# Dialog-to-Actions: Building Task-Oriented Dialogue System via Action-Level Generation

Yuncheng Hua
Xiangyu Xi*
Meituan Group, China
devin.hua@monash.edu
xixy10@foxmail.com

Zheng Jiang
Meituan Group, China
zjiang@seu.edu.cn

Guanwei Zhang
Meituan Group, China
zhangguanwei@meituan.com

Chaobo Sun
Meituan Group, China
sunchaobo@meituan.com

Guanglu Wan
Meituan Group, China
wanguanglu@meituan.com

Wei Ye
Peking University, China, China
wye@pku.edu.cn

## ABSTRACT

End-to-end generation-based approaches have been investigated and applied in task-oriented dialogue systems. However, in industrial scenarios, existing methods face the bottlenecks of reliability (e.g., domain-inconsistent responses, repetition problem, etc) and efficiency (e.g., long computation time, etc). In this paper, we propose a task-oriented dialogue system via action-level generation. Specifically, we first construct dialogue actions from large-scale dialogues and represent each natural language (NL) response as a sequence of dialogue actions. Further, we train a Sequence-to-Sequence model which takes the dialogue history as the input and outputs a sequence of dialogue actions. The generated dialogue actions are transformed into verbal responses. Experimental results show that our light-weighted method achieves competitive performance, and has the advantage of reliability and efficiency.

## CCS CONCEPTS

• **Computing methodologies** → **Discourse, dialogue and pragmatics**.

## KEYWORDS

Task-Oriented Dialogue System, Action-Level Generation, Dialog-to-Actions

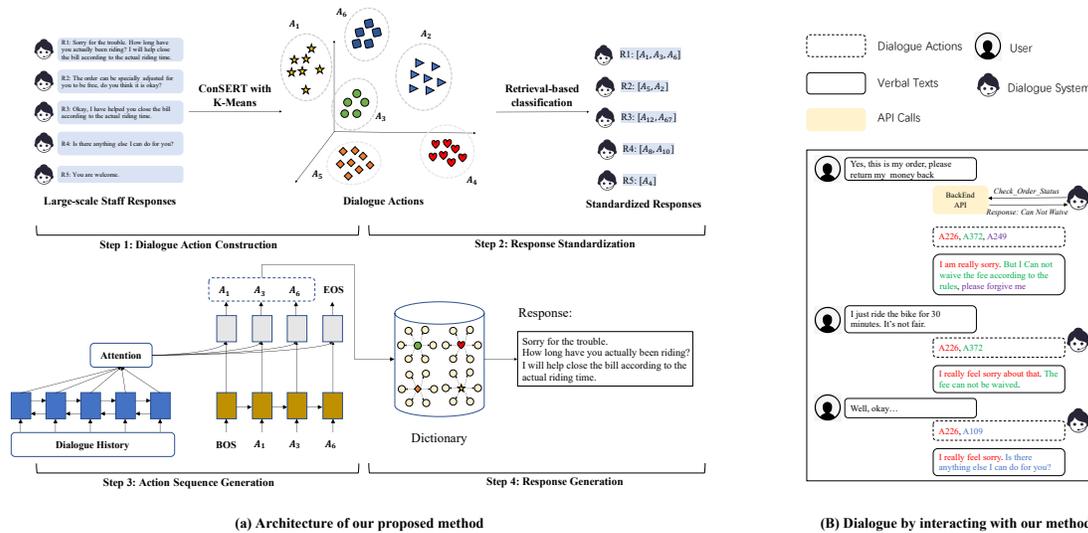*Yuncheng Hua and Xiangyu Xi contributed equally to this research.

## 1 INTRODUCTION

Recently, the end-to-end generation-based methods that directly output appropriate NL responses or API calls have been deeply investigated in task-oriented chatbots [2, 5, 9, 16, 22], and have been proven valuable for real-world business, especially after-sale customer services [1, 7, 8, 13, 14, 19, 21, 24, 25]. Based on the large-scale pre-trained language models [10, 11], generation-based methods have the advantage of simpler architecture and anthropomorphic interaction. Despite the significant progress, we find these token-level generation methods suffer from the following two limitations in practical scenarios.

**1. The token-level generation methods have limited reliability, which is essential for industrial task-oriented dialogue systems.** Due to the pre-trained language models' characteristics, the models may generate responses that are learned from the pre-training corpus. In certain cases, such responses are meaningless and not semantically incoherent with the current business domain, interrupting online interaction. Worse still, the models occasionally get stuck in generating repetitive responses across multiple turns (e.g., repeatedly enquiring the users for the same information). Above issues are also widely observed by other researchers [4, 12] and practitioners.[1]

**2. The token-level generation methods may fail to meet the efficiency requirement of the industrial systems, especially with large decoding steps.** The long computation time of the token-level generation models leads to unacceptable response latency of online dialogue systems, especially when the model generates a sentence of length that exceeds a threshold (e.g., 1,544 ms of T5 for a sentence of 30 words, as Figure 3 shows). Owing to the latency problem, a large number of service requests may be suspended or blocked during the peak period. Also, the computation resources (e.g., GPUs) required by the aforementioned systems might be unaffordable for small companies.

To address the above two problems, in this paper, we propose a task-oriented dialogue system based on the action-level generation method. Inspired by Xi et al. [19], we represent responses with **Dialogue Actions**, i.e., a class of the responses with unique and identical semantic meaning that can be automatically obtained by clustering. While Xi et al. [19] directly treats a whole response as a specific dialogue action, we split one response into multiple

---

[1]https://github.com/microsoft/DialoGPT/issues/45

**Figure 1: The system architecture and dialogue sample. In (b), the dialogue action and corresponding utterance segment are marked by the same color (e.g., "A226" and "I am really sorry").**

segments [6] and each segment can be mapped to a dialogue action. In this way, each response is represented as a sequence of dialogue actions. Given the dialogue context, a Seq2Seq model with an action-level recurrent decoder is used to generate the sequence of dialogue actions. Further, a frequency-based sampling method is used to compose the final response, based on the generated sequence of dialogue actions. Since the core component of our approach is the generation model which takes the *dialogue context* as the inputs and outputs *actions*, our method is named as D̲ialog-T̲o-A̲ctions (abbr. **DTA**). Compared with existing token-level generation-based systems, our DTA has the advantage of 1) reliability, since the generated natural language responses derive from the predefined dialogue actions; 2) efficiency, since the decoding space (i.e., dialogue actions) and the decoding steps are much smaller.

## 2 FRAMEWORK DESCRIPTION

### 2.1 Overview

We follow the workflow employed in the previous end-to-end task-oriented dialogue systems [2, 9], where the system takes the dialogue history as input, and generates a text string that either serves as a verbal staff response to the user or API calls (e.g., information inquiry, action execution, etc). When an API is invoked, the information returned from the API will be incorporated into the system's next response. A dialogue sample following such system interaction life cycle can be found in Figure 1 (b).

The key idea of our work is to generate dialogue actions and then compose a verbal response. To do so, we first construct dialogue actions from large-scale dialogues (Step 1) and represent each response as a sequence of dialogue actions (Step 2), as Figure 1 (a) shows. A Seq2Seq model with an action-level recurrent decoder is utilized to generate dialogue actions (Step 3), and the generated actions are further used to compose the verbal response (Step 4). We exemplify using the after-sale customer service of electric bike

rental business, where users and staffs communicate online through text messages. The technical details are introduced as follows.

### 2.2 Step 1: Dialogue Action Construction

A dialogue action refers to a cluster of utterances or utterance fragments that share identical semantic meaning and represent a common communicative intention, for instance making a request or querying information. Xi et al. [19] views a group of utterances with identical semantic information as dialogue action and selects a response corresponding to a specific staff action. However, the oversimplified setting, i.e., abstracting a whole utterance into an action, leads to relatively limited expressiveness and scalability. To make the responses more targeted and flexible, we construct dialogue actions based on utterance segments (of staff) rather than utterances. Specifically, each utterance is divided into multiple segments by a rule-based approach [6]. Further, following Xi et al. [19], we exploit a two-stage method to cluster the segments. Specifically, ConSERT [20] is utilized to generate representations for each utterance segment, and K-means is then applied to cluster the segments. We choose the number of clusters $K$ empirically to balance the purity and the number of the clusters, and treat each cluster of segments as a dialogue action (e.g., $A_1$ and $A_2$ in Figure 1 (a)).

### 2.3 Step 2: Response Standardization

Response standardization aims to standardize the responses (from the large-scale dialogues) by mapping each response to a sequence of dialogue actions. Following Yu et al. [23], we exploit a retrieval-based method, which retrieves clustered segments that are most similar to the given input utterance segment and label the input based on the corresponding clusters. As Figure 2 shows, given an input segment $x$, we use BM25 to recall top $k$ segments $\{u_1, ..., u_k\}$ from all clustered segments. Further, we exploit a BERT-based text similarity computation model $S$ to rerank the $k$ segments and select

the segment $\hat{u}$ with the highest similarity to $x$, denoting as:

$$\hat{u} = \underset{u_i \in \{u_1, ..., u_k\}}{\arg\max} \; S(x, u_i) \tag{1}$$

where $S(x, u_i)$ refers to the similarity between $x$ and $u_i$. $x$ is then annotated with the dialogue action $A_i$ that $\hat{u}$ belongs to.
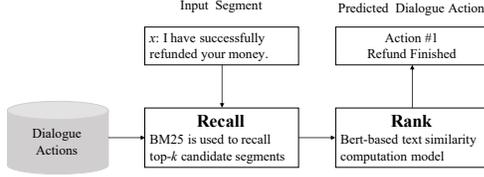


**Figure 2: The workflow of response standardization.**

Furthermore, we record the correspondence between the dialogue actions and utterance segments, as well as the frequencies of utterance segments in the dialogues. Specifically, we employ a key-value dictionary $\mathfrak{D}$, in which a key refers to a dialogue action $A_i$ while its value is a nested dictionary where the mapping relationship between the unique segments $\{x_1, ..., x_n\}$ and $A_i$ and the segments' occurrence frequencies $\{f_1, ..., f_n\}$ are recorded. The dictionary $\mathfrak{D}$ is used for composing the verbal response (Step 4).

## 2.4 Step 3: Action Sequence Prediction

Given a dialogue $\mathcal{D} = \{U_1, S_1, ..., U_T, S_T\}$ as a set of utterances exchanged between user ($U_i$) and staff ($S_i$) alternatively, through above carefully-designed steps, each staff utterance $S_i$ is represented as an action sequence $\mathcal{A}_{S_i}$. At the $m$-th turn, given a dialogue history $\mathcal{H}_m = \{U_{m-w}, S_{m-w}, ..., S_{m-1}, U_m\}$, we propose a Seq2Seq model to produce staff responses $S_m$. We first resort to the Seq2Seq model to output an action sequence of $\mathcal{A}_{S_m} = (A_1^{S_m}, A_2^{S_m}, ..., A_k^{S_m})$, where $k$ denotes the length of the action sequence, and then use the action sequence to form verbal response (in § 2.5).

**Encoder** Given the dialogue history $\mathcal{H}_m$, we sequentially concatenate all the utterances in $\mathcal{H}_m$ and each staff utterance's corresponding action sequence $\mathcal{A}_{S_i}$, forming a token sequence $(w_1, ..., w_n)$. The Bi-LSTM model is used to encode the token sequence into a sequence of continuous representations $H$:

$$H = (h_1, ..., h_n) = \text{BiLSTM}(w_1, ..., w_n) \tag{2}$$

**Decoder** Considering the efficiency requirement and small decoding space, we use the Luong attention method and employ LSTM as the decoder to calculate hidden state $s_t$ at time-step $t$ as follows:

$$s_t = \text{LSTM}(s_{t-1}, g(y_{t-1}), c_{t-1}) \tag{3}$$

where $y_{t-1}$ denotes the probability distribution over dialogue action space at step $t$-1 and $g(y_{t-1})$ denotes the action has the highest probability. After obtaining hidden state $s_t$ and context vector $c_t$, we generate probability distribution at time-step $t$ as follows:

$$y_t = \text{Softmax}(W_d[s_t; c_t]) \tag{4}$$

where $W_d$ is weight parameter. Given the ground-truth label $y_t$ at time-step $t$, we use $p(y_t|y_{<t}, \mathcal{H}_m)$ to denote the cross-entropy loss at step $t$ where $y_{<t}$ denotes the previously-generated actions.

The optimization objective is defined as:

$$L_{Gen} = -\sum_{\mathcal{D} \in C} \sum_{\mathcal{H}_m \in \mathcal{D}} \sum_{t=1}^{l_m} p(y_t|y_{<t}, \mathcal{H}_m) \tag{5}$$

where $C$ denotes the set of dialogues and $l_m$ denotes the length of the dialogue action sequence at the $m$-th turn of dialogue $\mathcal{D}$.

## 2.5 Step 4: Response Generation

Based on the action sequence generated in § 2.4, we compose the verbal response by selecting an utterance segment for each action and combining the segments sequentially. Considering the segments with higher frequencies are more likely to be the formal utterance that staff commonly use, we sample the segments from $\mathfrak{D}$ (built in § 2.3) following the principle that the higher the frequency, the more likely the segment to be selected. By doing this, we ensure the quality as well as the diversity of the verbal responses.

## 3 EXPERIMENTS

## 3.1 Experimental Settings

*3.1.1 Dataset.* We perform an experiment with a Chinese online after-sale customer service of electric bike rental business. In this scenario, the users may finish riding but forgot to lock the bike, and thus request the staff to remotely lock the bike and reduce the fees. The staff is required to judge whether the fee can be reduced by checking the status of the order via the back-end APIs. We collect the user-staff dialogues from the logs of online services for a week. The data statistics are shown in Table 1. The dialogues are randomly split into train, dev, and test sets with a ratio of 8:1:1. We construct 1,420 dialogue actions (§ 2.2), and each API call is treated as a dialogue action. The dataset will be released online.

| STAT TYPE | VALUE |
|---|---|
| Dialogs | 8,363 |
| Total turns | 55,576 |
| Avg. turns per dialog | 6.65 |
| Dialogue Actions | 1,420 |

**Table 1: Data statistics.**

*3.1.2 Baselines.* To evaluate the effectiveness and efficiency of our method, we compare it with the following state-of-the-art baselines: (1) **LSTM** which exploits a classical LSTM-based sequence-to-sequence architecture [15]; (2) **Transformer** which uses Transformer [17] as encoder and decoder; (3) **CDiaGPT** which is a GPT model pre-trained on a large-scale Chinese conversation dataset [18]; (4) **T5** which is a Text-to-Text Transfer Transformer (T5) model [11] pretrained with the CLUE Corpu.

*3.1.3 Evaluation Metrics.* To comprehensively evaluate the effectiveness of different models, we perform both offline evaluation (i.e., on the dataset) and online evaluation (i.e., online A/B testing).

**Offline Evaluation** The models take a specific ground-truth conversation history (i.e., context) as input and generate a response. Following Byrne et al. [2], we report the BLEU-4 score of each model in the test set. Considering the API calls are highly important, we observe the generated API calls in each turn and report the macro Precision (P), Recall (R), and F1-Score (F1) of API calls.

**Online Evaluation** Following Xi et al. [19], we deploy the models online and perform A/B testing. For each model, 120 dialogues are randomly sampled. The annotators who possess domain knowledge are required to perform the satisfaction assessment by grading each dialogue as "Low", "Medium", or "High" satisfaction degree. [2]

## 3.2 Main Results

The offline evaluation and online evaluation are shown in Table 3 and 2 respectively, from which we have the following observations: (1) Large-scale pre-trained language models significantly improve the performance of token-level generation models. For example, compared with the plain LSTM model, T5 achieves an absolute improvement of 26.20% for the BLEU-4 score and 4.98% for F1. (2) Compared with the CDiaGPT and the T5 models, our light-weighted DTA achieves competitive performance in the offline evaluation, and earns the highest satisfaction rating in the online evaluation, verifying the effectiveness of our proposed method.

| Model | Low | Medium | High |
|---|---|---|---|
| LSTM | 30.00 | 29.17 | 40.83 |
| Transformer | 27.50 | 28.33 | 44.17 |
| CDiaGPT | 12.50 | 13.33 | 74.17 |
| T5 | 5.83 | 15.83 | 78.33 |
| DTA | 8.33 | 12.50 | 79.17 |

**Table 2: Statistical Results of Online Evaluation (%)**

| Model | BLEU-4 | P | R | F1 |
|---|---|---|---|---|
| LSTM | 20.62 | 66.16 | 78.84 | 71.72 |
| Transformer | 26.21 | 62.59 | 75.90 | 64.74 |
| CDiaGPT | 42.54 | **71.18** | 86.43 | 77.11 |
| T5 | **46.83** | 70.91 | 87.25 | 76.70 |
| DTA | 44.82 | 68.03 | **90.80** | **77.74** |

**Table 3: Statistical Results of Offline Evaluation (%).**

## 3.3 In-Depth Analysis

*3.3.1 Effect of Efficiency Issue.* To investigate the effect of efficiency issue, we collect each model's computation time for processing the test samples under the same infrastructure (i.e., Tesla v100, 32GB RAM size, etc). Considering the computation time is highly correlated with the decoding steps, we first divide the generated responses into 10 subsets based on the response length, and then calculate the average computation time of each subset, as Figure 3 shows. We can observe that: (1) Despite the comparable performance, DTA has a significant advantage in computation efficiency over other models (e.g., 3.37 ms of DTA v.s. 1265.56 ms of CDiaGPT v.s. 2470.69 ms of T5 for subset "[50, 59]"). (2) DTA outperforms other models more significantly with longer responses. The reason is that the decoding steps of the token-level generation models are identical to the response length, while DTA performs action-level

decoding. The above observation verifies that our system provides an effective solution to build online dialogue services with limited computation resources.
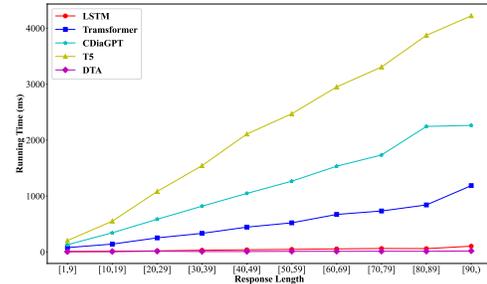


**Figure 3: The computation time of different models.**

*3.3.2 Effect of Reliability Issue.* We investigate the effect of reliability issue by quantitatively inspecting the repetition problem. Specifically, we calculate the Jaccard index of responses of each turn and the previous turns [3]. The average Jaccard index of each model is shown in Table 4, from which we can observe that: (1) The Jaccard index of human response is the smallest, indicating that existing models have room for improvement in terms of the repetition problem. (2) Our method has a much smaller Jaccard index than CDiaGPT and T5. The action sequence generation, together with the sampling strategy, can effectively alleviate the repetition problem.

| Model | Jaccard Index |
|---|---|
| Human Response | 0.129 |
| CDiaGPT | 0.214 |
| T5 | 0.207 |
| DTA | 0.142 |

**Table 4: Jaccard Index of different models.**

A concrete online dialogue is shown in Figure 1 (b), where the CDiaGPT generates exactly the same responses. Though DTA generates similar action sequences (e.g., combinations of Actions A226, A372, A249 and A109), there are much fewer cases where the action sequence exactly repeats previous turns. The small differences in action sequences can lead to large changes in verbal responses. Besides, the sampling mechanism ensures that the same action in different turns corresponds to different segments (e.g., A226 in three turns), which further enables DTA with better diversity.

## 4 CONCLUSION

In this paper, we propose a task-oriented dialogue system via action-level generation. An effective framework is proposed to build the generation model from the large-scale dialogues with minimum manual effort. The experimental analyses demonstrate our system's capability of tackling the reliability and the efficiency problems encountered with the existing end-to-end generation methods. In the future, we are interested in exploring an integrated system that unifies the discrete modules in DTA in an end-to-end architecture.

---

[2]The grading criteria can be summarized as follows: (i) Score "Low" denotes that Chatbot can not handle user requirements correctly. (ii) Score "Middle" denotes that Chatbot can handle user requirements correctly, but may generate a disfluent or incomplete response. (iii) Score "High" denotes that Chatbot can handle user requirements correctly and complete the conversation perfectly.

## 5 PRESENTER BIOGRAPHY

Presenter: Yuncheng Hua. He is an algorithm engineer at Meituan, focusing on researching and building dialogue systems.

## 6 COMPANY PORTRAIT

Meituan is China's leading shopping platform for locally found consumer products and retail services including entertainment, dining, delivery, travel and other services.

## REFERENCES

[1] Anish Acharya, Suranjit Adhikari, Sanchit Agarwal, Vincent Auvray, Nehal Belgamwar, Arijit Biswas, Shubhra Chandra, Tagyoung Chung, Maryam Fazel-Zarandi, Raefer Gabriel, Shuyang Gao, Rahul Goel, Dilek Hakkani-Tur, Jan Jezabek, Abhay Jha, Jiun-Yu Kao, Prakash Krishnan, Peter Ku, Anuj Goyal, Chien-Wei Lin, Qing Liu, Arindam Mandal, Angeliki Metallinou, Vishal Naik, Yi Pan, Shachi Paul, Vittorio Perera, Abhishek Sethi, Minmin Shen, Nikko Strom, and Eddie Wang. 2021. Alexa Conversations: An Extensible Data-driven Approach for Building Task-oriented Dialogue Systems. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations.* Association for Computational Linguistics, Online, 125–132. https://doi.org/10.18653/v1/2021.naacl-demos.15
[2] Bill Byrne, Karthik Krishnamoorthi, Saravanan Ganesh, and Mihir Kale. 2021. TicketTalk: Toward human-level performance with end-to-end, transaction-based dialog systems. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers).* Association for Computational Linguistics, Online, 671–680. https://doi.org/10.18653/v1/2021.acl-long.55
[3] Luciano da F Costa. 2021. Further generalizations of the Jaccard index. *arXiv preprint arXiv:2110.09619* (2021).
[4] Zihao Fu, Wai Lam, Anthony Man-Cho So, and Bei Shi. 2021. A theoretical analysis of the repetition problem in text generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 12848–12856.
[5] Youngsoo Jang, Jongmin Lee, and Kee-Eung Kim. 2021. GPT-Critic: Offline Reinforcement Learning for End-to-End Task-Oriented Dialogue Systems. In *International Conference on Learning Representations.*
[6] Meixun Jin, Mi-Young Kim, Dongil Kim, and Jong-Hyeok Lee. 2004. Segmentation of Chinese long sentences using commas. In *Proceedings of the Third SIGHAN Workshop on Chinese Language Processing.* 1–8.
[7] Sungjin Lee, Qi Zhu, Ryuichi Takanobu, Zheng Zhang, Yaoqin Zhang, Xiang Li, Jinchao Li, Baolin Peng, Xiujun Li, Minlie Huang, and Jianfeng Gao. 2019. ConvLab: Multi-Domain End-to-End Dialog System Platform. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations.* Association for Computational Linguistics, Florence, Italy, 64–69. https://doi.org/10.18653/v1/P19-3011
[8] Feng-Lin Li, Minghui Qiu, Haiqing Chen, Xiongwei Wang, Xing Gao, Jun Huang, Juwei Ren, Zhongzhou Zhao, Weipeng Zhao, Lei Wang, et al. 2017. Alime assist: An intelligent assistant for creating an innovative e-commerce experience. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management.* 2495–2498.
[9] Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, Peng Xu, Feijun Jiang, Yuxiang Hu, Chen Shi, and Pascale Fung. 2021. Bitod: A bilingual multi-domain dataset for task-oriented dialogue modeling. *arXiv preprint arXiv:2106.02787* (2021).
[10] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
[11] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* 21, 140 (2020), 1–67.
[12] Abigail See, Stephen Roller, Douwe Kiela, and Jason Weston. 2019. What makes a good conversation? how controllable attributes affect human judgments. *arXiv preprint arXiv:1902.08654* (2019).
[13] Shuangyong Song, Chao Wang, Haiqing Chen, and Huan Chen. 2021. An Emotional Comfort Framework for Improving User Satisfaction in E-Commerce Customer Service Chatbots. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Papers.* Association for Computational Linguistics, Online, 130–137. https://doi.org/10.18653/v1/2021.naacl-industry.17
[14] Kai Sun, Seungwhan Moon, Paul Crook, Stephen Roller, Becka Silvert, Bing Liu, Zhiguang Wang, Honglei Liu, Eunjoon Cho, and Claire Cardie. 2021. Adding Chit-Chat to Enhance Task-Oriented Dialogues. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics:*

*Human Language Technologies.* Association for Computational Linguistics, Online, 1570–1583. https://doi.org/10.18653/v1/2021.naacl-main.124
[15] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems* 27 (2014).
[16] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239* (2022).
[17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
[18] Yida Wang, Pei Ke, Yinhe Zheng, Kaili Huang, Yong Jiang, Xiaoyan Zhu, and Minlie Huang. 2020. A Large-Scale Chinese Short-Text Conversation Dataset. In *NLPCC.* https://arxiv.org/abs/2008.03946
[19] Xiangyu Xi, Chenxu Lv, Yuncheng Hua, Wei Ye, Chaobo Sun, Shuaipeng Liu, Fan Yang, and Guanglu Wan. 2022. A Low-Cost, Controllable and Interpretable Task-Oriented Chatbot: With Real-World After-Sale Services as Example. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Madrid, Spain) *(SIGIR '22).* Association for Computing Machinery, New York, NY, USA, 3398–3402. https://doi.org/10.1145/3477495.3536331
[20] Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. ConSERT: A Contrastive Framework for Self-Supervised Sentence Representation Transfer. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers).* Association for Computational Linguistics, Online, 5065–5075. https://doi.org/10.18653/v1/2021.acl-long.393
[21] Zhao Yan, Nan Duan, Peng Chen, Ming Zhou, Jianshe Zhou, and Zhoujun Li. 2017. Building task-oriented dialogue systems for online shopping. In *Thirty-First AAAI Conference on Artificial Intelligence.*
[22] Yunyi Yang, Yunhao Li, and Xiaojun Quan. 2021. Ubar: Towards fully end-to-end task-oriented dialog system with gpt-2. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 14230–14238.
[23] Dian Yu, Luheng He, Yuan Zhang, Xinya Du, Panupong Pasupat, and Qi Li. 2021. Few-shot Intent Classification and Slot Filling with Retrieved Examples. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.* Association for Computational Linguistics, Online, 734–749. https://doi.org/10.18653/v1/2021.naacl-main.59
[24] Qi Zhu, Zheng Zhang, Yan Fang, Xiang Li, Ryuichi Takanobu, Jinchao Li, Baolin Peng, Jianfeng Gao, Xiaoyan Zhu, and Minlie Huang. 2020. ConvLab-2: An Open-Source Toolkit for Building, Evaluating, and Diagnosing Dialogue Systems. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations.* Association for Computational Linguistics, Online, 142–149. https://doi.org/10.18653/v1/2020.acl-demos.19
[25] Xiaoming Zhu. 2019. Case ii (part a): Jimi's growth path: Artificial intelligence has redefined the customer service of jd. com. In *Emerging Champions in the Digital Economy.* Springer, 91–103.