

Implementation of GGN and SCQP methods via first-order Taylor approximations

Alejandro Astudillo

Ph.D. researcher

MECO Research Team

August 4, 2021



Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\min_w \phi_0(c_0(w))$$

$$s.t. \quad g_i(w) = 0, \quad i = 1, \dots, p$$

$$\phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q.$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\min_w \phi_0(c_0(w))$$

$$s.t. \quad g_i(w) = 0, \quad i = 1, \dots, p$$

$$\phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q.$$

QP from SQP method

$$\min_{s_k} \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k$$

$$s.t. \quad \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p$$

$$\nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q.$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\min_w \phi_0(c_0(w))$$

$$\text{s.t. } g_i(w) = 0, \quad i = 1, \dots, p$$

$$\phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q.$$

QP from SQP method

$$\min_{s_k} \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k$$

$$\text{s.t. } \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p$$

$$\nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q.$$

Exact Hessian of the Lagrangian

$$B_k = \nabla_w^2 \mathcal{L}(w_k, \lambda_k, \mu_k).$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\min_w \phi_0(c_0(w))$$

$$\text{s.t. } g_i(w) = 0, \quad i = 1, \dots, p$$

$$\phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q.$$

QP from SQP method

$$\min_{s_k} \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k$$

$$\text{s.t. } \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p$$

$$\nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q.$$

SCQP Hessian approximation

$$B_k^{SCQP}(w, \mu) := \frac{\partial c_0(w)^T}{\partial w} \nabla_c^2 \phi_0(c_0(w)) \frac{\partial c_0(w)}{\partial w} + \sum_{i=1}^q \mu_i \frac{\partial c_i(w)^T}{\partial w} \nabla_c^2 \phi_i(c_i(w)) \frac{\partial c_i(w)}{\partial w} \succeq 0.$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\begin{aligned} \min_w & \phi_0(c_0(w)) \\ \text{s.t.} & \quad g_i(w) = 0, \quad i = 1, \dots, p \\ & \quad \phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

QP from SQP method

$$\begin{aligned} \min_{s_k} & \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k \\ \text{s.t.} & \quad \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p \\ & \quad \nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

SCQP Hessian approximation

$$\begin{aligned} B_k^{SCQP}(w, \mu) & := \frac{\partial c_0(w)^T}{\partial w} \nabla_c^2 \phi_0(c_0(w)) \frac{\partial c_0(w)}{\partial w} \\ & \quad + \sum_{i=1}^q \mu_i \frac{\partial c_i(w)^T}{\partial w} \nabla_c^2 \phi_i(c_i(w)) \frac{\partial c_i(w)}{\partial w} \succeq 0. \end{aligned}$$

First-order Taylor approximation

$$f_{lin}(w, \bar{w}) = f(\bar{w}) + \nabla f(\bar{w})^T (w - \bar{w})$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\begin{aligned} \min_w \quad & \phi_0(c_0(w)) \\ \text{s.t.} \quad & g_i(w) = 0, \quad i = 1, \dots, p \\ & \phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

QP from SQP method

$$\begin{aligned} \min_{s_k} \quad & \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k \\ \text{s.t.} \quad & \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p \\ & \nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

SCQP Hessian approximation

$$\begin{aligned} B_k^{SCQP}(w, \mu) := & \frac{\partial c_0(w)}{\partial w}^T \nabla_c^2 \phi_0(c_0(w)) \frac{\partial c_0(w)}{\partial w} \\ & + \sum_{i=1}^q \mu_i \frac{\partial c_i(w)}{\partial w}^T \nabla_c^2 \phi_i(c_i(w)) \frac{\partial c_i(w)}{\partial w} \succeq 0. \end{aligned}$$

First-order Taylor approximation

$$f_{lin}(w, \bar{w}) = f(\bar{w}) + \nabla f(\bar{w})^T (w - \bar{w})$$

Trick in evaluation with the same symbolic variable

$$\left. \begin{aligned} \tilde{f}(w) &:= f_{lin}(w, w) = f(w), \\ \nabla \tilde{f}(w) &:= \nabla f(w), \\ \nabla^2 \tilde{f}(w) &:= 0 \end{aligned} \right\} \text{lin}(f(w))$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\begin{aligned} \min_w \quad & \phi_0(c_0(w)) \\ \text{s.t.} \quad & g_i(w) = 0, \quad i = 1, \dots, p \\ & \phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

QP from SQP method

$$\begin{aligned} \min_{s_k} \quad & \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k \\ \text{s.t.} \quad & \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p \\ & \nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

SCQP Hessian approximation

$$\begin{aligned} B_k^{SCQP}(w, \mu) := & \frac{\partial c_0(w)}{\partial w}^T \nabla_c^2 \phi_0(c_0(w)) \frac{\partial c_0(w)}{\partial w} \\ & + \sum_{i=1}^q \mu_i \frac{\partial c_i(w)}{\partial w}^T \nabla_c^2 \phi_i(c_i(w)) \frac{\partial c_i(w)}{\partial w} \succeq 0. \end{aligned}$$

First-order Taylor approximation

$$f_{lin}(w, \bar{w}) = f(\bar{w}) + \nabla f(\bar{w})^T (w - \bar{w})$$

Trick in evaluation with the same symbolic variable

$$\left. \begin{aligned} \tilde{f}(w) &:= f_{lin}(w, w) = f(w), \\ \nabla \tilde{f}(w) &:= \nabla f(w), \\ \nabla^2 \tilde{f}(w) &:= 0 \end{aligned} \right\} \text{lin}(f(w))$$

Ready to be solved with SCQP

$$\begin{aligned} \min_w \quad & \phi_0(\text{lin}(c_0(w))) \\ \text{s.t.} \quad & \text{lin}(g_i(w)) = 0, \quad i = 1, \dots, p \\ & \phi_i(\text{lin}(c_i(w))) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

Problem formulation

Optimization problem with
convex-over-nonlinear constraints

$$\begin{aligned} \min_w \quad & \phi_0(c_0(w)) \\ \text{s.t.} \quad & g_i(w) = 0, \quad i = 1, \dots, p \\ & \phi_i(c_i(w)) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

QP from SQP method

$$\begin{aligned} \min_{s_k} \quad & \frac{1}{2} s_k^T B_k s_k + \nabla \psi_0(w_k)^T s_k \\ \text{s.t.} \quad & \nabla g_i(w_k)^T s_k + g_i(w_k) = 0, \quad i = 1, \dots, p \\ & \nabla \psi_i(w_k)^T s_k + \psi_i(w_k) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

SCQP Hessian approximation

$$\begin{aligned} B_k^{SCQP}(w, \mu) := & \frac{\partial c_0(w)}{\partial w}^T \nabla_c^2 \phi_0(c_0(w)) \frac{\partial c_0(w)}{\partial w} \\ & + \sum_{i=1}^q \mu_i \frac{\partial c_i(w)}{\partial w}^T \nabla_c^2 \phi_i(c_i(w)) \frac{\partial c_i(w)}{\partial w} \succeq 0. \end{aligned}$$

First-order Taylor approximation

$$f_{lin}(w, \bar{w}) = f(\bar{w}) + \nabla f(\bar{w})^T (w - \bar{w})$$

Trick in evaluation with the same symbolic variable

$$\left. \begin{aligned} \tilde{f}(w) &:= f_{lin}(w, w) = f(w), \\ \nabla \tilde{f}(w) &:= \nabla f(w), \\ \nabla^2 \tilde{f}(w) &:= 0 \end{aligned} \right\} \text{lin}(f(w))$$

Ready to be solved with GGN

$$\begin{aligned} \min_w \quad & \phi_0(\text{lin}(c_0(w))) \\ \text{s.t.} \quad & \text{lin}(g_i(w)) = 0, \quad i = 1, \dots, p \\ & \underline{\text{lin}}(\phi_i(c_i(w))) \leq 0, \quad i = 1, \dots, q. \end{aligned}$$

Results

Tunnel-following NMPC for a Robot Manipulator

$$\min \int_0^{t_f} \left\| \begin{bmatrix} \dot{\theta}(\tau) - \dot{\theta}_{ref}(\theta(\tau)) \\ e_T(q(\tau), \theta(\tau)) \\ x(\tau) \\ u(\tau) \\ \nu(\tau) \end{bmatrix} \right\|_{\mathbf{W}}^2 + \alpha_l l(\tau) + \alpha_{l_O} l_O(\tau) d\tau$$

s.t. $0 = f_x(\dot{x}(\tau), x(\tau), u(\tau)),$
 $\dot{\zeta}(\tau) = f_\zeta(\zeta(\tau), \nu(\tau)),$
 $\underline{r}(\tau) \leq r(x(\tau), u(\tau)) \leq \bar{r}(\tau),$
 $\underline{q}(\tau) \leq q(\tau) \leq \bar{q}(\tau),$
 $\|e_p(q(\tau), \theta(\tau))\|^2 \leq \rho^2 + l(\tau),$
 $\|e_O(q(\tau), \theta(\tau))\|^2 \leq \rho_O^2 + l_O(\tau),$
 $l(\tau) \geq 0,$
 $l_O(\tau) \geq 0,$
 $0 \leq \theta(\tau) \leq \bar{\theta},$
 $x(0) = \hat{x}_0,$
 $(x(t_f), \zeta(t_f)) \in \varepsilon.$

Results

Tunnel-following NMPC for a Robot Manipulator

$$\min \int_0^{t_f} \left\| \begin{bmatrix} \dot{\theta}(\tau) - \dot{\theta}_{ref}(\theta(\tau)) \\ e_T(q(\tau), \theta(\tau)) \\ x(\tau) \\ u(\tau) \\ \nu(\tau) \end{bmatrix} \right\|_{\mathbf{W}}^2 + \alpha_l l(\tau) + \alpha_{l_O} l_O(\tau) d\tau$$

s.t.

$$0 = f_x(\dot{x}(\tau), x(\tau), u(\tau)),$$
$$\dot{\zeta}(\tau) = f_\zeta(\zeta(\tau), \nu(\tau)),$$
$$\underline{r}(\tau) \leq r(x(\tau), u(\tau)) \leq \bar{r}(\tau),$$
$$\underline{q}(\tau) \leq q(\tau) \leq \bar{q}(\tau),$$
$$\|e_p(q(\tau), \theta(\tau))\|^2 \leq \rho^2 + l(\tau),$$
$$\|e_O(q(\tau), \theta(\tau))\|^2 \leq \rho_O^2 + l_O(\tau),$$
$$l(\tau) \geq 0,$$
$$l_O(\tau) \geq 0,$$
$$0 \leq \theta(\tau) \leq \bar{\theta},$$
$$x(0) = \hat{x}_0,$$
$$(x(t_f), \zeta(t_f)) \in \varepsilon.$$

Results

Tunnel-following NMPC for a Robot Manipulator

$$\min \int_0^{t_f} \left\| \begin{bmatrix} \dot{\theta}(\tau) - \dot{\theta}_{ref}(\theta(\tau)) \\ e_T(q(\tau), \theta(\tau)) \\ x(\tau) \\ u(\tau) \\ \nu(\tau) \end{bmatrix} \right\|_{\mathbf{W}}^2 + \alpha_l l(\tau) + \alpha_{l_O} l_O(\tau) d\tau$$

s.t.

$$0 = f_x(\dot{x}(\tau), x(\tau), u(\tau)),$$

$$\dot{\zeta}(\tau) = f_\zeta(\zeta(\tau), \nu(\tau)),$$

$$\underline{r}(\tau) \leq r(x(\tau), u(\tau)) \leq \bar{r}(\tau),$$

$$\underline{q}(\tau) \leq q(\tau) \leq \bar{q}(\tau),$$

$$\|e_p(q(\tau), \theta(\tau))\|^2 \leq \rho^2 + l(\tau),$$

$$\|e_O(q(\tau), \theta(\tau))\|^2 \leq \rho_O^2 + l_O(\tau),$$

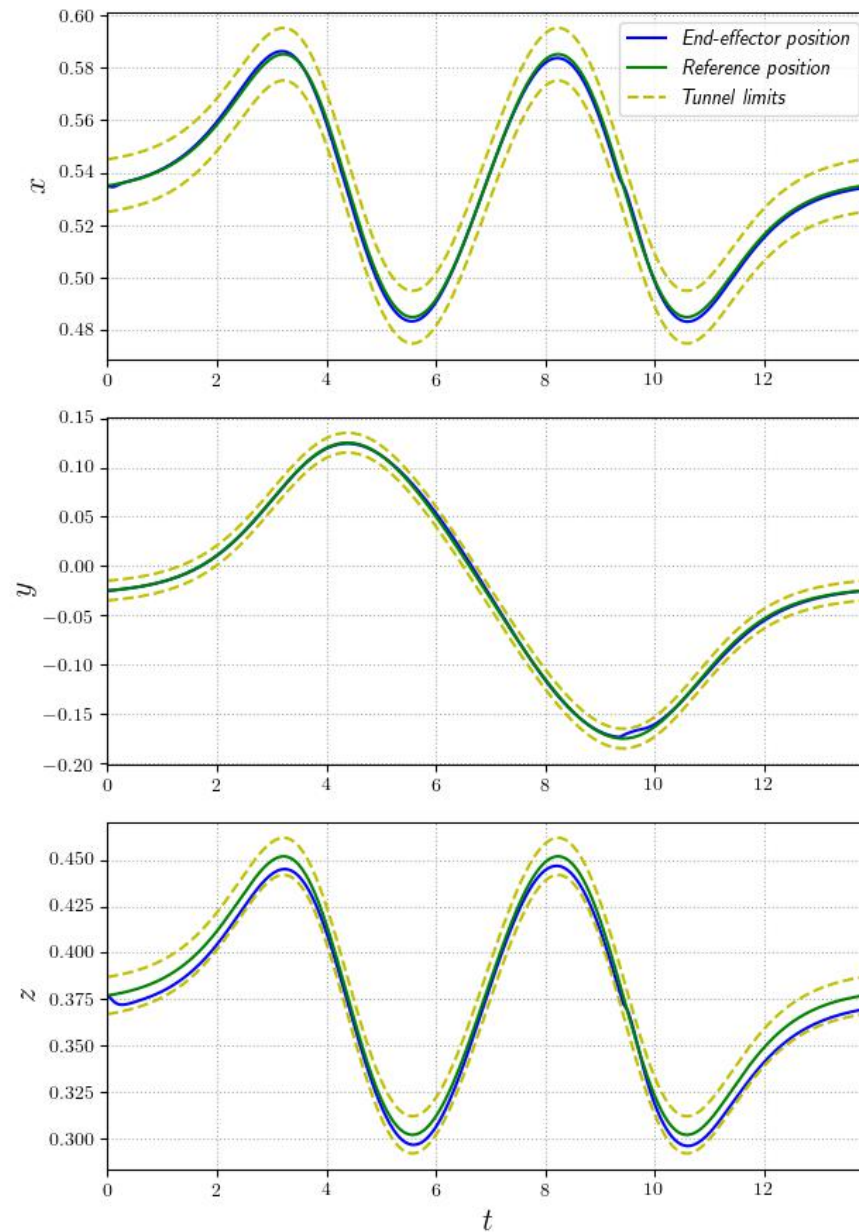
$$l(\tau) \geq 0,$$

$$l_O(\tau) \geq 0,$$

$$0 \leq \theta(\tau) \leq \bar{\theta},$$

$$x(0) = \hat{x}_0,$$

$$(x(t_f), \zeta(t_f)) \in \varepsilon.$$

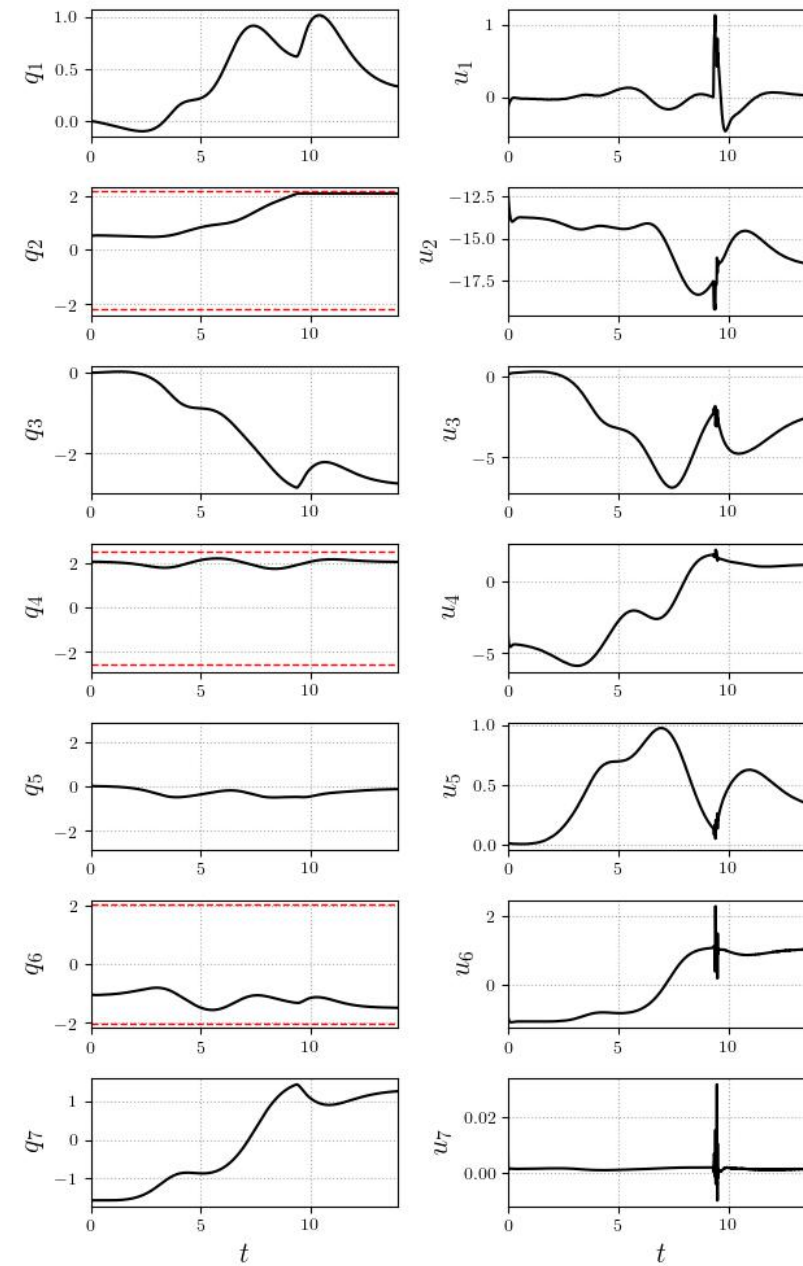
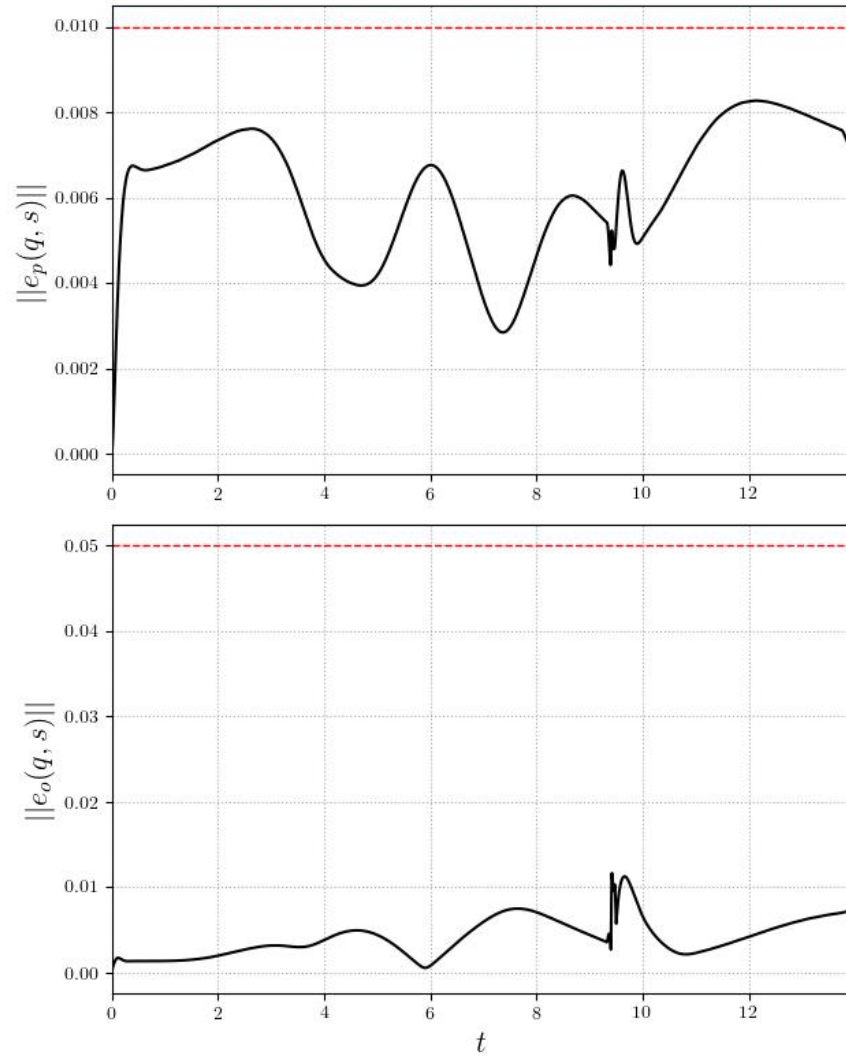


Thank you!

e-mail: alejandro.astudillovigoya@kuleuven.be

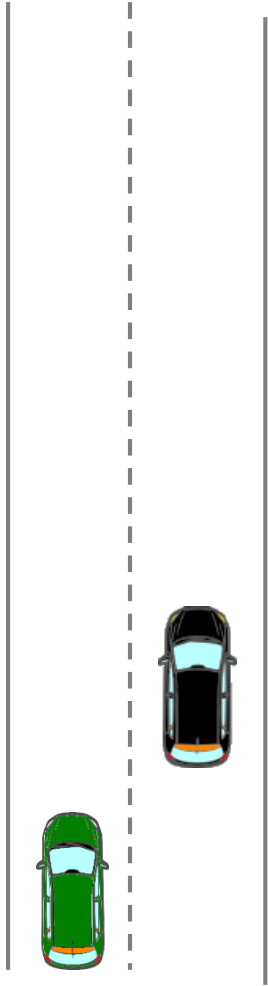


Extra



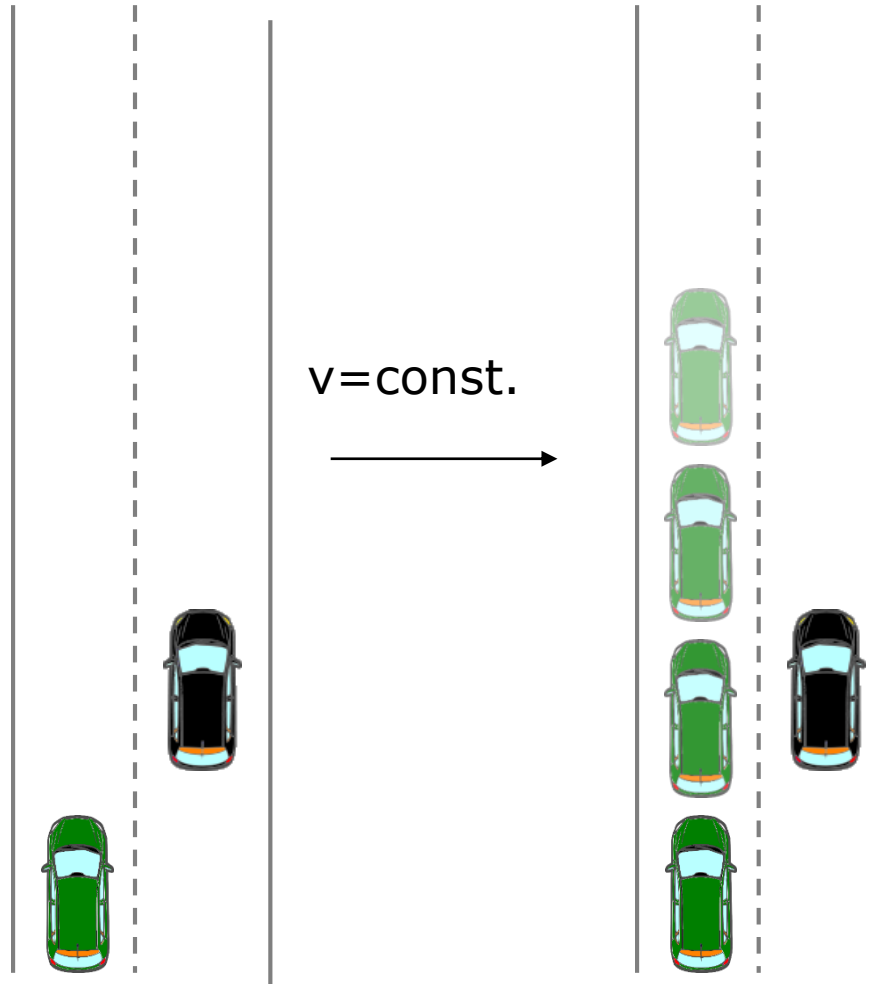
Interaction aware trajectory planning for automated vehicles

The Problem



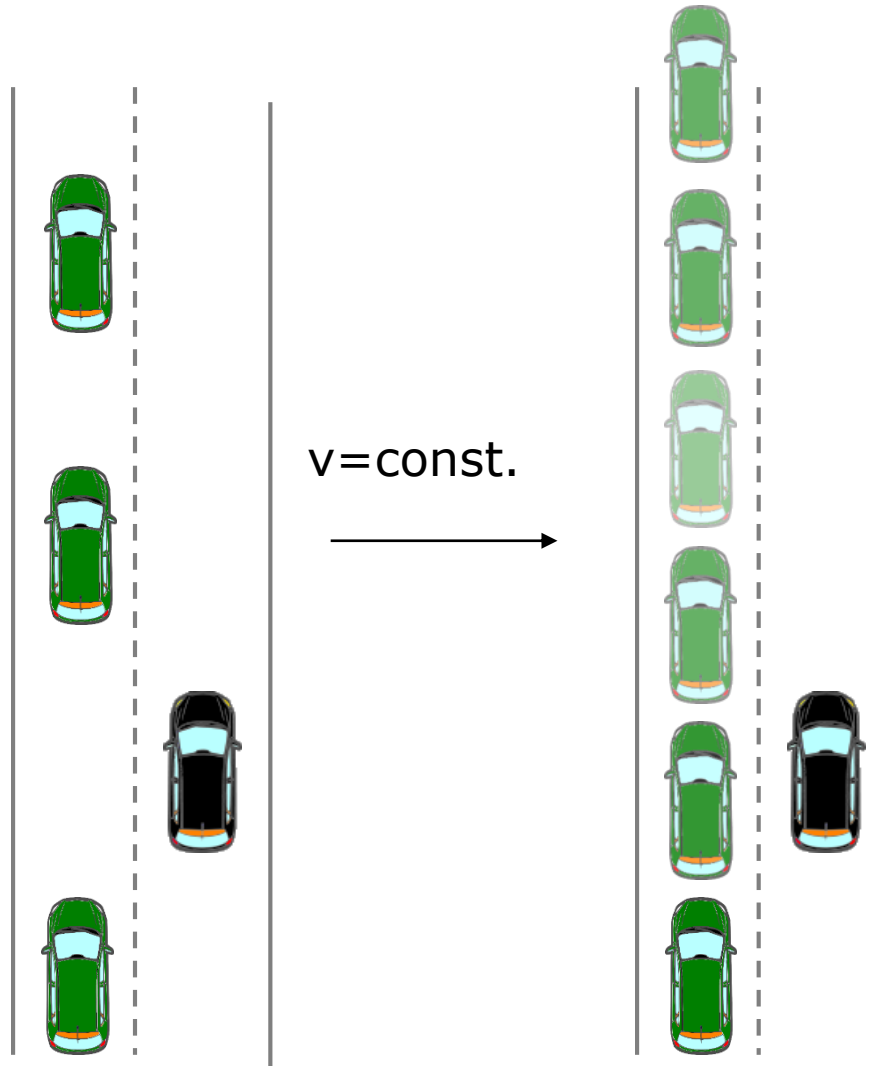
Interaction aware trajectory planning for automated vehicles

The Problem



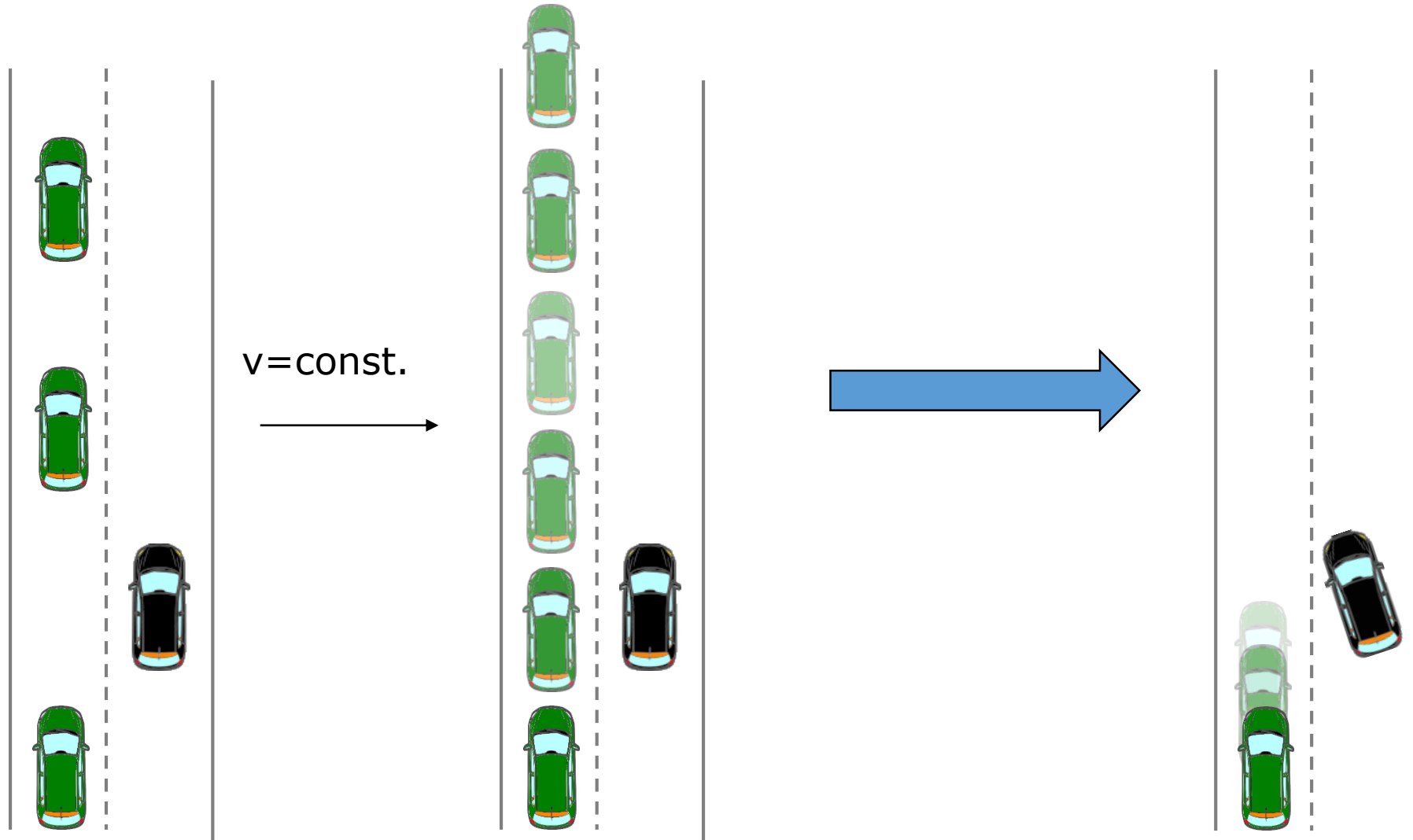
Interaction aware trajectory planning for automated vehicles

The Problem

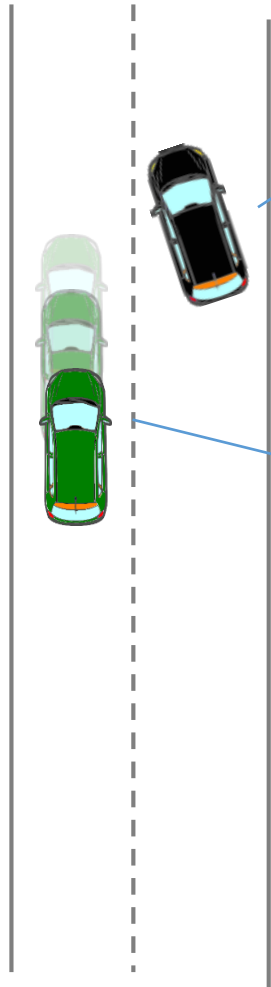


Interaction aware trajectory planning for automated vehicles

The Problem



Interaction aware trajectory planning for automated vehicles



$$F(x, y)$$

$$G(x, y)$$

$$f(x, y)$$

$$g(x, y)$$

Dynamic model

$$\vec{x} = \begin{bmatrix} x \\ y \\ \psi \\ v \end{bmatrix} \quad f = \begin{bmatrix} \dot{x} \\ \dot{y} \\ \dot{\psi} \\ \dot{v} \end{bmatrix} = \begin{bmatrix} v \cos(\psi + \beta) \\ v \sin(\psi + \beta) \\ \frac{v}{l} \tan(\delta) \cos(\beta) \\ a_x \end{bmatrix}$$

$$\begin{aligned} & \min F(x, y) \\ & \text{s.t. } G(x, y) \leq 0 \\ & \quad \gamma \in \{ \arg \min f(x, y); g(x, y) \leq 0 \} \end{aligned}$$

MPCC sup \approx
min $\|F^1(x,y)\|_2^2 + P^2(x,y) + \underbrace{J^T z \cdot w}_{\text{circled}}$

x, y, z, w

s.t. $M \cdot \bar{s}_0 + G(x,y) \leq 0$

$$0 = \underbrace{\nabla_y F(x,y) + \nabla_y g(x,y) \cdot z}$$

$$:= H(x,y,z)$$

$0 \leq z \perp w \geq 0$

$$0 = w + g(x,y)$$

Interaction aware trajectory planning for automated vehicles

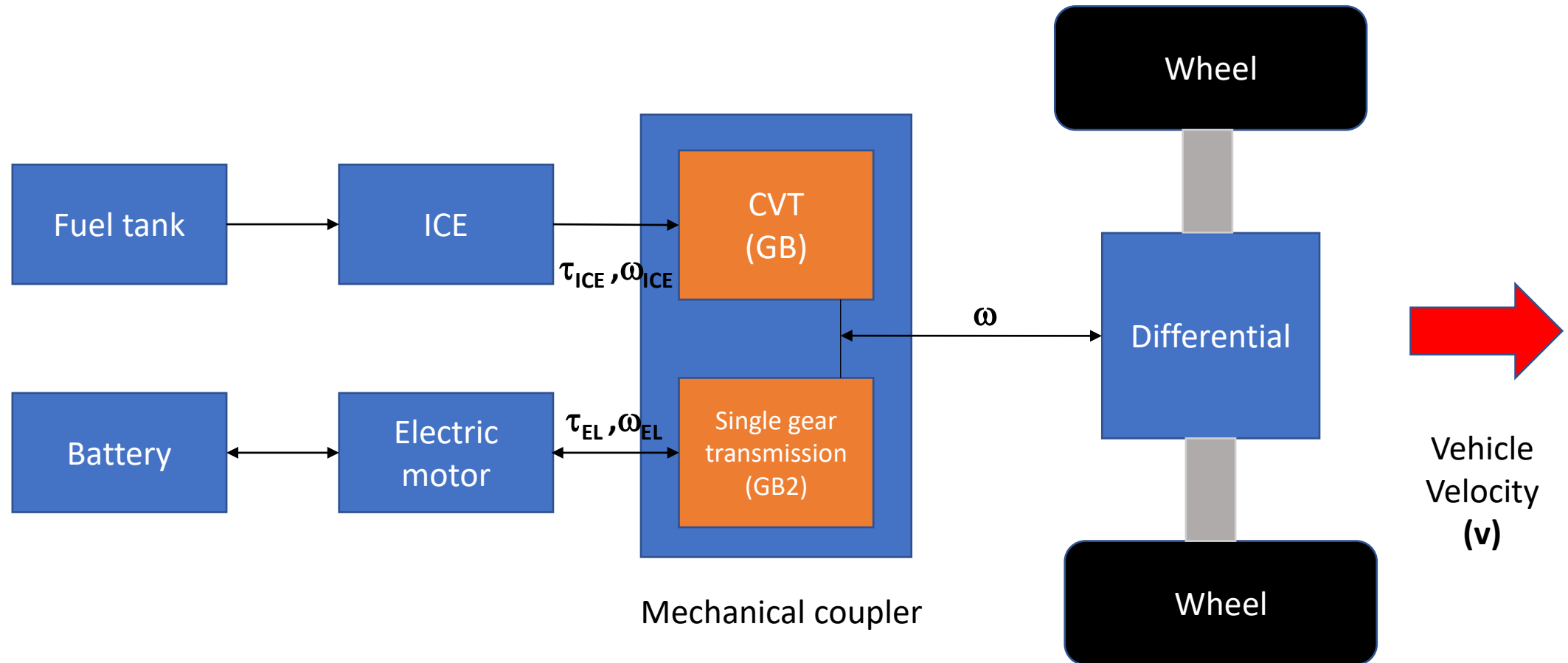
Video

Hybrid powertrain optimization

Manuel Antonio Perez Estevez
mperezestevez@unibz.it

Freiburg, 04/08/2021

Parallel hybrid powertrain simplified model



How to control the system to minimize greenhouse emissions?

Application problem → MPC

Model

$$\frac{dv(t)}{dt} = \frac{\frac{\tau_{ICE}(t)}{r} GB(t) + \frac{\tau_{EL}(t)}{r} GB_2(t) - \frac{1}{2} \rho_{air} A_f C_d v(t) - (\mu \cos(\alpha) + \sin(\alpha)) (g (m_v + N_p N_s m_{cell} + \frac{P_p}{\rho_{ICE}}))}{m_v + N_p N_s m_{cell} + \frac{P_p}{\rho_{ICE}}}$$

$$\frac{dSOC(t)}{dt} = -100 \frac{GB_2(t) v(t) \tau_{EL}(t)}{Cap_{cell} N_p N_s r \eta_{EL}(t) u(t) 3600}$$

$$\eta_{ICE}(t) = \eta_0 (a_1 + b_1 \left(\frac{v(t)}{n_p r} GB(t) \right) + c_1 \left(\frac{v(t)}{n_p r} GB(t) \right)^2 + d_1 \left(\frac{v(t)}{n_p r} GB(t) \right)^3) (a_2 + b_2(A) + c_2(A)^2 + d_2(A)^3)$$

Where:

$$A = \frac{v(t) GB(t) \tau_{ICE}(t)}{P_e(t) r}$$

$$P_e(t) = P_p \left(a_3 \left(\frac{v(t)}{n_p r} GB(t) \right) + b_3 \left(\frac{v(t)}{n_p r} GB(t) \right)^2 + c_3 \left(\frac{v(t)}{n_p r} GB(t) \right)^3 \right) + \varepsilon$$

$$\eta_{EL}(t) = LT(\tau_{EL}, v(t))$$

$$u(t) = u_{nom} (a_4 + b_4(SOC(t)) + c_4(SOC(t))^2 + d_4(SOC(t))^3 + e_4(SOC(t))^4 + f_4(SOC(t))^5) \dots - R \frac{v(t) GB_2(t) \tau_{EL}(t)}{N_p N_s r u(t) \eta_{EL}(t)}$$

Differential states

$$x(t) = \begin{bmatrix} v(t) \\ SOC(t) \end{bmatrix}$$

Algebraic states

$$z(t) = \begin{bmatrix} \eta_{ICE}(t) \\ \eta_{EL}(t) \\ u(t) \end{bmatrix}$$

Control actions

$$a(t) = \begin{bmatrix} \tau_{ICE}(t) \\ \tau_{EL}(t) \\ GB(t) \end{bmatrix}$$

Optimization problem

$$\min_{x,z,a} \frac{Coef_1}{r HV} \int_0^{t_{N-1}} \frac{v(t)GB(t)\tau_{ICE}(t)}{(\eta_{ICE}(t) + \varepsilon)} dt + (SOC_{target} - SOC_N)Cap_{cell}N_pN_s u_{nom}Coef_2$$

subject to

$$f_{impl}(x(t), z(t), a(t)) = 0$$

Constraints:

$$v(t) = v_{ref}(t)$$

$$\tau_{ICE}(t) \geq 0$$

$$0 \leq \eta_{EL}(t) \leq 1$$

$$2.5 \leq u(t) \leq 4.3$$

$$-200 \leq \tau_{EL}(t) \leq 180$$

$$0 \leq \eta_{ICE}(t) \leq 1$$

$$0 \leq SOC(t) \leq 100$$

$$\frac{\tau_{ICE}(t)v(t)GB(t)}{r} \leq P_p * 1.15$$

$$0 \leq GB(t) \leq 7$$

Differential states

$$x(t) = \begin{bmatrix} v(t) \\ SOC(t) \end{bmatrix}$$

Algebraic states

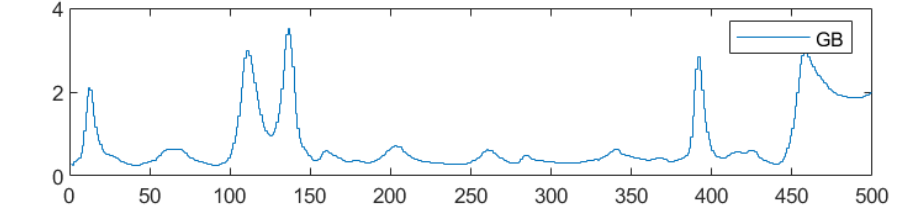
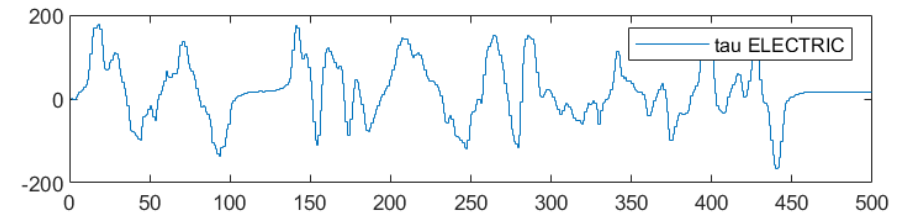
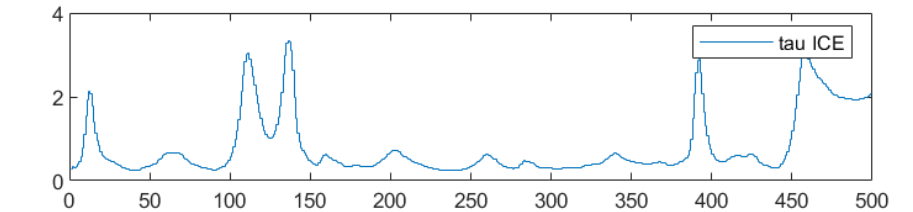
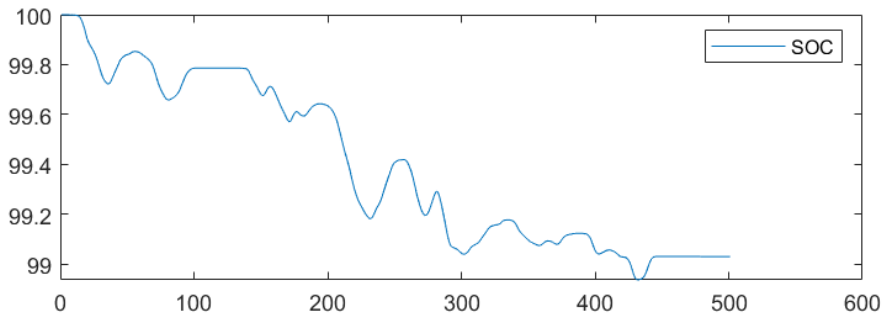
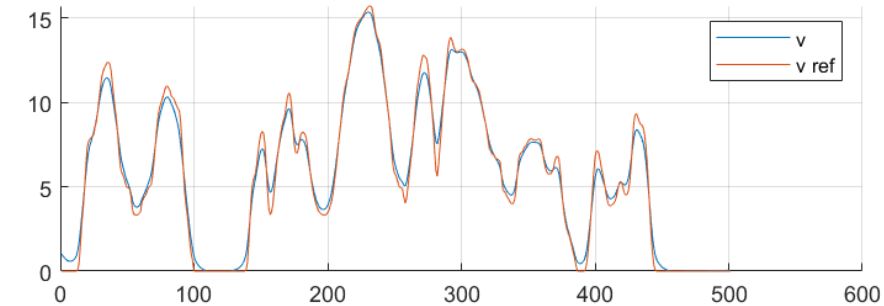
$$z(t) = \begin{bmatrix} \eta_{ICE}(t) \\ \eta_{EL}(t) \\ u(t) \end{bmatrix}$$

Control actions

$$a(t) = \begin{bmatrix} \tau_{ICE}(t) \\ \tau_{EL}(t) \\ GB(t) \end{bmatrix}$$

Preliminary results

$$\min_{x,z,a} \frac{Coef_1}{r HV} \int_0^{t_{N-1}} \left(\frac{v(t)GB(t)\tau_{ICE}(t)}{(\eta_{ICE}(t) + \epsilon)} + \beta (v(t) - v_{ref})^2 \right) dt$$



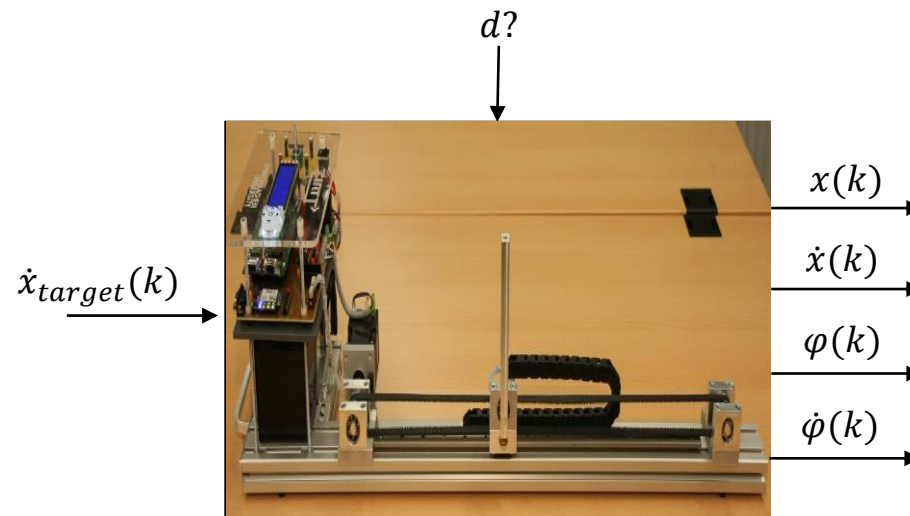
Control demonstrator cartpole with RL (TD3)

Manipulated Variable

\dot{x}_{target} Target velocity of speed controlled step motor

Challenges

Continuous state space
Continuous action space



$dT = 0.01 \text{ sec}$

Controlled Variable

x Cart position

\dot{x} Cart velocity

φ Pendulum angle

$\dot{\varphi}$ Angular velocity of pendulum

Control Goal

Translates to

Swing up

-1

Maintain upright position

-1

Stay within available space

Barrier?

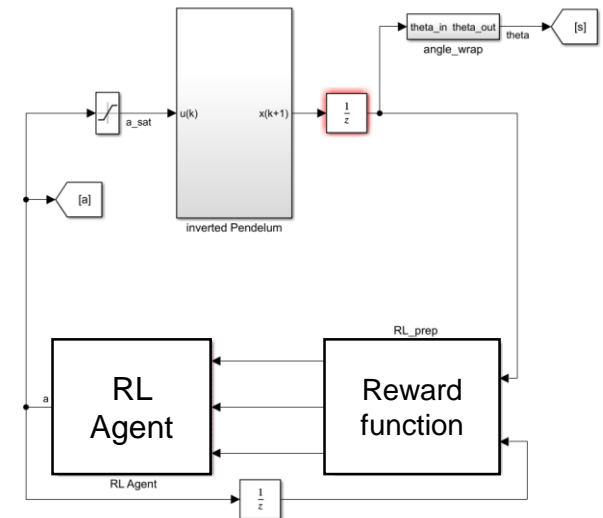
Control Problem (Maintain upright position)

If stays in boundaries : $Reward = 1 - x^2 + \cos(\theta)^2 - \dot{x}^2 - \dot{\theta}^2$, else $Reward = -100$

$$\text{s.t. } \frac{d}{dt} \begin{pmatrix} x \\ \dot{x} \\ \varphi \\ \dot{\varphi} \end{pmatrix} = \begin{pmatrix} \frac{K \cdot \dot{x}_{target} - \dot{x}}{T} \\ \dot{\varphi} \\ \frac{m \cdot g \cdot L \cdot \sin(\varphi)}{2 \cdot J} - \dot{x} \cdot \cos(\varphi) - \mu \cdot \dot{\varphi} \end{pmatrix} \quad \begin{matrix} -0.25 \text{ m} < x_{min} < 0.25 \text{ m} \\ -12^\circ < \varphi < 12^\circ \end{matrix}$$

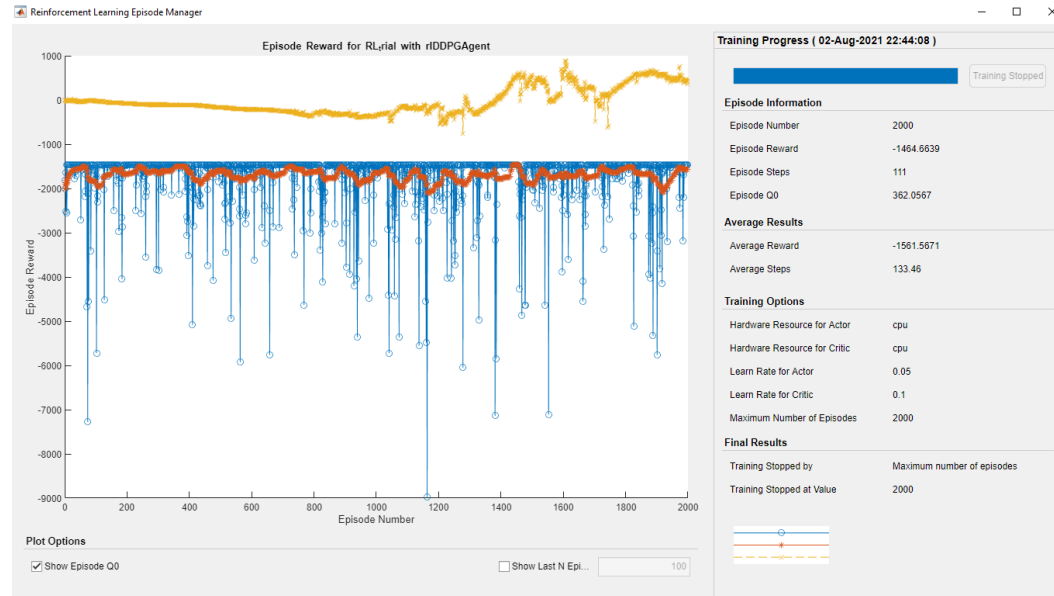
Procedure

- Matlab RL Toolbox
 - *Training time > 3 h, not converged*
- *Baseline3: Pythorch implementation RL*
 - *training time 10 Min*
 - *train cardpole gym environmet (exercise 8)*
 - *train IRT gym env -> not successful yet*
 - *export agent Neuronal Network to Matlab to control demonstrater -> to do*

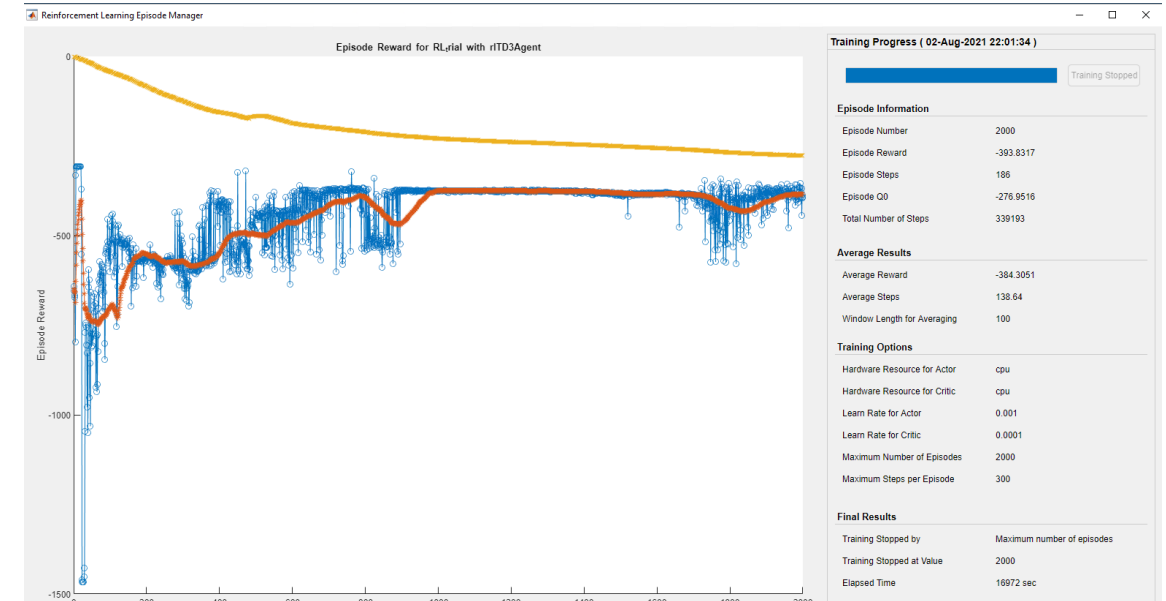


Training in Simulink

DDPG



TD3



Some results (1)

$$\text{Reward} = 1$$

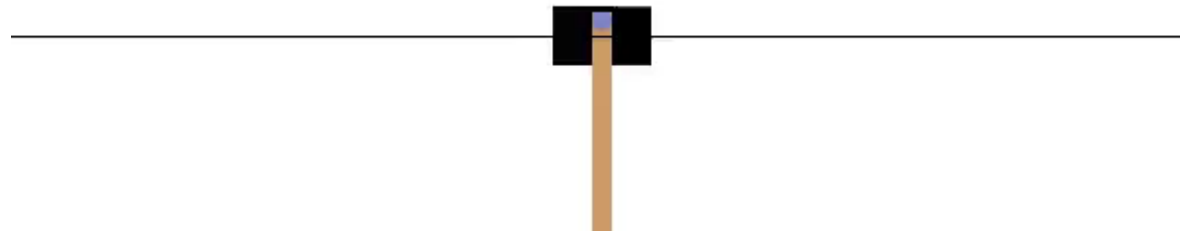


$$\text{Reward} = 1 - x^2 + \cos(\theta)^2 - \dot{x}^2 - \dot{\theta}^2$$



Some results (2)

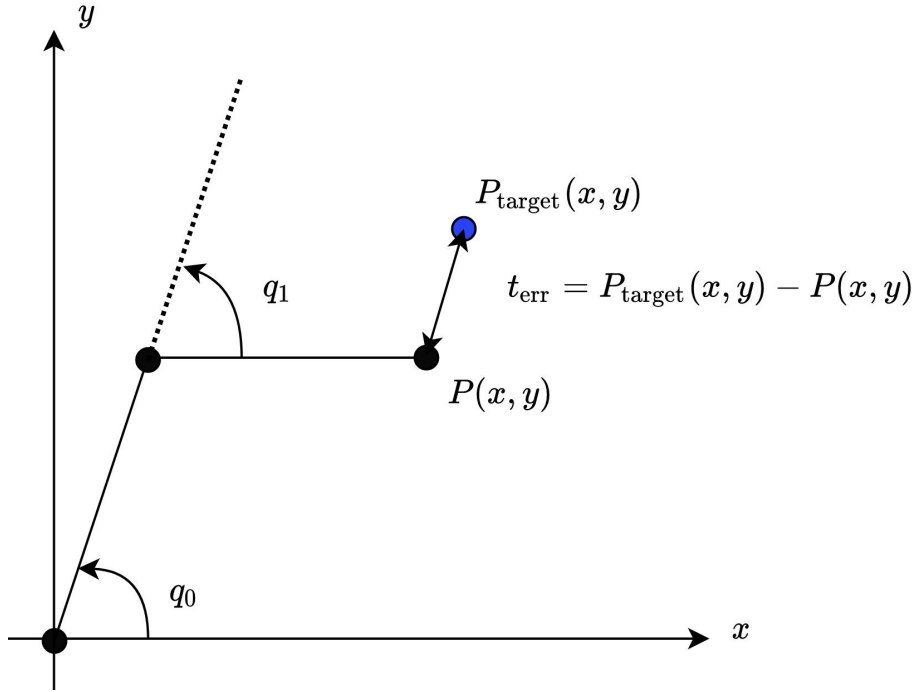
$$\text{Reward} = 1 - x^2 + \cos(\theta)^2 - \dot{x}^2 - \dot{\theta}^2$$
$$s_0 = (0, 0, \pi, 0)^T$$



Learning point to point motions for robot manipulators using RL

Johan & Ruan

Problem formulation



State space:

$$s = [q_0, q_1, \dot{q}_0, \dot{q}_1]$$

$$q_0, q_1 \in (-\pi, \pi]$$

$$\dot{q}_0, \dot{q}_1 \in [-8, 8]$$

Action space:

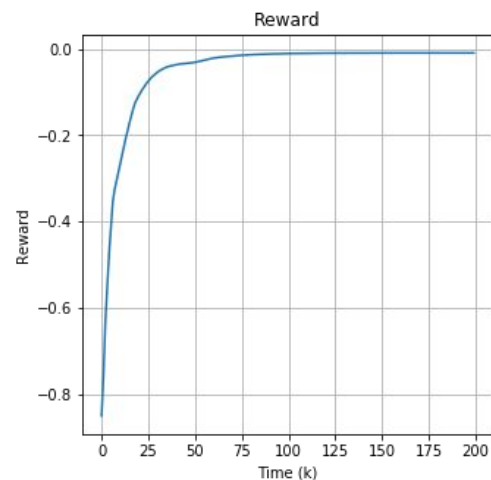
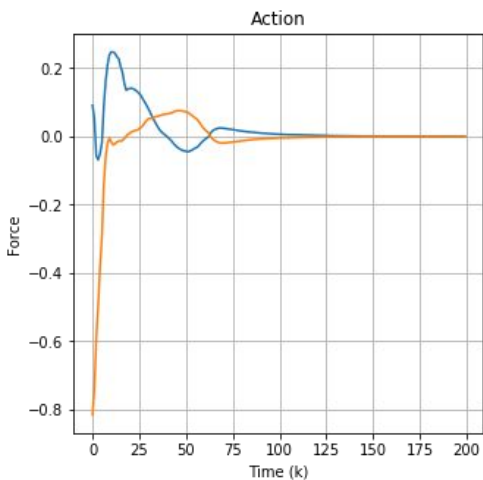
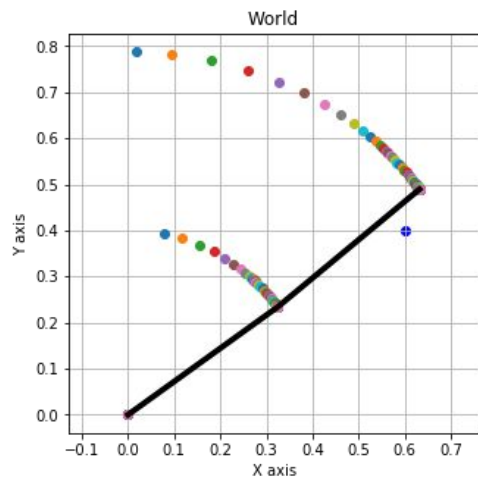
Continuous: $\ddot{q}_0, \ddot{q}_1 \in (-1.0, 1.0]$

Reward:

$$r = -(\|t_{\text{err}}\|_2^2 + 0.5\|a\|_2^2)$$

Results

Continuous: Soft Actor Critic



References:

- Soft Actor Critic algorithm: <https://stable-baselines.readthedocs.io/en/master/modules/sac.html>

Mobile robot parking with obstacles using RL

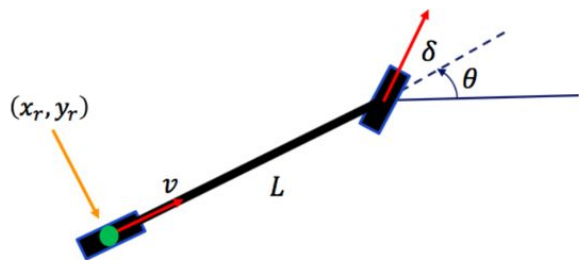
Federico Ulloa Rios

Santiago Iregui Rincón

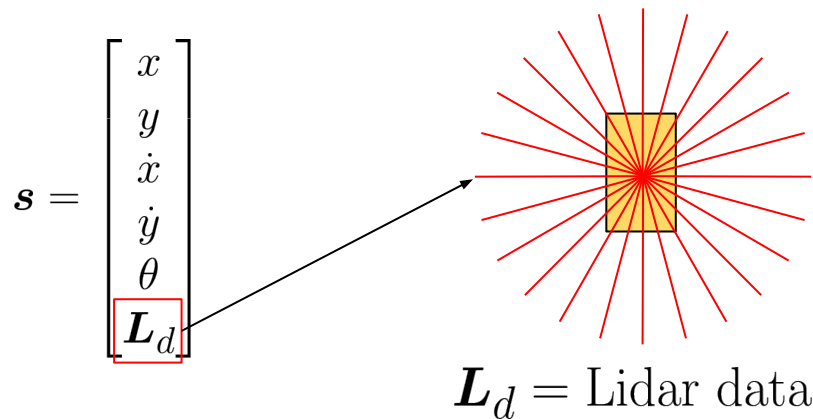
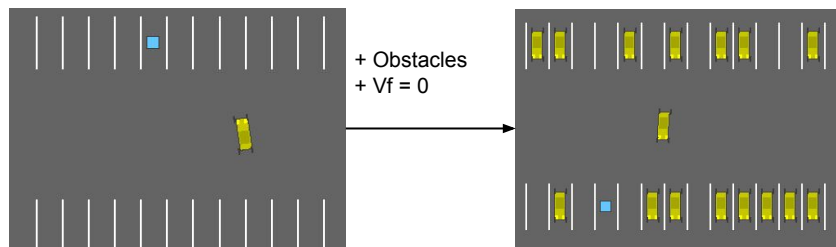


Problem Formulation

Deterministic dynamics (Bicycle model)



Modified GYM Environment



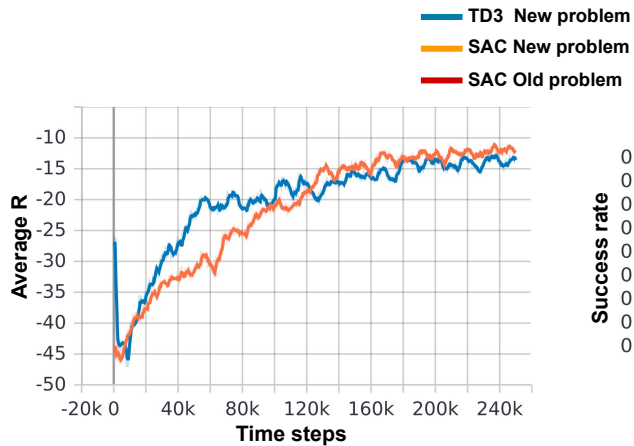
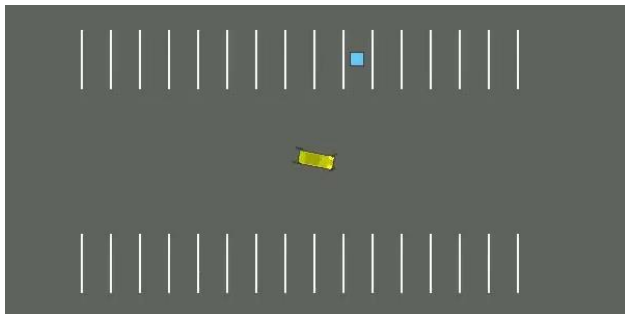
$L_d = \text{Lidar data}$

$$a = \begin{bmatrix} \alpha \\ \delta \end{bmatrix} \quad R = -(|s_r - s_t| \cdot w)^p - 5 * \psi$$

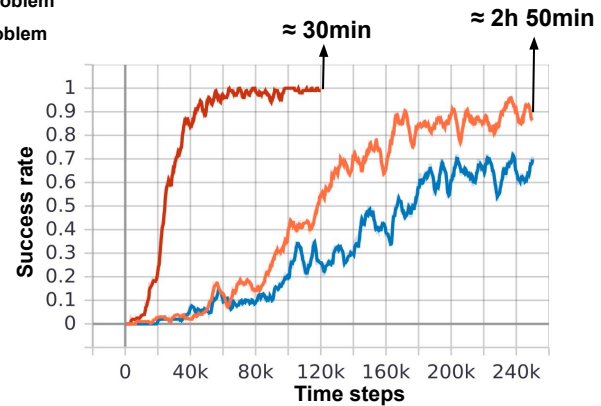
- Continuous action and state space
- Off-policy learning actor critic methods (SAC & TD3)

Results

SAC Old problem

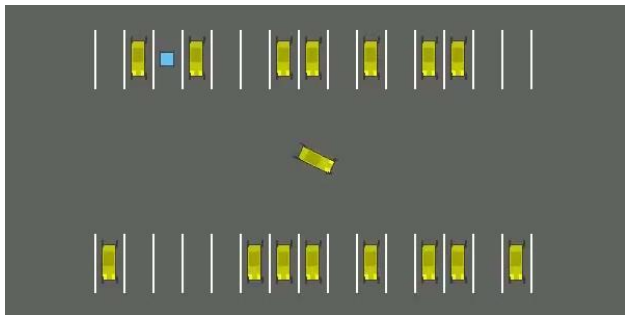


Training time @ i7-8850H & NVIDIA Quadro P600

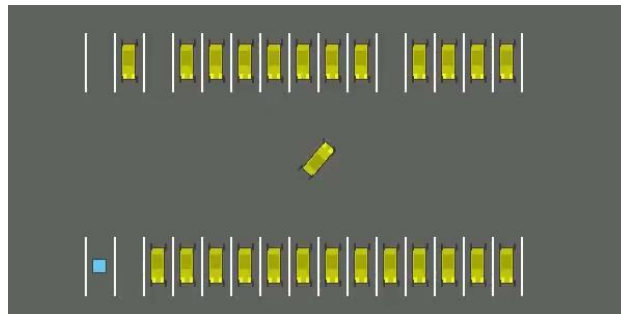


≈ 30min

≈ 2h 50min



SAC New problem



Model Predictive Control and Reinforcement Learning

Quadrotor Control using MPC

Date: 04.08.2021

Sourish Pramanick
Niket Ahuja
Nayana Koneru

Problem Formulation

$$s = [x, y, z, \phi, \theta, \psi, u, v, w, p, q, r]^T$$

$$u = [f_t, \tau_x, \tau_y, \tau_z]^T$$

$$C(s, a) = \sum_{k=0}^{N-1} (s_k - s_{ref,k})^T Q (s_k - s_{ref,k}) + a^T R a$$

$$\dot{\phi} = p + r[c(\phi)t(\theta)] + q[s(\phi)t(\theta)]$$

$$\dot{\theta} = q[c(\phi)] - r[s(\phi)]$$

$$\dot{\psi} = r \frac{c(\phi)}{c(\theta)} + q \frac{s(\phi)}{c(\theta)}$$

$$\dot{p} = \frac{I_y - I_z}{I_x} r q + \frac{\tau_x + \tau_{wx}}{I_x}$$

$$\dot{q} = \frac{I_z - I_x}{I_y} p r + \frac{\tau_y + \tau_{wy}}{I_y}$$

$$\dot{r} = \frac{I_x - I_y}{I_z} p q + \frac{\tau_z + \tau_{wz}}{I_z}$$

$$\dot{u} = r v - q w - g[s(\theta)] + \frac{f_{wx}}{m}$$

$$\dot{v} = p w - r u + g[s(\phi)c(\theta)] + \frac{f_{wy}}{m}$$

$$\dot{w} = q u - p v + g[c(\theta)c(\phi)] + \frac{f_{wz} - f_t}{m}$$

$$\dot{x} = w[s(\phi)s(\psi) + c(\phi)c(\psi)s(\theta)] - v[c(\phi)s(\psi) - c(\psi)s(\phi)s(\theta)] + u[c(\psi)c(\theta)]$$

$$\dot{y} = v[c(\phi)c(\psi) + s(\phi)s(\psi)s(\theta)] - w[c(\psi)s(\phi) - c(\phi)s(\psi)s(\theta)] + u[c(\theta)s(\psi)]$$

$$\dot{z} = w[c(\phi)c(\theta)] - u[s(\theta)] + v[c(\theta)s(\phi)]$$

$$\min_{s,a} \sum_{k=0}^{N-1} c(s_k, a_k) + E(s_N)$$

$$s_0 = \bar{s}_0$$

$$s_{k+1} = f(s_k, a_k)$$

$$-\pi \leq \phi \leq \pi$$

$$-\frac{\pi}{2} \leq \theta \leq \frac{\pi}{2}$$

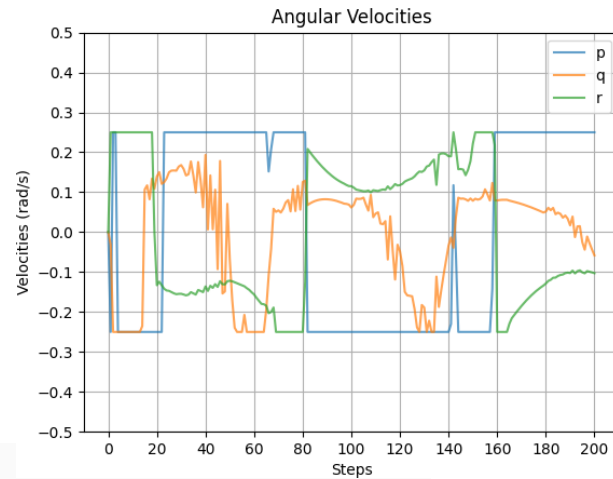
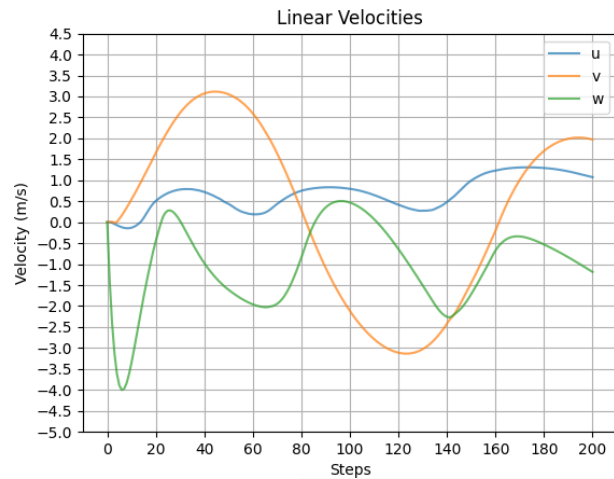
$$-\pi \leq \psi \leq \pi$$

$$-0.25 \leq \dot{\phi} \leq 0.25$$

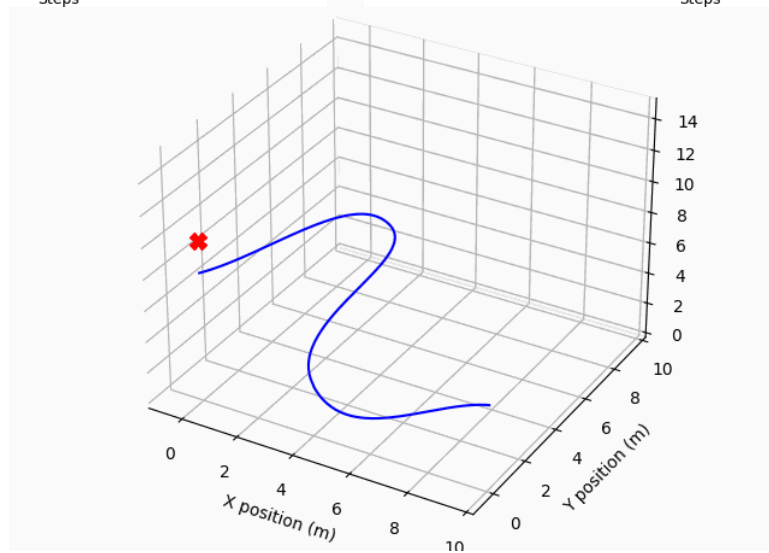
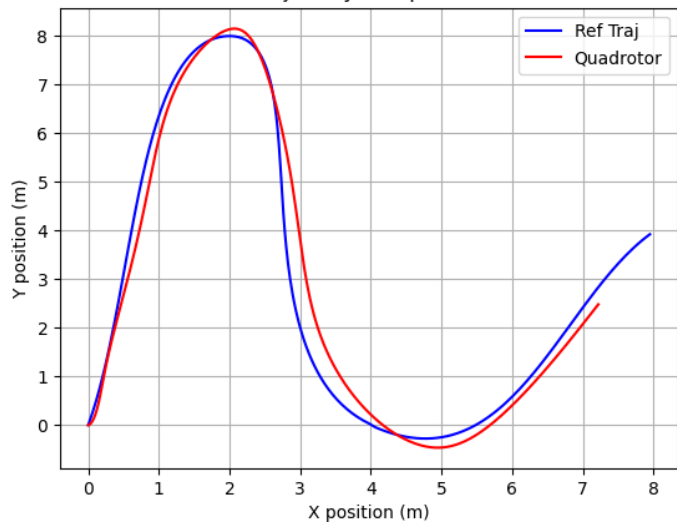
$$-0.25 \leq \dot{\theta} \leq 0.25$$

$$-0.25 \leq \dot{\psi} \leq 0.25$$

Results



Trajectory Comparison



Reinforcement Learning for Airborne Wind Energy Power Optimization

Jochem De Schutter and Jasper Hoffmann

Systems Control and Optimization Laboratory, ALU Freiburg

August 4, 2021



- ▶ system dynamics ($n_s = 3, n_a = 1$):

$$\dot{\psi} = g_k v_a \delta + \dot{\phi} \cos \theta$$

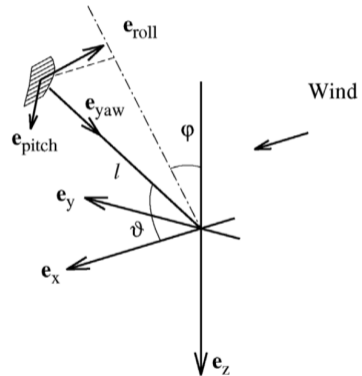
$$\dot{\phi} = -\frac{v_a}{l \sin \theta} \sin \psi$$

$$\dot{\theta} = -\frac{v_w}{l} \sin \theta + \frac{v_a}{l} \cos \theta$$

- ▶ economic cost ("max. pulling force")

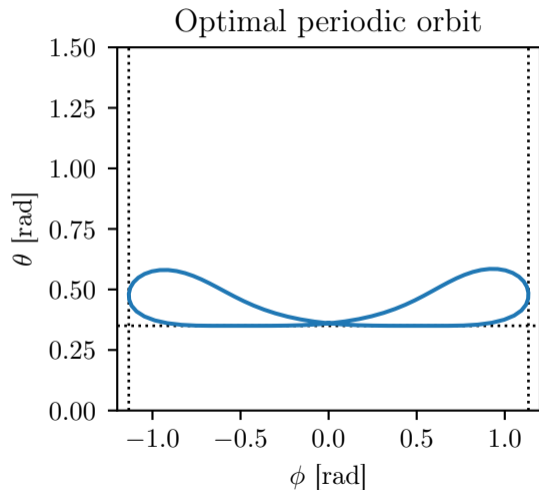
$$c(s, a) := -C_R v_a^2 + \alpha \dot{\psi}^2$$

$$v_a := v_w E \cos \theta$$

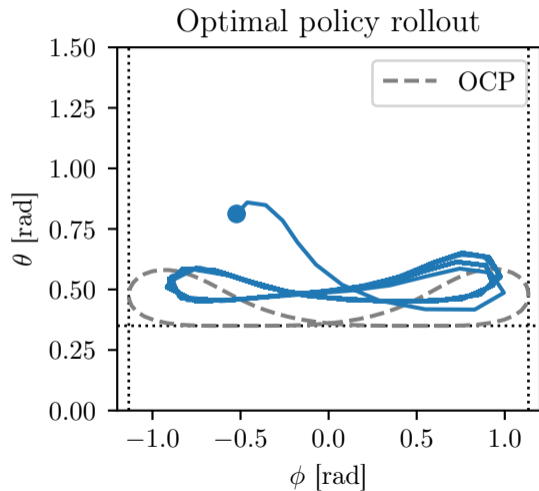




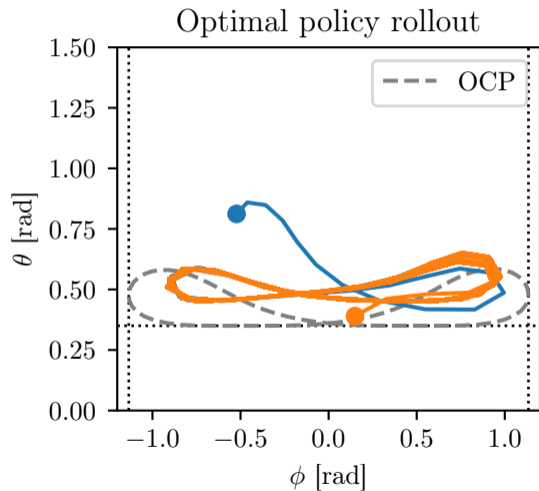
$$\begin{aligned} & \text{minimize}_{s(\cdot), a(\cdot), T} \quad \frac{1}{T} \int_0^T c(s(t), a(t)) dt \\ & \text{subject to} \quad \dot{s} = f(s(t), a(t)), \quad t \in [0, T], \\ & \quad \quad \quad s(0) = s(T), \\ & \quad \quad \quad h(s(t), a(t)) \geq 0, \quad t \in [0, T] \end{aligned}$$



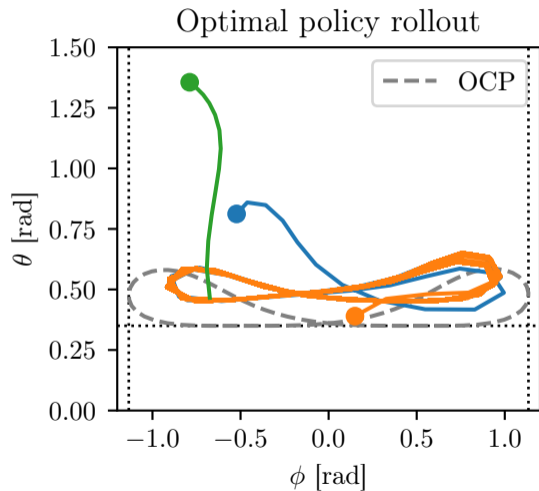
- ▶ Penalty-based constraint formulation (gym)
- ▶ Proximal policy optimization (stable-baselines3)
- ▶ 10M time-steps, $t_{\text{learn}} \approx 5\text{h}$.
- ▶ $|\bar{c}_{\text{RL}}| = 21.5\text{kN} \approx 23.8\text{kN} = |\bar{c}_{\text{OCP}}|$



- ▶ Penalty-based constraint formulation (gym)
- ▶ Proximal policy optimization (stable-baselines3)
- ▶ 10M time-steps, $t_{\text{learn}} \approx 5\text{h}$.



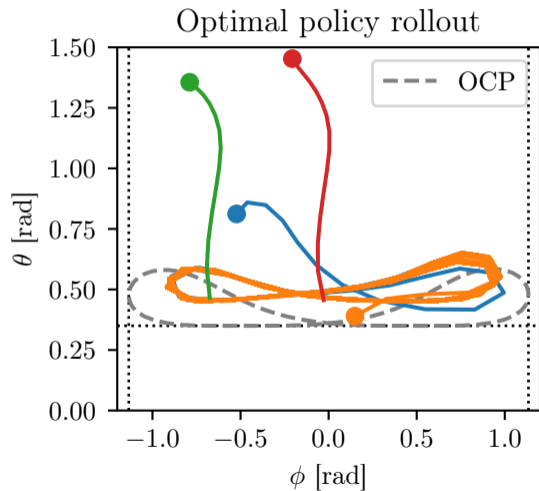
- ▶ Penalty-based constraint formulation (gym)
- ▶ Proximal policy optimization (stable-baselines3)
- ▶ 10M time-steps, $t_{\text{learn}} \approx 5\text{h}$.



Reinforcement learning



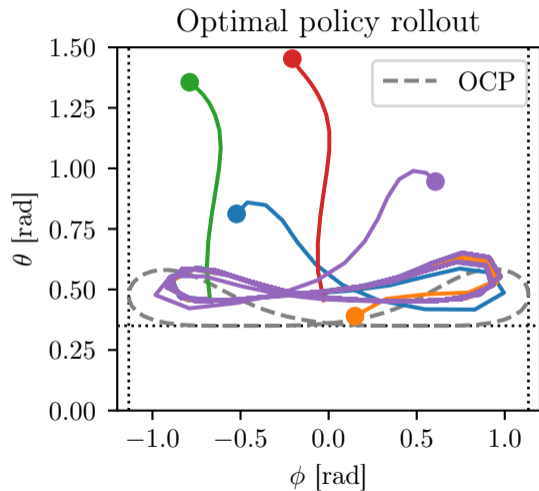
- ▶ Penalty-based constraint formulation (gym)
- ▶ Proximal policy optimization (stable-baselines3)
- ▶ 10M time-steps, $t_{\text{learn}} \approx 5\text{h}$.



Reinforcement learning



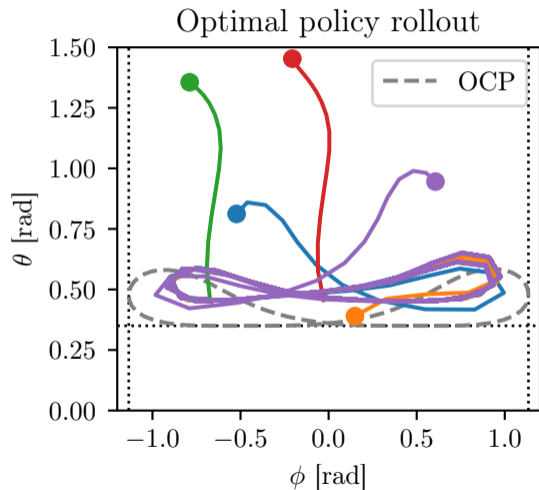
- ▶ Penalty-based constraint formulation (gym)
- ▶ Proximal policy optimization (stable-baselines3)
- ▶ 10M time-steps, $t_{\text{learn}} \approx 5\text{h}$.



- ▶ Penalty-based constraint formulation (gym)
- ▶ Proximal policy optimization (stable-baselines3)
- ▶ 10M time-steps, $t_{\text{learn}} \approx 5\text{h}$.

Future ideas:

- ▶ compare to economic MPC $\pi(s)$
- ▶ learn stochastic wind model



Point-Mass Racing Game

MPC and RL Course 2021

Christian Leininger, Rudolf Reiter

Background Idea

- **Game: Two cars: The ego car has better starting position, the opponent car has a higher velocity limit.**
- **Goal: Preventing the opponent car from overtaking, while still racing and avoiding obstacles.**
- **Design idea: Restrict the action space of an RL - policy to be safe by choosing the actions to influence only the cost of an MPC, that satisfies constraints.**
- **Basically the MPC “hides” inside the environment**
- **The MPC cost consists of a parameterizable part and a fixed part**
 - Parameterizable part: Linear plane in the position states of the vehicle
 - Fixed part: Quadratic cost to keep close to the middle part
- **The opponent agent is simulated by just the MPC policy with a fixed cost function**

Implementation

- **MPC:**

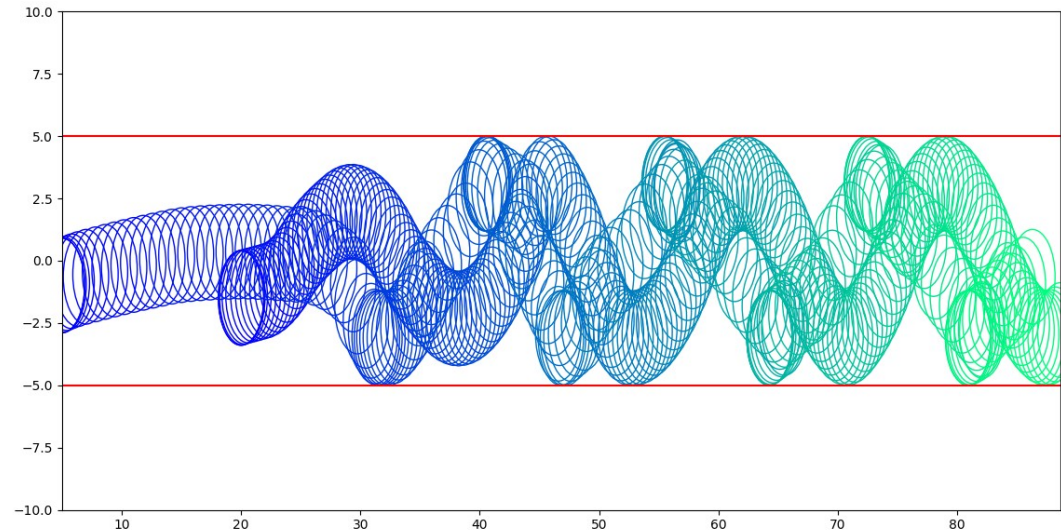
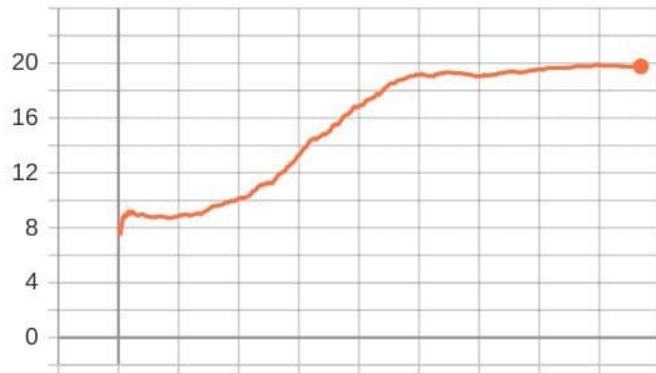
- “Hidden” in the RL environment
- CasADi, to optimize for a trajectory
- Perfect tracking is pretended
- Trajectory followed for some steps until the RL policy changes the cost function
- Circular shaped cars (point-masses) to keep it simple
- Slack variables to stay feasible
- Includes simple state estimation of the other agents trajectory (Not known to each other!)

- **RL:**

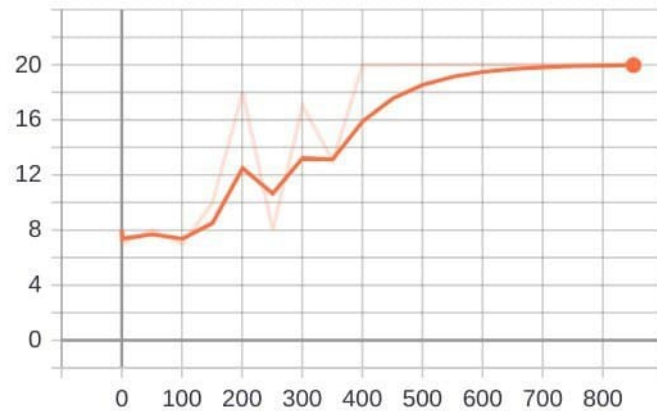
- TDQ Truncated Mixture of Continuous Distributional Quantile Critics
- Computes quantiles of Q-value Function
- Control over and underestimation

Results

Aver_reward



eval_reward



Works as expected!

Generalizes even to minor parameter changes. (positions, bounds, dimension)

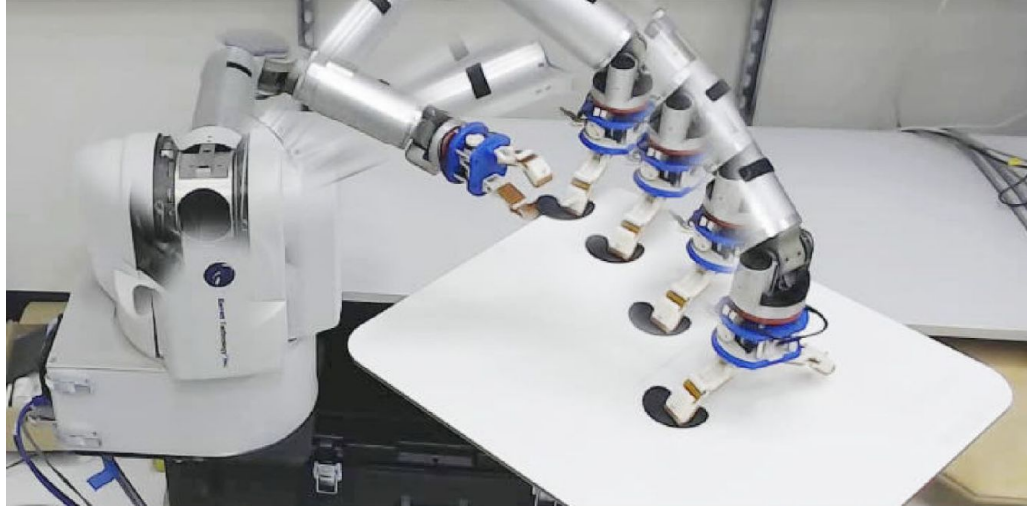
“Guaranteed safety” (to some degree)

Stability might be an issue for a tracking controller

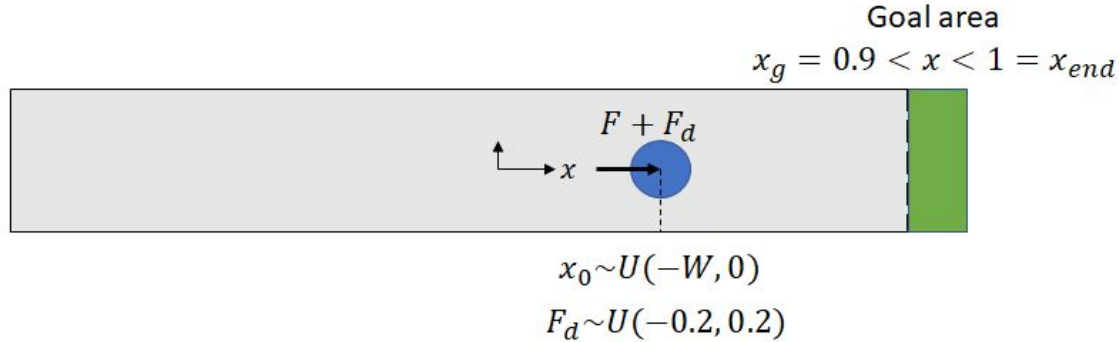
**Thank you for
your attention!
Questions?**

Non-prehensile Table Top Manipulation

Flavia Acerbo, Tommaso Sartor, Ajay Sathya



Problem Formulation



$$s = x$$

$$a = F \in [0, 10]$$

$$r_t = \begin{cases} 100, & \text{if } x_g < x < x_{end} \\ -1000, & \text{if } x > x_{end} \\ x - x_g & \text{else} \end{cases}$$

A disk lies flat on a surface (0.7 friction) at a random initial position. The disk can be pushed on the x-axis with an impulse (continuous force amplitude) at each time step.

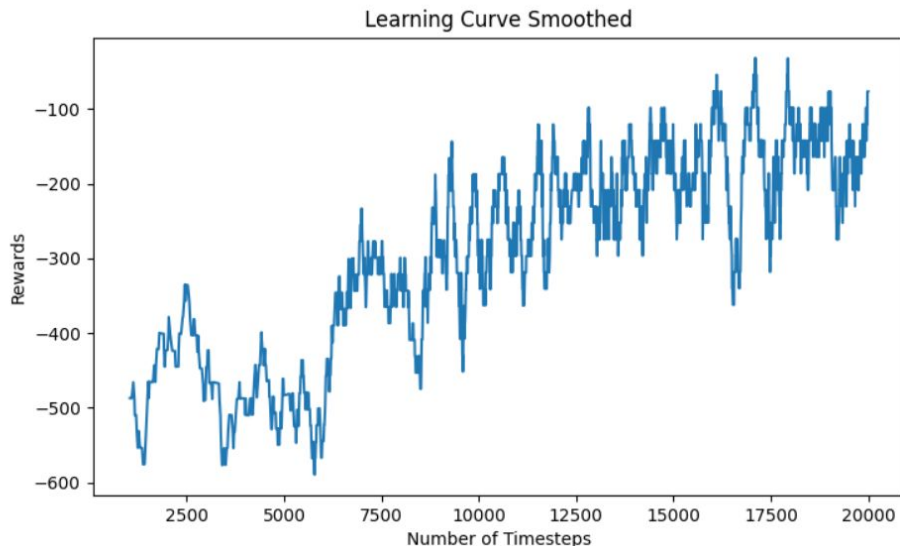
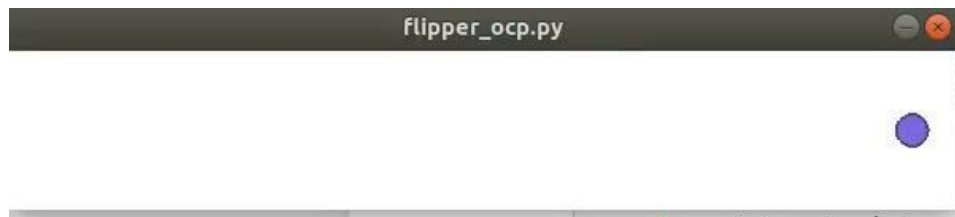
We observe the next position after the disk has come to rest.

A stochastic disturbance modifies the force that is applied to disk.

The goal is to make the disk reach the border of the surface (where it is possible to be grasped), without making it fall.

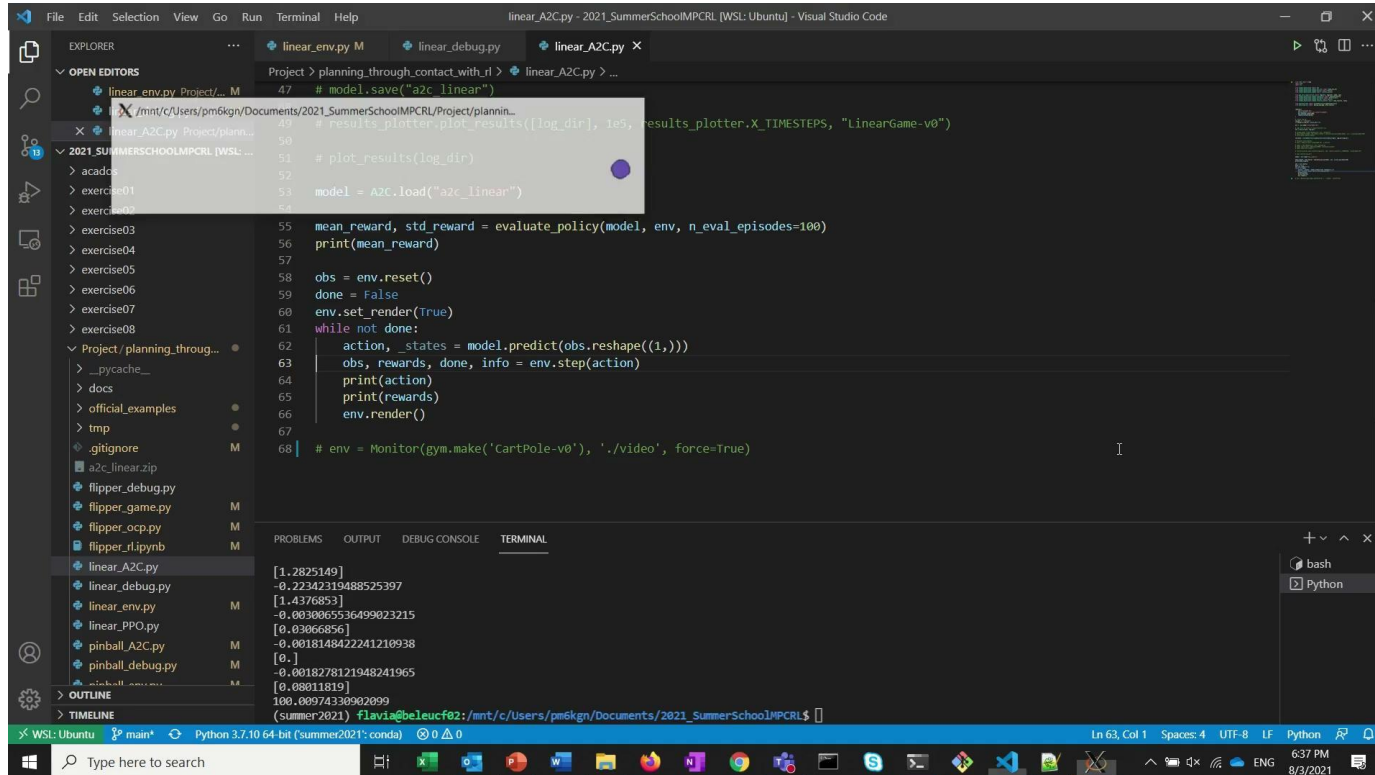
Results

- Implemented two approaches: MPC and RL.
- A2C algorithm: it learns to take big actions when away from the border and more cautious actions when close (where disturbance can be more influential)
- Using MPC, the controller failed only in case of large model discrepancy.
- Expected reward:
 - MPC policy 98.45
 - A2C policy: 99.41
 - Random policy: -871



Extra slides

A2C rollout



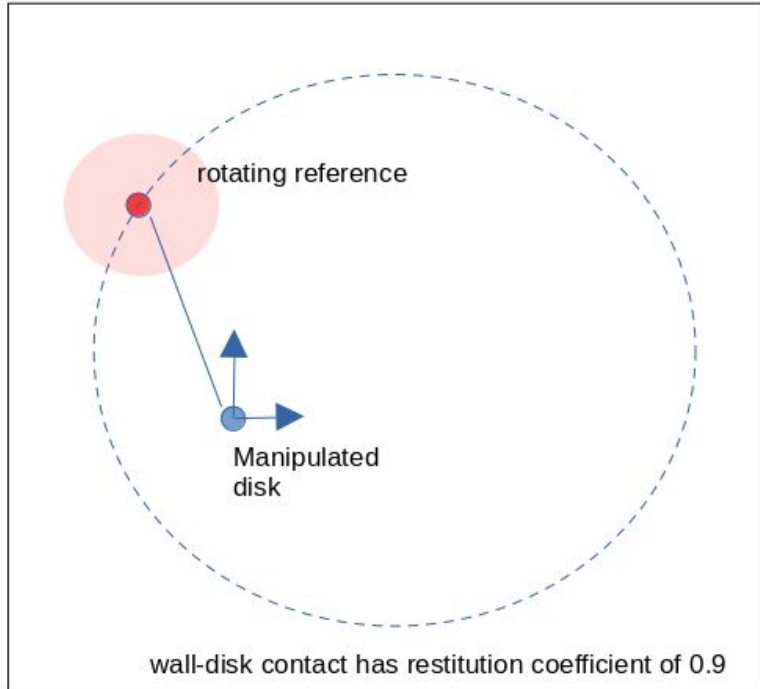
```
File Edit Selection View Go Run Terminal Help linear_A2C.py - 2021_SummerSchoolMPCRL [WSL: Ubuntu] - Visual Studio Code
EXPLORER
  OPEN EDITORS
    linear_env.py Project/... M
    linear_A2C.py Project/plann...
    2021_SUMMERSCHOOLMPCRL [WSL: ...
      acadco
      exercise01
      exercise02
      exercise03
      exercise04
      exercise05
      exercise06
      exercise07
      exercise08
      Project/planning_through...
        _pycache_
        docs
        official_examples
        tmp
        .gitignore M
        a2c_linear.zip
        flipper_debug.py
        flipper_game.py M
        flipper_ocp.py M
        flipper_rl.ipynb M
        linear_A2C.py
        linear_debug.py
        linear_env.py M
        linear_PPO.py
        pinball_A2C.py M
        pinball_debug.py M
        pinball_env.py M
      OUTLINE
      TIMELINE
  WSL: Ubuntu main Python 3.7.10 64-bit (Summer2021: conda) 0 0 0 Ln 63, Col 1 Spaces: 4 UTF-8 LF Python 6:37 PM 8/3/2021

Project > planning_through_contact_with_rl > linear_A2C.py > ...
47 # model.save("a2c_linear")
48
49 # results_plotter.plot_results([log_dir], test_results_plotter.X_TIMESTEPS, "LinearGame-v0")
50
51 # plot_results(log_dir)
52
53 model = A2C.load("a2c_linear")
54
55 mean_reward, std_reward = evaluate_policy(model, env, n_eval_episodes=100)
56 print(mean_reward)
57
58 obs = env.reset()
59 done = False
60 env.set_render(True)
61 while not done:
62     action, _states = model.predict(obs.reshape((1,)))
63     obs, rewards, done, info = env.step(action)
64     print(action)
65     print(rewards)
66     env.render()
67
68 # env = Monitor(gym.make('CartPole-v0'), './video', force=True)
```

Terminal Output:

```
[1.2825149]
-linear_debug.py -0.22342319488525397
-linear_env.py [1.4376853]
-linear_PPO.py -0.0030065536499023215
[0.03066856]
-pinball_A2C.py -0.0018148422241210938
[0.]
-pinball_debug.py -0.0018278121948241965
[0.08011819]
100.00974330902099
(Summer2021) flavia@beleucf02:/mnt/c/Users/pm6kgn/Documents/2021_SummerSchoolMPCRL $
```

Tracking time-varying reference



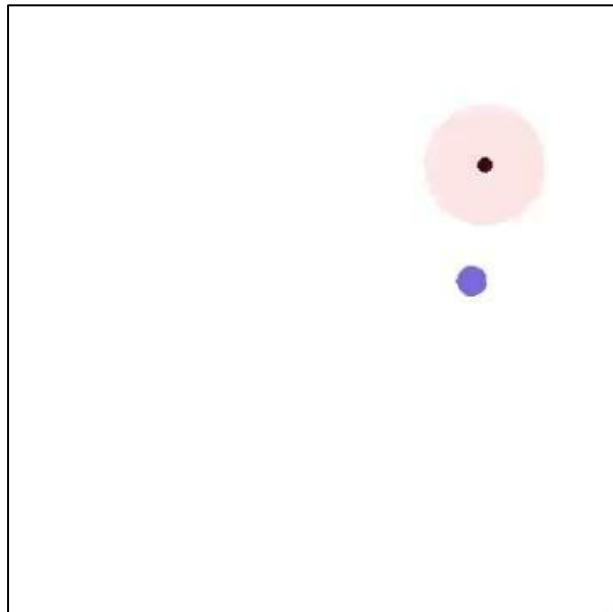
Reward is a function of distance between reference and manipulated disk

Episode is **completed** if disk is inside reference shaded area for a given number of steps

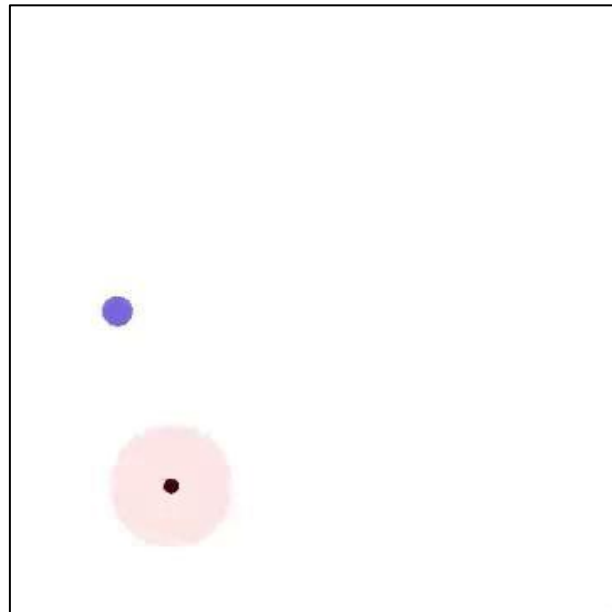
Action space is a discrete 9x9 space
4 level of force in two plus no-op for each cartesian axis,
a pair of force f_x , f_y is applied at each step

Chasing reference - videos

heuristic



a2c policy



RL approaches to the game 2048

Project by Hoang Dang & Mario Kantz

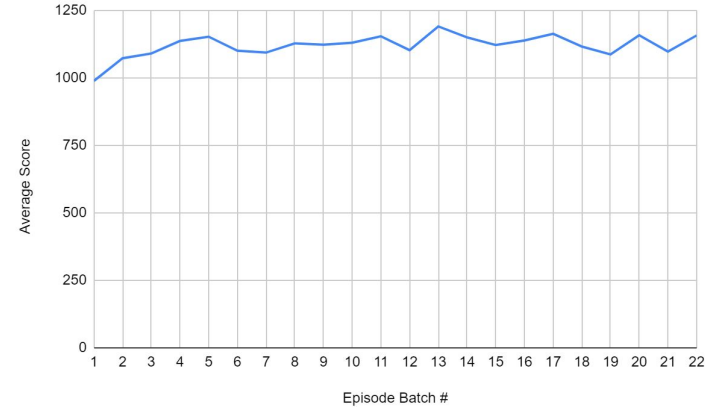
The Game 2048

- Solitaire game on 4x4 square grid
- Player can shove all tiles into a cardinal direction, merging equal tiles to their sum
- Score is the sum of values of merged tiles
- Between each player action, another tile is spawned (90% for 2, 10% for 4)
- Player may only shove into a direction if that moves or merges at least one tile
- Game ends when no more moves are possible
- Termination is thus guaranteed after at most 131,070 player moves

		2	4
		4	8
	2	16	32
	2	2	16

Preliminary Results

- We trained a DQN agent like in exercise 7 on an OpenAI gym domain of the game
- 16, 50, 50, 4 nodes
- Training occurred in blocks of 1000 episodes on a running index n counting completed blocks
- $\epsilon_n = 0.1 / \sqrt{1+n}$, to be able in theory to reach deeper parts of the game tree
- Performance capped at only $\sim +10\%$ of random performance after ~ 5000 episodes



Q-LEARNING ON-POLICY FROM REINFORCEMENT LEARNING IN MPC

ALVARO JAVIER FLOREZ MARTINEZ
ALEJANDRO ASTUDILLO VIGOYA

KU LEUVEN

MODEL PREDICTIVE CONTROL AND REINFORCEMENT LEARNING
SUMMER COURSE
FACULTY OF ENGINEERING, UNIVERSITY OF FREIBURG

PROBLEM DESCRIPTION

A parametrized MPC will be use to approximate the optimal policy and the value functions

$$Q_{\theta}(s, a) = \underset{\mathbf{x}, \mathbf{u}}{\text{minimize}} \quad \lambda_{\theta}(x_0) + \gamma^N V_{\theta}(x_N) + \sum_{k=0}^{N-1} \gamma^k \ell_{\theta}(x_k, u_k)$$

subject to

$$x_0 = s \quad u_0 = a,$$
$$x_{k+1} = f_{\theta}(x_k, u_k),$$
$$h_{\theta}(x_k, u_k) \leq 0$$

$$\pi_{\theta}(s) = \arg \min_a Q_{\theta}(s, a), \quad V_{\theta}(s) = \min_a Q_{\theta}(s, a)$$

Sensitivity of fully converged MPC

$$\nabla_{\theta} Q_{\theta}(s, a) = \nabla_{\theta} \mathcal{L}_{\theta}(s, y^*)$$

Q-Learning

$$\delta_k = \ell(s_k, a_k) + \gamma \min_{a_{k+1}} Q_{\theta}(s_{k+1}, a_{k+1}) - Q_{\theta}(s_k, a_k)$$

$$\theta \leftarrow \theta + \alpha \delta_k \nabla_{\theta} Q_{\theta}(s_k, a_k)$$

PROBLEM DESCRIPTION

A parametrized MPC will be used to approximate the optimal policy and the value functions

$$\underset{\mathbf{x}, \mathbf{u}, \sigma}{\text{minimize}} \quad V_0 + \frac{\gamma^N}{2} x_N^T S_N x_N + \sum_{k=0}^{N-1} f^T \begin{bmatrix} x_k \\ u_k \end{bmatrix} + \sum_{k=0}^{N-1} \frac{1}{2} \gamma^k \left(\|x_k\|^2 + \|u_k\|^2 + \omega^T \sigma_k \right)$$

subject to

$$\begin{aligned} x_0 &= s, \\ \mathbf{x}_{k+1} &= \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \\ \begin{bmatrix} 0 \\ -1 \end{bmatrix} + \underline{\mathbf{x}} - \sigma_k &\leq \mathbf{x}_k \leq \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \bar{\mathbf{x}} + \sigma_k, \\ \sigma_k &\geq 0, \\ -1 &\leq u_k \leq 1 \end{aligned}$$

where

$$\theta = [V_0, \underline{\mathbf{x}}, \bar{\mathbf{x}}, \mathbf{b}, f, \mathbf{A}, \mathbf{B}]$$

Q-Learning

$$\begin{aligned} \delta_k &= \ell(s_k, a_k) + \gamma V_\theta(s_{k+1}) - Q_\theta(s_k, a_k) \\ \theta &\leftarrow \theta + \alpha \delta_k \nabla_\theta Q_\theta(s_k, a_k) \end{aligned}$$

RESULTS

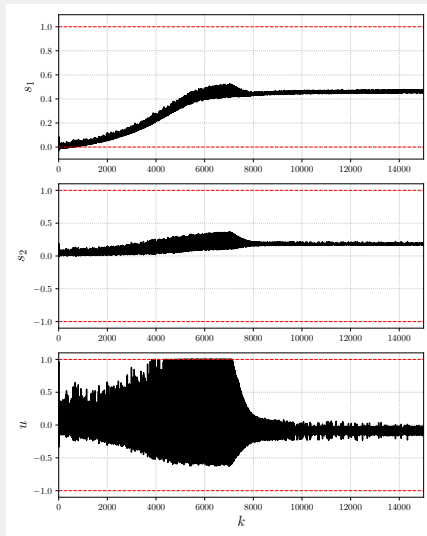


Figure 1: States and input

RESULTS

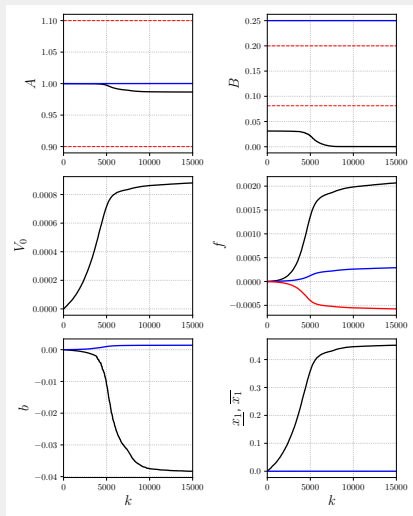
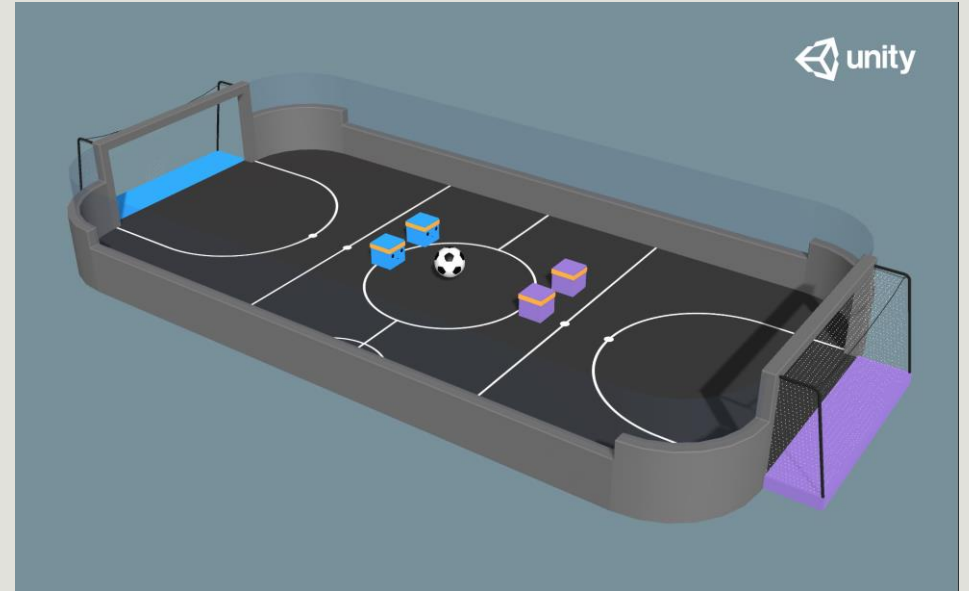


Figure 2: Parameters

How to train a computer to play advanced air hockey?



BASIC GAME AUTOMATIZATION USING UNITY ENVIRONMENT
AND REINFORCEMENT LEARNING

First Approach: Self-Made DQN-Agent

DQN: Lecture code from exercise 8 adapted to the unity environment and game scenario

State Space: 336 states, but only 264 accessible

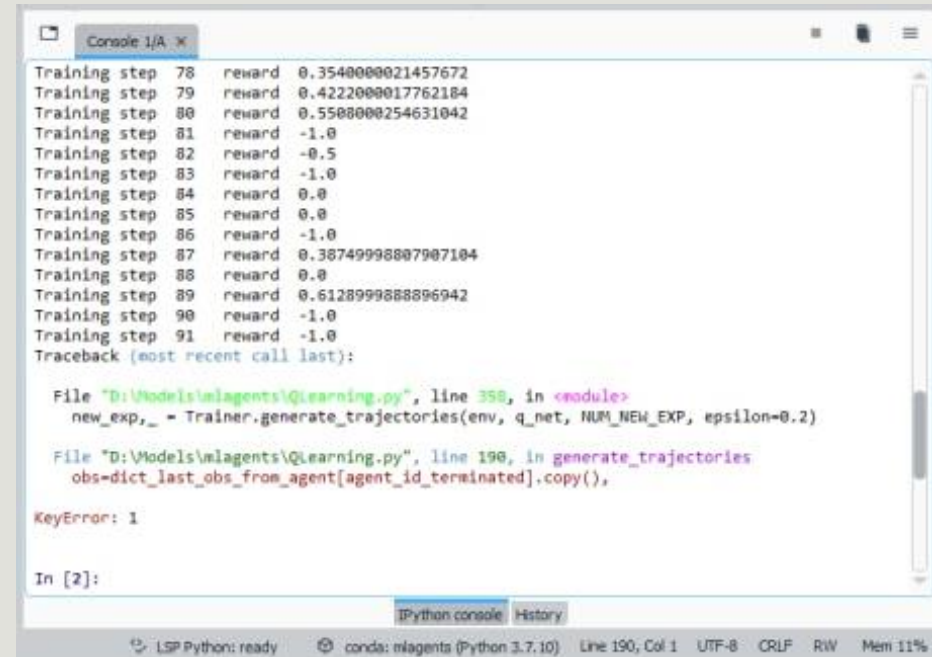
Action Space: 3x3 actions per agent

Training Idea: Train agent using DQN with the opponent agent as random noise (actions)

Results:

Training: Unsuccessful -> Error in the environment and python interface

Unity environment, not as straight forward to use in combination with python as expected



```
Console 1/A x
Training step 78 reward 0.3540000021457672
Training step 79 reward 0.4222000017762184
Training step 80 reward 0.5580000254631042
Training step 81 reward -1.0
Training step 82 reward -0.5
Training step 83 reward -1.0
Training step 84 reward 0.0
Training step 85 reward 0.0
Training step 86 reward -1.0
Training step 87 reward 0.38749998807907104
Training step 88 reward 0.0
Training step 89 reward 0.6128999888896942
Training step 90 reward -1.0
Training step 91 reward -1.0
Traceback (most recent call last):
  File "D:\Models\mlagents\Qlearning.py", line 358, in <module>
    new_exp,_ = Trainer.generate_trajectories(env, q_net, NUM_NEW_EXP, epsilon=0.2)
  File "D:\Models\mlagents\Qlearning.py", line 190, in generate_trajectories
    obs=dict_last_obs_from_agent[agent_id_terminated].copy(),
KeyError: 1

In [2]:
Python console History
LSP Python: ready conda: mlagents (Python 3.7.10) Line 190, Col 1 UTF-8 CRLF RW Mem 11%
```

Final Results: POCA vs. PPO

MA-POCA (MultiAgent Posthumous Credit Assignment): Novel multi-agent trainer that trains a centralized critic, a neural network that acts as a „coach“ for a whole group of agents
=> Team rewards are included, the agents learn cooperatively

PPO (Proximal Policy Optimization)
=> Only single rewards are accounted for, the agent learn only for themselves and team efforts not additionally accounted for

For the Presentation we hopefully have results on a second screened PC (remote work is today still not the best option)



KU LEUVEN

(Optimal) Quadrotor Control

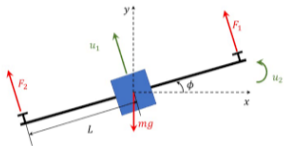
MPC vs RL

Shamil Mamedov, Mathias Bos

KU Leuven

August 2021

1 Problem Formulation



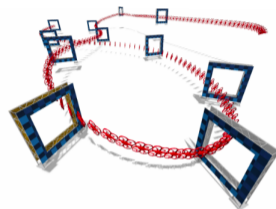
$$m\ddot{x} = -(F_1 + F_2) \sin(\phi)$$

$$m\ddot{y} = (F_1 + F_2) \cos(\phi) - mg$$

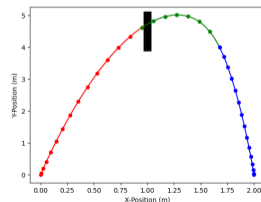
$$I\ddot{\phi} = (F_1 - F_2)L$$

states: $[x, y, \phi, \dot{x}, \dot{y}, \dot{\phi}]$

control actions: $[F_1, F_2]$



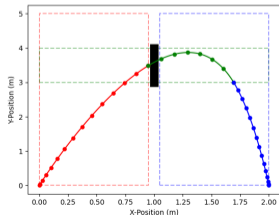
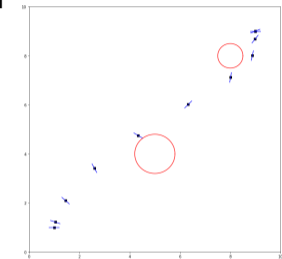
Y. Song*, M. Steinweg*, E. Kaufmann, D. Scaramuzza,
"Autonomous Drone Racing with Deep Reinforcement
Learning", 2021



2 Results

Optimal control

minimize T



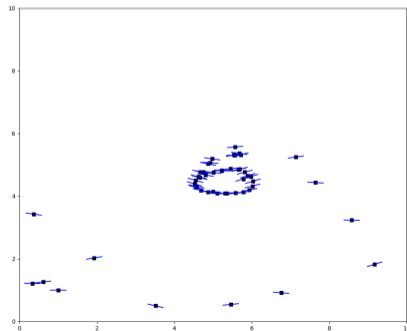
Reinforcement Learning

Stable baselines 3: SAC

Reward:

- negative distance to goal
- or $1/\text{dist}$
- penalty on hitting sides

'Future work': other reward choices...



ILQR-controlled inverted pendulum

Yizhen Wang

Microsystems engineering

$$\min_{\{s_N\}, \{a_{N-1}\}} \sum_{k=0}^{N-1} \frac{1}{2} \begin{bmatrix} s_k \\ a_k \end{bmatrix}^T \begin{bmatrix} Q & 0 \\ 0 & R \end{bmatrix} \begin{bmatrix} s_k \\ a_k \end{bmatrix} + \frac{1}{2} s_N^T P s_N$$

subject to,

$$s_0 - \bar{s}_0 = 0, \\ s_{k+1} - f_k(s_k, a_k) = 0, \quad k = 0, \dots, N-1$$

where,

$$s = \begin{bmatrix} p \\ \theta \\ v \\ \omega \end{bmatrix}, \quad a = [F]$$

For OCP,

N=20

Q=diag(10³, 10⁴, 10⁻², 10⁻²)

For LQR,

N=4

Q=diag(10², 10³, 10⁻², 10⁻²)

