



On Deep Reinforcement Learning for Spacecraft Guidance

Kirk Hovell* and Steve Ulrich†

Carleton University, Ottawa, Ontario K1S 5B6, Canada

This paper introduces a novel technique, named deep guidance, that leverages deep reinforcement learning, a branch of artificial intelligence, that enables guidance strategies to be learned rather than designed. The deep guidance technique consists of a learned guidance strategy that feeds velocity commands to a conventional controller to track. Control theory is combined with deep reinforcement learning in order to lower the learning burden and facilitate the transfer of the trained system from simulation to reality. In this paper, a proof-of-concept spacecraft pose tracking and docking scenario is considered, in simulation and experiment, to test the feasibility of the proposed approach. Results show that such a system can be trained entirely in simulation and transferred to reality with comparable performance.

I. Introduction

AUTONOMOUS spacecraft rendezvous and docking operations have become an active research area in recent decades. Applications, such as on-orbit servicing, assembly, and debris capture¹ require the capability for a chaser spacecraft to autonomously and safely maneuver itself in proximity to a potentially uncooperative target object. A common strategy is pose tracking, i.e., synchronizing the translational and rotational motion of the chaser with respect to the target, such that there is no relative motion between the two objects. Only then can the chaser perform its final approach and capture or dock with the target. Guidance and control algorithms have been developed for this purpose. For example, a guidance and control scheme for capturing a tumbling debris with a robotic manipulator was developed by Aghili.² Wilde et al.³ developed inverse dynamics models to generate guidance paths and included experimental validation. Ma et al.⁴ applied feed-forward optimal control for orienting a chaser spacecraft at a constant relative position with respect to a tumbling target. A variety of dual quaternion approaches have been explored.⁵⁻⁷ Pothen and Ulrich^{8,9} used the Udwalla-Kalaba equation to formulate the close-range rendezvous problem and included experimental validation. Lyapunov vector fields have been used to command docking with an uncooperative target spacecraft by Hough and Ulrich.¹⁰

The currently-developed guidance and control techniques presented above are hand-crafted to solve a particular task and require significant engineering effort. As more complex tasks are introduced, the engineering effort needed to hand-craft solutions may become infeasible. For example, developing a guidance law for a chaser spacecraft to detumble a piece of tumbling space debris when the two objects are connected via flexible tethers does not have a clear solution that can be hand-crafted. Motivated by more difficult guidance tasks, this paper introduces a new approach that builds upon a branch of artificial intelligence, called deep reinforcement learning, to augment the guidance capabilities of spacecraft for difficult tasks.

Reinforcement learning is based on the idea of an agent trying to choose actions in order to maximize the rewards it accumulates over a period of time. The agent uses a *policy* that, when given an input, returns a suggested action to take. At each timestep, a scalar reward, which may be positive or negative, is given to the agent and corresponds to task completion. Through trial-and-error, the agent attempts to learn a policy that maps inputs to the actions that maximize the expected reward. Therefore, by selecting an appropriate reward scheme, complex behaviours can be learned by the agent without being explicitly programmed. For a tethered space debris detumbling task, for example, rewards could be given for reducing the angular velocity

*PhD Candidate, Department of Mechanical and Aerospace Engineering, 1125 Colonel By Drive. Student Member AIAA.

†Associate Professor, Department of Mechanical and Aerospace Engineering, 1125 Colonel By Drive. Senior Member AIAA.

of the debris and penalties could be given for fuel usage, forcing the agent to learn a fuel-efficient detumbling policy. The engineering effort is reduced to specifying the reward system rather than the complete logic required to complete the task and is the main appeal of reinforcement learning. Neural networks have become a popular choice for representing the policy in reinforcement learning as they are universal function approximators.¹¹ When neural networks are used within reinforcement learning the technique is called *deep reinforcement learning*. The core concepts in reinforcement learning have been around for decades, but have only become useful in recent years due to the rapid rise in computing power. Many notable papers have been published this decade using deep reinforcement learning to solve previously-unsolvable tasks. In 2015, Mnih et al.¹² applied deep reinforcement learning to play many Atari 2600¹³ games at a superhuman level by training a policy to select button presses as a function of the screen pixels. Silver et al.^{14,15} used deep reinforcement learning to master the game of Go in 2016, a decade earlier than expected.

Training deep reinforcement learning policies on physical robots is time consuming, expensive, and leads to significant wear-and-tear on the robot because, even with state-of-the-art learning algorithms, the task may have to be repeated hundreds or thousands of times before learning succeeds. Training a policy onboard a spacecraft is not viable, due to fuel, time, and compute limitations. An alternative is to train the policy in a simulated environment and to transfer the trained policy to the physical robot. If the simulated model is sufficiently accurate and the task is not highly dynamic, policies trained in simulation may be directly transferrable to a real robot.¹⁶ For example, Tai et al.¹⁷ trained a room-navigating robot in simulation and deployed it to reality with success. While good results were obtained, this technique is unlikely to generalize well to more difficult or dynamic tasks due to the effect known as the simulation-to-reality gap.^{16,18,19} It states: since the simulator the policy is trained within cannot perfectly model the dynamics of the real world, a policy trained in simulation will often fail to perform well in a real-world environment due to policy-overfitting of the simulated dynamics. Efforts to get around this problem often take advantage of domain randomization,²⁰ i.e., randomizing the environmental parameters for each simulation to force the policy to become robust to environmental changes. Domain randomization has been used successfully in a quadrepedal robot recovery task²¹ and in the manipulation of objects with a robotic hand.^{20,22} However, domain randomization significantly increases the training time. Other efforts to solve the simulation-to-reality problem involve continuing to train the policy once deployed to experiment. A drifting car²³ policy was partially trained in simulation and then fine-tuning training was performed once experimental data was collected. A quadrotor stabilization task²⁴ summed the policy output with a conventional PD controller to guide the learning process and help with the simulation-to-reality transfer.

Deep reinforcement learning has also been applied to aerospace applications, though mostly in simulation. A simulated fleet of wildfire surveillance aircraft used deep reinforcement learning to command the flight-path of the aircraft.²⁵ Deep reinforcement learning has also been applied to spacecraft map generation while orbiting small bodies,²⁶ spacecraft orbit control in unknown gravitational fields,²⁷ and spacecraft orbital transfers.²⁸ Others have used reinforcement learning to train a policy that performs guidance and control for pinpoint planetary landing.^{29,30} Neural networks have been trained to approximate offline-generated optimal guidance paths for pinpoint planetary landing, such that the neural network approximates an optimal guidance algorithm that can be executed in real-time.^{31,32}

Inspired by deep reinforcement learning's ability to have behaviour be learned rather than hand-crafted, and motivated by the need for new simulation-to-reality transfer techniques, this paper introduces a novel technique that allows for reinforcement learning's use on real spacecraft. The proposed technique builds off of the planetary landing work,^{29,30} where neural networks were responsible for outputting guidance commands that were fed to a conventional controller. There, the neural networks were trained to approximate hand-crafted optimal guidance trajectories. Here, deep reinforcement learning is used to train a guidance policy whose trajectories are fed to a conventional controller to track. Using reinforcement learning allows for an unbiased guidance policy to be discovered by the agent instead of being shown many hand-crafted guidance trajectories to mimic. The proposed technique is in contrast to typical reinforcement learning research, where the policy is responsible for learning both the guidance and control logic. By restricting the policy to learn only the guidance portion, we harness the high-level, unbiased, task-solving abilities of reinforcement learning while deferring the control aspect to the well-established control theory community. Control theories have been developed that are able to perform trajectory tracking well under dynamic uncertainty,³³⁻³⁵ and can therefore handle model discrepancies between simulation and reality. Harris et al.³⁶ recently presented a strategy that uses reinforcement learning to switch between a set of available controllers depending on the system state. The authors³⁶ suggest that reinforcement learning should not be tasked with learning guidance

and control since control theory already has great success. By combining deep reinforcement learning for guidance with conventional control theory, the policy is prevented from learning a controller that overfits the error-prone simulated dynamics. We call our deep reinforcement learning guidance strategy *deep guidance*. In light of the above, the novel contributions of this paper are:

1. The *deep guidance* technique, that combines deep reinforcement learning guidance with a conventional controller.
2. Experimental demonstrations showing this deep guidance strategy can be trained in simulation and deployed to reality without any fine-tuning.
3. The first, to the best of the authors' knowledge, experimental demonstration of artificial intelligence commanding the motion of a spacecraft platform.

It should be noted that although the authors were motivated by difficult guidance tasks, such as detumbling tethered space debris, in this initial paper a proof-of-concept task—a simple spacecraft pose tracking and docking scenario—is considered. Demonstrating the deep guidance technique on a simple task will highlight its potential for use on more difficult tasks.

This paper is organized as follows: Sec. II presents background on deep reinforcement learning and the specific learning algorithm used in this paper, Sec. III presents the novel guidance concept developed by the authors, Sec. IV describes the pose tracking scenario considered, Sec. V presents numerical simulations demonstrating the effectiveness of the technique, Sec. VI presents experimental results, and Sec. VII concludes this paper.

II. Deep Reinforcement Learning

The goal of deep reinforcement learning is to discover a policy, π_θ , represented by a feed-forward neural network and whose subscript denotes it has trainable weights θ , that maps states, $\mathbf{x} \in \mathcal{X}$, to actions, $\mathbf{a} \in \mathcal{A}$ that, when executed, maximizes the expected rewards received over one episode (in reinforcement learning, one simulation is called one episode). The action is obtained by feeding the state to the policy, as follows

$$\mathbf{a} = \pi_\theta(\mathbf{x}). \quad (1)$$

While many deep reinforcement learning algorithms are currently available, this work uses the Distributed Distributional Deep Deterministic Policy Gradient algorithm³⁷ (D4PG). D4PG, released in early 2018, was selected because it operates in continuous state and action spaces, it has a deterministic output, it can be trained in a distributed manner to take advantage of multi-CPU machines, and it achieves state-of-the-art performance.

II.A. D4PG Algorithm

The D4PG³⁷ algorithm is an actor-critic algorithm, implying there is a policy neural network, $\pi_\theta(\mathbf{x})$, that maps states to actions *and* a value neural network, $Z_\phi(\mathbf{x}, \mathbf{a})$, with trainable weights ϕ , that maps a state-action pair to a probability distribution of the expected rewards for the remainder of the episode. The total rewards expected to be received from a given state, \mathbf{x} , and taking action \mathbf{a} from the policy, $\pi_\theta(\mathbf{x})$, is given by

$$J(\theta) = \mathbb{E}\{Z_\phi(\mathbf{x}, \pi_\theta(\mathbf{x}))\} \quad (2)$$

where $J(\theta)$ is the expected rewards from the given state as a function of the policy weights, θ , and \mathbb{E} denotes the expectation. The objective of reinforcement learning is then to find a policy $\mathbf{a} = \pi_\theta(\mathbf{x})$ that maximizes $J(\theta)$ through systematically adjusting θ .

To maximize Eq. (2), we must first establish the policy and value neural networks, shown in Fig. 1. The policy network, shown in Fig. 1a, accepts the system state \mathbf{x} as the input, shown with three elements x_1 , x_2 , and x_3 in the figure, and outputs a commanded action, $\mathbf{a} = \pi_\theta(\mathbf{x})$, which may be multidimensional. Each arrow in the network corresponds to a trainable weight that is parameterized by θ . The value network, shown in Fig. 1b, accepts the state and action vectors as inputs and outputs a value distribution, $Z_\phi(\mathbf{x}, \mathbf{a})$. In other words, it returns a probability distribution of how many discounted rewards are expected to be

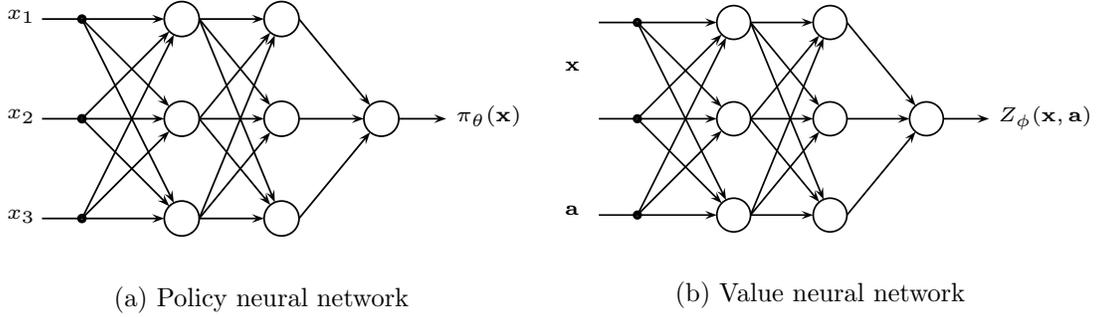


Figure 1. Policy and value neural networks used in the D4PG algorithm.

received from taking action \mathbf{a} in state \mathbf{x} and continuing to follow the policy until the end of the episode. The value network has trainable weights ϕ .

Simulated state, action, and reward data are obtained from running episodes and placing the data in a large *replay buffer*. Batches of simulated data are randomly sampled from the replay buffer and used to train the policy and value neural networks. From a given batch of state, action, next state, and reward data, the value network can be trained using supervised learning. Gradient-descent is used to minimize the cross entropy loss function given by

$$L(\phi) = \mathbb{E}[-Y \log(Z_\phi(\mathbf{x}, \mathbf{a}))] \quad (3)$$

where Y is the target value distribution, i.e., the new best prediction of the true value distribution based on the sampled data, as first introduced by Bellmare et al.³⁸ It is calculated using

$$Y = \sum_{n=0}^{N-1} \gamma^n r_n + \gamma^N Z_{\phi'}(\mathbf{x}_N, \pi_{\theta'}(\mathbf{x}_N)) \quad (4)$$

where r_n is the reward received at timestep n , γ is the discount factor for future rewards, i.e., future rewards are weighted lower than current rewards, and $Z_{\phi'}(\mathbf{x}_N, \pi_{\theta'}(\mathbf{x}_N))$ is the value distribution evaluated N timesteps into the future. The N -step return³⁹ is used, where N data points into the future are included for a more accurate prediction of Y . It should also be noted that the symbols θ' and ϕ' indicate that these are not the true weights of the policy and value networks but rather an exponential moving average of them, calculated by

$$\theta' = (1 - \epsilon)\theta' + \epsilon\theta \quad (5)$$

$$\phi' = (1 - \epsilon)\phi' + \epsilon\phi \quad (6)$$

with $\epsilon \ll 1$. Having a copy of the policy and value networks with smoothed weights has been empirically shown to have a stabilizing effect on the learning process.¹²

Equation (4) recursively calculates an updated prediction for the value distribution for a given state according to the reward data received. Then, by minimizing the loss function in Eq. (3) through adjusting ϕ , the value network slowly approaches these updated predictions, which are then used in Eq. (4). This recursive process led Sutton and Barto⁴⁰ to write “*we learn a guess from a guess*”.

With the training procedure for the value network outlined, the policy network must now be trained. The goal is to adjust the policy parameters, θ , in the direction of increasing the expected rewards for a given state, $J(\theta)$. Since neural networks are differentiable, the chain rule can be used to compute

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial J(\theta)}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \theta} \quad (7)$$

where $\frac{\partial J(\theta)}{\partial \mathbf{a}}$ is computed from the value network and $\frac{\partial \mathbf{a}}{\partial \theta}$ is computed from the policy network. In a sense, we are differentiating through the value network into the policy network. More formally:

$$\nabla_{\theta} J(\theta) = \mathbb{E}[\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \mathbb{E}[\nabla_{\mathbf{a}} Z_{\phi}(\mathbf{x}, \mathbf{a})] |_{\mathbf{a}=\pi_{\theta}(\mathbf{x})}] \quad (8)$$

describes how the policy parameters, θ , should be updated to increase the expected rewards when the policy π_θ is used.

To implement the algorithm, a number of agents run independent episodes using the most up-to-date version of the policy. Gaussian exploration noise is applied to the chosen action to force exploration, and is the basis for discovering new strategies. The action is obtained through

$$\mathbf{a}_t = \pi_\theta(\mathbf{x}_t) + \mathcal{N}(0, 1)\sigma \quad (9)$$

where $\mathcal{N}(0, 1)$ is the standard normal distribution and σ is the exploration noise standard deviation. The collected data are the state, \mathbf{x}_t , action, \mathbf{a}_t , reward, \mathbf{r}_t , and next state, \mathbf{x}_{t+N} , and are placed into a replay buffer that stores the most recent R data points. Asynchronously, a learner randomly samples batches of data from the replay buffer and uses them to train the value network one step using Eq. (3) and then trains the policy network one step using Eq. (8). Over time, the accuracy of the value network and the average performance of each agent is expected to increase. The D4PG algorithm is summarized in Algorithm 1.

Algorithm 1: D4PG³⁷

Input: Batch size M , n-step return length N , number of actors K , replay buffer size R , exploration rate σ , discount factor γ , policy and value network learning rates α and β

Learner

- 1 Initialize policy network weights θ and value network weights ϕ randomly
- 2 Initialize smoothed policy and target value network weights $\theta' = \theta$ and $\phi' = \phi$
- 3 Launch K actors and copy policy weights θ to each actor
- repeat**
- 4 Sample a batch of M data points from the replay buffer
- 5 Compute the target value distribution used to train the value network

$$Y = \sum_{n=0}^{N-1} \gamma^n r_n + \gamma^N Z_{\phi'}(\mathbf{x}_N, \pi_{\theta'}(\mathbf{x}_N))$$
- 6 Update value network weights ϕ by minimizing the loss function $L(\phi) = \mathbb{E}[-Y \log(Z_\phi(\mathbf{x}, \mathbf{a}))]$ using learning rate β
- 7 Compute policy gradients using $\nabla_\theta J(\theta) = \mathbb{E}[\nabla_\theta \pi_\theta(\mathbf{x}) \mathbb{E}[\nabla_{\mathbf{a}} Z_\phi(\mathbf{x}, \mathbf{a}) |_{\mathbf{a}=\pi_\theta(\mathbf{x})}]$ and update policy weights via $\theta = \theta + \nabla_\theta J(\theta)\alpha$
- 8 Update the smoothed network weights slowly in the direction of the main policy and value network weights $\theta' = (1 - \epsilon)\theta' + \epsilon\theta$ and $\phi' = (1 - \epsilon)\phi' + \epsilon\phi$ for $\epsilon \ll 1$
- 9 **until** *acceptable performance*

Actor

- repeat**
 - 10 From the given state, use the policy to obtain an action and add exploration noise with σ as the standard deviation $\mathbf{a}_t = \pi_\theta(\mathbf{x}_t) + \mathcal{N}(0, 1)\sigma$
 - 11 Step environment forward one timestep using action \mathbf{a}_t
 - 12 Record $(\mathbf{x}_t, \mathbf{a}_t, r_t = \sum_{n=0}^{N-1} \gamma^n r_n, \mathbf{x}_{t+N})$ and store in the replay buffer
 - 13 At the end of each episode, obtain the most up-to-date version of the policy π_θ from the Learner and reset the environment.
 - 14 **until** *acceptable performance*
-

III. Deep Guidance

Using deep reinforcement learning for spacecraft guidance may allow for difficult tasks to be accomplished through learning an appropriate behaviour rather than hand-crafting such a behaviour. In order for the reinforcement learning algorithm presented in Sec. II.A to be trained in simulation and deployed to reality without any fine-tuning on the spacecraft, it is proposed that the learning algorithm cannot be responsible for the entire guidance, navigation, and control stack (GNC). This is in order to prevent the policy from overfitting the simulated dynamics and being unable to handle the transition to a dynamically-uncertain real

world, i.e., the simulation-to-reality problem.¹⁶ For this reason, the authors present a system, called deep guidance, which uses deep reinforcement learning as a guidance system along with a conventional controller. Conventional controllers are able to handle dynamic uncertainties and modelling errors that typically plague reinforcement learning policies that attempt to learn the entire GNC routine. It is assumed that perfect navigation is available. A block-scheme diagram of the proposed system is shown in Fig. 2. The learned

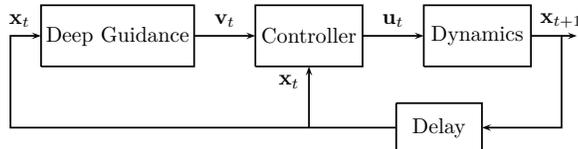


Figure 2. Proposed deep guidance strategy.

“Deep Guidance” block has the current state \mathbf{x}_t as its input and the desired velocity \mathbf{v}_t as its output. The desired velocity is fed to a conventional controller, that also receives the current state \mathbf{x}_t , and calculates a control effort \mathbf{u}_t . The control effort is executed on the dynamics that generate a new state, \mathbf{x}_{t+1} , one timestep later.

During training of the deep guidance model, an ideal controller is assumed, as shown in Fig. 3. This

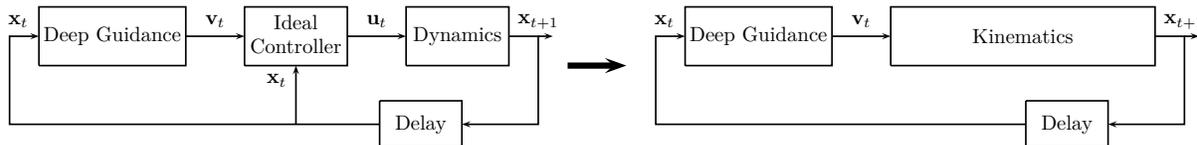


Figure 3. Proposed deep guidance with an ideal controller for training purposes in simulation.

ensures the guidance model does not overfit to any specific controller, making this approach controller-independent. Since the ideal controller perfectly commands the dynamics to move at the desired velocity \mathbf{v}_t , the ideal controller and the dynamics model may be combined into a single kinematics model.

Once trained, any controller may be combined with the deep guidance system for use on a real robot. Since the deep guidance system only experiences an ideal controller during training, it is possible that the guidance model will not experience any non-ideal states which may be encountered with a non-ideal controller, which may harm guidance performance. For this reason, Gaussian noise may be applied to the output of the kinematics to force the system into undesirable states during training that may be encountered by a non-ideal controller.

The proposed system is trained and tested on a spacecraft pose tracking and docking task detailed in the following section.

IV. Problem Statement

This section presents the simulated spacecraft pose tracking and docking environment the deep guidance system will be trained within. Though the deep guidance technique may allow for more complex behaviours to be learned, the simple task considered here is presented as a proof-of-concept. Similarly to other spacecraft researchers with experimental validation,^{3,41–43} the problem uses double-integrator planar dynamics such that the available experimental facility can validate the simulation work. Including a full orbital model was deemed to be outside the scope of this paper, since the goal of this work is to demonstrate the simulation-to-reality ability of the proposed deep guidance approach. Furthermore, a double-integrator model is representative of proximity operations in orbit over small distances and timescales.⁴⁴

A chaser spacecraft exists in a planar environment and is tasked with approaching and docking with a target spacecraft. An example initial scenario is shown in Fig. 4. The chaser and target spacecraft start at an initial position with zero velocity. The chaser spacecraft is given some time to maneuver to the safe hold point in front of the target. Then, the chaser spacecraft is tasked with approaching the target such that the two may dock.

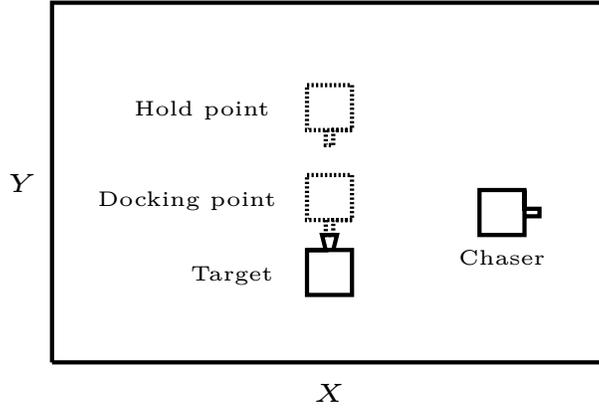


Figure 4. Sample spacecraft pose tracking and docking task.

IV.A. Kinematics and Dynamics Models

During training, a kinematics model is used which approximates a dynamics model and an ideal controller, as shown in Fig. 3. The deep guidance policy network accepts the system state and outputs the commanded action, which in this case is the velocity, \mathbf{v}_t ,

$$\mathbf{v}_t = \pi_{\theta}(\mathbf{x}_t) \quad (10)$$

which is then numerically integrated using the Scipy⁴⁵ Adams/Backward differentiation formula methods in Python to obtain \mathbf{x}_{t+1} .

To evaluate the learning performance, the trained policy is periodically evaluated on an environment with full dynamics and a controller, as shown in Fig. 2. In other words, it is “deployed” to another simulation for evaluation in much the same way that it will be deployed to an experiment in Sec. VI. The deep guidance policy outputs the desired velocity as in Eq. (10) that is fed to a simple proportional velocity-controller of the form

$$\mathbf{u}_t = K_p (\mathbf{v}_t - \dot{\mathbf{x}}_t) \quad (11)$$

where $K_p = 2$ and \mathbf{u}_t is the control effort. The K_p value is 20 times smaller for the angular motion controller. The control effort is executed on double-integrator dynamics using Newton’s Laws

$$\ddot{x} = \frac{F_x}{m} \quad (12)$$

$$\ddot{y} = \frac{F_y}{m} \quad (13)$$

$$\ddot{\psi} = \frac{\tau}{I} \quad (14)$$

where F_x is the force applied in the X direction, F_y is the force applied in the Y direction, τ is the torque applied about the Z axis, m is the chaser spacecraft mass, I is its moment of inertia, \ddot{x} is the acceleration in X , \ddot{y} is its acceleration in Y , and $\ddot{\psi}$ is the angular acceleration about Z . The accelerations are integrated twice to obtain the position at the following timestep. The following subsection discusses the reward function used to incentivize the desired behaviour.

IV.B. Reward Function

To calculate the reward given to the agent at each timestep, a reward field, f , is generated according to

$$f(\mathbf{x}_t) = -|\text{goal} - \mathbf{x}_t| \quad (15)$$

where goal is the goal state. The reward field is zero at the goal state and becomes negative linearly as the chaser moves away from the goal state. The goal state is set to be the hold point in Fig. 4 for the first 60 or 90s (depending on the task) of each episode and then the goal state is switched to the docking point for

the remainder of the episode. The rotational reward field is scaled down by a factor of ten compared to the translational reward field to ensure that rotation does not dominate the reward field.

The reward given to the agent, $r_t(\mathbf{x}, \mathbf{a})$, depends on the action taken in a given state. Therefore, the difference in the reward field between the current and previous timestep is used to calculate the reward given to the agent.

$$r_t = f(\mathbf{x}_t) - f(\mathbf{x}_{t-1}) \quad (16)$$

Using this approach, a positive reward will be given to the agent if it chooses an action that moves the chaser closer to the goal state and a negative reward otherwise. The learning algorithm details are presented in the following subsection.

IV.C. Learning Algorithm Details

The policy and the value neural networks have 400 neurons in their first hidden layer and 300 neurons in their second layer—Fig. 1 is not to scale. In the value network, the action is fed directly into the second layer of the network, as this was empirically shown to be beneficial.^{37,46} The value network therefore has 139,651 trainable weights and the policy network has 124,003 trainable weights. Each neuron in both hidden layers uses a rectified linear unit as its nonlinear activation function, shown below:

$$g(y) = \begin{cases} 0 & \text{for } y < 0 \\ y & \text{for } y \geq 0. \end{cases} \quad (17)$$

In the output layer of the policy network, a $g(y) = \tanh(y)$ nonlinear activation function is used to force the commanded velocity to be bounded. The output layer of the value neural network uses the softmax function, shown below, to force the output value distribution to indeed be a valid probability distribution.

$$g(y_i) = \frac{e^{y_i}}{\sum_{k=1}^K e^{y_k}} \quad \forall i = 1, \dots, B \quad (18)$$

for each element y_i and for B bins in the value distribution. Drawing from the original value distribution paper,³⁸ $B = 51$ bins are used in this work. The value distribution bins are evenly spaced on the empirically-determined interval $[-1000, 100]$, as this is the range of accumulated rewards encountered during this pose tracking and docking task.

The policy and value networks are trained using the Adam stochastic optimization routine⁴⁷ with a learning rate of $\alpha = \beta = 0.0001$. The replay buffer R can contain 10^6 transition data points and a mini batch size $M = 256$ is used. The target networks are updated on each training iteration with $\epsilon = 0.001$. The noise standard deviation applied to the action commanded by the policy to force exploration during training is $\sigma = \frac{1}{3} [\max(\mathbf{a}) - \min(\mathbf{a})] (0.9999)^E$, where E is the episode number; having a standard deviation of one third the action range empirically leads to good exploration of the action space. The noise standard deviation is decayed exponentially as more episodes are performed to narrow the action search area, at a rate that halves roughly every 7,000 episodes. Ten actors are used, $K = 10$, such that simulated data using the most up-to-date version of the policy is being collected by ten actors simultaneously. A discount factor of $\gamma = 0.99$ was used. The Tensorflow⁴⁸ machine learning framework was used to generate and train the neural networks.

During training, a kinematics environment is used, as shown in Fig. 3, and uses the state vector:

$$\mathbf{x}_t = \begin{bmatrix} x & y & \psi & e_x & e_y & e_\psi \end{bmatrix}^T \quad (19)$$

where e_x , e_y , and e_ψ are the errors between the chaser's current state and the goal state.

Every five training episodes, the current policy is “deployed” and run in a full dynamics environment with a proportional velocity-controller to evaluate its performance, as shown in Fig. 2. During deployment, $\sigma = 0$ in Eq. (9) such that no exploration noise is applied to the deep guidance velocity commands. The state vector becomes

$$\mathbf{x}_t = \begin{bmatrix} x & y & \psi & \dot{x} & \dot{y} & \dot{\psi} & e_x & e_y & e_\psi \end{bmatrix}^T. \quad (20)$$

When the full dynamics state vector is fed to the policy, however, the velocity elements are omitted as the policy network was not trained with velocity input.

V. Simulation Results

To test the deep guidance approach, two variations of the spacecraft pose tracking and docking task are studied. The first uses a stationary target while the second uses a rotating target. Each scenario is run both with constant initial conditions and with randomized initial conditions. The 30 cm cube satellite platform has a uniform simulated mass of 10 kg.

V.A. Docking with a Stationary Target

The chaser and target nominal initial conditions are $[3 \text{ m}, 1 \text{ m}, 0 \text{ rad}]$ and $[1.85 \text{ m}, 0.6 \text{ m}, \frac{\pi}{2} \text{ rad}]$ for the chaser and target, respectively. The hold point is 1.0 m offset from the front-face of the target and the docking point is 0.5 m offset. The simulation is run for 120 seconds with a 0.2 second timestep. For the first 60 seconds, the chaser is incentivized to move towards the hold point and is incentivized to move to the docking point for the remaining time. The commanded velocity bounds are $\pm 0.05 \text{ m/s}$ and $\pm \frac{\pi}{18} \text{ rad/s}$.

Results of the first spacecraft pose tracking and docking task are shown in Fig. 5, and are the result of 15 hours of training on an Intel® i7-8700K CPU. The learning curve, shown in Fig. 5a, plots the total rewards

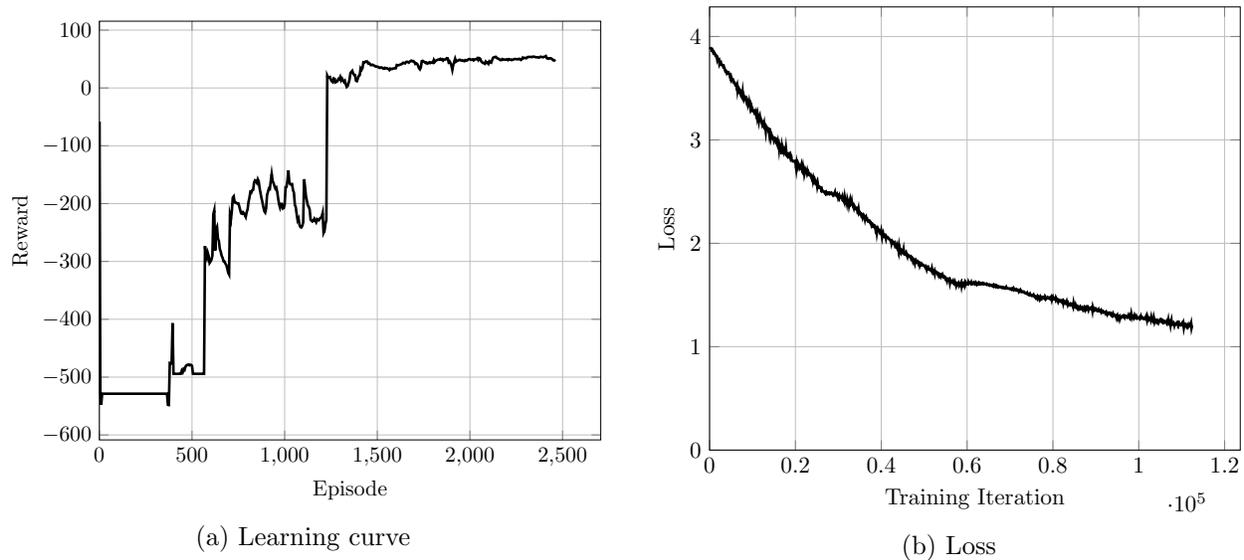


Figure 5. Training performance for chaser pose tracking and docking for a stationary target with constant initial conditions.

received on each episode, increases as expected during training. Here, it shows how the deep guidance strategy successfully learned the desired behaviour after roughly 2,000 episodes (12 hours), and then continued with minor improvements for the remainder of the training run. The loss function, calculated using Eq. (3) and shown in Fig. 5b, decreases as anticipated, indicating that the value network output distribution is approaching the target values calculated using Eq. (4), on average. Sample trajectories during training are shown in Fig. 6. The pose tracking and docking task was successfully learned. A video of the sample trajectories can be found at <https://youtu.be/JkmYxhQ30vI>.

Next, the chaser and target initial conditions were randomized at the beginning of each episode to force the deep guidance system to generalize across a range of initial states and not simply master a single trajectory. The initial states were randomized around the nominal ones according to a normal distribution with a standard deviation of 0.3 m for position and $\frac{\pi}{2}$ rad for attitude. Figure 7 shows the learning curve and associated loss function during training. The learning curve, shown in Fig. 7a, shows that the agent successfully learned a more general guidance strategy in the presence of randomized initial conditions. The learning curve appears more noisy than the learning curve shown in Fig. 5a because each episode may have slightly more or less rewards available depending on the initial conditions. Learned sample trajectories are shown in Fig. 8 and a video can be found at <https://youtu.be/8ReScA3y1KE>.

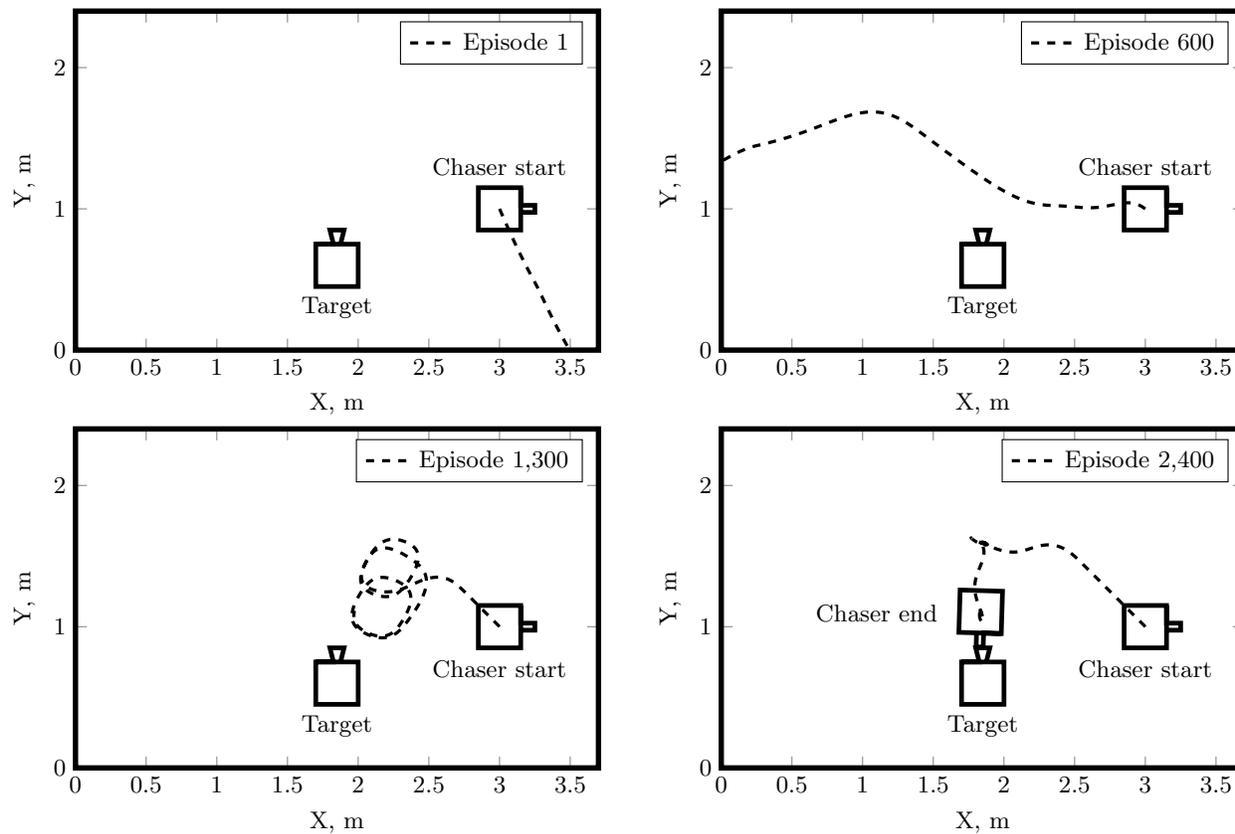


Figure 6. Visualization of chaser trajectories during the training process at various episodes.

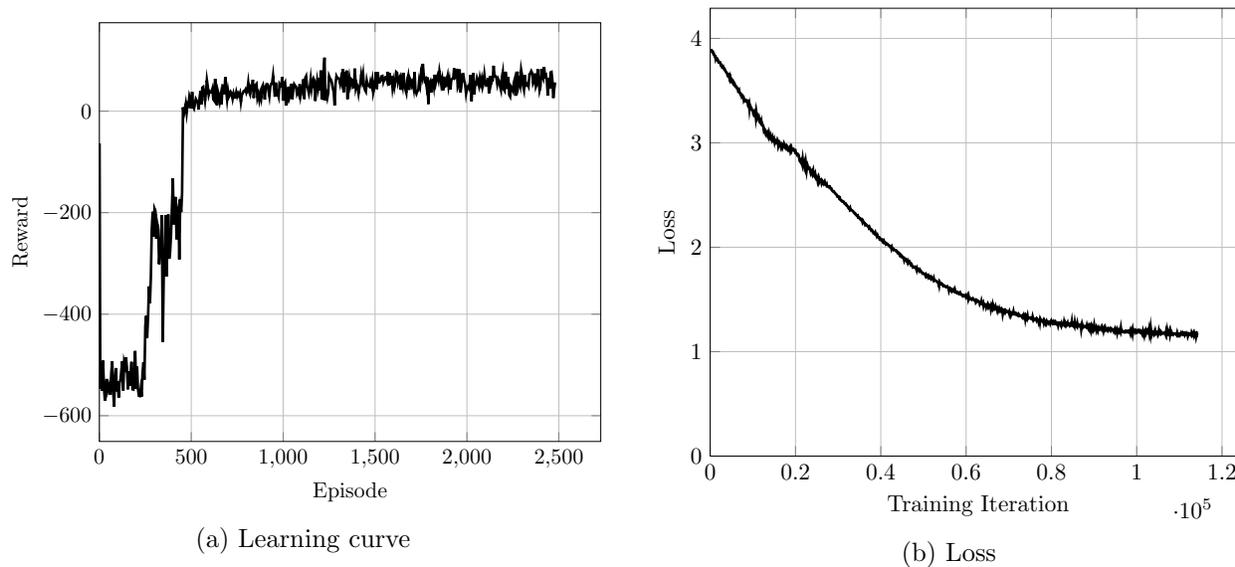


Figure 7. Training performance for chaser pose tracking and docking for a stationary target with randomized initial conditions.

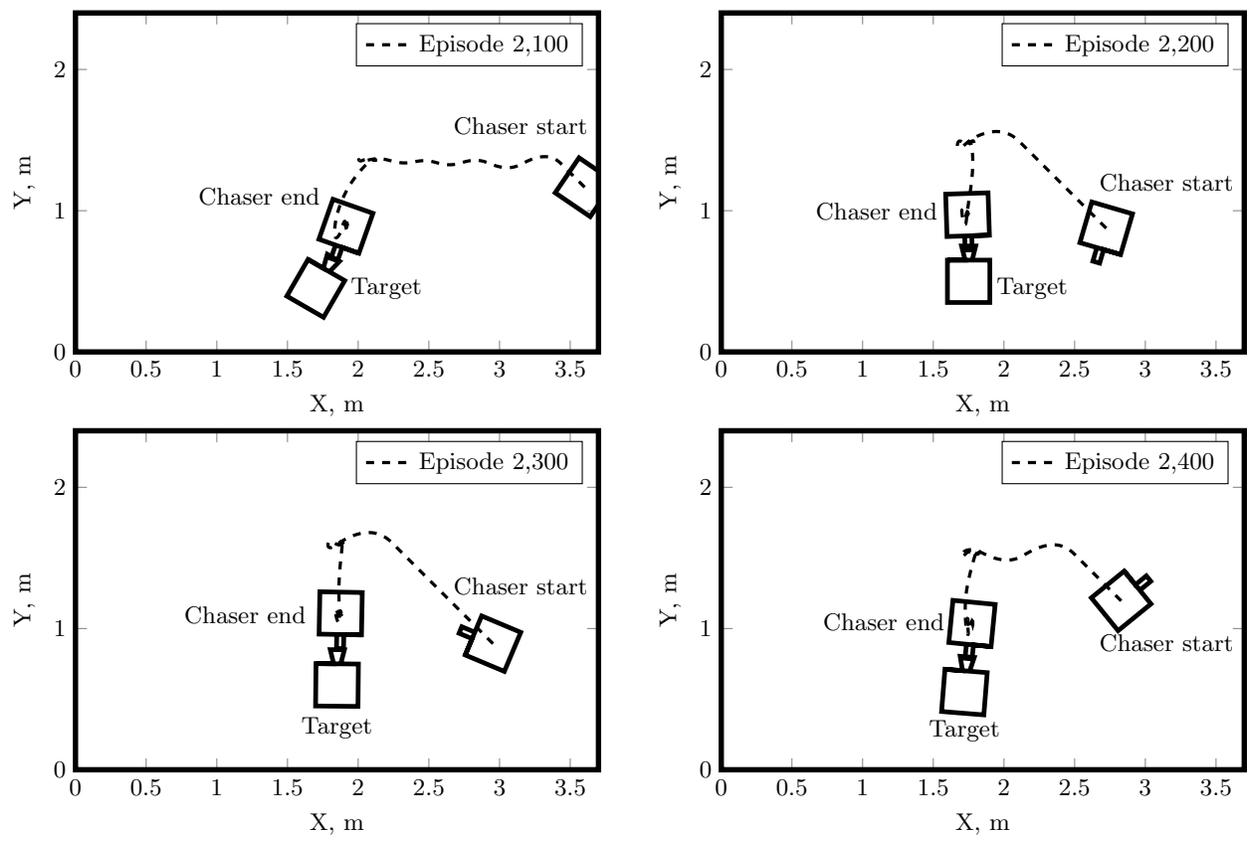


Figure 8. Examples of learned chaser trajectories with randomized initial conditions.

V.B. Docking with a Spinning Target

This subsection presents the second scenario that the deep guidance system was trained on. That is, a spacecraft pose tracking and docking task in the presence of a spinning target. All learning parameters are identical to those presented in Sec. IV.C demonstrating the generality of the proposed deep guidance approach. The target spacecraft is given a constant counter-clockwise angular velocity of $\omega = \frac{\pi}{45}$ rad/s. The velocity bounds on the chaser are increased to ± 0.1 m/s. The simulation is run for 180 seconds with a 0.2 second timestep. For the first 90 seconds, the chaser is incentivized to track the moving hold position, and afterwards it is incentivized to track the moving docking point. The initial conditions are [3 m, 1 m, 0 rad] and [1.85 m, 1.2 m, 0 rad] for the chaser and target, respectively.

The hold point is reduced to 0.9 m offset from the front-face of the target and the docking point is 0.5 m offset. Since the target is rotating, the hold and docking points are inertially moving with time. Constant initial conditions were first used, which produced the learning curve and loss function shown in Fig. 9.

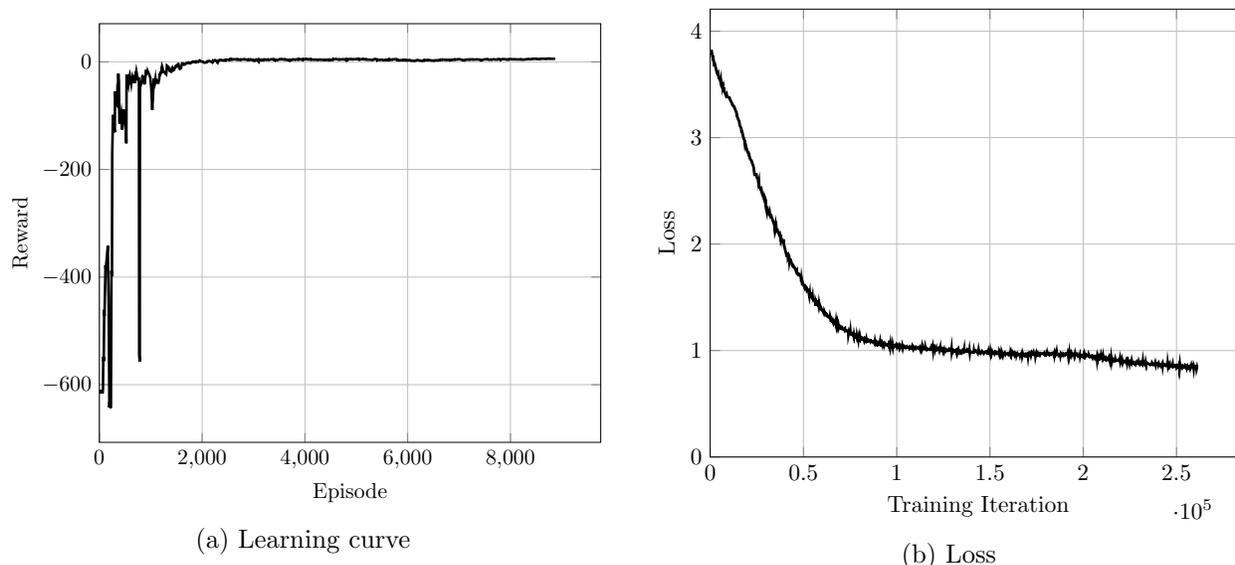


Figure 9. Training performance for chaser docking with a rotating target with constant initial conditions.

The learning curve in Fig. 9a shows that a deep guidance policy was successfully learned. Sample trajectories, shown in Fig. 10, show the motion of the chaser. After 2,000 episodes, a reasonably good policy was achieved. That is, the chaser rotates around the target and then makes its final approach for docking, as intended. A video of the sample trajectories can be found at <https://youtu.be/NWu5N81hUj0>.

The final simulation presented randomizes the initial conditions of the chaser and target at the beginning of each episode in the same manner presented in Sec. V.A. Results of this training run are shown in Fig. 11. Again, a deep guidance policy was successfully learned, as demonstrated by the increasing learning curve in Fig. 11a (and the decreasing loss function in Fig. 11b). The learning curve has a noisy steady state due to the initial condition randomness affecting each episode's total available rewards. Sampled successful trajectories are shown in Fig. 12 and a video of the trajectories can be found at https://youtu.be/ry_v_0zpQuI.

The deep guidance system presented in this paper was successfully trained on a simulated spacecraft pose tracking and docking task. It allows the designer to easily specify a reward function to convey the desired behaviour to the agent instead of specifying hand-crafted guidance trajectories. Although the guidance trajectories learned in this paper by the deep guidance system would be trivial to hand-craft, the purpose of this paper is to introduce the deep guidance technique. It is expected that the deep guidance technique will unlock the ability for more difficult learned guidance strategies in future work. All code used has been open-sourced and can be accessed at www.github.com/Kirkados/AIAA_GNC2020_D4PG.

The trained guidance policy is tested in experiment in the following section.

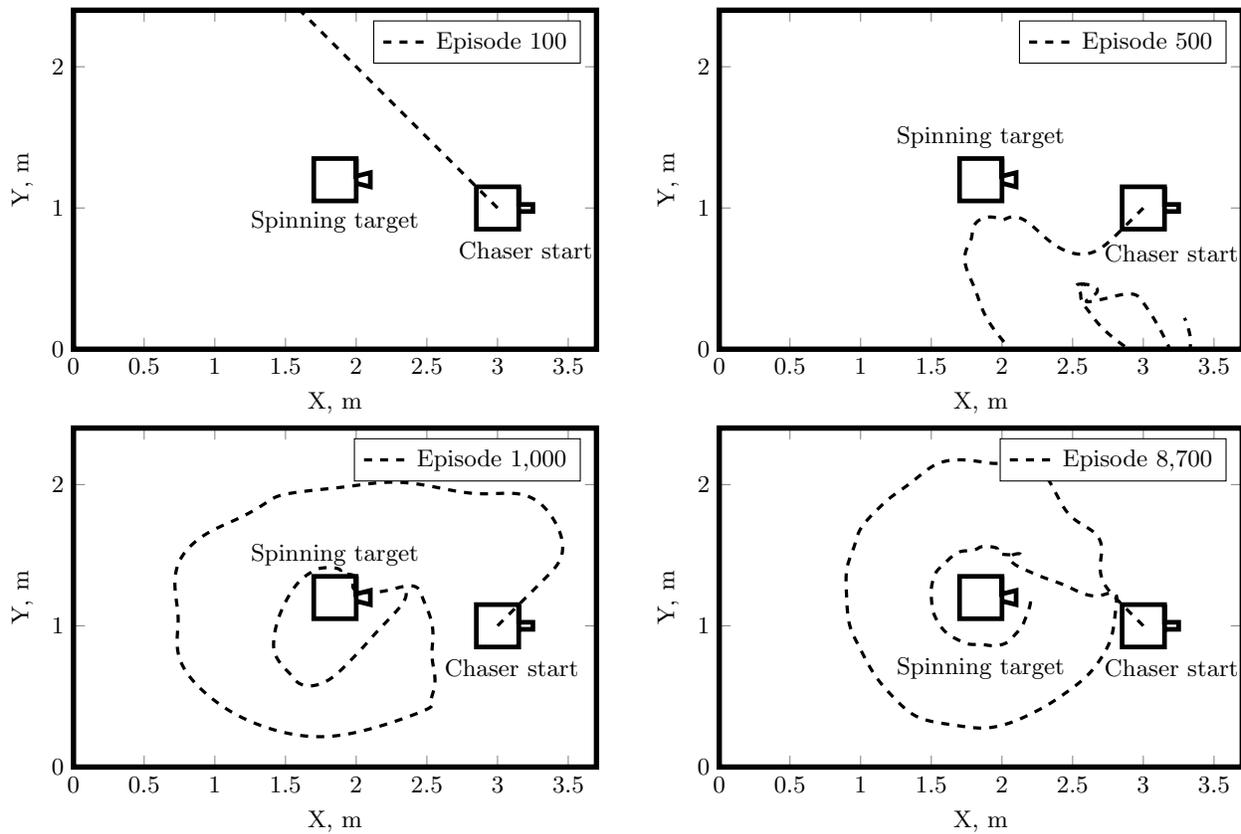


Figure 10. Visualization of chaser trajectories during the training process at various episodes.

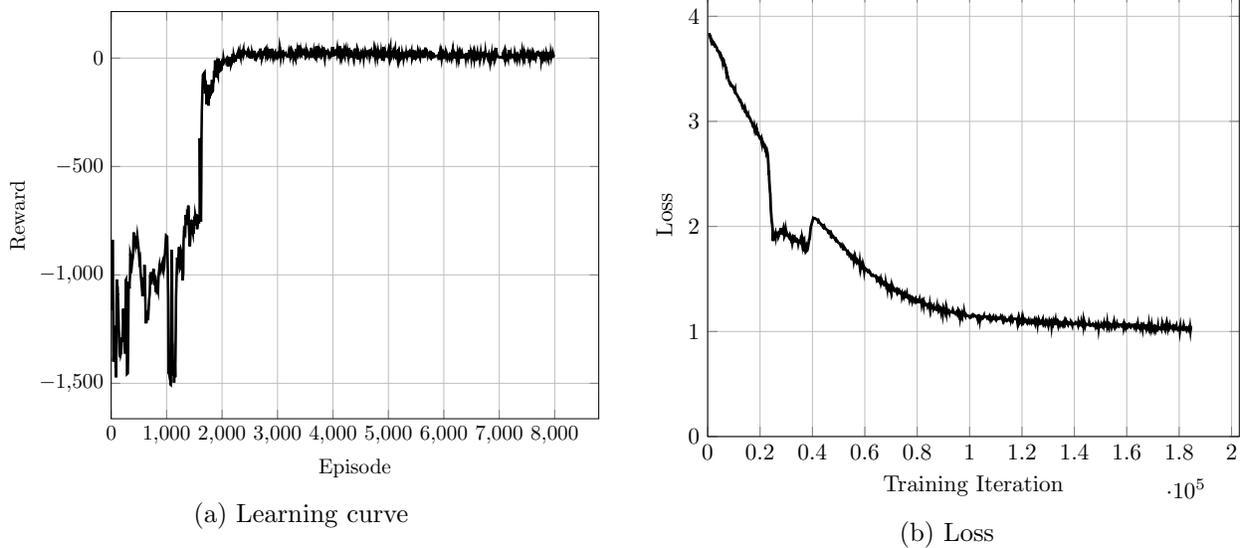


Figure 11. Training performance for chaser pose tracking and docking for a rotating target with randomized initial conditions.

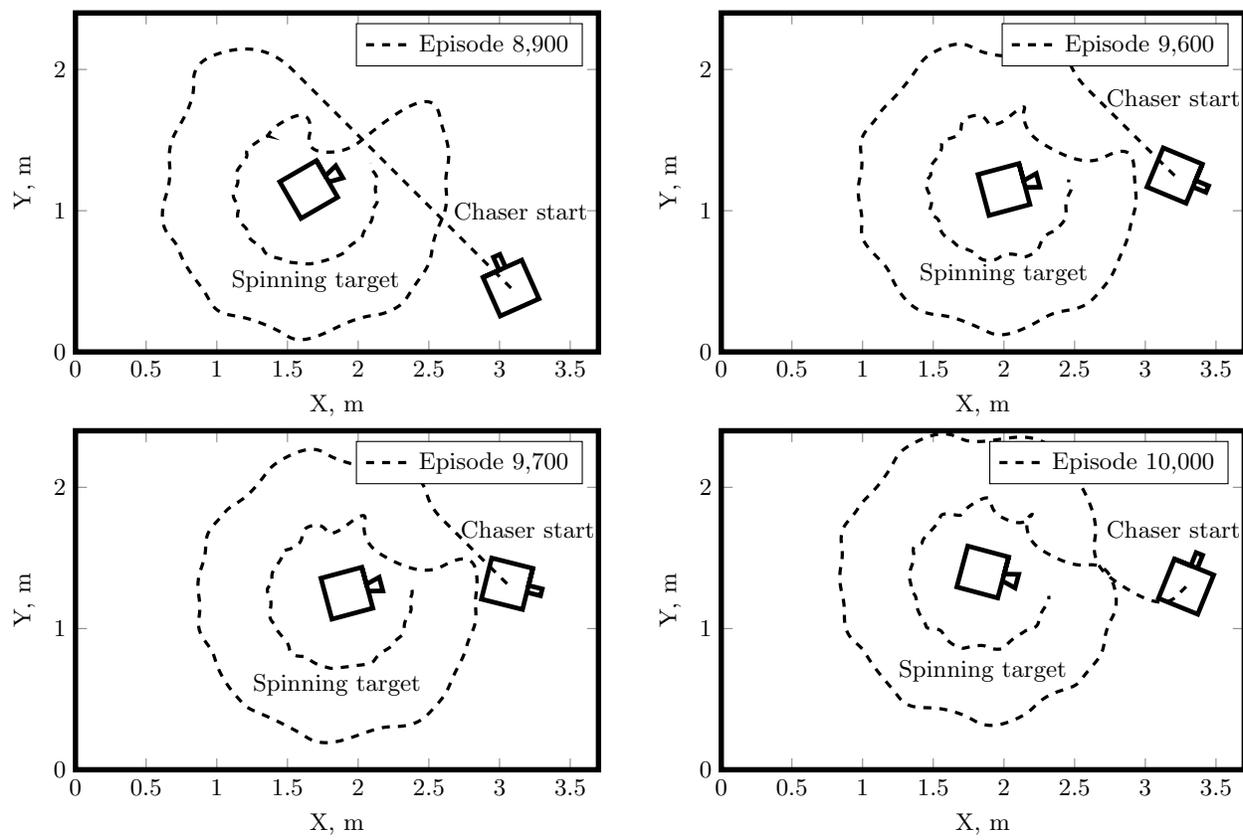


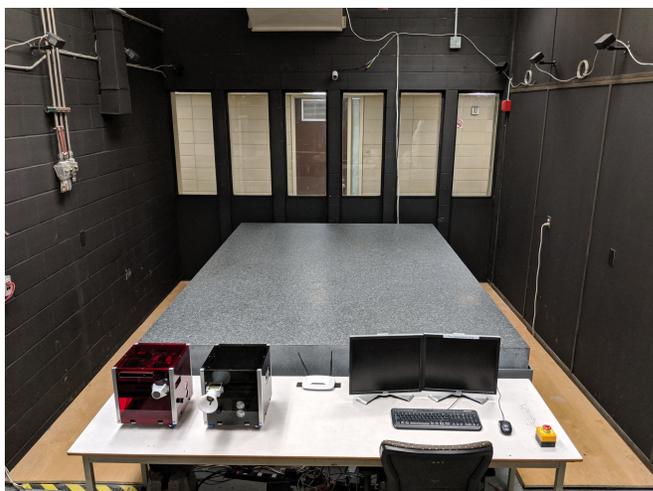
Figure 12. Examples of learned chaser trajectories with randomized initial conditions.

VI. Experimental Validation

To validate the numerical simulations, and to test if the proposed deep guidance strategy can overcome the simulation-to-reality gap, experiments are performed in a laboratory environment at Carleton University. A planar gravity-offset testbed is used, where two spacecraft platforms are positioned on a flat granite surface. Air bearings are used to provide a near friction-free planar environment. The experimental facility is discussed, followed by the experimental setup and results.

VI.A. Experiment Facility

Experiments were conducted at Carleton University's Spacecraft Robotics and Control Laboratory, using the Spacecraft Proximity Operations Testbed (SPOT). Specifically, SPOT consists of two air-bearing spacecraft platforms operating in close proximity on a 2.4 m \times 3.7 m granite surface. The use of air bearings on the platforms reduces the friction to a negligible level. Due to surface slope angles of 0.0026° and 0.0031° along both directions, residual gravitational accelerations of 0.439 and 0.525 mm/s² perturb the dynamics of the floating platforms along the X and Y directions, respectively. Both platforms have dimensions of 0.3 \times 0.3 \times 0.3 m and are actuated by expelling compressed air at 550 kPa (80 psi) through eight miniature air nozzles distributed around each platform, thereby providing full planar control authority. Each thruster generates approximately 0.25 N of thrust and is controlled at a frequency of 500 Hz by a pulse-width modulation scheme using solenoid valves. Pressurized air for the thrusters and the air bearing flotation system is stored on-board in a single air cylinder at 31 MPa (4,500 psi). The structure consists of an aluminum frame with four corner rods on which three modular decks are stacked. To protect the internal components, the structure is covered with semi-transparent acrylic panels. Figure 13a shows the SPOT laboratory facility and Fig. 13b shows two SPOT platforms in a proximity operations configuration.



(a) An overview of the SPOT facility.



(b) SPOT platforms.

Figure 13. The Spacecraft Proximity Operations Testbed (SPOT).

The motion of both platforms is measured in real-time through four active LEDs on each platform which are tracked by an 8-camera PhaseSpace[®] motion capture system. This provides highly accurate ground-truth position and attitude data. All motion capture cameras are connected to a PhaseSpace[®] server, which is connected to a ground station computer. The ground station computer wirelessly communicates ground truth information to the platforms' onboard computers, which consist of Raspberry Pi-3s running the Raspbian Linux operating system. Based on the position and attitude data the platforms wirelessly receive, they can perform feedback control, by calculating the required thrust to maneuver autonomously and actuating the appropriate solenoid valves to realize this motion. The ground station computer also receives real-time telemetry data (i.e., any signals of interest, as specified by the user) from all on-board computers, for post-experiment analysis purposes.

A MATLAB/Simulink[®] numerical simulator that recreates the dynamics and emulates the different

on-board sensors and actuators is first used to design and test the upcoming experiment. Once the performance in simulations is satisfactory, the control software is converted into C/C++ using Embedded Coder[®], compiled, and then executed on the platforms' Raspberry Pi-3 computers.

An Nvidia Jetson TX2 Module is used to run the trained deep guidance policy neural network in real-time. It accepts the current system state and returns a guidance velocity signal to the Raspberry Pi-3 that executes a control law to track the commanded velocity.

VI.B. Setup

The simulations presented in Sec. V were chosen such that they are replicable experimentally. Both pose tracking and docking tasks are attempted in experiment. The final parameters, θ , of the deep guidance policies trained with randomized initial conditions are exported for use on the chaser SPOT platform. Initial conditions are similar to those used previously in simulations.

The platforms remain in contact with the table until a strong lock has been acquired on the LEDs by the motion capture system. Following this, the platforms begin to float and maneuver to the desired initial conditions. Then, the target remains stationary or begins rotating while the chaser platform uses the deep guidance policy trained in simulation to guide itself towards the hold point and finally the docking point on the target.

It should be noted that a significant number of discrepancies exist between the simulated environment the policy was trained within and the experimental facility. The simulated environment did not account for the discrete thrusters and their limitations, the control thrust allocation strategy, signal noise, system delays, friction, air resistance, centre of mass offsets, thruster plume interaction, and table slope. In addition, the spacecraft mass used to evaluate the training was 10 kg whereas the experimental spacecraft platforms have a mass of 16.9 kg. Significant discrepancies exist between the simulated training environment and the experimental environment, which is an excellent test for the simulation-to-reality capabilities of the proposed deep guidance technique.

VI.C. Results

A trajectory of the experiment with a stationary target is shown in Fig. 14a. It shows the the deep guidance successfully outputs velocity commands that brings the chaser to the hold point and then additional velocity commands to move towards the target docking port. A video of experiment 1 can be found at: <https://youtu.be/rk6sY1D7XWU>.

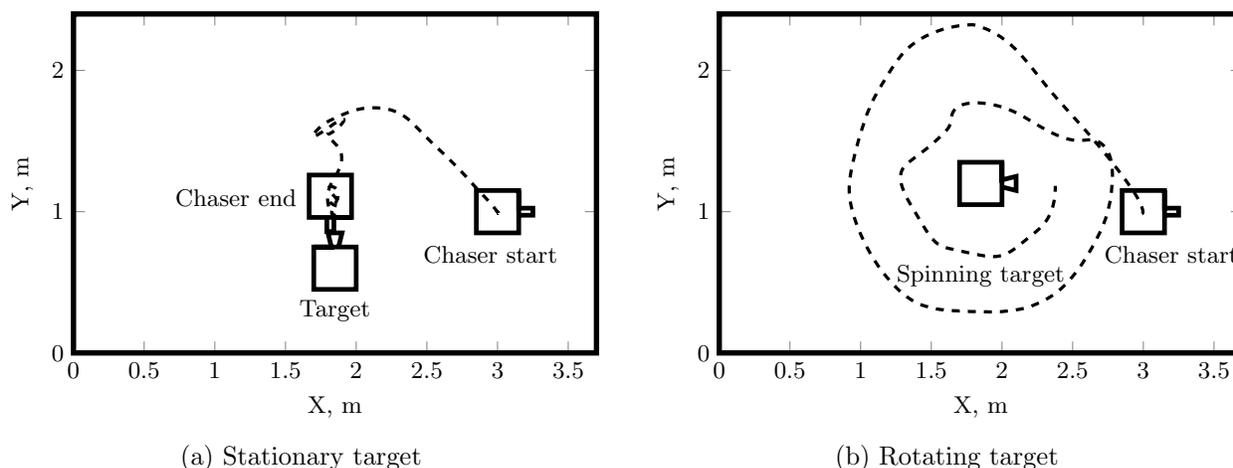


Figure 14. Experimental trajectories of deep guidance in Carleton's Spacecraft Proximity Operations Testbed.

A trajectory of the experiment with a rotating target is shown in Fig. 14b. The learned deep guidance technique successfully commands an appropriate velocity signal for the chaser to complete the task. A video of the motion can be found at <https://youtu.be/25aIkRv9GdY>

The deep guidance technique was successfully trained exclusively in simulation and deployed to an experimental facility and achieved similar performance to that during training. Combining the neural-network

guidance with conventional control allowed the trained system to handle unmodeled effects present in the experiment. The proposed technique successfully demonstrates the deep guidance technique as a viable solution to the simulation-to-reality problem present in deep reinforcement learning.

VII. Conclusion

This paper introduced deep reinforcement learning to the guidance problem for spacecraft robotics. Through training a guidance policy to accomplish a goal, complex behaviours can be learned rather than hand-crafted. The simulation-to-reality gap dictates that policies trained in simulation often do not transfer well to reality. To avoid this, the authors restrict the policy to output a guidance signal, which a conventional controller is tasked with tracking. Conventional control can handle the modelling errors that typically plague reinforcement learning. This paper tests this learned guidance technique, which the authors call *deep guidance*, on a simple problem. That is, spacecraft pose tracking and docking. Numerical simulations show that the proximity operation tasks can be successfully learned using the deep guidance technique. The trained policies are then deployed to two experiments, with comparable results to those in simulation even though the simulated environment did not model all effects present in the experimental facility. Future work will further explore the generality of the technique and its use on more difficult problems.

Acknowledgments

This research was financially supported in part by the Natural Sciences and Engineering Research Council of Canada under the Postgraduate Scholarship-Doctoral PGSD3-503919-2017 award.

References

- ¹Flores-Abad, A., Ma, O., Pham, K., and Ulrich, S., “A Review of Space Robotics Technologies for On-orbit Servicing,” *Progress in Aerospace Sciences*, Vol. 68, 2014, pp. 1–26.
doi:10.1016/j.paerosci.2014.03.002
- ²Aghili, F., “A Prediction and Motion-planning Scheme for Visually Guided Robotic Capturing of Free-floating Tumbling Objects with Uncertain Dynamics,” *IEEE Transactions on Robotics*, Vol. 28, No. 3, 2012, pp. 634–649.
doi:10.1109/TRO.2011.2179581
- ³Wilde, M., Ciarcià, M., Grompone, A., and Romano, M., “Experimental Characterization of Inverse Dynamics Guidance in Docking with a Rotating Target,” *Journal of Guidance, Control, and Dynamics*, Vol. 39, No. 6, 2016, pp. 1173–1187.
doi:10.2514/1.g001631
- ⁴Ma, Z., Ma, O., and Shashikanth, B. N., “Optimal Approach to and Alignment with a Rotating Rigid Body for Capture,” *Journal of the Astronautical Sciences*, Vol. 55, No. 4, 2007, pp. 407–419.
doi:10.1007/BF03256532
- ⁵Dong, H., Hu, Q., and Akella, M. R., “Dual-Quaternion-Based Spacecraft Autonomous Rendezvous and Docking Under Six-Degree-of-Freedom Motion Constraints,” *Journal of Guidance, Control, and Dynamics*, Vol. 41, No. 5, 2018, pp. 1150–1162.
doi:10.2514/1.g003094
- ⁶Filipe, N. and Tsiotras, P., “Adaptive Position and Attitude Tracking Controller for Satellite Proximity Operations using Dual Quaternions,” *Journal of Guidance, Control, and Dynamics*, Vol. 38, No. 4, 2015, pp. 566–577.
doi:10.2514/1.G000054
- ⁷Gui, H. and Vukovich, G., “Finite-time Output-Feedback Position and Attitude Tracking of a Rigid Body,” *Automatica*, Vol. 74, 2016, pp. 270–278.
doi:10.1016/j.automatica.2016.08.003
- ⁸Pothen, A. A. and Ulrich, S., “Close-range Rendezvous with a Moving Target Spacecraft Using Udwadia-Kalaba Equation,” *American Control Conference*, American Automatic Control Council, 2019, pp. 3267–3272.
- ⁹Pothen, A. A. and Ulrich, S., “Pose Tracking Control for Spacecraft Proximity Operations Using the Udwadia-Kalaba Framework,” *AIAA Guidance, Navigation, and Control Conference*, Orlando, FL, 2020.
- ¹⁰Hough, J. and Ulrich, S., “Lyapunov Vector Fields for Thrust-Limited Spacecraft Docking with an Elliptically-Orbiting Uncooperative Tumbling Target,” *AIAA Guidance, Navigation, and Control Conference*, Orlando, FL, 2020.
- ¹¹Hornik, K., Stinchcombe, M., and White, H., “Multilayer Feedforward Networks are Universal Approximator,” *Neural Networks*, Vol. 2, No. 5, 1989, pp. 359–366.
doi:10.1016/0893-6080(89)90020-8
- ¹²Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D., “Human-Level Control Through Deep Reinforcement Learning,” *Nature*, Vol. 518, No. 7540, 2015, pp. 529–533.
doi:10.1038/nature14236
- ¹³Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M., “The Arcade Learning Environment: An Evaluation Platform

for General Agents,” *Journal of Artificial Intelligence Research*, Vol. 47, 2013, pp. 253–279.

doi:10.1613/JAIR.3912

¹⁴Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., and Hassabis, D., “Mastering the Game of Go with Deep Neural Networks and Tree Search,” *Nature*, Vol. 529, No. 7587, 2016, pp. 484–489.

doi:10.1038/nature16961

¹⁵Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., Van Den Driessche, G., Graepel, T., and Hassabis, D., “Mastering the Game of Go Without Human Knowledge,” *Nature*, Vol. 550, No. 7676, 2017, pp. 354–359.

doi:10.1038/nature24270

¹⁶Kober, J., Bagnell, J. A., and Peters, J., “Reinforcement Learning in Robotics: A Survey,” *The International Journal of Robotics Research*, Vol. 32, No. 11, 2015, pp. 1238–1274.

doi:10.1177/0278364913495721

¹⁷Tai, L., Paolo, G., and Liu, M., “Virtual-to-real Deep Reinforcement Learning: Continuous Control of Mobile Robots for Mapless Navigation,” *IEEE International Conference on Intelligent Robots and Systems*, Vancouver, Canada, 2017, pp. 31–36.

doi:10.1109/IROS.2017.8202134

¹⁸Sünderhauf, N., Brock, O., Scheirer, W., Hadsell, R., Fox, D., Leitner, J., Upcroft, B., Abbeel, P., Burgard, W., Milford, M., and Corke, P., “The Limits and Potentials of Deep Learning for Robotics,” *Computing Research Repository*, 2018.

arXiv:1804.06557

¹⁹Bousmalis, K., Irpan, A., Wohlhart, P., Bai, Y., Kelcey, M., Kalakrishnan, M., Downs, L., Ibarz, J., Pastor, P., Konolige, K., Levine, S., and Vanhoucke, V., “Using Simulation and Domain Adaptation to Improve Efficiency of Deep Robotic Grasping,” *Computing Research Repository*, 2017.

arXiv:1709.07857

²⁰OpenAI, “Learning Dexterous In-Hand Manipulation,” *Computing Research Repository*, 2018.

arXiv:1808.00177

²¹Lee, J., Hwangbo, J., and Hutter, M., “Robust Recovery Controller for a Quadrupedal Robot using Deep Reinforcement Learning,” *Computing Research Repository*, 2019.

arXiv:1901.07517

²²OpenAI, Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., Schneider, J., Tezak, N., Tworek, J., Welinder, P., Weng, L., Yuan, Q., Zaremba, W., and Zhang, L., “Solving Rubik’s Cube with a Robot Hand,” *Computing Research Repository*, 2019, pp. 1–51.

arXiv:1910.07113

²³Cutler, M. and How, J. P., “Autonomous Drifting Using Simulation-aided Reinforcement Learning,” *IEEE International Conference on Robotics and Automation*, Stockholm, Sweden, 2016, pp. 5442–5448.

doi:10.1109/ICRA.2016.7487756

²⁴Hwangbo, J., Sa, I., Siegwart, R., and Hutter, M., “Control of a Quadrotor with Reinforcement Learning,” *IEEE Robotics and Automation Letters*, Vol. 2, No. 4, 2017, pp. 2096–2103.

doi:10.1109/LRA.2017.2720851

²⁵Julian, K. D. and Kochenderfer, M. J., “Distributed Wildfire Surveillance with Autonomous Aircraft Using Deep Reinforcement Learning,” *Journal of Guidance, Control, and Dynamics*, Vol. 42, No. 8, 2019, pp. 1768–1778.

doi:10.2514/1.G004106

²⁶Chan, D. M. and Agha-Mohammadi, A. A., “Autonomous Imaging and Mapping of Small Bodies Using Deep Reinforcement Learning,” *IEEE Aerospace Conference*, Big Sky, MT, 2019, pp. 1–12.

doi:10.1109/AERO.2019.8742147

²⁷Willis, S., Izzo, D., and Hennes, D., “Reinforcement Learning for Spacecraft Maneuvering Near Small Bodies,” *Advances in the Astronautical Sciences*, Vol. 158, Napa, CA, 2016, pp. 1351–1368.

²⁸Lafarge, N. B., Miller, D., Howell, K. C., and Linares, R., “Guidance for Closed-Loop Transfers using Reinforcement Learning with Application to Libration Point Orbits,” *AIAA Guidance, Navigation, and Control Conference*, Orlando, FL, 2020, pp. 1–5.

²⁹Scorsoglio, A., Furfaro, R., Linares, R., and Gaudet, B., “Image-based Deep Reinforcement Learning for Autonomous Lunar Landing,” *AIAA Guidance, Navigation, and Control Conference*, Orlando, FL, 2020.

³⁰Gaudet, B. and Furfaro, R., “Adaptive Pinpoint and Fuel Efficient Mars Landing Using Reinforcement Learning,” *IEEE/CAA Journal of Automatica Sinica*, Vol. 1, No. 4, 2014, pp. 397–411.

doi:10.1109/JAS.2014.7004667

³¹Furfaro, R., Simo, J., Gaudet, B., and Wibben, D. R., “Neural-based Trajectory Shaping Approach for Terminal Planetary Pinpoint Guidance,” *AAS/AIAA Astrodynamics Specialist Conference*, Hilton Head, SC, 2013, pp. 1–18.

³²Sánchez-Sánchez, C. and Izzo, D., “Real-time Optimal Control via Deep Neural Networks: Study on Landing Problems,” *Journal of Guidance, Control, and Dynamics*, Vol. 41, No. 5, 2018, pp. 1122–1135.

doi:10.2514/1.G002357

³³Ulrich, S., Saenz-Otero, A., and Barkana, I., “Passivity-based Adaptive Control of Robotic Spacecraft for Proximity Operations Under Uncertainties,” *Journal of Guidance, Control, and Dynamics*, Vol. 39, No. 6, 2016, pp. 1441–1450.

doi:10.2514/1.G001491

³⁴Jafarnejadsani, H., Sun, D., Lee, H., and Hovakimyan, N., “Optimized L1 Adaptive Controller for Trajectory Tracking of an Indoor Quadrotor,” *Journal of Guidance, Control, and Dynamics*, Vol. 40, No. 6, 2017, pp. 1415–1427.

doi:10.2514/1.G000566

³⁵Huang, P., Wang, D., Meng, Z., Zhang, F., and Guo, J., “Adaptive Postcapture Backstepping Control for Tumbling Tethered Space Robot-Target Combination,” *Journal of Guidance, Control, and Dynamics*, Vol. 39, No. 1, 2016, pp. 150–156. doi:10.2514/1.G001309

³⁶Harris, A., Teil, T., and Schaub, H., “Spacecraft Decision-Making Autonomy using Deep Reinforcement Learning,” *AAS/AIAA Space Flight Mechanics Meeting*, 2019, pp. 1–19.

³⁷Barth-Maron, G., Hoffman, M. W., Budden, D., Dabney, W., Horgan, D., TB, D., Muldal, A., Heess, N., and Lillicrap, T., “Distributed Distributional Deterministic Policy Gradients,” *International Conference on Learning Representations*, Vancouver, Canada, 2018. arXiv:1804.08617

³⁸Bellemare, M. G., Dabney, W., and Munos, R., “A Distributional Perspective on Reinforcement Learning,” *Computing Research Repository*, 2017. arXiv:1707.06887

³⁹Mnih, V., Badia, A., Mirza, M., Graves, A., and Lillicrap, T., “Asynchronous Methods for Deep Reinforcement Learning,” *International Conference on Machine Learning*, New York, NY, 2016. arXiv:1602.01783

⁴⁰Sutton, R. S. and Barto, A. G., *Reinforcement Learning: An Introduction*, Cambridge: MIT Press, 1998.

⁴¹Zappulla, R., Park, H., Virgili-Llop, J., and Romano, M., “Real-Time Autonomous Spacecraft Proximity Maneuvers and Docking Using an Adaptive Artificial Potential Field Approach,” *IEEE Transactions on Control Systems Technology*, Vol. 27, No. 6, 2019, pp. 2598–2605. doi:10.1109/TCST.2018.2866963

⁴²Ciarcià, M., Grompone, A., and Romano, M., “A Near-optimal Guidance for Cooperative Docking Maneuvers,” *Acta Astronautica*, Vol. 102, 2014, pp. 367–377. doi:10.1016/j.actaastro.2014.01.002

⁴³Mammarella, M., Capello, E., Park, H., Guglieri, G., and Romano, M., “Tube-based Robust Model Predictive Control for Spacecraft Proximity Operations in the Presence of Persistent Disturbance,” *Aerospace Science and Technology*, Vol. 77, 2018, pp. 585–594. doi:10.1016/j.ast.2018.04.009

⁴⁴Saulnier, K., Pérez, D., Huang, R. C., Gallardo, D., Tilton, G., and Bevilacqua, R., “A Six-degree-of-freedom Hardware-in-the-loop Simulator for Small Spacecraft,” *Acta Astronautica*, Vol. 105, No. 2, 2014, pp. 444–462. doi:10.1016/j.actaastro.2014.10.027

⁴⁵Oliphant, T. E., “Python for Scientific Computing,” *Computing in Science & Engineering*, Vol. 9, No. 3, 2007, pp. 10–20. doi:10.1109/MCSE.2007.58

⁴⁶Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D., “Continuous Control with Deep Reinforcement Learning,” *Computing Research Repository*, 2016. arXiv:1509.02971

⁴⁷Kingma, D. P. and Ba, J., “Adam: A Method for Stochastic Optimization,” *International Conference on Learning Representations*, San Diego, CA, 2015. arXiv:1412.6980

⁴⁸Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jozefowicz, R., Jia, Y., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Schuster, M., Monga, R., Moore, S., Murray, D., Olah, C., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X., “TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems,” *Computing Research Repository*, 2016, Software available from tensorflow.org. arXiv:1603.04467