# MSDF: A General Open-Domain Multi-skill Dialog Framework

Yu Zhao, Xinshuo Hu, Yunxin Li, Baotian Hu$^{(\boxtimes)}$, Dongfang Li, Sichao Chen, and Xiaolong Wang

Harbin Institute of Technology, Shenzhen, China
`hubaotian@hit.edu.cn, xlwangsz@hit.edu.cn`

**Abstract.** Dialog systems have achieved significant progress and have been widely used in various scenarios. The previous researches mainly focused on designing dialog generation models in a single scenario, while comprehensive abilities are required to handle tasks under various scenarios in the real world. In this paper, we propose a general **Multi-Skill Dialog Framework**, namely **MSDF**, which can be applied in different dialog tasks (e.g. knowledge grounded dialog and persona based dialog). Specifically, we propose a transferable *response generator* pre-trained on diverse large-scale dialog corpora as the backbone of MSDF, consisting of BERT-based encoders and a GPT-based decoder. To select the response consistent with dialog history, we propose a *consistency selector* trained through negative sampling. Moreover, the flexible copy mechanism of external knowledge is also employed to enhance the utilization of multiform knowledge in various scenarios. We conduct experiments on knowledge grounded dialog, recommendation dialog, and persona based dialog tasks. The experimental results indicate that our MSDF outperforms the baseline models with a large margin. In the Multi-skill Dialog of 2021 Language and Intelligence Challenge, our general MSDF won the 3rd prize, which proves our MSDF is effective and competitive.
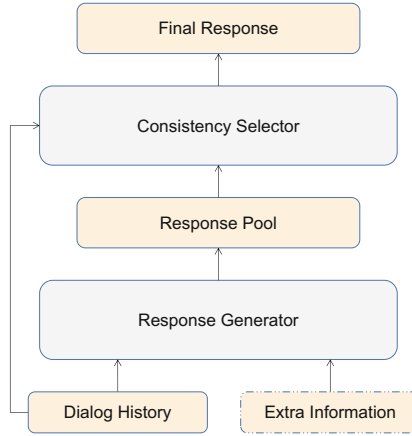
**Keywords:** Multi-skill dialog · Knowledge grounded dialog · Conversational recommendation · Persona based dialog

## 1 Introduction

Propelled by the acquisition of large-scale dialog corpora and the advance of pre-training technology, dialog generation models have made great progress, and dialog systems have been applied in various scenarios, such as chit-chat, knowledge grounded dialog, and conversational recommendation. However, most of the previous works only focused on modeling the dialog system within a single scenario, which is difficult to be applied in other scenarios directly. It does not

---

Y. Zhao and X. Hu—Contribute equally to this work.

**Fig. 1.** Overall Multi-Skill Dialog Framework (MSDF). The response generator first generates diverse responses through a sampling-based decoding algorithm conditioned on the dialog history and optional extra information (e.g. knowledge, user profile and/or machine persona). Then, the consistency selector selects the most contextually consistent response as the final response.

satisfy the requirements of practical application in the real world, where the dialog model needs to generate responses in various scenarios. To handle different tasks in various scenarios, multiple dialog skills are requested, such as knowledge utilization, commodities recommendation, and persona understanding skills. It is of great necessity to model a general multi-skill dialog framework that can be flexibly applied in various scenarios.

How to model the general multi-skill dialog systems that can effectively use information from diverse sources still remains challenging. On the one hand, the model needs to use various information (e.g. structured and unstructured knowledge, persona, conversation topics), and the previous works usually design complex models [6,11] to utilize specific information in a single scenario, which results in lacking universality. On the other hand, the previous works used complicated data processes and training processes [9] to optimize models on specific dialog corpus, thus, the model is difficultly transferred to other scenarios.

In this work, we propose the general multi-skill dialog framework **MSDF** to address the above problems, which consists of a pre-trained dialog *response generator* and a *consistency selector*. As depicted in Fig. 1, MSDF generates responses in two stages: 1) generating diverse responses as the candidate *response pool*, via *dialog history* (and *extra information*); 2) selecting the *final response* by consistency selector. Specifically, we first pre-train a universal and transferable encoder-decoder based model on various diverse dialog corpora, including chitchat, knowledge dialog, and recommendation dialog, to obtain a general model with multiple coarse-grained skills enhanced. Then, we apply multiple identical encoders to encode different source information and equip the decoder with

the multi-source information fusion module, which can be flexibly transferred to different application scenarios. Moreover, a dialog history can be mapped into multiple acceptable responses, which is also known as the one-to-many mapping [13], especially in the open domain. Thus, we introduce a BERT-based consistency selector to choose the most contextually consistent response with dialog history from the response pool, which is trained via negative sampling to distinguish consistent and inconsistent responses with dialog history.

Experiments are conducted to evaluate the performance of our MSDF in knowledge grounded dialog, recommendation dialog, and persona dialog tasks. The experimental results indicate that our MSDF outperforms the baseline models with a large margin in terms of F1 and BLEU scores. In the Multi-skill Dialog of 2021 Language and Intelligence Challenge, our MSDF won the third prize, with 6th rank in the human evaluation of the finals. Both automatic and human evaluation results indicate that our MSDF is effective and competitive.

Our contributions are as bellows.

– We propose a general multi-skill dialog framework that solves various tasks in different scenarios, namely MSDF, consisting of a dialog *response generator* with multi-source information encoders and a *consistency selector*.
– We pre-train an encoder-decoder based dialog generation model on various types of large-scale open-domain dialog corpora, which can be effectively transferred into our MSDF to solve different tasks.
– The experimental results indicate that our MSDF outperforms the baseline models with a large margin in knowledge dialog, recommendation dialog, and persona dialog tasks, demonstrating the effectiveness of MSDF.
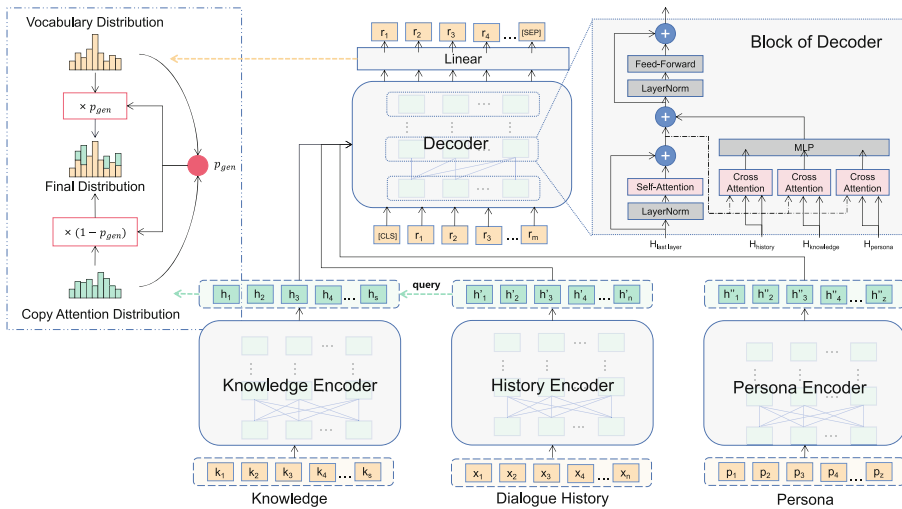
## 2   Related Work

The Multi-Skill Dialog of 2021 Language and Intelligence Challenge focuses on three kinds of dialog generation tasks, including the knowledge grounded dialog, recommendation dialog, and persona dialog. The previous works always modeled dialog systems to solve a single task in a scenario. For knowledge grounded dialogs, the model generates responses conditioned on dialog history and the given external knowledge. [8] employed posterior probability distribution of knowledge to guide the learning of prior distribution during training. [6] proposed the sequential decision making method to select knowledge used in response. And [12] modeled the knowledge selection as walking over knowledge graphs. For the persona dialogs, the model is required to understand and utilize the given persona to generate personalized responses. [9] improved the performance of persona understanding through persona negative sampling. [22] trained a VAE (Variational Auto-encoder) model to produce persona-related topic words jointly generating responses. For recommendation dialogs, the model needs to make recommendations through conversations, which is usually separated to two tasks: recommendation and dialog generation. [4,27] incorporated the common sense knowledge graph to improve the user profile understanding and recommendation dialog generation. [28] incorporated topic planning to enhance the

recommendation dialog. Following the work of [3], we also propose a general multi-skill dialog framework that can handle various tasks.

# 3 Method

## 3.1 Multi-skill Dialog

The multi-skill dialog task aims to construct a general dialog generation model with multiple skills to solve various tasks in various scenarios. In our work, we focus on the dialog modeling task with multiple skills, including three sub-tasks: 1) knowledge dialog, 2) recommendation dialog, and 3) persona dialog. The descriptions of the three sub-tasks are as follows[1] (Fig. 2).



**Fig. 2.** Responses generation process in MSDF. The top-right corner is the detail of decoder block, and the top-left corner shows the process of copy mechanism.

**Knowledge Dialog.** The inputs of the knowledge dialog include dialog goals $\mathbf{G} = [\mathbf{g}_1, \mathbf{g}_2, \ldots, \mathbf{g}_r]$, knowledge information $\mathbf{K} = [\mathbf{k}_1, \mathbf{k}_2, \ldots, \mathbf{k}_s]$, and dialog history $\mathbf{H} = [\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n]$, where $\mathbf{g}_i$ is a description of dialog target or a conversation topic, $\mathbf{k}_i$ is the form of a triple or an unstructured sentence, and $\mathbf{x}_i$ consists of a sequence of words. The knowledge grounded dialog generation task requires the model to use appropriate knowledge to generate the response $\mathbf{R} = [r_1, r_2, \ldots, r_m]$, where $r_i$ denotes the $i$-th word in response, $m$ denotes the length of response.

---

[1] Details in: https://aistudio.baidu.com/aistudio/competition/detail/67.

**Recommendation Dialog.** The inputs of the recommendation dialog include user profiles $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \ldots, \mathbf{u}_t]$, conversation scenario $\mathbf{S}$, dialog targets $\mathbf{G}$, knowledge information $\mathbf{K}$, and dialog history $\mathbf{H}$, where $\mathbf{u}_i$ is an aspect of the user profile in a key-value form. The recommendation dialog generation task requires the model to make recommendations based on the user profile and the scenario through conversation.

**Persona Dialog.** The inputs of the persona dialog include persona information of machine $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \ldots, \mathbf{p}_z]$, and dialog history $\mathbf{H}$, where $\mathbf{p}_i$ is a sentence describing the persona. The persona dialog generation task requires the model to generate natural, fluent and informative machine response in line with dialogue history and persona information.

### 3.2   Overall Framework

We propose a general multi-skill dialog framework MSDF, mainly consisting of a response generator and a history-response consistency selector, as shown in Fig. 1. MSDF generates the response in two steps: response generator is applied to generate diverse responses as the response pool, and then consistency selector chooses the most consistent response from the response pool as the final output.

**Response Generator.** The previously proposed generation models, such as GPT2 [14], PLATO2 [1], and BART [7], have shown their great generation performance obtaining from large-scale corpora. To obtain a universal pre-trained dialog model that can be flexibly transferred to the three sub-tasks, we pre-train an encoder-decoder-based model, which consists of a BERT-based encoder and a GPT-based decoder, and the optimization objective is to generate a target response conditioned on dialog history. The six large-scale dialog datasets, including LCCC [19], Weibo [17], Douban [21], DuConv [20], KdConv [26] and DuRecDial [11], are used to pre-train the model. Specifically, we process all the data to a general dialog form: for the multi-turn conversations, we split them into history-response pairs; for the conversations with extra knowledge or user profiles, we ignore them during pre-training. Moreover, to accelerate the pre-training process, we use the parameters of DialoGPT [24] to initialize both the encoder and decoder. We find that the model pre-trained with all data tends to generate short responses, thus, we additionally select the half of data with longer responses to pre-train it in the latter part of pre-training phase. In each sub-task, we duplicate the encoder of pre-trained model to encode specific information, and equip the decoder with an information fusion module.

Denote $\mathbf{H}_{l-1}^{D} \in \mathbb{R}^{L_d \times d}$ as the output of $(l-1)$-th decoder block, where $L_D$ is the length of decoding sequence and $d$ is the hidden size of decoder blocks, the calculating process in each decoder block is described as follows:

$$\mathbf{H}_l^{SA} = SelfAttention\left(LayerNorm(\mathbf{H}_{l-1}^{D})\right) + \mathbf{H}_{l-1}^{D}, \qquad (1)$$

$$\mathbf{H}_l^{FA} = Fusion\left(\mathbf{H}_l^{SA}, \mathbf{H}_{history}, \mathbf{H}_{knowledge}, \mathbf{H}_{persona}\right), \qquad (2)$$

$$\mathbf{H}_l^D = FFN\left(LayerNorm(\mathbf{H}_l^{FA})\right) + \mathbf{H}_l^{FA}, \tag{3}$$

where $\mathbf{H}_{history} \in \mathbb{R}^{n \times d}$, $\mathbf{H}_{knowledge} \in \mathbb{R}^{s \times d}$ and $\mathbf{H}_{persona} \in \mathbb{R}^{z \times d}$ denotes the encoded information, $\mathbf{H}_l^{SA} \in \mathbb{R}^{L_d \times d}$ denotes the result of self-attention, and $\mathbf{H}_l^{FA} \in \mathbb{R}^{L_d \times d}$ denotes the fused multi-source information. The $Fusion$ function is calculated in the following way:

$$\mathbf{A}_l(\mathbf{H}^E) = softmax\left(\frac{\left(LayerNorm(\mathbf{H}_l^{SA})\mathbf{W}_l^Q\right)(\mathbf{H}^E\mathbf{W}_l^K)^T}{\sqrt{d}}\right)(\mathbf{H}^E\mathbf{W}_l^V), \tag{4}$$

$$\mathbf{H}_l^{FA} = [\mathbf{A}_l(\mathbf{H}_{history}); \mathbf{A}_l(\mathbf{H}_{knowledge}); \mathbf{A}_l(\mathbf{H}_{persona})]\mathbf{W}_l^P + \mathbf{H}_l^{SA}, \tag{5}$$

where [;] indicates concatenation on the last dimension, $\mathbf{W}_l^Q$, $\mathbf{W}_l^K$ and $\mathbf{W}_l^V$ are learnable projection matrices and $\mathbf{W}_l^P \in \mathbb{R}^{3d \times d}$ is a learnable attention fusing matrix (we get $\mathbf{W}_l^P \in \mathbb{R}^{2d \times d}$ when there are only two cross-attention to be fused).

Finally, the hidden state $\mathbf{H}_l^D \in \mathbb{R}^{L_D \times d}$ of the $l$-th decoder layer is output by the $FFN$ layer and residual layer:

$$\mathbf{H}_l^D = FFN(LayerNorm(\mathbf{H}_l^{FA})) + \mathbf{H}_l^{FA}. \tag{6}$$

With the great performance of hybrid pointer-generator network [15] in summarization, we utilize the attention-based copy mechanism to generate knowledge-enhanced response. Decoder hidden states are put into linear language model to get original vocabulary distribution $\mathbf{P}_{vocab} \in \mathbb{R}^{L_D \times L_V}$,

$$\mathbf{P}_{vocab} = softmax(\mathbf{H}^D\mathbf{W}^{LM}), \tag{7}$$

where $W^{LM} \in \mathbb{R}^{d \times L_V}$ is the learnable language model head and $L_V$ denotes the vocabulary length. And then attention-based copy mechanism is utilized to generate extra knowledge enhanced response. We obtain knowledge copy attention $\mathbf{A}_{copy} \in \mathbb{R}^{L_D \times L_K}$, via cross-attention of decoded hidden states and encoded knowledge (or persona) hidden state:

$$\mathbf{A}_{copy} = softmax((\mathbf{H}^D\mathbf{W}^Q)(\mathbf{H}_{knowledge}\mathbf{W}^K)^T), \tag{8}$$

where $\mathbf{W}^Q$ and $\mathbf{W}^K$ are learnable projection matrices, and $L_K$ denotes the length of the knowledge sequence. Then generation probability $\mathbf{p}_{gen} \in [0, 1]$ can be calculated:

$$\mathbf{p}_{gen} = sigmoid([\mathbf{A}_{copy}\mathbf{H}_{knowledge}; \mathbf{H}^D]\mathbf{W}^{mlp}), \tag{9}$$

where $W^{mlp} \in R^{2d \times 1}$ is learnable matrices. We obtain the following probability distribution over the merged vocabulary to predict word w:

$$\mathbf{P}(w) = \log(\mathbf{p}_{gen}\mathbf{P}_{vocab}(w) + (1 - \mathbf{p}_{gen})\mathbb{I}(w_i = w)A_{copy}(w)), \tag{10}$$

where $\mathbb{I}(\cdot)$ is an indicator function, and $A^{copy}(w)$ is the element at $w$-th column in the copy attention matrix.

**Consistency Selector.** Inspired by the BERT next sentence prediction [5], we propose a history-response consistency selector, since the response should be the next sentence by dialog history. The consistency selector is a binary classifier to distinguish the consistent and inconsistent response with dialog history. We construct the consistency selector by a pre-trained model RoBBERTa [10], plus a linear head layer. We first concatenate the dialog history and response as the pre-trained model input to get context representation, and then put the first token (usually known as [CLS] token) to the linear head to get a binary classification score. The consistency selector is trained on positive and negative sampling examples from training data, where inconsistent responses are randomly sampled. During inference, we get the positive score from the consistency selector output as the consistency score.

### 3.3   Data Processing

We separate different input resources into four categories: dialog history, knowledge information, persona information, and current reply (or previously generated response during inference). The data preprocessing and reprocessing are demonstrated as follows.

The knowledge dialog generation model consists of a dialog history encoder, a knowledge encoder and a decoder. During data preprocessing of DuConv, we reformat the knowledge graph to a pseudo unstructured knowledge sentences, by concatenating the subject, predicate and object, and join all knowledge sentence with special token "[SEP]". Since the dialog topics are limited in movies and film stars and there are at most two topics during a dialog, we introduce special tokens to format the subject in the knowledge graph, including "[movie1]", "[movie2]", "[star1]", and "[star2]", which will be restored in data reprocessing after response generated. And we also add speaker tokens "[speaker1]" and "[speaker2]" to distinguish speakers in the dialog history. We also reprocess personal pronouns and figures due to the knowledge, such as outcome date of movies, and birthday of stars.

The recommendation dialog generation model consists of a dialog history encoder, a knowledge encoder, a persona encoder and a decoder. To facilitate optimization, we make the knowledge encoder and persona encoder share parameters. The knowledge graph in DuRecDial is also reformatted in the same way as DuConv, with all the candidate recommendation goal planning concatenated at the end. Since the recommendation goal is labeled with each response, we introduce a new special token "[goal]" as a separator, and join the golden goal with response e.g. "[goal]问User性别[goal]我该称呼您是先生还是女士" ("[goal] Asking about the user's gender [goal] Should I call you Mr. or MS"). This golden goal prefix in the response performs as the semantic guidance, which is typically like conditional generation. The dialog situation is viewed as extra knowledge concatenated before knowledge information. We also replace the user name with a new special token "[uname]" and add speaker tokens. Personal pronouns and figures will be corrected and goal information will be removed during data reprocess.

The persona dialog generation model consists of a dialog history encoder, a persona encoder and a decoder. We trained a simple Word2Vec by open-source gensim implementation[2], on CPC sentences to estimate the similarity between the response and persona by cosine similarity of average word vector, and randomly drop training examples with persona information similarity less than 0.7. It is observed that the generator prefers to generate longer sentences and there are too many responses whose topic is about the user job, so we drop some examples with too long responses and "工作" ("work") in dialog utterances. Considering that conversations in CPC are more like chit-chat than task-oriented dialog, we abandon the copy mechanism in this generator.

## 4   Experiment

### 4.1   Experimental Settings

**Dataset.** For pre-training, we use the Weibo dataset [17], Douban multi-round conversation [21], LCCC dataset [19], Emotional Conversational Dataset [25], retrieval-assisted conversational dataset [2], and Kdconv dataset [26]. For fine-tuning, we use DuConv [20] for knowledge grounded dialog, DuRecDial [11] for the conversational recommendation, and Chinese persona chat (CPC) for persona dialog. Statistics of all dataset are summarized in Table 1.

**Table 1.** Statistics of all datasets. [†] denotes the datasets for fine-tuning, and the others are for pre-training.

|  | Train | Dev | Test |
|---|---|---|---|
| DuConv[†] [20] | 19858 | 2000 | 5000 |
| DuRecDial[†] [11] | 6618 | 946 | 4645 |
| Chinese Persona Chat[†] | 23000 | 1500 | 3000 |
| Weibo dataset [17] | 3103764 | 443394 | 886790 |
| Douban multi-round conversation [21] | 5000000 | 25001 | 1186 |
| LCCC dataset [19] | 11987759 | 20000 | 10000 |
| Emotional Conversational Dataset [25] | 899207 | 110000 | 110000 |
| retrieval-assisted conversational dataset [2] | 5498480 | 107332 | 156706 |
| Kdconv dataset [26] | 3000 | 300 | 2751 |

**Implementation Details.** For the response generator, the encoders and decoder settings follow the DialoGPT [24], where the hidden size is 768 for 12 layers, the maximum input length is 512, and there are 12 heads in multi-head attention. For all dropout layers, the dropout rate is set to 0.1. We adopt

---

cross-entropy loss as our loss function, and the parameters are saved in term of the minimum cross-entropy loss on the development datasets. The consistency selector is implemented alike a NLI model as described in Sect. 3.2, which is also optimized through cross-entropy loss function.

## 4.2 Evaluation

The competitive baseline models are compared with the proposed MSDF on different quantitative metrics. We use BLEU-1 and BLEU-2 to evaluate the n-gram lexical similarity, F1 to evaluate the Chinese character level similarity, and DISTINCT1 and DISTINCT-2 to evaluate the diversity of generated responses. The total SCORE for the automatic evaluation is calculated by averaging all the F1/BLEU1/BLEU2 scores for the subtasks. To exhibit the improvement of our MSDF, we also implement and test four baseline models: Seq2Seq [18], HRED [16], DialoGPT [24] and BERT-GPT [23]. All these baseline models only take dialog history as inputs, without aquiring extra information. Since test datasets are not released, we evaluate all the models on LUGE platform[3]. The performance of baseline models and our MSDF is presented in Table 2, including ablation experiments.

With respect to the human evaluation, we refer readers to the competition leaderboard[4] for the single-turn or multi-turn dialog evaluation results.

**Table 2.** Automatic evaluation of test set B on baseline models and our MSDF (ranked by total score). All the results of F1, BLEU, and DISTINCT are average scores from three sub-tasks.

|  | SCORE | F1 | BLEU1/2 | DISTINCT1/2 |
|---|---|---|---|---|
| Seq2Seq [18] | 0.522 | 22.33 | 0.202/0.096 | 0.038/0.100 |
| HRED [16] | 0.565 | 23.66 | 0.220/0.108 | 0.038/0.105 |
| DialoGPT [24] | 0.373 | 17.05 | 0.141/0.061 | **0.079/0.313** |
| BERT-GPT [23] | 0.573 | 24.51 | 0.215/0.113 | 0.061/0.214 |
| MSDF | **0.934** | **38.62** | **0.333/0.215** | 0.057/0.183 |
| -consistency selector | 0.872 | 36.00 | 0.314/0.198 | 0.051/0.173 |
| -extra knowledge | 0.705 | 28.98 | 0.271/0.145 | 0.049/0.162 |

## 4.3 Discussion

Our MSDF outperforms the baseline models with a large margin, even though without extra information, which strongly presents the effectiveness. The consistency selector preferred to select the most common response. It reduces the

---

variance of generating performance and significantly improves automatic evaluation results, without considering the limited diversity of responses in human evaluation. Besides, the attention-based copy mechanism is of great importance for generating knowledge-enhanced response. According to our observation, DialoGPT and BERT-GPT still talk rubbish after fine-tuning, despite resulting in higher DISTINCT scores. In the competition, our MSDF got 9-th rank in automatic evaluation and 6-th rank in human evaluation. It increased by 3 ranks in human evaluation compared to the automatic evaluation. We attribute this to the multi-skill pre-training, from which our model could generate more human-like responses, in spite of the mismatch with the golden references in automatic evaluation.

## 5   Conclusion

This paper describes our general multi-skill dialog framework MSDF, consisting of a response generator and a dialog history consistency selector. We first pre-train the basic encoder-decoder on diverse datasets and then fine-tune it on the specific dataset to construct the response generator with a strong specific skill. We won the third prize in Multi-task Dialog of 2021 Language and Intelligence Challenge. Experiments are also conducted on several baseline models. The vast improvement over the baseline models indicates that our framework is effective and competitive. In future work, we will experiment with our MSDF on more tasks to evaluate the comprehensive skills of our framework and further improve its performance.

## References

1. Bao, S., et al.: Plato-2: towards building an open-domain chatbot via curriculum learning. arXiv preprint arXiv:2006.16779 (2020)
2. Cai, D., Wang, Y., Bi, W., Tu, Z., Liu, X., Shi, S.: Retrieval-guided dialogue response generation via a matching-to-generation framework. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 1866–1875 (2019)
3. Cao, Y., Bi, W., Fang, M., Tao, D.: Pretrained language models for dialogue generation with multiple input sources. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings, pp. 909–917 (2020)

4. Chen, Q., et al.: Towards knowledge-based recommender dialog system. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 1803–1813 (2019)

5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186 (2019)

6. Kim, B., Ahn, J., Kim, G.: Sequential latent knowledge selection for knowledge-grounded dialogue. In: International Conference on Learning Representations (2019)

7. Lewis, M., et al.: Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 7871–7880 (2020)

8. Lian, R., Xie, M., Wang, F., Peng, J., Wu, H.: Learning to select knowledge for response generation in dialog systems. In: IJCAI International Joint Conference on Artificial Intelligence, p. 5081 (2019)

9. Liu, Q., et al.: You impress me: dialogue generation via mutual persona perception. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 1417–1427 (2020)

10. Liu, Y., et al.: Roberta: a robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692 (2019)

11. Liu, Z., Wang, H., Niu, Z.Y., Wu, H., Che, W., Liu, T.: Towards conversational recommendation over multi-type dialogs. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp. 1036–1049 (2020)

12. Moon, S., Shah, P., Kumar, A., Subba, R.: Opendialkg: explainable conversational reasoning with attention-based walks over knowledge graphs. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 845–854 (2019)

13. Ni, J., Young, T., Pandelea, V., Xue, F., Adiga, V., Cambria, E.: Recent advances in deep learning-based dialogue systems. arXiv preprint arXiv:2105.04387 (2021)

14. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language models are unsupervised multitask learners. OpenAI blog **1**(8), 9 (2019)

15. See, A., Liu, P.J., Manning, C.D.: Get to the point: summarization with pointer-generator networks. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1073–1083 (2017)

16. Serban, I.V., Sordoni, A., Bengio, Y., Courville, A., Pineau, J.: Building end-to-end dialogue systems using generative hierarchical neural network models. In: Thirtieth AAAI Conference on Artificial Intelligence (2016)

17. Shang, L., Lu, Z., Li, H.: Neural responding machine for short-text conversation. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 1577–1586 (2015)

18. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. Adv. Neural. Inf. Process. Syst. **27**, 3104–3112 (2014)

19. Wang, Y., Ke, P., Zheng, Y., Huang, K., Jiang, Y., Zhu, X., Huang, M.: A large-scale Chinese short-text conversation dataset. In: CCF International Conference on Natural Language Processing and Chinese Computing, pp. 91–103. Springer (2020)

20. Wu, W., Guo, Z., Zhou, X., Wu, H., Zhang, X., Lian, R., Wang, H.: Proactive human-machine conversation with explicit conversation goal. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 3794–3804 (2019)
21. Wu, Y., Wu, W., Xing, C., Zhou, M., Li, Z.: Sequential matching network: a new architecture for multi-turn response selection in retrieval-based chatbots. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 496–505 (2017)
22. Xu, M., et al.: A neural topical expansion framework for unstructured persona-oriented dialogue generation. arXiv preprint arXiv:2002.02153 (2020)
23. Zeng, G., et al.: Meddialog: a large-scale medical dialogue dataset. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 9241–9250 (2020)
24. Zhang, Y., et al.: Dialogpt: large-scale generative pre-training for conversational response generation. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pp. 270–278 (2020)
25. Zhou, H., Huang, M., Zhang, T., Zhu, X., Liu, B.: Emotional chatting machine: Emotional conversation generation with internal and external memory. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
26. Zhou, H., Zheng, C., Huang, K., Huang, M., Zhu, X.: Kdconv: A chinese multi-domain dialogue dataset towards multi-turn knowledge-driven conversation. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 7098–7108 (2020)
27. Zhou, K., Zhao, W.X., Bian, S., Zhou, Y., Wen, J.R., Yu, J.: Improving conversational recommender systems via knowledge graph based semantic fusion. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 1006–1014 (2020)
28. Zhou, K., Zhou, Y., Zhao, W.X., Wang, X., Wen, J.R.: Towards topic-guided conversational recommender system. In: Proceedings of the 28th International Conference on Computational Linguistics. pp. 4128–4139 (2020)