# Visual-Based Forklift Learning System Enabling Zero-Shot Sim2Real Without Real-World Data

Koshi Oishi[*,1], Teruki Kato[1], Hiroya Makino[1], and Seigo Ito[1]

*Abstract*— **Forklifts are used extensively in various industrial settings and are in high demand for automation. In particular, counterbalance forklifts are highly versatile and are employed in diverse scenarios. However, efforts to automate these processes are lacking, primarily owing to the absence of a safe and performance-verifiable development environment. This study proposes a learning system that combines a photorealistic digital learning environment with a 1/14-scale robotic forklift environment to address this challenge. Inspired by the training-based learning approach adopted by forklift operators, we employ an end-to-end vision-based deep reinforcement learning approach. The learning is conducted in a digitalized environment created from CAD data, making it safe and eliminating the need for real-world data. In addition, we safely validate the method in a physical setting using a 1/14-scale robotic forklift with a configuration similar to that of a real forklift. We achieved a 60% success rate in pallet loading tasks in real experiments using a robotic forklift. Our approach demonstrates zero-shot sim2real with a simple method that does not require heuristic additions. This learning-based approach is considered a first step towards the automation of counterbalance forklifts.**

## I. INTRODUCTION

Forklifts are essential in various industries, including factories, logistics centers, ports, and construction sites. In particular, counterbalance forklifts are central to the industry owing to their robustness and power, generating significant demand for automation. Recent advancements have led to the automation of reach-type forklifts, considering their low power requirements and high maneuverability [1]. In contrast, automating counterbalance forklift operations requires advanced controllers that can leverage the versatility and overcome the limited maneuverability of forklifts. As humans learn these tasks through training, applying deep reinforcement learning (DRL) to automation is a natural progression [2], [3], [4]. However, research on its application to forklifts remains limited, likely owing to the risks of conducting experiments in real environments and the lack of training datasets. Therefore, a learning system that can overcome these challenges is required.

DRL has issues such as the dangers of real-world training and low sample efficiency. Sim2real, in which a policy is trained in a digital environment and transferred to the real
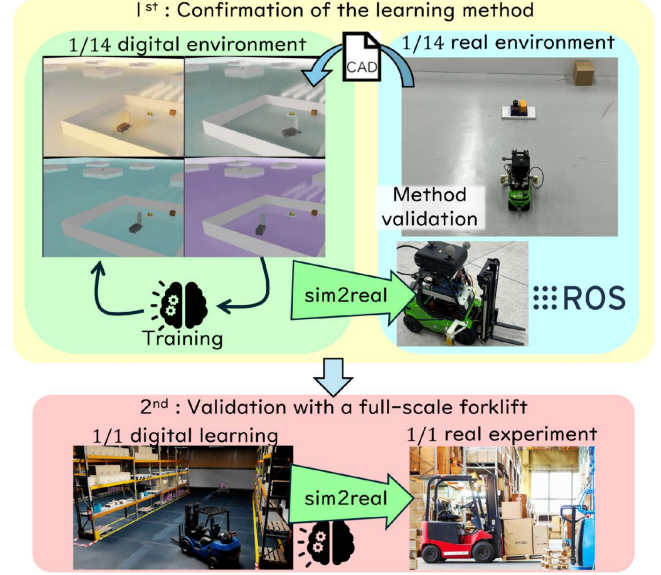


Fig. 1: Proposed concept. To implement DRL-based control on a full-scale forklift, we first test sim2real in a smaller environment. Sim2real is applied to the full-scale forklift after confirming the safety of the method.

world, is a common strategy for these issues [5]. However, learning in a digital environment creates a new challenge known as the domain gap. Various methods have been proposed to bridge this gap. However, as collecting failure data poses significant challenges, methods that do not rely on real-world data are preferable for forklifts. Notable studies without real-world data include those focusing on soccer and mountain climbing with bipedal and quadrupedal robots, respectively [6], [7]. These studies used domain randomization, which randomizes the environment during training [8], [9]. However, the use of visual information in these studies was limited. Moreover, as opposed to these small-scale robots, the first validation of a sim2real method using large machines such as forklifts carries risks, even though safety is verified in the digital environment.

We propose a practical learning system for automating counterbalance forklifts in pallet loading tasks involving the pallet approach and loading decision. In the pallet approach, a policy is optimized using end-to-end DRL to control the forklift based on the same visual and velocity information as that of humans operators. A loading decision policy to perform a loading task is acquired through supervised learning. Our concept involves training and validation in 1/14-scale

[1] The authors are with Toyota Central R&D Labs., Inc., 41-1 Yokomichi, Nagakute, Aichi, Japan.

[*] Corresponding author. e1616@mosk.tytlabs.co.jp

digital and real environments and subsequent extension of the learning method to full scale after confirming the safety in these environments, as shown in Fig. 1. Furthermore, to address the domain gap in which real-world data is lacking, we use domain randomization with NVIDIA Omniverse Isaac Sim (Isaac Sim), which is a highly advanced photorealistic digital environment [10]. Constructing this digital environment requires only CAD data, thereby eliminating the need for real-world data in training. Moreover, the 1/14-scale forklift used in the real environment features front-wheel drive and rear-wheel steering, similar to the full-scale version, and its lift function operates hydraulically. We developed the 1/14-scale forklift using a robot operating system (ROS) [11]. This study represents the first step towards the learning-based control of real counterbalance forklifts. The contributions of our research are as follows:

- We propose a forklift learning system using a photorealistic digital environment that allows for safe research and validation.
- We introduce an end-to-end DRL method for vision-based forklift control that focuses on the pallet approach, which does not require real-world data.
- We develop a dataset construction method for the pallet loading decision policy using a photorealistic digital environment.
- We demonstrate zero-shot sim2real using a 1/14-scale forklift.

The remainder of this paper is organized as follows. Related works are introduced in Section II. Section III presents the target tasks and learning setup for the proposed system. Our proposed learning system that uses a photorealistic simulator is described in Section IV. Section V outlines the demonstration in the real environment. Finally, concluding remarks and future directions are presented in Section V.

## II. RELATED WORK

Research on the automation of counterbalance forklifts has been limited to date. Walter et al. [12] and Teller et al. [13] focused on non-learning-based methods that use multiple LiDAR sensors for environmental perception. However, to the best of our knowledge, few studies have explored real-world implementations of the automation of counterbalance forklifts. Hadwiger et al. [14] conducted DRL using visual inputs on a counterbalance forklift in a digital environment. However, their method did not address the domain gap issue sufficiently and was only evaluated digitally.

End-to-end DRL, which acquires policies through training in a manner similar to human learning, can achieve human-like performance [2], [3]. Trained human operators play a central role in operating counterbalance forklifts, contributing to the adoption of the end-to-end DRL method. Common tasks in DRL applications include navigation and dexterous manipulation [5], [15]. Forklifts require both navigation and dexterous operation because the forks are directly attached to the wheels. In previous research, such tasks were addressed by using a robotic car for imprecise navigation and a robotic arm for dexterous manipulation, whereas our task requires

achieving dexterous navigation solely through the wheels [16].

DRL has been studied extensively in sim2real scenarios because of the dangers associated with real-world experimentation and the low sample efficiency of such experiments [5], [8], [17], [18]. Domain randomization is an effective method for addressing the domain gap in sim2real without relying on real-world data, and it has been widely used in many studies [8], [19], [20], [21]. Among these, the research closest to ours is the study on soccer, in which both navigation and dexterous manipulation were required [6]. Similar to our method, this research progressed to the stage of performing soccer tasks using low-cost onboard cameras [22]. However, NeRF was employed to achieve vision-based tasks, thereby negating the advantage of not using real-world data, which was a key strength in previous studies [23]. Furthermore, although NeRF is effective for static visuals, it is not well suited to moving objects such as pallets or loads [24].

Recent advancements have made it possible to create photorealistic digital environments in which human training can be performed using VR headsets [25]. Mittal et al. [10] and Yu et al. [26] proposed photorealistic learning systems using Isaac Sim. Although these studies used photorealistic digital environments to generate datasets for supervised learning, they did not report on vision-based DRL. Our study leverages the photorealism of Isaac Sim and domain randomization to address the zero-shot sim2real of vision-based DRL, which does not require real-world data and can adapt to dynamic environments.

## III. VISION-BASED FORKLIFT CONTROL VIA DRL

This section first describes the targeted tasks, followed by the learning setup employed by our learning system. The goal is to develop a forklift controller that can perform forklift tasks. This study emphasizes the overall learning framework, including the practical implementation; therefore, we use simple methods for training.

### A. Target task

The target task is the pallet loading operation by a forklift. It involves randomly positioning the forklift within the visible range of the pallet and then executing the loading process. We divide this task into two phases: the approach to the pallet and the lifting operation. Different policies are applied to each phase. We use a simple decision policy to determine whether to lift the pallet. Therefore, the focus is on the learning method related to the forklift-specific approach to the pallet. The observation data used for this task include visual and velocity information, similar to that used by a human operator. In addition, forklift operators lean to the sides to check the alignment between the forks and pallets. Therefore, we installed two cameras on the left and right sides of the forklift to capture images of the forks and pallets.

### B. Policy design for pallet approach

The approach policy is designed to output the throttle and steering of the forklift based on the visual and velocity
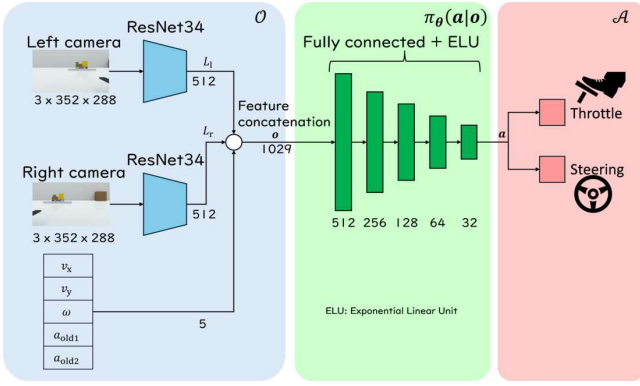
Fig. 2: Approach policy design.

inputs obtained from the camera images and velocity data, respectively. The camera images, which are $352 \times 288$ RGB images, are first resized to $224 \times 224$. These resized images are then converted into feature vectors, $L_l \in \mathbb{R}^{512}$ and $L_r \in \mathbb{R}^{512}$, using a ResNet pretrained on ImageNet [27], [28]. Subsequently, these feature vectors are used as inputs for the policy. This method was proposed as Resnet as representation for reinforcement learning [29]. No special processing is required because we use a standard pretrained ResNet without domain adaptation. In addition, assuming that speed sensors are installed, we utilize the two-dimensional speed information $\boldsymbol{v} := [v_x, v_y]^\top$ and yaw rate $\omega$. Finally, we include the past actions $a_{\text{old}1}$ and $a_{\text{old}2}$. Thus, the observation space $\mathcal{O}$ is a 1029-dimensional vector space.

The policy network $\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{o})$ outputs the normalized throttle and steering commands for the forklift. Here, $\boldsymbol{a} \in \mathcal{A} = [-1,1]^2$ represents the actions and $\boldsymbol{\theta}$ denotes the parameters of the policy network. The network $\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{o})$ consists of five fully connected layers with exponential linear units as activation functions. The overall structure is shown in Fig. 2.

*C. Training method*

We use the on-policy reinforcement learning algorithm proximal policy optimization (PPO) to train the approach policy $\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{o})$ [30]. The following subsections describe the loss and reward functions employed in our method.

*1) Loss function:* The actor and critic networks share a common network, which is optimized using the following single loss function:

$$L_t(\boldsymbol{\theta}) = \mathbb{E}_t \left[ L_t^{\text{PPO}}(\boldsymbol{\theta}) + c_1 L_t^{\text{value}}(\boldsymbol{\theta}) - c_2 H\left(\pi_{\boldsymbol{\theta}}(\cdot|\boldsymbol{o}_t)\right) \right.$$
$$\left. + c_3 L_t^{\text{bound}}(\boldsymbol{\theta}) \right], \tag{1}$$

where $L_t^{\text{PPO}}(\boldsymbol{\theta})$ denotes the loss of the PPO policy, $L_t^{\text{value}}(\boldsymbol{\theta})$ denotes the state-value loss, $H_t\left(\pi_{\boldsymbol{\theta}}(\cdot|\boldsymbol{o}_t)\right)$ denotes the entropy, $L_t^{\text{bound}}(\boldsymbol{\theta})$ denotes the boundary loss, and $c$ denotes the weight. Moreover, $L_t^{\text{PPO}}(\boldsymbol{\theta})$ is determined by

$$L_t^{\text{PPO}}(\boldsymbol{\theta}) = -\mathbb{E}_t \left[ \min\left(\rho_t(\boldsymbol{\theta})\hat{A}_t, \text{clip}(\rho_t(\boldsymbol{\theta}), 1-\varepsilon, 1+\varepsilon)\hat{A}_t\right) \right],$$
$$\rho_t(\boldsymbol{\theta}) = \frac{\pi_{\boldsymbol{\theta}}(\boldsymbol{a}_t|\boldsymbol{o}_t)}{\pi_{\boldsymbol{\theta}_{\text{old}}}(\boldsymbol{a}_t|\boldsymbol{o}_t)}, \tag{2}$$

where the clip function smooths the gradients by limiting the loss changes and $\varepsilon$ denotes the clipping parameter. In addition, $\hat{A}_t$ is computed using generalized advantage estimation [30]. Furthermore, $L_t^{\text{value}}(\boldsymbol{\theta})$ is obtained by

$$L_t^{\text{value}}(\boldsymbol{\theta}) = \mathbb{E}_t \left[ \left(V_{\boldsymbol{\theta}}(\boldsymbol{s}_t) - V_t^{\text{target}}\right)^2 \right], \tag{3}$$

where $V_{\boldsymbol{\theta}}(\boldsymbol{s}_t)$ denotes the state-value function, $V_t^{\text{target}}$ denotes the target state-value function, and $\boldsymbol{s}_t$ denotes the state of the environment, including privileged information. In addition, entropy is introduced to encourage exploration, as follows:

$$H(\pi_{\boldsymbol{\theta}}(\cdot|\boldsymbol{o}_t)) = -\sum_{\boldsymbol{a}} \pi_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{o}_t) \log \pi_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{o}_t). \tag{4}$$

Finally, the boundary loss $L_{\text{bound}}$ is introduced to prevent the actions from taking excessively large values:

$$L_t^{\text{bound}}(\boldsymbol{\theta}) = \|\boldsymbol{\mu}(\boldsymbol{o}_t)\|, \tag{5}$$

where $\boldsymbol{\mu}(\boldsymbol{o}_t)$ denotes the mean vector of the actions output by the policy $\pi_{\boldsymbol{\theta}}(\cdot|\boldsymbol{o}_t)$.

*2) Reward function:* Our reward function consists of positive and penalty rewards for desirable and undesirable states, respectively. In addition, calculating these rewards uses states of privileged information that are available in the digital environment. The deviation from a reference trajectory is used for the positive rewards. The initial position is the location of the forklift at the start of the task, and the terminal position is the pallet. The reference trajectory is based on the approximation of a clothoid curve, which remains fixed throughout the task. The positive reward $R_+$ is defined as follows:

$$R_+ = \alpha_1 \frac{1}{r_d} + \alpha_2 \frac{1}{r_{cd}} + \alpha_3 \frac{1}{r_{c\psi}} + \alpha_4 r_g, \tag{6}$$

where $r_d$ and $r_{cd}$ are the distances from the center of the forks to the pallet and clothoid curve, respectively, $r_{c\psi}$ is the difference between the orientation of the forks and the tangent to the clothoid curve, $r_g \in \{0,1\}$ is a special reward that is assigned when the forks reach the pallet position, and $\alpha$ represents the weight. The components of $R_+$ are illustrated in Fig. 3.

As opposed to general navigation tasks, our task requires inserting the forks without moving the pallet. Therefore, the penalty rewards are defined as follows:

$$R_- = \alpha_5 r_p + \alpha_6 r_v + \alpha_7 r_a + \alpha_8 r_{\text{ini}}, \tag{7}$$

where

$$r_p = \begin{cases} -1 & \textbf{if } \|\boldsymbol{v}_p\| > 0.01, \\ 0 & \textbf{otherwise}, \end{cases}$$

$$r_v = \begin{cases} -(\|\boldsymbol{v}\| - 0.07)^2 & \textbf{if } \|\boldsymbol{v}\| > 0.07, \\ 0 & \textbf{otherwise}, \end{cases}$$

$$r_a = -\|\boldsymbol{a} - \boldsymbol{a}_{\text{old}}\|^2,$$

$$r_{\text{ini}} = \begin{cases} -1 & \textbf{if } \|\boldsymbol{v}\| < 0.05 \textbf{ and } r_d > 0.3, \\ 0 & \textbf{otherwise}, \end{cases}$$

where $\boldsymbol{v}_p$ is the velocity of the pallet and $\boldsymbol{a}_{\text{old}}$ is the previous action. $r_p$ is the penalty for contacting the pallet;

Fig. 3: State used for calculating $R_+$.

TABLE I: Domain randomization targets

| Item | Range |
|------|-------|
| Observed speed | $\pm 10\%$ of value |
| Action | $\pm 10\%$ of command |
| Floor color | $\pm 20\%$ RGB |
| Pallet stand color | $\pm 20\%$ RGB |
| Pallet and load color | $\pm 20\%$ RGB |
| Light intensity | 100–100000 lm |
| Light temp. | 2000–7500 K |

a negative reward is assigned if the pallet starts to move. $r_v$ and $r_a$ are penalties that are designed to bridge the gap between the digital and real environments. Specifically, $r_v$ prevents excessive speed, whereas $r_a$ prevents the actions from becoming erratic. $r_{ini}$ is a penalty to prevent the forklift from freezing at the initial position. Therefore, the reward function can be summarized as follows:

$$R = R_+ + R_-. \tag{8}$$

## IV. FORKLIFT LEARNING SYSTEM

Learning-based controls often have black-box characteristics, which makes it challenging to validate them in high-powered, contact-heavy industrial machinery. Therefore, we constructed digital and real environments on a 1/14 scale. The digital environment was created using Isaac Sim, which supports the development of photorealistic environments [10]. We leveraged the extensive randomization capabilities and parallel learning features of the tool to address the domain gap challenges of reinforcement learning. The real environment was constructed using a ROS [11].

### A. Digital and real environment

We constructed digital and real environments for the task involving approaching a pallet within a 1.8 m square space. The appearances of the digital and real environments are shown in Fig. 4.

The forklift was randomly placed within the green triangle shown in Fig. 4a and made to approach the pallet. This task relied on the onboard cameras mounted on the forklift for visual input. Because the forklift does not appear in the onboard cameras, we simplified its digital model to reduce the implementation complexity. The objects in the digital environment were created using CAD data and specifications, without any real-world image capturing. As reflections are known to affect visual navigation in real environments

negatively, fluorescent lights were modeled and placed as shown in Fig. 4c [31]. The 1.8 m space was enclosed by white walls in the digital and real environments. The critical dimensions, namely the sizes of the pallet and forklift, are shown in Fig. 4f. The specifications of the forklift are discussed in the following section. Isaac Sim allowed for various randomizations within this digital environment. The information that was randomized in this study is listed in Table I. Examples of the floor and lighting randomization are shown in Fig. 5.

### B. Real forklift

For the real-world environment, we modified a 1/14-scale forklift from LESU, which is a China-based manufacturer [32]. The dimensions of this forklift were $100 \times 310$ mm, and the forks measured $75 \times 10$ mm. Similar to counterbalance forklifts, this forklift featured front-wheel drive and rear-wheel steering. In addition, the lift and tilt functions of the fork were powered by hydraulics.

All computers installed on this forklift were replaced with a Raspberry Pi 4 single-board computer to implement the ROS. Raspberry Pi received action $\boldsymbol{a}$, calculated by the ground control computer based on the approach policy $\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{o})$, and transmitted these commands to the respective servos and motors via the servo driver board, as shown in Fig. 6. Two USB cameras were connected to the Raspberry Pi; the captured images were sent to the ground control computer as observation. The cameras were mounted on both sides of the forklift, as shown in Fig. 4f, and were equipped with OV5640 sensors and 60-degree field-of-view lenses.

Velocity information was obtained using a motion capture system. Note that only velocity information was used in our method. The position data were only used for experimental evaluation. All control processes were executed at 15 Hz.

### C. Digital forklift

The forklift in the digital environment was implemented using a simplified design, as mentioned in Section IV-A. The drivetrain components were operated using the Articulations-based controllers of Isaac Sim. The parameters of these controllers were configured to match the speed response of the real forklift to the command inputs closely.

### D. DRL on Isaac Sim

We used "OmniIsaacGymEnvs," which is a parallel learning framework that is based on "rl_games" and is available in Isaac Sim [33], [34]. This tool enables the construction of multiple environments on the same ground plane, thereby improving the sampling efficiency. The randomization of the colors for the pallets and other objects was handled independently for each environment for domain randomization. We synchronized the task resets and shared the values for lighting and floor color randomization across all environments to simplify the implementation.

(a) Digital environment      (b) Side view (digital)      (c) Ceiling (digital)

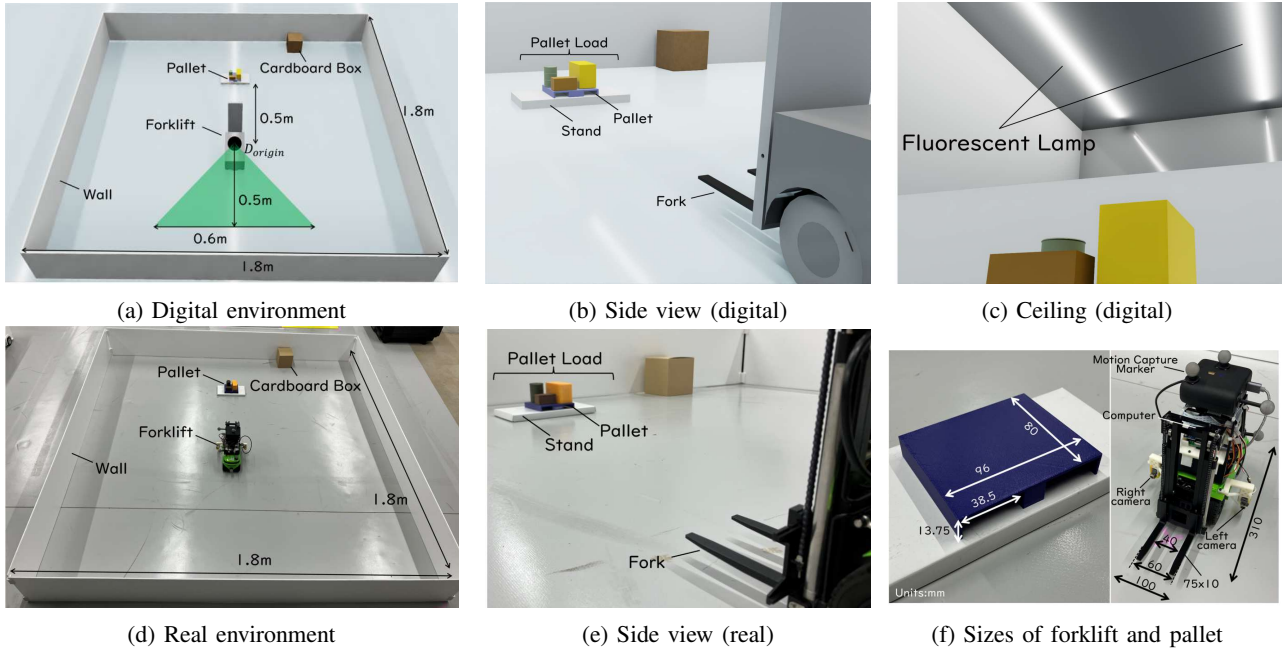(d) Real environment      (e) Side view (real)      (f) Sizes of forklift and pallet

Fig. 4: Digital and real environments. (a), (b), and (c) represent the digital environment. The green triangle in (a) indicates the initial position of the forklift, which is randomly determined at the start of the task. $D_{origin}$ denotes the origin. (d) and (e) represent the real environment. (f) represents the size of the real forklift and pallet. Units are in mm.
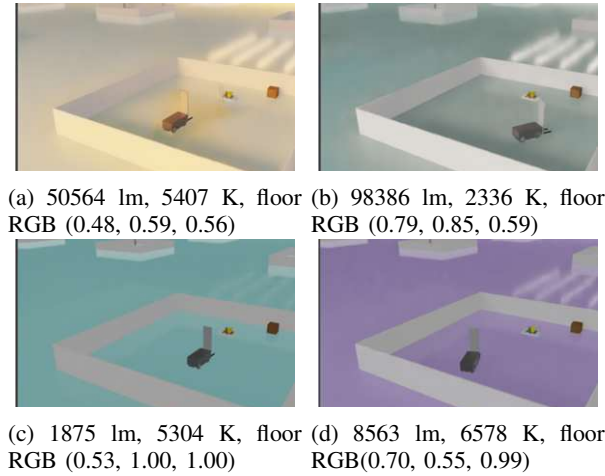


(a) 50564 lm, 5407 K, floor RGB (0.48, 0.59, 0.56)    (b) 98386 lm, 2336 K, floor RGB (0.79, 0.85, 0.59)

(c) 1875 lm, 5304 K, floor RGB (0.53, 1.00, 1.00)    (d) 8563 lm, 6578 K, floor RGB(0.70, 0.55, 0.99)
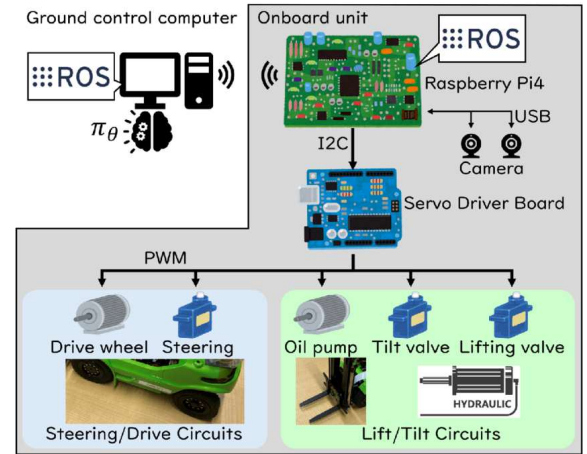
Fig. 5: Example of domain randomization.



Fig. 6: Real forklift control system.

### E. Decision policy for loading

To perform loading, verifying that the approach to the pallet has been successful is crucial. In this study, we designed a decision policy that activates when the forklift is stationary. The decision policy determines whether to lift the forks after the forklift approaches the pallet using the approach policy obtained through DRL. This policy was trained using binary classification through supervised learning. The training dataset was generated by collecting camera images from successful and unsuccessful attempts within the Isaac Sim environment, as shown in Fig. 7. Consequently, no real-world data were used to create this dataset.

## V. REAL DEMONSTRATION

This section presents the results of deploying our method in a real environment. As our method employs zero-shot sim2real, the policies trained in the digital environment were transferred directly to the real environment.

### A. Setup

We trained the policies in the digital environment described in Section IV using the method outlined in Section III. A ground control computer with 96 GB of RAM, an Intel® Xeon™ Gold 6146 3.2 GHz processor, and an NVIDIA RTX A6000 was used. The task was performed
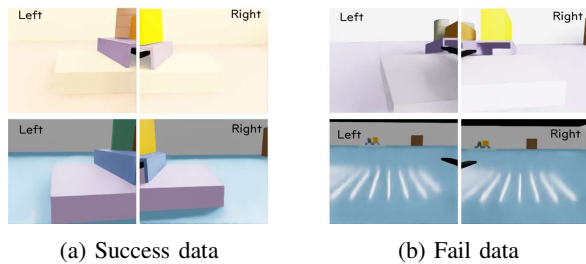
(a) Success data      (b) Fail data

Fig. 7: Dataset for decision policy.

TABLE II: Demonstration results

|  | Ours | Man A | Man B | Man C |
|---|---|---|---|---|
| Success rate (%) | 60 | 50 | 40 | 90 |
| Time (s) | 6.5 | 7.7 | 17.2 | 24.5 |

by positioning the forklift within a triangular area, similar to that during the training in the digital environment, as shown in Fig. 4. In this demonstration, we repeatedly positioned the forklift at an initial position, directed it to approach the pallet, and manually returned it to a new starting position after each task.

The decision policy for determining whether to lift the forks was triggered 3 s after the forklift came to a stop. In the case of a successful approach decision, the forklift followed a sequence of actions: lifting the forks, reversing, and lowering the forks to unload the pallet. Successful completion of the task was dependent on executing the sequence without dropping the pallet load.

We compared the results of our method with those of human-operated trials to demonstrate the difficulty of the task. Human operators controlled the forklift using a joystick, relying on the same camera images and velocity information as the approach policy. The operators were provided with a brief explanation and a 3-min practice session before the experiment. As hydraulic operation requires experience, the experiment was considered successful if the operator could insert the forks at least two-thirds into the pallet.

*B. Results*

We performed the task 10 times with our learned policies and the three human operators. The success rates and average times of the pallet approach required for success are shown in Table II. The average success rate for human operators was 60% and the completion time was 16.5 s. Although Man C achieved a success rate of 90%, he employed a cautious driving technique, stopping repeatedly and checking the camera images frequently. This indicates that the task was not straightforward. Our method achieved a 60% success rate in the pallet approach task at the highest speed. Figure 8 shows the position and velocity of the forklift forks with our method during the 10 trials based on the motion capture system. This result suggests that our approach policy learned efficient behavior, accelerating initially and decelerating near the pallet. Furthermore, the success rate of the loading decisions was 90%. A decision error occurred once when
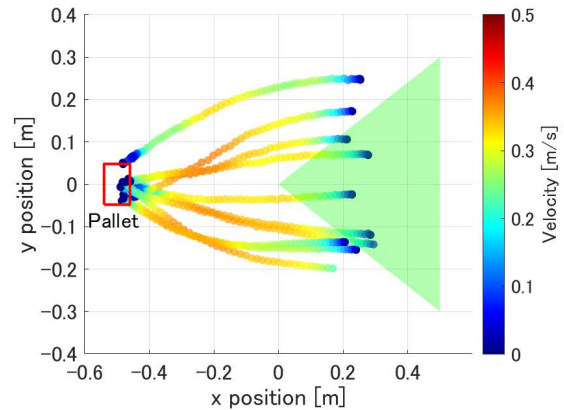


Fig. 8: Trajectory and velocity of the center of the fork. The green triangle indicates the range of randomized initial positions during digital training. As the trajectory represents the path of the fork, the center of the forklift body is positioned 0.2 m further back.
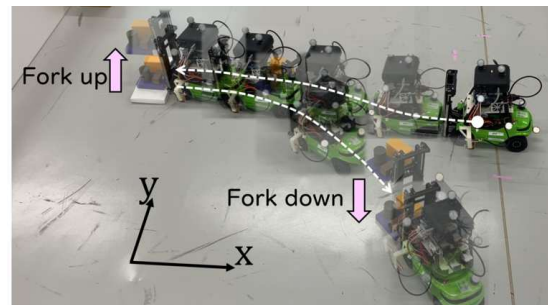


Fig. 9: Demonstration of forklift task.

the forks were inserted only on one side. Snapshots of a successful attempt are presented in Fig. 9.

Our method could produce a controller for forklift tasks using only visual and velocity data, similar to human operation, although the performance was not ideal. Furthermore, the experiment was safe and easy to conduct.

## VI. CONCLUSIONS AND FUTURE WORK

We have proposed a learning system for a counterbalance forklift to perform pallet loading tasks using only cameras and velocity data. In addition, we developed a 1/14-scale forklift with the same configuration as the 1/1-scale one. The forklift allows us to validate the trained policies for safety. Moreover, our method eliminates reliance on real-world data by utilizing a photorealistic environment and domain randomization. Our approach policy achieved a 60% success rate through efficient behavior in loading tasks in real-world experiments. In addition, the forklift could perform loading tasks based on the success or failure of its approach to the pallet using a decision policy learned from a dataset constructed with a digital environment. These policies were implemented through zero-shot sim2real, requiring no heuristic adjustments. Based on these results, this study presented a practical method for applying learning-based control

to counterbalance forklifts, offering significant potential to advance automation across various industries.

Future work will extend our training approach and environment to accommodate more diverse scenarios. In addition, we will incorporate insights from existing studies to improve the accuracy further and validate these results on a full-scale forklift.

## REFERENCES

[1] Linde Material Handling, "Automation for your warehouse." [Online]. Available: https://www.linde-mh.com/en/Solutions/Intralogistics-Automation

[2] 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, pp. 529–533, 2015.

[3] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *Journal of Machine Learning Research*, vol. 17, no. 39, pp. 1–40, 2016.

[4] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza, "Champion-level drone racing using deep reinforcement learning," *Nature*, vol. 620, no. 7976, pp. 982–987, 2023.

[5] W. Zhao, J. Queralta, and T. Westerlund, "Sim-to-real transfer in deep reinforcement learning for robotics: a survey," in *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, 2020, pp. 737–744.

[6] T. Haarnoja, B. Moran, G. Lever, S. Huang, D. Tirumala, J. Humplik, M. Wulfmeier, S. Tunyasuvunakool, N. Siegel, R. Hafner, M. Bloesch, K. Hartikainen, A. Byravan, L. Hasenclever, Y. Tassa, F. Sadeghi, N. Batchelor, F. Casarini, S. Saliceti, C. Game, N. Sreendra, K. Patel, M. Gwira, A. Huber, N. Hurley, F. Nori, R. Hadsell, and N. Heess, "Learning agile soccer skills for a bipedal robot with deep reinforcement learning," *Science Robotics*, vol. 9, 2024.

[7] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Science Robotics*, vol. 7, 2022.

[8] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 23–30.

[9] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 3803–3810.

[10] M. Mittal, C. Yu, Q. Yu, J. Liu, N. Rudin, D. Hoeller, J. Yuan, R. Singh, Y. Guo, H. Mazhar, A. Mandlekar, B. Babich, G. State, M. Hutter, and A. Garg, "Orbit: A unified simulation framework for interactive robot learning environments," *IEEE Robotics and Automation Letters*, vol. 8, pp. 3740–3747, 2023.

[11] M. Quigley, K. Conley, B. 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*, vol. 3. Kobe, Japan, 2009, p. 5.

[12] M. Walter, S. Karaman, E. Frazzoli, and S. Teller, "Closed-loop pallet manipulation in unstructured environments," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2010.

[13] S. Teller, M. Walter, M. Antone, A. Correa, R. Davis, L. Fletcher, E. Frazzoli, J. Glass, J. How, A. Huang, J. Jeon, S. Karaman, B. Luders, N. Roy, and T. Sainath, "A voice-commandable robotic forklift working alongside humans in minimally-prepared outdoor environments," in *2010 IEEE International Conference on Robotics and Automation(ICRA)*, 2010, pp. 526–533.

[14] S. Hadwiger and T. Meisen, "Autonomous load carrier approaching based on deep reinforcement learning with compressed visual information," in *2022 5TH International Conference on Artificial Intelligence for Industries, AI4I*, 2022, pp. 48–53.

[15] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine, "D4rl: Datasets for deep data-driven reinforcement learning," *arXiv preprint arXiv:2004.07219*, 2020.

[16] F. Schmalstieg, D. Honerkamp, T. Welschehold, and A. Valada, "Learning hierarchical interactive multi-object search for mobile manipulation," *IEEE Robotics and Automation Letters*, vol. 8, pp. 8549–8556, 2023.

[17] V. Wiberg, E. Wallin, A. Fälldin, T. Semberg, M. Rossander, E. Wadbro, and M. Servin, "Sim-to-real transfer of active suspension control using deep reinforcement learning," *Robotics and Autonomous Systems*, vol. 179, 2024.

[18] E. Su, C. Jia, Y. Qin, W. Zhou, A. Macaluso, B. Huang, and X. Wang, "Sim2real manipulation on unknown objects with tactile-based reinforcement learning," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 9234–9241.

[19] D. Yarats, I. Kostrikov, and R. Fergus, "Image augmentation is all you need: Regularizing deep reinforcement learning from pixels," in *International Conference on Learning Representations*, 2021.

[20] F. Sadeghi and S. Levine, "CAD$^2$RL: Real single-image flight without a single real image," in *Robotics: Science and Systems XIII*, 2017.

[21] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, " Sim-To-Real via Sim-To-Sim: Data-Efficient Robotic Grasping via Randomized-To-Canonical Adaptation Networks ," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 12 619–12 629.

[22] D. Tirumala, M. Wulfmeier, B. Moran, S. Huang, J. Humplik, G. Lever, T. Haarnoja, L. Hasenclever, A. Byravan, N. Batchelor, N. sreendra, K. Patel, M. Gwira, F. Nori, M. Riedmiller, and N. Heess, "Learning robot soccer from egocentric vision with deep reinforcement learning," in *8th Annual Conference on Robot Learning*, 2024.

[23] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[24] A. Byravan, J. Humplik, L. Hasenclever, A. Brussee, F. Nori, T. Haarnoja, B. Moran, S. Bohez, F. Sadeghi, B. Vujatovic, and N. Heess, "Nerf2real: Sim2real transfer of vision-guided bipedal motion skills using neural radiance fields," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 9362–9369.

[25] NVIDIA, "Isaac Sim - Robotics Simulation and Synthetic Data Generation." [Online]. Available: https://developer.nvidia.com/isaac-sim

[26] Q. Yu, M. Moghani, K. Dharmarajan, V. Schorp, W. C. H. Panitch, J. Liu, K. Hari, H. Huang, M. Mittal, K. Goldberg, and A. Garg, "Orbit-surgical: An open-simulation framework for learning surgical augmented dexterity," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 15 509–15 516.

[27] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[28] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.

[29] R. M. Shah and V. Kumar, "Rrl: Resnet as representation for reinforcement learning," in *the 38th International Conference on Machine Learning*, vol. 139, 2021, pp. 9465–9476.

[30] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[31] F. Xue, Y. Chang, T. Wang, Y. Zhou, and A. Ming, "Indoor obstacle discovery on reflective ground via monocular camera," *International Journal of Computer Vision*, vol. 132, no. 3, pp. 987–1007, 2024.

[32] LESU Model Technology Co., Ltd., "Lesu - manufacturer specialized in manufacturing simulation models." [Online]. Available: https://www.rclesu.com/

[33] D. Makoviichuk and V. Makoviychuk, "rl-games: A high-performance framework for reinforcement learning," May 2021. [Online]. Available: https://github.com/Denys88/rl_games

[34] NVIDIA, "Omniverse isaac gym reinforcement learning environments for isaac sim," 2022. [Online]. Available: https://github.com/isaac-sim/OmniIsaacGymEnvs