

# Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning Supplement

Abhishek Das<sup>1,\*</sup>; Satwik Kottur<sup>2,\*</sup>; José M.F. Moura<sup>2</sup>, Stefan Lee<sup>3</sup>, Dhruv Batra<sup>1,4</sup>

<sup>1</sup>Georgia Institute of Technology, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>Virginia Tech, <sup>4</sup>Facebook AI Research

[visualdialog.org](http://visualdialog.org)

## Abstract

*This supplement is organized as follows:*

- *Sec. 1 performs the following ablation study – how well would different models perform if they were provided human dialog as input? Specifically, how good are the different models on the task of human-dialog-based image retrieval?*
- *Sec. 2 describes details about hyperparameter selection and training.*
- *Finally, Sec. 3 presents qualitative results on image guessing task and generated language.*

*It contains 2 figures and 2 tables:*

- *Fig. 1 presents image-retrieval evaluation.*
- *Fig. 2 shows image guessing results for our model.*
- *Tab. 1 highlights example Q-BOT-A-BOT interactions for various models.*
- *Tab. 2 presents qualitative comparison of answers generated by various models.*

## 1. Dialog-based Image Retrieval

In Section 6 of the main paper, we evaluated how well models learn to cooperate at image guessing by setting up a retrieval task on the `test` split of VisDial v0.5. Specifically, we allow models to interact for 10 rounds where after each round of dialog exchange, Q-BOT predicts a feature representation  $\hat{y}_t$ , and we sort the entire `test` set in ascending order of distance to this prediction and compute the rank of the true image. Fig. 1a (reproduced from the main paper) shows mean percentile rank of the true image across rounds.

In this section, we conduct ablations to study how well different models perform when provided human dialog as input. Specifically, we study the performance of different models on the task of human-dialog-based image retrieval.

The task of human-dialog-based image retrieval allows us to probe two aspects of our models. First, we study how well RL trained models retain an understanding of human dialog from their SL-pretraining. Second, we can estimate a ‘human upper bound’ on the image guessing task. Specifically, when provided with a human-human dialog, A-BOT has ‘perfect’ perception and Q-BOT always asks consistent and non-repeating questions.

In addition to testing both models from the main paper (SL-pretrained and RL-full-QAf), we also train a supervised model, SL-pretrained-fc7only, specifically for this task of human-dialog-based image retrieval for comparison to those models. Specifically, a model trained to regress to image features given a human dialog. In practice, this amounts to finetuning  $f$  in SL-pretrained Q-BOT for image regression only, instead of a question-generation + fc7 regression multitask objective.

Fig. 1b shows results of this experiment as mean percentile rank of true image across rounds of human dialog. We note two observations. First, all models improve at the task of image retrieval when provided human-dialog, suggesting that language generation is a difficult task, and if we are provided ‘good’ dialog – self-consistent, accurate, and free from repetition – all models can improve at image guessing.

Second, we observe SL-pretrained-fc7only performs significantly better than the other models. This is unsurprising as this model is *explicitly trained* to regress from human dialog and lacks the additional utterance generation loss found in the other models. The strong performance of this model is encouraging because it implies that if we can generate better dialog, we can expect to do better at image retrieval.

\*The first two authors (AD, SK) contributed equally.

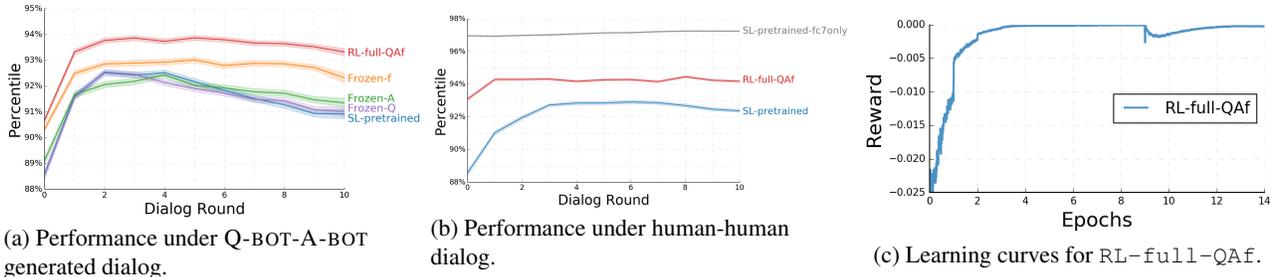


Figure 1: Guessing game evaluation – (a) based on Q-BOT-A-BOT dialog, and (b) based on ground-truth human dialog. (c) plots reward vs. epochs for RL-full-QAf.

## 2. Implementation Details

In this section, we provide additional details about the architecture, hyperparameters, and training procedure for the policy networks described in Section 4.1.

**Model Specifics.** Both Q-BOT and A-BOT policies are modeled via Hierarchical Recurrent Encoder-Decoder networks [2, 5, 6] and no parameters are shared between them.

All LSTMs (fact embedding, history encoder, question/answer decoder, and question encoder) have 2 layers and 512 dimensional hidden states. Word embeddings are 200 dimensional and are shared across question, history and decoder LSTMs. All models are implemented in Torch [1].

**Curriculum Training.** We use the penultimate layer activations from VGG-16 [7] as image representations and models are pretrained for 15 epochs on VisDial [2]. Afterwards they are transitioned to policy-gradient training by a fixed curriculum. Supervised training is performed for the first  $K$  rounds of a dialog and policy-gradient updates are performed for the final  $10 - K$ . We anneal  $K$  down by 1 every epoch (starting from 9). When  $K$  is zero, ground-truth captions are swapped for captions generated by an off-the-shelf captioning model [3], so that an infinite number of dialogs can be simulated without any human supervision. Models are trained for 5 epochs on these generated captions. Fig. 1c plots reward vs. epochs during policy learning. We use Adam [4] with a learning rate of  $10^{-3}$ , and clamp gradients to  $[-5, 5]$  to avoid explosion for all settings.

## 3. Qualitative Results

We provide additional qualitative results in Tab. 1, Tab. 2 and Fig. 2. Tab. 1 compares Q-BOT-A-BOT interactions for SL-pretrained and RL-full-QAf alongside the source images, their captions and human dialogs.

We notice that SL-pretrained agents often repeat questions –

*“what color are his shoes? black;  
what color are his shoes? black;  
what color ...”* (see first row).

RL-full-QAf agents do not fall into these loops as often; moreover, when they do repeat themselves, it is later in the dialog and they often sensibly recover from these loops –

*“can you see any trees? yes, there are lot of trees in background;  
can you see any buildings in background? no, I can not;  
can you see any trees? ...;  
does it look like they are in park?”* (row 1).

We also find that SL-pretrained agents tend to produce generic, ‘safe’ responses – *“I can’t tell; I don’t know; Can’t see”*, while RL-full-QAf interactions are noticeably more diverse and image-discriminative – *i.e.*,

*“what are they wearing?  
they are all wearing snow pants and jackets”* (row 1)  
*“where is man located?  
looks like classroom of some sort”* (row 9).

Notice that the human dialog does not mention location or classroom, despite these facts being visually discriminative.

In Sec. 6 of the main paper, we quantitatively evaluate how well A-BOT mimics human dialog (using the retrieval metrics from [2]). Recall that Frozen-Q-multi outperforms the other approaches on VisDial answering metrics. In Tab. 2, we present qualitative examples of responses generated by SL-pretrained, RL-full-QAf and Frozen-Q-multi A-BOT for ground-truth human-asked questions. We again find that SL-pretrained responses are predictably ‘safe’, while RL-full-QAf responses are more diverse and informative; however, this sometimes leads to RL-full-QAf being inaccurate. Frozen-Q-multi strikes a balance between response specificity and human-ness –

*“how old are men?  
SL-pretrained: I can’t tell,  
RL-full-QAf: they look like teenagers,  
Frozen-Q-multi: middle-aged”* (row 2);

*“what colors are umbrella?  
SL-pretrained: white,  
RL-full-QAf: rainbow colors,  
Frozen-Q-multi: many different colors”* (row 1).

Fig. 2 explores how Q-BOT's image guess changes through rounds of dialog. The left column shows the true image and caption and the middle column shows the text interaction between Q-BOT and A-BOT for select rounds. To show how Q-BOT's prediction changes through the dialog we show the five nearest neighbors of the prediction in the right column, with the nearest neighbor outlined in red. In the middle column we show where this nearest neighbor image ranks among other images in distance to the true feature, also showing other similarly distant images for context.

## References

[1] Torch. <http://torch.ch/>. 2  
[2] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura,

D. Parikh, and D. Batra. Visual Dialog. In *CVPR*, 2017. 2, 4, 5  
[3] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 2  
[4] D. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015. 2  
[5] I. V. Serban, A. Sordoni, Y. Bengio, A. Courville, and J. Pineau. Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models. In *AAAI*, 2016. 2  
[6] I. V. Serban, A. Sordoni, R. Lowe, L. Charlin, J. Pineau, A. Courville, and Y. Bengio. A Hierarchical Latent Variable Encoder-Decoder Model for Generating Dialogues. *arXiv preprint arXiv:1605.06069*, 2016. 2  
[7] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 2

Image + Caption	Human-Human dialog [2]	SL-pretrained Q-BOT-A-BOT dialog	RL-full-QAf Q-BOT-A-BOT dialog
	<p>What are their genders? 1 man 1 woman            Are they both adults? Yes            Do they wear goggles? Looks like sunglasses            Do they have hats on? Man does            Are there any other people? No            What color is man's hat? Black            Is it snowing now? No            What is woman wearing? Blue jacket and black pants            Are they smiling? Yes            Do you see trees? Yes</p>	<p>How many people? 2            Male or female? I can't tell            What are they wearing? I can't tell, they are far away            What color are skis? I can't tell            Are they wearing goggles? Yes            Are they wearing goggles? Yes</p>	<p>Are people male or female or both? I think 1 is male, can't see other 1 ...            What are they wearing? They are all wearing snow pants and jackets            Is it snowing? It does not appear to be snowing at moment            Can you see any trees? Yes, there are lot of trees in background            Can you see any buildings in background? No, I can not            Can you see any trees? Yes, there are lot of trees in background            Can you see any buildings in background? No, I do not see any buildings            Can you see any trees? Yes, there are lot of trees in background            Can you see any buildings in background? No, I do not see any buildings            Does it look like they are in park? It does not appear to be</p>
A couple of people in the snow on skis.			
	<p>How old are men? They look to be 30's            Are they in restaurant? I can't tell, its close up            Do men look they might be related? Maybe            Are they each holding slice of pizza in 1 hand and giving thumbs up with other? 1 holding slice, other holding box...            What kind of pizza do they have? cheese            Does box have pizza company name on it? No            Are there any pictures or graphics on box? No            Are men slender or heavy set? Slender            Do either of men wear glasses? Yes            1 or both? 1</p>	<p>How old are men? They look like they are in their 20's            Are they male or female? 1 male and 1 female            What are they wearing? T shirt and shorts            What kind of pizza is it? Looks like pepperoni            Does it look good? Yes it does            Can you tell what kind of pizza it is? Looks like pepperoni            Are there any toppings on pizza? Pepperoni            Are there any toppings on pizza? Pepperoni            Are there any toppings on pizza? Looks like cheese, but I can't tell...            Are there any other toppings on pizza? No</p>	<p>How old are men? They look like they are in their 20's            Are they wearing chef's hat? No            Do they have long hair? No            What kind of pizza is it? Looks like pepperoni            Does it look good? Yes            Are there any other people around? No            Are they inside or outside? Inside            Can you tell what time of day it is? No            Are there any other people around? No            What color are walls? White</p>
Two men giving thumbs up holding large slices of pizza.			
	<p>How old is man? In his teens            Is he wearing helmet? No            Is he wearing any protective gear? No            Do either of his hair? Dirty blonde or dark brown            What color is his shirt? Black with red letters            Is he wearing shorts? Yes            What color are his shorts? Denim            Are there other people visible? No            What color are his shoes? I can't see them            Is this at skate park? Yes</p>	<p>How old is man? I can't see his face, I don't know            What is he wearing? shorts and t shirt            What color is his shirt? white            what color are his pants? black            what color is his shirt? blue            What color are his shoes? Black            What color are his shoes? Black</p>	<p>What color is skateboard? It is hard to tell but I think it's black            Is he wearing helmet? Yes            Is he wearing knee pads? No, no knee pads            Is it day or night? It is daytime            Is it sunny? It is bright out, but I see no shadows, so I can not            Can you see any other people? There are 2 other people in background            Is it day or night? It is daytime            Is it sunny? Yes, it is sunny            What color is skateboard? It is hard to tell but I think it's black            Is he wearing helmet? Yes, he is wearing helmet</p>
A man riding a skateboard down the side of a ramp.			
	<p>What color is plane? White and blue            Is there writing on plane? Yes it says lufthansa            What color is ground? Grey tarmac            Any people? No people visible            Can you see airport? Yes            What color is his hair? Grey            Is there sunshine? No it's kind of overcast            Any other planes? No other planes            Any grass? No grass            Any trees? No tree</p>	<p>What color is plane? White            Is it large or small? Large            Any people? No            Any writing on plane? Yes            What does it say? UNK            What color is plane? White            Any trees? Yes</p>	<p>What color is airplane? White with red and blue stripes            Can you tell what airline it is? No, it's too far away            What time of day is it? It is daytime            Can you see any people? No            Can you see any buildings? Yes, in background            Is it day or night? It is daytime            Can you see sky? Yes            Are there clouds in sky? Yes, few            Does it look like it's going to rain? Yes it does            Does it look like it's going to rain? Yes it does</p>
An airplane sitting on an asphalt landing strip.			
	<p>What color is dog? Brown and white            What breed is dog? Look like britany spaniel            Is dog wearing collar? Not able to see            Can you tell what is cooking? Look like beans            Is oven stainless steel? Yes            Does it have see through glass? Yes            How many racks in oven? 2            Can you see knobs? No            Does oven have wooden handle? Not able to see            Are beans in pot? Glass dish</p>	<p>What kind of dog is it? I can't tell            What color is it? White            Is it large or small? Large            What kind of food is it? I can't tell            Is it large or small? Small            What color is it? White            Is it large or small? Large            What color are walls? White            What kind of floor? Can't see floor            What kind of floor? Can't see floor</p>	<p>What kind of dog is it? Terrier of some kind            What color is it? White and grey            Can you tell what kind of food it is? No            Are there people around? No people            Is it inside or outside? Inside            Can you tell what time of day it is? No I can't            What color is room? I can't see walls            Do you see any windows? No            Do you see any windows? No            Do you see any windows? No</p>
A dog watching food cook in an oven.			
	<p>Are there more people in this photo or just man? No            Just tennis player man in this photo? Yes            Is it outdoor court? Can't tell, it's close up            What is man wearing? White shirt and shorts            Can you see net? No            Can you see bottom of court? Yes            Can you see ball in photo? Yes            Is it close to man's racquet? Yes            Is man about to take swing? Yes            Is he left handed or right handed? Right</p>	<p>How old is man? I can't see his face, I don't know            What color is his hair? Brown            What is he wearing? Shorts and t shirt            What color is his shirt? White            What color is his shirt? White</p>	<p>How old is man? He looks to be in his 20's            What color is his hair? His hair is brown            Is he wearing glasses? No he is not            Is he wearing hat? No he is not            Is he wearing glasses? No he is not            Are there other people around? No it's just him            Can you see net? No I can't            What color is court? Court is green            Is it sunny? Yes it is            Do you see any trees? No I do not</p>
A man prepares to swing at the tennis ball.			
	<p>Is giraffe adult? I cannot tell, I only see face            Is grass at higher altitude than its feet? No            Is rock big? Yes            How is weather? Sunny            What color is rock? Tan color            Is grass green or dry? Green but looks little dry            Does giraffe look healthy? Yes            Does giraffe look happy? It looks tired            Is giraffe laying down? Yes            Is giraffe's mouth closed? Yes</p>	<p>Is there only 1 giraffe? Yes            Are there any other animals? No            Are there any people? No            Are there any other animals? No            Is it day or night? Day            Is it sunny? Yes            What time of day it is? I can't tell            Are there any other animals? No            Are there any trees? No            Are there any other animals? No</p>	<p>Is there only 1 giraffe? Yes            Are there any other animals? No            Are there any people? No            Does this look like zoo? Can't tell            Is there fence? No            Is it day or night? Day            Is it sunny? Yes            What time of day is it? Morning            Are there any trees? No            Can you see sky? No</p>
A giraffe resting its head in the grass near a rock.			
	<p>Are there any people in picture? No, there are no people            Is cheese white? Yes, cheese is white            Are crackers white? No, cracker is brownish            Are there lots of greens? No, only 1 sprig            Was picture taken in kitchen? I cannot tell, it is too zoomed in            How many crackers are there? There is 1            Is there design on plate? No, plate is plain white            Is there good lighting? Yes, lighting is decent            Is plate on table? I cannot tell I can only see plate as background            Does cracker look delicious? Absolutely</p>	<p>Are there any people? No            What color is plate? White            Is plate on table? I can't tell            What color is plate? White            Is plate on table? I can't tell            Is there anything else on table? No            Is there anything else on table? No            Is there anything else on table? No            Is there anything else on table? No</p>	<p>Is plate large or small? Large            Is it on table? I can't tell            What kind of donut? Plain with sugar on it            What kind of chips? Potato chips            Any people? No            Is this inside or outside? I can't tell, it's close up            What color is table? White            Any chairs? No            What color is table? I can't tell            Is it large or small? I can't see plate</p>
A cracker with shredded cheese and leafy greens on a white plate.			
	<p>How old is this man? He looks like late teens            Does he have short hair or long hair? Medium            What color is his hair? Brown            Is he wearing glasses? No            Does he have facial hair? No            How is he dressed? He's white dress shirt with red tie and matching red vest over it looks like school uniform            Is he larger man or is he skinny? He's skinny            Is this picture taken outside? Inside            Can you tell if it's daytime? I'm not sure            Can you tell what room he is standing in? He's in hallway</p>	<p>How old is man? He looks like he's in his 20s            T shirt and jeans            What color is his hair? Black            What is he wearing? T shirt and jeans            What color is his hair? Black            What is he wearing? T shirt and jeans            What color is his hair? Black            Is it short or long? Short            What color is his hair? Black            Is it short or long? Short</p>	<p>Where is man located? Looks like classroom of some sort            What is man wearing? Black t shirt and jeans            Are there any other people in picture? Yes, quite few people in background            What color is table? Picture is black and white, but it's wood table            Is there anything else on table? Not that I can see            What are other people doing? They are standing in front of him, but I don't know what they            Are there any windows? Not that I can see            What color are walls? I can't see walls            What time of day do you think it is? I can't tell            What is man wearing? Black t shirt and jeans</p>
A man making the live long and prosper sign from star trek.			

Table 1: Selected examples of Q-BOT-A-BOT interactions for SL-pretrained and RL-full-QAf. RL-full-QAf interactions are diverse, less prone to repetitive and safe exchanges (“can’t tell”, “don’t know” etc.), and more image-discriminative.

Image + Caption	Human question [2]	Human answer [2]	SL-pretrained	RL-full-QAf	Frozen-Q-multi
	<b>Round 1:</b> What colors are umbrella? <b>Round 2:</b> What kind of tower is it? <b>Round 4:</b> Can you see any clouds? <b>Round 9:</b> Is there any writing on tower?	It is blue, red, yellow and green I can't really tell Very few There isn't any writing it's just gray	White Not sure I can not see sky No	Rainbow colors It's small building Sky is not visible Not that I can see	Many different colors Small I No No, there is no writing
Multi-color umbrella floating in the air next to grey tower.					
	<b>Round 2:</b> How many people? <b>Round 5:</b> What are they wearing? <b>Round 8:</b> How old are men?	I see 3 Men are in dress pants and shirt and woman in jeans tank 30	Too many to count I can't tell I can't tell	2 Jackets They look like teenagers	3 Uniforms Middle aged
People sitting outside along a concrete wall on a sunny day.					
	<b>Round 2:</b> What color is cat? <b>Round 4:</b> Are there any appliances? <b>Round 9:</b> What else is there?	He is mostly cream colored with some tan markings Just microwave that I can see Knife block and some recipe binders	Black Stove Nothing else	White Stove, stove Toilet, toilet, toilet and top of toilet paper	Black and white No Wall
A kitchen scene with a cat sitting on the counter in front of a microwave.					
	<b>Round 1:</b> How old is person? <b>Round 6:</b> Where do you think he is? <b>Round 10:</b> What else can you tell me about this image?	I think around 17 It looks like neighborhood road He's very skinny boy	Can't tell I don't know There is snow on ground	Teenager At park There are lot of people in background	20s In park I see trees and sky
A person on a skateboard performing a trick.					
	<b>Round 5:</b> Are there any people? <b>Round 8:</b> Does lamb appear to be healthy? <b>Round 9:</b> What season does it seem to be?	There aren't any people there Lamb looks to be very healthy It looks to be spring	There are no people Yes, it does Summer	No people Yes it does It looks like summer	No Yes Spring
A small brown lamb walking across a lush green field.					

Table 2: Selected examples of answers generated by SL-pretrained, RL-full-QAf and Frozen-Q-multi A-BOT for ground-truth human-asked questions. Although differences are subtle, SL-pretrained answers are usually 'safe' and repetitive, RL-full-QAf answers diverse and sometimes unfaithful, and Frozen-Q-multi answers hit a balance in between.

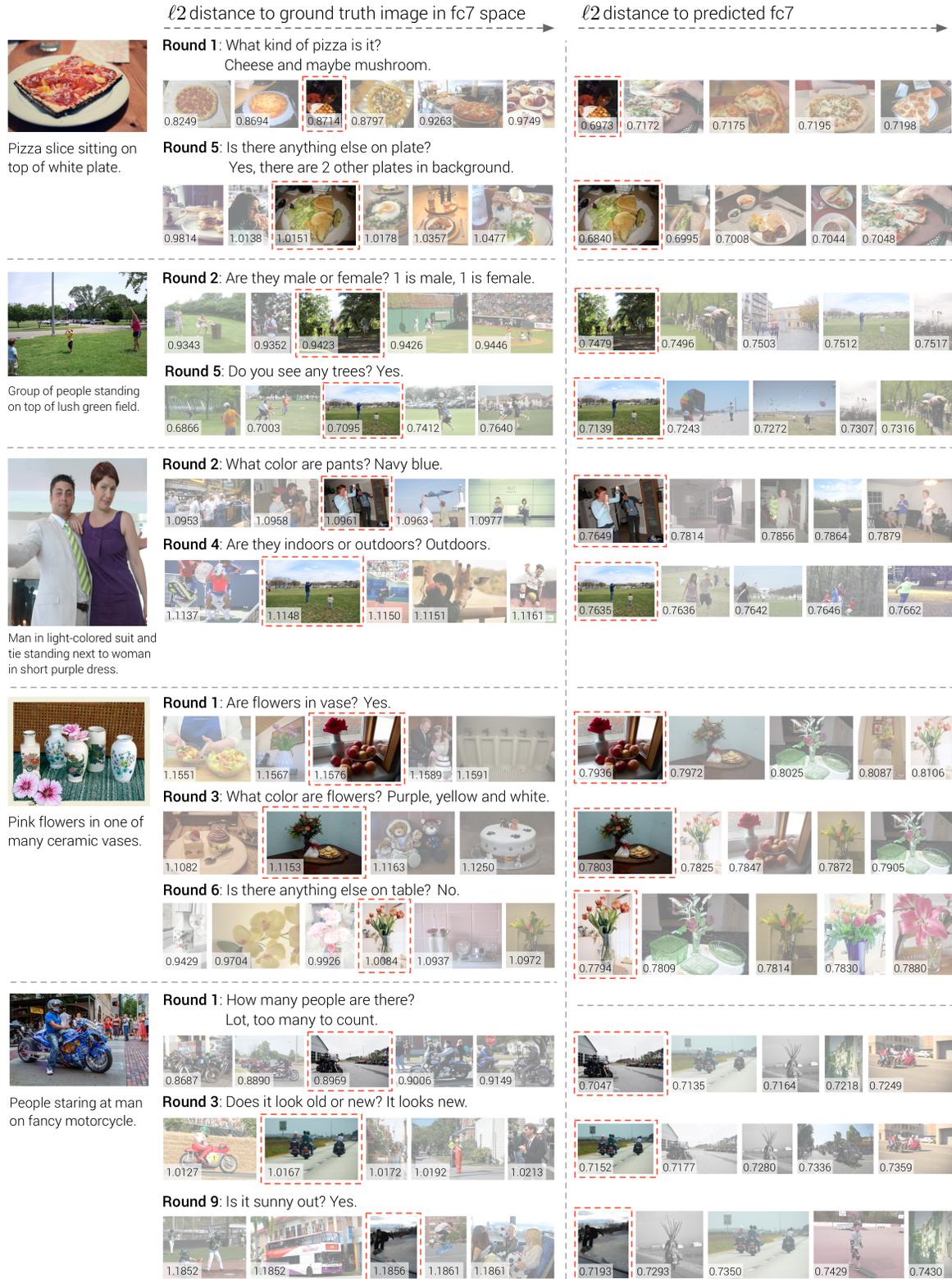


Figure 2: Qualitative results on predicted fc7-based image retrieval. Left column shows true image and caption, middle column shows dialog exchange, and a list of images sorted by their distance to the ground-truth image, right column shows a list of images sorted by their distance to predicted fc7. The image predicted by Q-BOT is highlighted in red. We can see that the predicted image is often semantically quite similar.