

RoboGSim: A Real2Sim2Real Robotic Gaussian Splatting Simulator

Xinhai Li^{1*}, Jialin Li^{2*}, Ziheng Zhang^{3†}, Rui Zhang⁴, Fan Jia³, Tiancai Wang³,
 Haoqiang Fan³, Kuo-Kun Tseng^{1‡}, Ruiping Wang^{2‡}
¹Harbin Institute of Technology, Shenzhen
²Institute of Computing Technology, Chinese Academy of Sciences
³MEGVII Technology ⁴Zhejiang University



Figure 1. **RoboGSim** is an efficient, low-cost interactive platform with high-fidelity rendering. It achieves demonstration synthesis with novel scenes, novel objects, and novel views, facilitating data scaling for policy learning. Additionally, it can perform the closed-loop simulation for safe, fair and realistic evaluation on different policy models.

Abstract

Efficient acquisition of real-world embodied data has been increasingly critical. However, large-scale demonstrations captured by remote operation tend to take extremely high costs and fail to scale up the data size in an efficient manner. Sampling the episodes under a simulated environment is a promising way for large-scale collection while existing simulators fail to high-fidelity modeling on texture and physics. To address these limitations, we introduce the RoboGSim, a real2sim2real robotic simulator, powered by 3D Gaussian Splatting and the physics engine. RoboGSim mainly includes four parts: Gaussian Reconstructor, Digi-

tal Twins Builder, Scene Composer, and Interactive Engine. It can synthesize the simulated data with novel views, objects, trajectories, and scenes. RoboGSim also provides an online, reproducible, and safe evaluation for different manipulation policies. The real2sim and sim2real transfer experiments show a high consistency in the texture and physics. We compared the test results of RoboGSim data and real robot data on both RoboGSim and real robot platforms. The experimental results show that the RoboGSim data model can achieve zero-shot performance on the real robot, with results comparable to real robot data. Additionally, in experiments with novel perspectives and novel scenes, the RoboGSim data model performed even better on the real robot than the real robot data model. This not only helps reduce the sim2real gap but also addresses the limitations of real robot data collection, such as its single-

*Equal contribution

†Project leader

‡Corresponding authors

source and high cost. We hope RoboGSim serves as a closed-loop simulator for fair comparison on policy learning. More information can be found on our project page <https://robogsim.github.io/>.

1. Introduction

Collecting large-scale manipulated data is of great importance for efficient policy learning. Some methods propose to capture the demonstrations as well as the actions through the remote operation [11, 37, 39]. While such operation relatively improves the collection efficiency, it tends to bring extremely large costs with the increasing data size. To solve this problem, some works [14, 34] attempt to generate the synthetic data under the simulated environment, which is further used to learn the manipulation policy. However, those Sim2Real approaches suffer from the large domain gap between simulated and real-world environments, making the learned policy invalid.

Recently, some works introduce the Real2Sim2Real (R2S2R) paradigm for robotic learning [3, 21]. The core insight is to perform realistic reconstruction via radiance field methods, such as NeRF [25] and 3D Gaussian Splatting (3DGS) [15], and insert learned representations into the simulator. Among those methods, the typical approach, Robo-GS [21], presents a Real2Sim pipeline and introduces a hybrid representation to generate digital assets enabling high-fidelity simulation. However, it lacks the demonstration synthesis on novel scenes, views, and objects, as well as verification as policy learning data. Moreover, it fails to perform closed-loop evaluation for different policies due to the misalignment between the latent representation, simulation, and real-world spaces.

In this paper, we develop a Real2Sim2Real simulator, called RoboGSim, for both high-fidelity demonstration synthesis and physics-consistent closed-loop evaluation. It mainly includes four parts: Gaussian Reconstructor, Digital Twins Builder, Scene Composer and Interactive Engine. Given the multi-view RGB image sequences and MDH [6] parameters of the robotic arm, Gaussian Reconstructor is built upon 3DGS [43] and reconstructs the scene and objects. Then, the Digital Twins Builder performs the mesh reconstruction and creates a digital twin in Isaac Sim. In Digital Twins Builder, we propose the layout alignment module to align the space between the simulation, real-world, and GS representation. After that, the Scene Composer combines the scene, robotic arm and objects in simulation, and renders the images from new perspective. Finally, in the Interactive Engine, RoboGSim works as the Synthesizer and Evaluator to performs the demonstration synthesis and closed-loop policy evaluation.

RoboGSim brings many advantages compared to existing (Real2)Sim2Real frameworks. It is the first neural sim-

ulator that unifies the demonstration synthesis and closed-loop evaluation. RoboGSim can generate realistic manipulated demonstrations with novel scenes, views, and objects for policy learning. It can also perform closed-loop evaluation for different policy networks, ensuring fair comparison under a realistic environment. In conclusion, our core contributions can be concluded as:

- **Realistic 3DGS-Based Simulator:** We develop a 3DGS-based simulator that reconstructs scenes and objects with realistic textures from multi-view RGB videos. RoboGSim is optimized for some challenging conditions like weak textures, low light, and reflective surfaces.
- **Digital Twin System:** We introduce the layout alignment module in the system. With the layout-aligned Isaac Sim, RoboGSim maps the physical interactions between objects and robotic arms from Real2Sim spaces.
- **Synthesizer and Evaluator:** RoboGSim can synthesize the realistic manipulated demonstrations with novel scenes, views, and objects for policy learning. It can also work as the Evaluator to perform model evaluation in a physics-consistent manner. The experiment results show that our generated data can achieve the same performance as real robot data, which in a way solves the sim2real gap problem. At the same time, in experiments with novel scenes and novel perspectives, our generated data is more effective than real robot data, even with the 2D Aug method. Our evaluator also partially verifying the performance of the real robot model in a closed-loop.

2. Related Work

2.1. Sim2Real in Robotics

The Real2Sim2Real approach fundamentally seeks to address the Sim2Real gap, which remains a persistent obstacle in the transformation from simulation to real world [8, 27]. In order to bridge the Sim2Real gap as much as possible, many feature-rich simulators have emerged in recent years, including [7, 23, 28, 35, 38]. To this end, various datasets and benchmarks have also been proposed for effective policy learning [12, 13, 16, 26].

Previous Sim2Real methods can be broadly classified into three categories: domain randomization, domain adaptation, and learning with disturbances [40]. Domain randomization methods are designed to expand the operational envelope of a robot in a simulator by introducing randomness. The simulation environment should be capable of migration of the aforementioned capabilities in real-world settings [1, 10, 14, 34]. Domain adaptation approaches aim to unify the feature space of simulated and real environments, facilitating the training and migration within the unified feature space [2, 19, 41]. The objective of learning methods introduce the disturbances into the simulated environment, in which the policy of robots is learned. It develops the ca-

capacity to operate effectively in the real world with noise and unpredictability [5, 36].

2.2. 3D Gaussian Splatting in Robotics

As a significant advancement in the field of 3D reconstruction, 3DGS [15] represents the scene as a large set of explicit Gaussian points and combines it with efficient rasterization to achieve high-fidelity real-time rendering, extending the capabilities of NeRF [25].

More recently, a number of studies have explored the use of 3DGS to perform manipulation tasks within embodied simulators and the real world. For example, Mani-Gaussian [22] introduces a dynamic GS framework alongside a Gaussian world model, which respectively represents Gaussian points implicitly and parameterizes them to model and predict future states and actions. Similarly, Gaussian-Grasper [42] utilizes RGB-D images as inputs and embeds semantic and geometric features into 3DGS through feature distillation and geometric reconstruction, thereby enabling language-guided grasping operations. To effectively transfer the knowledge learned in simulation to the real world and reduce the Sim2Real gap, recent works [18, 21, 29] based on 3DGS have appeared. Among them, the most similar to ours are Robo-GS [21] and SplatSim [29]. Robo-GS achieves manipulable robotic arm reconstruction by binding Gaussian points, grids, and pixels, with a primary focus on high-fidelity Real2Sim transfer; however, it provides limited discussion on the Sim2Real phase. SplatSim reconstructs both the robotic arm and objects in the scene and simultaneously verifies the feasibility of the method for Sim2Real tasks. However, it lacks discussions on generating digital twin assets of the objects, which are critical for achieving accurate manipulation.

3. Methods

3.1. Overall Architecture

As shown in Fig. 2, RoboGSim mainly includes four parts: Gaussian Reconstructor, Digital Twins Builder, Scene Composer, and Interactive Engine. Given multi-view images and MDH parameters of the robotic arm, Gaussian Reconstructor (Sec. 3.2) reconstructs scenes and objects using 3DGS, segments the robotic arm, and builds an MDH kinematic drive graph structure to enable accurate motion modeling of the robotic arm. Digital Twin Builder (Sec. 3.3) involves mesh reconstruction of the scene and objects. Through layout alignment, the asset data flow can be interconnected, facilitating the subsequent evaluation in Interactive Engine. Scene Composer (Sec. 3.4) achieves the synthesis of novel objects, scenes, and views. Interactive Engine (Sec. 3.5) synthesizes novel view/scene/object images for policy learning. It can also evaluate the policy networks in a closed-loop manner. Moreover, we can collect

the manipulated data in simulation using VR/Xbox equipment of the real world.

3.2. Gaussian Reconstructor

We employ the 3DGS method to reconstruct static scenes, followed by point cloud segmentation of the robotic arm’s joints. Subsequently, we utilize the MDH dynamic model to control the point clouds corresponding to each joint, facilitating the dynamic rendering of the robotic arm.

3D Gaussian Splatting (3DGS) [15] employs a set of multi-view images as input to achieve high-fidelity scene reconstruction. 3DGS represents the scene as a set of Gaussians and utilizes a differentiable rasterization rendering method to enable real-time rendering.

Specifically, for a scene $\mathcal{G} = \{g_i\}_{i=1}^N$ represented by N Gaussians, each Gaussian can be represented as $g_i = (\mu_i, \Sigma_i, o_i, c_i)$. Here, $\mu \in \mathbb{R}^3$, $\Sigma \in \mathbb{R}^{3 \times 3}$, $o \in \mathbb{R}$ and $c \in SH(4)$ denote the mean, covariance matrix, opacity and color factor, represented by spherical harmonic coefficients, respectively.

During the rendering process, the final color value C of the pixel can be obtained through a rendering method, similar to alpha-blending [15]. It utilizes a sequence of N ordered Gaussians that overlap with the pixel. Such process can be expressed as follows:

$$C = \sum_{i \in N} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (1)$$

$$\alpha_i = o_i \cdot \exp\left(\frac{1}{2} \delta_i^\top \Sigma_{2D}^{-1} \delta_i\right) \quad (2)$$

where α_i is the opacity of the i -th Gaussian. $\delta_i \in \mathbb{R}^2$ denotes the offset between 2D Gaussian center and current pixel. $\Sigma_{2D} \in \mathbb{R}^{2 \times 2}$ represents the 2D covariance matrix.

Modified Denavit-Hartenberg (MDH) [6] convention is a parameterized model to describe the kinematic chain of a manipulator. Each joint and link in the kinematic chain is characterized by a set of parameters. In MDH, a transformation matrix can be constructed for each link, achieving an accurate representation of the manipulator’s pose at each stage of motion. Let x_i, y_i, z_i denote the coordinates of the origin for the i -th joint. For a manipulator, the i -th joint configuration can be represented as:

$$\Theta = \{\beta_i, a_i, d_i, \theta_i\} \quad (3)$$

where β_i represents the twist angle, which is the rotation around the x -axis from the $(i - 1)$ -th joint to the i -th joint. a_i denotes the link length, measuring the distance along the x -axis from z_{i-1} to z_i . d_i is the link offset, indicating the displacement along the z -axis from x_{i-1} to x_i . θ_i represents the joint angle, rotation around the z -axis from x_{i-1} to x_i .

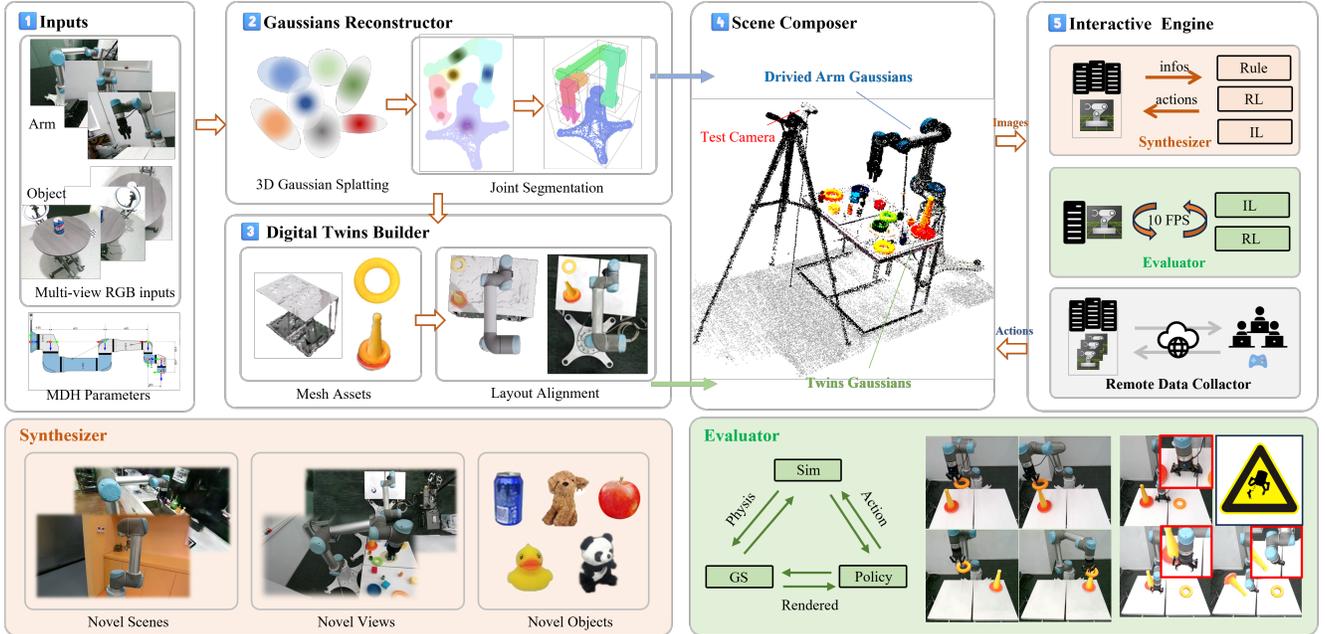


Figure 2. **Overview of the RoboGSim Pipeline:** (1) Inputs: multi-view RGB image sequences and MDH parameters of the robotic arm. (2) Gaussian Reconstructor: reconstruct the scene and objects using 3DGS, segment the robotic arm and build an MDH kinematic drive graph structure for accurate arm motion modeling. (3) Digital Twins Builder: perform mesh reconstruction of both the scene and objects, then create a digital twin in Isaac Sim, ensuring high fidelity in simulation. (4) Scene Composer: combine the robotic arm and objects in the simulation, identify optimal test viewpoints using tracking, and render images from new perspectives. (5) Interactive Engine: (i) The synthesized images with novel scenes/views/objects are used for policy learning. (ii) Policy networks can be evaluated in a close-loop manner. (iii) The embodied data can be collected by the VR/Xbox equipment.

The transformation matrix for each link T_i , using MDH parameters, can be written as:

$$T_i = \begin{bmatrix} \cos \theta_i & -\sin \theta_i \cos \beta_i & \sin \theta_i \sin \beta_i & a_i \cos \theta_i \\ \sin \theta_i & \cos \theta_i \cos \beta_i & -\cos \theta_i \sin \beta_i & a_i \sin \theta_i \\ 0 & \sin \beta_i & \cos \beta_i & d_i \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (4)$$

By sequentially multiplying these transformation matrices, we can obtain the final transformation matrix from the base to the end effector. We segment each joint and then treat all Gaussian points within a joint as a point mass. We further move all Gaussian points within a joint according to T_i , achieving kinematic-driven Control of the Gaussian points.

3.3. Digital Twins Builder

Digital twins should not only map real-world assets but also involve coordinate alignment. Through Real2Sim layout alignment and Sim2GS sparse keypoint alignment, we can digitize the real world, enabling the flow of digital assets between the real, simulated, and GS representation. This facilitates the conversion of digital assets in all directions, achieving comprehensive asset flooding.

3D Assets Generation: We employ two methods to generate 3D object assets. For real-world objects, we capture high-quality multi-view images of the objects using a

turntable and extract matching features with GIM [33] to address issues such as lack of texture and reflections. We then integrate the COLMAP pipeline [32] to obtain the initial SFM point cloud, which is subsequently used for reconstruction by 3DGS. Moreover, for novel objects available on the web, we initially employ generative 3D reconstruction methods [17, 20] to procure 3D gaussians and textured meshes of the objects. Subsequently, we utilize the method in GaussianEditor [4] that applies the diffusion model [31] to facilitate object reconstruction in 3DGS.

Layout Alignment: As shown in Fig. 2, since we follow the local coordinate system of the robotic arm, the world coordinates and Isaac Sim are axis-aligned. We first measure the real-world scene to align the size of the imported table scene in Isaac Sim. In the GS scene, a downward-facing camera is placed 1.6 meters above the base joint to render a segmentation map. For coordinate alignment, we place a downward-facing camera 1.6 meters above the base joint in Isaac Sim. By comparing the rendered scene from the BEV, front and side view segmentation, with the views from Isaac Sim, we adjust the shift to achieve layout alignment.

Sim2GS Alignment: Given the MDH-based transformation matrices T_i^{gs} and simulated transformation matrices

T_i^{sim} , there exists a transformation matrix $T_{gs(i)}^{sim}$ such that:

$$T_{gs(i)}^{sim} = T_i^{sim} \cdot T_{gs}^i \quad (5)$$

To compute the average transformation matrix T_{gs}^{sim} , we use the weighted sum and apply normalization:

$$T_{gs}^{sim} = \frac{\sum_{i=1}^6 w_i \cdot T_{gs(i)}^{sim}}{\left\| \sum_{i=1}^6 w_i \cdot T_{gs(i)}^{sim} \right\|} \quad (6)$$

where w_i is the weight of each joint.

For the target object T_{obj}^{sim} in Isaac Sim, we can transform it into the GS coordinate system using the following formula:

$$T_{obj}^{gs} = T_{sim}^{gs} \cdot T_{obj}^{sim} \quad (7)$$

Camera Localization: To transform the real-world coordinate system into the GS coordinate, we apply the localization approach from GS-SLAM [24]. For a pre-trained GS model, $\mathcal{G} = \{g_i\}_{i=1}^N$, we froze the attributes of 3DGS and optimize the external camera parameters T_C^W .

In camera localization, only the current camera pose is optimized without updates to the map representation. For monocular cases, we minimize the following photometric residual:

$$\mathcal{L}_{pho} = \left\| I(\mathcal{G}, T_C^W) - \bar{I} \right\|_1, \quad (8)$$

where $I(\mathcal{G}, T_C^W)$ represents rendering Gaussians \mathcal{G} from T_C^W , and \bar{I} is the observed image.

3.4. Scene Composer

Scene Editing: To merge the point cloud into the robotic arm scene, the transformation $T[R|t]$ of the marked point is first calculated. Then the coordinates of the point cloud in the new scene are projected into the arm coordinate based on the transformation. Expanding the 3D Covariance Σ in 3DGS into scale s and rotation quaternion q by:

$$\Sigma = q s s^T q^T \quad (9)$$

The ratio r of the transformation can be isolated and extracted as an independent component:

$$r = \sqrt{(RR^T)_{(0,0)}} \quad (10)$$

we can further use it to normalize the rotation matrix R :

$$R_{norm} = \frac{R}{r} \quad (11)$$

The scale attribute s of the Gaussian points is adjusted:

$$s = s + \log(r) \quad (12)$$

Apply the Transformation T to Gaussian point coordinates

$$\mu' = R\mu + t \quad (13)$$

$$\Sigma' = R_{norm} \Sigma R_{norm}^T \quad (14)$$

Object Editing: The transformation here can extend the transformation from the scene editing mentioned above. However, the difference is that the target object’s coordinate center is given by Eq. 7. The coordinate transformation for its Gaussian points can be represented:

$$\mu' = R(\mu - \mu_0) + \mu_0 + t \quad (15)$$

3.5. Interactive Engine

Our interactive engine can work as: Synthesizer and Evaluator. As Synthesizer, it produces large volumes of data with low-cost for downstream policy learning. As Evaluator, it can perform safe, real-time, and reproducible evaluation.

Synthesizer: We use the engine to generate numerous training trajectories, including robotic arm movements and target object trajectories. These trajectories drive the GS to generate massive and photorealistic simulated datasets for policy learning. This diverse data includes novel view renderings, scene combinations, and object replacements.

Evaluator: For trained models, testing directly on physical devices may pose safety risks or incur high costs for reproduction. Therefore, we convert the predicted trajectories into GS-rendered results to efficiently and rapidly evaluate the model’s prediction quality. Specifically, the Isaac Sim [28] outputs an initial state of the target object and robotic arm, and GS renders according to the status. The rendered images are then fed to the policy to predict the next frame’s action. The predicted action is passed to the simulation for kinematic inverse parsing, collision detection, and other physical interactions. Then, Isaac sim sends the parsed six-axis relative pose to the GS renderer, which then sends the rendered result as feedback to the policy network. This serves as visual feedback for predicting the next action, and the process iterates until the task is finished.

4. Experiments

Since there is no benchmarks available for Real2Sim2Real, we construct the following four groups of proxy experiments to comprehensively evaluate the performance of RoboGSim under simulation and real-world. We use UR5 robot arm for all experiments. The robot arm rendering is partially built upon the codebase of Robo-GS [21].

Real2Sim Novel Pose Synthesis verifies whether the robot arm pose captured in the real world can be effectively utilized to achieve precise control in the simulator.

Sim2Real Trajectory Replay checks whether the trajectories collected in the simulator can be accurately reproduced by the real-world robot arm.

RoboGSim as Synthesizer demonstrates the ability of RoboGSim to generate high-fidelity demonstrations with novel scenes, views, and objects, aligning with real world.

	Grasp Suc.			Place Suc.		
	TV	NV(MD)	NS	TV	NV	NS
Real	100%	30%	40%	90%	0%	20%
Real+2D AUG	80%	100%	60%	80%	0%	60%
RoboGSim	100%	70%	100%	90%	0%	90%

Table 1. Performance on ring-toss task in real world. TV denotes test view, NV is novel view and NS is the novel scene. MD means minor deviation.

RoboGSim as Evaluator shows that RoboGSim can effectively perform closed-loop evaluation for policy networks.

Method	Grasp Suc.	Place Suc.
Real-to-Real	100%	90%
Real-to-RoboGSim	100%	30%
Sim-to-Real	80%	0%
RoboGSim-to-Real	100%	90%

Table 2. Cross validation between real world and simulation. ‘‘Sim-to-Real’’ means that testing the VLA model trained with the Isaac sim data on the real world robot arm.

4.1. Real2Sim Novel Pose Synthesis

The objective of the novel pose synthesis is to validate the performance of Real2Sim reconstruction, with a particular focus on the accuracy of the robotic arm’s movements and the fidelity of the image texture. The static scene is reconstructed using the initial pose of the robotic arm from the first frame of GT. The trajectory collected from the real robotic arm is used as the driving force, and we employ the kinematic control for novel pose rendering. As shown in Fig. 3, the results demonstrate that our reconstruction accurately captures both the texture and the physical dynamics of the robotic arm, highlighting the fidelity achieved by RoboGSim. To compare with the video sequence driven by the real robot under the new viewpoint, RoboGSim achieves a 31.3 PSNR and 0.79 SSIM rendering result, while ensuring real-time rendering with 10 FPS.

4.2. Sim2Real Trajectory Replay

To verify whether the trajectories from Isaac Sim can perfectly align with the real machine and RoboGSim, we designed an experiment where the trajectory is collected using Isaac Sim, and then the trajectory is used to drive GS to render a Coke-grasping scene, while the same trajectory is used to drive the real machine to grasp a Coke can. As shown in Fig. 4, the comparison reveals a strong alignment between the simulated policy and the actual physical behavior of the robotic arm, highlighting the effectiveness of the Sim2Real transfer in our system. These results suggest that our simulation can reliably model real-world dynamics, facilitating successful policy transfer from simulation to the real world.

4.3. RoboGSim as Synthesizer

In this section, we use the vision-language-action (VLA) model to validate the effectiveness of synthetic data by RoboGSim. We use the LLAMA3-8B [9] as the LLM and CLIP [30] as the vision encoder. Two-layer MLP is used as the projection network. The VLA model is trained on 8xA100 (80GB) for 1 epoch. The training process is divided into three stages: (1) Pre-training with only the connector enabled, using the LAION-558K dataset. (2) Training with LLM unfrozen using the LLaVA665K dataset. (3) Supervised Finetuning (SFT) with robotic image-action data and the CLIP weight is frozen.

By using the real machine distribution to guide the RoboGSim distribution, we aim to improve the model’s success rate. We perform the experiments on a challenging ring-toss task (see Fig. 7), which is divided into two sub-tasks: picking up the ring and tossing it onto the target. The accuracy requirement for the Z-axis when picking up the ring is within 5mm. For real data, 1,000 samples are collected manually. For a fair comparison, we used 1,000 synthetic samples generated by RoboGSim. During testing, each model was tested 10 times, with three attempts allowed per trial for grasping. If all three attempts failed, the trial was marked as unsuccessful.

As shown in Tab. 1, We compared three models: one trained with real machine data, one trained with real machine data plus 2D AUG, and one trained with RoboGSim. The comparison was made in terms of test view, novel view, and novel scene. The results show that in the test view, the generated data from RoboGSim can achieve zero-shot capability, with performance comparable to the real machine data (both at 90%). In the novel scene case, RoboGSim performed much better than real machine data, reaching 90% compared to the 60% of real machine data. In the novel view experiment, RoboGSim also had less bias compared to the real machine data model. We also compared the effect of adding 2D AUG to the real machine data. After adding AUG, the performance in the test view dropped (90% \rightarrow 80%), but in the novel scene, it improved (40% \rightarrow 60%). However, in the novel view, the bias increased. Pure 2D AUG lacks spatial awareness, and its performance is far worse than that of RoboGSim, which has spatial intelligence. It should be noted that manual collection takes a

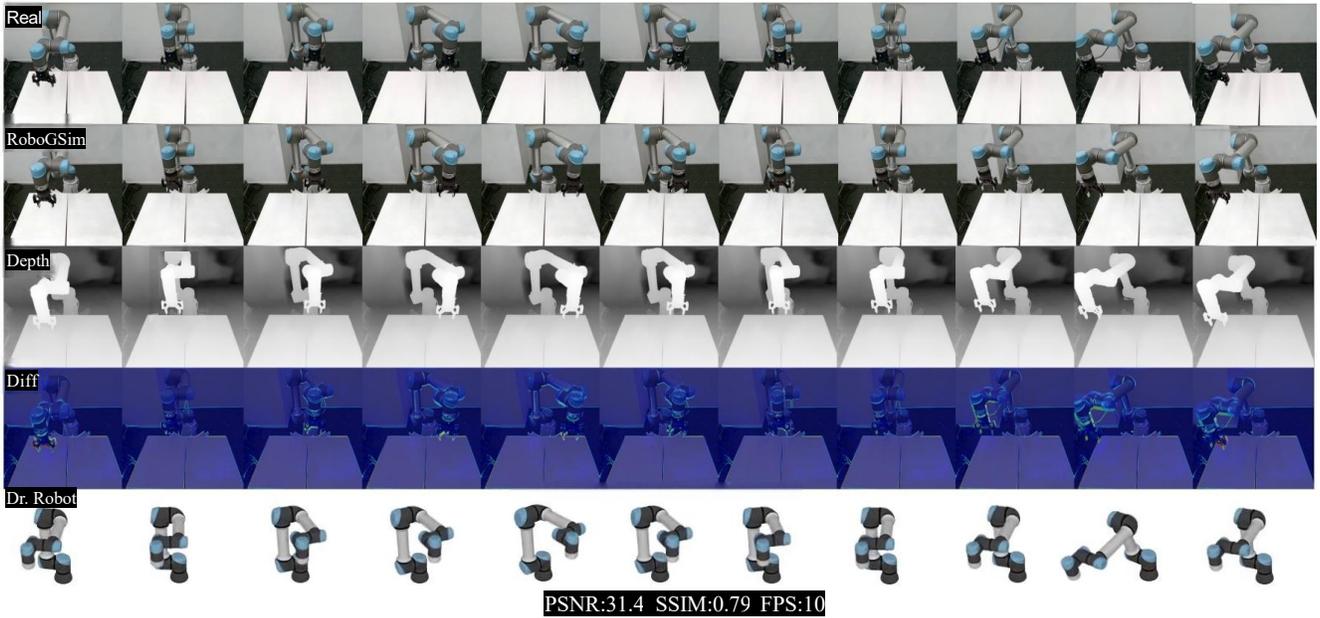


Figure 3. **Real2Sim Novel Pose Synthesis:** "Real" represents the capture of the real robotic arm from a new viewpoint. "RoboGSim" shows the rendering of the novel pose from the new viewpoint driven by the real recorded trajectory. "Depth" shows the rendering depth by GS. "Diff" is the difference calculated between the Real and the rendered RGB images. We compute the pixel distance of the same point between the Real and RoboGSim, which is 7.37.

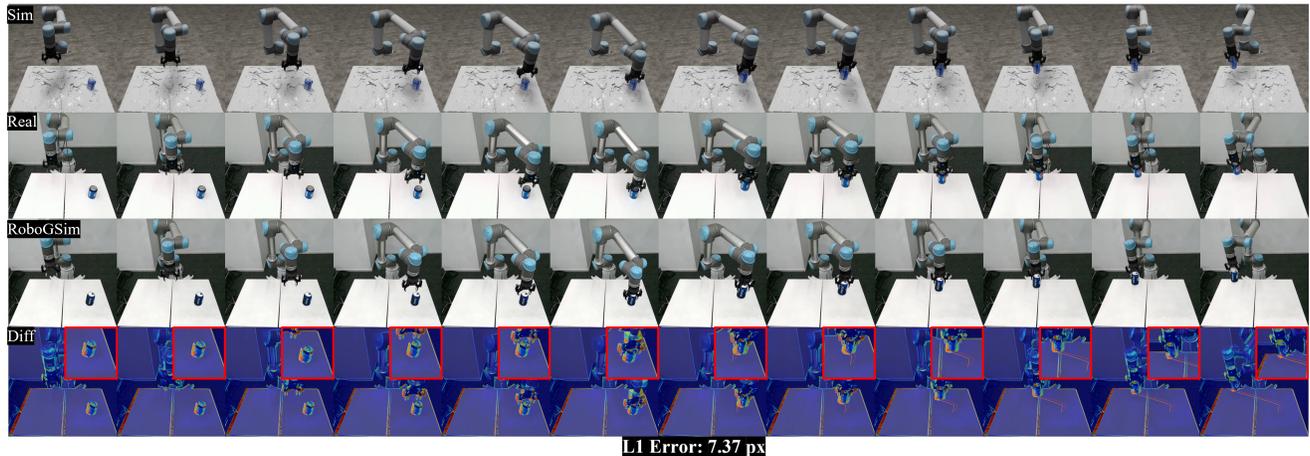


Figure 4. **Sim2Real Trajectory Replay:** The "Sim" row displays the video sequence collected from Isaac Sim. "Real" represents the demonstration driven by the trajectory in simulation. "RoboGSim" is the GS rendering result driven by the same trajectory. "Diff" indicates the differences between Real and the rendered results.

total of 40 hours while RoboGSim only requires 4 hours for synthesis. It is promising to further scale up the data size of synthesis for further performance improvements. Fig. 7 shows the visualization of some success and failure cases. Moreover, we also illustrate some more qualitative analysis for novel scene synthesis. As shown in Fig. 5, we display the results of the physical migration of the UR5 robot arm to new scenes, including a factory, a shelf, and two out-

door environments. The high-fidelity multi-view renderings demonstrate that RoboGSim enables the robot arm to operate seamlessly across diverse scenes.

4.4. RoboGSim as Evaluator

Realistic closed-loop evaluation is crucial for validation and comparison of policy networks. In this section, we mainly explore the effectiveness of using RoboGSim as an Evalu-

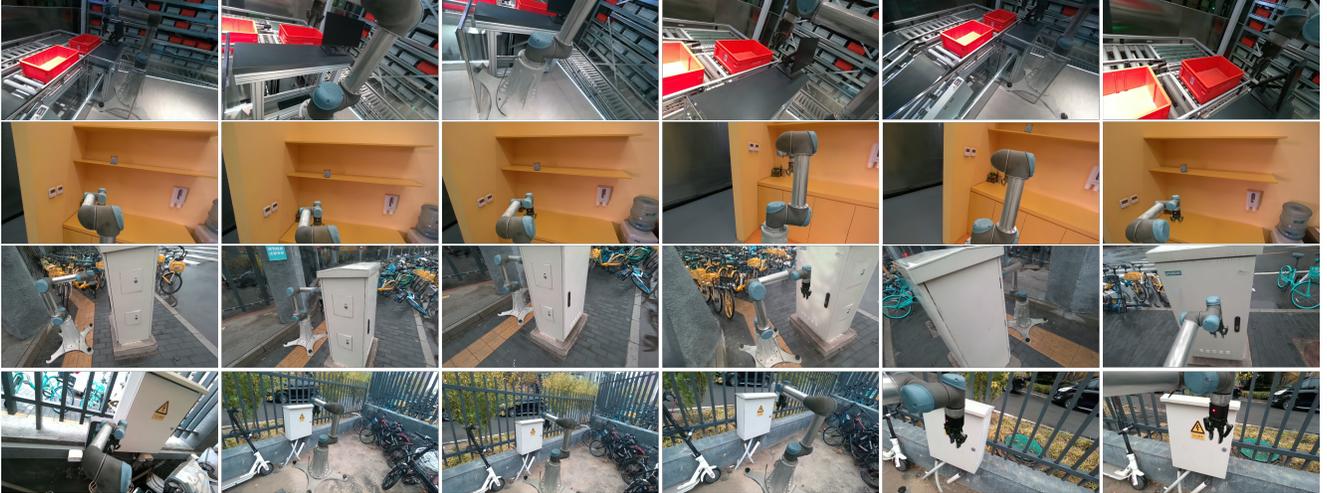


Figure 5. **Novel Scene Synthesis:** We show the results of the physical migration of the robot arm to new scenes, including a factory, a shelf, and two outdoor environments. The high-fidelity multi-view renderings demonstrate that RoboGSim enables the robot to operate seamlessly across diverse scenes.

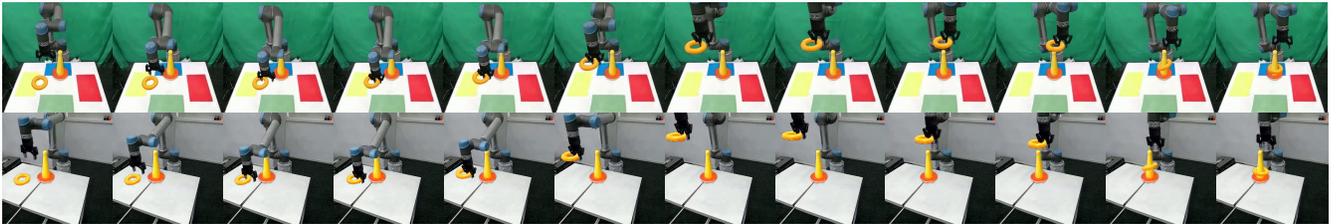


Figure 6. **RoboGSim as Synthesizer:** Rendering of the same training set in novel view and novel scenes.

Method	$L_1 \downarrow$	PSNR \uparrow
3DGS	0.01381	34.19939
2DGS	0.01798	32.36417
PGSR	0.01925	31.72386
MipNeRF360	0.02618	23.51348

Table 3. 3DGS achieves better static reconstruction than 2DGS

ator. It aims to show its high consistency with real-world inference. Given the well-trained VLA model, we deploy it for both real-world robots and RoboGSim simulation. As shown in Fig. 9, our closed-loop simulator RoboGSim can reproduce results similar to those from the real world. For similar bad cases, our RoboGSim can avoid the issues existing in the real world, like violations and collisions. Therefore, our evaluator provides a fair, safe, and efficient evaluation platform for policy.

5. Conclusion and Discussion

In this paper, we built a Real2Sim2Real simulator, based on 3DGS. We also introduce the digital twin system with spa-

tial alignment to enables 3D asset flow. With novel view-point, object, trajectory and scene, our RoboGSim engine can generate high-fidelity synthesized data. Additionally, due to our precise spatial alignment, RoboGSim can serve as evaluator that allows real-time online policy evaluation. Despite its great progress, the current version of RoboGSim has several limitations. It can only simulate rigid objects and the lighting for synthesized objects is not yet fully unified with the robotic arm. Moreover, generating geometrically consistent object meshes remains challenging, which is often key to completing complex manipulation tasks. In the near future, we will explore more advanced mesh extraction methods, further expand the task categories and establish the benchmarks to comprehensively evaluate the performance across diverse scenarios.

6. Acknowledgements

The work was supported by the National Science and Technology Major Project of China (2023ZD0121300).

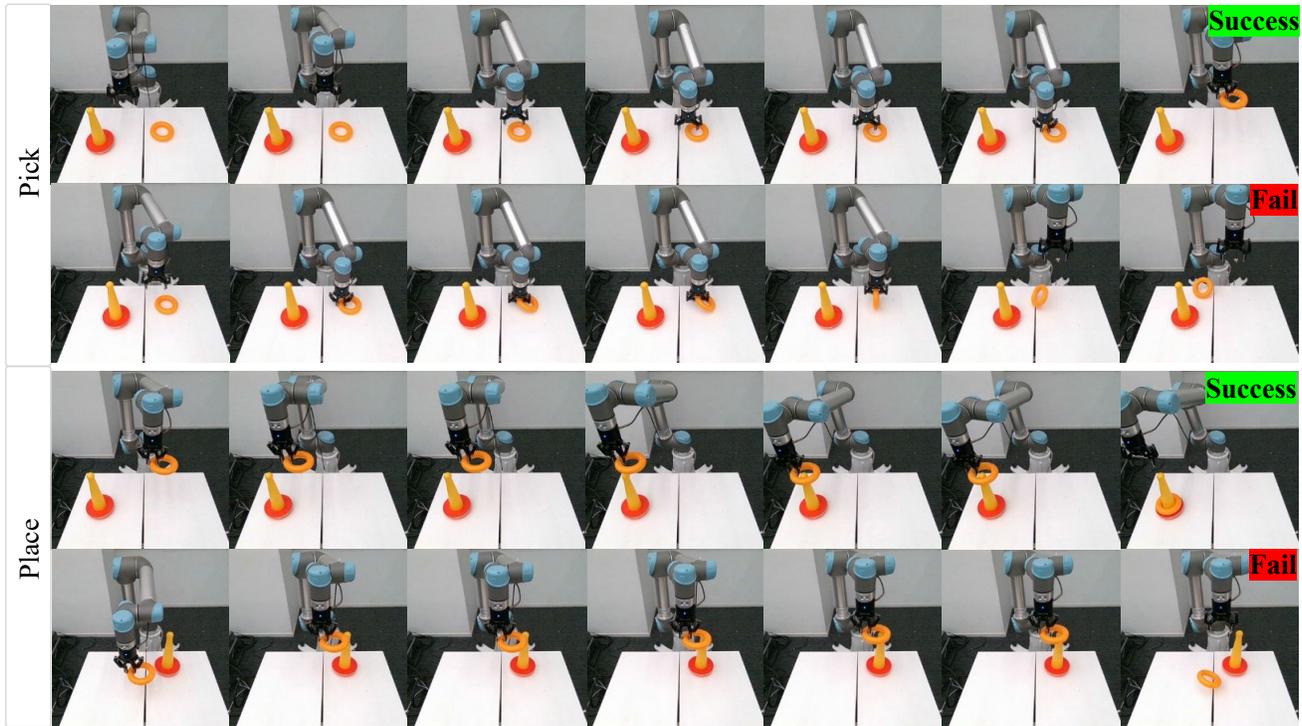


Figure 7. **RoboGSim as Synthesizer:** The first two rows show real robot videos captured from the test viewpoint, illustrating successful and failed cases of the VLA model on the Pick task. The last two rows display real robot videos captured from the test viewpoint, showing successful and failed cases of the VLA model on the Place task.

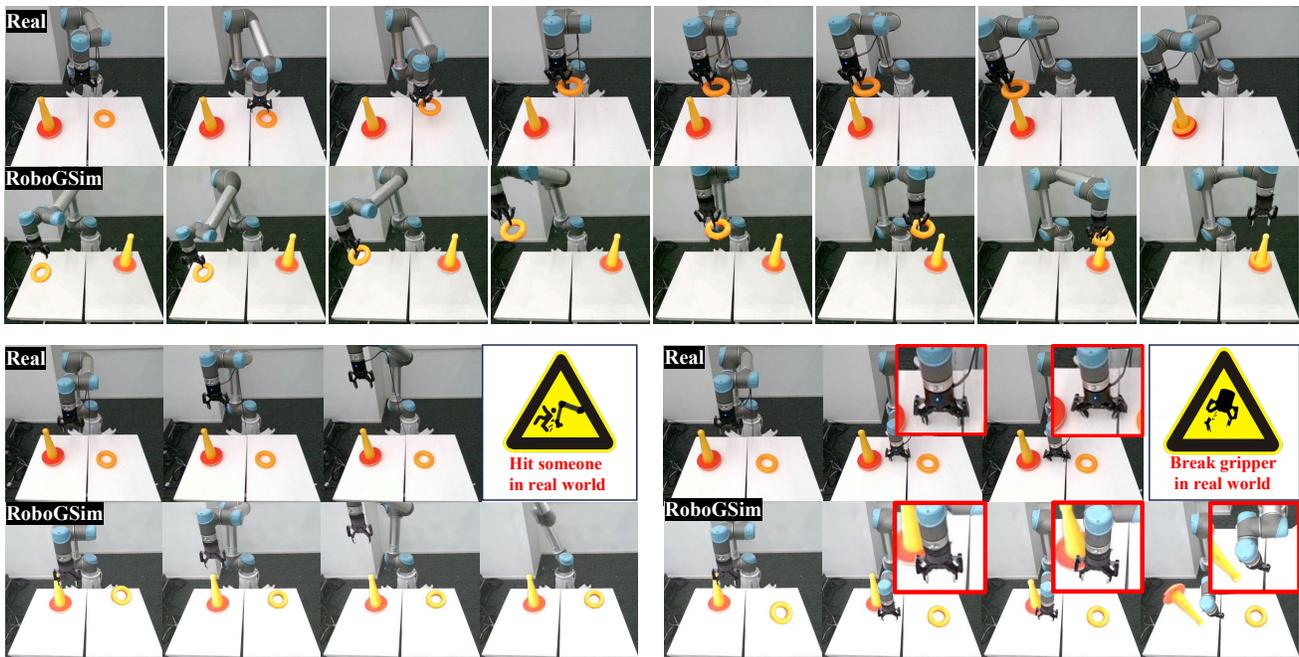


Figure 8. **RoboGSim as Evaluator:** The first two rows, labeled "Real" and "RoboGSim", show the footage captured from the real robot and RoboGSim, respectively. They are both driven by the trajectory generated by the same VLA network. In the third row, the left side shows the real-world inference where the robot arm exceeds its operational limits, resulting in a manual shutdown. The right side shows an instance where a wrong decision from the VLA network, causes the robotic arm to collide with the table. The fourth row presents the simulation results from RoboGSim, which can avoid dangerous collisions.

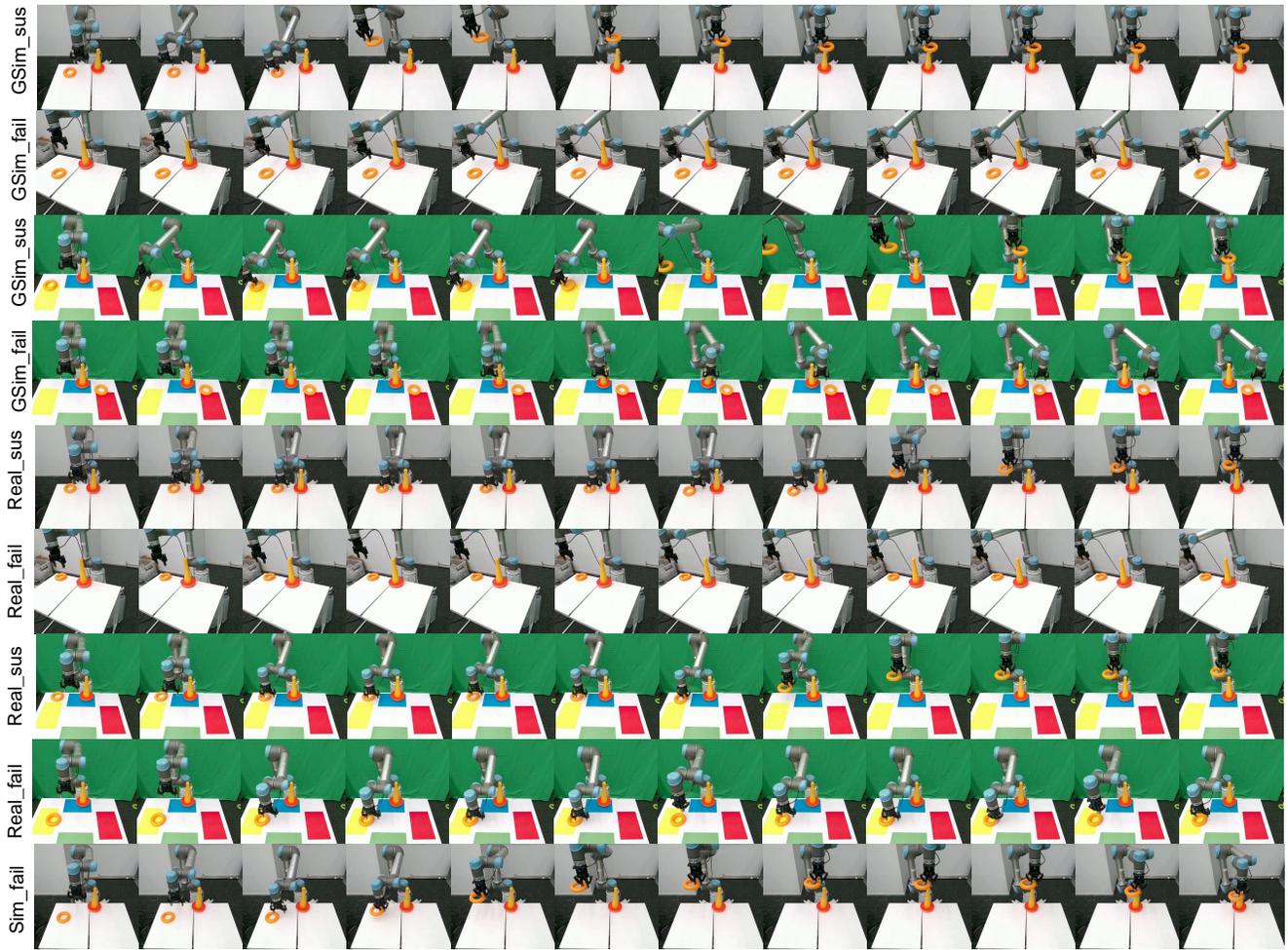


Figure 9. **Exp on Real robot:** The first row, *GSIm_sus*, represents successful cases of RoboGSim data under the test view. The second row, *GSIm_fail*, represents failure cases of RoboGSim data under the novel view. The third row, *GSIm_sus*, represents successful cases of RoboGSim data under the novel scene. The fourth row, *GSIm_fail*, represents failure cases of RoboGSim data under the novel scene. The fifth row, *Real_sus*, represents successful cases of real-world data under the test view. The sixth row, *Real_fail*, represents successful cases of real-world data under the novel view. The seventh row, *Real_sus*, represents successful cases of real-world data under the novel scene. The eighth row, *Real_fail*, represents failure cases of real-world data under the novel scene. The ninth row, *Sim_fail*, represents failure cases of Isaac Sim data under the test view.

References

- [1] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020. 2
- [2] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3722–3731, 2017. 2
- [3] Arunkumar Byravan, Jan Humplik, Leonard Hasenclever, Arthur Brussee, Francesco Nori, Tuomas Haarnoja, Ben Moran, Steven Bohez, Fereshteh Sadeghi, Bojan Vujatovic, and Nicolas Heess. Nerf2real: Sim2real transfer of vision-guided bipedal motion skills using neural radiance fields. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9362–9369, 2023. 2
- [4] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21476–21485, 2024. 4
- [5] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action dif-

- fusion. *The International Journal of Robotics Research*, 2024. 3
- [6] Peter I Corke. A simple and systematic approach to assigning denavit–hartenberg parameters. *IEEE transactions on robotics*, 23(3):590–594, 2007. 2, 3
- [7] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>, 2016–2021. 2
- [8] Konstantinos Dimitropoulos, Ioannis Hatzilygeroudis, and Konstantinos Chatzilygeroudis. A brief survey of sim2real methods for robot learning. In *International Conference on Robotics in Alpe-Adria Danube Region*, pages 133–140. Springer, 2022. 2
- [9] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 6
- [10] Ioannis Exarchos, Yifeng Jiang, Wenhao Yu, and C Karen Liu. Policy transfer via kinematic domain randomization and adaptation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 45–51. IEEE, 2021. 2
- [11] Zipeng Fu, Tony Z. Zhao, and Chelsea Finn. Mobile ALOHA: Learning bimanual mobile manipulation using low-cost whole-body teleoperation. In *8th Annual Conference on Robot Learning*, 2024. 2
- [12] Haoran Geng, Ziming Li, Yiran Geng, Jiayi Chen, Hao Dong, and He Wang. Partmanip: Learning cross-category generalizable part manipulation policy from point cloud observations. *arXiv preprint arXiv:2303.16958*, 2023. 2
- [13] Minh Heo, Youngwoon Lee, Doohyun Lee, and Joseph J. Lim. Furniturebench: Reproducible real-world benchmark for long-horizon complex manipulation. In *Robotics: Science and Systems*, 2023. 2
- [14] Johann Huber, François Héli on, Hippolyte Watrelot, Faiz Ben Amar, and St ephane Doncieux. Domain randomization for sim2real transfer of automatically generated grasping datasets. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4112–4118. IEEE, 2024. 2
- [15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimk uhler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 2, 3
- [16] Vikash Kumar, Rutav Shah, Gaoyue Zhou, Vincent Moens, Vittorio Caggiano, Abhishek Gupta, and Aravind Rajeswaran. Robohive: A unified framework for robot learning. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. 2
- [17] Xinhai Li, Huaibin Wang, and Kuo-Kun Tseng. Gaussdiffusion: 3d gaussian splatting for denoising diffusion probabilistic models with structured noise. *arXiv preprint arXiv:2311.11221*, 2023. 4
- [18] Ruoshi Liu, Alper Canberk, Shuran Song, and Carl Vondrick. Differentiable robot rendering. In *8th Annual Conference on Robot Learning*, 2024. 3
- [19] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. 2
- [20] Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023. 4
- [21] Haozhe Lou, Yurong Liu, Yike Pan, Yiran Geng, Jianteng Chen, Wenlong Ma, Chenglong Li, Lin Wang, Hengzhen Feng, Lu Shi, et al. Robo-gs: A physics consistent spatial-temporal model for robotic arm with hybrid representation. *arXiv preprint arXiv:2408.14873*, 2024. 2, 3, 5
- [22] Guanxing Lu, Shiyi Zhang, Ziwei Wang, Changliu Liu, Jiwen Lu, and Yansong Tang. Manigaussian: Dynamic gaussian splatting for multi-task robotic manipulation. In *European Conference on Computer Vision*, pages 349–366. Springer, 2025. 3
- [23] Viktor Makovychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance GPU based physics simulation for robot learning. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. 2
- [24] Hidenobu Matsuki, Riku Murai, Paul H.J. Kelly, and Andrew J. Davison. Gaussian splatting slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18039–18048, 2024. 5
- [25] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: representing scenes as neural radiance fields for view synthesis. *Commun. ACM*, 65(1):99–106, 2021. 2, 3
- [26] Mayank Mittal, Calvin Yu, Qinxu Yu, Jingzhou Liu, Nikita Rudin, David Hoeller, Jia Lin Yuan, Ritvik Singh, Yunrong Guo, Hammad Mazhar, Ajay Mandlikar, Buck Babich, Gavriel State, Marco Hutter, and Animesh Garg. Orbit: A unified simulation framework for interactive robot learning environments. *IEEE Robotics and Automation Letters*, 8(6): 3740–3747, 2023. 2
- [27] Jean-Baptiste Mouret and Konstantinos Chatzilygeroudis. 20 years of reality gap: a few thoughts about simulators in evolutionary robotics. In *Proceedings of the genetic and evolutionary computation conference companion*, pages 1121–1124, 2017. 2
- [28] NVIDIA. Isaac sim. <https://developer.nvidia.com/isaac/sim>, 2024. Software. 2, 5
- [29] Mohammad Nomaan Qureshi, Sparsh Garg, Francisco Yandun, David Held, George Kantor, and Abhishesh Silwal. SplatSim: Zero-shot sim2real transfer of rgb manipulation policies using gaussian splatting. *arXiv preprint arXiv:2409.10161*, 2024. 3
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervi-

- sion. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6
- [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 4
- [32] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 4
- [33] Xuelun Shen, Zhipeng Cai, Wei Yin, Matthias Müller, Zijun Li, Kaixuan Wang, Xiaozhi Chen, and Cheng Wang. Gim: Learning generalizable image matcher from internet videos. In *The Twelfth International Conference on Learning Representations*, 2024. 4
- [34] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter 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)*, pages 23–30. IEEE, 2017. 2
- [35] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012. 2
- [36] Jingkang Wang, Yang Liu, and Bo Li. Reinforcement learning with perturbed rewards. In *Proceedings of the AAAI conference on artificial intelligence*, pages 6202–6209, 2020. 3
- [37] David Whitney, Eric Rosen, Daniel Ullman, Elizabeth Phillips, and Stefanie Tellex. Ros reality: A virtual reality framework using consumer-grade hardware for ros-enabled robots. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–9, 2018. 2
- [38] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. SAPIEN: A simulated part-based interactive environment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [39] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning Fine-Grained Bimanual Manipulation with Low-Cost Hardware. In *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, 2023. 2
- [40] Wenshuai Zhao, Jorge Peña Queraltá, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 737–744, 2020. 2
- [41] Liming Zheng, Wenxuan Ma, Yinghao Cai, Tao Lu, and Shuo Wang. Gpdan: Grasp pose domain adaptation network for sim-to-real 6-dof object grasping. *IEEE Robotics and Automation Letters*, 8(8):4585–4592, 2023. 2
- [42] Yuhang Zheng, Xiangyu Chen, Yupeng Zheng, Songen Gu, Runyi Yang, Bu Jin, Pengfei Li, Chengliang Zhong, Zengmao Wang, Lina Liu, Chao Yang, Dawei Wang, Zhen Chen, Xiaoxiao Long, and Meiqing Wang. Gaussiangrasper: 3d language gaussian splatting for open-vocabulary robotic grasping. *IEEE Robotics and Automation Letters*, 9(9): 7827–7834, 2024. 3
- [43] Licheng Zhong, Hong-Xing Yu, Jiajun Wu, and Yunzhu Li. Reconstruction and simulation of elastic objects with spring-mass 3d gaussians. In *European Conference on Computer Vision*, pages 407–423. Springer, 2025. 2