

WeChat AI’s Submission for DSTC9 Interactive Dialogue Evaluation Track

Zekang Li¹²⁴, Zongjia Li²³, Jinchao Zhang², Yang Feng^{1*}, Jie Zhou²

¹Key Laboratory of Intelligent Information Processing
Institute of Computing Technology, Chinese Academy of Sciences

²WeChat AI, Tencent Inc, China

³School of EECS, Peking University

⁴University of Chinese Academy of Sciences

{lizakang19g, fengyang}@ict.ac.cn, zongjiali@pku.edu.cn
{dayerzhang, withtomzhou}@tencent.com

Abstract

We participate in the DSTC9 Interactive Dialogue Evaluation Track (Gunasekara et al. 2020) sub-task 1 (Knowledge Grounded Dialogue) and sub-task 2 (Interactive Dialogue). In sub-task 1, we employ a pre-trained language model to generate topic-related responses and propose a response ensemble method for response selection. In sub-task2, we propose a novel Dialogue Planning Model (DPM) to capture conversation flow in the interaction with humans. We also design an integrated open-domain dialogue system containing pre-process, dialogue model, scoring model, and post-process, which can generate fluent, coherent, consistent, and human-like responses. We tie 1st on human ratings and also get the highest Meteor, and Bert-score in sub-task 1, and rank 3rd on interactive human evaluation in sub-task 2.

Introduction

Our WeChat AI team participates in the DSTC9 Interactive Dialogue Evaluation Track sub-task 1 and sub-task 2. Sub-task 1 is knowledge-grounded dialogue generation on the Topical-Chat dataset, which is evaluated in a static manner. Sub-task 2 aims to extend dialog models beyond datasets and interactively evaluates dialogue systems with real users on DialPort (Zhao, Lee, and Eskenazi 2016). We mainly focus on improving the topical relevance of responses in sub-task 1 and improving topic depth, consistency, human-likeness in sub-task 2.

In the sub-task 1 (Knowledge Grounded Dialogue), our model architecture is built on the GPT2 model (Radford et al. 2019). We mainly focus on exploring better decoding methods and ensemble methods. In the decoding stage, sampling-based methods are likely to generate more diverse

and human-like responses compared to beam-search. But sampling-based methods always suffer from less topical relevance. To improve the topical relevance of responses, we propose a metric-based ensemble method for response selection.

In the sub-task 2 (Interactive Dialogue), the ultimate goal of the open domain dialogue system is to interact with real users effectively. Recently, due to pre-training with large scale transformer models, open-domain dialogue systems (Zhang et al. 2019b; Bao et al. 2020; Smith et al. 2020; Adwardana et al. 2020) have achieved great success. However, based on our real interaction experience with the systems, they can be improved on flexibility, topic depth, and consistency.

As to flexibility, current dialogue models are often lost in the dialogue when real users change the topic by accident because there is little topic change in the training data. Therefore, we propose a simple but efficient data augmentation method by constructing more flexible training data.

As to topic depth, as shown in Figure 2, we consider the whole dialogue process as many vector operations and propose the Dialogue Planning Model to capture the topic flow in the dialogue, which can improve topic depth effectively in the interaction with real users.

As to consistency, many dialogue models are likely to ignore the information in history when generating responses, such as opinion conflicts, personal information conflicts, and so on, which severely harms the dialogue consistency. To improve dialogue consistency, we employ a natural language inference model to detect the responses that conflict with dialogue history.

Related Work

Knowledge-Grounded Open-Domain Dialogue.

Knowledge-grounded open-domain dialogue is an important step towards a human-like dialogue system. Recent works mainly obtain knowledge from knowledge graphs (Zhou et al. 2018; Tuan, Chen, and Lee 2019), from unstructured

*Joint work with Pattern Recognition Center, WeChat AI, Tencent Inc, China. Yang Feng is the corresponding author. This work was done when Zekang Li and Zongjia Li was interning at Pattern Recognition Center, WeChat AI, Tencent.
Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

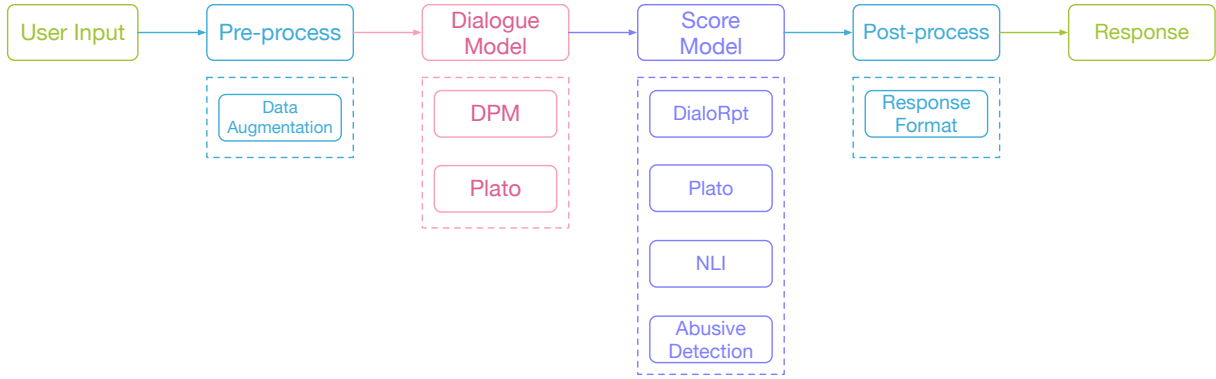


Figure 1: Overview of the interactive dialogue system. It contains 4 modules: Pre-process, Dialogue Model, Score Model, and Post-process. Pre-process mainly contains data augmentation. Dialogue Model includes two dialogue generation model: DPM and Plato, and generates response candidates. Score Model aims to select the most appropriate response through overall score by four scoring model. Post-process is used to make the response more human-like.

text (Dinan et al. 2018; Li et al. 2019; Lian et al. 2019; Kim, Ahn, and Kim 2020), and from visual information (Li et al. 2020; Huber et al. 2018). In the sub-task 1, we focus on using unstructured text as knowledge. Rather than learning from scratch like most recent work, we utilize a pre-trained language model and propose a response ensemble method to generate responses with more knowledge relevance.

Pre-trained Language Model.

Pre-trained Language models have brought about much improvement on various NLP tasks. GPT2 (Radford et al. 2019) and BERT (Devlin et al. 2018) are representative uni-directional and bi-directional language models. Based on GPT2, DialoGPT (Zhang et al. 2019b) is trained for dialogue response generation using Reddit comments. Meena (Adiwardana et al. 2020) utilize more social media data to make the chatbot more human-like. Besides, Blender finetunes the pre-trained model on human-annotated conversations. Furthermore, towards more diverse, human-like responses, Plato (Bao et al. 2020) introduces discrete latent variable and curriculum learning in the training process. To improve dialogue coherence and topic depth, we introduce a dialogue flow method upon the GPT2 model.

Our Method

For sub-task 1 (Knowledge Grounded Dialogue), we finetune the large pre-trained language model on the Topical-Chat dataset to generate topic-related knowledge grounded response candidates. Besides, we propose a response ensemble method to improve the topical relevance of response. For sub-task 2 (Interactive Dialogue), we build a dialogue system, which consists of four modules: pre-process, dialogue model, score model, and post-process, as shown in Figure 1.

Sub-task 1: Knowledge-Grounded Dialogue Generation

This task is to generate a response to a fixed dialogue context given the topic-related facts. Formally, let C and R denote

the dialogue context and response respectively, and K denote the topical-related facts. The probability to generate the response can be computed as:

$$P(R|C, K; \theta) = \prod_{i=1}^N P(R_i|C, K, R_{<i}; \theta) \quad (1)$$

where θ is the learnable parameter.

Model. The model architecture is based on gpt2-large (Radford et al. 2019). For fine-tuning, we concatenate the fact K , the dialogue context C , and the golden response R as the input sequence, and optimize the model by minimizing the following loss:

$$\mathcal{L} = - \sum_{i=1}^N \log(P(R_i|C, K, R_{<i}; \theta)) \quad (2)$$

To generate more diverse and human-like responses, we employ top-p sampling method (Holtzman et al. 2019) rather than greedy decoding and beam search.

Algorithm 1 Metric-based ensemble method for response selection.

- 1: **Input:** Response candidates $\{r_i | i = 1, 2, \dots, N\}$.
 - 2: **Output:** The most topic-related response r .
 - 3: Select a metric, such as Meteor Bert-score, and BLEU;
 - 4: Initialize metric score $M = \{m_i = 0 | i = 1, 2, \dots, N\}$;
 - 5: **for each** $i \in [1, N]$ **do**
 - 6: **for each** $j \in [1, i) \cup (i, N]$ **do**
 - 7: $m_i = m_i + Metric(r_i, r_j)$
 - 8: **end for**
 - 9: $m_i = m_i / (N - 1)$
 - 10: **end for**
 - 11: **Return:** $r_{argmax(M)}$.
-

Response Ensemble Method. Sampling-based decoding methods always bring about more diversity but suffer from less topical relevance. To improve the topical relevance, we propose a metric-based ensemble method to select the most

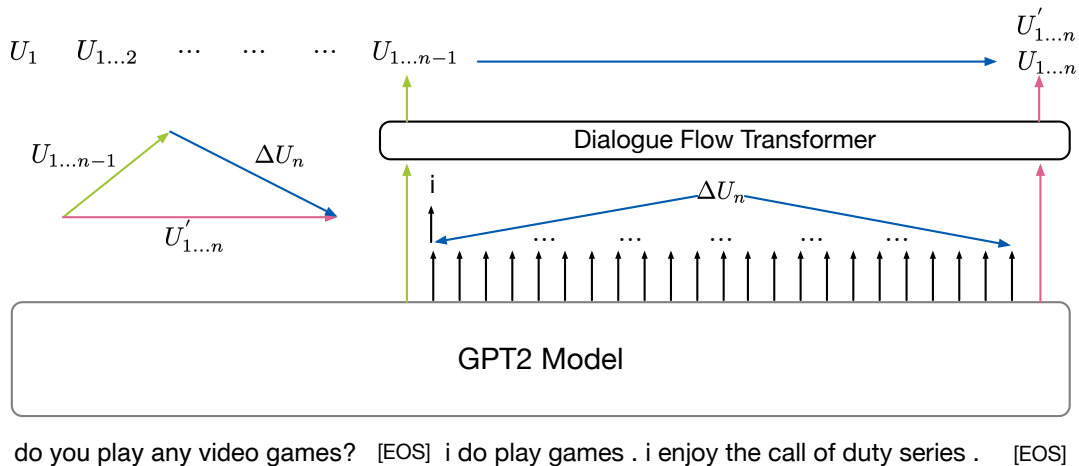


Figure 2: Dialogue planning model (DPM). We consider the whole dialogue process as many vector operations. $U_{1...n}$ denotes the representation of $1 \sim n$ utterance encoded by GPT2 model. The representation of utterance u_n conditioned on $u_{<n}$ can be calculated by $U'_{1...n} - U_{1...n-1}$, where $U'_{1...n}$ is the predicted representation by the dialogue flow transformer block.

Sampled Text	Meteor-self	Meteor-GT
Example 1: the warriors played the nba finals at the cow palace because the oakland arena was booked.	0.652	0.592
Example 2: the golden state warriors played the home games in the 1975 nba finals at the cow palace.	0.590	0.591
Example 3: the cow palace was the place to watch games in 1975.	0.341	0.540
Example 4: the golden state warriors played at the cow palace because the oakland arena was booked.	0.657	0.808
Example 5: the golden state warriors played in 1975 at the cow palace because the oakland arena was booked.	0.734	0.811
Ground-truth: in 1975 the golden state warriors had to play at the cow palace because their arena was booked.	-	-

Table 1: Response examples for the response ensemble method. As shown in the table, the sampled texts always lack of different parts of information. The response ensemble method can select the most integrated response. Meteor-self means the Meteor score with the other texts and Meteor-GT means the Meteor score with the ground-truth.

topical-relevant response from the generated response candidates, as shown in Algorithm 1. The method can improve the information integrity as well as the reference-based metric performance. As shown in Table 1, the 5 sampled texts lack of different parts of information. Through the method, we can select the most integrated and appropriate response.

Sub-task 2: Interactive Dialogue

There are many challenges to building an effective dialogue system, such as improving dialogue consistency, dialogue topic depth, and human-likeness. In this part, we will describe our methods to build and improve the dialogue system.

System Overview. As shown in Figure 1, our interactive open-domain dialogue system contains four modules: Pre-process, Dialogue Model, Score Model, and Post-process.

Pre-process. During the interaction, real users always change the dialogue topic by accident and have different preferences on topic depth. Therefore, we propose a sim-

ple but efficient data augmentation method to handle these problems. For different topic depth, given dialogue A , we randomly cut out some utterances from the end and get a new dialogue C . For the topic changes, given dialogue C and dialogue B , we concatenate C and B as a new dialogue D . In the training stage, we randomly sample dialogue C and D according to a fixed probability.

Dialogue Model. To generate diverse, informative, fluent responses, in the Dialogue model module, we employ two models: Dialogue Planning Model (DPM) and Plato (Bao et al. 2020) which is a large-scale pre-trained dialogue generation model with discrete latent variable. These two models generate many response candidates for the following scoring modules¹.

Dialogue Planning Model (DPM). As shown in Figure 2, to improve the dialogue coherence and topic depth, we

¹We use the public released Plato model (<https://github.com/PaddlePaddle/Knover/tree/master/plato-2>).

Models	Meteor	Bert-score	USR	Human Rating
Our model	0.142	0.869	-	-
+ bert-score ensemble	0.147	0.876*	4.34	4.13
+ meteor ensemble	0.160*	0.873	4.51	4.21
Sub-task 2 system	0.070	0.843	3.86	4.28*

Table 2: Automatic and human evaluation results on the test set provided by the organizers in DSTC9 Interactive Dialogue Evaluation Track sub-task 1. Note that * denotes that it ranks 1st over all submissions in the competition.

Models	FED	Human Rating
<i>Baseline</i>		
Transformer	3.69	3.60
DialoGPT	4.72	3.87
<i>Our system</i>		
Our System	4.61	4.08

Table 3: Automatic and human evaluation results on interactive dialogues provided by the organizers in DSTC9 Interactive Dialogue Evaluation Track sub-task 2.

design the Dialogue Planning Model based on the gpt2 model. Particularly, we design a dialogue flow transformer block (FLOW) upon the GPT2 model. Formally, given a dialogue containing n utterances $u = [u_1, u_2, \dots, u_n]$, suppose $U_{1\dots n}$ denote the representation of $1 \sim n$ utterance encoded by GPT2 model. We consider dialogue process as many vector operations as shown in Figure 2.

$$\Delta U_n = U'_{1\dots n} - U_{1\dots n-1} \quad (3)$$

where ΔU_n can be considered as the representation of utterance u_n conditioned on $u_{<n}$ and $U'_{1\dots n}$ is the predicted representation by the dialogue flow transformer block.

To train the Dialogue Planning Model, we design three tasks: Dialogue Flow Prediction, Response Generation, Bag-of-Words Prediction.

Dialogue Flow Prediction is to predict the representation of $1 \sim n$ utterances $U'_{1\dots n}$ based on $U_1, U_{1\dots 2}, \dots, U_{1\dots n-1}$.

$$U'_{1\dots n} = FLOW(U_1, U_{1\dots 2}, \dots, U_{1\dots n-1}) \quad (4)$$

We train the Dialogue Flow Prediction task by minimizing Mean Squared Error:

$$\mathcal{L}_{flow} = MSELoss(U_{1\dots n}, U'_{1\dots n}) \quad (5)$$

Response Generation is to generate dialogue response using utterances $u_{<n}$ and predicted representation of utterance u_n . Specifically, when generating each token, we concatenate ΔU and the gpt2 output hidden states. We optimize Response Generation task by minimizing the following loss:

$$\mathcal{L}_{gen} = - \sum_{i=1}^N \log(P(u_n^i | u_{<n}, u_n^{<i}, \Delta U_n; \theta)) \quad (6)$$

Bag-of-Words Prediction task is to predict the words in an utterance using ΔU_n , which can be considered as a topical

constraint. This task can be optimized by minimizing the following loss:

$$\mathcal{L}_{bow} = - \sum_{i=1}^N \log(u_n^i | \Delta U_n) \quad (7)$$

The overall loss to train Dialogue Planning Model can be computed as follows:

$$\mathcal{L} = \mathcal{L}_{flow} + \mathcal{L}_{gen} + \mathcal{L}_{bow} \quad (8)$$

Score Model. Score Model consists of four scoring models: DialoRpt, Plato, NLI, Abusive Detection. DialoRPT (Gao et al. 2020) is a large-scale dialog ranking model based on the human feedback of dialogue responses. Plato denotes the response selection model in Bao et al. (2020). NLI means natural language inference, which is useful for non-consistency detection. We exploit RoBERTa-large-mnli (Liu et al. 2019) to predict whether the response conflicts with dialogue history. Abusive Detection is to detect some abusive words.

Plato always provides relatively close scores (e.g. 1e-4) for top 10 responses, while DialoRPT gives scores with a larger gap. In our interactive experiments, DialoRPT always choose some responses that are not so coherent with the context, which is may be due to that it is trained based on the human likert-score in the forum. Therefore, we first use remove the responses with abusive words, and exploit plato to filter relatively coherent response candidates and then use NLI and DialoRPT to choose the best one. The whole score process is shown in Algorithm 2.

Algorithm 2 Scoring method for interactive dialogue response selection.

- 1: **Input:** Response candidates $R = \{r_i | i = 1, 2, \dots, N\}$, dialogue history.
 - 2: **Output:** The most appropriate response r .
 - 3: Detect abusive words and remove the response with abusive words.
 - 4: Score by Plato and select top-10 responses.
 - 5: For each response, use the NLI model to detect conflict with dialogue history. Remove the response with conflict.
 - 6: Score by DialoRPT and select top-1 response.
 - 7: **Return:** $r_{argmax(M)}$.
-

Post-process. To respond with more human-like responses, after selecting the most appropriate response through the

scoring model, post-process mainly formats the response more human-like, such as uppercase problem on special entities.

Experiments and Analysis

Experiment Settings

For sub-task 1, we conduct experiments on the Topical-chat dataset (Gopalakrishnan et al. 2019), which contains dialogues with topical knowledge. The dataset contains 9058 dialogues for training, 565 dialogues for freq validation, and 565 dialogues for freq test. We only train our model on the training data without any other data and select the model based on the performance on the freq validation data. We employ 60 response candidates for the response ensemble.

For sub-task 2, we initialize the Dialogue Planning Model with GPT2-large² and fine tune it on the BST dataset, which contains four human annotated conversations datasets: ConvAI2 (Zhang et al. 2018; Dinan et al. 2020), Empathetic Dialogues (Rashkin et al. 2018), Wizard of Wikipedia (Dinan et al. 2018), Blended Skill Talk (Smith et al. 2020). We train the model on 8 Tesla V100 GPUs for about 24 hours.

Metrics

For sub-task 1, we employ the following metrics to evaluate our model.

Bert-score: A reference-based evaluation metric that uses a pre-trained BERT model to greedily match each word in the generated response with the ground-truth response (Zhang et al. 2019a).

Meteor: A reference-based evaluation metric which is designed as an improvement on BLEU (Papineni et al. 2002) using a harmonic mean of precision and recall (Banerjee and Lavie 2005).

USR: An unsupervised and reference-free evaluation metric for response evaluation (Mehri and Eskenazi 2020b).

Human Ratings: Human evaluation is carried out on Amazon Mechanical Turk with the annotation questionnaire used in FED score (Mehri and Eskenazi 2020a). There are 100 context-response pairs sampled and each one is labeled by 3 annotators.

For sub-task 2, we submit our system on DialPort and collect dialogs through conversations with real users. The organizers mainly use human evaluation as well as the FED score to evaluate our dialogue system.

FED: A reference-free evaluation metric (Mehri and Eskenazi 2020a) which is designed to evaluate the interactive dialogue with real users.

Human Ratings: The human evaluation is carried out on Amazon Mechanical Turk with the annotation questionnaire used in the FED score (Mehri and Eskenazi 2020a). There are 200 dialogs evaluated for each system.

Experiment Results

Sub-task 1: The evaluation results on test set for sub-task 1 is shown in Table 2. Our model with bert-score ensemble

reaches the best bert-score and with meteor ensemble gets the best meteor score over all submissions. We also use the interactive system in sub-task 2 for this task, which tie 1st on human ratings.

Sub-task 2: As shown in Table 3, on human ratings, our system significantly outperforms baseline system Transformer and DialoGPT, which is provided by the organizer. However, on the FED score, our system is a little lower than DialoGPT. We consider that the FED score is based on DialoGPT, so the DialoGPT system gets a better FED score.

Conclusion

In this paper, we introduce the WeChat AI’s submission for DSTC9 Interactive Dialogue Evaluation Track sub-task 1 and sub-task 2. In sub-task 1, our model is based on GPT2 and we propose a simple but efficient ensemble method for knowledge-grounded dialogue. Our method achieved the highest Bert-score, Meteor, and human ratings using different systems in the competition. In sub-task 2, to improve topic depth and dialogue coherence, we propose the Dialogue Flow Model and we build an integrated open-domain dialogue system containing four modules: pre-process, dialogue model, scoring model, and post-process for generating more human-like responses. Our system significantly outperforms the baseline methods and ranks 3rd over all submissions.

Acknowledgement

We sincerely thank the anonymous reviewers for their thorough reviewing and valuable suggestions. This work is supported by National Key R&D Program of China (NO. 2018YFC0825201 and NO. 2017YFE0192900).

References

- Adiwardana, D.; Luong, M.-T.; So, D. R.; Hall, J.; Fiedel, N.; Thoppilan, R.; Yang, Z.; Kulshreshtha, A.; Nemade, G.; Lu, Y.; et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- Banerjee, S.; and Lavie, A. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 65–72.
- Bao, S.; He, H.; Wang, F.; Wu, H.; Wang, H.; Wu, W.; Guo, Z.; Liu, Z.; and Xu, X. 2020. Plato-2: Towards building an open-domain chatbot via curriculum learning. *arXiv preprint arXiv:2006.16779*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dinan, E.; Logacheva, V.; Malykh, V.; Miller, A.; Shuster, K.; Urbanek, J.; Kiela, D.; Szlam, A.; Serban, I.; Lowe, R.; et al. 2020. The second conversational intelligence challenge (convai2). In *The NeurIPS’18 Competition*, 187–208. Springer.

²<https://github.com/huggingface/transformers>

- Dinan, E.; Roller, S.; Shuster, K.; Fan, A.; Auli, M.; and Weston, J. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. *arXiv preprint arXiv:1811.01241* .
- Gao, X.; Zhang, Y.; Galley, M.; Brockett, C.; and Dolan, B. 2020. Dialogue Response Ranking Training with Large-Scale Human Feedback Data. In *EMNLP*.
- Gopalakrishnan, K.; Hedayatnia, B.; Chen, Q.; Gottardi, A.; Kwatra, S.; Venkatesh, A.; Gabriel, R.; Hakkani-Tür, D.; and AI, A. A. 2019. Topical-Chat: Towards Knowledge-Grounded Open-Domain Conversations. In *INTERSPEECH*, 1891–1895.
- Gunasekara, C.; Kim, S.; D’Haro, L. F.; Rastogi, A.; Chen, Y.-N.; Eric, M.; Hedayatnia, B.; Gopalakrishnan, K.; Liu, Y.; Huang, C.-W.; Hakkani-Tür, D.; Li, J.; Zhu, Q.; Luo, L.; Liden, L.; Huang, K.; Shayandeh, S.; Liang, R.; Peng, B.; Zhang, Z.; Shukla, S.; Huang, M.; Gao, J.; Mehri, S.; Feng, Y.; Gordon, C.; Alavi, S. H.; Traum, D.; Eskenazi, M.; Beirami, A.; Eunjoon, Cho; Crook, P. A.; De, A.; Geramifard, A.; Kottur, S.; Moon, S.; Poddar, S.; and Subba, R. 2020. Overview of the Ninth Dialog System Technology Challenge: DSTC9.
- Holtzman, A.; Buys, J.; Du, L.; Forbes, M.; and Choi, Y. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751* .
- Huber, B.; McDuff, D.; Brockett, C.; Galley, M.; and Dolan, B. 2018. Emotional dialogue generation using image-grounded language models. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, 1–12.
- Kim, B.; Ahn, J.; and Kim, G. 2020. Sequential Latent Knowledge Selection for Knowledge-Grounded Dialogue. *arXiv preprint arXiv:2002.07510* .
- Li, Z.; Li, Z.; Zhang, J.; Feng, Y.; Niu, C.; and Zhou, J. 2020. Bridging Text and Video: A Universal Multimodal Transformer for Video-Audio Scene-Aware Dialog. *arXiv preprint arXiv:2002.00163* .
- Li, Z.; Niu, C.; Meng, F.; Feng, Y.; Li, Q.; and Zhou, J. 2019. Incremental transformer with deliberation decoder for document grounded conversations. *arXiv preprint arXiv:1907.08854* .
- Lian, R.; Xie, M.; Wang, F.; Peng, J.; and Wu, H. 2019. Learning to select knowledge for response generation in dialog systems. *arXiv preprint arXiv:1902.04911* .
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692* .
- Mehri, S.; and Eskenazi, M. 2020a. Unsupervised evaluation of interactive dialog with dialogpt. *arXiv preprint arXiv:2006.12719* .
- Mehri, S.; and Eskenazi, M. 2020b. USR: An Unsupervised and Reference Free Evaluation Metric for Dialog Generation. *arXiv preprint arXiv:2005.00456* .
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language models are unsupervised multitask learners .
- Rashkin, H.; Smith, E. M.; Li, M.; and Boureau, Y.-L. 2018. Towards empathetic open-domain conversation models: A new benchmark and dataset. *arXiv preprint arXiv:1811.00207* .
- Smith, E. M.; Williamson, M.; Shuster, K.; Weston, J.; and Boureau, Y.-L. 2020. Can You Put it All Together: Evaluating Conversational Agents’ Ability to Blend Skills. *arXiv preprint arXiv:2004.08449* .
- Tuan, Y.-L.; Chen, Y.-N.; and Lee, H.-y. 2019. DyKgChat: Benchmarking dialogue generation grounding on dynamic knowledge graphs. *arXiv preprint arXiv:1910.00610* .
- Zhang, S.; Dinan, E.; Urbanek, J.; Szlam, A.; Kiela, D.; and Weston, J. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243* .
- Zhang, T.; Kishore, V.; Wu, F.; Weinberger, K. Q.; and Artzi, Y. 2019a. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675* .
- Zhang, Y.; Sun, S.; Galley, M.; Chen, Y.-C.; Brockett, C.; Gao, X.; Gao, J.; Liu, J.; and Dolan, B. 2019b. Dialogpt: Large-scale generative pre-training for conversational response generation. *arXiv preprint arXiv:1911.00536* .
- Zhao, T.; Lee, K.; and Eskenazi, M. 2016. Dialport: Connecting the spoken dialog research community to real user data. In *2016 IEEE Spoken Language Technology Workshop (SLT)*, 83–90. IEEE.
- Zhou, H.; Young, T.; Huang, M.; Zhao, H.; Xu, J.; and Zhu, X. 2018. Commonsense Knowledge Aware Conversation Generation with Graph Attention. In *IJCAI*, 4623–4629.