

# Dexterous Manipulation from Images: Autonomous Real-World RL via Substep Guidance

Kelvin Xu<sup>\*1</sup> Zheyuan Hu<sup>\*1</sup> Ria Doshi<sup>1</sup> Aaron Rovinsky<sup>1</sup> Vikash Kumar<sup>2</sup> Abhishek Gupta<sup>3</sup> Sergey Levine<sup>1</sup>  
<sup>1</sup> UC Berkeley <sup>2</sup> Meta AI Research <sup>3</sup> University of Washington



Fig. 1: Filmstrip of the final learned brush skill. Our agent is able to learn to grasp, in-hand reorient, and brush a surface using a kitchen brush. After around  $\sim 18$  hours of unattended learning, our system successfully performs each of these sub-tasks with  $\geq 80\%$  success.

**Abstract**—Complex and contact-rich robotic manipulation tasks, particularly those that involve multi-fingered hands and underactuated object manipulation, present a significant challenge to any control method. Methods based on reinforcement learning offer an appealing choice for such settings, as they can enable robots to learn to delicately balance contact forces and dexterously reposition objects without strong modeling assumptions. However, running reinforcement learning on real-world dexterous manipulation systems often requires significant manual engineering. This negates the benefits of autonomous data collection and ease of use that reinforcement learning should in principle provide. In this paper, we describe a system for vision-based dexterous manipulation that provides a “programming-free” approach for users to define new tasks and enable robots with complex multi-fingered hands to learn to perform them through interaction. The core principle underlying our system is that, in a vision-based setting, users should be able to provide high-level intermediate supervision that circumvents challenges in teleoperation or kinesthetic teaching which allows a robot to not only learn a task efficiently but also to autonomously practice. Our system includes a framework for users to define a final task and intermediate sub-tasks with image examples, a reinforcement learning procedure that learns the task autonomously without interventions, and experimental results with a four-finger robotic hand learning multi-stage object manipulation tasks directly in the real world, without simulation, manual modeling, or reward engineering.

## I. INTRODUCTION

Control methods face significant challenges when dealing with robotic manipulation tasks that are complex and involve multiple contact points, especially when using hands with multiple fingers and manipulating objects that are underactuated. Reinforcement learning (RL) offers an appealing choice for such settings, as it in principle enables a robot to learn to adeptly apply contact forces and manipulate objects without strong modeling assumptions, directly from real-world experience. However, running RL on real-world robotic platforms raises a number of practical issues that lie outside the standard RL formulation, such as difficulties with reward specification, state estimation and the practicalities of autonomous training. Addressing such issues

typically requires significant manual engineering or human intervention. This has led researchers to study alternative solutions, such as transfer from simulation [1]–[3], imitation learning [4]–[6], or use of cumbersome instrumentation, such as motion capture [2], [7]. Even when these issues can be overcome, effective real-world reinforcement learning typically requires considerable reward engineering [8], complex reset mechanisms or scripts [9], and other manually designed components. Each of these solutions erodes the original benefits of autonomy and ease of use that RL should in principle provide. Solving even the simplest tasks with RL requires considerable domain and robotics expertise to program reward functions and reset mechanisms for autonomous operation. Thus, in order to allow for learning-based dexterous manipulation systems to reach their full potential in terms of practicality, accessibility, and scalability, it is critical to limit the assumptions on manual engineering while still providing enough supervision for reinforcement learning to be tractable.

In this work, we propose a robotic learning system that can learn to control high-dimensional multi-fingered robotic hands from raw visual observations, without the need for extensive engineering for every new task. In the absence of simulation, manual reward shaping, and hand-designed state estimation instrumentation, we aim to enable RL to be as autonomous as possible. The robot should be able to practice the task for a long period of time without human intervention, and the task itself should be specified in a way that does not require per-task programming or human-in-the-loop supervision. To this end, we propose a system that autonomously practices a sequence of sub-skills based on high-level milestone specifications provided by the user that break up a complex task into more manageable sub-problems. The milestone specifications consist of snapshots of critical states illustrated by posing the robot and objects in the scene. For multi-fingered hands, such examples are significantly easier to provide than full demonstrations, and our system can use them to learn reward functions that provide sufficient shaping for RL in the real-world without per-task engineering or specific reward design. The system uses

<sup>\*</sup>Both authors contributed equally  
<https://sites.google.com/view/dexterous-avail/>

multi-camera visual observations to localize and manipulate objects, with policies learned end-to-end from pixels and no motion capture. By sequencing the sub-skills appropriately and introducing very simple physical instrumentation (in our experiments, tethering the object to prevent it from falling out of reach), the robot can learn dexterous behaviors by practicing for up to 48 hours fully autonomously. The milestone decomposition makes both reward inference and autonomous practicing significantly easier, enabling real-world learning of complex tasks. Our experiments (see Fig. 1 for an example) show that this approach can learn skills that involve basic grasping, in-hand reorientation, and object manipulation, through significant amounts of practicing, entirely from images and without task-specific reward engineering.

## II. RELATED WORK

Prior work has studied control of complex hands using trajectory optimization [10], [11], policy search [12]–[14], simulation to real-world transfer [3], [15], [16], and real-world reinforcement learning [17], [18]. In contrast to our work, the majority of this prior work has assumed access to compact state representations or accurate simulators and object models. Closer to the system we describe in this paper is prior work on learning visuomotor policies for dexterous manipulation [19]–[21]. However, with the exception of some work we discuss below, prior systems on RL for dexterous manipulation typically require assumptions on manually designed rewards, or ground truth object state observations. These assumptions hinder the application of RL in more real-world settings.

An important consideration in our system is the ability to specify a task without manual reward engineering, by using intermediate milestone examples. Previously studied methods for task specification include having humans provide demonstrations for imitation learning [4], [22], [23], using inverse RL [24]–[26], active settings where users can provide corrections [27]–[29], or ranking-based preferences [30], [31]. While some prior work [32], [33] also uses subgoals, these are firstly restricted to reaching only particular goal states rather than more abstract milestones and are only applied in much simpler simulated problems with perfect state estimation. Motivated by the goal of broader applicability, we do not assume access to expert demonstrations (e.g., via teleoperation or kinesthetic teaching), which can themselves be difficult to provide for high-dimensional systems [34], [35]. For example, providing kinesthetic demonstrations for a full hand-arm robotic system requires very challenging coordination and several simultaneous demonstrators, and is incompatible with vision-based policy learning (as the demonstrator is in the scene, and often occludes the robot or objects). In contrast, we utilize sparse images of intermediate outcomes that can be obtained simply by positioning the robot and object in particular states and build on the VICE framework [36] for reward inference. Our focus is not on devising a new *algorithm* for learning rewards, but on leveraging existing components, such as VICE, to build a complete, autonomous, robotic system that can enable

TABLE I: A comparison between the assumptions of AVAIL and prior autonomous RL methods.

Method	No Hand Engineered Reward	Multi-Task	Vision	High-DoFs
R3L [38]	✓	×	✓	×
MTRF [8]	×	✓	×	✓
Ours	✓	✓	✓	✓

scalable RL with a dexterous manipulator in the real world. Therefore, although some of the building blocks of our system are based on prior work, their combination and the capabilities they enable (learning image-based dexterous manipulation in the real world) are novel.

The most closely related robotic RL systems that have been previously proposed are R3L [37] and MTRF [8]. Our assumptions regarding lightweight instrumentation and vision-based autonomous learning most resemble those of Zhu et al. (2020) [37] (R3L). However, our work tackles a considerably more challenging setting: while Zhu et al. (2020) [37] studied a 3-finger claw mounted on a fixed base, we show that our method can control a 4-fingered hand on a 7 DoF arm. This is done by leveraging a multi-task RL formulation that builds on ideas from MTRF [8] instead of the novelty-based resets in R3L [38], which scale poorly in higher dimensional settings. In contrast to these prior works, our focus is on providing a framework that can enable vision-based learning of object manipulation skills with high-dimensional hands via a lightweight milestone-based task specification mechanism. Part of this requires an automated RL system that can run continuously for 48 hours, though unlike MTRF [8], we still employ lightweight physical instrumentation (by tethering the object to prevent it from falling outside of grasping range). We instead focus on the separate challenges of visual perception and reward specification, avoiding the need for manual reward engineering of the sort used by MTRF and completely circumventing the requirement for motion capture that was crucial in MTRF. We summarize these key system-level differences in Table I.

## III. ROBOTIC PLATFORM AND PROBLEM OVERVIEW

We first present an overview of our robot platform, describing the hardware and task setup, as well as the observation and action space. Then, we provide an overview of our problem setting, focusing on the practical goals of our system. We provide complete details related to our robotic platform in our project website<sup>1</sup> along with details of a simulated analogue that we employ for analysis and ablation experiments.

Our robotic system consists of a custom-built, 4-finger, 16-DoF robot hand, mounted on a 7-DoF Sawyer robotic arm. The arm and hand assembly are positioned over a tabletop surface (Fig. 3, left image). Our policy, which we operate at 8Hz, directly controls each joint position in addition to the Cartesian position and orientation of the arm, resulting in a 22-dimensional action space and 29-dimensional state

<sup>1</sup><https://sites.google.com/view/dexterous-avail>

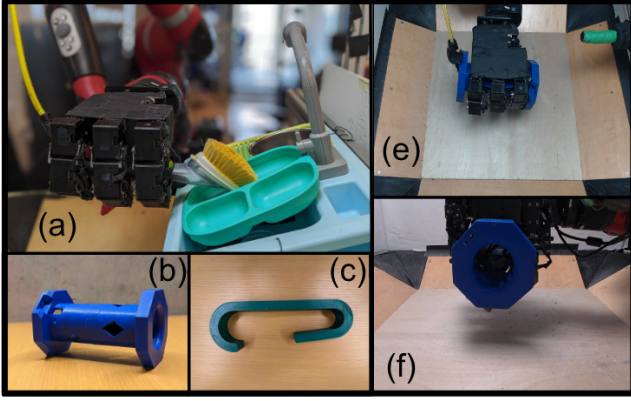


Fig. 3: An overview of our experimental platform: (a) Our robot consists of a 16 DoF four-finger hand mounted on a 7 DoF Sawyer arm; (a, b, c) Objects the robot manipulates in our experiments: a kitchen brush, a cylindrical hose connector, and a hook that must be attached to a handle; (e, f) Observations for the robot come from two monocular RGB cameras.

space. The system is designed to operate for upwards of 48 hours in contact-rich environments without breakage. In addition to the robot’s own joint encoders, two RGB image observations are provided to the robot via two low-cost web cameras and resized to  $84 \times 84$ . We discuss additional details on our project website.

Our tasks consist of manipulation behaviors such as reaching, grasping, in-hand and mid-air reorienting, and inserting. We consider three tasks (shown in Fig 3) for interacting with several different objects: inserting a hose into a connector on the side of the arena, hooking a rope onto a fixture, and cleaning a surface with a kitchen brush. Successfully completing each of these tasks requires correctly sequencing a series of sub-skills. For example, in order to complete our kitchen task, the robot must grasp and reorient that brush palm down without dropping it before making contact with the surface. Constructing and tuning both manual rewards and state estimation systems for each of these tasks separately would typically require laborious human engineering. For each of these tasks however, the only supervision we assume is to allow the user to place the robot and object in the desired position and capture a set of image “snapshots”. We describe how we use this supervision to drive reward inference and task selection via autonomous reinforcement learning in the sections below.

#### IV. PROBLEM FORMALISM AND USER ASSUMPTIONS

In this section, we formalize our problem setting and supervision assumptions. Consider first the Markov decision process (MDP) defined by the tuple  $(\mathcal{S}, \mathcal{A}, p_{dyn}, \rho, \gamma, R)$ , where  $\mathcal{S}$  denotes the state space,  $\mathcal{A}$  denotes the action spaces,  $p_{dyn} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}_{\geq 0}$  denotes the environment dynamics,  $R : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$  denotes the reward function,  $\rho : \mathcal{S} \mapsto \mathbb{R}_{\geq 0}$  denotes the initial state distribution and  $\gamma \in [0, 1)$  denotes the discount factor. The typical objective of episodic RL is to optimize the discounted return  $J(\pi) = \mathbb{E}_{\tau \sim \pi} [\sum_{t=0}^T \gamma^t R(s_t, a_t)]$  with respect to the policy  $\pi$ , where  $\tau = \{(s_i, a_i)\}_{i=0}^{T-1}$  is obtained by sampling  $s_0 \sim \rho(\cdot)$ ,  $a_t \sim \pi(\cdot | s_t)$  and  $s_{t+1} \sim p(\cdot | s_t, a_t)$ .

A principal concern of our work is to ask the question of how best to instantiate RL systems in the real world with minimal per-task engineering, instrumentation and intervention. Standard RL assumes a reward function  $R$  that in practice must often be hand engineered and tuned per-task by a user. This challenge is particularly acute in the dexterous manipulation setting where the desired behavior can often itself be composed of a sequence of complex “sub-tasks” (e.g., grasping, re-orienting, etc) with different objects that would need to be instrumented separately. In addition, independent of being challenging to learn, these sub-tasks must be appropriately sequenced in order to complete the task but also to allow the agent to continue to practice in the event of failure. This necessitates the provision of more fine-grained guidance via user supervision, while carefully balancing the cost of providing such supervision. Furthermore, most RL algorithms assume that the environment is episodic and resets are provided for free. This is not true when considering large scale autonomous operation.

To provide fine-grained supervision both for reward inference and autonomous practicing, we propose a method where a user supplies the robot with a set of sub-problems to practice. These sub-problems are defined by “milestone” examples, which constitute a graph structure:

**Definition IV.1 (Milestones graph).** We assume the user provides a set of outcome images that can be summarized by a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  of cardinality  $|\mathcal{K}|$  indexed by  $z$ , where each vertex  $v \in \mathcal{V}$  is composed of a set of  $M$  outcome images  $\{s_i^z\}_{i=1}^M$ . Each set of outcome images characterizes a semantically meaningful sub-task to be solved. In addition, upon accomplishing a sub-task, a directed edge  $(v, v') \in \mathcal{E}$ , or equivalently a binary label, is provided which indicates which sub-task is to be practiced next.

Consistent with the goal of having the agent continuously practice (i.e., not get stuck), we assume there are no sink nodes in the provided graph.<sup>2</sup> Then, instead of optimizing a single-task objective  $J(\pi)$ , we instead optimize all of the sub-tasks in the milestone graph simultaneously, resulting in a multi-task RL problem. Concretely, we learn a set of  $K$  policies  $\pi_z$  indexed by a categorical variable  $z$  (one for each milestone), optimizing a set of MDPs,  $\mathcal{M} \equiv (\mathcal{S}, \mathcal{A}, p_{dyn}(s_{t+1}|a_t, s_t), \{R_z\}_{z=0}^{K-1}, p_{task}(z|s))$ , where we have introduced a per milestone reward  $R_z$  and task predictor  $p_{task}$ . This leads to the following objective:

$$J_{MT}(\{\pi_i\}_{i=0}^{K-1}) = \sum_{i=0}^{\infty} \left[ \mathbb{E}_{\substack{s_0^i = s_T^{i-1} \\ \tau \sim \pi_{z_i}^i}} \left[ \sum_{t=0}^T \gamma^t R_z(s_t^i, a_t^i) \right] \right] \quad (1)$$

$$z_{i+1} \sim p_{task}(z_{i+1} | s_T^i). \quad (2)$$

<sup>2</sup>This assumption conceivably could be lifted by providing handling of safety-critical states to avoid irreversible sinks, [39], [40]. For simplicity, we leave addressing safety issues for future work.



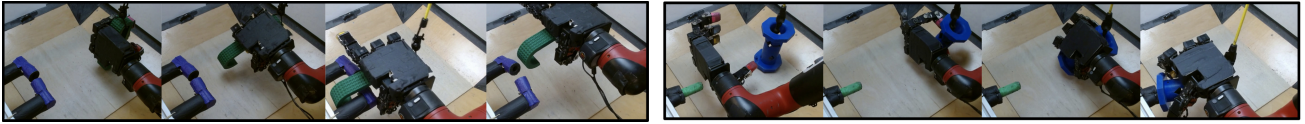


Fig. 5: Filmstrip of the final learned hooking (left) and insertion behavior (right). Using the user-provided milestones, our robot learns a set of skills that allows it to autonomously practice hooking and unhooking (left, right two images) and recover from failure (e.g., after dropping the hook) by regrasping and reorienting the hook (left, left two images). Similarly, using the user-provided milestones, our robot learns a set of skills (e.g., grasp, insert), which together enable successful insertion (right, rightmost image) as well as the stages needed to practice autonomously (right, left three images). After 36 hours of unattended training, our system hooks onto the handle with around a 95% success rate and successfully inserts with around 80% success rate.

### C. Algorithm Summary and Implementation Details

To summarize, given the milestone graph provided by the user, our system, AVAIL (Autonomy ViA mILestones), proceeds as follows. First, AVAIL performs supervised learning of the next task transitions provided by the user as described in Sec. V-B to learn  $p_{\text{task}}^{\phi}(z|s)$ . Next, during training, our approach chooses the most probable task  $z$  using an observed state, which is then used to collect experience using the corresponding policy  $\pi_z$ . We train a set of separate policies  $\pi_z$  for each of the  $K$  sets of example images, with separate critics  $Q_z$  and replay buffers  $\mathcal{B}_z$ . We parameterize each policy  $\pi_z$  as a deep neural network, and train each policy using the soft actor-critic algorithm (SAC) [44] using rewards  $R_z$  that are inferred via the multi-task VICE [36] algorithm trained in the loop. Finally, rather than resampling the task every step, we do so every  $T = 100$  steps.

## VI. EXPERIMENTAL EVALUATION

Our experiments whether AVAIL can learn complex manipulation skills in the real world with visual milestones. We evaluate our approach on three real world manipulation tasks that require successfully sequencing a set of skills and performing complex coordinated finger motions, in addition to a simulation evaluation for comparison with prior methods. Supplementary videos are available on the project website: <https://sites.google.com/view/dexterous-avail/>.

### A. Real-World Task Descriptions

We begin by describing our tasks. The objects and workspace in our experiments can be seen in Fig. 3. The environments are mostly uninstrumented, except for a passive tether that prevents the object from falling out of reach. Further task details can be found on the project website.

*a) Using a kitchen brush.:* Our first task requires the robot to use a two-sided cleaning brush to scrape a plate (see Fig. 1). This involves grasping and reorienting the brush with the fingers to face the plate, which is challenging due to the need to balance and rotate it. The palm-down manipulation of the brush is particularly challenging, as it requires balancing it so that it doesn't fall and rotating it around its long axis via a coordinated finger gait. Milestones include grasping the brush, scraping the surface, reorienting it, and making contact with the plate.

*b) Hose connector insertion.:* The goal of the second task is to attach a cylindrical hose connector to a peg connector, which requires reaching, grasping, reorienting, and performing the insertion as the task milestones. While the task is simpler in terms of dexterity, it requires visual

perception to carefully insert the connector onto the peg connector (see Fig. 5, left).

*c) Rope hooking.:* The third task requires attaching a hook to a handle ( Fig. 5, right). The robot must grasp, reorient, and hook/unhook the object for its milestones. This task requires visually servoing the hook over a handle.

### B. Real-world evaluation

To evaluate our system, we save the policies at regular intervals and evaluate their performance after training, so as not to interrupt the training process. For all tasks, the evaluation metric for each milestone is a binary success measurement based on the distance of the hand and object to the desired pose. We provide more details on our evaluation setup on our project website.

**Real-world skill learning:** The learning curves for real-world training are provided in Fig. 6. We observe that AVAIL automatically provides a degree of scaffolding by successfully learning skills early on in training (blue curve in left plot, blue and orange curve in center right plot), which corresponds to the robot being able to continuously regrasp and reorient the object. By the end of training, the robot successfully performs all milestones with a  $> 80$  success rate. No additional instrumentation is required beyond changing the object and fixture. Upon specifying the visual milestones, the robot is capable of completely unattended learning for approximately 48 hours of robot time.

**Real-world task graph learning:** To understand the effect of our learned task graph, we compare to a hand-designed task graph on our insertion task. This hand-designed task graph encodes a heuristic strategy where each of the tasks are practiced sequentially. The comparative success rates can be seen in Fig. 8. We find that our learned task graph outperforms this heuristic based task graph in terms of sample efficiency. We provide additional analysis on our real world rope hooking task on our project website. We find that our approach is robust to errors in our learned task graph, although future work could likely improve overall sample efficiency by improving task graph training.

### C. Simulated Comparison

Finally, we compare AVAIL to prior autonomous RL method on the (DHandValvePickup-v0) simulation domain developed in prior work [8]. We first compare to a standard state-of-the-art RL algorithm, soft actor-critic [44] (which we denote as SAC) using a reward learned from example images of a successful pickup or sparse reward. Note that prior work has used this approach for real-world

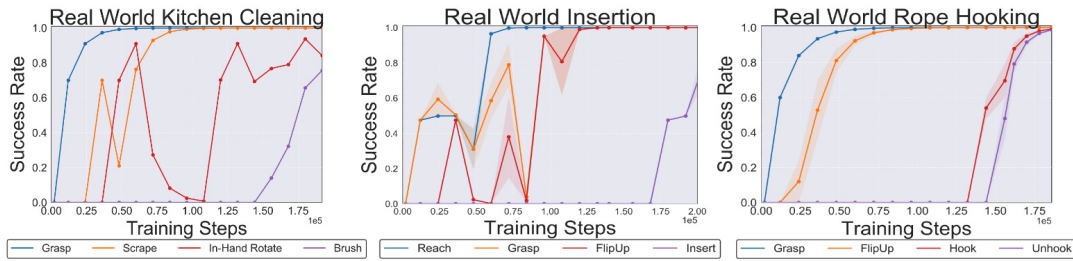


Fig. 6: Success rates of different milestones across our real world dexterous manipulation tasks. Our method is able to perform all the kitchen milestones with around 80% success rate (left), successfully perform pipe insertion (middle, red curve) with around 80% success, and nearly perfectly hook and unhook the rope (right, red/purple curve). Overall, success on other milestones improves earlier in training (blue, orange curves), equipping the robot with skills to autonomously retry the task.

robotics tasks [45]. Next, we compare to a forward backward controller [46], which can be seen as providing two milestones: one to pick up the object, and one to place it back on the table. Finally, we compare to R3L [37], where we provide the “forward” policy with a set of pickup goals and follow Zhu et al. [37] by interleaving training of the forward policy with a “perturbation controller” trained with an intrinsic reward based on random network distillation [38].

**Simulated comparative analysis.** We evaluate the performance of each method by sampling a random initial position of the object in the workspace and running the learned policy. We evaluate the task success over the course of training in Fig. 7. Prior methods do not make successful progress on this task, due to the combination of reset-free training and the lack of a shaped reward. Without any handling of the reset-free setting, both variants of SAC fail to progress. While R3L in principle can handle the autonomous setting by perturbing the state between trials, the large high-dimensional task simply provides too many ways for the purely novelty-seeking controller to modify the environment with a meaningful reset. The forward backward controller (red), which can be seen as an instantiation of our approach with two milestones, is the only prior method that succeeds in making progress.<sup>3</sup> Overall, this suggests that the improved performance can be achieved through providing more granular milestones.

## VII. DISCUSSION AND FUTURE WORK

We proposed a method for multi-task learning for dexterous manipulation from high dimensional image observations. Our method, AVAIL, constructs a task graph from a modest number of user-provided milestone examples. This task graph illustrates how to practice and reset the task, and provides guidance to the learning process in lieu of more standard manual reward shaping. While the milestone examples require human effort to provide, we expect in many cases that this effort is significantly lower than providing full demonstrations. Much like a teacher or coach might instruct a student not just by telling them the *goal* of a task but how they should go about practicing it, the milestone examples serve to provide guidance to the agent for how it should go about learning the desired behavior. Our experiments show that this approach effectively produces a learning process

<sup>3</sup>We additionally compared the forward backward controller to our approach in the real world, but found it did not make progress on our tasks. We provide details of this analysis on our project website.

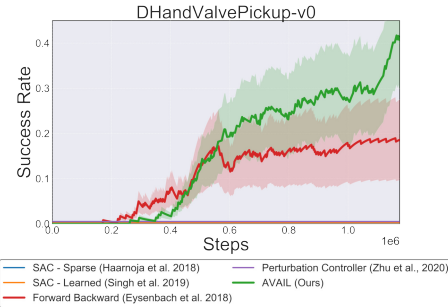


Fig. 7: Success rates of each method averaged over 5 seeds for the full task on the DHandValvePickup-v0 domain. Novelty based resets (purple) fail to make progress in this high DoF control problem. Compared to methods with fewer degree of supervision,  $K = 0, 1, 2$  milestones (blue, orange, red), our results illustrate the benefits of milestone supervision.

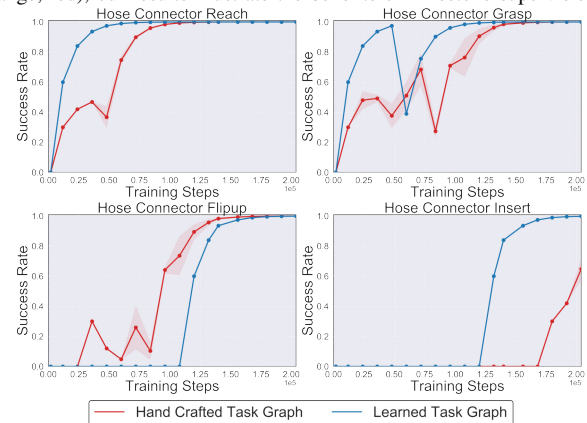


Fig. 8: A comparison of the success rate of each task on our real world insertion task. We find that using a learned task graph results in faster convergence on our real world robotic task compared to a hand-coded heuristic based task graph. We find the robot begins to consistently perform the final insertion task 25% faster than a hand-coded task graph.

where the agent first practices the easier tasks, and then builds up the more complex tasks on top of them, all the while learning autonomously without resetting.

One limitation of our current system is that we do employ some physical instrumentation, by tethering the object so that it doesn’t fall out of reach. This is because learning to pick up the object from any location was still too challenging for our method, as the range of possible scenarios was too large. Developing more capable reinforcement learning methods that can address this is an important direction for future work.

## VIII. ACKNOWLEDGEMENT

This research project was partially supported by the Office of Naval Research, with computing resources donated by Microsoft Azure.

## REFERENCES

- [1] 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.
- [2] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 3803–3810.
- [3] OpenAI, M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. W. Pachocki, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba, "Learning dexterous in-hand manipulation," *CoRR*, vol. abs/1808.00177, 2018. [Online]. Available: <http://arxiv.org/abs/1808.00177>
- [4] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, "A survey of robot learning from demonstration," *Robotics and autonomous systems*, vol. 57, no. 5, pp. 469–483, 2009.
- [5] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne, "Imitation learning: A survey of learning methods," *ACM Computing Surveys (CSUR)*, vol. 50, no. 2, pp. 1–35, 2017.
- [6] T. Osa, J. Pajarinen, G. Neumann, J. A. Bagnell, P. Abbeel, and J. Peters, "An algorithmic perspective on imitation learning," *Foundations and Trends® in Robotics Vol. 7, No. 1-2 (2018) 1–179*, 2018.
- [7] V. Kumar, A. Gupta, E. Todorov, and S. Levine, "Learning dexterous manipulation policies from experience and imitation," *arXiv preprint arXiv:1611.05095*, 2016.
- [8] A. Gupta, J. Yu, T. Z. Zhao, V. Kumar, A. Rovinsky, K. Xu, T. Devlin, and S. Levine, "Reset-free reinforcement learning via multi-task learning: Learning dexterous manipulation behaviors without human intervention," *ICRA*, 2021.
- [9] A. Nagabandi, K. Konolige, S. Levine, and V. Kumar, "Deep dynamics models for learning dexterous manipulation," in *Conference on Robot Learning*, 2020, pp. 1101–1112.
- [10] I. Mordatch, Z. Popović, and E. Todorov, "Contact-invariant optimization for hand manipulation," in *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation*. Eurographics Association, 2012.
- [11] V. Kumar, Y. Tassa, T. Erez, and E. Todorov, "Real-time behaviour synthesis for dynamic hand-manipulation," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2014, pp. 6808–6815.
- [12] J. Kober and J. Peters, "Policy search for motor primitives in robotics," *Advances in neural information processing systems*, vol. 21, 2008.
- [13] M. Posa, C. Cantu, and R. Tedrake, "A direct method for trajectory optimization of rigid bodies through contact," *The International Journal of Robotics Research*, vol. 33, no. 1, pp. 69–81, 2014.
- [14] A. Rajeswaran, V. Kumar, A. Gupta, J. Schulman, E. Todorov, and S. Levine, "Learning complex dexterous manipulation with deep reinforcement learning and demonstrations," in *Robotics: Science and Systems (RSS)*, 2018, pp. 1101–1112.
- [15] K. Lowrey, S. Kolev, J. Dao, A. Rajeswaran, and E. Todorov, "Reinforcement learning for non-prehensile manipulation: Transfer from simulation to physical system," in *2018 IEEE International Conference on Simulation, Modeling, and Programming for Autonomous Robots (SIMPAP)*. IEEE, 2018, pp. 35–42.
- [16] A. Allshire, M. Mittal, V. Lodaya, V. Makoviychuk, D. Makoviichuk, F. Widmaier, M. Wüthrich, S. Bauer, A. Handa, and A. Garg, "Transferring dexterous manipulation from gpu simulation to a remote real-world triferger," *arXiv preprint arXiv:2108.09779*, 2021.
- [17] H. Van Hoof, T. Hermans, G. Neumann, and J. Peters, "Learning robot in-hand manipulation with tactile features," in *2015 IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*. IEEE, 2015, pp. 121–127.
- [18] H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar, "Dexterous manipulation with deep reinforcement learning: Efficient, general, and low-cost," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 3651–3657.
- [19] D. Jain, A. Li, S. Singhal, A. Rajeswaran, V. Kumar, and E. Todorov, "Learning deep visuomotor policies for dexterous hand manipulation," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 3636–3643.
- [20] P. Mandikal and K. Grauman, "Learning dexterous grasping with object-centric visual affordances," *arXiv preprint arXiv:2009.01439*, 2020.
- [21] I. Akinola, J. Varley, and D. Kalashnikov, "Learning precise 3d manipulation from multiple uncalibrated cameras," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 4616–4622.
- [22] S. Ross, N. Melik-Barkhudarov, K. S. Shankar, A. Wendel, D. Dey, J. A. Bagnell, and M. Hebert, "Learning monocular reactive UAV control in cluttered natural environments," in *2013 IEEE International Conference on Robotics and Automation*, 2013.
- [23] S. Reddy, A. D. Dragan, and S. Levine, "Sqil: Imitation learning via reinforcement learning with sparse rewards," *arXiv preprint arXiv:1905.11108*, 2019.
- [24] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *AAAI*. AAAI Press, 2008.
- [25] M. Wulfmeier, P. Ondruska, and I. Posner, "Maximum entropy deep inverse reinforcement learning," *arXiv preprint arXiv:1507.04888*, 2015.
- [26] N. D. Ratliff, J. A. Bagnell, and M. Zinkevich, "Maximum margin planning," in *Machine Learning, Proceedings of the Twenty-Third International Conference ICML, 2006*.
- [27] D. P. Losey and M. K. O'Malley, "Including uncertainty when learning from human corrections," in *Conference on Robot Learning*. PMLR, 2018, pp. 123–132.
- [28] Y. Cui and S. Niekum, "Active reward learning from critiques," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 6907–6914.
- [29] J. D. Co-Reyes, A. Gupta, S. Sanjeev, N. Altieri, J. DeNero, P. Abbeel, and S. Levine, "Guiding policies with language via meta-learning," *CoRR*, vol. abs/1811.07882, 2018. [Online]. Available: <http://arxiv.org/abs/1811.07882>
- [30] V. Myers, E. Biyik, N. Anari, and D. Sadigh, "Learning multimodal rewards from rankings," in *Conference on Robot Learning*. PMLR, 2022, pp. 342–352.
- [31] D. S. Brown, W. Goo, and S. Niekum, "Better-than-demonstrator imitation learning via automatically-ranked demonstrations," in *Conference on robot learning*. PMLR, 2020, pp. 330–359.
- [32] O. Nachum, S. S. Gu, H. Lee, and S. Levine, "Data-efficient hierarchical reinforcement learning," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/e6384711491713d29bc63fc5eeb5ba4f-Paper.pdf>
- [33] A. Levy, R. P. Jr., and K. Saenko, "Hierarchical actor-critic," *CoRR*, vol. abs/1712.00948, 2017. [Online]. Available: <http://arxiv.org/abs/1712.00948>
- [34] B. Akgun, M. Cakmak, J. W. Yoo, and A. L. Thomaz, "Trajectories and keyframes for kinesthetic teaching: A human-robot interaction perspective," in *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, 2012, pp. 391–398.
- [35] V. Villani, F. Pini, F. Leali, and C. Secchi, "Survey on human–robot collaboration in industrial settings: Safety, intuitive interfaces and applications," *Mechatronics*, vol. 55, pp. 248–266, 2018.
- [36] J. Fu, A. Singh, D. Ghosh, L. Yang, and S. Levine, "Variational inverse control with events: A general framework for data-driven reward definition," *arXiv preprint arXiv:1805.11686*, 2018.
- [37] H. Zhu, J. Yu, A. Gupta, D. Shah, K. Hartikainen, A. Singh, V. Kumar, and S. Levine, "The ingredients of real-world robotic reinforcement learning," *arXiv preprint arXiv:2004.12570*, 2020.
- [38] Y. Burda, H. Edwards, A. J. Storkey, and O. Klimov, "Exploration by random network distillation," in *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. [Online]. Available: <https://openreview.net/forum?id=H1IJnR5Ym>
- [39] J. Garcia and F. Fernández, "A comprehensive survey on safe reinforcement learning," *Journal of Machine Learning Research*, vol. 16, no. 1, pp. 1437–1480, 2015.
- [40] K. Srinivasan, B. Eysenbach, S. Ha, J. Tan, and C. Finn, "Learning to be safe: Deep rl with a safety critic," *arXiv preprint arXiv:2010.14603*, 2020.
- [41] K. Li, A. Gupta, A. Reddy, V. H. Pong, A. Zhou, J. Yu, and S. Levine, "Mural: Meta-learning uncertainty-aware rewards for outcome-driven reinforcement learning," in *Proceedings of the 38th*

- International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 6346–6356. [Online]. Available: <https://proceedings.mlr.press/v139/i21g.html>
- [42] T. Hiraoka, T. Imagawa, T. Hashimoto, T. Onishi, and Y. Tsuruoka, “Dropout q-functions for doubly efficient reinforcement learning,” *arXiv preprint arXiv:2110.02034*, 2021.
- [43] I. Kostrikov, D. Yarats, and R. Fergus, “Image augmentation is all you need: Regularizing deep reinforcement learning from pixels,” *arXiv preprint arXiv:2004.13649*, 2020.
- [44] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel *et al.*, “Soft actor-critic algorithms and applications,” *arXiv preprint arXiv:1812.05905*, 2018.
- [45] A. Singh, L. Yang, K. Hartikainen, C. Finn, and S. Levine, “End-to-end robotic reinforcement learning without reward engineering,” *arXiv preprint arXiv:1904.07854*, 2019.
- [46] B. Eysenbach, S. Gu, J. Ibarz, and S. Levine, “Leave no trace: Learning to reset for safe and autonomous reinforcement learning,” *arXiv preprint arXiv:1711.06782*, 2017.