

# Automated Trajectory Planning: A Cascaded Deep Reinforcement Learning Approach for Low-Thrust Spacecraft Orbit-Raising

## Supplementary Material

### 1 Transfer Scenarios

The transfer scenarios presented in the paper not only differ in terms of the initial and final orbit but also vary due to the initial spacecraft parameters. These spacecraft parameters play a significant role in trajectory optimization problems. To ensure fair comparisons with other results in the literature, we standardized these parameters based on their values. As a result, we considered three Geostationary Transfer Orbit (GTO) initial orbit scenarios (GTO-1, GTO-2, GTO-3) and two Super-Geostationary Transfer Orbit (Super-GTO and Super-GTO-2) scenarios, as outlined in Table 1 in the main text. The details of these spacecraft parameters are presented in Table 1.

Scenarios	Impulse $I_{sp}$	Efficiency $\lambda$	Power P	Mass m	Thrust F	Thrust during shadows
GTO-1 to GEO	1800s	55%	5kW	1200kg	0.31N	No
GTO-2 to GEO	3300s	65%	5kW	450kg	0.20N	No
GTO-3 to GEO	3000s	-	-	2000kg	0.35N	Yes
Super-GTO to GEO	3300s	65%	10kW	1200kg	0.40N	No
Super-GTO-2 to NRHO	1500s	-	-	1000kg	1N	Yes

Table 1: Spacecraft initial parameters across all transfer scenarios

### 2 CDRL Framework

In this section we will provide the additional details and calculations for our cascaded deep reinforcement learning framework.

#### 2.1 State Elements

As introduced in the main paper in Section 4.3, our Cascaded Deep Reinforcement Lerning (CDRL) aprroach utilized the five 'he' elements as a state vector which are shown as follows:

$$\mathbf{s} = [h \ h_X \ h_Y \ e_X \ e_Y]^T \quad (1)$$

where  $h$  denotes the magnitude of angular momentum, while  $h_X$  and  $h_Y$  denote the component of specific angular momentum along the Earth-centered inertial reference frame and the components of the eccentricity vector  $e_x$  and  $e_y$  in a non-inertial reference frame obtained after a 2-1 Euler rotation sequence. The initial values of these elements for all scenarios are stated in Tab 2.

#### 2.2 Convergence Parameters

In Section 4.3, we examine five pivotal terminal conditions at each time step, which are integral spacecraft orbital elements: eccentricity ( $e$ ), semi-major axis ( $a_{sm}$ ), inclination ( $i$ ), right ascension of the ascending node ( $\Omega$ ), and argument of periapsis ( $\omega$ ). These orbital elements are derived from the state ( $he$ ) elements, as defined in Eq. 1. The conversion of these state elements to orbital elements for convergence assessment proves advantageous, particularly in

Scenarios	Position	State Parameters				
		<b>h (km<sup>2</sup>/s)</b>	<b>h<sub>x</sub> (km<sup>2</sup>/s)</b>	<b>h<sub>y</sub> (km<sup>2</sup>/s)</b>	<b>e<sub>x</sub></b>	<b>e<sub>y</sub></b>
GTO-1 - GEO	Initial	67288.41965	0	-32107.258	0.7306	0
	Target	29640.2292	0	0	0	0
GTO-2 - GEO	Initial	67246.21689	0	-30529.143	0.7310	0
	Target	29640.2292	0	0	0	0
GTO-3 - GEO	Initial	68196.3519	0	-8311.044	0.725	0
	Target	29640.2292	0	0	0	0
Super-GTO - GEO	Initial	70532.88179	0	-26991.765	0.8705	0
	Target	29640.2292	0	0	0	0
Super-GTO-2 - NRHO	Initial	101055.564	8993.124	-44988.209	0.6342	0.1422
	Target	384073.943	43240.356	-111563.294	-0.0209	-0.1249

Table 2: State parameters (he) defined for all transfer scenarios

27 GEO transfer scenarios where only three orbital elements (eccentricity ( $e$ ),  $a_{sm}$ , and inclination ( $i$ ), as detailed in Table  
 28 2 of the main text) are required for convergence. This transformation from five to three state elements significantly  
 29 mitigates the non-linearity of the problem, enhancing convergence and optimizing results. This transformation is  
 30 also noteworthy in NRHO transfer scenarios, where the convergence of these three orbital parameters substantially  
 31 influences the remaining two orbital parameters (RAAN, argp) as well.

32 The computation for converting  $he$  elements to orbital state elements is delineated as follows:

33 The first condition checks the eccentricity of the spacecraft's orbit. The magnitude of eccentricity is calculated  
 34 from the  $e_x$  and  $e_y$  parameters, which are part of the state vector. Here  $e_{tol}$  denotes the tolerance value for eccentricity.

$$e_{tar} \leq \left[ e = \sqrt{e_x^2 + e_y^2} \right] \leq e_{tar} + e_{tol} \quad (2)$$

35 The second condition checks the semi-major axis of the spacecraft's orbit, as shown in Eq. (??). The semi-major  
 36 axis ( $a_{sm}$ ) is calculated using the  $h$ ,  $e_x$ , and  $e_y$  parameters from the state vector, as well as the gravitational parameter  
 37  $\mu$ . The value of  $a_{sm}^{tar}$  is the desired target value for the semi-major axis, and  $a_{sm}^{tol}$  is the tolerance value for the semi-major  
 38 axis.

$$a_{sm}^{tar} - a_{sm}^{tol} \leq \left[ a_{sm} = \frac{h^2}{\mu(1 - \sqrt{(e_x^2 + e_y^2)})} \right] \leq a_{sm}^{tar} + a_{sm}^{tol} \quad (3)$$

39 The third condition checks the inclination angle of the spacecraft's orbit, as shown in Eq. (??), where  $i_{tol}$  denotes  
 40 the tolerance value for the inclination angle.

$$i_{tar} \leq \left[ i = \sqrt{\frac{h_x^2 + h_y^2}{h}} \right] \leq i_{tar} + i_{tol} \quad (4)$$

41 The fourth condition assesses the right ascension of the ascending node (RAAN), denoted as ( $\Omega$ ). The calculation  
 42 involves first determining the vertical component  $h_z$ , which is then utilized to find the value of  $\Omega$ . The computation  
 43 steps for finding  $\Omega$  and establishing the tolerance ranges,  $\Omega_{tol}$ , are detailed as follows:

$$h_z = \sqrt{h^2 - h_x^2 - h_y^2} \quad n = \text{cross} \left( \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} h_x \\ h_y \\ h_z \end{bmatrix} \right) \quad (5)$$

$$\Omega = \begin{cases} \arccos \left( \frac{n(1)}{\|n\|} \right), & \text{if } n(2) \geq 0 \\ 360 - \arccos \left( \frac{n(1)}{\|n\|} \right), & \text{otherwise} \end{cases} \quad (6)$$

$$\Omega_{tar} - \Omega_{tol} \leq \Omega \leq \Omega_{tar} \quad (7)$$

The fifth condition evaluates the argument of periapsis ( $\omega$ ). Unlike other parameters, this parameter cannot be directly computed from the  $he$  elements. Instead, additional calculations are required to first calculate the rotation matrix and then it determine the value of ( $\omega$ ). The subsequent calculations, along with the corresponding tolerance settings, are outlined as follows:

$$\zeta = \arctan\left(\frac{h_x}{h_z}\right); \quad \eta = -\frac{h_y}{h}; \quad \eta_{\cos} = \sqrt{\frac{h^2 - h_y^2}{h^2}}; \quad (8)$$

$$R_\zeta = \begin{bmatrix} \cos(\zeta) & 0 & -\sin(\zeta) \\ 0 & 1 & 0 \\ \sin(\zeta) & 0 & \cos(\zeta) \end{bmatrix}; \quad R_\eta = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \eta_{\cos} & \eta \\ 0 & -\eta & \eta_{\cos} \end{bmatrix}; \quad (9)$$

$$R_{ECI \text{ to } R} = R_\eta \cdot R_\zeta; \quad \mathbf{e}_{ECI} = \mathbf{R}_{ECI \text{ to } R}^T \begin{bmatrix} e_x \\ e_y \\ 0 \end{bmatrix}; \quad (10)$$

$$\omega = \begin{cases} \arccos\left(\frac{\mathbf{n} \cdot \mathbf{e}_{ECI}}{\|\mathbf{n}\| \cdot \|\mathbf{e}_{ECI}\|}\right), & \text{if } e_{ECI}(3) \geq 0 \\ 360 - \arccos\left(\frac{\mathbf{n} \cdot \mathbf{e}_{ECI}}{\|\mathbf{n}\| \cdot \|\mathbf{e}_{ECI}\|}\right), & \text{otherwise} \end{cases} \quad (11)$$

$$\omega_{tar} \leq \omega \leq \omega_{tar} + \omega_{tol} \quad (12)$$

The tolerance ranges and target values for all of these parameters in Eq.(2,3,4,7,12) are presented in Table 2 in the main text.

## 2.3 Reward Weights

We introduced the gradient-aided reward function in our approach, elaborated in Section 4.3 of the main text. The reward function incorporates user-defined weights, which are presented in Table 3. These weights are represented by  $W = [w_1; w_2; w_3]$ , where  $W[:, 1]$  corresponds to the weightage of  $a_{sm}$ ,  $W[:, 2]$  pertains to the weightage of  $e$ , and  $W[:, 3]$  reflects the weightage of  $i$  in Eq. 16 (main text). These weight assignments are aligned with the first, second, and third columns of Table 3, respectively. We also fixed the value of  $\tau$  in Eq. 16 (main text) as 0.5. The rewards weights are shown as follows:

	GTO to GEO			Super-GTO to GEO			Super-GTO-2 to NRHO		
	$W[:, 1]$	$\mathbf{W}[:, 2]$	$\mathbf{W}[:, 3]$	$W[:, 1]$	$W[:, 2]$	$\mathbf{W}[:, 3]$	$W[:, 1]$	$W[:, 2]$	$\mathbf{W}[:, 3]$
$\mathbf{w}_1$	1e <sup>3</sup>	2e <sup>3</sup>	3e <sup>2</sup>	1e <sup>3</sup>	4e <sup>3</sup>	3e <sup>2</sup>	8e <sup>5</sup>	8e <sup>5</sup>	8e <sup>5</sup>
$\mathbf{w}_2$	1e <sup>-2</sup>	1.9e <sup>-7</sup>	3e <sup>-5</sup>	1e <sup>-2</sup>	1.9e <sup>-9</sup>	3e <sup>-5</sup>	3e <sup>1</sup>	3e <sup>1</sup>	3e <sup>1</sup>
$\mathbf{w}_3$	5e <sup>2</sup>	7e <sup>2</sup>	3e <sup>2</sup>	5e <sup>2</sup>	2e <sup>3</sup>	3e <sup>2</sup>	1e <sup>2</sup>	1e <sup>2</sup>	1e <sup>2</sup>

Table 3: Weights ( $W = [w_1; w_2; w_3]$ ) used in calculating the reward function.

## 3 Implementation Details

In this section, we present the dynamic model assumptions in Section 3.1, eclipse model assumptions in Section 3.2, and hyperparameter settings in Section 3.3. The information is intended to offer comprehensive details to facilitate result replication by any interested party.

### 3.1 Modeling Assumptions

We assume constant spacecraft thrust, excluding periods in Earth's shadow. Modeling assumptions for the spacecraft's dynamic model include neglecting orbital perturbations and radiation damage. Specifically, we focus on onboard thrust as the sole force, ignoring additional perturbations and assuming constant thrust during the Sun-lit trajectory. These simplifications enable a fair comparison with existing sequential and DRL approaches in the literature. Incorporating orbital perturbations or radiation damage can be done by adding corresponding terms to the model.

67 In our modeling approach, we maintain the assumption of constant spacecraft thrust, with exceptions only during  
 68 the spacecraft's passage through the Earth's shadow. Our dynamic model incorporates certain assumptions to stream-  
 69 line the analysis. Firstly, we neglect the influence of orbital perturbations in this paper. While the force terms in  
 70 Eq 8 (main text) can encompass various forces acting on the spacecraft (thrust, J2 perturbation, gravitational forces)  
 71 we specifically consider the force to be solely due to onboard thrust. Secondly, we disregard the impact of radiation  
 72 damage in this paper. In actuality, the spacecraft's solar arrays may experience degradation when traversing the Van  
 73 Allen belts, leading to a reduction in available thrust. However, we simplify our model by assuming constant thrust  
 74 during the Sun-lit portion of the trajectory, overlooking the radiation damage effects. These modeling assumptions  
 75 are made for the purpose of facilitating a fair comparison with the sequential approach [2] and a previously studied  
 76 Deep Reinforcement Learning (DRL) approach [1] found in the literature. It is important to note that if one wishes to  
 77 incorporate orbital perturbations into the problem, additive terms representing those perturbations need to be included  
 78 in Eq 8 (main text). Similarly, for those interested in considering the effect of radiation damage, an artificial neural  
 79 network-based radiation damage prediction can be integrated into the framework.

### 80 **3.2 Eclipse Model Assumptions**

81 As the spacecraft undergoes multiple revolutions around the Earth to reach its final orbit employing all-electric propul-  
 82 sion, it is highly likely to traverse the Earth's shadow. During these shadow passages, the spacecraft has the option to  
 83 utilize onboard batteries to power the thrusters or switch them off and coast. This study assumes coasting during the  
 84 spacecraft's passage through the Earth's shadow in GEO transfer scenarios.

85 To identify the regions where the spacecraft enters the Earth's shadow, a shadow model is required. In this work,  
 86 we employ the cylindrical eclipse model. The cylindrical Earth shadow model assumes that the shadow cast by Earth is  
 87 cylindrical in shape and remains fixed in space without movement. The conditions to determine whether the spacecraft  
 88 is in eclipse are defined as follows:

$$X_I < 0, \quad (13)$$

$$\sqrt{Y_I^2 + Z_I^2} < R_E \quad (14)$$

89 where  $X_I$ ,  $Y_I$ , and  $Z_I$  represent the components of the Cartesian position vector of the spacecraft in the Inertial  
 90 frame, and  $R_E$  is the radius of the Earth. The equations to convert the spacecraft's state vector, as utilized in this work,  
 91 to Cartesian coordinates are discussed in [2].

<b>Learning rate</b>	$3 \exp -4$
<b>Discount factor</b>	0.99
<b>Buffer size</b>	$1 \exp 6$
<b>Time Penalty <math>\tau</math></b>	0.5

Table 4: Hyper-parameters settings used in CDRL Training.

### 92 **3.3 CDRL Parameters settings**

93 We conducted experiments on an Intel(R) Xeon(R) CPU E5-1620 v4 operating at a frequency of 3.5GHz with 8 cores,  
 94 coupled with the NVIDIA GeForce GTX 1080 graphics processing unit (GPU), and 32GB of random access memory  
 95 (RAM) to meet the computational requirements. The implementation of our Cascaded Deep Reinforcement Learning  
 96 (DRL) algorithms was carried out in Python 3.7 using the PyTorch framework, and built up over the Soft Actor-Critic  
 97 implementation by stable-baselines. This hardware/software configuration significantly contributed to the efficient and  
 98 effective development and training of our models. The hyperparameters for our actor, critic, and target critic models  
 99 are presented in Table 5

100 The actor network outputs the mean ( $\mu$ ) and variance ( $\sigma$ ) for each action, with the state values as input. Utilizing  
 101 these mean and variance values, the actor network generates a Gaussian distribution and samples actions from it.  
 102 Additionally, it produces the log probabilities of the sample distribution to calculate the entropy. The other hyper-  
 103 parameters utilized in the CDRL training are presented in table 4

Layers	Actor		Critic		Target Critic	
	Size	Activation Fun.	Size	Activation Fun.	Size	Activation Fun.
<b>Input layer</b>	6	ReLU	8	ReLU	8	ReLU
<b>Hidden 1</b>	256	ReLU	256	ReLU	256	ReLU
<b>Hidden 2</b>	256	ReLU	256	ReLU	256	ReLU
<b>output</b>	2	ReLU	1	Linear	1	Linear

Table 5: Network parameters used in CDRL.

## 4 Additional Training Results

In this section, we present the convergence scores and converging time plots during the training of CDRL networks for two-body scenarios. It's important to note that in the main text, we exclusively showcased these convergence plots for three-body scenarios due to space limitations. As illustrated in Fig. 1 and Fig. 2, it is evident that as the training score increases, the episodic time for the converged episodes decreases. We continued the training of DRL networks until we achieved convergence in episodes along with stable high scores.

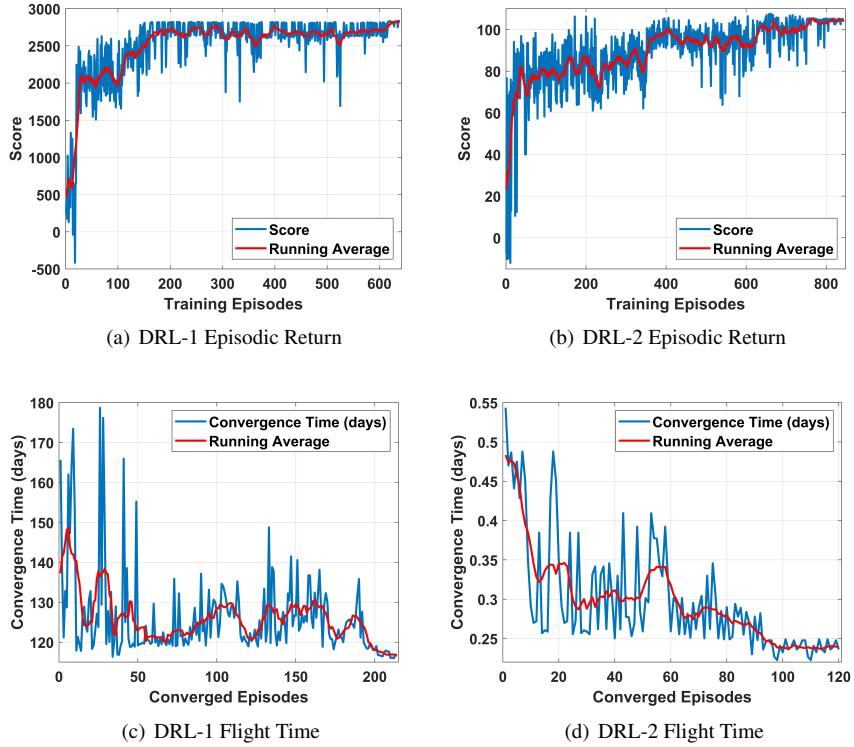


Figure 1: Episodic return and convergence times over training steps for GTO to GEO transfer scenario. DRL-1 and DRL-2 represent the two cascaded DRL agents. Episodic return is for all training episodes, whereas the flight time is for converged episodes during the training.

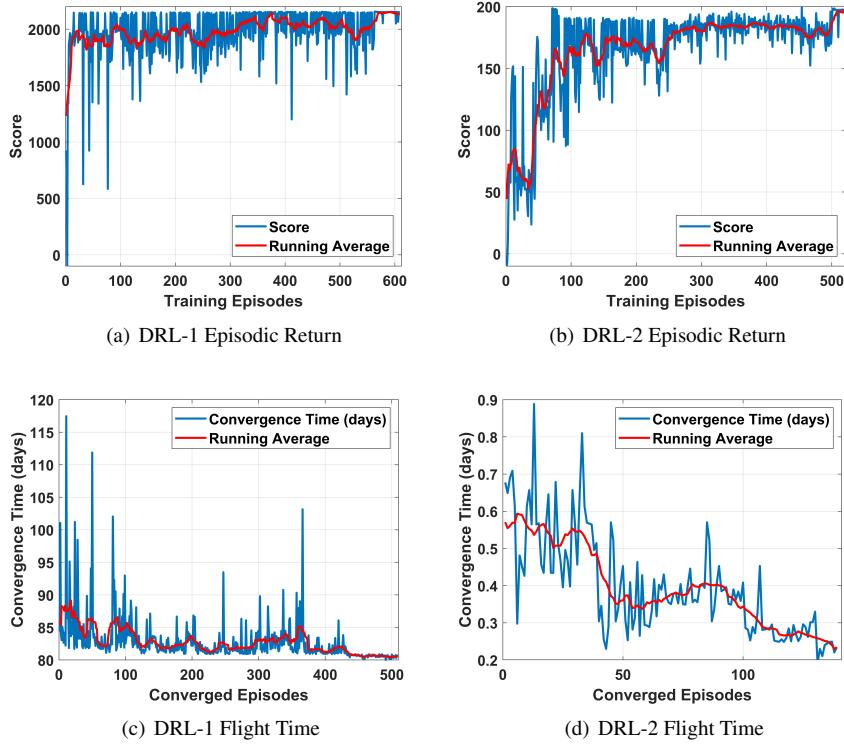


Figure 2: Episodic return and convergence times over training steps for Super-GTO to GEO transfer scenario. DRL-1 and DRL-2 represent the two cascaded DRL agents. Episodic return is for all training episodes, whereas the flight time is for converged episodes during the training.

## References

- [1] Hyeokjoon Kwon, Snyoll Oghim, and Hyochoong Bang. Autonomous guidance for multi-revolution lowthrust orbit transfer via reinforcement learning. *AAS 21*, 315, 2021.
- [2] Suwat Sreesawet and Atri Dutta. Fast and robust computation of low-thrust orbit-raising trajectories. *Journal of Guidance, Control, and Dynamics*, 41(9):1888–1905, 2018.