

# PlotQA: Reasoning over Scientific Plots

Nitesh Methani \*

Pritha Ganguly\*

Mitesh M. Khapra

Pratyush Kumar

Department of Computer Science and Engineering  
Robert Bosch Centre for Data Science and AI (RBC-DSAI)  
Indian Institute of Technology Madras, Chennai, India  
{nmethani, prithag, miteshk, pratyush}@cse.iitm.ac.in

## Appendices

The supplementary material is organised as follows: Section 1 presents the detailed statistics of the PlotQA dataset. Section 2 describes the methodology of knowledge graph creation from structured data. In Section 3, we provide sample plots from the PlotQA dataset. In Section 4, we list all the 74 question templates that were formulated from the crowd sourced questions. In Section 5, we further analyze our proposed pipeline and discuss some of the errors that occur in each stage.

### 1. PlotQA Statistics

In this section, we present the detailed statistics of the PlotQA dataset. In particular, we report the dataset distribution according to the different plot types (vertical bar (vbar), horizontal bar (hbar), line and dot-line) and different question templates (structural understanding, data-retrieval, and reasoning) present in the PlotQA dataset. Table 1 presents the question distribution, categorized by their template type and detailed statistics of the different plots present in different dataset splits.

### 2. Construction of knowledge graph from structured data

Following [1], we convert the semi-structured table into a knowledge graph which has two types of nodes *viz.* row nodes and entity nodes. The rows of the table become row nodes, whereas the cells of each row become the entity nodes in the graph. Directed edges exist from the row nodes to the entity nodes of that column and the corresponding table column header act as edge-labels. An example of knowledge graph of the semi-structured table given in Figure 1a is shown in Figure 1b. For reasoning on the knowledge graph, we adopted the same methodology as given in [1]. The questions are converted to a set of candidate logical

forms by applying compositional semantic parsing. Each of these logical forms is then ranked using a log-linear model and the highest ranking logical form is applied to the knowledge graph to get the final answer.

Years	Brazil	Iceland	Kazakhstan	Thailand
1996	13.174405	7.895492	19.037112	8.821224
1997	11.680978	7.642265	14.860660	9.322298
1998	9.304022	5.221005	15.246865	8.174856
1999	11.370439	5.453609	9.747697	9.685002

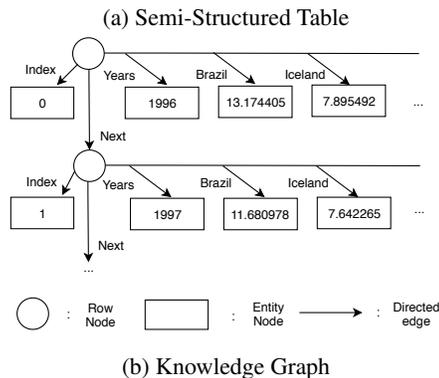


Figure 1: An example of the knowledge graph constructed from the semi-structured table.

### 3. Samples from the PlotQA dataset

Few examples of the {plot, question, answer} triplets from the PlotQA dataset are shown in Figure 3. For each of the plots, most of the question templates discussed in section 4 are applicable but depending on the context of the plot, their language varies from its surface form.

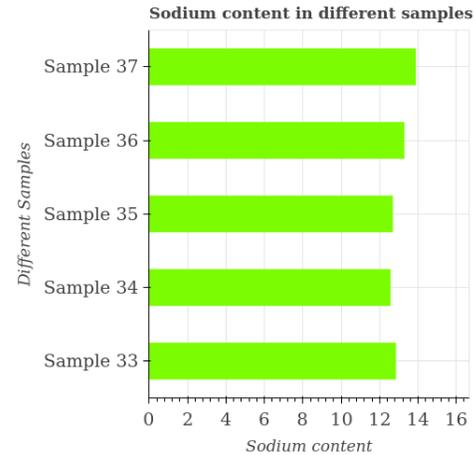
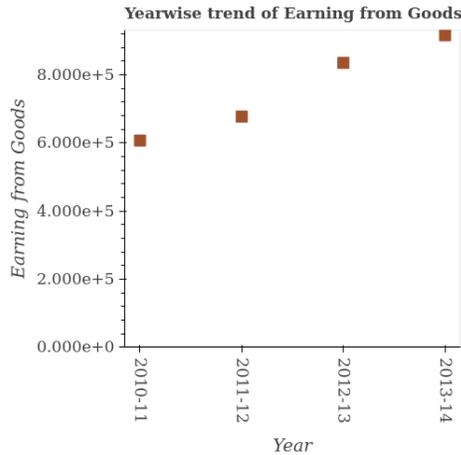
### 4. Question Templates

In this section, we present the 74 question templates which we have used for the question generation. Note that,

\*The first two authors have contributed equally

Dataset Split	Plot Types				Question Types			Answer Types		
	vbar	hbar	line	dot-line	Structural	Data-Retrieval	Reasoning	Yes/No	Fixed vocab.	Open vocab.
Train	52,463	52,700	25,897	26,010	871,782	2,784,041	16,593,656	784,115	3,095,774	16,369,590
Validation	11,249	11,292	5,547	5,571	186,994	599,573	3,574,081	167,871	600,424	3,592,353
Test	11,242	11,292	5,549	5,574	186,763	596,359	3,559,392	167,727	667,742	3,507,045

Table 1: Detailed Statistics for different splits of the PlotQA dataset.

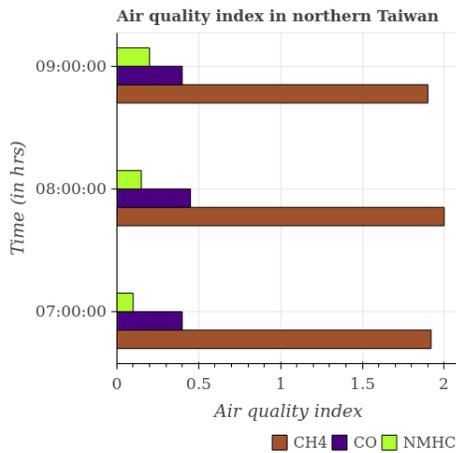


**Q1:** What is the difference between the highest and the second highest amount of earnings from goods?

**A:**  $0.9e + 5$

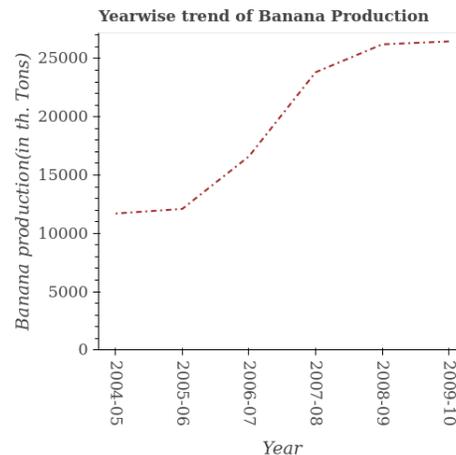
**Q2:** What is the ratio of the sodium content in Sample 37 to that in Sample 33?

**A:** 1.086



**Q3:** What is the average air quality index value for NHMC per hour?

**A:** 0.167



**Q4:** Does the amount of banana production monotonically increase over the years?

**A:** Yes

Figure 2: Sample {plot, question, answer} triplet present in the PlotQA dataset.

not all question templates are applicable to each and every type of plot. Also depending on the context of the plot, the question varies from the template's surface form.

### 1. Structural Understanding :

1. Does the graph contain any zero values?

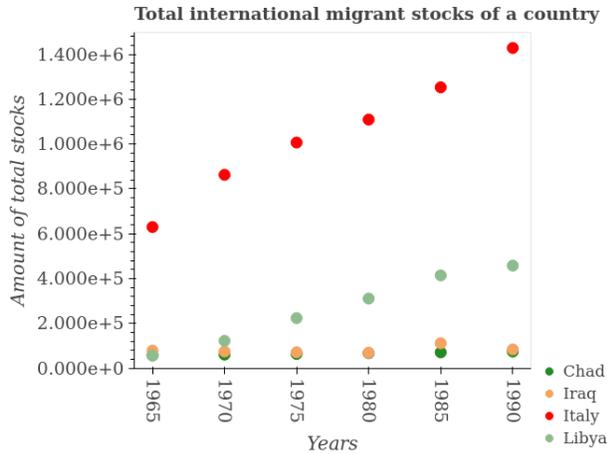
2. Does the graph contain grids ?

3. Where does the legend appear in the graph ?

4. How many legend labels are there?

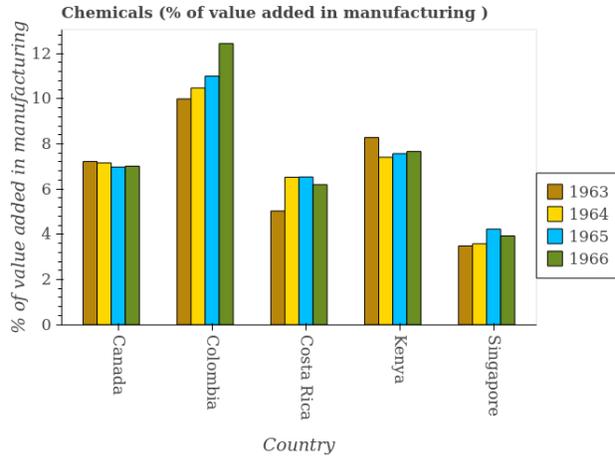
5. How are the legend labels stacked?

6. How many <plural form of X\_label> are there in



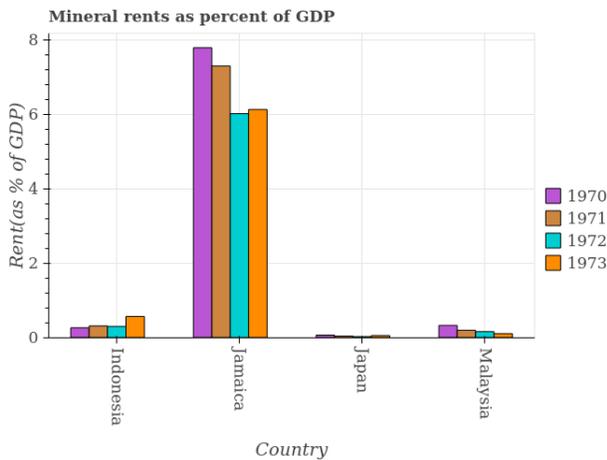
**Q5:** What is the difference between two consecutive major ticks on the Y-axis?

**A:** 2.000e + 5



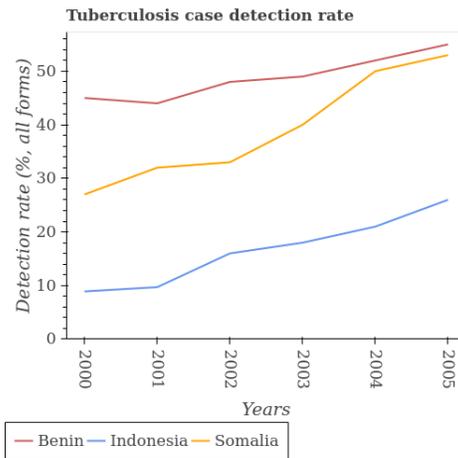
**Q6:** In how many cases, is the number of bars for a given year not equal to the number of legend labels?

**A:** 0



**Q7:** In how many countries, is the mineral rent (as % of GDP) in 1970 greater than the average mineral rent (as % of GDP) in 1970 taken over all countries?

**A:** 1



**Q8:** What is the total tuberculosis detection rate in Indonesia?

**A:** 101

Figure 3: Few more sample {plot, question, answer} triplet present in the PlotQA dataset.

- the graph?
- How many <figure-type>s are there?
  - How many different colored <figure-type>s are there?
  - How many groups of <figure-type>s are there?
  - Are the number of bars on each tick equal to the number of legend labels?
  - Are the number of bars in each group equal?
  - How many bars are there on the  $i^{th}$  tick from the left?
  - How many bars are there on the  $i^{th}$  tick from the right?

- How many bars are there on the  $i^{th}$  tick from the top?
- How many bars are there on the  $i^{th}$  tick from the bottom?
- Are all the bars in the graph horizontal?
- How many lines intersect with each other?
- Is the number of lines equal to the number of legend labels?

## 2. Data Retrieval :

- What does the  $i^{th}$  bar from the left in each group represent?

2. What does the  $i^{th}$  bar from the right in each group represent?
3. What does the  $i^{th}$  bar from the top in each group represent?
4. What does the  $i^{th}$  bar from the bottom in each group represent?
5. What is the label of the  $j^{th}$  group of bars from the left?
6. What is the label of the  $j^{th}$  group of bars from the top?
7. Does the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  monotonically increase over the  $\langle plural\ form\ of\ X\_label \rangle$  ?
8. What is the difference between two consecutive major ticks on the Y-axis ?
9. Are the values on the major ticks of Y-axis written in scientific E-notation ?
10. What is the title of the graph ?
11. Does  $\langle legend\_label \rangle$  appear as one of the legend labels in the graph ?
12. What is the label or title of the X-axis ?
13. What is the label or title of the Y-axis ?
14. In how many cases, is the number of  $\langle figure\_type \rangle$  for a given  $\langle X\_label \rangle$  not equal to the number of legend labels ?
15. What is the  $\langle Y\_value \rangle$  in/of  $\langle i^{th}\ X\_tick \rangle$  ?
16. What is the  $\langle Y\_value \rangle$  of the  $i^{th}$   $\langle legend\_label \rangle$  in  $\langle i^{th}\ X\_tick \rangle$  ?
17. Does the  $\langle Y\_label \rangle$  monotonically increase over the  $\langle plural\ form\ of\ X\_label \rangle$  ?
18. Is the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  strictly greater than the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  over the  $\langle plural\ form\ of\ X\_label \rangle$  ?
19. Is the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  strictly less than the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  over the  $\langle plural\ form\ of\ X\_label \rangle$  ?
5. Across all  $\langle plural\ form\ of\ X\_label \rangle$ , what is the maximum  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  ?
6. Across all  $\langle plural\ form\ of\ X\_label \rangle$ , what is the minimum  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  ?
7. In which  $\langle singular\ form\ of\ X\_label \rangle$  was the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  maximum ?
8. In which  $\langle singular\ form\ of\ X\_label \rangle$  was the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  minimum ?
9. What is the sum of  $\langle title \rangle$  ?
10. What is the difference between the  $\langle Y\_label \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  and  $\langle j^{th}\ x\_tick \rangle$  ?
11. What is the average  $\langle Y\_label \rangle$  per  $\langle singular\ form\ of\ X\_label \rangle$  ?
12. What is the median  $\langle Y\_label \rangle$  ?
13. What is the total  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  in the graph?
14. What is the difference between the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  and that in  $\langle j^{th}\ x\_tick \rangle$  ?
15. What is the difference between the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  and the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  in  $\langle j^{th}\ x\_tick \rangle$  ?
16. What is the average  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  per  $\langle singular\ form\ of\ X\_label \rangle$  ?
17. In the year  $\langle i^{th}\ x\_tick \rangle$ , what is the difference between the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  and  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  ?
18. What is the difference between the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  and  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  ?
19. In how many  $\langle plural\ form\ of\ X\_label \rangle$ , is the  $\langle Y\_label \rangle$  greater than  $\langle N \rangle$  units ?
20. Do a majority of the  $\langle plural\ form\ of\ X\_label \rangle$  between  $\langle i^{th}\ x\_tick \rangle$  and  $\langle j^{th}\ x\_tick \rangle$  (inclusive/exclusive) have  $\langle Y\_label \rangle$  greater than  $\langle N \rangle$  units ?

### 3. Reasoning :

1. Across all  $\langle plural\ form\ of\ X\_label \rangle$ , what is the maximum  $\langle Y\_label \rangle$  ?
2. Across all  $\langle plural\ form\ of\ X\_label \rangle$ , what is the minimum  $\langle Y\_label \rangle$  ?
3. In which  $\langle X\_label \rangle$  was the  $\langle Y\_label \rangle$  maximum ?
4. In which  $\langle X\_label \rangle$  was the  $\langle Y\_label \rangle$  minimum ?
21. What is the ratio of the  $\langle Y\_label \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  to that in  $\langle j^{th}\ x\_tick \rangle$  ?
22. Is the  $\langle Y\_label \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  less than that in  $\langle j^{th}\ x\_tick \rangle$  ?
23. In how many  $\langle plural\ form\ of\ X\_label \rangle$ , is the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  greater than  $\langle N \rangle$  units?
24. What is the ratio of the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  in  $\langle i^{th}\ x\_tick \rangle$  to that in  $\langle j^{th}\ x\_tick \rangle$  ?

25. Is the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  in  $\langle i^{th} x\_tick \rangle$  less than that in  $\langle j^{th} x\_tick \rangle$  ?
26. Is the difference between the  $\langle Y\_label \rangle$  in  $\langle i^{th} x\_tick \rangle$  and  $\langle j^{th} x\_tick \rangle$  greater than the difference between any two  $\langle plural\ form\ of\ X\_label \rangle$  ?
27. What is the difference between the highest and the second highest  $\langle Y\_label \rangle$  ?
28. Is the sum of the  $\langle Y\_label \rangle$  in  $\langle i^{th} x\_tick \rangle$  and  $\langle (i + 1)^{th} x\_tick \rangle$  greater than the maximum  $\langle Y\_label \rangle$  across all  $\langle plural\ form\ of\ X\_label \rangle$  ?
29. What is the difference between the highest and the lowest  $\langle Y\_label \rangle$  ?
30. In how many  $\langle plural\ form\ of\ X\_label \rangle$ , is the  $\langle Y\_label \rangle$  greater than the average  $\langle Y\_label \rangle$  taken over all  $\langle plural\ form\ of\ X\_label \rangle$  ?
31. Is the difference between the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  in  $\langle i^{th} x\_tick \rangle$  and  $\langle j^{th} x\_tick \rangle$  greater than the difference between the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  in  $\langle i^{th} x\_tick \rangle$  and  $\langle j^{th} x\_tick \rangle$  ?
32. What is the difference between the highest and the second highest  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  ?
33. What is the difference between the highest and the lowest  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  ?
34. In how many  $\langle plural\ form\ of\ X\_label \rangle$ , is the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  greater than the average  $\langle Y\_label \rangle$  of/in  $\langle legend\_label \rangle$  taken over all  $\langle plural\ form\ of\ X\_label \rangle$  ?
35. Is it the case that in every  $\langle singular\ form\ of\ X\_label \rangle$ , the sum of the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  and  $\langle legend\_label2 \rangle$  is greater than the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label3 \rangle$  ?
36. Is the sum of the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  in  $\langle i^{th} x\_tick \rangle$  and  $\langle j^{th} x\_tick \rangle$  greater than the maximum  $\langle Y\_label \rangle$  of/in  $\langle legend\_label2 \rangle$  across all  $\langle plural\ form\ of\ X\_label \rangle$  ?
37. Is it the case that in every  $\langle singular\ form\ of\ X\_label \rangle$ , the sum of the  $\langle Y\_label \rangle$  of/in  $\langle legend\_label1 \rangle$  and  $\langle legend\_label2 \rangle$  is greater than the sum of  $\langle Y\_label \rangle$  of  $\langle legend\_label3 \rangle$  and  $\langle Y\_label \rangle$  of  $\langle legend\_label4 \rangle$  ?

## 5. Some failure cases

In this section, we visualize the output of each stage in our multistage pipeline and show some interesting failure cases with the help of an example.

- **VED Stage:** Although the bounding boxes predicted by Faster R-CNN coupled with Feature Pyramid Network

(FPN) fit reasonably well at an IOU of 0.5, it is not acceptable as the values extracted from these bounding boxes will lead to incorrect table generation and subsequent QA. Example: In Figure 4b, consider the bar representing the “Indoor User Rating” value for “Vodafone”. The overlap between the ground-truth box (blue) and the predicted box (red) is higher than 0.5 but the values extracted from the detected box is 4.0 as opposed to the actual value which is 3.73. Another interesting failure case is shown in Figure 6b. There are multiple overlapping data points and the model is able to detect only one of the points. This leads to incomplete table generation as shown in Figure 7b where the values for Liberia for the years 2008, 2009 and 2010 could not be extracted. This small error might be acceptable for other VQA tasks but for PlotQA these small errors will escalate to multiple incorrect answers.

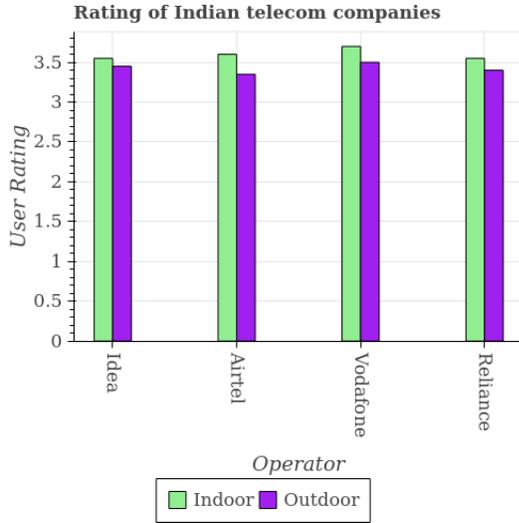
- **OCR stage:** A slight misalignment in the bounding boxes predicted by VED module causes significant errors while extracting the text. Example: In Figure 4b, consider the box enclosing the legend label “Indoor”. The rightmost edge of the predicted bounding box is drawn over the letter “r”, which makes the OCR module incorrectly recognize the text as “Indoo”. A similar error is made while performing OCR on the X-axis title which is read as “Dperator” instead of “Operator”. Consider another example where there are misaligned bounding boxes on axes tick-labels as shown in Figure 6b. The values extracted are 200B, -2009 and -100 as opposed to the ground-truth values 2008, 2009 and 100. This slight error leads to incorrect column name in the subsequent generated tables (Figure 5b and Figure 7b) and incorrect answers to all the questions pertaining to these labels as shown in Table 2 and Tale 3.

- **SIE stage:** Figure 5a shows the oracle table which is generated by using the ground-truth annotations and Figure 5b shows the table generated after passing the plot image through the different stages of our proposed multistage pipeline. It is evident from the generated table that the errors propagated from the VED and the OCR stage has lead to an incorrect table generation.

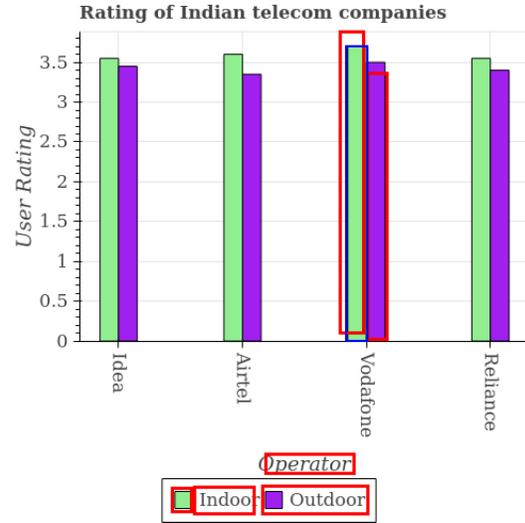
- **QA stage:** In Table 2 and Table 3 we compare the answer predictions made by different models with the ground-truth answer on randomly sampled questions. Note that, our proposed model combines the complementary strengths of both, QA-as-classification and QA-as-multistage pipeline, models.

## References

- [1] P. Pasupat and P. Liang. Compositional semantic parsing on semi-structured tables. In *ACL*, 2015.



(a) Input plot image



(b) Few examples of the predicted bounding boxes

Figure 4: Errors made by the VED stage (highlighted in red).

Operator	Indoor	Outdoor
Idea	3.55	3.45
Airtel	3.62	3.38
Vodafone	3.73	3.51
Reliance	3.57	3.41

(a) Oracle table generated using ground-truth annotations

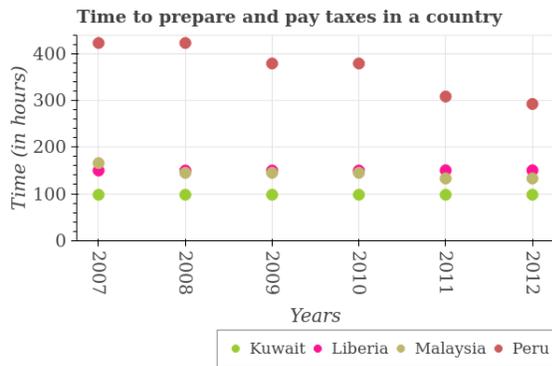
Doperator	Indoo	Outdoor
Idea	3.53	3.45
Airtel	3.60	3.40
Vodafone	4.0	3.35
Reliance	3.62	3.40

(b) Generated Semi-Structured Table

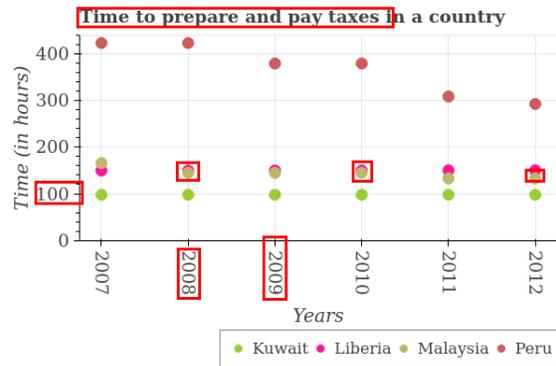
Figure 5: Errors made by the OCR and SIE stage (highlighted in red). Note that most of these errors have been propagated from the VED stage.

Question	Ground Truth	QA as classification	Multistage Pipeline	Our Model
Q1. What is the average indoor user rating per operator?	3.62	500	3.54	3.54
Q2. What is the total user rating for Vodafone in the graph?	7.24	5	7.35	7.35
Q3. What is the label or title of the X-axis?	Operator	Years	Dperator	Dperator
Q4. What is the indoor user rating of Airtel?	3.62	0	3.60	3.60
Q5. How many groups of bars are there?	4	4	4	4
Q6. What is the ratio of indoor user rating of Airtel to that of the outdoor user rating of Vodafone?	1.03	No	3.35	3.35
Q7. What is the difference between the highest and lowest outdoor user rating?	0.13	1.5	0.10	0.10
Q8. Does "Indoor" appear as one of the legend-labels in the graph?	Yes	Yes	1.0	Yes
Q9. For how many operators are the indoor user rating greater than the average outdoor user rating taken over all operators?	4	4	3.4	4
Q10. Is the sum of outdoor user rating in Reliance and Idea greater than the maximum outdoor user rating across all operators?	Yes	Yes	6.85	Yes

Table 2: Answers predicted by different models on the sample questions.



(a) Input plot image



(b) Few examples of the predicted bounding boxes

Figure 6: Errors made by the VED stage (highlighted in red).

Years	Kuwait	Liberia	Malaysia	Peru
2007	98	150	166	424
2008	98	150	145	424
2009	98	150	145	380
2010	98	150	145	380
2011	98	150.5	133	309
2012	98	150.5	133	293

(a) Oracle table generated using ground-truth annotations

Years	Kuwait	Liberia	Malaysia	Peru
2007	100-	150	165	420
200B	100-	-	155	420
-2009	100-	-	155	380
2010	100-	-	153	383
2011	100-	155	-	310
2012	100-	155	-	295

(b) Generated Semi-Structured Table

Figure 7: Errors made by the OCR and SIE stage (highlighted in red). Note that most of these errors have been propagated from the VED stage.

Question	Ground Truth	QA as classification	Multistage Pipeline	Our Model
<b>Q1.</b> How are the legend-labels stacked?	horizontal	horizontal	horizontal	horizontal
<b>Q2.</b> What is the time (in hours) required to prepare and pay taxes in Kuwait in 2008?	98	100	100-	100-
<b>Q3.</b> What is the difference between in time (in hours) required to prepare and pay taxes in Kuwait in 2009 and the time (in hours) required to prepare and pay taxes in Liberia in 2008?	-52	-0.5	100-	100-
<b>Q4.</b> What is the difference between the highest and the lowest time (in hours) required to prepare and pay taxes in Peru?	131	500	125	125
<b>Q5.</b> Is it the case that every year the sum of time (in hours) required to prepare and pay taxes in Peru and Kuwait is greater than the sum of the time (in hours) required to prepare and pay taxes in Liberia and Malaysia?	Yes	Yes	320	Yes

Table 3: Answers predicted by different models on the sample questions.