# AutoEval: Autonomous Evaluation of Generalist Robot Manipulation Policies in the Real World

Zhiyuan Zhou<sup>1</sup>, Pranav Atreya<sup>1</sup>, You Liang Tan<sup>2</sup>, Karl Pertsch<sup>1</sup>, Sergey Levine<sup>1</sup>

Abstract-Scalable and reproducible policy evaluation has been a long-standing challenge in robot learning. Evaluations are critical to assess progress and build better policies, but evaluation in the real world, especially at a scale that would provide statistically reliable results, is costly in terms of human time and hard to obtain. Evaluation of increasingly generalist robot policies requires an increasingly diverse repertoire of evaluation environments, making the evaluation bottleneck even more pronounced. To make real-world evaluation of robotic policies more practical, we propose AutoEval, a system to autonomously evaluate generalist robot policies around the clock with minimal human intervention. Users interact with AutoEval by submitting evaluation jobs to the AutoEval queue, much like how software jobs are submitted with a cluster scheduling system, and AutoEval will schedule the policies for evaluation within a framework supplying automatic success detection and automatic scene resets. We show that AutoEval can nearly fully eliminate human involvement in the evaluation process, permitting around the clock evaluations, and the evaluation results correspond closely to ground truth evaluations conducted by hand. To facilitate the evaluation of generalist policies in the robotics community, we provide public access (https://auto-eval.github.io) to multiple AutoEval scenes in the popular BridgeData robot setup with WidowX robot arms. In the future, we hope that AutoEval scenes can be set up across institutions to form a diverse and distributed evaluation network.

#### I. INTRODUCTION

Robot foundation models promise to drastically change the robot learning "workflow": instead of training policies for individual tasks or environments, these models are trained across a range of scenes, tasks, and robot embodiments [1]-[9], providing generalist policies that can solve new tasks in new settings. This shift to generalist training necessitates an analogous shift in how these policies are evaluated. While traditional evaluations for single-task policies typically involve a few dozen policy rollouts that are practical to do by hand, robot foundation models may require hundreds of rollouts across a variety of tasks and scenes to obtain an accurate assessment of their generalist capabilities. For instance, a comprehensive evaluation of the recently introduced OpenVLA model [4] against its baselines required more than 2,500 rollouts across four robot setups and three institutions, and a total of more than 100 hours of human labor for resetting scenes, rolling out policies, and recording success rates. Evaluations during the course of model development and design ablations may compound this effort multiple times over. Prior works have tried to address this evaluation bottleneck by building realistic simulated environments for evaluation [10], but the gap between simulation and the real

world can render results unreliable, and many tasks like cloth or liquid manipulation are challenging to simulate at sufficient fidelity. In this work we aim to develop a system for robot policy evaluation that combines the *reliability* of real world evaluations, with the *scalability* required for the evaluation of generalist robot policies.

A key bottleneck for the scalability of real-world robot evaluations is the human operator time required to conduct the evaluation, reset the scene, and score policy success. If we can reduce required human involvement to a minimum, we can drastically increase the throughput of real robot evaluations by running evaluations around the clock.

Our central contribution is the development of an autonomous evaluation system, AutoEval, that can evaluate user-supplied policies in the real world around the clock. We demonstrate that AutoEval can scale to diverse evaluation environments by instantiating it in three automated evaluation environments for table-top manipulation tasks in the BridgeData V2 environment [11]. Our experiments show that the two aspects of evaluation that typically rely most on human effort, scene resets and success determination, can both be automated with high fidelity with learning based methods, yielding evaluation results that correlate well with ground truth human evaluations. AutoEval drastically increases the evaluation throughput, enabling 500 evaluation episodes per 24-hour period. We also find that AutoEval provides a more reliable policy performance estimate than prior simulated evaluation approaches or offline metrics, while at the same time supporting a wider range of hardto-simulate tasks like cloth manipulation.

We open-source our code and a detailed step-by-step guide for setting up new automated evaluation platforms. Additionally, we open access (https://auto-eval. github.io) to multiple of our Bridge-AutoEval cells, enabling researchers from other institutions to evaluate their policies on our Bridge-AutoEval systems. We hope that this takes a step towards democratizing robotics research and enabling fair comparisons of robot policies on unified evaluation setups.

# II. AUTONOMOUS EVALUATION OF ROBOT POLICIES IN THE REAL WORLD

During robot evaluations, the policy is typically asked to perform the same task multiple times, while applying randomizations to the initial state of the robot and the environment, to get a statistically significant estimate of the policy's performance under the initial state distribution  $\rho(s)$ .

<sup>&</sup>lt;sup>1</sup> University of California, Berkeley, <sup>2</sup> NVIDIA



Fig. 1. We introduce AutoEval, a system for scalable, automated real robot evaluation of generalist robot policies. Automated evaluation results closely match human-run evaluations, while providing a more reliable performance signal than prior simulated evaluation approaches or offline metrics. AutoEval reduces human supervision time for evaluation by more than 99%. We provide public access to our AutoEval cells to facilitate standardization and ease of policy benchmarking.

The output of the policy evaluation is an evaluation score ranging from 0 to 1, representing the success probability.

We present an overview of our AutoEval system in algorithm 1. At its core, it follows the same structure as a conventional, human-run evaluation, running multiple trials with intermittent resets and performance scoring. However, AutoEval introduces multiple learned modules that automatically perform the tasks that typically require a *human* evaluator. Namely, AutoEval consists of three key modules: (1) a success classifier, that evaluates a policy's success on a given task, (2) a reset policy, that resets the scene back to a state from the initial state distribution upon completion of a trial, and (3) programmatic safety measures and fault detections that prevent robot damage and call for human intervention when necessary. All three components are implemented via flexible, *learned* models, and can thus be easily adapted to automate the evaluation of a wide range of robot tasks. Algorithm 1 Autonomous Policy Evaluation Loop (AutoEval)

- 1: **Input:** Task T, policy  $\pi$ , initial state distribution  $\rho(s)$ , success classifier  $C_T$ , reset policy  $\pi_T$ , reset classifier  $C_{\rho(s)}$
- 2: **Output:** Estimated probability of success for task T
- 3: for each trial do
- 4: Start State: Start from initial state  $s_0 \sim \rho(s)$
- 5: **Policy Rollout:** Rollout policy  $\pi$  for K timesteps
- 6: Success Check: Assign success label using  $C_T(s_K)$
- 7: **Reset Scene:** Rollout reset policy  $\pi_T$  to return to  $\rho(s)$
- 8: **Failure:** If unable to reset or robot unhealthy, notify human operator to help
- 9: end for

Next, we provide details on the design and training of each component of our AutoEval system.

Success classifier. The success classifier  $C_T : S \rightarrow \{0, 1\}$ serves to approximate the ground truth task-success T:  $S \rightarrow \{0,1\}$  that maps image states to a binary success label. Instead of hand-crafting a task-specific success rule as done in prior work [12], [13], AutoEval trains a learned success classifier C<sub>T</sub>, a recipe which can be easily applied to a wide range of robot tasks. Concretely, we collect a small set of example images of success and failure states. We use approximately 1000 images, which takes less than 10 minutes to collect by tele-operating the robot and saving the frames in the trajectory. We then fine-tune a pre-trained vision-language model (VLM) for the task of binary success detection. Given a language prompt, e.g., "Is the drawer open? Answer yes or no", and an image observation, the model is trained to predict whether the task was successfully completed. We use a pre-trained VLM to obtain a classifier that is robust to small perturbations of the environment without needing to collect a large number of example images for fine-tuning. In practice, we use the Paligemma VLM [14] for training the success classifier. More detailed information is provided in Appendix L.

**Reset policy.** The reset policy  $\pi_T(a|s)$  "undoes" what the evaluation policy  $\pi$  did during the evaluation rollout, returning the scene and robot to a state from the initial state distribution  $\rho(s)$ . Again, instead of relying on taskspecific "hardware resets" like springs or magnets, our aim with AutoEval is to design a system that can be flexibly applied to a range of robot tasks. We thus use a learned policy for resetting the scene. To learn a reset policy, we manually collect a small set of approximately 100 highquality demonstrations trajectories that reset the scene from plausible end-states of both successful and failed policy rollouts. In practice, this data collection takes typically less than two hours. We then fine-tune a generalist robot policy with behavioral cloning to act as a reset policy. Starting from a generalist policy checkpoint ensures that the reset policy is more robust, and fewer reset demonstrations are required to obtain reliable resets.



Fig. 2. Bridge-AutoEval cell: our robot setup for autonomous policy evaluation in the real world. It consists of a WidowX 250 6-DoF robot arm and a Logitech C920 HD RGB-camera. The scenes reproduce popular evaluation tasks from the BridgeData [11] robot dataset.



Fig. 3. Three scenes in the Bridge-AutoEval experiments: sink, drawer, and cloth. In total we support five tasks for autonomous evaluation: two pick-and-place task in sink, two drawer tasks in drawer, and one deformable cloth manipulation task in cloth.

# III. BRIDGE-AUTOEVAL: OPEN-SOURCE AUTOMATED EVAL PLATFORM

In this section, we describe an instantiation of our automated evaluation system for multiple environments and tasks from the BridgeData V2 dataset [11], [15]. BridgeData is a diverse manipulation dataset containing 60k+ manipulation demonstrations with a WidowX 6DoF robot arm, that span 13 different skills and 24 environments. State-of-the-art generalist manipulation policies like OpenVLA [4], RT2-X [1], CrossFormer [7], and Pi0 [8], [16] are all trained on BridgeData or a super set of it [17], and therefore policy evaluations on this setup are a natural testbed for scalable evaluation approaches for generalist policies.

We built three Bridge-AutoEval cells as shown in Figure 3, which we call the drawer scene, the sink scene, and the cloth scene. Each scenes support evaluation of one to two manipulation tasks: drawer supports evaluating "open the drawer" and "close the drawer"; sink supports evaluating pick-and-place tasks "put the eggplant in the blue sink" and "put the eggplant in the yellow basket"; cloth support the deformable object manipulation task "fold the cloth from top right to bottom left".

One contribution of our work is that we make two of our Bridge-AutoEval cells **publicly available** at https:// auto-eval.github.io, so other researchers can schedule evaluations for their policies. We hope that over time, this can contribute to making evaluations in robotics more reproducible and comparable. To make this practical, we provide a public web UI to access our Bridge-AutoEval cells and monitor the evaluation progress, as shown in Figure 7. Users can choose the scene and task on which they want to perform evaluation, and provide the IP address and port for a local "policy server", that serves the policy they want to evaluate. Our Bridge-AutoEval system will automatically queue the jobs for evaluation, and execute evaluation jobs around the clock from all users. At the end of a policy evaluation, AutoEval provides users with downloadable rollout data and a detailed performance report of the autonomous evaluation, which contains rollout videos, success rates, episode durations, and frequencies of motor resets or required human interventions. Figure 8 shows part of an example report.

# IV. EXPERIMENTAL RESULTS

# A. Experimental Setup

**Tasks.** We evaluate policies on the five Bridge V2 [11] evaluation tasks described in section III: opening and closing a drawer, placing a plastic eggplant in a sink and a basket, and folding a piece of cloth. Details in Appendix D.

**Policies.** We run evaluations with six recently released generalist robot policies from the robotics community: **Open-VLA** [4], **Octo** [3], **Open-** $\pi_0$ , **MiniVLA**, **SuSIE**, **SuSIE-LL**. Detailed description of these policies in Appendix E.

**Comparisons.** We compare multiple approaches for scalable evaluation of generalist policies. Concretely, we compare our approach, AutoEval, to prior work on simulated evaluation of robot manipulation policies, SIMPLER [10]. SIMPLER builds realistic simulated versions of real-world environments and then evaluates policies purely in simulation. For our experiments, we reuse the existing SIMPLER environment for the Bridge sink environment, and build a new SIMPLER simulation environment for the drawer scene (See Figure 9). Deformable objects such as the cloth in our cloth scene are hard to simulate [18], [19] and the simulator of SIMPLER, Maniskill [20], does not support simulating deformable objects so we do not evaluate the cloth scene in simulation. In addition, we compare to using mean-squared error on a validation set ("val-MSE") as a scalable approach for offline evaluation of robot policies.

**Metrics.** Human-run real-world evaluations is the gold standard for robotic policy evaluation. Following [10] we use two metrics to measure how closely the respective evaluation results match those of human-run evaluations: (1) **Pearson correlation** [21], which measures the linear consistency between two random variables, and is a widely used statistical tool for assessing correlation, with scores nearing 1 indicating high correlation. (2) **MMRV** (Mean Maximum Rank Violation) [10], which measures the consistency of policy *ranking* and, as described in [10], can be more robust to noise on the evaluation results. Low MMRV scores indicate closely matching evaluation results. Details in Appendix F.

#### B. AutoEval Closely Matches Human Evaluation Results

For each evaluation method, we run 50 evaluation rollouts for each policy in each of our five tasks (except "val-MSE",



Fig. 4. Correlation of scalable evaluation approaches to oracle humanrun evaluations. AutoEval closely matches human evaluations, achieving high correlation and low MMRV score (plotted in the figure is 1 - MMRVfor clarity). In contrast, SIMPLER simulated evaluations and validation MSE do not correlate as well with human evaluations.



Fig. 5. Qualitative visualization of AutoEval evaluation rollouts on three of our tasks. After the policy execution is done, the success classifier determines whether the rollout was successful. Then, the reset policy returns the environment into a state from the initial state distribution for the next evaluation. Our evaluations cover representative robot manipulation tasks: pick-place, articulate and deformable object manipulation.

which does not require rollouts). We report results in fig. 4, with a detailed breakdown of results per task, policy, and evaluation method in Appendix, table I to table III. Similar to prior work [10], we find that simple validation MSE is a poor evaluation metric for robot policies: it actually negatively correlates with real robot performance and thus does not provide a reliable performance estimate. We find that SIM-PLER evaluations in simulation provide a better performance signal, but lack reliability. Concretely, our results show that SIMPLER occasionally matches real-world performance well (e.g., for the "open drawer" task), but in other cases not accurately reflects the policy's performance. For example for Open- $\pi_0$  in the "put eggplant to sink" task, the policy performs very poorly in simulated evaluations, but achieves high success rate in the real world. Intuitively, different policies may suffer differently from the remaining sim-toreal gap in SIMPLER evaluations. As a result, SIMPLER's effectiveness is policy dependent and it cannot provide a *reliable* policy evaluation.

In contrast, we find that our approach, AutoEval, closely matches the results of oracle human-run evaluations, with an average Pearson score of 0.942 and MMRV of 0.015 (plotted as 1-MMRV in Figure 4 of 0.985). In particular, an MMRV score close to zero indicates that it rarely disrupts the ranking of policies. Intuitively, since evaluations are still run in the



Fig. 6. Visualization of a 24 hour AutoEval evaluation run with ~850 total evaluation episodes. AutoEval is able to run autonomously over extended periods of time and only required a total of 3 human interventions over a 24 hour period. On average, the evaluation throughput of AutoEval is on par with that of human evaluations, but saves **99%+** human operator time.

real world, there is no sim-to-real gap that could negatively affect policy performance. In practice, we find that success detector and reset policy work reliably during evaluation. We show qualitative examples of autonomous evaluation rollouts in fig. 5, and further examples in appendix I. Importantly, we find that AutoEval drastically reduces the human effort required to run real robot evaluations, cutting the human evaluator time for robot evaluations by > 99% compared to conventional, human-run evaluations.

### C. AutoEval Robustly Runs Over Long Time Spans

For this investigation, we performed a long-running evaluation session over the course of 24 hours, repeatedly interleaving the evaluation of various policy checkpoints, using the "open drawer" and "close drawer" tasks. In Figure 6, we present the evaluation throughput over the 24 hours, as well as the number of human interventions needed over the span of the whole evaluation. We present evaluation throughput in terms of the number of valid evaluation steps taken per minute Over the course of a day, a single AutoEval cell is able to run 60,000 evaluation steps (roughly 850 episodes on the drawer scene), with an average speed of 42 evaluation steps per minute. The AutoEval throughput varies in Figure 6 because of the different inference speed of different policies. The average AutoEval speed, shown in dotted blue line, is slightly lower but on par with the average evaluation speed of a human evaluator performing manual resets of the environment and recording success rates. Even though AutoEval has a slightly lower throughput, AutoEval runs autonomously and only required a total of three human interventions in the span of 24 hours to reset the scene or robot. Every time a human operator needed to intervene, they simply needed to check and reset the objects' position in the scene, and potentially move the robot arm into reset position if a motors failed and the robot fell on the table. Afterward, the human operator can make AutoEval resume simply with the press of a button. Assuming that each human reset operation takes 1 minute, 24h of autonomous evaluation only costs 3 minutes of human time, compared to  $\approx 16$  hours if a human evaluator wanted to run the same number of trials by hand. This means that AutoEval can reduce human time involvement by >99%.

In addition, we evaluate AutoEval's performance over two months of continuous operation. appendix Q shows that AutoEval yield consistent results over long time periods.

#### Appendix

# A. AutoEval Web UI Interface

Figure 7 illustrates the web UI for submitting evaluation jobs to Bridge-AutoEval and Figure 8 illustrates the resulting report provided to users upon completion of the automated evaluation.



Fig. 7. Web UI for submitting evaluation jobs to the Bridge-AutoEval cells. Users choose a desired task and provide the IP address for a policy server they host for evaluation, and can monitor the evaluation through the UI.



Fig. 8. Excerpt from an AutoEval result report, provided to the user upon completion of the automated evaluation. Users can see the per-episode success rate, rate of evaluation progress, instances of automatic recovery from motor failures, and qualitative rollout videos as well as classifier result plotted with initial and final frames to obtain a wholistic understanding of the policy's performance.

#### B. RELATED WORK

Generalist robot policies. Recent advancements in robot foundation models have demonstrated significant progress [1], [3], [4], [6], [8], [16], [22]–[27], fueled by large-scale robot datasets [5], [11], [17], [28]. These models are trained to perform diverse tasks (e.g., pick-and-place,

cloth folding) [4], [8], [11], [16], adapt to various scenes with different backgrounds and distractors [29], [30], and control multiple robot embodiments (e.g., quadrupeds, manipulator arms, drones) [7], [31]. With the increase in capabilities of these generalist robot policies, evaluation becomes ever more time-consuming, because measuring model performance needs evaluations of a variety of different skills and scenes. For example, reporting results for [4] required a few thousand evaluation trials and more than 100 hours of human labor. Evaluation trials needed during development probably compounded this number several times. This makes development and comprehensive evaluation of generalist robot policies increasingly challenging, calling for an evaluation method that is much more scalable.

Robot policy evaluation in the real world. Evaluating robot policies in a fair, comprehensive, and reproducible way is challenging. Robotic methods and systems today are mostly tested in custom settings at the institution where the method is developed. Cross-institution evaluation encounters difficulties with different hardware, task definitions, and performance measures [32]. To address this, multiple works have proposed real robot setups that have reproducible components (such as 3-D printed objects or cheap robot hardware) that are meant to be replicated across institutions [11], [33]-[38]. In addition to robot manipulators, there have also been efforts for standardized hardware in other robot embodiments [39], [40]. However, the sensitivity of policies to environmental factors like lighting, camera angles, and robot type makes it hard to accurately reproduce real robot setups across institutions, even when the same set of objects and hardware are used. Others have built evaluation systems that are hosted at a central location to compare different approaches. Some take the form of live competitions [32], [41]–[44], while others are hosted at research institutions and open to the public [45], [46]. However, these evaluations all require human involvement to supervise the policy evaluation or to reset the scene, making it expensive in terms of human time and therefore significantly limiting the number of real robot evaluations benchmark participants can perform. In addition, the live competitions are logistically challenging and hard to operate continually. These reproducibility and scalability constraints will become even more apparent as the capabilities of robot policies expand to more scenes, tasks, and embodiments. Our approach, AutoEval, can substantially improve the throughput of real robot evaluations by replacing parts of the evaluation pipeline traditionally completed by humans with specialized learned components, thus enabling robots to "evaluate themselves" 24/7. Notably, [47] proposed a setup for remote, autonomous policy evaluation in the real world as part of their Real Robot Competition, but they focused on evaluations in a single environment, engineered to require no resets and allow for scoring with task-specific, hand-defined rules. In contrast, our AutoEval system is designed for evaluation of generalist policies by enabling autonomous evaluation on a wider range of tasks (e.g., pickplace, articulate object & cloth manipulation) via learned reset and scoring modules. While the goal of this work is not to build a comprehensive benchmark for robot foundation models, which requires evaluations spanning many tasks, scenes, and embodiments, we demonstrate that our system can be used to automate evaluations for a diverse set of tasks and provide a step-by-step guide to set up new automated evaluations within a few hours. We hope that by reproducing this recipe at other institutions, the robotics community will over time be able to construct a comprehensive evaluation benchmark for generalist policies.

Evaluation in simulation. While human-run evaluations in the real world are the gold standard used by most prior works, they require extensive human effort and do not scale well as the capabilities of models increase. As a result, simulation has been a popular tool for high-throughput evaluation in robot learning research [48]-[63]. However, there are still discrepancies between these simulators and the real world, making simulated evaluation different from realworld evaluation. First of all, real-world physics of contacts, collisions, and friction are hard to simulate accurately [64]-[70]. Even if the physics simulation is perfect, not all physical parameters can be precisely measured in the real world to be replicated in simulation (e.g., friction coefficients and actuation delays) [71], [72]. Policies that interact with real-world objects usually exhibit different behavior than they do on their simulated counterparts. Secondly, policies need to deal with real world factors such as noisy and delayed sensory inputs that do not play a big part in simulation. Finally, the visual difference such as texture and lighting between simulated images and real-world observations makes the two types of evaluation quite different [73], [74]. Recent works have tried to reduce the visual discrepancy by building realistic simulators for policy evaluation [10], [75]. SIMPLER [10] constructs high-fidelity replicas of real robot evaluation scenes and demonstrates strong correlation of simulated rollouts to human-run rollouts in the corresponding real robot environments. However, gaps between simulation and the real world remain, and our evaluations show that they can affect different models to varying degrees, leading to inconsistent policy performance rankings between simulation and real world evaluation. Additionally, there is a large number of tasks, like cloth or liquid manipulation, that are challenging to simulate at sufficient fidelity to enable simulated evaluation. In contrast, our approach performs evaluations on real robot systems and thus provides a more reliable signal for policy performance, including on tasks that are hard to simulate, while retaining scalability by minimizing the need for human intervention.

**Autonomous robot operations.** Multiple prior works have identified the need for human supervision and oversight as a key limiting factor in robot learning [22], [29], [76]–[79]. While these works typically focus on autonomous policy *improvement* instead of autonomous policy *evaluation*, they share many challenges around robot resets and success detection. Thus, many of the techniques we employ for learning reset policies and success detectors are inspired by prior work in autonomous robot learning, and even some of the metrics are shared, e.g., measuring the frequency of human interven-

tion [80]. However, to our knowledge, our work is the first to explore the design of a general system for autonomous evaluation of generalist policies. While most robot learning researchers are (painfully) aware of the cost of evaluations, existing efforts toward automating real robot evaluations have been limited to task-specific solutions that often involve instrumenting the environment, e.g., with spring-driven or scripted reset mechanisms [13], [22], [81]. In contrast, we provide a task-agnostic, scalable approach for automating robot evaluations with flexible, learned components based on generalizable and broadly applicable foundation models.

#### C. SIMPLER environments

Figure 9 illustrates the two SIMPLER simulated environments we use in Section IV.



Fig. 9. SIMPLER [10] simulated evaluation scenes for the tested environments. Simulated evaluation is fast and cheap, but can struggle from visual and physics discrepancies between simulation and the real world.

#### D. Experimental Details

For our evaluations on the robot tasks, all tasks are performed using a WidowX 6-DoF robot arm. During humanrun evaluations, success is counted when the drawer is completely closed or opened at least 1.5cm, respectively, if the eggplant is fully inside the sink or basket at the end of the episode, and if the cloth is folded to at least a quarter of the way diagonally. We randomize the initial position of the eggplant, drawer, and the cloth at the beginning of each episode.

# *E.* Generalist Robot Policies Evaluated in Experiments Section.

Here we provide more details on all the generalist robot policies we evaluate in Section IV. OpenVLA [4], a 7B parameter vision-language-action model (VLA) pre-trained on the Open X-Embodiment dataset [17], Octo [82], a 27M parameter transformer policy, also pre-trained on Open X-Embodiment, **Open-** $\pi_0$  [83], an open-source reproduction of the 3B parameter  $\pi_0$  VLA [8] (the original  $\pi_0$  was not available in open-source at the time of writing), pre-trained on the Bridge V2 dataset, MiniVLA [84], a 3B parameter VLA pre-trained on the Bridge V2 dataset [11], SuSIE [85], a hierarchical policy that combines a image diffusion subgoal predictor with a small diffusion low-level policy, pre-trained on Bridge V2, and SuSIE-LL, which directly executes the goal-conditioned behavioral cloning low-level policy from SuSIE. This set of policies is a representative sample of current state-of-the-art generalist policies. All policies contain



Fig. 10. AutoEval evaluation scores remain consistent over 8 hours of autonomous evaluations. After 8 hours, the WidowX robot's motors overheat and evaluation scores start to drift. As a result, we pause evaluation for 20 min every 6 hours to let the motors cool off. Error bars show 95% confidence intervals.

the Bridge V2 dataset as part of their training data, and we evaluate the publicly released checkpoints for all models.

#### F. MMRV Score Details

MMRV is computed as follows: given N policies  $\pi_{1..N}$ and their respective success rates  $R_{A,1..N}$ ,  $R_{B,1..N}$  estimated via two evaluation procedures A and B, we compute:

$$\begin{aligned} \text{RankViolation}(i,j) &= |R_{A,i} - R_{A,j}| \\ &\cdot \mathbf{1}[(R_{B,i} < R_{B,j}) \neq (R_{A,i} < R_{A,j}) \\ \\ \text{MMRV}(R_A, R_B) &= \frac{1}{N} \sum_{i=1}^{N} \max_{1 \le j \le N} \text{RankViolation}(i,j). \end{aligned}$$

For each tested evaluation approach we compute MMRV with reference to human-run "oracle" evaluations, where low MMRVs indicate closely matching evaluation results.

#### G. Are AutoEval results consistent across time?

We test the consistency of AutoEval evaluations, i.e., Auto-Eval's ability to produce comparable performance estimates across multiple iterations of evaluating the same policy. To test this, we run the Open- $\pi_0$  policy through a sequence of 9 evaluations on the "open drawer" task, each consisting of 50 individual trials, or a total of 450 trials. Using AutoEval, the full evaluation takes ~11 hours. We report the results of this evaluation in fig. 10. We find that AutoEval produces consistent evaluation results across long periods of time. Concretely, for the first 7 evaluation runs, or a total of 350 evaluation episodes, AutoEval performance evaluation are within the margins of what might be considered the natural variance of robot evaluations  $(\pm 10\%)$ . We see a regression in performance after approximately 8 hours of continuous operation, which we attribute to an overheating of the motors of our rather affordable WidowX robot (;\$3500) after many hours of operation. To mitigate the effects of such overheating in practice, we pause autonomous evaluations for 20 minutes every 6 hours to let the motors cool off before resuming evaluations.

#### H. Safety During Extended Autonomous Robot Operations

To ensure that the robots can autonomously and safely operate for a long time, we take several measures to ensure the safety of the robot and to preserve the scene. First, we set safety boundaries for the robot such that the policy cannot go beyond certain xyz axis (e.g. beyond the view of the camera) so that it does not run into objects unintentionally. Second, since the WidowX robot arms do not natively support impedance control, we limit the maximum effort on each of the robot joints, so that ineffective policies do not press too hard against objects and cause motor failure. The common robot failure is due to joint failure when interacting and colliding with the objects in the scene, hence we constantly monitor and log the joint effort values, safely software reboot the joints when joint errors are detected during each trial.

#### I. Visualizations of AutoEval Rollouts

Figure 11 presents evaluation trajectories in the five different Bridge-AutoEval tasks. The actual language commands fed to the evaluated policies are:

- 1) "Close the drawer"
- 2) "Open the drawer"
- 3) "Put the eggplant in the yellow basket"
- 4) "Put the eggplant in the blue sink"
- 5) "fold the cloth from top right to bottom left"

# J. Evaluation on Bridge-SIMPLER [10]

In AutoEval, we introduced a new **Drawer Scene** to the existing SIMPLER [10] setup for the WidowX robot. The scene was visually matched with the AutoEval's Drawer setup, and overlaid with the same background to ensure consistency. A 3D model of the drawer, with exact dimensions matching the real-world setup, was also created. This scene introduced two evaluation tasks: "*open and close the drawer*". To add variability to each evaluation trial, we randomized the end effector's initial pose, the drawer's initial pose, and the lighting conditions in the background.

In addition to the Drawer Scene, AutoEval includes a **Sink Setup**, which closely resembles the existing SIMPLER [10] Sink Scene. In SIMPLER, the task here is the "move the eggplant to the basket" task. We also introduced a reverse task, "move eggplant to the sink," effectively making the scene reset-free. This allows for both forward and reverse tasks in the same environment.

With these two scenes and four tasks, we conducted 50 runs for each scene across five different generalist policies. The detailed results are shown in table I

# K. Detailed Evaluation Results on Bridge-AutoEval

In table II to table III, we provide detailed evaluation results for our comparison of different scalable evaluation approaches across the five Bridge V2 evaluation tasks.

#### L. Success Classifier in Bridge-AutoEval cells

To train success classifiers for the Bridge-AutoEvalscenes, we finetune the Paligemma VLM to act as a classifier. We manually collect a dataset of roughly 1000 images for each



Fig. 11. Samples of autonomous policy evaluation trials with AutoEval on the five tasks. The classifier result of each task is visualized on the right hand side, and the reset policy is not shown.

Policy	Drawer Scene		Sink Scene		
	Open	Close	Eggplant	Eggplant	
	Drawer	Drawer	to Basket	to Sink	
OpenVLA	2/50	27/50	1/50	0/50	
Open $\pi_0$	34/50	25/50	46/50	6/50	
Octo	3/50	0/50	6/50	3/50	
SuSIE-LL	1/50	0/50	0/50	0/50	
SUSIE	0/50	41/50	7/50	0/50	
MiniVLA	30/50	23/50	11/50	2/50	

TABLE I

EVALUATION RESULTS ON SIMPLER [10] FOR DRAWER AND SINK Scene on four tasks and six different policies.

Policy	Drawer		Sink		
	Open	Close	То	То	Fold
	Drawer	Drawer	Basket	Sink	Cloth
OpenVLA	40/50	46/50	1/50	0/50	13/50
Open $\pi_0$	29/50	46/50	7/50	47/50	12/50
Octo	1/50	5/50	0/50	0/50	4/50
SuSIE-LL	0/50	1/50	0/50	0/50	0/50
SuSIE	1/50	18/50	0/50	0/50	9/50
MiniVLA	33/50	49/50	38/50	0/50	8/50

TABLE II

AUTOEVAL RESULTS ON FIVE BRIDGE-AUTOEVAL TASKS ACROSS SIX DIFFERENT GENERALIST POLICIES.

scene, and manually label them. We form VQA questions with the labels, and finetune the base 3B parameter VLM with quantized LoRA using a learning rate of 2e - 5, batch size of 4 for 80 iterations.

For each evaluation scene, approximately 1000 image frames are collected to fine-tune the VLM. The corresponding language prompts are:

- 1) Sink Scene: "is the eggplant in the sink or in the basket? answer sink or basket or invalid"
- 2) Drawer Scene: "is the drawer open? answer yes or no"
- 3) Cloth Tabletop Scene: "is the blue cloth folded or unfolded? answer yes or no"

#### M. Reset Policy in Bridge-AutoEval cells

To train reset policies for the Bridge-AutoEvalcells, we finetune the generalist OpenVLA policy with LoRA, with batch size 64 and learning rate  $10^{-4}$  for 1000 iterations. For each scene, we collect roughly 100 demonstration trajectories via teleoperation, and train with a standard behavior cloning loss. We also use a scripted policy for one of our reset policy - "Close the Drawer" task, where the reset success rate is not sensitive to variation in scene.

# N. Step-by-step AutoEval Construction Guide

Overall, we find that the construction of an AutoEval cell for a new task can be completed within 1-3h of human effort, and less than 5 hours total, including model training time for success classifiers and reset policy. This is compared to tens of hours of human evaluation time that can be saved even

Policy	Drawer		Sink		
	Open	Close	То	То	Fold
	Drawer	Drawer	Basket	Sink	Cloth
OpenVLA	40/50	46/50	1/50	0/50	12/50
Open $\pi_0$	24/50	45/50	7/50	47/50	3/50
Octo	0/50	0/50	0/50	0/50	2/50
SuSIE-LL	0/50	0/50	0/50	0/50	0/50
SuSIE	2/50	13/50	0/50	0/50	10/50
MiniVLA	32/50	49/50	38/50	0/50	8/50
TABLE III					

GROUND TRUTH HUMAN EVALUATION RESULTS FOR THE FIVE BRIDGE-AUTOEVAL TASKS ACROSS SIX DIFFERENT GENERALIST POLICIES.

within a single typical research project. Below, we provide a step-by-step guide for creating an AutoEval setup for a new evaluation tasks.

- Train Reset Policy: Start by collecting approximately 100 high-quality robot demonstrations of resetting behavior from sensible final states of policy rollouts. Try to cover a diverse set of "reset start states", including those that failed the original task, to obtain a robust reset policy. Once you collected the dataset, fine-tune a generalist policy like OpenVLA [4], e.g., using LoRA fine-tuning, on your reset demonstration dataset. If you find that the reset is unreliable and fails often, consider collecting more reset demonstrations and re-training. For a small set of tasks, *scripted* reset policies are feasible too, but they require task-specific setup, which we won't discuss here.
- 2) Train Task Success Classifier: Collect approximately 1000 images of success and failure states. An easy way to collect such data is by "waving" the robot arm near the final state of the task at hand to quickly collect a number of "nearby" states for both successful and unsuccessful task execution. Then fine-tune a vision-language model like Paligemma [14] on this dataset and test the performance of your classifier by scoring a few final states of real policy rollouts. If you find that the classifier still mis-labels these states, you can put the states with corrected labels back into your training dataset and re-train until the accuracy of the classifier is high.
- 3) Setup Safety and Robustness Measures: If your robot does not have an integrated p-stop that prevents forceful collisions with the environment, consider implementing a limit on the motor current to prevent damage of the robot. Additionally, it can be useful to implement workspace bounds that prevent the robot from e.g. throwing objects out of a bin in a large majority of cases. Finally, train a reset success detector that checks whether the reset policy was successful by mirroring the steps above for training the success detector, but this time for the inverse task. We also recommend implementing an "on-call" system that sends push notifications (e.g. via Slack) to operators



Fig. 12. Analyzing 50 AutoEval runs on the sink scene: the main failure modes is false positive results because the reset policy failed to reset the scene.

whenever the robot reports an irrecoverable issue or the reset policy fails for  $N \approx 3$  times in a row.

4) Connect to AutoEval Environment Server: Connect your evaluation setup to the AutoEval environment server. This allows external policies to interact with the real robot evaluation scene via a simple REST API service. This also allows for automated scheduling, evaluation, and logging of results and rollouts within the AutoEval framework. Upon acceptance of this manuscript, we will release the code for our AutoEval software framework, to make reproduction easy.

#### O. Analyzing AutoEval Failure Modes

While our previous experiments show that AutoEval closely matches the results of human-run evaluations, we observe that over extended periods of operation errors occur occasionally. To better understand the sources of these errors and help the design of future autonomous evaluation cells, we perform a detailed analysis of all failures occurring in a 50 episodes AutoEval run on the "put eggplant in blue sink" task with the Open- $\pi_0$  policy. We visualize the outcome of our analysis in fig. 12. While many episodes experienced motor failure because of harsh contact with the scene, AutoEval handles such failure automatically by re-running those trials, and only report evaluation trials that do not contain motor failures. We find that for only three out of 50 trials, the autonomous evaluation fails, since the episodes mistakenly get classified as successes and the reset policy fails.

One key takeaway from this failure analysis is that our Bridge-AutoEval setup is already very reliable with few errors, and that most room for improvement is in improving the *efficiency* by reducing the number of motor failures during evaluation, e.g. by implementing a more compliant robot controller that prevents harsh environment interactions.

# P. Bridge-AutoEval Deployment Details

As described in section III, we open access to our Bridge-AutoEval cells to the research community. The two different AutoEval cells accepts and executes jobs in parallel. While the two WidowX robots will accept evaluation jobs 24/7, we enforce a 20 minute rest period every 6 hours where the robot will torque off and let the motors cool off (see fig. 10 for why this is necessary). The reset period will only happen between evaluation jobs.

Since we host the reset policies for the four tasks in Bridge-AutoEval 24/7, we optimize for lightweight policies (as compared to the fine-tuned OpenVLA reset policy we use in section IV). For the two tasks on Drawer, we use scripted reset policy; for the two tasks on Sink, we fine-tune MiniVLA [84] on the same demos. We find that all reset policies have success rate > 95%.

Additionally, while the AutoEval system is developed for evaluating generalist manipulation policies, the infrastructure that supports autonomous robot operation can in practice be used for other purposes as well, such as autonomous data collection [29] or online reinforcement learning finetuning [86], [87].

#### Q. Evaluation Results Reproducible Across Months

We find that AutoEval reproduces results even after more than 2 months of continued use, demonstrating its robustness to aging effects. We compare AutoEval results that are two months apart for two policies on three tasks as shown in table IV. During the two months, AutoEval operated continuously for a rough total of 200 hours. table IV shows that all evaluations perform similarly when evaluated two months apart, and the reset policy and success classifiers still have accuracies 96% and 96% respectively. We attribute such robustness to (1) safety controllers (Appendix A) limiting robot joint efforts to prevent high-force contact and damages, and (2) foundation model pre-training (Paligemma VLM, OpenVLA) making policies and detectors resilient to minor scene changes. Over two months, we have observed minimal "aging" – e.g. there are light scratches on the drawer upon close examination, but they are not visible in the 256x256 pixel policy image observations and does not impact the drawer physics.

Policy	Open	Close	Eggplant
	Drawer	Drawer	to Basket
OpenVLA (old)	39/50	48/50	4/50
OpenVLA	40/50	46/50	1/50
Avg. $\Delta$ Success	+2%	-2%	-6%
Open $\pi_0$ (old)	30/50	45/50	9/50
Open $\pi_0$	29/50	46/50	7/50
Avg. $\Delta$ Success	-2%	+2%	-4%

TABLE IV

AUTOEVAL RESULTS OBTAINED TWO MONTHS APART: RESULTS REMAIN HIGHLY CONSISTENT ACROSS TWO DIFFERENT POLICIES ON THREE DIFFERENT TASKS. ALL CORRELATE WELL TO HUMAN EVALUATIONS.

#### R. Initial States in Bridge-AutoEval Cells

We find that our learned reset policy is able to reset to a consistent distribution of initial states. As an example, we plot the centroids of all eggplant initial positions for three representative AutoEval runs of the "Eggplant to Basket" task in fig. 13 (50 trials each). Qualitatively, we find that the reset distributions of other tasks are similarly overlapping, and also roughly cover the task distribution.



Fig. 13. Initial state distribution for 3 different AutoEval runs is consistent. Red dots show the centroid position of the eggplant. Each run uses 50 trials.

#### S. Improving AutoEval with Additional Human Involvement

Though AutoEval results highly correlate with ground truth human-run evaluations, it is not perfect (as shown in fig. 12). Additional human effort, when available, can further improve AutoEval's accuracy. The easiest and most effective way to apply extra human effort can be spent going through the evaluation report after AutoEval finishes to remove the runs where the reset policy fails, and relabel the success manually. Going through 50 trials of AutoEval roughly takes 1-2 minutes of human time. This enables ground-truth judgment of evaluations runs while still saving the majority of human time required in robot evaluations.

# T. Limitations

AutoEval environment creation time. Our current approach for creating new environments for automated evaluation requires some up-front human effort to train the reset policy and success classifier. In our experience, the complete process only takes a few hours for a new scene and is quickly outweighed by the time savings of autonomous evaluation, but future work can explore more efficient ways of constructing reset policies and success classifiers to further reduce the effort for setting up a new scene for autonomous evaluation. We also expect that future improvements to vision foundation models and generalist policies will make the training of robust success classifiers and reset policies easier, possibly to the point where we can "train" these modules simply by providing a handful of examples in context.

**Mobile manipulation tasks.** Our experiments capture a set of robot manipulation tasks that are reflective of the kind of tasks commonly used for evaluating generalist robot policies today, where the primary focus is on table-top manipulation tasks. We believe that our approach will transfer well to a wide range of other single-arm and bi-manual table top manipulation tasks. However, mobile robot tasks, particularly mobile manipulation tasks, may pose new challenges e.g. with regards to robust resets at room scale, success estimation under partial observability, and operational safety, all of which pose important directions for future work.

**Binary success metrics.** AutoEval evaluation currently only supports binary success estimates (did the policy succeed at a task or fail). When humans run evaluations, they can provide a more granular assessment of the policy's performance, including task progress scores and a qualitative analysis of the policy's proficiency. While AutoEval users can obtain similar assessments from re-watching the logged evaluation videos, this is a time-consuming process. In future work, it would be exciting to investigate whether more granular performance analysis can be provided in an automatic evaluation framework, e.g., by querying powerful video summarization models.

# ACKNOWLEDGMENT

We would like to thank Kyle Stachowicz and Mitsuhiko Nakamoto for valuable advice and discussions. Pranav is supported by the NSF Graduate Research Fellowship.

#### REFERENCES

- A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, *et al.*, "Rt-2: Visionlanguage-action models transfer web knowledge to robotic control," *arXiv preprint arXiv:2307.15818*, 2023.
- [2] K. Bousmalis, G. Vezzani, D. Rao, C. Devin, A. X. Lee, M. Bauza, T. Davchev, Y. Zhou, A. Gupta, A. Raju, *et al.*, "Robocat: A selfimproving foundation agent for robotic manipulation," *arXiv preprint arXiv*:2306.11706, 2023.
- [3] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo, T. Kreiman, Y. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine, "Octo: An open-source generalist robot policy," in *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- [4] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, *et al.*, "Openvla: An open-source vision-language-action model," *arXiv preprint arXiv:2406.09246*, 2024.
- [5] D. Shah, A. Sridhar, A. Bhorkar, N. Hirose, and S. Levine, "Gnm: A general navigation model to drive any robot," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 7226–7233.
- [6] A. Sridhar, D. Shah, C. Glossop, and S. Levine, "Nomad: Goal masked diffusion policies for navigation and exploration," in 2024 *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 63–70.
- [7] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine, "Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation," *arXiv preprint arXiv:2408.11812*, 2024.
- [8] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, *et al.*, "π<sub>0</sub>: A visionlanguage-action flow model for general robot control," *arXiv preprint arXiv*:2410.24164, 2024.
- [9] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei, "Voxposer: Composable 3d value maps for robotic manipulation with language models," *arXiv preprint arXiv:2307.05973*, 2023.
- [10] X. Li, K. Hsu, J. Gu, K. Pertsch, O. Mees, H. R. Walke, C. Fu, I. Lunawat, I. Sieh, S. Kirmani, S. Levine, J. Wu, C. Finn, H. Su, Q. Vuong, and T. Xiao, "Evaluating real-world robot manipulation policies in simulation," *arXiv preprint arXiv:2405.05941*, 2024.
- [11] H. R. Walke, K. Black, T. Z. Zhao, Q. Vuong, C. Zheng, P. Hansen-Estruch, A. W. He, V. Myers, M. J. Kim, M. Du, *et al.*, "Bridgedata v2: A dataset for robot learning at scale," in *Conference on Robot Learning*. PMLR, 2023, pp. 1723–1736.
- [12] S. Gu, E. Holly, T. Lillicrap, and S. Levine, "Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates," in 2017 IEEE international conference on robotics and automation (ICRA). IEEE, 2017, pp. 3389–3396.
- [13] A. Nagabandi, K. Konolige, S. Levine, and V. Kumar, "Deep dynamics models for learning dexterous manipulation," in *Conference on Robot Learning*. PMLR, 2020, pp. 1101–1112.
- [14] L. Beyer, A. Steiner, A. S. Pinto, A. Kolesnikov, X. Wang, D. Salz, M. Neumann, I. Alabdulmohsin, M. Tschannen, E. Bugliarello, *et al.*, "Paligemma: A versatile 3b vlm for transfer," *arXiv preprint arXiv:2407.07726*, 2024.
- [15] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and S. Levine, "Bridge data: Boosting generalization of robotic skills with cross-domain datasets," *arXiv preprint arXiv:2109.13396*, 2021.

- [16] K. Pertsch, K. Stachowicz, B. Ichter, D. Driess, S. Nair, Q. Vuong, O. Mees, C. Finn, and S. Levine, "Fast: Efficient action tokenization for vision-language-action models," *arXiv preprint arXiv:2501.09747*, 2025.
- [17] O. X.-E. Collaboration, A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Singh, A. Brohan, A. Raffin, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Ichter, C. Lu, C. Xu, C. Finn, C. Xu, C. Chi, C. Huang, C. Chan, C. Pan, C. Fu, C. Devin, D. Driess, D. Pathak, D. Shah, D. Büchler, D. Kalashnikov, D. Sadigh, E. Johns, F. Ceola, F. Xia, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Schiavi, H. Su, H.-S. Fang, H. Shi, H. B. Amor, H. I. Christensen, H. Furuta, H. Walke, H. Fang, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Kim, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu, J. Wu, J. Luo, J. Gu, J. Tan, J. Oh, J. Malik, J. Tompson, J. Yang, J. J. Lim, J. Silvério, J. Han, K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, K. Zhang, K. Majd, K. Rana, K. Srinivasan, L. Y. Chen, L. Pinto, L. Tan, L. Ott, L. Lee, M. Tomizuka, M. Du, M. Ahn, M. Zhang, M. Ding, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. D. Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, P. R. Sanketi, P. Wohlhart, P. Xu, P. Sermanet, P. Sundaresan, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Martín-Martín, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Moore, S. Bahl, S. Dass, S. Song, S. Xu, S. Haldar, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Dasari, S. Belkhale, T. Osa, T. Harada, T. Matsushima, T. Xiao, T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, V. Jain, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Wang, X. Zhu, X. Li, Y. Lu, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Xu, Y. Wang, Y. Bisk, Y. Cho, Y. Lee, Y. Cui, Y. hua Wu, Y. Tang, Y. Zhu, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Xu, and Z. J. Cui, "Open X-Embodiment: Robotic learning datasets and RT-X models," https://arxiv.org/abs/2310.08864, 2023.
- [18] H. Ha and S. Song, "Flingbot: The unreasonable effectiveness of dynamic manipulation for cloth unfolding," in *Conference on Robot Learning*. PMLR, 2022, pp. 24–33.
- [19] X. Lin, Y. Wang, J. Olkin, and D. Held, "Softgym: Benchmarking deep reinforcement learning for deformable object manipulation," in *Conference on Robot Learning*. PMLR, 2021, pp. 432–448.
- [20] T. Mu, Z. Ling, F. Xiang, D. Yang, X. Li, S. Tao, Z. Huang, Z. Jia, and H. Su, "Maniskill: Generalizable manipulation skill benchmark with large-scale demonstrations," arXiv preprint arXiv:2107.14483, 2021.
- [21] K. Pearson, "Vii. note on regression and inheritance in the case of two parents," *proceedings of the royal society of London*, vol. 58, no. 347-352, pp. 240–242, 1895.
- [22] D. Kalashnikov, J. Varley, Y. Chebotar, B. Swanson, R. Jonschkowski, C. Finn, S. Levine, and K. Hausman, "Mt-opt: Continuous multi-task robotic reinforcement learning at scale," *arXiv preprint* arXiv:2104.08212, 2021.
- [23] K. Ehsani, T. Gupta, R. Hendrix, J. Salvador, L. Weihs, K.-H. Zeng, K. P. Singh, Y. Kim, W. Han, A. Herrasti, *et al.*, "Imitating shortest paths in simulation enables effective navigation and manipulation in the real world," *arXiv preprint arXiv:2312.02976*, 2023.
- [24] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar, "Roboagent: Towards sample efficient robot manipulation with semantic augmentations and action chunking," 2023.
- [25] S. Liu, L. Wu, B. Li, H. Tan, H. Chen, Z. Wang, K. Xu, H. Su, and J. Zhu, "Rdt-1b: a diffusion foundation model for bimanual manipulation," arXiv preprint arXiv:2410.07864, 2024.
- [26] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen, *et al.*, "Palm 2 technical report," *arXiv preprint arXiv:2305.10403*, 2023.
- [27] S. Ye, J. Jang, B. Jeon, S. Joo, J. Yang, B. Peng, A. Mandlekar, R. Tan, Y.-W. Chao, B. Y. Lin, *et al.*, "Latent action pretraining from videos," *arXiv preprint arXiv:2410.11758*, 2024.
- [28] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis, *et al.*, "Droid: A large-scale in-the-wild robot manipulation dataset," *arXiv preprint arXiv:2403.12945*, 2024.
- [29] Z. Zhou, P. Atreya, A. Lee, H. Walke, O. Mees, and S. Levine, "Autonomous improvement of instruction following skills via foundation models," arXiv preprint arXiv:2407.20635, 2024.
- [30] L. Fu, H. Huang, G. Datta, L. Y. Chen, W. C.-H. Panitch, F. Liu,

H. Li, and K. Goldberg, "In-context imitation learning via next-token prediction," *arXiv preprint arXiv:2408.15980*, 2024.

- [31] J. Yang, C. Glossop, A. Bhorkar, D. Shah, Q. Vuong, C. Finn, D. Sadigh, and S. Levine, "Pushing the limits of crossembodiment learning for manipulation and navigation," *arXiv preprint arXiv:2402.19432*, 2024.
- [32] K. Van Wyk, J. Falco, and E. Messina, "Robotic grasping and manipulation competition: Future tasks to support the development of assembly robotics," in *Robotic Grasping and Manipulation: First Robotic Grasping and Manipulation Challenge, RGMC 2016, Held in Conjunction with IROS 2016, Daejeon, South Korea, October 10–12, 2016, Revised Papers 1.* Springer, 2018, pp. 190–200.
- [33] B. Yang, J. Zhang, D. Jayaraman, and S. Levine, "Replab: A reproducible low-cost arm benchmark platform for robotic learning," *ICRA*, 2019.
- [34] M. Heo, Y. Lee, D. Lee, and J. J. Lim, "Furniturebench: Reproducible real-world benchmark for long-horizon complex manipulation," *arXiv* preprint arXiv:2305.12821, 2023.
- [35] J. Luo, C. Xu, F. Liu, L. Tan, Z. Lin, J. Wu, P. Abbeel, and S. Levine, "Fmb: a functional manipulation benchmark for generalizable robotic learning," *The International Journal of Robotics Research*, p. 02783649241276017, 2023.
- [36] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar, "Benchmarking in manipulation research: The ycb object and model set and benchmarking protocols," *arXiv preprint arXiv:1502.03143*, 2015.
- [37] S. Tyree, J. Tremblay, T. To, J. Cheng, T. Mosier, J. Smith, and S. Birchfield, "6-dof pose estimation of household objects for robotic manipulation: An accessible dataset and benchmark," in 2022 *IEEE/RSJ International Conference on Intelligent Robots and Systems* (*IROS*). IEEE, 2022, pp. 13 081–13 088.
- [38] J. Leitner, A. W. Tow, N. Sünderhauf, J. E. Dean, J. W. Durham, M. Cooper, M. Eich, C. Lehnert, R. Mangels, C. McCool, *et al.*, "The acrv picking benchmark: A robotic shelf picking benchmark to foster reproducible research," in 2017 IEEE international conference on robotics and automation (ICRA). IEEE, 2017, pp. 4705–4712.
- [39] L. Paull, J. Tani, H. Ahn, J. Alonso-Mora, L. Carlone, M. Cap,

Y. F. Chen, C. Choi, J. Dusek, Y. Fang, *et al.*, "Duckietown: an open, inexpensive and flexible platform for autonomy education and research," in 2017 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2017, pp. 1497–1504.

- [40] D. Pickem, P. Glotfelter, L. Wang, M. Mote, A. Ames, E. Feron, and M. Egerstedt, "The robotarium: A remotely accessible swarm robotics research testbed," in 2017 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2017, pp. 1699–1706.
- [41] E. Krotkov, D. Hackett, L. Jackel, M. Perschbacher, J. Pippine, J. Strauss, G. Pratt, and C. Orlowski, "The darpa robotics challenge finals: Results and perspectives," *The DARPA robotics challenge finals: Humanoid robots to the rescue*, pp. 1–26, 2018.
- [42] N. Correll, K. E. Bekris, D. Berenson, O. Brock, A. Causo, K. Hauser, K. Okada, A. Rodriguez, J. M. Romano, and P. R. Wurman, "Analysis and observations from the first amazon picking challenge," *IEEE Transactions on Automation Science and Engineering*, vol. 15, no. 1, pp. 172–188, 2016.
- [43] H. Kitano, M. Asada, Y. Kuniyoshi, I. Noda, and E. Osawa, "Robocup: The robot world cup initiative," in *Proceedings of the first international* conference on Autonomous agents, 1997, pp. 340–347.
- [44] Earth Rover Challenge Team, "The Earth Rover Challenge," 2025, accessed: 2025-02-01. [Online]. Available: https://sites.google.com/ view/the-earth-rover-challenge
- [45] G. Zhou, V. Dean, M. K. Srirama, A. Rajeswaran, J. Pari, K. Hatch, A. Jain, T. Yu, P. Abbeel, L. Pinto, *et al.*, "Train offline, test online: A real robot learning benchmark," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 9197–9203.
- [46] S. Yenamandra, A. Ramachandran, K. Yadav, A. Wang, M. Khanna, T. Gervet, T.-Y. Yang, V. Jain, A. W. Clegg, J. Turner, *et al.*, "Homerobot: Open-vocabulary mobile manipulation," *arXiv preprint arXiv:2306.11565*, 2023.
- [47] S. Bauer, M. Wüthrich, F. Widmaier, A. Buchholz, S. Stark, A. Goyal, T. Steinbrenner, J. Akpo, S. Joshi, V. Berenz, *et al.*, "Real robot challenge: A robotics competition in the cloud," in *NeurIPS 2021 Competitions and Demonstrations Track.* PMLR, 2022, pp. 190–204.
- [48] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas,

D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, *et al.*, "Deepmind control suite," *arXiv preprint arXiv:1801.00690*, 2018.

- [49] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison, "Rlbench: The robot learning benchmark & learning environment," *IEEE Robotics* and Automation Letters, vol. 5, no. 2, pp. 3019–3026, 2020.
- [50] Y. Lee, E. S. Hu, and J. J. Lim, "IKEA furniture assembly environment for long-horizon complex manipulation tasks," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021. [Online]. Available: https://clvrai.com/furniture
- [51] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone, "Libero: Benchmarking knowledge transfer for lifelong robot learning," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [52] S. Nasiriany, A. Maddukuri, L. Zhang, A. Parikh, A. Lo, A. Joshi, A. Mandlekar, and Y. Zhu, "Robocasa: Large-scale simulation of everyday tasks for generalist robots," *arXiv preprint arXiv:2406.02523*, 2024.
- [53] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard, "Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks," *IEEE Robotics and Automation Letters (RA-L)*, vol. 7, no. 3, pp. 7327–7334, 2022.
- [54] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, *et al.*, "Isaac gym: High performance gpu-based physics simulation for robot learning," *arXiv preprint arXiv:2108.10470*, 2021.
- [55] S. Tao, F. Xiang, A. Shukla, Y. Qin, X. Hinrichsen, X. Yuan, C. Bao, X. Lin, Y. Liu, T.-k. Chan, *et al.*, "Maniskill3: Gpu parallelized robotics simulation and rendering for generalizable embodied ai," *arXiv preprint arXiv:2410.00425*, 2024.
- [56] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, M. Deitke, K. Ehsani, D. Gordon, Y. Zhu, *et al.*, "Ai2-thor: An interactive 3d environment for visual ai," *arXiv preprint arXiv:1712.05474*, 2017.
- [57] X. Puig, E. Undersander, A. Szot, M. D. Cote, T.-Y. Yang, R. Partsey, R. Desai, A. W. Clegg, M. Hlavac, S. Y. Min, *et al.*, "Habitat 3.0: A co-habitat for humans, avatars and robots," *arXiv preprint arXiv:2310.13724*, 2023.
- [58] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and

S. Levine, "Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning," in *Conference on robot learning*. PMLR, 2020, pp. 1094–1100.

- [59] O. Ahmed, F. Träuble, A. Goyal, A. Neitz, Y. Bengio, B. Schölkopf, M. Wüthrich, and S. Bauer, "Causalworld: A robotic manipulation benchmark for causal structure and transfer learning," *arXiv preprint* arXiv:2010.04296, 2020.
- [60] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine, W. Ai, B. Martinez, *et al.*, "Behavior-1k: A human-centered, embodied ai benchmark with 1,000 everyday activities and realistic simulation," *arXiv preprint arXiv:2403.09227*, 2024.
- [61] C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. Vainio, C. Gokmen, G. Dharan, T. Jain, *et al.*, "igibson 2.0: Objectcentric simulation for robot learning of everyday household tasks," *arXiv preprint arXiv:2108.03272*, 2021.
- [62] A. Mandlekar, S. Nasiriany, B. Wen, I. Akinola, Y. Narang, L. Fan, Y. Zhu, and D. Fox, "Mimicgen: A data generation system for scalable robot learning using human demonstrations," *arXiv preprint* arXiv:2310.17596, 2023.
- [63] O. Mees, L. Hermann, and W. Burgard, "What matters in language conditioned robotic imitation learning over unstructured data," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 11205–11212, 2022.
- [64] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in 2012 IEEE/RSJ international conference on intelligent robots and systems. IEEE, 2012, pp. 5026–5033.
- [65] A. Juliani, "Unity: A general platform for intelligent agents," arXiv preprint arXiv:1809.02627, 2018.
- [66] E. Coumans, "Bullet physics simulation," in ACM SIGGRAPH 2015 Courses, 2015, p. 1.
- [67] J. Lee, M. X. Grey, S. Ha, T. Kunz, S. Jain, Y. Ye, S. S. Srinivasa, M. Stilman, and C. Karen Liu, "Dart: Dynamic animation and robotics toolkit," *The Journal of Open Source Software*, vol. 3, no. 22, p. 500, 2018.
- [68] NVIDIA, "Physx," 2020. [Online]. Available: https://developer.nvidia. com/physx-sdk

- [69] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, et al., "Sapien: A simulated part-based interactive environment," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11097–11107.
- [70] G. Authors, "Genesis: A universal and generative physics engine for robotics and beyond," December 2024. [Online]. Available: https://github.com/Genesis-Embodied-AI/Genesis
- [71] P. Huang, X. Zhang, Z. Cao, S. Liu, M. Xu, W. Ding, J. Francis, B. Chen, and D. Zhao, "What went wrong? closing the sim-to-real gap via differentiable causal discovery," in *Conference on Robot Learning*. PMLR, 2023, pp. 734–760.
- [72] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE, 2017, pp. 23–30.
- [73] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. VanderBilt, M. Wallingford, *et al.*, "Robothor: An open simulation-to-real embodied ai platform," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3164–3174.
- [74] J. Zhang, L. Tai, P. Yun, Y. Xiong, M. Liu, J. Boedecker, and W. Burgard, "Vr-goggles for robots: Real-to-sim domain adaptation for visual control," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 1148–1155, 2019.
- [75] Z. Li, T.-W. Yu, S. Sang, S. Wang, M. Song, Y. Liu, Y.-Y. Yeh, R. Zhu, N. Gundavarapu, J. Shi, *et al.*, "Openrooms: An open framework for photorealistic indoor scene datasets," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 7190–7199.
- [76] M. Ahn, D. Dwibedi, C. Finn, M. G. Arenas, K. Gopalakrishnan, K. Hausman, B. Ichter, A. Irpan, N. Joshi, R. Julian, S. Kirmani, I. Leal, E. Lee, S. Levine, Y. Lu, S. Maddineni, K. Rao, D. Sadigh, P. Sanketi, P. Sermanet, Q. Vuong, S. Welker, F. Xia, T. Xiao, P. Xu, S. Xu, and Z. Xu, "Autort: Embodied foundation models for large scale orchestration of robotic agents," 2024.
- [77] L. Pinto and A. Gupta, "Supersizing self-supervision: Learning to

grasp from 50k tries and 700 robot hours," in 2016 IEEE international conference on robotics and automation (ICRA). IEEE, 2016, pp. 3406–3413.

- [78] A. S. Chen, H. Nam, S. Nair, and C. Finn, "Batch exploration with examples for scalable robotic reinforcement learning," *IEEE Robotics* and Automation Letters, vol. 6, no. 3, pp. 4401–4408, 2021.
- [79] T. Lampe, A. Abdolmaleki, S. Bechtle, S. H. Huang, J. T. Springenberg, M. Bloesch, O. Groth, R. Hafner, T. Hertweck, M. Neunert, *et al.*, "Mastering stacking of diverse shapes with large-scale iterative reinforcement learning on real robots," *arXiv preprint arXiv:2312.11374*, 2023.
- [80] J. M. Beer, A. D. Fisk, and W. A. Rogers, "Toward a framework for levels of robot autonomy in human-robot interaction," *Journal of human-robot interaction*, vol. 3, no. 2, p. 74, 2014.
- [81] D. B. D'Ambrosio, S. Abeyruwan, L. Graesser, A. Iscen, H. B. Amor, A. Bewley, B. J. Reed, K. Reymann, L. Takayama, Y. Tassa, *et al.*, "Achieving human level competitive robot table tennis," *arXiv preprint arXiv:2408.03906*, 2024.
- [82] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, *et al.*, "Octo: An open-source generalist robot policy," *arXiv preprint arXiv:2405.12213*, 2024.
- [83] A. Z. Ren, "Open-pi-zero: An open-source hardware project," https: //github.com/allenzren/open-pi-zero, 2024, accessed: 2025-01-31.
- [84] S. Belkhale and D. Sadigh, "Minivla: A better vla with a smaller footprint," 2024. [Online]. Available: https://github.com/ Stanford-ILIAD/openvla-mini
- [85] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine, "Zero-shot robotic manipulation with pretrained imageediting diffusion models," *arXiv preprint arXiv:2310.10639*, 2023.
- [86] M. Nakamoto, Y. Zhai, A. Singh, M. S. Mark, Y. Ma, C. Finn, A. Kumar, and S. Levine, "Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning," arXiv preprint arXiv:2303.05479, 2023.
- [87] Z. Zhou, A. Peng, Q. Li, S. Levine, and A. Kumar, "Efficient online reinforcement learning fine-tuning need not retain offline data," arXiv preprint arXiv:2412.07762, 2024.