# Supplementary Material
# Describing Unseen Videos via Multi-Modal Cooperative Dialog Agents

Ye Zhu[1], Yu Wu[2,3], Yi Yang[2], and Yan Yan[1]⋆

[1] Texas State University, San Marcos, USA
{ye.zhu, tom_yan}@txstate.edu
[2] ReLER, University of Technology Sydney, Australia
yu.wu-3@student.uts.edu.au,yi.yang@uts.edu.au
[3] Baidu Research, Beijing, China

We mainly provide additional qualitative experimental results and analysis in this document.

**Internal Dialog Generations.** Selected examples of the generated dialog interactions between two agents are shown in Fig. 1 and Fig. 2. The corresponding round of the ground truth dialog is also given as comparisons. During the training process, the question and answer decoders are trained to imitate the ground truth dialog by providing the ground truth words as sequence input labels, but the entire model is optimized in the final video description task. As shown in the figures, the generated dialog is not always natural and precise in the sense of natural language understanding, although the generated questions and answers sometimes do contain keywords that are later used in the final descriptions (*e.g.*, kitchen, fridge, couch, dish). *Q-BOT* also tends to ask repeated questions in the internal dialog generation process, especially with small starting round number.

**Video Descriptions.** More qualitative results about the final video description task are presented in Fig. 2. We observe that *Q-Basic* tends to generate similar descriptions for different inputs. It is a common and existing problem in the related research fields that combine natural language processing and computer vision [2, 1]. Other methods (*i.e.*, *Q-Strong*, *A-BOT* and our *Q-BOT* with cooperative learning) alleviate the problem and include more specific details in the generated descriptions. However, the sentence structures are not as varied as the human (GT) performance, which leads to another diversity problem that researchers seek to solve [3, 5, 4]. The qualitative results show that our *Q-BOT* with cooperative learning is able to describe some details that are not presented in the initial inputs (*e.g.*, the refrigerator in the first example), which proves that *Q-BOT* obtains useful information from the dialog interactions with *A-BOT*.

**Further Analysis.** We further analyze different aspects of this work. Firstly, we believe this task setting is of great potential to provide a new perspective that brings more secure and reliable AI systems. The task setting with implicit information sources is not only limited to the proposed video description task, but also applicable in a wider context such as scene understanding and segmentation. Secondly, under the overall framework of multi-modal learning

---

⋆ Corresponding author.

| Input of Q-BOT | Video descriptions | Ground Truth Dialog | Generated Dialog |
|---|---|---|---|



Input of Q-BOT:

Q1: is there a person in the video?
A1: yes, there is one person in the video.
Q2: is it a male or female?
A2: the person is male.

Video descriptions:

Human(GT): a man is vacuuming the kitchen floor. when he finishes he picks up some garments and leaves the kitchen.

Q-Basic: a man walks into a room and takes off his shirt. he then walks away.

Q-Strong: a man walks into a room and takes off his jacket. he then takes a bite of a sandwich and then walks away.

A-BOT: a man is vacuuming the floor. he then goes to the living room and grabs a vacuum cleaner.

**Ours: a man is in the kitchen. he is holding a cup of coffee. he takes a drink from a cup and walks out of the room.**

Ground Truth Dialog:

Q3: does more than one person appear in the video?
A3: there is only one person.
Q4: is he in the kitchen?
A4: yes, did you see the machine he uses.
Q5: is he holding something?
A5: yes, he is holding a handheld vacuum.
Q6: is he cleaning?
A6: yes he is, he is vacuuming the floor.
Q7: does he pick up anything else?
A7: he does very good, he picked up a garment.
Q8: does he walk out the kitchen?
A8: yes he walks out after picking up the garment.
Q9: where does he pick up the clothes from?
A9: towards the back of the kitchen behind the table.
Q10: does he come back into the kitchen after taking the clothes?
A10: no he does not.

Generated Dialog:

Q3: is there one person?
A3: yes, there is.
Q4: is there any person in pics?
A4: yes, there is one person.
Q5: is there anything in the kitchen?
A5: yes.
Q6: is there anything prepared?
A6: yes, i think so.
Q7: is there anything crossed?
A7: no, there is one person.
Q8: is there any person in pics?
A8: yes, i think so.
Q9: is there any person in pics?
A9: yes, i think so.
Q10: is there any person in pics?
A10: yes.

**Fig. 1.** Qualitative example with complete generated internal dialog.

with implicit information sources, there remain multiple future directions that worth being further exploited. For example, given the fact that *Q-BOT* tends to ask similar questions during the internal dialog process, we can consider incorporating internal supervision to motivate *Q-BOT* to ask more diverse and informative questions. In the meanwhile, we observe from the experiments that the performance does not necessarily improve when given more chances to freely ask/answer questions during the internal dialog interactions for a given video, indicating the fact that *Q-BOT* might have the potential to acquire the necessary information in less than 10 rounds of question-answer interactions. One possible future direction inspired by the observation is to encourage more efficient dialog interactions between the two agents.
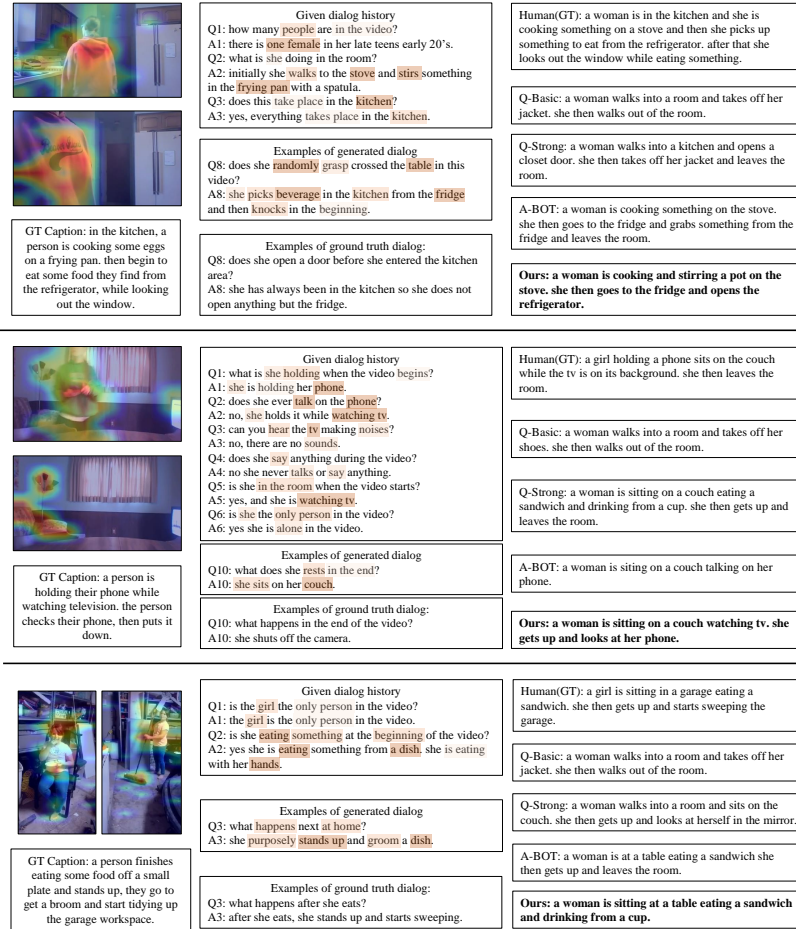
**Row 1**

GT Caption: in the kitchen, a person is cooking some eggs on a frying pan. then begin to eat some food they find from the refrigerator, while looking out the window.

Given dialog history
Q1: how many people are in the video?
A1: there is one female in her late teens early 20's.
Q2: what is she doing in the room?
A2: initially she walks to the stove and stirs something in the frying pan with a spatula.
Q3: does this take place in the kitchen?
A3: yes, everything takes place in the kitchen.

Examples of generated dialog
Q8: does she randomly grasp crossed the table in this video?
A8: she picks beverage in the kitchen from the fridge and then knocks in the beginning.

Examples of ground truth dialog:
Q8: does she open a door before she entered the kitchen area?
A8: she has always been in the kitchen so she does not open anything but the fridge.

Human(GT): a woman is in the kitchen and she is cooking something on a stove and then she picks up something to eat from the refrigerator. after that she looks out the window while eating something.

Q-Basic: a woman walks into a room and takes off her jacket. she then walks out of the room.

Q-Strong: a woman walks into a kitchen and opens a closet door. she then takes off her jacket and leaves the room.

A-BOT: a woman is cooking something on the stove. she then goes to the fridge and grabs something from the fridge and leaves the room.

**Ours: a woman is cooking and stirring a pot on the stove. she then goes to the fridge and opens the refrigerator.**

**Row 2**

GT Caption: a person is holding their phone while watching television. the person checks their phone, then puts it down.

Given dialog history
Q1: what is she holding when the video begins?
A1: she is holding her phone.
Q2: does she ever talk on the phone?
A2: no, she holds it while watching tv.
Q3: can you hear the tv making noises?
A3: no, there are no sounds.
Q4: does she say anything during the video?
A4: no she never talks or say anything.
Q5: is she in the room when the video starts?
A5: yes, and she is watching tv.
Q6: is she the only person in the video?
A6: yes she is alone in the video.

Examples of generated dialog
Q10: what does she rests in the end?
A10: she sits on her couch.

Examples of ground truth dialog:
Q10: what happens in the end of the video?
A10: she shuts off the camera.

Human(GT): a girl holding a phone sits on the couch while the tv is on its background. she then leaves the room.

Q-Basic: a woman walks into a room and takes off her shoes. she then walks out of the room.

Q-Strong: a woman is sitting on a couch eating a sandwich and drinking from a cup. she then gets up and leaves the room.

A-BOT: a woman is siting on a couch talking on her phone.

**Ours: a woman is sitting on a couch watching tv. she gets up and looks at her phone.**

**Row 3**

GT Caption: a person finishes eating some food off a small plate and stands up, they go to get a broom and start tidying up the garage workspace.

Given dialog history
Q1: is the girl the only person in the video?
A1: the girl is the only person in the video.
Q2: is she eating something at the beginning of the video?
A2: yes she is eating something from a dish. she is eating with her hands.

Examples of generated dialog
Q3: what happens next at home?
A3: she purposely stands up and groom a dish.

Examples of ground truth dialog:
Q3: what happens after she eats?
A3: after she eats, she stands up and starts sweeping.

Human(GT): a girl is sitting in a garage eating a sandwich. she then gets up and starts sweeping the garage.

Q-Basic: a woman walks into a room and takes off her jacket. she then walks out of the room.

Q-Strong: a woman walks into a room and sits on the couch. she then gets up and looks at herself in the mirror.

A-BOT: a woman is at a table eating a sandwich she then gets up and leaves the room.

**Ours: a woman is sitting at a table eating a sandwich and drinking from a cup.**

**Fig. 2.** More qualitative results for the unseen video description task. We provide the input static images to *Q-BOT* on the left. The given dialog history, examples of generated dialog and the ground truth dialog are presented in the middle. Comparisons of the final descriptions among different methods are shown on the right.

# References

1. Alamri, H., Cartillier, V., Das, A., Wang, J., Cherian, A., Essa, I., Batra, D., Marks, T.K., Hori, C., Anderson, P., et al.: Audio visual scene-aware dialog. In: CVPR (2019)
2. Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J.M., Parikh, D., Batra, D.: Visual dialog. In: CVPR (2017)
3. Jain, U., Lazebnik, S., Schwing, A.G.: Two can play this game: visual dialog with discriminative question generation and answering. In: CVPR (2018)
4. Wang, L., Schwing, A., Lazebnik, S.: Diverse and accurate image description using a variational auto-encoder with an additive gaussian encoding space. In: NeurIPS (2017)
5. Wu, Q., Wang, P., Shen, C., Reid, I., Van Den Hengel, A.: Are you talking to me? reasoned visual dialog generation through adversarial learning. In: CVPR (2018)