

LeRoP: A Learning-Based Modular Robot Photography Framework

Hao Kang¹, Jianming Zhang², Haoxiang Li³, Zhe Lin², TJ Rhodes², Bedrich Benes¹

Abstract—We introduce a novel framework for automatic capturing of human portraits. The framework allows the robot to follow a person to the desired location using a Person Re-identification model. When composing is activated, the robot attempts to adjust its position to form the view that can best match the given template image, and finally takes a photograph. A template image can be predicted dynamically using an off-the-shelf photo evaluation model by the framework, or selected manually from a pre-defined set by the user. The template matching-based view adjustment is driven by a deep reinforcement learning network. Our framework lies on top of the Robot Operating System (ROS). The framework is designed to be modular so that all the models can be flexibly replaced based on needs. We show our framework on a variety of examples. In particular, we tested it in three indoor scenes and used it to take 20 photos of each scene: ten for the pre-defined template, ten for the dynamically generated ones. The average number of adjustment was 11.20 for pre-defined templates and 12.76 for dynamically generated ones; the average time spent was 22.11 and 24.10 seconds respectively.

I. INTRODUCTION

Getting a good picture requires getting into a fine location, proper light, and right timing. While getting a great picture requires talent and artistic preparation, there is a possible way of getting a well-taken image automatically. For example, a system could follow the rule of thirds, empty area, keeping eyes in sight, and other commonly-accepted conventions [1].

Some previous work has been done in this direction. For example, Kim et al. [2] introduced a robot that can move and capture photographs according to composition lines of the human target. Robot Photographer of Luke [3] can randomly walk in an unstructured environment, and take photographs of humans basing on heuristic composition rules. To get a good picture, we need to move the robot to a certain location. At the same time, a good sense of composition needs to be taught to the robot. However, the robustness of geometric-based robot motion planning relies heavily on sensor accuracy. Moreover, traditional rule-based designs are usually rigid in handling various composition situations.

Therefore, we propose a flexible learning-based view adjustment framework for taking indoor portraits that we call: LeRoP. The objective of our work is to create a framework that can train a robot to automatically move and capture the best view of a person. We implemented the LeRoP utilizing

a photo evaluation model to propose good views, and using a Deep Reinforcement Learning (DRL) model to adjust the robot position and orientation towards the best view to capture. Additionally, the framework is interactive to the user. For example, our robot with the framework can be triggered to follow the target to a photograph location before searching for the best view. We ingeniously utilized a 360 camera as a supplement to the main camera (for photography). The omnidirectional view helps avoid repetitive rotations in tracking and selects views for the DRL template matching procedure. The modular design allows the photo evaluation aesthetic model to be swapped flexibly basing on photo style preferences. The DRL model can adapt to different view adjustment methods by simply re-training the network with a suitable reward function and apply to different hardware settings.

An example in Figure 1 shows our framework at work. The user first selects *Tracking Mode* to lead the robot to a user desired location. The user then selects *Composing Mode* to start the autonomous composing process. Once *Composing Mode* is activated, the robot observes the scene with its 360° camera for hunting the best view (template). Once the template is detected, it moves into a location matches the template that allows for taking a high-resolution portrait with the second camera.

We tested it on a robot system built on *Turtlebot* shown in Figure 2, and our robot can (a) interact with the user, (b) identify and follow the user, (c) propose well-composed template dynamically or use supplemental pre-defined template, and (d) adjust position to match the template and capture the portrait.

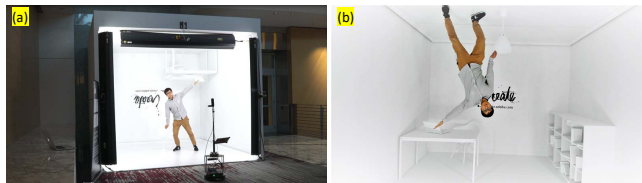


Fig. 1. Our LeRoP robot at work. Photo (a) is the third-person view of the working scene. Photo (b) is the final capture (the photo has been rotated by 180° for better visualization).

We claim the following contributions:

- 1) A template matching based DRL solution with its synthetic virtual training environment for robot view adjustment.
- 2) A method to utilize an omnidirectional camera to support tracking and DRL view selection.

¹Hao Kang and Bedrich Benes are with the Purdue University, 610 Purdue Mall, West Lafayette, IN 47901, USA {kang235, bbenes}@purdue.edu

²Jianming Zhang, Zhe Lin and TJ Rhodes are with the Adobe Research, 345 Park Ave., San Jose, CA 95110, USA {jianmzha, zlin, trhodes}@adobe.com

³Haoxiang Li is with the AIBee, 350 Cambridge Avenue, Suite 250, Palo Alto, CA 94306, USA hxli@aibee.com

- 3) An interactive modular robot framework design that supports automatically capture high quality human portraits.

II. RELATED WORK

Autonomous cameras for both virtual and real photography have long been explored. We refer the reader to the review [4] that summarizes autonomous cameras from the viewpoint of camera planning, controlling, and target selecting. Galvane et al. [5]–[7] also provided several studies on automatic cinematography and editing for virtual environments.

One of the earliest robotic cameras was introduced in [8], [9]. This camera allows intelligent framing of subjects and objects in a TV studio upon verbal request and script information. Byers et al. [10] developed probably the first robot photographer that can navigate with collision avoidance and frame image using face detection and predefined rules. They expanded their work by discussing their observations, experiences, and further plans with the robot in follow up studies [11], [12]. Four principles from photography perspective [1] were applied to their system and to many later robot photographer systems. In particular they implemented the rule of thirds, the empty-space rule, the no-middle rule, and the edge rule.

The follow-up work [13] extended the work of Byers et al. [10]–[12] by making the robot photographer interactive with users. A framework introduced in [3] described RGB-D data-based solutions with Microsoft Kinect and the capability of detecting direction via human voice recognition was added in [2]. Campbell et al. [14] introduced a mobile robot system that can automate group-picture-framing by applying optical flow and motion parallax techniques. The photo composing in most of these studies rely on heuristic composition rules or similar techniques. Such settings usually do not adequately account for the correlation between the foreground and the background, as well as the light and color effects of the entire image. They also lack a generalizable framework for varied application scenarios.

Recent studies use UAV technology for drone photography trajectory planning [15]–[17] and drone photography system redesign [18], [19]. Although these studies can make it easier for users to obtain high-quality photos, the methods usually only provide an auxiliary semi-automated photo taking process. Users still need to use subjective composition principles to get satisfactory photos.

One crucial part in autonomous camera photography is view selection (*i.e.* framing). Besides the composition principles [1] that are widely used on robot photography [20]–[23] for view selection, there are many recent efforts in photo quality evaluation and aesthetics analysis using computational methods [24]–[28]. These methods use machine learning algorithms to learn aesthetic models for image aesthetic evaluation. The recent advances in deep learning further elevate the research in this direction. Large datasets containing photo rankings basing on aesthetics and attributes are curated [29]–[31] and they allow training deep neural

networks to predict photo aesthetics levels and composition quality [31]–[36]. Compared with aesthetic evaluation methods based on pre-set rules, these data-driven approaches can handle more general and complex scenes where these rules cannot simply be applied. Therefore, we embrace the deep learning based aesthetic models for view selection in our system.

The success in Altari 2600 video games [37] and AlphaGo [38] showed the power of DRL in solving decision-making problems. DRL also enables the control policy end-to-end learning for virtual agents [39], [40], and real robots [41], [42]. Our study can be simplified as a visual servoing problem of robot special navigation. Several DRL driven visual servoing [43], [44] and navigation [45], [46] scenarios were well studied to allow their autonomous agents to interact with the environments. The previous research proved the capability and aptness of using DRL to learn optimal behaviors effectively and efficiently.

III. SYSTEM OVERVIEW

We designed the framework on the robot to: interact with the user, track the user to the desired location, and adjust the position that can take a well-composed portrait. The framework will be systematically discussed from both the hardware and software perspectives.

A. Hardware

The hardware of the entire framework consists of eight major components shown in Figure 2 and the devices with corresponding models are listed in Table I.

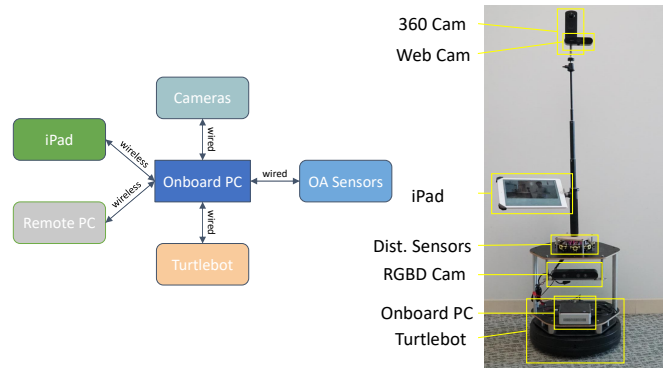


Fig. 2. Hardware system overview.

B. Software

The software framework is built on top of *Robot Operating System* (ROS) [47], and the architecture is shown in Figure 3.

The *Core Node* runs on the on-board PC, and it is in charge of communication between nodes on controlling different hardware and software components.

The *Kinematic Node* controls the linear and angular motion of the robot, and avoids collision with obstacles.

The *Camera Node* provides the framework with vision ability such as real-time video streaming and photo shooting.

TABLE I
THE HARDWARE CONFIGURATION.

Device	Model
360 Camera	Ricoh Theta S
Web Camera	Logitech Brio 4K
RGB-D Camera	Orbbec Astra
Tablet Computer	iPad
Distance Sensors	TeraRanger Tower
Mobile Base	Yujin Turtlebot 2
Onboard PC	CPU: Intel Core i7-5557U@3.1GHz RAM: 8GB GPU: Intel Iris Graphics 6100
Remote PC	CPU: Intel Xeon E5-2603v3@1.60GHz RAM: 16GB GPU: Nvidia Geforce GTX 1080

The Kinematic and Camera nodes both reside on the on-board PC. The vision is essential to the robot photographer, relating to two major modes: *Following Mode* and *Photographing Mode*. The two modes can be activated and switched through the Interaction Node which is presented as an iPad application.

The application has a *Graphical User Interface (GUI)* that takes touch gestures and human poses as input, and gives graphical results and voice prompts as output (Section IV-C). The two modes are implemented with two separate nodes. The *Tracker Node* (Section IV-A) is deployed to the remote PC, which analyzes real-time video stream using Person Detection and Re-Identification Neural Network models to identify and follow the user. The tracking is made with the help of wide-view panorama images provided by the 360 camera and the depth information supported by the RGB-D camera. The *Composer Node* also resides on the remote PC. *Composer Node* (Section IV-B) utilizes a Deep Neural Network (DNN) Composition model to determine: the best target view (template), adjust the robot towards the target view with a Deep Reinforcement Learning (DRL) Template Matching Model, and finally shoot the target and select the best photo from candidates with a DNN Best Frame Selection model.

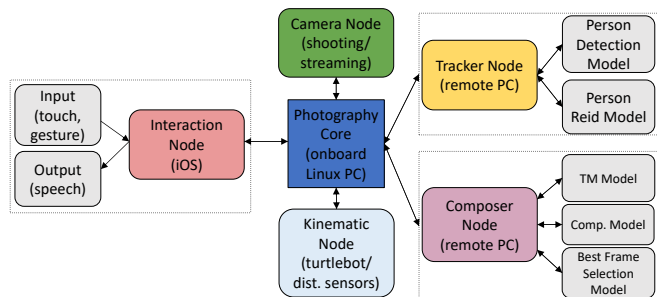


Fig. 3. Software framework overview.

IV. IMPLEMENTATION

A. The Tracker

The *Tracker Module* is responsible for the spatial tracking of the user. It allows the robot to generate instructions for the *Kinematic Node* to follow the user. The architecture of the *Tracker Module* is presented in Figure 4.

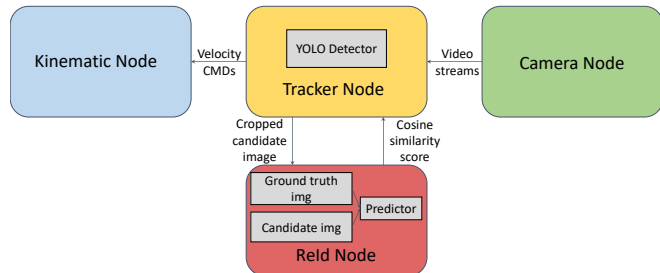


Fig. 4. The tracker module architecture.

Our tracker module utilizes the images from the 360 and RGB-D cameras. The panorama image provides a 360° view for omni-directionally searching and identifying the user around the robot. When the user is located, the robot rotates until the user is centered in the RGB-D camera view. The depth information is used to retrieve the distance between the robot and the target. The distance value determines the linear tracking velocity of the robot. A velocity smoother [48] is used to control the robot’s acceleration. Tracking is not activated when the user is within “operating zone”, which is within 0.5 meters around the robot. When obstacles less than 0.5 meters are detected by OA sensors, the linear velocity of the robot decreases to zero. When the tracking target is missing from the view, the robot stops and waits for the target to appear or the tracker to be reset.

The *Tracker Node* uses YOLO [49] to generate candidate person bounding boxes for the input panoramic images. The candidate person images are cropped out and broadcast to the *ReId Node* that uses the person ReId model from [50] to compare the reference person image with each candidate person image, and predicts cosine similarity scores. The candidate person most similar to the ground truth person is considered to be the target with a goal drop threshold (0.80). The reference person image is initially set when the tracker is activated by the user, and continuously updated when the score of a candidate is higher than the threshold (0.95). The processing time for the tracker node is about 0.04 seconds for a 1440 × 360 input image. The velocity commands generated by the *Tracker Node* are broadcast and received by the *Kinematic Node* for robot motions.

B. The Composer

After the robot tracks the user to a desired photography location, photo composing can be activated to automatically adjust the robot position and take photos for the user based on (a) the static pre-defined templates (Figure 7a) or (b) dynamically proposed well-composed views (Figure 9d). The pipeline for the photo composing process is described in Figure 6.

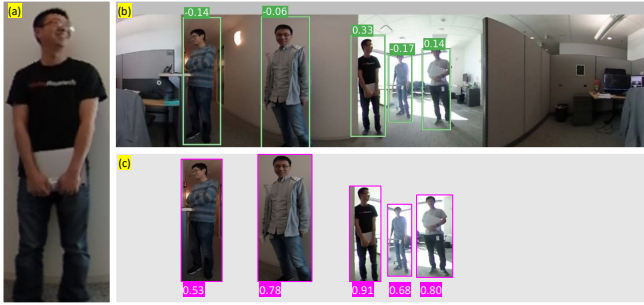


Fig. 5. *Person ReID* results. The ground truth person image in (a) is the query person to track. Part (b) demonstrates the proposed bounding boxes of candidates, and their similarity scores using the joint ReID model [51]. The reference person images in (c) are generated with YOLO [49]. The similarity scores in (c) are predicted with the triplet-loss ReID model [50]. Both methods can provide correct result. Note that the scores in (b) and (c) are not normalized to the same scale.



Fig. 6. The pipeline of photo composing.

1) *Template Generation*: A template contains information that is used to guide the robot composing a user satisfied photo. We use three input data from a template: (a) the location of the person in the photo, (b) the size of the person in the photo, and (c) the pose of the person in the photo. The template can be chosen manually or generated automatically.

There is a set of pre-defined templates varying in location, size, and pose for the user to manually pick from the system. Figure 7a demonstrates a pre-defined template with a cartoon avatar. With this template set, the robot moves around to compose the final photo that matches the template (Section IV-B.2) for the user (Figure 7b). The pre-defined template is not necessarily a cartoon image; the system supports any single person photo with proper aspect ratio.

The system also provides a dynamic template generation solution for photo composing. The novel modular solution enables autonomous photographing with the robot. The solution requires the panorama photos from the 360 camera and the final capture with the high-quality webcam. An example panorama photo is shown in Figure 8a. The panorama photo is cropped and remapped with the method described in [52] to form a collection of candidate templates (Figure 8b). The candidate templates are generated with different levels of distance and yaw angles. Each candidate template is guaranteed to contain the person target. The remapping de-warps the images with the camera parameters of the webcam to make sure the view is reachable from webcam.

The panorama photo processing is performed on the remote PC, and the procedures of finding and using the best template is shown in Figure 9. Candidate templates (Figure 9a) are passed through an off-the-shelf photo evaluation model (Figure 9b). The View Evaluation Net presented in [31] is used in the system on the remote PC, which evaluates and scores photos basing on compositions

(Figure 9c). The candidate template with the highest score is chosen (Figure 9d). Once the robot finishes Template Matching (Section IV-B.2), the final shot is taken with the webcam (Figure 9e). The system is designed to be modular, so that the photo evaluation model can be swapped adaptively with needs.



Fig. 7. Pre-defined templates can be used for photo composing. (a) demonstrates an example of a pre-defined template. (b) shows the final capture with the template from (a).

2) *Template Matching*: The template matching is a process of making the position, size, and pose of the target person in the webcam view as similar as possible to the person in the template by moving the robot to the appropriate location. The similarity can be estimated by comparing the distance of human pose key-points between the current webcam view and the template. OpenPose [53] is used to extract the pose key-points and Figure 10 demonstrates a simple example of template matching procedures. The template is shown in the bottom right corner in each camera view. Figure 10a is the initial camera view of the robot. The robot first turns right. The target is centered in the camera view as shown in Figure 10b. The robot then moves forward. The target becomes bigger as shown in Figure 10c. The robot moves forward again. The size of the target increases as Figure 10d presents. The robot eventually turns right to reach the end state. The target matches the template with small error below the threshold as shown in Figure 10e.

The goal of template matching is to get the robot to capture the desired view with the fixed webcam. The input is the visual information from the webcam video stream. The output is a sequence of robot motor actions. The template matching can be simplified as a problem of robot special navigation. Deep Reinforcement Learning (DRL) has been successfully applied to the robot navigation problem in recent studies [54], [55]. The advantage of using DRL in our framework is that such settings are more adaptive to changes than rule-based and geometric-based solutions. New policies can be simply re-trained with a tuned reward function.

In typical reinforcement learning settings, the agent receives (s_t, a_t, r_t, s_{t+1}) at each time t when interacting with the environment, where s_t is the current state, a_t is the action the agent makes according to its policy π , r_t is the reward produced by the environment based on s_t and a_t , and s_{t+1} is the next state after transitioning through the environment. The goal is to maximize the expected cumulative return from each state s_t in $\mathbb{E}_\pi[\sum_{i \geq 0} \gamma^i r(s_{t+i}, a_{t+i})]$, where γ is the discount factor in $(0, 1]$.

For the template matching, the actions are robot discrete linear and angular velocities. The reward function can be set as an exponential function in Equation 1.

$$r = e^{-\alpha \|v-v'\|}, \quad (1)$$

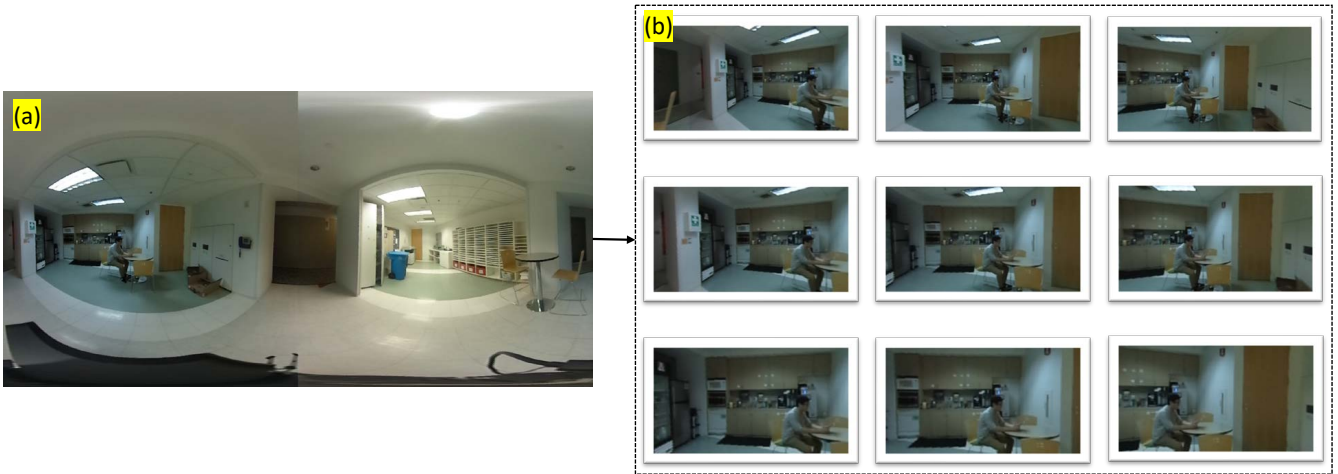


Fig. 8. The panorama photo in (a) is cropped and de-warped to form a collection of candidate templates in (b).

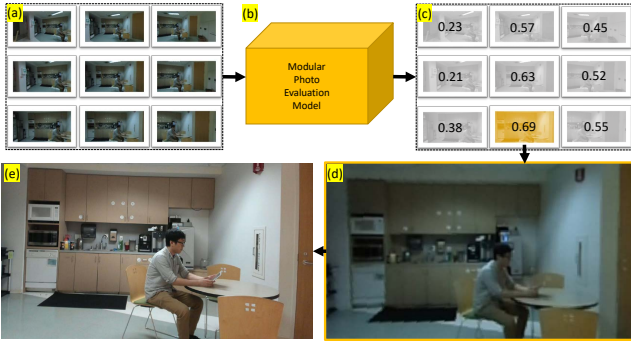


Fig. 9. Finding the best template. Part (a) presents a collection of candidate templates. The candidate templates are passed through a modular photo evaluation network in (b). Scores are retrieved from the network for each candidate template in (c). The template with the highest score is chosen to be the template in (d). The final photo is composed with the webcam in (e).

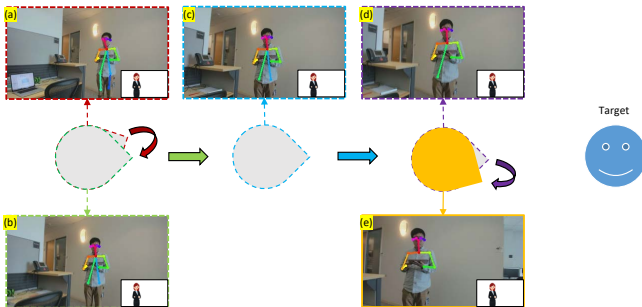


Fig. 10. A step-by-step template matching example.

where v and v' represent the current target keypoint vector and the goal template keypoint vector extracted with OpenPose [53] respectively, and α is a constant factor ($2.5e-03$) to adjust the L^2 distance scale. The DRL training and experiment is discussed in Section V.

The architecture of the composer module is demonstrated in Figure 12. The *Observer* receives a webcam image from the *Camera Node*, and extracts out human pose key-points as an observation. The *Composer Node* passes the observation

to the *Action Parser Node*, which uses a DRL model to predict the next robot action based on the current observation. When an action is decided, the *Action Parser Node* sends it back to the *Composer Node*. The *Responder* parses the robot action into corresponding velocity commands and broadcasts out for execution. The composer module runs on the remote PC.

3) *The Final Photo Acquisition*: After the current webcam view matches the template, the robot starts capturing photos. The webcam shutter can be triggered with a special pose if a pre-defined pose template is manually selected in the beginning of composing (Figure 11). The pose triggering is implemented by comparing the similarity of pose key-point [53] coordinates between the person in template and the person in webcam view. A set of candidate photos are taken during the photo acquisition. The user can manually select the best photo among the candidates, or a best frame selection model such as [56] can automatically decide the best photo for the user based on frame or face scores.



Fig. 11. Triggering camera shutter with a pose. The template in (a) shows a pre-defined pose template. The photo in (b) shows the target making the pose to trigger the webcam shutter.

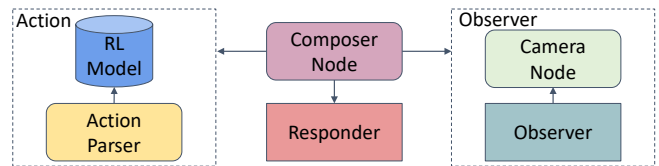


Fig. 12. The composer module architecture.

C. Interaction Node

The *Interaction Node* is deployed on an iPad as an iOS application written in Objective-C. It provides the GUI to the user for Human-Robot Interaction. The application has two modes: the *Following Mode*, and the *Photographing Mode*. When the *Following Mode* is activated, the *Tracker Module* is executed. When the *Photographing Mode* is activated, the *Composer Module* is executed. The touch screen takes user input and gives the user on screen graphical output or text-to-speech prompts. The details of the Interaction Node is shown in the complementary video.

V. EXPERIMENTS

The virtual training environment for Robot Template Matching View Adjustment was set up with synthetic images cropped and processed from panorama photos sampled in real scenes. We trained the view adjustment network with two implementations of DRL. The model was optimized by using action memory and adaptive velocity. We tested the robot photographer in 5 real scenes.

A. Data Preparation

A Ricoh Theta S 360 camera was used to take panorama photos on a 5×5 grid mat (Figure 13). One photo was taken at each grid point. The distance between adjacent grid points was 20 cm. In each scene, 1-5 printed QR codes were placed. A total of 25 photos were collected at each scene; and a total of 15 indoor scenes were selected.

Figure 14 indicates the procedures of setting up the virtual training environment. Each 360° photo (Figure 14a) was cropped every 15° into 24 "rotation images" (Figure 14b). For each scene, there are 600 images. A person in a short video (Figure 14c) was cropped out as a collection of frames with minor differences on poses. 30 person videos were processed. One random frame of the cropped person was attached to the location of the QR code for each rotation image (Figure 14d). The frames of a same person were used for the same QR code in a scene. The scale of the person frame was determined by the size of the QR code in the photo. OpenAI Gym [57] virtual environment was set up with 600 synthesis images for each scene (Figure 14e). The robot rotation action was simulated with the 24 rotation images from the same 360° photo. The robot translation action was simulated with the crops of adjacent 360° photo with the same yaw angle. Nearest neighbour snapping was used in translation simulation. Equation 2 shows the location calculation of a robot translation action.

$$\begin{bmatrix} x \\ y \end{bmatrix} = \lfloor \begin{bmatrix} x' \\ y' \end{bmatrix} + \delta \begin{bmatrix} \sin \theta \\ \cos \theta \end{bmatrix} \rfloor \quad (2)$$

where $[x', y']^T$ present the robot initial location coordinate, and $[x, y]^T$ present the robot location coordinate after a translation action. θ is the yaw angle, and δ is the signed step scalar which presents the size and direction (forward and backward) of one translation action.

B. Training

The DRL model was trained with Advantage Actor Critic (A2C) method, which is a synchronous deterministic variant of Asynchronous Advantage Actor Critic (A3C) [58], and Actor Critic using Kronecker-Factored Trust Region (ACKTR) [59] method respectively using a PyTorch implementation [60]. The environment observation was set to be a vector that contains the selected pose key-point coordinates extracted with OpenPose [53]. The robot translation and rotation actions were expressed as one-hot vector array.

In the training, we noticed that the average total rewards of A2C and ACKTR were always on the same level. However, ACKTR tended to converge at an earlier stage than A2C (Figure 15 and Figure 16). We experimentally improved the average total rewards by using two methods: (a) adaptive velocity and (b) action memory. The robot actions were reflected as linear and angular velocities in one hot vector array. If we increase the array size by setting different levels of velocities to allow the agent to sample, a slightly higher average total reward can be reached for both A2C and ACKTR as Figure 15 shows. Also, if we add previous actions made by the agent to the current observation as memory, both A2C and ACKTR can get higher average total rewards as Figure 16 shows.

C. Testing

We tested the robot photographer at three indoor scenes, and used it to take 20 photos at each scene (ten using pre-defined template and ten using dynamically generated template). Two of the 60 tests failed to conduct the final capture within 30 actions (1 using pre-defined template, 1 using dynamically generated template) and thus are removed from the test samples. The linear velocity of the robot was set to be no more than 0.15 m/s; and the angular velocity was set to be no more than 0.5 rad/s. For dynamically generated template tests, the number of actions has mean $\bar{X} = 12.76$, standard deviation $SD = 4.67$; and the composing time has $\bar{X} = 24.10$, $SD = 8.57$. For pre-defined template tests, the number of action has $\bar{X} = 11.20$, $SD = 4.10$; and the composing time has $\bar{X} = 22.11$, $SD = 7.47$. The resulting photo evaluation was omitted, and the related user study refers to the off-the-shelf composition model study [31].

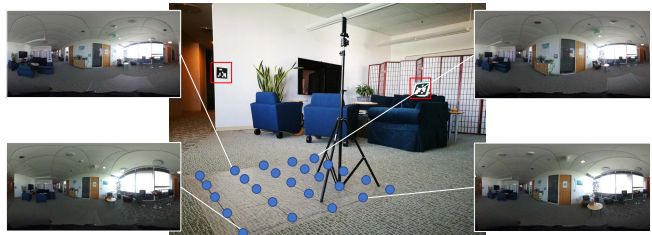


Fig. 13. Training data collection with a 360 camera on a 5×5 grid mat.

VI. CONCLUSIONS

We have developed a novel learning based modular framework for robot photography. The framework allows the robot

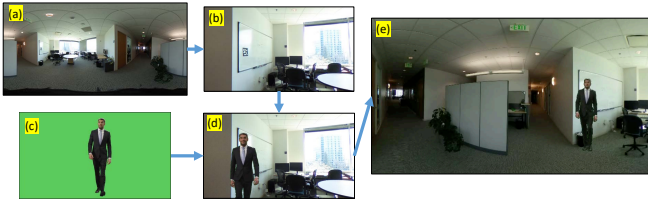


Fig. 14. Virtual training environment setup.

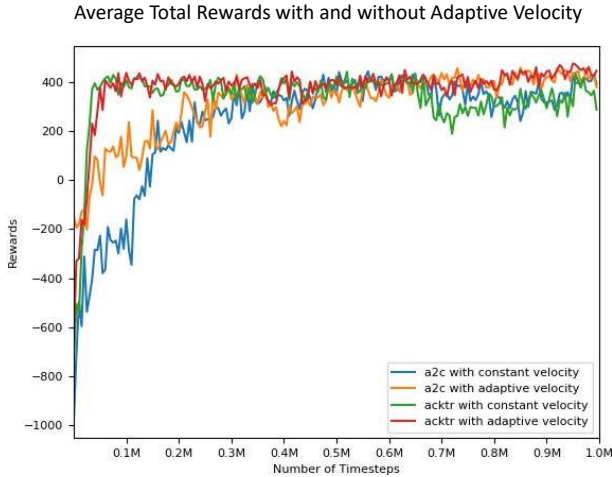


Fig. 15. Improvements on average total rewards in training. The figure shows a comparison of average total rewards with constant velocity and with 3-level adaptive velocity using A2C and ACKTR.

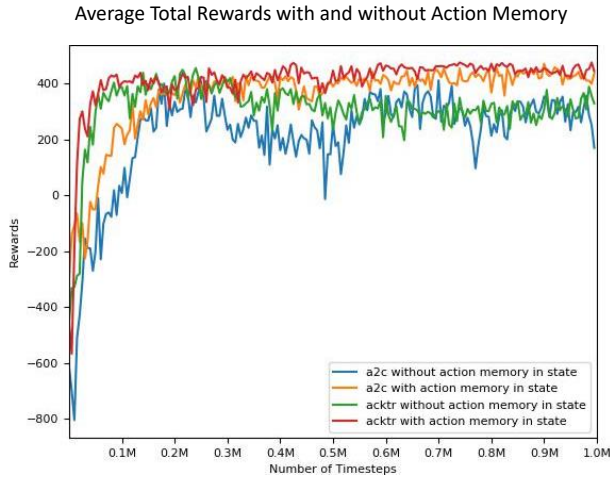


Fig. 16. Improvements on average total rewards in training. The figure shows a comparison of average total rewards with and without previous 5 actions as memory using A2C and ACKTR.

to take well-composed photographs of a person. The robot photographer has a GUI displayed on an attached iPad that gives voice prompts to the user. Our robot can track the user to a desired location. Then, it starts to adjust its position to form the view in a webcam that matches the best template portrait to capture. The best template is searched by using a modular photo evaluation aesthetic model on cropped images of a panorama photo. The view adjustment is driven by a

DRL model basing on template matching. A synthetic virtual environment for navigation training solution was provided.

The system has several *limitations*. It has limited on-board computation power, so that it relies on a powerful remote PC to run DNN models. The Turtlebot that serves as the basis of our system is relatively simple with few degrees of freedom.

Future work includes testing the solution by using a more complex robot. Also, our system currently supports only a single person portrait. New policies would need to be re-trained to get better support on taking group photos. In future work, we also would like to test different photo evaluation aesthetic models, and extend the work to outdoor scenes.

REFERENCES

- [1] T. Grill and M. Scanlon, *Photographic composition*. Amphoto Books, 1990.
- [2] M.-J. Kim, T.-H. Song, S.-H. Jin, S.-M. Jung, G.-H. Go, K.-H. Kwon, and J.-W. Jeon, "Automatically available photographer robot for controlling composition and taking pictures," in *IROS 2010*. IEEE, 2010, pp. 6010–6015.
- [3] M. Zabaraukas and S. Cameron, "Luke: An autonomous robot photographer," in *ICRA 2014*. IEEE, 2014, pp. 1809–1815.
- [4] J. Chen and P. Carr, "Autonomous camera systems: A survey," in *Workshops at the 28th AAAI Conference on AI*, 2014, pp. 18–22.
- [5] Q. Galvane, M. Christie, R. Ronfard, C.-K. Lim, and M.-P. Cani, "Steering behaviors for autonomous cameras," in *Proc. of Motion on Games*, ser. MIG '13. ACM, 2013, pp. 71:93–71:102.
- [6] Q. Galvane, R. Ronfard, M. Christie, and N. Szilas, "Narrative-driven camera control for cinematic replay of computer games," in *Proc. of the 7th International Conference on Motion in Games*, ser. MIG '14. ACM, 2014, pp. 109–117.
- [7] Q. Galvane, M. Christie, C. Lino, and R. Ronfard, "Camera-on-rails: Automated computation of constrained camera paths," in *Proc. of the 8th ACM SIGGRAPH Conference on Motion in Games*, ser. MIG '15. ACM, 2015, pp. 151–157.
- [8] C. S. Pinhanez and A. P. Pentland, *Intelligent studios: Using computer vision to control TV cameras*. Perceptual Computing Section, Media Laboratory, Massachusetts Institute of Technology, 1995.
- [9] C. S. Pinhanez and A. F. Bobick, "Intelligent studios modeling space and action to control tv cameras," *Applied Artificial Intelligence*, vol. 11, no. 4, pp. 285–305, 1997.
- [10] Z. Byers, M. Dixon, K. Goodier, C. M. Grimm, and W. D. Smart, "An autonomous robot photographer," in *IROS 2013*, vol. 3. IEEE, 2003, pp. 2636–2641.
- [11] W. D. Smart, C. M. Grimm, M. Dixon, and Z. Byers, "(not) interacting with a robot photographer," in *Proc. of the AAAI Spring Symposium on Human Interaction with Auto. Syst. in Complex Environments*, 2003.
- [12] Z. Byers, M. Dixon, W. D. Smart, and C. M. Grimm, "Say cheese! experiences with a robot photographer," *AI magazine*, vol. 25, no. 3, p. 37, 2004.
- [13] H. Ahn, D. Kim, J. Lee, S. Chi, K. Kim, J. Kim, M. Hahn, and H. Kim, "A robot photographer with user interactivity," in *IROS 2006*. IEEE, 2006, pp. 5637–5643.
- [14] J. Campbell and P. Pillai, "Leveraging limited autonomous mobility to frame attractive group photos," in *Proc. of the ICRA 2005*, April 2005, pp. 3396–3401.
- [15] C. Gebhardt, B. Hepp, T. Nageli, S. Stevšić, and O. Hilliges, "Airways: Optimization-based planning of quadrotor trajectories according to high-level user goals," in *Proc. of the 2016 CHI Conference on Human Factors in Computing Systems*. ACM, 2016, pp. 2508–2519.
- [16] C. Richter, A. Bry, and N. Roy, "Polynomial trajectory planning for aggressive quadrotor flight in dense indoor environments," in *Robotics Research*. Springer, 2016, pp. 649–666.
- [17] M. Roberts and P. Hanrahan, "Generating dynamically feasible trajectories for quadrotor cameras," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 4, p. 61, 2016.
- [18] H. Kang, H. Li, J. Zhang, X. Lu, and B. Benes, "Flycam: Multitouch gesture controlled drone gimbal photography," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3717–3724, Oct 2018.

- [19] Z. Lan, M. Shridhar, D. Hsu, and S. Zhao, "Xpose: Reinventing user interaction with flying cameras," in *Robotics: Science and Systems*, 2017.
- [20] R. Gadge and K. Karlapalem, "Aesthetic guideline driven photography by robots," in *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, vol. 22, no. 3, 2011, p. 2060.
- [21] B. Gooch, E. Reinhard, C. Moulding, and P. Shirley, "Artistic composition for image creation," in *Rendering Techniques 2001*. Springer, 2001, pp. 83–88.
- [22] C. Cavalcanti, H. Gomes, R. Meireles, and W. Guerra, "Towards automating photographic composition of people," in *Proc. of the IASTED International Conference on Visualization, Imaging, and Image Processing*, 2006, pp. 25–30.
- [23] S. Banerjee and B. L. Evans, "In-camera automation of photographic composition rules," *IEEE Transactions on Image Processing*, vol. 16, no. 7, pp. 1807–1820, 2007.
- [24] R. Datta, D. Joshi, J. Li, and J. Z. Wang, "Studying aesthetics in photographic images using a computational approach," in *ECCV 2006*. Springer, 2006, pp. 288–301.
- [25] S. Dhar, V. Ordonez, and T. L. Berg, "High level describable attributes for predicting aesthetics and interestingness," in *CVPR 2011*. IEEE, 2011, pp. 1657–1664.
- [26] Y. Ke, X. Tang, and F. Jing, "The design of high-level features for photo quality assessment," in *CVPR 2006*, vol. 1. IEEE, 2006, pp. 419–426.
- [27] W. Luo, X. Wang, and X. Tang, "Content-based photo quality assessment," in *ICCV 2011*. IEEE, 2011, pp. 2206–2213.
- [28] M. Nishiyama, T. Okabe, I. Sato, and Y. Sato, "Aesthetic quality classification of photographs based on color harmony," in *CVPR 2011*. IEEE, 2011, pp. 33–40.
- [29] N. Murray, L. Marchesotti, and F. Perronnin, "Ava: A large-scale database for aesthetic visual analysis," in *CVPR 2012*. IEEE, 2012, pp. 2408–2415.
- [30] S. Kong, X. Shen, Z. Lin, R. Mech, and C. Fowlkes, "Photo aesthetics ranking network with attributes and content adaptation," in *ECCV 2016*. Springer, 2016, pp. 662–679.
- [31] Z. Wei, J. Zhang, X. Shen, Z. Lin, R. Mech, M. Hoai, and D. Samaras, "Good view hunting: Learning photo composition from dense view pairs," in *CVPR 2018*. IEEE, 2018, pp. 5437–5446.
- [32] X. Lu, Z. Lin, X. Shen, R. Mech, and J. Z. Wang, "Deep multi-patch aggregation network for image style, aesthetics, and quality estimation," in *Proc. ICCV 2015*. IEEE, 2015, pp. 990–998.
- [33] X. Lu, Z. Lin, H. Jin, J. Yang, and J. Z. Wang, "Rapid: Rating pictorial aesthetics using deep learning," in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 457–466.
- [34] L. Marchesotti, N. Murray, and F. Perronnin, "Discovering beautiful attributes for aesthetic image analysis," *International journal of computer vision*, vol. 113, no. 3, pp. 246–266, 2015.
- [35] L. Mai, H. Jin, and F. Liu, "Composition-preserving deep photo aesthetics assessment," in *Proc. of CVPR 2016*. IEEE, 2016, pp. 497–506.
- [36] L. Kang, P. Ye, Y. Li, and D. Doermann, "Convolutional neural networks for no-reference image quality assessment," in *Proc. of CVPR 2014*. IEEE, 2014, pp. 1733–1740.
- [37] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [38] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al., "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, p. 484, 2016.
- [39] X. B. Peng, G. Berseth, K. Yin, and M. van de Panne, "Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning," *ACM Transactions on Graphics (Proc. SIGGRAPH 2017)*, vol. 36, no. 4, 2017.
- [40] J. H. Libin Liu, "Learning basketball dribbling skills using trajectory optimization and deep reinforcement learning," *ACM Transactions on Graphics*, vol. 37, no. 4, August 2018.
- [41] S. Gu, E. Holly, T. Lillicrap, and S. Levine, "Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3389–3396.
- [42] J. Hwangbo, I. Sa, R. Siegwart, and M. Hutter, "Control of a quadrotor with reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, pp. 2096–2103, 2017.
- [43] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [44] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 421–436, 2018.
- [45] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. J. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre, K. Kavukcuoglu, D. Kumaran, and R. Hadsell, "Learning to navigate in complex environments," *CoRR*, vol. abs/1611.03673, 2016. [Online]. Available: <http://arxiv.org/abs/1611.03673>
- [46] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, "Target-driven visual navigation in indoor scenes using deep reinforcement learning," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, May 2017, pp. 3357–3364.
- [47] M. Quigley, K. Conley, B. P. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, "Ros: an open-source robot operating system," in *ICRA Workshop on Open Source Software*, 2009.
- [48] I. Kostrikov, "Pytorch implementations of reinforcement learning algorithms," <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr>, 2018.
- [49] J. Redmon and A. Farhadi, "Yolo9000: Better, faster, stronger," in *CVPR 2017*, July 2017, pp. 6517–6525.
- [50] A. Hermans, L. Beyler, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [51] T. Xiao, S. Li, B. Wang, L. Lin, and X. Wang, "Joint detection and identification feature learning for person search," in *CVPR 2017*. IEEE, 2017, pp. 3376–3385.
- [52] K. Mo, "Panorama image viewer and cropper," <https://github.com/daerduoCarey/PanoramaImageViewer>, 2018.
- [53] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *CVPR*, 2016.
- [54] J. Zhang, J. T. Springenberg, J. Boedecker, and W. Burgard, "Deep reinforcement learning with successor features for navigation across similar environments," in *IROS 2017*, Sept 2017, pp. 2371–2378.
- [55] K. Mo, H. Li, Z. Lin, and J.-Y. Lee, "The adobeindoornav dataset: Towards deep reinforcement learning based real-world indoor robot visual navigation," *arXiv preprint arXiv:1802.08824*, 2018.
- [56] J. Ren, X. Shen, Z. Lin, and R. Mch, "Best Frame Selection in a Short Video," Adobe Research, Tech. Rep., 2018.
- [57] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *CoRR*, vol. abs/1606.01540, 2016. [Online]. Available: <http://arxiv.org/abs/1606.01540>
- [58] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International conference on machine learning*, 2016, pp. 1928–1937.
- [59] Y. Wu, E. Mansimov, R. B. Grosse, S. Liao, and J. Ba, "Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation," in *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017, pp. 5279–5288.
- [60] Yujin, "Yocs velocity smoother," <https://github.com/yujinrobot/yujinocs>, 2018.