


# Visual Reference Resolution using Attention Memory for Visual Dialog

## Supplementary Document

### A More Qualitative Results on MNIST Dialog

Input image	#	Questions	Answers
	1	How many digits in a yellow background are there ?	two
	2	How many digits with a stroke are there among them ?	two
	3	How many 5's are there among them ?	one
	4	What is the color of it ?	violet
	5	What is the style of the digit at the left of it ?	flat
	6	What is the number of the digit ?	1
	7	What is the background color of the digit ?	white
	8	What is the number of the digit at the left of it ?	1


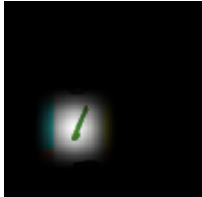



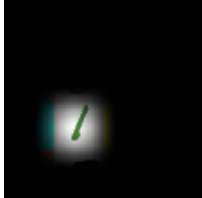



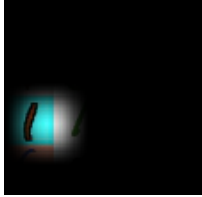
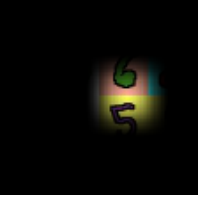
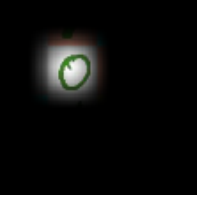
Retrieved attention from memory	Final attention	Manually modified retrieved attention	Final attention	Current question
				Q6
Predicted answer: 1		Predicted answer: 2		
				Q7
Predicted answer: white		Predicted answer: yellow		
				Q8
Predicted answer: 1		Predicted answer: 0		

Figure 1: Qualitative results of proposed model. Given input image and dialog at top row, retrieved attention (column 1), final attention (column 2) and predicted answers are presented for last three questions. Retrieved attentions focus on reference of ambiguous expressions and final attentions focus on region of target object based on relationship in question. Additionally, we manually modify retrieved attention by assigning high probability to randomly chosen single location (column 3) and show final attention obtained with modified retrieved attention (column 4). With modified retrieved attentions, final attentions also change accordingly. Note that misaligned retrieved attentions are corrected in final attention as depicted in Q8.



## B Qualitative Results on Visual Dialog (VisDial)






Input image	
	
Caption: <i>Cat sitting in small bowl on wood flooring indoors</i>	
QA pair, predicted answer and rank of GT answer	Attended image
<b>Q1:</b> <i>What color is the bowl ?</i> GT answer: <i>White</i> Predicted answer: <i>White</i> Rank of GT: <i>1</i>	
<b>Q2:</b> <i>Do you see any people?</i> GT answer: <i>No</i> Predicted answer: <i>No, just the cat</i> Rank of GT: <i>2</i>	
<b>Q3:</b> <i>What color is the cat ?</i> GT answer: <i>Grey, white, and black</i> Predicted answer: <i>Grey, black and white</i> Rank of GT: <i>6</i>	
<b>Q4:</b> <i>Is the cat wearing any collar?</i> GT answer: <i>No</i> Predicted answer: <i>No</i> Rank of GT: <i>1</i>	

Figure 3: Qualitative results of dialog in VisDial. Given the image and the caption at the top, a sequence of questions are presented in each row with the attention from the proposed model, the GT answer, the predicted answer of the model and the rank of the GT answer. The attention map is concentrated on the reference of question while it is distributed over the entire image when the reference of the question is not present in the image as in Q2 and Q4.

Input image



Caption:

*A hot dog covered in mustard and cheese sits next to French fries*

QA pair, predicted answer and rank of GT answer

Attended image

**Q1:** *Are hot dogs in bun ?*

GT answer: *Yes, there's only one though*

Predicted answer: *Yes*

Rank of GT: 9

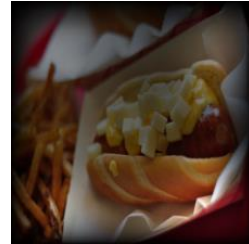


**Q2:** *Are they on plate ?*

GT answer: *No, it's in cardboard container*

Predicted answer: *No, it's in cardboard container*

Rank of GT: 1



**Q3:** *Are they stake fries ?*

GT answer: *No, they are shoestring fries*

Predicted answer: *No, they are shoestring fries*

Rank of GT: 1



**Q4:** *Are they on table ?*

GT answer: *They seem to be*

Predicted answer: *Yes*

Rank of GT: 17

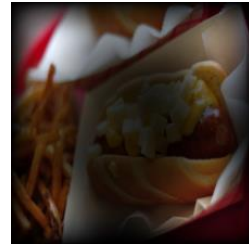


Figure 4: Qualitative results of another dialog in VisDial.









Dialog Information	Input image	Attended image
<p>Caption: <i>A large bear standing upright with mountains in the background</i></p> <p>Previous QA: <i>Is this the only bear here ? / yes</i></p> <p>Current question: <i>What color is it's fur ?</i></p> <p>GT answer: <i>Brown</i></p> <p>Predicted answer: <i>Brown</i></p> <p>Rank of GT: <i>1</i></p>		
<p>Caption: <i>A train that is on a large rail way</i></p> <p>Previous QA: <i>Is the train moving ? / No it is stopped</i></p> <p>Current question: <i>What color is the train ?</i></p> <p>GT answer: <i>It is white and red with some blue on it</i></p> <p>Predicted answer: <i>It is white and red with some blue on it</i></p> <p>Rank of GT: <i>1</i></p>		
<p>Caption: <i>An airplane parked in the middle of a runway</i></p> <p>Previous QA: <i>Can you see the airport ? / No</i></p> <p>Current question: <i>Is it a sunny day ?</i></p> <p>GT answer: <i>Yes</i></p> <p>Predicted answer: <i>Yes</i></p> <p>Rank of GT: <i>1</i></p>		
<p>Caption: <i>A zebra standing next to a wire beside a chain link fence</i></p> <p>Previous QA: <i>Is the zebra in a zoo ? / Maybe, can't tell</i></p> <p>Current question: <i>Is the zebra young or old ?</i></p> <p>GT answer: <i>Grown I think</i></p> <p>Predicted answer: <i>Young</i></p> <p>Rank of GT: <i>5</i></p>		

Figure 5: More qualitative results of questions in different dialogs.