

Towards Fast Adaptation of User Simulation to New Domains

DSTC8 Track Proposal

1 Task Description

Recent advances in meta-learning have allowed to train neural networks to adapt to a new task after seeing only a few examples of the task. Most progress on the topic has been made in computer vision [1, 2] but results on some natural language processing tasks are encouraging [3, 4]. This proposal aims to leverage these new techniques for neural multi-domain dialogue systems [5–8].

Neural dialogue systems require very large datasets to learn to output consistent and grammatically-correct sentences [9–11]. In general, about 100k dialogues are necessary to learn a good language model in the chit-chat domain [9]. With this challenge, our goal is to investigate whether sample complexity can decrease with time, *i.e.*, if a dialogue system that was trained on a large corpus can learn to converse about a new domain given a much smaller in-domain corpus. This is a step towards scaling up the design of neural dialogue models.

We propose a task in two steps: 1) train a general-purpose dialogue system on a large corpus, and 2) adapt this model to new domains with little in-domain data. The first step allows to learn a good language model and some general conversational skills. For this step, we suggest to use the publicly available Reddit dataset [12]. We provide pre-processing code for this data so that all participants work on the same corpus. The task is to rapidly adapt these conversational skills to new domains. For this part, we have collected a corpus of goal-oriented domain-specific dialogues. These datasets are described in the following section.

2 Resources

Reddit Dataset We constructed a corpus of dialogues from Reddit submissions and comments spanning November 2017 through October 2018. Content is selected from a curated list of one thousand high-traffic subreddits. Our extraction and filtering methodology is based on that used in the DSTC7 sentence generation task [13], the key difference being we sample at most two threads per submission. The corpus consists of five million training dialogues, with an additional one million dialogues reserved for validation.

Goal-Oriented Corpus We collected 40 355 goal-oriented dialogues via crowd-sourcing. These dialogues span 51 domains and 242 tasks. For each dialogue, we paired two crowd-workers, one had the role of being a bot, and the other one was the user. Both were given a domain and a task. Examples of domains are: bus schedule, apartment search, alarm setting, banking and event reservation. We defined several tasks per domain. An example of task for the bus schedule domain is: *Inform the user that the bus stop they are asking about has been moved two blocks north* on the bot side, and *Ask if a certain bus stop is currently operational* on the user side. Statistics of the dataset are provided in Table 1. We provide a training/validation/test split for this dataset. We will also release an analysis of the complexity of the domains to provide some insights on the adaptation challenges.

Note that all entities were invented by the crowd-workers (for instance, the address of the bus stop) and the goal of this challenge is to automate the user utterances and not the bot utterances. Samples from this dataset are given in the Appendix.

| | Mean | Minimum | Maximum |
|--------------------------------|-------|---------|---------|
| Number of tasks per domain | 4.8 | 3 | 11 |
| Number of dialogues per domain | 791.9 | 288 | 1990 |
| Number of dialogues per task | 167.6 | 32 | 285 |
| Number of turns per dialogue | 11.9 | 9 | 57 |

Table 1: Statistics of the goal-oriented dataset

Baseline Model We will provide code for a baseline meta-learner that participants could build upon. We will also release evaluation results for this model at the beginning of the challenge.

3 Evaluation

The automatic evaluation of dialogue generation is a highly debated topic [14–16]. For the last DSTC, some tracks proposed to use metrics such as BLEU [17], ROUGE [18], etc. for development and then performed human evaluation to make the final ranking of the submitted models. For this challenge, we do not advise using word-overlap metrics. One reason for this decision is that tasks in the goal-oriented corpus do not specify entity names, so the users and the bots came up with names. Therefore, word-overlap metrics will potentially measure a lot of noise. We leave it to the appreciation of the participants to use these metrics for development but we will only perform human evaluation to rank the submitted models.

The challenge will be divided into domain generalization and task generalization. We will hold out a subset of the domains and tasks, respectively. We will provide a small fixed size training set as input to the meta learner and a test set. The test set will only contain unfinished dialogues. Examples of input and output are given in the Appendix. Participants are asked to submit a single response for each unfinished dialogue. Humans will rate responses for appropriateness, informativeness and utility [19] *with respect to the task the user is performing in the dialogue*. We will additionally evaluate up to one model per team where third-party input data has been used. The code for performing this evaluation will be released and available for participants to use during development.

4 Organizers

The organizers of this challenge are all affiliated with Microsoft Research in Montréal: Hannes Schulz, Adam Atkinson, Kaheer Suleman, Mahmoud Adada, Shikhar Sharma, Wendy Tay, Layla El Asri. Their email addresses are all formatted as: first.last@microsoft.com (for Layla El Asri: layla.elasri@microsoft.com).

References

- [1] S. Ravi and H. Larochelle. “Optimization as a model for few-shot learning”. In: *Proceedings of the International Conference on Learning Representations*. 2017.
- [2] C. Finn, P. Abbeel, and S. Levine. “Model-agnostic meta-learning for fast adaptation of deep networks”. In: *Proceedings of the International Conference on Machine Learning*. 2017.
- [3] J. Gu, Y. Wang, Y. Chen, K. Cho, and V. O. Li. “Meta-Learning for Low-Resource Neural Machine Translation”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 2018.
- [4] T. Munkhdalai, X. Yuan, S. Mehri, and A. Trischler. “Rapid adaptation with conditionally shifted neurons”. In: *Proceedings of the International Conference on Machine Learning*. 2018.
- [5] B. Peng, X. Li, L. Li, J. Gao, A. Çelikyilmaz, S. Lee, and K. Wong. “Composite Task-Completion Dialogue Policy Learning via Hierarchical Deep Reinforcement Learning”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 2017.
- [6] A. Rastogi, D. Hakkani-Tür, and L. P. Heck. “Scalable Multi-Domain Dialogue State Tracking”. In: *Proceedings of Automatic Speech Recognition and Understanding*. 2017.
- [7] P. Budzianowski, T.-H. Wen, B.-H. Tseng, I. Casanueva, S. Ultes, O. Ramadan, and M. Gašić. “MultiWOZ - A Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 2018.
- [8] I. Shalyminov, A. Eshghi, and O. Lemon. “Multi-Task Learning for Domain-General Spoken Disfluency Detection in Dialogue Systems”. In: *Proceedings of the Workshop Series on the Semantics and Pragmatics of Dialogue*. 2018.
- [9] O. Vinyals and Q. V. Le. “A Neural Conversational Model”. In: *arXiv:1506.05869* (2015).
- [10] J. Li, W. Monroe, A. Ritter, M. Galley, J. Gao, and D. Jurafsky. “Deep Reinforcement Learning for Dialogue Generation”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 2016.
- [11] T.-H. Wen, Y. Miao, P. Blunsom, and S. Young. “Latent Intention Dialogue Models”. In: *Proceedings of the International Conference on Machine Learning*. 2017.
- [12] J. Baumgartner. *Pushshift website. Reddit dataset*. <https://files.pushshift.io/reddit/>. Accessed 15-April-2018. 2018.
- [13] K. Yoshino, C. Hori, J. Perez, L. F. D’Haro, L. Polymenakos, C. Gunasekara, W. S. Lasecki, J. Kummerfeld, M. Galley, C. Brockett, J. Gao, B. Dolan, S. Gao, T. K. Marks, D. Parikh, and D. Batra. “The 7th dialog system technology challenge”. In: *arXiv preprint* (2018).
- [14] C.-W. Liu, R. Lowe, I. Serban, M. Noseworthy, L. Charlin, and J. Pineau. “How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 2016.

- [15] S. Sharma, L. E. Asri, H. Schulz, and J. Zumer. “Relevance of Unsupervised Metrics in Task-Oriented Dialogue for Evaluating Natural Language Generation”. In: *arXiv:1706.09799* (2017).
- [16] J. Novikova, O. Dušek, A. C. Curry, and V. Rieser. “Why We Need New Evaluation Metrics for NLG”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. 2017.
- [17] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. “BLEU: a Method for Automatic Evaluation of Machine Translation”. In: *Proceedings of the Conference of the Association for Computational Linguistics*. 2002.
- [18] C.-Y. Lin. “ROUGE: A Package for Automatic Evaluation of summaries”. In: *Proceedings of the ACL workshop on Text Summarization Branches Out*. 2004.
- [19] M. Galley, C. Brockett, X. Gao, B. Dolan, and J. Gao. “End-to-End conversation Modeling: DSTC7 Task 2 Description”. In: *DSTC7 workshop (forthcoming)*.

Appendix

Goal-Oriented Dialogues

List of Domains

Calendar update, pizza ordering, movie listings, event reservation, weather checking, flight booking, contact information update, restaurant reservation, playlist editing, looking up information, shopping, store details, sports information, quote of the day, how to, prompt generation, library request, banking, restaurant picking, name suggestion, vacation ideas, city information, music suggestion, agreement, pet advice, apartment search, Guinness record check, geography, alarm setting, contact management, phone settings, appointment reminder, home, policy, decision making, catalogue lookup, ski, bus schedule, insurance, information retrieval, sorting, scam lookup, time zone, play schedule, game rules, wedding planning, status check, gift ideas, tourism, hotel reservation and phone plan.

Samples of Dialogues

Domain: Ski

User Task: You want to know if there are good ski hills an hour’s drive from your current location.

Bot Task: Tell the user that there are no ski hills in their immediate location.

Bot Hello how may I help you?

User Is there any good ski hills an hour’s drive from my current location?

Bot I’m sorry to inform you that there are no ski hills in your immediate location

User Can you help me find the nearest?

Bot Absolutely! It looks like you’re about 3 hours away from Bear Mountain. That seems to be the closest.

User Hmm.. sounds good

Bot Alright! I can help you get your lift tickets now!When will you be going?

User Awesome! please get me a ticket for 10pax

Bot You’ve got it. Anything else I can help you with?

User None. Thanks again!

Bot No problem!

Example of input/output for this dialogue:

Input: domain = Ski, user task = You want to know if there are good ski hills an hour's drive from your current location., bot task = Tell the user that there are no ski hills in their immediate location, dialogue history = Hello how may I help you?

Output: user response = Is there any good ski hills an hour's drive from my current location?

Domain: Game rules

User Task: Ask how to win at the text adventure Zork. If the Bot cannot help you, ask them what games it can help you with.

Bot Task: Tell the user that you are programmed to help them play games, not win them. If the user starts talking about other games, go along with them and try to help.

Bot Hello how may I help you?

User I want to know how I can win playing Zork?

Bot I am programmed to help you play games, not win them

User What games can you help me with?

Bot Any game

User Can you specify a certain game?

Bot Zork

User What are the rules of Zork?

Bot You have to follow through a cave adventure

User What is the ending of the game?

Bot The user comes out of the cave

Example of input/output for this dialogue:

Input: domain = game rules, user task = Ask how to win at the text adventure Zork. If the Bot cannot help you, ask them if what games it can help you with., bot task = Tell the user that you are programmed to help them play games, not win them. If the user starts talking about other games, go along with them and try to help., dialogue history = Hello how may I help you?; I want to know how I can win playing Zork?; I am programmed to help you play games, not win them

Output: user response = What games can you help me with?

Domain: Music suggestion

User Task: Tell the bot that you like a certain genre of music, but don't prompt it to give you suggestions

Bot Task: Provide suggestions, even if the user doesn't ask for them (or tells you they don't need any)

Bot Hello how may I help you?

User Hello. I like classic rock music.

Bot Dream on Aerosmith

User No, I'm not a huge fan of that song. Is there anything else available?

Bot Sweet Home Alabama

User Well, I do like Lynyrd Skynyrd, but that song is so cliché. Do you have anything else by the same band?

Bot Free bird

User Yes, that is the perfect song for how I am feeling right now!

Bot Great, Is there anything else?

User No, but I really appreciate your assistance. Thank you so much!

Bot You're welcome

Example of input/output for this dialogue:

Input: domain = music suggestion, user task = Tell the bot that you like a certain genre of music, but don't prompt it to give you suggestions, bot task = Provide suggestions, even if the user doesn't ask for them (or tells you they don't need any)., dialogue history = Hello how may I help you?; Hello. I like classic rock music.; Dream on Aerosmith; No, I'm not a huge fan of that song. Is there anything else available?; Sweet Home Alabama;

Output: user response = Well, I do like Lynyrd Skynyrd, but that song is so cliché. Do you have anything else by the same band?

Domain: Information retrieval

User Task: Describe a movie where these guys have to climb a mountain and then throw a thing into the mountain and ask the bot if they know what the movie is called.

Bot Task: Ask the user to elaborate as much as possible about the movie they are asking about, and then suggest they're thinking of Lord of the Rings.

Bot Hello how may I help you?

User I'm trying to remember this movie. Can you help?

Bot I would be happy to help. Can you tell me what it is about?

User In it these guys have to climb up a mountain. Then they throw something into the mountain.

Bot Do you know what genre the movie is, such as comedy, horror, sci-fi, or fantasy?

User Not sure about that.

Bot Do you any of the actor or character names in the movie?

User No... All I can remember is the stuff about the mountain.

Bot Is it a new movie, or an older movie?

User I don't know.

Bot It sounds like the movie you are talking about could be Lord of the Rings. Does that sound right?

User It could be that. Thanks for the help.

Example of input/output for this dialogue:

Input: domain = information retrieval, user task = Describe a movie where these guys have to climb a mountain and then throw a thing into the mountain and ask the bot if they know what the movie is called., bot task = Ask the user to elaborate as much as possible about the movie they are asking about, and then suggest they're thinking of Lord of the Rings., dialogue history = Hello how may I help you?; I'm trying to remember this movie. Can you help?; I would be happy to help. Can you tell me what it is about?; I would be happy to help. Can you tell me what it is about?

Output: user response = In it these guys have to climb up a mountain. Then they throw something into the mountain.