictions are close to the gold answers which indicates that the model was able to understand the questions, but failed to output the exact correct numbers. Unlike TaPas [20] and VisionTaPas where the model predicts the operation and then applies it on the selected cells, Chart-T5 generates the final answer as text in an autoregressive manner which may result in hallucinations.

- **Logical Questions with Nested Operations** While our pretraining tasks include many questions with nested operations, our model is still struggling with such questions in our ChartQA dataset (see Figure 4.6b for examples). Currently, we pretrain the model to

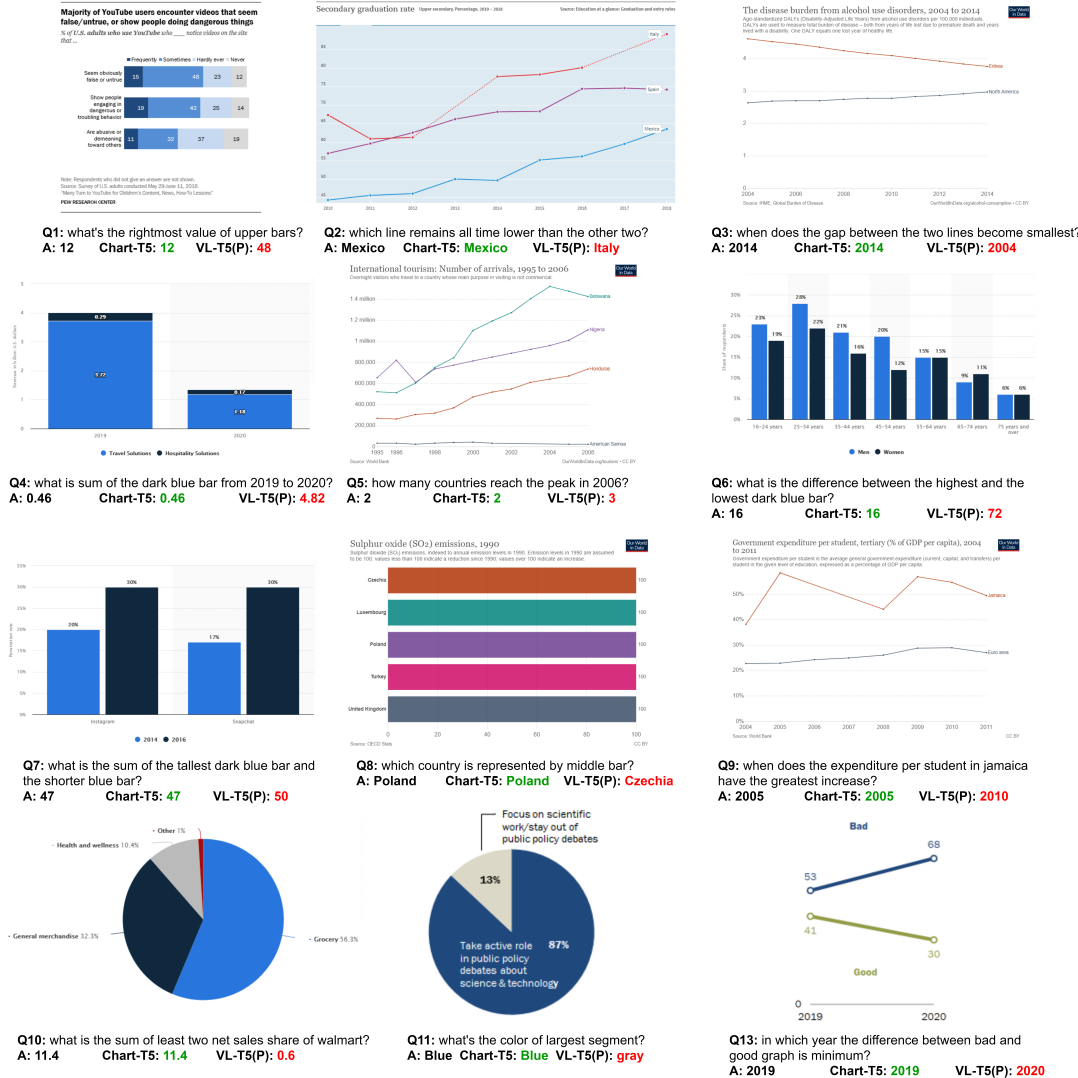


Figure 4.4: Sample outputs of the Chart-T5 and VL-T5 (Pretrained) models on our new ChartQA test set. Answers in green are correct and answers in red are incorrect.

predict the final answer directly which can be considered as a high-level task. To further analyze and eliminate this limitation, we are planning to decouple the mathematical expression prediction phase (cells selection, operations prediction, and their orders (e.g., Ratio(Sum(3, 4)), 50)) from the aggregation phase which executes the mathematical

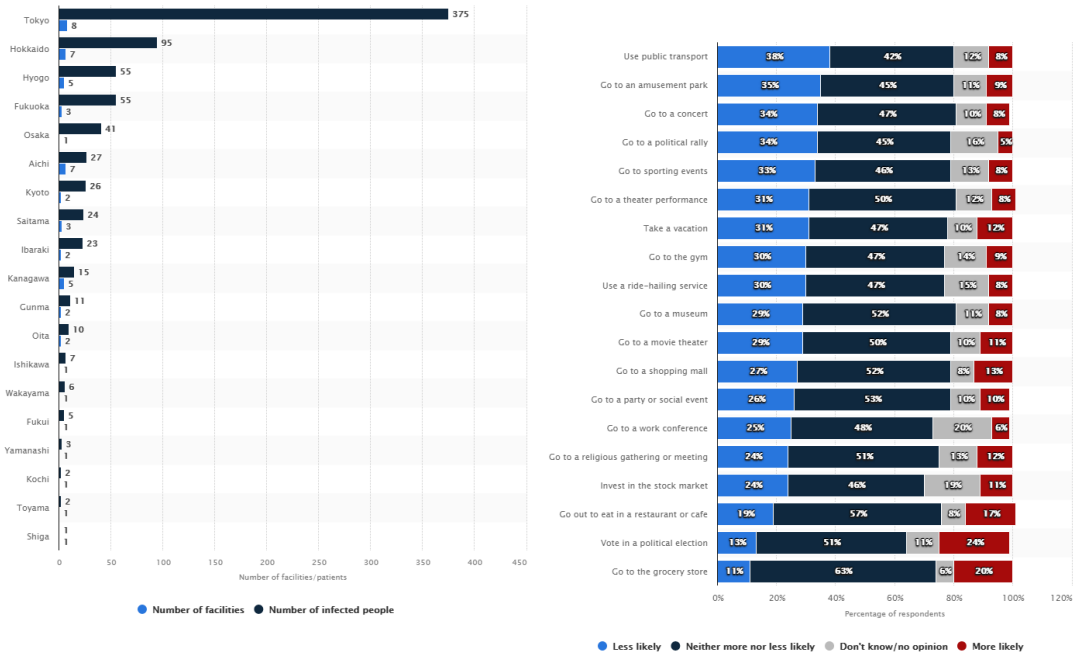


Figure 4.5: Some charts with many data points from the ChartQA-M dataset.

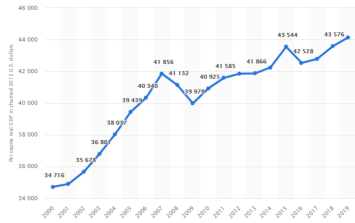
expression to produce the final answer. Splitting our high-level tasks into multiple low-level ones can help us better analyze the limitations of the model and improve its performance on the questions with nested operations.

- Chart Data Extraction** Our new task setup assumes that the Chart-T5 model have access to the underlying data and visual information (*e.g.*, bounding boxes) of the chart elements. In contrast, a lot of the real-world charts are stored in bitmap format on the Web, so our current approach does not work on them. In the future, we are planning to extend our previous chart data extraction approach to automatically extract the visual information as well as the underlying data of the chart elements which we can utilize to construct the visual data table of the chart image.

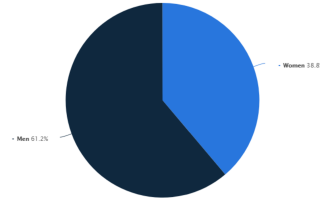
• **Input Representation** In order to feed the visual data table to the input of our Chart-T5 model, we currently flatten the visual data table into sequences and concatenate all their tokens in order. Such input linearization approach makes it more challenging for the model to understand the structure of the visual data table (*e.g.*, associating the column header with their corresponding data cells in the same column). Moreover, although the visual data table improves the input representation to our model by relating the visual features and the data values, it drops important semantic and spatial relations between the chart elements (*e.g.*, containment relation between the bars and the plot area). In the future, we are planning to develop better chart representations using semantic knowledge graphs [66] that better conveys the chart underlying structure. We will also utilize special embeddings [61] to encode the graph structure in the input to our model.

4.6 Discussion

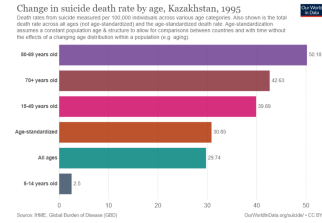
In this chapter, we presented our transfer learning approach for chart question answering. First, we discussed our new input representation, the visual data table, that relates the visual features of the chart elements with their underlying data values. Second, we also presented the Chart-T5 model and how we preprocess its input. Third, we discussed our three pretraining tasks: (i) Visual Attribute Retrieval (ii) Numerical Reasoning (iii) Visual Reasoning, and how they teach the model to reason over the structural and visual properties of the chart elements. Fourth, we presented the results of our Chart-T5 model on the



Q1: what is the per capita real gross domestic product of montana in the year 2007 (in chained 2012 us dollars)?
 A: 41856
 Output: 41867

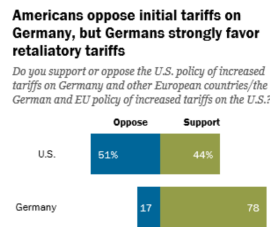


Q2: work out the ratio of the bigger segment to the smaller one?
 A: 1.577
 Output: 1.6875

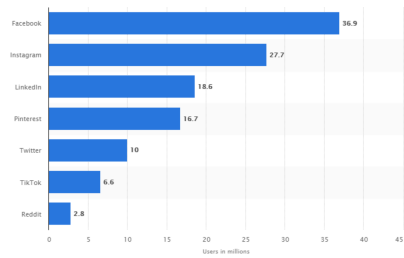


Q3: what is the biggest difference in the age between the highest suicidal age ground and the lowest one?
 A: 48.68
 Output: 47.68

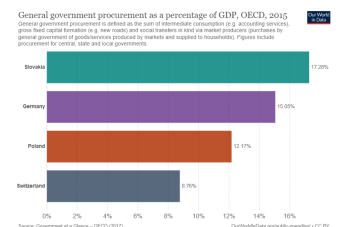
(a) Hallucinations in the model's outputs



Q1: divide average of green bars by the smallest value of the blue bar, what's the result (round to one decimal place)?
 A: 3.6
 Output: 0.063



Q2: among facebook, instagram, linkedin, pinterest and twitter users percent, what is the average minus the median?
 A: 3.38
 Output: 0.245



Q3: add the three largest government procurement in gdp and multiply it by the smallest government procurement in gdp across all the countries
 A: 389.82
 Output: 0.094

(b) Questions with Nested Operations

Figure 4.6: Example of errors from the Chart-T5 model.

ChartQA dataset and showed that it outperforms the previous best performing model, VL-T5 (pretrained on PlotQA). Finally, we analyzed the results of the Chart-T5 model by highlighting its limitations and showing the possible areas for improvements.

5 Conclusions and Future Work

In this chapter, we first present our conclusions. Then, we discuss our directions for future work.

5.1 Conclusions

In this thesis work, we have two main contributions: (i) We present the ChartQA dataset (ii) We propose new approaches with transformer-based models that achieve the state-of-the-art performance.

We introduced ChartQA, a new large-scale chart question answering benchmark with 9.8K human-written questions focusing on visual and logical reasoning on 4.6K chart images. Our chart images were crawled from four different online sources (Statista, OWID, OECD, and PEW) to maximize the variations in the visual styles. We have also augmented the dataset with 23.1K machine-generated questions from the Statista human-written chart summaries using the T5 model. Our dataset covers the most common three chart types: bar, line, and pie charts. Moreover, the topic distribution in our data is quite diverse (politics, economy, health, society, .etc). Unlike previous template-based

questions, our human-authored questions have more language variations and may exhibit the informal, intricate, and nuanced nature of language. Most of our questions are either compositional or visual and compositional which makes our dataset more challenging.

In this thesis, we also introduce new approaches that combine visual features and underlying data values from a chart to answer questions. We study two task setups: (i) Gold Data Table Provided (ii) Gold Data Table Not Provided. In the first setup, our models achieve near perfect results on the previous datasets (FigureQA, DVQA, PlotQA) and state-of-the-art results on the ChartQA dataset. In the second task setup, we automatically extract the underlying data table from the chart image to use it in our models. We extend the ChartOCR data extraction approach to automatically output the fully structured data table from a chart image. Although the performance slightly decreases, our models still achieve the state-of-the-art results on the DVQA, PlotQA, and ChartQA datasets.

Finally, we also study transfer learning techniques to address some of the major challenges in our dataset, ChartQA, and boost the performance. We introduce a new input representation, visual data table, that relates the visual features with their underlying data values from the data table. We also introduce three pretraining tasks: (i) Visual Attribute Retrieval (ii) Numerical Reasoning (iii) Visual Reasoning. Although our model, Chart-T5, improves the performance on the ChartQA dataset, especially the human-authored questions, our evaluation also reveals several unique challenges that emerge from the

visual and reasoning questions with nested mathematical/logical operations.

We hope that our benchmark and models will serve as a starting point for other researchers to address the challenges in our dataset.

5.2 Future Work

While our newly created benchmark is considered the first ChartQA benchmark with real-world charts and human-authored questions, it is relatively small compared to the previous large-scale synthetic datasets. Moreover, our real-world charts were crawled from only four online sources which may limit the visual styles of our charts and the generalizability of our models. In the future, we are planning to develop a larger dataset with more charts from additional online sources to increase the variations in the visual styles, and human-authored questions.

While our models achieve impressive results on the ChartQA dataset as well as the previous datasets, they still have some areas for improvement. First, the input representation to our models does not encode the full underlying structure of the charts. Although our new input format, visual data table, makes the first step by associating the visual features with the data values, it still drops important spatial and semantic relations between the chart elements. In addition, flattening the visual data table in the input to our models may affect and limit the model’s understanding of the input structure. In the future, we are planning to explore better chart representations including semantic graph representations

[66] that can exploit the relations among the question, chart visual features, and data values. We will also explore various ways to encode the structure of the input graphs using special types of embeddings [61]. Second, our models' performance decreases when we use the extracted data table instead of the gold one mainly due to the errors of our chart data extraction approach. Hence, we plan to develop an end-to-end deep learning chart data extraction approach that could help improve the performance and generalize well to different chart styles. Third, we will explore more pretraining tasks and techniques [71, 48] that can further improve the numerical reasoning skills of the Chart-T5 model, especially on the questions with nested mathematical/logical operations. Finally, we will generalize our pretraining tasks to include other chart types (*e.g.*, scatter plots) and focus on other downstream tasks (*e.g.*, chart summarization [29] in addition to ChartQA).

Bibliography

- [1] Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Dhruv Batra, and Devi Parikh. Vqa: Visual question answering, 2016.
- [2] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv*, abs/1607.06450, 2016.
- [3] Jeonghun Baek, Geewook Kim, Junyeop Lee, Sungrae Park, Dongyoon Han, Sangdoon Yun, Seong Joon Oh, and Hwalsuk Lee. What is wrong with scene text recognition model comparisons? dataset and model analysis. In *International Conference on Computer Vision (ICCV)*, 2019.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>.
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [6] R. Chaudhry, S. Shekhar, U. Gupta, P. Maneriker, P. Bansal, and A. Joshi. Leafqa: Locate, encode attend for figure question answering. In *2020 IEEE Winter*

- Conference on Applications of Computer Vision (WACV)*, pages 3501–3510, 2020. doi: 10.1109/WACV45572.2020.9093269.
- [7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning, 2020.
- [8] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. In *ICML*, 2021.
- [9] Minseok Cho, Reinald Kim Amplayo, Seung won Hwang, and Jonghyuck Park. Adversarial tableqa: Attention supervision for question answering on tables. *ArXiv*, abs/1810.08113, 2018.
- [10] J. Choi, Sanghun Jung, Deok Gun Park, J. Choo, and N. Elmqvist. Visualizing for the non-visual: Enabling the visually impaired to use visualization. *Computer Graphics Forum*, 38, 2019.
- [11] Pradeep Dasigi, Matt Gardner, Shikhar Murty, Luke Zettlemoyer, and Eduard Hovy. Iterative search for weakly supervised semantic parsing. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2669–2680, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1273. URL <https://aclanthology.org/N19-1273>.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xi-aohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.

- [14] Julian Eisenschlos, Syrine Krichene, and Thomas Müller. Understanding tables with intermediate pre-training. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 281–296, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.27. URL <https://aclanthology.org/2020.findings-emnlp.27>.
- [15] Siddhant Garg, Thuy Vu, and Alessandro Moschitti. TANDA: transfer and adapt pre-trained transformer models for answer sentence selection. *CoRR*, abs/1911.04118, 2019. URL <http://arxiv.org/abs/1911.04118>.
- [16] Mor Geva, Ankit Gupta, and Jonathan Berant. Injecting numerical reasoning skills into language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 946–958, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.89. URL <https://aclanthology.org/2020.acl-main.89>.
- [17] Moonsu Han, Minki Kang, Hyunwoo Jung, and Sung Ju Hwang. Episodic memory reader: Learning what to remember for question answering from streaming data. *CoRR*, abs/1903.06164, 2019. URL <http://arxiv.org/abs/1903.06164>.
- [18] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. doi: 10.1109/ICCV.2017.322.
- [19] Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend, 2015.
- [20] Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. TaPas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.398. URL <https://www.aclweb.org/anthology/2020.acl-main.398>.
- [21] Enamul Hoque, Vidya Setlur, Melanie Tory, and Isaac Dykeman. Applying pragmatics principles for interaction with visual analytics. *IEEE Transactions on Visualization and Computer Graphics*.
- [22] Mohit Iyyer, Wen-tau Yih, and Ming-Wei Chang. Search-based neural structured learning for sequential question answering. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,

- pages 1821–1831, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1167. URL <https://aclanthology.org/P17-1167>.
- [23] Daekyoung Jung, Wonjae Kim, Hyunjoo Song, Jeong-in Hwang, Bongshin Lee, Bohyoung Kim, and Jinwook Seo. *ChartSense: Interactive Data Extraction from Chart Images*, page 6706–6717. Association for Computing Machinery, New York, NY, USA, 2017. ISBN 9781450346559. URL <https://doi.org/10.1145/3025453.3025957>.
- [24] Daekyoung Jung, Wonjae Kim, Hyunjoo Song, Jeongin Hwang, B. Lee, B. H. Kim, and Jinwook Seo. Chartsense: Interactive data extraction from chart images. *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, 2017.
- [25] Kushal Kafle, Mohammed Yousefhussien, and Christopher Kanan. Data augmentation for visual question answering. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 198–202, Santiago de Compostela, Spain, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-3529. URL <https://aclanthology.org/W17-3529>.
- [26] Kushal Kafle, Scott Cohen, Brian L. Price, and Christopher Kanan. DVQA: understanding data visualizations via question answering. *CoRR*, abs/1801.08163, 2018. URL <http://arxiv.org/abs/1801.08163>.
- [27] Kushal Kafle, Robik Shrestha, Brian L. Price, Scott Cohen, and Christopher Kanan. Answering questions about data visualizations using efficient bimodal fusion. *CoRR*, abs/1908.01801, 2019. URL <http://arxiv.org/abs/1908.01801>.
- [28] Samira Ebrahimi Kahou, Adam Atkinson, Vincent Michalski, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. *CoRR*, abs/1710.07300, 2017. URL <http://arxiv.org/abs/1710.07300>.
- [29] Shankar Kanthara, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq R. Joty. Chart-to-text: A large-scale benchmark for chart summarization. *ArXiv*, abs/2203.06486, 2022.
- [30] 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*, CHI ’20, page 1–13, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450367080.

doi: 10.1145/3313831.3376467. URL <https://doi.org/10.1145/3313831.3376467>.

- [31] 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*, pages 1–13, 2020.
- [32] Dae Hyun Kim, Vidya Setlur, and Maneesh Agrawala. Towards understanding how readers integrate charts and captions: A case study with line charts. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–11, 2021.
- [33] Hei Law and Jia Bin Deng. Cornernet: Detecting objects as paired keypoints. *International Journal of Computer Vision*, 128:642–656, 2019.
- [34] Matan Levy, Rami Ben-Ari, and Dani Lischinski. Classification-regression for chart comprehension. *CoRR*, abs/2111.14792, 2021. URL <https://arxiv.org/abs/2111.14792>.
- [35] Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenetorp, and Sebastian Riedel. PAQ: 65 Million Probably-Asked Questions and What You Can Do With Them. *Transactions of the Association for Computational Linguistics*, 9:1098–1115, 10 2021. ISSN 2307-387X. doi: 10.1162/tacl_a_00415. URL https://doi.org/10.1162/tacl_a_00415.
- [36] Xiaoya Li, Xiaofei Sun, Yuxian Meng, Junjun Liang, Fei Wu, and Jiwei Li. Dice loss for data-imbalanced NLP tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 465–476, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.45. URL <https://aclanthology.org/2020.acl-main.45>.
- [37] Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020.
- [38] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. URL <http://arxiv.org/abs/1405.0312>.
- [39] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *CoRR*, abs/1908.02265, 2019. URL <http://arxiv.org/abs/1908.02265>.

- [40] Junyu Luo, Zekun Li, Jinpeng Wang, and Chin-Yew Lin. Chartocr: Data extraction from charts images via a deep hybrid framework. *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1916–1924, 2021.
- [41] Mateusz Malinowski and Mario Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NIPS*, 2014.
- [42] Ahmed Masry and Enamul Hoque. Integrating image data extraction and table parsing methods for chart question answering. *Chart Question Answering Workshop, in conjunction with the Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–5, 2021.
- [43] Gonzalo Gabriel Méndez, Miguel A. Nacenta, and Sebastien Vandenheste. *IVoLVER: Interactive Visual Language for Visualization Extraction and Reconstruction*, page 4073–4085. Association for Computing Machinery, New York, NY, USA, 2016. ISBN 9781450333627. URL <https://doi.org/10.1145/2858036.2858435>.
- [44] 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 (WACV)*, March 2020.
- [45] Richard Meyes, Melanie Lu, Constantin Waubert de Puiseau, and Tobias Meisen. Ablation studies in artificial neural networks, 2019. URL <https://arxiv.org/abs/1901.08644>.
- [46] Thomas Müller, Francesco Piccinno, Massimo Nicosia, Peter Shaw, and Yasemin Altun. Answering conversational questions on structured data without logical forms. *ArXiv*, abs/1908.11787, 2019.
- [47] Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Nick Schoelkopf, Riley Kong, Xiangru Tang, Murori Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, and Dragomir Radev. Fetaqa: Free-form table question answering. *arXiv preprint arXiv:2104.00369*, 2021.
- [48] Kuntal Kumar Pal and Chitta Baral. Investigating numeracy learning ability of a text-to-text transfer model. In *EMNLP*, 2021.
- [49] Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1470–1480, Beijing, China,

- July 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1142. URL <https://aclanthology.org/P15-1142>.
- [50] Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1470–1480, Beijing, China, July 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1142. URL <https://www.aclweb.org/anthology/P15-1142>.
- [51] Jorge Poco and Jeffrey Heer. Reverse-engineering visualizations: Recovering visual encodings from chart images. *Comput. Graph. Forum*, 36(3):353–363, jun 2017. ISSN 0167-7055. doi: 10.1111/cgf.13193. URL <https://doi.org/10.1111/cgf.13193>.
- [52] Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018.
- [53] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [54] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- [55] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. *CoRR*, abs/1606.05250, 2016. URL <http://arxiv.org/abs/1606.05250>.
- [56] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *CoRR*, abs/1506.01497, 2015. URL <http://arxiv.org/abs/1506.01497>.
- [57] Ankit Rohatgi. Webplotdigitizer: Version 4.4. <https://automeris.io/WebPlotDigitizer>, 2020.
- [58] Adam Santoro, David Raposo, David G. T. Barrett, Mateusz Malinowski, Razvan Pascanu, Peter W. Battaglia, and Timothy P. Lillicrap. A simple neural network module for relational reasoning. *CoRR*, abs/1706.01427, 2017. URL <http://arxiv.org/abs/1706.01427>.

- [59] M. Savva, Nicholas Kong, Arti Chhajta, Li Fei-Fei, Maneesh Agrawala, and J. Heer. Revision: automated classification, analysis and redesign of chart images. *Proceedings of the 24th annual ACM symposium on User interface software and technology*, 2011.
- [60] David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical reasoning abilities of neural models. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1gR5iR5FX>.
- [61] Martin Schmitt, Leonardo F. R. Ribeiro, Philipp Dufter, Iryna Gurevych, and Hinrich Schütze. Modeling graph structure via relative position for text generation from knowledge graphs. In *Proceedings of the Fifteenth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-15)*, pages 10–21, Mexico City, Mexico, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.textgraphs-1.2. URL <https://aclanthology.org/2021.textgraphs-1.2>.
- [62] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 464–468, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2074. URL <https://aclanthology.org/N18-2074>.
- [63] N. Siegel, Zachary Horvitz, Roie Levin, S. Divvala, and Ali Farhadi. Figureseer: Parsing result-figures in research papers. In *ECCV*, 2016.
- [64] 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)*, pages 3275–3284, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.264. URL <https://www.aclweb.org/anthology/2020.emnlp-main.264>.
- [65] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019.
- [66] Damien Teney, Lingqiao Liu, and Anton van den Hengel. Graph-structured representations for visual question answering. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3233–3241, 2017.

- [67] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL <http://arxiv.org/abs/1706.03762>.
- [68] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [69] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019.
- [70] Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. LUKE: Deep contextualized entity representations with entity-aware self-attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6442–6454, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.523. URL <https://aclanthology.org/2020.emnlp-main.523>.
- [71] Peng-Jian Yang, Ying Chen, Yuechan Chen, and Daniel Matthew Cer. Nt5?! training t5 to perform numerical reasoning. *ArXiv*, abs/2104.07307, 2021.
- [72] Yi Yang, Wen-tau Yih, and Christopher Meek. WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2013–2018, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1237. URL <https://aclanthology.org/D15-1237>.
- [73] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alexander J. Smola. Stacked attention networks for image question answering. *CoRR*, abs/1511.02274, 2015. URL <http://arxiv.org/abs/1511.02274>.
- [74] Pengcheng Yin, Zhengdong Lu, Hang Li, and Kao Ben. Neural enquirer: Learning to query tables in natural language. In *Proceedings of the Workshop on Human-Computer Question Answering*, pages 29–35, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/W16-0105. URL <https://aclanthology.org/W16-0105>.
- [75] Pengcheng Yin, Graham Neubig, Wen tau Yih, and Sebastian Riedel. Tabert: Pretraining for joint understanding of textual and tabular data, 2020.
- [76] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-

language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5579–5588, June 2021.

- [77] Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from natural language using reinforcement learning, 2017.

A Appendix

A.1 Pretraining Tasks Templates

In this section, we present our pretraining tasks templates that we used to generate the queries. Also, we would like to highlight that not all the queries are applicable to every chart type. Some queries may be applicable to only one type (e.g. pie charts).

- Visual Attributes Retrieval
 1. position: first bar from the top/left in the second group from the top/left
 2. position: first bar from the bottom/right in the second group from the bottom/right
 3. position: second bar from the bottom/right in the first group from the bottom/right
 4. position: second bar from the right/bottom in the first group from the left/top
 5. position: topmost/leftmost bar

6. position: bottommost/rightmost bar
7. position: second bar from the top/left
8. position: second bar from the right/bottom
9. position: leftmost topmost bar
10. position: leftmost bottommost bar
11. position: rightmost topmost bar
12. position: rightmost bottommost bar
13. position: leftmost <color> data
14. position: rightmost <color> data
15. position: second from the left <color> data
16. position: second from the right <color> data
17. color: which legend represented by <color>
18. color: what is the color of <legend>

- Numerical Reasoning

1. Which one is greater, <x1> or <x2>?
2. Divide the sum of largest and lowest values by <n>
3. When did line <legend - label> peak?

4. What is the difference between maximum and minimum of <legend - label>?
5. Sum pie segments above <value>
6. What is the sum of top three values?
7. What is the median/mode of <legend - label>?
8. What is the negative peak of <legend - label>?
9. What is the largest/smallest value of <legend - label>?
10. Which two x-axis labels of <legend - label> sums up to <value>?
11. What is the sum of the second highest and second lowest value of <legend - label>?
12. Which x-axis label is second highest for <legend - label>?
13. What is the sum of two middle values of <legend - label>?
14. Which two x-axis labels of <legend - label> have a difference of <value> ?
15. What is the average of <legend - label> from <x - label - 1> to <x - label - 2>?
16. What is the average of the highest and lowest value of <legend - label - 1>?

17. What is the sum of the average of <legend - label - 1> and average of <legend - label - 2>?
18. What is the sum/difference of the maximum of <legend - label - 1> and minimum of <legend - label - 2>?
19. Which x-axis label has the maximum/minimum difference between <legend - label - 1> and minimum of <legend - label - 2>?
20. Which x-axis label witnessed the smallest value of <legend - label>?
21. Which label contains largest/smallest values across all labels?
22. Sum up the medians of all the data series in this chart
23. What is the average of all values above <value>?
24. What is the sum of the largest and smallest difference between <legend - label - 1> and <legend - label - 2>?
25. What is the maximum/minimum difference between <legend - label - 1> and <legend - label - 2>?
26. What is the ratio of the largest to the smallest pie segment?
27. What is the ratio of the two largest/smallest segments?

- Visual Reasoning

1. What is the difference between the leftmost and rightmost bars?

2. What is the sum of the bars in the second group from the left?
3. What is the sum of the bars in the first group from the right?
4. What is the ratio between the two leftmost bars?
5. What is the difference between the rightmost $\langle \text{color} - 1 \rangle$ bar and leftmost $\langle \text{color} - 2 \rangle$ bar?
6. What is the average of $\langle \text{color} \rangle$ bars values?
7. How many $\langle \text{color} \rangle$ bars are larger than $\langle N \rangle$?
8. What is the average of the bars in the second group from the right?
9. How many bars in the leftmost group have a value over $\langle N \rangle$?
10. What does the $\langle \text{color} \rangle$ represent?
11. What is the median value of the $\langle \text{color} \rangle$ bars/line?
12. What is the average of the $\langle \text{color} - 1 \rangle$ sum and $\langle \text{color} - 2 \rangle$ sum?
13. What is the average of the $\langle \text{color} - 1 \rangle$ median and $\langle \text{color} - 2 \rangle$ median?
14. What is the least difference between the $\langle \text{color} - 1 \rangle$ and $\langle \text{color} - 2 \rangle$ bars/line?
15. What is the ratio between the leftmost and rightmost bar in the first group from the left?

16. What is the maximum value in the <color> bars/line?
17. What is the minimum value in the <color> bars/line?
18. What is the sum of <color> bars/line?
19. What is the difference between the maximum values of the two leftmost bar groups?
20. Sum of the first <color - 1> and last <color - 2> bars/line points?
21. Difference between the two lowest <color> bars.
22. Add largest and smallest <color> line/bar values and divide by 2
23. What is the value of <color> line/bars in <x - axis - label>?
24. Sum/Average of <color - 1> and <color - 2> values in <x - axis - label>? Sum of highest points in <color - 1> and <color - 2> lines/bars
25. Which color has the highest/smallest values?
26. How many values are equal in <color - 1> line/bar?
27. Sum two rightmost values of <color> graph
28. Product of two smallest values in the graph
29. Sum of lowest and median values of <color> graph/bars
30. When did <color> line reached the peak

31. What is the average of the rightmost three points of <color> line
32. How many <color> data points are above <value>?
33. What's the ratio of the largest and the third/second-largest <color> bar?
34. Is the sum of lowest value of <color - 1> and <color - 2> bar greater than largest value of <color - 3> bar?
35. Is the median value of <color - 1> bars greater than the median value of <color - 2> bars?
36. Is the median of all the <color - 1> bars greater than the largest value of <color - 2> bar?
37. What's the product of <color> bars in India and Japan?
38. Is the sum of the two middle bars greater than the sum of top and bottom bars?
39. What's the ratio of the <x - axis - label - 1> <color - 1> bar and the <x - axis - 2> <color - 2> bar?
40. Is the total of all <color - 1> bars greater than the total of all <color - 2> bars?
41. Take the sum of the two smallest <color - 1> bars and smallest <color - 2> bars, deduct the smaller value from the larger value, what's the result?
42. What is the sum/average of two smallest/largest <color> bars?

43. What is the ratio of `<color - 1>` and `<color - 2>` segments?

44. What segment is represented by `<color>`?