

H2R: A Human-to-Robot Data Augmentation for Robot Pre-training from Videos

Anonymous CVPR submission

Paper ID 00076

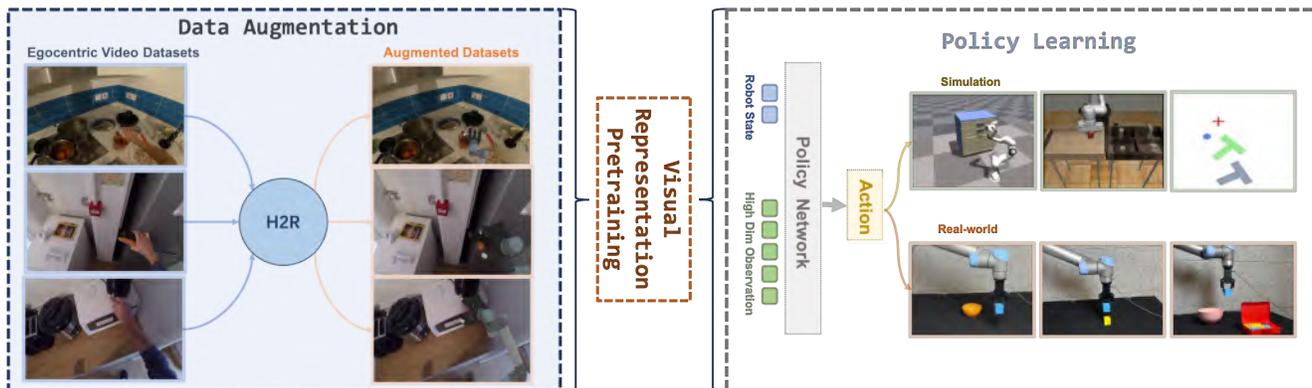


Figure 1. **H2R** is a data augmentation technique designed to enhance robot pre-training by converting first-person human hand operation videos into robot-centric visual data. By bridging the visual domain gap, H2R improves pre-trained visual encoders for downstream robot policies (imitation/reinforcement learning), validated across simulation benchmarks and real-world robotic tasks.

Abstract

001 Large-scale pre-training using videos has proven effective for robot learning, as it enables the model to acquire
 002 task knowledge from first-person human operation data that reveals how humans perform tasks and interact with their
 003 environment. However, the models pre-trained on such data can be suboptimal for robot learning due to the significant
 004 visual gap between human hands and those of different robots. To remedy this, we propose H2R, a simple data
 005 augmentation technique for robot pre-training from videos, which extracts the human hands from first-person videos
 006 and replaces them with those of different robots to generate new video data for pre-training. Specifically, we start by
 007 detecting the 3D position and key points of human hands, which serves as the basis for generating robots in the simulation
 008 environment that exhibit similar motion postures. Then, we calibrate the intrinsic parameters of the simulator camera
 009 to match the camera in the original video and render the images of generated robots. Finally, we overlay these images
 010 onto the original video to replace human hands. Such a procedure bridges the visual gap between the human hand
 011 and the robotic arm and produces an augmented dataset for pre-training.

022 We conduct extensive experiments on a variety of robotic tasks, ranging from standard simulation benchmarks to robotic real-world tasks,
 023 with varying pre-training strategies, video datasets, and policy learning methods. The experimental results show that H2R can improve the representation capability of visual encoders pre-trained by various methods. In imitation learning, H2R consistently enhances the average success rate across different pre-training methods, with improvements ranging from 0.9% to 10.2%. The effect of this improvement is highly stable. In reinforcement learning, most pre-training methods show improvements. Our real-world evaluations across diverse manipulation tasks demonstrate that H2R-enhanced visual representations consistently outperform baseline models, achieving success rate improvements ranging from 6.7% to 15% across all model-task configurations.

1. Introduction

039 Pre-training of generalizable robotic features for object manipulation and motion navigation constitutes a crucial ob-
 040 ject in robot learning. In this paper, we propose H2R, a data augmentation technique designed to enhance robot pre-training by converting first-person human hand operation videos into robot-centric visual data. By bridging the visual domain gap, H2R improves pre-trained visual encoders for downstream robot policies (imitation/reinforcement learning), validated across simulation benchmarks and real-world robotic tasks.

042 jective within the realm of robotics. Inspired by the remark- 095
043 able accomplishments of large scale pre-training in com- 096
044 puter vision [18, 28, 34, 46, 68] and natural language pro- 097
045 cessing [1, 12, 41, 47, 65], many efforts have been devoted 098
046 to harness large-scale data to construct generalizable rep- 099
047 resentations in the robotics field [2, 3, 8]. Nevertheless, 100
048 when it comes to robot manipulation, the process of col-
049 lecting demonstrations is not only labor-intensive but also
050 highly costly [3, 15, 16, 19, 29, 32, 33, 35, 37]; at the same
051 time, there exist many large-scale egocentric video datasets
052 showing how human perform manipulation and navigation,
053 which can serve as a can serve as a cheap alternative of
054 demonstrations for the pre-training of generalizable visual
055 features for robotics.

056 Recent works [46, 63, 69] analyze such egocentric hu-
057 man video datasets such as Ego4D [23], SSV2 [20], and
058 Epic Kitchens [9] with the aim of gleaning prior knowl-
059 edge about object manipulation and enabling the acquisi-
060 tion of general and robust feature representations. How-
061 ever, the gap in visual representations between the human
062 arm and the robotic arm remain largely unaddressed in prior
063 work and can hinder the transferability of models trained on
064 egocentric datasets to robotic systems. Specifically, when
065 utilizing the robot expert data to fine-tune the pre-trained
066 robotic representations for downstream robotic tasks, the
067 model has to learn to bridge the visual gap between the first
068 person human hand and the robots in addition to acquiring
069 nuanced task-specific skills demonstrated in the robot ex-
070 pert data. This would result in increased complexity during
071 the fine-tuning process and suboptimal performance.

072 To mitigate this issue, we propose H2R (as shown in
073 Figure 1), a simple data augmentation method that con-
074 verts videos of **H**uman hand operations into that of **R**obotic
075 arm manipulation. H2R consists of two major procedures:
076 the first part is to generate the robotic arm’s movements to
077 imitate the human hand movements in a video, followed
078 by the second stage that overlays the robotic arm’s move-
079 ments onto the human hand’s movements in the video.
080 Specifically, in the first part, we employ state-of-the-art 3D
081 hand reconstruction model HaMeR [50] to accurately detect
082 the position and posture of the human hand in egocentric
083 videos. Then, we simulate the same robot state in simula-
084 tors to obtain the mask of robot hands. While in the second
085 stage, we use the Segment Anything Model [36] to auto-
086 matically separate human hand from background, and use
087 the inpainting model LaMa [58] to fill the removed hand
088 mask. After that, we align the camera intrinsic parameters
089 of the images detected in HaMeR with those in the simu-
090 lator, and then achieve pixel-level matching between the
091 robotic arm images in the simulators and the human hand
092 images in the egocentric video. Finally, we overlay the
093 robotic arm images captured by the simulator’s camera onto
094 the areas where the human hands are removed. Through

such a process, H2R explicitly reduces the gap between hu-
man and robot hands by creating realistic robotic arm move-
ments that visually mimic human hand actions. It allows
the model to learn the task-specific actions demonstrated by
the human hand, but with robotic arm visual representations
that are more suitable for robotic systems.

For pre-training, we used the SSV2 dataset with 62,500
videos, from which 16 keyframes were randomly sampled
per video for MAE [28] and R3M [46]. Additionally, we ex-
tracted 117,624 action clips from 2,486 videos in the Ego4D
dataset for MPI [68]. Specific settings are detailed in the ex-
perimental section.

We demonstrate the effectiveness of H2R by integrating
it into a holistic policy learning framework. We trained stan-
dard MAE [28], R3M [46], and MPI [68] vision encoders on
egocentric videos obtained by the proposed H2R. We then
freeze the encoder model as a feature extractor and train
both an Reinforcement Learning (RL) policy by employ-
ing mainstream RL learning methods such as PPO [56] and
Imitation Learning (IL) policy with behavior clone and Dif-
fusion Policy [6]. Finally, for the RL policies, we evaluate
them on MVP [51], a closed-loop benchmark, and compare
with results where the encoders are trained on the original
egocentric video data. We observe a significant improve-
ment of the training stability, which brings more effective-
ness in the context of RL policy learning seeing the unstable
nature of the bare RL training. For the IL policies, the BC
policies are trained and tested on Robomimic [45], while the
Diffusion Policy models are trained and tested on their own
baseline. Both of the BC and Diffusion policies showed a
significant improvement on the success rate and stability on
the downstream tasks.

Through extensive real-world experiments, we validate
the effectiveness of H2R in real-world robotic manipulation
tasks. We employ Diffusion Policy [6] (DP) and Equivari-
ant Diffusion Policy [64] (eq-dp) as policy frameworks for
downstream training, integrating pre-trained visual repre-
sentation models MAE and R3M into the policy networks.
The results demonstrate that H2R significantly enhances the
performance of both MAE and R3M-based policies.

Our paper provides three contributions:

- We propose a data-centric pipeline, H2R, to mitigate the
gap between human and robot hands when utilizing large-
scale egocentric video datasets to pre-train generalizable
visual features for robots.
- We apply H2R to SSV2 and Ego4D datasets and train
a visual encoder that is more suitable for robotic tasks.
Built upon this, we yield a robust robotic manipulation
policy through RL and IL training on robot expert data.
- We demonstrate the effectiveness of H2R through exten-
sive experiments on closed-loop benchmarks.

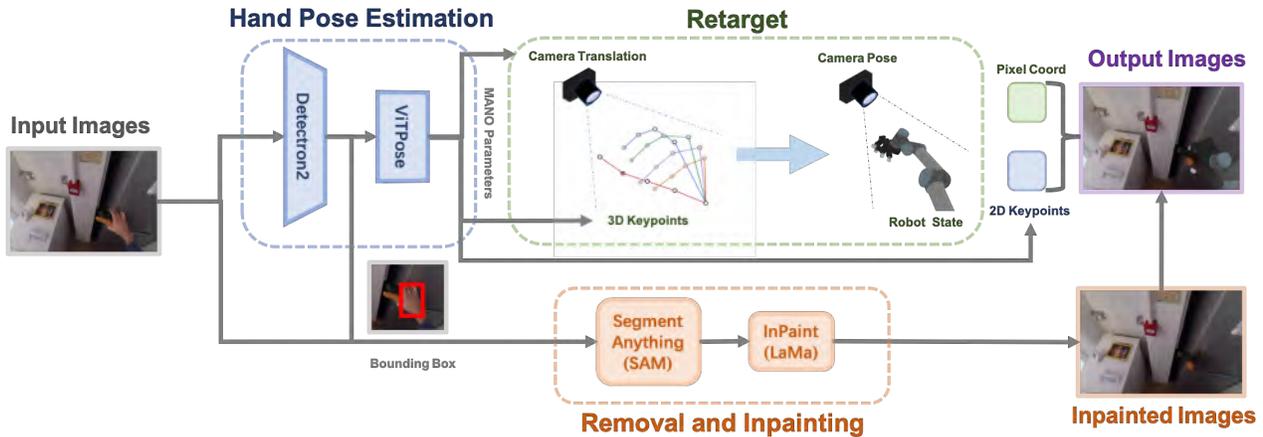


Figure 2. **H2R Pipeline.** H2R involves replacing human hands with robot arms by first using the HaMeR model to detect hand poses and camera parameters. The human hand is then removed using the SAM, and the inpainting model LaMa fills in the gap. A robot hand is constructed based on the detected pose and keypoints, with the camera perspective adjusted to match the original image. Finally, the robot hand is overlaid onto the image, ensuring accurate alignment with the human hand.

2. Related Work

2.1. Robot Policy Learning

Training robot policies [6, 8, 40, 43, 63, 64] in a data-driven manner have been adopted by the robotics community as well as the machine learning community. This serves as a paradigm to automatically yield models for performing robotic tasks including grasping, manipulation, locomotion, navigation, and other complex tasks. [39, 59]. Currently, policy learning methods can be classified into two types: imitation learning (IL)-based [6, 45, 67, 70] and reinforcement learning (RL)-based [24, 56].

IL-based methods [6, 45, 67, 70] train robot policies based on successful demonstration of task execution within the dataset. Supervised by behavior cloning [17, 60] objective along with other auxiliary objectives, the policy predicts a sequence of future actions based on current and past observations. To deal with the non-markovian transition of robot configuration under scenario such as stationary process, ACT [70] employs a temporal fusion of sequence predicted at multiple time steps and thus mitigates the related confounder problem. To deal with the multi-modality nature of robot motion, diffusion models are adopted [6, 67]. For IL-based methods, data diversity contributes largely to the generalizability of model.

On the other hand, RL-based methods [66] resort for the RL paradigm of learning an optimal policy by defining a reward function. These methods formulate robotic manipulation tasks as MDP processes and apply RL algorithms such as PPO [56], SAC [24], and more. Typically, RL for robotics tasks are realized by researchers via RL training in simulator, sim2real transfer, and policy deploying on real robots such as legged robots or aerial drones for locomo-

tion [30, 38, 71], robot arms and dexterous hands for manipulation [62, 66], mobile robots for navigation [7].

For both IL and RL-based methods, a strong feature extractor backbone serves as a cornerstone for learning a robust policy. Therefore, Well-conceived data-centric pipeline is crucial and contributes to the backbone training.

2.2. Visual Encoder Pretraining for Robotics

Researchers investigated visual representation [48] under various perspectives such as model architecture [13, 25], training objective [26, 27], dataset [11, 42, 55, 57], and more. PVR-Control [49] demonstrates the effectiveness of visual representation which surpasses the state representation under the investigated scenarios. RPT [53] explores tokenized representation of transformer and trains the corresponding encoder through masked token-prediction.

Unsupervised training methods such as Masked Auto-Encoder (MAE) [27] and contrastive learning [4, 5] are employed by researchers [46, 51] for training video encoder and enhancing generalizability. Specifically, MVP [51] introduced video representation for downstream RL tasks while R3M [46] combines time-contrastive learning and video-language alignment. To effectively perform language-guided robotic tasks, researchers of Voltron [34] utilize MAE [27] and contrastive learning [4, 5] for low-level control and high-level planning, respectively.

2.3. Data Quality to Learning Method

Yielding a universal visual representation through a data-centric fashion is crucial for the visual encoder along with the policy to generalize to in-domain and even out-of-domain scenarios [66]. Data-centric analysis indicates the importance of data regarding to the pretraining of visual

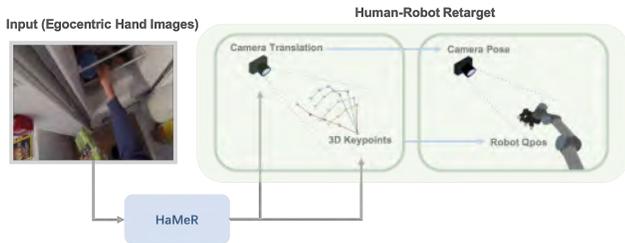


Figure 3. **Human-Robot retarget.** We adopt the HaMeR model to extract hand keypoints and obtain camera parameters, corresponding coordinate systems are constructed, and this information is used to adjust the robot hand’s pose and camera perspective, enabling precise hand pose retargeting.

representation such as data distribution and inclusion of out-of-domain data [10], task and domain adaptation [44], etc. Mirage [52] realizes domain transfer of policy between different robots through a pre-processing process within the image space of training data. Our paradigm differs from Mirage from the perspective that we further investigate a more generalized form for robotic representation learning in a data-centric way which results in a more robust policy training.

3. H2R: Human-to-Robot Data Augmentation

In this section, we describe H2R, a data augmentation pipeline for robot pre-training from videos, our key insight is to remove the hands in every frame and replace them with robotic arms. Figure 2 shows H2R pipeline.

Our proposal is to replace the human hands in every frame with that of a robot, generating an augmented dataset D_{aug} . This approach aims to mitigate the visual gap between human hands and robots, facilitating the transfer of knowledge for easier adaptation of models trained on ego-centric data to robotic tasks. In particular, we hope that the vision encoder trained on the augmented dataset D_{aug} would outperform that trained on the original dataset D in downstream robotic tasks.

3.1. Pipeline of H2R

3D Hand Pose Estimation. In order to overlay the human hands in the image with different robots, we firstly need an efficient and accurate model to detect the hand information. Recent HaMeR [50], a state-of-the-art 3D hand detection and reconstruction model. we detect the position of the hand and its key points, and the internal and external parameters of the rendering camera of the RGB image. Such position information of the identified hand is then used to remove the hands from the image.

Human Arm and Hand Remove. We leverage the Segment Anything Model (SAM) [36] to automatically separate the human hand from the background using hand pose infor-

mation detected by HaMeR. Even though there is only hand information provided but no arm information Thus, SAM could detect both hand and arm as a single object and separate it from background, showing good robustness under the varying conditions of clothing and occlusion. Finally, a state-of-the-art inpainting model LaMa [58] is used to fill the removed hand mask. After this step, we obtain the RGB images with the human hand removed for the later stage of adding robot hands.

Robotic Arm and End Effector Construction. The final step involves constructing the robot arm and end effector, then overlaying it onto the generated images from the previous stage (as shown in Figure 3). For the robotic arm reconstruction, Since HaMeR does not provide information about the arm keypoints, we initially set the target robot to a neutral pose and then adjust the missing joint point information. For the robotic end effector reconstruction For the dexterous hand, the angles of each joint are determined by the angles formed by the corresponding three keypoints, while for the gripper, the degree of opening and closing is determined by the distance between the corresponding fingers.

Simulator Camera Position Alignment. The visual bias introduced by the camera perspective is significantly larger than the action retargeting itself; thus, we leverage the hand keypoints and camera parameters from HaMeR to adjust the camera pose in the simulator. Specifically, the two coordinate systems C_H and C_S can be uniquely determined by the human hand and the robot arm, and the camera’s position in C_H can be used in C_S to ensure the same perspective of the camera. We build the coordinate system ${}^W\mathbf{I}_H$ based on the hand keypoints:

$${}^W\mathbf{I}_H = \{{}^w\mathbf{i}_{H,x}, {}^w\mathbf{i}_{H,y}, {}^w\mathbf{i}_{H,z}\} \quad (1)$$

Where ${}^w\mathbf{i}_{H,x}$, ${}^w\mathbf{i}_{H,y}$, ${}^w\mathbf{i}_{H,z}$ are unit vectors along the x-axes, y-axes and z-axes of the human hand coordinate system. With the keypoints get in HaMeR, we build the three axis of coordinates with the following functions:

$$\begin{aligned} {}^w\mathbf{i}_{H,x} &= {}^w\mathbf{i}_{0,9} \\ {}^w\mathbf{i}_{H,y} &= {}^w\mathbf{i}_{0,9} \times {}^w\mathbf{i}_{0,13} \\ {}^w\mathbf{i}_{H,z} &= {}^w\mathbf{i}_{H,x} \times {}^w\mathbf{i}_{H,y} \end{aligned} \quad (2)$$

Where ${}^w\mathbf{i}_{0,9}$, ${}^w\mathbf{i}_{0,13}$ are unit vectors along middle finger and ring finger. Similarly, To construct the mapping from hand pose to robot arms, we need to get another coordinate system ${}^W\mathbf{I}_S$ in the simulator:

$${}^W\mathbf{I}_S = \{{}^w\mathbf{i}_{S,x}, {}^w\mathbf{i}_{S,y}, {}^w\mathbf{i}_{S,z}\} \quad (3)$$

The method of determining the axis of coordinates is the same:

$$\begin{aligned} {}^w\mathbf{i}_{S,x} &= {}^w\mathbf{i}_{0,2} \\ {}^w\mathbf{i}_{S,y} &= {}^w\mathbf{i}_{0,2} \times {}^w\mathbf{i}_{0,3} \\ {}^w\mathbf{i}_{S,z} &= {}^w\mathbf{i}_{S,x} \times {}^w\mathbf{i}_{S,y} \end{aligned} \quad (4)$$

Where $\mathbf{i}_{0,2}, \mathbf{i}_{0,3}$ are unit vectors along robot fingers that correspond to human middle and ring fingers. We build the following two coordinate transformation matrix to construct the mapping:

$$\begin{aligned} {}^W_H \mathbf{R} &= \begin{pmatrix} {}^W \mathbf{I}_H & \mathbf{key}_0 \\ \mathbf{O} & 1 \end{pmatrix} \\ {}^W_S \mathbf{R} &= \begin{pmatrix} {}^W \mathbf{I}_S & \mathbf{ee_pos} \\ \mathbf{O} & 1 \end{pmatrix} \end{aligned} \quad (5)$$

Where $\mathbf{key}_0, \mathbf{ee_pos}$ are the positions of human wrist and robot wrist. After obtaining the two coordinate systems, we need to determine the position of the camera in the simulator (${}^W \mathbf{cam}_{sim}$) and the position of the camera in the real world (${}^H \mathbf{cam}_{Real}$), thus we can ensure we get the same pose of the human hand and robot arms

$$\begin{aligned} {}^H \mathbf{cam}_{Real} &= {}^W_H \mathbf{R}^{-1} \times {}^W \mathbf{cam}_{Real} \\ {}^S \mathbf{cam}_{sim} &= {}^H \mathbf{cam}_{Real} \\ {}^W \mathbf{cam}_{sim} &= {}^W_S \mathbf{R} \times {}^H \mathbf{R}^{-1} \times {}^W \mathbf{cam}_{Real} \end{aligned} \quad (6)$$

Robot Hand Rendering and Copy-paste. After setting the action, the segmentation mask of the robot arm is obtained by shooting with the camera. The result of the HaMeR model contains the pixel coordinates of the key points of the human hand. By calculating the pixel coordinates of the corresponding links in the robot hand, the robot hand can be copy-pasted to the original image based on the corresponding relationship, ensuring that it is pixel-level aligned with the original human hand in the image (see Figure 4).

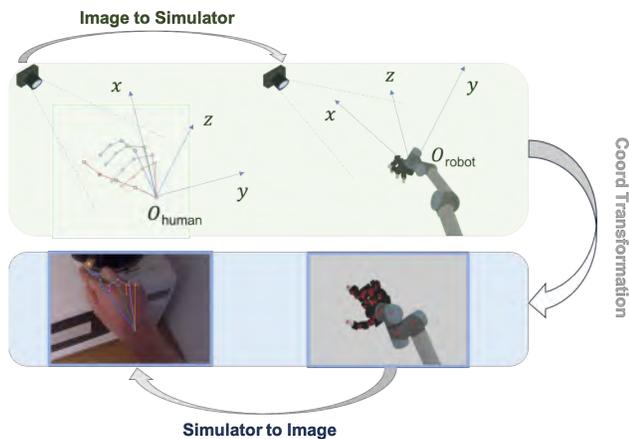


Figure 4. **Robot hand rendering and copy-paste.** The HaMeR model provides hand keypoints and camera parameters, which are used to align the simulator’s camera pose with the original view. The robot arm is then rendered in the simulator, and by matching the pixel coordinates of the arm’s links, it is overlaid onto the original image with pixel-level alignment to the human hand.

3.2. Model Training

Encoder Pre-training. We adopt the MAE [28, 63], R3M [46], and MPI [68] frameworks for pre-training, each employing a Vision Transformer (ViT) Base [14] model as the visual encoder. The SSv2 dataset [21] is used for MAE and R3M training, whereas the Ego4D dataset [22] is employed for MPI training. For the MAE and R3M pre-training methods, in addition to pre-training on the H2R data and raw data, we also applied a simple CutMix baseline to demonstrate the effectiveness of using the robotic arm to cover the human hand, which overlays a fixed set of specific images of robotic arms with grippers onto the original images, ensuring that the overlaid images cover the human hands as much as possible, without exceeding the detected bounding box. Our H2R is different from such baseline by employing robot hand construction to better match the pose of the hand and arm in the images. Based on the type of robotic arm used in CutMix, we categorize the augmented set into three types: CutMix1 represents the UR5 robotic arm, CutMix2 refers to the Franka robotic arm, and CutMix3 combines both the UR5 and Franka robotic arms.

Policy Training. Finally, we employ several existing policy training methods to fine-tune the pre-trained model for evaluations on downstream robotic tasks. We reuse their original implementations to ensure that any performance improvements are solely attributable to our data augmentation approach. For RL models, we evaluate downstream tasks using the PixMC [63] benchmark and employ PPO [56] for policy learning. Additionally, we utilize Robomimic [45] and Diffusion Policy [6] for evaluating IL models. The Robomimic baseline is primarily used for BC policies, and we test three tasks with the Robomimic datasets. For the Diffusion Policy, we specifically evaluate the push task to assess the robustness of our method across different approaches.

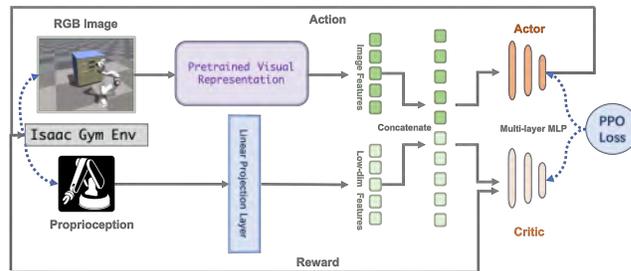


Figure 5. **RL-driven policy training pipeline.** We propose a training pipeline for RL-driven policy learning, designed to evaluate performance across various simulation benchmarks.

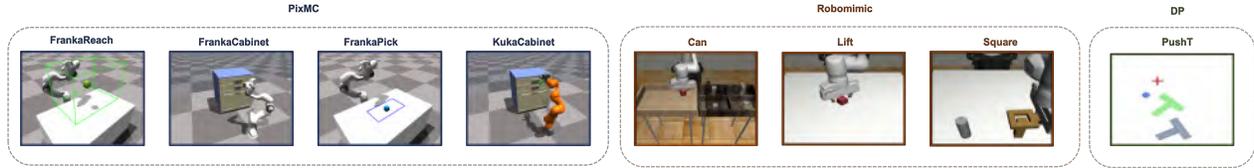


Figure 6. **Simulation benchmark.** We choose 4 tasks from the PixMC and 3 tasks from the Robomimic, covering a range of robotic manipulation skills. We also include the PushT task, designed for the Diffusion Policy framework, as an additional benchmark to evaluate performance in a different task setup.

347 4. Simulation Experiment

348 4.1. Experiment Setup

349 **Simulation Benchmark.** For evaluations in simulation, we
 350 select a total of 8 simulation benchmarks across different
 351 environments, which are PixMC [63], Robomimic [45], and
 352 Diffusion Policy [6]. In particular, for PixMC, we select
 353 **FrankaReach**, **FrankaCabinet**, **FrankaPick**, and **Kuka-**
 354 **Cabinet** to assess the robot’s ability to interact with ob-
 355 jects. For Robomimic, we include tasks such as **Move-**
 356 **Can**, **Square**, and **Lift**, where the robot performs actions
 357 like moving or lifting objects. We also use the **PushT** task,
 358 designed for the Diffusion Policy framework, which evalu-
 359 ates a robot’s ability to push an object to a target location.
 360 These simulation tasks, visualized in Figure 6, span a range
 361 of manipulation skills, providing a comprehensive evalua-
 362 tion of robot performance. For each pre-training method
 363 (MAE [28, 63], R3M [46], MPI [68]), we evaluate the per-
 364 formance of pre-trained encoders with H2R in reinforc-
 365 ement learning and imitation learning. For the PixMC, most
 366 tasks involve motion control of robotic arms, we primar-
 367 ily use reinforcement learning (RL) methods to validate the
 368 effectiveness of H2R. However, for more complex simula-
 369 tion tasks, such as Robomimic benchmark, experimental re-
 370 sults tend to be more sensitive to the reward mechanisms in
 371 reinforcement learning. Therefore, to avoid the impact of
 372 reward on task success rates, we adopt imitation learning
 373 methods (IL) combined with a pre-trained visual encoder
 374 for testing. This approach will evaluate the effectiveness of
 375 H2R in bridging the gap between human hand and robotic
 376 arm visual perception.

377 **Pre-training Dataset.** We select SSV2 (Something-
 378 Something V2) [21] and Ego4D [22] as the primary datasets
 379 for our experiments. SSV2 contains 220,847 video clips
 380 of human actions with everyday objects, designed to help
 381 models understand fine-grained hand gestures. Ego4D is a
 382 large-scale egocentric dataset with 3,670 hours of video col-
 383 lected from 923 participants worldwide, aimed at advancing
 384 first-person visual perception. We use the SSV2 dataset in
 385 the MAE and R3M methods, and the Ego4D dataset in the
 386 MPI method. For MAE and R3M, we select 62,500 videos
 387 from the SSV2 dataset and randomly sample 16 keyframes
 388 from each video. For MPI, we extract 117,624 action clips

389 from 2,486 videos in the Ego4D dataset, each clip consist-
 390 ing of three frames (start, middle, and end).

391 **Evaluation Metric.** We repeat each experiment three times
 392 with different seeds and report the averaged results. For
 393 tasks in PixMC, we train the models using reinforcement
 394 learning for 2,000 steps and report the final success rate.
 395 For tasks in Robomimic, we train for 200 steps and report-
 396 ing the mean success rate. For tasks in RLbench, rlbench
 397 For the PushT task, we train the Diffusion Policy model for
 398 200 epoches and report the success rate in the simulation en-
 399 vironment. The training hyperparameters used in this work
 400 are identical to those described in the original paper.

401 4.2. Results.

402 **Reinforcement Learning.** From Table 1, we observe that
 403 the improvement brought by H2R in reinforcement learning
 404 shows more variation depending on the task. Some tasks see
 405 an improvement, while others experience a decline. How-
 406 ever, on average, across all tasks, there is still an overall
 407 improvement. Additionally, the performance with CutMix
 408 data is particularly better with R3M, while the use of H2R
 409 data yields excellent results with MAE. For example, when
 410 using the MAE pretraining method, the use of our H2R
 411 data results in a 29.7% improvement in the average suc-
 412 cess rate of the tasks. On the other hand, encoders trained
 413 with CutMix data show improvements ranging from 18.0%
 414 to 21.4%. When using the R3M pretraining method, the im-
 415 provement in average success rate with H2R data is smaller,
 416 but the performance boost with CutMix data is more pro-
 417 nounced. Finally, when using the MPI pretraining method,
 418 the use of H2R data results in a modest reduction in the av-
 419 erage success rate.

420 **Imitation Learning.** From Table 2, we observe that the en-
 421 coder trained on H2R data shows consistent improvements
 422 across various tasks compared to the encoder trained on the
 423 original data, with the average success rate improvement on
 424 all tasks ranging from 0.9% to 10.2%. Especially for the
 425 more challenging MoveCan task, it can improve the suc-
 426 cess rate by 25.5%. Additionally, while encoders trained
 427 on the relatively simple CutMix data show improvement on
 428 tasks in Robomimic, their performance in the PushT task
 429 remains slightly worse than the encoders trained on original
 430 data. These results demonstrate the effectiveness of using

	FrankaReach	FrankaCabinet	FrankaPick	KukaCabinet	Average
MAE	97.5	88	0	90	46.9
MAE+CutMix1	93.5 (-4.0%)	90.5 (+2.5%)	0.0 (0.0%)	80.5 (-9.5%)	66.1 (+19.2%)
MAE+CutMix2	96.5 (-1.0%)	100.0 (+12.0%)	0.0 (0.0%)	63.0 (-27.0%)	64.9 (+18.0%)
MAE+CutMix3	98.5 (+1.0%)	90.5 (+2.5%)	0.0 (0.0%)	84.0 (-6.0%)	68.3 (+21.4%)
MAE+H2R	96.0 (-1.5%)	92.0 (+4.0%)	31.5 (+31.5%)	87.0 (-3.0%)	76.6 (+29.7%)
R3M	63	0	0	0	15.8
R3M+CutMix1	95.5 (+32.5%)	99.5 (+99.5%)	0.0 (0.0%)	1.0 (+1.0%)	49.0 (+33.2%)
R3M+CutMix2	98.5 (+35.5%)	97.5 (+97.5%)	0.0 (0.0%)	0.0 (0.0%)	49.0 (+33.2%)
R3M+CutMix3	97.5 (+34.5%)	85.5 (+85.5%)	0.0 (0.0%)	0.0 (0.0%)	45.8 (+30.0%)
R3M+H2R	8.5 (-54.5%)	0.0 (0.0%)	0.0 (0.0%)	81.0 (+81.0%)	17.9 (+2.1%)
MPI	83.5	0	0	58	35.4
MPI+H2R	88.0 (+4.5%)	20.0 (+20.0%)	0.0 (0.0%)	0.0 (-58.0%)	27 (-8.4%)

Table 1. **Reinforcement learning experiment result.** We report the success rate (%) over RL-based tasks for MAE, R3M, and MPI.

	MoveCan	Square	Lift	Average	PushT
MAE	54	25.5	94.5	58	59.2
MAE+CutMix1	72.0 (+18.0%)	30.0 (+4.5%)	95.0 (+0.5%)	65.7 (+7.7%)	37.5 (-21.7%)
MAE+CutMix2	58.0 (+4.0%)	36.0 (+10.5%)	90.0 (-4.5%)	61.3 (+3.3%)	40.0 (-19.2%)
MAE+CutMix3	78.0 (+24.0%)	32.0 (+9.3%)	92.0 (-2.5%)	67.3 (+2.7%)	42.0 (-17.2%)
MAE+H2R	79.5 (+25.5%)	29.5 (+4.0%)	95.5 (+1.0%)	68.2 (+10.2%)	64.5 (+5.3%)
R3M	59.5	20.5	85	55	15
R3M+CutMix1	69.5 (+10.0%)	30.0 (+9.5%)	91.0 (+6.0%)	63.5 (+8.5%)	19.0 (+4.0%)
R3M+CutMix2	66.0 (+6.5%)	26.0 (+5.5%)	83.0 (-2.0%)	58.3 (+3.3%)	17.0 (+2.0%)
R3M+CutMix3	68.0 (+8.5%)	26.0 (+5.5%)	84.0 (-1.0%)	59.3 (+4.3%)	14.0 (-1.0%)
R3M+H2R	61.5 (+2.0%)	37.5 (+17.0%)	85.0 (0.0%)	61.3 (+6.3%)	22.0 (+7.0%)
MPI	58	21	96	58.3	62.7
MPI+H2R	62.5 (+4.5%)	24.5 (+3.5%)	94.5 (-1.5%)	60.5 (+2.2%)	63.8 (+0.9%)

Table 2. **Imitation learning experiment result.** We report the success rate (%) over IL-based tasks for MAE, R3M, and MPI.

431 the robotic arm to cover the human hand in video data, as
432 well as the effectiveness of H2R in imitation learning.

433 5. Real World Experiment

434 5.1. Experiment Setup

435 **Real-world Tasks.** To validate the effectiveness of H2R in
436 downstream manipulation tasks, we implement three real-
437 world manipulation tasks using a UR5 [61] robotic arm with
438 a Robotiq [54] Gripper integration. The single RealSense
439 L515 [31] camera is used to obtain visual observations in
440 the real world. Realsense L515 is set above and behind the
441 robotic arm, which provides a similar viewpoint to the hu-
442 man video data used in the pretrained visual model. Our
443 real-world setup are shown in Figure 7. We provide detailed
444 descriptions of the three implemented manipulation tasks as
445 follows.

1. **Pick and Place:** Grasp a cube and place it into a bowl. 446
2. **Stack Cubes:** Stack a blue cube atop a yellow cube. 447
3. **Pick from Box:** Retrieve a cube from a box, place it into 448
a bowl, and then close the box lid. 449

All the tasks are visualized in Appendix. Pick and Place 450
task is the simplest of three tasks but still requires precise 451
cube recognition, grasping, and placement within a design- 452
ated bowl area. By contrast, Stack Cubes task demands 453
higher positional accuracy during placement, 454

requiring precise identification of the yellow cube’s loca- 455
tion for successful stacking. Pick from Box task combines 456
grasping with articulated object manipulation, necessitating 457
longer-horizon planning and higher precision. For instance, 458
the robot must avoid the lid while retrieving the cube and se- 459
lect optimal contact points to close the lid post-placement, 460
challenging its ability to learn from high-dimensional visual 461
inputs. 462

Policy Model	Tasks	MAE	MAE+H2R	R3M	R3M+H2R
Diffusion Policy	Pick and Place	45	65(+20%)	40	50(+10%)
	Stack Cubes	50	55(+5%)	55	70(+15%)
	Pick from Box	55	50(-5%)	45	65(+20%)
	Average	50	56.7(+6.7%)	46.7	61.7(+15%)
Equivariant Diffusion Policy	Pick and Place	55	70(+15%)	60	70(+10%)
	Stack Cubes	50	50(0%)	65	75(+10%)
	Pick from Box	55	75(+25%)	50	70(+20%)
	Average	53.3	65(+11.7%)	58.3	75(+13.3%)

Table 3. **Real-World success rate.** We report the success rate (%) over real-world tasks for MAE, R3M. Percentage changes due to H2R are shown in parentheses, with blue indicating improvement and red indicating degradation.

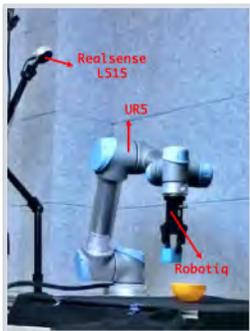


Figure 7. Real-world experiment scene.

Policy Training Process Details. For dataset collection, we collect expert demonstrations through human teleoperation, comprising 30 demonstrations per task. For downstream policy training, we select the Diffusion Policy (DP) [6] and Equivariant Diffusion Policy (eq-dp) [64] as policy frameworks. We apply upstream pre-trained MAE and R3M visual representation models to downstream policy learning, selecting four configurations for comparison: MAE, MAE+H2R, R3M, and R3M+H2R. Pretrained models are incorporated as frozen vision encoders in the policy network to evaluate their effectiveness. We use a single RGB camera image as the high-dimensional observation space and the robot proprioception as the low-dimensional observation space. Both are combined as input observations to the policy network. Policy is trained for 300 epochs using the collected data for each task.

5.2. Experiment Results

We evaluate the success rates of each model-task combination in real-world deployments. The results, as shown in Table 3, demonstrate that H2R significantly enhance the performance of visual encoders across diverse robotic tasks. H2R augmentation improves MAE-based policies in 6 out of 7 task configurations, with the largest gains in Pick from Box (+25% for Equivariant Diffusion) and Pick

and Place (+20% for Diffusion Policy). Across all tasks, H2R augmentation consistently enhanced R3M-based policies, with the most notable improvements observed in geometrically complex scenarios such as Pick from Box, where R3M+H2R paired with Equivariant Diffusion Policy achieved a 20% success rate increase. These results highlight the potential of our approach to enhance visual encoders for real-world robotic applications, even in complex and dynamic environments.

6. Conclusion

We proposed H2R, a data augmentation technique that bridges the visual gap between human hand demonstrations and robotic arm manipulations by replacing human hands in first-person videos with robotic arm movements. Using 3D hand reconstruction and image inpainting models, H2R generates synthetic robotic arm manipulation sequences, making them more suitable for robot pre-training. Experiments across simulation benchmarks and real-world tasks demonstrate consistent improvements in success rates for encoders trained with various pre-training methods (MAE, R3M, MPI), highlighting its effectiveness and generalizability. H2R enables efficient transfer of task knowledge from human demonstrations to robotic systems, reducing reliance on costly robot-specific data collection.

511

References

- 512 [1] Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-
513 son, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri,
514 Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu,
515 Jonathan H. Clark, Laurent El Shafey, Yanping Huang,
516 Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark
517 Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Ke-
518 fan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernan-
519 dez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan
520 Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks,
521 Michele Catasta, Yong Cheng, Colin Cherry, Christopher A.
522 Choquette-Choo, Aakanksha Chowdhery, Clément Crepy,
523 Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin,
524 Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxi-
525 aoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia,
526 Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven
527 Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea
528 Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Itty-
529 cheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy,
530 Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine
531 Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li,
532 Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Fred-
533 erick Liu, Marcello Maggioni, Aroma Mahendru, Joshua
534 Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado,
535 John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie
536 Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan
537 Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Cas-
538 tro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Ren-
539 nee Shelby, Ambrose Slone, Daniel Smilkov, David R. So,
540 Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasude-
541 van, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui
542 Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu,
543 Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao
544 Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny
545 Zhou, Slav Petrov, and Yonghui Wu. Palm 2 technical re-
546 port, 2023. 2
- 547 [2] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen
548 Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding,
549 Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence,
550 Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakr-
551 ishnan, Kehang Han, Karol Hausman, Alexander Herzog,
552 Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan
553 Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal,
554 Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu,
555 Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kan-
556 ishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar,
557 Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait
558 Singh, Radu Soricut, Huang Tran, Vincent Vanhoucke, Quan
559 Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin
560 Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu,
561 and Brianna Zitkovich. Rt-2: Vision-language-action mod-
562 els transfer web knowledge to robotic control, 2023. 2
- 563 [3] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen
564 Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakr-
565 ishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Ju-
566 lian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally
567 Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalash-
nikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey
Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor
Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta,
Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka
Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin
Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clay-
ton Tan, Huang Tran, Vincent Vanhoucke, Steve Vega, Quan
Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu,
and Brianna Zitkovich. Rt-1: Robotics transformer for real-
world control at scale, 2023. 2
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Ge-
offrey Hinton. A simple framework for contrastive learning
of visual representations. *arXiv preprint arXiv:2002.05709*,
2020. 3
- [5] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad
Norouzi, and Geoffrey Hinton. Big self-supervised mod-
els are strong semi-supervised learners. *arXiv preprint
arXiv:2006.10029*, 2020. 3
- [6] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric
Cousineau, Benjamin Burchfiel, and Shuran Song. Diffu-
sion policy: Visuomotor policy learning via action diffusion.
arXiv preprint arXiv:2303.04137, 2023. 2, 3, 5, 6, 8
- [7] Reinis Cimurs, Il Hong Suh, and Jin Han Lee. Goal-driven
autonomous exploration through deep reinforcement learn-
ing. *IEEE Robotics and Automation Letters*, 7(2):730–737,
2022. 3
- [8] Embodiment Collaboration, Abby O’Neill, Abdul Rehman,
Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta,
Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim
Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex
Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky,
Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov,
Anikait Singh, Animesh Garg, Aniruddha Kembhavi, An-
nie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma,
Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan
Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard
Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles
Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu,
Cheng Chi, Chenguang Huang, Christine Chan, Christo-
pher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei
Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak
Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman,
Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan
Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao,
Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gau-
rav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng,
Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen
Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi,
Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hi-
roki Furuta, Homanga Bharadhwaj, Homer Walke, Hongjie
Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel
Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn
Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jay Vakil,
Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao,
Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan
Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jing-
pei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey
Hejna, Jonathan Booyer, Jonathan Tompson, Jonathan Yang,

- 626 Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, 684
627 Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, 685
628 Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, 686
629 Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin 687
630 Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty El- 688
631 lis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Ku- 689
632 nal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, 690
633 Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei- 691
634 Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca 692
635 Weihs, Magnum Chen, Marion Lepert, Marius Memmel, 693
636 Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, 694
637 Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, 695
638 Mingtong Zhang, Mingyu Ding, Minh Ho, Mohan Kumar 696
639 Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, 697
640 Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suen- 698
641 derhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi 699
642 Shafullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pan- 700
643 nag R Sanketi, Patrick "Tree" Miller, Patrick Yin, Paul 701
644 Wohllhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre 702
645 Sermanet, Pieter Abbeel, Priya Sundareshan, Qiuyu Chen, 703
646 Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto 704
647 Mart'in-Mart'in, Rohan Bajjal, Rosario Scalise, Rose Hen- 705
648 drix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Men- 706
649 donca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bus- 707
650 tamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry 708
651 Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, 709
652 Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Hal- 710
653 dar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, 711
654 Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen 712
655 Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel 713
656 Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, 714
657 Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Ma- 715
658 sushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, 716
659 Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor 717
660 Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent 718
661 Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, 719
662 Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, 720
663 Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yan- 721
664 song Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen 722
665 Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yix- 723
666 uan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, 724
667 Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin 725
668 Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang 726
669 Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, 727
670 Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and 728
671 Zipeng Lin. Open x-embodiment: Robotic learning datasets 729
672 and rt-x models, 2024. 2, 3 730
- [9] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, 731
673 Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide 732
674 Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. 733
675 Scaling egocentric vision: The epic-kitchens dataset. In 734
676 *Proceedings of the European conference on computer vision* 735
677 (*ECCV*), pages 720–736, 2018. 2 736
- [10] Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Ab- 737
679 hinav Gupta. An unbiased look at datasets for visuo-motor 738
680 pre-training. In *7th Annual Conference on Robot Learning*, 739
681 2023. 4 740
682
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, 741
and Li Fei-Fei. Imagenet: A large-scale hierarchical image 742
database. In *2009 IEEE conference on computer vision and 743*
683 *pattern recognition*, pages 248–255. Ieee, 2009. 3 744
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina 745
Toutanova. Bert: Pre-training of deep bidirectional 746
transformers for language understanding. *arXiv preprint 747*
684 *arXiv:1810.04805*, 2018. 2 748
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, 749
Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, 750
Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl- 751
vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image 752
is worth 16x16 words: Transformers for image recognition 753
at scale. In *9th International Conference on Learning Rep- 754*
685 *resentations, ICLR 2021, Virtual Event, Austria, May 3-7,* 755
686 *2021. OpenReview.net*, 2021. 3 756
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, 757
Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, 758
Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl- 759
vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image 760
is worth 16x16 words: Transformers for image recognition 761
at scale, 2021. 5 762
- [15] Jiafei Duan, Yi Ru Wang, Mohit Shridhar, Dieter Fox, and 763
Ranjay Krishna. Ar2-d2: Training a robot without a robot. 764
In *Conference on Robot Learning*, pages 2838–2848. PMLR, 765
2023. 2 766
- [16] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, 767
Chenxi Wang, Junbo Wang, Haoyi Zhu, and Cewu Lu. 768
Rh20t: A comprehensive robotic dataset for learning diverse 769
skills in one-shot, 2023. 2 770
- [17] Pete Florence, Corey Lynch, Andy Zeng, Oscar Ramirez, 771
Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, 772
Igor Mordatch, and Jonathan Tompson. Implicit behavioral 773
cloning, 2021. 3 774
- [18] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan 775
Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, 776
Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. 777
Dat- 778
acomp: In search of the next generation of multimodal 779
datasets. *Advances in Neural Information Processing Sys-* 780
781 *tems*, 36, 2024. 2 782
- [19] Jensen Gao, Annie Xie, Ted Xiao, Chelsea Finn, and Dorsa 783
Sadigh. Efficient data collection for robotic manipulation via 784
compositional generalization, 2024. 2 785
- [20] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michal- 786
ski, Joanna Materzynska, Susanne Westphal, Heuna Kim, 787
Valentin Haenel, Ingo Freund, Peter Yianilos, Moritz 788
Mueller-Freitag, et al. The "something something" video 789
database for learning and evaluating visual common sense. 790
In *Proceedings of the IEEE international conference on com-* 791
683 *puter vision*, pages 5842–5850, 2017. 2 792
- [21] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michal- 793
ski, Joanna Materzyńska, Susanne Westphal, Heuna Kim, 794
Valentin Haenel, Ingo Freund, Peter Yianilos, Moritz 795
Mueller-Freitag, Florian Hoppe, Christian Thureau, Ingo Bax, 796
and Roland Memisevic. The "something something" video 797
database for learning and evaluating visual common sense, 798
2017. 5, 6 799
- [22] Kristen Grauman, Andrew Westbury, Eugene Byrne, 800
Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson 801

- 742 Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Mar- 800
743 tin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar 801
744 Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, 802
745 Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant 803
746 Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien 804
747 Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichten- 805
748 hofer, Adriano Fragomeni, Qichen Fu, Abrahm Gebrese- 806
749 lasie, Cristina Gonzalez, James Hillis, Xuhua Huang, Yifei 807
750 Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kot- 808
751 tur, Anurag Kumar, Federico Landini, Chao Li, Yanghao 809
752 Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, 810
753 Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will 811
754 Price, Paola Ruiz Puentes, Merey Ramazanova, Leda Sari, 812
755 Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, 813
756 Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma 814
757 Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Cran- 815
758 dall, Dima Damen, Giovanni Maria Farinella, Christian Fue- 816
759 gen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, 817
760 Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, 818
761 Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, 819
762 Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo 820
763 Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around 821
764 the world in 3,000 hours of egocentric video, 2022. 5, 6 822
[23] Kristen Grauman, Andrew Westbury, Eugene Byrne, 823
765 Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson 824
766 Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: 825
767 Around the world in 3,000 hours of egocentric video. In *Pro- 826*
768 *ceedings of the IEEE/CVF Conference on Computer Vision 827*
769 *and Pattern Recognition*, pages 18995–19012, 2022. 2 828
[24] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and 829
770 Sergey Levine. Soft actor-critic: Off-policy maximum 830
771 entropy deep reinforcement learning with a stochastic 831
772 actor, 2018. cite arxiv:1801.01290Comment: ICML 832
773 2018 Videos: sites.google.com/view/soft-actor-critic Code: 833
774 github.com/haarnoja/sac. 3 834
[25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 835
775 Deep residual learning for image recognition. In *Proceed- 836*
776 *ings of the IEEE Conference on Computer Vision and Pattern 837*
777 *Recognition (CVPR)*, 2016. 3 838
[26] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross 839
778 Girshick. Momentum contrast for unsupervised visual repre- 840
779 sentation learning. In *Proceedings of the IEEE/CVF Confer- 841*
780 *ence on Computer Vision and Pattern Recognition (CVPR)*, 842
781 2020. 3 843
[27] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr 844
782 Dollár, and Ross Girshick. Masked autoencoders are scalable 845
783 vision learners. In *Proceedings of the IEEE/CVF Conference 846*
784 *on Computer Vision and Pattern Recognition (CVPR)*, pages 847
785 16000–16009, 2022. 3 848
[28] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr 849
786 Dollár, and Ross Girshick. Masked autoencoders are scalable 850
787 vision learners. In *Proceedings of the IEEE/CVF conference 851*
788 *on computer vision and pattern recognition*, pages 16000– 852
789 16009, 2022. 2, 5, 6 853
[29] Alexander Herzog, Kanishka Rao, Karol Hausman, Yao 854
790 Lu, Paul Wohlhart, Mengyuan Yan, Jessica Lin, Montser- 855
791 rat Gonzalez Arenas, Ted Xiao, Daniel Kappler, Daniel 856
792 Ho, Jarek Rettinghouse, Yevgen Chebotar, Kuang-Huei 857
793 Lee, Keerthana Gopalakrishnan, Ryan Julian, Adrian Li, 858
794 Chuyuan Kelly Fu, Bob Wei, Sangeetha Ramesh, Khem 859
795 Holden, Kim Kleiven, David Rendleman, Sean Kirmani, Jeff 860
796 Bingham, Jon Weisz, Ying Xu, Wenlong Lu, Matthew Ben- 861
797 nice, Cody Fong, David Do, Jessica Lam, Yunfei Bai, Benjie 862
798 Holson, Michael Quinlan, Noah Brown, Mrinal Kalakrish- 863
799 nan, Julian Ibarz, Peter Pastor, and Sergey Levine. Deep rl at 864
scale: Sorting waste in office buildings with a fleet of mobile 865
manipulators, 2023. 2 866
[30] Kevin Huang, Rwik Rana, Alexander Spitzer, Guanya Shi, 867
and Byron Boots. DATT: Deep adaptive trajectory tracking 868
for quadrotor control. In *7th Annual Conference on Robot 869*
Learning, 2023. 3 870
[31] Intel Corporation. Intel@realsense™lidar camera 1515. 871
[https://www.intelrealsense.com/lidar- 872](https://www.intelrealsense.com/lidar-camera-1515/)
[camera-1515/](https://www.intelrealsense.com/lidar-camera-1515/), 2025. Accessed: 2025-02-01. 7 873
[32] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Freder- 874
ik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. 875
Bc-z: Zero-shot task generalization with robotic imitation 876
learning, 2022. 2 877
[33] Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Ben- 878
jamin Swanson, Rico Jonschkowski, Chelsea Finn, Sergey 879
Levine, and Karol Hausman. Mt-opt: Continuous multi-task 880
robotic reinforcement learning at scale, 2021. 2 881
[34] Siddharth Karamcheti, Suraj Nair, Annie S. Chen, Thomas 882
Kollar, Chelsea Finn, Dorsa Sadigh, and Percy Liang. 883
Language-driven representation learning for robotics. In 884
Robotics: Science and Systems (RSS), 2023. 2, 3 885
[35] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Bal- 886
akrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush 887
Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang 888
Chen, Kirsty Ellis, Peter David Fagan, Joey Hejna, Masha 889
Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree 890
Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, 891
Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Mem- 892
mel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Al- 893
bert Zhan, Kevin Black, Cheng Chi, Kyle Beltran Hatch, 894
Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pan- 895
nag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, 896
Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Hee- 897
won Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, 898
Rohan Baijal, Mateo Guaman Castro, Daphne Chen, Qi- 899
yuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul Fos- 900
ter, Jensen Gao, David Antonio Herrera, Minho Heo, Kyle 901
Hsu, Jiaheng Hu, Donovan Jackson, Charlotte Le, Yun- 902
shuang Li, Kevin Lin, Roy Lin, Zehan Ma, Abhiram Mad- 903
dukuri, Suvir Mirchandani, Daniel Morton, Tony Nguyen, 904
Abigail O’Neill, Rosario Scalise, Derick Seale, Victor Son, 905
Stephen Tian, Emi Tran, Andrew E. Wang, Yilin Wu, An- 906
nie Xie, Jingyun Yang, Patrick Yin, Yunchu Zhang, Os- 907
bert Bastani, Glen Berseth, Jeannette Bohg, Ken Gold- 908
berg, Abhinav Gupta, Abhishek Gupta, Dinesh Jayaraman, 909
Joseph J Lim, Jitendra Malik, Roberto Martín-Martín, Sub- 910
ramanian Ramamoorthy, Dorsa Sadigh, Shuran Song, Jia- 911
jun Wu, Michael C. Yip, Yuke Zhu, Thomas Kollar, Sergey 912
Levine, and Chelsea Finn. Droid: A large-scale in-the-wild 913
robot manipulation dataset, 2024. 2 914
[36] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, 915

- 858 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-
859 head, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and
860 Ross Girshick. Segment anything, 2023. 2, 4
- [37] Alex X. Lee, Coline Devin, Yuxiang Zhou, Thomas Lampe,
861 Konstantinos Bousmalis, Jost Tobias Springenberg, Arunk-
862 umar Byravan, Abbas Abdolmaleki, Nimrod Gileadi, David
863 Khosid, Claudio Fantacci, Jose Enrique Chen, Akhil Raju,
864 Rae Jeong, Michael Neunert, Antoine Laurens, Stefano Sal-
865 iceti, Federico Casarini, Martin Riedmiller, Raia Hadsell,
866 and Francesco Nori. Beyond pick-and-place: Tackling
867 robotic stacking of diverse shapes, 2021. 2
- [38] Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen
870 Koltun, and Marco Hutter. Learning quadrupedal locomotion
871 over challenging terrain. *Science Robotics*, 5(47):eabc5986,
872 2020. 3
- [39] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter
873 Abbeel. End-to-end training of deep visuomotor policies,
874 2016. 3
- [40] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz,
875 and Deirdre Quillen. Learning hand-eye coordination for
876 robotic grasping with deep learning and large-scale data col-
877 lection. *The International Journal of Robotics Research*, 37
878 (4-5):421–436, 2018. 3
- [41] Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt
881 Jordan, Samir Gadre, Hritik Bansal, Etash Guha, Sedrick
882 Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muen-
883 nighoff, Reinhard Heckel, Jean Mercat, Mayee Chen, Suchin
884 Gururangan, Mitchell Wortsman, Alon Albalak, Yonatan
885 Bitton, Marianna Nezhurina, Amro Abbas, Cheng-Yu Hsieh,
886 Dhruva Ghosh, Josh Gardner, Maciej Kilian, Hanlin Zhang,
887 Rulin Shao, Sarah Pratt, Sunny Sanyal, Gabriel Ilharco,
888 Giannis Daras, Kalyani Marathe, Aaron Gokaslan, Jieyu
889 Zhang, Khyathi Chandu, Thao Nguyen, Igor Vasiljevic,
890 Sham Kakade, Shuran Song, Sujay Sanghavi, Fartash
891 Faghri, Sewoong Oh, Luke Zettlemoyer, Kyle Lo, Alaaeldin
892 El-Nouby, Hadi Pouransari, Alexander Toshev, Stephanie
893 Wang, Dirk Groeneveld, Luca Soldaini, Pang Wei Koh, Je-
894 nia Jitsev, Thomas Kollar, Alexandros G. Dimakis, Yair Car-
895 mon, Achal Dave, Ludwig Schmidt, and Vaishaal Shankar.
896 Datacomp-lm: In search of the next generation of training
897 sets for language models. *arXiv preprint arXiv:2406.11794*,
898 2024. 2
- [42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays,
900 Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence
901 Zitnick. Microsoft coco: Common objects in context. In
902 *Computer Vision – ECCV 2014*, pages 740–755, Cham,
903 2014. Springer International Publishing. 3
- [43] Jianlan Luo, Charles Xu, Jeffrey Wu, and Sergey Levine.
904 Precise and dexterous robotic manipulation via human-in-
905 the-loop reinforcement learning, 2024. 3
- [44] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason
906 Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre
907 Berges, Tingfan Wu, Jay Vakil, Pieter Abbeel, Jitendra Ma-
908 lik, Dhruv Batra, Yixin Lin, Oleksandr Maksymets, Aravind
909 Rajeswaran, and Franziska Meier. Where are we in the
910 search for an artificial visual cortex for embodied intelli-
911 gence? In *Advances in Neural Information Processing Sys-*
912 *tems*, pages 655–677. Curran Associates, Inc., 2023. 4
- [45] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiri-
913 any, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio
914 Savarese, Yuke Zhu, and Roberto Martín-Martín. What mat-
915 ters in learning from offline human demonstrations for robot
916 manipulation, 2021. 2, 3, 5, 6
- [46] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea
921 Finn, and Abhinav Gupta. R3m: A universal visual
922 representation for robot manipulation. *arXiv preprint*
923 *arXiv:2203.12601*, 2022. 2, 3, 5, 6
- [47] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,
924 Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo
925 Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anad-
926 kat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Bal-
927 com, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
928 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-
929 Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko,
930 Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman,
931 Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai,
932 Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
933 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fo-
934 tis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Ja-
935 son Chen, Mark Chen, Ben Chess, Chester Cho, Casey
936 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Cur-
937 rier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah
938 Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve
939 Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,
940 Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
941 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo
942 Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogi-
943 neni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon,
944 Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross,
945 Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han,
946 Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke,
947 Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,
948 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost
949 Huizinga, Shantanu Jain, Shawn Jain, Joanne Jiang, Angela
950 Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto,
951 Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser,
952 Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,
953 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina
954 Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros,
955 Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, An-
956 drew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen
957 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee,
958 Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel
959 Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa
960 Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Mal-
961 facini, Sam Manning, Todor Markov, Yaniv Markovski,
962 Bianca Martin, Katie Mayer, Andrew Mayne, Bob Mc-
963 Grew, Scott Mayer McKinney, Christine McLeavey, Paul
964 McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob
965 Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin,
966 Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu,
967 Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Rei-
968 ichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard
969 Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub
970 Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Gi-
971 ambattista Parascandolo, Joel Parish, Emy Parparita, Alex

- 974 Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Fil- 1031
975 ipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde 1032
976 de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, 1033
977 Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, 1034
978 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya 1035
979 Ramesh, Cameron Raymond, Francis Real, Kendra Rim- 1036
980 bach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, 1037
981 Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sas- 1038
982 try, Heather Schmidt, David Schnurr, John Schulman, Daniel 1039
983 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, 1040
984 Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, 1041
985 Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, 1042
986 Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Fe- 1043
987 lipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie 1044
988 Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, 1045
989 Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick 1046
990 Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Val- 1047
991 lone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, 1048
992 Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, 1049
993 Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welin- 1050
994 der, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, 1051
995 Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren 1052
996 Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, 1053
997 Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech 1054
998 Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, 1055
999 Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William 1056
1000 Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. 2 1057
- [48] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. 1058
1001 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, 1059
1002 Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, 1060
1003 Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell 1061
1004 Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael 1062
1005 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Je- 1063
1006 gou, Julien Mairal, Patrick Labatut, Armand Joulin, and Pi- 1064
1007 otr Bojanowski. DINOv2: Learning robust visual features 1065
1008 without supervision. *Transactions on Machine Learning Re- 1066*
1009 *search*, 2024. 3 1067
- [49] Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, 1068
1011 and Abhinav Gupta. The unsurprising effectiveness of pre- 1069
1012 trained vision models for control. In *Proceedings of the 1070*
1013 *39th International Conference on Machine Learning*, pages 1071
1014 17359–17371. PMLR, 2022. 3 1072
- [50] Georgios Pavlakos, Dandan Shan, Ilija Radosavovic, Angjoo 1073
1016 Kanazawa, David Fouhey, and Jitendra Malik. Reconstruct- 1074
1017 ing hands in 3d with transformers, 2023. 2, 4 1075
- [51] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, 1076
1019 Jitendra Malik, and Trevor Darrell. Real-world robot learn- 1077
1020 ing with masked visual pre-training, 2022. 2, 3 1078
- [52] Ilija Radosavovic, Baifeng Shi, Letian Fu, Ken Goldberg, 1079
1022 Trevor Darrell, and Jitendra Malik. Robot learning with sen- 1080
1023 sorimotor pre-training. In *7th Annual Conference on Robot 1081*
1024 *Learning*, 2023. 4 1082
- [53] Ilija Radosavovic, Baifeng Shi, Letian Fu, Ken Goldberg, 1083
1026 Trevor Darrell, and Jitendra Malik. Robot learning with sen- 1084
1027 sorimotor pre-training. In *Conference on Robot Learning*, 1085
1028 pages 683–693. PMLR, 2023. 3 1086
- [54] Robotiq Inc. Adaptive grippers. [https://robotiq.](https://robotiq.com/products/adaptive-grippers) 1087
1030 [com/products/adaptive-grippers](https://robotiq.com/products/adaptive-grippers), 2025. Ac- 1088
1031 cessed: 2025-02-01. 7 1032
- [55] Christoph Schuhmann, Richard Vencu, Romain Beaumont, 1033
1034 Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo 1035
1036 Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: 1037
1038 Open dataset of clip-filtered 400 million image-text pairs. 1039
1039 *ArXiv*, abs/2111.02114, 2021. 3 1040
- [56] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Rad- 1041
1042 ford, and Oleg Klimov. Proximal policy optimization algo- 1043
1044 rithms. *CoRR*, abs/1707.06347, 2017. 2, 3, 5 1045
- [57] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang 1046
1047 Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: 1048
1049 A large-scale, high-quality dataset for object detection. In 1049
1050 *2019 IEEE/CVF International Conference on Computer Vi- 1051*
1051 *sion (ICCV)*, pages 8429–8438, 2019. 3 1052
- [58] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, 1053
1054 Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, 1055
1056 Naejin Kong, Harshith Goka, Kiwoong Park, and Victor 1057
1058 Lempitsky. Resolution-robust large mask inpainting with 1059
1059 fourier convolutions, 2021. 2, 4 1060
- [59] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Woj- 1061
1062 ciech Zaremba, and Pieter Abbeel. Domain randomization 1063
1064 for transferring deep neural networks from simulation to the 1064
1065 real world, 2017. 3 1066
- [60] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral 1067
1068 cloning from observation. *arXiv preprint arXiv:1805.01954*, 1068
1069 2018. 3 1069
- [61] Universal Robots A/S. Ur5 - collaborative robots. <https://www.universal-robots.com/products/ur5-robot/>, 2025. Accessed: 2025-02-01. 7 1070
- [62] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, 1071
1072 Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Im- 1073
1074 proving dexterous grasping policy learning via geometry- 1074
1075 aware curriculum and iterative generalist-specialist learning. 1075
1076 *arXiv preprint arXiv:2304.00464*, 2023. 3 1077
- [63] Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra 1078
1079 Malik. Masked visual pre-training for motor control. *arXiv 1079*
1080 *preprint arXiv:2203.06173*, 2022. 2, 3, 5, 6 1080
- [64] Jingyun Yang, Zi ang Cao, Congyue Deng, Rika Antonova, 1081
1082 Shuran Song, and Jeannette Bohg. Equibot: Sim(3)- 1082
1083 equivariant diffusion policy for generalizable and data effi- 1083
1084 cient learning, 2024. 2, 3, 8 1084
- [65] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie 1085
1086 Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud 1085
1087 transformers with masked point modeling. In *Proceedings 1086*
1088 *of the IEEE/CVF conference on computer vision and pattern 1087*
1089 *recognition*, pages 19313–19322, 2022. 2 1088
- [66] Zhecheng Yuan, Tianming Wei, Shuiqi Cheng, Gu Zhang, 1089
1090 Yuanpei Chen, and Huazhe Xu. Learning to manipulate any- 1089
1091 where: A visual generalizable framework for reinforcement 1090
1092 learning. *arXiv preprint arXiv:2407.15815*, 2024. 3 1091
- [67] Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, 1092
1093 Muhan Wang, and Huazhe Xu. 3d diffusion policy: Gen- 1092
1094 eralizable visuomotor policy learning via simple 3d repre- 1093
1095 sentations. In *Proceedings of Robotics: Science and Systems 1094*
1096 *(RSS)*, 2024. 3 1095
- [68] Jia Zeng, Qingwen Bu, Bangjun Wang, Wenke Xia, Li Chen, 1096
1097 Hao Dong, Haoming Song, Dong Wang, Di Hu, Ping Luo, 1097
1098 1098

- 1089 Heming Cui, Bin Zhao, Xuelong Li, Yu Qiao, and Hongyang
1090 Li. Learning manipulation by predicting interaction, 2024. 2,
1091 5, 6
- 1092 [69] Jia Zeng, Qingwen Bu, Bangjun Wang, Wenke Xia, Li Chen,
1093 Hao Dong, Haoming Song, Dong Wang, Di Hu, Ping Luo,
1094 et al. Learning manipulation by predicting interaction. *arXiv*
1095 *preprint arXiv:2406.00439*, 2024. 2
- 1096 [70] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea
1097 Finn. Learning Fine-Grained Bimanual Manipulation with
1098 Low-Cost Hardware. In *Proceedings of Robotics: Science*
1099 *and Systems*, Daegu, Republic of Korea, 2023. 3
- 1100 [71] Ziwen Zhuang, Zipeng Fu, Jianren Wang, Christopher Atke-
1101 son, Sören Schwertfeger, Chelsea Finn, and Hang Zhao.
1102 Robot parkour learning. In *Conference on Robot Learning*
1103 *(CoRL)*, 2023. 3

1104 **Appendix**1105 **CLIP-based Evaluation of Augmentation Effectiveness.**

1106 To quantitatively assess the effectiveness of H2R in bridging
1107 the visual gap between human and robot hands, we employ
1108 CLIP to measure the similarity between images and
1109 action descriptions before and after augmentation. Specifically,
1110 we generate two text prompts for an image:

1111 • **Human-centric prompt:** “A human is [action].”

1112 • **Robot-centric prompt:** “A robot is [action].”

1113 where [action] describes the task being performed (e.g.,
1114 “picking up a cube”) in the image. We compute CLIP
1115 similarity scores between the image and their respective
1116 prompts, where higher robot-prompt scores indicate better
1117 visual alignment with robotic manipulation.

	Img1	Img2	Img3	Img4	Img5	Img6
ori	31.4	23.7	31.2	32.0	27.5	29.7
aug	28.4	28.8	32.7	29.6	28.7	27.2

Table 4. **CLIP similarity scores.** Higher values indicate better alignment between images and robot-centric descriptions.

1118 Visual examples in Figure 11 demonstrate how our aug-
1119 mentation successfully adapts human motions to robotic
1120 kinematics. The CLIP similarity scores in Table 4 con-
1121 firm that the augmented images maintain comparable align-
1122 ment with robot-centric descriptions. CLIP scores occa-
1123 sionally decrease for some tasks, likely due to minor ar-
1124 tifacts in hand-object interaction synthesis. However, as
1125 demonstrated in the following subsections, these discrep-
1126 ancies do not impede downstream policy performance, sug-
1127 gesting that H2R prioritizes functionally relevant visual fea-
1128 tures over pixel-perfect realism.

1129 **Failure Cases Analysis of the Real-world Experiment** In
1130 addition to evaluating the success rate, we performed a de-
1131 tailed analysis of all failure cases by decomposing each task
1132 into distinct operational phases, as shown in Table 5. We
1133 divided three real-world tasks into multiple stages based on
1134 the complete motion sequence of a robotic arm. For each
1135 task, we classified the failure cases according to the furthest
1136 phase achieved, where later stages correspond to higher task
1137 completion levels.

1138 In our real-world evaluation, we show the frequency dis-
1139 tribution of task-specific failure cases in Figure 8. we find
1140 that regardless of the task or model used in the experiments,
1141 Case 1 constitutes a significant proportion of the failure
1142 cases. A major reason for this is the policy’s inability to
1143 accurately locate the position to interact with the target ob-
1144 ject. This misalignment can be attributed to various factors
1145 such as camera noise, environmental lighting changes, ob-
1146 ject occlusions, or the model’s limited adaptability to new
1147 environments. To delve deeper into these issues, we present

tasks	Pick and Place	Stack Cubes	Pick from Box
case1	Picking failure	Picking failure	Picking failure
case2	Placing failure	Stacking failure	Placing failure
case3	Success	Success	Closing failure
case4	/	/	Success

Table 5. **Cases of each real-world task.** We list 3 cases for Pick and Place task, 3 cases for Stack Cubes task and 4 cases for Pick from Box task.

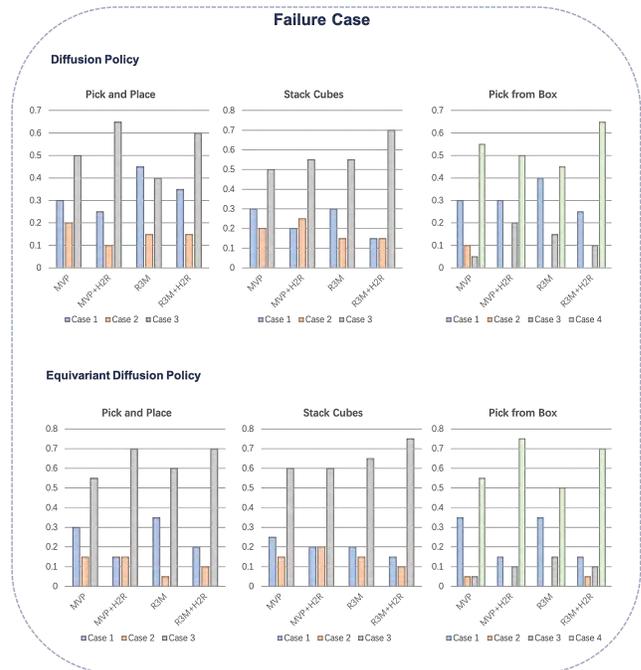


Figure 8. **Failure case analysis.** We divided each task into 3-4 cases to enable a detailed analysis of execution failure causes and finer-grained evaluation.

1148 detailed examples of typical failure cases in Figure 9. These
1149 factors can prevent the end-effector from correctly identi-
1150 fying and approaching the target location, leading to task
1151 failure. We also observe that H2R-augmented visual repre-
1152 sentation models not only improve overall success rates in
1153 real-world tasks but also significantly reduce the occurrence
1154 of Case 1 failures across most of the tasks, which indicates
1155 that even in failed attempts, the robot demonstrates higher
1156 partial-task completion.

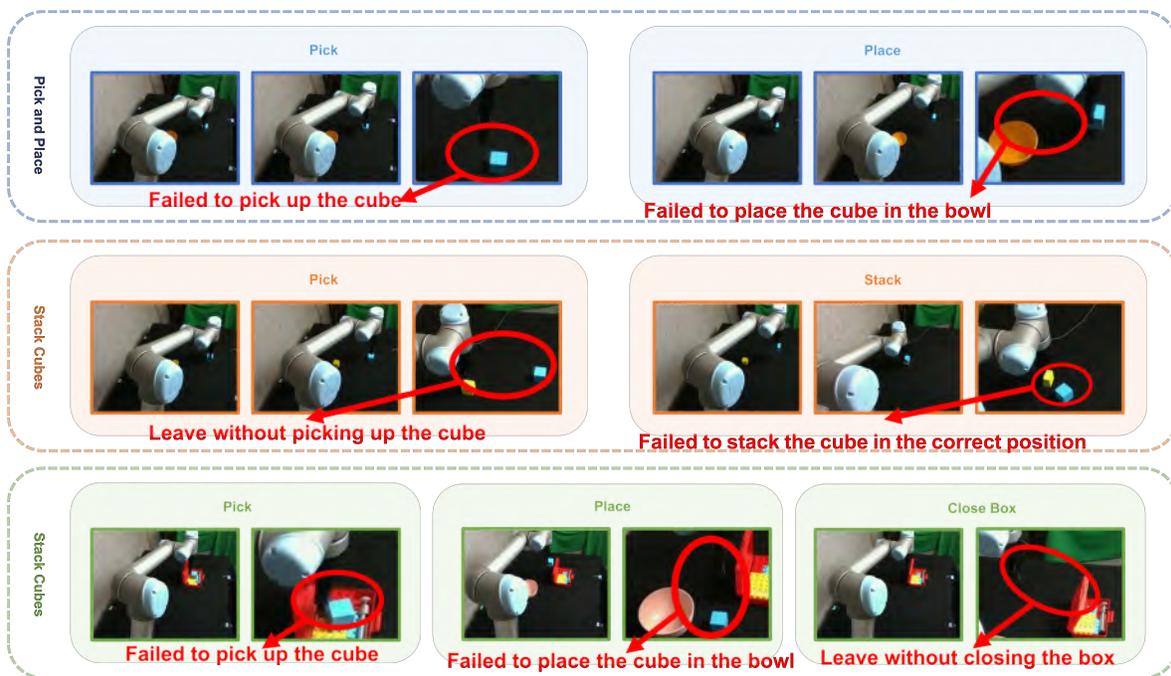


Figure 9. **Failure case visualization.** We provided first-person visualizations from the robot’s viewpoint for every failure scenario.



Figure 10. **Real-world task.** Illustration of three real-world manipulation tasks ranging from simple to complex.



Figure 11. **H2R samples.** Visual comparison between original human data (top) and our augmented data (bottom).