

Compressed Context for Schema Guided Zero Shot Dialog

Anonymous Authors

Affiliation

Address Line 1

Abstract

Creating Task Oriented Dialog (TOD) Systems that generalize well to unseen domains has been a challenging research area in dialog. Most systems use the dialog history as context, and as the number of turns in a dialog increases, the context becomes too long and in many cases contain repetitive and unnecessary information. Also, the metrics used to evaluate dialog systems mainly focus on dialog state tracking (DST) and response generation, and seem to neglect system actions. In this paper, we propose a novel TOD system that uses a compressed context consisting of the latest DST and the last user utterance. Using a compressed context allows us to feed additional information like the schema, list of system actions, user actions and service results to the model, and still use language models like GPT-2. We propose a two step training process, where in the first pass we calculate the cross entropy loss on the context and target, and in the second pass we calculate the loss only on the target. To get a better understanding of the system actions, we propose two new metrics called Average Action Accuracy (AAA) and Joint Action Accuracy (JAA). Experimental results on the Schema Guided Dialogue (SGD) dataset show that our model outperforms the state of the art models in terms of zero shot generalization.

1 Introductions

- Basic Tod intro
- Challenge of zero-shot TOD
- Existing methods
- Our Method
- Pros of our methodology
- Summary of results

TOD systems interact with users in the form of dialog using natural language, to accomplish user tasks. The system needs to understand user needs and provide the best possible response to the user. The task of extracting user intent and goals from conversations by filling belief slots is called Dialog State Tracking (DST) (Wang, Liu, and Zhao 2016).

Using the DST and dialog history, the system needs to decide what actions to take and then convey that action in the form of natural language to the user.

Traditional TOD systems were built using a pipeline approach, where each component was created separately and then integrated together. However, with the adaptation of large pretrained language models (Devlin et al. 2019; Radford et al. 2019), researchers have moved towards end-to-end systems (Hosseini-Asl et al. 2020; Peng et al. 2021; Lee et al. 2020; Yang, Li, and Quan 2020; Jeon and Lee 2021; Sun et al. 2022; Yang et al. 2022; Noroozi et al. 2020), and have formulated the problem as a cascaded generation problem (Su et al. 2021), which is to sequentially generate the DST, system action and response.

2 Methodology

Steps: Pretraining(CE Loss), Training (loss on target only)

Context: Dialog History like SimpleTOD DST instead of dialog history

Additional info Schema List of system actions User actions Service Results

- Summary and importance of approach
- Base model (GPT-2)
- Prompt (context, additional info)
- Training (pretraining, training)

As shown in Figure ??, our model receives the DST and user utterance of the previous turn, and additional information as the context and auto-regressively generates the target, which consists of the updated DST, system action and response.

2.1 Training

We train the model in two steps, where in the first step we calculate the cross entropy loss on the context and target, and in the second step we calculate the loss only on the target. The goal of the first step is for the model to understand the general structure of the text and the second step is to fine tune the model to generate the target text.

3 Experimental Setup

3.1 Datasets

The Schema Guided Dialogue (SGD) dataset is a large scale dataset for task oriented dialogue that consists of over 16K multi domain dialogs between a human and a virtual assistant covering 16 domains. The dataset also provides a schema for each domain that provides a textual description of the domain, list of slots and list of intents. A slot contains a name, textual description, and possible values for categorical slots and an intent contains a name, textual description, optional slots and result slots.

3.2 Evaluation Metrics

To have a fair comparison with other methods, we use the evaluation script from SGD dataset and report the Active Intent Accuracy, Requested Slot F1, Average Goal Accuracy and Joint Goal Accuracy.

However, the metrics in the SGD dataset are geared towards DST, and do not take into account the system actions and response. To evaluate task completion we report Inform, Success, Average Action Accuracy (AAA) and Joint Action Accuracy (JAA). For evaluating response generation we report the ROUGE-2 (Lin and Och 2004) score and GLEU (Wu et al. 2016) instead of BLEU as it performs better on individual sentence pairs. The combined score is calculated as suggested in (Mehri, Srinivasan, and Eskenazi 2019) with $(\text{Inform} + \text{Success}) \times 0.5 + \text{GLEU}$.

The metric Inform measures whether a system has provided a correct entity and Success measures whether it has answered all the requested information. To calculate inform, from the ground truth system actions, we filter actions by action type inform (Inform, Inform Count) and check if they are predicted correctly. To calculate success, we filter actions by slot names that are in the requested slots and check if the action slot values are predicted correctly. AAA and JAA are similar to the goal metrics in SGD, but are calculated from system actions.

4 Results

- Best Results
- Comparison Table with other methods
- Ablation Study
- Ablation table of our model to show importance of each component

We compare our approach with other popular approaches that have been proposed for task oriented dialog and report the results in Table 1.

4.1 Ablation Study

To better understand the effect of different components in our approach, we present an ablation study in Table ??.

Describe effect of different components.

5 Related Works

5.1 Supervised End to End Models

Pretrained language models like BERT (Devlin et al. 2019), GPT-2 (Radford et al. 2019) and T5 (Raffel et al. 2019) have been used extensively in the literature for End to End models for TOD systems (Hosseini-Asl et al. 2020), (Peng et al. 2021), (Lee et al. 2020), (Yang, Li, and Quan 2020), (Jeon and Lee 2021), (Sun et al. 2022), (Yang et al. 2022), (Noroozi et al. 2020). In these models, the context consists of user and system utterance, whereas in our model we use the last user utterance and the previous state DST as context. Moreover, most of these models have the best performance in supervised settings and do not have the primary focus on zero-shot generalization.

5.2 Zero Shot End to End Models

Zero Shot DST models (Feng, Wang, and Li 2020), (Zhao et al. 2022) incorporate schema as part of the context and generalize well for DST, however these models do not focus on system actions and response generation. (Noroozi et al. 2020)

References

- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *ArXiv abs/1810.04805*.
- Feng, Y.; Wang, Y.; and Li, H. 2020. A sequence-to-sequence approach to dialogue state tracking. In *Annual Meeting of the Association for Computational Linguistics*.
- Hosseini-Asl, E.; McCann, B.; Wu, C.-S.; Yavuz, S.; and Socher, R. 2020. A simple language model for task-oriented dialogue. *ArXiv abs/2005.00796*.
- Jeon, H., and Lee, G. G. 2021. Dora: Toward policy optimization for task-oriented dialogue system with efficient context. *Comput. Speech Lang.* 72:101310.
- Lee, H.; Jo, S.; Kim, H.; Jung, S.; and Kim, T.-Y. 2020. Sumbt+larl: End-to-end neural task-oriented dialog system with reinforcement learning. *ArXiv abs/2009.10447*.
- Lin, C.-Y., and Och, F. 2004. Looking for a few good metrics: Rouge and its evaluation. In *Ntcir workshop*.
- Mehri, S.; Srinivasan, T.; and Eskenazi, M. 2019. Structured fusion networks for dialog. *arXiv preprint arXiv:1907.10016*.
- Noroozi, V.; Zhang, Y.; Bakhturina, E.; and Kornuta, T. 2020. A fast and robust bert-based dialogue state tracker for schema guided dialogue dataset. *ArXiv abs/2008.12335*.
- Peng, B.; Li, C.; Li, J.; Shayandeh, S.; Lidén, L.; and Gao, J. 2021. Soloist: Building task bots at scale with transfer learning and machine teaching. *Transactions of the Association for Computational Linguistics* 9:807–824.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language models are unsupervised multitask learners.

Model	Setting	Avg act Acc	Joint act Acc	Avg goal Acc	Joint goal Acc	Inform	Intent Acc	Req Slots F1	Response BLEU	Response ROUGE	Success	Avg u.act Acc	Joint u.act Acc	Combined
Our	all	56.25	46.07	73.53	45.59	66.46	88.31	95.61	16.22	18.52	58.63	74.08	65.63	78.76
	seen	58.41	48.64	75.7	49.87	68.7	89.09	95.97	23.12	18.52	61.05	75.97	67.99	88
	unseen	55.56	45.31	73.15	44.85	65.14	88.19	95.51	14.29	18.52	58.07	73.5	64.94	75.89

Table 1: Main Results

Raffel, C.; Shazeer, N. M.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *ArXiv abs/1910.10683*.

Su, Y.; Shu, L.; Mansimov, E.; Gupta, A.; Cai, D.; Lai, Y.-A.; and Zhang, Y. 2021. Multi-task pre-training for plug-and-play task-oriented dialogue system. *arXiv preprint arXiv:2109.14739*.

Sun, H.; Bao, J.; Wu, Y.; and He, X. 2022. Bort: Back and denoising reconstruction for end-to-end task-oriented dialog. In *NAACL-HLT*.

Wang, B.; Liu, K.; and Zhao, J. 2016. Inner attention based recurrent neural networks for answer selection. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1288–1297. Berlin, Germany: Association for Computational Linguistics.

Wu, Y.; Schuster, M.; Chen, Z.; Le, Q. V.; Norouzi, M.; Macherey, W.; Krikun, M.; Cao, Y.; Gao, Q.; Macherey, K.; Klingner, J.; Shah, A.; Johnson, M.; Liu, X.; Łukasz Kaiser; Gouws, S.; Kato, Y.; Kudo, T.; Kazawa, H.; Stevens, K.; Kurian, G.; Patil, N.; Wang, W.; Young, C.; Smith, J.; Riesa, J.; Rudnick, A.; Vinyals, O.; Corrado, G.; Hughes, M.; and Dean, J. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation.

Yang, Y.; Ding, H.; Liu, Q.; and Quan, X. 2022. Ubarv2: Towards mitigating exposure bias in task-oriented dialogs. *ArXiv abs/2209.07239*.

Yang, Y.; Li, Y.; and Quan, X. 2020. Ubar: Towards fully end-to-end task-oriented dialog systems with gpt-2. In *AAAI Conference on Artificial Intelligence*.

Zhao, J.; Gupta, R.; Cao, Y.; Yu, D.; Wang, M.; Lee, H.; Rastogi, A.; Shafraan, I.; and Wu, Y. 2022. Description-driven task-oriented dialog modeling. *ArXiv abs/2201.08904*.