



Towards fully autonomous orbit management for low-earth orbit satellites based on neuro-evolutionary algorithms and deep reinforcement learning

Alexander Kyuroson^{*}, Avijit Banerjee, Nektarios Aristeidis Tafanidis, Sumeet Satpute, George Nikolakopoulos

Robotics and Artificial Intelligence Group, Department of Computer, Electrical and Space Engineering, Luleå University of Technology, 971 87 Luleå, Sweden

ARTICLE INFO

Recommended by T. Parisini

Keywords:

Deep reinforcement learning
Satellite constellation
Orbit management
Robotics

ABSTRACT

The recent advances in space technology are focusing on fully autonomous, real-time, long-term orbit management and mission planning for large-scale satellite constellations in Low-Earth Orbit (LEO). Thus, a pioneering approach for autonomous orbital station-keeping has been introduced using a model-free Deep Policy Gradient-based Reinforcement Learning (DPGRL) strategy explicitly tailored for LEO. Addressing the critical need for more efficient and self-regulating orbit management in LEO satellite constellations, this work explores the potential synergy between Deep Reinforcement Learning (DRL) and Neuro-Evolution of Augmenting Topology (NEAT) to optimize station-keeping strategies with the primary goal to empower satellite to autonomously maintain their orbit in the presence of external perturbations within an allowable tolerance margin, thereby significantly reducing operational costs while maintaining precise and consistent station-keeping throughout their life cycle. The study specifically tailors DPGRL algorithms for LEO satellites, considering low-thrust constraints for maneuvers and integrating dense reward schemes and domain-based reward shaping techniques. By showcasing the adaptability and scalability of the combined NEAT and DRL framework in diverse operational scenarios, this approach holds immense promise for revolutionizing autonomous orbit management, paving the way for more efficient and adaptable satellite operations while incorporating the physical constraints of satellite, such as thruster limitations.

1. Introduction

Recent advancements in satellite miniaturization and their utilization in various fields such as communication and Earth observation have created the need for long-term orbit management and mission planning for LEO satellite constellations (Banerjee et al., 2023). To address these demands, Machine Learning (ML) (Li et al., 2020b) and Deep Learning (DL) (Li et al., 2020a) have gained exponential attention within the aerospace community for their applications in autonomous and real-time Guidance, Navigation and Control (GNC) systems (Izzo et al., 2019) for future deployment over the past decades. Such pioneering direction represents a departure from conventional methods, offering real-time onboard capabilities that pave the way for more autonomous and resilient decision-making with long-duration missions (Li et al., 2019). Notably, the application of ML and DL in GNC systems is expected to allow for the execution of diverse periodic operations such as maneuver planning for station-keeping, ensuring the satellite formations stability in the presence of orbital perturbations, thereby enabling precise control to actively maintain the

desired relative position of satellite within the constellation. Such advancements enable more fuel-efficient maneuvers with more effective and independent control, marking a significant leap forward in space exploration.

After the final orbital insertion to relocate the satellite from its parking orbit to its desired orbit, continuous perturbations (Viswanathan et al., 2022) necessitate periodic corrections to ensure the satellite remains within an acceptable tolerance range of its intended trajectory. To address this, defining a trajectory tracking problem (Sankaranarayanan et al., 2023) involves using the nominal satellite trajectory, unaffected by perturbations, as the reference path. The goal is to determine the necessary maneuver plan or the control input to keep the satellite aligned with this reference trajectory. However, such station-keeping maneuvers become extremely complex given that the satellite is part of a constellation and its relative state, i.e., position and velocity, must be maintained with respect to other satellites within the formation (Smith et al., 2021). Furthermore, other factors, such as fuel

^{*} Correspondence to: Luleå University of Technology, 971 87 Luleå, Sweden.

E-mail addresses: akyuroson@gmail.com (A. Kyuroson), aviban@ltu.se (A. Banerjee), nektaf@ltu.se (N.A. Tafanidis), sumsat@ltu.se (S. Satpute), geonik@ltu.se (G. Nikolakopoulos).

<https://doi.org/10.1016/j.ejcon.2024.101052>

Received 14 May 2024; Accepted 10 June 2024

Available online 15 June 2024

0947-3580/© 2024 The Author(s). Published by Elsevier Ltd on behalf of European Control Association. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

consumption and the duration of orbital correction, must be considered as they will impact operational capabilities.

DRL has been shown as a promising framework to not only autonomously control the movements of a satellite while maintaining it within its designated orbit but also reduce operational costs and maximize mission returns (Harris et al., 2019). It must be noted that DRL algorithms can continuously learn and optimize control strategies by receiving feedback from the environment such as the position of satellite and its orbital dynamics, thereby making optimal state-dependent decisions accordingly (Cai et al., 2022). DRL algorithms can be designed to consider various parameters such as gravitational forces, orbital perturbations, and other dynamic factors affecting the position of satellite. Such state-dependent learning processes allow the satellite to adapt to dynamic and unstable environmental conditions and execute maneuvers that optimize station-keeping, reducing the need for constant manual interventions.

Furthermore, due to the size of trained policies, which are based on Multi-Layered Perceptron (MLP) network architectures and therefore are computationally efficient, DRL agents can be deployed on platforms with limited computing power, to achieve autonomous capabilities such as collision avoidance (Zhang et al., 2016), path-planning, trajectory optimization (Sullivan & Bosanac, 2020) as well as station-keeping (Miller & Linares, 2019). Expanding on this concept, a complex model-free data-driven policy based on Proximal Policy Optimization (PPO) was proposed for propulsion-less maneuvers that leverage differential drag states to solve multi-satellite constellation problems (Smith et al., 2021). Moreover, similar DRL algorithm was used to perform station-keeping maneuvers for a satellite operating near a Sun-Earth L_2 Southern quasi-halo trajectory in an ephemeris model while utilizing Bayesian optimization for guided parameter selection to achieve a policy with high accuracy in a deterministic fashion (Bonasera et al., 2023).

Given the advantages mentioned above of DRL for use in orbit management, this work investigates the feasibility of utilization of DRL in combination with NEAT (Peng et al., 2018; Stanley & Miikkulainen, 2002) for station-keeping, which can lead to more adaptive and efficient control strategies, while enhancing the ability of satellite to autonomously manage its position within a tolerance range of its orbital parameters, thereby reducing operational costs and ensuring precise and consistent station-keeping throughout the mission. Therefore, the main contributions of this study are as follows: (a) A novel satellite environment for real-time simulation of orbital station-keeping that includes non-linear orbital dynamics and perturbations. (b) Employing NEAT for topology and parametric optimization as well as the creation of an expert agent for comparison with optimal policy achieved by DRL in the proposed environment under mission-specific constraints. (c) Utilization of model-free DPGR algorithms in conjunction with LEO satellites with low-thrust constraints for maneuvers to the desired orbit, while incorporating dense reward scheme and reward shaping based on domain knowledge. (d) A preliminary comparative analysis of Deep Deterministic Policy Gradient (DDPG) (Lillicrap et al., 2019), Twin Delayed Deep Deterministic Policy Gradient (TD3) (Fujimoto et al., 2018) and NEAT to further improve learning while fine-tuning the system constraints and address catastrophic inference.

The remainder of this article is structured as follows. Section 2 outlines the related theories to implement the proposed framework. Thereafter, Section 3 presents a detailed description of the implemented DRL and NEAT framework and discusses the agent and environmental design for solving LEO satellite station-keeping. Furthermore, a detailed evaluation of achieved simulation results is presented in Section 4. Finally, we conclude this article by discussing the achieved results and future work in Section 5.

2. Background

In this Section, the orbital dynamics, encompassing both the equations of motion and the perturbations affecting the satellite, are presented. Moreover, a detailed description of the DRL algorithms that are employed for orbital station-keeping within the satellite orbital environment is provided, followed by an overview of evolutionary algorithms, specifically NEAT, for its use for Reinforcement Learning (RL).

2.1. Orbital dynamics

The orbital environment model for the satellite, used in both training as well as testing of the RL agent, is described in the Earth Centered Inertial (ECI) (J200) frame of reference and is as follows (Hu et al., 2023; Vallado, 2001):

$$\ddot{\mathbf{r}}(t) = -\frac{\mu}{\|\mathbf{r}\|^3} \mathbf{r} + \mathbf{a}_t + \mathbf{a}_{3b} + \mathbf{a}_{drag} + \mathbf{a}_g + \mathbf{a}_{srp}, \quad (1)$$

where $\mathbf{a}_t = R_{\mathbb{H}}^{\mathbb{I}} \mathbf{U}$, and the satellite position $\mathbf{r}(t) = [r_x, r_y, r_z]^T$ and velocity $\mathbf{v}(t) = [v_x, v_y, v_z]^T$ are represented in ECI reference frame defined as $\mathbb{I} = \{X, Y, Z\}$ as depicted in Fig. 1. Additionally, μ denotes the Earth's gravitational parameter, and $\|\mathbf{r}\|$ indicates the distance of the satellite from the center of the Earth. The satellite is orbiting around the Earth, primarily under the influence of the gravitational field. However, its motion is affected under the influence of various perturbation factors such as non-spherical gravity \mathbf{a}_g due to Earth's oblateness, atmospheric drag \mathbf{a}_{drag} , influences of the third body \mathbf{a}_{3b} such as the Sun and Moon as well as the solar radiation pressure \mathbf{a}_{srp} . In the context of station-keeping for the LEO region, the non-spherical Earth's gravitational field and the air drag are the primary influential factors. To account for the effects of disturbances while maintaining the station-keeping objectives, controlled actuation, \mathbf{a}_t , is considered provided by the onboard low-thrust electric propulsion system. The thrust force is expressed as $\mathbf{U} = [F_R, F_T, F_N]^T$ in the Local-Vertical Local-Horizontal (LVLH) reference frame $\mathbb{H} = \{R, T, N\}$, where $R_{\mathbb{H}}^{\mathbb{I}}$ denotes the corresponding transformation matrix from \mathbb{H} to \mathbb{I} .

External perturbations continuously deviate the satellite from its desired orbit. A brief description of the perturbation models used to construct the orbital environment is described in the following part.

Earth's Gravitational Perturbation: The gravitational field can be described as a potential, whose gradient provides acceleration due to gravity. The higher degree spherical harmonics gravity accelerations are given as (Vallado, 2001):

$$\mathbf{a}_g = \nabla \frac{\mu}{r} \sum_{n=2}^{\infty} \sum_{m=0}^n H_1 H_2, \quad (2)$$

where $H_1 = P_{nm}(\sin \varphi) (C_{nm} \cos m\lambda + S_{nm} \sin m\lambda)$, $H_2 = \left(\frac{R_E}{r}\right)^n$, and n and m are the degree and order of spherical harmonics, respectively. R_E denotes the Earth's radius, r is the geocentric distance of the satellite, and P_{nm} is the Legendre polynomial with argument $\sin \varphi$. Moreover, φ and λ denote the respective latitude and longitude of the satellite and C_{nm} and S_{nm} are the harmonic coefficients of the potential. The Earth's gravity potential can be modeled accurately, using the higher degree and order Earth's Gravity Model (EGM). The EGM96 gravity model provides the data for spherical harmonic coefficients complete to the degree and order 360 (NASA & NIMA, 2004).

Atmospheric Drag: Given that the satellite is operating in LEO, the acceleration drag experienced is calculated according to Eq. (3), where $\mathbf{v}_{rel} = \mathbf{v} - (0, 0, \Omega_{Earth})^T \times \mathbf{r}$ is the relative velocity of the satellite and earth, and Ω_{Earth} , is the Earth's rotation. Furthermore, C_d is the drag coefficient of the satellite, while the atmospheric density at the altitude of the spacecraft is denoted by ρ , the total mass of the spacecraft is indicated by M , and the cross-sectional area in nominal attitude is given by A_d .

$$\mathbf{a}_{drag} = -0.5\rho \left(\frac{A_d C_d}{M}\right) \|\mathbf{v}_{rel}\|^2 \left(\frac{\mathbf{v}_{rel}}{\|\mathbf{v}_{rel}\|}\right). \quad (3)$$

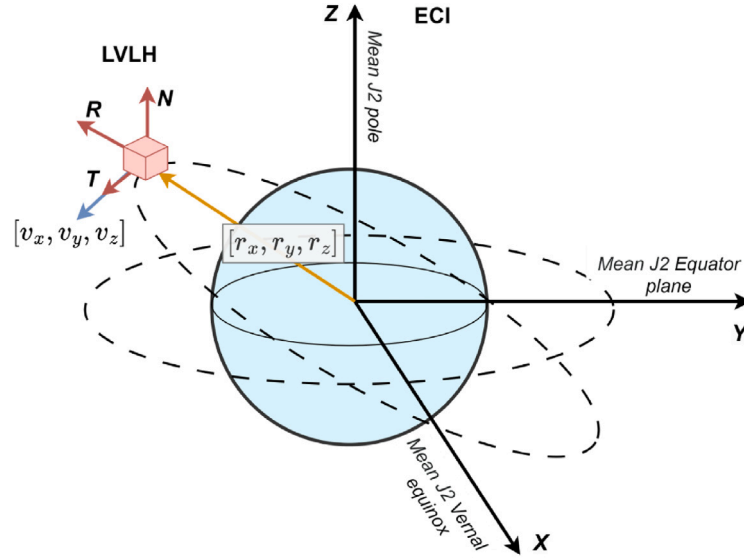


Fig. 1. The ECI and LVLH reference frames.

Third Body Perturbation: The third-body gravity perturbations are exerted by the Sun and Moon. Thus, the gravity acceleration due to the third body effects can be expressed as (Vallado, 2001):

$$\mathbf{a}_{3b} = \mu_{\odot} \left(\frac{\mathbf{r}_{\odot s}}{r_{\odot s}^3} - \frac{\mathbf{r}_{\odot \oplus}}{r_{\odot \oplus}^3} \right) + \mu_{\ominus} \left(\frac{\mathbf{r}_{\ominus s}}{r_{\ominus s}^3} - \frac{\mathbf{r}_{\ominus \oplus}}{r_{\ominus \oplus}^3} \right) \quad (4)$$

where μ_{\odot} and μ_{\ominus} corresponds to the gravitational parameter of Sun and Moon, respectively. Moreover, $\mathbf{r}_{\oplus s}$ and $\mathbf{r}_{\odot s}$ are the respective position vectors from the satellite to the Sun and Moon, and $\mathbf{r}_{\odot \oplus}$ and $\mathbf{r}_{\ominus \oplus}$ denote the respective position vector from the Earth to the Sun and Moon. The Eq. (4) can be expressed as $\mathbf{a}_{3b} = \mathbf{a}_{\text{moon}} + \mathbf{a}_{\text{sun}}$, where the required position vectors of the celestial bodies are calculate via SPICE Toolkit.

Solar Radiation Pressure: For the Solar Radiation Pressure (SRP) induced acceleration, the cannonball disturbance model and dual conical shadow model (Gill & Montenbruck, 2013) are considered. The SRP induced acceleration is given by:

$$\mathbf{a}_{srp} = \left(\frac{S C_r A_s}{c \cdot M} \right) \frac{-\mathbf{r}_{ss}}{\|\mathbf{r}_{ss}\|^3}, \quad (5)$$

where $S \in \mathbb{R}$ is the solar flux at 1 Au, which is calculated from the conical shadow model, $C_r \in \mathbb{R}$ denotes the surface reflectivity coefficient, A_s is the exposed area, c denotes the speed of light, and M is the total mass of the respective satellite.

2.2. Deep reinforcement learning

DRL can be described as a learning paradigm, where an agent attempts to learn an optimal policy, π , by interacting with its environment, \mathcal{E} ; see Fig. 2. Based on received observations, x_t , at each time-step, t , the agent takes an action, $a_t \in \mathbb{R}^n$, which will result in a scalar reward feedback, r_t . The resulting policy maps each state to a probability distribution for a given action, $\pi : S \rightarrow \mathcal{P}(\mathcal{A})$. By utilizing the Markov Decision Process (MDP), the DRL can be modeled based on state space, S , action space \mathcal{A} , initial state distribution $\mathcal{P}(s_1)$, transition dynamics $\mathcal{P}(s_{t+1}|s_t, a_t)$, a reward function $\mathcal{R}(s_t, a_t)$ and a discount factor $\gamma \in [0, 1]$. To obtain the optimal policy by maximizing the expected return from the initial distribution, $\mathcal{J} = \mathbb{E}_{r_t, s_t \sim \mathcal{E}, a_t \sim \pi} [\mathcal{R}_1]$ is the main goal in DRL framework.

2.2.1. Deep Deterministic Policy Gradient (DDPG)

DDPG can be described as an off-policy RL algorithm that has adopted an actor-critic architecture as shown in Fig. 3. Moreover, the actor and critic networks are denoted as μ and Q , respectively, while

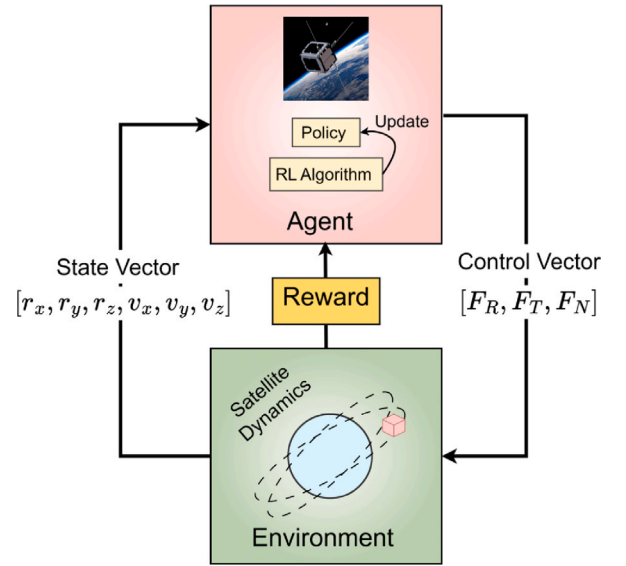


Fig. 2. The layout of the proposed DRL framework.

their corresponding target networks are represented by μ' and Q' . By combining value-based and policy-based methods, it can concurrently learn a Q -function and a policy based on the sampled batches from its experience pool (Lillicrap et al., 2019). By directly taking specific actions based on the current state, the actor-network learns a deterministic policy. Moreover, experience replay is employed to stabilize training, and target networks are used to compute more stable Q -value estimates. Critic Network updates its parameters, θ^Q , based on the Temporal Difference (TD) method, where the loss function for the critic network is given by:

$$L(\theta^Q) = \frac{1}{N} \sum_i (Q(s_i, a_i | \theta^Q) - y_i)^2, \quad (6)$$

where $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{\mu'})) | \theta^Q$ is the target value for the critic network and γ the discount factor.

2.2.2. Twin delayed DDPG

To address the overestimation of Q -value, which can be observed in actor-critic methods such as DDPG, TD3 adopts an improved clipping

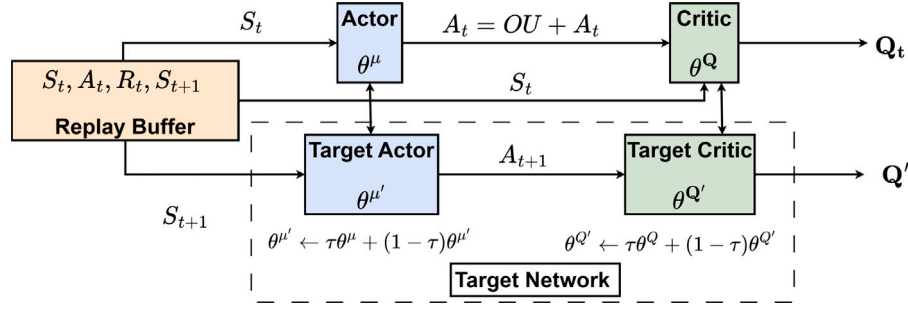


Fig. 3. DDPG algorithm architecture flowchart.

variant of double Q -learning (Van Hasselt et al., 2016) by utilizing two critic networks that have access to the same pool of sampled experiences (Fujimoto et al., 2018). To update the critic networks, the minimum Q -value is selected based on the following expression:

$$L = \left(r + \gamma \min_{i=1,2} Q^i(s_{t+1}, a_{t+1} | \theta_i^{Q'}) - Q(s_{t+1}, a_{t+1} | \theta_i^Q) \right)^2, \quad (7)$$

where θ_1^Q and θ_2^Q represent the corresponding parameters of critic networks and $\theta_1^{Q'}$ and $\theta_2^{Q'}$ are the parameters of the target networks. Moreover, $r(s_t, a_t)$ is the reward feedback based on the deterministic action $a_t = \mu(s_t | \theta^\mu)$. To minimize the loss function, Eq. (7) is used to update both θ_1^Q and θ_2^Q . Furthermore, by adding exploration noise, ϵ , to the action, the value function can be updated more smoothly while reducing the Q -function error exploitation. It must be noted that based on the domain-specific problem, ϵ can be either represented by uncorrelated Gaussian distribution (Fujimoto et al., 2018), $\mathcal{N}(0, \sigma_a)$, where σ_a represents the standard deviation observed in n steps, or temporally correlated noise such as Ornstein-Uhlenbeck (OU) process (Lillicrap et al., 2019). The regularization of a_{t+1} by adding clipped noise, ϵ , can facilitate the smoothing of the target policy and increase the exploration ability of the algorithm. The TD3 employs the delayed target network update method by using State Action Reward State Action (SARSA)-style regularization technique to prevent accumulations of error and to ensure small TD error during the training. The soft update rule for the target network is given by $\theta_i^{Q'} \leftarrow \tau \cdot \theta_i^Q + (1-\tau)\theta_i^{Q'}$, where θ_i^Q , $\theta_i^{Q'}$ and τ represent the parameters of the main and the target network and the updating rate, respectively.

2.3. Neuro-evolutionary topology optimization

NEAT can be described as an evolutionary algorithm, where both the weights and topologies of Neural Networks (NNs) are evolved such that for any given control problem (Stanley & Miikkulainen, 2004), it can discover an optimal policy without any reliance on indirect inference based on the value function. Furthermore, it has been shown that NEAT can efficiently survey policy space for the environments where the gradients are complex and difficult to calculate (Peng et al., 2018). Combined with its capability to mitigate the impact of initialization as well as initial topology configuration of NNs on learning performance, NEAT has been considered to be vital for the optimization of domain-specific MLPs in DRL framework (Whiteson et al., 2005). However, due to its initialization approach, where extremely rudimentary NNs without any hidden layers are utilized to evolve into highly sophisticated networks gradually, it has been deemed sample inefficient (Peng et al., 2018). Therefore, other variations of NEAT such as Hypercube-based Neuro-Evolution of Augmenting Topology (HyperNEAT) (Risi & Stanley, 2011) and adaptive HyperNEAT (Risi & Stanley, 2010) were proposed to address some of the shortcomings of NEAT.

The topology optimization via NEAT is achieved by incrementally augmenting the initial population of generated NNs via mutation operators in addition to adding nodes and links after each generation assessment (Stanley & Miikkulainen, 2002). Thus, discovering the most optimal topology with the minimum required number of nodes and their related weights for a given environment. Furthermore, NEAT utilizes speciation by measuring the number of excess, E , and disjoint, D , genes between a given pair of genomes to calculate their compatibility distance, δ , thus avoiding any topological override between different populations for any given species (Peng et al., 2018). The compatibility distance can be expressed as $\delta = c_1 \frac{E}{N} + c_2 \frac{D}{N} + c_3 \overline{W}$, where c_1 , c_2 , and c_3 are the associated weights to adjust the importance of each factor while \overline{W} and N represent the average weight differences of matching genes and the number of genes, respectively. During the reproduction phase, NEAT leverages explicit fitness sharing, thereby preventing the takeover of the entire population by any single species to allow possible future evolution and crossovers (Stanley & Miikkulainen, 2004). The adjusted fitness, f'_i , for given organism, i , is given by:

$$f'_i = \frac{f_i}{\sum_{j=1}^n sh(\delta(i, j))}, \quad (8)$$

where $sh(\delta(i, j)) \in \{0, 1\}$ and represents the sharing function that depending on the compatibility threshold, δ_i , reduce $\sum_{j=1}^n sh(\delta(i, j))$ to the number of organisms in a given species. It must be noted that based on f'_i , the number of offsprings are adjusted such that the lowest performing species are eliminated before generating the next generation (Risi & Stanley, 2011).

3. Methodology

In this Section, the proposed orbital environment and its related constraints are presented. Moreover, a brief description of DRL algorithms as well as NEAT and their related hyperparameters are provided. The reward function design for addressing the station-keeping problem is given in detail, and its subsequent components are further analyzed.

3.1. Station-keeping in DRL framework

Fig. 2 illustrates the proposed architecture of the satellite environment and its utilization in conjunction with DRL framework for autonomous orbital management, specifically optimal station-keeping. To determine the optimal maneuvers required for driving the satellite back to its nominal reference trajectory based on its current state, $\mathbf{x}(t)$, the desired nominal trajectory without any perturbations is calculated and utilized as the target state in DRL algorithm to determine the optimal control input, \mathbf{U} , for maintaining the satellite on its reference trajectory. Furthermore, the reward feedback, \mathcal{R} , is proposed based on the convergence towards the nominal state and allows the optimization of the policy, π , during the training phase by interacting with the simulated environment that encompasses the satellite dynamics.

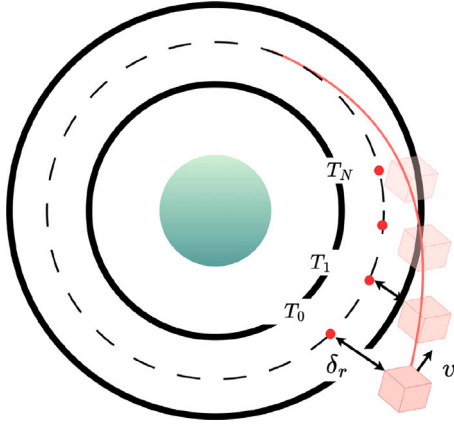


Fig. 4. Satellite trajectory tracking, where δ_r is the error between the reference and actual position at time-step T_n .

3.1.1. Experiment setup

The station-keeping problem can be described as the optimal trajectory tracking problem and is given by:

$$\min_U V(s(t)) = \int_{\tau=t}^{t_f} \|(s(\tau) - s_d(\tau))\|_2^Q + \|U(\tau)\|_2^R d\tau, \quad (9)$$

where t_f denotes a specific final time, while $s_d(\tau)$ and $s(\tau)$ represent the desired and actual states of the spacecraft at any given time instance τ , respectively. Furthermore, Eq. (9) is subjected to $\dot{s}(t) = As(t) + \frac{1}{m}BU + D$, where $U(t) \in [\mathbf{F}_{min}, \mathbf{F}_{max}]$ represents the thrust magnitude with its minimum and maximum thresholds. In the proposed environment, the low-thrust constraints are applied such that the maximum thrust magnitude is 18 mN. The time-step T_n is set to 10 s with the thrust duration also fixed at 10 s. Moreover, the state initialization is randomized such that the initial velocity and position of the satellite have been affected by the perturbations without any prior interference with position error, $\delta_r \leq 45$ km, as illustrated in Fig. 4.

3.1.2. Action & observation space

The state observation for a given satellite is determined based on satellite dynamics and the perturbations present in its environment, as depicted in Fig. 2. The observation space, $\mathbf{x}(t)$, contains both the velocity vector, $\mathbf{v}(t)$, as well as the position vector, $\mathbf{r}(t)$, of the satellite in a given time-step and is obtained from the orbital propagator. Both velocity and position vectors are included in observation space for the DRL algorithm due to the station-keeping problem formulation, where both the position and velocity of the satellite must be tightly tracked to maintain a nominal orbit.

Furthermore, the action vector, \mathbf{a}_t , is the control input, which is expressed as $\mathbf{a}_t = [\mathbf{F}_R, \mathbf{F}_T, \mathbf{F}_N]$ in LVLH reference frame and corresponds to the exerted forces from the satellite thrusters. The normalized thrust vector elements with their corresponding range $\in [-1, 1]$ are denoted by \mathbf{F}_R , \mathbf{F}_T and \mathbf{F}_N . By considering six onboard thrusters, which are aligned with all three axes of LVLH frame of reference, the actuation command for the combined thrusters located on the same axis can be described as acceleration and deceleration along a given axis for a given state.

3.1.3. Reward function design

The reward function and its design is a crucial element of DRL framework as it directly impacts the learning outcomes of a given agent (Hu et al., 2023). Therefore, a well-formulated reward function is required to enable the agent to efficiently explore the environment for optimal policy convergence. To comply with the standard reward formulation based on policy gradient frameworks of RL and MDP, a policy given as π_θ can be found such that the agent must optimize the parameter θ by maximizing the expected accumulated reward

Table 1

Hyperparameters for the DRL algorithms.

Parameters	DDPG values	TD3 values
Learning rate critic	0.00025	0.0001
Learning rate actor	0.000025	0.00001
Discount factor	0.99	0.99
Update interval	2	2
Target network update rate	0.001	0.001
Batch size	64	64
OU noise (σ)	0.15	0.15
OU noise (ζ)	0.0	0.0
OU noise (ϕ)	0.2	0.2

$J(\theta) = \mathbb{E}_{s \sim \rho^\pi, a \sim \pi_\theta} [r(s, a)]$. According to the policy gradient theorem, the gradient of $J(\theta)$ with respect to θ can be expressed as following:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim \rho^\pi, a \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^\pi(s, a)], \quad (10)$$

where Q^π represents the state-action value function. To realize fully autonomous station-keeping, state-specific thresholds for both position and velocity are defined such that extreme divergence from optimal orbit will result in the termination of the simulation with an extremely high penalty. Furthermore, the reward function, \mathcal{R}_{tot} , has been decomposed into multiple segments to achieve fuel-optimal maneuvers to reach the nominal trajectory within the least feasible steps. Therefore, a dense reward scheme is proposed to facilitate the agent in its exploration endeavor while minimizing the samples required to reach an optimal policy. The total reward function is expressed as a combination of linear and exponential functions and is formulated as follows:

$$\mathcal{R}_{tot}(s, a) = \mathcal{R}_{state} + \mathcal{R}_{position} + \mathcal{R}_{bound} + \mathcal{R}_{action} - \mathcal{R}_{step}, \quad (11)$$

where \mathcal{R}_{state} represents the unified state reward based on $\{s_t, s_{t+1}, s_{target}\}$ and is given by:

$$\mathcal{R}_{state} = -\ln(\|s_{t+1} - s_{target}\|) + \ln(\|s_t - s_{target}\|). \quad (12)$$

Moreover, $\mathcal{R}_{position}$ is a compounded reward to assess the approximation to the nominal position by exponentially rewarding the agent trajectory for given actions. This will result in closer proximity to the desired orbit and can be expressed as:

$$\mathcal{R}_{position} = e^{\left(20 - \frac{\|\delta_r\|}{\psi}\right)}, \quad (13)$$

where $\psi = 10^3$ is a scaling factor and $\delta_r = \bar{\mathbf{r}}_{t+1} - \bar{\mathbf{r}}_{target}$ represents the error between the reference and actual position as shown in Fig. 4. Eq. (13) provides an auxiliary reward when the agent crosses the maximum tolerable trajectory threshold, $\delta_r \leq 20$ km. Furthermore, to minimize the consumption of propellant, the agent is rewarded based on the magnitude of control input, which is given by:

$$\mathcal{R}_{action} = \beta \cdot (1 - \|\bar{\mathbf{a}}\|), \quad (14)$$

where $\beta = 10^2$ is an empirical weight chosen such that the fuel consumption is scaled with respect to $\mathcal{R}_{position}$, thereby resulting in less fuel utilization when the satellite is within the desired boundary threshold, $\delta_r \leq 1$ km. To ensure that the exploration is performed in an optimal and sample-efficient fashion, the \mathcal{R}_{bound} is utilized to incentivize any exploration where the satellite exceeds its initial divergence outside of selected boundaries. Therefore, \mathcal{R}_{bound} is solely activated when the achieved error between the nominal and actual position in a new state, $\bar{\mathbf{r}}_{t+1} \geq 35$ km, and is calculated as follows:

$$\mathcal{R}_{bound} = -\left(10 \cdot \psi + e^{\left(\frac{\|\delta_r\|}{\psi} - 40\right)}\right). \quad (15)$$

To minimize the number of steps while achieving the desired orbit, \mathcal{R}_{step} is used to encourage the agent to converge to the solution by minimizing its effort and can be expressed by $\mathcal{R}_{step} = \gamma \cdot \kappa$, where $\gamma = 5 \cdot 10^{-4}$ is an empirically selected weight and κ represents the current step.

3.1.4. Terminal constraints

The terminal constraints play a significant role in the DRL training process, as it prevents sample-inefficient training. Therefore, the terminal position error is defined as follows:

$$T_{constraints} = \begin{cases} True & \|\bar{r}_{t+1} - \bar{r}_{target}\| \geq 50 \text{ km} \\ False & Otherwise. \end{cases} \quad (16)$$

Given that the terminal criterion is reached, the training episode is halted, and the agent is penalized and receives a -10^4 penalty. Furthermore, in any successful episode, where $\|\delta_r\| \leq 1 \text{ km}$, the agent is rewarded with an additional 10^4 points.

3.1.5. Exploration noise

The exploration of the environment by the agent is facilitated by incorporating noise into the control inputs generated by the actor-network as depicted in Fig. 3. This strategy increases the efficacy of DRL based control by effectively expanding the agent exploration into previously unseen state–action space during the training process (Lillicrap et al., 2019). Furthermore, the addition of noise creates inherent robustness in the generated model for real-world scenarios, given the existence of noise in the acquired data during live missions.

The OU process, which models the velocity of Brownian particles in the presence of friction, is utilized during the training phase for both DDPG and TD3 as temporally correlated noise. The OU process is formulated as follows:

$$\delta a = \phi(\zeta - a) + \sigma \mathcal{N}(0, 1), \quad (17)$$

where ϕ , ζ , and σ represent the decay rate, long-term mean, and variation of the generated noise, respectively (Lillicrap et al., 2019). The associated values for each parameter of OU process for both DRL algorithms are provided in Table 1.

3.1.6. Hyperparameters

Both DDPG and TD3 have been considered as DRL algorithm candidates to evaluate the proposed satellite environment and obtain an optimal policy for the station-keeping problem. The main driving factor for selecting these algorithms is their wide utilization in application with both observation and action space in the continuous domain Lillicrap et al. (2019). The topology for actor and critic networks for both algorithms was selected after experimenting with networks of various sizes. A similar network topology is used throughout all training and evaluation phases to maintain neutrality and accurately assess the capabilities of both algorithms. Specifically, the actor and critic networks are constructed with two deep, fully connected layers comprising 400 and 300 neurons, respectively, activated by the ReLU function. Furthermore, the output layer utilizes the hyperbolic tangent function, tanh, as its activation mechanism, enabling the clipping of the thrust vector, \mathbf{a} , within its predetermined bounds. Details regarding the remaining hyperparameters for both algorithms are provided in Table 1.

3.2. Station-keeping with NEAT

Similar to the proposed DRL framework depicted in Fig. 2, NEAT algorithm is used to leverage the proposed satellite environment for autonomous orbit management, specifically addressing the station-keeping problem with defined constraints as stated briefly in the previous section.

3.2.1. Experiment setup

A similar environment and experimental setting as DRL framework is deployed, where the main goal is evolution and optimization of the feature network, thereby determining the optimal topology for the defined problem to study the correlation between state and action space.

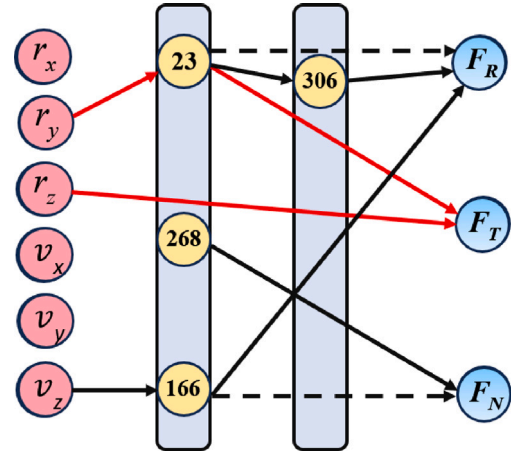


Fig. 5. NNs topology from NEAT for the best-achieved fitness with two hidden layers, where yellow nodes indicate clusters of neurons while the color of edges, red and black, indicate a negative or positive correlation, respectively.

3.2.2. Hyperparameters

NEAT has been employed to obtain an optimal policy that satisfies similar constraints and criteria as DRL algorithms. Furthermore, the initial number of hidden layers was set to 2 so that similar MLP as DDPG and TD3 actor-network is created to further investigate the selected NNs architecture and possibly discover other optimal topology for the station-keeping problem for low-thrust satellites. It must be noted the termination condition was selected such that given the maximum population, the average fitness of the species must be $\geq 10^4$. The remaining hyperparameters to configure NEAT are provided in Table 2.

4. Results

In this section, the results from both DRL algorithms as well as NEAT for station-keeping are presented. Furthermore, the component-wise residual for the state vector is provided to analyze the behavior of policy generated by both methods. It must be noted that both the environmental framework and the agents have been implemented in the Python 3.8 environment, while Pytorch is used for the implementation of DRL algorithms.

4.1. Evaluation of reward

The average episodic reward for both DDPG and TD3 during the training process is shown in Fig. 8. As Fig. 8 illustrates, although both DRL algorithms are capable of solving the station-keeping task, DDPG has the most unreliable performance due to catastrophic inference, which is one of the primary challenges of DRL algorithms given non-stationary data stream. To combat these issues, several strategies

Table 2
Hyperparameters for NEAT algorithm.

Parameters	Values
Population size (p)	50
Max generations (n)	100
Add node rate (m_n)	0.02
Add connection rate (m_c)	0.5
Weight init mean	0.0
Weight init std	1.0
Weight mutate rate (m_w)	0.3
Weight mutate power	0.81
Activation functions	ReLU & tanh
Aggregation functions	mean, max & sum
Crossover rate (c)	0.1

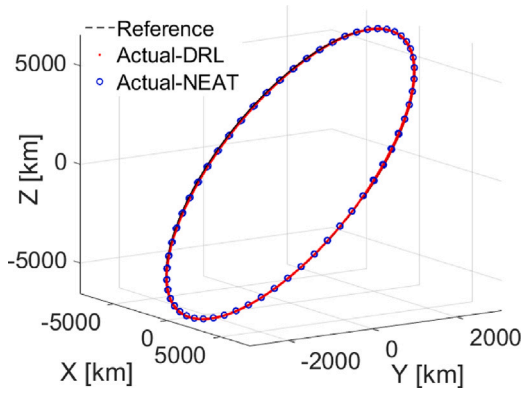


Fig. 6. Correction of satellite trajectory utilizing DRL and NEAT.

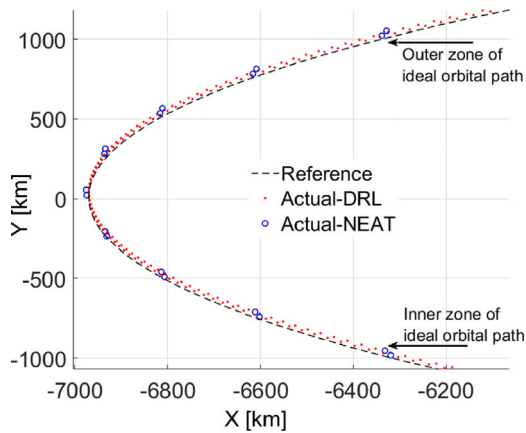


Fig. 7. Close-up view of a segment of satellite trajectory based on DRL and NEAT.

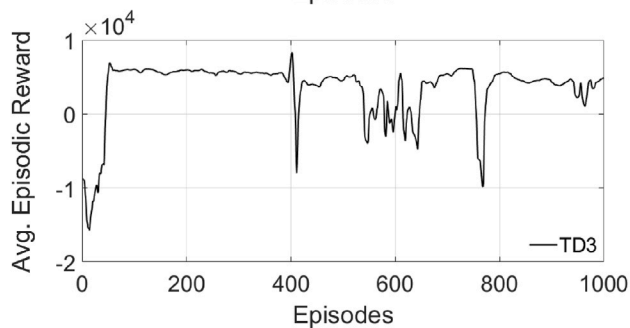
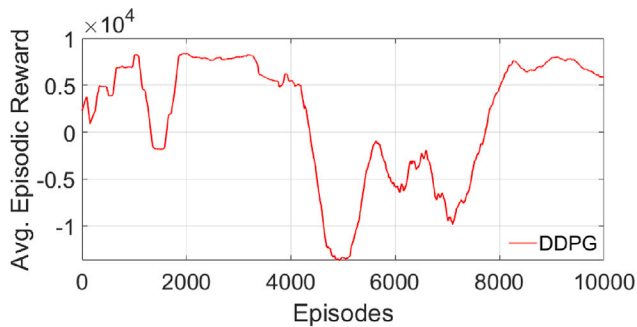


Fig. 8. Average episodic training reward for DDPG and TD3.

were devised for continual learning that align with the main reasons that TD3 was adopted due to its stable learning and faster convergence to the solution.

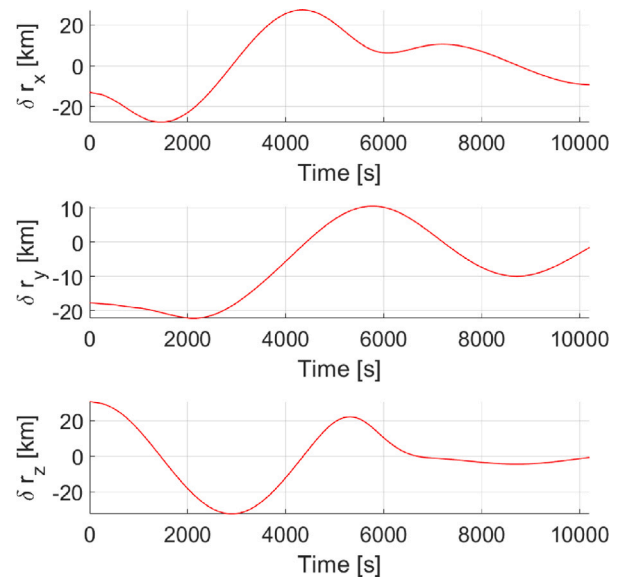


Fig. 9. Component-wise trajectory error for TD3.

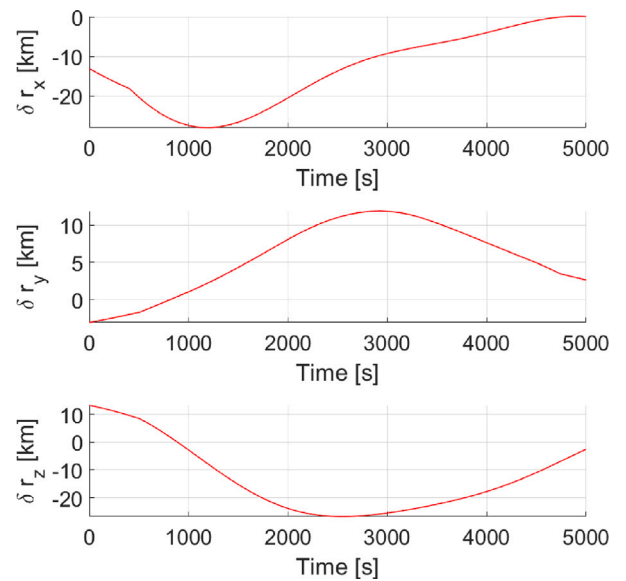


Fig. 10. Component-wise trajectory error for NEAT.

4.2. Evaluation of DRL & NEAT

The optimal policy from TD3 is utilized to compare the trajectory and velocity error between the DRL framework and NEAT. It must be noted that optimized NNs topology devised by NEAT closely resembles the MLP architecture proposed for TD3 (Fujimoto et al., 2018) as shown in Fig. 5. Moreover, it has been shown that in the context of station-keeping, some state elements are not vital for achieving the nominal trajectory and velocity; see Fig. 5.

Figs. 6 and 7 illustrate the nominal trajectory of the satellite without perturbations as well as trajectory correction performed by TD3 and NEAT for station-keeping for a given satellite in the presence of perturbations. It must be noted that the target orbit for this study is considered to be a circular orbit with a radius of 600 km, and an inclination of 70 deg. Furthermore, Figs. 9 and 10 depict the residual between the target and the actual position of satellite for TD3 and NEAT as the correctional maneuver is performed. Similarly, the error between the desired and the actual velocity of the satellite for TD3 is shown

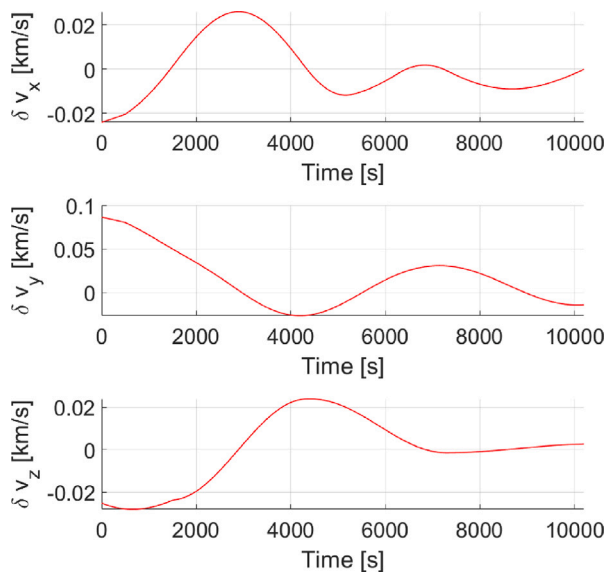


Fig. 11. Component-wise velocity error for TD3.

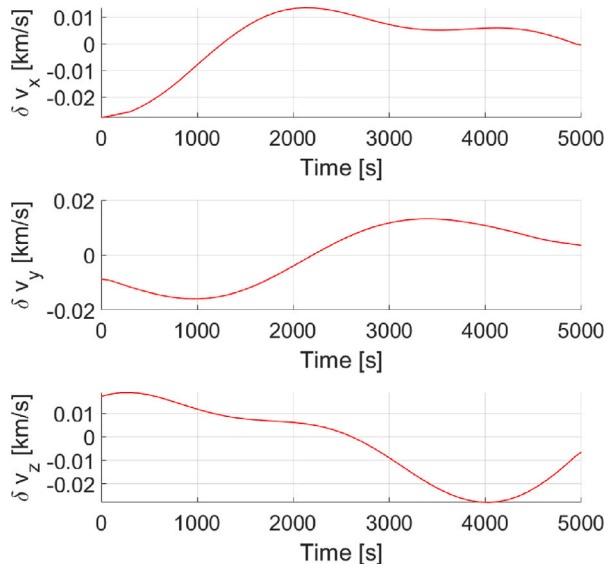


Fig. 12. Component-wise velocity error for NEAT.

in Fig. 11. NEAT algorithm produces similar velocity residuals as illustrated in Fig. 12, which indicates that both algorithms can perform trajectory maneuvers such that the nominal trajectory is reached in the presence of perturbations.

5. Conclusions

This article introduces a novel approach towards achieving a fully autonomous orbital station-keeping for LEO satellites through the amalgamation of NEAT and DRL. The quest for complete autonomy becomes essential in response to the pressing need for an autonomous decentralized orbit management strategy for vast satellite constellations. Towards this, the proposed methodology provides a framework for continuous, precise, and adaptable positioning of satellites without consistent human intervention. By synergizing DRL and NEAT while integrating dense reward schemes, the proposed approach strives to optimize the orbital deviation to empower satellites to independently maintain their orbits despite external perturbations, staying within

acceptable tolerance margins. Furthermore, this approach considers the physical limitations imposed by thruster constraints, thereby paving the way for more efficient and adaptable satellite operations. Thus, it showcases the potential of AI-driven solutions in advancing autonomous control systems for complex real-world applications, departing from conventional methods and paving the way for more efficient and resilient satellite operations. Population-based optimization, prioritized replay buffer, and other variations of NEAT can be investigated to optimize the results further. Moreover, other exploration strategies such as Generalized State-Dependent Exploration (g-SDE) (Raffin et al., 2020) and colored noise (Eberhard et al., 2023) can be leveraged to improve the generalization of achieved policy by introducing smoother state exploration.

CRedit authorship contribution statement

Alexander Kyuroson: Data curation, Formal analysis, Investigation, Methodology, Software, Validation, Visualization, Writing – original draft, Writing – review & editing. **Avijit Banerjee:** Conceptualization, Funding acquisition, Supervision, Writing – original draft, Writing – review & editing. **Nektarios Aristeidis Tafaoidis:** Data curation, Investigation, Visualization, Writing – original draft. **Sumeet Satpute:** Conceptualization, Funding acquisition, Investigation, Supervision. **George Nikolakopoulos:** Conceptualization, Funding acquisition, Supervision, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgment

This work has been partially funded by the European Space Agency (ESA) open Invitations to Tender (ITT) and innovation research grant in OPTACOM project, in collaboration with OHB Sweden under Grant Contract no: OPC-OSE-CC-0536.

References

- Banerjee, A., Mukherjee, M., Satpute, S., & Nikolakopoulos, G. (2023). Resiliency in space autonomy: a review. *Current Robotics Reports*, 4(1), 1–12.
- Bonasera, S., Bosanac, N., Sullivan, C. J., Elliott, L., Ahmed, N., & McMahon, J. W. (2023). Designing Sun–Earth L2 halo orbit stationkeeping maneuvers via reinforcement learning. *Journal of Guidance, Control, and Dynamics*, 46(2), 301–311.
- Cai, Y., Zhang, E., Qi, Y., & Lu, L. (2022). A review of research on the application of deep reinforcement learning in unmanned aerial vehicle resource allocation and trajectory planning. In *2022 4th international conference on machine learning, big data and business intelligence* (pp. 238–241). <http://dx.doi.org/10.1109/MLBDBI58171.2022.00053>.
- Eberhard, O., Hollenstein, J. J., Pinneri, C., & Martius, G. (2023). Pink noise is all you need: Colored noise exploration in deep reinforcement learning. In *International conference on learning representations*.
- Fujimoto, S., Hoof, H., & Meger, D. (2018). Addressing function approximation error in actor-critic methods. In *International conference on machine learning* (pp. 1587–1596). PMLR.
- Gill, E., & Montenbruck, O. (2013). *Satellite orbits: Models, methods and applications*. Springer, ISBN: 978-3-540-67280-7, <http://dx.doi.org/10.1007/978-3-642-58351-3>.
- Harris, A., Teil, T., & Schaub, H. (2019). Spacecraft decision-making autonomy using deep reinforcement learning. In *29th AAS/AIAA space flight mechanics meeting* (pp. 1–19).
- Hu, J., Yang, H., Li, S., & Zhao, Y. (2023). Densely rewarded reinforcement learning for robust low-thrust trajectory optimization. *Advances in Space Research*.
- Izzo, D., Märten, M., & Pan, B. (2019). A survey on artificial intelligence trends in spacecraft guidance dynamics and control. *Astrodynamics*, 3, 287–299.
- Li, H., Baoyin, H., & Topputo, F. (2019). Neural networks in time-optimal low-thrust interplanetary transfers. *IEEE Access*, 7, 156413–156419.
- Li, H., Chen, S., Izzo, D., & Baoyin, H. (2020). Deep networks as approximators of optimal low-thrust and multi-impulse cost in multitarget missions. *Acta Astronautica*, 166, 469–481.

- Li, B., Huang, J., Feng, Y., Wang, F., & Sang, J. (2020). A machine learning-based approach for improved orbit predictions of LEO space debris with sparse tracking data from a single station. *IEEE Transactions on Aerospace and Electronic Systems*, 56(6), 4253–4268.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., & Wierstra, D. (2019). Continuous control with deep reinforcement learning. arXiv: 1509.02971.
- Miller, D., & Linares, R. (2019). Low-thrust optimal control via reinforcement learning. In *29th AAS/AIAA space flight mechanics meeting, vol. 168* (pp. 1817–1834). American Astronautical Society Ka'anapali, Hawaii.
- NASA, & NIMA (2004). EGM96 the NASA GSFC and NIMA joint geopotential model. URL: <https://cddis.nasa.gov/926/egm96/egm96.html>.
- Peng, Y., Chen, G., Singh, H., & Zhang, M. (2018). NEAT for large-scale reinforcement learning through evolutionary feature learning and policy gradient search. In *Proceedings of the genetic and evolutionary computation conference* (pp. 490–497). New York, NY, USA: Association for Computing Machinery, ISBN: 9781450356183, <http://dx.doi.org/10.1145/3205455.3205536>.
- Raffin, A., Kober, J., & Stulp, F. (2020). Smooth exploration for robotic reinforcement learning. In *Conference on robot learning*.
- Risi, S., & Stanley, K. O. (2010). Indirectly encoding neural plasticity as a pattern of local rules. In *International conference on simulation of adaptive behavior* (pp. 533–543). Springer.
- Risi, S., & Stanley, K. O. (2011). Enhancing ES-hyperNEAT to evolve more complex regular neural networks. In *Proceedings of the 13th annual conference on genetic and evolutionary computation* (pp. 1539–1546).
- Sankaranarayanan, V. N., Banerjee, A., Satpute, S., Roy, S., & Nikolakopoulos, G. (2023). Adaptive control for a payload carrying spacecraft with state constraints. *Control Engineering Practice*, 135, Article 105515.
- Smith, B., Abay, R., Abbey, J., Balage, S., Brown, M., & Boyce, R. (2021). Propulsionless planar phasing of multiple satellites using deep reinforcement learning. *Advances in Space Research*, 67(11), 3667–3682.
- Stanley, K. O., & Miikkulainen, R. (2002). Evolving neural networks through augmenting topologies. *Evolutionary Computation*, 10(2), 99–127.
- Stanley, K. O., & Miikkulainen, R. (2004). Competitive coevolution through evolutionary complexification. *Journal of Artificial Intelligence Research*, 21(1), 63–100.
- Sullivan, C. J., & Bosanac, N. (2020). Using reinforcement learning to design a low-thrust approach into a periodic orbit in a multi-body system. In *AIAA scitech 2020 forum* (p. 1914).
- Vallado, D. A. (2001). *Fundamentals of astrodynamics and applications: vol. 12*, Springer Science & Business Media.
- Van Hasselt, H., Guez, A., & Silver, D. (2016). Deep reinforcement learning with double Q-learning. In *Proceedings of the AAAI conference on artificial intelligence, vol. 30* (pp. 2094–2100).
- Viswanathan, V. K., Papadimitriou, A., Banerjee, A., Mansouri, S. S., & Nikolakopoulos, G. (2022). Exogenous disturbance estimation for autonomous navigation around small celestial bodies. In *2022 IEEE 61st conference on decision and control* (pp. 3760–3766). IEEE.
- Whiteson, S., Stone, P., Stanley, K. O., Miikkulainen, R., & Kohl, N. (2005). Automatic feature selection in NeuroEvolution. In *Proceedings of the 7th annual conference on genetic and evolutionary computation* (pp. 1225–1232).
- Zhang, T., Kahn, G., Levine, S., & Abbeel, P. (2016). Learning deep control policies for autonomous aerial vehicles with MPC-guided policy search. In *2016 IEEE international conference on robotics and automation* (pp. 528–535). IEEE.