

# Predictive RL for Continuous Trajectory Tracking

Georg Schäfer<sup>1,2</sup> Jakob Rehr<sup>1</sup> Stefan Huber<sup>1</sup> Simon Hirlaender<sup>2</sup>

<sup>1</sup> Josef Ressel Centre for Intelligent and Secure Industrial Automation,  
Salzburg University of Applied Sciences, Salzburg, Austria

<sup>2</sup> Paris Lodron University of Salzburg, Salzburg, Austria

## 1. The Industrial Problem

- » Industrial cyber-physical systems require **smooth, highly precise control**.
- » Standard deep reinforcement learning agents are often **purely reactive**, leading to **lag and overshoots** when tracking targets.
- » Controllers like Model Predictive Control (MPC) **anticipate future states** but demand high online computational resources.
- » **Our Goal:** Combine the **foresight of MPC** with the **fast inference** of deep reinforcement learning.

## 2. Methodology & Experimental Setup

- » **Proximal Policy Optimization** applied to a **1-DOF Quanser Aero 2** testbed, representing non-linear industrial dynamics.
- » Reactive Baseline State:  $s_t = (\theta_t, \dot{\theta}_t, r_t)$
- » **Predictive State:**  $s_t = (\theta_t, \dot{\theta}_t, r_t, \dot{r}_t, r_{t+1:\Delta t}, \dots, r_{t+N:\Delta t})$
- » with  $\theta$ : pitch,  $\dot{\theta}$ : vel.,  $r$ : target,  $\dot{r}$ : target vel.,  $N$ : # future targets,  $\Delta t$ : step interval.
- » **Evaluation:** 10 simulation runs per configuration evaluating absolute mean deviation on a fixed trajectory, followed by **zero-shot sim-to-real transfer** to the physical hardware.

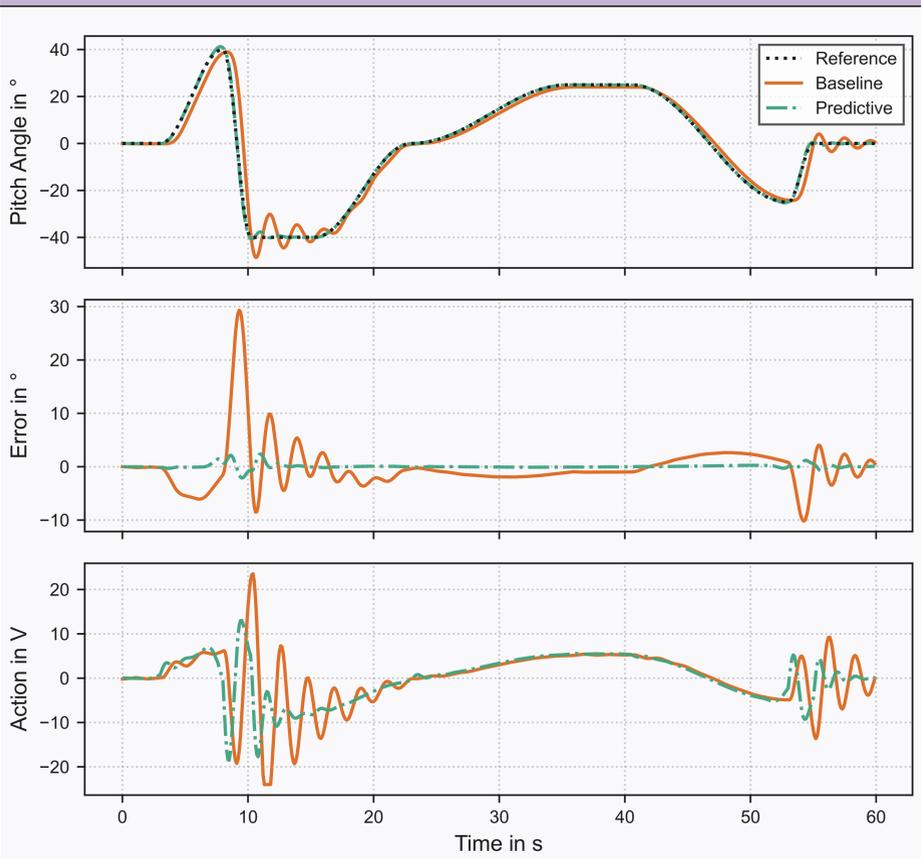
## 3. Results: Mean Absolute Deviation (°)

$\dot{r}$	$N$	$\Delta t$	Worst	Best	Mean	Std	Hardware
No	0	–	3.05	2.51	2.73	0.18	2.34
No	1	0.1	1.88	1.49	1.70	0.13	2.09
No	1	1.0	0.95	0.78	0.85	0.06	<b>1.11</b>
No	4	0.5	1.02	0.49	0.70	0.14	1.21
Yes	0	–	0.90	0.41	0.52	0.15	1.23
Yes	1	0.1	0.74	0.34	0.48	0.10	1.43
Yes	1	1.0	0.50	0.22	0.34	0.08	1.14
Yes	4	0.5	0.53	0.19	<b>0.31</b>	0.10	<b>1.11</b>

## 4. Conclusion & Outlook

- » Adding a **1.0s prediction horizon** cut real-world tracking error in half (from **2.34° to 1.11°**).
- » A **single future target** performed equally well as the most **dense predictive state**. This simplicity acts as a **regularizer**, achieving peak performance without needing velocity information.
- » **Outlook:** Training on the physical hardware and benchmarking its inference speed and energy efficiency directly against MPC.

### Simulation (Baseline vs. $\dot{r}$ , $N = 4$ , $\Delta t = 0.5s$ )



### Physical System (Baseline vs. $\dot{r}$ , $N = 4$ , $\Delta t = 0.5s$ )

