# Decentralized Aerial Manipulation of a Cable-Suspended Load using Multi-Agent Reinforcement Learning

Jack Zeng<sup>1\*</sup>, Andreu Matoses Gimenez<sup>1</sup>, Eugene Vinitsky<sup>2</sup>, Javier Alonso-Mora<sup>1</sup>, Sihao Sun<sup>1\*</sup>

<sup>1</sup>Delft University of Technology, <sup>2</sup>NYU Tandon School of Engineering

\*Corresponding authors: jack-zeng@hotmail.com, s.sun-2@tudelft.nl

**Abstract:** This paper presents the first decentralized method to enable real-world 6-DoF manipulation of a cable-suspended load using a team of Micro-Aerial Vehicles (MAVs). Our method leverages multi-agent reinforcement learning (MARL) to train an outer-loop control policy for each MAV. Unlike state-of-the-art controllers that utilize a centralized scheme, our policy does not require global states, inter-MAV communications, nor neighboring MAV information. Instead, agents communicate implicitly through load pose observations alone, which enables high scalability and flexibility. It also significantly reduces computing costs during inference time, enabling onboard deployment of the policy. In addition, we introduce a new action space design for the MAVs using linear acceleration and body rates. This choice, combined with a robust low-level controller, enables reliable sim-to-real transfer despite significant uncertainties caused by cable tension during dynamic 3D motion. We validate our method in various realworld experiments, including full-pose control under load model uncertainties, showing setpoint tracking performance comparable to the state-of-the-art centralized method. We also demonstrate cooperation amongst agents with heterogeneous control policies, and robustness to the complete in-flight loss of one MAV. Videos of experiments: https://autonomousrobots.nl/paper\_websites/ aerial-manipulation-marl

**Keywords:** Aerial Manipulation, Multi-Agent Reinforcement Learning, Micro Aerial Vehicles

# 1 Introduction

Autonomous Micro Aerial Vehicles (MAVs) offer great capability for transporting slung loads to dangerous and remote locations [1]. While a single low-cost MAV has limited payload capacity, collaborative teams of MAVs can transport significantly heavier loads. In addition, by connecting each MAV with the load at different points using tethers, the full pose of the load can be controlled by changing the position of the MAVs, yielding a cooperative cable-suspended manipulation solution, which shows great potential for aerial-based construction, inspection, and resecuring [2, 3, 4, 5, 6].

To coordinate and control MAV fleets, the state-of-the-art method [6] employs a centralized framework that accurately captures the strong dynamical coupling between the MAVs and the suspended load. This ensures safety and stability while addressing the significant underactuation inherent to cable-suspended systems, preventing actuator saturations and reciprocal collisions. However, using centralized control strategies for such systems suffers from critical drawbacks: computational complexity tends to scale exponentially with the number of agents for many approaches, rendering real-time control infeasible for larger teams with a centralized scheme [6, 7]. In addition, dependence on global state information and centralized communication is often impractical due to limits on sensors and communication bandwidth. A plausible solution, decentralization, remains an open challenge

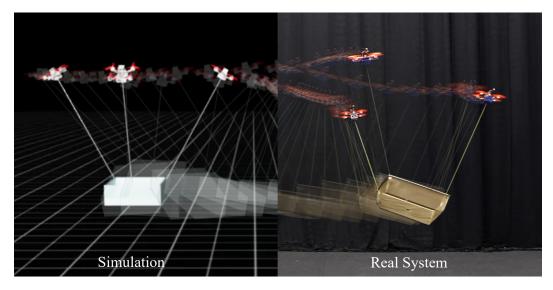


Figure 1: Multi-MAV lifting system performing full-pose control of a cable-suspended load. Left: simulation environment used to train the decentralized outer-loop control policy. Right: policy transferred to the real system.

to effectively coordinate MAV fleets due to partial observability, limited communication bandwidth, and decision-making under strong dynamical coupling between agents while co-manipulating an object.

In this work, we present the first decentralized algorithm to achieve a real-world demonstrated full-pose manipulation of a cable-suspended payload using a team of MAVs. Our method leverages multi-agent reinforcement learning (MARL) and **does not require any inter-agent communication**. Instead, each agent only takes their own state and identity, the load pose, and the target load pose as observations. We train the policy through MARL in a centralized training with decentralized execution (CTDE) paradigm using multi-agent proximal policy optimization (MAPPO) [8]. Each MAV learns to communicate implicitly through the load pose information. To fill the sim-to-real gap in this highly dynamic cooperative task, we design the action space of the reinforcement learning (RL) policy as reference linear accelerations and body rates of the MAV and combine the RL policy with a low-level controller based on incremental nonlinear dynamic inversion (INDI) [9, 10, 11]. The low-level controller follows the linear acceleration command with the body rate reference as the feedforward commands, ensuring agile and smooth control maneuvers during the cooperative manipulation.

Our method enables zero-shot transfer of the policy from simulation to real-world deployment to achieve full-pose control accuracy comparable to the state-of-the-art centralized controller [6], and is deployed fully onboard. In addition, experiments with real MAVs demonstrate that our method remains robust under load model uncertainties, operates effectively in heterogeneous agent settings where one MAV uses a different controller, and remains functional even when one of the MAVs completely fails.

# Our core contributions are as follows:

- The first method to achieve fully decentralized and onboard-deployed cooperative aerial manipulation in experiments with real MAVs, without any inter-agent communication.
- A novel action space design for MAVs manipulating a cable-suspended load, together with a robust low-level controller, enabling successful zero-shot sim-to-real transfer.
- First demonstration of robust full-pose control of the cable-suspended load under heterogeneous conditions and even under complete in-flight failure of an MAV.

# 2 Related works

Cooperative aerial manipulation of a cable-suspended load typically embraces a centralized paradigm to consider the cable-load-MAVs system as a whole and requires global state observations to ensure safety and performance. Early research on multi-MAV cable-suspended load problems often relied on model simplifications, such as assuming a quasi-static regime to ignore dynamic coupling effects [12, 13, 14, 15], which cannot address force-related constraints and perform dynamic motions. Another class of methods leverages system flatness [16] and dynamic equations to account for dynamic coupling effects. An example is the cascaded scheme, which employs an outer-loop geometric controller to generate the commanded wrench for the load, distributes it as desired cable tensions, and executes it through inner-loop controllers of MAVs [3, 17, 18, 19]. The outer-loop controller can be replaced by various approaches, such as inverse dynamics control [20], linear quadratic regulator [4], and nonlinear model predictive control (NMPC) [5]. Recent work [6] leverages whole-body dynamics and NMPC to generate reference trajectories followed by an adaptive low-level controller, showing high agility and accuracy.

However, these centralized methods require exponentially higher computational budgets and communication burdens with the number of agents involved. Therefore, decentralized controllers, such as distributed MPC [21, 22] have been proposed and tested in simulation to address the problem with the computational issues. But these methods still require reliable inter-agent data transfer to obtain real-time states from other agents, which does not fundamentally solve the problems with limited communication bandwidth.

Multi-agent reinforcement learning has been extensively studied for complex multi-agent systems, including cooperative scenarios [23, 24, 25]. Beyond achieving expert-level performance in video games [26, 27], MARL has been successfully applied to robotics, enabling decentralized control of multiple agents. For instance, researchers have leveraged MARL to develop cooperative strategies in robot football [28, 29], as well as multi-robot object manipulation with quadrupedal robots, including pushing [30] and cable-based towing [31]. Unlike our approach, these manipulation methods [30, 31] rely on neighboring agent information through communication or onboard perception. In many cases, MARL is employed to optimize high-level task objectives while relying on mid- and low-level controllers for motor and sub-task execution, capitalizing on RL's ability to optimize a long-horizon task-level objective [32].

Recent work by [33] demonstrates MARL's potential for cooperative object manipulation using simulated humanoids, relying solely on object bounding box information without explicit inter-agent communication. However, their approach depends on handcrafted reward functions that guide the humanoids toward predefined grasping points and walking behaviors. In MAV applications, MARL has been explored for tasks like swarming [34], but challenges remain due to the platform's agility, instability, and reliance on high-frequency, low-latency control [35]. Recently, MARL has shown potential for training multi-MAV lifting systems using global state observations [36]. However, a significant challenge remains to address the sim-to-real gap and partial observability, especially for the multi-MAV lifting system, where dynamic uncertainties are substantial due to complex aerodynamic disturbances and unknown cable tensions.

Our method effectively bridges this gap by leveraging multi-agent reinforcement learning (MARL) to achieve the first real-world demonstration of decentralized aerial manipulation, operating without global state observations or inter-agent communication. Furthermore, the method is deployed entirely onboard, enabled by its computational efficiency.

## 3 Methods

An overview of the full approach is shown in Figure 2. Our method utilizes MARL to train an outer-loop control policy, which generates reference accelerations and body rates for the low-level controller in real-time based on local observations of the ego-MAV state, its robot ID, payload- and goal pose. The low-level controller, including an INDI attitude controller, tracks these references

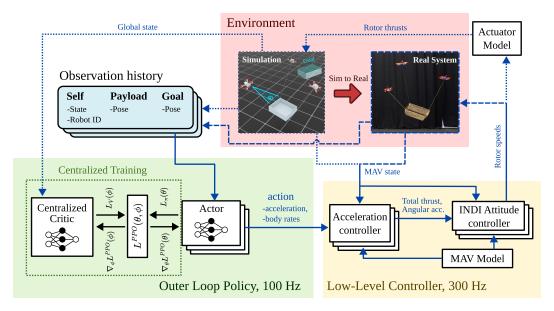


Figure 2: Overview of our method. Dotted lines indicate components only for training; dashed lines indicate those only for real-system deployment; solid lines for both. The training process involves the centralized critic (which observes the privileged global state), direct access to MAV states, and the actuator model that maps rotor speeds to thrust forces. Shared actors make decisions based on local observations, without access to other agents' states. The output actions, namely acceleration and body rates, are tracked by a robust model-based low-level controller based on INDI.

based on the MAV model and accelerometer measurements. The privileged full state is observed by the centralized critic during training, which is discarded at execution time. Collected experience is shared across actors to update the parameters of a shared policy. This enables training to be centralized while execution remains decentralized, allowing each agent to run the policy independently onboard after zero-shot transfer from simulation to the real world.

We model cooperative aerial manipulation as a decentralized partially observable Markov decision process (Dec-POMDP) [37] with a shared reward function. A Dec-POMDP is defined by  $\langle \mathcal{I}, \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ , where  $\mathcal{I}$  denotes the set of agents with the total number of agents being equal to N,  $\mathcal{S}$  is the environment state,  $\mathcal{A} = \{a_i\}_{i=1}^N$  is the joint action space of all agents,  $\mathcal{O} = \{o_i\}_{i=1}^N$  is each agent's partial observation of the environment,  $\mathcal{P}: \mathcal{S} \times \mathcal{A} \to \mathcal{S}$  is the transition model,  $R: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$  is the shared reward function and  $\gamma$  is the discount factor. At each timestep t, the current state  $\mathbf{s}_t \in \mathcal{S}$  transitions to a new state  $\mathbf{s}_{t+1}$  based on the joint action  $\mathbf{a}_t \in \mathcal{A}$  and the transition function  $\mathcal{P}$ . Each agent i then receives the shared reward as feedback from the environment.

Our approach employs the CTDE paradigm [38], utilizing privileged global state information during training for the asymmetric centralized critic while relying solely on local observations for policy execution. Each agent i has a policy  $\pi_i:\omega_i(o_i)\to a_i$  that maps its local observation, processed through its observation function  $\omega_i$ , to an action  $a_i$ . We implement parameter sharing across agents (i.e.,  $\pi_i=\pi_j \ \forall i,j$ ), thus reducing  $\pi_i$  to a homogeneous policy  $\pi$ . The set of observation functions for all agents can be denoted as  $\Omega=\{\omega_i\}_{i=1}^N$ . The final decentralized partially observable problem is thus defined by the tuple  $\langle \mathcal{I}, \mathcal{S}, \mathcal{A}, \mathcal{O}, \Omega, \mathcal{P}, \mathcal{R}, \gamma \rangle$ 

Observations and rewards The state of each MAV is given by  $\boldsymbol{x}_i = \begin{bmatrix} \boldsymbol{p}_{M,i}, \, \boldsymbol{R}_{M,i}, \, \boldsymbol{v}_{M,i}, \, \boldsymbol{\omega}_{M,i} \end{bmatrix}$ , where  $\boldsymbol{p}_{M,i} \in \mathbb{R}^3$  denotes the MAV's position,  $\boldsymbol{R}_{M,i} \in \mathbb{R}^9$  is the vector composed of elements of its rotation matrix,  $\boldsymbol{v}_{M,i} \in \mathbb{R}^3$  and  $\boldsymbol{\omega}_{M,i} \in \mathbb{R}^3$  denote its linear and angular velocities. We use the subscript i to denote the i-th MAV. The state of the load is given by  $\boldsymbol{x}_L = \begin{bmatrix} \boldsymbol{p}_L, \, \boldsymbol{R}_L, \, \boldsymbol{v}_L, \, \boldsymbol{\omega}_L \end{bmatrix}$  where  $\boldsymbol{p}_L \in \mathbb{R}^3$  denotes the load's position,  $\boldsymbol{R}_L \in \mathbb{R}^9$  is the vector composed of elements of its

rotation matrix,  $v_L \in \mathbb{R}^3$  and  $\omega_L \in \mathbb{R}^3$  denote its linear and angular velocities. The state of the goal relative to the payload is denoted by  $\boldsymbol{x}_G = \begin{bmatrix} \boldsymbol{d}_G, \, \boldsymbol{R}_G \end{bmatrix}$  where  $\boldsymbol{d}_G \in \mathbb{R}^3$  and  $\boldsymbol{R}_G \in \mathbb{R}^9$  represent the goal position relative to the current load position and the vector composed of elements of its relative rotation matrix from the current load orientation to the goal orientation respectively. All quantities are described in the inertial world frame  $\mathcal{F}_I$ . The global state that is observable to the centralized critic during training is then denoted as:

$$s = \begin{bmatrix} \boldsymbol{x}_L, \, \boldsymbol{x}_G, \, \boldsymbol{x}_{M,1}, \, \boldsymbol{x}_{M,2}, \, \cdots, \, \boldsymbol{x}_{M,N} \end{bmatrix}$$
 (1)

Where N is the total number of MAVs. The local policies' observation space only includes the load pose, relative goal terms, their own respective MAV state, and a one-hot vector  $e_i$  indicating their identity to enable role differentiation among homogeneous agents, as the policy network parameters are shared across all MAVs. The observation space for the i-th MAV is described as:

$$\boldsymbol{o}_i = \left[ \boldsymbol{p}_L, \, \boldsymbol{R}_L, \, \boldsymbol{x}_G, \, \boldsymbol{x}_{M,i}, \, \boldsymbol{e}_i \right] \tag{2}$$

As the problem is partially observable, we use a history of observations by stacking the current and last 2 observations of the policy [39]. For a more detailed discussion on the history length, we refer the readers to Appendix A.7.

We train the policies using MAPPO [8], a model-free MARL algorithm that extends PPO [40] with CTDE. The reward at time t, denoted as  $r_t$ , is defined as:

$$r_t = r_t^{\text{pos}} + r_t^{\text{ori}} + r_t^{\text{down}} + r_t^{\text{act}} + r_t^{\text{br}} + r_t^{\text{thrust}}$$
(3)

Where  $r_t^{\mathrm{pos}}$  and  $r_t^{\mathrm{ori}}$  are rewards to track the goal position and orientation for the load,  $r_t^{\mathrm{down}}$  encourages the MAVs to aim their (proxy) downwash away from the load for stability against aerodynamic disturbances,  $r_t^{\mathrm{act}}$  and  $r_t^{\mathrm{br}}$  penalize action changes from the last time step and large body-rate outputs respectively for smoother flight,  $r_t^{\mathrm{thrust}}$  penalizes outputting large thrusts which encourages energy efficiency. For a detailed reward formulation, we refer the readers to Appendix A.8.

Action space and low-level controller To balance reliable sim-to-real transfer with sufficient control authority, the choice of action space is critical. Prior work in single MAV control demonstrates that high-level outputs (e.g., position or velocity) enhance robustness to disturbances and sim-to-real gaps but limit performance, whereas low-level outputs (e.g., snap) improve tracking precision at the cost of larger transfer discrepancies [41, 42]. To address this trade-off, we propose a mid-level action space in desired accelerations and body rates (ACCBR). This approach preserves adequate control capability while also being robust against uncertain disturbances and model mismatches from the cable-suspended load.

The low-level controller converts the acceleration reference  $a_{i,ref}$  from the outer-loop policy to the thrust direction command through the following acceleration controller:

$$\boldsymbol{z}_{i,\text{des}} = \frac{\boldsymbol{a}_{i,\text{ref}} - \boldsymbol{g} - \boldsymbol{f}_{i,\text{ext}}/m_i}{\|\boldsymbol{a}_{i,\text{ref}} - \boldsymbol{g} - \boldsymbol{f}_{i,\text{ext}}/m_i\|}, \quad \boldsymbol{f}_{i,\text{ext}} = m_i \boldsymbol{a}_{i,\text{filtered}} - \boldsymbol{f}_{i,\text{filtered}}$$
(4)

where external forces  $f_{\text{ext}}$ , primarily due to the cable tensions, are estimated using the MAV mass  $m_i$ , filtered accelerometer measurements  $a_{i,\text{filtered}}$  and collective thrust  $f_{i,\text{filtered}}$  computed from a classical quadratic thrust model and filtered rotor speed feedbacks [11]. The desired attitude command and the policy output body-rate command are then sent to the INDI attitude controller to generate rotor speed commands. We refer readers to [9, 10, 11] for further details on INDI.

**Training setup** We train our method completely in simulation and achieve zero-shot transfer to real-world experiments. The simulation environment is built using NVIDIA's Isaac Lab [43], and the MARL algorithms are modified from [44]. Training was conducted on a consumer-grade RTX 3090 GPU and completed in 17 hours. The network architecture is a 4-layer MLP of size [1024, 512, 256, 128] for both the shared policies and the centralized critic. The inputs to the network are normalized stacked observation histories with history size H = 3. For a complete overview of training details, network and agent parameters, we refer the readers to Appendix A.9.

The MAVs with a cable-suspended load spawn uniformly between -1 and 1 in xy, 0.5 and 1.5 in z, with a random heading. The goal position is sampled from the same range, but also allows pitch and roll of  $\pm 45^{\circ}$ . Despite sampling of the goal is limited to the predefined sets, the policy is still able to generalize and reach goal poses outside of it during execution.

# 4 Experiments and Results

# 4.1 Real-world experiments

Setpoint tracking Our real-world experiments demonstrate agile pose control of three MAVs with a cablesuspended load, tracking a 2 m displacement with  $(30^{\circ}, -20^{\circ}, -90^{\circ})$  attitude commands. We compare our decentralized method with the stateof-the-art centralized NMPC approach [6] in Figure 3. Despite being fully decentralized, our method achieves comparable tracking performance with positional and attitude RMSEs of 0.52 m (vs 0.45 m) and 22.93° (vs 16.24°), respectively. Note that RMSE comparisons favor NMPC as it tracks a reference trajec-

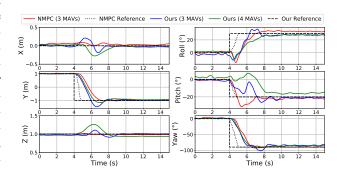


Figure 3: Time series of pose tracking results comparing our method and a centralized NMPC method [6]. Our method also includes a setup with 4 MAVs.

tory while we only track target poses, resulting in a larger RMSE in the transient area of the step command. The time-to-target (error  $< 0.10\,\text{m}/10^\circ$ ) is 6.84 s for NMPC vs. 8.36 s for ours, and the final displacement (RMSE) is  $0.05\,\text{m}/4.02^\circ$  for NMPC vs.  $0.04\,\text{m}/5.78^\circ$  for ours. We also show successful pose control with 4 MAVs (without cable slack), resulting in tracking RMSEs of 0.92 m and  $42.67^\circ$ . The increased error, compared to the 3 MAV case, may be due to the system becoming overconstrained, which introduces more complex coordination and (cable) dynamics [18]. In terms of computational efficiency, we run the NMPC and our method onboard a Raspberry Pi 5 (2.4 GHz quad-core ARM Cortex-A76). Our method inferences in **6 ms** at 100 Hz, versus NMPC's **78 ms** at 10 Hz. Crucially, while NMPC's computation time grows exponentially with agent count, e.g., 174 ms and 267 ms for 5 and 6 agents respectively, our agent-independent approach maintains **constant computation time regardless of team size**.

**Robustness against load model mismatch** To evaluate robustness, we add objects (0.216 kg, 15.4% of load mass) to the load, including four freely movable items that dynamically perturb both mass distribution and center of mass. Despite no inertia randomization during training, the system maintains strong tracking performance (0.63 m vs 0.60 m position RMSE; 26.93° vs 26.49° attitude RMSE. The low-level feedback controller automatically compensates for these disturbances, demonstrating inherent robustness to model uncertainties. Experimental results are shown in Figure 4B.

**Heterogeneous agents** Although our policy is trained under the assumption of homogeneous agents, it remains effective when deployed with heterogeneous agents. In this experiment, we let the load hover at a fixed point. Then we hacked one of the MAVs by replacing its RL policy with a model-based controller [11], and provided it with different setpoints to observe the behavior of the other two MAVs controlled by the RL policy. Specifically, we commanded the hacked MAV to move outwards on the y-axis by 0.7 m to pull the load away from the reference; we then commanded the hacked MAV to move inwards by 0.3 m to push it closer to the other two MAVs. Figure 4A provides a snapshot of the experiments. Since the policy is conditioned solely on the load pose and not on the states of the other agents, the two remaining MAVs utilizing the policy can compensate for load pose deviations from the reference. In contrast, the fully observable policy fails in these conditions due to the dependence on the states of all agents. Time series can be found in Appendix A.3.

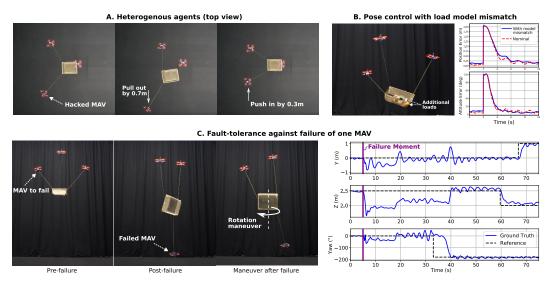


Figure 4: Real-world experiments. (A) Snapshot of the test with heterogeneous agents in which one MAV is manually controlled (hacked) to pull out and push in, and the other two MAVs counteract the interference of the hacked MAV. (B) Snapshot of the test where additional load is added to the original load, and the pose error with and without such model mismatch. (C) Snapshot of the case where one MAV fails in flight and the remaining two MAVs manage to control the load.

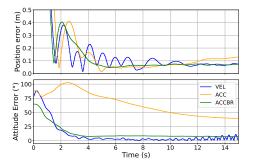
In-flight failure of one MAV The effectiveness of our method with a heterogeneous agent setup and robustness against load model uncertainties also offers strong fault tolerance in the case of agent failure. In this experiment, we deliberately turned off the hacked MAV (one of the two on the same side). As a result, the load was controlled by the remaining two MAVs. Note that with only two MAVs, the load orientation around the line joining the remaining two attachment points becomes unactuated. Even worse, the failed MAV hangs underneath the load, leading to additional disturbances to the post-failure system. Despite that, our method allows the other two MAVs to effectively control the remaining 5 DoFs of the load. We show that the system is still able to yaw by -180° and is also able to maintain position control by flying 0.5 meters down along the z-axis and maneuvering along the y-axis by 1 meter. The tracking results and snapshots of the setup after the failure are seen in Figure 4C. As in the heterogeneous agent case, the remaining agents can compensate for the missing MAV since the policy operates independently of other agents' states, thereby avoiding unstable behavior in out-of-distribution scenarios. In contrast, the fully observable policy fails under these conditions due to its reliance on the states of all agents. Time series illustrating both scenarios can be found in Appendix A.4.

# 4.2 Comparison among different action and observation spaces

We compare our selected observation and action spaces with alternatives in simulation for safety. The Agilicious flight stack is used with the Gazebo simulator [45] and RotorS [46] plugins, which add sensor noise, aerodynamic disturbances, and system latencies in a ROS environment. All policies are trained for 1 billion environment steps (10 h) and evaluated 10 times in Gazebo.

Action space We compare the ACCBR action space with three alternatives: velocity (VEL), linear acceleration (ACC), and collective thrust with body rates (CTBR). The ACCBR, VEL, and ACC outputs all utilize the same low-level controllers, which compensate for disturbances such as aero-dynamic forces and cable tension. In contrast, CTBR outputs feed directly into the INDI attitude controller without additional disturbance compensation.

The RMSE results in Table 1 demonstrate that the VEL action space achieves the best performance, followed by ACCBR, while ACC fails to track the load orientation accurately. Notably, the widely used CTBR approach [47, 32] fails to learn effectively. Since CTBR directly commands collective



Action space	Pos RMSE	Att RMSE
ACCBR	$0.64 \pm 0.00$	$33.87 \pm 0.91$
CTBR*	NaN	NaN
ACC	$0.54 \pm 0.00$	$87.89 \pm 1.85$
VEL	$0.56 \pm 0.06$	$\textbf{25.74} \pm \textbf{1.49}$
*Not able to take off		

Figure 5: Positional and attitude errors comparing different action spaces at test time in the Gazebo environment.

Table 1: Pose tracking RMSEs of different action spaces at test time in the Gazebo environment

thrust without leveraging the proposed low-level controller's disturbance compensation, we hypothesize that the unpredictable cable forces exerted on each MAV make the learning process prohibitively difficult, as there are no cable force sensors mounted for both training and evaluations.

However, while VEL achieves lower RMSE, Figure 5 shows it causes **hazardous oscillations**. ACCBR offers more stable hovering despite higher initial errors, making it safer and preferable for stability-critical tasks like inspection or delivery.

**Observation space** To benchmark the decentralized policy's performance, we compare three observation space cases: (1) the fully observable case with global state  $s = [x_L, x_G, x_1, x_2, x_3]$ , (2) an augmented partial observability case where each MAV i also receives the load twist and other MAVs' positions ("Partial augmented")  $o_i = [x_L, x_G, p_{j_1}, p_{j_2}, x_i, e_i]$  with  $p_{j_1}, p_{j_2}$  representing the neighboring agents' positions, and (3) the partially observable case. For partially observable cases, we include observation histories (H = 3) to improve state estimation and decision-making under uncertainty [39]. Figure 6 reveals comparable convergence across all configurations, indicating that load pose alone serves as a sufficient statistic for implicit MAV coordination, while the full global state contains redundant elements.

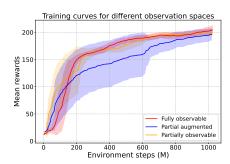


Figure 6: Training curves of fully observable, partial augmented, and partially observable observation spaces.

# 5 Conclusion

We introduced a decentralized method using MARL that allows for full-pose control of a cable-suspended load using three MAVs without any inter-MAV communication or neighboring MAV information. The policy is computationally tractable and executes entirely onboard. We proposed a novel action space of accelerations and body rates (ACCBR) along with a robust low-level controller and showcase zero-shot transfer from simulation to real-world deployment. Extensive testing with real MAVs shows that the setpoint tracking performance of our method is comparable to that of the state-of-the-art centralized NMPC [6], despite being fully decentralized and having significantly lower computation time. Our method demonstrates robustness against unknown disturbances, heterogeneous agents, and even the complete in-flight failure of one MAV. We attribute this resilience to two key factors: 1) closed-loop reference tracking by the low-level controller, which maintains stability despite perturbations, 2) decentralized policy independence, where local agents operate without dependence on neighboring states, preventing cascading failures. Our work shows promising results to enable scalable and robust cooperative aerial manipulation with minimal onboard sensing and no internal communications required.

## Acknowledgments

The authors would like to thank Dr. Yunlong Song, Dr. Dennis Benders, and Shlok Deshmukh for the insightful discussions, and Maurits Pfaff and Kseniia Khomenko for their help with the experiments. This work is funded by the Dutch Research Council (NWO) under Grant 20256 for the project "Accurate Aerial Manipulation under Uncertainties" and by the European Union under the ERC grant INTERACT, 101041863. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.

## References

- [1] E. N. Barmpounakis, E. I. Vlahogianni, and J. C. Golias. Unmanned aerial aircraft systems for transportation engineering: Current practice and future challenges. *International Journal of Transportation Science and Technology*, 5(3):111–122, 2016.
- [2] K. Sreenath and V. Kumar. Dynamics, control and planning for cooperative manipulation of payloads suspended by cables from multiple quadrotor robots. *rn*, 1(r2):r3, 2013.
- [3] T. Lee. Geometric control of quadrotor uavs transporting a cable-suspended rigid body. *IEEE Transactions on Control Systems Technology*, 26(1):255–264, 2017.
- [4] J. Geng, P. Singla, and J. W. Langelaan. Load-distribution-based trajectory planning and control for a multilift system. *Journal of Aerospace Information Systems*, 19(5):366–381, 2022.
- [5] G. Li and G. Loianno. Nonlinear model predictive control for cooperative transportation and manipulation of cable suspended payloads with multiple quadrotors. In 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 5034–5041. IEEE, 2023.
- [6] S. Sun, X. Wang, D. Sanalitro, A. Franchi, M. Tognon, and J. Alonso-Mora. Agile and cooperative aerial manipulation of a cable-suspended load. arXiv preprint arXiv:2501.18802, 2025.
- [7] L. Bakule and M. Papik. Decentralized control and communication. *Annual Reviews in Control*, 36(1):1–10, 2012.
- [8] C. Yu, A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. Bayen, and Y. Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in neural information processing systems*, 35:24611–24624, 2022.
- [9] E. J. Smeur, Q. Chu, and G. C. De Croon. Adaptive incremental nonlinear dynamic inversion for attitude control of micro air vehicles. *Journal of Guidance, Control, and Dynamics*, 39(3): 450–461, 2016.
- [10] E. Tal and S. Karaman. Accurate tracking of aggressive quadrotor trajectories using incremental nonlinear dynamic inversion and differential flatness. *IEEE Transactions on Control Systems Technology*, 29(3):1203–1218, 2020.
- [11] S. Sun, A. Romero, P. Foehn, E. Kaufmann, and D. Scaramuzza. A comparative study of non-linear mpc and differential-flatness-based control for quadrotor agile flight. *IEEE Transactions on Robotics*, 38(6):3357–3373, 2022.
- [12] J. Fink, N. Michael, S. Kim, and V. Kumar. Planning and control for cooperative manipulation and transportation with aerial robots. In *Robotics Research: The 14th International Symposium ISRR*, pages 643–659. Springer, 2011.

- [13] N. Michael, J. Fink, and V. Kumar. Cooperative manipulation and transportation with aerial robots. *Autonomous Robots*, 30:73–86, 2011.
- [14] M. Manubens, D. Devaurs, L. Ros, and J. Cortés. Motion planning for 6-d manipulation with aerial towed-cable systems. In *Robotics: science and systems (RSS)*, page 8p, 2013.
- [15] D. Sanalitro, H. J. Savino, M. Tognon, J. Cortés, and A. Franchi. Full-pose manipulation control of a cable-suspended load with multiple uavs under uncertainties. *IEEE Robotics and Automation Letters*, 5(2):2185–2191, 2020.
- [16] K. Sreenath, T. Lee, and V. Kumar. Geometric control and differential flatness of a quadrotor uav with a cable-suspended load. In *52nd IEEE conference on decision and control*, pages 2269–2274. IEEE, 2013.
- [17] G. Li, R. Ge, and G. Loianno. Cooperative transportation of cable suspended payloads with mavs using monocular vision and inertial sensing. *IEEE Robotics and Automation Letters*, 6 (3):5316–5323, 2021.
- [18] G. Li, X. Liu, and G. Loianno. Rotortm: A flexible simulator for aerial transportation and manipulation. *IEEE Transactions on Robotics*, 40:831–850, 2023.
- [19] K. Wahba and W. Hönig. Efficient optimization-based cable force allocation for geometric control of a multirotor team transporting a payload. *IEEE Robotics and Automation Letters*, 9 (4):3688–3695, 2024.
- [20] C. Masone and P. Stegagno. Shared control of an aerial cooperative transportation system with a cable-suspended payload. *Journal of Intelligent & Robotic Systems*, 103(3):40, 2021.
- [21] J. Wehbeh, S. Rahman, and I. Sharf. Distributed model predictive control for uavs collaborative payload transport. In 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 11666–11672. IEEE, 2020.
- [22] B. Wang, R. Huang, and L. Zhao. Auto-multilift: Distributed learning and control for cooperative load transportation with quadrotors. *arXiv* preprint arXiv:2406.04858, 2024.
- [23] L. Busoniu, R. Babuska, and B. De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 38(2):156–172, 2008.
- [24] K. Zhang, Z. Yang, and T. Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pages 321–384, 2021.
- [25] J. K. Gupta, M. Egorov, and M. Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In Autonomous Agents and Multiagent Systems: AAMAS 2017 Workshops, Best Papers, São Paulo, Brazil, May 8-12, 2017, Revised Selected Papers 16, pages 66–83. Springer, 2017.
- [26] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *nature*, 575(7782):350–354, 2019.
- [27] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Dębiak, C. Dennison, D. Farhi, Q. Fischer, S. Hashme, C. Hesse, et al. Dota 2 with large scale deep reinforcement learning. arXiv preprint arXiv:1912.06680, 2019.
- [28] S. Liu, G. Lever, Z. Wang, J. Merel, S. A. Eslami, D. Hennes, W. M. Czarnecki, Y. Tassa, S. Omidshafiei, A. Abdolmaleki, et al. From motor control to team play in simulated humanoid football. *Science Robotics*, 7(69):eabo0235, 2022.

- [29] Z. Li, F. Bjelonic, V. Klemm, and M. Hutter. Marladona-towards cooperative team play using multi-agent reinforcement learning. *arXiv preprint arXiv:2409.20326*, 2024.
- [30] Y. Feng, C. Hong, Y. Niu, S. Liu, Y. Yang, W. Yu, T. Zhang, J. Tan, and D. Zhao. Learning multi-agent loco-manipulation for long-horizon quadrupedal pushing. *arXiv* preprint *arXiv*:2411.07104, 2024.
- [31] W.-T. Chen, M. Nguyen, Z. Li, G. N. Sue, and K. Sreenath. Decentralized navigation of a cable-towed load using quadrupedal robot team via marl. *arXiv preprint arXiv:2503.18221*, 2025.
- [32] Y. Song, A. Romero, M. Müller, V. Koltun, and D. Scaramuzza. Reaching the limit in autonomous racing: Optimal control versus reinforcement learning. *Science Robotics*, 8(82): eadg1462, 2023.
- [33] J. Gao, Z. Wang, Z. Xiao, J. Wang, T. Wang, J. Cao, X. Hu, S. Liu, J. Dai, and J. Pang. Coohoi: Learning cooperative human-object interaction with manipulated object dynamics. *Advances in Neural Information Processing Systems*, 37:79741–79763, 2024.
- [34] S. Batra, Z. Huang, A. Petrenko, T. Kumar, A. Molchanov, and G. S. Sukhatme. Decentralized control of quadrotor swarms with end-to-end deep reinforcement learning. In *Conference on robot learning*, pages 576–586. PMLR, 2022.
- [35] J. Xing, A. Romero, L. Bauersfeld, and D. Scaramuzza. Bootstrapping reinforcement learning with imitation for vision-based agile flight. In 8th Annual Conference on Robot Learning.
- [36] B. Xu, F. Gao, C. Yu, R. Zhang, Y. Wu, and Y. Wang. Omnidrones: An efficient and flexible platform for reinforcement learning in drone control. *IEEE Robotics and Automation Letters*, 9(3):2838–2844, 2024.
- [37] F. A. Oliehoek and C. Amato. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
- [38] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch. Multi-agent actorcritic for mixed cooperative-competitive environments. Advances in neural information processing systems, 30, 2017.
- [39] M. J. Hausknecht and P. Stone. Deep recurrent q-learning for partially observable mdps. In *AAAI fall symposia*, volume 45, page 141, 2015.
- [40] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [41] E. Kaufmann, L. Bauersfeld, and D. Scaramuzza. A benchmark comparison of learned control policies for agile quadrotor flight. In 2022 International Conference on Robotics and Automation (ICRA), pages 10504–10510. IEEE, 2022.
- [42] J. Eschmann, D. Albani, and G. Loianno. Learning to fly in seconds. *IEEE Robotics and Automation Letters*, 2024.
- [43] M. Mittal, C. Yu, Q. Yu, J. Liu, N. Rudin, D. Hoeller, J. L. Yuan, R. Singh, Y. Guo, H. Mazhar, et al. Orbit: A unified simulation framework for interactive robot learning environments. *IEEE Robotics and Automation Letters*, 8(6):3740–3747, 2023.
- [44] A. Serrano-Muñoz, D. Chrysostomou, S. Bøgh, and N. Arana-Arexolaleiba. skrl: Modular and flexible library for reinforcement learning. *Journal of Machine Learning Research*, 24(254): 1–9, 2023.

- [45] N. Koenig and A. Howard. Design and use paradigms for Gazebo, an open-source multi-robot simulator. In *Proceedings of the 2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, volume 3, pages 2149–2154. IEEE, 2004.
- [46] F. Furrer, M. Burri, M. Achtelik, and R. Siegwart. Rotors—a modular gazebo may simulator framework. *Robot Operating System (ROS) The Complete Reference (Volume 1)*, pages 595–625, 2016.
- [47] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza. Champion-level drone racing using deep reinforcement learning. *Nature*, 620(7976):982–987, 2023.
- [48] P. Foehn, E. Kaufmann, A. Romero, R. Penicka, S. Sun, L. Bauersfeld, T. Laengle, G. Cioffi, Y. Song, A. Loquercio, et al. Agilicious: Open-source and open-hardware agile quadrotor for vision-based flight. *Science robotics*, 7(67):eabl6259, 2022.
- [49] J. Mayer, J. Westermann, J. P. G. H. Muriedas, U. Mettin, and A. Lampe. Proximal policy optimization for tracking control exploiting future reference information. *arXiv* preprint arXiv:2107.09647, 2021.
- [50] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [51] M. Cusumano-Towner, D. Hafner, A. Hertzberg, B. Huval, A. Petrenko, E. Vinitsky, E. Wijmans, T. Killian, S. Bowers, O. Sener, et al. Robust autonomy emerges from self-play. *arXiv* preprint arXiv:2502.03349, 2025.

# A Appendix

#### A.1 Limitations

Our method requires pose measurement of the load, which is not often practical beyond lab environments. In our experiment, we require an external motion capture system to provide high-frequency load pose measurement. For future real-world outdoor deployment, onboard sensing (e.g., a downward-facing camera for load pose estimation and SLAM for MAV localization) would be necessary. This would introduce new challenges, such as observation delays, imperfect state estimates, sensor noise, and different reference frames for the load and MAVs that require alignment and synchronization. Additionally, our current framework does not address obstacle avoidance, as we assume collision-free paths to the goal—an unrealistic assumption in unstructured environments. Future work will focus on integrating a robust perception stack and obstacle avoidance capabilities.

## A.2 Experimental setup

**Real-world evaluation setup** We evaluate our method in real-world experiments. Our experiment includes 3 MAVs built based on the Agilicious [48] flight stack. Each MAV is connected to a basket-shaped payload with 1-meter cables at three distinct locations. The MAVs weigh 0.6kg, and the payload weighs 1.4 kg. We conduct the experiment in an indoor flight space with motion capture systems. We attach motion capture markers to the MAVs and the payload to measure their positions and orientations and distribute them to each MAV through ROS at 100 Hz. The trained policy and low-level controllers are deployed onboard each MAV. The policy is inferred at 100 Hz to send acceleration and body-rate commands. The low-level controller is executed at 300 Hz to generate rotor speed commands.

#### A.3 Heterogeneous agents time series

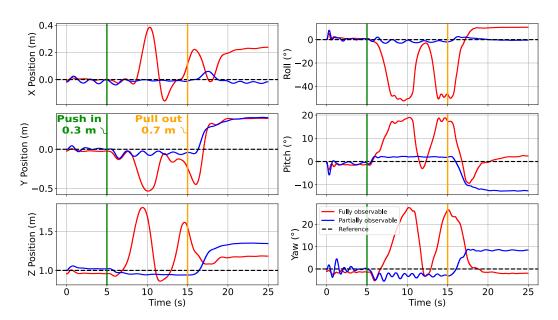


Figure 7: Time series of the load pose in the heterogeneous agents scenario, comparing the performance of the partially observable policy and the fully observable policy. The time points at which control commands are issued to push the load inward by 0.3 m relative to the desired policy position, or to pull it outward by 0.7 m, are indicated in green (push-in) and orange (pull-out), respectively.

Figure 7 compares the performance of partially observable and fully observable policies in the heterogeneous agents scenario. The partially observable policy, being independent of other agents'

states, allows the unaffected MAVs to compensate for the hacked agent, maintaining system stability. In contrast, the fully observable policy—which relies on neighboring agents' states—performs worse, exhibiting larger tracking errors (0.42 m vs. 0.28 m in position, 30.08 degrees vs. 8.88 degrees in attitude) and large oscillations during the inward push.

# A.4 In-flight failure of one MAV time series

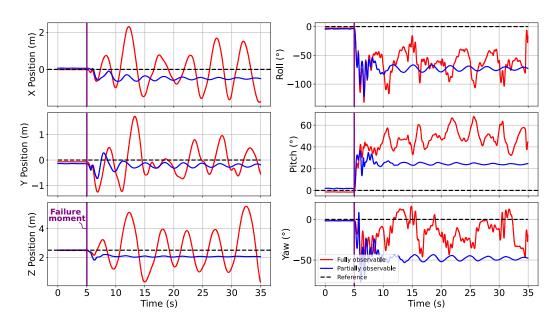


Figure 8: Time series of load pose in the in-flight failure of one MAV case without sending any commands, comparing a partially observable policy vs a fully observable policy. The thick purple line indicates the moment the MAV fails.

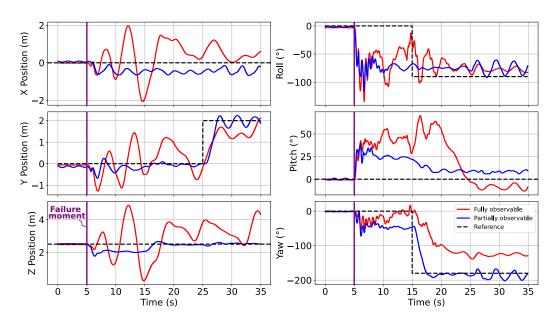


Figure 9: Time series of load pose in the in-flight failure of one MAV case, comparing a partially observable policy vs a fully observable policy. An attitude command is sent after 10 seconds and a positional command after 20 seconds. The thick purple line indicates the moment the MAV fails.

Figures 8 and 9 show the tracking performance of the partially observable and fully observable policies following an in-flight failure of one MAV. Figure 8 corresponds to the scenario in which no additional command inputs are issued, whereas Figure 9 corresponds to the scenario in which new attitude and position commands are introduced at  $t=15\,\mathrm{s}$  and  $t=25\,\mathrm{s}$ . In both scenarios, the partially observable policy successfully compensates for the MAV failure. In contrast, the fully observable policy exhibits strong oscillatory behavior, causing the suspended MAV to repeatedly crash to the ground. When new pose commands are sent, the fully observable policy fails to track them accurately, whereas the partially observable policy is still able to track 5 DoF. This results in larger tracking errors for the fully observable policy, which incurs position and attitude root-mean-square errors of 1.50 m and 73.37 degrees, respectively, compared to 0.67 m and 50.31 degrees for the partially observable policy. The robustness of the partially observable policy is attributed to its independence from the states of neighboring agents, which helps prevent cascading failures.

# A.5 Trajectory tracking

Although our method is not trained for trajectory tracking, we evaluate its trajectory tracking capabilities against that of the centralized NMPC [6] in Figure 10. The reference trajectory is a figure-eight trajectory with a maximum velocity of 1 m/s and a maximum acceleration of 0.5 m/s<sup>2</sup>. It is worth noting that our method only considers the reference pose information, while the NMPC also takes velocity information from the reference trajectory into account. For future specialized trajectory tasks, incorporating higher-order information such as velocity, as well as future reference points [49] into the observations would significantly improve tracking performance and make for a fairer comparison. Nonetheless, our method is able to successfully track the figure-eight trajectory, albeit with a high tracking error. Our method achieves positional and attitude RMSEs of 0.82 m (vs 0.10 m), and 18.22 degrees (vs 4.80 degrees).

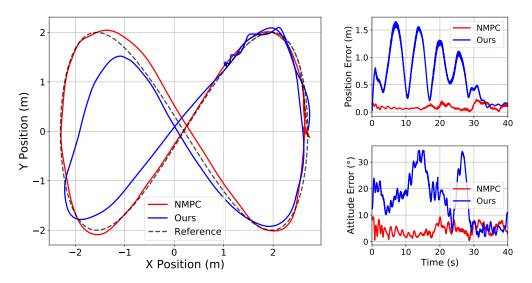


Figure 10: Comparison of our method, which is not trained for trajectory tracking, against the centralized NMPC in [6]. **Left**: top view of the flight path of the center of mass of the load while tracking a figure-eight trajectory with a maximum velocity of 1 m/s and maximum acceleration of 0.5 m/s<sup>2</sup>. **Right**: position (top) and attitude (bottom) tracking errors time series.

#### A.6 Performance without centralized critic

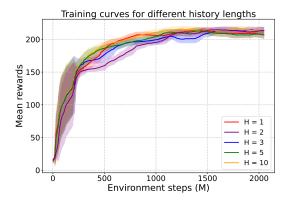
To assess the impact of using a centralized critic with access to privileged global state information, we compare its performance against a policy trained with a shared local critic. The local



critic has access only to local observations, which are the same as those available to the actor. The training curves in Figure 11 show that the setup with the local critic fails to converge to the same performance as with the centralized critic, and even collapses at the end. Specifically, the policy with the local critic fails to learn the position and orientation rewards effectively. We hypothesize that access to global state information allows the centralized critic to produce more accurate value estimates, which can indirectly support more effective credit assignment during learning [50], thereby improving task performance.

## A.7 Performance with different history lengths

We compare the performance of the partially observable policy with different history lengths in the observation space. H=1 means that the history only contains the observations of the current timestep (no previous observations). All policies are trained on a limited budget of 2 billion environment steps and are evaluated in the Gazebo environment.



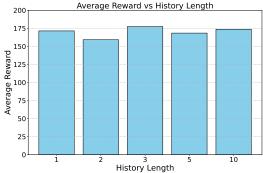


Figure 12: Training curves comparing different history lengths for the partially observable policy.

Figure 13: Mean reward of policies with different history lengths over 10 runs at test time in the Gazebo environment.

Figures 12 and 13 show that including historical observations has little impact on performance. We hypothesize that the load's pose—even without historical data—contains enough information to estimate the other agents' states, enabling implicit communication among the MAVs. Further investigation into the role of history in more complex scenarios, such as those with higher noise or additional MAVs, is left for future work.

## A.8 Reward function formulation

The reward function components are formulated as:

$$r_{t}^{\text{pos}} = \lambda_{1} \exp\left(-\lambda_{2} \|\boldsymbol{p}_{G} - \boldsymbol{p}_{L}\|\right),$$

$$r_{t}^{\text{ori}} = \lambda_{3} \exp\left(-\lambda_{4} \theta(\boldsymbol{q}_{G}, \boldsymbol{q}_{L})\right),$$

$$r_{t}^{\text{down}} = \lambda_{5} \left(1 - \exp\left(-\lambda_{6} \cdot \min_{i} \|f_{\text{int}}(\boldsymbol{p}_{M,i}, \boldsymbol{t}_{i}) - \boldsymbol{p}_{L}\|\right)\right),$$

$$r_{t}^{\text{act}} = \lambda_{7} \exp\left(-\|(\boldsymbol{a}_{t} - \boldsymbol{a}_{t-1})/N\|^{2}\right),$$

$$r_{t}^{\text{br}} = \lambda_{8} \exp\left(-\|\boldsymbol{\omega}_{t}/N\|\right),$$

$$r_{t}^{\text{thrust}} = \lambda_{9} \exp\left(-\max(\boldsymbol{T}_{t}/T_{\text{max}})\right),$$
(5)

Here  $p_G$  and  $p_L$  denote the goal and load positions respectively.  $\theta(q_G, q_L)$  denotes the quaternion error magnitude function which is calculated using the quaternion representation of the goal orientation  $q_G$ , and the load orientation  $q_L$ . The error is calculated by taking the norm of the axis-angle representation of the quaternion difference  $q_G \otimes q_L^*$ , where  $q_L^*$  is the conjugate of  $q_L$ .

The function  $f_{\text{int}}(p_{M,i}, t_i)$  computes the intersection point between two elements: the line defined by the i-th MAV's position  $p_{M,i}$  and its thrust direction  $t_i$ , and the plane containing the payload. This payload plane is characterized by its normal vector  $n = \ell_x \times \ell_y$ , where  $\ell_x$  and  $\ell_y$  represent arbitrary vectors spanning the load's local x-y plane. From all such intersection points computed for each MAV, the operator min selects the closest one to the payload position, corresponding to the most significant downwash effect.

The intersection calculation expands to:

$$f_{\text{int}}(\mathbf{p}_{M,i}, \mathbf{t}_i) = \mathbf{p}_{M,i} + \left(\frac{d - \mathbf{n} \cdot \mathbf{p}_{M,i}}{\mathbf{n} \cdot \mathbf{t}_i}\right) \mathbf{t}_i$$
 (6)

where  $d = n \cdot p_L$  defines the payload plane's offset from the origin through the payload position  $p_L$ .

The amount of MAVs is denoted by N, and a represents the control command, and  $\omega$  the body rate part of the control command.  $T \in \mathbb{R}^{4N}$  is the vector containing the rotor thrusts from each MAV, which is then normalized by the maximum thrust output  $T_{\max}$ .  $\lambda_1, \lambda_2 \cdots \lambda_9$  are different positive hyperparameters. All components are normalized by the simulation frequency. The chosen hyperparameters are shown in Table 2.

## A.9 Training configuration

The inputs to the network are normalized stacked observation histories with history size H=3. We also implement a form of advantage filtering [51] where 50% of the samples with the lowest advantage magnitude are dropped. This approach prioritizes learning from the most informative state transitions—specifically the underexplored extremes of the data distribution where actions have a clearly better or worse outcome—thereby improving data efficiency during training. For a complete overview of the network and agent parameters, we refer the readers to Table 3.

For setups with more than 3 MAVs, the mass of the load is sampled from a uniform distribution between 1.0 and 1.8 kg (the mass of the real payload is 1.4 kg). For the 3-MAV setup, the cables are modeled as rigid rods of 1 meter in length, connected to both the payload and the MAVs via ball joints. When using more than 3 MAVs, the system becomes overconstrained, which can lead to cable slack [12]. To address this, the cables are instead modeled as three rigid segments linked by ball joints.

The episodes have a duration of 20 seconds, where a single goal pose is given to encourage stable hovering of the payload. The episode times out after 20 seconds, in which case the return is bootstrapped using the value function estimate, or it terminates earlier if:

- any MAV or the payload is too close to the ground,
- the angle between the payload and the cable exceeds a certain threshold,

- the angle between the cable and the MAV exceeds a certain threshold,
- cables collide with each other,
- MAVs collide with each other,
- any rigid body is outside a specified bounding box,
- any of the cable tensions are below a specified threshold. ( $> 3~\mathrm{MAVs})$

**Reward function weights** The reward function weights shown in Table 2 are based on iterative tuning in simulation and real-world experiments.

Reward weight	Value
$\lambda_1$	1.5
$\lambda_2$	1.5
$\lambda_3$	1.5
$\lambda_4$	1.5
$\lambda_5$	0.5
$\lambda_6$	3.0
$\lambda_7$	0.5
$\lambda_8$	0.5
$\lambda_9$	0.5

Table 2: Reward function weights

**Hyperparameters of MAPPO** The hyperparameters of MAPPO are shown in table 3. The names of the parameters are based on the SKRL [44] learning library.

Hyperparameter	Value
number of envs	4096
rollouts	128
learing epochs	5
mini batches	4
discount factor	0.99
gae lambda	0.95
learning rate actor	5e-4
learning rate critic	1e-4
state preprocessor	RunningStandardScaler
shared state preprocessor	RunningStandardScaler
value preprocessor	RunningStandardScaler
grad norm clip	1.0
ratio clip	0.1
value clip	0.1
entropy loss scale	0.001
value loss scale	1.0
kl threshold	0.0

Table 3: MAPPO hyperparameters based on SKRL [44] learning library