

Deep Reinforcement Learning for Controls – Application to the Path Following Problem

Johannes Ultsch, Jonathan Brembeck, Ricardo de Castro

German Aerospace Center (DLR) - Institute of System Dynamics and Control (SR)

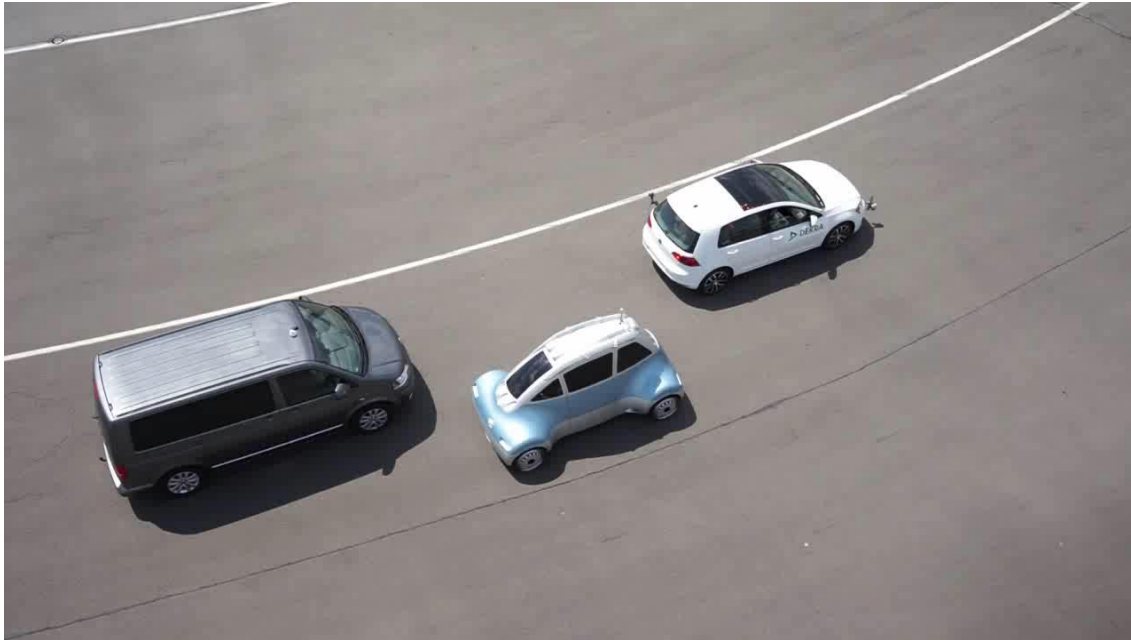
14th ESA Workshop on
Avionics, Data, Control and Software Systems
20-22 October 2020



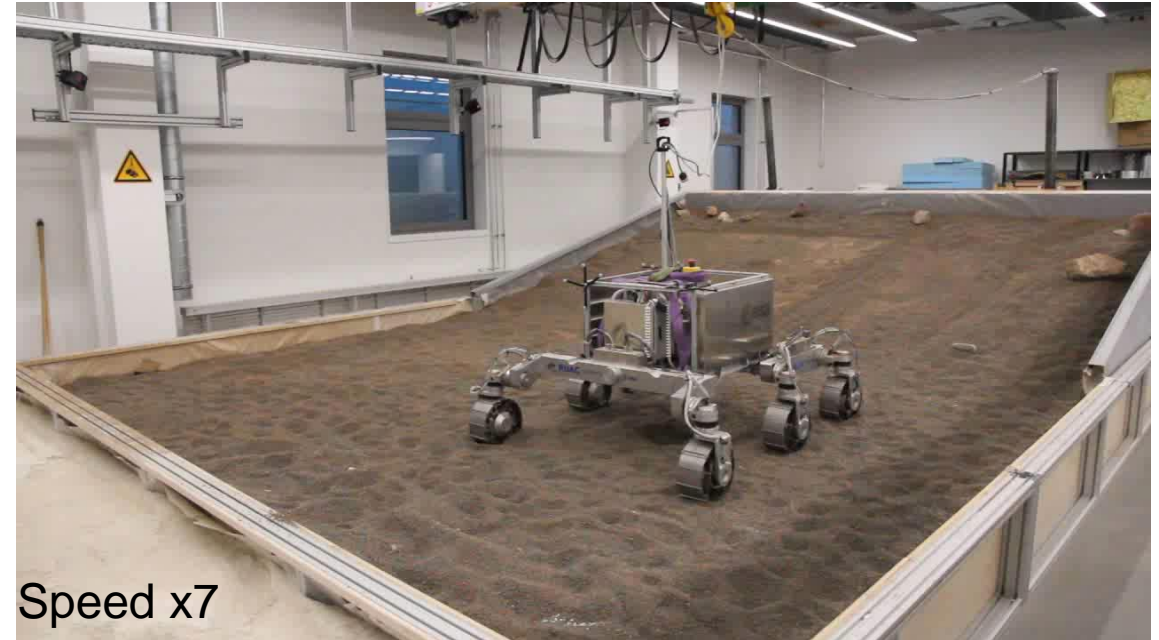
The ROboMObil - Synergies: Transport & Space

Synergies of operation modes

ROboMObil

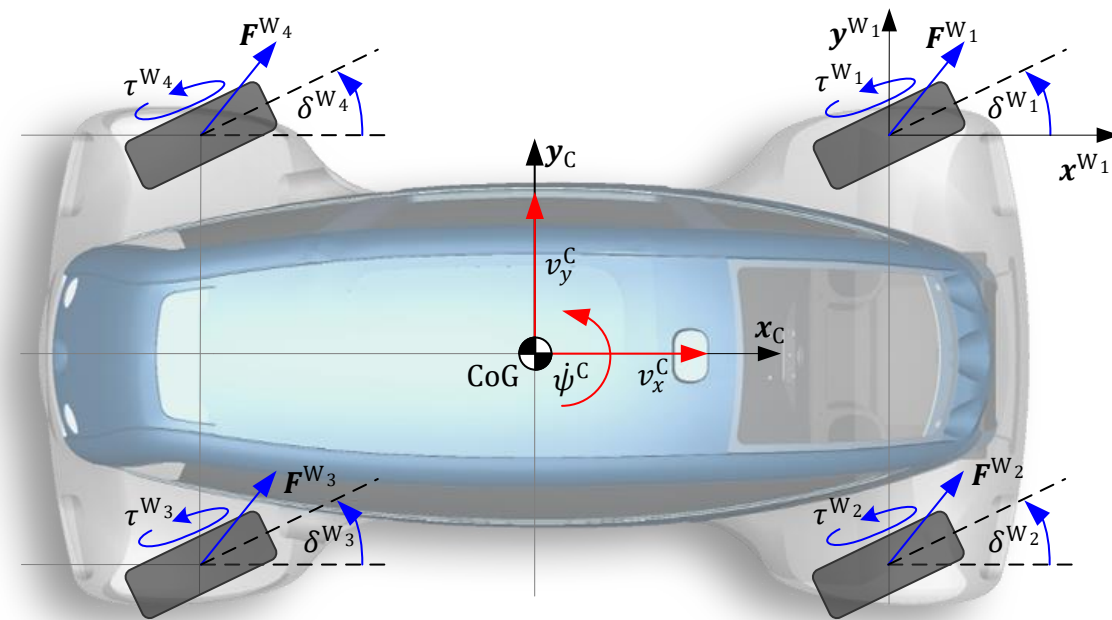


ExoMars rover

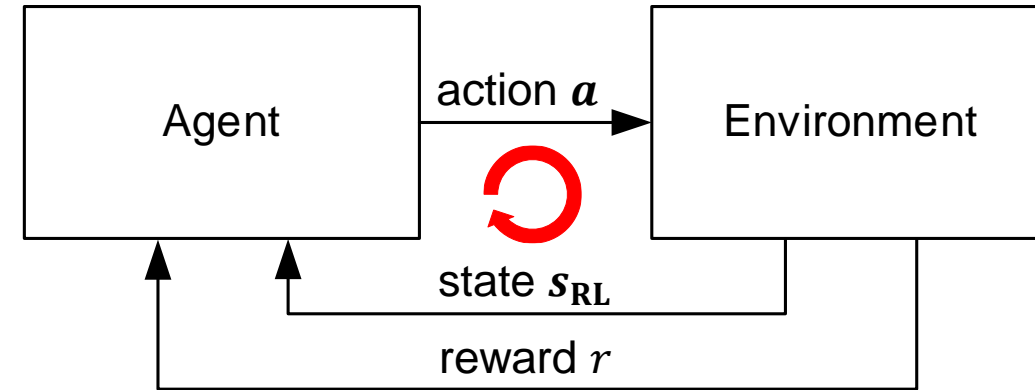


The ROboMObil - The Robotic Electric Vehicle

- Electric X-by-Wire vehicle
- 13 kWh Li-ion battery @ 350V
- 4 Steer-by-wire actuators
- 4 In-wheel traction motors
- 4 Semi-active dampers
- 2 Electro-hydraulic brake-by-wire actuators for friction brake
- Control inputs:
 - Desired steering angle: $\delta^{W_1}, \delta^{W_2}, \delta^{W_3}, \delta^{W_4}$
 - Desired traction torque: $\tau^{W_1}, \tau^{W_2}, \tau^{W_3}, \tau^{W_4}$



The Reinforcement Learning (RL) Approach



- Goal:

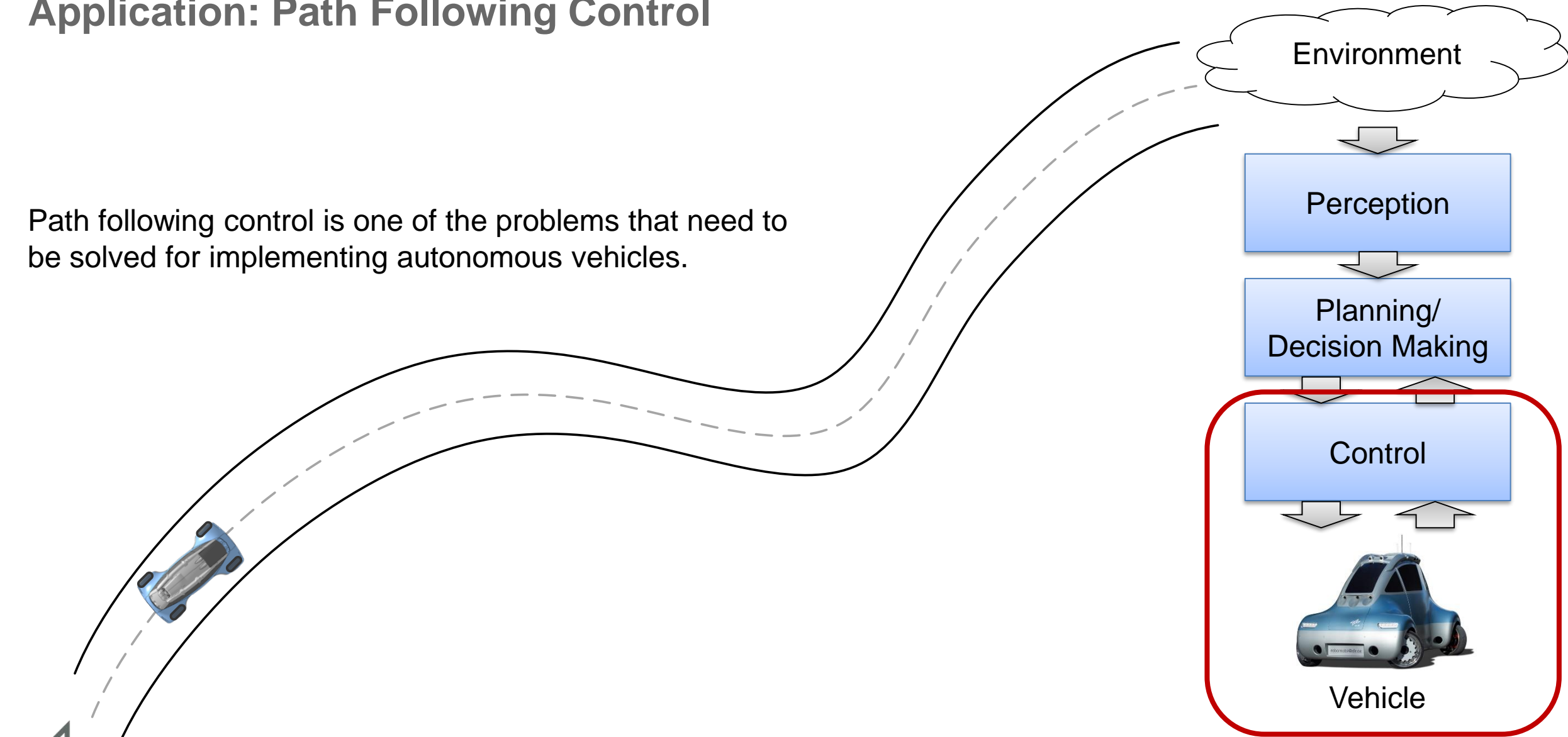
Find a policy $\pi(\mathbf{a}_k | \mathbf{s}_{\text{RL},k})$, which maximizes the *expected discounted return* G_k :

$$G_k = r_{k+1} + \gamma r_{k+2} + \gamma^2 r_{k+3} + \dots = \sum_{i=0}^{\infty} \gamma^i r_{k+i+1} \text{ mit } \gamma \in [0,1]$$

- The policy is adapted based on interaction with the system (*training*)
 - **Offline** training (simulation): safe, fast and scalable
 - **Online** training (real world): no reality gap
- Central design degrees of freedom
 - Selection of the **reward function** r which encodes the control goal
 - Selection of **RL algorithm** and hyperparameters
 - Selection of the **feedback state** \mathbf{s}_{RL}

Application: Path Following Control

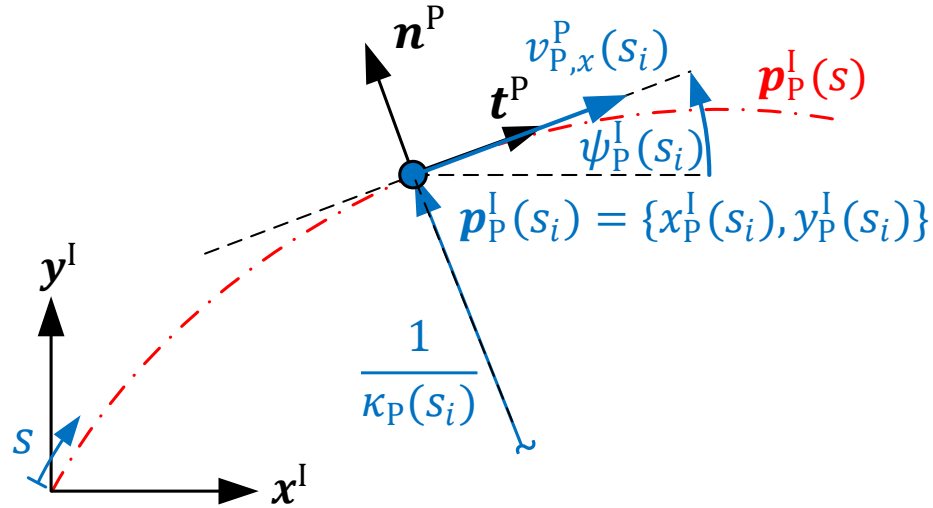
Path following control is one of the problems that need to be solved for implementing autonomous vehicles.



Motion Demand Representation

- Control goal: The vehicle should follow a predetermined motion demand λ as close as possible
- The motion demand is parametrized by the arc length s :

$$\lambda(s) = [x_P^I(s), y_P^I(s), \psi_P(s), \kappa_P(s), v_{P,x}^P(s)]^T$$



- In contrast to trajectory tracking control there is no explicit time dependency in the motion demand related to path following control

The Path Following Problem

- The reference point s^* on the path is defined as:

$$s^* = \arg \min_s \underbrace{\|\mathbf{p}_P(s) - \mathbf{p}_C\|_2}_{e(s)}$$

- We use the time independent path interpolation from [1] to efficiently calculate the reference point on the path

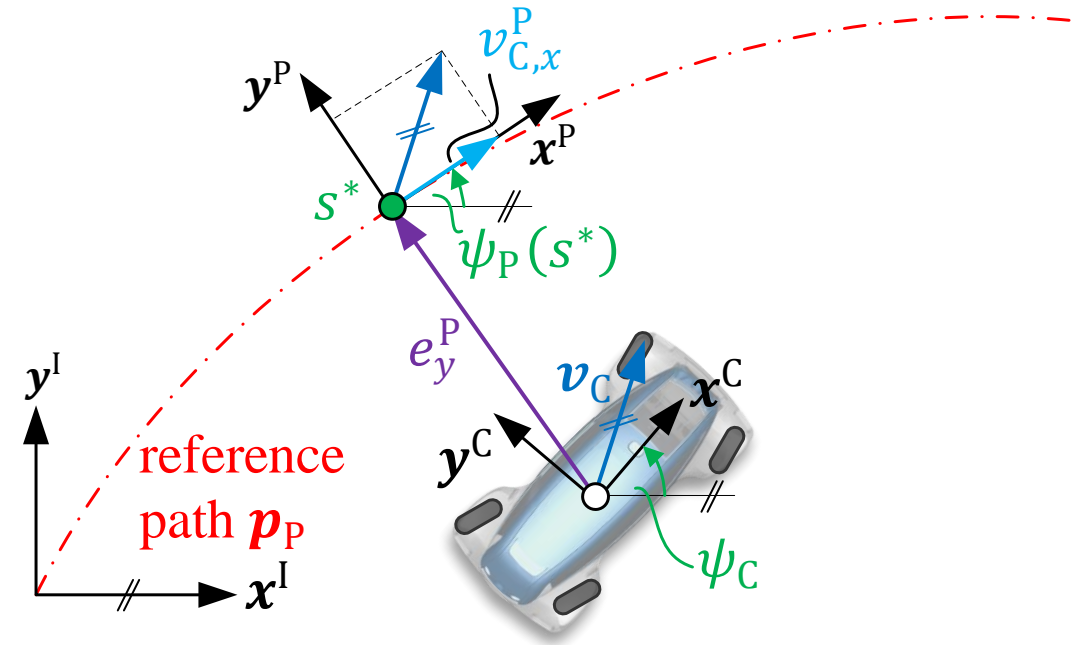
- Control goal:

- $\min e_y^P$ with the position error:
- $\min e_{v_x}^P$ with the velocity error:
 - $\min e_\psi$ with the orientation error:

$$e_y^P = \underbrace{y_P^P}_{=0} - y_C^P$$

$$e_{v_x}^P = v_{P,x}^P - v_{C,x}^P$$

$$e_\psi = \psi_P - \psi_C$$



[1] P. Ritzer, C. Winter and J. Brembeck, „Advanced path following control of an overactuated robotic vehicle,“ in Proceedings of the IEEE Intelligent Vehicles Symposium, 2015.

Application to the ROboMObil Architecture [2]

- RL-PFC:

- $\mathbf{a} = [\delta^f, \delta^r, \tau]$

- with $\tau = \tau^{W_1} = \tau^{W_2} = \tau^{W_3} = \tau^{W_4}$

- The control allocator distributes (δ^f, δ^r) according to Ackermann to all 4 wheels

RL-PFC TV:

- Torque vectoring (TV)

- $\mathbf{a} = [\delta^f, \delta^r, \tau^{W_1}, \tau^{W_2}, \tau^{W_3}, \tau^{W_4}]$

- $\mathbf{s}_{\text{RL},k} = \left[\underbrace{e_{y,k}^P, e_{v_x,k}^P, e_{v_y,k}^P, e_{\psi,k}^P, e_{y,k-1}^P, e_{v_x,k-1}^P, e_{v_y,k-1}^P, e_{\psi,k-1}^P}_{\text{position, velocity and orientation error}}, \underbrace{\kappa_{P,k}, \kappa_{P,k-1}}_{\text{curvature}}, \underbrace{\delta_{c,k}^f, \delta_{c,k}^r, \delta_{c,k-1}^f, \delta_{c,k-1}^r}_{\text{steering angles}} \right]$

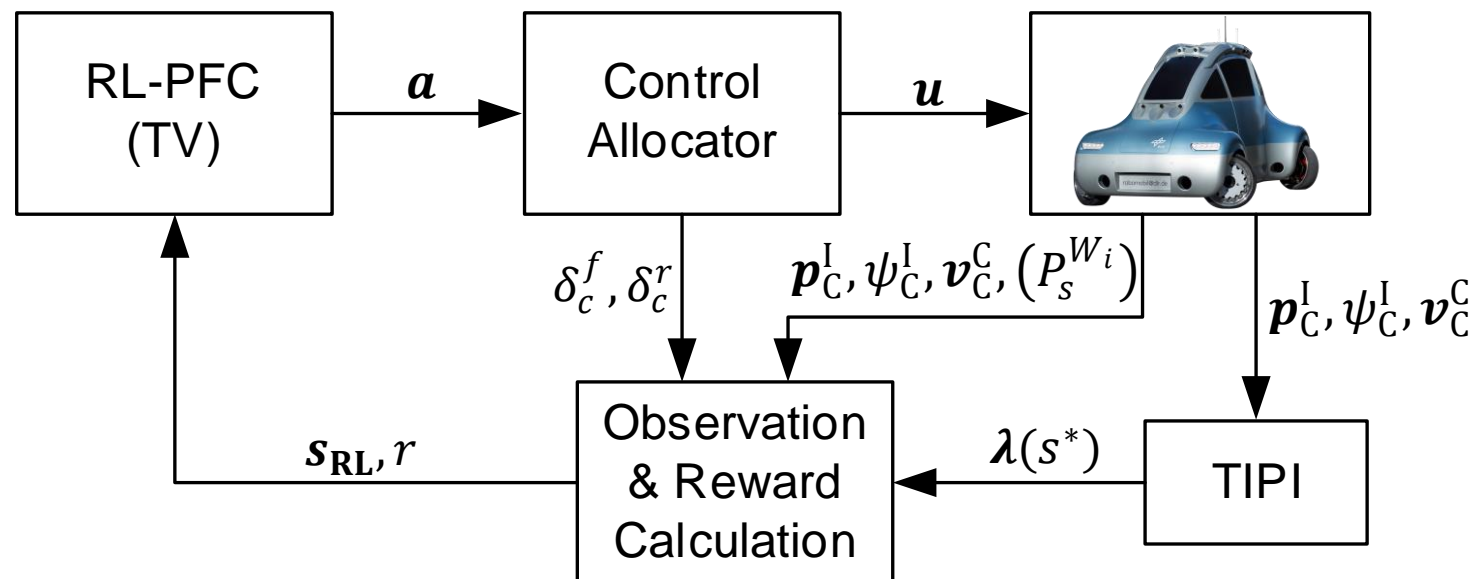
position, velocity and orientation error

curvature

steering angles

at timestep k and $k - 1$

$$\mathbf{u} = [\delta^{W_1}, \delta^{W_2}, \delta^{W_3}, \delta^{W_4}, \tau^{W_1}, \tau^{W_2}, \tau^{W_3}, \tau^{W_4}]$$



[2] J. Ultsch, J. Brembeck and R. de Castro, „Learning-Based Path Following Control for an Over-Actuated Robotic Vehicle,“ in VDI-AUTOREG, Düsseldorf, 2019.



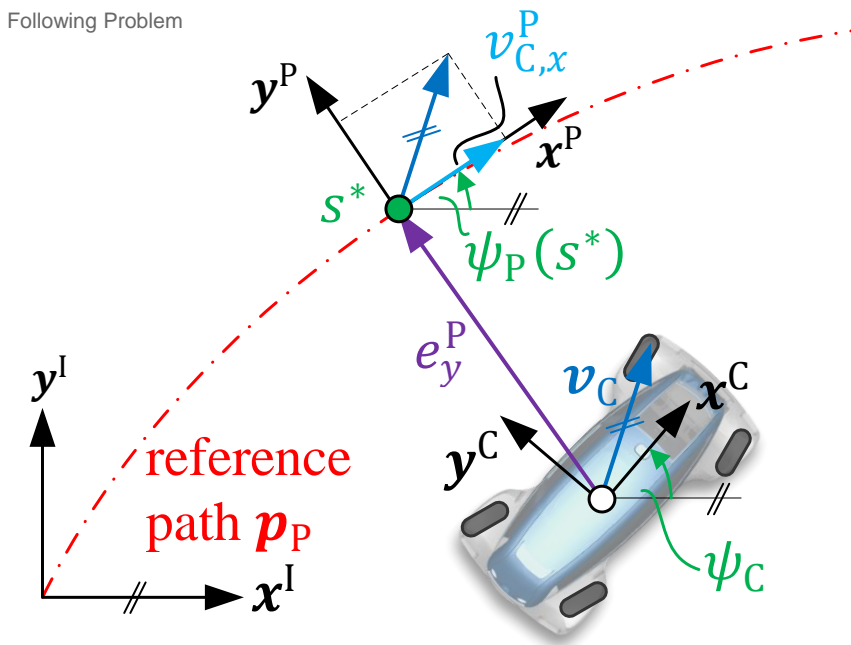
Reward Function Design (1)

Hierarchical control goals:

1. a) Minimize the position error:
 b) Minimize the velocity error:
 c) Minimize the orientation error:
2. Generate smooth steering reference:

$$\begin{aligned} \min e_y^P \\ \min e_{v_x}^P \\ \min e_\psi \end{aligned}$$

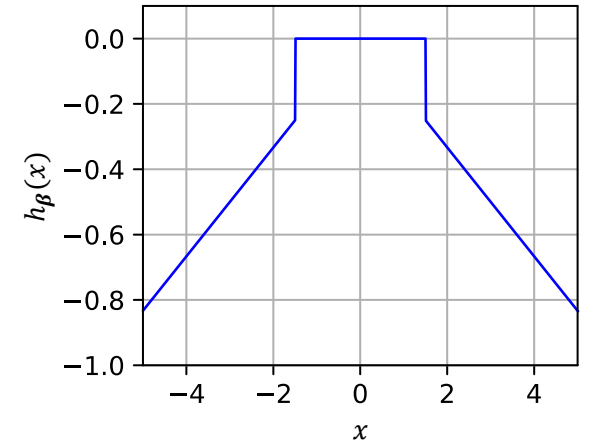
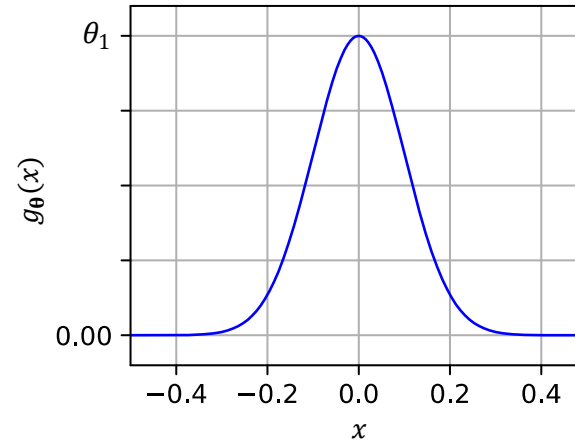
$$\min \Delta\delta^i \quad \text{with} \quad \Delta\delta^i(k) = \delta^i(k) - \delta_c^i(k-1); \forall i \in \{f, r\}$$



Used auxiliary functions:

$$\bullet g_\theta(x) = \theta_1 e^{\frac{-x^2}{2\theta_2}}$$

$$\bullet h_\beta(x) = \begin{cases} 0, & \text{if } |x| < \beta_1 \\ -\beta_2 \cdot |x|, & \text{else} \end{cases}$$



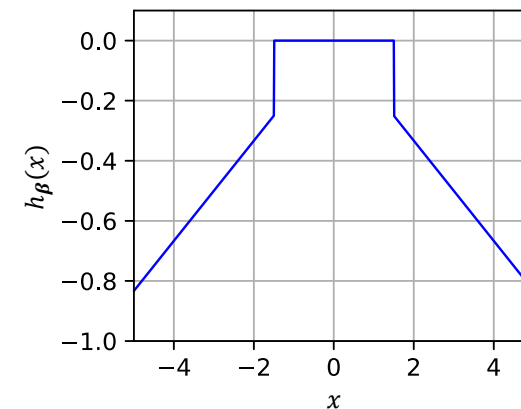
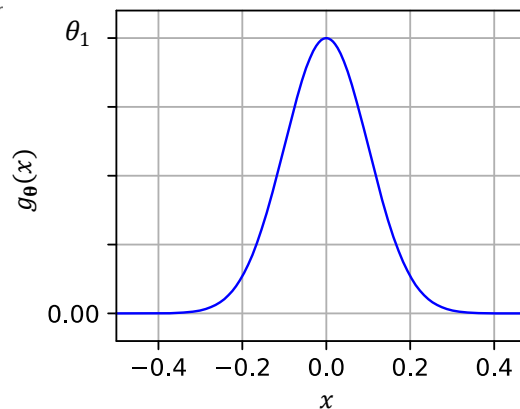
Reward Function Design (2)

$$r_{\text{RL-PFC}}(\cdot) = \underbrace{g_{\theta_e}(e_y^P) \left(1 + g_{\theta_\psi}(e_\psi^P) + g_{\theta_v}(e_{v_x}^P) \right)}_{\text{1. goal: error minimization}} + \underbrace{g_{\theta_{\Delta\delta}}(e_y^P) \left(h_{\beta_{\Delta\delta}}(\Delta\delta^f) + h_{\beta_{\Delta\delta}}(\Delta\delta^r) \right)}_{\text{2. goal: smooth steering}}$$



Reward Function Design (2)

RL-PFC:



$= 0$ for $e_y^P \rightarrow \infty$
 $= 1$ for $e_y^P = 0$

$$r_{\text{RL-PFC}}(\cdot) = g_{\theta_e}(e_y^P) \left(1 + g_{\theta_\psi}(e_\psi^P) + g_{\theta_v}(e_{v_x}^P) \right) + g_{\theta_{\Delta\delta}}(e_y^P) \left(h_{\beta_{\Delta\delta}}(\Delta\delta^f) + h_{\beta_{\Delta\delta}}(\Delta\delta^r) \right)$$

$= 0$ for $e_\psi^P \rightarrow \infty$
 $= 1$ for $e_\psi^P = 0$

$= 0$ for $e_{v_x}^P \rightarrow \infty$
 $= 1$ for $e_{v_x}^P = 0$

$= 0$ for $|\Delta\delta^i| < \beta_{\Delta\delta,1}$
 < 0 else

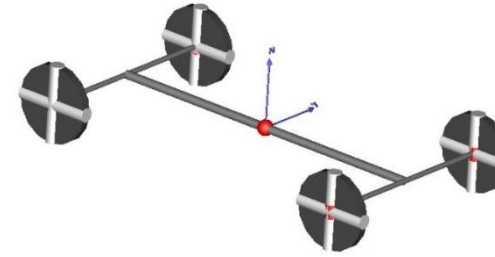
RL-PFC-TV:

$$r_{\text{RL-PFC-TV}}(\cdot) = r_{\text{RL-PFC}}(\cdot) - d_1 \sum_{i=1}^4 \left(\frac{P_s^{W_i}}{d_2} \right)^2$$

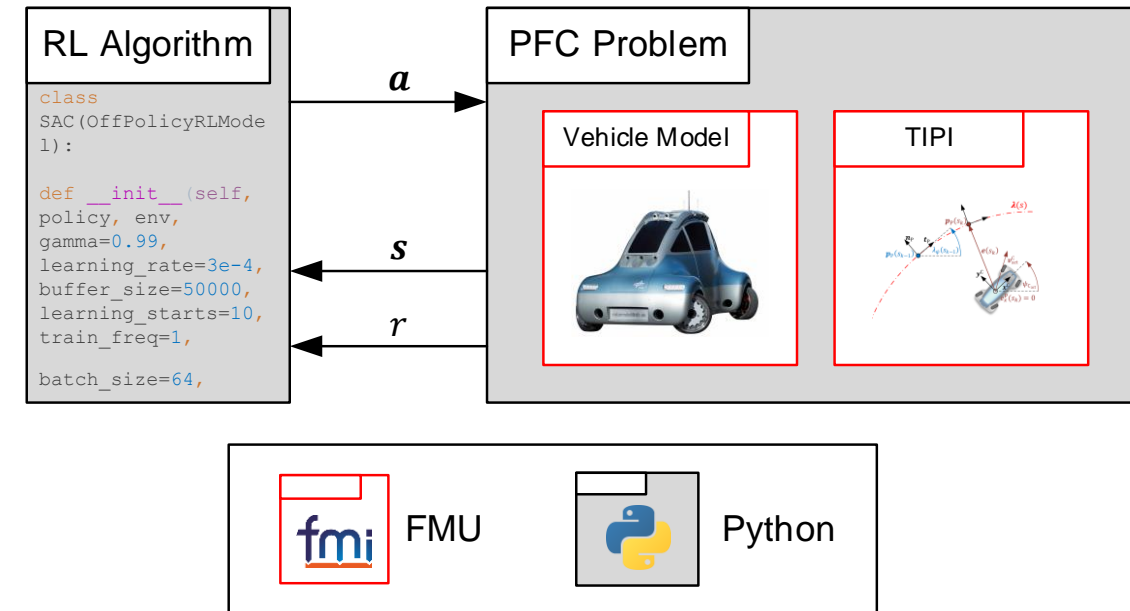
with slip-based losses: $P_s^{W_i} = F^{W_i} \cdot v_s^{W_i}$



Implementation Details

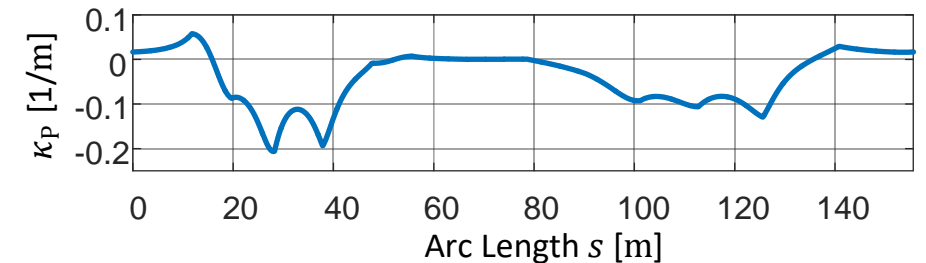
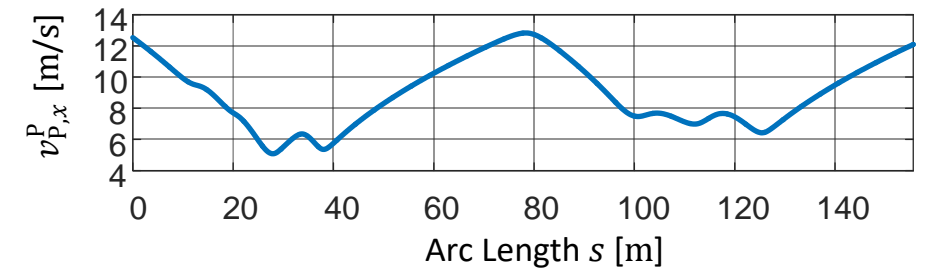
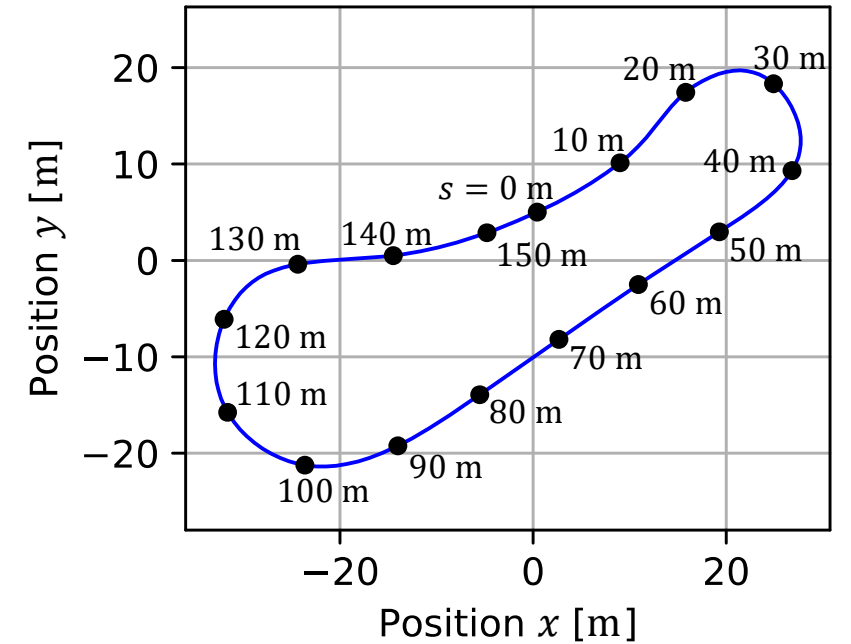
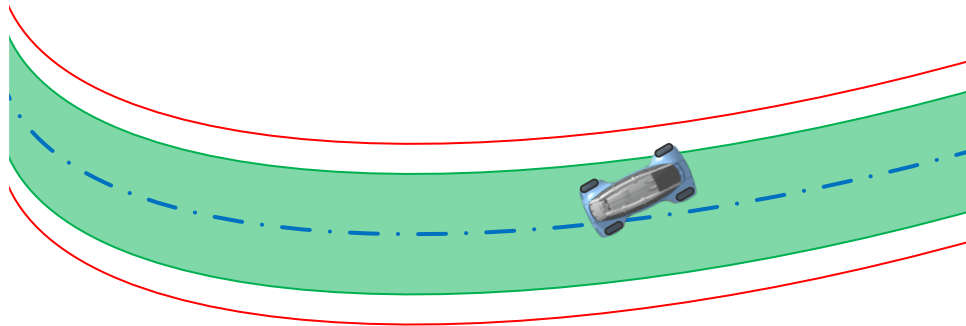


- **Aims of the training model:**
 - Fast simulation
 - Good approximation of the dominant dynamics
- The training model covers:
 - Planar two track model
 - Approximated actuator dynamics
 - Slip-based tire model
- Modular approach including 2 FMUs
- RL algorithm Soft Actor-Critic (SAC):
 - Robust to hyperparameter variations
 - High efficiency



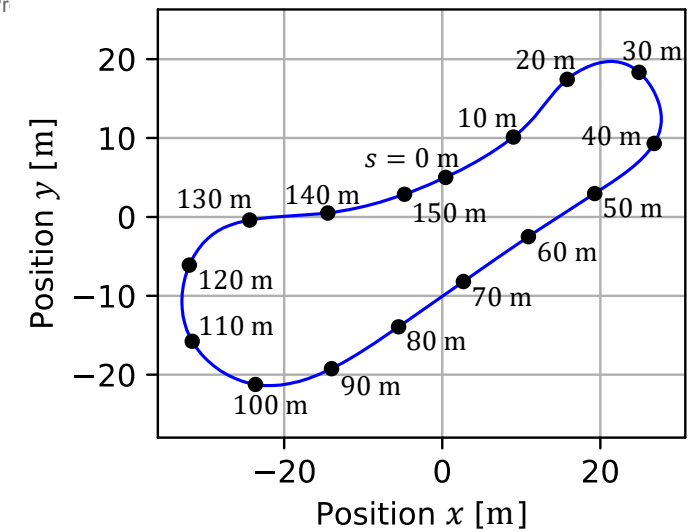
Trainingssetup

- Challenging training path (up to $5.5 \frac{\text{m}}{\text{s}^2}$ lateral acceleration, tight curve)
- Random initialization to support exploration
- Early stopping of episodes to avoid unnecessary exploration



Simulative Assessment on Training Path

Error		RL-PFC TV	RL-PFC	ref PFC
Position: e_y^P [cm]	RMSE	1.41	1.60	2.39
	max	4.19	3.55	9.99
Velocity: $e_{v_x}^P$ [m/s]	RMSE	0.165	0.174	0.0965
	max	0.413	0.526	0.383
Orientation: e_ψ [°]	RMSE	1.39	0.679	0.964
	max	3.39	1.38	2.66

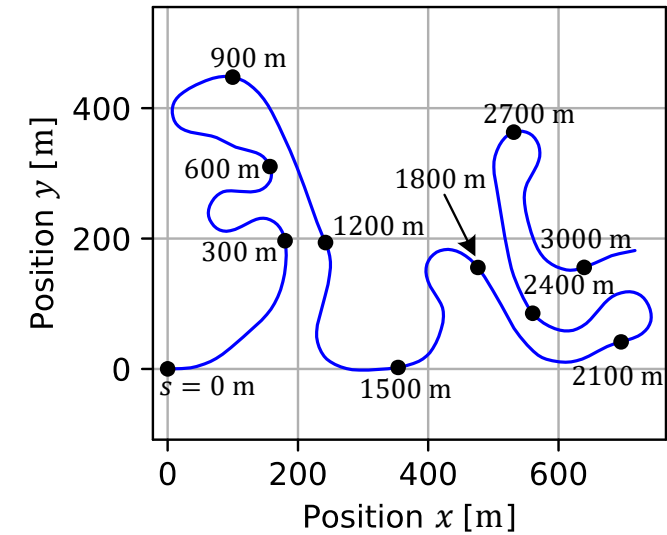


- ref PFC: 3 outer loop PD controllers combined with an optimization based control allocator (online optimization) [2]
- RL-PFC TV and RL-PFC exhibit superior position tracking
- RL-PFC exhibits superior orientation tracking
- RL-PFC controllers are designed to primary optimize position tracking
 - Velocity tracking is degraded in order to optimize position tracking

- RMSE: $\sqrt{\frac{1}{s_{\max}} \int_0^{s_{\max}} (e(s))^2 ds}$

Simulative Assessment on Unknown Path

Error		RL-PFC TV	RL-PFC	ref PFC
Position: e_y^P [cm]	RMSE	1.47	1.54	0.172
	max	2.59	3.51	0.579
Velocity: $e_{v_x}^P$ [m/s]	RMSE	0.130	0.142	0.00652
	max	0.311	0.314	0.0140
Orientation: e_ψ [°]	RMSE	0.866	0.813	0.121
	max	2.29	1.39	0.354



- ref PFC controller used for comparison performs better
- But: Both (model-free) RL-based controllers are able to follow the unknown path with an maximum position error in the cm range
- Training should cover a wide range of scenarios/ situations relevant in application

Thank you for your attention

