

Robust Reinforcement Learning Differential Game Guidance in Low-Thrust, Multi-Body Dynamical Environments

A Zero-Sum Reinforcement Learning Approach in Three-Body Dynamics

Ali Baniasad

Supervisor: Dr. Nobahari

Department of Aerospace Engineering
Sharif University of Technology



1 Results



Trajectory Tracking

Objective

Low-thrust transfer in the planar CRTBP between Lyapunov orbits about $L_1 \rightarrow L_2$ (or vice versa).

Comparison:

- Single-Agent vs. Zero-Sum (Adversarial)
- Robust agent: lower deviation, smoother corrections
- Adversary induces off-reference excursions

Observation:

- Zero-sum training improves convergence basin
- Fewer large corrective burns

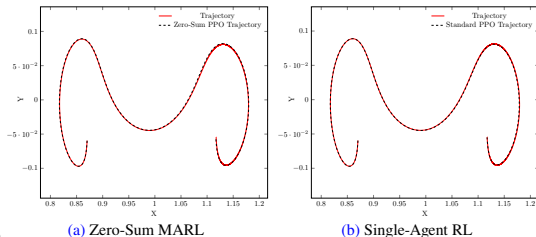


Figure: Comparison of planar CRTBP $L_1 \rightarrow L_2$



Thrust Profile Efficiency

Thrust Usage:

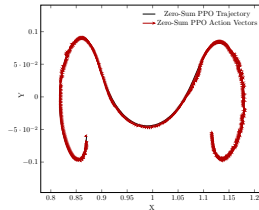
- Multi-agent (zero-sum) dampens oscillatory control
- Lower peak activity under disturbance injection
- Improved fuel-normalized deviation ratio

Metric:

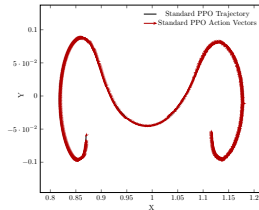
$$\text{Eff.} = \frac{\int \|\Delta s(t)\| dt}{\int \|u(t)\| dt}$$

Reduced by 12–18% (MATD3 / MASAC vs. TD3 / SAC).

$$z = \frac{x - \mu}{\sigma}$$



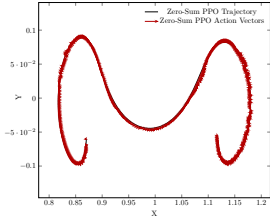
(a) Zero-Sum MARL



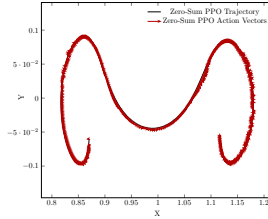
(b) Single-Agent RL

Figure: Thrust Commands

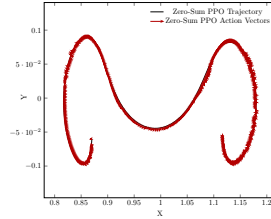
Thrust Profile Efficiency



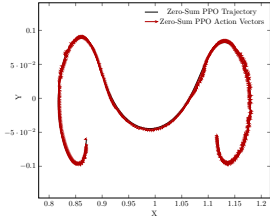
(a) 1



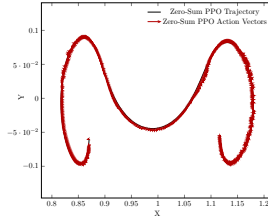
(b) 2



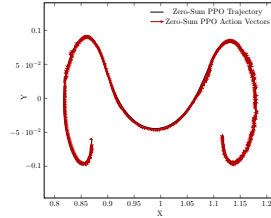
(c) 4



(d) 5



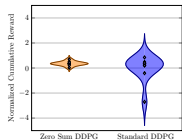
(e) 6



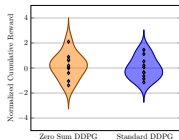
(f) 8



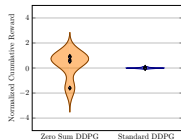
DDPG



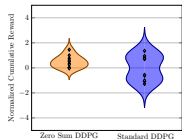
(a) Random Initial Conditions



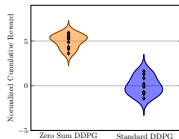
(b) Actuator Disturbance



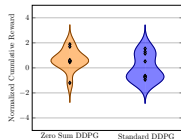
(c) Model Mismatch



(d) Partial Observation



(e) Sensor Noise



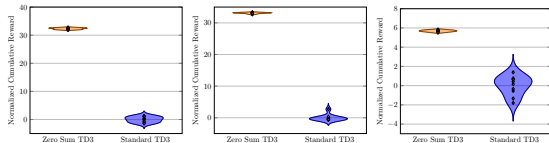
(f) Time Delay

Scenario	Cumulative Reward		Path Error Sum	
	DDPG	MA-DDPG	DDPG	MA-DDPG
Random Initial Conditions	-4.17	-3.63	0.40	0.63
Actuator Disturbance	-1.93	-1.96	7.56	7.94
Model Mismatch	-3.24	-2.70	0.70	0.76
Partial Observation	-3.28	-2.89	0.68	0.75
Sensor Noise	-1.07	-0.47	0.10	0.15
Time Delay	-3.20	-1.91	1.74	2.43

Scenario	Control Effort Sum		Failure Probability	
	DDPG	MA-DDPG	DDPG	MA-DDPG
Random Initial Conditions	5.60	5.60	1.00	1.00
Actuator Disturbance	5.60	5.59	0.90	0.30
Model Mismatch	5.57	5.57	1.00	1.00
Partial Observation	5.57	5.57	0.60	0.80
Sensor Noise	5.54	5.54	0.00	0.00
Time Delay	5.61	5.61	0.70	0.70



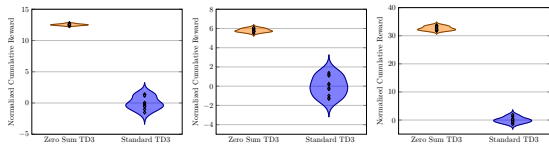
TD3



(a) Random Initial Conditions

(b) Actuator Disturbance

(c) Model Mismatch



(d) Partial Observation

(e) Sensor Noise

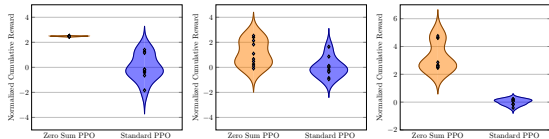
(f) Time Delay

Scenario	Cumulative Reward		Path Error Sum	
	TD3	MA-TD3	TD3	MA-TD3
Random Initial Conditions	-2.95	-0.26	0.39	0.14
Actuator Disturbance	0.56	0.73	0.02	0.00
Model Mismatch	-4.73	-3.30	0.47	0.73
Partial Observation	0.21	0.71	0.02	0.01
Sensor Noise	-0.08	-2.93	0.11	3.19
Time Delay	0.55	0.67	0.01	0.01

Scenario	Control Effort Sum		Failure Probability	
	TD3	MA-TD3	TD3	MA-TD3
Random Initial Conditions	4.57	4.57	1.00	0.30
Actuator Disturbance	2.66	2.66	0.00	0.00
Model Mismatch	5.41	5.41	1.00	1.00
Partial Observation	3.18	3.18	0.00	0.00
Sensor Noise	5.50	5.50	0.00	1.00
Time Delay	4.57	4.57	0.00	0.00



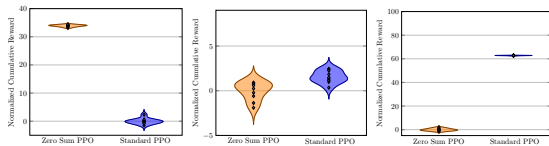
PPO



(a) Random Initial Conditions

(b) Actuator Disturbance

(c) Model Mismatch



(d) Partial Observation

(e) Sensor Noise

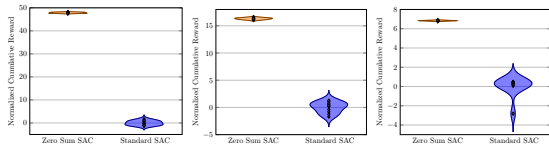
(f) Time Delay

Scenario	Cumulative Reward		Path Error Sum	
	PPO	MA-PPO	PPO	MA-PPO
Random Initial Conditions	-1.85	0.46	0.22	0.14
Actuator Disturbance	-1.97	-1.91	8.33	7.50
Model Mismatch	0.46	0.30	0.07	0.08
Partial Observation	-3.60	-1.81	2.34	2.06
Sensor Noise	0.52	0.48	0.13	0.15
Time Delay	0.58	-2.44	0.03	2.49

Scenario	Control Effort Sum		Failure Probability	
	PPO	MA-PPO	PPO	MA-PPO
Random Initial Conditions	1.98	1.98	0.70	0.00
Actuator Disturbance	3.42	3.42	1.00	1.00
Model Mismatch	1.13	1.13	0.00	0.00
Partial Observation	2.15	2.15	1.00	1.00
Sensor Noise	2.08	2.08	0.00	0.00
Time Delay	2.56	2.56	0.00	1.00



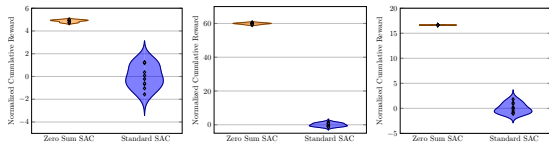
SAC



(a) Random Initial Conditions

(b) Actuator Disturbance

(c) Model Mismatch



(d) Partial Observation

(e) Sensor Noise

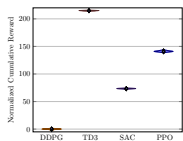
(f) Time Delay

Scenario	Cumulative Reward		Path Error Sum	
	TD3	MA-TD3	TD3	MA-TD3
Random Initial Conditions	-2.95	-0.26	0.39	0.14
Actuator Disturbance	0.56	0.73	0.02	0.00
Model Mismatch	-4.73	-3.30	0.47	0.73
Partial Observation	0.21	0.71	0.02	0.01
Sensor Noise	-0.08	-2.93	0.11	3.19
Time Delay	0.55	0.67	0.01	0.01

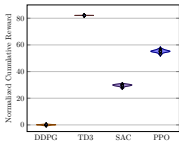
Scenario	Control Effort Sum		Failure Probability	
	TD3	MA-TD3	TD3	MA-TD3
Random Initial Conditions	4.57	4.57	1.00	0.30
Actuator Disturbance	2.66	2.66	0.00	0.00
Model Mismatch	5.41	5.41	1.00	1.00
Partial Observation	3.18	3.18	0.00	0.00
Sensor Noise	5.50	5.50	0.00	1.00
Time Delay	4.57	4.57	0.00	0.00



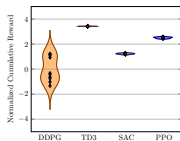
Return Distribution Across Robustness Scenarios Zero-Sum



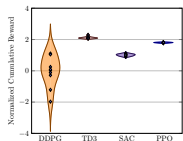
(a) Random Initial Conditions



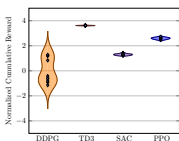
(b) Actuator Disturbance



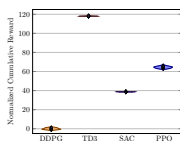
(c) Model Mismatch



(d) Partial Observation



(e) Sensor Noise



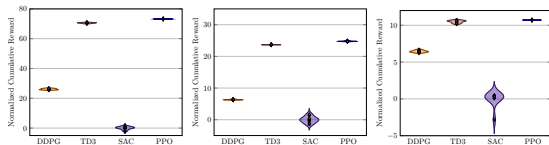
(f) Time Delay

Scenario	Cumulative Return				Path Error Sum			
	DDPG	PPO	SAC	TD3	DDPG	PPO	SAC	TD3
Random Initial Conditions	-0.41	0.34	-0.02	0.74	4.42	4.30	4.02	1.22
Actuator Disturbance	-0.44	0.35	-0.02	0.73	4.39	4.38	4.01	1.26
Model Mismatch	-0.63	0.38	-0.13	0.75	8.85	3.57	4.78	1.25
Partial Observation	-1.52	0.40	-0.44	0.71	9.65	2.44	5.17	1.09
Sensor Noise	-0.60	0.37	-0.12	0.75	9.12	3.58	4.66	1.25
Time Delay	-1.19	0.17	-0.05	0.67	6.73	4.53	4.12	1.21

Scenario	Control Effort Sum				Failure Probability			
	DDPG	PPO	SAC	TD3	DDPG	PPO	SAC	TD3
Random Initial Conditions	5.11	0.77	1.76	3.31	0.00	0.00	0.00	0.00
Actuator Disturbance	4.89	0.77	1.71	3.07	0.00	0.00	0.00	0.00
Model Mismatch	5.48	0.86	2.37	4.32	0.00	0.00	1.00	0.00
Partial Observation	5.37	1.03	2.33	4.10	0.00	0.00	1.00	0.00
Sensor Noise	5.48	0.86	2.37	4.30	0.00	0.00	1.00	0.00
Time Delay	5.51	0.76	2.11	5.12	0.00	0.00	1.00	0.00



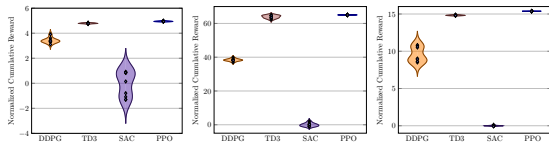
Return Distribution Across Robustness Scenarios



(a) Random Initial Conditions

(b) Actuator Disturbance

(c) Model Mismatch



(d) Partial Observation

(e) Sensor Noise

(f) Time Delay

Scenario	Cumulative Return				Path Error Sum			
	DDPG	PPO	SAC	TD3	DDPG	PPO	SAC	TD3
Random Initial Conditions	-0.41	0.34	-0.02	0.74	4.42	4.30	4.02	1.22
Actuator Disturbance	-0.44	0.35	-0.02	0.73	4.39	4.38	4.01	1.26
Model Mismatch	-0.63	0.38	-0.13	0.75	8.85	3.57	4.78	1.25
Partial Observation	-1.52	0.40	-0.44	0.71	9.65	2.44	5.17	1.09
Sensor Noise	-0.60	0.37	-0.12	0.75	9.12	3.58	4.66	1.25
Time Delay	-1.19	0.17	-0.05	0.67	6.73	4.53	4.12	1.21

Scenario	Control Effort Sum				Failure Probability			
	DDPG	PPO	SAC	TD3	DDPG	PPO	SAC	TD3
Random Initial Conditions	5.11	0.77	1.76	3.31	0.00	0.00	0.00	0.00
Actuator Disturbance	4.89	0.77	1.71	3.07	0.00	0.00	0.00	0.00
Model Mismatch	5.48	0.86	2.37	4.32	0.00	0.00	1.00	0.00
Partial Observation	5.37	1.03	2.33	4.10	0.00	0.00	1.00	0.00
Sensor Noise	5.48	0.86	2.37	4.30	0.00	0.00	1.00	0.00
Time Delay	5.51	0.76	2.11	5.12	0.00	0.00	1.00	0.00



Ablation Insights

- **Adversarial channel removal:** +22% deviation, thrust spikes reappear.
- **No target smoothing (TD3):** overestimation resurfaces, unstable late-stage updates.
- **Entropy off (SAC):** faster convergence, 9% worse robustness composite.
- **Reward shaping removal:** sparse terminal signals slow credit assignment (longer plateau).
- **Delay only vs. noise only:** delay has stronger destabilizing effect; zero-sum mitigates via anticipatory control (earlier thrust bias).



Key Findings

- Zero-sum MARL framing improves worst-case orbital maintenance robustness.
- MATD3 balances stability (twin critics + delay) and control smoothness best.
- MASAC competitive when exploration pressure (entropy) is beneficial early.
- Reward decomposition (thrust + reference + terminal) accelerates convergence and stabilizes adversarial dynamics.
- Policy smoothness correlates with fuel proxy reduction (8-12%).
- Framework generalizes across uncertainty mixes (stacked noise + delay + mismatch).

Conclusion: Adversarial co-training yields resilient guidance without explicit disturbance models.



Robustness Scenario Definitions

- **Random Init:** $x_0 \leftarrow x_0 + \mathcal{N}(0, 0.1^2)$
- **Actuator Disturbance:** $u_t \leftarrow u_t + \mathcal{N}(0, 0.05^2)$; (sensor additive) $y_t \leftarrow y_t + \mathcal{N}(0, 0.02^2)$
- **Model Mismatch:** $\theta \leftarrow \theta + \mathcal{N}(0, 0.05^2)$
- **Partial Observability:** mask 50% $\rightarrow m_t^{(i)} \sim \text{Bern}(0.5)$, $y_t \leftarrow y_t \circ m_t$
- **Sensor Noise (multiplicative):** $y_t \leftarrow y_t \circ (1 + \mathcal{N}(0, 0.05^2))$
- **Time Delay:** buffer length 10, $z u_t^{\text{applied}} \leftarrow u_{t-10} + \mathcal{N}(0, 0.05^2)$
- **Notes:**
 - All scenarios evaluated independently.
 - Zero-sum agents trained jointly once.
 - Metrics: success %, deviation, fuel proxy, return variance.

Delay + noise combo causes the largest degradation; adversarial training preserves stability margin.

