

Model Predictive Control and Reinforcement Learning

– Introduction –

Joschka Boedecker and Moritz Diehl

University Freiburg

Aim of this course



Understanding the main concepts of model predictive control (MPC) and reinforcement learning (RL) and their similarities and differences.

Applying the methods to practical optimal control problems in hands-on exercises and project work.

Agenda



Week 1 at Faculty of Engineering (HS 0-26 - Buil. 101)			
	Wed 4-Oct	Thu 5-Oct	Fri 6-Oct
09:00-10:30	Welcome + <i>Lecture 1</i> Dynamic Systems and Simulation (MD)	<i>Lecture 3</i> Dynamic Programming and Optimal Control (MD)	<i>Lecture 5</i> Monte Carlo RL, Temporal Difference and Q-Learning (JB)
10:30-11:00	Coffee Break	Coffee Break	Coffee Break
11:00-12:20	<i>Exercise 1</i> Dynamic Systems and Simulation	<i>Exercise 3</i> Dynamic Programming and Optimal Control	<i>Exercise 5</i> Q-learning
12:20-12:30		<i>Project Brainstorming</i>	<i>Project Brainstorming</i>
12:30-14:00	Lunch	Lunch	Lunch
14:00-15:30	<i>Lecture 2</i> Numerical Optimization (MD)	<i>Lecture 4</i> Deep Learning (JB)	<i>Lecture 6</i> RL with Function Approximation (JB)
15:30-16:00	Coffee Break	Coffee Break	Coffee Break
16:00-17:30	<i>Exercise 2</i> Numerical Optimization	<i>Exercise 4</i> PyTorch	<i>Exercise 6</i> DQN
		Break	
19:00-		Social Gathering	

Week 2 at Historical University (HS 1015 - Buil. KG I)				
Mon 9-Oct	Tue 10-Oct	Wed 11-Oct	Thu 12-Oct	Fri 13-Oct
Motivating Discussion + <i>Lecture 7</i> NMPC (MD)	<i>Lecture 9</i> Policy gradient and Actor-Critic methods (MD, JB)	<i>Lecture 11</i> Introduction MPC and RL Framework (SG)	<i>Lecture 13</i> When to use RL in MPC? (SG)	Project work
Coffee Break	Coffee Break	Coffee Break	Coffee Break	Coffee Break
<i>Exercise 7</i> NMPC (Acados)	<i>Exercise 9</i> Actor-Critic	<i>Exercise 11</i> MPCRL, tba	<i>Exercise 12</i> MPCRL, tba	Project work / Project presentations
<i>Project Brainstorming</i>	<i>Project Brainstorming</i>			
Lunch	Lunch	Lunch	Lunch	Lunch
<i>Lecture 8</i> Transformers (JB)	<i>Lecture 10</i> MPPI (HB) / Model-based RL (JB)	<i>Lecture 12</i> Safety and Stability in MPCRL (SG)	Project work	Project presentations
Coffee Break	Coffee Break	Coffee Break	Coffee Break	Coffee Break
<i>Exercise 8</i> Transformer in practice (forecasting)	<i>Exercise 10</i> MPPI	Project brainstorming and kick-off	Project work	
Break	From 17 to 18: Project Announcements	Break		
Aperitif at Waldsee		Workshop Dinner		

Team



Moritz Diehl
(Lecturer MPC)



Joschka Boedecker
(Lecturer RL)



Andrea Ghezzi
(Tutor)



Jasper Hoffmann
(Tutor)



Sebastien Gros
(NTNU Trondheim)



Yuan Zhang
(Tutor)

Discussion



- ▶ What do you know about Model Predictive Control?
- ▶ What are characteristics of Reinforcement Learning?
- ▶ What are differences to Supervised Learning?



Characteristics of MPC & Reinforcement Learning



- ▶ Both are frameworks to solve sequential decision making problems
- ▶ Both automatically design controllers based on desired outcomes (reward / stage cost, constraints)
- ▶ Closed-loop system visits different regions of the state space than uncontrolled system

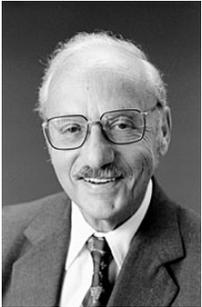
MPC

- ▶ System identification precedes control implementation, model fixed during execution
- ▶ Typically convex stage costs
- ▶ Constraints imposed explicitly
- ▶ Online optimization over prediction horizon, expensive
- ▶ Usually combined with state estimator

RL

- ▶ Controller directly learned from data, trial-and-error, exploration and exploitation trade off
- ▶ Both shaped/concave and 0-1/sparse rewards
- ▶ Constraints are imposed via penalties
- ▶ Typically parametrized controller, cheap online execution
- ▶ Usually, history included in definition of the state

MPC & Reinforcement Learning have a long history



Linear Programming ... MPC

- ▶ Linear Programming (LP) developed by **G. Dantzig** in 1947
- ▶ was extended to Quadratic Programming (QP), Nonlinear Programming (NLP), Integer Programming (IP), ... in field of **mathematical optimisation**
- ▶ Online solution of LP, QP, NLP, IP used for many planning problems and increasingly for industrial control problems in form of MPC



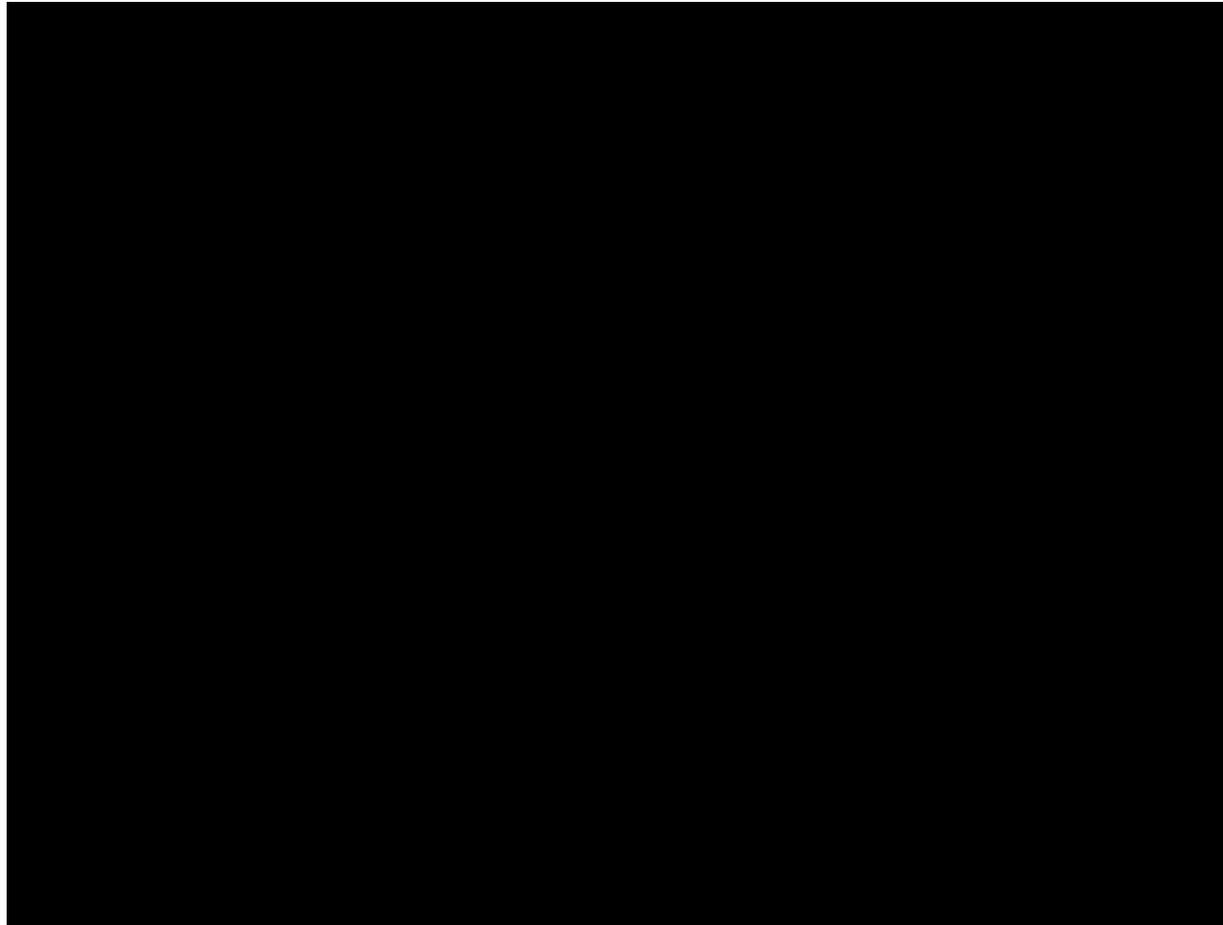
Dynamic Programming ... RL

- ▶ Dynamic Programming developed by **R. Bellman** in 1950s
- ▶ was extended to approximate dynamic programming, Monte Carlo Tree Search, Q-learning, policy search ... in field of **machine learning**
- ▶ Reinforcement Learning techniques are increasingly applied to solve difficult planning and decision making problems with impressive results e.g. in computer games and robotics.

Some Applications of RL



Learning to Play Atari Games from Pixel Input



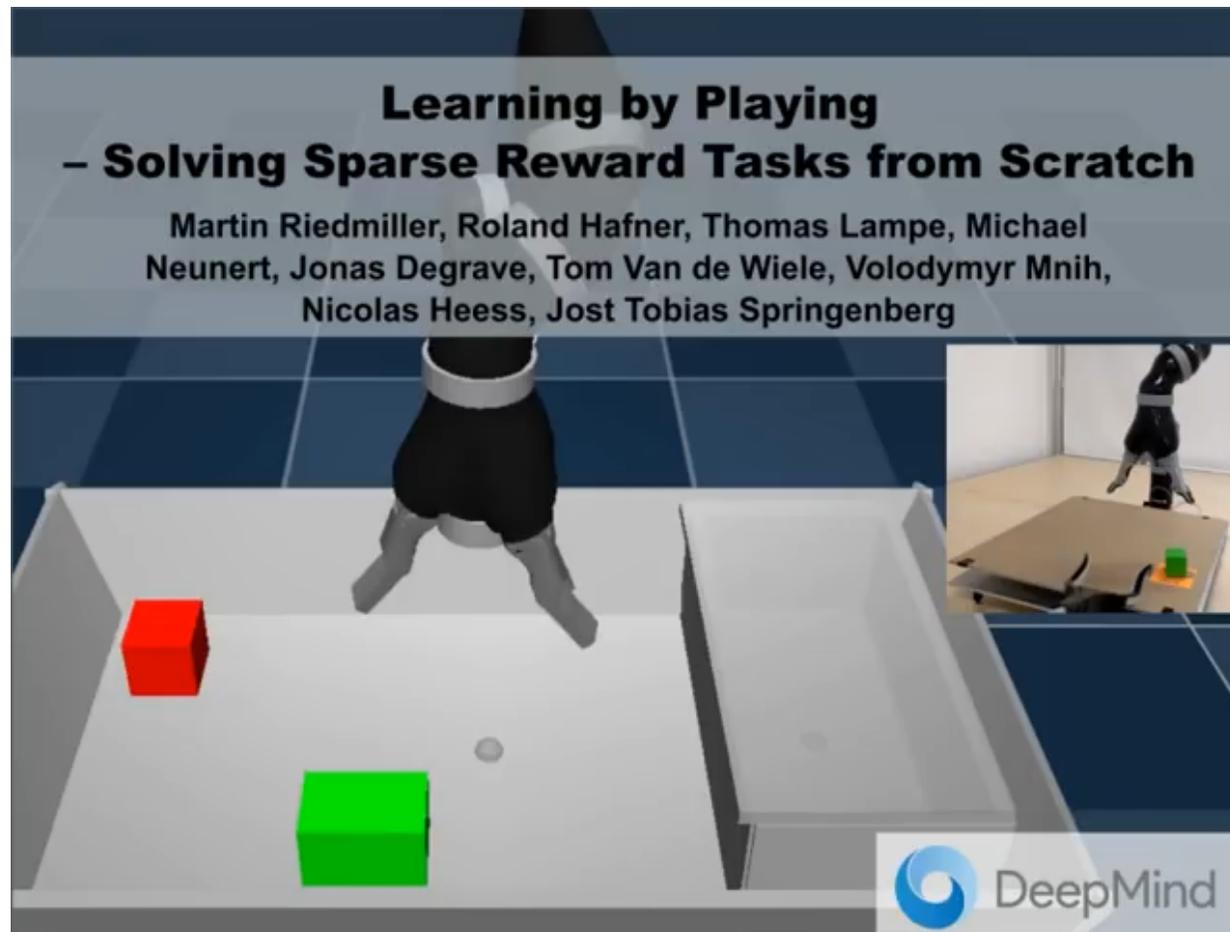
[Mnih et al., 2015]

Learning to Play the Game of Go Better Than Any Human



[Silver et al., 2016]

Learning Difficult Robot Manipulation Tasks from Scratch



[Riedmiller et al., 2018]



Approximate Real-Time Optimal Control Based on Sparse Gaussian Process Models

Joschka Boedecker, Jost Tobias Springenberg, Jan Wülfing, Martin Riedmiller

University of Freiburg

Department of Computer Science
Machine Learning Lab
Prof. Dr. Martin Riedmiller



**UNI
FREIBURG**



Deep Inverse Q-learning with Constraints



UNI
FREIBURG



Gabriel Kalweit^{*,1}, Maria Huegle^{*,1}, Moritz Werling² and
Joschka Boedecker¹

¹ University of Freiburg, Germany and ² BMWGroup, Germany



Some Applications of MPC



Time-Optimal Point-To-Point Motions



Fast oscillating systems (cranes, plotters, wafer steppers, ...)

Control aims:

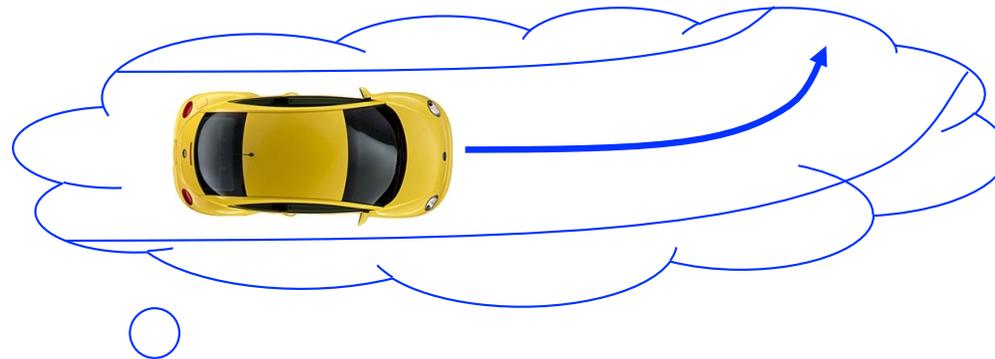
- reach end point as fast as possible
- do not violate constraints
- no residual vibrations

Idea: formulate as embedded optimization problem
in form of Model Predictive Control (MPC)



Model Predictive Control (MPC)

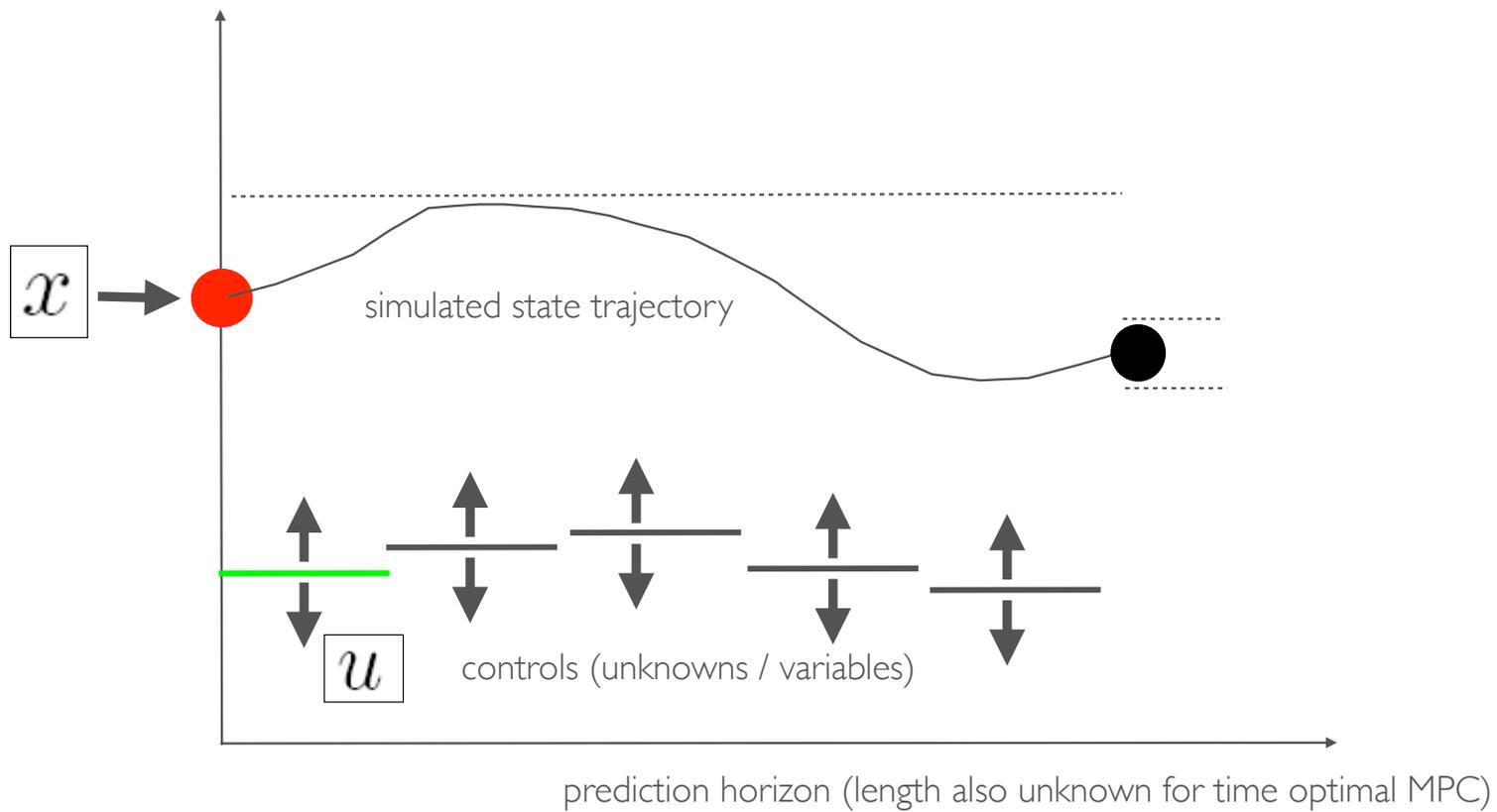
Always look a bit into the future



Example: driver predicts and optimizes, and therefore slows down before a curve

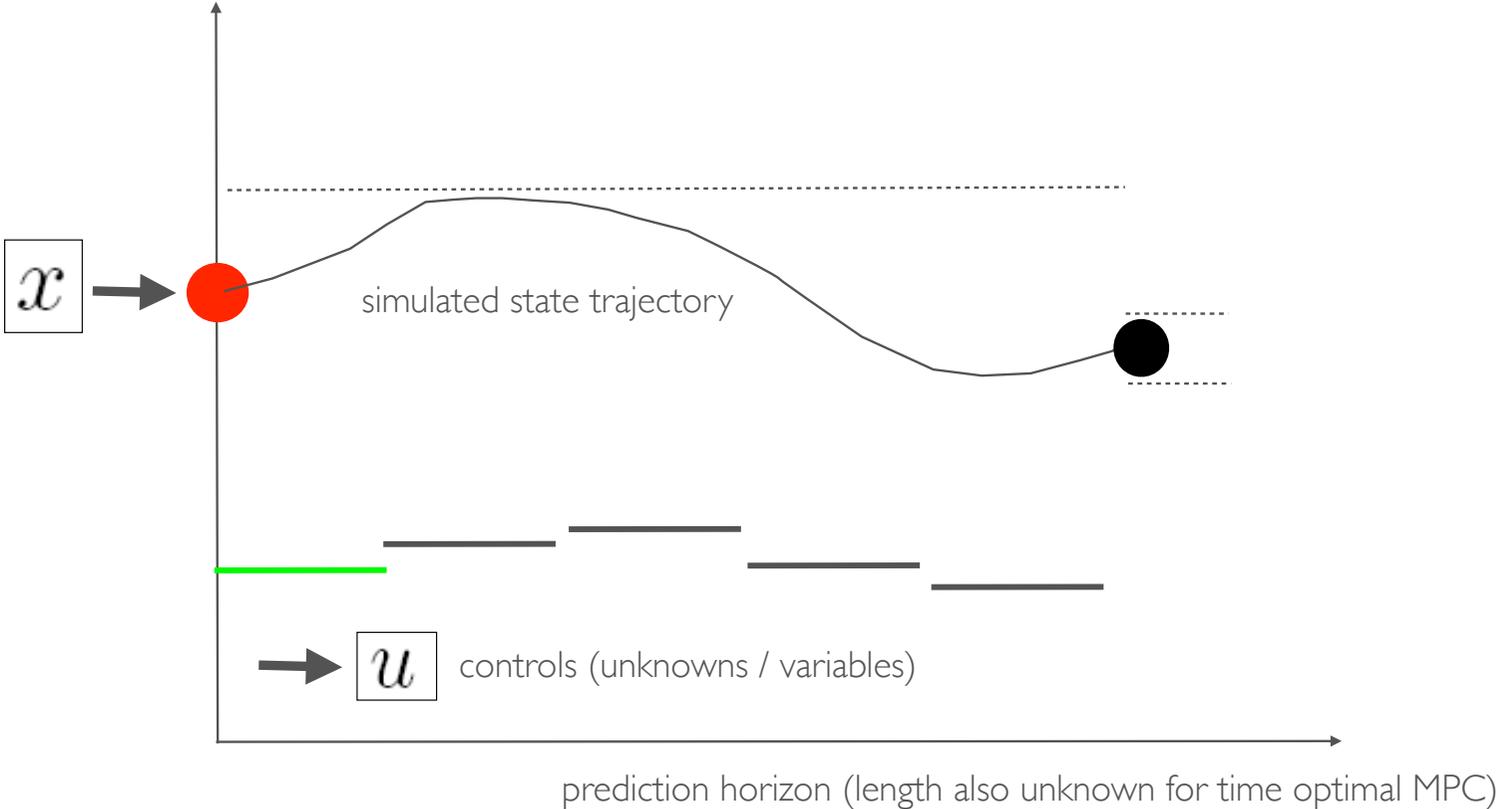
Optimal Control Problem in MPC

For given system state \mathbf{x} , which controls \mathbf{u} lead to the best objective value without violation of constraints?

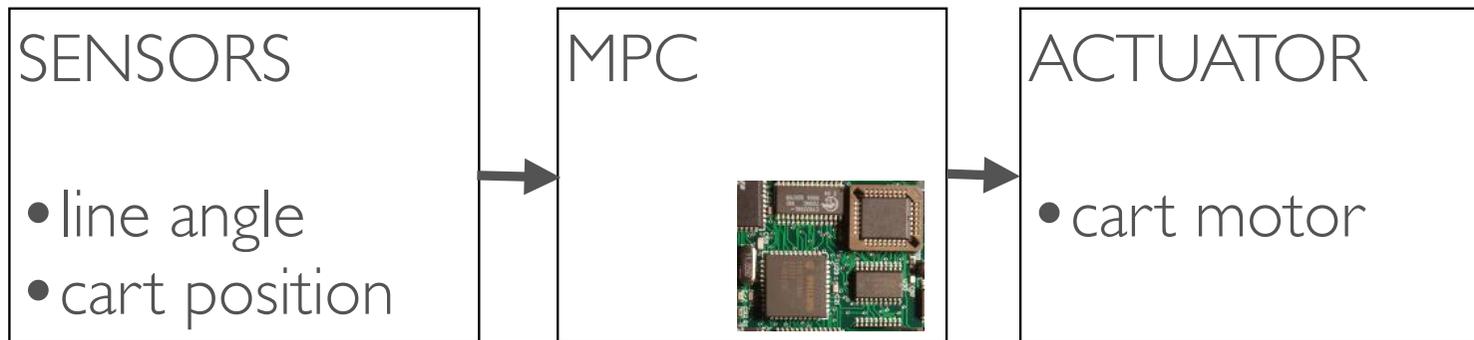


Optimal Control Problem in MPC

For given system state x , which controls u lead to the best objective value without violation of constraints ?



Time Optimal MPC of a Crane



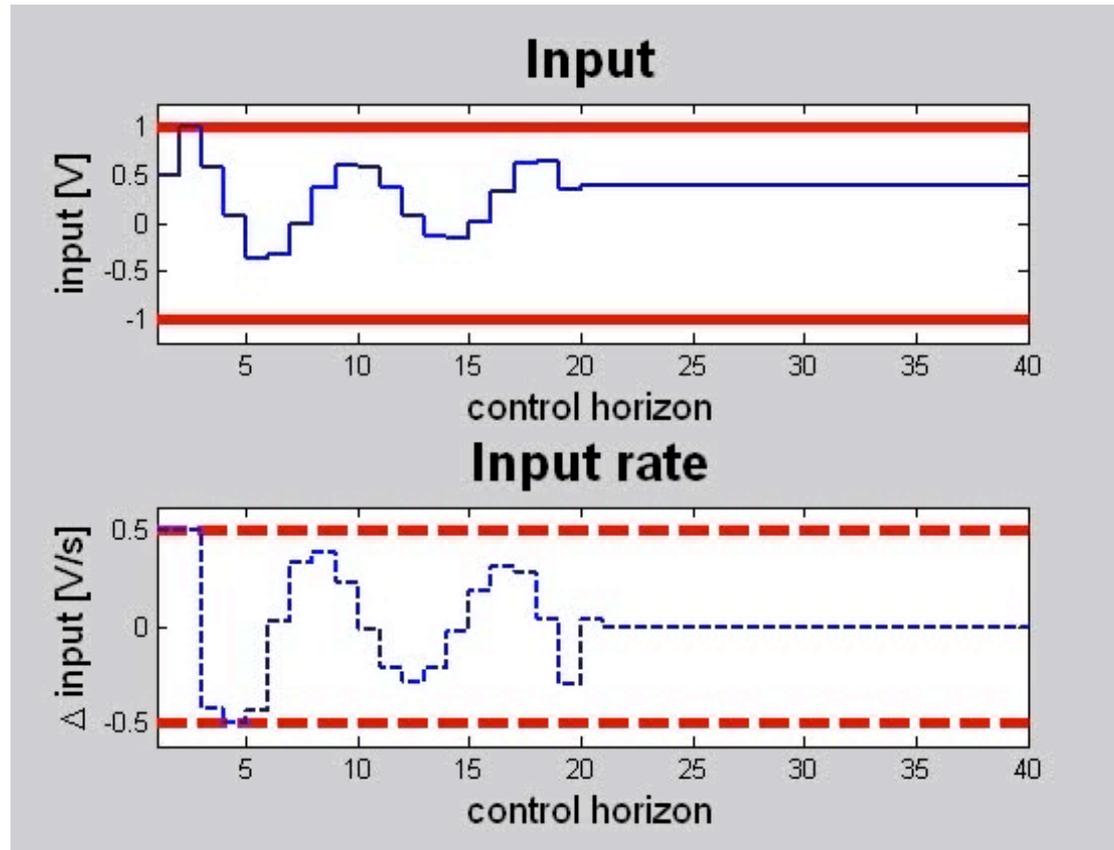
Hardware: xPC Target. Software: qpOASES [Ferreau, D., Bock, 2008]

Time Optimal MPC of a Crane

Univ. Leuven [Vandenbrouck, Swevers, D.]



Optimal Solutions in qpOASES Varying in Time



Time Optimal MPC in Industry: 25cm step, 100nm accuracy

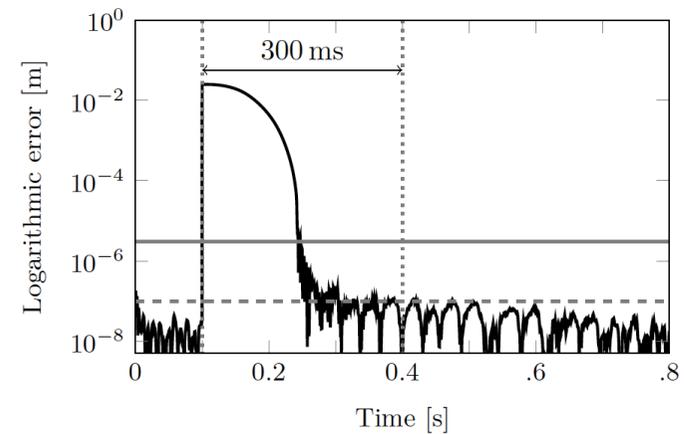
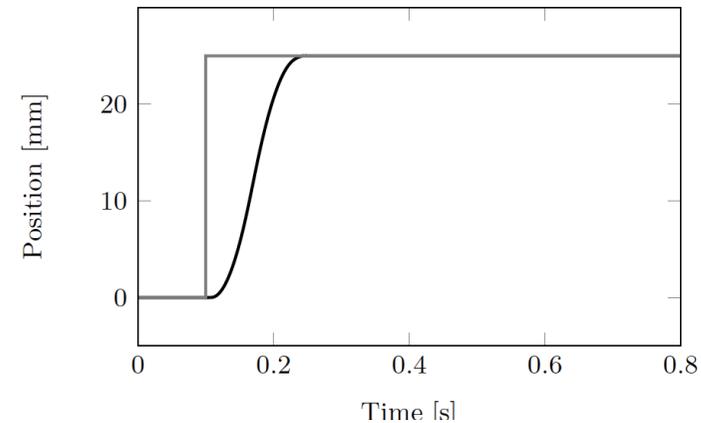


TOMPC at 250 Hz (+PID with 12 kHz)

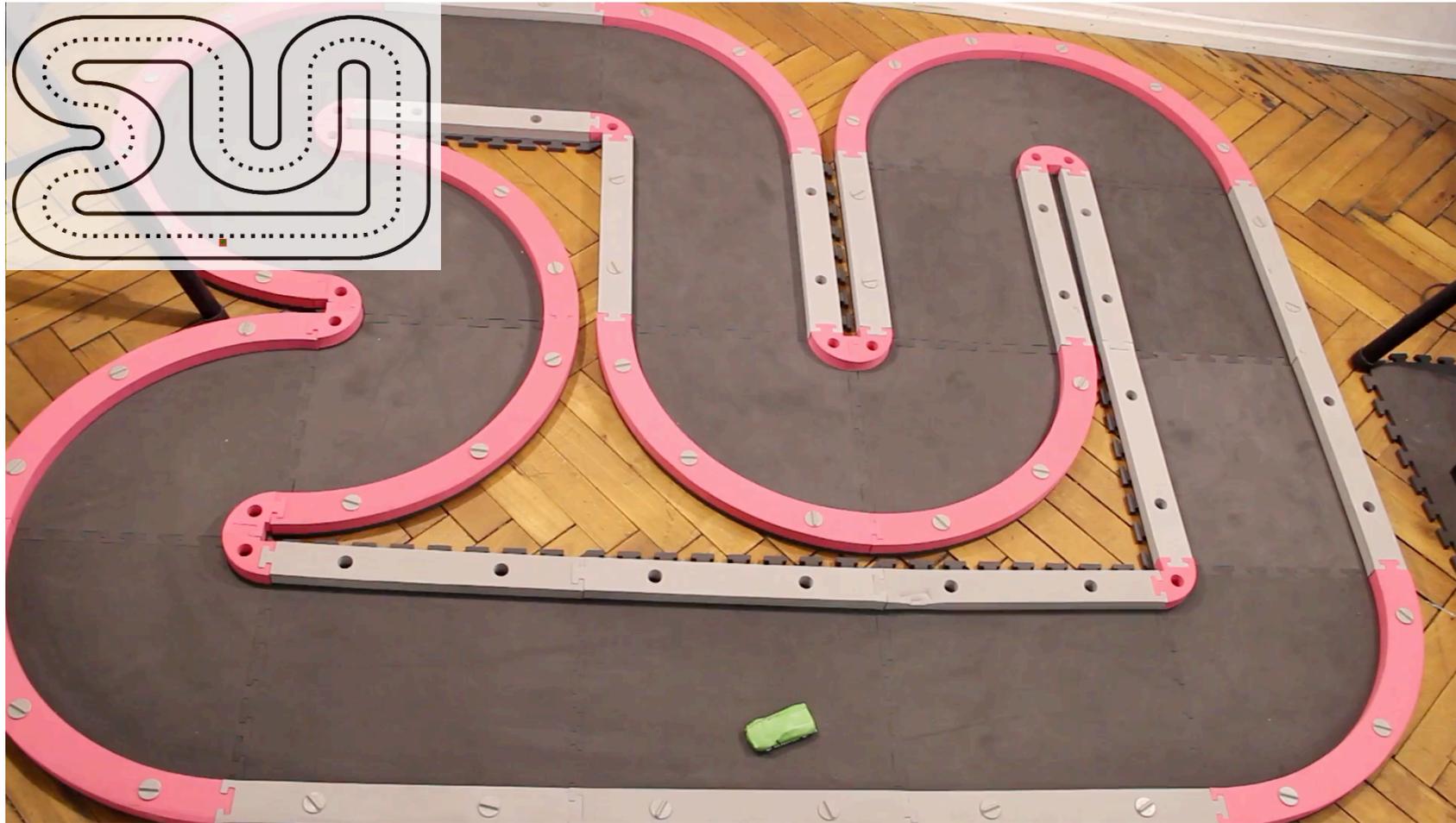
Lieboud's results after 1 week at ETEL:

- 25 cm step in 300 ms
- 100 nm accuracy

equivalent to: „fly 2,5 km with MACH15,
stop with 1 mm position accuracy“



Model Predictive Control of the Freiburg Race Cars



acados coupled into ROS, optimization every 10ms

[Kloeser et al., submitted]

Safe Motion Planning at Bosch via the Convex Inner Approximation Method [Schöls et al, 2020]

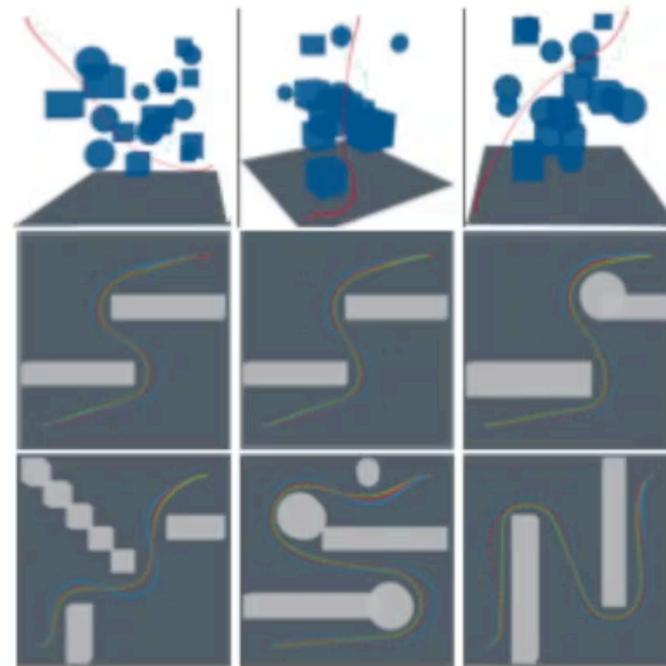
An NMPC Approach using Convex Inner Approximations for Online Motion Planning with Guaranteed Collision Freedom

Tobias Schoels^{1,2}, Luigi Palmieri², Kai O. Arras², and Moritz Diehl¹

Abstract—Even though mobile robots have been around for decades, trajectory optimization and continuous time collision avoidance remains subject of active research. Existing methods trade off between path quality, computational complexity, and kinodynamic feasibility. This work approaches the problem using a model predictive control (MPC) framework, that is based on a novel convex inner approximation of the collision avoidance constraint. The proposed Convex Inner Approximation (CIAO) method finds a dynamically feasible and collision free trajectory in few iterations, typically one, and preserves feasibility during further iterations. CIAO scales to high-dimensional systems, is computationally efficient, and guarantees both kinodynamic feasibility and continuous-time collision avoidance. Our experimental evaluation shows that the approach outperforms state of the art baselines in terms of planning efficiency and path quality. Furthermore real-world experiments show its capability of unifying trajectory optimization and tracking for safe motion planning in dynamic environments.

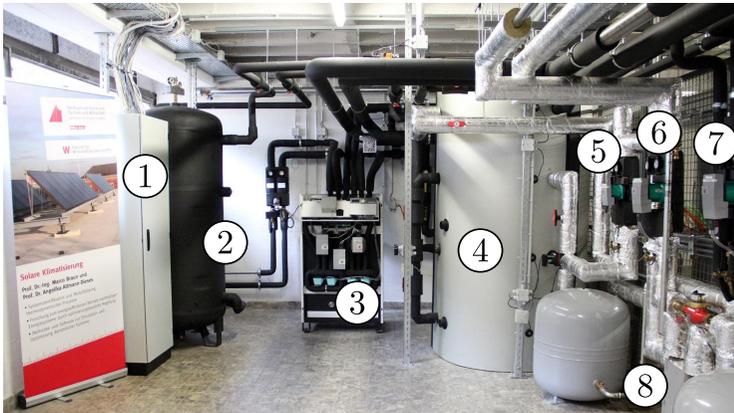
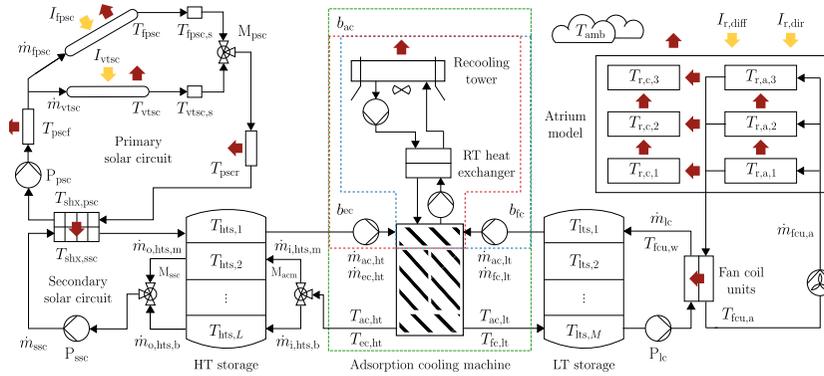
I. INTRODUCTION

Several existing mobile robotics applications (e.g. intralogistic and service robotics) require robots to operate in dynamic environments among other agents, such as humans or other autonomous systems. In these scenarios the reactive



Nonlinear Mixed-Integer Control of a Solar Adsorptive Cooling Machine

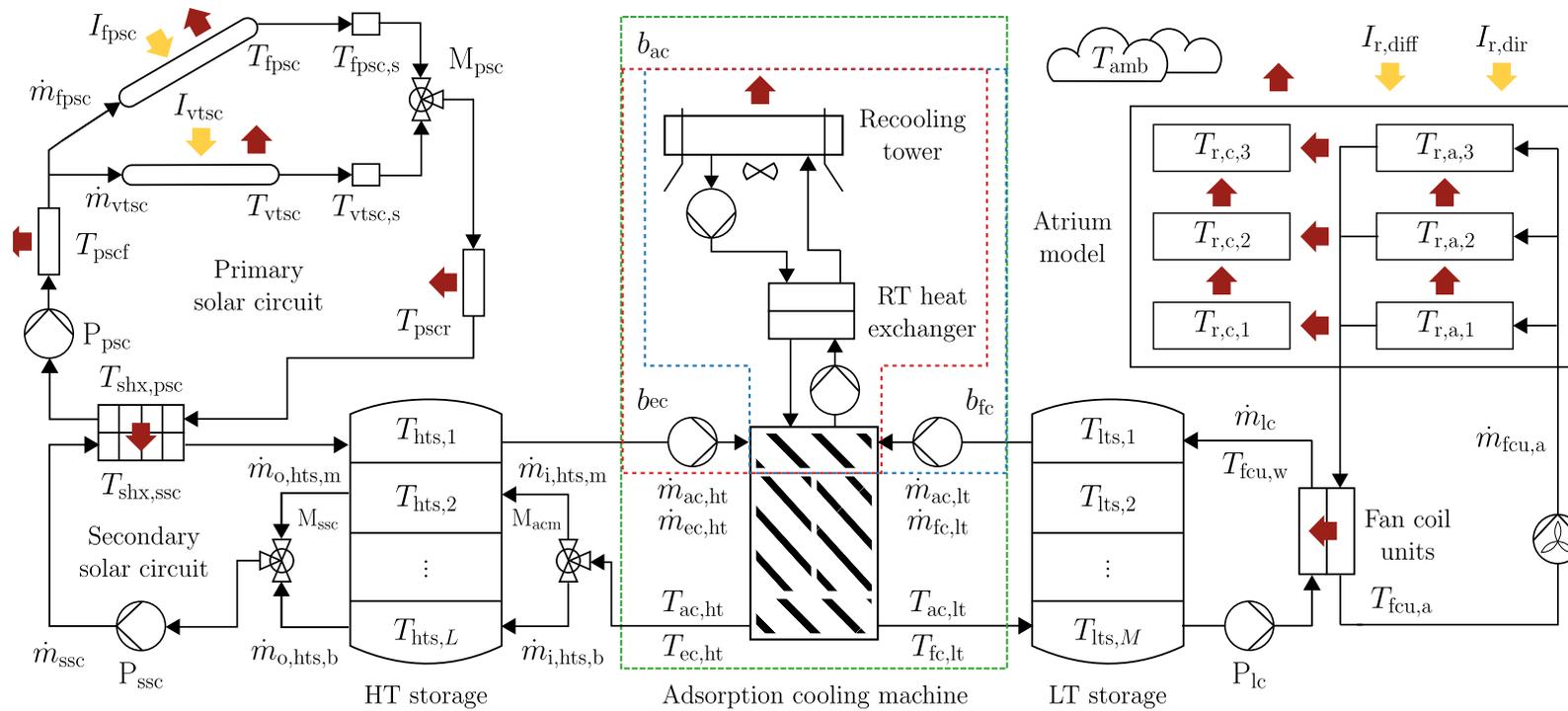
[Bürger et al., 2019]



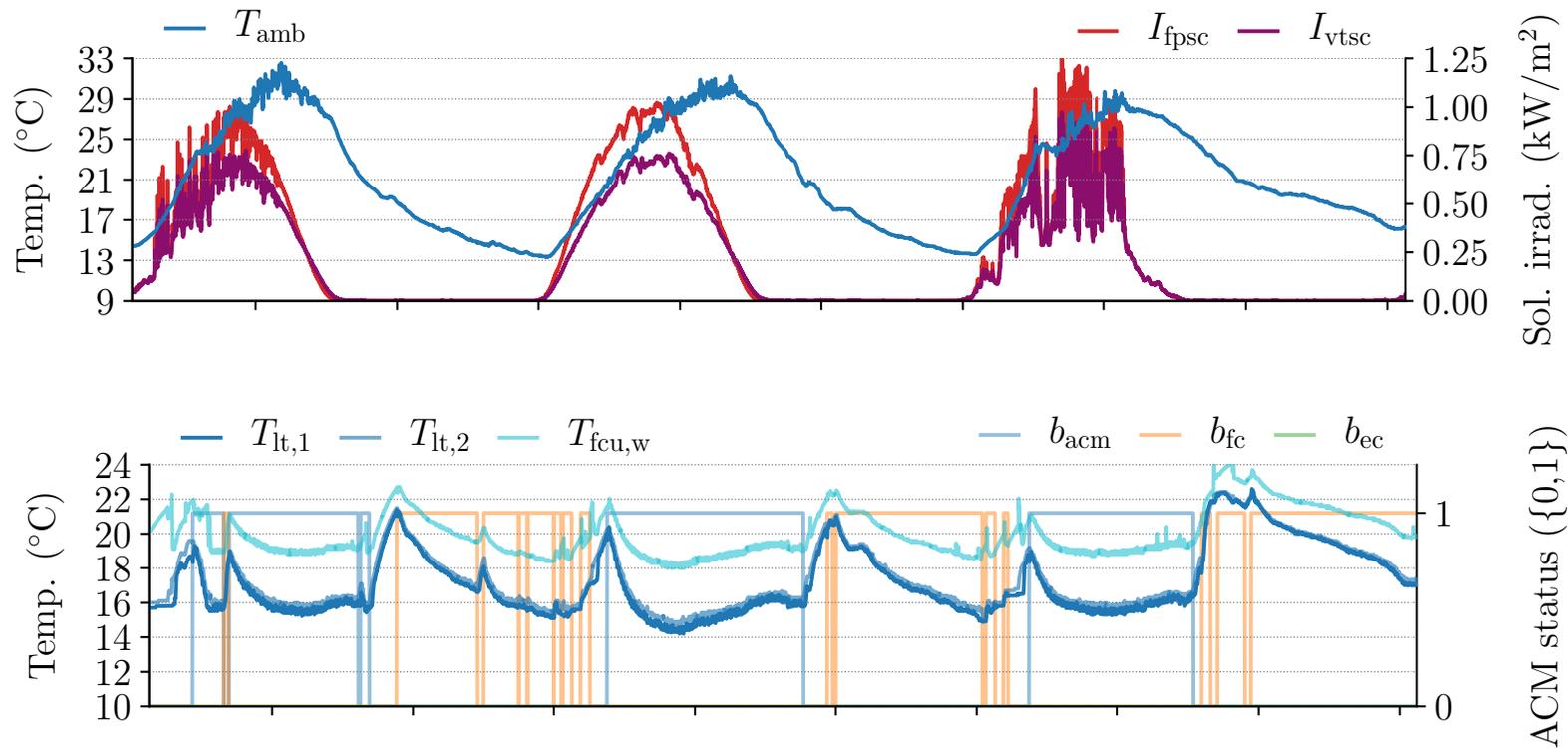
After discretisation of PDE components, obtain nonlinear ODE with 39 states, 6 continuous and 2 binary inputs.

Predict 24 hours. Aim: minimise electricity consumption.

Model Overview

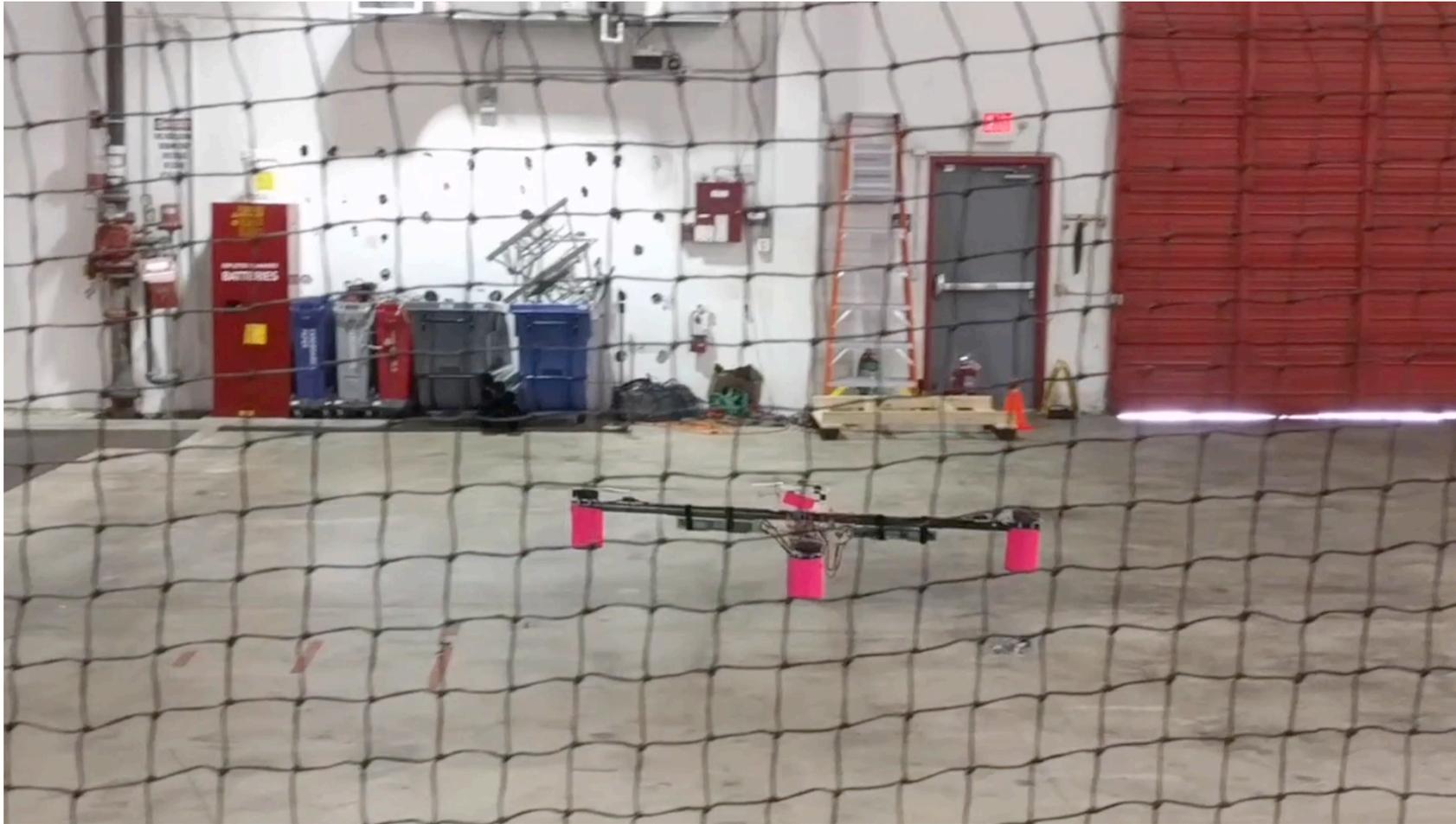


Experimental MPC Results from Sept 14-17, 2019



Every 2 minutes, a new optimization problem is solved, using a real-time algorithm based on CasADi, IPOPT [Wächter and Biegler 2006], and Pycombina [Bürger et al, 2019], an implementation of the combinatorial integral approximation (CIA) method [Sager 2009].

Human sized quadcopter control (Nonlinear MPC) at Kitty Hawk, California, using acados



[Zanelli, Horn, Frison, D., 2018]

Electrical Compressor Control at ABB (Norway)



- work of Dr. Joachim Ferreau and Dr. Thomas Besselmann, ABB
- nonlinear MPC with qpOASES and ACADO, 1ms sampling time
- first tests at 48 MW Drive
- currently, 15% of Norwegian Gas Exports are controlled by Nonlinear MPC

Joachim Ferreau (email from 7.3.2016):

The NMPC installations in Norway (actually 5 compressors at two different sites) are doing fine since last autumn – roughly 80 billion NMPC instances solved by now. In addition, they have proven to work as expected when handling external voltage dips.

eco4wind: MPC for wind turbine control

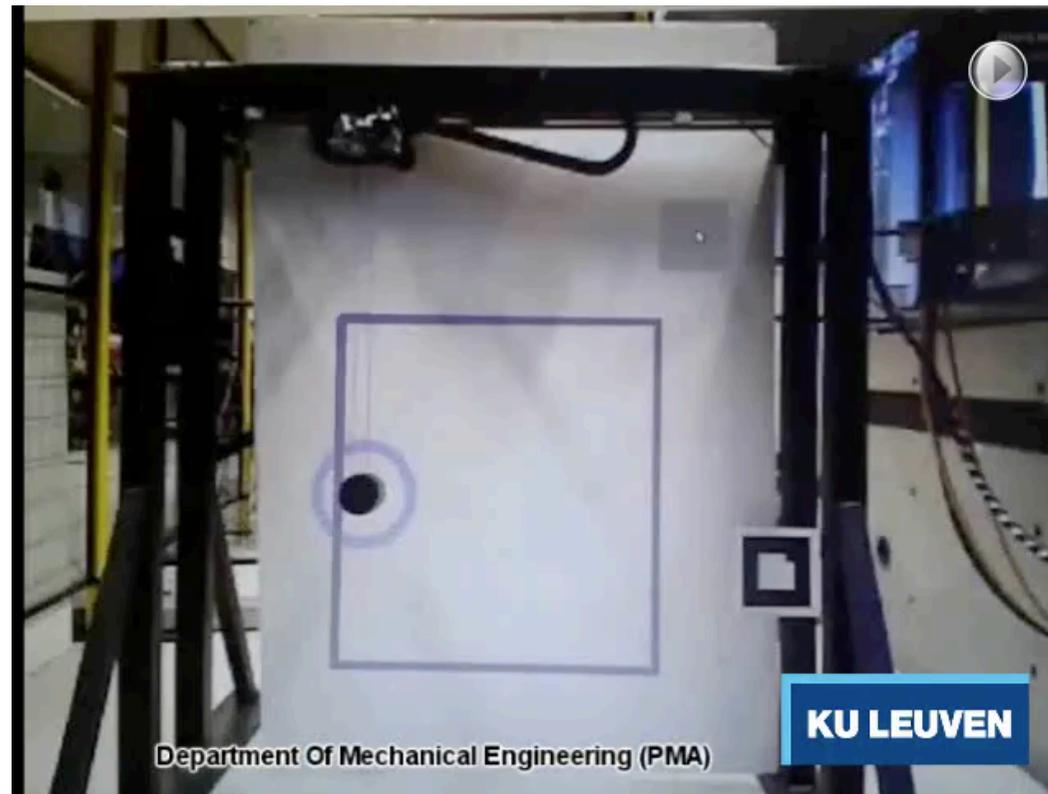


Industrial partners: IAV, SENVION

Nonlinear MPC with about 40 states based on ACADO code generation with QP solver
HPIPM running on industrial hardware at IAV

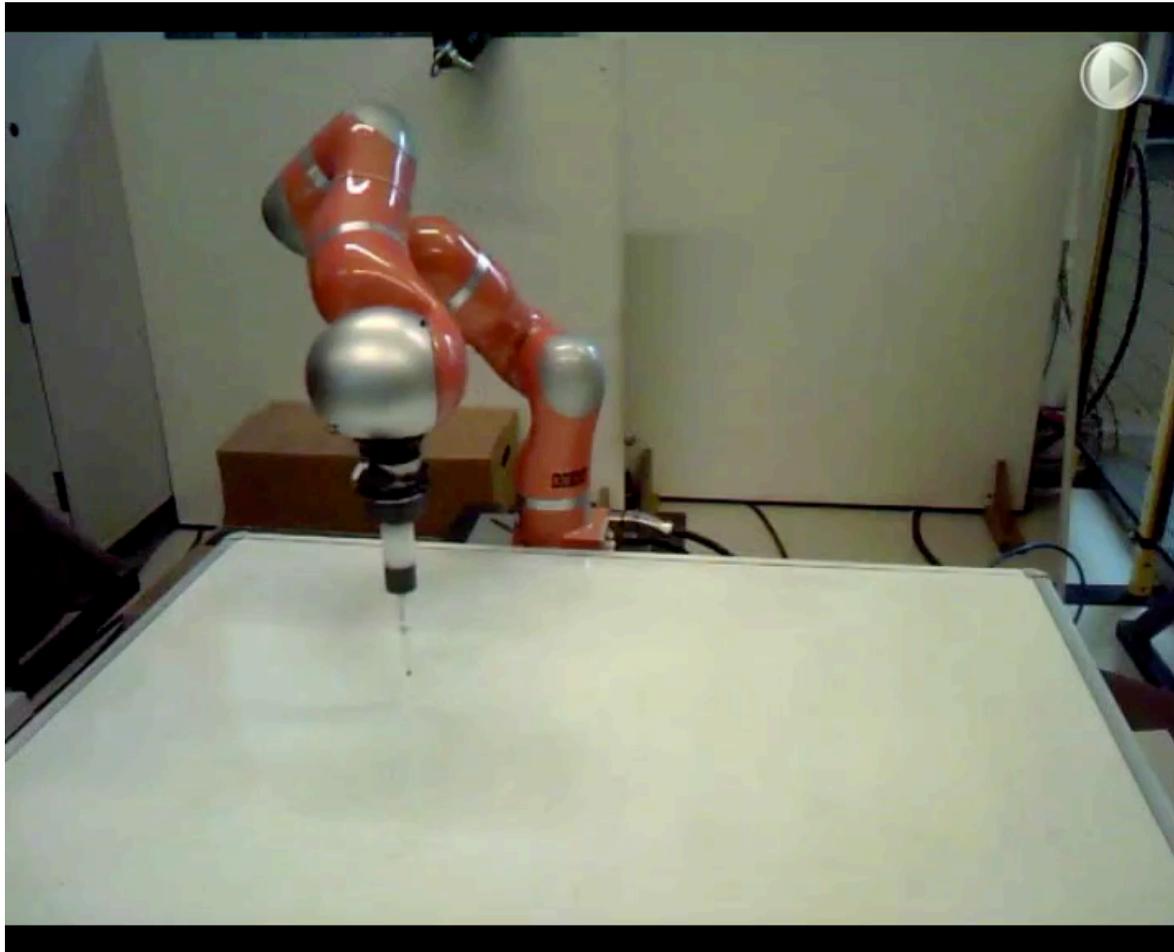
Time Optimal “drawing” by crane

Univ. Leuven [Wannes Van Loock et al.,] (CasADi)

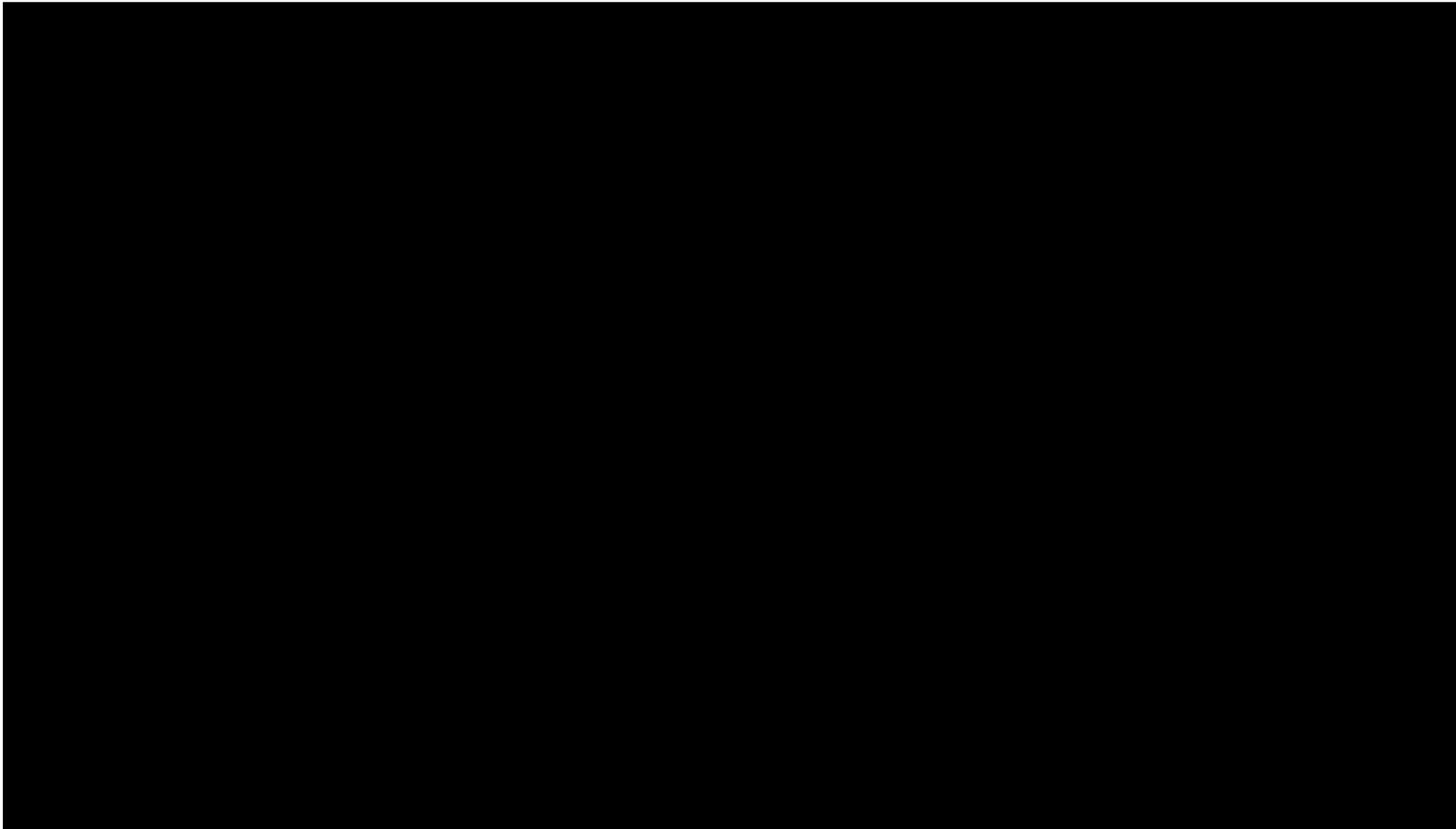


Time-optimal “hand writing” by robot

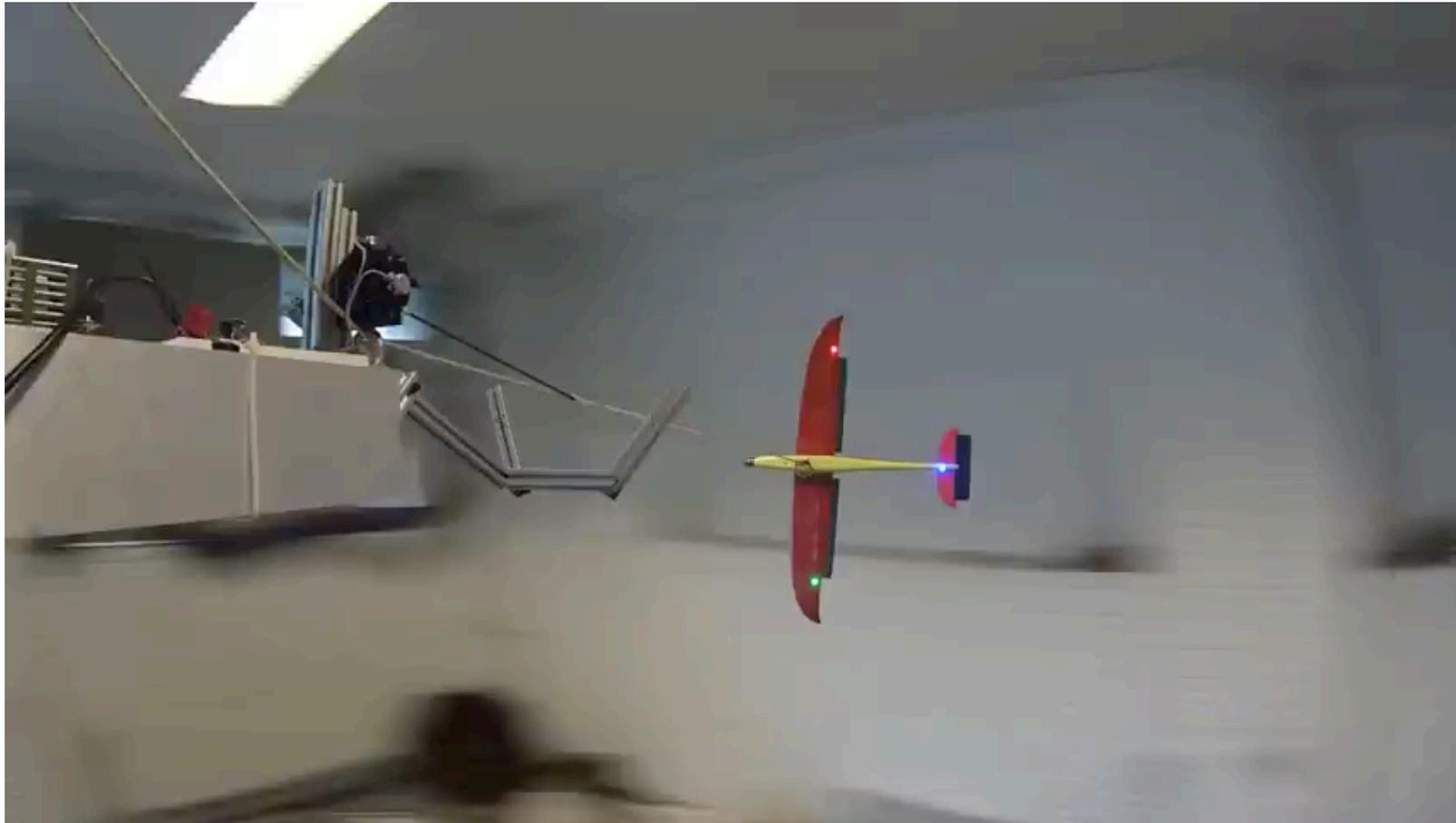
Univ. Leuven [Debrouwere, Swevers] using [Verscheure et al, IEEE TAC 2009]



Predictive control of flight carousel (in Freiburg)



Flight carousel (in Leuven, by M. Vukov)



Nonlinear MPC and Moving Horizon Estimation (MHE)

Closed loop experiments with NMPC & NMHE



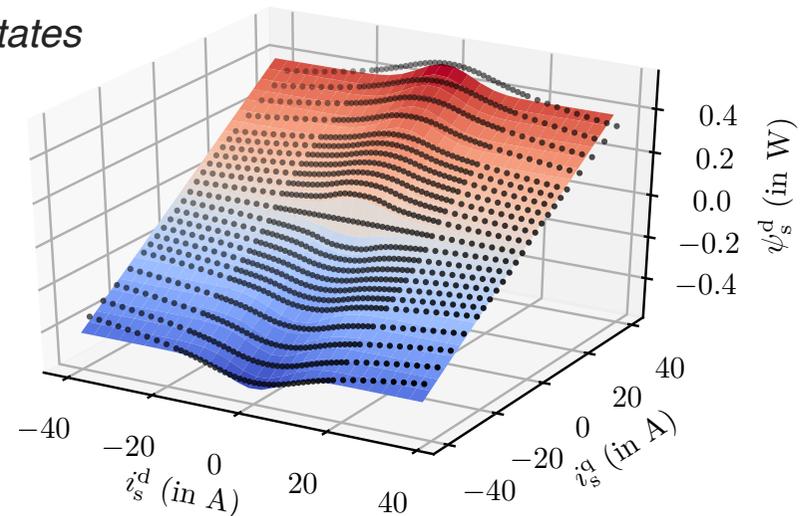
RSM: Control Oriented Differential Algebraic Equation Model

$$\begin{aligned} \frac{d}{dt} \psi_s &= u_s - R_s i_s - \omega J \psi_s + v, \\ 0 &= \psi_s - \Psi_s(i_s), \end{aligned}$$

- differential algebraic equation (DAE)
- currents i_s as implicitly defined *algebraic states*
- analytical flux map approximations:

$$\Psi_s^q(i_s^d, i_s^q, \theta_q) = \frac{c_0^q}{\sqrt{2\pi\sigma_d^2}} \exp\left(-\gamma\left(i_s^d, \sigma_d\right)\right) \operatorname{atan}\left(c_1^q i_s^q\right) + c_2^q i_s^q,$$

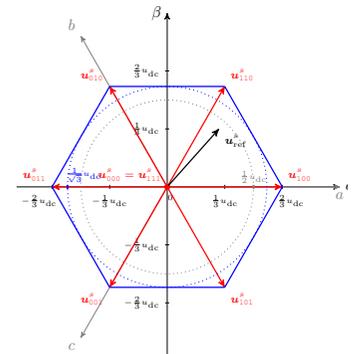
$$\gamma(x, y) := \frac{1}{2} \left(\frac{x}{y}\right)^2$$



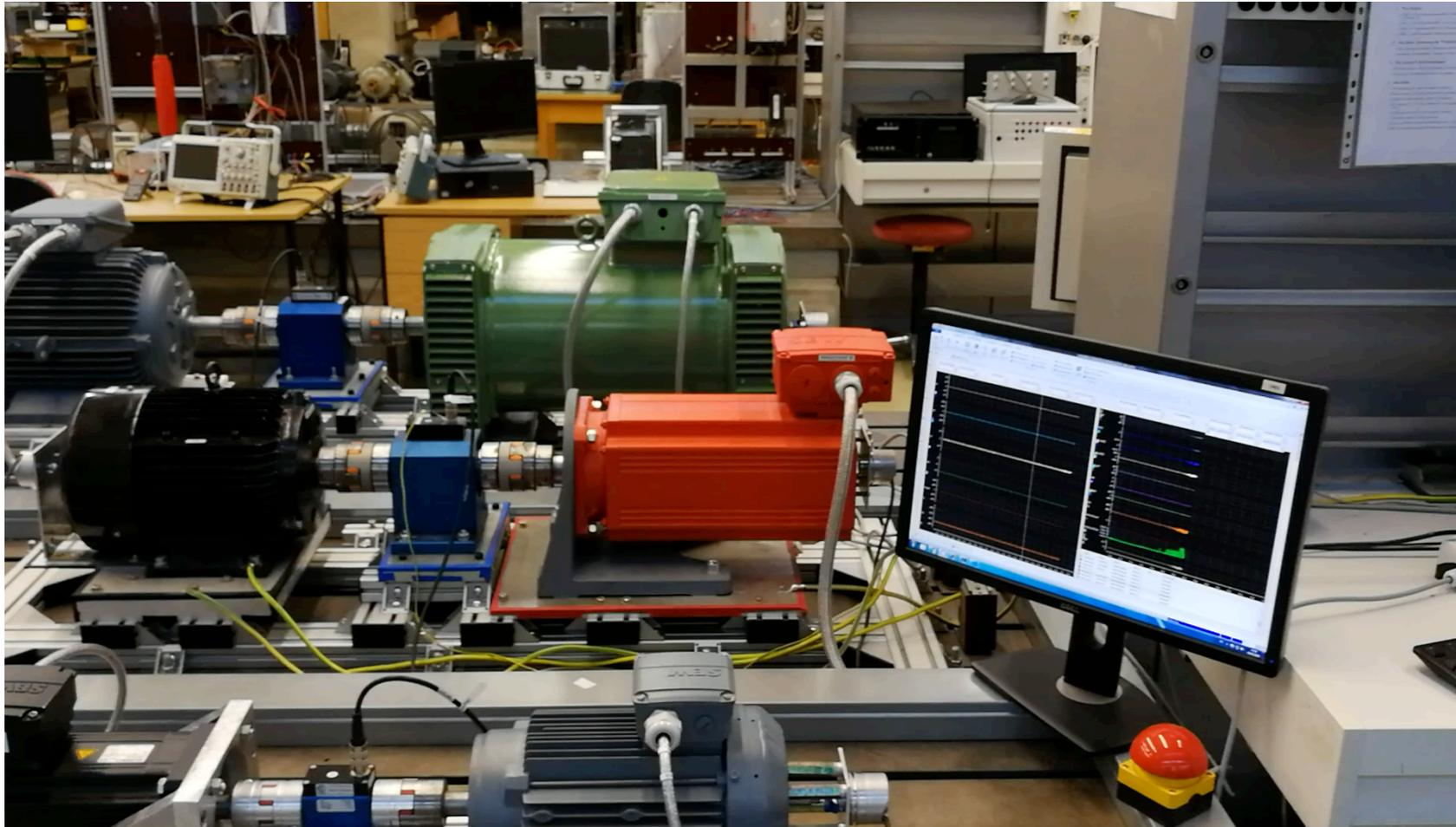
(b) fitted grey box flux model - q-component

Optimization Problem resulting from Direct Multiple Shooting:

$$\begin{aligned}
 & \min_{\substack{\psi_0, \dots, \psi_N \\ u_0, \dots, u_{N-1}}} \frac{T_h}{2N} \sum_{i=0}^{N-1} \left\| \begin{matrix} \psi_i - \bar{\psi} \\ u_i - \bar{u} \end{matrix} \right\|_W^2 + \frac{1}{2} \|\psi_N - \bar{\psi}\|_{W_N}^2 \\
 & \text{s.t.} \quad \psi_0 - \psi_e = 0, \\
 & \quad g(\psi_i, u_i, \omega_e, v_e) - \psi_{i+1} = 0, \quad i = 0, \dots, N-1, \\
 & \quad u_i^\top u_i \leq \left(\frac{u_{dc}}{\sqrt{3}} \right)^2, \quad i = 0, \dots, N-1, \\
 & \quad \hat{C}u_i \leq \hat{c}, \quad i = 0, \dots, N-1,
 \end{aligned}$$



RSM: Video from NMPC Experiments at TU Munich



CS-NMPC significantly better than PI Controller

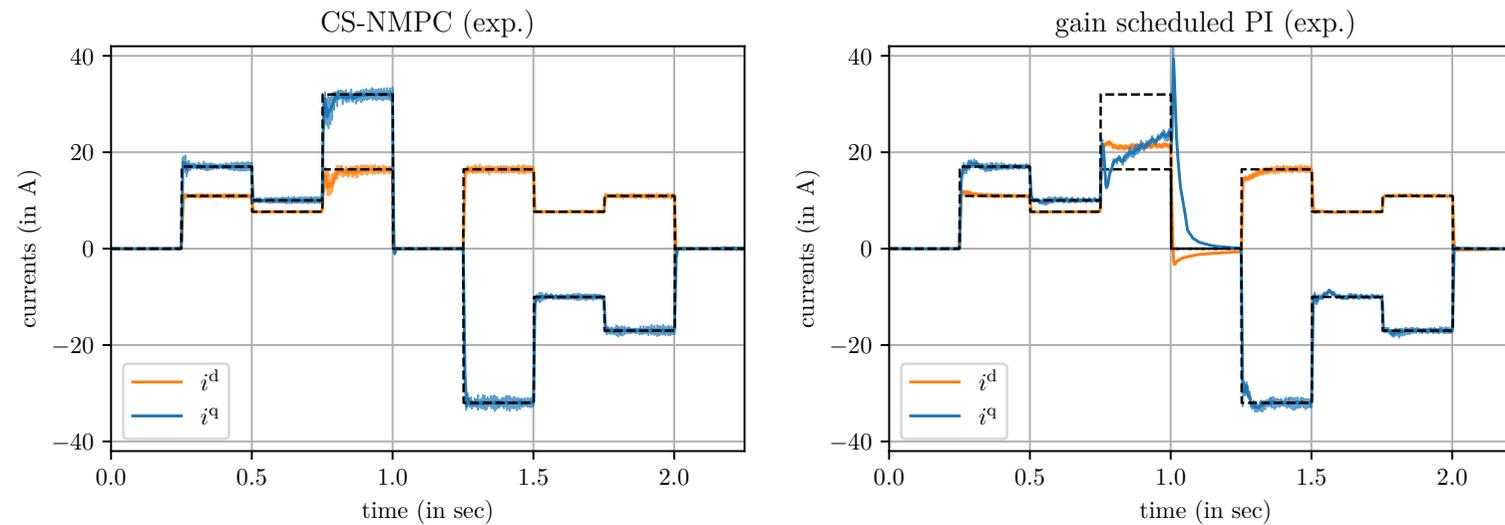


Figure 8: Current steps at $157 \frac{\text{rad}}{\text{s}}$ (experiment): results obtained using the proposed CS-NMPC controller (left) and gain-scheduled PI controller (right). The CS-NMPC controller outperforms the PI controller, especially when the input constraints become active (e.g., between $t = 0.75$ s and $t = 1.00$ s). At the same time, as it can be seen especially between $t = 1.25$ s and $t = 1.50$ s, a faster transient can be achieved, even when the constraints are active only for a short time.

Agenda



Week 1 at Faculty of Engineering (HS 0-26 - Buil. 101)			
	Wed 4-Oct	Thu 5-Oct	Fri 6-Oct
09:00-10:30	Welcome + <i>Lecture 1</i> Dynamic Systems and Simulation (MD)	<i>Lecture 3</i> Dynamic Programming and Optimal Control (MD)	<i>Lecture 5</i> Monte Carlo RL, Temporal Difference and Q-Learning (JB)
10:30-11:00	Coffee Break	Coffee Break	Coffee Break
11:00-12:20	<i>Exercise 1</i> Dynamic Systems and Simulation	<i>Exercise 3</i> Dynamic Programming and Optimal Control	<i>Exercise 5</i> Q-learning
12:20-12:30		<i>Project Brainstorming</i>	<i>Project Brainstorming</i>
12:30-14:00	Lunch	Lunch	Lunch
14:00-15:30	<i>Lecture 2</i> Numerical Optimization (MD)	<i>Lecture 4</i> Deep Learning (JB)	<i>Lecture 6</i> RL with Function Approximation (JB)
15:30-16:00	Coffee Break	Coffee Break	Coffee Break
16:00-17:30	<i>Exercise 2</i> Numerical Optimization	<i>Exercise 4</i> PyTorch	<i>Exercise 6</i> DQN
		Break	
19:00-		Social Gathering	

Week 2 at Historical University (HS 1015 - Buil. KG I)				
Mon 9-Oct	Tue 10-Oct	Wed 11-Oct	Thu 12-Oct	Fri 13-Oct
Motivating Discussion + <i>Lecture 7</i> NMPC (MD)	<i>Lecture 9</i> Policy gradient and Actor-Critic methods (MD, JB)	<i>Lecture 11</i> Introduction MPC and RL Framework (SG)	<i>Lecture 13</i> When to use RL in MPC? (SG)	Project work
Coffee Break	Coffee Break	Coffee Break	Coffee Break	Coffee Break
<i>Exercise 7</i> NMPC (Acados)	<i>Exercise 9</i> Actor-Critic	<i>Exercise 11</i> MPCRL, tba	<i>Exercise 12</i> MPCRL, tba	Project work / Project presentations
<i>Project Brainstorming</i>	<i>Project Brainstorming</i>			
Lunch	Lunch	Lunch	Lunch	Lunch
<i>Lecture 8</i> Transformers (JB)	<i>Lecture 10</i> MPPI (HB) / Model-based RL (JB)	<i>Lecture 12</i> Safety and Stability in MPCRL (SG)	Project work	Project presentations
Coffee Break	Coffee Break	Coffee Break	Coffee Break	Coffee Break
<i>Exercise 8</i> Transformer in practice (forecasting)	<i>Exercise 10</i> MPPI	Project brainstorming and kick-off	Project work	
Break	From 17 to 18: Project Announcements	Break		
Aperitif at Waldsee		Workshop Dinner		