

# Service humanoid robotics: a novel interactive system based on bionic-companionship framework

Jiaji Yang<sup>Corresp., 1</sup>, Eysin Chew<sup>1</sup>, Pengcheng Liu<sup>2</sup>

<sup>1</sup> Cardiff School of Technologies, Cardiff Metropolitan University, Cardiff, Cardiff, United Kingdom

<sup>2</sup> The Department of Computer Science, University of York, York, United Kingdom

Corresponding Author: Jiaji Yang

Email address: JYang@cardiffmet.ac.uk

At present, industrial robotics focuses more on motion control and vision, whereas humanoid service robotics (HSRs) are increasingly being investigated and researched in the field of speech interaction. The problem and quality of human-robot interaction (HRI) has become a widely debated topic in academia. Especially when HSRs are applied in the hospitality industry, some researchers believe that the current HRI model is not well adapted to the complex social environment. HSRs generally lack the ability to accurately recognize human intentions and understand social scenarios. This study proposes a novel interactive framework suitable for HSRs. The proposed framework is grounded on the novel integration of Trevarthen's (2001) companionship theory and neural image captioning (NIC) generation algorithm. By integrating image-to-natural interactivity generation and communicating with the environment to better interact with the stakeholder, thereby changing from interaction to a bionic-companionship. Compared to previous research a novel interactive system is developed based on the bionic-companionship framework. The humanoid service robot was integrated with the system to conduct preliminary tests. The results show that the interactive system based on the bionic-companionship framework can help the service humanoid robot to effectively respond to changes in the interactive environment, for example give different responses to the same character in different scenes.

# Service Humanoid Robotics: A Novel Interactive System Based on Bionic-Companionship Framework

Jiaji Yang<sup>1</sup>, Eysin Chew<sup>1</sup>, Pengcheng Liu<sup>2</sup>

<sup>1</sup> EUREKA Robotics Lab, Cardiff School of Technologies, Cardiff Metropolitan University, Cardiff, CF5 2YB, United Kingdom.

<sup>2</sup> The Department of Computer Science, University of York, York YO10 5GH, United Kingdom.

Corresponding Author:

Jiaji Yang<sup>1</sup>

Western Avenue, Cardiff, CF5 2YB

Email address: JYang@cardiffmet.ac.uk

## Abstract

At present, industrial robotics focuses more on motion control and vision, whereas humanoid service robotics (HSRs) are increasingly being investigated and researched in the field of speech interaction. The problem and quality of human-robot interaction (HRI) has become a widely debated topic in academia. Especially when HSRs are applied in the hospitality industry, some researchers believe that the current HRI model is not well adapted to the complex social environment. HSRs generally lack the ability to accurately recognize human intentions and understand social scenarios. This study proposes a novel interactive framework suitable for HSRs. The proposed framework is grounded on the novel integration of Trevarthen's (2001) companionship theory and neural image captioning (NIC) generation algorithm. By integrating image-to-natural interactivity generation and communicating with the environment to better interact with the stakeholder, thereby changing from interaction to a bionic-companionship.

Compared to previous research a novel interactive system is developed based on the bionic-companionship framework. The humanoid service robot was integrated with the system to conduct preliminary tests. The results show that the interactive system based on the bionic-companionship framework can help the service humanoid robot to effectively respond to changes in the interactive environment, for example give different responses to the same character in different scenes.

## Introduction

Humanoid service robots (HSRs) have seen a sharp rise in adoption recently and are seen as one of the major technologies that will drive the service industries in the next decade (Harris et al., 2018). An increasing number of researchers are committed to investigating HSRs to help humans complete repetitive or high-risk service and interactive tasks such as serving patients with infectious diseases, delivering meals and so on. Delivery robots, concierge robots, and chat robots have been increasingly used by travel and hospitality companies (Ivanov, 2019). Although the contribution of these achievements mainly comes from the rapid development of robotics engineering, Ivanov et al. (2019) indicated that future research focus will gradually shift from robotics engineering to human-robot interaction (HRI), thus opening up interdisciplinary research directions for researchers.

In the early days, Fong et al. (2003) proposed that in order to make robots perform better, the robot needs to be able to use human skills (perception, cognition, etc.) and benefit from human advice and expertise. This means that robots that rely solely on self-determination have limitations in performing tasks. The authors further propose that the collaborative work between humans and robots will be able to break this constraint, and research on human-robot interaction has begun to emerge. Fong et al. (2003) believe that to build a collaborative control system and complete human-robot interaction, four key problems must be solved. (1) The robot must be able to detect limitations (what can be done and what humans can do), determine whether to seek help, and identify when it needs to be resolved. (2) The robot must be self-reliant and secure. (3) The system must support dialog. That is, robots and humans need to be able to communicate with each other effectively. However, dialog is restricted at present. Through collaborative control, dialog should be two-way and require a richer vocabulary. (4) The system must be adaptive. Although most of the current humanoid service robots already support dialog and can complete simple interactive tasks, as propounded in the research, such dialog in the present time remains limited and “inhuman.” In the process of interacting with robots, humans always determine the state of the robot (the position of the robot or the action the robot is doing) through vision, and then communicate with the robot through a dialog system. However, HSR cannot perform this yet as they do not seem to fully satisfy the two-way nature of dialog. Therefore, this research responds to the current gap and attempts to differ from the current HRI research. This research attempts to introduce deep learning into the existing dialog system of HSR, thus advancing the field.

With the continuous development of humanoid robots, more and more humanoid robots are used in the service industry, especially the hospitality industry. Human-Robot Interaction (HRI) has become a hot potato by more and more researchers (Yang & Chew, 2020). However, with the deepening of research, some researchers found that when humans interact with humanoid service robots (HSRs), humans hope that HSRs should have the ability and interest to interact with the dynamic thoughts and enthusiasm of the partner's relationship, and can recognize the environment, blended with what others think is meaningful and the emotions to express sympathy (Yang & Chew, 2020). This coincides with Trevarthen Companionship Theory (Trevarthen, 2001), so the concept of human robot companion (HRC) was proposed this research. The earlier concept of the robot companion is mentioned by Dautenhahn et al. (2005): HSRs need to have a high degree of awareness and sensitivity to social environment. Through the review of the above literatures, it is proposed to establish an interactive and companion framework for HSRs using deep learning and neural image caption generation, thus advance the current field of HSRs to tackle with bionic-interactive tasks of the service industry and further evolve from conventional HRI to Human and Robot Companion (HRC) (See Tab. 1).

This study proposes that the introduction of visual data into the current HRI model of HSRs enables HSRs to have a high level of sensitivity to the social dynamic environment while interacting with humans, thereby enhancing the current HRI model to HRC. With the continuous development of deep learning, some researchers have recently realized the transformation of static pictures or videos from conventional camera input into text descriptions (Li et al., 2020; Hu et al., 2020; Luo, 2019). This deep learning algorithm model is called neural image capturing (NIC). This research attempts to adapt and integrate NIC into HSRs and propose a novel framework (bionic-companionship framework) to enhance the traditional HRI experience. This framework aims to improve the current HRI interaction mode in the field of HSRs to a higher level of HRC (Yang & Chew, 2021). The bionics in this research refers to the humanoid service robot imitating all the tastes of life, trying to adapt to the seven emotions of ancient human nature (joy, anger, sadness, fear, love, disgust, liking) and six biological wills (life, death, eyes, ears, mouth, nose) (Chew et al., 2021). The system proposed in this study combines visual intelligence and Speech Intelligence, and imitates human behavior in social activities, which is in line with the concept of robot bionics proposed by researchers such as Chew et al (2021). Therefore, this study believes that the proposed system is a bionic system.

## Related works

With the continuous development of HRI research, industrial robots have been able to interact with humans accurately and self-adaptively. Some advanced control systems (Zhang et al., 2020) and algorithms (Tang et al., 2020) have been proposed as Industrial robots provide reliable support for completing interactive tasks in an industrial environment. However, as HSRs began to enter the service industry, some research cases began to discover that there are still problems with the interaction of HSRs in the social environment. Caleb Solly (2018) believed that users can also

help robots when robots help users; meanwhile, users can give feedback to optimize the system. The feedback reflects not only the optimization of the robot system but also the satisfaction of customers. Chung's (2018) study indicated that hotels in the hospitality industry want to collect customer feedback in real-time to immediately disseminate positive feedback and respond to unsatisfactory customers while they are still on the scene. Guests want to inform their experience without affecting their privacy. Stakeholders in the hospitality industry hope that intelligent robots can interact more with users. Besides, Rodriguez (2015) concluded that the optimal distance between users and robots is 69.58 cm. Specifically, interaction with a certain greeting mode can attract users to maintain a longer interaction time; robots with the active search are more attractive to participants. The interaction time is longer than that of passively searching robots, suggesting that robots should be designed to keep at a certain distance from humans and consider adding the ability to allow robots to actively identify customers and attract them.

Research suggests that the current interactive system used by HSRs lacks the ability to process and adapt to dynamic social environments. The dynamic social environment here refers to the same human behavior and language often expressing different meanings in different social situations, such as In different situations, the handshake may require two completely different interactive messages to respond. Therefore, this research proposes the concept of HRC to develop a new interactive mode to solve the current problems faced by HRI in the hospitality industry. For a more detailed comparison of HRI and HRC, please refer to the video in the appendix link (<https://youtu.be/fZmV4MKeYtQ>).

## **Review of Neural Image Captioning**

The challenge of generating natural language descriptions from visual data has been extensively researched in the field of computer vision. However, early research has mainly focused on generating natural language descriptions from video-type visual data (Gerber, 1996; Mitchell et al, 2012). These systems convert complex visual data into natural languages using rule-based systems. However, because the rules are artificially designed, these systems are sufficiently robust, bionic, and have been shown to be beneficial in limited applications such as traffic scenarios (Vinyals et al., 2015). In the past decade, various researchers, inspired by the successful use of sequence-to-sequence training with neural networks for machine translation, proposed a method for generating image descriptions based on recurrent neural networks (RNNs) (Cho et al., 2014; Sutskever et al., 2014). In fact, this method of replacing the encoder in the encoder-decoder framework in machine translation with image features transforms the original complex task of generating image data caption into a simple process of "translating" the image into a sentence (Cho et al., 2014). Furthermore, Donahue et al. (2014) used long short-term memory (LSTM) for end-to-end large-scale visual learning processes. In addition to images, Donahue et al. (2014) also applied LSTM to videos, allowing their models to generate video descriptions. Vinyals et al. (2015) and Kiros et al. (2014) initially proposed the structure of a currently popular neural image generation algorithm based on the combination of a convolutional neural network (CNN) image recognition model and a natural language processing (NLP) structured model. Moreover, the

neural image captioning algorithm based on the attention mechanism has also attracted extensive attention in the field of computer vision. Denil et al. (2012) proposed a real-time target tracking and attention recognition model driven by sight data. Tang et al. (2014) proposed an attention-generation model based on deep learning. From the perspective of visual neuroscience, the model requires object-centric data collection for model generation. Subsequently, Mnih et al. (2014) proposed a new recurrent neural network model, which can adaptively select specific areas or locations to extract information from images or videos and process the selected area at high resolution. As the algorithm has increasingly mature, the application of the algorithm in related fields has also been breaking through recently, such as the caption generation of car images (Chen et al., 2017), the description generation of facial expressions (Kuznetsova et al., 2014), and educational NAO robots driven by image caption generation for video Q&A games for children's education (Kim, 2015). Recent research on image caption generation also shows that the accuracy and reliability of the technology have increased (Ding et al., 2019). In addition, reinforcement learning to automatically correct image caption generation networks have also been proposed (Fidler, 2017). These deep learning-based studies have undoubtedly laid a foundation for the possible NIC integration with HSRs as proposed in this study. The novel integration led to the possibility for humanoid robots to interact with humans while recognizing the social environment in real time, thereby improving the interactive service quality of the HSRs.

### ***Neural Image Caption Generation Algorithm 'Crash Into' Robot***

An increasing number of studies have been conducted on HRI combined with image caption generation algorithm. Kim et al. (2015) used the structure of a convolutional neural network (CNN) combined with RNN + deep concept hierarchies (DCH) to design and develop an educational intelligent humanoid robot system for play video games with children. In this study, CNN was used to extract and pre-process cartoons with educational features, and RNN and DCH were used to convert the collected video features into Q&A about cartoons. During the game, after watching the same cartoon, the child and the robot ask and answer questions based on the content of the cartoon. The research results show that such a system can interact effectively with children. However, for HRIs, such simple and limited-structured Q&A conditions cannot satisfy all the interaction scenarios required. Cascianelli et al. (2018) used a gull-gated recurrent unit (GRU) encoder-decoder architecture to develop a human-robot interface that provides interactive services for service robots. This research solves a problem called natural language video description (NLVD). The authors also compared the performance when using LSTM and GRU with two different algorithms to solve these problems. They demonstrated that the GRU algorithm runs faster and consumes less memory. This type of model may be more suitable for HSRs. Although the research model is competitive on public datasets, the experimental results on the designed datasets show that the model suffers from significant overfitting. This proves that in the actual model training process, a specific training dataset for HSR interaction should be established, and other methods (such as transfer learning) should be considered to improve the generalization ability of the model for interactive tasks. Luo et al. (2019, June) created a

description template to add various image features collected by the robot, such as face recognition and expression, to the generated description. Compared with the previous models, their interaction is slightly more natural and closer to the human description. However, Luo et al. use the model to provide limited services to industry managers, hard to generalize, and not for developing an entire HRI framework.

Like the research on robot vision language, research on robot vision action is in its infancy. Yamada et al. (2016) used RNNs to enable robots to learn commands online from humans and respond with corresponding behaviors. This research furthermore provides a reference and direction for humanoid robots to use deep learning to obtain online learning capabilities for human commands. Inspired by the above study, the rationale and hypothesis proposed in the present research are that the description generated by the neural image captions can drive HSRs to perform appropriate behaviors, and HSRs can even obtain online learning capabilities of interacting with surrounding people through studying and analyzing social environments. Tremblay et al. (2018) and Nguyen et al. (2018) believe that non-experts often lack the rationality of task descriptions when issuing instructions to robots. They use deep learning to allow robots to automatically generate human-readable instructions' descriptions according to the surrounding social environment. In addition, Nguyen et al. (2018) also used visual data to make humanoid robots imitate and learn human actions under corresponding commands so that the robot can learn how to complete the corresponding tasks only through visual data; however, social robots cannot complete precise control of movements when they imitate movements of visual data.

## **Contribution to the Knowledge: The Bionic-Companionship Framework with NIC for HSRs**

The contribution of the present study is the novel investigation and design of the bionic-companionship framework for HSRs, adapting and integrating neural image caption generation algorithms and bionic humanoid robots, to be validated in a lab-controlled environment and real-life exploration. The new HRC framework is anticipated to enhance HRI to reach a new state, making it possible for HSRs to become bionic companions of humans.

This study proposes adapting and integrating deep learning techniques to one of the world's most advanced HSRs so that robots can autonomously and in a timely fashion convert pictures or data information captured by robotic visions and sensors into texts or sentences in order to respond and communicate more naturally with humans. The conceptual model of the proposed system consists of various modules, as shown in Figures 1, 2. The contributions of this research are summarized as follows:

- (1) In order to solve the current problems of HSRs in the hospitality industry, a new interactive concept - HRC is proposed.
- (2) A novel bionic interaction framework is designed based on the proposed HRC.

(3) A system that can be used on HSRs is developed based on the bionic interaction framework, and the system has been tested and verified. The preliminary results prove that the system can enable HSRs to handle dynamic social environments.

### ***Humanoid service robot used in research***

The design and investigation of this HRC framework involves using the Canbot U05E humanoid robot (see Fig. 1 for the high-level design, Figs 2–5 for further details) (Canbot, 2021). The robot's 22-degree-of-freedom motion joints enable it to perform a variety of simulated movements, such as raising the head, turning the head, raising the arm, shaking the crank, shaking hands, leaning back, walking, and turning, and based on the proposed framework, it can acquire natural human behaviors and, as a result, efficiently interact with humans. In addition, Canbot U05E's advanced vision system and sensors can collect more complete environmental data for the proposed design and make the novel framework more robust. The robot is designed to imitate the human's seven senses, providing strong support for the concept and implementation of the bionic partner designed in this study.

### ***Bionic-Companionship Framework***

In this study, we review the previous works on this topic and research gaps in the literature and describe a novel humanoid service robot and human interaction framework with neural image subtitles as its core (details are shown in Fig. 2). The framework uses the structure of the NIC algorithm to better realize the interaction of HSRs from HRI to the direction of bionic-companionship. According to the initial descriptions of robot companions, as in the studies by Turkle (2006) and (Kim et al., 2015), the proposed framework should provide HSRs with more natural interactions and a more sensitive understanding of the environment, and hence, the design of the framework is divided into two subsystems (see the dotted red).

### ***Image/Video Description Generation System***

These subsystems are the core modules of the entire interactive framework. HSRs collect visual data of the surrounding environment through equipped visual sensors (such as HD or 3D cameras) and sensors (such as tactile and radar). The type of visual data collected depends on the complexity of the interactive task to be completed by HSRs. It is generally considered that more complex interactive tasks require the use of continuous images or real-time videos. The system uses the latest neural image generation algorithm structure and CNN to perform feature extraction on the pictures and video data of the surrounding social environment, and converts the data into feature vector sequences that can be used by RNN. Finally, the RNN completes the process of generating an interactive description from the visual data. HSRs use a speech synthesis system



that converts these descriptions into voices to communicate with humans. This process is different from the past mode of using HSRs human sensing sensors and setting fixed interactive feedback; the innovation of this system is that HSRs can automatically and naturally generate interactive feedback. This means that the change in the scene during the interaction will cause a continuous change in the interaction feedback, and this change is not preset by humans. In addition, in further conversation interactions, human voice response and social environment data will be coordinated by HSRs and produce continuous conversation interaction behavior.

### ***Command-Robot Behavior System***

For HSRs, simple conversation interactions are insufficient. HSRs should generate corresponding motions based on visual and human behavior data. For example, when humans wave to a robot, the robot should also actively respond. The hypothesis of this study is to classify or cluster description text generated from visual data and use these classified description texts to control the motions of HSRs in response to complex interactive tasks. For example, when the description generated by neural image captions is "Hello", then HSRs will automatically determine whether 'Hello' matches a category that requires interactive motion and performs corresponding motions such as waving.

## **Pilot testing, Preliminary Results, and Discussion**

In the present study, we designed and integrated a classic NIC model on the HSR and performed a preliminary evaluation.

### ***Introduction to HSR-NIC Model***

The structure of the HSR-NIC algorithm used in this study was adapted and enhanced from the model structure proposed by Mao et al. (2014) who used a classic encoder-decoder structure. In this study, the encoder uses the Xception pre-trained CNN to convert the input image into a feature vector. The word sequence is then input into the LSTM after a layer of word embedding layer, and finally, an add operation is performed on the word features output by the LSTM and the image features extracted by the trained CNN. These are then input into a decoder composed of a single-layer fully connected layer, which generates the probability distribution of the next word using a softmax layer. The LSTM introduced by the model can solve the long-term dependency problem in the traditional RNN, thereby improving the accuracy of the model. The dense representation of word embedding can reduce the amount of calculations involved in the model; it also enables the model to capture similar relationships between words. In addition, the model used in this study also introduces a dropout layer with a probability of 50% to increase the robustness

of the model. The teacher forcing mechanism was used during model training to accelerate the model training process. The optimizer used in the research is Adam, which has the advantages of making the model converge more quickly and automatically adjusting the learning rate with learning. The variables of the model are updated by minimizing the cross-entropy loss between the probability distribution of the predicted result and the probability distribution of the true result and back-propagation. The model structure diagram as follow (Fig. 3):

### **Model Forward Propagation Process**

The training process of the image captioning task can be described as follows: For a picture in the training set, its corresponding description is a sequence that represents the words in the sentence. For model  $\theta$ , given input image  $I$  from the HSR's vision, the probability of the model generating sequence is expressed as

$$P(S|I;\theta) = \prod_{t=0}^N P(S_t|S_0, S_1, \dots, S_{t-1}, I; \theta) \quad (1)$$

The logarithm of the likelihood function is used to obtain the log-likelihood function:

$$\log P(S|I;\theta) = \sum_{t=0}^N \log P(S_t|S_0, S_1, \dots, S_{t-1}, I; \theta) \quad (2)$$

The training objective of the model is to maximize the sum of the log-likelihoods of all training samples:

$$\theta^* = \arg_{\max_{\theta}} \sum_{(I,S)} \log P(S|I;\theta) \quad (3)$$

where  $(I, S)$  is the training sample. This method of maximum likelihood estimation is equivalent to empirical risk minimization using the log-loss function. Therefore, in the forward propagation process of this research model, the image feature vector  $I_v$  is extracted from the image using the CNN, and a two-dimensional vector of shape (batch size, 2048) is the output.

$$I_v = CNN_{\theta_c}(I) \quad (4)$$

The extracted image features need to be encoded by a fully connected layer into the context feature vector  $C$  that can be matched with word features. The word feature vector is the output  $O_t$  of the LSTM over the time step. The input word of LSTM passes through a word-embedding layer to generate a dense vector representation  $W(s)$ .

$$C = W_{\theta}(I_v), O_t = LSTM_{\theta}(W(s)) \quad (5)$$

Finally, word feature  $O_t$  and context feature  $C$  are together input into a decoder composed of a single fully connected layer after the softmax calculation generates the probability distribution of the next word  $P(S_i|I;\theta)$ .

$$P(S_i|I;\theta) = \text{softmax}(W_\theta(C + O_t)) \quad (6)$$

The loss function is expressed as

$$L = \sum_{t=1}^T y^{(t)} \log p^{(t)} + (1-y^{(t)}) \log (1-p^{(t)}) \quad (7)$$

### Training Dataset

For the present study, we use Flickr 8k (Rashtchian et al., 2010) as the training dataset. This is a new benchmark collection for sentence-based image descriptions and searches. It consists of 8,000 images. Each image was paired with five different captions. These captions provide content descriptions of the objects and events in the picture. The images do not contain any well-known people or locations but depict random scenes and situations. Examples of datasets are shown in Fig. 4. The Flickr 8k dataset not only contains images of animals and objects, but also of some social scenes. These data can help robots to better understand natural, day-to-day scenes.

### The Process of Humanoid Service Robot Generating Image Captions

To explore the feasibility of the bionic-companionship framework, preliminary tests were conducted on a real humanoid service robot (Canbot U05E). The process of generating image captions by a humanoid service robot is divided into four steps, as shown in Fig 5.

**Step 1.** The HSR-NIC API is responsible for controlling the robot to call the high-definition camera to collect surrounding environment information (the data collection in this study is focused on HSR capture images). The collected data will be sent to the local host service program through the HTTP protocol and wait for a response from the HSR.

**Step 2.** The HSR-NIC localhost server program receives the data, and the requests perform preliminary processing and cleaning of the data (image) and send the data (image) to the HSR-NIC model server program to wait for the calculation result (the generated caption description).

**Step 3.** The HSR-NIC model server program analyzes the image data according to the training parameters saved before, generates the descriptive caption, and returns it to the local server.

**Step 4.** The HSR-NIC local server program sends the caption description to the robot application through the HTTP protocol, and the robot application controls the robot to respond according to the caption description, such as speech synthesis and motion control.

### ***Preliminary Test Results and Limitation***

In this study, we conducted a preliminary test on a humanoid service robot integrated with the NIC algorithm. The results of the preliminary test were found to be promising.

With the discuss of the last chapter, the research will integrate the NIC into the HSRs to make the HSRs take advantage of the change of the surrounding environment interact with the human better. Therefore, the system proposed by this research will combine qualitative analysis and quantitative analysis to initially validate the performance of the system.

This study introduces the cross-entropy loss curve of the last 50 epochs of the model as the evaluation metric for quantitative analysis. As shown in the Fig. 6, the model finally converges to the minimum loss value of 2.65 in the training set and 2.71 in the validation set, which proves that the model has no over-fitting and under-fitting, and has generalization ability. Since the loss value is calculated from the sum of the difference between the probability value of each predicted word in the predicted description and the true value, the loss value will be affected by the sentence length of the predicted description. In related work, researchers (Li et al., 2020; Hu et al., 2020) used some more reliable evaluation methods to evaluate the performance of the model, including the BLUE4 (Papineni et al., 2002) and CIDEr (Vedantam et al., 2015). These evaluation metrics are usually used in the field of machine translation instead of manual evaluation. Since the tasks handled by the NIC model can be regarded as translated from images/scenes into English, the evaluation metrics can also be applied to the evaluation of NIC. This study will use qualitative analysis to replace quantitative analysis of metrics such as BLUE4 and CIDEr, so as to further evaluate the preliminary performance of HSR after the integrated NIC model.

As shown in Fig. 7 and Fig. 8, the researcher conducted two sets of tests in three different scenarios with HSR. In the first set of tests, the researcher wore a hat and changed scenarios. In the second set of tests, the researcher did not wear a hat, and the scene switching method was the same as in the first set. It can be seen from the experimental results that the humanoid robot can complete the perception of scene switching through this algorithm and generate a rough description of the scene. In the first set of tests, most of the content described was accurate. The robot equipped with the NIC algorithm can effectively identify ‘man’, ‘black shirt’, , and ‘sitting on a bench’. However, in the second group of tests, there were many errors in the recognition results. This could be attributed to the researcher’s long hair. Interestingly, researchers with long hair are easily identified as women or children. This indicates that the accuracy of the NIC algorithm still has room for improvement.

In addition, in order to test the performance of the system in a dynamic environment. The researcher conducted the test in a real environment (As Fig. 8). The researcher selected six real environments as the test data and let the robot generate interactive information. Among the six real interactive environments, there are three scenes that can be more accurately recognized by the robot and produce corresponding descriptions. The description information can correspond to the test environment, and the corresponding part of the description has been highlighted with the same color in the Fig. 8. Some of the objects, facilities, and human movements in these scenes can be accurately predicted, such as sidewalk, traffic, bench, building, building, etc. However, in the other three environments, the robot did not give an accurate description. The researchers believe that this may be due to the fact that the training set does not contain objects in these three environments, causing the model to fail to learn how to express the 'unfamiliar environment'.

In general, as per the results of the two experimental sets, it was proven that the robot equipped with the NIC algorithm can capture the changes in the surrounding environment and generate different feedbacks according to the changes. The results also demonstrate the feasibility of the proposed bionic-companionship framework. Although there is still a gap between the prediction results of the algorithm and the real communication scene, the researcher believes that special data collection for some specific interaction scenarios and model training for these specific data can be effective in addressing this gap. Future research directions will mainly focus on improving the accuracy of algorithms and achieving more human-like interactions. (The detailed process is shown in the [HSR-NIC demo video](#)). In addition, the researcher believes that the scene understanding of static images is the basis for dealing with dynamic environments. Some researches have mentioned that the introduction of related algorithms of object detection into NIC can identify and generate descriptions of scenes in dynamic environments. This is also the current research limitation of this research and the research challenges that will be faced in the future.

## Conclusions

This study presents a review of neural image generation algorithms and application cases in the field of robotics, and proposes a novel humanoid service robot and human interaction framework based on the bionic-companionship theory. The subsystems of the bionic-companionship framework are designed and introduced in detail. Preliminary tests also initially proved that the framework could increase the sensitivity of HSRs to changes in the surrounding environment. The proposed framework will contribute to further development from HRI to HRC. Future work will focus on implementing each of the subsystems in the framework and applying the framework to HSRs to verify its performance.

## References

455 CANBOT. (2020). Retrieved 4 June 2020, from [https://www.canbotrobots.com/html/yy-](https://www.canbotrobots.com/html/yy-detail.html)  
456 detail.html  
457 Caleb-Solly P, Dogramadzi S, Huijnen CA, Hvd Heuvel (2018) Exploiting ability for human  
458 adaptation to facilitate improved human-robot interaction and acceptance. *Inf Soc*  
459 34(3):153–165  
460 Cascianelli, S., Costante, G., Ciarfuglia, T. A., Valigi, P., & Fravolini, M. L. (2018). Full-GRU  
461 natural language video description for service robotics applications. *IEEE Robotics and*  
462 *Automation Letters*, 3(2), 841-848.  
463 C. Trevarthen (2017). Play with infants: The impulse for human story-telling, In, Tina Bruce,  
464 Pentti Hakkarainen and Milda Bredikyte (Eds.) *The Routledge International Handbook of*  
465 *Play in Early Childhood*. Abingdon: Taylor & Francis/Routledge, Chapter 15.  
466 [http://www.becera.org.uk/BECERA%202017/CT%20ON%20PLAYRoutledge%20Handboo](http://www.becera.org.uk/BECERA%202017/CT%20ON%20PLAYRoutledge%20Handbook%202017.pdf)  
467 [k%202017.pdf](http://www.becera.org.uk/BECERA%202017/CT%20ON%20PLAYRoutledge%20Handbook%202017.pdf)  
468 Chen, L., He, Y., & Fan, L. (2017). Let the robot tell: describe car image with natural language  
469 via LSTM. *Pattern Recognition Letters*, 98, 75-82.  
470 Chew E., Lee PL., Hu S. & Yang J. (2021). Investigating the First Robotic Nurses: Humanoid  
471 Robot Nightingale and Partners for COVID-19 Preventive Design, the Seventh edition of the  
472 International Workshop on New Trends in Medical and Service Robots Conference,  
473 University Hospital Basel, Basel, Switzerland, 7-9 June 2021. <https://www.mesrob2021.org>  
474 (Springer-indexed)  
475 Cho, Kyunghyun, van Merriënboer, Bart, Gulcehre, Caglar, Bougares, Fethi, Schwenk, Holger,  
476 and Bengio, Yoshua. (2014). Learning phrase representations using RNN encoder-decoder  
477 for statistical machine translation. In *EMNLP*  
478 Chung MJY, Cakmak M (2018) “how was your stay?” : Exploring the use of robots for  
479 gathering customer feedback in the hospitality industry. In: 2018 27th IEEE International  
480 Symposium on Robot and Human Interactive Communication (RO-MAN), IEEE, pp 947 –  
481 954  
482 Dautenhahn, K ; Woods, S ; Kaouri, C ; Walters, M.L ; Kheng Lee Koay & Werry, I., (2005).  
483 What is a robot companion - friend, assistant or butler? (2005) *IEEE/RSJ International*  
484 *Conference on Intelligent Robots and Systems*, pp.1192–1197.  
485 Denil, M., Bazzani, L., Larochelle, H., & de Freitas, N. (2012). Learning where to attend with  
486 deep architectures for image tracking. *Neural computation*, 24(8), 2151-2184.  
487 Ding, S., Qu, S., Xi, Y., Sangaiah, A.K. and Wan, S., 2019. Image caption generation with high-  
488 level image features. *Pattern Recognition Letters*, 123, pp.89-95.  
489 Donahue, Jeff, Hendriks, Lisa Anne, Guadarrama, Segio, Rohrbach, Marcus, Venugopalan,  
490 Subhashini, Saenko, Kate, and Darrell, Trevor. (2014). Long-term recurrent convolutional  
491 networks for visual recognition and description. *arXiv:1411.4389v2*.  
492 Fidler, S. (2017). Teaching machines to describe images with natural language feedback. In  
493 *Advances in Neural Information Processing Systems* (pp. 5068-5078).

- 494 Fong, T., Thorpe, C., & Baur, C. Collaboration (2003) dialogue, human-robot interaction.  
495 In *Robotics Research*, pp. 255-266. Springer, Berlin, Heidelberg.
- 496 Gerber, R. & Nagel, N.-H., (1996). Knowledge representation for the generation of quantified  
497 natural language descriptions of vehicle traffic in image sequences. *Proceedings of 3rd*  
498 *IEEE International Conference on Image Processing*, 2, pp.805–808 vol.2.
- 499 Gui, L. Y., Zhang, K., Wang, Y. X., Liang, X., Moura, J. M., & Veloso, M. (2018, October).  
500 Teaching robots to predict human motion. In *2018 IEEE/RSJ International Conference on*  
501 *Intelligent Robots and Systems (IROS)* (pp. 562-567). IEEE.
- 502 Harris, Karen, Austin Kimson, and Andrew Schwedel (2018), “Why the Automation Boom Could  
503 Be Followed by a Bust,” *Harvard Business Review* (March 13),  
504 <https://hbr.org/2018/03/why-the-automation-boom-could-be-followed-by-a-bust>.
- 505 Hu, X., Yin, X., Lin, K., Wang, L., Zhang, L., Gao, J., & Liu, Z. (2020) Vivo: Surpassing human.  
506 performance in novel object captioning with visual vocabulary pre-training. *arXiv preprint*  
507 *arXiv:2009.13682*.
- 508 Ivanov, S. (2019). Ultimate transformation: How will automation technologies disrupt the travel,  
509 tourism and hospitality industries? *Zeitschrift für Tourismuswissenschaft*, 11 (1), 25–43.
- 510 Ivanov, S., Gretzel, U., Berezina, K., Sigala, M., & Webster, C. (2019). Progress on robotics in  
511 hospitality and tourism: A review of the literature. *Journal of Hospitality and Tourism*  
512 *Technology*. <https://doi.org/10.1108/JHTT-08-2018-0087>.
- 513 Kim, K. M., Nan, C. J., Ha, J. W., Heo, Y. J., & Zhang, B. T. (2015, September). Pororobot: A  
514 deep learning robot that plays video Q&A games. In *2015 AAAI Fall Symposium Series*.
- 515 Luo, R. C., Hsu, Y. T., & Ye, H. J. (2019, June). Multi-Modal Human-Aware Image Caption  
516 System for Intelligent Service Robotics Applications. In *2019 IEEE 28th International*  
517 *Symposium on Industrial Electronics (ISIE)* (pp. 1180-1185). IEEE.
- 518 Li, X., Yin, X., Li, C., Zhang, P., Hu, X., Zhang, L., ... & Gao, J. (2020) Oscar: Object-semantics.  
519 aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*,  
520 pp. 121-137. Springer, Cham.
- 521 Mitchell, M., Han, X., Dodge, J., Mensch, A., Goyal, A., Berg, A., Yamaguchi, K., Berg, T.,  
522 Stratos, K., & Daumé III, H. (2012). Midge: Generating image descriptions from  
523 computer vision detections. In *Proceedings of the 13th conference of the European chapter*  
524 *of the association for computational linguistics*(pp. 747–756). Association for  
525 *Computational Linguistics*.
- 526 Nguyen, A., Kanoulas, D., Muratore, L., Caldwell, D. G., & Tsagarakis, N. G. (2018, May).  
527 Translating videos to commands for robotic manipulation with deep recurrent neural  
528 networks. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*  
529 (pp. 1-9). IEEE.[19]
- 530 Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. (2002). Bleu: a method for automatic evaluation  
531 of machine translation. In *Proceedings of the 40th annual meeting of the Association for*  
532 *Computational Linguistics*, pp. 311-318.

- P. Kuznetsova, V. Ordonez, T.L. Berg, Y. Choi, (2014), TREETALK: Composition and compression of trees for image descriptions, *Trans. Assoc. Comput.Ling.* 2 (1) 351–362.
- R. Kiros, R. Salahutdinov, R. Zemel, (2014), Multimodal neural Language models, in: *International Conference on Machine Learning*, pp. 595–603.
- Robert K. Yin (2014) *Case Study Research Design and Methods* (5th ed.). Thousand Oaks, CA: Sage. 282 pages.
- Rodriguez-Lizundia E, Marcos S Zalama (2015) A bellboy robot: Study of the effects of robot behaviour on user engagement and comfort. *Int J Human-Comp Stud* 82:83–95
- Rashtchian, C., Young, P., Hodosh, M., & Hockenmaier, J. (2010) Collecting image annotations using amazon’s mechanical turk. In *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk*, pp. 139-147.
- R. C. Luo, Y. Hsu, Y. Wen and H. Ye, "Visual Image Caption Generation for Service Robotics and Industrial Applications," (2019) *IEEE International Conference on Industrial Cyber Physical Systems (ICPS)*, Taipei, Taiwan, 2019, pp. 827-832, doi: 10.1109/ICPHYS.2019.8780171.
- Sutskever, Ilya, Vinyals, Oriol, and Le, Quoc VV. (2014). Sequence to sequence learning with neural networks. In *NIPS*, pp. 3104– 3112
- Tang, X., Zhang, Q., & Hu, L. (2020). An EKF-based performance enhancement scheme for stochastic nonlinear systems by dynamic set-point adjustment. *IEEE Access*, 8, 62261-62272.
- Tremblay, J., To, T., Molchanov, A., Tyree, S., Kautz, J., & Birchfield, S. (2018, May). Synthetically trained neural networks for learning human-readable plans from real-world demonstrations. In *2018 IEEE International Conference on Robotics and Automation (ICRA)* (pp. 1-5). IEEE.
- Trevarthen, C. (2001). Intrinsic motives for companionship in understanding: Their origin, development, and significance for infant mental health. *Infant Mental Health Journal: Official Publication of The World Association for Infant Mental Health*, 22(1 - 2), 95-131.
- Turkle, S. (2006). A nascent robotics culture: New complicities for companionship. *American Association for Artificial Intelligence Technical Report Series AAAI*.
- Vedantam, R., Lawrence Zitnick, C., & Parikh, D. Cider (2015) Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* pp. 4566-4575.
- V. Mnih, N. Hees, A. Graves, K. Kavukcuoglu, (2014), Recurrent models of visual attention, *NIPS*
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 3156-3164).
- Vedantam, R., Lawrence Zitnick, C., & Parikh, D. Cider (2015) Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* pp. 4566-4575.



573 Yang, J. & Chew, E. (2020) A Systematic Review for Service Humanoid Robotics Model in  
 574 Hospitality, Springer International Journal for Social Robotics.  
 575 <https://doi.org/10.1007/s12369-020-00724-y> (Impact Factor 3.168; indexed by JCSR &  
 576 Scopus, within the top 20 of Google Robotics publication)  
 577 Yang, J. & Chew, E. (2021) The Novel Design Model for Robotic Waitress, The International  
 578 Journal for Social Robotics (Impact Factor 3.168; WoS & Scopus-indexed, ranked within the  
 579 top 20 of the Top Publications Google Scholar citations)  
 580 Y. Tang, N. Srivastava, R.R. Salakhutdinov, (2014), Learning generative models with vi- sual  
 581 attention, in: NIPS, pp. 1808–1816.  
 582 Yamada, T., Murata, S., Arie, H., & Ogata, T. (2016). Dynamical integration of language and  
 583 behavior in a recurrent neural network for human–robot interaction. *Frontiers in*  
 584 *neurorobotics*, 10, 5.  
 585 Zhang, Q. C., Hu, L., & Gow, J. (2020). Output feedback stabilization for mimo semi-linear  
 586 stochastic systems with transient optimisation. *International Journal of Automation and*  
 587 *Computing*, 17(1), 83-95.

**Table 1** (on next page)

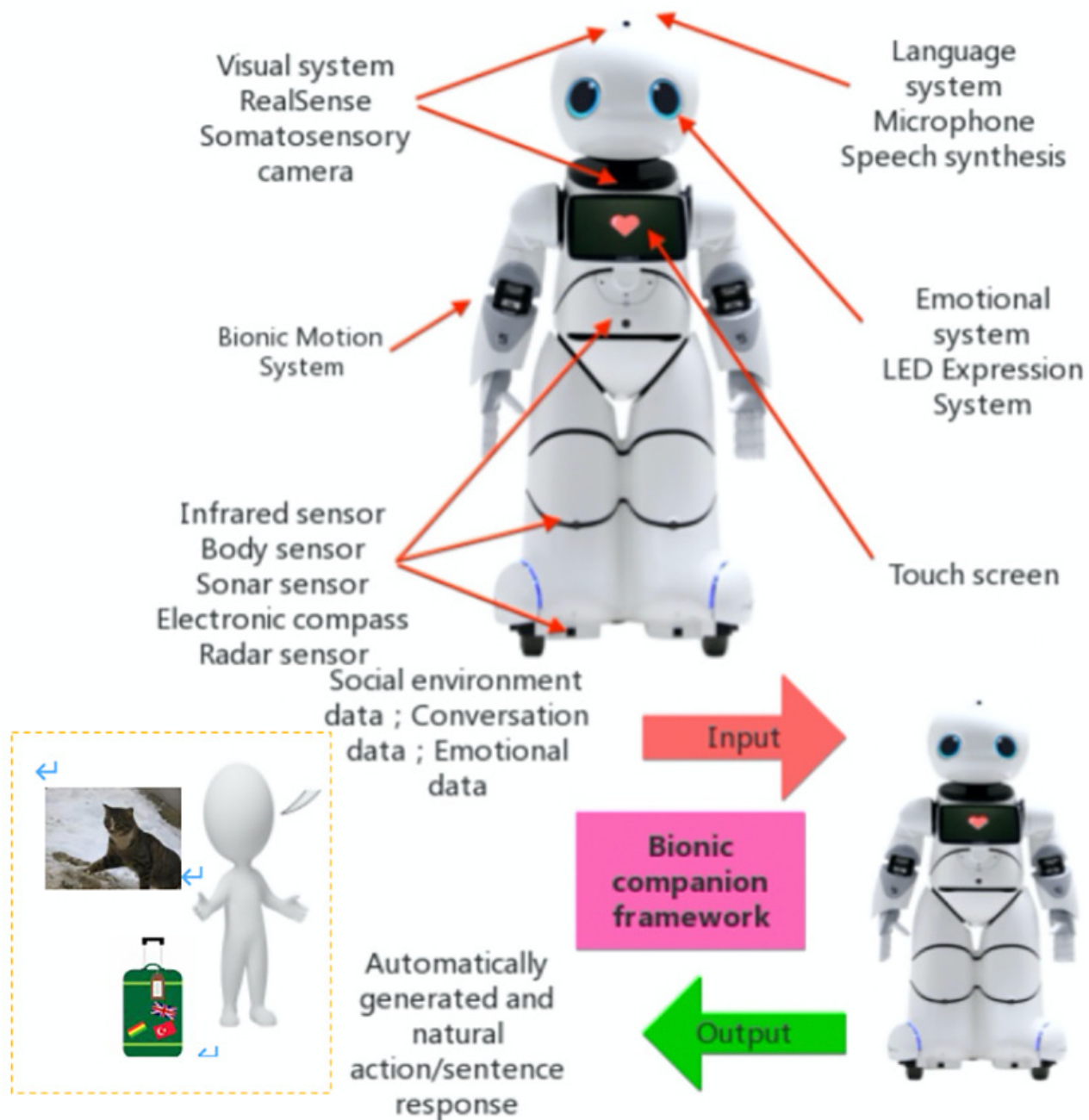
Scenario Based Comparison for HRI and HRC

| Scenario:   | HRI  | HRC   |
|---|--|---|
| <p><i>Scenario 1: Hospitality:</i></p> <p><b>The enhancement from HRI to HRC:</b></p> <ol style="list-style-type: none"> <li><b>1) Compared with HRI, robots in HRC can recognize the environment (luggage) and changes in customers' appearance (red shirt), which is in line with the proposed concept of companion should become dynamic in the theory of Companionship.</b></li> <li><b>2) More enthusiastic and bionic interaction capabilities (automatically detect whether they are regular customers, and greet enthusiastically)</b></li> <li><b>3) The robot in HRC remember the customers' past orders and provide meal recommendation for well-being.</b></li> </ol> | <p>When you enter a hotel, you see a reception area dominated by robots. When you approach the reception area, the HSR will say "Welcome to Hotel XYZ, please follow the instructions to check-in on my display screen". After completing the check-in, the robot will tell you the room number and issue you a room card, you go to your room, change a suit and prepare to go downstairs to eat. When you go back to the reception, the robot says 'welcome, please follow the instructions to place an order on my display'. You choose a few dishes that look good on the screen of the robot, but when the food comes up you don't seem to be satisfied with the taste...</p> | <p>When you enter a hotel, you can a reception area thoughtfully served by robots. The robots also see you and wave to you, 'Welcome Jack, you have a nice luggage, I can help you to check-in. What else can I do for you?' After completing the check-in, the robot will tell you the room number and issues you a room card. You go to your room and change to a red shirt to go downstairs to eat. When you go back to the reception, the robot says, "Welcome Jack, you wear nice red shirt, what can I do for you?" You choose a few dishes that look good on the robot's screen, but the robot tells you that 'According to your past order and diet preferences, these meals may not be suitable for you. Feel free to change it to a less cholesterol dishes with special house promotion and I recommend you to take this quality wine as a treat to have a healthy eating while enjoying your stay with us.'</p> |
| <p><i>Scenario 2: Health care:</i></p> <p><b>The enhancement from HRI to HRC:</b></p> <ol style="list-style-type: none"> <li><b>1) Robots have dynamic thinking and real-time neural image captioning ability: able to deal with</b></li> </ol>   | <p>You bought a robot at home to monitor your health. The robot obtains some of your health indicators (such as temperature, blood pressure, etc.) through some external devices. When there is a problem with your indicators, the robot can give you corresponding suggestions or help you contact a doctor. One day</p>   | <p>You bought a home care robot to monitor your health. The robot obtains some of your health indicators (such as pulses, blood pressure, etc.) through some external devices. When there is a problem with your indicators, the robot can give you corresponding suggestions or help you contact a doctor. One day you suddenly fainted at home for some reason. The robot discovered your real-</p>   |

|  |  |  |
|--|--|--|
| <p><b>emergencies and a quick decision making from what it sees the environment in real-time.</b></p> <p><b>2) Robots has been improved from conventional smart Q&amp;A and interactions to new concept of bionic companionship.</b></p> | <p>you suddenly fainted at home for some reasons, but because you did not aim at the detection device connected to the robot, the robot did not find your condition. Fortunately, your neighbor found you fainted at home. . .</p> | <p>time condition through the deep learning vision system and contacted the your family member or hospital in time, subject to what the robot sees, e.g. fainted human with lots of blood or motion (call hospital for emergency); fainted human with conscious and free speech (call family members).</p> |
|--|--|--|

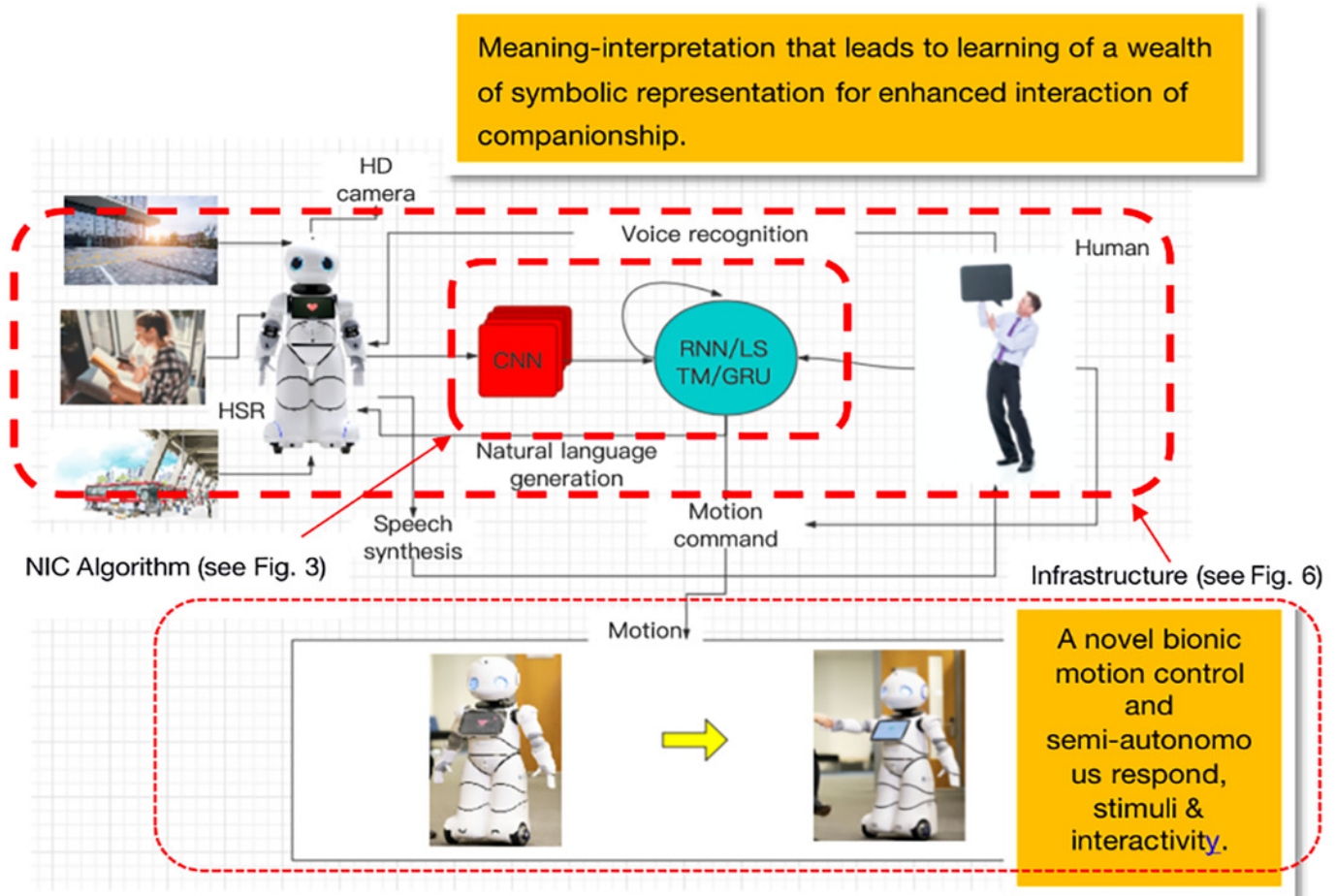
# Figure 1

HSR capabilities with the proposed high-level of HRC conceptual model



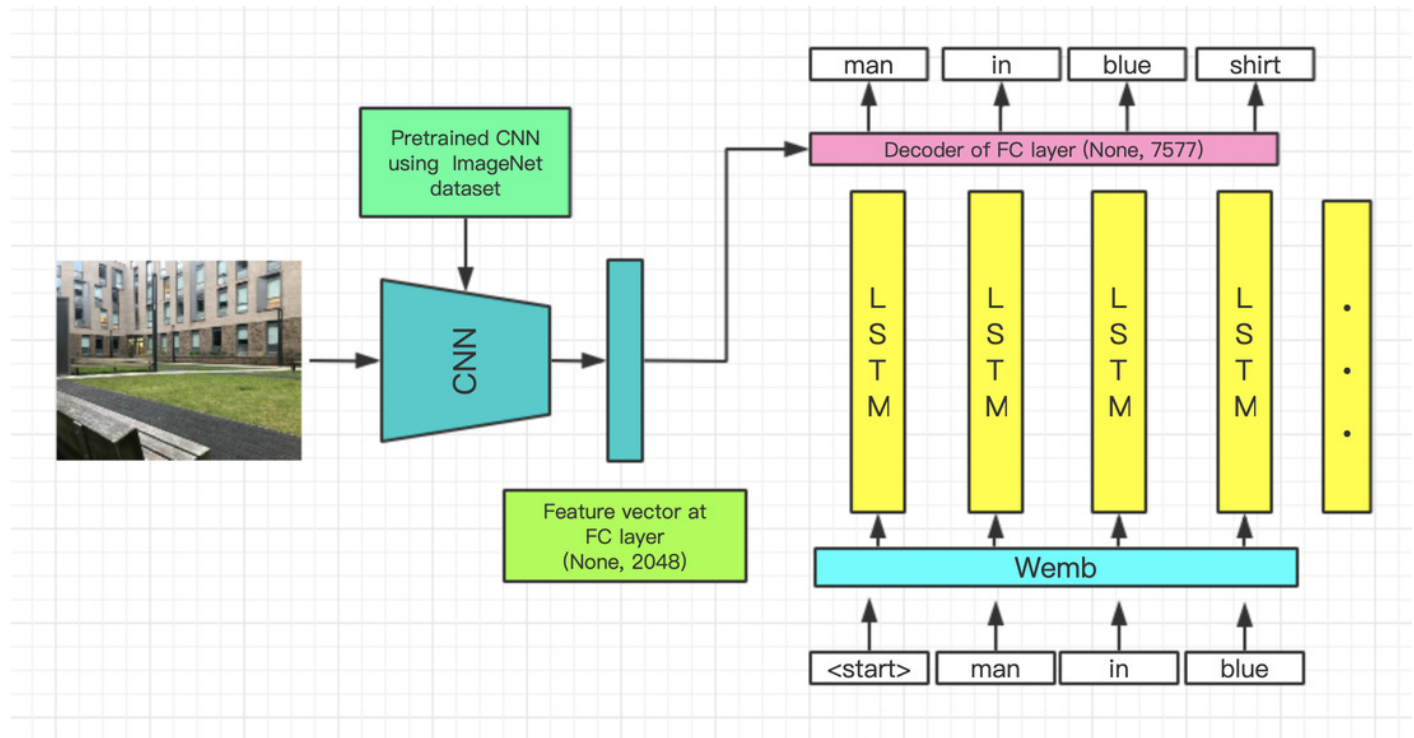
# Figure 2

## Bionic-companionship framework design



# Figure 3

Neural image captioning model structure for HSR





# Figure 4

Samples Flickr 8k (Rashtchian et al., 2010) training data set



0. A black dog is running after a white dog in the snow .
1. Black dog chasing brown dog through snow
2. Two dogs chase each other across the snowy ground .
3. Two dogs play together in the snow .
4. Two dogs running through a low lying body of water .

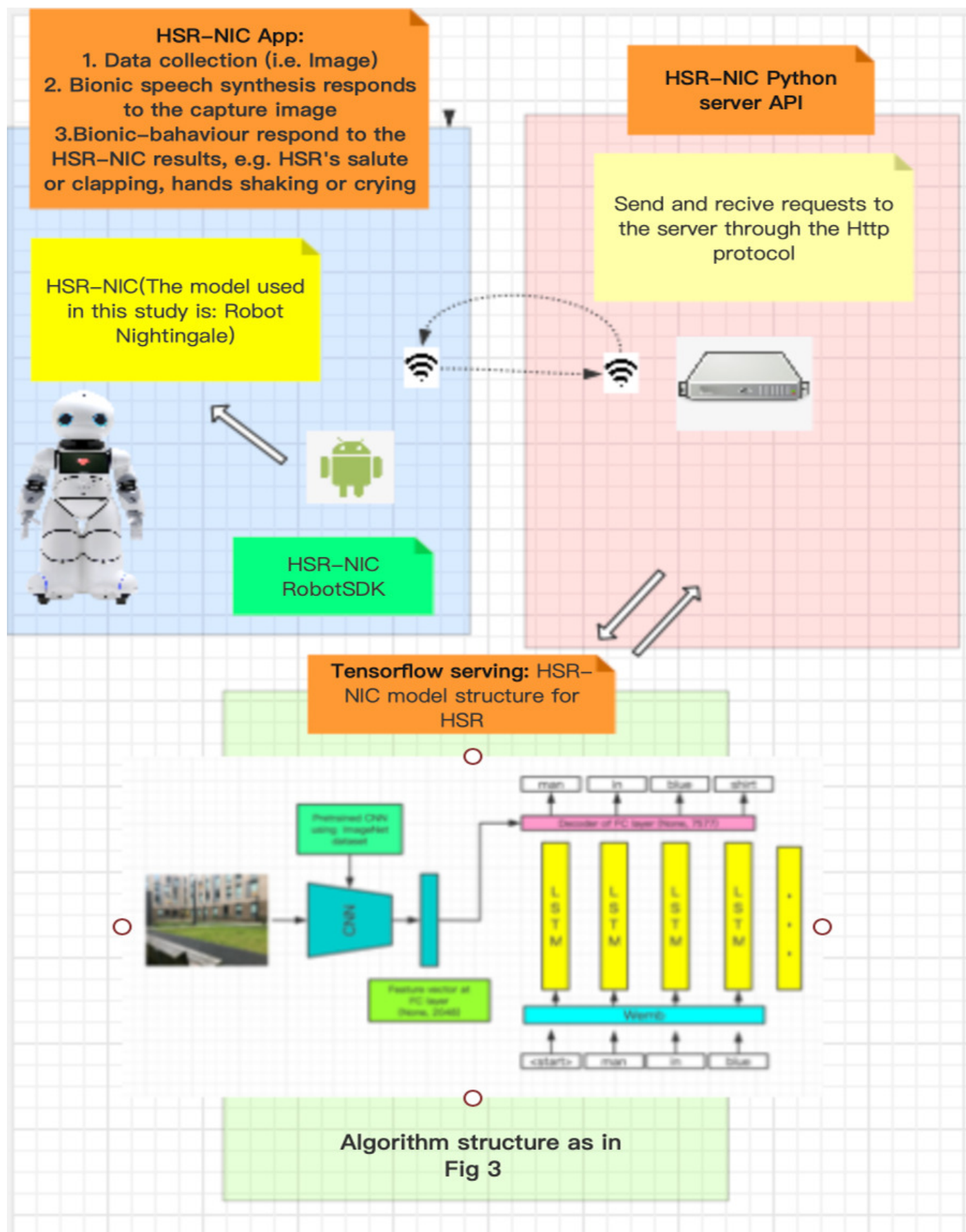


0. A man and a dog are sitting on a park bench with several people walking down the path that runs alongside it .
1. A man and a dog on a park bench in the foreground , with a group of walking people in the distance .
2. A man sits on a bench in the park with a dog , while others walk nearby .
3. People sit on a park bench while others jog on a path .
4. The man sits with his dog on a bench in a park .



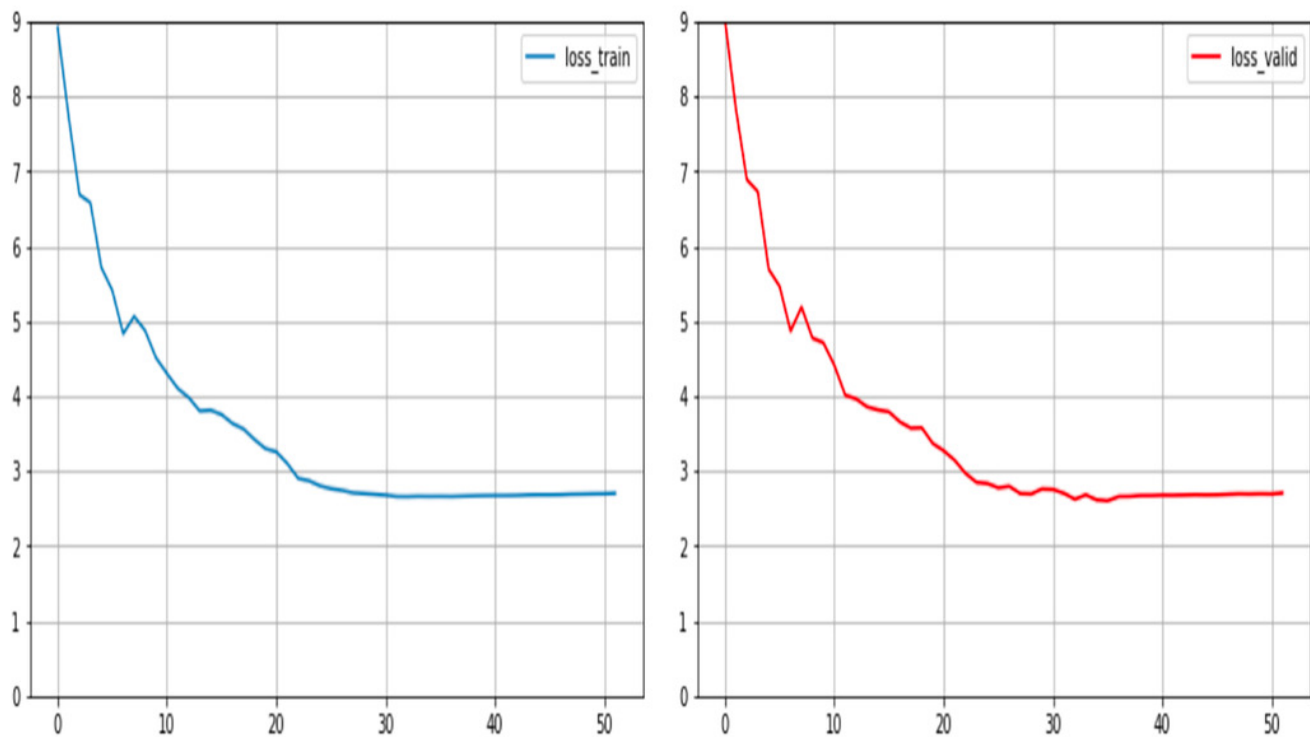
# Figure 5

The Infrastructure of the Humanoid Service Robot Generating Neural Image Captions with as part of the Bionic Companionship Framework



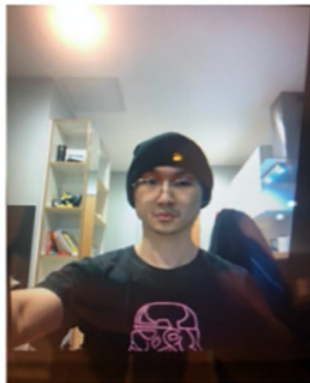
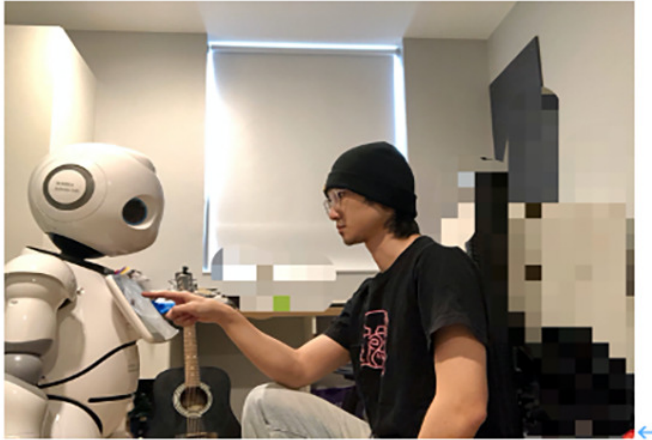
# Figure 6

Loss curve of NIC model on training set and validation set

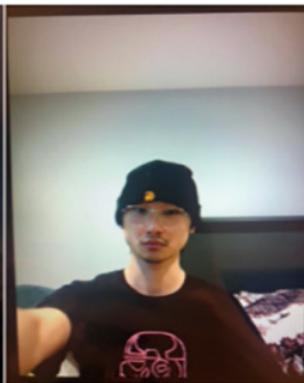


# Figure 7

A Series of Preliminary Testing Results Captured from Canbot U05E and Bionic-Companionship Preliminary Framework



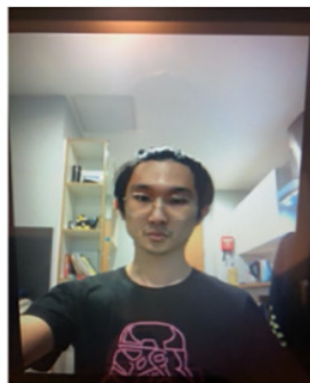
Man in black shirt and jeans is sitting on the ground with Animal



Man in black shirt and jeans is sitting on the ground



Man in black shirt and sunglasses is sitting on bench



Two girls are playing on the grass



Two girls are sitting on bed laughing



Two children are playing on the grass

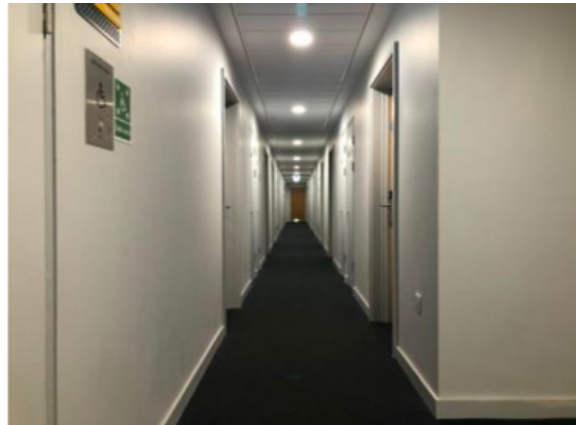
# Figure 8

Real social environment test examples





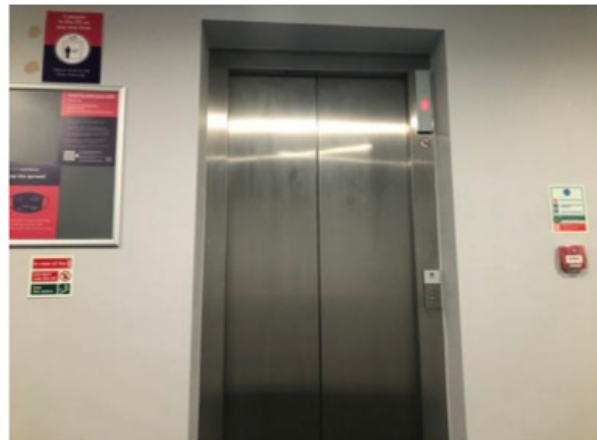
man in helmet is sitting on the sidewalk with traffic



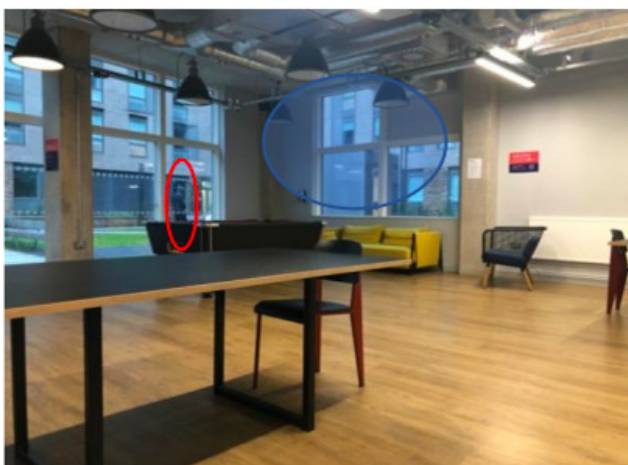
man is sitting on bench with his feet in the air



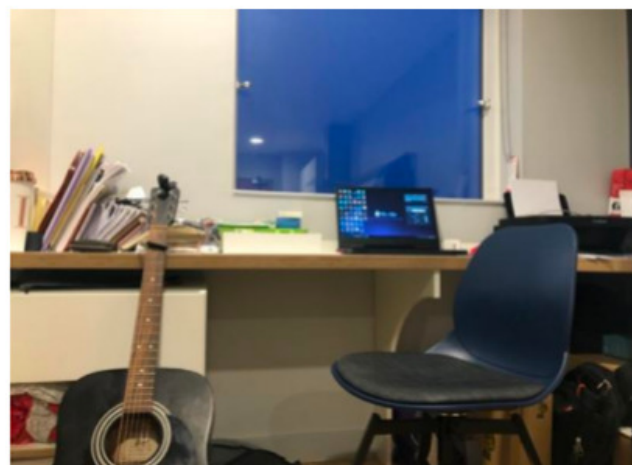
man in blue shirt is sitting on bench with his feet up in front of building



man sitting on the ground with his feet up in the air



man in black shirt is standing in front of large building



man is standing in front of crowd