Carnegie Mellon University
School of Computer Science

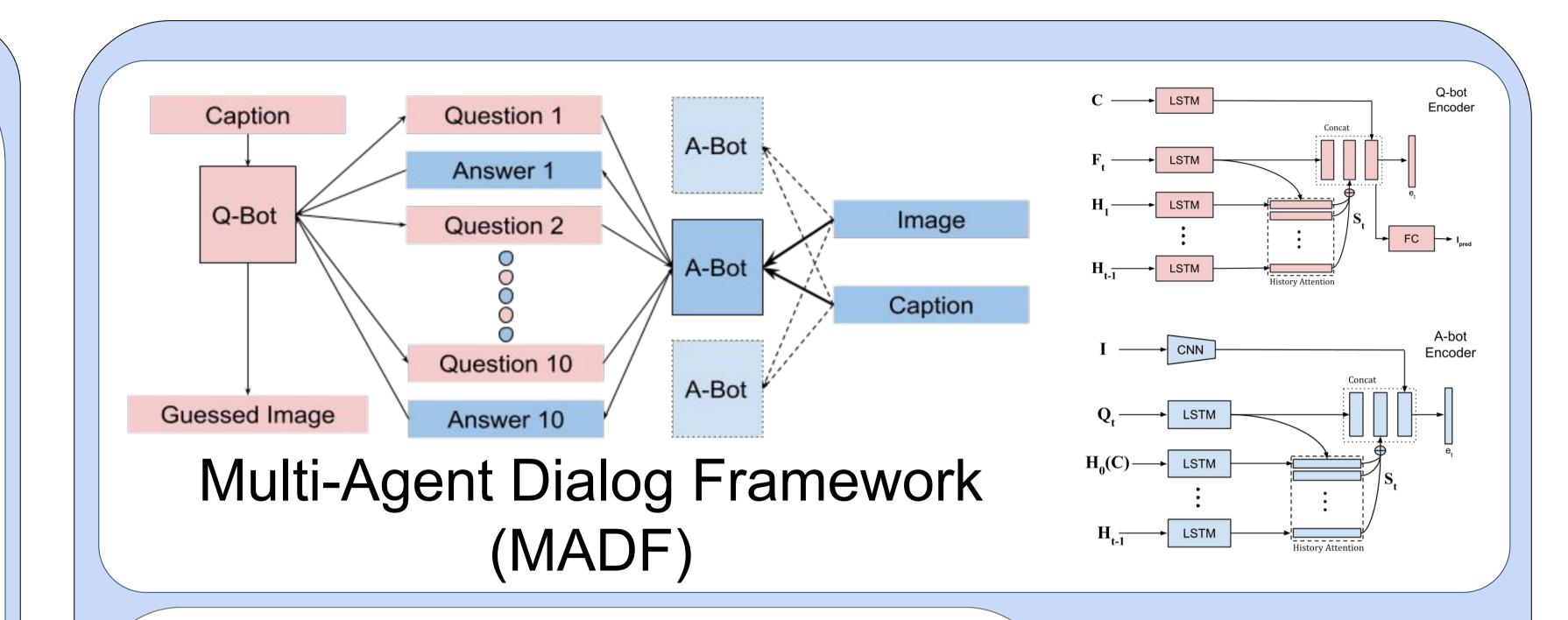Akshat Agarwal*, Swaminathan G.*, Vasu Sharma*, Katia Sycara

Github Code
https://goo.gl/gc6dGZ

Paper
https://goo.gl/pYWUV6

# Mind Your Language
# Learning Visually Grounded Dialog in a Multi-Agent Setting

## MOTIVATION



Goal-oriented dialog involves agents conversing with each other to achieve a particular goal, like transferring information

- Interpretability of these conversations desirable for transparency, motivating the use of Natural Language (NL)
- For AI, exchanging info by communicating in NL is inherently suboptimal
- Humans adhere to NL because they have to interact with an entire community, and having a private language for each person would be inefficient
- Hence, we propose a multi-agent dialog framework where each agent interacts with and learns from multiple agents, resulting in coherent and interpretable dialog

## VISUAL DIALOG TASK

- Formulated as a conversation between two collaborative agents, a Question (Q-) Bot and an Answer (A-) Bot
- A-Bot given an image and a caption, while Q-Bot is given only a caption - both agents share a common objective, which is for Q-Bot to form an accurate mental representation of the unseen image
- Facilitated by exchange of 10 pairs of questions and answers between the two agents, using a shared common vocabulary

## BACKGROUND



First, agents are trained in isolation via supervision for 15 epochs from VisDial dataset, using Max Likelihood Estimation loss wrt ground truth QA - results in repetitive responses

Then, they are smoothly transitioned to RL via a curriculum

1. Both agents interact with each other and learn by self-play
2. Q-Bot observes $\{c, q_1, a_1, ..., q_{10}, a_{10}\}$, A-Bot also observes $I$
3. Action: Predict words sequentially until a stop token is encountered (or max length reached)
4. Reward: Incentivizing information gain from each round of QA, measured using the predicted image embedding $y_t$

$$r_t(s_t^Q, (q_t, a_t, y_t)) = l(y_{t-1}, y^{gt}) - l(y_t, y^{gt})$$

5. No motivation to stick to rules and conventions of English language, making the RL optimization problem ill-posed



### Multi-Agent Dialog Framework (MADF)

- We solve the problem using our multi-agent setup where 1 Q-Bot communicates with 1 of multiple A-Bots (or vice-versa) for a batch of training, then chooses another A-Bot and repeats
- Much harder for the agents to deviate from natural language since coming up with a new language for each pair would be inefficient



**Algorithm 1** Multi-Agent Dialog Framework (MADF)

```
1:  procedure MULTIBOTTRAIN
2:      while train_iter < max_train_iter do                          ▷ Main Training loop over batches
3:          Qbot ← random_select (Q_1, Q_2, Q_3....Q_q)
4:          Abot ← random_select (A_1, A_2, A_3....A_a)               ▷ Either q or a is equal to 1
5:          history ← (0, 0, ...0)                                    ▷ History initialized with zeros
6:          fact ← (0, 0, ...0)                                       ▷ Fact encoding initialized with zeros
7:          Δimage_pred ← 0                                           ▷ Tracks change in Image Embedding
8:          Qz_1 ← Ques_enc(Qbot, fact, history, caption)
9:          for t in 1:10 do                                         ▷ Have 10 rounds of dialog
10:             ques_t ← Ques_gen(Qbot, Qz_t)
11:             Az_t ← Ans_enc(Abot, fact, history, image, ques_t, caption)
12:             ans_t ← Ans_gen(Abot, Az_t)
13:             fact ← [ques_t, ans_t]                                ▷ Fact encoder stores the last dialog pair
14:             history ← concat(history, fact)                       ▷ History stores all previous dialog pairs
15:             Qz_t ← Ques_enc(Qbot, fact, history, caption)
16:             image_pred ← image_regress(Qbot, fact, history, caption)
17:             R_t ← (target_image − image_pred)^2 − Δimage_pred
18:             Δimage_pred ← (target_image − image_pred)^2
19:         end for
20:         Δ(w_Qbot) ← 1/10 Σ_{t=1}^{10} ∇_{θ_Qbot} [G_t log p(ques_t, θ_Qbot) − Δimage_pred]
21:         Δ(w_Abot) ← 1/10 Σ_{t=1}^{10} G_t ∇_{θ_Abot} log p(ans_t, θ_Abot)
22:         w_Qbot ← w_Qbot + Δ(w_Qbot)                               ▷ REINFORCE and Image Loss update for Qbot
23:         w_Abot ← w_Abot + Δ(w_Abot)                               ▷ REINFORCE update for Abot
24:     end while
25: end procedure
```



The little girl is standing with skis on her feet

## RESULTS

| | Metric | N | Supervised | RL 1Q,1A | RL 1Q,3A | RL 3Q,1A |
|---|---|---|---|---|---|---|
| 1 | Q-Bot Relevance | 8 | 2.5 | 2.75 | 2 | 2.75 |
| 2 | Q-Bot Grammar | 8 | **2.25** | 2.875 | 2.5 | 2.375 |
| 3 | A-Bot Relevance | 12 | 2.5 | 2.583 | 2.25 | **1.67** |
| 4 | A-Bot Grammar | 12 | 1.92 | 3.5 | **1.83** | 2.25 |
| 5 | Overall Coherence | 20 | 2.8 | 3.05 | **2.3** | 1.85 |

Quantitative Metrics (below) and Human Evaluations (above; lower is better; 20 evaluators). RL 1Q,3A refers to dialog system trained with 1 Q-Bot, 3 A-Bots

Overall **Dialog Coherence** of RL-1Q,3A and 3Q,1A systems ranked much better according to humans

- Multiple A-Bots interacting with Q- Bot improves relevance, and vice versa
- The grammar improves for both bots in both 1Q,3A and 3Q,1A settings
- Having multiple A-Bots to interact with exposes the Q-Bot to diverse dialog, leading to more stable updates with lower bias

| Model | MRR | Mean Rank | R@1 | R@5 | R@10 |
|---|---|---|---|---|---|
| Answer Prior (Das et al., 2016) | 0.3735 | 26.50 | 23.55 | 48.52 | 53.23 |
| MN-QIH-G (Das et al., 2016) | 0.5259 | 17.06 | 42.29 | 62.85 | 68.88 |
| HCIAE-G-DIS (Lu et al., 2017) | 0.547 | 14.23 | 44.35 | 65.28 | 71.55 |
| Frozen-Q-Multi (Das et al., 2017) | 0.437 | 21.13 | N/A | 53.67 | 60.48 |
| CoAtt-GAN (Wu et al., 2017) | 0.5578 | 14.4 | **46.10** | **65.69** | 71.74 |
| SL(Ours) | **0.610** | **5.323** | 34.74 | 57.67 | **72.68** |
| RL - 1Q,1A(Ours) | 0.459 | 7.097 | 16.04 | 54.69 | 72.34 |
| RL - 1Q,3A(Ours) | 0.601 | 5.495 | 34.83 | 57.47 | 72.48 |
| RL - 3Q,1A(Ours) | 0.590 | 5.56 | 33.59 | 57.73 | 72.61 |

We outperform all previous architectures in MRR, Mean Rank and R@10: **consistently good responses**