# Collective Robot Reinforcement Learning with Distributed Asynchronous Guided Policy Search

Ali Yahya[1]    Adrian Li[1]    Mrinal Kalakrishnan[1]    Yevgen Chebotar[2]    Sergey Levine[3]

*Abstract*— In principle, reinforcement learning and policy search methods can enable robots to learn highly complex and general skills that may allow them to function amid the complexity and diversity of the real world. However, training a policy that generalizes well across a wide range of real-world conditions requires far greater quantity and diversity of experience than is practical to collect with a single robot. Fortunately, it is possible for multiple robots to share their experience with one another, and thereby, learn a policy collectively. In this work, we explore distributed and asynchronous policy learning as a means to achieve generalization and improved training times on challenging, real-world manipulation tasks. We propose a distributed and asynchronous version of Guided Policy Search and use it to demonstrate collective policy learning on a vision-based door opening task using four robots. We show that it achieves better generalization, utilization, and training times than the single robot alternative.
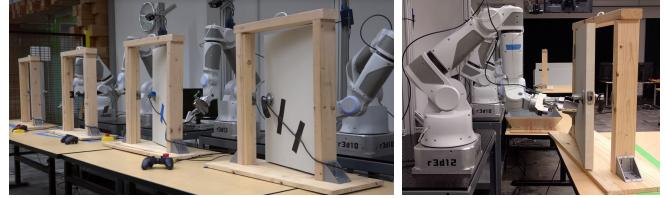
Fig. 1. Multiple robots collaborating to learn the door opening skill. Our system allows the robots to operate continuously to collect a large amount of diverse experience, while the policy is simultaneously trained with a replay buffer of the latest trajectory samples.

## I. INTRODUCTION

Policy search techniques show promising ability to learn feedback control policies for robotic tasks with high-dimensional sensory inputs through trial and error [1, 2, 3, 4]. Most successful applications of policy search, however, rely on considerable manual engineering of suitable policy representations, perception pipelines, and low-level controllers to support the learned policy. Recently, deep reinforcement learning (RL) methods have been used to show that policies for complex tasks can be trained end-to-end, directly from raw sensory inputs (like images [5, 6]) to actions. Such methods are difficult to apply to real-world robotic applications because of their high sample complexity. Methods based on Guided Policy Search (GPS) [7], which convert the policy search problem into a supervised learning problem, with a local trajectory-centric RL algorithm acting as a teacher, reduce sample complexity and thereby help make said applications tractable. However, training such a policy to generalize well across a wide variety of real-world conditions requires far greater quantity and diversity of experience than is practical to collect with a single robot.

Fortunately, it is possible for multiple robots to share their experience with one another, and thereby, learn a policy collectively. In this work, we explore distributed and asynchronous policy learning (also known hereafter in this work as collective policy learning) as a means to achieve generalization and improved training times on challenging,

real-world manipulation tasks. Collective policy learning presents a number of unique challenges. These challenges can be broadly categorized as utilization challenges and synchronization challenges. On one hand, we would like to maximize robot utilization — the fraction of time that the robots spend collecting experience for learning. On the other hand, each robot must allocate compute and bandwidth to process and communicate its experience to other robots, and the system as a whole needs to synchronize the assimilation of each robot's experience into the collective policy.

The main contribution of this work is a system for collective policy learning. We address the aforementioned utilization and synchronization challenges with a novel distributed and asynchronous variant of Guided Policy Search. In our system, multiple robots practice the task simultaneously, each on a distinct instance of the task, and jointly train a single policy whose parameters are maintained centrally by a parameter server. To maximize utilization, each robot continues to practice and optimize its own local policy while the single global policy is trained from a buffer of previously collected experience. For high-dimensional policies such as those based on neural networks, the increase in utilization that is conferred by asynchronous training is significant. Consequently, this approach dramatically brings down the total amount of time required to learn complex visuomotor policies using GPS, and makes this technique scalable to more realistic applications which require greater data diversity.

We evaluate our approach in simulation and on a real-world door opening task (shown in Figure 1), where both the pose and the appearance of the door vary across task instances. We show that our system achieves better generalization, utilization, and training times than the single robot alternative.

## II. RELATED WORK

Robotic motor skill learning has shown considerable promise for enabling robots to autonomously learn complex

[1]Ali Yahya, Adrian Li, and Mrinal Kalakrishnan are with X, Mountain View, CA 94043, USA. {alive,alhli,kalakris}@x.team

[2]Yevgen Chebotar is with the Department of Computer Science, University of Southern California, Los Angeles, CA 90089, USA. This research was conducted during Yevgen's internship at X.

[3]Sergey Levine is with Google Brain, Mountain View, CA 94043, USA.

motion skills [1, 2, 3, 4]. However, most successes in robotic motor skill learning have involved significant manual design of representations in order to enable policies to generalize effectively. For example, the well-known dynamic movement primitive representation [8] has been widely used to generalize learned skills by adapting the goal state, but it inherently restricts the learning process to trajectory-centric behaviors.

Enabling robotic learning with more expressive policy classes that can represent more complex strategies has the potential of eliminating the need for the manual design of representations. Recent years have seen improvement in the generalizability of passive perception systems, in domains such as computer vision, natural language processing, and speech recognition through the use of deep learning techniques [9]. These methods combine deep neural networks with large datasets to achieve remarkable results on a diverse range of real-world tasks. However, the requirement of large labeled datasets has limited the application of such methods to robotic learning problems. While several works have extended deep learning methods to simulated [5, 6] and real-world [7, 10] robotic tasks, the kind of generalization exhibited by deep learning in passive perception domains has not yet been demonstrated for robotic skill learning. This may be due to the fact that robotic learning experiments tend to use relatively small amounts of data in constrained domains, with a few hours of experience collected from a single robot in each experiment.

A central motivation behind our work is the ability to apply deep learning to robotic manipulation by making it feasible to collect large amounts of on-policy experience with real physical platforms. While this may seem impractical for small-scale laboratory experiments, it becomes much more realistic when we consider a possible future where robots are deployed in the real-world to perform a wide variety of skills. The challenges of asynchrony, utilization, and parallelism, which we aim to address in this work, are central for such real-world deployments. The ability of robotic systems to learn more quickly and effectively by pooling their collective experience has long been recognized in the domain of cloud robotics, where it is typically referred to as collective robotic learning [11, 12, 13, 14]. Our work therefore represents a step toward more practical and powerful collective learning with distributed, asynchronous data collection.

Distributed systems have long been an important subject in deep learning [15]. While distributed asynchronous architectures have previously been used to optimize controllers for simulated characters [16], our work is, to the best of our knowledge, the first to experimentally explore distributed asynchronous training of deep neural network policies for real-world robotic control. In our work, we parallelize both data collection and neural network policy training across multiple machines.

## III. PRELIMINARIES ON GUIDED POLICY SEARCH

In this section, we define the problem formulation and briefly summarize theb guided policy search (GPS) algorithm, specifically pointing out computational bottlenecks that can be alleviated through asynchrony and parallelism. A more complete description of the theoretical underpinnings of the method can be found in prior work [7]. The goal of policy search methods is to optimize the parameters $\theta$ of a policy $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$, which defines a probability distribution over robot actions $\mathbf{u}_t$ conditioned on the system state $\mathbf{x}_t$ at each time step $t$ of a task execution. Let $\tau = (\mathbf{x}_1, \mathbf{u}_1, \ldots, \mathbf{x}_T, \mathbf{u}_T)$ be a trajectory of states and actions. Given a task cost function $l(\mathbf{x}_t, \mathbf{u}_t)$, we define the trajectory cost $l(\tau) = \sum_{t=1}^{T} l(\mathbf{x}_t, \mathbf{u}_t)$. Policy optimization is performed with respect to the expected cost of the policy:

$$J(\theta) = E_{\pi_\theta}[l(\tau)] = \int l(\tau) p_{\pi_\theta}(\tau) d\tau,$$

where $p_{\pi_\theta}(\tau)$ is the policy trajectory distribution given the system dynamics $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$:

$$p_{\pi_\theta}(\tau) = p(\mathbf{x}_1) \prod_{t=1}^{T} p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) \pi_\theta(\mathbf{u}_t|\mathbf{x}_t).$$

Most standard policy search methods aim to directly optimize this objective, for example by estimating the gradient $\nabla_\theta J(\theta)$. However, this kind of direct model-free method can quickly become intractable for very high-dimensional policies, such as the large neural network policies considered in this work [4]. An alternative approach is to train the deep neural network with supervised learning, using a simpler local policy optimization method to produce supervision. To that end, guided policy search introduces a two-step approach for learning high-dimensional policies by combining the benefits of simple, efficient trajectory-centric RL and supervised learning of high-dimensional, nonlinear policies. Instead of directly learning the policy parameters with reinforcement learning, a trajectory-centric algorithm is first used to learn simple local controllers $p_i(\mathbf{u}_t|\mathbf{x}_t)$ for trajectories with various initial conditions, which might correspond, for instance, to different poses of a door for a door opening task. We refer to these controllers as *local policies*. In this work, we employ time-varying linear-Gaussian controllers of the form $p_i(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(\mathbf{K}_t\mathbf{x}_t + \mathbf{k}_t, \mathbf{C}_t)$ to represent these local policies, following prior work [7].

After optimizing local policies, the controls from these policies are used to create a training set for learning a complex high-dimensional *global policy* in a supervised manner. Hence, the final global policy generalizes to the initial conditions of multiple local policies and can contain thousands of parameters, which can be efficiently learned with supervised learning. Furthermore, while trajectory optimization might require the full state $\mathbf{x}_t$ of the system to be known, it is possible to only use the observations $\mathbf{o}_t$ of the full state for training a global policy $\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$. This allows the global policy to predict actions from raw observations at test time [7].

In this work, we will examine a general asynchronous framework for guided policy search algorithms, and we will show how this framework can be instantiated to extend two prior guided policy search methods: BADMM-based guided policy search [7] and mirror descent guided policy search

(MDGPS) [17]. Both algorithms share the same overall structure, with alternating optimization of the local policies via trajectory-centric RL, which in the case of our system is either a model-based algorithm based on LQR [18] or a model-free algorithm based on $PI^2$ [3], and optimization of the global policy via supervised learning through stochastic gradient descent (SGD). The adaptation of $PI^2$ to guided policy search is described in detail in a companion paper [19]. The difference between the two methods is the mechanism that is used to keep the local policies close to the global policy. This is extremely important, since in general not all local policies can be reproduced effectively by a single global policy.

*a) BADMM-based GPS:* In BADMM-based guided policy search [7], the alternating optimization is formalized as a constrained optimization of the form

$$\min_{\theta, p_1, \ldots, p_N} \sum_{i=1}^{N} E_{\tau \sim p_i}[l(\tau)] \text{ s.t. } p_i(\mathbf{u}_t|\mathbf{x}_t) = \pi_\theta(\mathbf{u}_t|\mathbf{x}_t) \, \forall \, \mathbf{x}_t, \mathbf{u}_t, i.$$

This is equivalent in expectation to optimizing $J(\theta)$, since $\sum_{i=1}^{N} E_{\tau \sim p_i}[l(\tau)] = \sum_{i=1}^{N} E_{\tau \sim \pi_\theta}[l(\tau)]$ when the constraint is satisfied, and $\sum_{i=1}^{N} E_{\tau \sim \pi_\theta}[l(\tau)] \approx E_{\mathbf{x}_1 \sim p(\mathbf{x}_1), \tau \sim \pi_\theta}[l(\tau)]$ when the initial states $\mathbf{x}_1^i$ are sampled from $p(\mathbf{x}_1)$. The constrained optimization is then solved using the Bregman ADMM algorithm [20], which augments the objective for both the local and global policies with Lagrange multipliers that keep them similar in terms of KL-divergence. These terms are denoted $\phi_i(\tau, \theta, \tau)$ and $\phi_\theta(p_i, \theta, \tau)$ for the local and global policies, respectively, so that the global policy is optimized with respect to the objective

$$\min_{\theta} \sum_{i=1}^{N} E_{\tau \sim p_i} \left[ \sum_{t=1}^{T} D_{\text{KL}}(\pi_\theta(\mathbf{u}_t|\mathbf{x}_t) \| p_i(\mathbf{u}_t|\mathbf{x}_t)) + \phi_\theta(p_i, \theta, \tau) \right], \tag{1}$$

and the local policies are optimized with respect to

$$\min_{p_i} E_{\tau \sim p_i(\tau)}[l(\tau)] \text{ s.t. } D_{\text{KL}}(p_i(\tau) \| \bar{p}_i(\tau)) < \epsilon, \tag{2}$$

where $\bar{p}_i$ is the local policy at the previous iteration. The constraint ensures that the local policies only change by a small amount at each iteration, to prevent divergence of the trajectory-centric RL algorithm, analogously to other recent RL methods [2, 21]. The derivations of $\phi_i(\tau, \theta, \mathbf{x}_t)$ and $\phi_\theta(p_i, \theta, \mathbf{x}_t)$ are provided in prior work [7].

*b) MDGPS:* In MDGPS [17], the local policies are optimized with respect to

$$\min_{p_i} E_{\tau \sim p_i(\tau)}[l(\tau)] \text{ s.t. } D_{\text{KL}}(p_i(\tau) \| \pi_\theta(\tau)) < \epsilon, \tag{3}$$

where the constraint directly limits the deviation of the local policies from the global policies. This can be interpreted as the generalized gradient step in mirror descent, which improves the policy with respect to the objective. The supervised learning step simply minimizes the deviation from the local policies, without any additional augmentation terms:

$$\min_{\theta} \sum_{i=1}^{N} E_{\tau \sim p_i} \left[ \sum_{t=1}^{T} D_{\text{KL}}(\pi_\theta(\mathbf{u}_t|\mathbf{x}_t) \| p_i(\mathbf{u}_t|\mathbf{x}_t)) \right]. \tag{4}$$

---

**Algorithm 1** Standard synchronous guided policy search
1: **for** iteration $k \in \{1, \ldots, K\}$ **do**
2:    Generate sample trajectories starting from each $\mathbf{x}_1^i$ by executing $p_i(\mathbf{u}_t|\mathbf{x}_t)$ or $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$ on the robot.
3:    Use samples to optimize each of the local policies $p_i(\mathbf{u}_t|\mathbf{x}_t)$ from each $\mathbf{x}_1^i$ with respect to Equation (2) or (3), using either LQR or $PI^2$.
4:    Optimize the global policy $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$ according to Equation (1) or (4) with SGD.
5: **end for**

---

This can be interpreted as the projection step in mirror descent, such that the overall algorithm optimizes the global policy subject to the constraint that the policy should lie within the manifold defined by the policy class, which in our case corresponds to neural networks.

Both GPS algorithms are summarized in Algorithm 1, and consist of two alternating phases: optimization of the local policies with trajectory-centric RL, and optimization of the global policy with supervised learning. Several different trajectory-centric RL algorithms may be used, and we summarize the ones used in our experiments below.

### A. Local Policy Optimization

The GPS framework is generic with respect to the choice of local policy optimizer. In this work, we consider two possible methods for local policy optimization:

*c) LQR with local models:* To take a model-based approach to optimization of time-varying linear-Gaussian local policies, we can observe that, under time-varying linear-Gaussian dynamics, the local policies can be optimized analytically using the LQR method, or the iterative LQR method in the case of non-quadratic costs [22]. However, this approach requires a linearization of the system dynamics, which are generally not known for complex robotic manipulation tasks. As described in prior work [18], we can still use LQR if we fit a time-varying linear-Gaussian model to the samples using linear regression. In this approach, the samples generated on line 1 of Algorithm 1 are used to fit a time-varying linear-Gaussian dynamics model of the form $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{F}_{\mathbf{x},t}\mathbf{x}_t + \mathbf{F}_{\mathbf{u},t}\mathbf{u}_t + \mathbf{f}_t, \mathbf{N}_t)$. As suggested in prior work, we can use a Gaussian mixture model (GMM) prior to reduce the sample complexity of this linear regression fit, and we can accommodate the constraints in Equation (2) or (3) with a simple modification that uses LQR within a dual gradient descent loop [18].

*d) $PI^2$:* Policy Improvement with Path Integrals ($PI^2$) is a model-free policy search method based on the principles of stochastic optimal control [3]. It does not require fitting of linear dynamics and can be applied to tasks with highly discontinuous dynamics or non-differentiable costs. In this work, we employ $PI^2$ to learn feedforward commands $\mathbf{k}_t$ of time-varying linear-Gaussian controllers as described in [19].

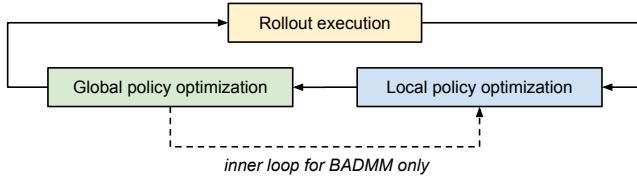The controls at each time step are updated according to

Fig. 2. Diagram of the training loop for synchronous GPS with a single replica. Rollout execution corresponds to line 2 in Algorithm 1, local policy optimization to line 3, and global policy optimization to line 4. In BADMM-based GPS, the algorithm additionally alternates between local and global policy optimization multiple times before executing new rollouts. This sequential version of the algorithm requires training to pause while performing rollouts, and vice versa.

the soft-max probabilities $P_{i,t}$ based on their cost-to-go $S_{i,t}$:

$$S_{i,t} = S(\tau_{i,t}) = \sum_{j=t}^{T} l(\mathbf{x}_{i,j}, \mathbf{u}_{i,j}), \ \ P_{i,t} = \frac{e^{-\frac{1}{\eta}S_{i,t}}}{\sum_{i=1}^{N} e^{-\frac{1}{\eta}S_{i,t}}},$$

where $l(\mathbf{x}_{i,j}, \mathbf{u}_{i,j})$ is the cost of sample $i$ at time $j$. In this way, trajectories with lower costs become more probable after the policy update. For learning feedforward commands, the policy update corresponds to a weighted maximum likelihood estimation of the new mean $\mathbf{k}_t$ and the noise covariance $\mathbf{C}_t$. In this work, we use relative entropy optimization [2] to determine the temperature $\eta$ at each time step independently, based on a KL-divergence constraint between policy updates.

Both the LQR model-based method and the PI$^2$ model-free algorithm require samples in order to improve the local policies. In the BADMM variant of GPS, these samples are always generated from the corresponding local policies. However, in the case of MDGPS, the samples can in fact be generated directly by the global policy, with the local policies only existing temporarily within each iteration for the purpose of policy improvement. In this case, new initial states can be sampled at each iteration of the algorithm, with new local policies instantiated for each one [17]. We make use of this capability in our experiments to train the global policy on a wider range of initial states in order to improve generalization.

## IV. Asynchronous Distributed Guided Policy Search

In synchronous GPS, rollout execution and policy optimization occur sequentially (see Figure 2). This training regime presents two challenges: (1) there is considerable downtime for the robot while the policies are being optimized, and (2) there are synchronization issues when extending the global policy optimization to use data collected across multiple robots.

To overcome these challenges, we propose a modification to GPS which is both asynchronous and distributed (see Figure 3). In our asynchronous distributed GPS method (ADGPS), the algorithm is decoupled into global and local worker threads. The global workers are responsible for continuously optimizing the global policy using a buffer of experience data, which we call the replay memory. The local workers execute the current controllers on their respective robots, adding the collected data to the replay memory. The
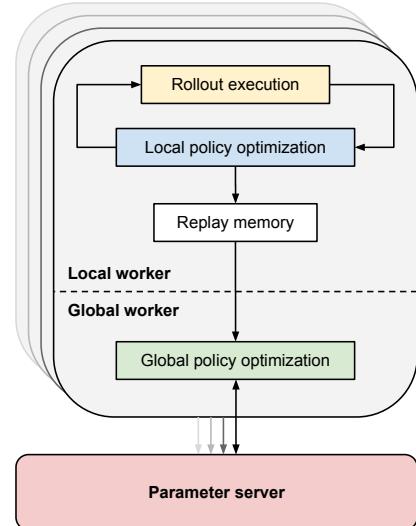


Fig. 3. The training loop for ADGPS with multiple replicas. Rollout execution and global policy optimization are decoupled via the replay memory. Multiple robots concurrently collect data and asynchronously update the parameter server, allowing maximal utilization of both computational and robot resources, as well as parallelization across multiple robots and servers.

local workers are also responsible for updating the local policies. Note, however, that updating the local policies is very quick when compared to global policy updates, since the local policy update requires either a small number of LQR backward passes, or simply a weighted average of the sample controls, if using the PI$^2$ method. This operation can be completed in just a few seconds, while global policy training requires stochastic gradient descent (SGD) optimization of a deep neural network, and can take hours.

The local and global worker threads communicate through the replay memory, which stores the rollouts and optimized trajectories from each local worker. Since the rollouts in this memory are not guaranteed to come from the latest policy, they are reweighted at every iteration using importance sampling. The global workers asynchronously read from the replay memory and apply updates to the global policy. By decoupling the local and global work, the robots can now continuously collect data by executing rollouts, while the global policy is optimized in the background. This system also makes it easy to add multiple robots into the training process, by adding additional local workers for every robot.

The global policy itself can be represented with any function approximator, but in our work, as in prior GPS methods, we use a deep neural network representation trained with stochastic gradient descent (SGD), which can itself be trained in a distributed manner. The global policy is stored on a parameter server [23], allowing multiple robots to concurrently collect data while multiple machines concurrently apply updates to the same global policy. By utilizing more robots, we are able to achieve much greater data diversity than would otherwise be realized with only a single robot, and by using multiple global worker threads, we can accelerate global policy training.

The replay memory may be either centralized or dis-

**Algorithm 2** Asynchronous distributed guided policy search (local worker)

1: **for** iteration $k \in \{1, \ldots, K\}$ **do**
2:     Generate sample trajectories starting from each $\mathbf{x}_1^i$ assigned to this worker, by executing either $p_i(\mathbf{u}_t|\mathbf{x}_t)$ or $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$ on the robot.
3:     Use samples to optimize each of the local policies $p_i(\mathbf{u}_t|\mathbf{x}_t)$ from each $\mathbf{x}_1^i$ with respect to Equation (2) or (3), using either LQR or PI$^2$.
4:     Append optimized trajectories $p_i(\mathbf{u}_t|\mathbf{x}_t)$ to replay memory $\mathcal{D}$.
5: **end for**

---

**Algorithm 3** Asynchronous distributed guided policy search (global worker)

1: **for** step $n \in \{1, \ldots, N\}$ **do**
2:     Randomly sample a mini-batch $\{\mathbf{x}_t^j\}$ from the replay memory $\mathcal{D}$, with corresponding labels obtained from the corresponding local policies $p_i(\mathbf{u}_t|\mathbf{x}_t^j)$, where $i$ is the instance from which sample $j$ was obtained.
3:     Optimize the global policy $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$ on this mini-batch for one step of SGD with respect to Equation (1) or (4).
4: **end for**

---

tributed. In our implementation of this system, each physical machine connected to each physical robot maintains its own replay memory. This is particularly convenient if we also run a single global worker thread on each physical machine, since it removes the need to transmit the high-bandwidth rollout data between machines during training. Instead, the machines only need to communicate model parameter updates to the centralized parameter server, which are typically much smaller than images or high-frequency joint angle and torque trajectories. In this case, the only centralized element of this system is the parameter server. Furthermore, since mini-batch gradient descent assumes uncorrelated examples within each batch, we found that distributed training actually improved stability when aggregating gradients across multiple robots and machines [24].

This entire system, which we call asynchronous distributed guided policy search (ADGPS), was implemented in the distributed machine learning framework TensorFlow [25], and is summarized in Algorithms 2 and 3. In our implementation, rollout execution and local policy optimization are still performed sequentially on the local worker as the optimization is a relatively cheap step; however, this is not strictly necessary and both steps could also be performed asynchronously. It is also possible to instantiate this system with varying numbers of global and local workers, or even a single centralized global worker. However, as discussed above, associating a single global worker with each local worker allows us to avoid transmitting the rollout data between machines, leading to a particularly efficient and convenient implementation.

## V. EXPERIMENTAL EVALUATION

Our experimental evaluation aims to answer two questions about our proposed asynchronous learning system: (1) does distributed asynchronous learning accelerate the training of complex, nonlinear neural network policies, and (2) does training across an ensemble of multiple robots improve the generalization of the resulting policies. The answers to these questions are not trivial: although it may seem natural that parallelizing experience collection should accelerate learning, it is not at all clear whether the additional bias introduced by asynchronous training would not outweigh the benefit of greater dataset size, nor that the amount of data is indeed the limiting factor.

### A. Simulated Evaluation

In simulation, we can systematically vary the number of robots to evaluate how training times scale with worker count, as well as study the effect of asynchronous training. We simulated multiple 7-DoF arms with parallelized simulators that each run in real time, in order to mimic rollout execution times that would be observed on real robots. The arms are controlled with torque control in order to perform a simple Cartesian reaching task that requires placing the end-effector at a commanded position. The robot state vector consists of joint angles and velocities, as well as its end-effector pose and velocity. We use a 9-DoF parameterization of pose, containing the positions of three points rigidly attached to the robot end-effector represented in the base frame. The Cartesian goal pose uses the same representation, and is fed to the global policy along with the robot state. The global policy must be able to place the end-effector at a variety of different target positions, with each instance of the task corresponding to a different target. We train the policy using 8 instances of the task, using 4 additional instances as a validation set for hyperparameter tuning, and finally test the global policy on 4 held-out instances. These experiments use the LQR variant of BADMM-based GPS.

We ran guided policy search with and without asynchrony, and with increasing numbers of workers from 1 to 8. Figure 4 shows the average costs across four test instances for each setting of the algorithm, plotted against the number of trials and wall-clock time, respectively. ADGPS-4 and ADGPS-8 denote 4 and 8 pairs of local and global workers, respectively, while AGPS is an asynchronous run with a single pair of workers. Note that asynchronous training does slightly reduce the improvement in cost per iteration, since the local policies are updated against an older version of the global policy, and the global policy is trained on older data. However, the iterations themselves take less time, since the global policy training is parallelized with data collection and local policy updates. This substantially improves the learning rate in terms of wall clock time. This is illustrated in Figure 5, which shows the relative improvement in wall-clock time (labeled as "speedup") compared to standard GPS, as well as the relative increase in sample complexity (labeled as "sample count") due to the slightly reduced policy improvement per iteration.
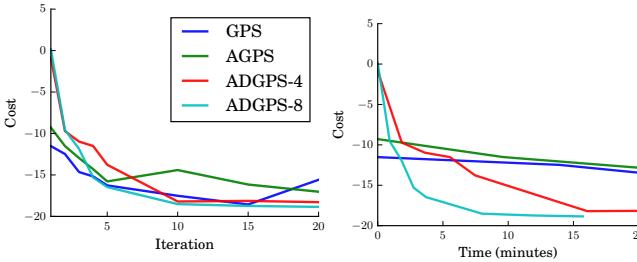
Fig. 4. Average costs of the 4 test instances used in the simulated reaching task, over number of iterations (*left*) as well as training duration (*right*). ADGPS-4 and ADGPS-8 denote 4 and 8 pairs of local and global workers, respectively, while AGPS is an asynchronous run with a single pair of workers. Note that asynchronous training does slightly reduce the improvement per iteration, but substantially improves training time when multiple workers are used.
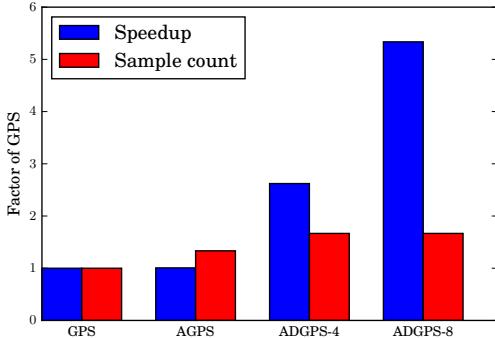


Fig. 5. Speedup in wall-clock training time and sample count comparison between GPS and the asynchronous variants, measured as the wall-clock time or sample count needed to reach a threshold cost value. Note that additional asynchronous workers incur only a modest cost in total sample count, while providing a substantial improvement in wall-clock training time.

### B. Real-World Evaluation

Our real-world evaluation is aimed at determining whether our distributed asynchronous system can effectively learn complex, nonlinear neural network policies, using visual inputs, and whether the resulting policies can generalize more effectively than policies learned on a single robot platform using the standard synchronous variant of GPS. To that end, we tackle a challenging real-world door opening task (Figure 6), where the goal is to train a single visuomotor policy that can open a range of doors with visual and mechanical differences in the handle (Figure 7), while also
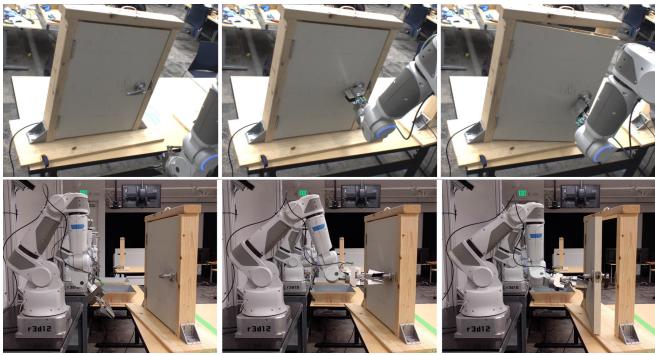


Fig. 6. Door task execution. Top: sample robot RGB camera images used to control the robot. Bottom: side view of one of the robots opening a door.



Fig. 7. Variation in door handles used in the experiment described in Section V-B. The three handles on the left are used during training, and the handle on the right is used for evaluation.

dealing with variations in the pose of the door with respect to the robot, variations in camera calibration, and mechanical variations between robots themselves.

We use four lightweight torque-controlled 7-DoF robotic arms, each of which is equipped with a two finger gripper, and a camera mounted behind the arm looking over the shoulder. The poses of these cameras are not precisely calibrated with respect to each robot. The input to the policy consists of monocular RGB images and the robot state vector as described in Section V-A. The robots are controlled at a frequency of 20Hz by directly sending torque commands to all seven joints. Each robot is assigned a specific door for policy training. The cost function is computed based on an IMU sensor attached to the door handle on the reverse side of the door. The desired IMU readings, which correspond to a successfully opened door, are recorded during kinesthetic teaching of the opening motion from human demonstration. We additionally add quadratic terms on joint velocities and commanded torques multiplied by a small constant to encourage smooth motions.

The architecture of the neural network policy we use is shown in Figure 8. Our architecture resembles prior work [7], with the visual features represented by feature points produced via a spatial softmax applied to the last convolutional response maps. Unlike in [7], our convolutional network includes multiple stages of pooling and skip connections, which allows the visual features to incorporate information at various scales: low-level, high-resolution, local features as well as higher-level features with larger spatial context. This allows the network to generate high resolution features while limiting the amount of computation performed at high resolutions, enabling evaluation of this deep model at camera frame rates.

*1) Policy pre-training:* We train the above neural network policy in two stages. First, the convolutional layers are pretrained with a proxy pose detection objective. To create data for this pretraining phase, we collect camera images while manually moving each of the training doors into various poses, and automatically label each image by using a geometry-based pose estimator based on the point pair feature (PPF) algorithm [27]. We also collect images of each robot learning the task with PI$^2$ (without vision), and label these images with the pose of the robot end-effector obtained from forward kinematics. Each pose is represented as a 9-DoF vector, containing the positions of three points rigidly attached to the object (or robot), represented in the world frame. The door poses are all labeled in the camera frame, which allows us to pool this data across robots into
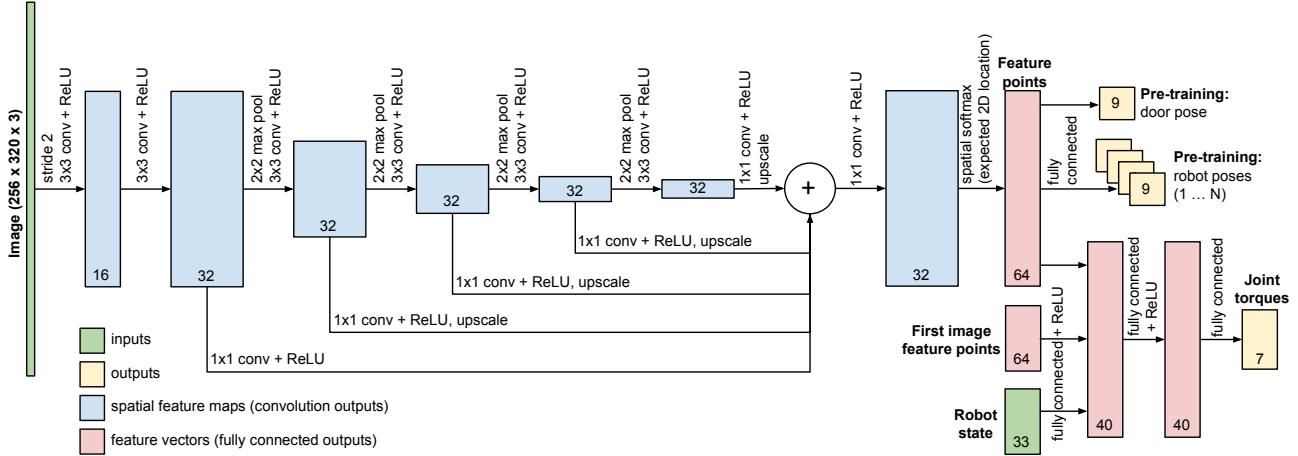
Fig. 8. The architecture of our neural network policy. The input RGB image is passed through a 3x3 convolution with stride 2 to generate 16 features at a lower resolution. The next 5 layers are 3x3 convolutions followed by 2x2 max-pooling, each of which output 32 features at successively reduced resolutions and increased receptive field. The outputs of these 5 layers are recombined by passing each of them into a 1x1 convolution, converting them to a size of 125x157 by using nearest-neighbor upscaling, and summation (similar to [26]). A final 1x1 convolution is used to generate 32 feature maps. The spatial soft-argmax operator [7] computes the expected 2D image coordinates of each feature. A fully connected layer is used to compute the object and robot poses from these expected 2D feature coordinates for pre-training the vision layers. The feature points for the current image are concatenated with feature points from the image at the first timestep as well as the 33-dimensional robot state vector, before being passed through two fully connected layers to produce the output joint torques.

a single dataset. However, since the robot endeffector poses are labeled in the base frame of each robot with an unknown camera offset, we cannot trivially train a single network to predict the pose label of any robot from the camera image alone. Hence, the pose of each robot is predicted using a separate output using a linear mapping from the feature points. This ensures that the 2-D image features learnt to predict the robot and door poses can be shared across all robots, while the 3-D robot pose predictions are allowed to vary across robots. The convolutional layers are trained using stochastic gradient descent (SGD) with momentum to predict the robot and door poses, using a standard Euclidean loss.

*2) Policy learning:* The local policy for each robot is initialized from its provided kinesthetic demonstration. We bootstrap the fully connected layers of the network by running four iterations of BADMM-based ADGPS with the $PI^2$ local policy optimizer. The pose of each door is kept fixed during bootstrapping. Next, we run 16 iterations of asynchronous distributed MDGPS with $PI^2$, where we randomly perturb each door pose at the start of every iteration. This sampling procedure allows us to train the global policy on a greater diversity of initial conditions, resulting in better generalization. The weights of the convolutional layers are kept frozen during all runs of GPS. In future work, it would be straightforward to also fine-tune the convolutional layers end-to-end with guided policy search as in prior work [7], but we found that we could obtain satisfactory performance without end-to-end training of the vision layers on this task.

*3) Results:* The trained policy was evaluated on a test door not seen during training time, by executing 50 trials per robot over a grid of translations and orientations. Figure 9 shows the results obtained using the policy after training. We find that all four robots are able to open the test door in most configurations using a single global policy, showing generalization over appearance and mechanical properties of

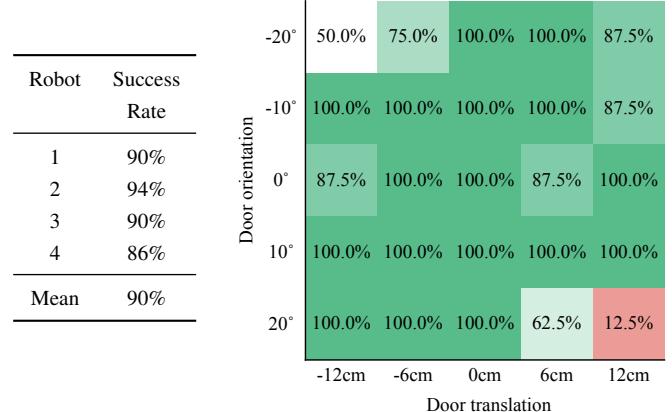| Robot | Success Rate |
|-------|--------------|
| 1 | 90% |
| 2 | 94% |
| 3 | 90% |
| 4 | 86% |
| Mean | 90% |



Fig. 9. Results from evaluating the visuomotor policy trained using ADGPS, using 50 trials per robot on a test door whose translations and orientations are sampled on a grid. *Left*: Success rates per robot averaged over the sampling grid. *Right*: Aggregate success rates across all robots for varying translations and orientations.

door handles, door position and orientation, camera calibration, and variations in robot dynamics. The lack of precise camera calibration per robot implies that the policy needs to visually track the pose of the door handle and the robot gripper, servoing it to the grasp pose. This is evident when watching the robot execute the learned policy (see video attachment): the initial motion of the robot brings the gripper into the view of the camera, after which the robot is able to translate and orient the gripper to grasp the handle before opening the door.

Furthermore, the trained policy was also evaluated with two test camera positions on the test door. The first camera position was arrived at by displacing the camera of one of the robots towards the ground by 5cm. The second position was arrived at by displacement that same camera away from the door by 4cm. The trained policy had success rates of 52% and 54% respectively with these two camera positions.

In comparison, a successful policy that was trained on only a single robot and a single door using GPS with PI$^2$ as in [19] fails to generalize to either an unseen door or different camera positions.

## VI. Discussion and Future Work

We presented a system for distributed asynchronous policy learning across multiple robots that can collaborate to learn a single generalizable motor skill, represented by a deep neural network. Our method extends the guided policy search algorithm to the asynchronous setting, where maximal robot utilization is achieved by parallelizing policy training with experience collection. The robots continuously collect new experience and add it to a replay buffer that is used to train the global neural network policy. At the same time, each robot individually improves its local policy to succeed on its own particular instance of the task. Our simulated experiments demonstrate that this approach can reduce training times, while our real-world evaluation shows that a policy trained on multiple instances of different doors can improve the generalization capability of a vision-based door opening policy.

Our method also assumes that each robot can execute the same policy, which implicitly involves the assumption that the robots are physically similar or identical. An interesting future extension of our work is to handle the case where there is a systematic discrepancy between robotic platforms, necessitating a public and private component to each policy. In this case, the private components would be learned locally, while the public components would be trained using shared experience and pooled across robots. This could allow distributed asynchronous training to extend even to heterogeneous populations of robots, where highly dissimilar robots might share globally useful experience, such as the statistics of natural images, while robot-specific knowledge about, for example, the details of low-level motor control, would be shared only with physically similar platforms.

## References

[1] R. Tedrake, T.W. Zhang, and H.S. Seung. Stochastic policy gradient reinforcement learning on a simple 3d biped. In *IROS*, 2004.

[2] J. Peters, K. Mülling, and Y. Altun. Relative entropy policy search. In *AAAI*, 2010.

[3] E. Theodorou, J. Buchli, and S. Schaal. A generalized path integral control approach to reinforcement learning. *Journal of Machine Learning Research*, 11:3137–3181, 2010.

[4] M.P. Deisenroth, G. Neumann, and J. Peters. A survey on policy search for robotics. *Foundations and Trends in Robotics*, 2(1-2):1–142, 2013.

[5] J. Schulman, S. Levine, P. Moritz, M. Jordan, and P. Abbeel. Trust region policy optimization. In *ICML*, 2015.

[6] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2016.

[7] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17 (1):1334–1373, 2016.

[8] J.A. Ijspeert, J. Nakanishi, and S. Schaal. Movement imitation with nonlinear dynamical systems in humanoid robots. In *ICRA*, 2002.

[9] Yann Lecun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 5 2015. ISSN 0028-0836. doi: 10.1038/nature14539.

[10] Thomas Lampe and Martin Riedmiller. Acquiring visual servoing reaching and grasping skills using neural reinforcement learning. In *IEEE International Joint Conference on Neural Networks (IJCNN 2013)*, Dallas, TX, 2013.

[11] M. Inaba, S. Kagami, F. Kanehiro, and Y. Hoshino. A platform for robotics research based on the remote-brained robot approach. *International Journal of Robotics Research*, 19(10), 2000.

[12] J. Kuffner. Cloud-enabled humanoid robots. In *IEEE-RAS International Conference on Humanoid Robotics*, 2010.

[13] B. Kehoe, A. Matsukawa, S. Candido, J. Kuffner, and K. Goldberg. Cloud-based robot grasping with the google object recognition engine. In *IEEE International Conference on Robotics and Automation*, 2013.

[14] B. Kehoe, S. Patil, P. Abbeel, and K. Goldberg. A survey of research on cloud robotics and automation. *IEEE Transactions on Automation Science and Engineering*, 12(2), April 2015.

[15] Martin Zinkevich, Markus Weimer, Lihong Li, and Alex J. Smola. Parallelized stochastic gradient descent. In *Neural Information Processing Systems (NIPS)*, 2010.

[16] Igor Mordatch, Kendall Lowrey, Galen Andrew, Zoran Popovic, and Emanuel V Todorov. Interactive control of diverse complex characters with neural networks. In *Advances in Neural Information Processing Systems*, pages 3132–3140, 2015.

[17] W. Montgomery and S. Levine. Guided policy search as approximate mirror descent. In *NIPS*, 2016.

[18] S. Levine and P. Abbeel. Learning neural network policies with guided policy search under unknown dynamics. In *NIPS*, 2014.

[19] Y. Chebotar, M. Kalakrishnan, A. Yahya, A. Li, S. Schaal, and S. Levine. Path integral guided policy search. In *Technical Report*, 2017.

[20] H. Wang and A. Banerjee. Bregman alternating direction method of multipliers. In *NIPS*, 2014.

[21] Christian Daniel, Gerhard Neumann, and Jan Peters. Hierarchical relative entropy policy search. In *AISTATS*, pages 273–281, 2012.

[22] W. Li and E. Todorov. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *ICINCO*, pages 222–229, 2004.

[23] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, et al. Large scale distributed deep networks. In *Advances in neural information processing systems*, pages 1223–1231, 2012.

[24] Jianmin Chen, Rajat Monga, Samy Bengio, and Rafal Jozefowicz. Revisiting distributed synchronous sgd. In *International Conference on Learning Representations Workshop Track*, 2016.

[25] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous distributed systems, 2015. URL http://tensorflow.org/.

[26] J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *NIPS*, 2014.

[27] S. Hinterstoisser, V. Lepetit, N. Rajkumar, and K. Konolige. Going further with point pair features. In *ECCV*, 2016.