

# Lecture 14: Outlook and Research Insights

Oliver Wallscheid



# Table of contents

- 1 Safe reinforcement learning
- 2 Real-world implementation with fast policy inference
- 3 Meta reinforcement learning

## Recap: optimal control and constraints

Real-world systems are always subject to certain state constraints  $\mathcal{X}$  and input limitations  $\mathcal{U}$ . Violating those can lead to safety issues.

$$v_k^* = \max_{\mathbf{u}_k} \sum_{i=0}^{N_p} \gamma^i r_{k+i+1}(\mathbf{x}_{k+i}, \mathbf{u}_{k+i}), \quad (14.1)$$

$$\text{s.t.} \quad \mathbf{x}_{k+i+1} = \mathbf{f}(\mathbf{x}_{k+i}, \mathbf{u}_{k+i}), \quad \mathbf{x}_{k+i} \in \mathcal{X}, \quad \mathbf{u}_{k+i} \in \mathcal{U}.$$

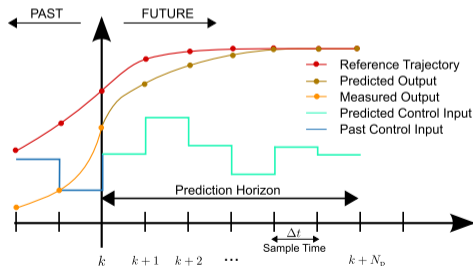


Fig. 14.1: MPC scheme (source: [www.wikipedia.org](http://www.wikipedia.org), by Martin Behrendt CC BY-SA 3.0)

# Application examples with safety-relevant constraints

Collaborative robot control (source: [www.wikipedia.org](http://www.wikipedia.org), CC BY-SA 4.0)



Autonomous car driving (source: [www.wikipedia.org](http://www.wikipedia.org), CC BY-SA 4.0)

Energy system control



Medication control (source: [www.wikipedia.org](http://www.wikipedia.org), CC BY-SA 4.0)

# Safety levels

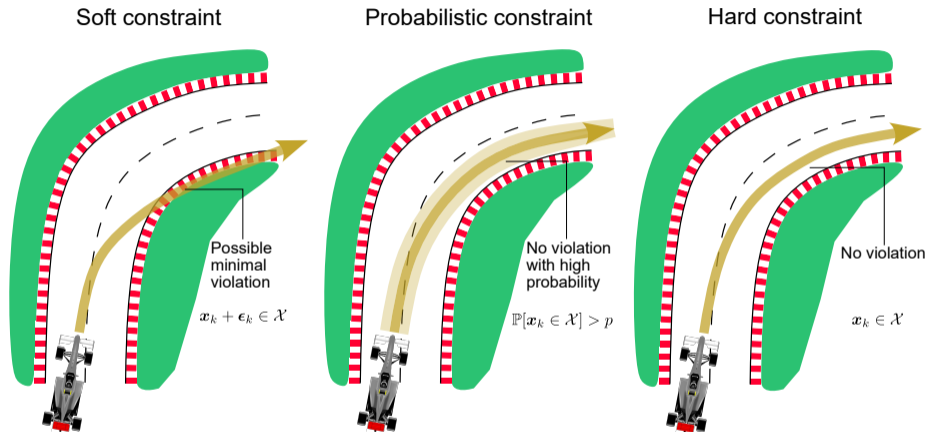
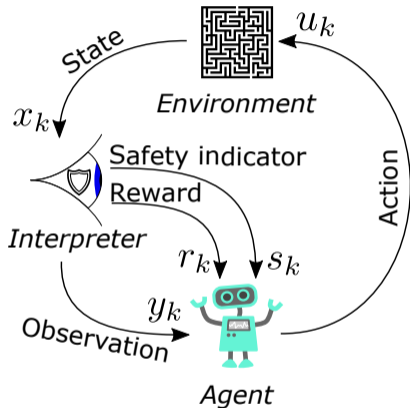
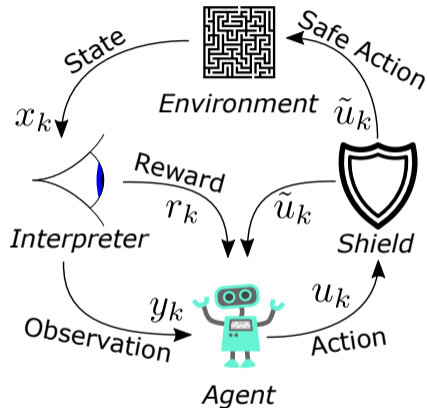


Fig. 14.2: Different levels of safety (derived from L. Brunke et al., *Safe Learning in Robotics: From Learning-Based Control to Safe Reinforcement Learning*, Annual Review of Control, Robotics, and Autonomous Systems, 2022)

# Bird's eye view on RL concepts integrating safety



(a) Safety critic: add a critic which indicates to which extent the current data sample fits to a safe situation



(b) Safety shield: use a priori or learned model knowledge of the environment to make predictions identifying actions leading to unsafe situations

# Achievable safety levels and model knowledge

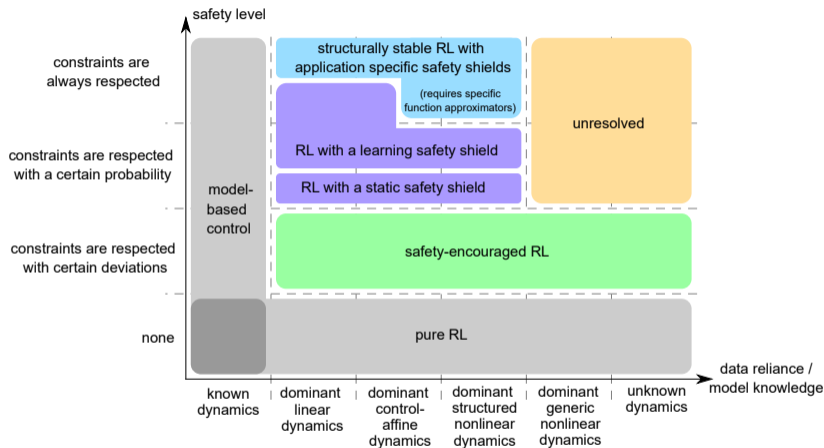
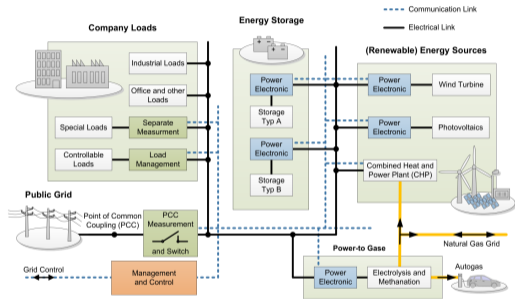
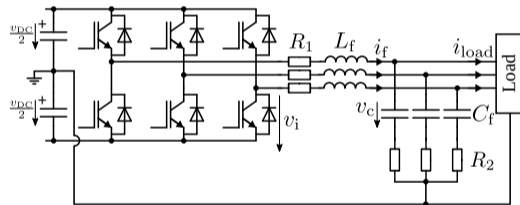


Fig. 14.3: Safety and model knowledge map (derived from L. Brunke et al., *Safe Learning in Robotics: From Learning-Based Control to Safe Reinforcement Learning*, Annual Review of Control, Robotics, and Autonomous Systems, 2022)

# Energy system control application



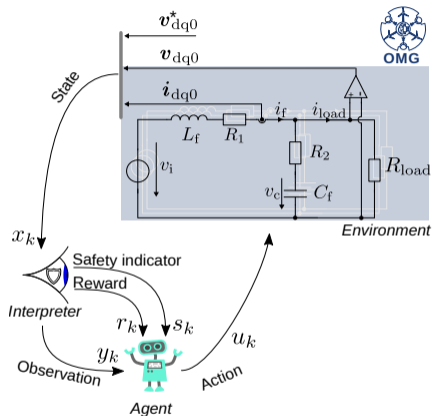
(a) Example microgrid that can be emulated in the LEA Microgrid Laboratory.



(b) Application under investigation: Three-phase grid-forming inverter disturbed by stochastic load



# Reference tracking with disturbance rejection



- ▶ Cont. state- and actionspace
- ▶ Deep deterministic policy gradient agent
- ▶ Grid-forming inverter
- ▶ Stochastic load acts as disturbance
- ▶ State per phase:  $\mathbf{x}_k = [i_f, v_C]$ ,  $v_i = v_{DC} \cdot u_k$
- ▶  $r_k = f(v_C, v^*, i_f) \in [1, -0.75]$
- ▶  $s_k = -1$ , if limit ( $i_f$  or  $v_C$ ) is exceeded

Fig. 14.4: Simulation setting with environment modeled using OpenModelica Microgrid Gym

# Reward design for grid-forming inverter

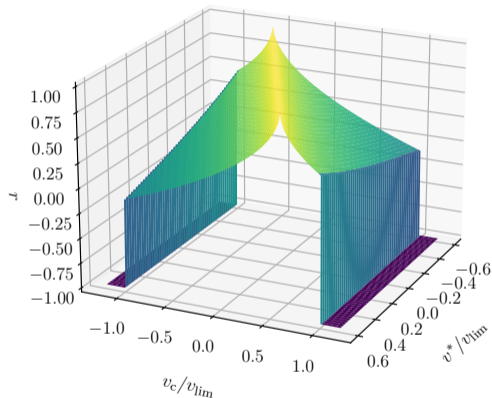


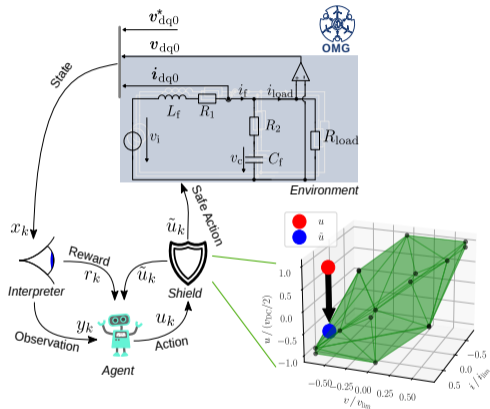
Fig. 14.5: Reward function 14.2 for different reference and measured voltages and currents below nominal current

- ▶ Three cases, depending on operation point

$$r = \begin{cases} \text{MRE}(v_C, v^*), & \text{A} \\ \text{MRE}(v_C, v^*) + f(i_f), & \text{B} \\ -1, & \text{C} \end{cases} \quad (14.2)$$

- ▶ **A**  $v_C \leq v_{\text{lim}} \wedge i_f \leq i_{\text{nom}}$
- ▶ **B**  $v_C \leq v_{\text{lim}} \wedge i_{\text{nom}} \leq i_f \leq i_{\text{lim}}$
- ▶ **C** otherwise
- ▶ Linear punishment term  $f(i_f)$

# Reference tracking with disturbance rejection using safety shield



- ▶ Safety shield: Ensure that action does not cause state limit violation in future system trajectories
- ▶ Such a state action pair is called feasible
- ▶ Calculation of **feasible set** requires a model
- ▶ Training data can be utilized to **identify** model
- ▶ Here, recursive least squares (RLS) is applied

Fig. 14.6: Safety shield based on feasible set

# Safety shield based on feasible set - proof of concept (1)

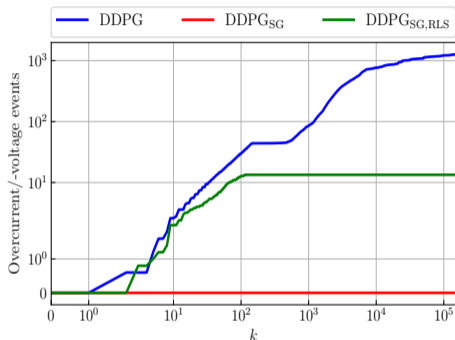
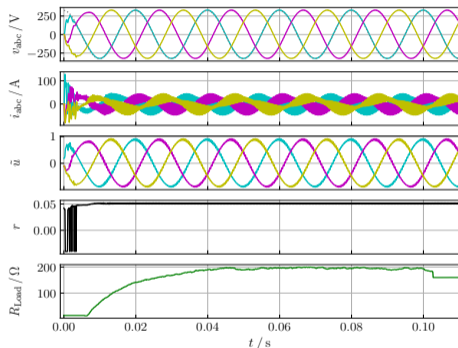


Fig. 14.7: Accumulated unsafe events (overcurrent/-voltage) per trainingstep  $k$

- ▶ Three different approaches
- ▶ **DDPG**: Agent without safety shield
- ▶ **DDPG<sub>SG</sub>**: Agent with safety shield using perfect a priori knowledge
- ▶ **DDPG<sub>SG,RLS</sub>**: Agent with safety shield without a priori knowledge, identifying model using RLS
- ▶ Five agents trained per approach
- ▶ Results in D. Weber et al., *Safe Reinforcement Learning-Based Control in Power Electronic Systems*, 2023

## Saftey shield based on feasible set - proof of concept (2)



- ▶ DDPG<sub>SG,RLS</sub> agent trained for 150000 steps
- ▶  $R_{Load}$  changes every step based on random process
- ▶ Additional events – load steps and drifts – triggered randomly

Fig. 14.8: Blackstart after training using DDPG<sub>SG,RLS</sub>

# Table of contents

- 1 Safe reinforcement learning
- 2 Real-world implementation with fast policy inference
- 3 Meta reinforcement learning

# Real-time implementation aspects (1)

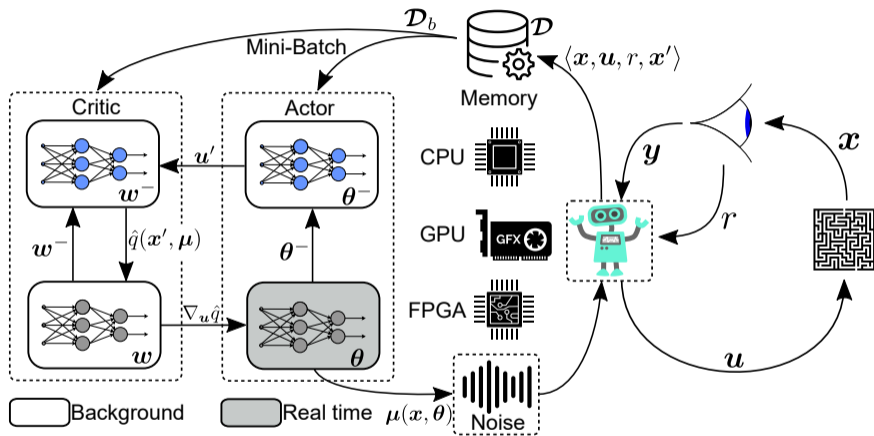
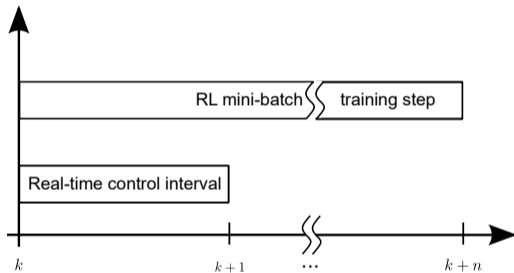
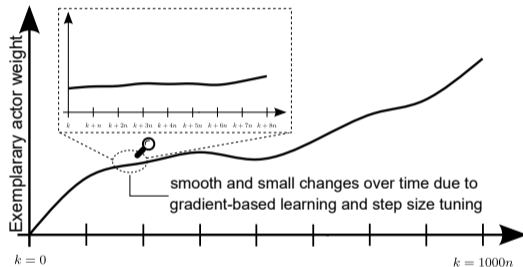


Fig. 14.9: DDPG implementation example (derivative work of Fig. 1.1 and [wikipedia.org](https://en.wikipedia.org), CC0 1.0)

## Real-time implementation aspects (2)



(a) Real-time control requirement vs. learning time



(b) Typical evolution of RL parameter weights during learning



# Application example: deep Q direct torque control

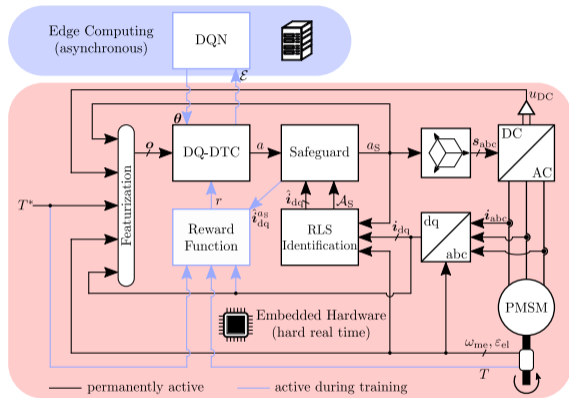


Fig. 14.10: Deep Q direct torque control schematic

- ▶ The DQ-DTC is basically a DQN
- ▶ Sampling time of the plant system is  $T_s = 50 \mu\text{s}$
- ▶ DQN inference, safeguarding and system identification must fit into  $T_s$
- ▶ Source: M. Schenke et al., *Finite-Set Direct Torque Control via Edge Computing-Assisted Safe Reinforcement Learning for a Permanent Magnet Synchronous Motor*, 2023

# Fast neural network inference

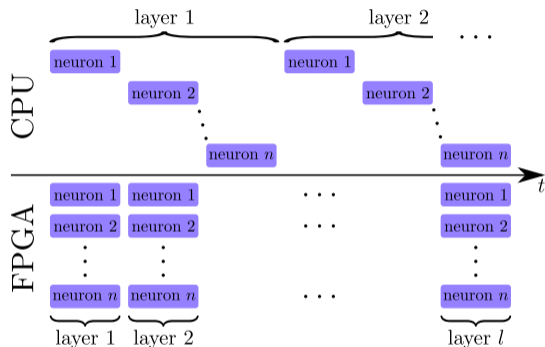


Fig. 14.11: Conceptual comparison of CPU and FPGA evaluation of a neural network

- ▶ Each neuron has the same job  
$$y_{n,l+1} = f(\mathbf{y}_l^\top \mathbf{w}_{n,l} + b_{n,l})$$
- ▶ CPU must evaluate each neuron sequentially
- ▶ FPGA can evaluate each neuron at the same time
- ▶ Maximum number of parallel computations is limited

# Edge reinforcement learning

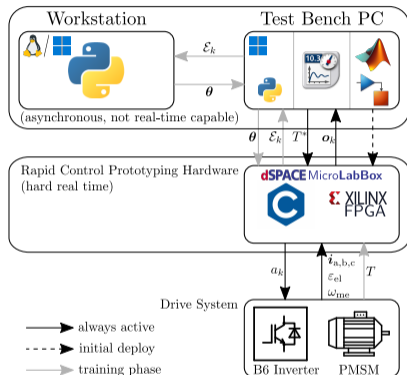


Fig. 14.12: Our edge reinforcement learning pipeline

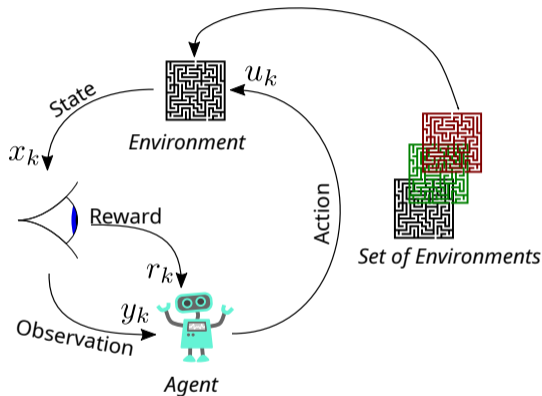
- ▶ Backward pass / learning steps are outsource to workstation
- ▶ Communication between test bench and workstation is based on TCP/IP
- ▶ Backward pass is generic and has no time constraints → low application effort

Youtube link: [Coffee machine vs. deep Q direct torque control](#)

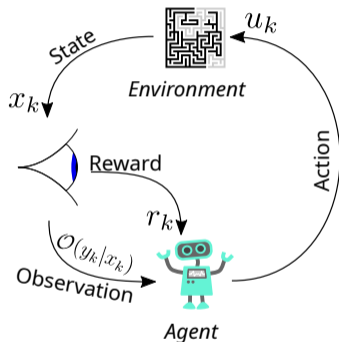
# Table of contents

- 1 Safe reinforcement learning
- 2 Real-world implementation with fast policy inference
- 3 Meta reinforcement learning

# Meta reinforcement learning - the setting (1)



(a) General problem class is similar, environments only differ in some characteristics, the agent could transfer learned behavior

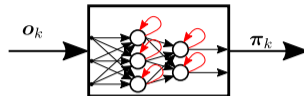


(b) Solution approach: treat the environment as partially observable, distinguishing details are not directly available

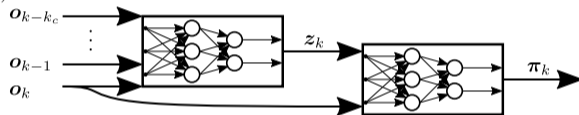
# Meta reinforcement learning - the setting (2)

- ▶ The agent must have some mechanism that allows adaptation to the specific environment
- ▶ This means, the distinguishing details must be extracted in some way
- ▶ Usually, they can be retrieved from a larger set of observations

a) Recurrent networks



b) Context networks



c) Expert knowledge

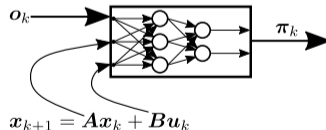


Fig. 14.13: Different concepts of meta learning

# Usage in electric drive control: classical agent

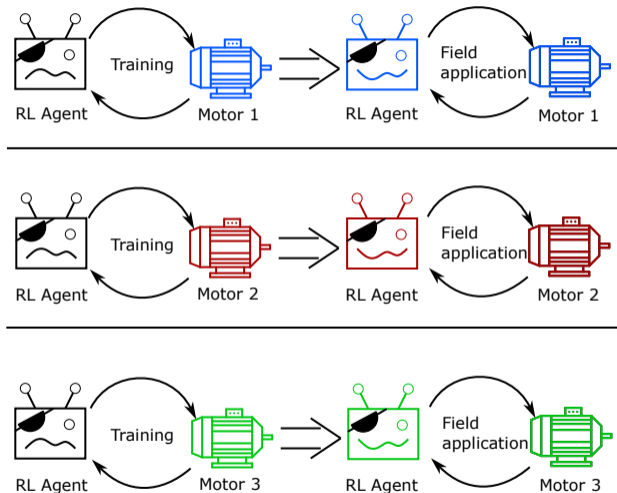


Fig. 14.14: Each agent must be trained individually → huge effort



# Usage in electric drive control: meta agent

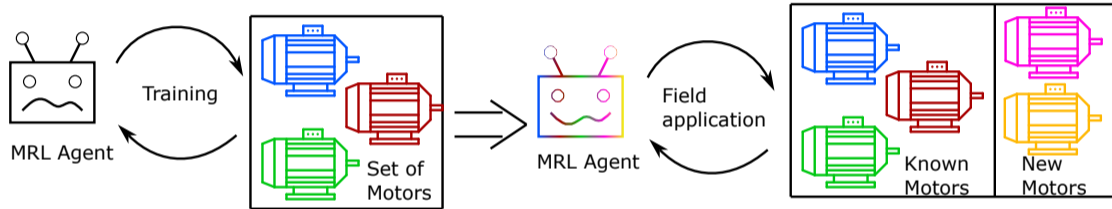


Fig. 14.15: One agent to control them all → effort is limited and independent of the number of controlled environments

# Our setup

- ▶ Make use of context network
- ▶ Generate context  $z$  with a fix set of observations  $\rightarrow z = \text{const.}$
- ▶ Source: D. Jakobeit et al., *Meta-Reinforcement Learning-Based Current Control of Permanent Magnet Synchronous Motor Drives for a Wide Range of Power Classes*, IEEE TPEL, 2023

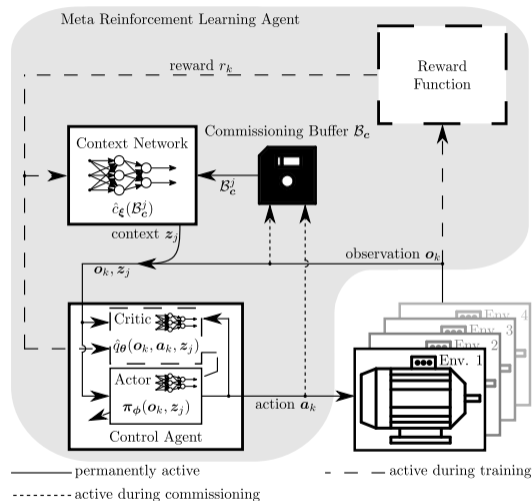
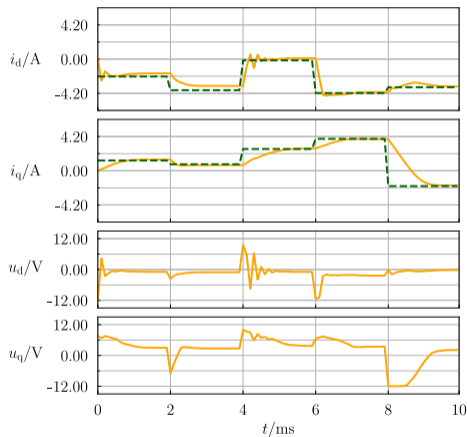
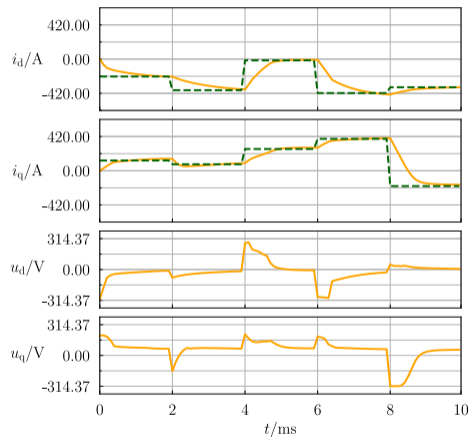


Fig. 14.16: A meta learning concept that we implemented successfully

# Evaluation on (very) different motors



(a) Current control on a PMSM with low rated power



(b) Current control on a PMSM with high rated power

- ▶ Application of RL on technical systems comes with many challenges, e.g.,
  - ▶ Safety limits,
  - ▶ Real-time / computational constraints,
  - ▶ Varying and/or partially unknown environments.
- ▶ Real-world implementations often require more than bare RL algorithms, e.g.,
  - ▶ Integration of available a priori expert knowledge,
  - ▶ Combination with model-based control engineering tools.
- ▶ Ideal integration of data-driven RL solutions together with expert-based control engineering parts is subject to many open research question.