

Predict-then-Decide: A Predictive Approach for Wait or Answer Task in Dialogue Systems

Zehao Lin, *Student Member, IEEE*, Shaobo Cui, Guodun Li, Xiaoming Kang, Feng Ji, Fenglin Li, Zhongzhou Zhao, Haiqing Chen, and Yin Zhang* , *Member, IEEE*

Abstract—Different people have different habits of describing their intents in conversations. Some people tend to deliberate their intents in several successive utterances, i.e., they use several consistent messages for readability instead of a long sentence to express their question. This creates a predicament faced by the application of dialogue systems, especially in real-world industry scenarios, in which the dialogue system is unsure whether it should answer the query of user immediately or wait for further supplementary input. Motivated by such an interesting predicament, we define a novel Wait-or-Answer task for dialogue systems. We shed light on a new research topic about how the dialogue system can be more intelligent to behave in this Wait-or-Answer quandary. Further, we propose a predictive approach named Predict-then-Decide (PTD) to tackle this *Wait-or-Answer* task. More specifically, we take advantage of a *decision* model to help the dialogue system decide whether to wait or answer. The decision of decision model is made with the assistance of two ancillary prediction models: a user prediction and an agent prediction. The user prediction model tries to predict what the user would supplement and uses its prediction to persuade the decision model that the user has some information to add, so the dialogue system should wait. The agent prediction model tries to predict the answer of the dialogue system and convince the decision model that it is a superior choice to answer the query of user immediately since the input of user has come to an end. We conduct our experiments on two real-life scenarios and three public datasets. Experimental results on five datasets show our proposed PTD approach significantly outperforms the existing models in solving this Wait-or-Answer problem.

I. INTRODUCTION

With the availability of large-scale dialogue corpora and the advances in deep learning and reinforcement learning, conversational artificial intelligence has made great progress. In recent years, many works try to improve the performance of data-driven dialogue systems from multiple dimensions, e.g. dialogue states [1], [2], dialogue generation [3], emotion integration [4], knowledge integration [5], [6].

In real-life scenarios, we observe that a large portion of dialogue system users often describe their intents in several successive utterances rather than a single utterance. This brings a critical dilemma in which the dialogue system is not sure

* Corresponding Author

Zehao Lin, Guodun Li and Yin Zhang are with the College of Computer Science and Technology, Zhejiang University, Hangzhou, Zhejiang, 310027, China. (e-mail: georgelin@zju.edu.cn, guodun.li@zju.edu.cn, and zhangyin98@zju.edu.cn).

Shaobo Cui, Xiaoming Kang, Feng Ji, Fenglin Li, Zhongzhou Zhao, and Haiqing Chen are with DAMO Academy, Alibaba Group. (e-mail: yuanchun.csb@alibaba-inc.com, kxm180043@alibaba-inc.com, zhongxiu.jf@alibaba-inc.com, fenglin.lf@alibaba-inc.com, zhongzhou.zhaozz@alibaba-inc.com and haiqing.chenhq@alibaba-inc.com).

Manuscript received March 12, 2021; revised July 14, 2021.

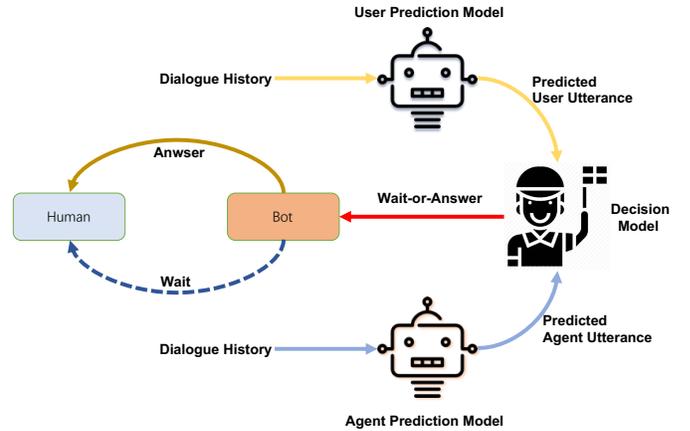


Fig. 1. An overview of the PTD framework. The user prediction model and agent prediction model predict user's and agent's future possible utterance, respectively, and assist the decision model in solving *Wait-or-Answer* problem.

whether it should **wait** for the further input from the user or simply **answer** the question right away. If dialogue systems can not solve wait-or-answer dilemma, users have to convey their intents in a single turn, which harms the user experience. Additionally, too early cut-in or delayed response of dialogue systems will puzzle the users and lead to conversation failure. This Wait-or-Answer dilemma becomes even more complicated and complex when it comes to multi-turn dialogue systems. Despite the surge of attention into the dialogue system models, very few research works have specially investigated the Wait-or-Answer problem. Some studies have been done on the incomplete user utterance task [7], [8], [9], however, in multi-turn dialogue environment, all utterances from users are complete sentences (i.e. all sentences are meaningful and unbroken in both syntax and semantics). This requires our model to be able to learn user intents and habits from context.

The most widely-applied approach to work around in industry [10], [11] is to extend dialogue systems' response time: the dialogue agents always wait for some time in case of user's further input. However, this *passive* extended-waiting-time strategy may cause incoherence of the conversation and lead to poor user experience. For example, a 3-second extended-waiting-time may seem short to some people but apparently too long for some other users. A comfortable waiting time varies with each individual, which makes it difficult to set the accurate waiting time. Recently, Liu et al. [12] address the similar problem under the semantic matching task and obtain the weights of the co-occurring utterances in the context to

determine whether the given context should be responded to.

Different from the previous strategy, we propose to address the Wait-or-Answer issue in a *predictive approach*, which purely relies on dialogue context information and can solve wait-or-answer task alone or as a supplement to time expanding and triggering methods. Specifically, we try to predict the user’s next action based on the current dialogue history: (1) to provide supplemental information (if so, the bot should wait). (2) to wait for the response from chatbot (if so, the bot should answer the user’s query). The most intuitive predictive approach is to apply a classification model, which we use is a classification model to predict whether the dialogue system should wait or answer based on the dialogue history. These kinds of methods only consider the information in past dialogue history but omit the user and agent’s possible future intention. Intuitively, suppose that we can predict what the user would supplement if the user wants to express further information, and what the dialogue system would answer if the user has completed his or her question and is waiting for an answer, the dialogue system has more confidence to decide whether to wait or answer.

Motivated by such intuitions, we propose a model named **Predict-then-Decide** (PTD). As shown in Figure 1, there is a decision model that controls whether the bot should answer the user query or wait for further information. Except for the decision model, there are two auxiliary prediction models: the user prediction and the agent prediction. The user prediction model persuades the decision model that the bot should wait for further input from users. While for the agent prediction model, it tries to convince the decision model that the bot should immediately answer users’ queries. As for the decision model, given the *suggestions* from these two prediction models, it makes its decision whether the bot should wait or answer.

At last, to evaluate the performance of PTD framework, we test several baselines including two rule-based systems and three supervised models, and three PTD variants on two real industry datasets collected by ourselves and three public datasets for better reliability. Experimental results and analysis show the improvements brought by PTD framework.

In summary, this paper makes the following contributions:

- This paper explicitly defines the Wait-or-Answer task, which is crucial to further enhance the capability of dialogue systems;
- We propose a novel framework, named Predict-then-Decide (PTD), to solve the Wait-or-Answer task, which uses two prediction models and one decision model to help dialogue systems decide whether to wait or answer;
- Experimental results on both real industry scenarios and public datasets demonstrate that our model significantly outperforms the baselines, which validates the benefits brought by our PTD framework. Our modified public datasets and code are released ¹ to the public for further research in both academia and industry.

II. PRELIMINARY

In this section we provide some background knowledge about Dialogue Systems in Section II-A. Besides, our proposed PTD framework involves both generative models (for prediction models in PTD) and classification models (for the decision model in PTD), we present some preliminary work about the generative and classification models in NLP in Section II-B.

A. Dialogue Systems

Researches on dialogue systems are mainly divided into two categories: task-oriented dialogue systems and chit-chat dialogue systems. Task-oriented dialogue systems [13], [14], [15] aim at solving tasks in specific domains with grounding knowledge while chit-chat bots [16], [17], [18] mainly concentrate on interacting with a human to provide reasonable responses and entertainment [19]. Recent years research on task-oriented dialogue systems mainly concentrates on dialogue states [20] and knowledge integration [6], [21] using pipeline or end to end models. Chit-chat bots focus on conversing with the human in open domains. Though chit-chat bots seem to perform totally different from task-oriented dialogue systems, actually as revealed in Yan et al. [22], nearly 80% utterances are chit-chat messages in the online shopping scenario and handling those queries is closely related to user experiences. Many studies have investigated how to apply neural networks to the components of dialogue systems or end-to-end dialogue frameworks [15], [23]. The advantage of deep learning is its ability to leverage large amounts of data from the internet, sensors, etc. The big conversation data and deep learning techniques like SEQ2SEQ [24] and attention mechanism [25] help the model understand the utterances, retrieve background knowledge and generate responses.

B. Generative and Classification Models

Dialogue Generation In general, two major approaches have been developed for dialogue divided by the reply types: generative methods such as sequence-to-sequence models, which generate proper responses during the conversation; and retrieval-based methods, which learn to select responses from the current conversation from a repository.

The generative method has been attracting more and more attention [26], [27]. The main reason is that, when compared to retrieval-based dialogue systems, generative models can sometimes produce more fluent and flexible replies, making them more user friendly in some cases. Unlike the retrieval method, Natural Language Generation (NLG) translate a communication goal selected by the dialogue manager into a natural language form [28]. It reflects the naturalness of a dialogue system, and thus the user experience. Another reason is that, in addition to the fluency and accuracy of responses, generative systems are far more flexible to use for common users than retrieval based systems.

Conventional template or rule-based approaches mainly contain a set of templates, rules, and hand-crafted heuristics designed by domain experts. This makes it labor-intensive yet rigid, motivating researchers to find more data-driven

¹Open Source Repository: <https://github.com/mumeblossom/PTD>

approaches [5], [21] to optimize a generation module from corpora, one of which, Semantically Controlled LSTM (SC-LSTM) [29], a variant of LSTM [30], gives semantic control on language generation with an extra component. As for the fully-data driven dialogue systems, SEQ2SEQ [24] based encoder-decoder frameworks and attention mechanisms [25] are still the most widely adopted [21], [5], [31] techniques. Transformer [32] based models, e.g. BART [33], T5 [34] and GPT [35], show its effectiveness compared with traditional CNN or RNN based models in generation tasks. They are able to utilize a large amount of natural language data from the internet by self-supervised learning and fine-tune themselves in downstream tasks.

Text Classification Text classification is a critical problem in all NLP tasks, and it has been widely investigated and studied [36], [37], [38], [39] in recent decades. Text classification can be done on multiple levels, including document classification [40], [41], sentence classification [42], emotion classification [43], and so on.

Though end-to-end methods play a more and more important role in dialogue system, the text classification modules [36], [37] remain very useful in many problems like emotion recognition [44], gender recognition [45], intent detection [46], etc. There have been several widely used text classification methods proposed, e.g. Recurrent Neural Networks (RNNs) and CNNs. Typically RNN is trained to recognize patterns across time, while CNN learns to recognize patterns across space. [47] proposed TextCNNs trained on top of pre-trained word vectors for sentence-level classification tasks and achieved excellent results on multiple benchmarks.

Besides RNNs and CNNs, a powerful network architecture called Transformer [32] is solely based on attention mechanisms and achieves promising performance in many NLP tasks. To make the best use of unlabeled data, Devlin et al. [48] introduce a new language representation model with two auxiliary pre-training tasks called BERT, which stands for Bidirectional Encoder Representations from Transformers.

III. THE WAIT-OR-ANSWER TASK

Why Do We Study The Wait-or-Answer Task? Conventional dialogue systems mainly concentrate on the accuracy and fluency of generated or retrieved answers. These kinds of dialogue systems, including most commercial chatbots, require users to strictly follow the designed conversation instructions. For example, users must be ready to finish all of the words they want to say in one breath and without pausing. This requires users describing their intents in a single sentence.

However, the aforementioned setting in existing dialogue systems does NOT hold in real-life settings. For instance, as shown in Figure 2, in a real-world scenario where a user requests information about a theater, the agent firstly starts the conversation with “Good morning. Vane Theater at your service.” (A1), then the user replies with three sentences, firstly “Hello” (U11), secondly “I’m thinking about watching a Chinese traditional opera with a foreign girl.” (U12) and thirdly “What’s on this weekend?” (U13). Generally speaking, users won’t speak all sentences without a breath. If the agent cuts

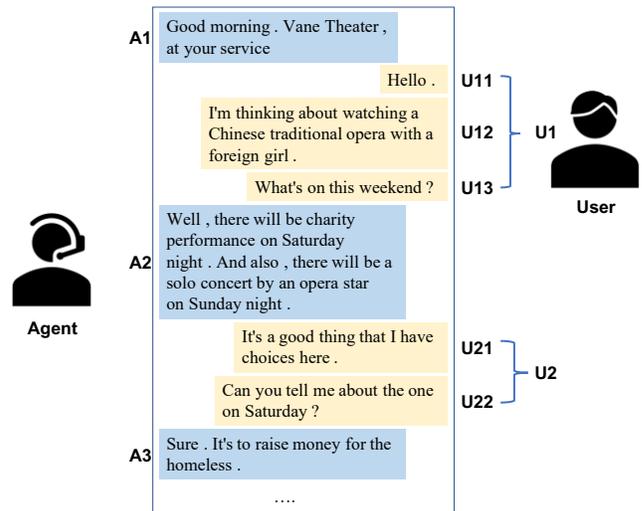


Fig. 2. A multi-turn dialogue fragment. In this case, a user sends split utterances in a turn, e.g. split U1 to {U11, U12 and U13}

in the wrong moment of the conversation, e.g. immediately replies to the user’s second statement, the agent has to guess what the user wants and omit the important information “on this weekend” in the third sentence. So in this case, the agent should wait for the user until he or she has finished sending his or her last message, otherwise, the pace of the conversation will be messed up. However, existing dialogue agents can not handle this scenario well and will reply to the user’s every utterance immediately or after a fixed time interval.

There are mainly two issues when applying existing dialogue agents to the real-life conversation: (1) Existing dialogue systems lack the capability of deciding to avoid generating bad responses based on semantically incomplete utterance, when they receive a short utterance from users as the start of a conversation. (2) Existing dialogue systems may cut into a conversation at an inappropriate point, which could confuse the user and mess up the pace of conversation and thus lead to nonsense interactions. In other words, the existing dialogue system can NOT catch the right moment to Answer or Wait.

As stated above, it is worthwhile to investigate this Wait-or-Answer task which would empower the dialogue system to enhance their ability to make appropriate decisions in the wait-or-answer dilemma.

Task Formulation The *passive* extended-waiting-time strategy may cause incoherence of the conversation and lead to poor user experience. So our task is to actively predict whether the agent should answer immediately. Under this premise, our problem is formulated as follows. There is a conversation history represented as a sequence of utterances: $X = \{x_1, x_2, \dots, x_m\}$, where each utterance x_i itself is a sequence of words $x_{i_1}, x_{i_2}, x_{i_3} \dots x_{i_n}$. In addition, each utterance has following additional labeled tags:

- (1). Turn id: which turn the utterance locates in.
- (2). Sub-turn id: the position of the utterance in its turn which may have more than one utterance.
- (3). Speaker id: who sends out the utterance. 0 means the user while 1 means the agent.

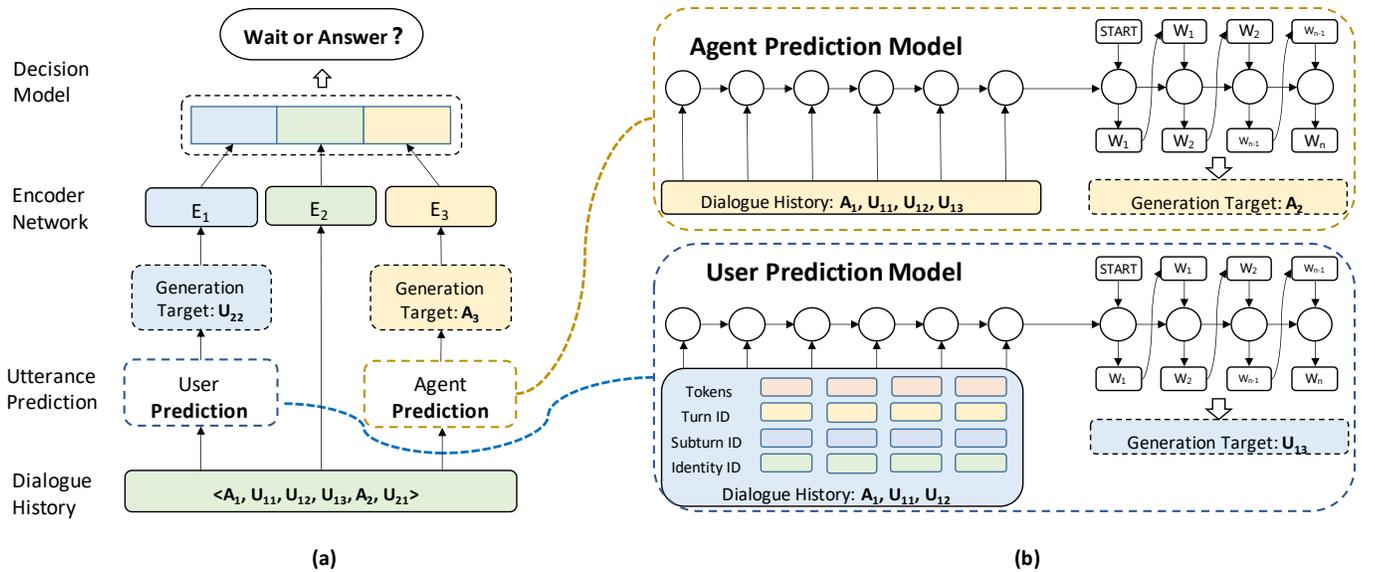


Fig. 3. (a) The PTD framework will complete *Wait-or-Answer* task using dialogue history and two trained prediction models' predictions.(b) Details of the agent and user prediction models trained using the same dialogues but different samples.

Now, given a dialogue history X and tags T , the goal of the model is to predict a label $Y \in \{0, 1\}$, the action the agent would take, where $Y = 0$ means the agent will wait for the user for the next message, and $Y = 1$ means the agent will reply immediately. Formally we will choose the action that maximizes the following probability:

$$Y = \arg \max_y P(y|X, T). \quad (1)$$

IV. THE PREDICT-THEN-DECIDE (PTD) FRAMEWORK

In this section, we firstly present the overview of our PTD framework in Section IV-A. Section IV-B and Section IV-C are about the detailed model structures of prediction model and decision model respectively. At last, we describe the overall training and inference phase in Section IV-D.

A. The Overview of PTD Framework

The PTD framework has two prediction models and one decision model, as shown in Figure 1. The decision model makes the final decision about whether the dialogue system should answer users' queries immediately or wait for the users' further information. We use two prediction models to assist the decision model. The user prediction model forecasts what the user might supplement. Then, the user prediction model uses this **simulated** query to convince the decision model to wait for the user's following input since the user does have some information to supplement. The agent prediction model, nevertheless, predicts the dialogue system's answer for users' present query. Then, the agent prediction model utilizes this **simulated** answer to make the decision model believe that the dialogue system should answer the user's queries immediately because the user has finished its input.

In fact, these two prediction models act as the world model [49], which creates a virtual environment to simulate the possible future dialogue to train the agent, for the decision

model. More specifically, the output of the user prediction model (simulated query) and the output of the agent prediction model (simulated answer) both function as the simulated experience. Peng et al. [50] first propose Deep Dyna-Q, incorporating into the dialogue agent a world model to mimic real user response and generate simulated experience. Compared to models that directly apply a classifier, such as TextCNN and BERT, to dialogue systems to solve the *Wait-or-Answer* problem, our proposed prediction models are better at learning semantic information, from both history and future possible utterances, by training on the corpus, and providing supplemental prediction to the decision model. Prediction models will also magnify the errors, providing negative feedback and making it easier for our PTD to understand which decision is better.

B. Prediction Model

The prediction model in PTD generates the next possible utterance given the dialogue history. There are two prediction models in our method: the user prediction model and the agent prediction model. The goal of the two prediction models is to learn the user's and agent's speaking style respectively and generate possible future utterances from different perspectives.

As shown in Figure 3 (a), prediction model itself is a sequence generation model. We use one-hot embedding to convert all words and the related tags to one-hot vectors $w_n \in \mathbf{R}^{|V|}$, where $|V|$ is the length of the vocabulary list. Then we extend each word x_{i_j} in utterance x_i by concatenating the token itself with turn tag, identity tag, and sub-turn tag. We adopt SEQ2SEQ as the basic architecture and LSTMs as the encoder and decoder networks. LSTMs will encode each extended word w_t as a continuous vector h_t at each time step t . The process can be formulated as: $h_t = \text{LSTMs}(h_{t-1})$. For the same piece of dialogue, we split it into different samples for different prediction models. As shown in Figure 2 and 3

(a), we use (A1, U11, U12) as dialogue history input and U13 as ground truth to train the user prediction model and use (A1, U11, U12, U13) as dialogue history and A2 as ground truth to train the agent prediction model.

Apart from the extra labels, the prediction models are trained in the conventional way: $h_t = \text{LSTMs}(h_{t-1})$, and the decoder is the similarly structured LSTMs but h_t will be fed to a Softmax with $W_v \in \mathbf{R}^{h \times |V|}$, $b_v \in \mathbf{R}^{|V|}$, which will produce a probability distribution p_t over all words, formally: $p_t = \text{Softmax}(W_v h_t + b_v)$.

Decoder will select the word based on p_t at each time step. The loss for prediction model is the sum of the negative log-likelihood of the correct word at all time steps:

$$L_{\text{PRE}} = - \sum_{t=1}^N \log(p_t), \quad (2)$$

where N is the length of the generated sentence. During inference, we also apply a beam search to improve generation performance. Finally, the trained agent prediction model and user prediction model are obtained.

C. Decision Model

The decision module is fundamentally a text classifier. In our settings, to fully utilize the dialogue history and latent semantic information, we rewrite the objective in Equation (1) as follows:

$$R' = \arg \max_y P(y|X, T, R_a, R_u), \quad (3)$$

where $R_a = \text{IG}_a(X, T)$ and $R_u = \text{IG}_u(X, T)$ represent the generated utterances from the agent prediction model and user prediction model respectively. $R' \in [0, 1]$ is a selection indicator where $R' = 1$ means selecting R_a whereas 0 means selecting R_u .

We adopt several architectures like Bi-GRUs, TextCNNs, and BERT as the base model of the decision module. Without loss of generality, here we illustrate how to build a decision model by taking TextCNNs as an example.

As shown in Figure 3, the three similarly structured TextCNNs following the work [47] take the inferred responses R_a , R_u and dialogue history X , tags T . For each raw word sequence x_1, \dots, x_n , we encode each word as one-hot vector $w_i \in \mathbf{R}^{|V|}$. By looking up a word embedding matrix $E \in \mathbf{R}^{|V| \times d}$, the input text is represented as an input matrix $Q \in \mathbf{R}^{l \times d}$, where l is the length of sequence of words and d is the dimension of word embedding features. The matrix is then fed into the similar structured CNNs using one layer convolution with max-over-time pooling to get the feature maps of X , R_a and R_u :

$$\begin{aligned} \hat{C}_x &= \text{TextCNNs}(X); \\ \hat{C}_a &= \text{TextCNNs}(R_a); \\ \hat{C}_u &= \text{TextCNNs}(R_u). \end{aligned} \quad (4)$$

Then we will have two possible dialogue paths, X with R_a and X with R_u , representing D_a and D_u :

$$\begin{cases} D_a = W_1[\hat{C}_x; \hat{C}_a] + b_1 & \text{Answer path: } X \text{ with } R_a. \\ D_u = W_2[\hat{C}_x; \hat{C}_u] + b_2 & \text{Wait path: } X \text{ with } R_u. \end{cases} \quad (5)$$

Then the decision model will predict which of these two paths (user path or agent path) should be taken. Namely, the decision model conducts a binary classification task:

$$P = \text{Softmax}(W_4(W_3[D_a; D_u] + b_3) + b_4), \quad (6)$$

where W_1 to W_4 and b_1 to b_4 are learnt parameters. At last we will get a two-dimensional probability distribution P , which indicates the most reasonable response. The loss function of the decision model is the negative log-likelihood of the probability of choosing the correct action:

$$L_{\text{DEC}} = - \sum_{i=1}^M \sum_{j=0}^1 Y_i(j) \log(P(j)), \quad (7)$$

where $j = \{0, 1\}$ is the action label, M is the number of samples and Y_i is an one-hot encoding of the ground-truth label of the i -th sample.

The decision modules based on Bi-GRU or BERT are implemented similarly to TextCNNs in Equation (4).

Algorithm 1 Training and inference of PTD Framework

```

1: procedure TRAIN-USER-PREDICTION( $X, T$ )
2:   for mini-batch( $X, T$ ) do
3:      $R_u \leftarrow \text{IG}_u(X, T)$ 
4:     minimize Equation (2) to optimize user prediction
      model
5:   end
6:   return trained user prediction model  $\text{IG}_u$ 
7: end procedure
8: procedure TRAIN-AGENT-PREDICTION( $X, T$ )
9:   for mini-batch( $X, T$ ) do
10:     $R_a \leftarrow \text{IG}_a(X, T)$ ;
11:    minimize Equation (2) to optimize agent prediction
      model
12:   end
13:   return trained agent prediction model  $\text{IG}_a$ 
14: end procedure
15: procedure TRAIN-DECISION( $X, T$ )
16:   for mini-batch( $X, T$ ) do
17:     $R_u \leftarrow \text{IG}_u(X, T)$ 
18:     $R_a \leftarrow \text{IG}_a(X, T)$ 
19:     $R' \leftarrow \text{decision model}(X, T, R_u, R_a)$ 
20:    minimize Equation (7) to optimize decision model
21:   end
22:   return trained decision model
23: end procedure
24: procedure INFERENCE( $X, T$ )
25:    $R_u \leftarrow \text{IG}_u(X, T)$ 
26:    $R_a \leftarrow \text{IG}_a(X, T)$ 
27:    $R' \leftarrow \text{decision model}(X, T, R_u, R_a)$ 
28: end procedure
29: TRAIN-USER-PREDICTION;
30: TRAIN-AGENT-PREDICTION;
31: TRAIN-DECISION;
32: INFERENCE;

```

TABLE I
REAL INDUSTRY DATASETS STATISTICS.

Datasets	E-COMM			After-Sales		
	Train	Valid	Test	Train	Valid	Test
Vocabulary Size		9948			3877	
Dialogues	1676	209	211	1796	224	226
Avg. Turns/Dialogue	12.50	11.79	12.03	10.14	8.76	11.26
Avg. User Sub-turns	1.40	1.41	1.40	1.56	1.59	1.55
Avg. Utterance Length	41.37	45.27	40.75	50.07	49.78	47.04
Avg. Agent's Utterances Length	28.78	29.39	28.83	24.49	24.40	24.13
Avg. User's Utterances Length	5.69	8.08	6.10	8.84	8.54	7.72
Agent Wait Samples Size	20976	2464	2538	18213	1962	2544
Agent Answer Samples Size	8348	1016	1027	10112	1150	1401

TABLE II
PUBLIC DATASETS STATISTICS. NOTE THAT THE STATISTICS ARE BASED ON THE MODIFIED DATASET DESCRIBED IN SECTION V-A3

Datasets	MultiWOZ			DailyDialog			CCPE		
	Train	Valid	Test	Train	Valid	Test	Train	Valid	Test
Vocabulary Size		2443			6219			4855	
Dialogues	8423	1000	1000	11118	1000	1000	398	49	52
Avg. Turns/Dialogue	6.32	6.97	6.98	4.09	4.21	4.03	9.7	9.96	9.92
Avg. Split User Turns	1.89	1.92	1.94	2.09	2.12	2.12	3.12	3.02	2.73
Avg. Utterance Length	10.54	10.7	10.56	8.71	8.54	8.75	7.93	8.02	7.78
Avg. Agent's Utterances Length	14.43	14.78	14.69	12.04	11.81	12.17	8.7	8.84	8.19
Avg. User's Utterances Length	6.18	6.28	6.17	5.91	5.87	5.96	7.61	7.66	7.56
Agent Wait Samples Size	47341	6410	6573	49540	4717	4510	8183	973	894
Agent Answer Samples Size	53249	6970	6983	41547	3846	3689	3455	436	464

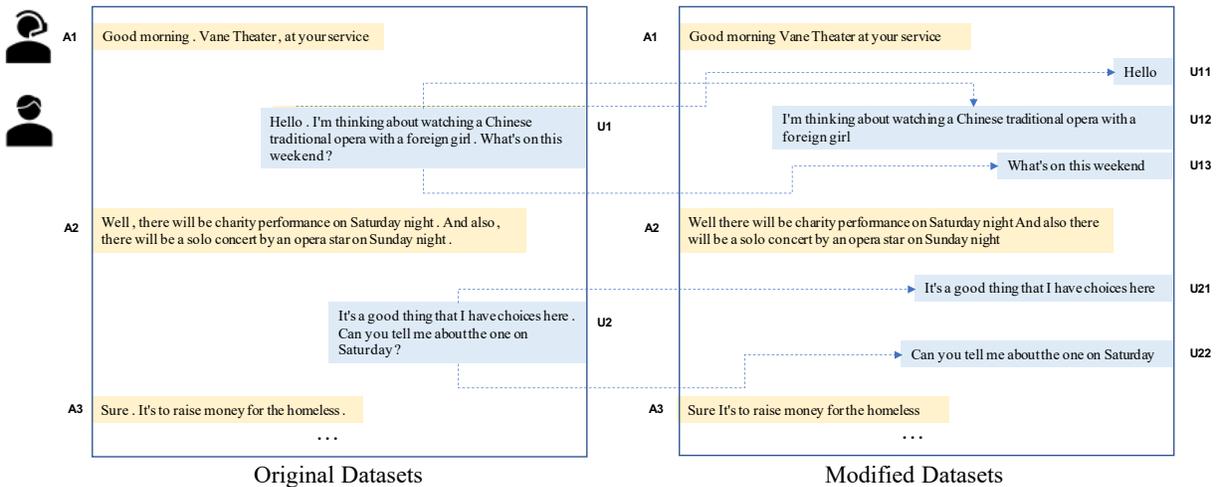


Fig. 4. Comparison of a piece of the original dialogue and Our modified dialogue following the section V-A3. Note that some elements task-oriented corpus like MultiWOZ slot values, knowledge base, and ontology content are not shown in Original Datasets and have been deleted in Modified Datasets.

D. Training and Inference of PTD Framework

As shown in Figure 3 and Algorithm 1, we show the procedure of the model’s training and inference. We first train user prediction model on the dialogue history with ground-truth as user’s utterance (For example, from $[A_1, U_{11}]$ to U_{12} in Figure 2) and agent prediction model on the dialogue history with ground-truth as agent’s utterance (For example, from $[A_1, U_{11}, U_{12}, U_{13}]$ to A_2). And then we infer predicted user and agent future possible utterances from user and agent prediction models as decision model’s training data. In this kind of design, the two prediction models will not only simulate future possible dialogue to support *wait* and *answer* action, but also magnify the distinction between decisions, e.g. the performance of user prediction model will be poor when the ground-truth is *answer* because user prediction model never learned how to speak like an agent, this will make decision model easier to distinguish which decision is better. At last, we feed the predicted utterances and original dialogue history together to train the decision model. During the inference procedure, we simply use two prediction models to predict possible user and agent’s utterance (U'_{22} and A'_3) and combined with dialogue histories $[A_1, U_{11}, U_{12}, U_{13}, A_2, U_{21}]$ for the decision model to decide whether the model should wait or answer right away.

V. EXPERIMENTAL SETUP

In this section, we firstly present the process of data construction in Section V-A, after which we give the evaluation metrics we use in Section V-B. Finally, we detail the training setup of baselines and our PTD models in Section V-C and Section V-D respectively.

A. Datasets Construction

1) *Realistic Scenarios*: We test our approach and baselines on two real scenarios. Table I shows the statistics of two real industry datasets.

- **E-COMM**. E-COMM is a dataset that contains real dialogue history between customers and service representatives from one of the biggest cross-border e-commerce. In a session of dialogue, both customers and customer service may have sub-turns. We merge the sub-turns of customers service in each turns into one long utterance.
- **After-Sales**. After-Sales is a dataset that contains real dialogue history between customers and customer service, talking about equipment maintenance and repairing topics. We take the same pre-processing method as E-COMM.

2) *Public Datasets*: As the proposed approach mainly concentrates on the interaction between human and computer, we select and modify three datasets of very different styles to evaluate the performance and generalization of our method. Two of them are task-oriented dialogue datasets. One is a large MultiWOZ 2.0² and the other is a smaller dataset Coached Conversational Preference Elicitation (CCPE)³, which has

many more turns per dialogue. The last dataset is a chit-chat dataset DailyDialog⁴. All datasets are collected during human-to-human conversations. We evaluate and compare the results with the baseline methods from multiple perspectives. Table II shows the statistics of datasets and details of datasets are described as follows:

- **MultiWOZ 2.0** [20]. MultiDomain Wizard-of-Oz dataset (MultiWOZ) is a fully-labeled collection of human-human written conversations. Compared with previous task-oriented dialogue datasets, e.g. DSTC 2 [51] and KVR [14], it is a much larger multi-turn conversational corpus and across several domains and topics.
- **DailyDialog** [52]. DailyDialog is a high-quality multi-turn dialogue dataset, which contains conversations about daily life. In this dataset, humans often first respond to the previous context and then propose their own questions and suggestions. In this way, people pay more attention to others’ words and are willing to continue the conversation. The speaker’s behavior will be more unpredictable and complex than in the task-oriented dialogue datasets.
- **CCPE** [53]. CCPE is a dataset consisting of 502 English dialogues. Though it seems much smaller than MultiWOZ 2.0 and DailyDialog, CCPE has 12,000 annotated utterances between a user and an assistant discussing movie preferences in natural language. It is collected using a Wizard-of-Oz methodology between two paid crowd-workers and focuses on the movie domain. We select this dataset to test if our model can run well on both large and small datasets.

3) *The Pipeline of Dataset Construction*: As the task we concentrate on, making a decision to wait or answer, is quite different from traditional dialogue systems, existing dialogue datasets will be unable to provide the information for training and testing. Thus we propose a fairly simple and general dataset construction method to directly rebuild over the existing public dialogue corpus.

We modify the datasets with the following steps:

- I. Delexicalisation**: For task-oriented dialogue, slot labels are important for navigating the system to complete a specific task. However, those labels and accurate values from ontology files will not benefit our task essentially. So we replace all specific values with a slot placeholder in the pre-processing step.
- II. Utterance segmentation**: Existing datasets concentrate on the dialogue content, combining multiple sentences into one utterance each turn when gathering the data. In this step, we split the combined utterance into multiple *complete* utterances according to original datasets. This makes task in line with user habits. And models can not make decisions simply based on the completeness of sentences. In this pre-processing procedure, we do not divide all user utterances but half part of all in train/development/test sets. All punctuation marks are eliminated after this procedure to improve the applicability in most situations including Spoken Dialogue Systems.
- III. Extra Labeling**: We add several labels, including turn tags, sub-turn tags, and role tags, to each split and original

²<http://dialogue.mi.eng.cam.ac.uk/index.php/corpus/>

³<https://research.google/tools/datasets/coached-conversational-preference-elicitation/>

⁴<http://yanran.li/dailydialog.html>

TABLE III

AN EXAMPLE OF THE PREDICTION MODEL’S GENERATION AND DECISION MODEL’S SELECTION. DECISION MODEL CHOOSES THE AGENT PREDICTION MODEL MEANS AGENT SHOULD ANSWER THE USER DIRECTLY.

Example	
Dialogue History	User: what is the address for pizza hut in cherry hinton
Ground-Truth	Agent: the address is [restaurant_address]
Prediction	Agent Prediction Model the address is [restaurant_address] can i help you with anything else User Prediction Model i am looking for a guesthouse
Decision Model Selection	Agent Prediction Model

TABLE IV

ACCURACY RESULTS ON TWO REAL SCENARIO DATASETS. BETTER RESULTS BETWEEN BASELINES AND CORRESPONDING PTD MODELS ARE IN BOLD.

Models	E-COMM		After-Sales	
	Accuracy	F1	Accuracy	F1
ATLU	39.16	35.74	43.46	32.91
PTSU	48.18	57.47	49.80	58.58
Bi-GRU	71.96	78.99	64.32	78.23
GRU-PTD	72.06	83.69	70.85	82.62
TextCNN	69.91	76.60	65.78	78.63
TextCNN-PTD	72.08	83.71	70.71	82.53
BERT	69.45	81.46	66.02	79.43
BERT-PTD	71.98	83.71	71.14	82.86

sentences in order to (1) label the speaker role and dialogue turns (2) mark the ground truth for supervised training and evaluate the baselines and our model.

Finally, we have the modified datasets which imitate the real-life human chatting behaviors. As shown in Figure 4, we compare one original dialogue (in this example, from DailyDialog) with our modified one. Our modified datasets and code are open-sourced to both academic and industrial communities.

B. Evaluation Metrics

In our *Wait-or-Answer* task, we define the *Answer* action of the agent as the positive samples and the *Wait* action is the negative action. As both the positive and negative actions are important in this task, so we choose the model with the accuracy metrics instead of precision or recall.

To compare with dataset baselines in multiple dimensions and test the model’s performance, we use the overall Bilingual Evaluation Understudy (BLEU) [54], which is the cumulative score for BLEU-1 to BLEU-4, to evaluate the prediction models’ generation performance. As for the decision model, we use the accuracy score as the main metrics to evaluate the wait-or-answer decision and select models. Though our aim is to obtain a model with best accuracy in distinguishing wait or answer, apart from BLEU and accuracy, we also present Precision, Recall and F1 to evaluate baselines and our models from multiple perspectives. Details are as follows:

- *Bilingual Evaluation Understudy (BLEU)* [54]. BLEU has been widely employed in evaluating sequence generation including machine translation, text summarization, and dialogue systems. BLEU calculates the n-gram precision which is the fraction of n-grams in the candidate text which is present in any of the reference texts.

- *Accuracy* The accuracy metric is the probability of whether the decision model can successfully classify the ground truth in the test dataset. The accuracy score in our experiments is the correct ratio in all samples.
- *Precision* also called positive predictive value is the fraction of relevant instances among the retrieved instances. In our case, we calculate precision by the ratio of correctly predicted answer actions in all predicted answer actions of the test dataset.
- *Recall* also known as sensitivity is the fraction of the total amount of relevant instances that are actually retrieved. In our case, we calculate recall by the ratio of correctly predicted answer actions in all answer actions of the test dataset.
- *F1 Score* Only consider the precision p or the recall r is difficult to determine which one is really better of not both p and r get a better score. F1 score considers both the precision p and the recall r of the test to compute the score. We calculate the F1 score by the harmonic mean of the precision and the recall.

C. Baselines and Their Training Setup

To make the best practice hyper-parameter settings adopted by each training set in baselines and our models. For baselines, we only rely on dialogue history to make decisions. We conduct experiments on the following baselines with fine-tuned parameters:

- **Gated Recurrent Units (GRU)** [55]: we use GRU to encode the input history as vector and convert it into a two-dimensional probability vector with softmax to predict whether the model should wait or answer. We test hidden size ranging from 200 to 600, dropout rate ranging from 0.2 to 0.8, batch size in [32, 64, 128, 256].
 - **TextCNN** [47]: we use TextCNN to encode the input history as vector and convert it into a two-dimensional probability vector with softmax to predict whether the model should wait or answer and search the best performance in batch size ranging in [32, 64, 128, 256], dropout rate from 0.3 to 0.7, kernel number, which is the number of convolution kernels of each size type, from 100 to 600, kernel size in [(1,2,3),(3,4,5),(5,6,7),(7,8,9)].
 - **BERT** [48]: we use BERT to encode the input history as vector and convert it into a two-dimensional probability vector with softmax to predict whether the model should wait or answer and test learning rate in [2e-5, 3e-5, 5e-5], training epochs in [2.0, 3.0, 4.0] and batch size in [16, 32].
- We also employ two rule-based baselines referred to the work in Liu et al. [12]:

TABLE V
RESULTS ON THREE PUBLIC DATASETS. BETTER RESULTS BETWEEN BASELINES AND CORRESPONDING PTD MODELS ARE IN **BOLD**.

Dataset	MultiWOZ				DailyDialog				CCPE			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
ATLU	52.57	60.89	22.18	32.52	57.71	54.67	35.21	42.84	61.27	35.24	15.95	21.96
PTSU	60.78	60.93	66.53	63.61	59.26	53.62	70.07	60.75	42.27	33.44	69.61	45.17
Bi-GRU	79.12	75.69	87.70	80.85	75.23	72.07	72.94	71.86	67.53	54.83	31.15	38.49
GRU-PTD	82.03	79.02	88.87	83.27	77.80	77.97	77.29	75.71	72.69	63.40	48.47	53.50
TextCNN	77.68	73.61	88.85	80.03	75.79	71.03	78.49	73.91	68.65	59.78	27.63	36.35
TextCNN-PTD	80.75	77.17	89.08	82.52	79.02	77.14	74.87	75.35	73.32	68.99	41.90	51.43
BERT	80.75	76.93	89.46	82.73	78.68	75.03	78.86	76.90	70.99	59.21	48.49	53.31
BERT-PTD	82.73	80.36	87.99	84.00	79.35	75.92	79.23	77.54	75.41	67.86	53.23	59.66

- **Active Triggering based on Longest Utterance (ATLU)**: It considers the lengths of utterances in the input history to make decisions. If the length of the last utterance is longer than others, the model should answer directly.
- **Passive Triggering based on Shortest Utterance (PTSU)**: Similar to ATLU, it compares the length of the last utterance with the lengths of all utterances in the input history to make decisions. The model should wait only if the length of the last utterance is shorter than others.

D. PTD Models and Their Training Setup

To test the performance of our proposed PTD framework, we apply our PTD framework in the baselines and obtain **GRU-PTD**, **TextCNN-PTD** and **BERT-PTD**. The detailed setting on public datasets is described as follows:

- **GRU-PTD**: for GRU-PTD on MultiWOZ, batch size is 32, hidden size is 300, dropout rate is 0.3. On DailyDialog, batch size is 64, hidden size is 500, and dropout rate is 0.5. On CCPE, batch size is 32, hidden size is 200, and dropout rate is 0.8.
- **TextCNN-PTD**: for TextCNN-PTD on MultiWOZ, batch size is 64, kernel number is 400, kernel size is (7,8,9), and dropout rate is 0.3. On DailyDialog, batch size is 32, kernel number is 400, kernel size is (5,6,7), and dropout rate is 0.5. On CCPE, batch size is 64, kernel number is 600, kernel size is (5,6,7), and dropout rate is 0.4.
- **BERT-PTD**: the maximum sequence length to 128, batch size is 32, and the number of training epochs is 3.0 to 4.0.

During training, we also adopt a learning rate decay factor as 0.5. All experiments employ the teacher-forcing scheme [56], feeding the gold target of last time. We also perform early stopping for decision model when its performance does not increase during 6 consecutive validation epochs. We test the hidden size in [32, 64, 128, 256] and set dropout rate in [0.1, 0.2]. The learning rate is initiated with 0.001 and the training batch is set to 64. The metrics results are coming from the best result settings for each dataset.

VI. EXPERIMENTAL RESULTS AND ANALYSIS

A. Experiment Results

A Case Showing How PTD Works For a better understanding of our PTD framework, we also present an example of how PTD acts in Table III. The User prediction model predicts

what the user might supplement based on the dialogue history: *i am looking for a guesthouse*. The agent prediction model, however, predicts what the dialogue system might answer: *the address is [restaurant_address] can i help you with anything else*. Based on both predictions, the decision model concludes that it's a better choice to take the supplementary input of agent. So the dialogue system decides to continue the dialogue rather than to wait.

Experimental Results To illustrate the benefits brought by our PTD framework, we present the comparison results between our PTD models⁵ and the baselines on real industry datasets in Table IV. We can observe that our PTD models achieve superior performance compared with their counterparts. Additionally, to evaluate our models and baselines from other multiple dimensions, we present the comparison results⁶ on accuracy, precision, recall, and F1 results in Table V, which also prove the benefits brought by our PTD framework. To better analyze the effects of prediction models in our PTD framework, we present different PTD models' performance scores with different types of prediction model in Table VI.

B. Experimental Analysis

Shortcomings of Traditional Rule-based Models As shown in Table IV, we can find out that though the ATLU and PTSU models are easy to implement, their results on all datasets are unsatisfactory and unstable. ATLU achieves 39.16 and 43.46 in terms of accuracy. PTSU has slightly better performance than ATLU. Both of two rule-based model have low performance in terms of F1 scores. This indicates that the rules based models have difficulty in this Wait-or-Answer task. The phenomenon that ATLU achieves 39.16 on E-COMM and 43.46 on After-Sales in terms of accuracy shows that the performance of rule-based models are unstable and highly rely on the habits of users in datasets. A significant difference between Wait-or-Answer task and incomplete utterance classification task is that all utterances in Wait-or-Answer task are complete. A long utterance does not mean the intention is complete and a short utterance does not represent an incomplete dialogue. Therefore, rule-based models rely on the surface features of utterances e.g., length of sentences and grammars, are not reliable on all situations in the Wait-or-Answer task.

⁵Without loss of generality, the prediction models in our PTD models: GRU-PTD, TextCNN-PTD, BERT-PTD all adopt an LSTM structure applied with the attention mechanism.

⁶These models are selected by the accuracy scores

TABLE VI

THE EFFECTS OF DIFFERENT TYPES OF PREDICTION MODELS ON THE TEXTCNN DECISION MODEL. THE PRE COLUMNS ARE THE BLEU SCORE OF PREDICTION MODELS GENERATED QUERIES OR ANSWERS. WoA (WAIT-OR-ANSWER) COLUMNS ARE DECISION MODEL’S (DEC) ACCURACY SCORE.

PRE. Type		MultiWOZ			DailyDialog			CCPE		
		PRE (BLEU)		DEC (Acc.)	PRE (BLEU)		DEC (Acc.)	PRE (BLEU)		DEC (Acc.)
		Agent	User	WoA	Agent	User	WoA	Agent	User	WoA
N/A		-	-	77.68	-	-	75.79	-	-	68.65
LSTM	Agent Prediction	11.77	0.80	80.04	4.51	0.61	76.37	15.71	0.00	70.04
	User Prediction	0.30	8.87		0.15	8.70		0.00	1.14	
LSTM + Attn.	Agent Prediction	12.47	0.72	80.75	19.19	0.60	79.02	23.86	0.00	73.32
	User Prediction	0.24	9.71		0.26	24.52		0.00	1.46	
LSTM w/ GLOVE + Attn.	Agent Prediction	13.37	0.67	80.38	19.01	0.67	78.56	19.56	0.00	71.62
	User Prediction	0.51	10.61		0.21	24.65		0.00	1.77	

Benefits Brought By PTD Framework From Table IV, we can see that our BERT-PTD model achieves the best performance on all datasets, e.g. BERT-PTD achieves 71.14 and 82.86 on After-Sales in terms of Accuracy and F1. Both are the highest performance on this dataset. Besides, the other two PTD models: GRU-PTD and TextCNN-PTD also significantly outperform their corresponding baselines, e.g. BERT-PTD achieves 71.98 on E-COMM in terms of Accuracy, which is the lowest score among all PTD models but still outperforms all other baseline models (the highest baseline model Bi-GRU achieves 71.96). Even the most rudimentary PTD model: GRU-PTD can beat all baselines (GRU, TextCNN, and BERT) in all these datasets.

The results on more evaluation metrics on public datasets in Table V also verify that our PTD framework is a more suitable choice for the Wait-or-Answer task. In most cases, the PTD models have better precision and recall scores, this shows PTD framework will decrease the possibility of presenting false positive and false negative results. We can also find that sometimes baseline will get higher recall but lower precision and F1, e.g. BERT gets 89.46 on MultiWOZ and TextCNN gets 78.49 on DailyDialog in terms of recall. This shows that baseline models can make false positive decisions more frequently. After analyzing some results samples, e.g. as shown in Table III, we can find that because of the supplement of generative models, PTD models are better at understanding the intuitions of users behind dialogue context. In contrast, the baseline models may make wrong decisions when meeting a situation with complete semantic and syntax dialogue context.

Above all, experimental results demonstrate that PTD model will improve the performance on the wait-or-answer task and decrease the possibility of making false positive and false negative decisions.

PTD’s Advantage on Small-scale Datasets One of the most crucial limits of the dialogue systems’ applications is the lack of high-quality datasets. In this case, we analyze the PTD’s performance on small-scale datasets. CCPE is relatively small-scale datasets, which consists of only 502 dialogues and significantly more average turns (9.7 in train set compared with 4.09 in DailyDialog and 6.32 in MultiWOZ). Besides, the numbers of positive (Agent Answer) samples and negative (Agent Wait) samples are more imbalanced. This makes it much more difficult to train a satisfactory model.

As shown in Table V, we can see that baselines: Bi-GRU,

TextCNNs, and BERT achieve accuracy scores of 67.53, 68.65, and 70.99. We can observe that baselines’ performance is significantly worse than that on large-scale datasets such as MultiWOZ. However, our PTD models: GRU-PTD, TextCNN-PTD, and BERT-PTD all achieve improved scores. As shown in Table VI, we can see that in small datasets, the prediction models’ BLEU scores are not worse than that on larger datasets like MultiWOZ. And all type of prediction models help decision models get significant improvement, and improvement is positively correlated with the prediction models’ performance. The LSTM with Attention-based prediction models gets the best generation scores and the best decision model results. In this case, we can conclude that on small-scale and imbalanced datasets, baselines have more difficulty in achieving satisfying results. However, our PTD models can learn more semantic information from the dialogue history with the user prediction model and the agent prediction model. In this way, our PTD models can achieve much more satisfying results than the baselines. This improvement can be explained by the fact that the prediction model in PTD can exploit the information in dialogue history more thoroughly.

Effects of Prediction Models in PTD Framework Another interesting issue about the PTD framework is the prediction models’ effects on the PTD framework. We investigate this issue by answering the following two questions:

(1). *Do the prediction models work as we expect?* We want to check out if the user prediction model can truly predict the user’s supplementary input and the agent prediction model can predict the dialogue system’s answer precisely. But, do they work as we expect? We conduct an experiment on the MultiWOZ dataset. As shown in Table VI, the LSTM based agent prediction model get the BLEU score at 11.77 on agent samples, in which the ground-truth is agents’ utterances, and the user prediction model gets the BLEU score at 0.3 on agent samples. Similar results are shown in other prediction models’ experiments. This phenomenon doesn’t mean that the user prediction model runs terrible. Actually, these results show that our user prediction model successfully behaves like a user. And its difficulty in generating agent utterance also meets our design. The example is also shown in Table III, in which the predicted agent utterance by user prediction model seems a high-quality fluent sentence and is also suitable for the scene. However, referring to the dialogue history, it is not a good choice since user in the last turn has said a sentence with

similar intention *what is the address for pizza hut in cherry hinton*, so the user prediction models' prediction *i am looking for a guesthouse* is not a good choice for decision, which means, the decision model prefers to answer user's utterance. From above we can conclude that contrasting results of the two prediction models work as we expect and help the decision model in *Wait-or-Answer* task.

(2). *Can better prediction model lead to better PTD models?* Another interesting question is that if the improvement of the prediction model can always boost the performance of PTD models. Take the DailyDialog as an example, we can see that with the enhancement of the attention mechanism and pre-trained GLOVE, the prediction models' performance increases⁷. The accuracy of the PTD models also increases: from 76.37 to 79.02, from 76.37 to 78.56. We can also observe the same phenomenon on MultiWOZ. From those results, we can conclude that there is a positive correlation between the performance of prediction models and the final performance of PTD models. From the above analysis, we can conclude that both the user prediction model and the agent prediction model can significantly enhance the decision models by predicting the dialogue interaction behavior. Recently, transformer based models show its effectiveness in generation tasks, e.g. BART [33], T5 [34] and GPT [35], they may further improve performance in PTD. We leave it as a future work.

VII. RELATED WORK

In this Section, we describe some work related to the wait-or-answer task, concentrating on response triggering recognition and predicting incomplete user utterance.

Coman et al. [57] implement an incremental Dialog State Tracker which is updated on a token basis to identify the point of maximal understanding in an ongoing utterance. DeVault et al. [58] propose a method for determining when a system has reached a point of maximal understanding of an ongoing user utterance to responsive overlap behaviors in dialogue systems, opening possibilities for systems to interrupt, acknowledge or complete a user's utterance while it is still in progress. Liu et al. [12] highlight the issues in previous models that in a supervised setting, the response timing information in the dialogues may be inaccurate in real life scenarios, such as customer service. MRTM address the inappropriate triggered responses problem very similar to wait-or-answer situation. MRTM is a self-supervised learning scheme leveraging the semantic matching relationships between the context and the response to train a semantic matching model and obtains the weights of the co-occurring utterances in the context through an asymmetrical self-attention mechanism. PTD adopts this similar settings that all training labels, e.g. turn tags and role tags are generated from dialogue history without human intervention. MRTM mainly concentrates on response selection models, while PTD explores to solve this task in a more general application scenario.

Some researches [7] have also investigated the incomplete utterance restoration and rewriting. Liu et al. [9] formulate

the incomplete utterance rewriting as a semantic segmentation task and propose a model for predicting the edit operations in parallel. Pan et al. [8] facilitate the study of incomplete utterance restoration for open-domain dialogue systems and propose a "pick-and-combine" model to restore the incomplete utterance from its context.

VIII. CONCLUSION

Conventional dialogue systems require that users must describe their intents in a single utterance, otherwise dialogue systems will answer immediately, which may cause misunderstanding or reply to the wrong question. In this paper, we explicitly define the aforementioned quandary as a novel Wait-or-Answer task. We further propose a framework named Predict-then-Decide (PTD) model to tackle with this Wait-or-Answer task. Our paper sheds light on the enhancement of the existing dialogue systems' ability to handle the wait-or-answer problem. We believe that the clearly-defined Wait-or-Answer task and PTD framework can provide an interesting topic for both academic and industrial NLP communities.

ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China (No. 62072399, No. 61402403, No. U19B2042), MoE Engineering Research Center of Digital Library, Chinese Knowledge Center for Engineering Sciences and Technology, Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies, and the Fundamental Research Funds for the Central Universities.

REFERENCES

- [1] C. Wu, A. Madotto, E. Hosseini-Asl, C. Xiong, R. Socher, and P. Fung, "Transferable multi-domain state generator for task-oriented dialogue systems," in *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*. Association for Computational Linguistics, 2019, pp. 808–819.
- [2] H. Le, R. Socher, and S. C. H. Hoi, "Non-autoregressive dialog state tracking," in *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [3] J. Li, W. Monroe, T. Shi, S. Jean, A. Ritter, and D. Jurafsky, "Adversarial learning for neural dialogue generation," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*. Association for Computational Linguistics, 2017, pp. 2157–2169.
- [4] T. Hasegawa, N. Kaji, N. Yoshinaga, and M. Toyoda, "Predicting and eliciting addressee's emotion in online dialogue," in *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, ACL 2013, 4-9 August 2013, Sofia, Bulgaria, Volume 1: Long Papers*. The Association for Computer Linguistics, 2013, pp. 964–972.
- [5] M. Ghazvininejad, C. Brockett, M. Chang, B. Dolan, J. Gao, W. Yih, and M. Galley, "A knowledge-grounded neural conversation model," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*. AAAI Press, 2018, pp. 5110–5117.
- [6] C. Wu, R. Socher, and C. Xiong, "Global-to-local memory pointer networks for task-oriented dialogue," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [7] M. Huang, F. Li, W. Zou, H. Zhang, and W. Zhang, "SARG: A novel semi autoregressive generator for multi-turn incomplete utterance restoration," *CoRR*, vol. abs/2008.01474, 2020.

⁷With the attention mechanism, the BLEU score increases from 4.51 to 19.19. With the attention mechanism and pre-trained GLOVE vector, the BLEU score increases from 4.51 to 19.01.

- [8] Z. F. Pan, K. Bai, Y. Wang, L. Zhou, and X. Liu, "Improving open-domain dialogue systems via multi-turn incomplete utterance restoration," 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 2019, Hong Kong, China, November 3-7, 2019*. Association for Computational Linguistics, 2019, pp. 1824–1833.
- [9] Q. Liu, B. Chen, J. Lou, B. Zhou, and D. Zhang, "Incomplete utterance rewriting as semantic segmentation," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*. Association for Computational Linguistics, 2020, pp. 2846–2857.
- [10] D. W. Oard, "Query by babbling: A research agenda," in *Proceedings of the First Workshop on Information and Knowledge Management for Developing Region*, ser. IKM4DR '12. New York, NY, USA: Association for Computing Machinery, 2012, p. 17–22.
- [11] R. Nordlie, "'user revelation'—a comparison of initial queries and ensuing question development in online searching and in human reference interactions," in *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, 1999, pp. 11–18.
- [12] C. Liu, J. Jiang, C. Xiong, Y. Yang, and J. Ye, "Towards building an intelligent chatbot for customer service: Learning to respond at the appropriate time," in *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*. ACM, 2020, pp. 3377–3385.
- [13] C. D. Manning and M. Eric, "A copy-augmented sequence-to-sequence architecture gives good performance on task-oriented dialogue," in *Proceedings of the 15th Conference of the European Chapter of the ACL, EACL 2017, Valencia, Spain, April 3-7, 2017, Volume 2: Short Papers*, 2017, pp. 468–473.
- [14] M. Eric, L. Krishnan, F. Charette, and C. D. Manning, "Key-value retrieval networks for task-oriented dialogue," in *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue, Saarbrücken, Germany, August 15-17, 2017*. Association for Computational Linguistics, 2017, pp. 37–49.
- [15] Z. Yan, N. Duan, P. Chen, M. Zhou, J. Zhou, and Z. Li, "Building task-oriented dialogue systems for online shopping," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, 2017, pp. 4618–4626.
- [16] R. Yan, "'chitty-chitty-chat bot': Deep learning for conversational AI," in *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*. ijcai.org, 2018, pp. 5520–5526.
- [17] A. W. Li, V. Jiang, S. Y. Feng, J. Sprague, W. Zhou, and J. Hoey, "Follow alice into the rabbit hole: Giving dialogue agents understanding of human level attributes," *arXiv preprint arXiv:1910.08293*, 2019.
- [18] B. Hancock, A. Bordes, P.-E. Mazare, and J. Weston, "Learning from dialogue after deployment: Feed yourself, chatbot!" in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3667–3684.
- [19] H. Chen, X. Liu, D. Yin, and J. Tang, "A survey on dialogue systems: Recent advances and new frontiers," *SIGKDD Explor. Newsl.*, vol. 19, no. 2, p. 25–35, Nov. 2017.
- [20] P. Budzianowski, T. Wen, B. Tseng, I. Casanueva, S. Ultes, O. Ramadan, and M. Gasic, "Multiwoz - A large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*. Association for Computational Linguistics, 2018, pp. 5016–5026.
- [21] Z. Lin, X. Huang, F. Ji, H. Chen, and Y. Zhang, "Task-oriented conversation generation using heterogeneous memory networks," 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 2019, Hong Kong, China, November 3-7, 2019*. Association for Computational Linguistics, 2019, pp. 4557–4566.
- [22] Z. Yan, N. Duan, P. Chen, M. Zhou, J. Zhou, and Z. Li, "Building task-oriented dialogue systems for online shopping," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*. AAAI Press, 2017, pp. 4618–4626.
- [23] Z. C. Lipton, X. Li, J. Gao, L. Li, F. Ahmed, and L. Deng, "Bbq-networks: Efficient exploration in deep reinforcement learning for task-oriented dialogue systems," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*. AAAI Press, 2018, pp. 5237–5244.
- [24] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, 2014, pp. 3104–3112.
- [25] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *Proceedings of the EMNLP 2015, Lisbon, Portugal*, 2015, pp. 1412–1421.
- [26] S. Liu, H. Chen, Z. Ren, Y. Feng, Q. Liu, and D. Yin, "Knowledge diffusion for neural dialogue generation," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*. Association for Computational Linguistics, 2018, pp. 1489–1498.
- [27] X. Zhao, W. Wu, C. Tao, C. Xu, D. Zhao, and R. Yan, "Low-resource knowledge-grounded dialogue generation," in *8th International Conference on Learning Representations, ICLR 2020*, 2020.
- [28] J. Gao, M. Galley, and L. Li, "Neural approaches to conversational AI," in *Proceedings of ACL 2018, Melbourne, Australia, July 15-20, 2018, Tutorial Abstracts*, Y. Artzi and J. Eisenstein, Eds. Association for Computational Linguistics, 2018, pp. 2–7.
- [29] T. Wen, M. Gasic, N. Mrksic, P. Su, D. Vandyke, and S. J. Young, "Semantically conditioned lstm-based natural language generation for spoken dialogue systems," pp. 1711–1721, 2015.
- [30] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, p. 1735–1780, Nov. 1997.
- [31] H. Chen, Z. Ren, J. Tang, Y. E. Zhao, and D. Yin, "Hierarchical variational memory network for dialogue generation," in *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*, 2018, pp. 1653–1662.
- [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, E. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [33] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "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*, 2020, pp. 7871–7880.
- [34] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.
- [35] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901.
- [36] M. Jiang, Y. Liang, X. Feng, X. Fan, Z. Pei, Y. Xue, and R. Guan, "Text classification based on deep belief network and softmax regression," *Neural Computing and Applications*, vol. 29, no. 1, pp. 61–70, 2018.
- [37] K. Kowsari, D. E. Brown, M. Heidarysafa, K. J. Meimandi, M. S. Gerber, and L. E. Barnes, "Hdltex: Hierarchical deep learning for text classification," in *16th IEEE International Conference on Machine Learning and Applications, ICMLA 2017, Cancun, Mexico, December 18-21, 2017*. IEEE, 2017, pp. 364–371.
- [38] S. Lai, L. Xu, K. Liu, and J. Zhao, "Recurrent convolutional neural networks for text classification," in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*. AAAI Press, 2015, pp. 2267–2273.
- [39] K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, 2019.
- [40] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola, and E. Hovy, "Hierarchical attention networks for document classification," in *Proceedings of the 2016 conference of the North American chapter of the ACL: human language technologies*, 2016, pp. 1480–1489.
- [41] L. M. Manevitz and M. Yousef, "One-class svms for document classification," *Journal of machine Learning research*, vol. 2, no. Dec, pp. 139–154, 2001.
- [42] A. Komninos and S. Manandhar, "Dependency based embeddings for sentence classification tasks," in *Proceedings of the 2016 conference of*

the North American chapter of the ACL: human language technologies, 2016, pp. 1490–1500.

- [43] R. Xia and Z. Ding, “Emotion-cause pair extraction: A new task to emotion analysis in texts,” in *Proceedings of the 57th ACL*. Florence, Italy: ACL, Jul. 2019, pp. 1003–1012.
- [44] Z. Song, X. Zheng, L. Liu, M. Xu, and X. Huang, “Generating responses with a specific emotion in dialog,” in *Proceedings of the 57th ACL*. Florence, Italy: ACL, Jul. 2019, pp. 3685–3695.
- [45] A. M. Hoyle, L. Wolf-Sonkin, H. Wallach, I. Augenstein, and R. Cotterell, “Unsupervised discovery of gendered language through latent-variable modeling,” in *Proceedings of the 57th ACL*. Florence, Italy: ACL, Jul. 2019, pp. 1706–1716.
- [46] G. Yan, L. Fan, Q. Li, H. Liu, X. Zhang, X.-M. Wu, and A. Y. Lam, “Unknown intent detection using Gaussian mixture model with an application to zero-shot intent classification,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics, Jul. 2020, pp. 1050–1060.
- [47] Y. Kim, “Convolutional neural networks for sentence classification,” in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*. ACL, 2014, pp. 1746–1751.
- [48] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” pp. 4171–4186, 2019.
- [49] D. Ha and J. Schmidhuber, “Recurrent world models facilitate policy evolution,” in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, 2018, pp. 2455–2467.
- [50] B. Peng, X. Li, J. Gao, J. Liu, and K.-F. Wong, “Deep dyna-q: Integrating planning for task-completion dialogue policy learning,” in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2182–2192.
- [51] M. Henderson, B. Thomson, and J. D. Williams, “The second dialog state tracking challenge,” in *Proceedings of the 15th SIGDIAL*. Philadelphia, PA, U.S.A.: ACL, Jun. 2014, pp. 263–272.
- [52] Y. Li, H. Su, X. Shen, W. Li, Z. Cao, and S. Niu, “DailyDialog: A manually labelled multi-turn dialogue dataset,” in *Proceedings of the 8th IJCNLP*. Taipei, Taiwan: Asian Federation of Natural Language Processing, Nov. 2017, pp. 986–995.
- [53] F. Radlinski, K. Balog, B. Byrne, and K. Krishnamoorthi, “Coached conversational preference elicitation: A case study in understanding movie preferences,” in *Proceedings of the 20th Annual SIGDial Meeting on Discourse and Dialogue*. Stockholm, Sweden: Association for Computational Linguistics, Sep. 2019, pp. 353–360.
- [54] K. Papineni, S. Roukos, T. Ward, and W. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th ACL, July 6-12, 2002, Philadelphia, PA, USA.*, 2002, pp. 311–318.
- [55] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, “Empirical evaluation of gated recurrent neural networks on sequence modeling,” in *NIPS 2014 Workshop on Deep Learning*, 2014.
- [56] G. Bekey and K. Y. Goldberg, *Neural Networks in Robotics*. Springer Science & Business Media, 1992, vol. 202.
- [57] A. C. Coman, K. Yoshino, Y. Murase, S. Nakamura, and G. Riccardi, “An incremental turn-taking model for task-oriented dialog systems,” in *Interspeech 2019, 20th Annual Conference of the International Speech Communication Association, Graz, Austria, 15-19 September 2019*. ISCA, 2019, pp. 4155–4159.
- [58] D. DeVault, K. Sagae, and D. R. Traum, “Can I finish? learning when to respond to incremental interpretation results in interactive dialogue,” in *Proceedings of the SIGDIAL 2009 Conference, The 10th Annual Meeting of the Special Interest Group on Discourse and Dialogue, 11-12 September 2009, London, UK*. The Association for Computer Linguistics, 2009, pp. 11–20.

Zehao Lin received his B.Sc. degree from Zhejiang University in 2015. He is currently working toward the Ph.D. degree in the College of Computer Science and Technology, Zhejiang University. His current research interests include natural Language processing, dialogue systems, and multimodal.

Shaobo Cui is a senior algorithm engineer in DAMO Academy, Alibaba Group. His current research interests include machine learning, mathematical optimization, natural language processing, and their applications. Previously, he received his M.Sc. degree from Tsinghua University.

Guodun Li received his B.Sc. degree from Hangzhou Dianzi University in 2020. He is currently working toward the M.Sc. degree in the College of Computer Science and Technology, Zhejiang University. His current research interests include natural language processing and their applications.

Xiaoming Kang received his M.Sc. degree from Fudan University in 2016. He is now a Senior Algorithm Engineer at Alibaba Group. His current research interests include natural language processing algorithms and their applications. He mainly focused on Dialogue System and Question answering.

Feng Ji received his B.Sc. degree from Tongji University and Ph.D. degree from Fudan University in 2003 and 2012 respectively. He was a senior algorithm engineer in Alibaba Group and now currently works in Tencent. His research interests include artificial intelligence, natural language understanding and generation, dialogue systems, information retrieval and recommendation systems.

Feng-Lin Li received his Ph.D. degree from University of Trento in 2016. He is currently working at DAMO Academy, Alibaba Group. His research interests include natural language processing and knowledge graph.

Zhongzhou Zhao received his M.Sc. degree in computer science from the Harbin Institute of Technology and the University of Pavia. He is currently a staff algorithm engineer at DAMO Academy, Alibaba Group. He is one of the early members of AliMe chatbot, and is currently responsible for the AI algorithm aspect of AliMe Avatar. His research interests include MRC, NLG, VQA and Multimodal Understanding.

Haiqing Chen received his B.S degree from Zhejiang University of technology in 2009. Now He is a senior algorithm expert of Alibaba cloud intelligence business group Damo Academic. His current research interests include natural language processing, question & answering, dialogue, machine learning and deep learning.

Yin Zhang received his Ph.D. degree from Zhejiang University in 2009. He is currently an associate professor with the College of Computer Science and Technology, Zhejiang University, China. His research interests include machine reading comprehension, question answering systems, multi-agent systems, and digital library.