

Proactive Interaction Framework for Intelligent Social Receptionist Robots

Yang Xue¹, Fan Wang¹, Hao Tian², Min Zhao³, Jiangyong Li¹, Haiqing Pan³ and Yueqiang Dong¹

Abstract—Proactive human-robot interaction (HRI) allows the receptionist robots to actively greet people and offer services based on vision, which has been found to improve acceptability and customer satisfaction. Existing approaches are either based on multi-stage decision processes or based on end-to-end decision models. However, the rule-based approaches require sedulous expert efforts and only handle minimal pre-defined scenarios. On the other hand, existing works with end-to-end models are limited to very general greetings or few behavior patterns (typically less than 10). To address those challenges, we propose a new end-to-end framework, the TransFormer with Visual Tokens for Human-Robot Interaction (TFVT-HRI)¹. The proposed framework extracts visual tokens of relative objects from an RGB camera first. To ensure the correct interpretation of the scenario, a transformer decision model is then employed to process the visual tokens, which is augmented with the temporal and spatial information. It predicts the appropriate action to take in each scenario and identifies the right target. Our data is collected from an in-service receptionist robot in an office building, which is then annotated by experts for appropriate proactive behavior. The action set includes 1000+ diverse patterns by combining language, emoji expression, and body motions. We compare our model with other SOTA end-to-end models on both offline test sets and online user experiments in realistic office building environments to validate this framework. It is demonstrated that the decision model achieves SOTA performance in action triggering and selection, resulting in more humanness and intelligence when compared with the previous reactive reception policies.

INTRODUCTION

Receptionist robots work in public areas such as lobbies and shopping malls, helping to guide visitors, post instructions, and answer questions [1], [2], [3]. As most existing receptionist robots can only passively answer users' calls, new visitors may have little motivation to initiate an interaction with the robot in many cases. This can be attributed to the lack of knowledge of the robot or low-expectations of its capability. While proactive reception robots can try to initiate interactions themselves, there are also risks of arousing antipathy in inappropriate proactive behavior cases.

To make proactive interaction more human-like and socially acceptable, it is desirable to correctly understand user intention. Previous works are based on empirical social rules and require the specification of scenarios at first, such as "A human is passing by," or "A human is waving a hand

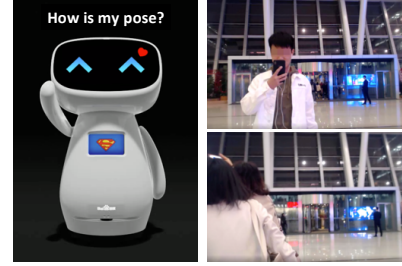


Fig. 1. **Left:** An example of proactive multi-modal behavior of the Xiaodu robot for users who are taking a photo ahead of it. Specifically, the robot is making a superman pose and saying "How is my pose". Here we use the avatar for better visualization. Check Fig. 3 for more information on Xiaodu robot. **Right:** Two different "photo-taking" scenarios in the served lobby.

towards the robot." To identify those scenarios, sophisticated sensor settings [4], [5], [6] are commonly required. Further, there are "micro signals" that are too complicated to be categorized. For instance, it is common for visitors to take a photo of the reception robot alone or take a group photo with the robot, shown in Fig. 1.

While prior researchers also propose to use an end-to-end framework in generating proactive greetings [7], [8], [9], their works suffer from the following limitations: First, the robot behaviors are restricted to very few types, typically less than 10, most of which are general words such as "Hello"[10], [11], "Excuse me", "May I Help You" [4]. Although those words can "break the ice" to some extent, they provide little information and shows no understanding of the scenario. Second, it is desirable to reason over time instead of depending on one single frame for proactive HRI. For instance, a visitor may hesitate and look around for a long time, which is a vital sign of seeking help. Third, the robot may be required to target a specific person or a particular group of persons instead of the scenario, especially when multiple persons are in sight, while most end-to-end decision-making is scenario-based.

To address aforementioned challenges, we propose a new end-to-end framework namely TransFormer with Visual Tokens for Human-Robot Interaction (TFVT-HRI). It has the following unique features: First, to reason over time while targeting specific persons, we use a *visual token extractor* to turn image-level signals into object-level tokens. Second, we use consecutive frames of a video clip, along with the temporal and spatial information of each token. Third, we use a transformer model to process the tokens. The attention over visual tokens enhances information exchange among different objects and different times to better interpret the

This work is under review of ICRA 2021

¹Yang Xue, Fan Wang, Jiangyong Li and Yueqiang Dong are from Baidu Natural Language Processing Department (e-mail: {xueyang02, wangfan04, lijyong01, dongyueqiang}@baidu.com;)

²Hao Tian is from Baidu Research, (e-mail: tianhao@baidu.com)

³Min Zhao and Haiqing Pan are from Baidu AI Interaction Design Lab, (e-mail: {zhaomin04, panhaiqing}@baidu.com)

¹Source code repository: <https://github.com/PaddlePaddle/PaddleRobotics>.

scenario. Fourth, to improve the generalization ability of our model, we employ natural language pretraining [12] as well as image pretraining in the representation layers. To train the decision model, we collect a rich dataset from in-service reception robots in office-building lobbies, where our staff imitate the daily visitors by implying different interaction intentions. We then invite experts to review those videos and label appropriate proactive behaviors. We collect over 1000+ multi-modal actions composed of language, emoji expression, and body movements. The action set includes not only simple greetings such as “Good morning, miss”, but also scenario-specific reactions such as “How is my pose” (Fig. 1), “Are you interested in playing a game with me?” and “Are you looking for some places?” It not only permits more human-like interactions but is also able to guide the visitors to further communications in depth.

To validate the proposed method, we first compare TFVT-HRT with the state-of-the-art action recognition method, fine-tuned R(2+1)D [13] model on the test set of our collected data, which shows that TFVT-HRT pushes on to frontier in this task. We then conduct a user experiment on the *Xiaodu* robot (Fig. 3), a receptionist robot serving in office buildings, museums, etc. It is shown that TFVT-HRT can generate human-like behaviors showing a deep understanding of the scenario and achieves a substantially higher score in overall comfort, naturalness, friendliness, and intelligence than prior wake word-based HRI system of the *Xiaodu* robot.

RELATED WORK

Social Rule based Systems. A large bunch of the previous works is based on social rules to initiate interaction with a human. Some of them make a straightforward classification of two classes: whether the user has the intention of interaction [4], [6]. The robot behavior includes a single pattern greeting only. Bergstrom et al. [5] utilize a laser range finder to collect the motion of the visitors around the robot, by which the visitors are classified into four groups depending on their trajectories, revealing the different level of interest in the robot. The robot behaviors are then specified based on the categorization. Heenan et al. [11] derive action set from Kendon’s theory of conducting interaction [14]. The classification is based on the user’s distance and head orientation, and the robot action is derived from Hall’s proxemics theory [15], which is also suggested by Zhao et al. [10]. Further, Zhao et al. [10] extend proactive HRI to progressive initiations including visual observation and eye contact. Those approaches inevitably require the pre-definition of scenarios.

Imitation Learning for Human-Robot Interaction. Imitation learning is frequently employed to enable robots to learn complex human skills. For HRI, it has been used to instruct a robot to express emotional body language [8], and teach a robot to play the role of a travel agent [9], showing that learning from human demonstrations can enable a robot to learn communication skills. However, imitation learning has not been widely used in the proactive HRI task yet.

Human Intention and Video Action Recognition. To understand the human intention, besides tracking the motion, speed, etc. of human pedestrians with laser and other sensors [4], [6], it is also possible to use video clips or images. The recently proposed pre-trained R(2+1)D model [13] aims to the elaborate classification of human action, which is a competitive approach in proactive action generation. However, to generate correct behaviors we still need an exhaustive specification of reacting rules after classification. Further, in the case of multiple visitors, designing rules can be too complicated. We provide a detailed comparison of our method with action recognition methods in our experiments.

Metrics to Evaluate Interaction Initiation. In previous studies, proactive HRI is mainly evaluated by system performance and users’ subjective evaluations. For system performances, Liao et al. [16] and Rashed [17] use the success rate of initiation of interaction; Shi et al. [18] evaluate the recognition accuracy of participants state. For users’ subjective evaluations, Zhao et al. [10] and Ozaki et al. [6] utilize questionnaire to study users’ experiences. Liao et al. [16] consider the feedbacks on satisfaction, likability, perceived interruption, perceived hedge, and social-agent orientation. Bergstrom et al. [5] and Shi et al. [18] mainly focus on the naturalness or appropriateness of the robot behavior. In this work, we will evaluate our framework from both offline datasets and online questionnaires.

METHODOLOGY

The proposed framework takes sequential RGB frames as the inputs and decides whether to initiate an interaction, the target to interact, and the action selection. The overview of the framework is illustrated in Fig. 2, which is composed of three modules: visual token extractor (Fig. 2(a)); multi-modal action encoder (Fig. 2(c)); and transformer-based decision model (Fig. 2(b)).

Visual Token Extractor

We denote the input image flow as $I^{(1)}, I^{(2)}, \dots, I^{(t)}$, where we use superscript (t) to represent the t -th frame. In this module, we first extract the region of interests by YOLOv4 [19]. YOLOv4 is a powerful object-detection model that detects a wide range of classes of objects. In our case, we are mainly interested in 6 categories that may be connected to human identity and action recognition in our scenarios, including “person”, “backpack”, “handbag”, “suitcase”, “tie”, and “cell phone”, while other categories are neglected.

To reduce the computation cost, we re-use the output of the backbone of YOLOv4, the CSPDarkNet53 [20], as the feature extractor. As the bounding boxes are in different sizes, we apply RoIAlign pooling to normalize the sizes of feature maps, inspired by the MaskRCNN [21], as well as global average pooling (GAP) to reduce the feature representation dimension. Through this process we are able to represent each pedestrian i in the image $I^{(t)}$ with a 512-dimensional feature vector $v_i^{(t)}$. We also add the embedding vector of 2D position information [22] and object classification ID, which

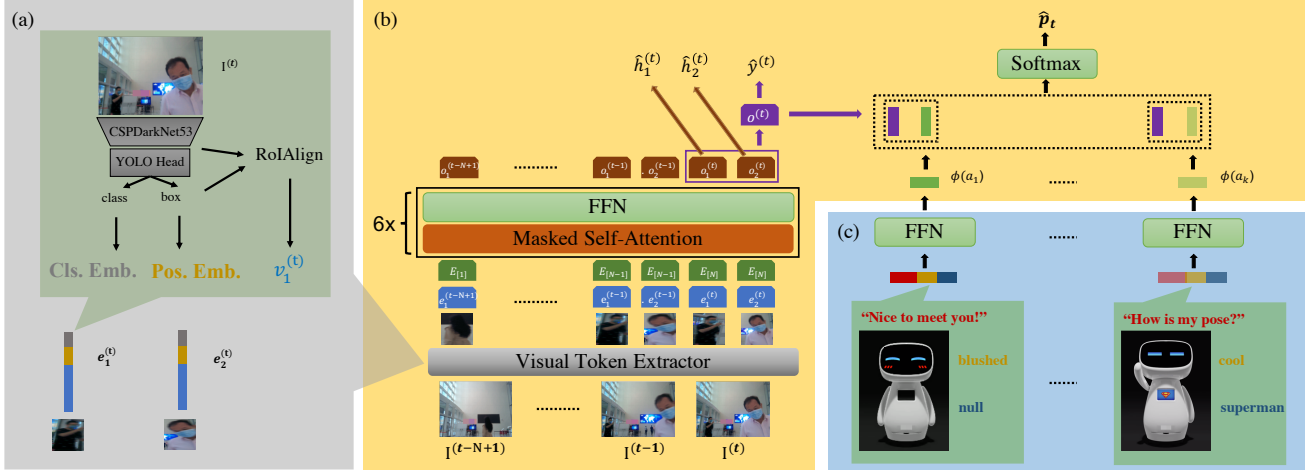


Fig. 2. Illustration of the TFVT-HRI framework, which is composed of three modules: (a) Visual token extractor that extracts visual tokens of objects in one frame. (b) Transformer-based decision model which encodes video clip to predict proactive behaviors (c) Multi-modal action encoder that encodes the natural language, expression, and body motion to form a representation vector.

gives the final visual token of (\oplus denotes “concatenation” and EMB represents “embedding vector”)

$$e_i^{(t)} = v_i^{(t)} \oplus EMB(POS(i, t)) \oplus EMB(CLASS(i, t)). \quad (1)$$

Multi-modal Action Encoder

The action is denoted as a combination of an utterance (u), an emoji facial expression (f), and a body motion (m). We specify a static set of distinct actions $a_k = (u_k, f_k, m_k), k \in [1, K]$ by collecting the actions from the annotated dataset. To better represent the action in semantic space, we utilize an open-source natural language pre-training model ERNIE [12] to process each utterance into a vector. Then facial expression and the motion are limited to several pre-defined patterns, such as “smiling” and “hand-shaking” (Fig. 3). For each type of body motion and expression, we use an embedding vector. The representation of three modalities are concatenated and further processed by an additional feed-forward neural layer (FFN) to yield the representation for k -th action, denoted by

$$\phi(a_k) = FFN(ERNIE(u_k) \oplus EMB(f_k) \oplus EMB(m_k)). \quad (2)$$

Transformer-based Decision Model

To enable the decision model to reason over time, normally, we could try to track each individual and encode the trajectory. Unfortunately, tracking one individual requires registration of visual tokens across different frames, which can introduce additional error to the decision. To this end, instead of tracking each individual, we present a transformer decision model. Transformer is widely used in natural language processing [23] by allowing the tokens to fully interact with each other through self-attention. In our case, it is not only able to capture the trajectory of each pedestrian by attending over time, but also possible to capture the interaction between the target and other individuals, e.g., talking with another person, holding a cell phone, carrying a suitcase. Also, to preserve consistency for training and inference, for each visual token $e_i^{(t)}$ that we consider, we allow it to

attend to $e_i^{(t')}$ only when $t' \leq t$, which implies masked self-attention in transformer blocks (MTRN for short). In our case, we use 6 transformer blocks.

For each frame $I^{(t)}$ received, we use visual token extractor to acquire a list of visual tokens $e_1^{(t)}, e_2^{(t)}, \dots$, from which we consider the top- M ($M = 20$ in the experiments) visual tokens. The priority of the selection goes as follows: the “person” class is selected in the first place; the larger the bounding box, the higher priority; in case there are less than M visual tokens, we add padding terms. In every time step, we use the latest N frames as the input, which gives $M \times N$ tokens in one sequence. For each token, we further add an embedding vector of the relative frame ID $E_{[i]}$. The encoding process can be represented by

$$\begin{aligned} o_1^{(t-N+1)}, \dots, o_M^{(t-N+1)}, \dots, o_M^{(t)} = & MTRN(\\ & e_1^{(t-N+1)} \oplus E_{[1]}, \dots, e_M^{(t-N+1)} \oplus E_{[1]}, \\ & e_1^{(t-N+2)} \oplus E_{[2]}, \dots, e_M^{(t-N+2)} \oplus E_{[2]}, \\ & \dots, \\ & e_1^{(t)} \oplus E_{[N]}, \dots, e_M^{(t)} \oplus E_{[N]}) \end{aligned} \quad (3)$$

Starting from Eq. 3, the prediction of the decision model is three-fold: 1. decide whether to initiate an interaction, which depend on $\hat{y}^{(t)}$; 2. predict the target to interact, which is dependent on $\hat{h}_i^{(t)}$; 3. select the action to be taken, which is dependent on $\hat{p}^{(t)}(a_k)$. We first take the max-pooling of all the positions in the last frame to represent the scenario, which is defined as

$$o^{(t)} = \text{MaxPooling}(o_1^{(t)}, o_2^{(t)}, \dots, o_M^{(t)}), \quad (4)$$

we then specify the three-fold output of the decision model as follows:

$$\hat{y}^{(t)} = \sigma(W_y \cdot o^{(t)} + b_y) \quad (5)$$

$$\hat{h}_i^{(t)} = \sigma(W_h \cdot o_i^{(t)} + b_h) \quad (6)$$

$$\begin{aligned} \hat{p}^{(t)}(a_1), \dots, \hat{p}^{(t)}(a_k) = & \text{Softmax}(\phi(a_1) \cdot o^{(t)}, \dots, \\ & \phi(a_k) \cdot o^{(t)}) \end{aligned} \quad (7)$$

Correspondingly, the annotation in training data include three parts: interaction trigger $y^{(t)} = 0/1$, where $y^{(t)} = 1$ indicate a proper opportunity to initiate the interaction; interaction target indicator $h_i^{(t)} = 0/1$, where $h_i^{(t)} = 1$ implies that the visual token is one of the target to be interacted; the action selection $\mathbf{I}_a^{(t)}$ which is a one-hot vector of length K . We use cross-entropy loss \mathcal{L}_{ce} for each of them, giving the loss function of Eq. 8.

$$\begin{aligned} \mathcal{L} = & \sum_t \mathcal{L}_{ce}(y^{(t)}, \hat{y}^{(t)}) + \sum_t y^{(t)} \cdot \mathcal{L}_{ce}(\mathbf{I}_a^{(t)}, \hat{\mathbf{p}}^{(t)}) \\ & + \sum_t y^{(t)} \sum_i \mathcal{L}_{ce}(h_i^{(t)}, \hat{h}_i^{(t)}) \end{aligned} \quad (8)$$

For inference, we use weighted sampling on $\hat{\mathbf{p}}^{(t)}$ for action selection; for target filtering, we first remove the visual tokens of classes other than “person”, then we use a threshold H , such that $h_i^{(t)} > H$ indicate a valid target. At executing the actions, the robot turns towards the centroid of all the valid targets. We argue that turning to the interaction target is essential to the user experience.

EXPERIMENTAL RESULTS

Data Collection

Hours-long videos were collected from two working *Xiaodu* robots (Fig. 3) in the office lobbies, and then they were annotated and processed into 5-second-long video clips as the training and test set. The office lobbies for our data collection are significantly different in the light condition. As shown in Fig. 3 (b) and (c), the B-lobby has a worse light condition than the A-lobby because of the backlighting environment. This makes the vision-based approaches more challenging in the B-lobby. After collected hours-long videos from these two lobbies, we used the following method to preprocess, annotate, and post-process them:

- 1) The raw videos were preprocessed by a multiple object tracking model [24], which marks the bounding boxes of tracked persons with a unique ID for each. The experts label suitable targets for initiating interactions with track IDs.
- 2) Experts watched the preprocessed videos and selected the suitable timestamps for initiating interactions. For each selected timestamp, the expert annotated the tracking IDs of targeting interaction objects, chose the facial expression and the body motion from a list that the *Xiaodu* robot permits, and filled the textbox with a proper utterance for greeting.
- 3) Annotated videos were post-processed into 5-second-long video clips in which positive examples were from experts’ annotations, and negative examples were randomly sampled from segments without annotations then filtered out in-interaction segments.

Statistically, for the training set, we collected and labeled 3900 positive cases (i.e., trigger label is 1) with 1000+ unique multi-modal actions (triplet of utterance to speak, emoji to display, and body action to execute), in which 1720 cases are from B-lobby, and 2180 cases are from A-lobby. For the

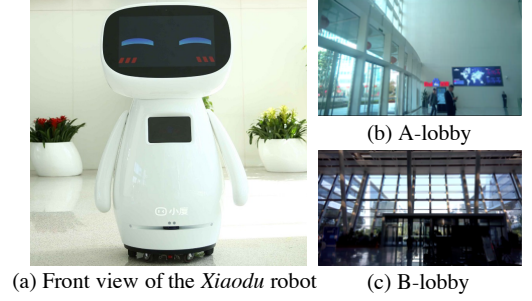


Fig. 3. (a) The *Xiaodu* robot, a receptionist robot that we used to conduct all the experiments. An RGB camera is on its head to record live video that serves as input to the model. Execution of the multi-modal action includes three modalities: its stereo speaker can play synthesized TTS given the utterance; a LED screen on the face can display over 30 emojis; for motions, it can perform hand-shaking, hand-waving, hug, moving, etc. (b) The A-lobby from the perspective of the robot. (c) The B-lobby from the perspective of the robot.

test set, from 12-hour videos, we collected and labeled 74 positive cases in A-lobby and 91 positive cases in B-lobby. However, when evaluating the performance on the test set, we used the full 135k+ negative cases instead of sampling to verify the model sufficiently.

Offline Evaluation

To evaluate the proactive HRI framework, it is desirable to verify whether the system predicts the correct timestamp to initiate interaction and whether the action selection is proper. However, although each positive case in the test set has a labeled action, we found that multiple actions can be suitable for one scenario. Therefore, evaluating the action selection is extremely challenging. In this part, we evaluated the capability of predicting proper triggering time only. The user experience research in the next part will evaluate whether the selected action is proper.

We used the state-of-the-art model called R(2+1)D [13] in a similar video action recognition task as our baseline. The R(2+1)D model was pre-trained on a large video dataset from social media. We fine-tuned it on our dataset in an end-to-end manner to predict the multi-modal action ID. To enable a fair comparison between our Transformer setting and R(2+1)D, we add a “NULL” action in the action set, such that selecting the null action is equal to not triggering an initiation. For R(2+1)D, the output is a $(K + 1)$ -way classification on the action space. In this way, the R(2+1)D can predict the triggering time without predicting $\hat{y}^{(t)}$. For our model, we also add a “NULL” action in the action set, this leads to three different inference modes:

- **Trigger-Only** By removing the null action, We rely only on $\hat{y}^{(t)}$ to decide whether to initiate interaction.
- **Actor-Only** By setting $\hat{y}^{(t)} = 1$ identically, we rely only on $\hat{\mathbf{p}}^{(t)}$ to select the null action to serve as the trigger.
- **Trigger-Actor** Initiation is only triggered when we both have $\hat{y}^{(t)} > H$ and avoid selecting the null action.

To find out the contribution of different elements in the visual token, we designed ablation studies by removing the RoIAlign pooling feature ($v_i^{(t)}$), the position embedding,

and the classification embedding in turn. We reported the performance of successful proactive interaction (proactively and appropriately interact with the user) in terms of precision and recall. The performance of the trigger is related to the threshold, so we include two more metrics: average precision (AP) and average recall (AR) for better evaluation. Furthermore, we use the F1 score to demonstrate overall performance. All the results are shown in Table I. Note that we select the threshold with the maximum F1 score to report the precision and recall. In the Trigger-Actor mode, we also use this threshold for the trigger. As we can see, the results show that the proposed transformer model outperforms the state-of-the-art end-to-end R(2+1)D model. Besides, the ablation study shows that the most important factor is the information from the RoIAlign pooling feature ($v_i^{(t)}$). Without this information, the framework suffers the largest F1 score drop.

User Experience Research

Because the data were collected under a passive HRI system, the distribution of the test set was different from that of the online experiments. To figure out whether the proposed framework performs well on a real robot, we implement this user experience research on a real robot in the A and B lobbies, comparing the existing wake word-based interaction system of the *Xiaodu* robot. On the Jetson AGX Xavier, the proposed TFVT-HRI model runs at 6.25fps, whereas the R(2+1)D model runs at 1.89fps. Since the R(2+1)D model suffers from 3x computation latency, and its performance is worse than the proposed framework, we do not deploy it on the real robot for this user experience research.

Participants. We recruited 30 new employees who have never interacted with the *Xiaodu* robot, but 26 of them used other similar AI devices (they share the same wake word “X” to start an interaction). Among the participants, there are 16 male and 14 female; their ages range from 21 to 35 years old ($M=27.1$, $SD=3.84$). All participants reported normal or corrected to normal vision and normal hearing. All participants volunteered to participate in the study and agreed to make audio and video recordings of the research process. At the end of the study, all participants were given

appropriate compensation.

Design. The experiment adopted a between-subjects design. Participants were randomly assigned to one of the two groups, the experimental group and the control group. In the experimental group, the robot was deployed with the proposed framework and would proactively initiate interaction. In the control group, the robot was waiting reactively for the wake word.

The dependent variables included the objective factor (success rate) and subjective factors (emotions, attitudes). For emotions, we utilize the Self-Assessment Manikin (SAM) technique [25] because of its high correlation with psychological response [26] and focus on the metrics of valence and arousal. For attitudes, the participants are asked to evaluate the level of overall comfort, naturalness, friendliness, and intelligence using a 7-point Likert questionnaire (from 1-7, higher scores indicated stronger agreement). The present study only focused on the initiating process of HRI. Participants were also instructed to give ratings based on their experience of initiating interaction.

A success case means that the participant use at least one function of the *Xiaodu* robot. For example, after deploying the proposed framework, the robot may greet then suggest one function to the participant. If the participant noticed it and response, this is a success case. For the existing wake word-based HRI, the participant should try to wake, the robot then following text instruction on the screen to make a success interaction case. A failure case means that the participant neither responds to the proactive calling of the robot nor wakes the robot and follows the instruction.

For the objective metric, success rate, 100% of participants had success interactions when using the proposed framework. Whereas the success rate was 80% on the existing wake word-based HRI system. For emotions, the results of descriptive statistics are shown in Table II. Independent-Samples T Test indicated that there was no significant difference in valence ($t(28) = 1.218$, $p = 0.233 > 0.05$) and arousal ($t(28) = 1.906$, $p = 0.067 > 0.05$). We combined the scatter diagram of the valence-arousal plain and affect words coordinates in Fig. 4. It informed us that two groups both have a positive emotion, EXCITED. The Levene’s Test indicated

TABLE I
PERFORMANCE OF MODEL VARIANTS AND RESULTS OF ABLATION STUDY.

Methods		A-Lobby					B-Lobby				
		Precision	Recall	AP	AR	F1	Precision	Recall	AP	AR	F1
R(2+1)D+ig65m		0.383	0.837	-	-	0.526	0.267	0.859	-	-	0.407
TFVT-HRI	Trigger-Only	0.905	0.851	0.718	0.927	0.877	0.583	0.778	0.418	0.949	0.667
	Actor-Only	0.894	0.868	-	-	0.881	0.485	0.876	-	-	0.624
	Trigger-Actor	0.905	0.851	-	-	0.877	0.585	0.775	-	-	0.667
TFVT-HRI w/o $v_i^{(t)}$	Trigger-Only	0.717	0.930	0.713	0.857	0.810	0.382	0.989	0.434	0.894	0.551
	Actor-Only	1.0	0.121	-	-	0.216	0.800	0.148	-	-	0.250
	Trigger-Actor	1.0	0.061	-	-	0.114	1.0	0.041	-	-	0.078
TFVT-HRI w/o Pos. Emb.	Trigger-Only	0.782	0.910	0.755	0.893	0.841	0.420	0.941	0.399	0.929	0.581
	Actor-Only	0.831	0.790	-	-	0.810	0.424	0.911	-	-	0.579
	Trigger-Actor	0.864	0.691	-	-	0.768	0.405	0.757	-	-	0.527
TFVT-HRI w/o Cls. Emb.	Trigger-Only	0.855	0.810	0.786	0.738	0.832	0.411	0.966	0.420	0.861	0.576
	Actor-Only	0.700	0.212	-	-	0.326	0.314	0.800	-	-	0.451
	Trigger-Actor	0.714	0.152	-	-	0.251	0.337	0.648	-	-	0.443

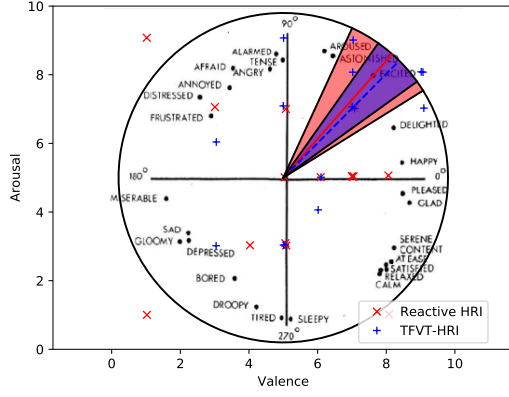


Fig. 4. A scatter diagram of the effects on the circumplex model [27]. In the diagram, the blue line and red line are the respective mean of the affect (i.e., the overall experience of feeling or emotion) on the proposed TFVT-HRI framework and the reactive HRI system; the blue and red sector indicate the variances of our model and passive method respectively. For checking users’ emotions easily, we merge the scatter diagram and affect words coordinates together. Clearly, both approaches arouse emotions close to EXCITED.

that all the attitudes variables except for the intelligence score met the assumptions of variance homogeneity at $p > 0.05$. Thus, we used the Independent-Samples T Test for Overall Comfort Level, Naturalness, Friendliness scores, and T’ Test for intelligence score. Compared to the control group, the experimental group reported significantly higher scores in Overall Comfort Level ($t(28) = 2.141, p = 0.041 < 0.05$), Naturalness ($t(28) = 2.354, p = 0.026 < 0.05$), Friendliness ($t(28) = 2.705, p = 0.012 < 0.05$), and Intelligence ($t'(24.679) = 2.225, p = 0.035 < 0.05$).

In summary, the proposed framework has a higher success rate of initiating an interaction than the existing wake word-based HRI system of the *Xiaodu* robot due to that the proposed framework can capture the intention of the participants and trigger multi-modal interaction signals (language, facial expression, and body language). However, no significant difference was found in the emotions of the two groups. The robot’s characters may have some positive impacts on users’ emotions, such as its cute appearance (Fig. 3), and child-like voices. However, the proposed framework leads to significant differences in feelings of overall comfort, naturalness, friendliness, and intelligence.

Showcases and the Discovered Social Rules

Some interaction cases are shown in Fig. 5. As we can see in Fig. 5 (a), when a group of visitors passed by, the *Xiaodu* robot raised its hand and greeted these people using voice. Here we emphasize that the *Xiaodu* robot started this action when the visitors were in the near field, matching the



Fig. 5. Showcases that the *Xiaodu* robot interacts with pedestrians following the commonly recognized social rules.

prior studies based on Hall’s proxemics theory [15]. In the second case (Fig. 5 (b)), a visitor were taking photos with the *Xiaodu* robot. The robot proactively blinked its eyes as responses. The third case (Fig. 5 (c)) shows that the *Xiaodu* robot successfully led the girl to a depth interaction. Overall, the TFVT-HRI framework substantially learned certain knowledge of social rules from the experts’ demonstrations.

CONCLUSIONS

In this work, we propose a new end-to-end framework for the proactive HRI task, TFVT-HRI. The proposed framework is featured with object-level modeling, reasoning over time and semantic space, and a relatively large multi-modal action space. Offline dataset validation and online user study are carried out. It is proved that the proposed model is better at the HRI task than other end-to-end models such as R(2+1)D, which is also proved to require less computational cost than R(2+1)D. The user study has revealed that the proposed framework achieves a substantially higher score in overall comfort, naturalness, friendliness, and intelligence than the existing wake word-based reactive HRI system on the *Xiaodu* robot. One possible extension of this work is to employ reinforcement learning online utilizing user feedback signals. Also, we could further extend the action space to continuous space to permit the robot to react with a broader range of behaviors, for which it is necessary to explore multi-modal generative models.

TABLE II
DESCRIPTIVE STATISTIC RESULTS OF USER QUESTIONNAIRES.

Methods	Participants	Emotion				Attitudes							
		Valence		Arousal		Comfort		Naturalness		Friendliness		Intelligence	
		M	SD	M	SD	M	SD	M	SD	M	SD	M	SD
TFVT-HRI	15	6.20	1.93	6.27	2.15	5.20	1.42	4.87	1.41	6.07	1.22	4.53	1.73
Reactive HRI	15	5.27	2.25	4.72	2.25	4.13	1.30	3.73	1.22	4.67	1.59	3.33	1.18

REFERENCES

- [1] R. Nisimura, T. Uchida, A. Lee, H. Saruwatari, K. Shikano, and Y. Matsumoto, "Aska: receptionist robot with speech dialogue system," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 2. IEEE, 2002, pp. 1314–1319.
- [2] T. Hashimoto, S. Hiramatsu, T. Tsuji, and H. Kobayashi, "Realization and evaluation of realistic nod with receptionist robot saya," in *ROMAN 2007-The 16th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 2007, pp. 326–331.
- [3] A. Niculescu, B. Van Dijk, A. Nijholt, and S. L. See, "The influence of voice pitch on the evaluation of a social robot receptionist," in *2011 International Conference on User Science and Engineering (i-USEr)*. IEEE, 2011, pp. 18–23.
- [4] Y. Kato, T. Kanda, and H. Ishiguro, "May i help you?-design of human-like polite approaching behavior," in *2015 10th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 2015, pp. 35–42.
- [5] N. Bergstrom, T. Kanda, T. Miyashita, H. Ishiguro, and N. Hagita, "Modeling of natural human-robot encounters," in *2008 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2008, pp. 2623–2629.
- [6] Y. Ozaki, T. Ishihara, N. Matsumura, T. Nunobiki, and T. Yamada, "Decision-making prediction for human-robot engagement between pedestrian and robot receptionist," in *2018 27th IEEE International Symposium on Robot and Human Interactive Communication (ROMAN)*. IEEE, 2018, pp. 208–215.
- [7] Y. Ozaki, T. Ishihara, N. Matsumura, and T. Nunobiki, "Can user-centered reinforcement learning allow a robot to attract passersby without causing discomfort?" in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 6986–6992.
- [8] N. T. V. Tuyen, S. Jeong, and N. Y. Chong, "Emotional bodily expressions for culturally competent robots through long term human-robot interaction," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 2008–2013.
- [9] M. Doering, D. F. Glas, and H. Ishiguro, "Modeling interaction structure for robot imitation learning of human social behavior," *IEEE Transactions on Human-Machine Systems*, vol. 49, no. 3, pp. 219–231, 2019.
- [10] M. Zhao, D. Li, Z. Wu, S. Li, X. Zhang, L. Ye, G. Zhou, and D. Guan, "Stepped warm-up—the progressive interaction approach for human-robot interaction in public," in *International Conference on Human-Computer Interaction*. Springer, 2019, pp. 309–327.
- [11] B. Heenan, S. Greenberg, S. Aghel-Manesh, and E. Sharlin, "Designing social greetings in human robot interaction," in *Proceedings of the 2014 conference on Designing interactive systems*, 2014, pp. 855–864.
- [12] Y. Sun, S. Wang, Y. Li, S. Feng, X. Chen, H. Zhang, X. Tian, D. Zhu, H. Tian, and H. Wu, "Ernie: Enhanced representation through knowledge integration," *arXiv preprint arXiv:1904.09223*, 2019.
- [13] D. Ghadiyaram, M. Feiszli, D. Tran, X. Yan, H. Wang, and D. Mahajan, "Large-scale weakly-supervised pre-training for video action recognition," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 038–12 047, 2019.
- [14] A. Kendon, *Conducting interaction: Patterns of behavior in focused encounters*. CUP Archive, 1990, vol. 7.
- [15] E. T. Hall, *The hidden dimension*. Garden City, NY: Doubleday, 1966, vol. 609.
- [16] Q. V. Liao, M. Davis, W. Geyer, M. Muller, and N. S. Shami, "What can you do? studying social-agent orientation and agent proactive interactions with an agent for employees," in *Proceedings of the 2016 ACM Conference on Designing Interactive Systems*, ser. DIS '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 264–275.
- [17] M. G. Rashed, "Observing people's behaviors in public spaces for initiating proactive human-robot interaction by social robots," Ph.D. dissertation, Graduate School of Science and Engineering, Saitama University, Japan, 2016.
- [18] C. Shi, M. Shiomi, T. Kanda, H. Ishiguro, and N. Hagita, "Measuring communication participation to initiate conversation in human–robot interaction," *International Journal of Social Robotics*, vol. 7, no. 5, pp. 889–910, 2015, we'd like to thank everyone who helped with this project. This research was supported by KAKENHI 25240042.
- [19] A. Bochkovskiy, C.-Y. Wang, and H. Liao, "Yolov4: Optimal speed and accuracy of object detection," *ArXiv*, vol. abs/2004.10934, 2020.
- [20] C.-Y. Wang, H. Liao, I.-H. Yeh, Y.-H. Wu, P.-Y. Chen, and J.-W. Hsieh, "Cspnet: A new backbone that can enhance learning capability of cnn," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1571–1580, 2020.
- [21] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask r-cnn," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [22] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," *ArXiv*, vol. abs/2005.12872, 2020.
- [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 5998–6008.
- [24] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," *2017 IEEE International Conference on Image Processing (ICIP)*, pp. 3645–3649, 2017.
- [25] M. M. Bradley and P. J. Lang, "Measuring emotion: The self-assessment manikin and the semantic differential," *Journal of Behavior Therapy and Experimental Psychiatry*, vol. 25, no. 1, pp. 49 – 59, 1994.
- [26] P. J. LANG, M. K. GREENWALD, M. M. BRADLEY, and A. O. HAMM, "Looking at pictures: Affective, facial, visceral, and behavioral reactions," *Psychophysiology*, vol. 30, no. 3, pp. 261–273, 1993.
- [27] J. Russell, "A circumplex model of affect," *Journal of Personality and Social Psychology*, vol. 39, pp. 1161–1178, 12 1980.