

Spacecraft Command and Control with Safety Guarantees using Shielded Deep Reinforcement Learning

Andrew Harris* and Hanspeter Schaub†

Increasingly complex space missions have motivated the development of autonomous command and control approaches which must handle high-dimensional, continuous observation and action spaces with hard-to-analyze behavior. Deep reinforcement learning (DRL) techniques are a rising area of research for dealing with such problems, but lack performance and safety guarantees which reduce their applicability for spacecraft operations. This work identifies promising strategies from the DRL literature for providing safety and performance guarantees using correct-by-construction shield synthesis, and techniques for identifying the robustness and stability of trained agents in a computationally expedient manner. Additionally, open-source spacecraft simulation tools derived from the AVS Basilisk astrodynamics simulation package are presented and discussed. Shielded learning agents are presented against a naive DRL agent approach for the command and control of a LEO ground observation mission and compared on the basis of performance, computational efficiency, and safety.

I. Introduction

Autonomous operation of spacecraft has emerged as a major priority among space agencies and private companies alike. Decades of development have yielded few missions that approach the goal of full autonomy. While support tools for operators become increasingly sophisticated, next-generation autonomy for space mission operations will require the introduction of artificial intelligence to further supplant or replace the role of human operators. Advances in artificial intelligence have demonstrated the capability for human-level reasoning on complex tasks; however, systematic constraints within the space domain limit their immediate applicability. Unlike terrestrial applications, which can benefit from on-the-ground debugging, spacecraft operational autonomy approaches must meet high bars of verification and validation to prevent mission failure.

This work focuses on three criteria to enable the verification and validation of future spacecraft autonomy approaches: providing safety guarantees for deep-learning based management agents, simulating spacecraft operations to high fidelity in a computationally efficient way, and automatically validating the performance of trained management agents.

*Research Assistant, Ann and H.J. Smead Department of Aerospace Engineering Sciences, University of Colorado Boulder, Boulder, CO, 80309 USA.

†Glenn L. Murphy Chair of Engineering, Smead Department of Aerospace Engineering Sciences, University of Colorado, 431 UCB, Colorado Center for Astrodynamics Research, Boulder, CO 80309-0431. AAS Fellow, AIAA Fellow.

These foci directly extend the applicability of prior work in the use of Deep Reinforcement Learning (DRL) to create autonomous, self-improving management agents to operate spacecraft without human input. Additionally, these tools should not impact the performance of the learning agent to preserve the advantages of learning-based approaches for creating decision-making agents.

As outlined in prior work [1, 2], machine learning techniques such as Deep Reinforcement Learning have the potential to improve upon current state-of-the-art approaches to operational autonomy by combining the best attributes of optimization- and rule-driven autonomy. At present, operators typically rely on either pre-defined transition criteria that must be carefully set by experts ([3, 4]) or utilize computationally intense planning and scheduling algorithms that depend on high-fidelity models and analysis conducted by experts [5–8]. In contrast, autonomy approaches based around learning agents are able to deal directly with high or infinite dimensional input/output spaces and non-linear and non-smooth problem dynamics [9] so long as a simulator can be constructed. Additionally, DRL-based decision agents encode their knowledge in deep neural networks, which can be rapidly evaluated in constant time once trained.

These benefits have led to a large and growing body of work on the use of machine learning “agents” for planning and scheduling. Multiple works ([10, 11]) have investigated the use of reinforcement learning to address UAV-based sensor tasking and health management problems. Related areas, such as sensor tasking for space situational awareness, have also been addressed using deep reinforcement learning [12, 13]. Finally, the recent success of the OpenAI Five at winning games which require high levels of strategy and planning [14] from self-play alone suggests that reinforcement learning techniques can tackle extremely complex, dynamic, continuous state and action spaces and still replicate or exceed human-level performance.

A small collection of other works in the application of machine learning techniques to spacecraft problems exists in the recent literature, mostly focusing on the application of learning approaches to control problems in uncertain environments. Several works [1, 15] consider reinforcement learning in the context of autonomous aerobraking controllers, with mixed results. Others explore machine learning techniques for asteroid proximity operations [16] or autonomous lunar landing [17]. Importantly, these approaches have focused on low-level control with reinforcement learning, an area that has been traditionally been addressed by conventional estimation and control techniques which can offer boundaries or guarantees on performance. In contrast, this work explicitly examines applications of reinforcement learning to high-level spacecraft planning and decision-making problems.

The existing reinforcement learning literature is primarily concerned with sample efficiency, which drives computational costs and therefore implementation difficulty. Safety and robustness of these ML-driven systems is a rising concern. While many real systems share common constraints on the learning process [18], unique features of the space mission life-cycle create additional factors and constraints that motivate this work. Unlike many other reinforcement learning domains, fairly accurate a-priori models of system behavior are well known for space systems [2], or are at least bounded by mission requirements. At the same time, statuses that would cause a mission to fail, such as a low-power

condition, must be avoided at all costs during agent execution and are desirable to avoid in training for similar reasons.

This work is arranged as follows. First, the spacecraft command and control problem is reformulated as a Partially-Observable Markov Decision Problem (POMDP) amicable to deep reinforcement learning techniques; given this problem formulation, the concept of a one-step sequential decision agent is introduced alongside plausible concepts of implementation. Next, a simulation framework for spacecraft-driven POMDPs is presented for the benefit of current and future research in this area. Next, the theory of shield synthesis is reviewed and applied to the spacecraft control problem to provide safety guarantees during learning and deployment. The efficacy of this strategy versus a naive learning agent is presented through simulations on the aforementioned simulation framework.

II. Problem Statement

A. Command and Control Framework

Traditional spacecraft operations planning and execution is a complex, multi-step process with many stakeholders which relies heavily on expert knowledge. For reference, a generic version of this paradigm is presented here. First, mission stakeholders specify mission objectives and a reference mission trajectory. Given this trajectory and a set of desired tasks, a set of activities or operational modes are defined and scheduled as spacecraft resources (power, fuel, compute time) and mission resources (observation/maneuver/communication windows) permit. Finally, these plans are converted into an action sequence, up-linked to a spacecraft, and executed by on-board software. Simultaneously, teams of human operators typically monitor mission execution and spacecraft health parameters and intervene when parameters fall outside of a defined specification, either directly by changing the current action sequence or indirectly by initiating a re-planning sequence. While this process or processes like it have been used successfully for decades, it relies heavily on human expertise to create priorities, construct action sequences, and verify spacecraft behavior. In the

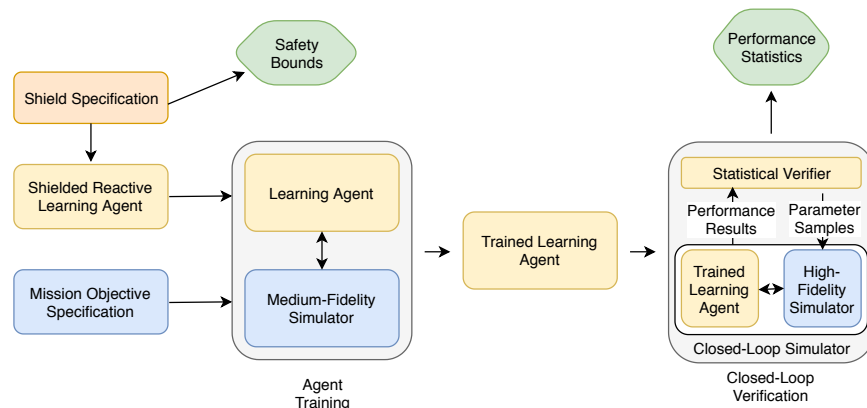


Fig. 1 End-to-End training pipeline for autonomous operations agents.

search for future autonomy approaches, it is desirable to both replicate existing capabilities in the realm of rule-based and optimization-oriented autonomy while improving their extensibility, robustness to un-modeled dynamics, and

computational burden. To provide a feasible scope, this work specifically considers the mission-level decision-making problem wherein sub-plans (“modes”) have already been identified, either by some other planning routine or by designers pre-flight. In this context, a decision-making agent must account not only for mission objectives, but also the constraints imposed by spacecraft hardware, orbital and attitude mechanics, and uncertainty regarding known or unknown environmental parameters.

A common framework for representing and addressing such problems are Partially-Observable Markov Decision Processes, which compactly represent the problems facing a software agent acting in an evolving environment according to some higher-level objective[19]. The mathematics of such processes, and challenges associated with them, are reviewed briefly here.

A model of several time-steps of a classical POMDP is presented in Fig. 2, and discussed further here. As in traditional Markov Decision Processes (MDPs), the state in a POMDP is updated by a transition function F , and at any given time can be computed as a function of the previous state and the most recent action taken by the considered agent(s):

$$s_k = F(s_{k-1}, a_{k-1}) \quad (1)$$

This state s_k is observed by the agent according to some observation function H :

$$o_k = H(s_k) \quad (2)$$

Given an observation o_k of the state, the agent then selects an action a_k to influence the future state according to some policy π :

$$a_k = \pi(o_k) \quad (3)$$

While these transition functions represent physical or software-defined process dynamics, the objective of an agent is ultimately motivated by a reward function R :

$$r_k = R(s_{k-1}, a_{k-1}, s_k) \quad (4)$$

Taken together, these components form a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{O}, F, H, \mathcal{R}.)$

The objective of a software agent within a POMDP is to select a policy π that maximizes its realized reward. While the general POMDP case places no restrictions on the nature of any of the transition functions or states, the consideration of infinite-dimensional, continuous state and action spaces can be extremely computationally intensive. For this reason, many applied autonomy approaches that leverage POMDPs perform some degree of discretization to their state or action space. Additionally, it is noted that POMDPs attempt to describe holistic, system-level problems within a unified

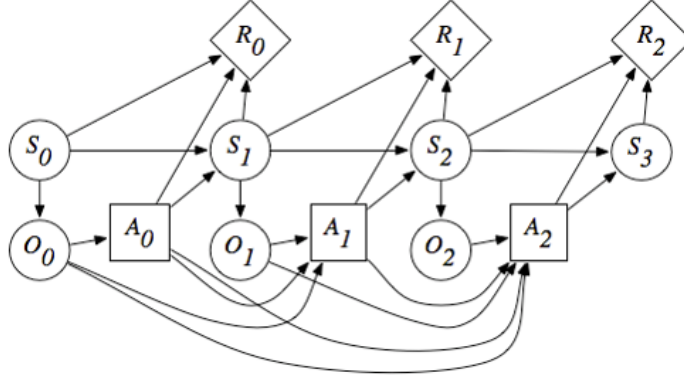


Fig. 2 Sequential Partially Observable Markov Decision process framework for representing decision problems.

framework that is theoretically related to but practically divorced from traditional estimation and control approaches. For these reasons, POMDP-based approaches to autonomy are most frequently studied in cases where traditional estimation and controls approaches are not readily tractable, including human-assisted machine decision-making [11] or multi-vehicle coordination problems [20].

A representative space orbit operations scenario is presented and defined as a partially-observable Markov decision process (POMDP) to be solved using Proximal Policy Optimization (PPO) For a spacecraft, the general high-level autonomy POMDP can be stated as follows. Given the constraints of orbital dynamics, on-board hardware, and pre-defined software behaviors, select the sequence of behaviors that best satisfies mission objectives. This approach is described in greater detail in [2].

B. Deep Reinforcement Learning

Astrodynamics and spacecraft-planning problems are typically considered in the context of continuous estimation and control, as many of the processes facing such systems are infinite-dimensional with well-understood, reasonably accurate models. Unfortunately, the high-level relationships between spacecraft actions and the satisfaction of mission objectives is less analytically tractable, and frequently mixes discrete reward states (such as whether a geological feature has been imaged) with continuous ones (such as the management of spacecraft power states). Reinforcement Learning (RL) techniques represent one class of algorithms for addressing decision processes which lack analytically tractable models; however, traditional RL techniques require discrete state and action spaces. Prior work [1] has shown that accurate discrete state spaces for spacecraft decision problems can be extremely large and therefore infeasible to explore. To address these shortcomings, Deep Reinforcement Learning (DRL) techniques use deep neural networks as function approximators in place of tables and can therefore learn on continuous state or action spaces without discretization.

Reinforcement Learning techniques are intended to solve general Markov Decision Processes (MDPs), which are simplified forms of POMDPs without the issue of observation functions. The goal of reinforcement learning is to find an optimal policy π^* that maximizes the expected future reward of the agent. The optimal policy, π^* , is the policy with

the largest expected sum of rewards or value function. The cumulative value, V , of a given state is provided by the discounted sum of the rewards from the current infinitely into the future and is given below:

$$V(s_0, s_1, s_2, \dots) = \sum_{t=0}^{\infty} \gamma^t r_t \quad 0 \leq \gamma < 1 \quad (5)$$

where γ is the reward discount factor. This term weights the importance of future rewards relative to the current reward. Given this framework, the optimal policy is that which maximizes the expected discounted future reward.

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \quad (6)$$

This leads to another expression for the cumulative value function referred to as Bellman's Equation:

$$V(s) = R(s) + \gamma \max_a \sum_{s'} p(s'|s, a) V(s') \quad (7)$$

where $p(s'|s, a)$ is the probability of the agent being in state, s' , after performing action, a , in state s .

This work uses Proximal Policy Optimization [21] as implemented by the `stable-baselines` Python package [22], an extended variant of Trust-Region Policy Optimization. PPO uses a loss function that penalizes the learning agent from dramatically changing its network weights from iteration to iteration using advantage estimation and a clipping function:

$$L^{CLIP}(\theta) = \hat{E}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)] \quad (8)$$

where θ represents the policy parameters, \hat{E}_t represents the learned expectation over timesteps, r_t is the ratio of the probability under the new and old policies, \hat{A}_t is the estimated advantage at t , and ϵ is a hyperparameter representing the clipping range. The advantage function is defined as the difference between the state-action value function $Q^\pi(s, a)$ and the state value function $V^\pi(s)$. This approach has been shown to produce faster, more reliable convergence than other results, and represents the state-of-the-art in model-free deep reinforcement learning.

C. Agent Implementation Frameworks

A major assumption in our formulation of the spacecraft control problem as a (PO)MDP shown in Eqn. 11 is the discretization of time, which—when combined with the mechanics of learning as described in Section II.B—results in decision-making agents that can only *react* to current observations, as shown in Fig. 3. Rather than utilizing a specific plan or strategy, all relevant planning and strategy information is encoded in the deep network utilized by the agent. In practice, evaluating neural networks is nearly constant-time and can be readily hardware-accelerated, making this implementation attractive for future on-board use where system information is readily available and humans are already

out of the loop.

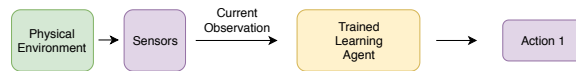


Fig. 3 Sequential decision-making agent architecture.

At the same time, many existing systems assume that discrete sets of actions will be periodically up-linked from the ground and lack the on-board processing power to evaluate a neural network. For these systems, an architecture which uses a ground-side simulator to propagate forward existing observations and actions is proposed as shown in Fig. 4. The incorporation of a simulator allows for the agent to make “future” decisions based on current knowledge and plan ahead. This architecture is also attractive for near-term implementation, as it allows human operators to verify and validate action sequences in advance of execution.

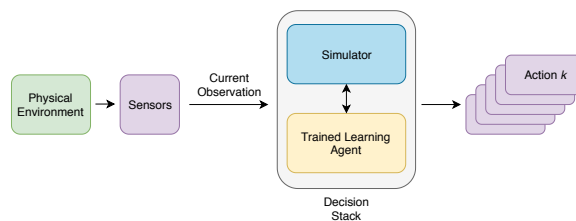


Fig. 4 Planning architecture using a sequential decision-making agent.

Examination of the properties and benefits of planning versus reactive agents is left outside the scope of this work, which focuses on establishing training and safety properties for DRL-based sequential decision-making agents for spacecraft command and control.

D. Safety Guarantees

Safety in the face of uncertain spacecraft performance, environmental parameters, and operating sequences is a critical requirement for future autonomy architectures. While some reinforcement learning techniques can bound their performance with respect to a reward function within an MDP, there are virtually none which can guarantee safety on their own. In practice, this is dealt with through reward engineering; unsafe action or state combinations are given large costs or penalties to achieved reward. This approach has several key disadvantages: many problems for which reinforcement learning is well-suited have complex environment/reward interactions, which makes manual reward engineering difficult; when reward engineering is feasible, it does not prevent the agent from taking unsafe actions in conditions outside the training set presented by its environment; finally, there is no quantifiable boundary or degree of safety provided through reward engineering. These shortcomings have motivated the search for alternative approaches to safety that can be combined with common DRL approaches.

Reactive synthesis is one category of techniques that can provide performance bounds and guarantees for controllers

on specified systems. In general, reactive synthesis algorithms operate on discrete, known, finite systems and attempt to produce behavior on such systems that satisfies a specification written in a temporal logic language, such as Linear Temporal Logic (LTL). Also described as “correct-by-construction” approaches, reactive synthesis algorithms only produce control policies that meet a given specification; if the specification cannot be met on the current system, no policy will be produced, allowing for designers to check feasibility before implementation. While powerful for addressing systems with discrete, finite, known dynamics, reactive synthesis approaches scale poorly with system and specification complexity, which limits their applicability for solving general spacecraft planning problems, which are difficult to discretize to sufficient fidelity[1].

Shielded learning techniques [23] combine common DRL approaches with reactive synthesis-based shields to combine the power of black-box optimization with formal guarantees of safety. Shielded learning depends on the construction of a coarse, finite-state safety MDP from the original MDP the learning agent is intended to solve that is conservative with respect to the original environment’s dynamics and the safety specification, yet limited enough that reactive synthesis can be applied to it. Next, a safety specification is created using Linear Temporal Logic which encapsulates all desired safety conditions and provided as an input to a reactive synthesis algorithm, such as a two-player game, which produces a discrete, state-dependent strategy. Finally, this strategy is implemented alongside the learning agent as shown in Figure 5; in this implementation, the shield accepts observations of the current system state and the action attempted by the learning agent, and permits the action only if it aligns with the shield’s strategy. This implementation architecture is applicable to both training and on-line use of the sequential decision agent, allowing it to provide safety boundaries during mission execution.

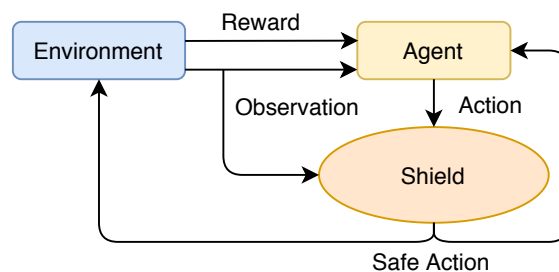


Fig. 5 Post-Posed shielded reinforcement learning framework.

An example of this transformation in practice is shown for a system with two safety-critical dimensions in Fig. 6. Mission designers first identify state combinations that represent mission failure, such as depleting the spacecraft’s battery or allowing reaction wheels to spin up beyond manufacturer’s specifications. In addition to the hard safety boundaries, operators and mission planners typically incorporate additional boundaries to act as margins of safety against actual failure; these are represented by the dashed lines labeled “operational boundary,” which are used to define “warning states.” While in this boundary, operators typically take immediate action to return the system to safe, nominal

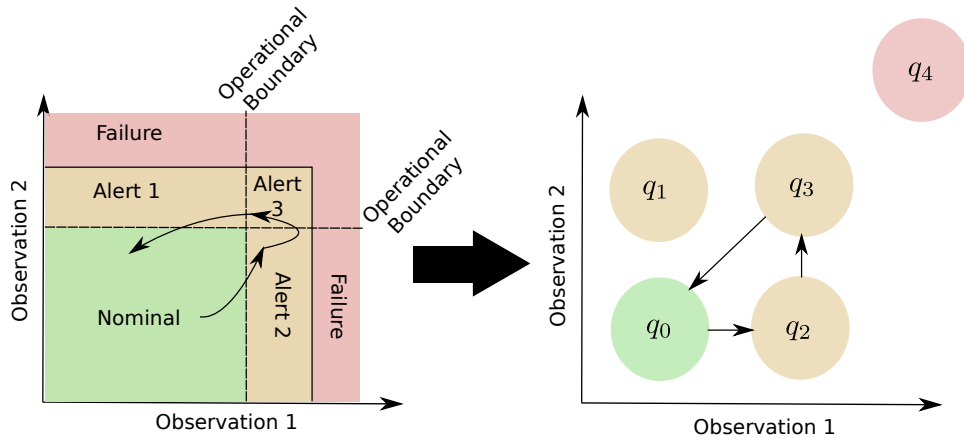


Fig. 6 Conversion from continuous states to a discrete safety MDP.

operating conditions. In this view, the system’s behavior can be plotted on a phase-plot, where individual samples of the system’s true trajectory are represented as curves in the observation variable space. The continuous, bounded system creates a natural framework for the construction of a safety MDP, wherein each warning state becomes a discrete state, including products of warning states. It is important that the safety MDP contain all information necessary for the system to operate safely, which may require the inclusion of states which are not themselves safety risks but which affect the performance of actions necessary for the safety of the system. This process results in a discrete “safety” MDP which exists in parallel with the continuous POMDP.

E. Verification Strategy

Training AI agents necessarily involves the use of models that may contain parametric uncertainties. A key concern with deep reinforcement learning approaches is in preventing over-fitting to certain model behaviors and parameters. To verify that a given AI operations agent is functional within its specification under plausible conditions, a robust verification strategy is required. Here, a two-part verification strategy is used: first, the agent is employed in a higher-fidelity model of the system to check for signs of over-fitting; next, a statistically informed robustness metric is used to allocate these high-fidelity simulation runs over the set of plausible parametric uncertainties.

The field of statistical verification theory is especially applicable to complex cyber-physical systems for which strict analytical guarantees on system performance are not available. Unlike classical model-checking techniques, which both depend on deterministic system models and provide only binary success or failure metrics, statistical verification techniques utilize samples of system performance generated by a simulator alongside a specification to automatically test system performance.

F. Simulation Framework

Both the training process for DRL-based operations agents and the verification framework require the ability to simulate a space mission to high fidelity. DRL techniques in particular can struggle when transferring from simulated to real experiences due to the “simulation gap,” as DRL agents can over-fit on specific attributes of low-fidelity simulators which do not generalize to the real world. In the same vein, verification techniques require the existence of high-fidelity, trusted simulation capability which adequately captures the behavior of the real system. For spacecraft, this requires the ability to simulate not only traditional astrodynamics components (orbital and attitude dynamics), but also the behavior of flight software components.

The Basilisk astrodynamics simulation package represents an ideal toolset for both of these applications. Specifically, Basilisk provides:

- 1) **High Fidelity Astrodynamics:** The Basilisk dynamics engine can simulate fully-coupled multi-body dynamics in tandem with GPU-accelerated orbital dynamics [24], allowing for the simulation of second- and third-order effects like attitude/orbit coupling, fuel slosh, and flexing panels.
- 2) **Flight Software Simulation/Integration:** Developed as a tool to aid flight software development by providing a flight-like environment for testing, Basilisk provides first-class support for the integration of flight software components.
- 3) **Computational Performance:** Compute-heavy code is written in C/C++ and is highly performant as a result; even with tasks like image generation in the loop, BSK-based simulations are thousands of times faster than real-time, allowing for rapid generation of samples for both DRL and verification algorithms.
- 4) **Integration with common ML/RL frameworks:** Basilisk is written with SWIG and provides a Python API for setting up, executing, and analyzing simulations, which allows it to be integrated with other common ML/RL packages (Tensorflow, Keras, gym, scikit-learn).

To facilitate the integration of Basilisk with other machine learning tools, a library of OpenAI gym environments which utilize Basilisk for spacecraft simulation has been created and opened to the public. This library supports common DRL frameworks such as OpenAI’s `baselines` and the `stable-baselines` fork.

III. Performance Comparison

To demonstrate the viability and applicability of deep and shielded learning techniques, this section applies them to a reference problem implemented using the Basilisk deep learning framework and compares their performance in both training (time to convergence, performance with respect to the reward function) and execution (qualitative evaluation of safety).

A. Reference Mission Operations Problem

For the purposes of this work, a scenario consisting of a single spacecraft conducting ground observations of Earth is considered. In general, the goal of the operations agent is to maximize both the time spent pointing at the ground and the accuracy of that ground-pointing mode; as such, a reward function which diminishes smoothly as the spacecraft attitude varies away from the ground-pointing reference is selected as

$$R_s = \frac{1}{1 + |\sigma_{\text{err}}|} \text{ if } a = \text{Science} \quad (9)$$

Pointing is accomplished through the use of three reaction wheels with randomized initial biases; attitude determination is accomplished using a truth-plus-noise simulation of an ideal attitude estimator. In addition to managing the science pointing mode, the spacecraft operations agent must also ensure that the system remains power-positive by pointing the spacecraft's body-fixed solar panel towards the sun. The spacecraft's power consumption is modeled using a simple net power process:

$$\dot{J} = W_{\text{in}} - W_{\text{out}} \quad (10)$$

where J is the total energy stored by the spacecraft's battery, W_{in} is the power produced by the solar panel which is assumed to follow a cosine law, and W_{out} is the constant load power drawn by the spacecraft; for the purposes of this work, the load power is assumed to be constant.

To further complicate the operations problem, reaction wheel saturation is also modeled. In LEO, a primary source of disturbance torques for spacecraft occur from interactions with planetary atmospheres. To this end, the spacecraft geometry is considered as a standard "box-and-wing" model with a large offset area representing the solar panel. Left uncorrected, reaction wheel speeds would increase to counteract the aerodynamic torques until they saturate, rendering the spacecraft uncontrollable. To desaturate the wheels, a set of RCS thrusters and a wheel desaturation algorithm are implemented as a third and final flight mode. This mode sets the attitude reference towards the sun, but periodically pulses the thrusters to reduce the wheel momentum. Importantly, this mode is constructed using a pre-existing desaturation algorithm that assumes a small body angular rate when computing thruster firing sequences; when entered before the attitude control system can stabilize the system, this mode produces destabilizing behavior. This type of constraint is representative of one the real-world challenges of incorporating strategies for autonomy around existing flight software stacks and operations procedures.

Table 1 Initial conditions for the real-valued MDP; \mathcal{U} represents a uniform distribution.

| Variable | Value |
|---------------|---|
| r_{eq} | 3396.19 km |
| a | $r_{eq} + 400.0km$ |
| e | $\mathcal{U}(0, 0.5)$ |
| i | $\mathcal{U}(-90^\circ, 90^\circ)$ |
| ω | $\mathcal{U}(0^\circ, 360^\circ)$ |
| Ω | $\mathcal{U}(0^\circ, 360^\circ)$ |
| ν | $\mathcal{U}(0^\circ, 360^\circ)$ |
| σ_{BN} | $\mathcal{U}(\mathbf{0}, \mathbf{1})$ |
| ω_{BN} | $\mathcal{U}(\mathbf{0} \text{ rad/s}, \mathbf{0.1} \text{ rad/s})$ |
| ω_{BN} | $\mathcal{U}(\mathbf{0}, \mathbf{0.1})$ |
| ω_{RW} | $\mathcal{U}(-600 \text{ RPM}, 600 \text{ RPM})$ |
| J_{stored} | $\mathcal{U}(5 \text{ W-Hr}, 10 \text{ W-Hr})$ |
| t_{mode} | 3 minutes |
| T_{max} | 540 modes |

Table 2 Safety MDP labelling parameters

| Observed Variable | Operational Limit | Safety Limit |
|-------------------|-------------------|--------------|
| $ \omega_{BN} $ | 0.05 rad/s | N/A |
| $ \omega_{RW} $ | 1,000 RPM | 1,500 RPM |
| J_{stored} | 5 W-Hr | 0 W-Hr |

$$P = \begin{cases}
 s & = \{\mathbf{r} \in \mathbb{R}^3, \dot{\mathbf{r}} \in \mathbb{R}^3, \boldsymbol{\sigma}_{BN} \in \mathbb{O}^3, \boldsymbol{\omega}_{BN} \in \mathbb{R}^3, \boldsymbol{\omega}_{RW} \in \mathbb{R}^3, \mathbf{J} \in \mathbb{R}^1\} \\
 o & = \{\boldsymbol{\sigma}_{BN} \in \mathbb{O}^3, \boldsymbol{\omega}_{BN} \in \mathbb{R}^3, \boldsymbol{\omega}_{RW} \in \mathbb{R}^3, \mathbf{J} \in \mathbb{R}^1\} \\
 a & = \{\text{Mission, Sun Pointing, Desaturation}\} \\
 T & = \{f_{\text{Mission}}, f_{\text{Sun Pointing}}, f_{\text{Desaturation}}\} \\
 R & = \{R_s, -50 \text{ if } J = 0 \text{ or } |\omega_{RW}| > 250 \frac{\text{rad}}{\text{s}}\}
 \end{cases} \quad (11)$$

The abstract MDP described by Equation 11 represents a command and control problem for a single spacecraft in LEO with hardware constraints and is used as a reference problem. For training, the initial conditions are drawn from uniform random distributions over a range of LEO orbits; similarly, the spacecraft's internal states are randomized to ensure coverage over this space. A summary of the MDP's parameters is shown in Table 1.

Additionally, the parameters of the safety MDP are listed in Table 2.

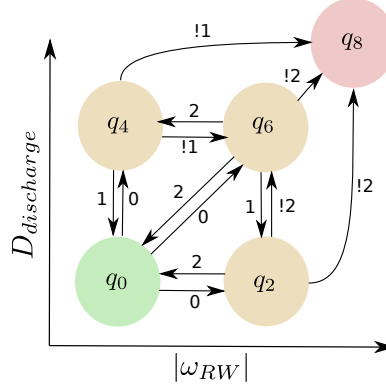


Fig. 7 Safety MDP constructed for the LEO attitude mode planning simulator. $D_{discharge}$ represents the depth of discharge and is inversely analogous to J . Modes relating to “tumble” states with large body rates are omitted for clarity.

B. Shield Construction

To apply the shielded learning technique to space mission operations, a simplified version of the mission POMDP is first constructed using a-priori knowledge. Here, “alert” states are defined using the operational limits found in Table 2. These limits are applied to transform the continuous-time, continuous-state system described by Equation 11 into a simplified, discrete MDP in the observed variables, represented graphically in Fig. 7. This MDP is stated as P_{disc} :

$$P = \begin{cases} s & = \{\omega_{BN} \in \{\text{nominal, high}\}, |\omega_{RW}| \in \{\text{nominal, alert, failure}\}, J \in \{\text{nominal, low, failure}\}\} \\ o & = \{q \in \{q_0, q_1, \dots, q_7, q_8\}\} \\ a & = \{\text{Mission, Sun Pointing, Desaturation}\} \\ T & = \{f_{\text{Mission}}, f_{\text{Sun Pointing}}, f_{\text{Desaturation}}\} \\ R & = \{\emptyset\} \end{cases} \quad (12)$$

While substantially smaller than the continuous state POMDP, the safety MDP encodes important information; for example, desaturation events are only feasible when the spacecraft is not in a tumbling state, and tumbling states themselves do not lead to failure unless the battery charge or wheel speed are already near the failure criteria. In addition, the various state combinations that lead to failure are lumped into q_8 for brevity; this permits the use of the simple LTL specification

$$\varphi = G(\neg\text{“fail”}) \quad (13)$$

which is represented using the Büchi automaton shown in Fig. 8, and can be understood in English as “globally never allow the state to reach the failure state.”

To solve this safety game, the game itself was implemented as a stochastic Markov game (smg) within the PRISM-

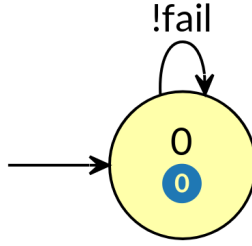


Fig. 8 The one-state Büchi automaton representing the safety specification for the system.

games solver. In this case, PRISM-games solves the safety game using Value Iteration [25]. PRISM-games then saves the shield strategy as a .adv file, which encodes the state-action strategy which maximizes the probability of remaining safe. For this work, the resulting strategy is memoryless and state-based, making it especially amicable to on-line implementation.

To use this adversary strategy, the `stable-baselines`[22] implementation of PPO2 was extended to conform to the post-posed shield framework shown in Fig. 5.

C. Training Results

To provide a comparison between the shielded and unshielded approaches to DRL-based spacecraft autonomy, three agents of each type were trained on the reference problem with separate, random seeds with identical network parameters, hyperparameters, and training durations. The resulting training curves are shown in Figure 9. Notably, convergence behavior is broadly similar between each initialization within each agent category, which indicates that the spacecraft problem is well-posed and does not suffer from the same stochastic convergence that other common DRL environments produce. Clearly, the shielded agents produce substantially better mean rewards at virtually every point in the training process, with the final shielded agents achieving more than twice the mean reward of the unshielded agents. This performance is the result of two benefits of shielding: first, the shielded agents do not spend as much time exploring regions of the state/action space related to failure, as the shield activations keep the agent away from these regions; second, the “safety” aspect of the shield prevents the agent from receiving a reward penalty associated with failure. These results show that the addition of shielding to learning processes for typical spacecraft decision problems to which safety is a core attribute can dramatically improve performance even during training.

D. Performance Results

To verify that the agents are indeed performing in a safe manner, a “simulator” consisting of the agent in a closed-loop interaction with the environment was set up and run multiple times for the best-performing shielded and unshielded agents. The resulting phase-plot diagrams of the agent’s behavior in the observed battery and wheel speed are demonstrated in Figure 10. The shielded learning agent is able to immediately recover after breaching the battery

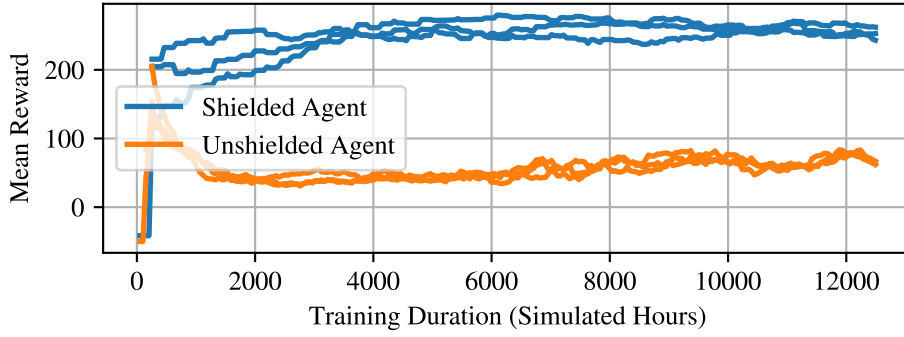
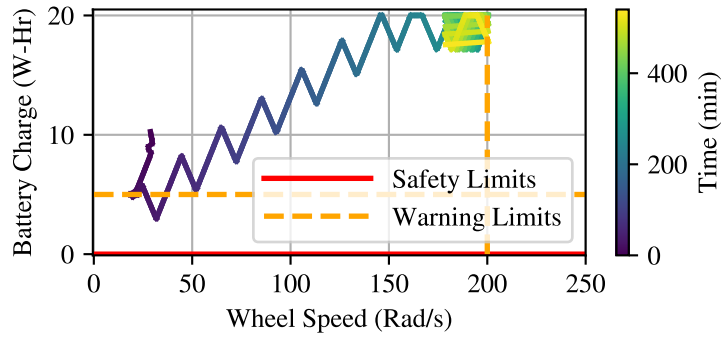
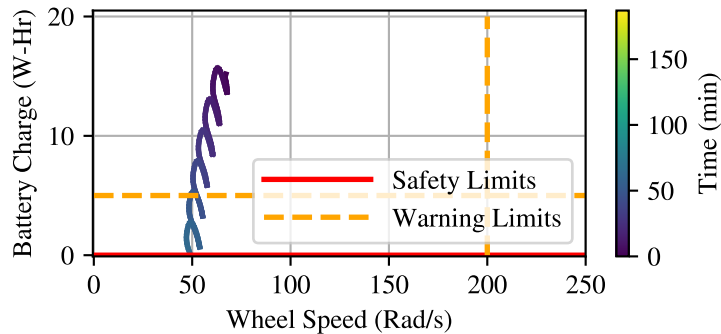


Fig. 9 Comparison of achieved mean reward during training versus quantity of training time.

charge warning limit, and remains bounded by the wheel speed limit while converging to a limit cycle in the upper-right of the nominal section of the phase space. On the other hand, the unshielded agent allows itself to run out of power relatively quickly over the simulation period and does not recover, indicating that it has converged to a local minima in the training space.



(a) Phase plot of the system observations for a run of the shielded agent.



(b) Unshielded Agent; note that the agent fails by depleting the spacecraft's battery.

Fig. 10 Observation phase plots for the shielded and unshielded agents.

IV. Conclusion

A methodology for considering spacecraft command and control problems as sequential decision problems suitable for the application of modern machine learning tools has been presented and extended using the Basilisk astrodynamics framework. In addition, the technique of reactive synthesis and shielded reinforcement learning has been reviewed and applied to a detailed reference spacecraft command and control problem. In comparison to naive approaches to reinforcement learning, the shielded learning approach produces sequential decision agents that both operate safely under prescribed limits and achieves quantitatively better performance versus the unshielded learning agent.

References

- [1] Harris, A., and Schaub, H., "Towards Reinforcement Learning Techniques for Spacecraft Autonomy," *42nd Annual AAS Guidance, Navigation and Control Conference*, , No. AAS 18-078, 2018, pp. 1–10.
- [2] Harris, A., Teil, T., and Schaub, H., "Spacecraft Decision-Making Autonomy Using Deep Reinforcement Learning," *29th AAS/AIAA Space Flight Mechanics Meeting, Hawaii*, , No. AAS 19-447, 2019, pp. 1–19.
- [3] Kubitschek, D. G., "Impactor Spacecraft Encounter Sequence Design for the Deep Impact Mission," *Jet Propulsion*, 2005, pp. 1–14. URL <https://smartech.gatech.edu/handle/1853/8031>.
- [4] Foster, C., Mason, J., Vittaldev, V., Leung, L., Beukelaers, V., Stepan, L., and Zimmerman, R., "Constellation Phasing with Differential Drag on Planet Labs Satellites," *Journal of Spacecraft and Rockets*, Vol. 55, No. 2, 2018. doi:10.2514/1.A33927, URL <https://arc.aiaa.org/doi/pdf/10.2514/1.A33927>.
- [5] Chien, S., Sherwood, R., Tran, D., Castano, R., Cichy, B., Davies, A., Rabideau, G., Tang, N., Burl, M., Mandl, D., Frye, S., Hengemihle, J., Agostino, J. D., Bote, R., Trout, B., Shulman, S., Ungar, S., Gaasbeck, J. V., Boyer, D., Systems, C., Griffin, M., and Mit, H.-h. B., "Autonomous Science on the EO-1 Mission," *Proceedings of International Symposium on Artificial Intelligence, Robotics and Automation in Space (i-SAIRAS)*, , No. May, 2003.
- [6] Chien, S., and Jonsson, A., "Automated Planning & Scheduling for Space Mission Operations JPL," , No. February, 2005.
- [7] Chien, S. A., Tran, D., Rabideau, G., Schaffer, S. R., Mandl, D., and Frye, S., "Timeline-Based Space Operations Scheduling with External Constraints," *Proceedings of the 20th International Conference on Automated Planning and Scheduling (ICAPS)*, , No. Icaps, 2010, pp. 34–41.
- [8] Choo, T. H., and Skura, J. P., "SciBox: A software library for rapid development of science operation simulation, planning, and command tools," *Johns Hopkins APL Technical Digest (Applied Physics Laboratory)*, Vol. 25, No. 2, 2004, pp. 154–161.
- [9] Mnih, V., Silver, D., and Riedmiller, M., "Atari Deep Reinforcement learning," *Nips*, 2013, pp. 1–9. doi:10.1038/nature14236.
- [10] Sutton, R. S., and Barto, A. G., "Reinforcement learning," *Learning*, Vol. 3, No. 9, 2012, p. 322. doi:10.1109/MED.2013.6608833, URL <https://books.google.com/books?id=CAFR6IBF4xYC{&}pgis=1{&}5Cnhttp://incompleteideas.net/sutton/book/the-book.html{&}5Cnhttps://www.dropbox.com/s/f4tnuhipchpkgoj/book2012.pdf>.

- [11] Julian, K. D., and Kochenderfer, M. J., “Autonomous Distributed Wildfire Surveillance using Deep Reinforcement Learning,” , No. January, 2018, pp. 1–16. doi:10.2514/6.2018-1589.
- [12] Linares, R., and Furfaro, R., “Dynamic Sensor Tasking for Space Situational Awareness via Reinforcement Learning,” *Proceedings of the Advanced Maui Optical and Space Surveillance Technologies Conference, Maui Economic Development Board*, 2016, pp. 1–10. URL <https://www.amostech.com/TechnicalPapers/2016/SSA-Algorithms/Linares.pdf>.
- [13] Linares, R., and Furfaro, R., “An Autonomous Sensor Tasking Approach for Large Scale Space Object Cataloging,” *Advanced Maui Optical and Space Surveillance Technologies Conference (AMOS)*, 2017, pp. 1–17. URL www.amostech.com.
- [14] OpenAI, “OpenAI Five,” <https://blog.openai.com/openai-five/>, 2018.
- [15] Cianciolo, A. D., Maddock, R. W., Prince, J. L., Bowes, A., Powell, R. W., White, J. P., Tolson, R., Shaughnessy, O., and Carrelli, D., “Autonomous Aerobraking Development Software : Phase 2 Summary,” 2018, pp. 1–16.
- [16] Gaudet, B., Furfaro, R., Process, M. D., Learning, R., Regulator, L. Q., Arizona, T., and Arizona, T., “Robust Spacecraft Hovering Near Small Bodies in,” *test*, , No. August, 2012, pp. 1–20. doi:10.2514/6.2012-5072.
- [17] Roberto Furfaro, I. B., “Deep Learning for Autonomous Lunar Landing,” *Proceedings of the 2018 AAS/AIAA Astrodynamics Specialist Conference, Snowbird UT*, 2018.
- [18] Dulac-Arnold, G., Mankowitz, D., and Hester, T., “Challenges of Real-World Reinforcement Learning,” 2019. URL <http://arxiv.org/abs/1904.12901>.
- [19] Cassandra, A. R., “A Survey of POMDP Applications,” *Uncertainty in Artificial Intelligence*, 1997, pp. 472–480.
- [20] Sample, E., Ahmed, N., and Campbell, M., “An Experimental Evaluation of Bayesian Soft Human Sensor Fusion in Robotic Systems,” , No. August, 2012, pp. 1–19. doi:10.2514/6.2012-4542.
- [21] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O., “Proximal Policy Optimization Algorithms,” *Arxiv*, 2017, pp. 1–12. doi:10.1016/j.jdeveco.2016.04.001, URL <http://arxiv.org/abs/1707.06347>.
- [22] Hill, A., Raffin, A., Ernestus, M., Gleave, A., Traore, R., Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., Schulman, J., Sidor, S., and Wu, Y., “Stable Baselines,” <https://github.com/hill-a/stable-baselines>, 2018.
- [23] Alshiekh, M., Bloem, R., Ehlers, R., Könighofer, B., Niekum, S., and Topcu, U., “Safe Reinforcement Learning via Shielding,” *ArXiv*, 2017, pp. 1–23. URL <http://arxiv.org/abs/1708.08611>.
- [24] Alcorn, J., Schaub, H., Piggott, S., and Kubitschek, D., “Simulating Attitude Actuation Options Using the Basilisk Astrodynamics Software Architecture,” *67 th International Astronautical Congress*, 2016.
- [25] Chen, T., Forejt, V., Kwiatkowska, M., Parker, D., and Simaitis, A., “PRISM-games: A model checker for stochastic multi-player games,” *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Vol. 7795 LNCS, 2013, pp. 185–191. doi:10.1007/978-3-642-36742-7_13.