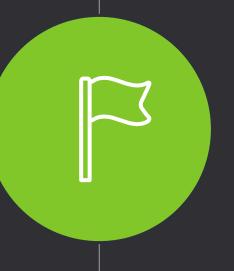


ADVERSARIAL LEARNING OF ROBUST AND SAFE CONTROLLERS FOR CYBER-PHYSICAL SYSTEMS

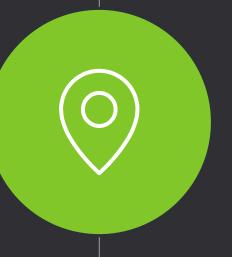
Master in Data Science and Scientific Computing

Francesco Franchina - 24 April 2020



PURPOSE

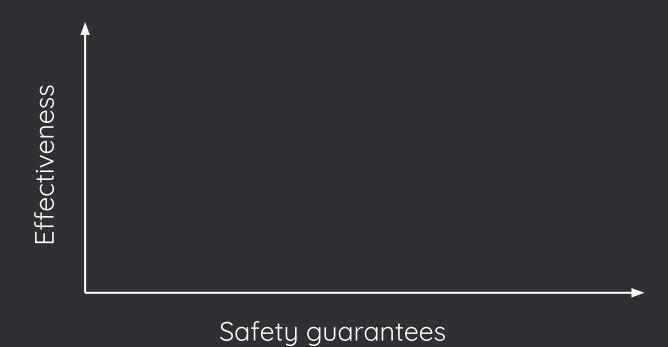
Propose a novel method for obtaining a safe controller after a training procedure



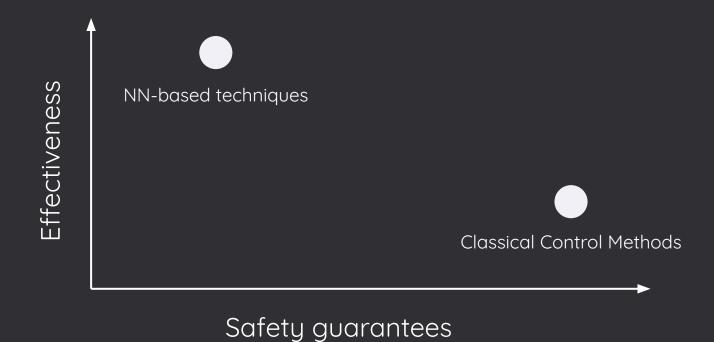
STARTING FROM

If it ain't broke, don't fix it!

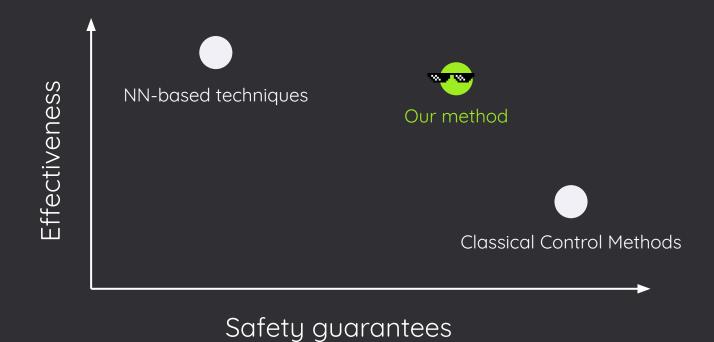
TRADEOFF



TRADEOFF



TRADEOFF



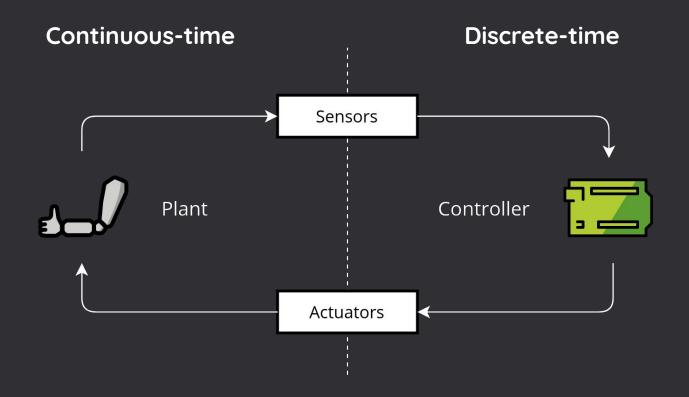
1 BACKGROUND CPS, GAN, STL

Cyber

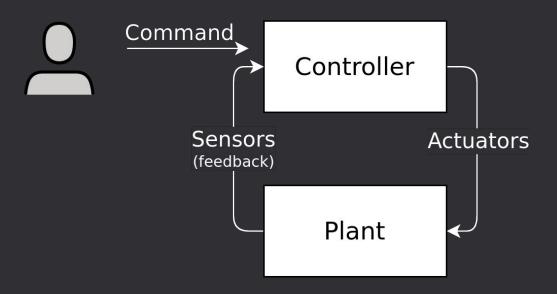
- Controller
- Algorithm
- Discrete-time

Cyber Physical Controller • Plant Algorithm • Nature's laws Discrete-time • Continuous-time

	Cyber	Physical		System
0	Controller	• Plant	0	Communication
0	Algorithm	Nature's laws	0	Sensors
0	Discrete-time	Continuous-time	0	Actuators



Closed-loop control



Generator

- Neural Network
- Learns a distribution
- Generative model

Generator

- Neural Network
- Learns a distribution
- Generative model

Discriminator

- Neural Network
- Classifier

Generator

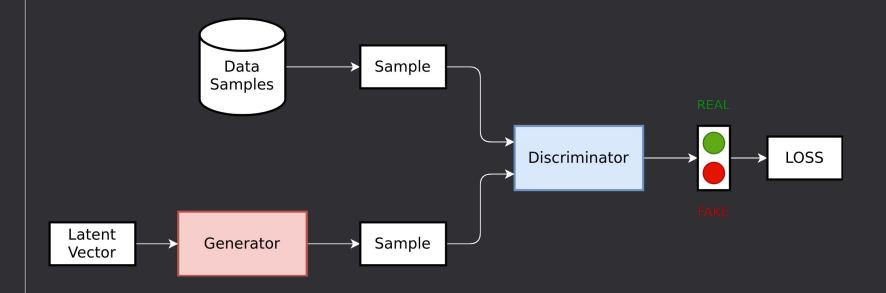
Discriminator

- Neural Network
- Learns a distribution
- Generative model

- Neural Network
- Classifier

Zero-sum game

Architecture



- Makes statements
- Operates on signals
- Monitoring

Makes statements

Operates on signals

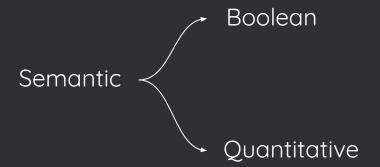
Monitoring

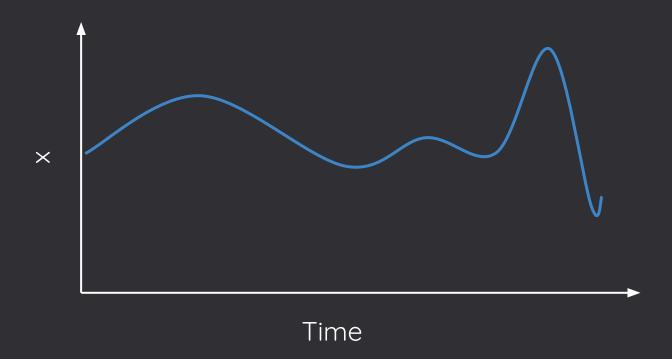
Syntax

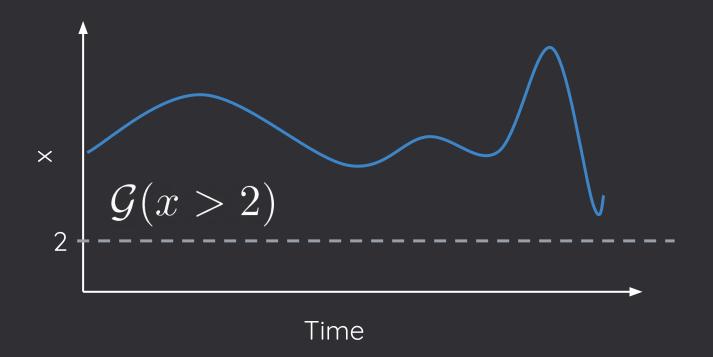
Semantic

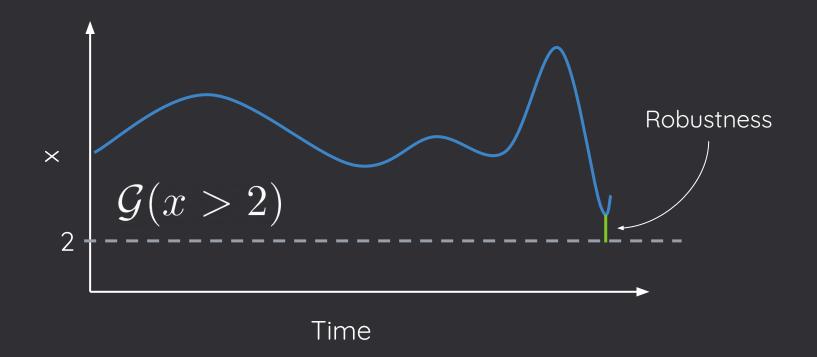
- Makes statements
- Operates on signals
- Monitoring

Syntax









2 OUR METHOD

Putting all the pieces together

Adversarial training of NNs

STL robustness as loss

Adversarial training of NNs

STL robustness as loss

 \prod

Physical Model

Differential equation

Adversarial training of NNs

STL robustness as loss



Physical Model

Differential equation



Agent

- Plant of CPS
- Senses environment
- Acts on the environment



Agent

- Plant of CPS
- Senses environment
- Acts on the environment



Environment

- Not controlled
- Can sense
- Can mutate



Agent α



Environment β





Agent α

IN: \mathbf{u}_{lpha}

OUT: $\xi_{lpha}(s) = \mathbf{o}_{lpha}$



Environment eta

IN: \mathbf{u}_{eta}

OUT: $\xi_{eta}(s) = \mathbf{o}_{eta}$



Simulator
$$f$$

$$\mathbf{s}_{i+1} = \mathbf{s}_i + f(\mathbf{s}_i, (\mathbf{u}_{\alpha})_i, (\mathbf{u}_{\beta})_i, t_i)$$

PHYSICAL MODEL - OUR METHOD

Simulator f

$$\mathbf{s}_{i+1} = \mathbf{s}_i + f(\mathbf{s}_i, (\mathbf{u}_{lpha})_i, (\mathbf{u}_{eta})_i, t_i)$$

Next state

Current state

Agent's action

Environment's Time action



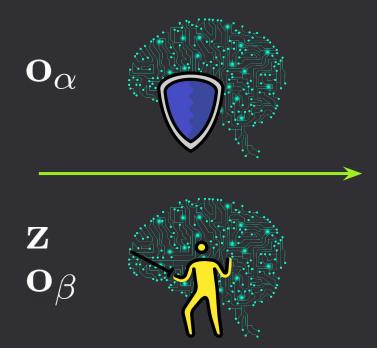
Defender



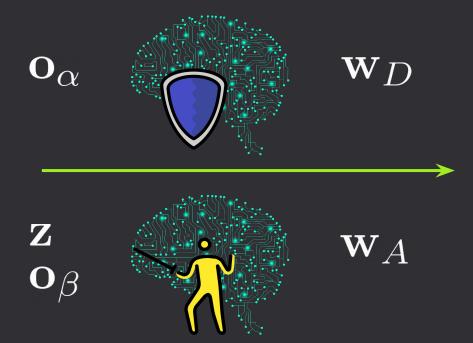
Attacker



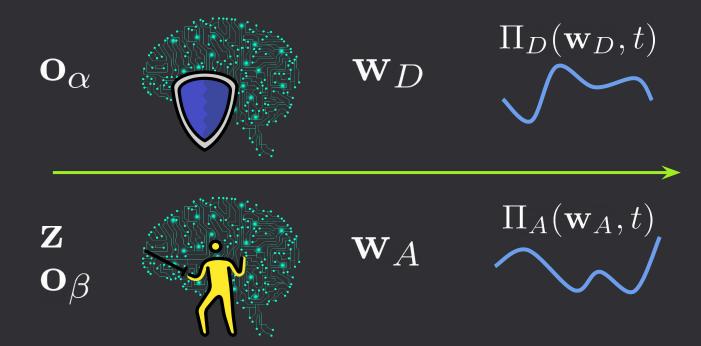
Planning



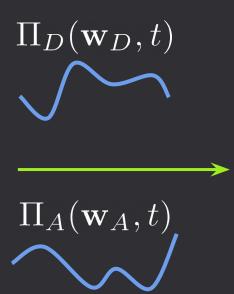
Planning

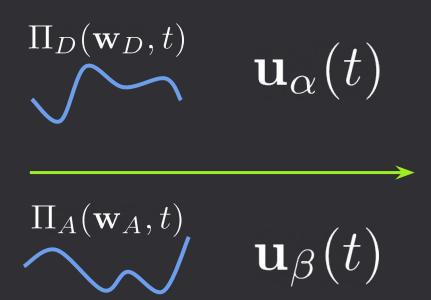


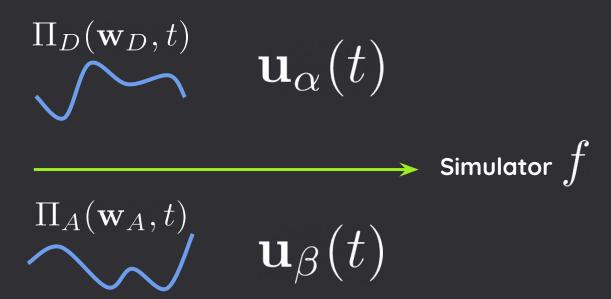
Planning

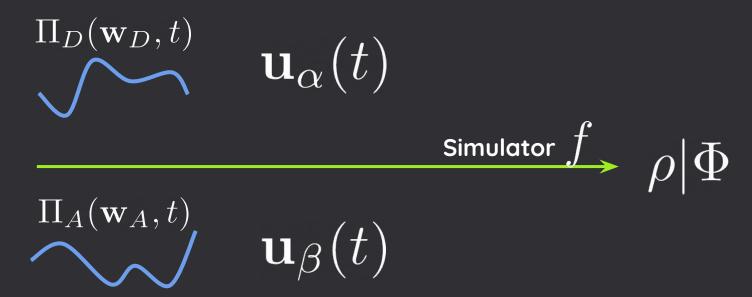


Planning $\overline{\Pi_D}(\mathbf{w}_D,\overline{t})$ \mathbf{w}_D $|\Pi_A(\overline{\mathbf{w}}_A,t)|$ \mathbf{W}_A

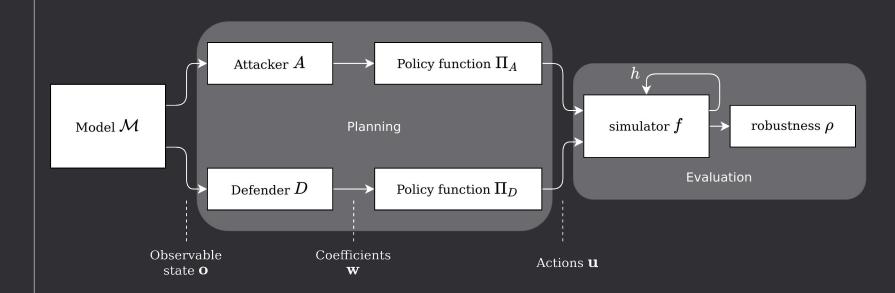




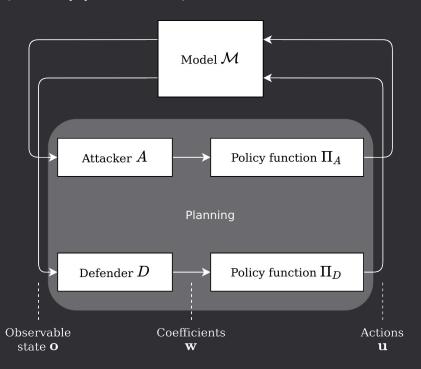




Big picture (for training)



Big picture (for application)



3 CASE STUDIES

Cruise control and Car platooning

Keep **always** the same speed.



Agent

- The car
- Can control its velocity
- Can measure the steepness
- Cannot steer

Agent

- The car
- Can control its velocity
- Can measure the steepness
- Cannot steer

Environment

- The road
- Cannot change
- Cannot have knowledge
- Goes straight

$$\Phi = \mathcal{G}(v_c \ge \tilde{v} - \varepsilon \land v_c \le \tilde{v} + \varepsilon)$$

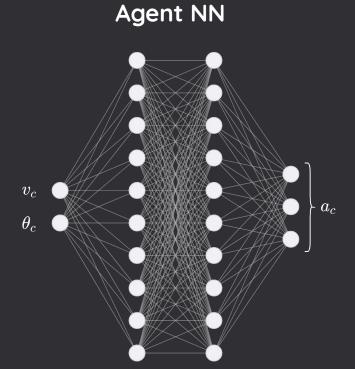
$$\Phi = \mathcal{G}(v_c \ge \tilde{v} - \varepsilon \land v_c \le \tilde{v} + \varepsilon)$$

$$m\frac{dv}{dt} = m\left(\frac{dv}{dt}\right)_{in} - \nu mg\cos\theta_c - mg\sin\theta_c$$

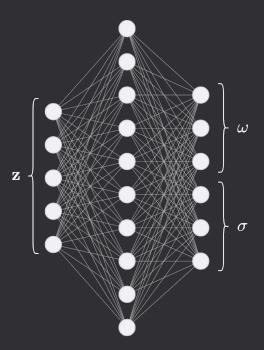
$$\Phi = \mathcal{G}(v_c \ge \tilde{v} - \varepsilon \land v_c \le \tilde{v} + \varepsilon)$$

$$m\frac{dv}{dt} = m\left(\frac{dv}{dt}\right)_{in} - \nu mg\cos\theta_c - mg\sin\theta_c$$

$$r(x) = \sum_{i=1}^{d} \omega_i \exp\left(-\frac{||x - \mu_i||^2}{\sigma_i^2}\right)$$

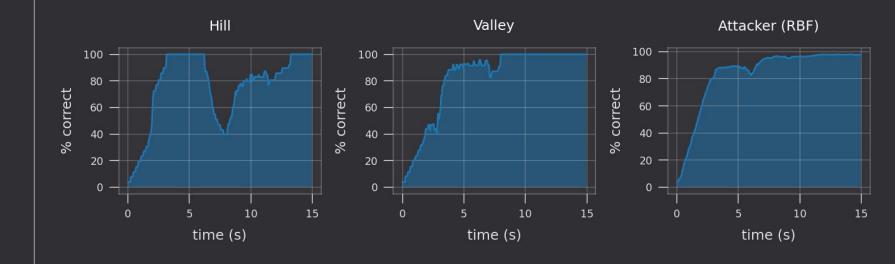


Environment NN

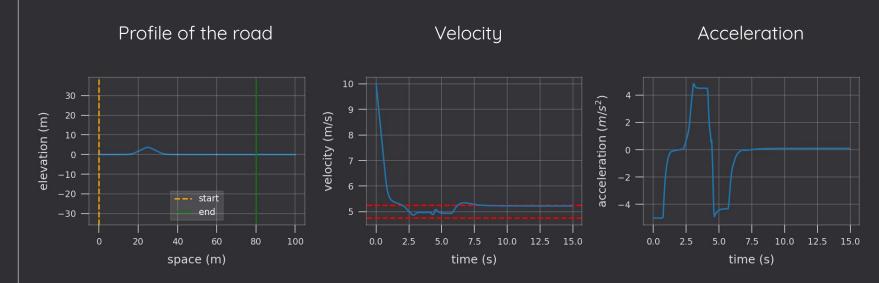


Percentage of correctness over time

10k simulations

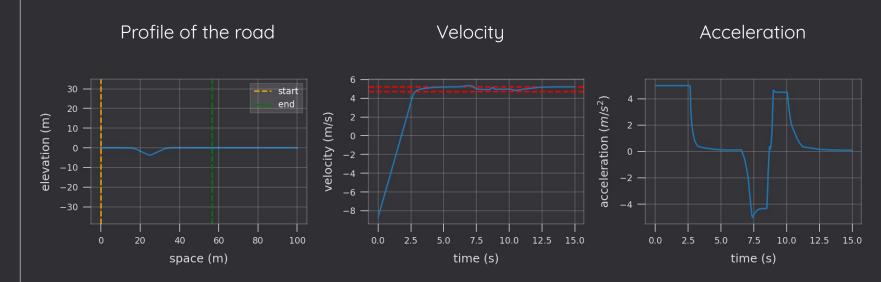


Hill test



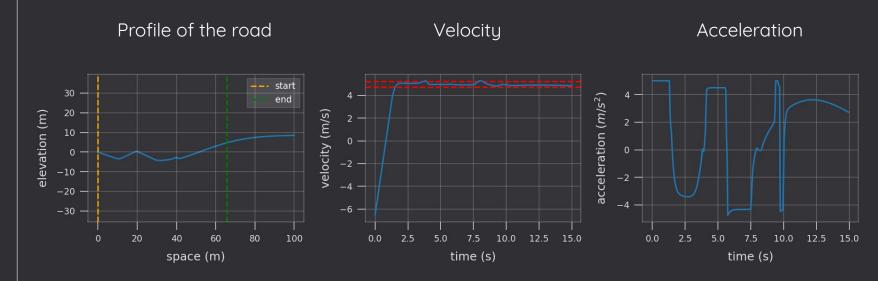
Initial velocity: 11.98 m/s

Valley test



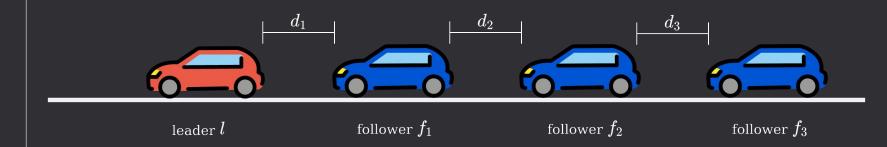
Initial velocity: -8.94 m/s

Attacker test

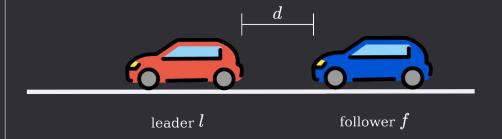


Initial velocity: -6.93 m/s

Keep **always** the distance within the safety range.



Keep always the distance within the safety range.



Divide et impera

Agent

- The follower
- Can control its velocity
- Can measure the distance
- Knows the other's velocity
- Goes only forward

Environment

- The leader
- Can control its velocity
- Can measure the distance
- Knows the other's velocity
- Goes only forward

$$\Phi = \mathcal{G}(d \le d_{max} \land d \ge d_{min})$$

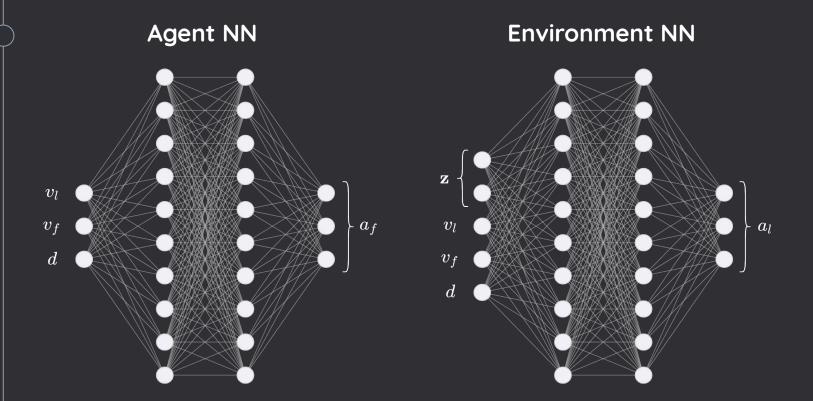
→ Agent

■ Environment

$$\Phi = \mathcal{G}(d \le d_{max} \land d \ge d_{min})$$

$$m\frac{dv}{dt} = m\left(\frac{dv}{dt}\right)_{in} - \nu mg$$

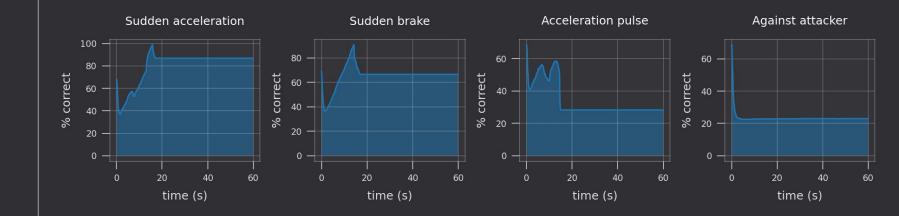
$$m\frac{dv}{dt} = m\left(\frac{dv}{dt}\right)_{in} - \nu mg$$



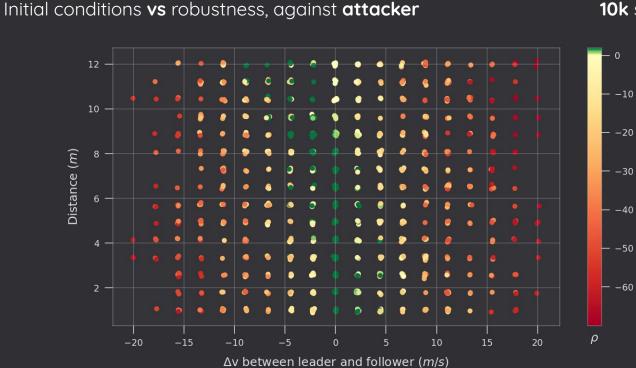
CAR PLATOONING RESULTS - CASE STUDIES

Percentage of correctness over time

10k simulations

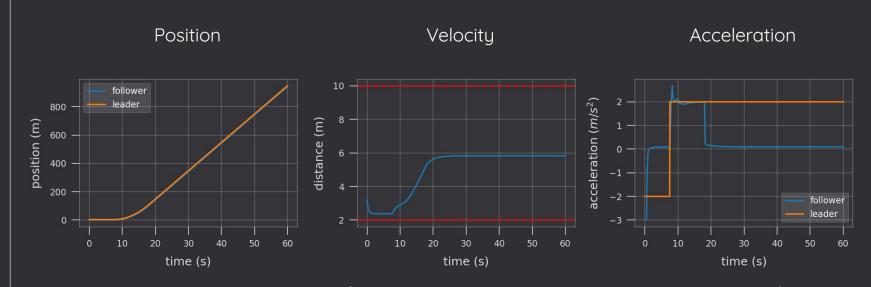


CAR PLATOONING RESULTS - CASE STUDIES



10k simulations

Sudden acceleration test

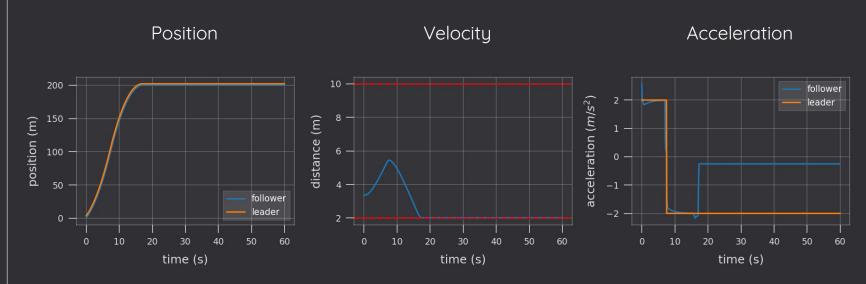


Leader's initial velocity: -0.03 m/s

Initial distance: 3.31 m

Follower's initial velocity: 2.30 m/s

Sudden brake test

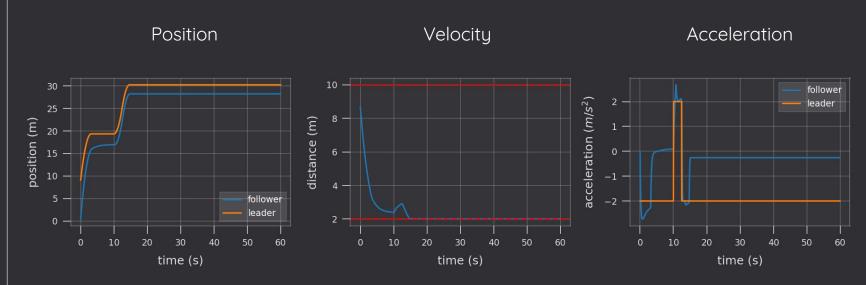


Leader's initial velocity: 6.72 m/s

Initial distance: 3.35 m

Follower's initial velocity: 6.65 m/s

Acceleration pulse test

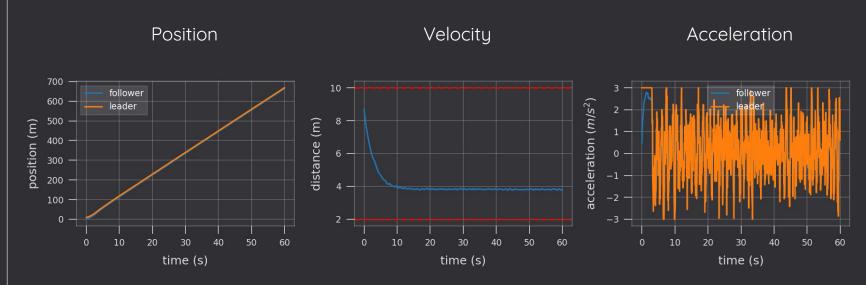


Leader's initial velocity: 6.70 m/s

Initial distance: 8.78 m

Follower's initial velocity: 8.90 m/s

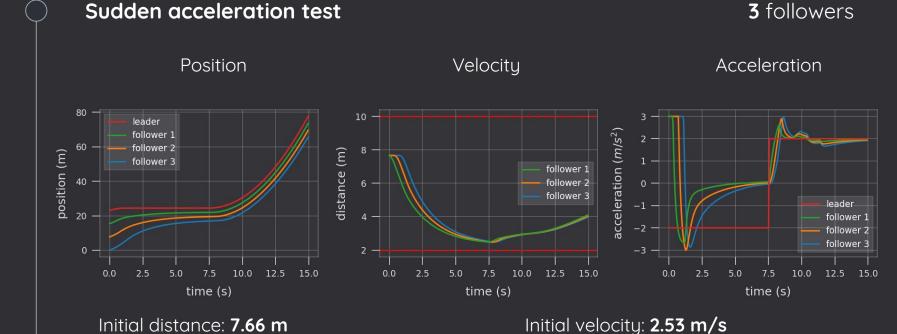
Attacker test

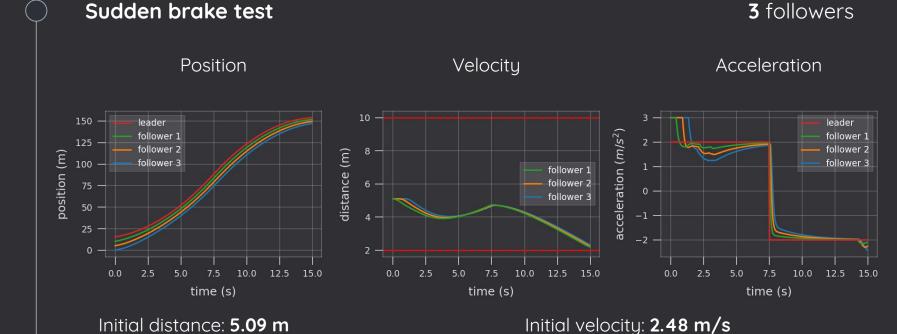


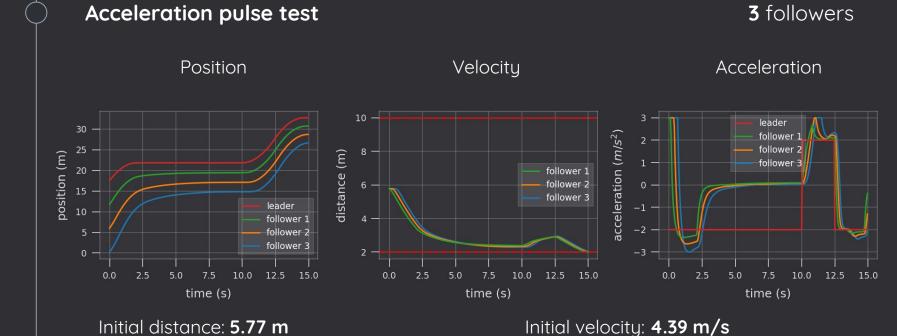
Leader's initial velocity: 11.98 m/s

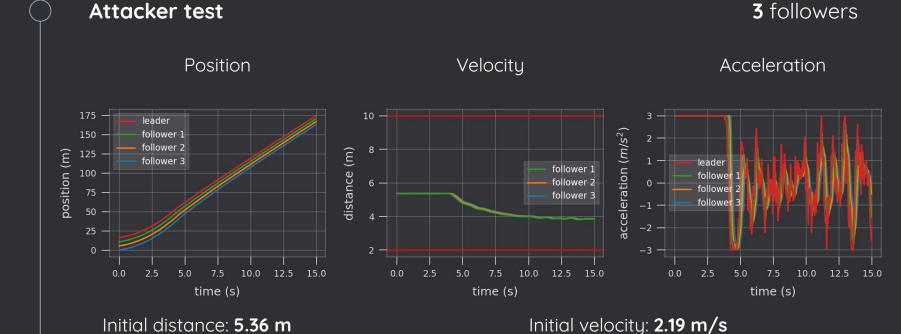
Initial distance: 11.98 m

Follower's initial velocity: 11.98 m/s









4

CONCLUSION

Recap and future work

RECAP - CONCLUSION

- We developed a method that
 - Makes use of adversarial training in a new way
 - Makes use of episodes instead of a training set
 - Trains a performant and safe controller
 - Produces interesting results

We developed a framework ready to be applied to new models.

FUTURE WORK - CONCLUSION

- Possible developments
 - Faster training using mini-batches
 - Approximation of non-differentiable models
 - Optimized sampling of the initial conditions



Thanks!

Coming soon on GitHub...