# Zero-Shot Skill Composition and Simulation-to-Real Transfer by Learning Task Representations

Zhanpeng He<sup>\*</sup>, Ryan Julian<sup>\*</sup>, Eric Heiden, Hejia Zhang, Stefan Schaal, Joseph J. Lim, Gaurav Sukhatme, and Karol Hausman

Abstract-Simulation-to-real transfer is an important strategy for making reinforcement learning practical with real robots. Successful sim-to-real transfer systems have difficulty producing policies which generalize across tasks, despite training for thousands of hours equivalent real robot time. To address this shortcoming, we present a novel approach to efficiently learning new robotic skills directly on a real robot, based on model-predictive control (MPC) and an algorithm for learning task representations. In short, we show how to reuse the simulation from the pre-training step of sim-to-real methods as a tool for foresight, allowing the sim-to-real policy adapt to unseen tasks. Rather than end-to-end learning policies for single tasks and attempting to transfer them, we first use simulation to simultaneously learn (1) a continuous parameterization (i.e. a skill embedding or latent) of task-appropriate primitive skills, and (2) a single policy for these skills which is conditioned on this representation. We then directly transfer our multiskill policy to a real robot, and actuate the robot by choosing sequences of skill latents which actuate the policy, with each latent corresponding to a pre-learned primitive skill controller. We complete unseen tasks by choosing new sequences of skill latents to control the robot using MPC, where our MPC model is composed of the pre-trained skill policy executed in the simulation environment, run in parallel with the real robot. We discuss the background and principles of our method, detail its practical implementation, and evaluate its performance by using our method to train a real Sawyer Robot to achieve motion tasks such as drawing and block pushing.

## I. INTRODUCTION

Reinforcement learning algorithms have been proven to be effective to learn complex skills in simulation environments in [1], [2] and [3]. However, practical robotic reinforcement learning for complex motion skills remains a challenging and unsolved problem, due to the high number of samples needed to train most algorithms and the expense of obtaining those samples from real robots. Most existing approaches to robotic reinforcement learning either fail to generalize between different tasks and among variations of single tasks, or only generalize by requiring collecting impractical amounts of real robot experience. With recent advancements in robotic simulation, and the widespread availability of large computational resources, a popular family of methods seeking to address this challenge has emerged, known as "sim-to-real" methods. These methods seek to offload most

\* These authors contributed equally to the paper.



Fig. 1. The Sawyer robot performing the reaching task in simulation (left) and real world (right)



Fig. 2. The Sawyer robot performing the box pushing task in simulation (left) and real world (right)

training time from real robots to offline simulations, which are trivially parallelizable and much cheaper to operate. Our method combines this "sim-to-real" schema with representation learning and model-predictive control (MPC) to make transfer more robust, and to significantly decrease the number of simulation samples needed to train policies which achieve families of related tasks.

The key insight behind our method is that the simulation used in the pre-training step of a simulation-to-real method can also be used online as a tool for foresight. It allows us to predict the behavior of a known policy on an unseen task. When combined with a latent-conditioned policy, where the latent actuates variations of useful policy behavior (e.g. skills), this simulation-as-foresight tool allows our method to use what the robot has already learned to do (e.g. the pre-trained policy) to bootstrap online policies for tasks it has never seen before. That is, given a latent space of useful behaviors, and a simulation which predicts the rewards for those behaviors on a new task, we can reduce the adaptation problem to intelligently choosing a sequence of latent skills

Zhanpeng He, Ryan Julian, Eric Heiden, Hejia Zhang, Stefan Schaal, Joseph Lim and Gaurav Sukhatme are with the Department of Computer Science, University of Southern California, Los Angeles, USA. {zhanpenh, rjulian, heiden, hejia, sschaal, limjj, gaurav}@usc.edu

Karol Hausman is with Google Brain, Mountain View, CA, USA karolhausman@google.com

which maximize rewards for the new task.

## II. BACKGROUND

Most simulation-to-real approaches so far have focused on addressing the "reality gap" problem. The reality gap problem is the domain shift performance loss induced by differences in dynamics and perception between the simulation (policy training) and real (policy execution) environments. Training a policy only in a flawed simulation generally yields control behavior which is not adaptable to even small variations in the environment dynamics. Furthermore, simulating the physics behind many practical robotic problems (e.g. sliding friction and contact forces) is an open problem in applied mathematics, meaning it is not possible to solve a completely accurate simulation for many important robotic tasks [4]. Rather than attempt to create an explicit alignment between simulation and real [5], or randomize our simulation training to a sufficient degree to learn a policy which generalizes to nearby dynamics [6], our method seeks to learn a sufficient policy in simulation, and adapt it quickly to the real world online during real robot execution.

Our proposed approach is based on four key components: reinforcement learning with policy gradients (RL) [7], variational inference [8], model-predictive control (MPC), and physics simulation. We use variational inference to learn a low-dimensional latent space of skills which are useful for tasks, and RL to simultaneously learn single policy which is conditioned on these latent skills. The precise long-horizon behavior of the policy for a given latent skill is difficult to predict, so we use MPC and an online simulation to evaluate latent skill plans in the in simulation before executing them on the real robot.

#### **III. RELATED WORK**

Learning skill representations to aid in generalization has been proposed in works old and new. Previous works proposed frameworks such as Associative Skill Memories [9] and probabilistic movement primitives [10] to acquire a set of reusable skills. Our approach is built upon [11], which learns a embedding space of skills with reinforcement learning and variational inference, and [12] which shows that these learned skills are transferable and composable on real robots. While [12] noted that predicting the behavior of latent skills is an obstacle to using this method, our approach addresses the problem by using model-predictive control to successfully complete unseen tasks with no fine-tuning on the real robot. Exploration is a key problem in robot learning, and our method uses latent skill representations to address this problem. Using learned latent spaces to make exploration more tractable is also studied in [13] and [14]. Our method exploits a latent space for task-oriented exploration: it uses model-predictive control and simulation to choose latent skills which are locally-optimal for completing unseen tasks, then executes those latent skills on the real robot.

Using reinforcement learning with model-predictive control has been explored previously. Kamthe *et al.* [15] proposed using MPC to increase the data efficiency of reinforcement algorithms by training probabilistic transition models for planning. In our work, we take a different approach by exploiting our learned latent space and simulation directly to find policies for novel tasks online, rather than learning and then solving a model.

Simulation-to-real transfer learning approaches include randomizing the dynamic parameters of the simulation [6], and varying the visual appearance of the environment [16], both of which scale linearly or quadratically the amount of computation needed to learn a transfer policy. Other strategies, such as that of Barrett et al. [17] reuse models trained in simulation to make sim-to-real transfer more efficient, similar to our method, however this work requires an explicit pre-defined mapping between seen and unseen tasks. Saemundson et al. [18] use meta-learning and learned representations to generalize from pre-trained seen tasks to unseen tasks, however their approach requires that the unseen tasks be very similar to the pre-trained tasks, and is few-shot rather than zero-shot. Our method is zero-shot with respect to real environment samples, and can be used to learn unseen tasks which are significantly out-of-distribution, as well as for composing learned skills in the time domain to achieve unseen tasks which are more complex than the underlying pre-trained task set.

Our work is closely related to simultaneous work performed by Co-Reyes *et al.* [19]. Whereas our method learns an explicit skill representations using pre-chosen skills identified by a known ID, [19] learn an implicit skill representation by clustering trajectories of states and rewards in a latent space. Furthermore, we focus on MPC-based planning in the latent space to achieve robotic tasks learned online with a real robot, while their analysis focuses on the machine learning behind this family of methods and uses simulation experiments.

# IV. METHOD

## A. Skill Embedding Algorithm

In our multi-task RL setting, we pre-define a set of lowlevel skills with IDs  $\mathcal{T} = \{1, \ldots, N\}$ , and accompanying, per-skill reward functions  $r^t(s, a)$ .

In parallel with learning the joint low-level skill policy  $\pi_{\theta}$  as in conventional RL, we learn an embedding function  $p_{\phi}$  which parameterizes the low-level skill library using a latent variable z. Note that the true skill identity t is hidden from the policy behind the embedding function  $p_{\phi}$ . Rather than reveal the skill ID to the policy, once per rollout we feed the skill ID t, encoded as s one-hot vector, through the stochastic embedding function  $p_{\phi}$  to produce a latent vector z. We feed this same value of z to the policy for the entire rollout, so that all steps in a trajectory are correlated with the same value of z.

$$\mathfrak{L}(\theta, \phi, \psi) = \max_{\pi} \mathbb{E}_{\pi(a, z|s, t)} \left[ \sum_{i=0}^{\infty} \gamma^{i} \hat{r}(s_{i}, a_{i}, z, t) \middle| s_{i+1} \right]$$
(1)

where

$$\hat{r}(s_i, a_i, z, t) = \alpha_1 \mathbb{E}_{t \in T} \left[ \mathbb{H} \left( p_\phi(z|t) \right) \right] + \alpha_2 \log q_\psi(z|s_i^H) + \alpha_3 \mathbb{H} \left( \pi_\theta(a_i|s_i, z) \right) + r^t(s_i, a_i)$$

To aid in learning the embedding function, we learn an inference function  $q_{\psi}$  which, given a state-only trajectory window  $s_i^H$  of length H, predicts the latent vector z which was fed to the low-level skill policy when it produced that trajectory. This allows us to define an augmented reward which encourages the policy to produce distinct trajectories for different latent vectors. We learn  $q_{\psi}$  in parallel with the policy and embedding functions, as shown in Eq. 1.

We add a policy entropy bonus  $\mathbb{H}(\pi_{\theta}(a_i|s_i, z))$ , which ensures that the policy does not collapse to a single solution for each skill. For a detailed derivation, refer to [11].

# B. Skill Embedding Criterion

In order for the learned latent space to be useful for completing unseen tasks, we seek to constrain the embedding distribution to satisfy two important properties:

- 1) **High entropy:** Each task should induce a distribution over latent vectors which is wide as possible, corresponding to many variations of a single skill.
- Identifiability: Given an arbitrary trajectory window, the inference network should be able to predict with high confidence the latent vector fed to the policy to produce that trajectory.

When applied together, these properties ensure that during training the policy is trained to encode high-reward controllers for many parameterizations of a skill (highentropy), while simultaneously ensuring that each of these latent parameterizations corresponds to a distinct variation of that skill. This dual constraint is the key for using model predictive control or other composing methods in the latent space as discussed in Sec. IV-C.



Fig. 3. Skill Embedding Algorithm and MPC

We train the policy and embedding networks using Proximal Policy Optimization [20], though our method may be used by any parametric reinforcement learning algorithm. We use the MuJoCo physics engine [21] to implement our Sawyer robot simulation environments. We represent the policy, embedding, and inference functions using multivariate Gaussian distributions whose mean and diagonal covariance are parameterized by the output of a multi-layer perceptron. The policy and embedding distributions are jointly optimized by the reinforcement learning algorithm, while we train the inference distribution using supervised learning and a simple cross-entropy loss.

## C. Using Model Predictive Control for Zero-Shot Adaptation

To achieve unseen tasks on a real robot with no additional training, we freeze the multi-skill policy learned in Sec. IV-A, and use a new algorithm which we refer to as a "composer." The composer achieves unseen tasks by choosing new sequences of latent skill vectors to feed to the frozen skill policy. Exploring in this smaller space is faster and more sample-efficient, because it encodes highlevel properties of tasks and their relations. Each skill latent induces a different pre-learned behavior, and our method reduces the adaptation problem to choosing sequences of these pre-learned behaviors–continuously parameterized by the skill embedding–to achieve new tasks.

Note that we use the simulation itself to evaluate the future outcome of the next action. For each step, we set the state of the simulation environment to the observed state of the real environment. This equips our robot with with the ability to predict the behavior of different skill latents. Since our robot is trained in a simulation-to-real framework, we can reuse the simulation from the pre-training step as a tool for foresight when adapting to unseen tasks. This allow us to select a latent skill online which is locally-optimal for a task, even if that task was seen not during training. We show that this scheme allows us to perform zero-shot task execution and composition for families of related tasks. This is in contrast to existing methods, which have mostly focused on direct alignment between simulation and real, or data augmentation to generalize the policy using brute force. Despite much work on simulation-to-real methods, neither of these approaches has demonstrated the ability to provide the adaptation ability needed for general-purpose robots in the real world. We believe our method provides a third path towards simulationto-real adaptation that warrants exploration, as a higher-level complement to these effective-but-limited existing low-level approaches.

We denote the new task  $t^{\text{new}}$  corresponding to reward function  $r^{\text{new}}$ , the real environment in which we attempt this task  $\mathcal{R}(s'|s, a)$ , and the RL discount factor  $\gamma$ . We use the simulation environment  $\mathcal{S}(f'|f, \dashv)$ , frozen skill embedding  $p_{\phi}(z|t)$ , and latent-conditioned skill policy  $\pi_{\theta}(a|s, z)$ , all trained in Sec. IV-A, to apply model-predictive control in the latent space as follows (Algorithm 1).

We first sample k candidate latents  $\mathcal{Z} = \{z_1, \ldots, z_k\}$ according to  $p(z) = \mathbb{E}_{t \sim p(t)} p_{\phi}(z|t)$ . We observe the state  $s_{\text{real}}$  of real environment  $\mathcal{R}$ .

For each candidate latent  $z_i$ , we set the initial state of the simulation S to  $s_{\text{real}}$ . For a horizon of T time steps, we sample the frozen policy  $\pi_{\theta}$ , conditioned on the candidate latent  $a_{j \in T} \sim \pi_{\theta}(a_j | s_j, z_i)$ , and execute the actions  $a_j$ the simulation environment S, yielding and total discounted



Fig. 4. Using model-predictive control with embedding functions and multi-task policy

reward  $R_i^{\text{new}} = \sum_{j=0}^T \gamma^j r^{\text{new}}(s_j, a_j)$  for each candidate latent. We then choose the candidate latent acquiring the highest reward  $z^* = \operatorname{argmax}_i R_i^{\text{new}}$ , and use it to condition and sample the frozen policy  $a_{l \in N} \sim \pi_{\theta}(a_j | s_j, z^*)$  to control the real environment  $\mathcal{R}$  for a horizon of N < T time steps.

We repeat this MPC process to choose and execute new latents in sequence, until the task has been achieved.

# Algorithm 1 MPC in Skill Latent Space

**Require:** A latent-conditioned policy  $\pi_{\theta}(a|s, z)$ , a skill embedding distribution  $p_{\phi}(z|t)$ , a skill distribution prior p(t), a simulation environment S(s'|s, a), a real environment  $\mathcal{R}(s'|s, a)$ , a new task  $t^{\text{new}}$  with associated reward function  $r^{\text{new}}(s, a)$ , an RL discount factor  $\gamma$ , an MPC horizon T, and a real environment horizon N.

```
while t^{new} is not complete do
   Sample \mathcal{Z} = \{z_1, \dots, z_k\} \sim p(z) = \mathbb{E}_{t \sim p(t)} p_{\phi}(z|t)
   Observe s_{\text{real}} from \mathcal{R}
   for z_i \in \mathcal{Z} do
       Set initial state of S to s_{real}
       for j \in \{1, ..., T\} do
           Sample a_j \sim \pi_{\theta}(a_j | s_j, z_i)
           Execute simulation s_{j+1} = \mathcal{S}(s_j, a_j)
       end for
       Calculate R_i^{\text{new}} = \sum_{j=0}^T \gamma^j r^{\text{new}}(s_j, a_j)
   end for
   Choose z^* = \operatorname{argmax}_{z_i} R_i^{\text{new}}
   for l \in \{1, ..., N\} do
       Sample a_l \sim \pi_{\theta}(a_l | s_l, z^*)
       Execute real environment s_{l+1} = \mathcal{R}(s_l, a_j)
   end for
end while
```

The choice of MPC horizon T has a significant effect on the performance of our approach. Since our latent variable encodes a skill which only partially completes the task, executing a single skill for too long unnecessarily penalizes a locally-useful skill for not being globally optimal. Hence, we set the MPC horizon T to not more than twice the number of steps that a latent is actuated in the real environment N.

#### V. EXPERIMENTS

We evaluate our approach by completing two sequencing tasks on a Sawyer robot: drawing a sequence of points and pushing a box along a sequential path. For each of the experiments, the robot must complete an overall task by sequencing skills learned during the embedding learning process. Sequencing skills poses a challenge to conventional RL algorithms due to the sparsity of rewards in sequencing tasks [22]. Because the agent only receives a reward for completing several correct complex actions in a row, exploration under these sequencing tasks is very difficult for conventional RL algorithms. By reusing the skills we have consolidated in the embedding space, we show a high-level controller can effectively compose these skills in order to achieve such difficult sequencing tasks.

#### A. Sawyer: Drawing a Sequence of Points

In this experiment, we ask the Sawyer Robot to move its end-effector to a sequence of points in 3D space. We first learn the low level policy that receives an observation with the robot's seven joint angles as well as the Cartesian position of the robot's gripper, and output incremental joint positions (up to 0.04 rads) as actions. We use the Euclidean distance between the gripper position and the current target is used as the cost function. We trained the policy and the embedding network on eight goal positions in simulation, forming a 3D rectoid enclosing the workspace. Then, we use the model-predictive control to choose a sequence latent vector which allows the robot to draw an unseen shape. For both simulation and real robot experiments, we attempted two unseen tasks: drawing a rectangle in 3D space (Figs. 5 and 7) and drawing a triangle in 3D space (Figs. 6 and 8).



Fig. 5. Gripper position plots for the unseen rectangle-drawing experiment in simulation. In this experiment, the unseen task is drawing a rectangle in 3D space.



Fig. 6. Gripper position plots in unseen triangle-drawing experiment in simulation. In this experiment, the unseen task is to move the gripper to draw a triangle.



Fig. 7. Gripper position plots for the triangle-drawing experiment on the real robot. In this experiment, the unseen task is to draw a triangle.



Fig. 8. Gripper position plots in unseen triangle-drawing experiment on the real robot. In this experiment, the unseen task it to move the gripper to draw a triangle.

B. Sawyer: Pushing the Box through a Sequence of Waypoints

In this experiment, we test our approach with a task that requires contact between the Sawyer Robot and an object. We ask the robot to push a box along a sequence of points in the table plane. We choose the Euclidean distance between the position of the box and the current target position as the reward function. The policy receives a state observation with the relative position vector between the robot's gripper and the box's centroid and outputs incremental gripper movements (up to  $\pm 0.03$  cm) as actions.

We first pre-train a policy to push the box to four goal locations relative to its starting position in simulation. We trained the low-level multi-task policy with four tasks in simulation: 20 cm up, down, left, and right of the box starting position. We then use the model-predictive control to choose a latent vectors and feed it with the state observation to frozen multi-task policy which controls the robot.

For both simulation and real robot experiments, we use the simulation as a model of the environments. In the simulation experiments, we use model-predictive controller to push the box to three points. In the real robot experiments, we ask the Sawyer Robot to complete two unseen tasks: pushing up-then-left and pushing left-then-down.



Fig. 9. Plot of block positions and gripper positions in simulation experiments. In the first experiment (left), the robot pushes the box to the right, up and then left. In the second experiment (right), the robot pushes the box to the left, then up, and then back to its starting position.



Fig. 10. Plot of block positions and gripper positions in real robot experiments. In experiment I (left), the robot pushes the box to the left, and then down. In experiment II (right), the robot push the box to the up, and then left.

#### VI. RESULTS

## A. Sawyer Drawing

In the unseen drawing experiments, we sampled k = 15 vectors from the skill latent distribution, and for each of them performed an MPC optimization with a horizon of T = 4 steps. We then execute the latent with highest reward for N = 2 steps on the target robot. In simulation experiments, the Sawyer Robot successfully draw a rectangle with by sequencing 54 latents (Fig. 2) and drew by sequencing a triangle with 56 latents (Fig. 3). In the real robot experiments, the Sawyer Robot successfully completed the unseen rectangle-drawing task by choosing 62 latents (Fig. 4) in 2 minutes of real time and completed the unseen triangle-drawing task by choosing 53 latents (Fig. 5) in less than 2 minutes.

## B. Sawyer Pusher Sequencing

In the pusher sequencing experiments, we sample k = 50 vectors from the latent distribution. We use an MPC optimization with a simulation horizon of T = 30 steps, and execute each chosen latent in the environment for N = 10 steps. In simulation experiments, the robot completed the unseen up-left task less than 30 seconds of equivalent real time and the unseen right-up-left task less than 40 seconds of equivalent real time. In the real robot experiments, the

robot successfully completed the unseen left-down task by choosing 3 latents over approximately 1 minute of real time, and the unseen push up-left task by choosing 8 latents in about 1.5 minutes of real time.

## C. Analysis

These experiment results show that our learned skills are composable to complete the new task. In comparison with performing a search as done in [12], our approach is faster in wall clock time because we perform the model prediction in simulation instead of on the real robot. Note that our approach can utilize the continuous space of latents, whereas previous search methods only use an artificial discretization of the continuous latent space. In the unseen box-pushing real robot experiment (Fig. 7, *Right*), the Sawyer robot pushes the box towards the bottom-right right of the workspace to fix an error it made earlier in the task. This intelligent reactive behavior was never explicitly trained during the pre-training in simulation process. This shows that by sampling from our latent space, the model-predictive controller successfully selects a skill that is not pre-defined during training process.

## VII. CONCLUSION

In this work, we combine task representation learning simulation-to-real training, and model-predictive control to efficiently acquire policies for unseen tasks with no additional training. Our experiments show that applying model predictive control to these learned skill representations can be a very efficient method for online learning of tasks. The tasks we demonstrated are more complex than the underlying pre-trained skills used to achieve them, and the behaviors exhibited by our robot while executing unseen tasks were more adaptive than demanded by the simple reward functions us. Our method provides a partial escape from the reality gap problem in simulation-to-real methods, by mixing simulation-based long-range foresight with locally-correct online behavior.

For future work, we plan to apply our model-predictive controller as an exploration strategy to learn a composer policy that uses the latent space as action space. We look forward to efficiently learning a policy on real robots with guided exploration in our latent space.

## ACKNOWLEDGEMENTS

The authors would like to thank Angel Gonzalez Garcia, Jonathon Shen, and Chang Su for their work on the garage<sup>1</sup> reinforcement learning for robotics framework, on which the software for this work was based. We also want to thank the authors of multiworld<sup>2</sup> for providing a well-tuned Sawyer Block Pushing simulation environment. This research was supported in part by National Science Foundation grants IIS-1205249, IIS-1017134, EECS-0926052, the Office of Naval Research, the Okawa Foundation, and the Max-Planck-Society. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding organizations.

#### REFERENCES

- V. Mnih, K. Kavukcuoglu, D. Silver, A. Rusu, J. Veness, M. Bellemare, A. Graves, M. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Humanlevel control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015. [Online]. Available: https://www.nature.com/nature/journal/v518/n7540/pdf/nature1423 6.pdf
- [2] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016. [Online]. Available: https://doi.org/10.1038/nature16961
- [3] A. Yahya, A. Li, M. Kalakrishnan, Y. Chebotar, and S. Levine, "Collective robot reinforcement learning with distributed asynchronous guided policy search," *CoRR*, vol. abs/1610.00673, 2016. [Online]. Available: http://arxiv.org/abs/1610.00673
- [4] N. Jakobi, P. Husbands, and I. Harvey, "Noise and the reality gap: The use of simulation in evolutionary robotics," in *Advances in Artificial Life*, F. Morán, A. Moreno, J. J. Merelo, and P. Chacón, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 1995, pp. 704–720.
- [5] A. A. Visser, N. Dijkshoorn, M. V. D. Veen, and R. Jurriaans, "Closing the gap between simulation and reality in the sensor and motion models of an autonomous ar.drone." [Online]. Available: http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.841.1746&re p=rep1&type=pdf
- [6] X. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," *CoRR*, vol. abs/1710.06537, 2017. [Online]. Available: http://arxiv.org/abs/1710.06537
- J. Schulman, P. Abbeel, and X. Chen, "Equivalence between policy gradients and soft q-learning," *CoRR*, vol. abs/1704.06440, 2017.
   [Online]. Available: http://arxiv.org/abs/1704.06440
- [8] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *CoRR*, vol. abs/1312.6114, 2013. [Online]. Available: http://arxiv.org/abs/1312.6114
- [9] P. Pastor, M. Kalakrishnan, L. Righetti, and S. Schaal, "Towards associative skill memories," in *Humanoids*, Nov 2012.
- [10] E. Rueckert, J. Mundo, A. Paraschos, J. Peters, and G. Neumann, "Extracting low-dimensional control variables for movement primitives," in *ICRA*, May 2015.
- [11] K. Hausman, J. Springenberg, Z. Wang, N. Heess, and M. Riedmiller, "Learning an embedding space for transferable robot skills," in *ICLR*, 2018. [Online]. Available: https://openreview.net/forum?id=rk07ZXZRb
- [12] R. C. Julian, E. Heiden, Z. He, H. Zhang, S. Schaal, J. Lim, G. S. Sukhatme, and K. Hausman, "Scaling simulation-to-real transfer by learning composable robot skills," in *International Symposium on Experimental Robotics*. Springer, 2018. [Online]. Available: https://ryanjulian.me/iser\_2018.pdf
- [13] S. Gu, E. Holly, T. Lillicrap, and S. Levine, "Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates," in *ICRA*. IEEE, 2017.
- [14] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, "Diversity is all you need: Learning skills without a reward function," Feb. 2018. [Online]. Available: http://arxiv.org/abs/1802.06070
- [15] S. Kamthe and М. P. Deisenroth, "Data-efficient reinforcement learning with probabilistic model predictive control," CoRR, vol. abs/1706.06491, 2017. [Online]. Available: http://arxiv.org/abs/1706.06491
- [16] F. Sadeghi and S. Levine, "CAD2RL: Real single-image flight without a single real image," in RSS, 2017.
- [17] S. Barrett, M. E. Taylor, and P. Stone, "Transfer learning for reinforcement learning on a physical robot," in *Ninth International Conference on Autonomous Agents and Multiagent Systems - Adaptive Learning Agents Workshop (AAMAS - ALA)*, May 2010. [Online]. Available: http://www.cs.utexas.edu/users/ailab/?AAMASWS10-barrett

<sup>&</sup>lt;sup>1</sup>https://github.com/rlworkgroup/garage

<sup>&</sup>lt;sup>2</sup>https://github.com/vitchyr/multiworld

- [18] S. Sæmundsson, K. Hofmann, and M. P. Deisenroth, "Meta reinforcement learning with latent variable gaussian processes," *CoRR*, vol. abs/1803.07551, 2018. [Online]. Available: http://arxiv.org/abs/1803.07551
- [19] J. D. Co-Reyes, Y. Liu, A. Gupta, B. Eysenbach, P. Abbeel, and S. Levine, "Self-Consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings," Jun. 2018. [Online]. Available: http://arxiv.org/abs/1806.02813
- [Online]. Available: http://arxiv.org/abs/1806.02813
   [20] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *CoRR*, vol. abs/1707.06347, 2017. [Online]. Available: http://arxiv.org/abs/1707.06347
- [21] "MuJoCo: A physics engine for model-based control," in IROS, 2012.
- [22] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba, "Hindsight experience replay," in *NIPS*, 2017.