# Quadrotor Morpho-Transition:
# Learning vs Model-Based Control Strategies

Ioannis Mandralis[1,2*], Richard M. Murray[2], Morteza Gharib[1]

*Abstract*— Quadrotor Morpho-Transition, or the act of transitioning from air to ground through mid-air transformation, involves complex aerodynamic interactions and a need to operate near actuator saturation, complicating controller design. In recent work, morpho-transition has been studied from a model-based control perspective, but these approaches remain limited due to unmodeled dynamics and the requirement for planning through contacts. Here, we train an end-to-end Reinforcement Learning (RL) controller to learn a morpho-transition policy and demonstrate successful transfer to hardware. We find that the RL control policy achieves agile landing, but only transfers to hardware if motor dynamics and observation delays are taken into account. On the other hand, a baseline MPC controller transfers out-of-the-box without knowledge of the actuator dynamics and delays, at the cost of reduced recovery from disturbances in the event of unknown actuator failures. Our work opens the way for more robust control of agile in-flight quadrotor maneuvers that require mid-air transformation. Video; Code.

## I. INTRODUCTION

Ground aerial robotic systems are ideally poised to increase the reliability and scope of autonomous robotic missions. Whilst robots specially adapted to single locomotion types like quadrotors or bipeds may achieve impressive performance in their respective domains, they suffer from fundamental disadvantages. Aerial robots face important limitations due to battery life and payload capacity, whilst ground robots have limited ability to explore the aerial domain. Combining ground and aerial locomotion modes thus promises to deliver increased versatility, better energy efficiency, expanded exploration capabilities, as well as increased fault tolerance.

In recent work, a novel ground aerial robot named ATMO (Aerially Transforming Morphobot) capable of *morpho-transition* was introduced [1]. This type of maneuver, depicted in Figure 1 requires smoothly transitioning from flying to driving locomotion by morphing near the ground and landing on dual-purpose wheel-thruster appendages with the largest possible tilt angle, i.e. as close as possible to drive configuration. Contrary to conventional quadrotor landings that involve landing by vertical, non-transforming descent, transforming prior to landing can enable landing and driving in crevices, morphing to enter through narrow tunnels, or

[1]Aerospace Engineering Department, California Institute of Technology, 1200 E California Blvd, Pasadena, CA, USA mgharib@caltech.edu

[2]Computing and Mathematical Sciences Department, Caltech, MA, California Institute of Technology, 1200 E California Blvd, Pasadena, CA, USA mgharib@caltech.edu USA. rmurray@caltech.edu

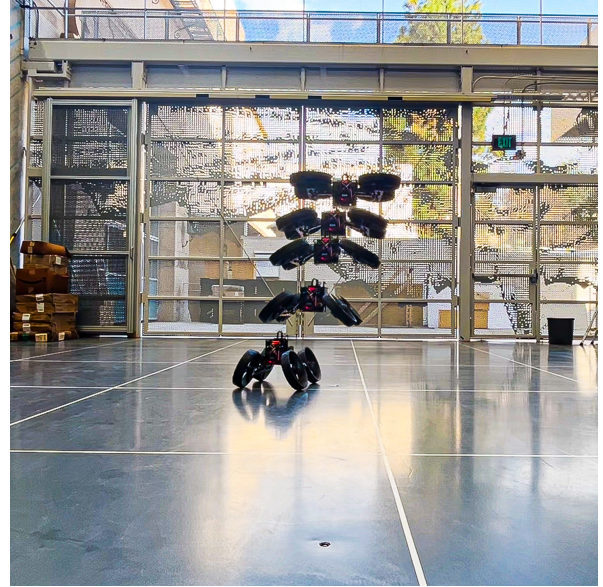∗ Corresponding author: imandralis@caltech.edu



Fig. 1. Snapshots of ATMO performing the Morpho-Transition maneuver using an end to end Reinforcement Learning control policy developed in this work. In this maneuver the robot starts in quadrotor flight mode and lands on its wheels by mid-air transformation.

enabling smooth ground aerial transition when ground clearance is required due to debris or otherwise rough terrain.

However, performing such a maneuver can be extremely challenging due to the need to stabilize mid-air while operating near actuator saturation. As ATMO changes its body posture to transform from quadrotor to ground rover mode, it passes through a configuration in which the thrust generated by the spinning rotor blades is not enough to counteract gravity. When attempting to hold position in this configuration, actuators are saturated making disturbance rejection very difficult. Thus, planning through this critical point to land with the largest possible tilt angle, combined with the need to take into account the ground contact upon landing, complicates the design of model based controllers.

This problem was first studied in [1] by using model predictive control with an adaptive cost function that varies with proximity to the ground and body shape to ensure successful landing on wheels. This method solved the problem of feasibility loss introduced at high tilt angles due to actuator saturation, but required extensive engineering work to achieve the reported results. The method was compared against off-the-shelf cascaded PID quadrotor controllers which were unable to solve the task due to the inherent differences between the dynamics of quadrotor flight and quadrotors in morpho-
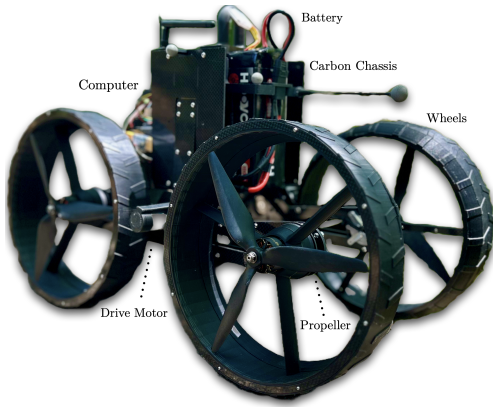
Fig. 2. A photo of ATMO, the Aerially Transforming Morphobot, in ground configuration with main components labeled. The robot can drive and fly by virtue of wheel-thruster actuators and a morphing mechanism to change between the two modes.

transition.

In this work we develop a competing RL method for achieving morpho-transition. The RL agent autonomously learns to stabilize ATMO, reject disturbances, and finds a feasible trajectory to solve the task from a single high level reward objective. Unlike the MPC method, where special care has to be taken when approaching the ground due to the hybrid switching from flight dynamics to ground dynamics, the RL method does not need to explicitly account for ground contact, enabling smooth control throughout the flight-drive operational envelope. The RL method operates at the RPM command level making this an end-to-end approach which bypasses the need for lower level body rate or attitude control loops which require extensive manual tuning, domain knowledge, and engineering effort.

### A. Related Work

Significant work in mid-air transformation control has been done for quadrotors which can move their thrusters in the plane normal to the direction of the thrust. This can be important for quadrotors, for example, which can use this adaptation to fit through narrow gaps, or to optimize flight agility and efficiency. In this space, a feedback linearizing LQR body rate controller is proposed in [2] with LQR gains adapted online to take into account the changes of center of gravity and inertia due to shape change. These are then fed into a control allocation matrix, which is a function of the thruster angles, that converts desired torque and collective thrust to thruster commands. This approach is beneficial because it fits directly into the classical cascaded PID loop used for quadrotor control. A competing approach uses an adaptive body rate controller that adapts online to changing actuator configuration [3].

Fewer works have investigated how rotorcraft with thrusters that tilt away from the vertical axis can be controlled. This type of system poses significant challenges to controllers due to the thrust not being aligned with the body $z$-axis resulting in out-of-plane thrust components

that are unaccounted by cascaded PID architectures designed for fixed frame quadrotors. In this domain, there have been developments such as passively morphing out-of-plane quadrotors which use mechanical springs to fold and pass through small gaps [4], [5]. However, the problem of continuous control of morphing with out-of-plane thrust, rather than hybrid switching, has only been considered using model predictive, or optimization-based control [6], [1]. These methods require extensive manual engineering as well as first principles modeling to produce reliable control performance.

An attractive alternative to the model-based techniques is Reinforcement Learning. Research in RL for flight control has mainly focused on learning how to fly quadrotors with the aim of pushing the envelope of performance and agility [7], [8]. Significant strides have been made in training RL controllers to directly output RPM commands, rather than control commands at a higher level of abstraction, removing the need for manually engineered lower level control loops [9], [10], [11], [12]. Reinforcement learning based landing and perching has been explored for quadrotors [13], [14], [15] but to the best of our knowledge, it has not been used for control of mid-air robotic transformation.

### B. Contribution

In this work we demonstrate successful transfer to hardware of an end-to-end RL controller and compare our results to an MPC controller for the challenging task of morpho-transition. We show that, using the RL controller, we can achieve controlled landings with tilt angles greater than $65°$, which is challenging to achieve with model-based controllers. We evaluate the performance of the RL method compared to the MPC method in [1] by examining the disturbance rejection characteristics under unknown observation delays, motor dynamics and recovery from partial actuator failure. Our contributions can be summarized as follows:

1) We develop a reward function that accurately encodes the problem of morpho-transition by taking advantage of the perfect information available in the simulator.
2) We propose a performance benchmark for the morpho-transition maneuver based on impact velocity and final distance to the goal after push disturbances of different magnitude and direction. We thoroughly characterize the controllers by exploiting massively parallel GPU simulation.
3) We extend the Isaac Lab simulation framework to include common quadrotor aerodynamics such as the thrust and moment coefficients, as well as motor dynamics and observation delays for the specialized tasks requiring aerial transformation. This code, along with the MPC method, is publicly available here.

## II. BACKGROUND

### A. Robot Description

ATMO [1] is an aerially transforming Morphobot capable of flying, driving, and transforming mid-air to land on its wheels. A photo of the robotic platform is shown in Figure
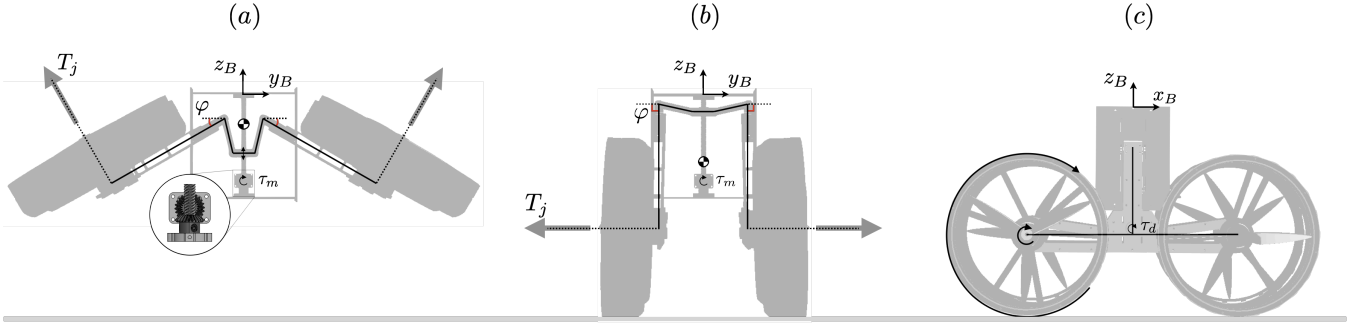
Fig. 3. (a) Front-view of ATMO in flight mode with thrust $T_j$ and the body frame depicted. (b) Front-view of ATMO in ground mode. (c) Side-view of ATMO showing the wheel rotation. The tilt angle $\varphi$ is depicted, as well as the motor torques required for morphing, $\tau_m$, and driving, $\tau_d$, respectively.

2, along with the main parts and components labeled. The design rationale and system specifications are described in [1]. We provide a brief summary below, for completeness.

ATMO can smoothly transition between quadcopter and ground rover mode by virtue of a closed loop kinematic chain, shown in Figure 4, which controls the sagittal plane tilt angle, $\varphi$, of four wheel-thruster pairs using a single brushed DC motor. The mechanism offers a high reduction allowing the robot posture to lock in place when un-actuated, and enables transformation with minimal actuation requirements. ATMO's mass of 5.5kg includes a high capacity battery pack for sustained operation, two belt-pulley systems on either side of the robot actuated by DC motors to enable differential drive steering as well as an onboard NVIDIA Jetson Orin Nano computer and CubeOrange PX4 enabled flight controller for sensor fusion and state estimation. All communication is achieved using ROS2 [16].

As ATMO morphs in the sagittal plane to transition from flight to drive modes the thrust produced by the spinning propellers is oriented away from the body-vertical axis resulting in lower overall thrust available to counteract gravity. At the critical tilt angle, $\varphi_c = 60°$,[1] the thruster actuators become fully saturated when trying to counteract gravity. In this configuration, the robot cannot hold its position and reject disturbances simultaneously. Furthermore, when flying in a morphed configuration, unequal values of thrust on either side of the H-frame rotor configuration may result in thrust components along the body $y_B$-axis, shown in Figure 3. This makes ATMO's dynamics fundamentally different to a fixed-frame quadrotor.

### B. Morpho-Transition Quadrotor Dynamics

Following [1], we modeled ATMO as 7 inertial components: the robot base, the left arm and right arm, and the four spinning propellers, as shown in Figure 3. The wheel joints are not considered for aerial maneuvers.

The robot generalized coordinates and velocities are defined as: $q = (p, \xi, \varphi)$ and $w = (v, \omega)$ where $p$ is the position of the robot center of mass, $\xi$ is the quaternion of

the robot base frame relative to the inertial world frame, $\varphi$ is the body tilt angle, $v$ is the velocity of the center of mass, and $\omega$ is the angular velocity of the robot relative to the world frame expressed in the body frame.

The control inputs are defined as $u = (u_{\text{aero}}, u_{\text{body}})$, where $u_{\text{aero}} \in [0, 1]^4$ denote the normalized propeller RPMs and $u_{\text{body}} \in [-1, 1]$ is the commanded tilting velocity. The thrust force and drag torque caused by the rotating propellers are assumed to be linear functions of the control inputs and act parallel to the rotor axes which are defined in the positive body $z$ axis: $T_j = c_T u_{\text{aero}}^j \hat{z}_j$ and $M_j = c_M c_T u_{\text{aero}}^j \hat{z}_j$. The thrust and moment coefficients $c_T, c_M$ are identified experimentally through thrust stand system identification experiments. $\hat{z}_j$ denotes the rotation axis of propeller $j$.

The body shape change occurs by spinning a central shaft with a DC motor encoder combination which results in linear motion of joint $A$, as shown in Figure 4. This in turn is translated to a rotation of point $E$, which corresponds to the intersection of the propeller axis with the link $CDE$. Thus a single actuator can transform the robot. This effect is modeled as a pure integrator for the tilt angle $\varphi$, reflecting the physical properties of the non-compliant morphing mechanism. The motor dynamics of the thrusters are modeled as first order linear systems with input $u_{\text{aero}}$ and output the normalized propeller RPMs $\Omega$.

Overall, ATMO's dynamic equations of motion can be written in the following form,

$$
\begin{cases}
M(q)\dot{w} + b(q, w) + g(q) = S(q)\Omega & (1) \\
\dot{\Omega} = T_m^{-1}(u_{\text{aero}} - \Omega) & (2) \\
\dot{\varphi} = \dot{\varphi}_{\max} u_{\text{body}}, & (3)
\end{cases}
$$

where $M(q)$ denotes the mass matrix, $b(q, w)$ encodes the coriolis terms, and $g(q)$ encodes the gravitational dynamics. The actuation matrix $S(q)$ maps the generalized control input into the generalized acceleration space. The thruster motor constant is $T_m$ and the maximum speed of the closed loop kinematic linkage is $\dot{\varphi}_{\max}$. The robot equations were derived using the open-source software package proNEu

---

[1]Computed as the point where the weight of the robot equals the maximum projection of the thrust in the vertical axis: $\cos \varphi_c = \frac{mg}{T_{\max}}$.

[17]. In state-space, the dynamics take the following form:

$$\dot{\boldsymbol{x}} = \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{u}), \tag{4}$$

$$\boldsymbol{x} = (\boldsymbol{p}, \boldsymbol{\xi}, \boldsymbol{v}, \boldsymbol{\omega}, \varphi, \boldsymbol{\Omega}) \in \mathbb{R}^{18}, \tag{5}$$

$$\boldsymbol{u} = (\boldsymbol{u}_{\text{aero}}, u_{\text{body}}) \in [0, 1]^5. \tag{6}$$

Second order aerodynamic effects like proximity effects, propeller flapping, or vehicle rotor dynamic couplings [18] are ignored.

### C. Model Predictive Controller

The model predictive controller for the morpho-transition maneuver is detailed in [1] and briefly summarized here. The controller is a trajectory tracking MPC that receives a desired reference $x_{\text{ref}}, u_{\text{ref}}$ and solves the following optimal control problem iteratively in a receding horizon control loop:

$$\begin{aligned} \underset{\boldsymbol{x}, \boldsymbol{u}}{\text{minimize}} \quad & \int_0^{t_f} L(\boldsymbol{x}, \boldsymbol{u}) dt \\ \text{subject to} \quad & \dot{\boldsymbol{x}} = \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{u}), \quad \forall t \in [0, t_f], \\ & \boldsymbol{u}_{\text{aero}} \in [0, 1]^4. \end{aligned}$$

The dynamics are discretized and enforced using multiple shooting [19] with $N$ collocation points and a look-ahead time $t_f$. A nonlinear program is solved using sequential quadratic programming in a real-time iteration scheme on-board the NVIDIA Jetson Orin Nano using the software packages CasADi and ACADOS [20], [21]. The controller runs at 150Hz. A look-ahead of $t_f = 2.0$ seconds and $N = 20$ collocation points are used. The cost function is given by,

$$L(\boldsymbol{x}, \boldsymbol{u}) = \alpha(\boldsymbol{x}) \underbrace{L_1(\boldsymbol{x}, \boldsymbol{u})}_{\text{flying}} + (1 - \alpha(\boldsymbol{x})) \underbrace{L_2(\boldsymbol{x}, \boldsymbol{u})}_{\text{transitioning}}, \tag{7}$$

and is a convex combination of quadratic costs $L_1, L_2$ specialized for flight and transition [1]. The cost is varied online according to the height and body shape through the scalar function $\alpha(\boldsymbol{x})$. The MPC controller assumes that the RPM commands $\boldsymbol{\Omega}$ are equal to the control commands i.e. $\boldsymbol{u}_{\text{aero}} = \boldsymbol{\Omega}$, thus operating in the reduced state space $\boldsymbol{x} = (\boldsymbol{p}, \boldsymbol{\xi}, \boldsymbol{v}, \boldsymbol{\omega})$. The tilt angle $\varphi$ is supplied from an external reference generator and tracked by a low level controller. For details on the implementation of the MPC algorithm, the reader is referred to [1].

### III. REINFORCEMENT LEARNING APPROACH

#### A. Simulation

To train a control policy for morpho-transition maneuvers, we implement morphing quadrotor dynamics as an extension of the state-of-the-art Isaac Lab simulation framework [22], [23]. The following section gives a brief overview of the dynamics implemented in the simulator. ATMO is simulated as a rigid body using a URDF (Unified Robot Description Format) representation with kinematic and inertial parameters generated from a CAD model. The simplified dynamics models for the closed loop kinematic linkage, aerodynamics
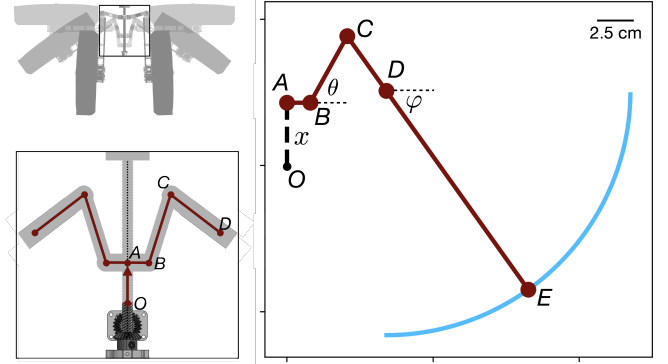


Fig. 4. A depiction of the closed loop kinematic linkage enabling symmetric morphing and transformation from air to ground mode.

models, and the motor dynamics are implemented in discrete time in the pre-physics step of the simulation, which performs the following operations:

$$\begin{cases} \boldsymbol{\Omega}[k+1] = \beta \boldsymbol{u}_{\text{aero}}[k] + (1 - \beta)\boldsymbol{\Omega}[k] & (8) \\ \boldsymbol{T}_j[k+1] = c_T \boldsymbol{\Omega}_j[k+1]\hat{\boldsymbol{z}}_j & (9) \\ \boldsymbol{M}_j[k+1] = c_M^j c_T \boldsymbol{\Omega}_j[k+1]\hat{\boldsymbol{z}}_j & (10) \\ \varphi_d[k+1] = \varphi[k] + \dot{\varphi}_{\text{max}} u_{\text{body}}[k] T_s. & (11) \end{cases}$$

The motor dynamics are discretized exactly, under the assumption of zero-order hold between simulation time-steps, leading to a first-order filter with constant $\beta = 1 - e^{-T_s/T_m}$. A nominal motor time constant of $T_m = 0.15$ is used, and the simulator time step is $T_s = 0.02$ seconds, leading to a filter constant $\beta = 0.12$. The filtered normalized motor RPMs are then used to determine the thrust and moment acting on the four propeller axes, $\boldsymbol{T}_j, \boldsymbol{M}_j$.

The desired tilt angle $\varphi_d$ is given to the simulator and tracked by two PD position controllers (one for each arm joint), with stiffness $k_p = 1 \times 10^{15}$ and damping $1 \times 10^5$. The actual tilt mechanism on the hardware, depicted in Figure 4, is a closed loop kinematic chain that is difficult to simulate [24], [25]. Instead, we simplify the actuation system to a "stiff" mechanism driven by a single velocity input. The large gains chosen for the PD controller make the actuator a pure velocity controller that ensures the mechanism is stiff - approximating its real world characteristics.

#### B. Task Definition & Domain Randomization

At the beginning of each episode ATMO spawns uniformly at a random $(x_0, y_0)$ position in a square of 2 meters side length around the origin, and an initial height, $z_0 \sim \mathcal{U}(1, 2)$ meters from the ground. The goal state is set as the origin. The initial orientation relative to the inertial frame is randomized by sampling a roll and pitch from $\mathcal{U}(-\frac{\pi}{6}, \frac{\pi}{6})$ and an initial yaw angle $\psi_0$ from $\mathcal{U}(0, 2\pi)$. Initial velocities, $\boldsymbol{v}_0$, and angular velocities, $\boldsymbol{\omega}_0$ are sampled uniformly from the distribution $\mathcal{U}(-0.1, 0.1)$. The initial tilt angle is sampled randomly as $\varphi_0 \sim \mathcal{U}(0, \pi/6)$. Each episode ends when the elapsed simulation time exceeds five seconds, when the robot
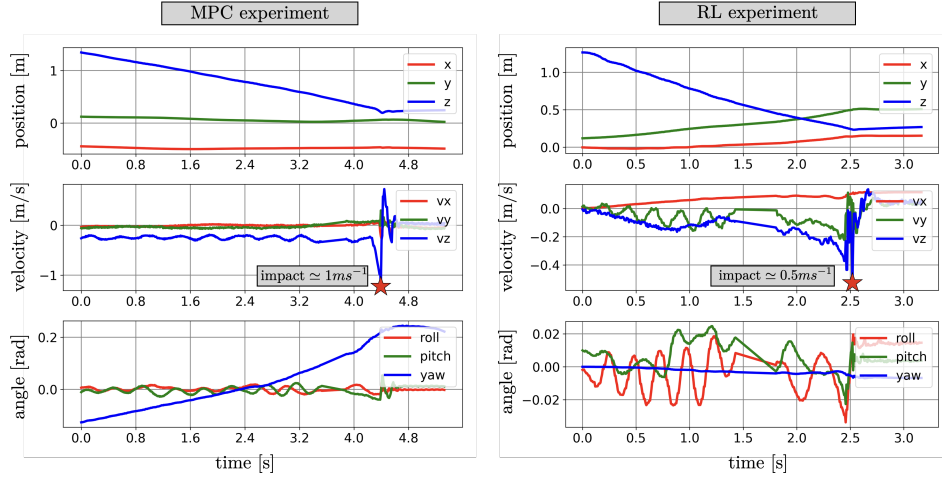
Fig. 5. Morpho-Transition maneuver performed on the hardware using the two controllers. The performance of the MPC controller is shown on the left, and the RL controller on the right. The key states during the maneuver are reported.

speed exceeds a maximum threshold of $2.0ms^{-1}$, or the $xy$ position of the robot deviates from the goal position by more than 1.5 meters.

We make the policy robust to disturbances by pushing the robot in random directions once per episode. At the beginning of each episode a push time, $t^*$, and push duration, $T^*$, are sampled from uniform distributions $t^* \sim \mathcal{U}(0, 0.8t_f), T^* \sim \mathcal{U}(0, 0.2)$, where $t_f$ is the episode length i.e. 5 seconds. The force and moment direction and intensity are sampled at the beginning of each episode and applied in the pre-physics step of the simulator for the given duration at the push time. The force and moment scales are set to $f^* = 0.20c_T$ and $\tau^* = 0.20c_Tc_M$. The maximum impulse and angular impulse imparted on the robot during training is thus: $J_{\text{disturbance}}^{\max} = 0.04c_T$ and $\Delta L_{\text{disturbance}}^{\max} = 0.04c_Tc_M$.

The motor time constant is randomized slightly around the nominal value $T_m = 0.15$ seconds at the beginning of each episode $T_m \sim \mathcal{U}(0.10, 0.20)$. The thrust and moment constants are assumed to have a $\pm 20\%$ uncertainty and are sampled as $c_T \sim \mathcal{U}(0.8c_T, 1.2c_T)$ and $c_M \sim \mathcal{U}(0.8c_M, 1.2c_M)$. Finally, the uncertainty of the tilt actuator dynamics is incorporated by randomizing the maximum tilt velocity parameter $\dot{\varphi}_{\max} \sim \mathcal{U}(0.8\frac{\pi}{8}, 1.2\frac{\pi}{8})$. We found that training without randomizing these parameters results in neither sim-to-sim nor sim-to-real transfer.

### C. Action and Observation Spaces

The actions are the outputs of the trained RL policy, $\pi(s)$. For the policy network we use three hidden layers of 128 units each with Exponential Linear Unit (ELU) activation functions between the layers. The network outputs are passed through a sigmoid layer to ensure the control actions are $u = \pi(s) \in [0, 1]^5$. The outputs of the policy network are interpreted as the 4 thruster control inputs $u_{\text{aero}}$ and the tilt angle velocity input $u_{\text{body}}$.

The inputs to the policy network are the policy observations,

$$s_\pi = (p, R, v, \omega, \varphi, u^-) \in \mathbb{R}^{19+5n}, \qquad (12)$$

where $u^-$ is an observation history of $n = 10$ previous time steps. For the rotation representation we found that it was necessary to use the full rotation matrix $R \in SO(3)$. Using a Quaternion representation of rotation did not result in sim-to-real transfer since quaternions double cover the space of rotations, meaning that the policy network must learn that $\xi$ and $-\xi$ represent the same rotation.

We use an asymmetric actor critic [26] scheme and give 14 additional privileged observations to the critic to improve value estimation. The critic observation space thus becomes,

$$s_Q = (s_\pi, f, \tau, t^*, T^*, t, J_c, \Omega^-) \in \mathbb{R}^{33+5n}, \qquad (13)$$

where $f, \tau$ are the disturbance force and moment, $t^*, T^*, t$ are the push time, push duration, and current time, and $J_c$ is the impulse acting on the robot due to ground contact forces. This is approximated as, $J_c = \int_t^{t+T_s} f_c dt \simeq [f_c(k+1) - f_c(t)] T_s$, where $f_c$ is the contact force between the ground and the robot, obtained from the Isaac Lab simulation environment. We gave the critic network access to the observed RPM values $\Omega^-$ which was shown in [11] to stabilize training. A small amount of uniform noise, with scales available in the Appendix, is added to the policy observations. No noise is added to the action history.

### D. Delays & Motor Dynamics

Finally, we found that for successful sim-to-real transfer, it was essential to add an observation delay of one simulation time step, or $20ms$, to the observations fed into the policy network during train time. This approximates the delays present due to inter-computer communication on the hardware. Without taking into account observation delays as well as motor dynamics the RL control policy could not be transferred to hardware.

### E. Reward Function

The reward function is designed to incentivize landing on the wheels in the vicinity of the goal state. The main term that produces this behavior is: $r_c^w = (c_w \land a)$. Here, $c_w$ is a
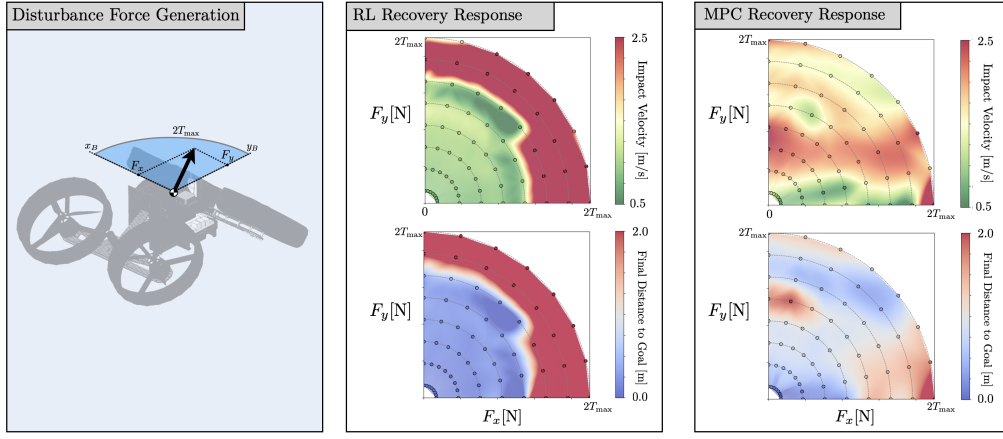
Fig. 6. Comparison between RL and MPC methods. The impact velocity and final distance to goal was recorded for 64 different push forces for each controller. The direction of the push force was varied in the $xy$ plane and the magnitude of the push force was varied between $0.60c_T$ and $8.0c_T$. The data points are shown in a scatter plot with an interpolated heat map overlayed.

Boolean variable active when any of the wheels make contact with ground and $a$ is a Boolean variable denoting that the robot is in the acceptance state, defined as the state where the $xy$ position of the robot is within 40 cm from the goal position.

To ensure that the robot has sufficiently rich reward information we added various reward shaping terms. First, we penalized the linear and angular velocities of the robot base frame, the action rate, and deviations from flat orientation. We rewarded the distance to the $xy$ coordinates of the goal position as well as descending at a constant rate of $0.5ms^{-1}$. We also rewarded high tilt angles as the robot approaches the ground. Finally, we took advantage of the ability to accurately measure contact times and forces in the simulator to penalize large impulse forces with the ground as well as to penalize undesirable thrust actions that occur while the robot is in contact with the ground. We impose a penalty of $-1$ unit for early termination. The full reward function expression is:

$$
\begin{aligned}
r(\boldsymbol{s}, \boldsymbol{u}) =\ & a_0||\boldsymbol{v}||^2 + a_1||\boldsymbol{\omega}||^2 + a_2|1 - q_a| \\
& + a_3||\boldsymbol{u} - \boldsymbol{u}^-||^2 + a_4 c_0||\boldsymbol{u}_{\mathrm{aero}}||^2 \\
& + a_5 J_c^2 + a_6(c_w \wedge a) \\
& + a_7\phi(d) + a_8\phi(z^2 + e_\varphi^2) \\
& + a_9\phi(e_{v_z}^2). \quad\quad (14)
\end{aligned}
$$

Here, $\phi(\cdot)$ is the function $\phi(x) = e^{-4x}$. $q_a$ denotes the axis of the quaternion, $c_0$ is a boolean variable that is equal to 1 when the first contact with the ground has occurred, $d$ is the distance to the $xy$ position of the goal state, and $e_\varphi = \varphi - \frac{\pi}{2}$. The coefficients $a_1 \dots a_9$ are given in the appendix.

*F. Training*

We use the Proximal Policy Optimization algorithm (PPO) [27] and exploit massively parallel training on GPU [24]. We use the RL-games implementation [22] with an initial learning rate of $\lambda = 1 \times 10^{-5}$. Our model was trained for 1000 policy updates after which the reward had converged to
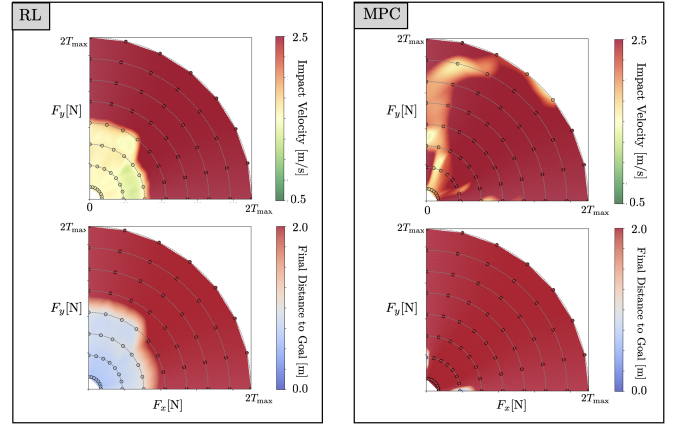


Fig. 7. Recovery characteristics under partial actuator failure. The thrust and moment coefficients of the thrusters are multiplied by $[0.8, 0.9, 0.85, 1.1]$ with order front right, back left, front left, back right.

a stable value. This took about 20 minutes on a GPU-enabled desktop computer.[2]

## IV. HARDWARE EXPERIMENTS

We performed different Morpho-Transition maneuvers in the CAST flight arena at Caltech, using the MPC and the RL controllers. Optitrack motion capture information was used to supply position and orientation estimates to an Extended Kalman Filter running onboard ATMO removing the need for GPS or vision based localization. ATMO was initialized on the ground and flown up to 1.25 meters using the PX4 quadrotor position controller [28]. At that point, the RL or MPC controllers were engaged and the morpho-transition landing maneuver was performed. The results of two representative experiments are shown in Figure 5, and the videos are available as supplementary materials.

Both controllers were successfully able to complete the task. The final tilt angle achieved by the MPC controller

[2]NVIDIA GeForce RTX 4070 with 16GB memory

was $\varphi = 60°$ and the RL controller achieved $\varphi = 65°$. The full maneuver can be seen in the Supplementary Video accompanying this paper. The RL controller achieves better impact velocity $0.5ms^{-1}$ as opposed to the $1.0ms^{-1}$ impact achieved by the MPC controller. It also manages to do this while operating past the actuator saturation tilt angle.

On the other hand, the RL controller exhibits larger oscillations in the roll angle evidencing slight instability during the descent. We believe this is due to inaccurate estimation of the system latency as well as the motor dynamics time constant. The effect of delays can be mitigated by inferencing the RL control policy directly on the flight controller microcontroller hardware rather than on the companion computer which adds latency due to the need to stream commands through the ROS2 network. Overall, the MPC method exhibits more stable angular dynamics but shows significant drift in yaw.

## V. RECOVERY FROM DISTURBANCES

To further characterize the differences between the two methods, we exploited the parallel simulation environment of Isaac Lab to examine the effect of the direction and magnitude of push disturbance on both controllers.

The MPC method is tasked with following a downward descending trajectory at 0.5 m/s and the RL agent performs actions according to the policy learned during the training process. The simulation was set up with $T_m = 0.15$ and $20ms$ observation delay. We systematically applied a push force, in the $xy$ plane of center-of-mass frame, at different angles and magnitudes in the first quadrant of the body $xy$ plane, as shown in the first panel of Figure 6. The angle was varied in increments of $\frac{\pi}{16}$ in the first quadrant of the body $xy$ plane. ATMO's symmetry does not require us to test push forces in other quadrants. The magnitude of the force was varied in increments of $0.925c_T$ from $0.6c_T$ to $8.0c_T$. Thus the maximum disturbance force that the robot recovery was tested for was of magnitude 2 times greater than the overall thrust capability of the robot ($T_{\max} = 4c_T$).

To characterize the ability of the controllers to recover from push disturbances we employ two metrics: 1) the impact velocity, and 2) the final distance from the goal. For each applied disturbance force the control policies were rolled out in simulation for 5 seconds and the impact velocity and final distance to the goal were recorded. Lower impact velocity and lower final distance to goal indicate that the controller was able to successfully recover from the push disturbance. Our results are shown in Figure 6. In the case of the RL method, the control policy performed very well up to between $5 - 6c_T$ disturbance force. In this region, there is a low impact velocity close to the commanded $0.5ms^{-1}$. The final distance to goal also stays close to zero. In other words, the RL agent is able to fully recover from push disturbances which are far beyond what it has been trained on. Once the threshold of $5 - 6c_T$ is surpassed however, the RL agent fails to land with an acceptable impact velocity or final distance to goal.

The MPC controller, on the other hand, exhibits more uniform performance. It is able to recover reasonably well from larger disturbances but it does not exhibit a uniform region of excellent disturbance recovery performance. Indeed, although the MPC controller often arrives at the goal position, it also often impacts the ground with unacceptable impact velocity.

## VI. RECOVERY FROM PARTIAL ACTUATOR FAILURE

Finally, we tested the performance of the two controllers under partial thruster failure. This can happen in experimental scenarios due to motor wear or low energy state. To simulate this scenario, we multiplied the nominal value of the thrust and moment coefficients by $[0.8, 0.9, 0.85, 1.1]$. This is a significant perturbation that the RL has not been trained on[3], and the MPC has no knowledge of. The results are shown in Figure 7. The RL controller is able to recover from considerable disturbance forces, even under the partial actuator loss. Naturally, the recovery region is much smaller than in the case where the thrusters are operating at nominal efficiency, Figure 6, but the RL agent is still able to recover from considerable applied impulses. In contrast, the MPC controller cannot recover at all, resulting in failure for almost all push tests considered. Note that it is possible to extend the MPC algorithm to take into account actuator failure as done in [29], which is expected to significantly improve performance. However, this would require explicit estimation of the actuator failure to trigger online changes in the optimal control problem being solved by the MPC.

## VII. CONCLUSIONS & FUTURE WORK

We have demonstrated successful transfer of an end-to-end deep RL controller to hardware for the task of quadrotor morpho-transition. To achieve this we developed a training approach that can generalize for the control of aerial robots in need of mid-air transformation control. The RL approach was compared to an MPC method. We found that end-to-end RL that has been trained on a well informed distribution of disturbance dynamics can reject small disturbances more reliably than the equivalent end-to-end MPC method. The MPC method performs worse at small disturbances but is able to recover from large disturbances where the RL controller fails. Moreover, the RL method is able to recover from partial actuator failure without explicit knowledge of the failure. The MPC method can in theory be extended to achieve this but would require explicit estimation of the actuator failure to trigger online changes in the optimization. This highlights the potential of RL to generate control policies that can generalize from partial sensor information to produce sophisticated control behaviors. In future work, we will seek to characterize the disturbance rejection capabilities of the two controllers in experiment, and investigate how an MPC approach that takes motor dynamics and observation delays into account compares to the end-to-end RL control policy.

---

[3]the RL was trained only on uniform changes in the thrust and moment coefficients.

## ACKNOWLEDGMENT

## REFERENCES

[1] I. Mandralis, R. Nemovi, A. Ramezani, R. M. Murray, and M. Gharib, "ATMO: an aerially transforming morphobot for dynamic ground-aerial transition," *Communications Engineering*, vol. 4, no. 1, Apr. 2025. [Online]. Available: http://dx.doi.org/10.1038/s44172-025-00413-6

[2] D. Falanga, K. Kleber, S. Mintchev, D. Floreano, and D. Scaramuzza, "The foldable drone: A morphing quadrotor that can squeeze and fly," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 209–216, 2019.

[3] A. Desbiez, F. Expert, M. Boyron, J. Diperi, S. Viollet, and F. Ruffier, "X-morf: A crash-separable quadrotor that morfs its x-geometry in flight," in *2017 Workshop on Research, Education and Development of Unmanned Aerial Systems (RED-UAS)*, 2017, pp. 222–227.

[4] N. Bucki and M. W. Mueller, "Design and control of a passively morphing quadcopter," in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 9116–9122.

[5] N. Bucki, J. Tang, and M. W. Mueller, "Design and control of a midair-reconfigurable quadcopter using unactuated hinges," 2021. [Online]. Available: https://arxiv.org/abs/2103.16632

[6] I. Mandralis, E. Sihite, A. Ramezani, and M. Gharib, "Minimum time trajectory generation for bounding flight: Combining posture control and thrust vectoring," in *2023 European Control Conference (ECC)*, 2023, pp. 1–7.

[7] J. Hwangbo, I. Sa, R. Siegwart, and M. Hutter, "Control of a quadrotor with reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, pp. 2096–2103, 2017.

[8] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza, "Champion-level drone racing using deep reinforcement learning," *Nature*, vol. 620, no. 7976, p. 982–987, Aug. 2023. [Online]. Available: http://dx.doi.org/10.1038/s41586-023-06419-4

[9] A. Molchanov, T. Chen, W. Hönig, J. A. Preiss, N. Ayanian, and G. S. Sukhatme, "Sim-to-(multi)-real: Transfer of low-level robust control policies to multiple quadrotors," 2019. [Online]. Available: https://arxiv.org/abs/1903.04628

[10] 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)*, 2022, pp. 10 504–10 510.

[11] J. Eschmann, D. Albani, and G. Loianno, "Learning to fly in seconds," 2023. [Online]. Available: https://arxiv.org/abs/2311.13081

[12] D. Zhang, A. Loquercio, J. Tang, T.-H. Wang, J. Malik, and M. W. Mueller, "A learning-based quadcopter controller with extreme adaptation," 2024. [Online]. Available: https://arxiv.org/abs/2409.12949

[13] R. Polvara, M. Patacchiola, S. Sharma, J. Wan, A. Manning, R. Sutton, and A. Cangelosi, "Autonomous quadrotor landing using deep reinforcement learning," 2017. [Online]. Available: https://arxiv.org/abs/1709.03339

[14] J. E. Kooi and R. Babuška, "Inclined quadrotor landing using deep reinforcement learning," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 2361–2368.

[15] R. Peter, L. Ratnabala, D. Aschu, A. Fedoseev, and D. Tsetserukou, "Lander.ai: Adaptive landing behavior agent for expertise in 3d dynamic platform landings," 2024. [Online]. Available: https://arxiv.org/abs/2403.06572

[16] S. Macenski, T. Foote, B. Gerkey, C. Lalancette, and W. Woodall, "Robot operating system 2: Design, architecture, and uses in the wild," *Science Robotics*, vol. 7, no. 66, May 2022. [Online]. Available: http://dx.doi.org/10.1126/scirobotics.abm6074

[17] M. Hutter, C. Gehring, and R. Siegwart, "proNEu: Derivation of analytical kinematics and dynamics," Autonomous Systems Lab, ETHZ, Tech. Rep., 2011.

[18] P.-J. Bristeau, P. Martin, E. Salaun, and N. Petit, "The role of propeller aerodynamics in the model of a quadrotor UAV," in *2009 European Control Conference (ECC)*. IEEE, Aug. 2009. [Online]. Available: http://dx.doi.org/10.23919/ECC.2009.7074482

[19] H. Bock and K. Plitt, "A multiple shooting algorithm for direct solution of optimal control problems *," *IFAC Proceedings Volumes*, vol. 17, no. 2, p. 1603–1608, Jul. 1984. [Online]. Available: http://dx.doi.org/10.1016/S1474-6670(17)61205-9

[20] R. Verschueren, G. Frison, D. Kouzoupis, J. Frey, N. van Duijkeren, A. Zanelli, B. Novoselnik, T. Albin, R. Quirynen, and M. Diehl, "acados - a modular open-source framework for fast embedded optimal control," *Mathematical Programming Computation*, vol. 14, pp. 147 – 183, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:204960656

[21] J. A. E. Andersson, J. Gillis, G. Horn, J. B. Rawlings, and M. Diehl, "CasADi: a software framework for nonlinear optimization and optimal control," *Mathematical Programming Computation*, vol. 11, no. 1, p. 1–36, Jul. 2018. [Online]. Available: http://dx.doi.org/10.1007/s12532-018-0139-4

[22] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, and G. State, "Isaac gym: High performance gpu-based physics simulation for robot learning," 2021. [Online]. Available: https://arxiv.org/abs/2108.10470

[23] M. Mittal, C. Yu, Q. Yu, J. Liu, N. Rudin, D. Hoeller, J. L. Yuan, R. Singh, Y. Guo, H. Mazhar, A. Mandlekar, B. Babich, G. State, M. Hutter, and A. Garg, "Orbit: A unified simulation framework for interactive robot learning environments," *IEEE Robotics and Automation Letters*, vol. 8, no. 6, pp. 3740–3747, 2023.

[24] N. Rudin, H. Kolvenbach, V. Tsounis, and M. Hutter, "Cat-like jumping and landing of legged robots in low gravity using deep reinforcement learning," *IEEE Transactions on Robotics*, vol. 38, no. 1, pp. 317–328, 2022.

[25] T. El-Agroudi, F. G. Maurer, J. A. Olsen, and K. Alexis, "In-flight attitude control of a quadruped using deep reinforcement learning," in *8th Annual Conference on Robot Learning*, 2024. [Online]. Available: https://openreview.net/forum?id=67tTQeO4HQ

[26] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, "Asymmetric actor critic for image-based robot learning," 2017. [Online]. Available: https://arxiv.org/abs/1710.06542

[27] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," 2017. [Online]. Available: https://arxiv.org/abs/1707.06347

[28] L. Meier, D. Honegger, and M. Pollefeys, "PX4: A node-based multithreaded open source robotics framework for deeply embedded platforms," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 6235–6240.

[29] F. Nan, S. Sun, P. Foehn, and D. Scaramuzza, "Nonlinear MPC for quadrotor fault-tolerant control," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 5047–5054, 2022.

## VIII. APPENDIX

### A. Reward function coefficients

| Coefficient | Value | Description |
|---|---|---|
| $a_0$ | $-0.10\,T_s$ | Velocity Penalty |
| $a_1$ | $-0.30\,T_s$ | Angular Velocity Penalty |
| $a_2$ | $-0.10\,T_s$ | Orientation penalty |
| $a_3$ | $-0.07\,T_s$ | Control Rate Penalty |
| $a_4$ | $-0.13\,T_s$ | Ground Thrust Penalty |
| $a_5$ | $-1.0$ | Contact Impulse Penalty |
| $a_6$ | $0.40$ | Contact in Acceptance Reward |
| $a_7$ | $0.30\,T_s$ | Distance to Goal Reward |
| $a_8$ | $0.40\,T_s$ | Morphing reward |
| $a_9$ | $0.30\,T_s$ | Descending reward |

### B. Noise scales

We added uniform noise of magnitude 0.005 to the position $p$ and rotation matrix $R$, 0.035 to the velocity $v$ and angular velocity $\omega$, and 0.018 to the tilt angle $\varphi$.