

1

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

Akshat Agarwal\*, Swaminathan Gurumurthy\*, Vasu Sharma\*, Katia Sycara

Adaptive Learning Agents Workshop International Conference on Autonomous Agents and Multiagent Systems (AAMAS) 2018





Interpretable, goal-oriented dialog between artificial agents

- Al needs to communicate with humans conveying visual and textual info
  - use cases: assistive systems for visually impaired, assistants (like Siri/Alexa)
- It will become common to have two agents communicate with each other towards a goal, say, reserving a table (like Google Duplex)
- We want these conversations to be interpretable to humans for sake of transparency and ease of debugging

## **Motivation**





Interpretable, goal-oriented dialog between artificial agents

- Humans adhere to natural language because they have to interact with an entire community
- Having a private language for each person would be inefficient
- Previous work on visual dialog showed a pair of agents adapting to each other start communicating in a private language to maximize the flow of information

We propose a multi-agent dialog framework (MADF) where each agent **interacts with and learns from multiple agents;** and show that it results in more coherent and human-interpretable dialog between agents, without compromising on task performance



- 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





 The VisDial dataset<sup>1</sup> contains ~80k images each with 10 pairs of questions and answers, collected from humans

**Carnegie Mellon** 

- 2. To elicit temporal continuity, grounding in the image and naturalistic conversations, workers were paired on AMT to chat in real time
- 3. One worker (Q) saw only the caption to a hidden image, and had to ask questions about the image to 'imagine the scene' better
- 4. The 2<sup>nd</sup> worker (A) saw image and caption, and had to answer the questions
- 5. 10 pairs of questions and answers exchanged for each image

<sup>1</sup> Das, Abhishek, et al. "Visual dialog." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Vol. 2. 2017.



- 1. The agents (Q-Bot and A-Bot) are pre-trained on the VisDial dataset using supervision<sup>1</sup>. They do not interact with each other in this phase
- 2. This is followed by making them interact and adapt to each other by reinforcement learning<sup>2</sup>
  - a. They are rewarded by the environment to maximize transfer of information with each QA pair
  - b. The transition from supervised to reinforcement learning is handled smoothly via a curriculum.

## Visual Dialog RL Framework





Note: The agents have been pretrained on the VisDial dataset before interacting as shown above

## **Agent Architectures**

## **Carnegie Mellon**



Fact ( $F_{t-1}$ ) is concatenation of a question ( $Q_{t-1}$ ) and its answer ( $A_{t-1}$ ) C is the caption History ( $F_1$ , $F_2$ ...  $F_{t-2}$ ) is a combination of all previous QA pairs for any particular image Attend over history using fact,  $F_{t-1}$ , to produce  $H_t^{(Q)}$ Concat  $F_{t-1}$ ,  $H_t^{(Q)}$  and C embeddings and pass through linear layers to get  $e_t^{(Q)}$  and  $y_t$ 

<sup>1</sup> Lu, Jiasen, et al. "Best of both worlds: Transferring knowledge from discriminative learning to a generative visual dialog model." Advances in Neural Information Processing Systems. 2017.

## **Agent Architectures**

## **Carnegie Mellon**



 $F_0(C)$  is the caption History  $(F_0, F_1, F_2...F_{t-1})$  is a combination of all previous QA pairs for any particular image Attend over history using question,  $Q_t$  to produce  $H_t^{(A)}$ Concat  $Q_t$ ,  $H_t^{(A)}$  and  $y_{gt}$  embeddings and pass through linear layer to get  $e_t^{(A)}$ 

<sup>1</sup> Lu, Jiasen, et al. "Best of both worlds: Transferring knowledge from discriminative learning to a generative visual dialog model." *Advances in Neural Information Processing Systems*. 2017.



- 1. Two agents, a Q-Bot and an A-Bot are first trained in isolation via supervision from the VisDial dataset for 15 epochs
- 2. Then smoothly transitioned to reinforcement learning via a curriculum
  - For the first K rounds of dialog for each image, agents are trained by supervision, and for remaining 10-K rounds they are made to interact and train via RL.
  - b. K starts at 9 and is reduced to 0 over 10 epochs

<sup>1</sup> Das, Abhishek, et al. "Visual dialog." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Vol. 2. 2017. <sup>2</sup> Das, Abhishek, et al. "Learning cooperative visual dialog agents with deep reinforcement learning." arXiv preprint arXiv:1703.06585 (2017).



- Both agents trained using a Maximum Likelihood Estimation (MLE) loss against the ground truth QA for every round of dialog<sup>1</sup>
- Q-Bot simultaneously minimizes Mean Squared Error (MSE) loss between the true and predicted image embeddings<sup>2</sup>
- 3. Problems:
  - a. MLE results in repetitive, 'safe' responses (*e.g. I don't know, I can't see*)
  - b. No interaction during training leads to unexpected responses during testing when they interact with each other and face out of distribution QA

<sup>1</sup> Das, Abhishek, et al. "Visual dialog." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Vol. 2. 2017. <sup>2</sup> Das, Abhishek, et al. "Learning cooperative visual dialog agents with deep reinforcement learning." arXiv preprint arXiv:1703.06585 (2017).



- 1. Both agents allowed to interact with each other and learn by self-play<sup>1</sup>
- 2. No ground-truth data except images and captions

- 3. Q-Bot observes  $\{c, q_1, a_1, ..., q_{10}, a_{10}\}$ , A-Bot observes  $\{I, c, q_1, a_1, ..., q_{10}, a_{10}\}$
- *I* : image, *c* : caption,  $q_i, a_j$ : i<sup>th</sup> dialog pair exchanged where i = [1,..10]
- 4. Action: Predict words sequentially until a stop token is encountered (or max length reached)
- 5. Reward: Incentivizing information gain from each round of QA, measured using the predicted image embedding  $y_t$

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

No motivation to stick to rules and conventions of English language<sup>2</sup>, making the RL optimization problem ill-posed

<sup>&</sup>lt;sup>1</sup> Das, Abhishek, et al. "Learning cooperative visual dialog agents with deep reinforcement learning." arXiv preprint arXiv:1703.06585 (2017). <sup>2</sup> Kottur, Satwik, et al. "Natural Language Does Not Emerge 'Naturally' in Multi-Agent Dialog." arXiv preprint arXiv:1706.08502 (2017).

## Multi-Agent Dialog Framework (MADF)





**Carnegie Mellon** 

We create either multiple Q-Bots to interact with a single A-Bot, OR multiple A-Bots to interact with a Q-Bot - and randomly pick one pair of agents to interact for each batch of images, and learn via REINFORCE

## **MADF Algorithm**

#### **Carnegie Mellon**



Algorithm 1 Multi-Agent Dialog Framework (MADF)								
1: procedure MULTIBOTTRAIN								
2:	while train_iter < max_train_iter do	Main Training loop over batches						
3:	$Qbot \leftarrow random\_select(Q_1, Q_2, Q_3, Q_q)$							
4:	$Abot \leftarrow random\_select (A_1, A_2, A_3A_a)$	$\triangleright$ Either q or a is equal to 1						
5:	$history \leftarrow (0, 0, 0)$	▷ History initialized with zeros						
6:	$fact \leftarrow (0, 0, 0)$	Fact encoding initialized with zeros						
7:	$\Delta image\_pred \leftarrow 0$	Tracks change in Image Embedding						
8:	$Qz_1 \leftarrow Ques\_enc(Qbot, fact, history, columnation)$	aption)						
9:	<b>for</b> t in 1:10 <b>do</b>	▷ Have 10 rounds of dialog						
10:	$ques_t \leftarrow Ques\_gen(Qbot, Qz_t)$							
11:	$Az_t \leftarrow Ans\_enc(Abot, fact, history, image, ques_t, caption)$							
12:	$ans_t \leftarrow Ans\_gen(Abot, Az_t)$							
13:	$fact \leftarrow [ques_t, ans_t]$	▷ Fact encoder stores the last dialog pair						
14:	$history \leftarrow concat(history, fact)$	History stores all previous dialog pairs						
15:	$Qz_t \leftarrow Ques\_enc(Qbot, fact, history)$	(a, caption)						
16:	$image\_pred \leftarrow image\_regress(Qbot$	, fact, history, caption)						
17:	$R_t \leftarrow (target\_image-image\_pred)$	$^{2}-\Delta image_{-}pred$						
18:	$\Delta image\_pred \leftarrow (target\_image\_pred)^2$							
19:	end for $1 - 10$							
20:	$\Delta(w_{Qbot}) \leftarrow \frac{1}{10} \sum_{t=1}^{10} \nabla_{\theta_{Qbot}} [G_t \log p(qu)]$	$(es_t,  heta_{Qbot}) - \Delta image\_pred]$						
21:	$\Delta(w_{Abot}) \leftarrow \frac{1}{10} \sum_{t=1}^{10} G_t \nabla_{\theta_{Abot}} \log p(ans)$	$_{t}, heta_{Abot})$						
22:	$w_{Qbot} \leftarrow w_{Qbot} + \Delta(w_{Qbot})$	▷ REINFORCE and Image Loss update for Qbot						
23:	$w_{Abot} \leftarrow w_{Abot} + \Delta(w_{Abot})$	▷ REINFORCE update for Abot						
24:	end while							
25: 0	end procedure							

## Quantitative Results for Answer Ranking



Model	MRR	Mean Rank	<b>R@1</b>	R@5	<b>R@10</b>
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

**Carnegie Mellon** 

About the Metrics: **Mean Rank and MRR** compute the average rank (and average of their reciprocals), respectively, assigned to the ground truth answer, over a set of 100 candidate answers for each question (provided in the VisDial dataset). **Recall@k** computes the percentage of answers with rank less than k

Used VisDial v0.9 dataset, with 83k train images + 40k test images. Results are reported on the test set

We outperform all previous architectures in MRR, Mean Rank and Recall @ 10, showing consistently good answers

While RL-1Q,1A performance drops (since it is being optimized for image estimation and not answer ranking, unlike SL), our multi-agent systems RL-1Q,3A and 3Q,1A **recover most of the performance gap** 

# Quantitative Results for Image Retrieval





**Carnegie Mellon** 

Fig. shows Image Retrieval Percentile Score (Y-axis) vs Dialog Round, from 1-10 (X-axis)

The score is calculated by ranking the Q-Bot's prediction of the image over all 40k images in the VisDial test set

- While the performance of SL decreases (because of the nature of LSTMs to forget), the performance of our RL systems (RL-1Q,3A and RL-3Q,1A) remain constant
- 2. This, combined with previous result, validate our hypothesis that our multi-agent dialog framework results in more coherent dialog without compromising on task performance

# Carnegie MellonQualitative Evaluations ofGenerated Language by Humans



2	Metric	Ν	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

- 20 evaluators shown one image and the 10 QA pairs generated for 4 systems, asked to rank on each metric
- 2. Overall Coherence of 1Q,3A and 3Q,1A ranked much better
- 3. Q-Bot more relevant in 1Q,3A and A-Bot more relevant in 3Q,1A which proves our hypothesis that exposure to multiple agents to interact with improves relevance

## **Sample Generations**



2 men in a grassy park playing frisbee			Human-Human Dialogue				SL			
		How many men? Are they playing frisbee? What color frisbee? Just white? Do men look skilled? Why not? They don't seem good? Is frisbee in air? Who has it? Is he bigger than other?		2 Yes White Yes No Playing Guess not, I said no No Person 1 no	are the how old do you s are they can you are sittin are they what co are the	boys playing volleyball ? do they look ? see any buildings ? male ? see any cars ? tell how many women ng there or playing ? combing leather ? lor are the men 's shoes? lor are the women ? boys well groomed ?	yes , a professional boy maybe in their late 20 's no , i don't see buildings yes , both men and women no , can only see the bunch of there are more people than 3 no , they 're wood women in they are black 'the snowboarder is white yes			
RL - 1Q, 1A			RL - 1Q, 3A			RL - 3Q, 1A				
do you recognize 1 ball ? that you could you are seeing controllers ? do you recognize the guy and age is the adult ? do you recognize them ? how many people are there ? do you do they fence famous ? do you see fire ? is the fire dog UNK ? do they have long manes ? can you see other people watching them ? balack and white, it looks you couldn't be late teens , his motorcycle but UNK tint and UNK motorcycle but UNK tint and UNK is the fire dog UNK ? balack and white, it looks Yes there is a people watching them ?			How old do the men appear? Is this at a beach? Do they have on bathing suits? How old are they? What color frisbee? Do they have a regular ball shirt With how old are they? Is there other people in the pic? How many of them are playing? What is the woman doing?	on?	Young ad W 1 of them Mid Yes, there is a man behind the sit		What color is umbrella? What are they wearing? What color is frisbee? What are they doing? Are they all holding rackets? Are there any other people? What color is the frisbee? Are there any other people? Are the people tall? Are they in a field?	Black with a blue stripe T shirts and jeans White Sitting on the beach, talking Yes Yes Creamy green Yes a lot Looks very tall no		



- Carnegie Mellon
  - 1. Understanding why task performance (image retrieval score) does not improve with dialog rounds (though intuitively, with more information having been exchanged, it should!)
  - 2. Improving the image embeddings used as ground truth to have richer information



- **Carnegie Mellon** 
  - 1. Through quantitative evaluations of answer ranking and image retrieval task performance, we show that our multi-agent systems generate interpretable dialog without compromising on task performance
  - 2. Through qualitative evaluations of dialog relevance and coherence by humans, we show that our multi-agent systems produce more coherent dialog
  - 3. This approach has also been validated in related works published in parallel:
    - a. Cao, Kris, et al. "Emergent Communication through Negotiation." arXiv preprint arXiv:1804.03980 (2018).
    - b. Lee, Jason, et al. "Emergent translation in multi-agent communication." arXiv preprint arXiv:1710.06922 (2017).





Github: https://goo.gl/gc6dGZ

# Thank You!

Questions?