

Robot Learning as an Empirical Science: Best Practices for Policy Evaluation

Hadas Kress-Gazit^{1,2}, Kunimatsu Hashimoto², Naveen Kuppuswamy², Paarth Shah²,
Phoebe Horgan², Gordon Richardson², Siyuan Feng², Benjamin Burchfiel²

¹Cornell University and ²Toyota Research Institute

Corresponding author: hadaskg@cornell.edu

Abstract: The robot learning community has made great strides in recent years, proposing new architectures and showcasing impressive new capabilities; however, the dominant metric used in the literature, especially for physical experiments, is “success rate”, i.e. the percentage of runs that were successful. Furthermore, it is common for papers to report this number with little to no information regarding the number of runs, the initial conditions, and the success criteria, little to no narrative description of the behaviors and failures observed, and little to no statistical analysis of the findings. In this paper we argue that to move the field forward, researchers should provide a nuanced evaluation of their methods, especially when evaluating and comparing learned policies on physical robots. To do so, we propose best practices for future evaluations: explicitly reporting the experimental conditions, evaluating several metrics designed to complement success rate, conducting statistical analysis, and adding a qualitative description of failures modes. We illustrate these through an evaluation on physical robots of several learned policies for manipulation tasks.

Keywords: Evaluation, Best practices, Metrics

1 Introduction

Recent years have seen significant advancements in machine learning, with successful deployments of machine learning models in the wild now becoming commonplace [1, 2, 3]. Robotics has also undergone significant changes with increasing adoption of data-driven machine-learning methods. These range from reinforcement learning (RL) [4], to deep RL [5] to the recent emergence of foundation models - general-purpose perception-action models [6, 7, 8] trained on large and diverse datasets and capable of performing in myriad in-the-wild domains.

As the field progresses, complex behaviors have been shown both in simulation and with physical hardware. However, while papers typically describe their architectural designs and training paradigms in depth, their evaluation criteria and process are often sparse in details that would help move the field forward. Specifically, the evaluation often solely focuses on “success rate” –the percentage of autonomous runs that were successful– with little description of the experimental conditions, number of evaluations, success criteria, performance, failure modes, and typically without any statistical analysis. This lack of detail and nuance makes it difficult to assess the true state of the field and impacts two communities of researchers; those who develop learning algorithms, and those who wish to use them. The former because it is not clear what the exact state of the art is and what the fruitful research directions are. The latter, who may wish to use policies as a black-box in some larger system, because they do not have a clear understanding of possible failure modes and under what conditions the algorithms were evaluated.

Contribution: We propose best practices for policy evaluation to improve the science of robot learning. These include suggestions for the experimental setup, different metrics, and the analysis of the results. While we focus on evaluation with physical robots, these best practices extend well to

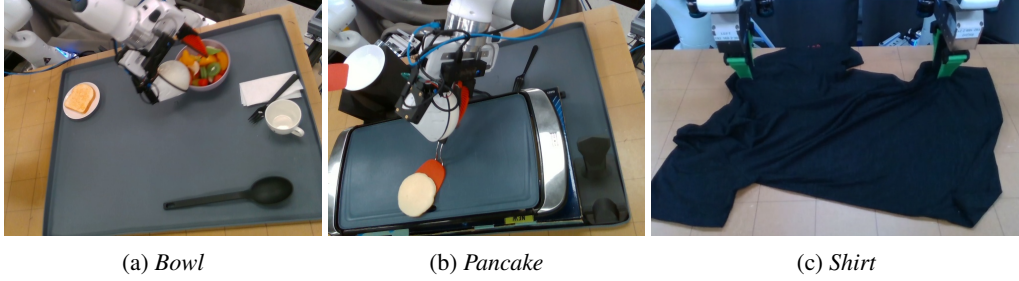


Figure 1: Tasks and robots

simulation. We illustrate our points through several examples of manipulation tasks performed by physical robots.

Data used in this paper: We illustrate our proposed best practices using data collected from physical robot evaluation runs. For each skill, we train several behavior cloning (BC) policies; all policies use the same set of human teleoperated demonstrations, but differ in architecture, observation space, and hyper parameters. While this paper considers a set of learned robot policies as case studies, we treat their implementation details as black-boxes.

- **Push Bowl** (*Bowl*, Fig. 1a, 6 policies): Starting with the end effector is in the air, a Franka Emika Panda arm pushes a bowl filled with (fake) fruit to a specific area on the table.
- **Flip and Serve Pancake** (*Pancake*, Fig. 1b, 3 policies): Two Franka Emika Panda arms perform the task together; first each arm grasps a spatula, then the right arm flips the pancake, then the left arm lifts the pancake and places it on a plate.
- **Fold Shirt** (*Shirt*, Fig. 1c, 2 policies): Two Franka Emika Panda arms perform the task together; they grasp a T-shirt, folds it three times, and center it on the table.

1.1 Related Work

Experiment design and analysis in robotics and machine learning: As disciplines mature, they begin forming best practices for evaluation and reporting, to ensure further progress. Such guidelines help communicate what information is useful to the community and set expectations, which is especially useful for people entering the field. Examples for best practices in fields that are close to robot learning include a primer on conducting experiments in human-robot interaction [9], and statistical analysis for evaluation of deep reinforcement learning [10, 11].

Metrics in robotics and related fields: The most common metric used in the robot learning literature is “success rate” (e.g. [8, 12, 13]). This is true both for simulation and physical experiments; however due to the fact that in simulation we have the full state of the world, there are more explicitly defined success metrics for simulated tasks. One common metric is distance from goal, especially in pick and place type tasks [14, 15, 16]. For physical experiments, especially work that evaluates different tasks, some works do not state explicitly what the criteria is, and others rely on human judgement (e.g. [17, 18]). Recent work has proposed a metric for behavior entropy, capturing a model’s ability to create diverse behaviors [19]. For more specific tasks, researchers have created quantitative metrics such as weight of material transferred during a scooping and a pouring task [20], time to fall, angle of rotation, torque applied for an in-hand manipulation task [21], and end state belonging to a predetermined set of states, evaluated manually [22]. Several works also capture sub-goal achievement, as we discuss in this paper. Some of these are in simulation where it is possible to automate sub-goal evaluation, especially if it is related to location or contact in space, and some manually evaluated for physical experiments (e.g. [23, 24]). More broadly, metrics have been proposed in adjacent engineering fields including: human-robot interaction studies [25], natural language processing [26], and computer vision [27, 28, 29, 30].

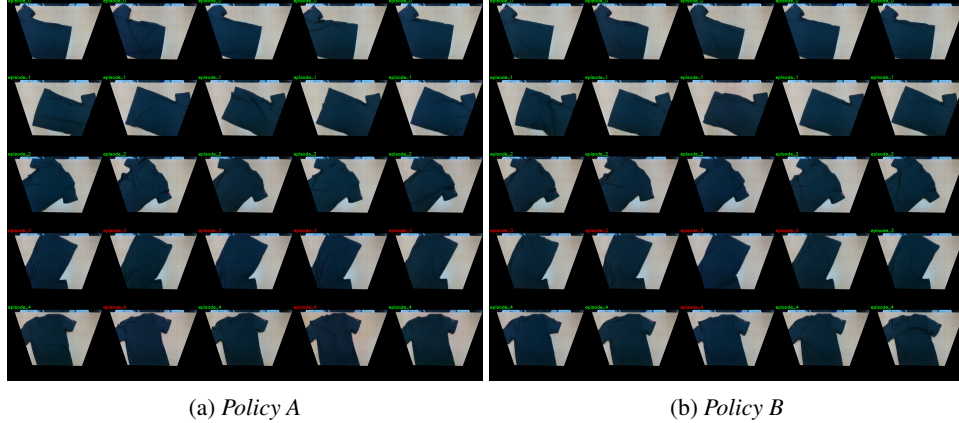


Figure 2: Initial conditions (ICs) for Example 1; Policy A on the left and Policy B on the right. Each row represents 5 evaluations with the same IC; the ICs are numbered 0-4 from the top row. The color of the text above each photo indicates the success; green is a successful rollout, red is a failure. We can see that the IC are visually consistent. Furthermore, for ICs 0,1,2 both policies always succeed, while for IC 3 they mostly fail.

2 Experiment Setup

Ultimately, the purpose of an experimental evaluation in the field of robot learning is to gain knowledge and insight about design choices of the learned policy (training, architecture, learning objective, hyper parameters, etc.). To gain the insight, policies would ideally be compared under identical conditions; however, as opposed to simulation, with a physical system it is impossible to replicate the exact same experimental conditions between evaluation runs. Furthermore, we know that the performance of learned policies is highly dependent on the experimental conditions. In this section, we discuss experiment best practices to mitigate this challenge and maximize the utility of valuable and expensive physical experiments. When employed, they create a (closer to) level playing field for the learned policies being compared, reduce confounding variables that decrease the signal-to-noise ratio of experiments, and allow for increased confidence in experimentally-drawn conclusions.

2.1 Success criteria

First and foremost, there must be a clear, detailed, and unambiguous definition of success. While this sounds obvious, many papers do not provide a definition of success. When success criteria are not explicit, evaluation may become biased; as an anecdote, we note a time when three of the authors watched a robot pour ice from a cup into a sink. The robot did pour the ice into the sink but then hesitated when lowering the cup, moving it for a few seconds above the table, ultimately placing it on the table and tipping it over. Two of the authors viewed the episode as a failure, while the third thought it was a success. This ambiguity allows inadvertent evaluation bias to distort the results.

2.2 Initial conditions

It is well known among robot learning practitioners that today’s learning-based robots are highly sensitive to their deployment environment. For non-interactive tasks, the most critical aspect of this are the initial conditions from which rollouts are performed - these include object types and locations in the environment, lighting conditions and camera locations, to name a few.

It is easy and commonplace to inadvertently change initial distributions during evaluation (e.g. object placements, environmental shifts such as sunlight, background changes) in ways that could significantly affect policy performance but are not immediately obvious to a person. Unfortunately, it is not common for empirical work to control for these details or explicitly describe the initial conditions used across evaluation runs.

Example 1 (Effect of initial condition). *Figure 2 shows the set of initial conditions (ICs) we used to evaluate policies A and B for shirt. While the overall success rate of policy A and B are similar (72% vs 80%), we can see that the success is highly dependent on the IC; if we were not careful in controlling for IC, we might arrive at a different conclusion. For example, if we used different ICs for each policy, and by chance use mainly ICs 1,2 for one policy and ICs 3,4 for the other, we may conclude that one policy is significantly more performant than the other.*

Example 2 (Detecting distribution shift by matching initial conditions). *Going back to Example 1, we performed an additional evaluation of the policies, using the same initial conditions, after a lab reorganization. We ran 8 evaluation runs of policy A, 2 each on IC 0,1,3,4 and got 0/8 (0%) successes. This is in contrast to the previous performance of A on the same initial conditions which was 13/20 (65%). By controlling for the initial conditions, we can detect a suspected distribution shift resulting in degradation of the policies.*

2.3 Experimental Process Best Practices

In many robot learning papers that describe new algorithmic approaches, such as new architectures, researchers perform ablation studies and the analysis of the new approach is in comparison to the ablations and other baseline approaches. For this comparison to be meaningful, and more importantly, to reduce as much as possible unintended bias in the evaluation, we recommend the following:

- Defining detailed success criteria ahead of time for overall success and semantic metrics, as described in Section 3.1. For example, “The robot picked up the cup with the ice, poured all the ice cubes into the sink, and placed the cup back on the table without tipping it over” is better than “The robot poured the ice into the sink.” We advocate for the behavior designer to be the person writing the criteria, but not the person evaluating - that way it is easier to detect ambiguous or incomplete descriptions.
- Reducing, as much as possible, the unintended variability in environmental conditions between different policy rollouts, especially in the initial conditions. This can be done by matching initial conditions using image overlays or markings in the scene, and ensuring policies are evaluated in the same session so lighting and other environmental conditions are more likely to be the same.
- A/B testing, i.e., interleaving policy rollouts when comparing different policies in a way that is blind to the evaluator. This means evaluating all of the policies within one session, as opposed to different policies in different sessions; this will mitigate unintended bias from the evaluator since they will not know which policy is running. See [31] for a compelling argument.
- Ensuring consistency across evaluators and separating the role of demonstrator and evaluator. People who work closely with the robot and gather demonstrations have a better understanding of how to set up the robot and environment for maximum success; separating the roles and ensuring the same person (or group of people) performs all the evaluations will create a more consistent assessment of the policies.

3 Evaluation metrics

In this section we describe several metrics that can provide different nuanced information regarding the behavior of a robot. We divide the discussion into two overall types of metrics: *semantic information*, and *performance metrics*; the former is a binary type of evaluation with Yes/No as the result, capturing “correctness” of the behavior. The latter provides a set of continuous numbers that can be thought of as a dense reward in a reinforcement learning setting, capturing the “quality” of the behavior. For both types of metrics, some are task-agnostic, for example trajectory smoothness, while other are task-specific; we argue that both are important for analyzing the behavior of learned models.

| Sub-goal | Policy A (Y/N) | Policy B (Y/N) | Policy C (Y/N) |
|-------------------------------|----------------|----------------|----------------|
| Overall success? | 15 / 3 | 11 / 6 | 4 / 19 |
| Robot collided with anything? | 0 / 18 | 2 / 15 | 0 / 23 |
| Right arm picked up spatula? | 18 / 0 | 17 / 0 | 23 / 0 |
| Left arm picked up spatula? | 18 / 0 | 17 / 0 | 23 / 0 |
| Robot flipped pancake? | 18 / 0 | 16 / 1 | 23 / 0 |
| Robot picked up pancake? | 15 / 3 | 12 / 5 | 5 / 18 |

Table 1: Partial rubric for *Pancake* evaluation of three policies. Of note is Policy C; while the overall success rate is low (17%), it was able to complete part of the task (picking up the spatulas and flipping the pancake) 100% of the time (as denoted in bold). Looking only at overall success would deem Policy C to be an unsuccessful policy; however, that is not the full picture, as it is able to accomplish more than half the task consistently.

3.1 Semantic information

Capturing a notion of “success” or “failure”, these metrics are designed as Yes/No questions; they include success rate (percentage of “Yes”), subgoal completion, and failures modes. While these metrics are task specific (even for success rate, one needs to define what success in the task means), we distinguish between metrics that are descriptive (Section 3.1.1) and metrics that are computable, for example those based on logic (Section 3.1.2).

3.1.1 Rubrics

To measure semantic task progress, which provides a more fine-grained signal for evaluating a policy and allows the evaluator to gather quantitative information regarding failures, we recommend explicitly defining a rubric for each task. Some of the rubric items can be task agnostic, for example “did the robot exhibit unexpected collisions”, but the majority will be task specific.

This rubric should be filled out by the evaluator during the evaluation. While this adds an extra burden on the evaluator, in practice, when running evaluations on physical robots the evaluator is there making sure the experiment is progressing as intended; we argue that having them also fill out the rubric is worth the time as it can provide useful information.

Example 3 (Pancake partial success). *Table 1 displays our rubric that the evaluator filled out when rolling out three different policies. Just looking at the overall task success rate makes policy C seem worthless; however, we can see that all policies were able to grasp the spatulas and flip the pancake, but policies B and C struggled with picking the pancake up (last row of Table 1) and placing it on the plate (overall task success).*

Recently there has been work on automating the detection of such semantic information, and especially detection of failures, through the use of learned models [32] or visual-language models [33]. While this is a promising direction, as those papers mention, this is not yet a reliable way to automate the assessment of semantic information for policy rollouts on physical systems.

3.1.2 Signal Temporal Logic

In contrast to rubrics which require a person to evaluate the results, researchers can create functions that automatically produce granular evaluation data from rollout information. This can take the form of predicate functions which are functions from the state of the system to $\mathcal{B} = \{True, False\}$.

Temporal logics [34] provide syntax and semantics for formulas that are richer than what can be captured by a predicate; for several temporal logics there exists automated evaluation procedures to assess the truth value of formulas. We propose the use of Signal Temporal Logic (STL) and both its Boolean (“Yes/No”) and quantitative semantics, also known as the *robustness metric* [35, 36](Section 3.2.1) as a way to create rich, computable metrics. We provide the syntax and semantics of STL in the appendix.

STL formulas can capture properties as simple as “maintain a distance from an obstacle”, and as complex as “the right arm should be on the right side of the table unless the left arm dropped

the spatula, in which case the right arm should pick up the spatula within 10 seconds of it being dropped” and other behaviors that include timing, conditionals, conjunctions and disjunctions.

We propose the use of STL because given a formula and state information from the robot and the environment, both semantic and performance metrics can be computed automatically, reducing the burden on the evaluator. The challenge is providing the state information; for physical robots we can use proprioceptive information and classifiers, as we show in Example 4. For simulation evaluation, the world state is known, therefore STL metrics can shine.

3.2 Performance metrics

Ultimately, the purpose of learning policies is to create autonomous robots; in many domains, robots will work with and around people. When working with people, a robot being “correct” is not enough. The manner in which the robot behaves may impact the interaction and the acceptance of the robot [37, 38]. Here we suggest performance metrics for learned policies.

3.2.1 Signal Temporal Logic robustness

As mentioned above, STL has two types of semantics: Boolean and quantitative. Given an STL formula and a trajectory X_t , we can evaluate whether the trajectory evaluates to *True* or *False* (Boolean) but we can also calculate *how close* the formula is to satisfying or violating the formula (robustness). A positive robustness indicates that the formula is *True*, negative that it is *False*. Thus, STL can be used to determine both correctness and quality of policies. The definition of the qualitative semantics can be found in [36]; in the following we use the RTAMT [39, 40] library to calculate the robustness of STL formulas on rollout data.

Example 4 (push bowl). *Assume we prefer the robot not touch the table when it is pushing the bowl. To automatically evaluate this property we look at two signals; the z coordinate of the robot end effector, as calculated by the Panda using forward kinematics (where the z axis is perpendicular to the table, $z = 0$ is the table and z is positive for end effector positions above the table) , and contact, a signal that has values of around 20 for no contact and values greater than 500 when there is contact. We create this signal from a sensor attached to the end effector. The STL specification capturing the required property is*

$$\Box((\text{contact} > 100) \rightarrow (z > 0.25))$$

which means that at all times, whenever the robot is making contact, the value of z must be greater than 0.25. Since this is an implication, the specification is satisfied if either the left hand side is False (i.e. there is no contact), or the right hand side is True (i.e. z is always larger than 0.25)

Figure 3 shows the robustness score of all the trajectories for the 6 policies. We can see two modes in these figures - trajectories with robustness of around 80 and trajectories with robustness around 0. The former corresponds to trajectories that do not make contact; the robustness corresponds to how close the left hand side of the implication is to being False. If the typical non-contact value is around 20, and our threshold is 100, we have 80 in robustness - contact can increase by 80 before we need to satisfy the right hand side. Trajectories around 0 correspond to trajectories in which the robot made contact. Points with negative values represent trajectories for which the minimum z value during contact was less than 0.25, thus violating the specification.

When comparing the policies, we can see that the brown pentagram (policy ID 5) is the most prone to not making contact at all (most runs with robustness ~ 80), while the orange square (policy ID 1) is the most prone to violating the z -value requirement (most runs with negative robustness).

In our example we can see two types on behaviors because the signals themselves have a different range of values. This might not be the case in all situations, but in all cases the robustness will give insight regarding the success of a policy with respect to a possibly complex objective, and crucially, quantify how close the policy is to succeeding (or failing).

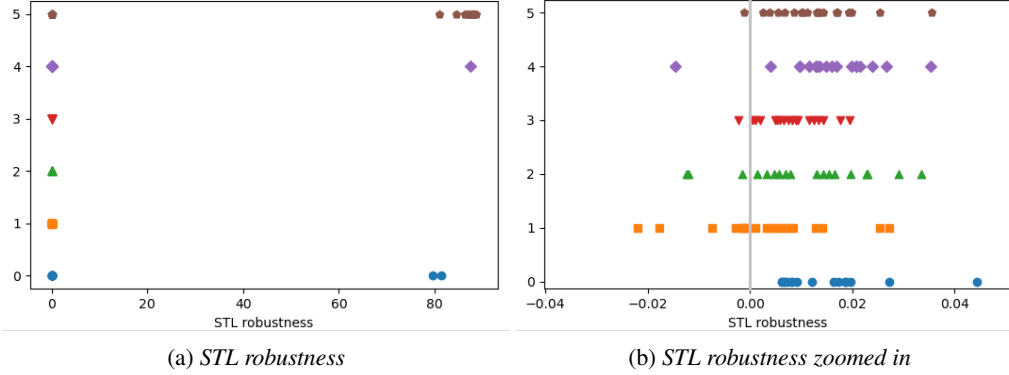


Figure 3: Robustness metric for Example 4 for the 6 policies (Y-axis). There are two types of behaviors; the points around 80 indicate that the robot did not make contact with the bowl, the points around 0 indicate that the robot did. In 3b we can see which trajectories violated the STL formula and by how much - points with negative robustness represent rollouts for which during contact the end effector z-coordinate was smaller than 0.25. We added a gray line at zero to help visualize.

3.2.2 Smoothness metrics

Smoothness of the robot trajectory may impact the human-robot interaction [37]; Refai et. al. [41] explored several possible smoothness metrics and conclude that the most appropriate one is SPectral ARC length (SPARC) [42, 43] computed over speed profiles. This metric is well suited for reaching tasks, and [42] discuss extensions to rhythmic movement. Looking at robot data, we can observe differences in the trajectory smoothness for different SPARC values; the more negative the value, the less smooth the trajectory is.

Example 5 (Bowl SPARC). Figure 4 shows the SPARC values for robot end effector trajectories before the first contact for all the rollouts in which the robot made contact with the bowl, i.e. we removed rollouts in which the robot does not touch the bowl. Looking, for example, at the policy represented by the green triangles (policy ID 2), we can see that the SPARC value ranges from -2.72 (smoothest) to -11.13 (least smooth). The rollouts that correspond to the two extreme values are both successful; however, the quality of the motion is different with the smooth trajectory making contact with the bowl almost instantly and completing the task in 4 seconds while the least smooth moves in a periodic manner above the bowl before making contact (air balling) and taking 15 seconds to finish the task.

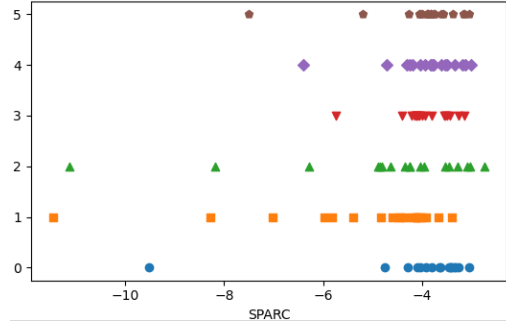


Figure 4: SPARC data for Bowl (Example 5). We consider rollouts in which the robot makes contact with the bowl and we calculate SPARC for the pre-contact trajectory. Each data point corresponds to a rollout. Different policies are color coded. More negative SPARC value corresponds to less smooth motion.

4 Analysis and Reporting

The final aspect of our recommendations for best practices addresses the reporting of the results and the insights gained from the experiment to the community. We break down the information into three categories: experimental parameters, statistical analysis, and failures. We provide an example

evaluation report in the appendix, Section 6.4; in this evaluation we also demonstrate how we can gain insights regarding differences between policies despite them having similar success rates.

4.1 Experiment parameters

We recommend that every experimental evaluation provide the following information: 1) a clear description of the semantic metrics, i.e. an explicit description of what is considered a success overall and in subgoals, 2) the number of evaluations performed across each condition and not just percentages (see Section 4.2), 3) the timing of the evaluations, i.e. were all the policies evaluated in an A/B fashion, were they evaluated in one session, were evaluations performed across different days/weeks, and 4) information regarding the initial conditions; visually as in Examples 1 and 8 or as in [44] (Fig.8), or in a narrative form. For the initial conditions we further recommend that researchers provide information regarding the relationship of the evaluation initial condition to the training initial conditions; were the evaluation ICs chosen as to try to capture in-distribution evaluation or were they chosen explicitly to evaluate out of distribution behavior?

4.2 Statistical Analysis

Providing a point estimate for success and performance metrics can be misleading [10]. Conducting a statistical analysis provides the community with a more nuanced understanding of the results. The statistical analysis can follow the frequentist approach, by providing interval estimates [10] or bounds [45], or the Bayesian approach [46] of estimating parameters of distributions; in the following we illustrate the Bayesian approach and the importance of communicating the experimental details for the analysis. [47] suggests best practices for reporting such statistical analysis.

Example 6 (Bayesian Analysis). *We estimate the success rate of policies A and B for the Pancake task. From Table 1, we see that the success rates of policy A and B are 83.3% (15/18) and 64.7% (11/17) respectively. In Bayesian analysis, we treat task success as a Bernoulli distribution with (unknown) parameter p , the probability of success, which is a random variable we try to estimate given the data. Given a prior on p , here a uniform distribution between 0 and 1, and the data, we can estimate the distribution of p ; the more the distributions of the policies overlap, the less confidence we have that one policy is better than the other. For this example, calculated with [48], we estimate a .11 probability that B actually performs better than A (Figure 6 in the appendix).*

Revisiting Example 6, we illustrate why it is important to explicitly state the number of evaluations and not just the percentage. Keeping the percentage the same but changing the number of evaluations results in very different conclusions (as shown in Figure 7 in the appendix). Reducing the number of evaluations makes any conclusions weaker, adding more evaluations makes them stronger. This is not surprising, but it emphasizes the need to provide additional information and analysis.

4.3 Failure modes

We recommend providing a detailed description of common and surprising failure modes encountered during the evaluation. This includes information regarding failure categories, narrative descriptions, frequency, and visual descriptions through images and videos (e.g. [44, 49]). This information is useful for several reasons; first, it sets expectations for researchers who are looking to incorporate learned models into their system as to what they can expect and what they need to be mindful of. Second, it provides a “gradient” for researchers in robot learning regarding the current state of the art and where more research is needed. Third, it provides a baseline for future progress; demonstrated change in the type or frequency of failures corresponds to progress in the field.

Example 7 (Failure description). *In the bowl example, we examine 6 different policies that were trained on identical training data. One common failure mode is for the robot to miss the bowl when it approaches to establish contact. From the rubric and the STL robustness (Figure 3) we see that policy 5 is the most prone to this error with 9 failures, while policies 1,2, and 3 did not exhibit this failure during evaluation. This type of granular analysis can motivate additional evaluation or insight into architectural, pretraining, or other design decisions.*

5 Discussion

In the following we offer some additional thoughts regarding the empirical nature of robot learning.

Purpose of using different metrics: Current learning techniques show a lot of promise for enabling robots to do tasks that were once considered impossible; however, they are not yet sufficiently performant for real deployment. To close this gap, it is critical to track progress and identify areas of deficit that need to be addressed. Similarly to the development of nuanced metrics in computer vision, such an approach to robot policy evaluation will be critical to advance robotics as well. Given the diversity of the approaches and tasks researchers consider, some metrics might make more sense than others. Furthermore, some metrics are going to be task specific, especially sub-goal achievement and STL specifications. We are advocating for researchers to choose the set of metrics that is best suited for their approach, task, and domain, but that provide more nuanced information than only success rate.

Simulation: We focus on physical experiments; however, all the best practices discussed are applicable to evaluation in simulation. Some aspects of the evaluation become straight forward, for example ensuring identical initial condition and environmental state (lighting, friction, etc) and automating all the metrics, at the expense of needing to reason about the sim-to-real gap. Closing this gap is an active area of research (e.g. [50]).

Releasing evaluation data: The robot learning community has been active in publicly releasing training data and models (e.g. [13]); however, while the data is well suited for ingestion by a machine learning algorithm, it is typically difficult to process manually. For example, it is not always clear what different arrays represent or how to decode images into videos that can be viewed by people. Furthermore, it is not common for researchers to publicly release evaluation data. While we propose several metrics, many useful metrics will be application dependent. For example automation applications may prefer approaches that exhibit consistency while human-interaction settings favor smoothness of motion. While we encourage researchers in the field to look beyond success rate in their evaluations, we also strongly advocate for open release of evaluation rollout data to enable additional nuanced posthoc analysis by the community.

Closing thoughts: Evaluating physical systems is difficult. It requires controlling numerous potentially confounding variables, balancing sample-size with diversity of conditions and comparisons given a fixed resource budget, contending with equipment malfunction, and managing human subjectivity and bias. Despite these challenges, rigorous evaluation is absolutely critical for an empirically driven field, such as robot learning, to sustain lasting progress. In this work, we advocate for a set of best practices designed to mitigate the above issues. While not panacea, we believe widespread adoption of these practices will significantly reduce noise and in turn improve the pace of progress in the field.

6 Appendix

6.1 STL syntax and semantics

Signal temporal logic (STL) is defined over continuous-valued signals X_t . The base element in the logic is a predicate μ and its associated predicate function $h(X)$. The truth value of μ at time t is determined as:

$$\mu ::= \begin{cases} True & h(X_t) \geq 0. \\ False & h(X_t) < 0 \end{cases}$$

In physical experiments, X would typically be robot state, in simulation it can also include the state of environment such as object locations.

An STL formula is recursively defines as [36]:

$$\varphi := \mu \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \varphi_1 \mathcal{U}_{\mathcal{I}} \varphi_2$$

where φ, φ_1 and φ_2 are STL formulas and $\mathcal{I} = [a, b]$ is an interval over which the formula is considered, where $0 \leq a < b$. In practice, since we consider finite trajectories, $b < \text{inf}$.

The logic contains Boolean operators: \neg negation, \wedge conjunction (“and”), and using those we can also construct \vee disjunction (“or”), \rightarrow implication, and \leftrightarrow “if and only if”. In addition, STL contains temporal operators: “Until” \mathcal{U} - the formula $\varphi_1 \mathcal{U}_{\mathcal{I}} \varphi_2$ evaluates to *True* if φ_1 is *True* until φ_2 becomes *True* in interval \mathcal{I} . Additional temporal operators are constructed using the Boolean operators and \mathcal{U} : “Always” $\Box_{\mathcal{I}}\varphi$ where φ must be *True* throughout the interval \mathcal{I} for the formula to be *True* and “Eventually” $\Diamond_{\mathcal{I}}\varphi$ where φ must be *True* at some point during the interval \mathcal{I} for the formula to be *True*. The semantics of STL can be found in [35, 36]. Section 3.2.1 provides a concrete example for such a metric; the sign of the robustness metric corresponds to the truth value of the Boolean semantics.

6.2 Visualizing initial conditions - additional example

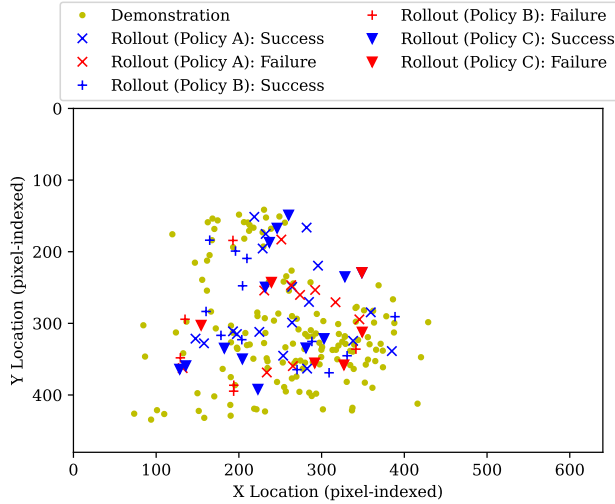


Figure 5: Initial location of the bowl in *Bowl*

Example 8 (Visualizing initial condition). In Figure 5 we visualize one aspect of the initial condition of a task - the initial location of the object that is to be manipulated. Figure 5 shows the initial locations of the bowl for *Bowl*. We used semantic segmentation to generate masks of the manipulant, then applied OpenCV [51]’s SimpleBlobDetector to filter out misclassified pixels, leaving a single blob of maximum size per image, and finally calculated the 2D centroid of the manipulant pixels in the image frame.

| | |
|-------------------------|-------------------|
| Number of demonstration | 154 |
| Type of Policy | Success / Failure |
| Policy A | 18 / 10 |
| Policy B | 12 / 6 |
| Policy C | 13 / 6 |

Table 2: Number of success and failure in *Bowl* for the policies shown in Figure 5

In Figure 5, we compare the initial conditions of the evaluations of three policies. Table 2 contains the respective success and failure numbers. While the success rate is similar, we can qualitatively see that the training distribution is not uniform and the evaluation conditions are not consistent between policies.

6.3 Bayesian analysis

The following figures are the distributions of p for Example 6. They illustrate how success percentage alone, which is fixed for all the figures, does not allow for in depth analysis of the effectiveness of one policy with respect to another.

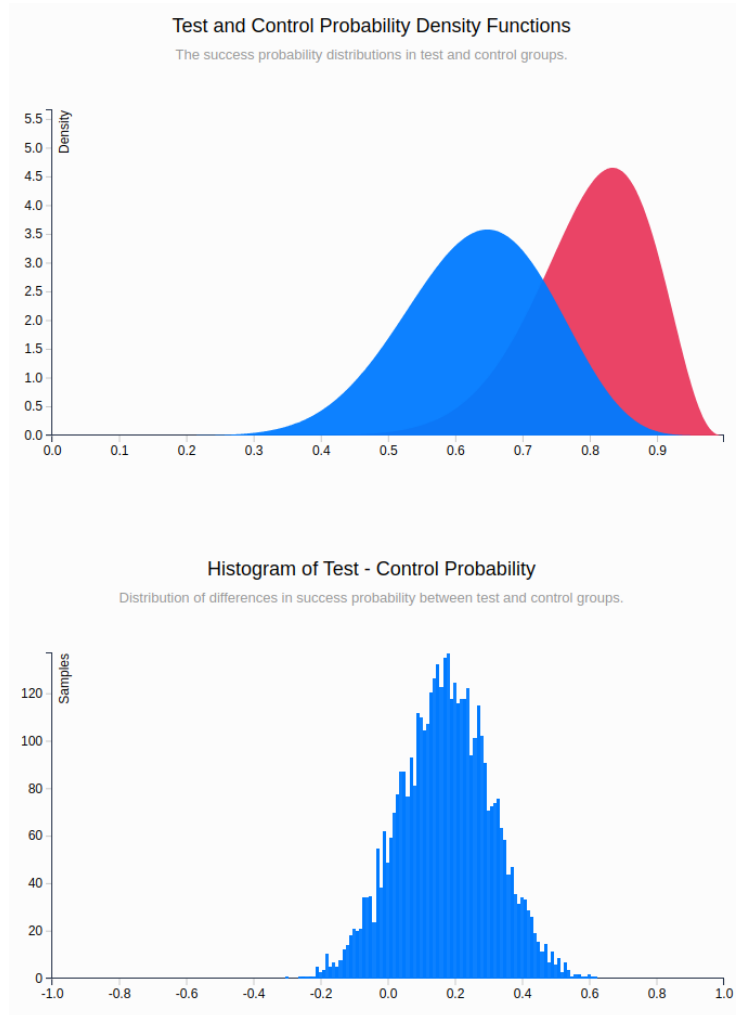


Figure 6: Estimated distributions of the success rate p for the *pancake* task (Example 6). Red is policy A, blue is policy B. The top shows the estimated distributions, the bottom a histogram of the difference between p . We can see that a difference of 0 is in the support of the histogram.

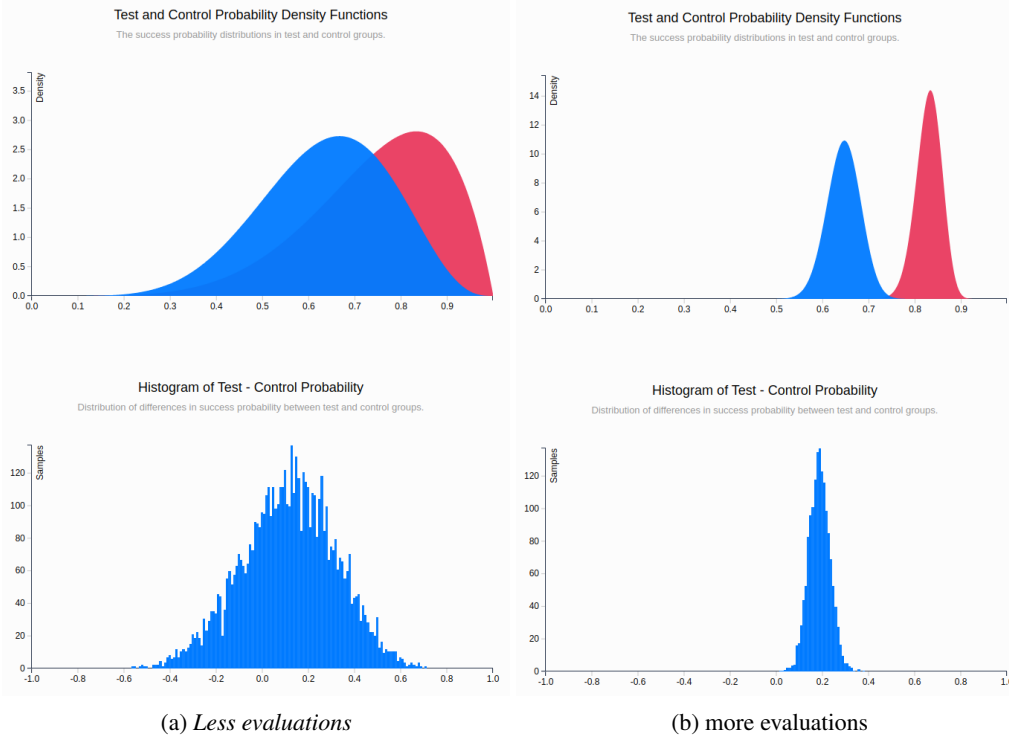


Figure 7: Estimated distributions of the success rate p for the *pancake* task (Example 6) when we use the same success rate but vary the number of evaluations. Red is policy A, blue is policy B. The top shows the estimated distributions, the bottom a histogram of the difference between p . In Figure 6 policy success/failure is 15/3 for policy A and 11/6 for policy B. Here, Figure (a) is 5/1 for policy A and 6/3 for policy B (slightly different success rate because it must be natural numbers), and Figure (b) is 150/30 for policy A and 110/60 for policy B. We can see that for (b), a difference of 0 between p is no longer in the support of the histogram, indicating that A outperforms B under the experimental conditions.

6.4 Example Evaluation Report

To illustrate the recommendations presented in this paper, we provide an example evaluation report for the comparison of two policies (policy A and policy B) learned through behavior cloning. The robot task is to pick up an energy bar and place it on a wooden tray. The robot is bimanual and can use either of its arms to perform the task. The details of the policies do not matter because we are not making a point about those policies but rather about how to report on their evaluation. If we were comparing the policies, in a discussion, we would draw conclusions regarding the policies based on the information presented below.

As shown below, even when we cannot draw conclusions regarding which policy is more successful, a detailed evaluation can lead to research insights and future directions.

6.4.1 Experiment parameters

The policies were evaluated in an A/B fashion (blind to the evaluator) across two days, the first day each policy was run 12 times, the second day (4 days later) each policy was run 8 times. The evaluator was not involved in the data collection or policy training. We created 10 different initial conditions (ICs) that were similar to initial conditions in the training set and we evaluated each policy twice on each IC (for a total of 20 evaluations per policy) by visually matching the initial conditions using an image overlay tool. The ICs are shown in Figure 8.

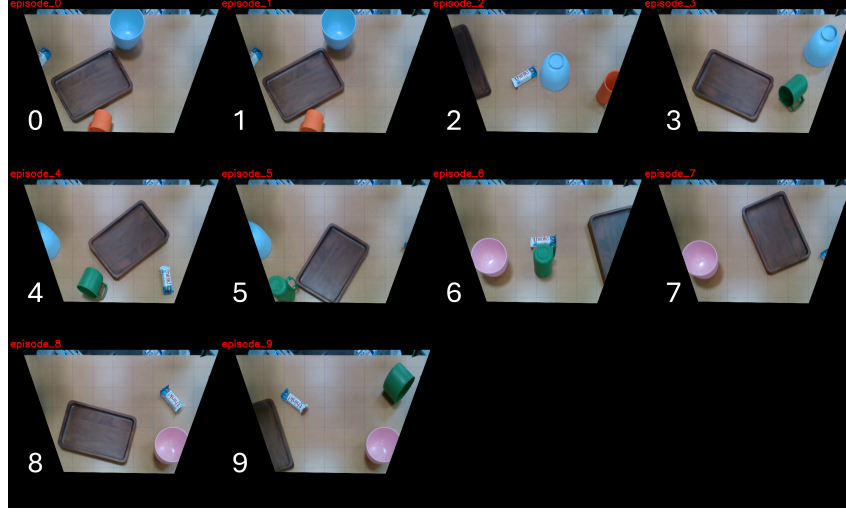


Figure 8: Initial location of the placing energy bar task

Success criteria: A run was considered successful if it ended with “Energy bar is on the wooden tray and the tray is on the table.” Note that we did not consider collisions with other items as failures.

6.4.2 Results

Success rate: Policy A succeeded in completing the task 13/20 times (0.65), while policy B succeeded 14/20 (0.7). Modeling the success of a policy as a binomial distribution, and assuming a uniform prior on the success probability, we cannot state that one policy outperforms the other. See Figure 11 for the posterior distributions.

Despite similar overall success, the policies performed differently and failed in different ways and on different ICs.

Performance: The following performance analysis contains only the *successful runs* for each of the policies. In Figures 9 and 10 we show the performance metrics for both policies and for both robot arms.

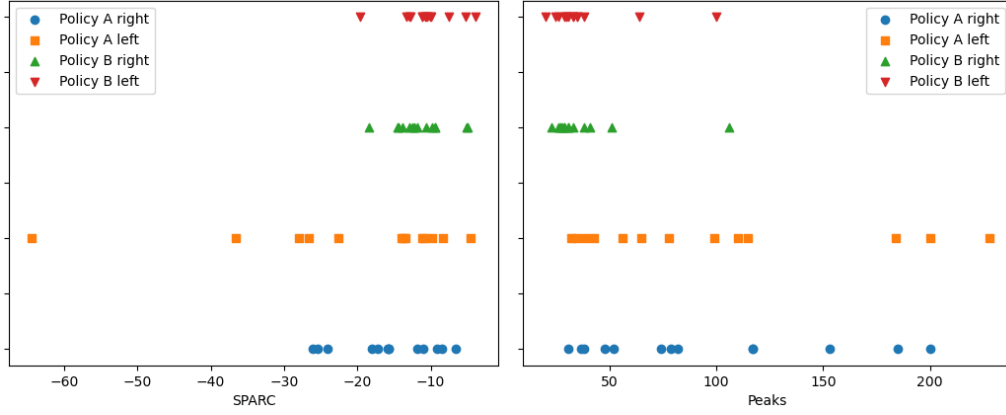
Smoothness of the trajectories: Figure 9 shows two metrics that relate to the smoothness of the trajectory. Figure 9a shows the SPARC metric—smaller absolute values correspond to smoother trajectories. Figure 9b shows the number of peaks in the velocity of the end effector; more peaks correspond to longer and less smooth trajectories.

From Figure 9, we can see that Policy B (triangles) produces smoother trajectories according to both smoothness metrics.

STL robustness example: We observed qualitatively that policy B’s grasps of the energy bar are more stable than those of policy A. To quantify this, we write an STL specification over two signals: z , the z -value of the end effector in a global frame centered on the table, and $gripper_diff$ which is the difference in gripper width between two consecutive time steps; when the gripper is closing, $gripper_diff > 0$ and when it is opening $gripper_diff < 0$. We plot the robustness of the following STL formula:

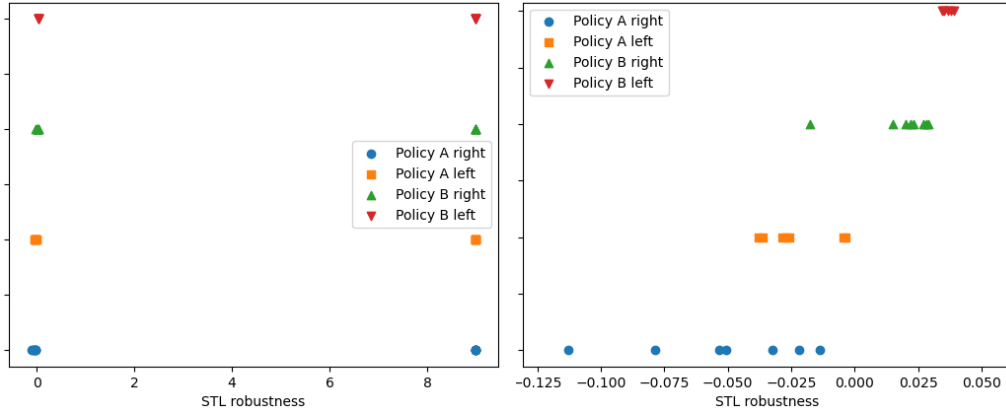
$$\Box((gripper_diff * 1000 > 9) \rightarrow (z < 0.25))$$

This formula states that always, if the difference between the gripper width in two consecutive time steps is greater than 0.009 (an observed change when the gripper is closing) then the end effector z -value should be less than 0.25. This formula is true and has positive robustness metric if either 1) the robot does not close its gripper (the left side of the implication is false) or 2) the end effector is close to the table when the gripper closes.



(a) *SPARC: smaller absolute value corresponds to smoother motion* (b) *Velocity peaks: more peaks corresponds to longer and less smooth motions*

Figure 9: Trajectory smoothness metrics for the example in Section 6.4 for the 2 policies and the two arms (Y-axis).



(a) *STL robustness*

(b) *STL robustness zoomed in*

Figure 10: STL Robustness metric for the example in Section 6.4 for the 2 policies and the two arms (Y-axis). There are two types of behaviors; the points around 9 indicate that the robot did not close its gripper, the points around 0 indicate that the robot did. In 10b we can see which trajectories violated the STL formula and by how much - points with negative robustness represent rollouts for which when the gripper was closing, the z-coordinate was larger than 0.25 indicating that the robot was attempting to grasp too high.

Figure 10 shows the robustness metric for both policies and both arms. We write the STL formula as $gripper_diff * 1000 > 9$ and not as $gripper_diff > 0.009$ so that we can observe the different modes in the robustness metric (gripper did not close resulting in a robustness of 9 vs gripper closed and a resulting robustness depending on the value of z). From the figure we can see that each policy and each arm had runs in which the gripper did not close. That is expected because each run had only one arm that grasped the energy bar. In the zoomed in figure (Figure 10b) we see that policy B (triangles) has higher robustness values, meaning that policy B tends to grasp lower than policy A. This is consistent with the failure modes we observed as described below.

Failure analysis: Both policies failed on IC 6; Policy A also failed on both runs of IC 4 and one run each of ICs 7,8,9, while policy B failed on both runs of ICs 2 and 9.

Policy A's main failure mode was the grasping of the energy bar, it either failed to grasp or dropped it prematurely (6/7 failures). Policy B, on the other hand, consistently picked up the energy bar (20/20

runs) but then either moved away from the tray or placed the energy bar in a different location (6/6 failures).

These different failures modalities may have different implications in a human-robot interaction scenario, so while the success rate is similar, one policy might be considered better by users of the robot, if this information is provided. Furthermore, for robot learning researchers, this information may provide insight on the algorithmic choices and future research directions.

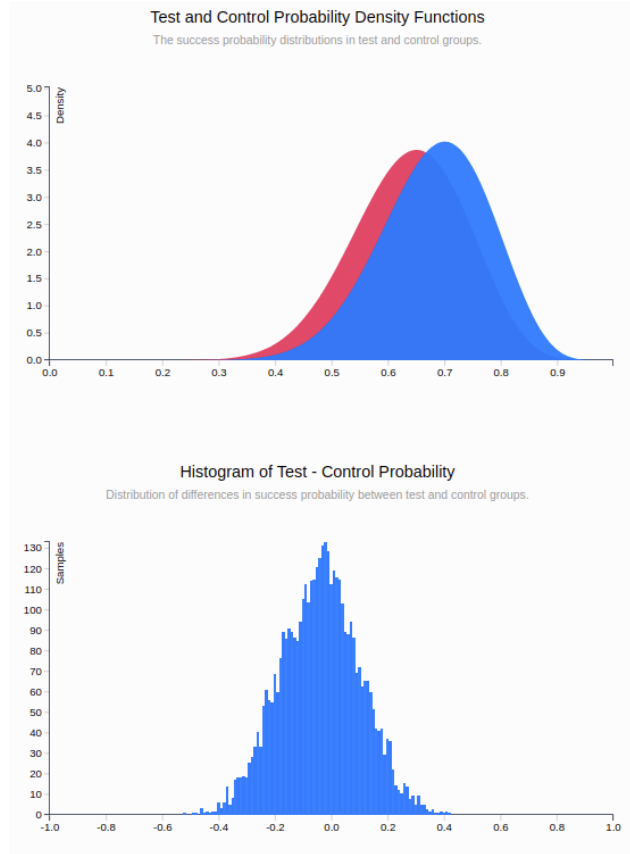


Figure 11: Estimated distributions of the success rate p for the energy bar task, calculated with [48]. Red is policy A, blue is policy B. The top shows the estimated distributions, the bottom a histogram of the difference between p . We can see that a difference of 0 is in the support of the histogram.

References

- [1] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment anything. *arXiv:2304.02643*, 2023.
- [2] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. URL <https://github.com/CompVis/latent-diffusion><https://arxiv.org/abs/2112.10752>.
- [3] OpenAI, :, J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, R. Avila, I. Babuschkin, S. Balaji, V. Balcom, P. Baltescu, H. Bao, M. Bavarian, J. Belgum, I. Bello, J. Berdine, G. Bernadett-Shapiro, C. Berner, L. Bogdonoff, O. Boiko, M. Boyd, A.-L. Brakman, G. Brockman, T. Brooks, M. Brundage, K. Button, T. Cai, R. Campbell, A. Cann, B. Carey, C. Carlson, R. Carmichael, B. Chan, C. Chang, F. Chantzis, D. Chen, S. Chen, R. Chen, J. Chen, M. Chen, B. Chess, C. Cho,

- C. Chu, H. W. Chung, D. Cummings, J. Currier, Y. Dai, C. Decareaux, T. Degry, N. Deutsch, D. Deville, A. Dhar, D. Dohan, S. Dowling, S. Dunning, A. Ecoffet, A. Eleti, T. Eloundou, D. Farhi, L. Fedus, N. Felix, S. P. Fishman, J. Forte, I. Fulford, L. Gao, E. Georges, C. Gibson, V. Goel, T. Gogineni, G. Goh, R. Gontijo-Lopes, J. Gordon, M. Grafstein, S. Gray, R. Greene, J. Gross, S. S. Gu, Y. Guo, C. Hallacy, J. Han, J. Harris, Y. He, M. Heaton, J. Heidecke, C. Hesse, A. Hickey, W. Hickey, P. Hoeschele, B. Houghton, K. Hsu, S. Hu, X. Hu, J. Huizinga, S. Jain, S. Jain, J. Jang, A. Jiang, R. Jiang, H. Jin, D. Jin, S. Jomoto, B. Jonn, H. Jun, T. Kaftan, Łukasz Kaiser, A. Kamali, I. Kanitscheider, N. S. Keskar, T. Khan, L. Kilpatrick, J. W. Kim, C. Kim, Y. Kim, H. Kirchner, J. Kiros, M. Knight, D. Kokotajlo, Łukasz Kondraciuk, A. Kondrich, A. Konstantinidis, K. Kosic, G. Krueger, V. Kuo, M. Lampe, I. Lan, T. Lee, J. Leike, J. Leung, D. Levy, C. M. Li, R. Lim, M. Lin, S. Lin, M. Litwin, T. Lopez, R. Lowe, P. Lue, A. Makanju, K. Malfacini, S. Manning, T. Markov, Y. Markovski, B. Martin, K. Mayer, A. Mayne, B. McGrew, S. M. McKinney, C. McLeavey, P. McMillan, J. McNeil, D. Medina, A. Mehta, J. Menick, L. Metz, A. Mishchenko, P. Mishkin, V. Monaco, E. Morikawa, D. Mossing, T. Mu, M. Murati, O. Murk, D. Mély, A. Nair, R. Nakano, R. Nayak, A. Neelakantan, R. Ngo, H. Noh, L. Ouyang, C. O’Keefe, J. Pachocki, A. Paino, J. Palermo, A. Pantuliano, G. Parascandolo, J. Parish, E. Parparita, A. Passos, M. Pavlov, A. Peng, A. Perelman, F. de Avila Belbute Peres, M. Petrov, H. P. de Oliveira Pinto, Michael, Pokorny, M. Pokrass, V. Pong, T. Powell, A. Power, B. Power, E. Proehl, R. Puri, A. Radford, J. Rae, A. Ramesh, C. Raymond, F. Real, K. Rimbach, C. Ross, B. Rotsted, H. Roussez, N. Ryder, M. Saltarelli, T. Sanders, S. Santurkar, G. Sastry, H. Schmidt, D. Schnurr, J. Schulman, D. Selsam, K. Sheppard, T. Sherbakov, J. Shieh, S. Shoker, P. Shyam, S. Sidor, E. Sigler, M. Simens, J. Sitkin, K. Slama, I. Sohl, B. Sokolowsky, Y. Song, N. Staudacher, F. P. Such, N. Summers, I. Sutskever, J. Tang, N. Tezak, M. Thompson, P. Tillet, A. Tootoonchian, E. Tseng, P. Tuggle, N. Turley, J. Tworek, J. F. C. Uribe, A. Vallone, A. Vijayvergiya, C. Voss, C. Wainwright, J. J. Wang, A. Wang, B. Wang, J. Ward, J. Wei, C. Weinmann, A. Welihinda, P. Welinder, J. Weng, L. Weng, M. Wiethoff, D. Willner, C. Winter, S. Wolrich, H. Wong, L. Workman, S. Wu, J. Wu, M. Wu, K. Xiao, T. Xu, S. Yoo, K. Yu, Q. Yuan, W. Zaremba, R. Zellers, C. Zhang, M. Zhang, S. Zhao, T. Zheng, J. Zhuang, W. Zhuk, and B. Zoph. Gpt-4 technical report, 2023.
- [4] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013. doi:10.1177/0278364913495721. URL <https://doi.org/10.1177/0278364913495721>.
- [5] C. Tang, B. Abbatematteo, J. Hu, R. Chandra, R. Martín-Martín, and P. Stone. Deep reinforcement learning for robotics: A survey of real-world successes, 2024. URL <https://arxiv.org/abs/2408.03539>.
- [6] 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, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source generalist robot policy. <https://octo-models.github.io>, 2023.
- [7] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran, V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In *arXiv preprint arXiv:2212.06817*, 2022.
- [8] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, P. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman, A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, L. Lee, T.-W. E. Lee, S. Levine, Y. Lu, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao, K. Reymann, M. Ryoo, G. Salazar, P. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut,

- H. Tran, V. Vanhoucke, Q. Vuong, A. Wahid, S. Welker, P. Wohlhart, J. Wu, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023.
- [9] G. Hoffman and X. Zhao. A primer for conducting experiments in human–robot interaction. *J. Hum.-Robot Interact.*, 10(1), oct 2020. doi:10.1145/3412374. URL <https://doi.org/10.1145/3412374>.
- [10] R. Agarwal, M. Schwarzer, P. S. Castro, A. C. Courville, and M. Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Information Processing Systems*, 34, 2021.
- [11] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger. Deep reinforcement learning that matters. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), Apr. 2018. doi:10.1609/aaai.v32i1.11694. URL <https://ojs.aaai.org/index.php/AAAI/article/view/11694>.
- [12] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar. RoboAgent: Generalization and Efficiency in Robot Manipulation via Semantic Augmentations and Action Chunking, 2023.
- [13] 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. Thompson, 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. Sundareshan, 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.
- [14] T. Yu, D. Quillen, Z. He, R. Julian, A. Narayan, H. Shively, A. Bellathur, K. Hausman, C. Finn, and S. Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning, 2021.
- [15] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning, 2019.
- [16] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. In *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018. doi:10.15607/RSS.2018.XIV.049.
- [17] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training, 2023.

- [18] Do as i can and not as i say: Grounding language in robotic affordances.
- [19] X. Jia, D. Blessing, X. Jiang, M. Reuss, A. Donat, R. Lioutikov, and G. Neumann. Towards diverse behaviors: A benchmark for imitation learning with human demonstrations. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=6pPYRXKPpw>.
- [20] G. Zhou, V. Dean, M. K. Srirama, A. Rajeswaran, J. Pari, K. Hatch, A. Jain, T. Yu, P. Abbeel, L. Pinto, C. Finn, and A. Gupta. Train offline, test online: A real robot learning benchmark, 2023.
- [21] H. Qi, A. Kumar, R. Calandra, Y. Ma, and J. Malik. In-Hand Object Rotation via Rapid Motor Adaptation. In *Conference on Robot Learning (CoRL)*, 2022.
- [22] K. Bousmalis, G. Vezzani, D. Rao, C. Devin, A. X. Lee, M. Bauza, T. Davchev, Y. Zhou, A. Gupta, A. Raju, A. Laurens, C. Fantacci, V. Dalibard, M. Zambelli, M. Martins, R. Pevceviciute, M. Blokzijl, M. Denil, N. Batchelor, T. Lampe, E. Parisotto, K. Žolna, S. Reed, S. G. Colmenarejo, J. Scholz, A. Abdolmaleki, O. Groth, J.-B. Regli, O. Sushkov, T. Rothörl, J. E. Chen, Y. Aytar, D. Barker, J. Ortiz, M. Riedmiller, J. T. Springenberg, R. Hadsell, F. Nori, and N. Heess. *RoboCat: A Self-Improving Foundation Agent for Robotic Manipulation*, 2023.
- [23] M. Heo, Y. Lee, D. Lee, and J. J. Lim. FurnitureBench: Reproducible Real-World Benchmark for Long-Horizon Complex Manipulation. In *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023. doi:10.15607/RSS.2023.XIX.041.
- [24] A. Ajay, S. Han, Y. Du, S. Li, A. Gupta, T. Jaakkola, J. Tenenbaum, L. Kaelbling, A. Srivastava, and P. Agrawal. *Compositional Foundation Models for Hierarchical Planning*, 2023.
- [25] A. Steinfeld, T. Fong, D. Kaber, M. Lewis, J. Scholtz, A. Schultz, and M. Goodrich. Common metrics for human-robot interaction. In *Proceedings of the 1st ACM SIGCHI/SIGART Conference on Human-Robot Interaction*, HRI '06, page 33–40, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595932941. doi:10.1145/1121241.1121249. URL <https://doi.org/10.1145/1121241.1121249>.
- [26] K. Blagec, G. Dorffner, M. Moradi, S. Ott, and M. Samwald. *A global analysis of metrics used for measuring performance in natural language processing*, 2022.
- [27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- [28] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- [29] G. Borgefors. Distance transformations in arbitrary dimensions. *Computer Vision, Graphics, and Image Processing*, 27(3):321–345, 1984. ISSN 0734-189X. doi:[https://doi.org/10.1016/0734-189X\(84\)90035-5](https://doi.org/10.1016/0734-189X(84)90035-5). URL <https://www.sciencedirect.com/science/article/pii/0734189X84900355>.
- [30] Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover’s distance as a metric for image retrieval. *International journal of computer vision*, 40:99–121, 2000.
- [31] V. Vanhoucke. I heart irreproducible research; better experimental protocols for real world research. URL <https://towardsdatascience.com/i-irreproducible-research-dbc48eb140ac>.
- [32] A. Inceoglu, E. E. Aksoy, and S. Sariel. Multimodal detection and identification of robot manipulation failures, 2023.

- [33] L. Guan, Y. Zhou, D. Liu, Y. Zha, H. B. Amor, and S. Kambhampati. "task success" is not enough: Investigating the use of video-language models as behavior critics for catching undesirable agent behaviors, 2024.
- [34] E. A. Emerson. Temporal and modal logic. In J. Van Leeuwen, editor, *Formal Models and Semantics*, Handbook of Theoretical Computer Science, pages 995–1072. Elsevier, Amsterdam, 1990. ISBN 978-0-444-88074-1. doi:<https://doi.org/10.1016/B978-0-444-88074-1.50021-4>. URL <https://www.sciencedirect.com/science/article/pii/B9780444880741500214>.
- [35] O. Maler and D. Nickovic. Monitoring temporal properties of continuous signals. In *Formal Techniques, Modelling and Analysis of Timed and Fault-Tolerant Systems*, pages 152–166. Springer, 2004.
- [36] A. Donzé and O. Maler. Robust satisfaction of temporal logic over real-valued signals. In K. Chatterjee and T. A. Henzinger, editors, *Formal Modeling and Analysis of Timed Systems*, pages 92–106, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg. ISBN 978-3-642-15297-9.
- [37] M. Rubagotti, I. Tusseyeva, S. Baltabayeva, D. Summers, and A. Sandygulova. Perceived safety in physical human–robot interaction—a survey. *Robotics and Autonomous Systems*, 151:104047, 2022. ISSN 0921-8890. doi:<https://doi.org/10.1016/j.robot.2022.104047>. URL <https://www.sciencedirect.com/science/article/pii/S0921889022000173>.
- [38] J. Urakami and K. Seaborn. Nonverbal cues in human–robot interaction: A communication studies perspective. *J. Hum.-Robot Interact.*, 12(2), mar 2023. doi:[10.1145/3570169](https://doi.org/10.1145/3570169). URL <https://doi.org/10.1145/3570169>.
- [39] D. Ničković and T. Yamaguchi. Rtamt: Online robustness monitors from stl. In *International Symposium on Automated Technology for Verification and Analysis*, pages 564–571. Springer, 2020.
- [40] RTAMT toolbox. URL <https://github.com/nickovic/rtamt>.
- [41] M. I. Mohamed Refai, M. Saes, B. L. Scheltinga, J. van Kordelaar, J. B. J. Bussmann, P. H. Veltink, J. H. Buurke, C. G. M. Meskers, E. E. H. van Wegen, G. Kwakkel, and B.-J. F. van Beijnum. Smoothness metrics for reaching performance after stroke. part 1: which one to choose? *Journal of NeuroEngineering and Rehabilitation*, 18(1):154, Oct 2021. ISSN 1743-0003. doi:[10.1186/s12984-021-00949-6](https://doi.org/10.1186/s12984-021-00949-6). URL <https://doi.org/10.1186/s12984-021-00949-6>.
- [42] S. Balasubramanian, A. Melendez-Calderon, A. Roby-Brami, and E. Burdet. On the analysis of movement smoothness. *Journal of NeuroEngineering and Rehabilitation*, 12(1):112, Dec 2015. ISSN 1743-0003. doi:[10.1186/s12984-015-0090-9](https://doi.org/10.1186/s12984-015-0090-9). URL <https://doi.org/10.1186/s12984-015-0090-9>.
- [43] Sparc. URL <https://github.com/siva82kb/SPARC>.
- [44] C. Chi, Z. Xu, C. Pan, E. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots, 2024.
- [45] J. A. Vincent, H. Nishimura, M. Itkina, P. Shah, M. Schwager, and T. Kollar. How generalizable is my behavior cloning policy? a statistical approach to trustworthy performance evaluation, 2024.
- [46] J. K. Kruschke. Bayesian estimation supersedes the t test. *Journal of Experimental Psychology: General*, 142(2):573–603, 2013. doi:[10.1037/a0029146](https://doi.org/10.1037/a0029146). URL <https://doi.org/10.1037/a0029146>.

- [47] J. K. Kruschke. Bayesian analysis reporting guidelines. *Nature Human Behaviour*, 5(10): 1282–1291, Oct 2021. ISSN 2397-3374. doi:[10.1038/s41562-021-01177-7](https://doi.org/10.1038/s41562-021-01177-7). URL <https://doi.org/10.1038/s41562-021-01177-7>.
- [48] Bayesian a/b test calculator. URL <https://making.lyst.com/bayesian-calculator/>.
- [49] P. Liu, Y. Orru, C. Paxton, N. M. M. Shafiullah, and L. Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024.
- [50] 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, 2024.
- [51] G. Bradski. The opencv library. *Dr. Dobb's Journal: Software Tools for the Professional Programmer*, 25(11):120–123, 2000.