

# A Reduced-Order Approach to Assist with Reinforcement Learning for Underactuated Robotics

Jérémy Augot<sup>1,2</sup>, Aaron J. Snoswell<sup>2</sup> and Surya P. N. Singh<sup>2</sup>

**Abstract**—Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## I. INTRODUCTION

As robots become more compliant and variable they become more capable, and more complex. Underactuated robots, for example, provide great design freedom, yet their adoption has been limited by the need for manual controller designs.

Deep Reinforcement Learning (DRL) methods offer automated tools for robotic control by identifying a policy that maximizes the expected sum of rewards [1]. This has been extended to continuous control domains via algorithms such as the Deep Deterministic Policy Gradient (DDPG) [2], Proximal Policy Optimization (PPO) [3], Soft Actor-Critic (SAC) [4], and Twin Delayed Deep Deterministic (TD3) [5]. The power of these methods comes, in part, from the high-dimensional, nonlinear function approximation of both the policy and the expected value by neural networks. The dimensionality of these rich representations also presents a challenge for sample collection, stability, and convergence [6]. Reducing the order seems a natural approach to addressing this challenge, but doing so directly requires carefully defining suitable (state-space) features as the aforementioned methods are sensitive to the chosen representation [7].

Interestingly, many robots are underactuated by design, or at times, by circumstance (e.g., due to motor saturation). Such systems achieve their tasks due to the inherent dependence of their actuation states. This implicit structure suggests that an (automatic) reduction of the actuation space may make these systems more compatible with (model-free) DRL methods; and would allow more variable underactuated systems. Towards this, we introduce a highly variable, single actuator robotic locomotion benchmark based around a passive compliant toy system (see also Figure 1).

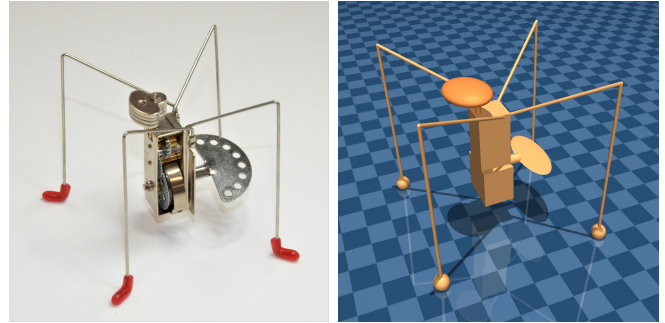


Fig. 1: The wind-up children’s toy ‘Katita’ (left) was the inspiration for our underactuated ‘Jitterbug’ continuous control task (right). In the simulated robot the wind-up spring is replaced with a controlled single degree-of-freedom motor. For scale, the blue checks on the simulated floor on the right are 1cm in size.

The contributions of this paper are ... We present a reduced order variant for this domain. We introduce the Jitterbug Problem as a diverse underactuated robotics task with dynamic, compliant and non-trivial motion in five scenarios of varying gradations of task complexity. And, we find that reduced order models can have similar (reward) performance, yet empirically have less training variance and slightly better learning rates.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a

<sup>1</sup>Ecole CentraleSupélec, Paris, France

<sup>2</sup>The Robotics Design Lab at The University of Queensland, Brisbane, Australia

nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

## II. RELATED WORK

Reinforcement Learning is an ....

Reinforcement Learning (RL) methods offer automated tools for controller design of such systems via trial and error [8], but have been limited by feature representations and forward models [9]. Deep Reinforcement Learning (DRL) algorithms for continuous control address this through the use of deep neural networks for both the policy ('actor') and value function ('critic') [2], [4]. Such systems have been ...

### A. Deep Reinforcement Learning

...and have shown promise in a host of challenging applications including visuomotor learning [10] and VR teleoperation [11].

### B. Reduced Order Approaches for Underactuated Robotics

Underactuated robots are fundamentally correlated in their actuation space [12]. When a model is available control strategies including from optimal control [13], direct collocation [14], and trajectory optimization [15] have been applied. Some novel robots, such as compliant legged robots or soft robots, tend to exhibit complex dynamics that do stymie dynamical system and model reduction approaches as an explicit physical model may not be available [16]. Learning the model empirically, such as via a neural network, typically centers around obtaining a model that forecasts the robot's next timestep which is then used subsequently for control, such as via Model Predictive Control (MPC) for a legged robot [17] or for manipulation with a soft-gripper [18].

Highly variable motion is not random. Given that actuation space for underactuated robots is structured, an autoencoder would be able to automatically discover some of these correlations [19], [20]. A concern, however, is that the reduced model is not always intuitive, which would complicate the subsequent synchronization with a model based control strategy. Model-free deep reinforcement learning (DRL) continuous control methods, however, are rather compatible with an autoencoder's reduced state description.

The use model reduction complements said DRL methods...

...This has been considered for reinforcement learning ...

This has been particularly effective in cases with very high dimensional inputs, such as from video, for low-dimensional independent constraints, such as obstacle locations [10], [21]....

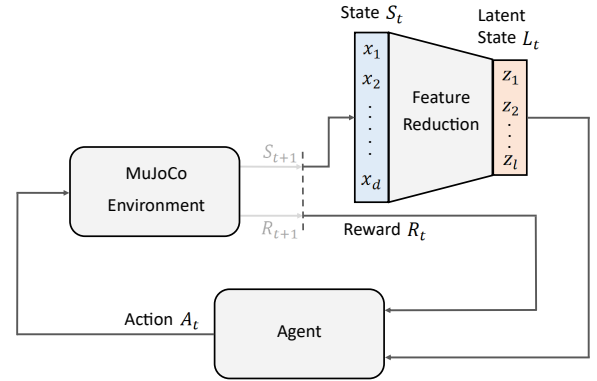


Fig. 2: The architecture of our system.

## III. METHOD

We use the intuition that an underacted system has correlations between states and thus there is a lower-dimensional representation that could describe this process.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna,

vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

#### IV. A NOVEL, UNDER-ACTUATED CONTROL BENCHMARK

We implemented our Jitterbug benchmark using the DeepMind Control Suite (DMC) framework [REF]. DMC is a framework and set of benchmark tasks for continuous control published by Google DeepMind in 2018. DMC benchmarks consist of a *domain* defining a robotic and environment model and *tasks* which are instances of that domain with specific MDP structure

DMC uses the robust Multi-Joint dynamics with Contact (MuJoCo) robotics physics engine for simulation [REF]. To aid comparison across tasks, DMC imposes constraints on rewards ( $R \in [0, 1]$ ) and episode length ( $H = 1000$ ). As such, for any DMC task, cumulative episode return  $\approx 1000$  indicates success.

DMC tasks are compatible subsets of the popular OpenAI Gym framework [REF], meaning many popular RL algorithm frameworks can be used with these benchmarks.

##### A. The Jitterbug Domain

Our Jitterbug model was inspired by the children’s toy Katita (Figure 1). We aimed to reproduce the physical dynamics of this toy while enabling control by replacing the wind-up spring with a single actuator of equivalent torque. Our Jitterbug model conforms to the dimensions and mass of the Katita, however we replace the wind-up spring with a controlled single degree-of-freedom motor. We retain the

(non-functional) wind up crank to more closely model the mass distribution of the physical Katita.

We used high-speed recording and visual tachometry to measure the Katita motor speed and leg vibration modes. By reverse-engineering the Katita gearbox we estimated the torque output of the drive spring and configured the MuJoCo actuator appropriately. We modelled the legs as rigid bodies with shoulder and elbow hinge joints. The hinge stiffness was manually tuned to reproduce the dominant leg vibration mode observed in our high-speed footage. The Jitterbug model density was set using standard values for stainless steel ( $7700 \text{ kg/m}^3$ ) for the body and tough plastic ( $1100 \text{ kg/m}^3$ ) for the feet. Figure 3 shows the physical composition of our simulated Jitterbug model.

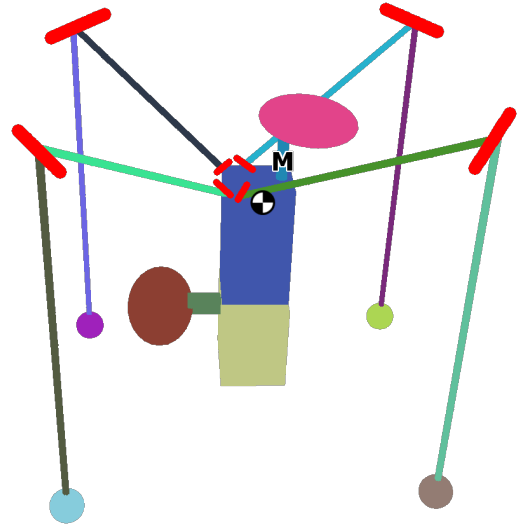


Fig. 3: Schematic representation of the Jitterbug model. Individual rigid bodies are in different colours and we highlight the position of the center of mass ( $\bullet$ ), hinge joints ( $\text{⚡}$ ) and the single motor (**M**).

Due to the importance of contact and stiff dynamics in the Jitterbug’s locomotion, we found it necessary to adjust MuJoCo’s default settings, selecting an integration timestep of 0.0002 and semi-implicit Euler integration. With these settings we qualitatively observed a close correspondence between the Katita and the simulated dynamics under constant motor actuation on the Jitterbug.

DMC supports the definition of physically-based camera models for to enable learning from raw pixels if desired. We defined several cameras for the Jitterbug domain including an overhead, tracking and ego-centric view.

##### B. The Jitterbug Task Suite

The Jitterbug dynamics naturally induce very high variance motion under a range of motor velocities. We defined a collection of five tasks of increasing difficulty based on the

Jitterbug domain. The tasks were designed with increasingly sparse reward signals to increase the difficulty.

For all tasks we choose  $\gamma = 0.99$  and consider a task solved when cumulative episode reward is  $\gtrapprox 900$ . In all tasks the Jitterbug is reset to a random pose near the origin at the start of an episode. All tasks have a single continuous action controlling the motor  $\mathcal{A} = [-1, 1]$  (larger/smaller values are clipped) and continuous state and observation spaces.

For each task, we report  $(\dim(\mathcal{S}), \dim(\mathcal{O}))$  and a brief description of the reward structure. Tasks are reported here ordered easiest to hardest.

- 1) *Move From Origin* (16, 15): The Jitterbug must move away from the origin in any direction. N.b. a sufficiently fast constant motor velocity is sufficient to solve this task.
- 2) *Face In Direction* (17, 16): The Jitterbug must rotate to face in a randomly selected yaw direction.
- 3) *Move In Direction* (20, 19): The Jitterbug is rewarded for velocity in a randomly selected direction in the X, Y plane.
- 4) *Move To Position* (19, 18): The Jitterbug must move to a randomly selected position in the X, Y plane.
- 5) *Move To Pose* (20, 19): The Jitterbug must move to a randomly selected position in the X, Y plane and rotate to face in a randomly selected yaw direction. N.b. Due to the multiplication of position and yaw reward components, this task has a *very* sparse reward signal!

In addition, for all tasks the Jitterbug must remain upright to achieve reward. Falling does not terminate the episode early as the leg dynamics are sufficiently springy that bouncing into the upright pose again can allow recovery from this condition (albeit at the loss of some reward). Indeed - we observed some learned strategies that appeared to utilize this mode of locomotion!

## V. EXPERIMENTS

### A. Characterising The Jitterbug Tasks

To verify feasibility, we hand-crafted heuristic policies that can solve each task. To characterise the difficulty of the Jitterbug task suite, we performed preliminary hyper-parameter tuning to select reasonable settings and trained several RL algorithms on the tasks.

Figure 7 reports training curves for example on- and off-policy algorithms. We contrast the performance of PPO (an on-policy method) and DDPG (an off-policy method). We also overlay the performance of our heuristic policies for comparison. Each figure shows the median and 10<sup>th</sup> - 90<sup>th</sup> percentile episode return across 10 different seeds.

Our selected hyper-parameters are reported in Table I. For all cases, we used fully-connected neural networks with hidden layers of size 350 and 250 with ReLU activation. Where applicable, we use separate networks for the actor and critic (i.e. no shared weights).

We ran additional experiments using TRPO, A2C and SAC and observed similar performance to the reported results.

Training curves for these algorithms are not included here for brevity.

TABLE I: Algorithm hyper-parameters. Bold items were changed from the defaults offered by the `stable_baselines` package.

Parameter	Value
<i>Shared</i>	
Optimizer	Adam [22]
<b>Learning Rate (<math>\alpha</math>)</b>	<b><math>1\text{E}^{-4}</math></b>
Network Architecture(s)	Fully Connected
Number of Hidden Layers	2
<b>Hidden Layer Sizes</b>	<b>[350, 250]</b>
<b>Activation Functions</b>	<b>ReLU</b>
<i>DDPG</i>	
<b>Batch Size</b>	<b>256</b>
Training Steps	50
Rollout and Evaluation Steps	100
<b>Replay Buffer Size</b>	<b><math>1\text{E}^6</math></b>
Soft Update Coefficient ( $\tau$ )	$1\text{E}^{-3}$
Parameter Noise	None
<b>Action Noise</b>	<b>Ornstein-Uhlenbeck</b> $\mu = 0.3, \sigma = 0.3, \theta = 0.15$
<i>PPO</i>	
<b>Steps / Environment / Update</b>	<b>256</b>
<b>Entropy Coefficient</b>	<b><math>1\text{E}^{-2}</math></b>
Value Function Coefficient	0.5
Max Gradient Norm	0.5
Bias-Variance Coefficient ( $\lambda$ )	0.95
Minibatches	4
Policy Clipping Range	0.2
Value Clipping Range	None
Surrogate Optimization Epochs	4

### B. Characterising Learned Policies

To verify the learned policies were sensible (i.e. to confirm the absence of ‘reward hacking’) we qualitatively and quantitatively investigated the exhibited behaviours.

We observed that a key difference between successful and unsuccessful trained policies seemed to be the ability to learn piecewise control functions. For example, for all tasks but *Move From Origin*, achieving high reward requires careful modulation of the reactive torque applied to the Jitterbug body by the motor counterweight. One way to achieve this (the method we use in our heuristic policies) is by pulsing the motor in different directions. We observed that successful policies learned to pulse the motor in short bursts in alternating directions (e.g.  $\sim 180^\circ$  at a time, see Figure 4), whereas unsuccessful policies would often drive the motor continuously. In doing so, the successful policies were able to achieve high cumulative episode return, and accomplish the high-level task encoded by the reward (Figure 5).

### C. Reduced-Order Training

*Autoencoder*: The input  $X$  is corrupted by adding random noise to each of the features. Random noise used:  $X_{\text{noisy}} =$

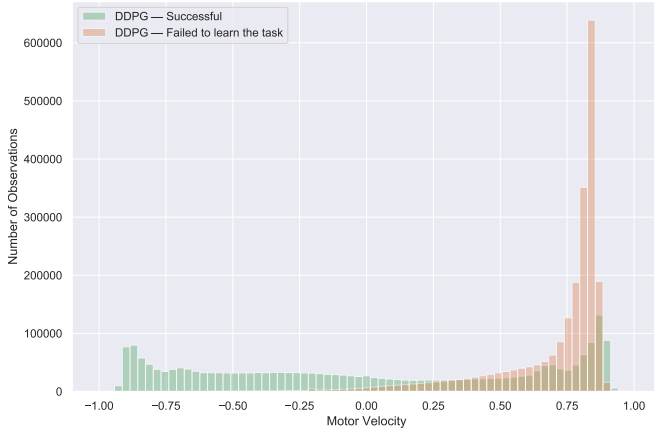


Fig. 4: Characterising policy behaviours for the *Move In Direction* task. We plot the distribution of motor velocities across many episodes for a successful policy (green) and unsuccessful policy (orange). The successful policy learns to pulse the motor in alternating directions (the same strategy used by our heuristic policy). In contrast, the unsuccessful policy gets stuck in a local minima where the motor is continuously driven in one direction.



Fig. 5: Heatmap showing Jitterbug position over 100 episodes before (left) and after (right) training DDPG on the task *Move In Direction*. In the second figure, to evaluate the agent, the target direction was fixed at  $+45^\circ$ .

$X + \mathcal{N}(0, 0.1)$ . The autoencoder is fed with the corrupted data  $X_{noisy}$  and is trained to reproduce the original data  $X$ . Loss used: Mean Squared Error between the output  $Y$  and the original data  $X$ .

Dimensions:

- Input/Output:  $d = 16$
- Latent: several cases -  $l = 12, 8, 4$

*Training the Autoencoder:* Gathering of the data used to train the autoencoder: random policy, i.e. an agent was run for 5M steps taking only random actions and each observed state was saved in a file.

*Training:* from this dataset, 80% of the data was used to train and 20% to test. Use: Once trained, only the encoder part of the autoencoder is used (see figure)

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero,

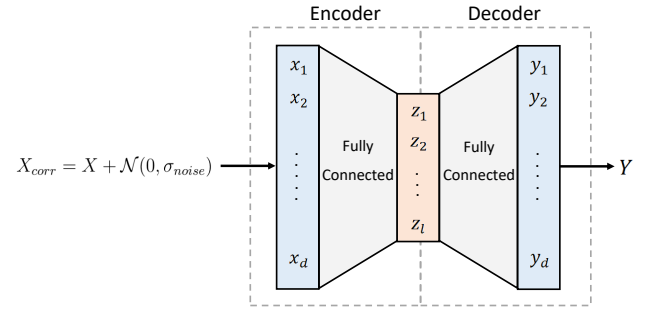


Fig. 6: WE used a De-Noising AutoEncoder as a means to learn a reduced-order state representation.

nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus



nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

## VI. DISCUSSION

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## VII. CONCLUSION

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero,

nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## APPENDIX

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## ACKNOWLEDGMENT

We thank Dr. Hanna Kurniawati at the ANU Robust Decision-making and Learning Laboratory for discussions and simulation assistance. This research was partly supported by an Australian Research Council Discovery Project (DP160100714). A. Snoswell is supported in part through and Australian Government Research Training Program Scholarship.

## REFERENCES

- [1] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [2] 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," *arXiv preprint arXiv:1509.02971*, 2015.
- [3] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [4] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," *arXiv preprint arXiv:1801.01290*, 2018.
- [5] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," *arXiv preprint arXiv:1802.09477*, 2018.
- [6] R. Islam, P. Henderson, M. Gombrokchi, and D. Precup, "Reproducibility of benchmarked deep reinforcement learning tasks for continuous control," *arXiv preprint arXiv:1708.04133*, 2017.
- [7] S. Bhatnagar, D. Precup, D. Silver, R. S. Sutton, H. R. Maei, and C. Szepesvári, "Convergent temporal-difference learning with arbitrary smooth function approximation," in *Advances in Neural Information Processing Systems*, pp. 1204–1212, 2009.

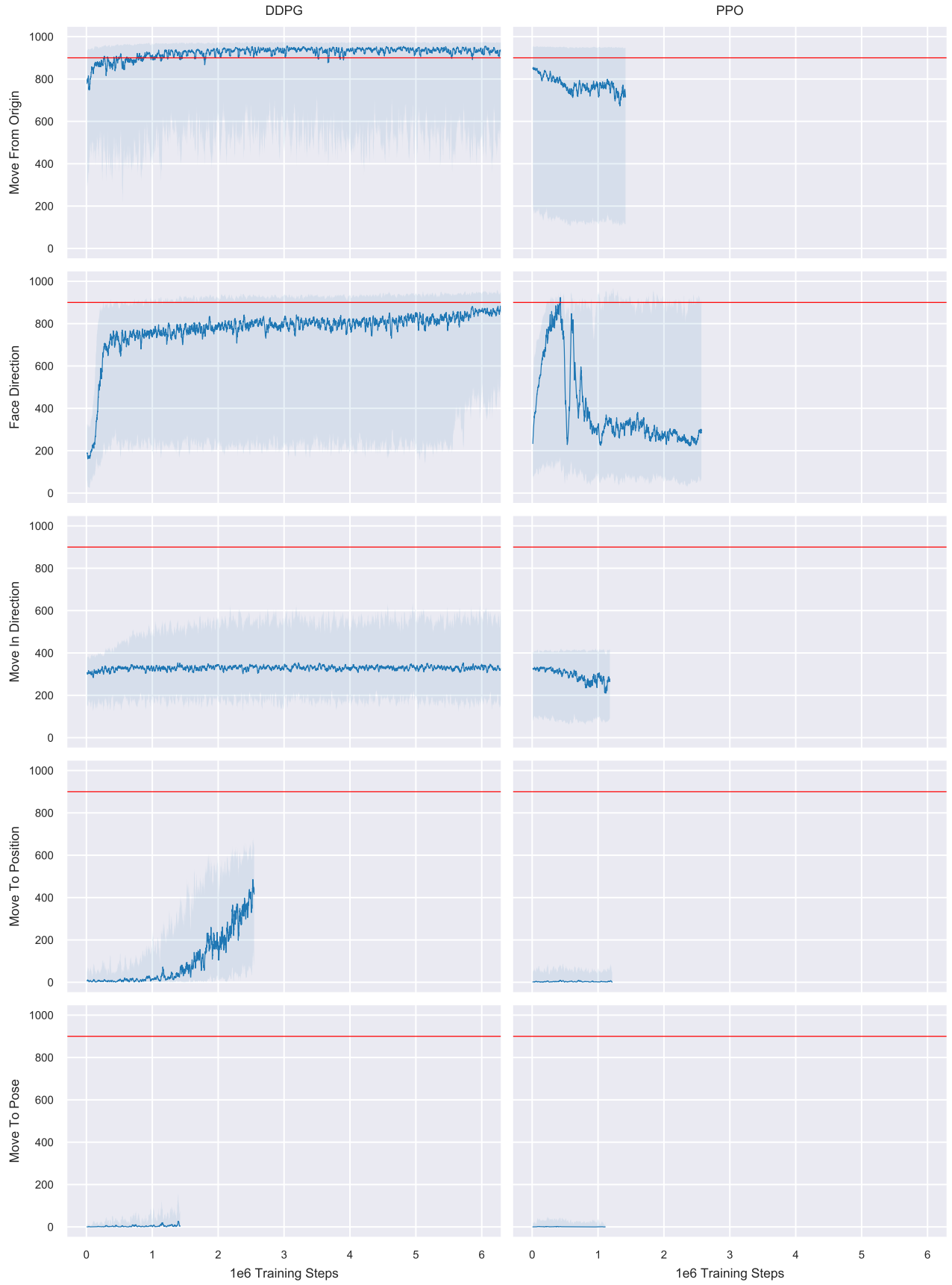


Fig. 7: Characterising the Jitterbug task suite. We compare the training progress of DDPG and PPO up to 6 million training steps (6000 episodes). We show median (solid line) and the 10<sup>th</sup> and 90<sup>th</sup> quartiles (shaded area) of per-episode episode return across 10 random seeds in each figure. A task is considered 'solved' if the trained agent consistently scores  $\gtrsim 900$  return per episode (red line). All plots are filtered with a  $20 \times 10^3$  step moving average filter.

- [8] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. 1998.
- [9] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel, "Benchmarking deep reinforcement learning for continuous control," in *International Conference on Machine Learning*, pp. 1329–1338, 2016.
- [10] C. Finn, X. Y. Tan, Y. Duan, T. Darrell, S. Levine, and P. Abbeel, "Deep spatial autoencoders for visuomotor learning," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 512–519, IEEE, 2016.
- [11] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel, "Deep imitation learning for complex manipulation tasks from virtual reality teleoperation," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1–8, IEEE, 2018.
- [12] R. Tedrake, "Underactuated robotics: Learning, planning, and control for efficient and agile machines: Course notes for mit 6.832," tech. rep., Massachusetts Institute of Technology, 2009.
- [13] J. T. Betts, *Practical methods for optimal control and estimation using nonlinear programming*, vol. 19. SIAM, 2010.
- [14] O. Von Stryk, "Numerical solution of optimal control problems by direct collocation," in *Optimal Control*, pp. 129–143, Springer, 1993.
- [15] M. Kalakrishnan, S. Chitta, E. Theodorou, P. Pastor, and S. Schaal, "Stomp: Stochastic trajectory optimization for motion planning," in *2011 IEEE international conference on robotics and automation*, pp. 4569–4574, IEEE, 2011.
- [16] K. Nakajima, H. Hauser, T. Li, and R. Pfeifer, "Information processing via physical soft body," *Scientific reports*, vol. 5, p. 10487, 2015.
- [17] A. Nagabandi, G. Yang, T. Asmar, R. Pandya, G. Kahn, S. Levine, and R. S. Fearing, "Learning image-conditioned dynamics models for control of underactuated legged millirobots," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4606–4613, IEEE, 2018.
- [18] T. Nishimura, K. Mizushima, Y. Suzuki, T. Tsuji, and T. Watanabe, "Thin plate manipulation by an under-actuated robotic soft gripper utilizing the environment," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1236–1243, IEEE, 2017.
- [19] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [20] A. Ng, "Sparse autoencoder," cs294a lecture notes, Stanford University, 2011.
- [21] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Thompson, S. Levine, and P. Sermanet, "Learning latent plans from play," *arXiv preprint arXiv:1903.01973*, 2019.
- [22] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.