# Towards VQA Models That Can Read

Amanpreet Singh[1], Vivek Natarajan, Meet Shah[1], Yu Jiang[1], Xinlei Chen[1],
Dhruv Batra[1,2], Devi Parikh[1,2], and Marcus Rohrbach[1]

[1]Facebook AI Research, [2]Georgia Institute of Technology

## 1. OCR and Answer Space Analysis

We perform the following analysis on TextVQA's validation set of size 5,734. We find that 44.9% (2575) of LoRRA's predicted answers are from OCR tokens (i.e., using the copy mechanism). The remaining 55.1% of predicted answers are from the SA. This shows that our approach does in fact rely heavily on what it reads in the image, and relies on its copy mechanism to generalize and produce answers that have never been seen before or are rare in the training data. While predicting answers from OCR tokens, the model gets the entire answer string correct 27% (696) of the time, and partially correct (i.e., matches one word in answer) 11% (284) of the time. The percentage of partially correct answers shows the possibilities of getting better results by using n-grams of OCR tokens or spelling correction for fixing incorrect OCR predictions. Similarly, while predicting from the answer space, the model gets the answer correct 22.4% (709/3,159).

We find that 30.6% (1759) of questions have their answers in OCR tokens. For these questions, LoRRA chooses to predict from OCR tokens, 68% (1,208/1,759) of the times and correctly answers 57.5% (697/1,208) of these. Similarly, 48% (2794) of questions have their answers in SA. For these questions, LoRRA chooses to predict from LA 66.75% (1,865/2,794) of the times and gets 38% (710/1,865) of these correct.

81% questions in TextVQA's validation set have more than or equal to 2 OCR tokens. Among these 4645 questions, LoRRA chooses to copy from OCR tokens 49.7% (2,309) and gets 24.3% (560/2,309) of these correct. This suggests that LoRRA doesn't randomly copy OCR token from list of available tokens.

## 2. TextVQA Examples and LoRRA Predictions

In Fig. 1, we show representative examples from our TextVQA dataset along with the predictions of Pythia+LoRRA. Each example shows the ground truth answer, the predictions from LoRRA, whether the answer prediction was from OCR tokens or the pre-determined answer space, and attention weights for each of the OCR tokens.
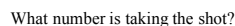
The examples indicate the following points:

- The model is able to successfully answer questions about times, dates, brands, cities and places, and is often able to correctly spell them even if the OCR tokens had them misspelled (as it is able to predict the correct answer from answer space). See Fig. 1k (short hand's hour), Fig. 1g (birthday date), Fig. 1r (no turn on red sign), Fig. 1s (city "london" correct from the large text), Fig. 1o (samsung).

- The model is able to successfully answer questions involving colors and spatial reasoning. See Fig. 1e (player on the right), Fig. 1f (location of 25 coin), Fig. 1c (location of banner).

- Model is able to distinguish objects based on spatial positioning, colors and relative positioning well. See Fig. 1q where the model needs to identify the correct sign based on multiple colors, or Fig. 1r where the model needs to identify the correct sign in the red circle and shows that the model is not biased toward "stop" as other VQA models are, or Fig. 1a where model needs to predict the correct number based on spatial reasoning from the choice of 7 or 14, or Fig. 1f where the model needs to predict the correct coin's value based on the position mentioned.

- The model is also able to reason about basic sizes (less, greater, smallest) and shapes (circle). See Fig. 1f where the model needs to figure out which one is top and silver coin, or Fig. 1k where the model needs to figure out which one is the shorter hand, or Fig. 1q where the model needs to figure out which one is lowest measurement among four.

- The model often predicts an answer from the answer space as informed by OCR tokens. See Fig. 1k where the Pythia model which doesn't use OCR predicts 3 in this case, but our approach predicts 4 which is the correct answer.

- The model often answers questions about cities with "new york". See Fig. 1j where the model predicts New York instead of San Francisco. We have observed this bias in other city related questions as well.

- For yes/no questions, even though "yes" is the more common answer, the model does predict "no" frequently. See Fig. 1m, Fig. 1l.

- Sometimes when the answer is not in the answer space, but the partial answer is in OCR tokens, the model predicts the partial answer which is closest to the actual answer. See Fig. 1e where the model predicts "fly" instead of "fly emirates", or Fig. 1g where the model predicts only the birthday date "19", instead of "may 19". By construction our model can only copy a single OCR token, but our TextVQA dataset contains Q/A pairs which require copying multiple OCR tokens in the right order. Explore this is an interesting direction for future work.

- The model sometimes gets seemingly simple questions wrong by predicting generic answers. See Fig. 1h where the model can't predict "embossed" even though it is in OCR tokens, or see Fig. 1b where the model predicts most common letter "g" in the answer space instead of predicting based on "a-2" in the OCR tokens.

- The model has a strong dependency on the quality of OCR tokens produced. If the OCR module missed some text in the image, the model's output can be wrong. See Fig. 1i or Fig. 1p where the OCR tokens do not contain the ground truth answer or see Fig. 1u where the OCR system is unable to correctly read "irig" the second time.

## 3. Interface Screenshots

We show three stages of data collection pipeline in Fig. 2, Fig. 3, Fig. 4 and Fig. 5. Fig. 2 and Fig. 3 shows the entrance and first stage of our task which is used to identify and remove images without text in them. Fig 4 shows the second stage of our task which is used to collect questions on images with text. Finally, the third stage interface is shown in 5 which is used to collect the answer for a question about an image.

What number is taking the shot?

Ground Truth: 14 | Prediction: 14 | From: OCR Tokens

(a)

What letter is written on the blue part of the sign?

Ground Truth: a | Prediction: g | From: Answer Space

(b)

What is the brand name on the blue banner at the top of the image?

Ground Truth: watsons | Prediction: watsons | From: OCR Tokens

(c)

What country is the brut from?

Ground Truth: france | Prediction: france | From: Answer Space

(d)

What is on the jersey of the player on the right?

Ground Truth: fly emirates | Prediction: fly | From: OCR Tokens

(e)

How much is the top silver coin worth?

Ground Truth: 25 | Prediction: 25 | From: OCR Tokens

(f)

When is dad's birthday?

Ground Truth: may 19 | Prediction: 19 | From: OCR Tokens

(g)

What edition is this?

Ground Truth: embossed | Prediction: second | From: Answer Space

(h)

What brand of soda on the table?

Ground Truth: dr pepper | Prediction: aaro | From: OCR Tokens

(i)

What city are the giants from?

Ground Truth: san francisco | Prediction: new york | From: Answer Space

(j)

What number is the small hand on?

Ground Truth: 4 | Prediction: 4 | From: Answer Space

(k)

Is the text about physics?

Ground Truth: yes | Prediction: no | From: Answer Space

(l)

Are all the books with same title?

Ground Truth: no | Prediction: no | From: Answer Space

(m)

What letter is written in yellow on the red sign?

Ground Truth: m | Prediction: m | From: Answer Space

(n)

What is the brand of this smartphone?

Ground Truth: samsung | Prediction: samsung | From: Answer Space
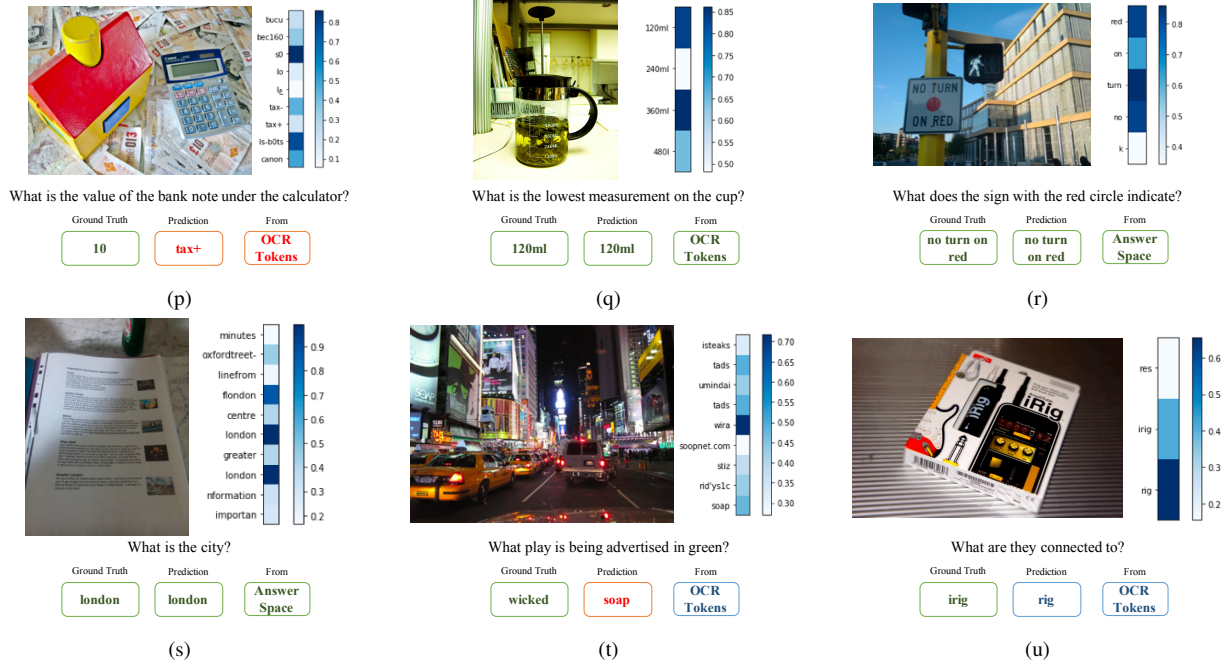
(o)

Figure 1: **TextVQA Examples and LoRRA's predictions on them.** We show multiple examples from TextVQA, ground truth answers, along with predictions from LoRRA, attention maps on OCR tokens and where LoRRA predicted the answer from (OCR tokens or predetermined answer space. Green, red, and blue boxes correspond to correct, incorrect, and partially correct answers, respectively. On the right side of each image, we show attention bars which depict attention weights (0-1) for each of the OCR tokens.
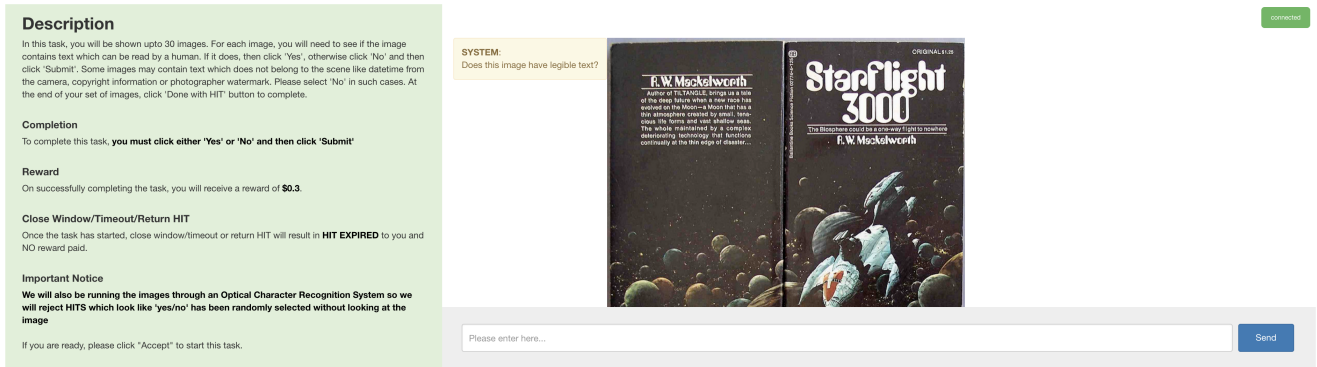


Figure 2: **Introduction page for our task.**

Figure 3: **Text detection main task.** First stage of our tasks is used to identify and remove images without text.
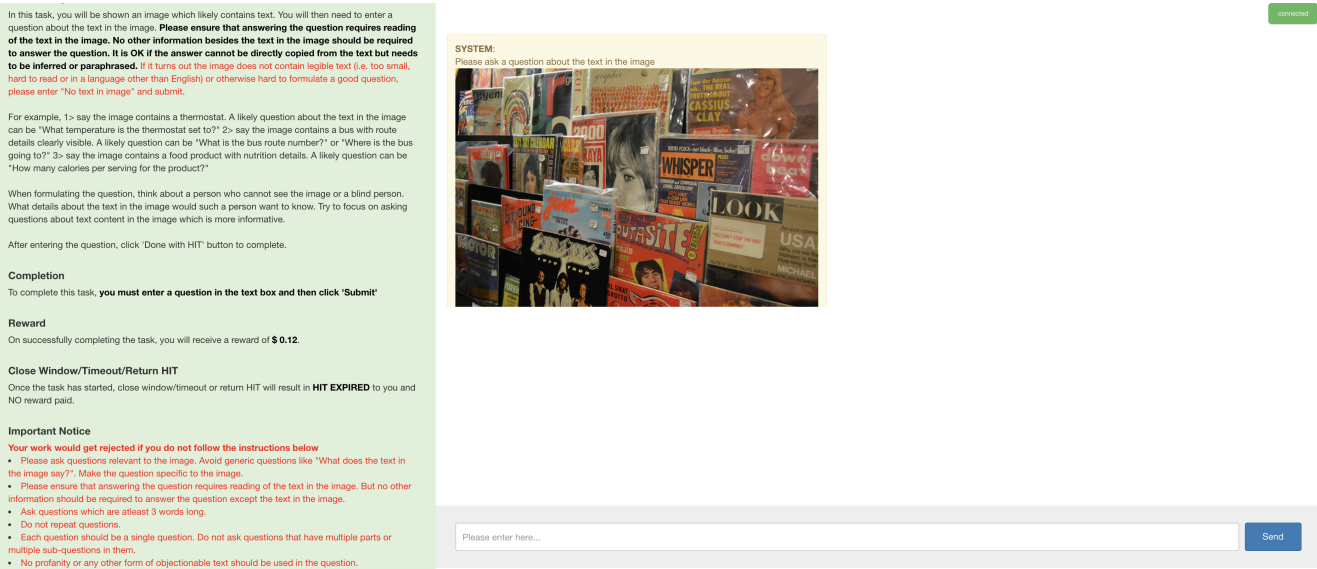


Figure 4: **Question task.** In the second stage, we ask workers to ask a question about an image whose answer requires reading text in the image. We provide instructions and rules to ensure that we get high quality questions.

**Description**

In this task, you will be shown upto 20 images. For each image, you will be shown a question about the text in the image. You will then need to answer the question based on the text in the image. **Please try to keep your answers brief. Your answers should be how most other people would answer the questions.**
Please use your browser shorcut to zoom into the image if necessary to see the text clearly

If it turns out that the input text is not a question, please select 'Not a question' checkbox and submit.

If it turns out that the image does not have text, please select 'No text in image' checkbox and submit.

If it turns out that answering the question does not require reading text in the image, please select 'Answering does not require reading text in the image' checkbox and submit.

If it turns out that answering the question requires information that you do not have, please select 'Unanswerable' checkbox and submit.

But please use this checkboxes only if you are absolutely certain. Make your best effort to answer the question and not use the checkboxes.

After entering the answer, click 'Done with HIT' button to complete.

**Completion**
To complete this task, **you must enter an answer to the question in the text box and then click 'Submit'**

**Close Window/Timeout/Return HIT**
Once the task has started, close window/timeout or return HIT will result in **HIT EXPIRED** to you and NO reward paid.

**Important Notice**
**Your work would get rejected if you do not follow the instructions below**
• Please ensure that you answer the question based on the image and the text in the image. If it looks like you have answered the question without looking at the image, you work will be rejected.
• Use the checkboxes only in case of the circumstances provided above. Try your best not to use it and provide a freeform answer in all other cases.
• Try to keep your answer as brief as possible.
• For yes/no questions, please just say 'yes/no'. Do not use terms like 'yeah/nah/ya' etc
• For numerical answers, please use digits unless you need to copy answer from the image text.
• If you need to speculate (e.g., "What just happened?"), provide an answer that most people would agree on.
• No profanity or any other form of objectionable text should be used in the question.

If you are ready, please click "Accept" to start this task.



Figure 5: **Answer task.** In the third stage, we ask workers to answer a question about the image.