# A    Reinforcement Learning Algorithm Details

Based on the supervisedly pretrained Q-Bot and A-Bot, RL is used to fine-tune the two models. Following previous works, we use REINFORCE algorithm to update policy parameters.

Q-Bot's policy is formulated as $\pi_Q(q_t|S_t;\theta_Q)$, where $S_t$ is the state, $q_t$ is the generated question and the policy is learned by the deep neural network parameterized by $\theta_Q$. Similarly, A-Bot's policy is formulated as $\pi_A(a_t|S_t;\theta_A)$. For policy $\pi_Q$ and $\pi_A$, the objective function of the policy gradient is given by: $J(\theta_Q,\theta_A) = E_{\pi_Q,\pi_A}\sum_{t=1}^{T} r_t(S_t,A_t)$. Following the REIN-FORCE, the gradient of $J(\theta_Q)$ can be written as: $\nabla J(\theta_Q) = E_{\pi_Q,\pi_A} r_t(\cdot)\nabla\theta_Q log\pi_Q(q_t|S_t)$. Similarly, the gradient of $J(\theta_A)$ can be written as: $\nabla J(\theta_A) = E_{\pi_Q,\pi_A} r_t(\cdot)\nabla\theta_A log\pi_A(a_t|S_t)$. Accordingly, Q-Bot and A-Bot are optimized with policy gradients while maximizing the reward $r_t$ as claimed in Eq. 8. In Algorithm 1, we exemplify the training procedure that uses RL at the beginning of the self-talking dialogue. To simplify, the target image feature is denoted by *tar_feat* and the predicted image feature is denoted by *pred_feat*.

# B    Training Details

**Supervise Learning.** We use the same Q-Bot as in previous work [8]. In concrete, the model consists of 4 components: 1) fact encoder, 2) history encoder, 3) question decoder, 4) feature regression network. Following previous settings, all encoders and decoder are 2-layer LSTM with 512-d hidden states, the feature regression network is a 1-layer MLP. Fact encoder and history encoder encodes the textual dialogue history. Based on the encoding of dialogue history, question decoder (a 2-layer LSTM) generates the question and feature regression Network (a single fully connected layer) outputs the predicted image feature. Q-Bot is optimized with cross-entropy loss and L2 loss that minimizes the distance between predicted image feature and target image feature. Following previous works, we pretrain Q-Bot for 20 epochs with a learning rate of 1e-3 that is iteratively decayed to a minimum of 5e-5. Adam is the optimizer and dropout rate is 0.5.

Among many attention-based A-Bot models, we use a classic attention-based A-Bot with History-Conditioned Image Attentive Encoder (HCIAE) [16] to verify the proposed method. According to the visual attention computation method in HCIAE, the current concerned visual contents **i-vc$_t$** in Eq. 2 can further be described as:

$$\mathbf{att_t} = softmax(\mathbf{w_a}^T tanh(\mathbf{W_H M_t^H} + (\mathbf{W_q m_t^q})\mathbf{1^T}), \qquad (9)$$

where $\mathbf{M_t^H}$ is history feature, $\mathbf{m_t^q}$ is question feature, $\mathbf{w_a}$, $\mathbf{W_H}$ and $\mathbf{W_q}$ are learnable parameters. A-Bot is optimized with cross entropy loss. Following their implementation, we pretrain A-Bot for 40 epochs with a learning rate of 4e-4, decayed by 0.75 after 10 epochs. We use Adam as the optimizer. Dropout rate is 0.5.

**Reinforcement Learning.** The $\gamma$ in the polarization operation in Eq. 4 is set to be 0.6 empirically. The reward coefficients in Eq. 8 are set as follows: $\alpha_O$=1, $\alpha_D$=1e-1, $\alpha_I$=1e-2. The coefficients vary because different rewards are at different magnitudes. We set the coefficients based on the consideration of balancing task success and dialogue quality. We give further analysis on coefficients setting in Appendix E.

RL smoothly begins according to a curriculum learning setting: SL is used for the first K rounds of dialogue and policy-gradient update works in the left 10 - K rounds. Following [24], we start at K = 9 and gradually anneal it to 4. We train the agents for 30 epochs with a batch size of 32. As learning rate is inconsistent in SL, we also use different settings in

---

**Algorithm 1** Training the Q-Bot and A-bot using REINFORCE with ECS-based rewards.

---

1: **for** Each Update **do**
2:     **# Initialize A-VC:**
3:     $a\text{-}vc_0 \leftarrow 0$
4:     $i\text{-}vc_0 \sim U(0,1)$
5:     $\triangle pred\_feat \leftarrow 0$
6:     **for** $t = 1$ to $T$ **do**
7:         $q_t, pred\_feat \leftarrow QBot(H_t)$
8:         **# Obtain I-VC:**
9:         $a_t, att_t \leftarrow ABot(I, H_t, q_t)$
10:        $i\text{-}vc_t = att_t$
11:        **# Update A-VC:**
12:        $a\text{-}vc_t \leftarrow OR(a\text{-}vc_{t-1}, polarize(i\text{-}vc_t))$
13:        **# Compute rewards:**
14:        $r_t^D \leftarrow D_{KL}(i\text{-}vc_t || i\text{-}vc_{t-1})$
15:        **if** $\sum_{j=1}^{N} a\text{-}vc_t^j - a\text{-}vc_{t-1}^j > 0$ **then**
16:            $r_t^I \leftarrow 1$
17:        **else** $r_t^I \leftarrow 0$
18:        $r_t^O \leftarrow (tar\_feat, pred\_feat)^2 - \triangle pred\_feat$
19:        $r_t \leftarrow \alpha_O r_t^O + \alpha_D r_t^D + \alpha_I r_t^I$
20:        $\triangle pred\_feat \leftarrow (tar\_feat, pred\_feat)^2$
21:        Evaluate $\nabla J(\theta_Q)$ and update QBot
22:        Evaluate $\nabla J(\theta_A)$ and update ABot
23:    **end for**
24: **end for**

---

RL. We optimize the Q-Bot with a learning rate of 1e-3 that is iteratively decayed up to a minimum of 5e-5 while learning rate of A-Bot is initialized with 1e-4 and decayed up to 5e-6. We clamp gradients to [-5, 5] to avoid explosion.

# C    Additional Qualitative Analysis on A-Bot

We find our A-Bot tends to generate descriptive responses that not only directly answer the question but also describe correlated visual information. We give the selected generated dialogue examples in Fig. 8(a). As can be seen, when is asked questions like *"what else?"*, our A-Bot gives detailed description about visual contents that have not been mentioned yet while the A-Bot optimized by single original reward tediously says *"no"*. As in the top example, our A-Bot constantly adds new information of the undisclosed image by answering *"yes , a glass of water , a coffee cup , a cup of some sort , and another plate"* to *"do you see any other items ?"* and answering *"a cell phone , part of what looks like a UNK , the rest of the paper , and a…"* to *"anything else seen besides the interest ?"*. The descriptive responses help Q-Bot have better knowledge of the image and encourage Q-Bot to generate new questions based on the given detailed answers – as in the bottom example in Fig. 8(a), after knowing that there are tennis court in the background, Q-Bot then asks *"what color is the surface ?"*.

In Fig. 8(b), we give the examples that replace the jointly-trained Q-Bot with the ground-

truth questions, to illustrate the adaptability of our A-Bot and make a direct comparison. As can be seen, our A-Bot still performs outstandingly as it tends to give detailed answers that include additional explanations or descriptions moreover the simple yes/no. In the first example, when the question is *"is this in a zoo?"*, our A-Bot answers *"no, it's in the wild"*, which corrects Q-Bot's belief on the undisclosed image, while the basic A-Bot only answers *"i don't think so"*. In the third example, when the question is *"can you see the leaves of the tree?"*, our A-Bot answers *"yes, they are green"* which gives additional information about the related visual contents. Informative and diverse answers like these are more human-like and are helpful to keep users engaged.

# D  Additional Qualitative Analysis on Comparing Methods

We give additional generated dialogue examples from comparing methods in Fig. 9. As what has been analyzed before, our method helps achieve less repetitive and more visual informative dialogues. As can be seen, dialogues generated by *RL* start to repeat at an early stage of the dialogue, dialogues from *RL-Diverse* are less repetitive but tend to repeat at successive rounds, while dialogues generated by *RL-Ours* can always converse on new visual contents. Besides, our dialogues valuably show consistency sometimes. For example, in the 2nd sample, after the A-Bot answers *"shirt and pants"*, the Q-Bot further asks *"What color is his shirt?"*. Also, in the 3rd sample, as the A-Bot mentions *"there is a candle"*, the Q-Bot then asks *"What color is the candle?"*

# E  Analysis on Different Settings of Reward-Coefficients

We explore the training effects brought by different settings of reward-coefficients, i.e. the coefficients of DR and IR. We show the results in Tab. 2, which includes indicators that measure the performances in three-aspects: 1) joint performance of Q-Bot and A-Bot by PMR, i.e. the Percentile Mean Rank in the image-guessing task; 2) Q-Bot performance by Unique Question; 3) A-Bot performance by MRR and NDCG.

In setting 1, we intuitively assign the reward coefficients to let weighted rewards be in the same order of magnitude. We also conduct experiments with other reward-coefficients settings, as in 2 and 3. We think different coefficient settings lead to different training results on Q-Bot and A-Bot. Setting 1 achieves the most balanced cooperative Q-Bot and A-Bot while other settings may help improve the performances on separate tasks, i.e., image-guessing and Visual Dialogue answering, which are also meaningful.

| settings | $\alpha_D$ | $\alpha_I$ | Unique question$^Q$ | MRR$^A$ | NDCG$^A$ | PMR$^{QA}$ |
|---|---|---|---|---|---|---|
| 1 | 0.1 | 0.01 | **8.11**$_{\pm 0.08}$ | 49.36$_{\pm 0.00}$ | 57.81$_{\pm 0.00}$ | 96.86$_{\pm 0.06}$ |
| 2 | 0.05 | 0.01 | 6.79$_{\pm 0.07}$ | 49.38$_{\pm 0.00}$ | **58.00**$_{\pm 0.00}$ | 96.57$_{\pm 0.08}$ |
| 3 | 0.01 | 0.001 | 5.43$_{\pm 0.00}$ | **49.42**$_{\pm 0.00}$ | 57.84$_{\pm 0.00}$ | **97.00**$_{\pm 0.00}$ |

Table 2: Performances on VisDial v1.0 val under different reward-coefficients settings. $\alpha_D$ and $\alpha_I$ are coefficients for DR and IR, respectively. The superscripts represent the indicator belongs to which agent or both.

| Image | RL | RL - Ours |
|---|---|---|
| | is this in color ?<br>A:yes<br>do you see any people ?<br>A:just a hand<br>is this inside or outside ?<br>A:i ca n't tell it 's a close up<br>what color is the plate ?<br>A:white<br>anything else in the room ?<br>A:no , just the hot dog<br>anything else in the room ?<br>A:no , just the hot dog<br>anything else in the room ?<br>A:no , just the hot dog<br>anything else in the room ?<br>A:no , just the food<br>anything else in the room ?<br>A:no , just the hot dog<br>anything else in the room ?<br>A:no , just the food | is the photo in color ?<br>A:yes<br>is it a professional photo ?<br>A:no<br>is it well lit ?<br>A:yes<br>is it daytime ?<br>A:yes<br>is it sunny ?<br>A:yes<br>do you see any people ?<br>A:yes<br>do you see any other items ?<br>A:yes , a glass of water , a coffee cup , a cup of some sort , and another plate<br>what color is the shirt ?<br>A:i ca n't see the top of the shirt<br>anything else seen besides the interest ?<br>A:a cell phone , part of what looks like a UNK , the rest of the paper , and a<br>what color is the table ?<br>A:i ca n't see the table |
| | is this in color ?<br>A:yes<br>how many people do you see ?<br>A:2<br>are they male or female ?<br>A:male<br>how old do they look ?<br>A:30 's<br>anything else in the room ?<br>A:no<br>anything else in the room ?<br>A:no<br>anything else in the room ?<br>A:no<br>anything else in the room ?<br>A:no<br>anything else in the room ?<br>A:no<br>anything else in the room ?<br>A:no | is the photo in color ?<br>A:yes<br>is it a professional photo ?<br>A:no<br>is it well lit ?<br>A:yes<br>is it daytime ?<br>A:yes<br>how old are the men ?<br>A:20 's<br>what race are the men ?<br>A:white<br>what else is in the background ?<br>A:tennis rackets , a tennis court , part of a tennis ball , part of a blue sky<br>what color is the man 's shirt ?<br>A:orange<br>anything else in the background ?<br>A:no , just a baseball glove , part of a tennis ball , part of a tennis player , the<br>what color is the surface ?<br>A:green |

(a)

| | Question | RL A-Bot | RL-Ours A-Bot |
|---|---|---|---|
| | is this photo in color ?<br>is this outdoors ?<br>is it sunny ?<br>can you see any trees ?<br>can you see grass ?<br>how many zebras are there ?<br>what color are they ?<br>can you see any dirt ?<br>is this in a zoo ?<br>is the field grassy ? | yes<br>yes , it is<br>*yes , it is*<br>no<br>yes , there is a lot of grass<br>there are 2<br>black and white<br>*yes , there is a lot of grass*<br>*i do n't think so*<br>yes | yes<br>yes , it is<br>*no , it 's cloudy*<br>no<br>yes , i can see grass<br>2<br>black and white<br>*yes , in the background*<br>*no , it 's in the wild*<br>yes |
| | is there any people in the picture ?<br>is there anything in the picture ?<br>is it dark ?<br>is it a big room ?<br>what color is the carpet ?<br>are the lights big ?<br>what color of the towel ?<br>is the vanity large ?<br>anything in the sink ?<br>anything else stand out ? | no<br>yes<br>no<br>*no*<br>*brown*<br>yes<br>white<br>yes<br>no<br>*no* | no<br>yes<br>no<br>*no, it 's small*<br>*multi colored*<br>yes<br>white<br>yes<br>no<br>*yes , there is a cabinet* |
| | is this photo in color ?<br>does the photo look professional ?<br>is it daytime ?<br>is it sunny ?<br>are there people ?<br>are there animals ?<br>can you see any signs ?<br>can you see the leaves of the tree ?<br>are there any cars ?<br>are there any buildings ? | yes<br>no<br>yes<br>yes<br>no<br>no<br>*yes*<br>*yes*<br>no<br>no | yes<br>no<br>yes<br>yes<br>no<br>no animals<br>*yes , but i ca n't read them*<br>*yes , they are green*<br>no cars<br>no |

(b)

Figure 8: Illustration of A-Bot's ability to generate descriptive responses when (a) answering jointly optimized Q-Bot, (b) answering ground truth questions. Comparing models are: 1)'RL': basic agents fine-tuned by RL using the original rewards; 2) RL-Ours: basic agents fine-tuned by RL with two ECS-based rewards in addition.

| Image + Caption | RL | RL - Diverse | RL - Ours |
|---|---|---|---|
| Caption: 2 dogs with a hat standing beside each other | what color are dogs ?1 is white and other is brown<br>are they big dogs ? medium sized<br>what color are they ?1 is white and other is black<br>anything else in room ?there is fence in background<br>anything else in room ?nothing else<br>anything else in room ?nothing else<br>anything else in room ? nothing else<br>anything else in room ?nothing else<br>anything else in room ?nothing else<br>anything else in room ?nothing else | what kind of dogs are they ?they are labs<br>what color are dogs ?1 is white and other is brown<br>what kind of dogs ?they look like labs<br>what color are they ?they are white<br>are they inside or outside ? outside<br>are there any trees around ?yes<br>are there any people in picture ?no<br>are there any trees ?yes<br>are there any buildings ?no<br>are there trees around ?there are trees in background | is the photo in color ?yes<br>is it a professional photo ?no<br>is it well lit ?yes<br>is it daytime ?yes<br>what color are the dogs ?1 is white and 1 is brown<br>are the dogs big or small ?big<br>do you see anything else besides the dogs ?no<br>do you see anything else besides interest ?no<br>what color is the floor ?brown<br>anything else in the background ?no |
| Caption: a man stands at a dining room table holding a knife to cut a cake | what color is his hair ?dark brown<br>is he wearing glasses ?no<br>what color is his hair ?dark brown<br>is he wearing glasses ?no<br>anything else in the room ?yes , a lot of glasses , coffee table<br>anything else in the room ?there is a small table in the background , and part of the table , part of the table ,<br>anything else in the room ?there is a small table in the corner of the table , part of the table , part of the<br>anything else in the room ?there is a small table behind the man , and the man 's hand is on the wall behind<br>anything else in the room ?there is a small table in the background , and the corner of the man , the man 's hand<br>anything else in the room ?there is a glass , a table , a table | how old is the man ?late 20 's early 30 's<br>what is he wearing ?a button up shirt and a tie<br>what color is his hair ?brown<br>is he wearing glasses ?no<br>what color is his hair ?brown<br>what color is his shirt ?white<br>what color is the table ?brown<br>can you see furniture around ?yes , a table , a table<br>can you see furniture around ?yes , there is a table , a table<br>can you see furniture around ?yes , there is a table , a table | is the photo in color ? yes<br>is the photo close up ? no<br>how old is the man ?20 's<br>what race is the man ? white<br>how old does he look ? 30<br>what is he wearing ?shirt and pants<br>is there anything else other than him ?yes , a bookshelf , shelves<br>what color is his shirt ? white<br>anything else in the background ?tables , tables , tables<br>what color is the floor ?i ca n't see the floor |
| Caption: an elderly woman posing with 2 girls and a boy in front of a birthday cake in the kitchen | how old does the woman look ?in her 30s<br>what is she wearing ?a white shirt<br>what kind of cake is it ?white cake<br>how many candles are on the cake ?0<br>anything else in the room ?yes , a table , cups<br>what color are the cups ?clear glass<br>what color are the cups ?clear glass<br>anything else in the room ?there are a few decorations in the background , but i ca n't tell what they are<br>anything else in the room ?there are a few decorations in the background , but i ca n't tell what they are<br>anything else in the room ?there is a table behind the cake , and the table is in the background , but i can | is the photo in color ?yes<br>how old is the girl ?maybe 5<br>what color is her hair ?brown<br>what is she wearing ?a white shirt<br>what color is her hair ?brown<br>what is she wearing ?a white shirt<br>what color is her hair ?brown<br>what color are the walls ?white<br>can you see furniture around ?yes , part of a table , a table<br>can you see furniture around ?yes , the kitchen table and the kitchen is in the kitchen , the kitchen is in the kitchen | how old is the woman ?i would say late 20 's early 30's<br>how old is the woman ? she looks like she 's in her 20 's<br>what color is her hair ? her hair is brown<br>what color is the cake ? the cake is white<br>what color is the girls hair ?her hair is brown<br>is there anything else on the table ?yes , there is a candle and a candle<br>what color is the candle ?the candles is yellow<br>is there anything else in the background ?yes , there are a lot of stuff in the background ,cabinets<br>what color is the floor ? the floor is not visible<br>anything else in the room ?there are a lot of things on the table , and the cabinets are the kitchen cabinets , the other |
| Caption: a large america airlines plane flying through the air | what color is the plane ?red , white , and blue<br>is there any writing on it ?yes<br>what does it say ?UNK<br>what color is the plane ?red , white , and blue<br>is there any writing on the plane ?yes<br>what does it say ?UNK<br>anything else in the background ?no<br>anything else in the room ?no<br>anything else in the room ?no<br>anything else in the room ?no | is the photo in color ?yes<br>is there any people ?no<br>is this indoors or outdoors ?outdoors<br>can you see the sky ?yes<br>can you see the clouds ?no<br>can you see trees ?no<br>can you see buildings ?no<br>can you see trees around ?no<br>can you see buildings around ?no<br>can you see trees around ?no | what color is the plane ? red and white<br>is it a passenger plane ? yes<br>can you see the windows on the plane ? yes<br>can you see the windows on the plane ? yes<br>can you see any people ?no<br>can you see the sky ?yes<br>can you see anything<br>else besides interest ?no<br>what color is the plane ? white and red<br>can you see anything else besides the plane ? no<br>can you see the ground ?no |

Figure 9: Additional generated dialogue Examples from comparing methods, i.e., 1)'RL': basic agents fine-tuned by RL using the original rewards; 2)'RL-Diverse': diverse agents fine-tuned by RL using the original rewards; 3) RL-Ours: basic agents fine-tuned by RL with ECS-based rewards in addition.