

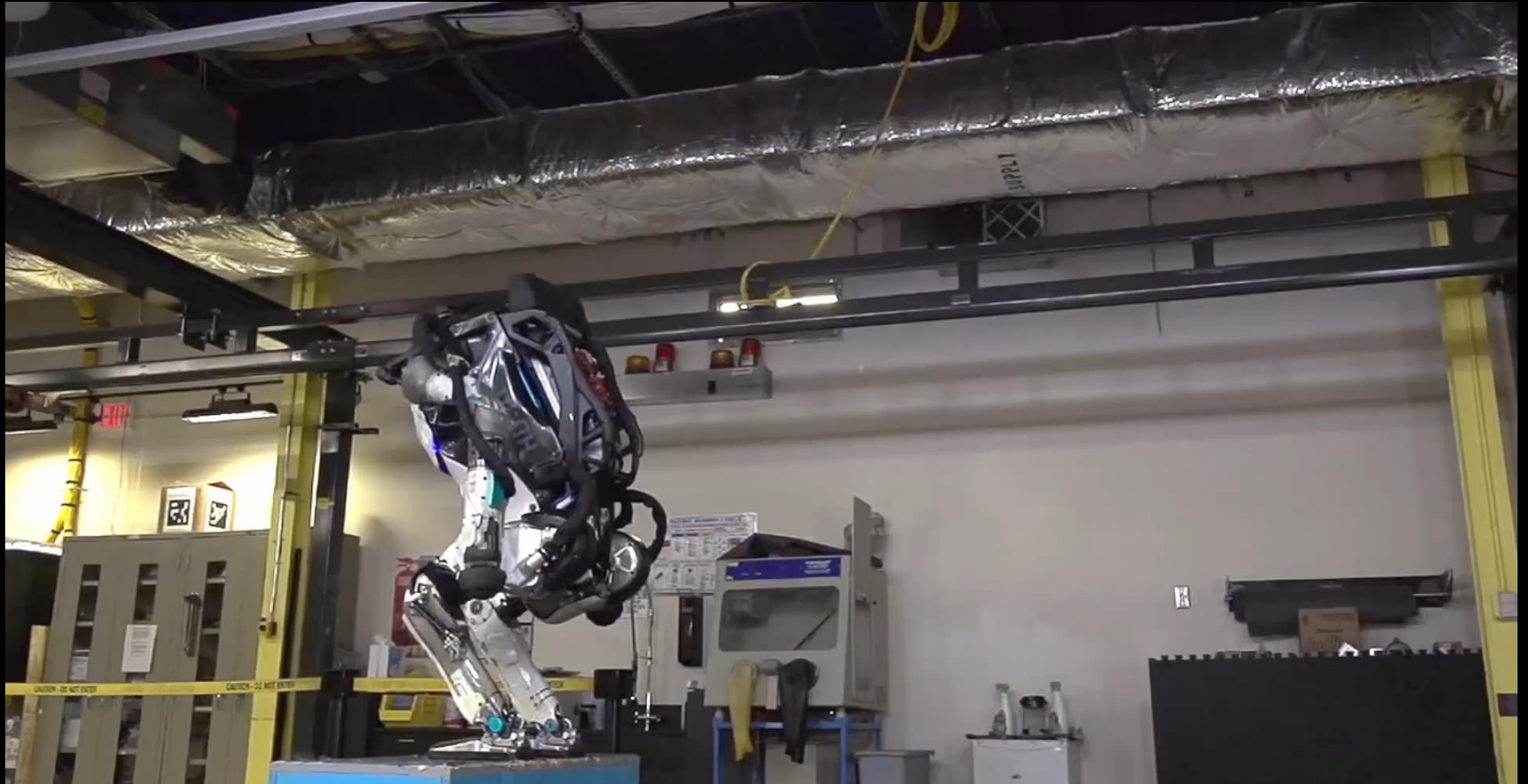
Robots Learning (Through) Interactions

May 7th 2020

Jens Kober



Advanced Robot Hardware + Manual Programming



<https://youtu.be/fRj34o4hN4I>

What Makes Tasks Hard?

- Complex dynamics
- Uncertainties and variations
 - Objects
 - Environment
 - Tasks
 - Human behaviour
- Occurring in all robot domains!

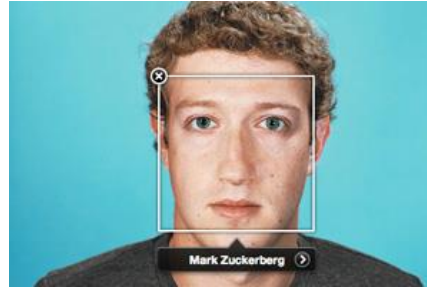
Learning to Interact



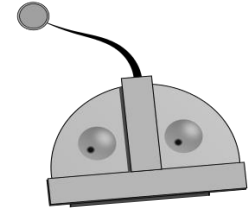
Why is Robotics Different?



IBM Research



sociable.co



europe1.fr



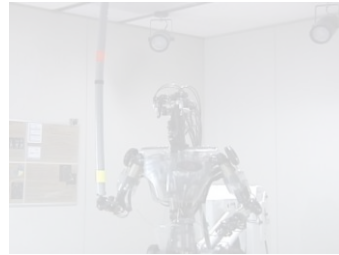
DeepMind

openclipart.org

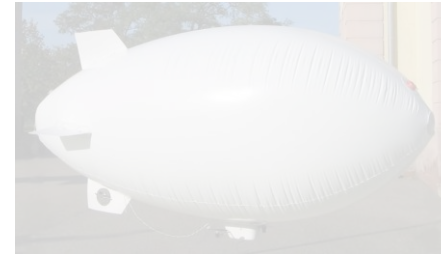
(Reinforcement) Learning in Robotics



Mahadevan & Connell, AI 1992



Schaal, NIPS 1996



Rottmann et al., IROS 2007



Kormushev & Calinon, IROS 2010

Domain appropriate methods needed!

- Safe with real robots
- Fast learning
 - Sample efficient – “small” data
 - Incorporate prior knowledge
 - Few open parameters
- Real time computations



Google

Naïve Deep RL

<https://youtu.be/vv85S4Z-ZG0>



Tim de Bruin

Imitation Learning

- Teacher demonstrates skill, student tries to mimic



testmyprep.com

Reinforcement Learning

- Practice, practice, practice...



cnn.com



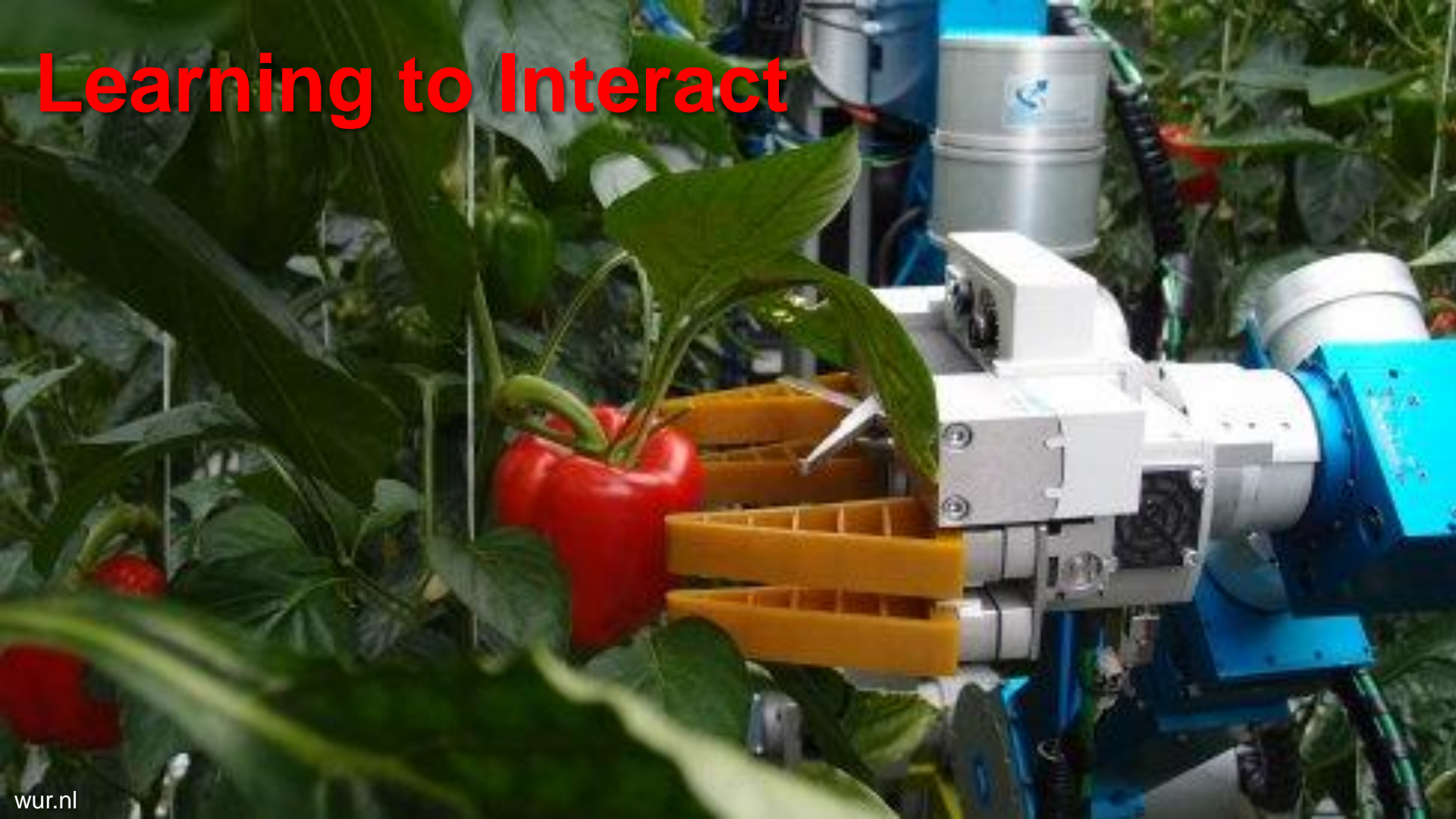
incaavalanche.com

How to Learn?

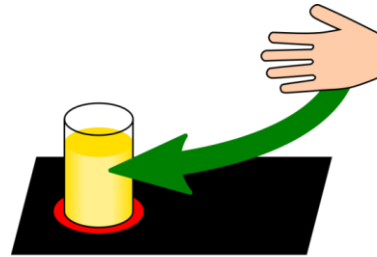
- Continued student-teacher interaction
 - Additional demonstrations
 - Intermittent feedback
- Largely missing in robot learning!
- Benefits
 - Speed-up
 - Complex tasks
 - Intuitive



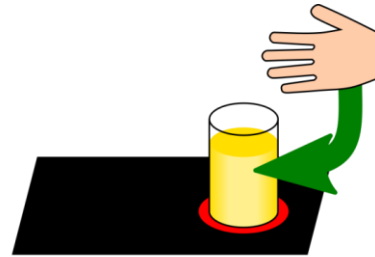
Learning to Interact



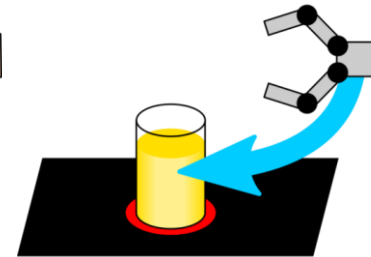
Imitation Learning



a) Demonstration 1

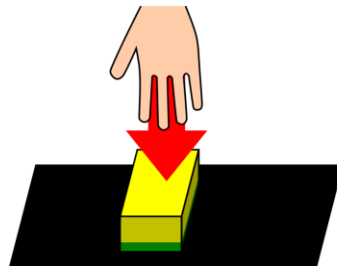


b) Demonstration 2

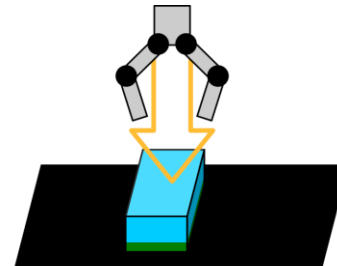


c) Generalization 1

- Combinatorial explosion
- Force \rightarrow Position? Position \rightarrow Force?



a) Demonstration



b) Generalization

https://youtu.be/t_ZoiKcEM0M



Kinesthetic teaching

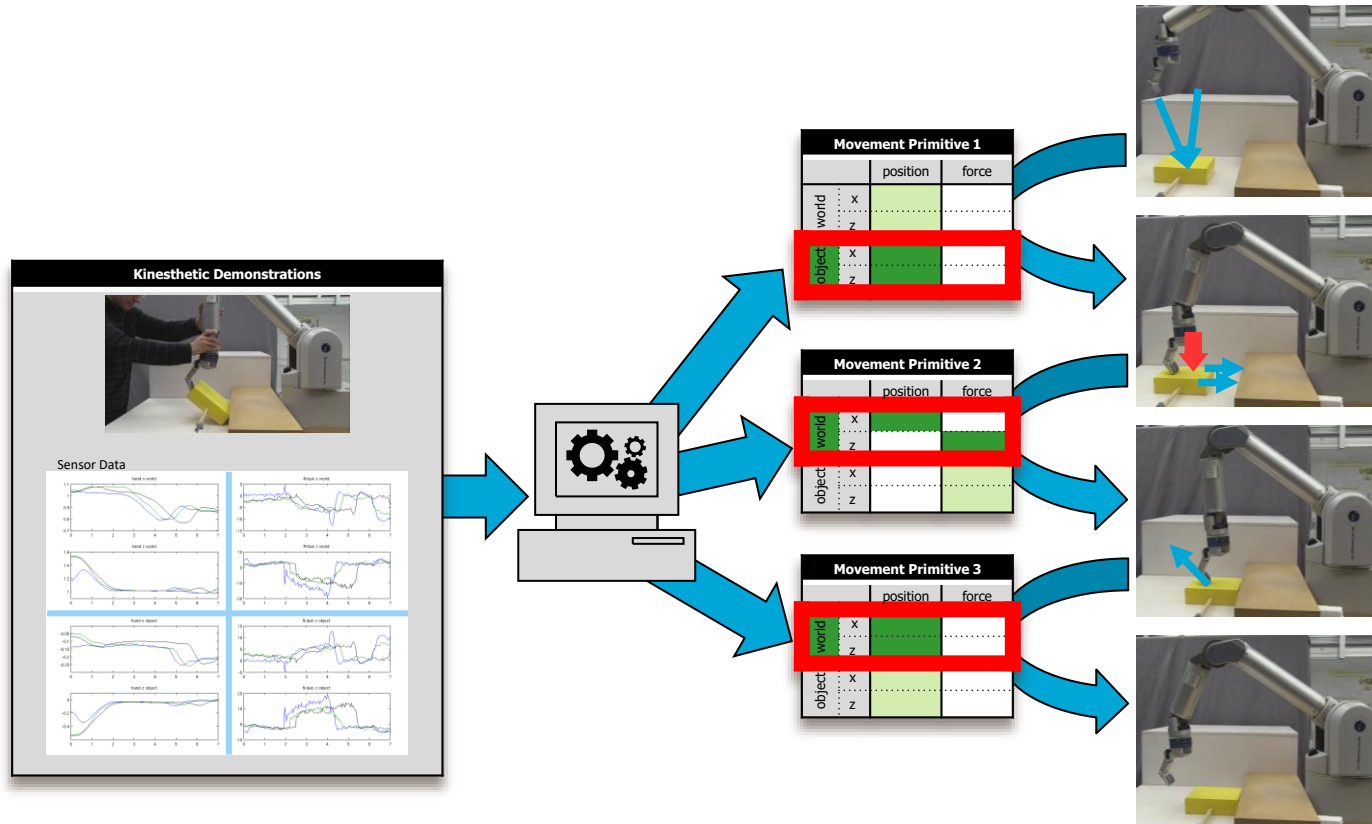


Skill reproduction

Kober, Gienger, & Steil, ICRA 2015

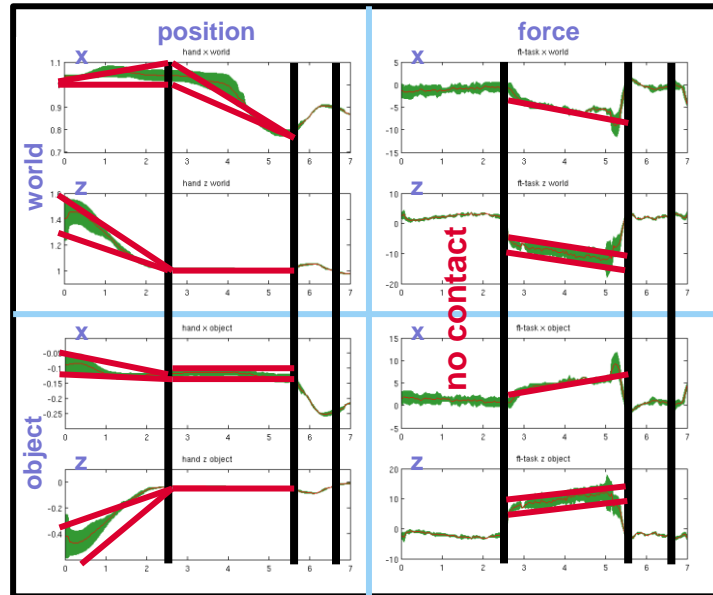
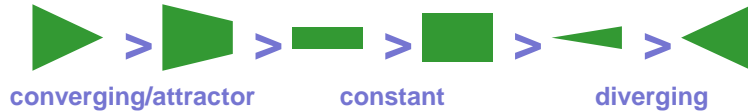
HR
Europe
Honda Research Institute

Primitives within a Sequence



Primitives within a Sequence

Score per MP/frame/component/modality based on statistics & contact



Scores: Movement Primitive 2

		position	force
world	x	90.0	20.0
	z	10.0	16.0
object	x	8.0	20.0
	z	10.0	16.0



Scores: Movement Primitive 2

		position	force
world	x		
	z		
object	x		
	z		

https://youtu.be/t_ZoiKcEM0M



HRI Europe
Honda Research Institute

<https://youtu.be/WCayQ8xliU8>



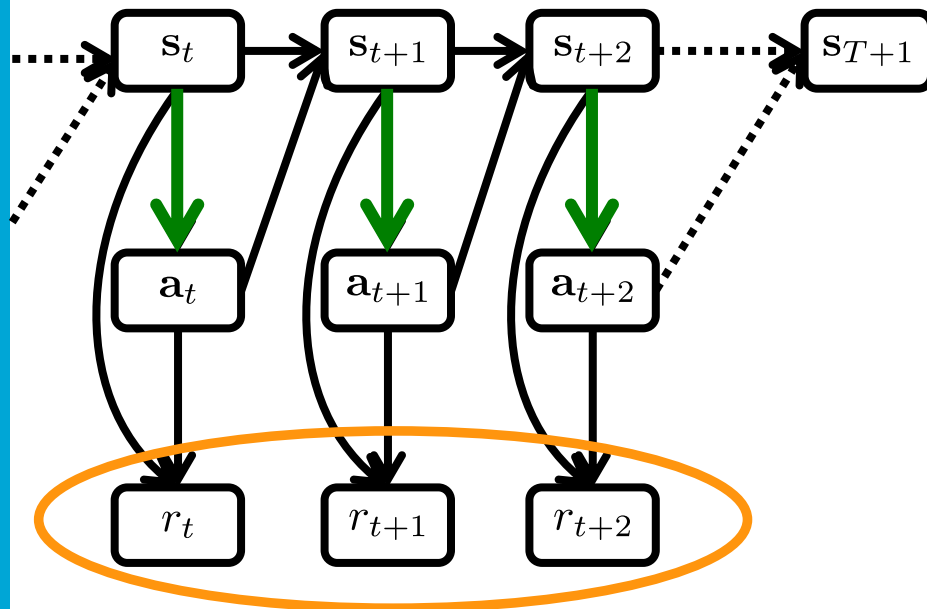
Simon Manschitz



Kinesthetic Demonstration

Reinforcement Learning

Sutton & Barto 1998



states
 $s_{1:T+1}$

actions
 $a_{1:T}$

rewards
 $r_{1:T}$

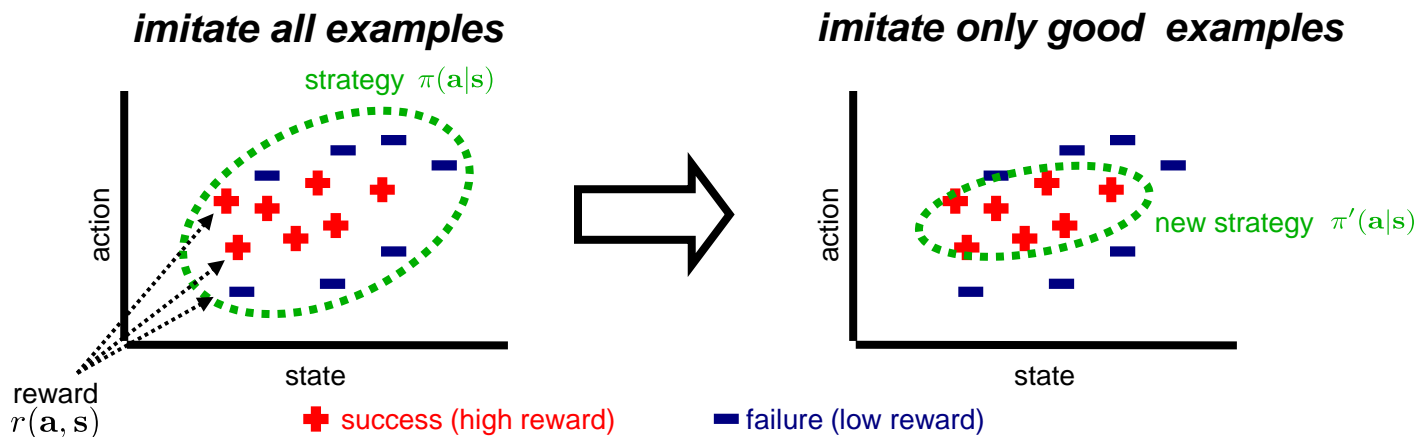
policy
 $\pi(a_t | s_t, t, \theta)$

policy parameters

Objective: maximize expected return $J(\theta)$

Reward-Weighted Imitation

- Maximize reward = optimize strategy
- One possibility: reward-weighted imitation



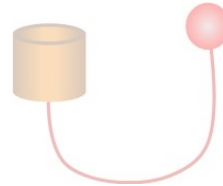
$$\theta' = \theta + E \left\{ \sum_{t=1}^T Q_t^{\text{sa}} \mathbf{W}_t^{\text{s}} \right\}^{-1} E \left\{ \sum_{t=1}^T Q_t^{\text{sa}} \mathbf{W}_t^{\text{s}} \epsilon_t \right\}$$

Reward-weighted Imitation

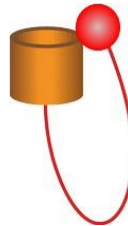
Trial 1
Return 0,1



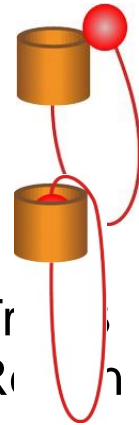
Trial 2
Return 0,3



Trial 3
Return 0,8



Trial 4
Return 0,8



Trial 4
Return 1,0

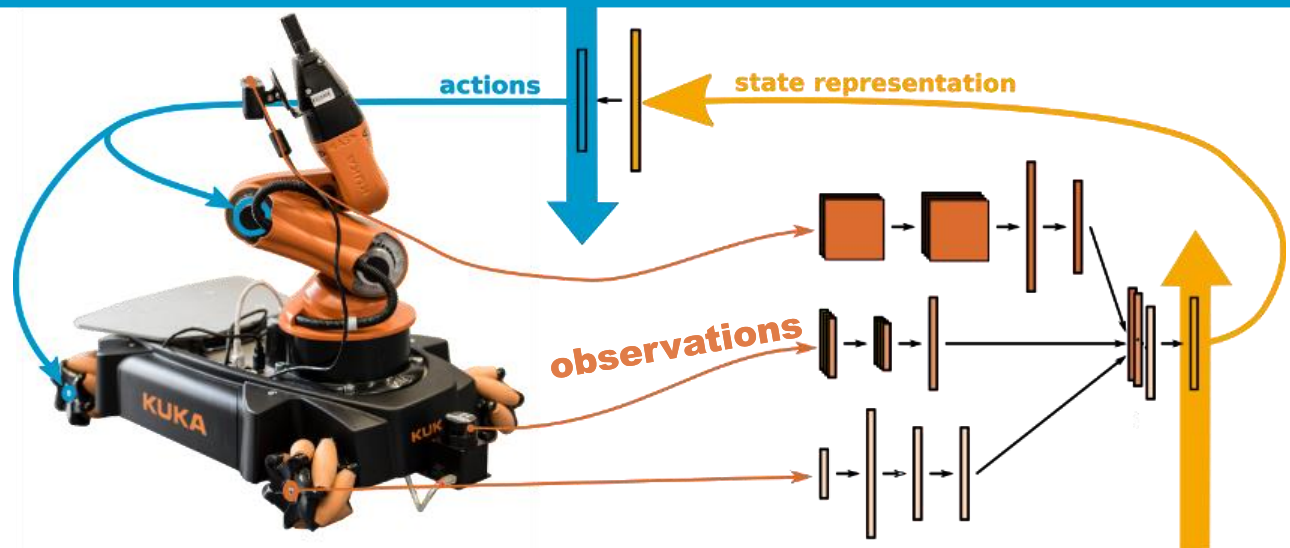


<https://youtu.be/cNyoMVZQdYM>



Reinforcement Learning:

- Learn behaviors using trial and error and a scalar reward function

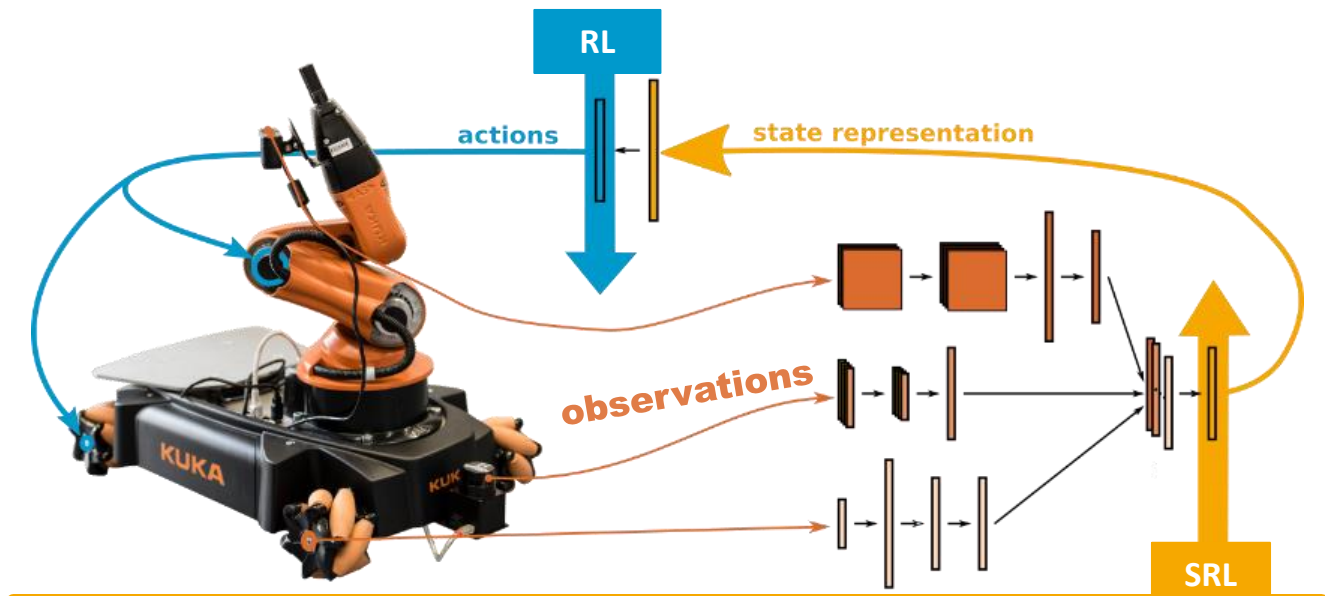


State Representation Learning:

- Cost functions acting on the state representation
- Prior knowledge about the world
- Force learning a more general representation



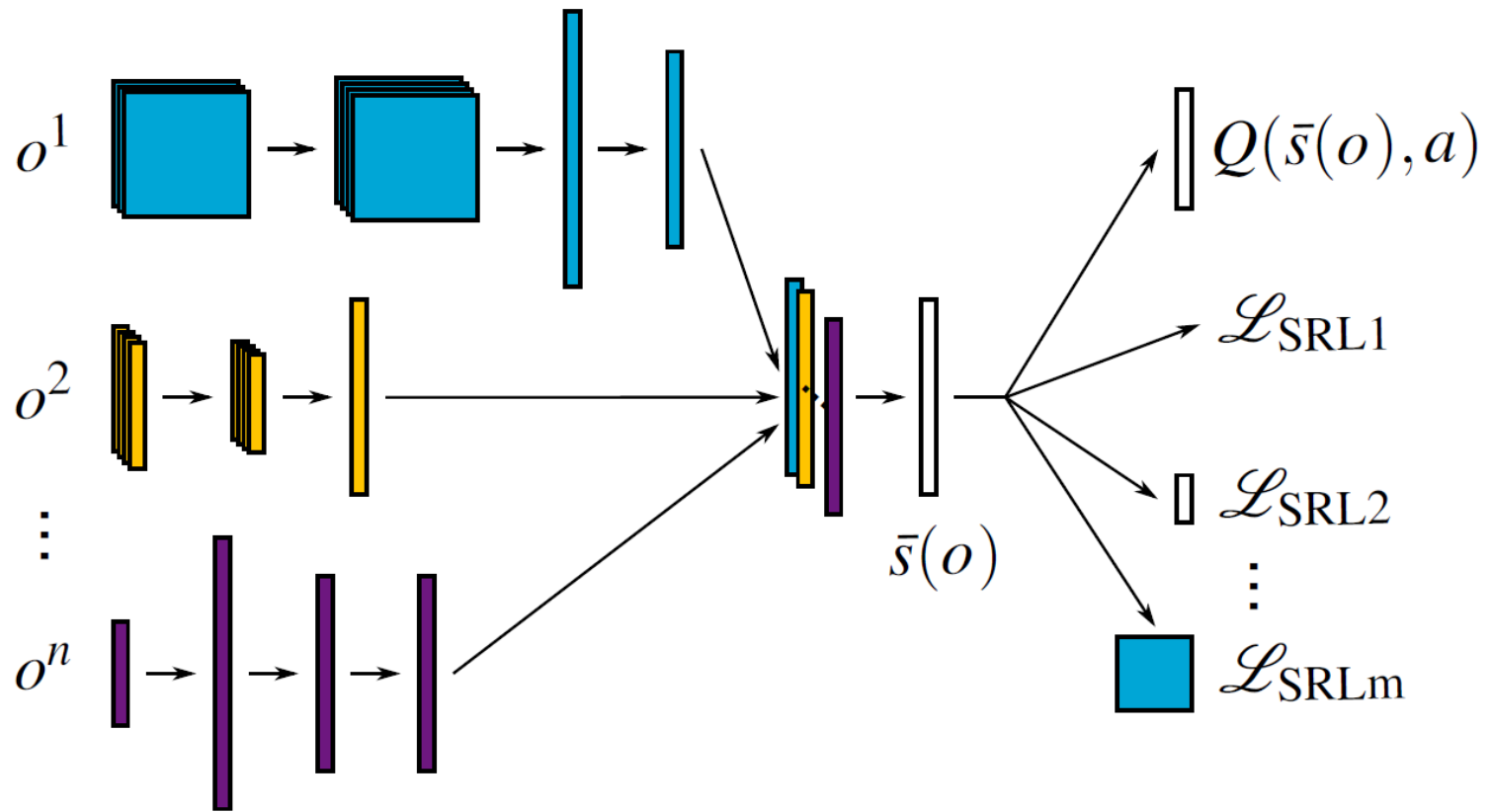
Tim de Bruin

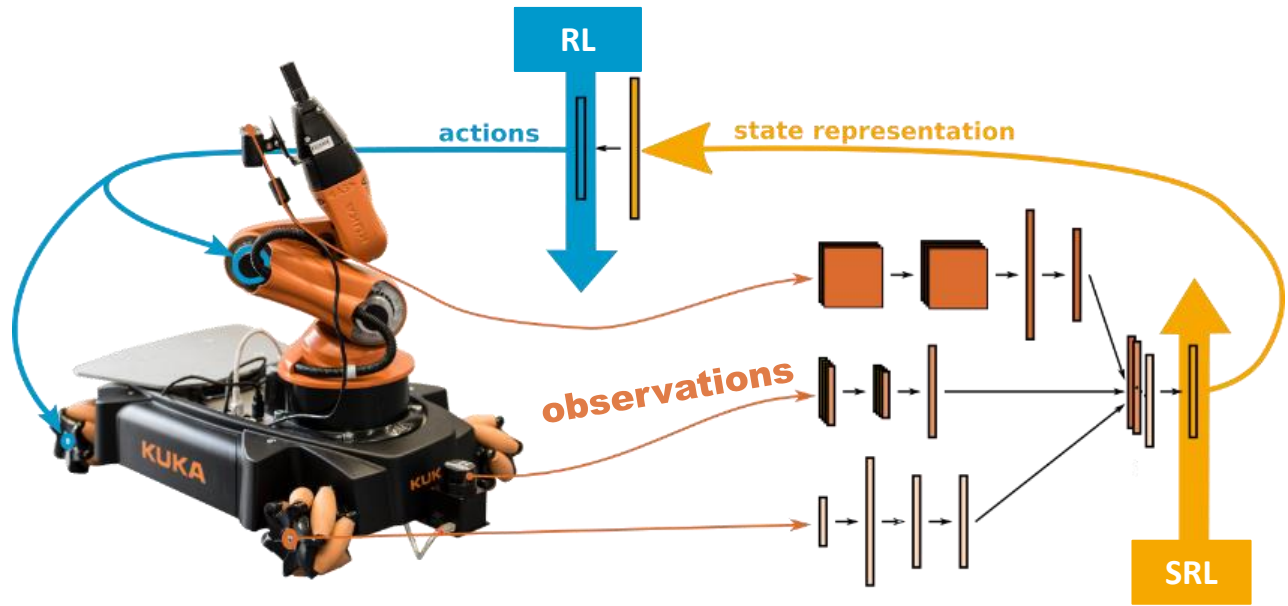


- Auto encoding
- Instantaneous reward prediction
- (Inverse) state dynamics
- Slowness and diversity
- Reinforcement Learning



Tim de Bruin

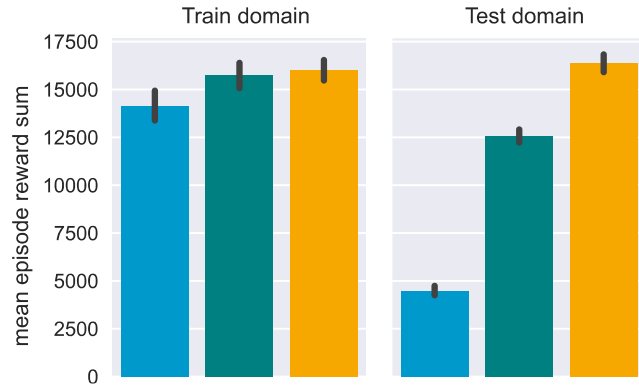




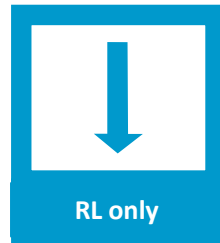
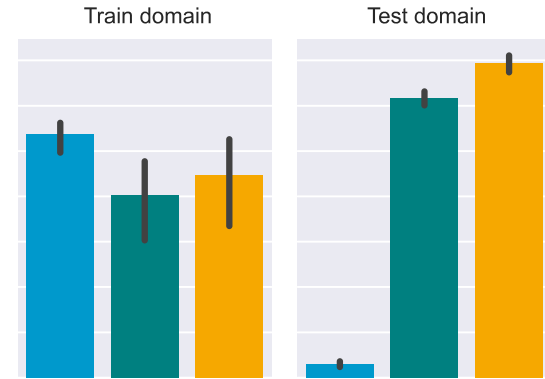
Experiment



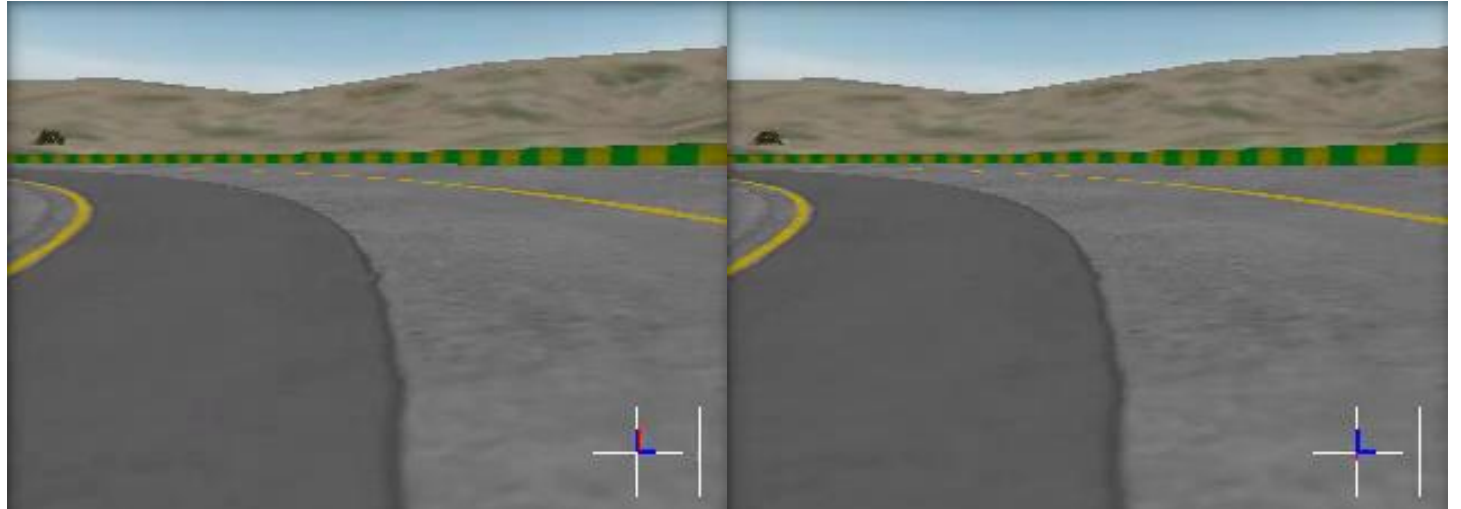
With pre-training



Without pre-training



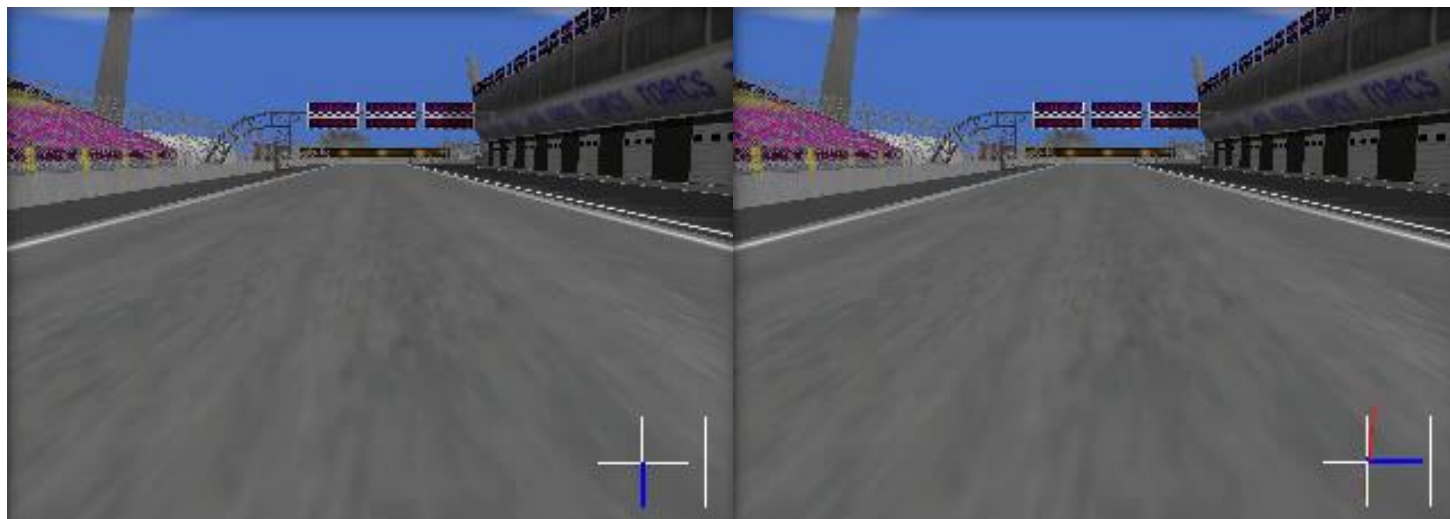
Train Track



RL only

RL + SRL

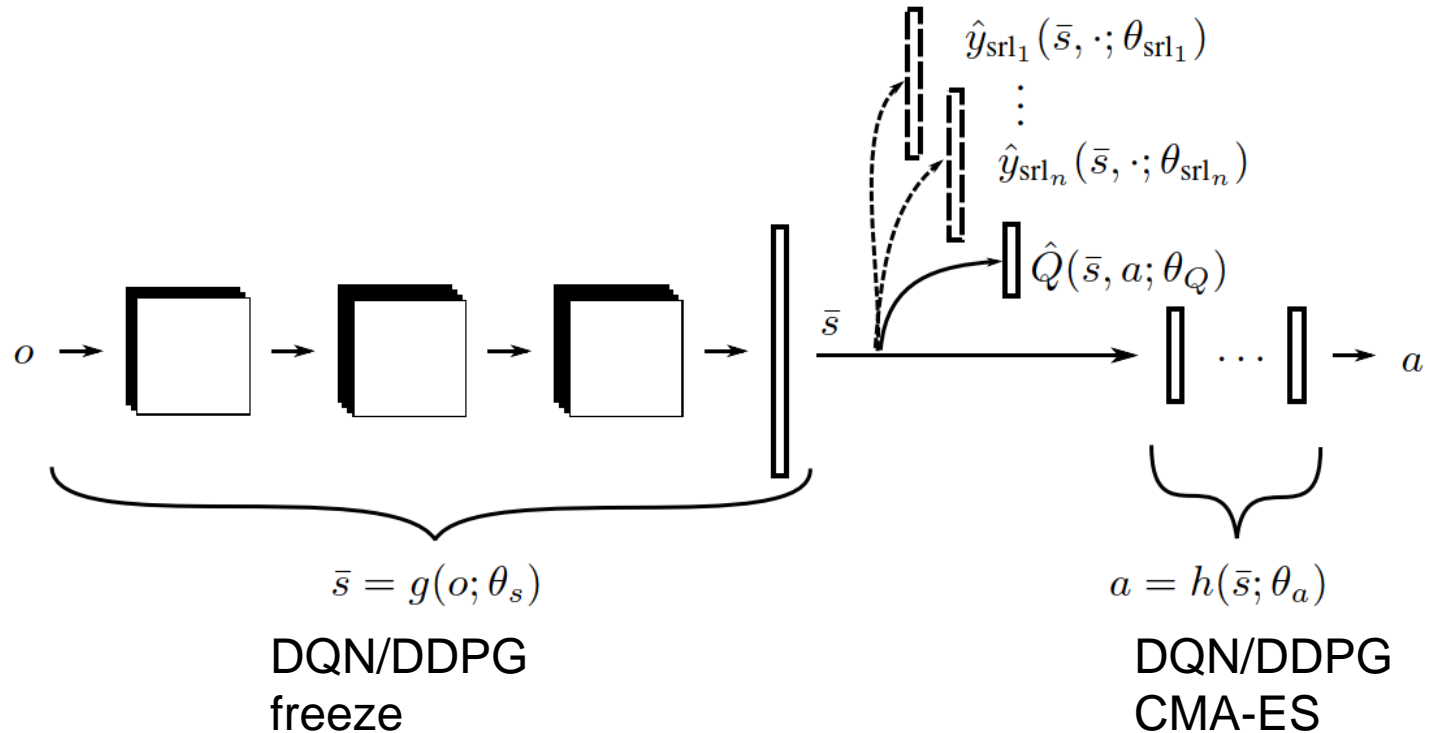
Test Track

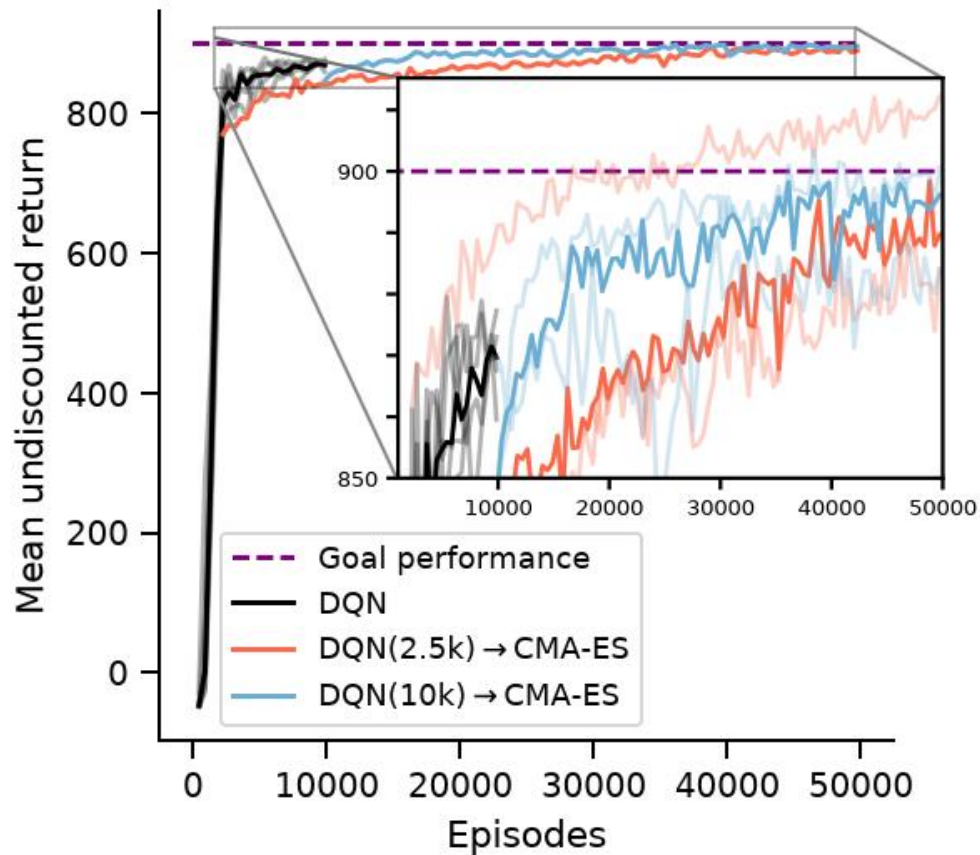
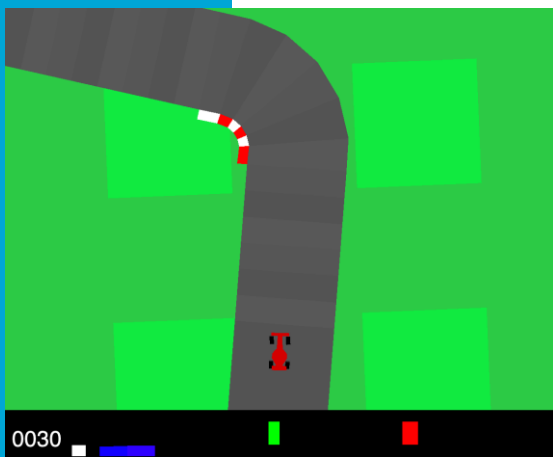


RL only

RL + SRL

Splitting SRL and RL





<https://youtu.be/D7zqglDkEq4>

*rotate box 180 degrees
counterclockwise*



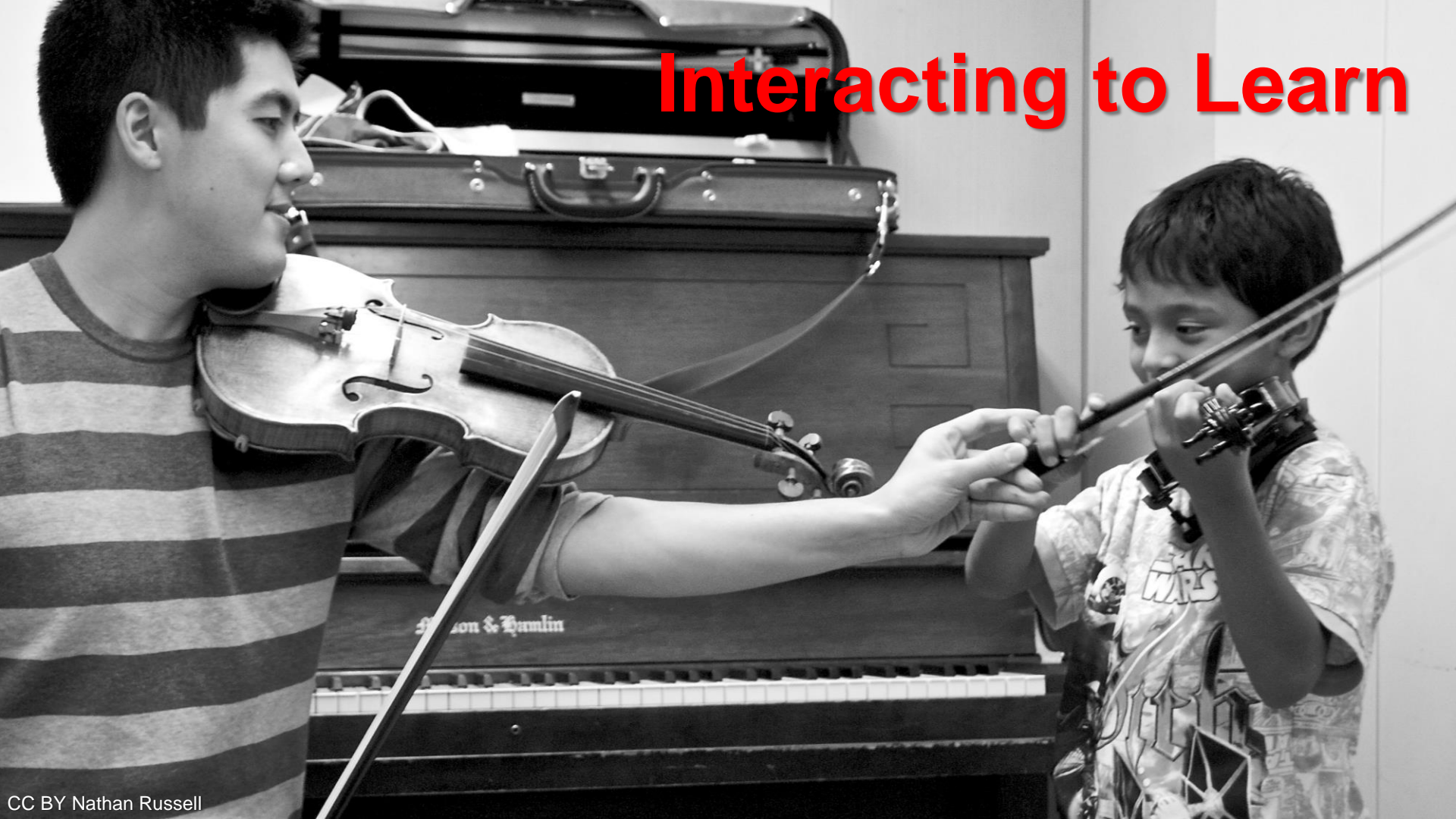
Linda
van der Spaa



latd.com

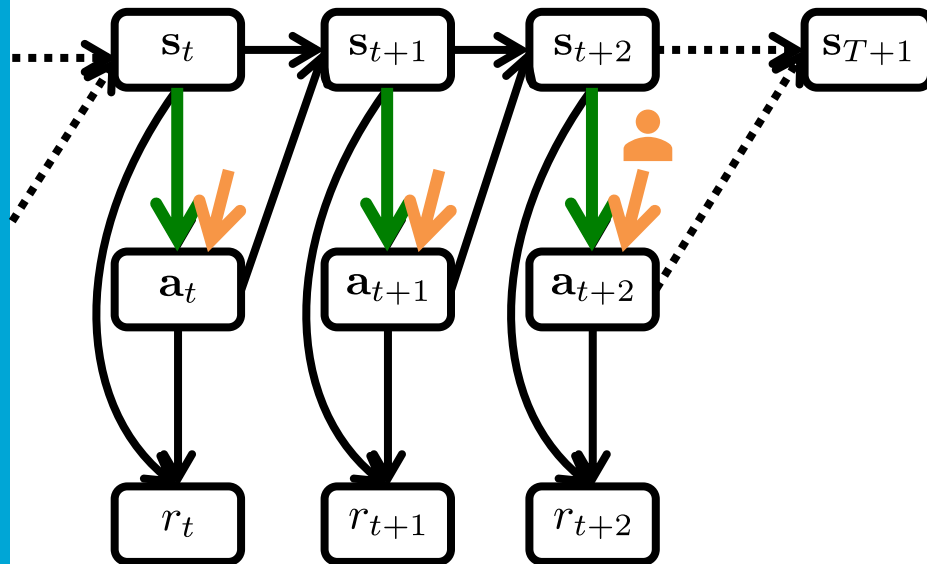
Questions so far?

Interacting to Learn



Reinforcement Learning: Exploration & Human Advice

Sutton & Barto 1998



states
 $s_{1:T+1}$

actions
 $a_{1:T}$

rewards
 $r_{1:T}$

policy
 $\pi(a_t | s_t, t, \theta)$

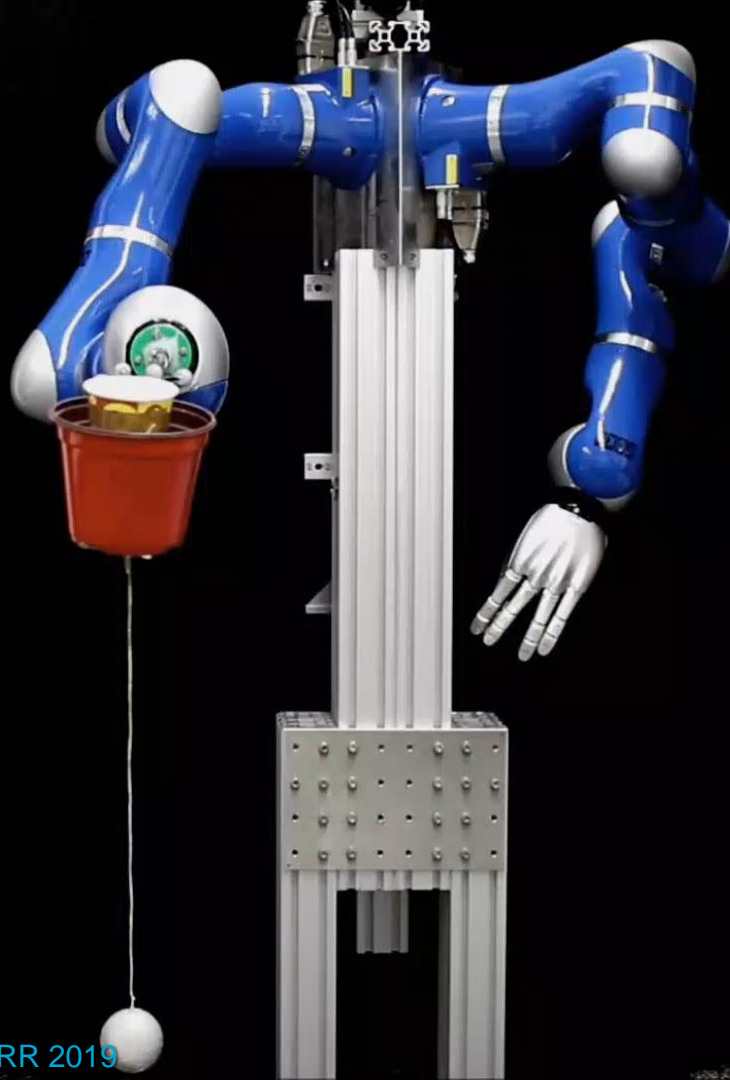
exploration
 ϵ

Model human advice as exploration

<https://youtu.be/ptsINZdum2s>



Carlos Celemin



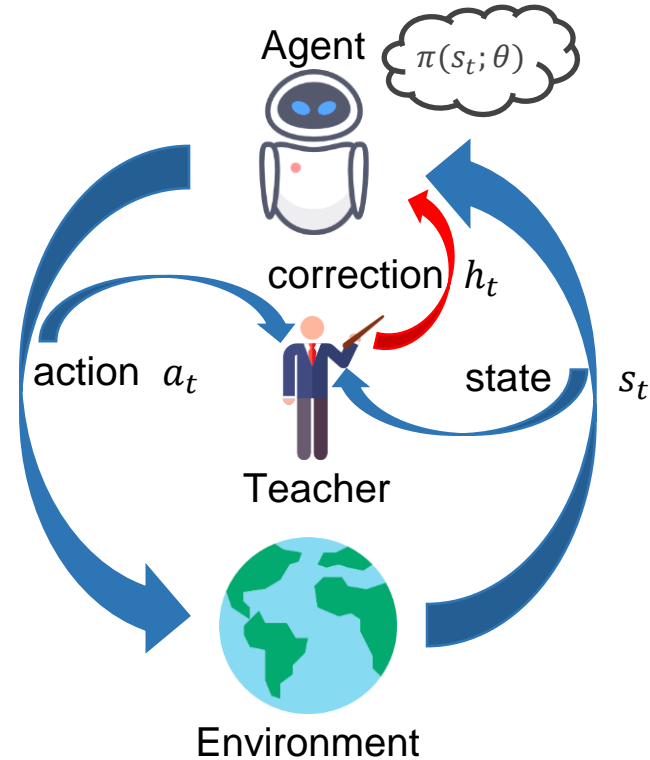
Back to Imitation Learning...

Interactive Learning: COACH

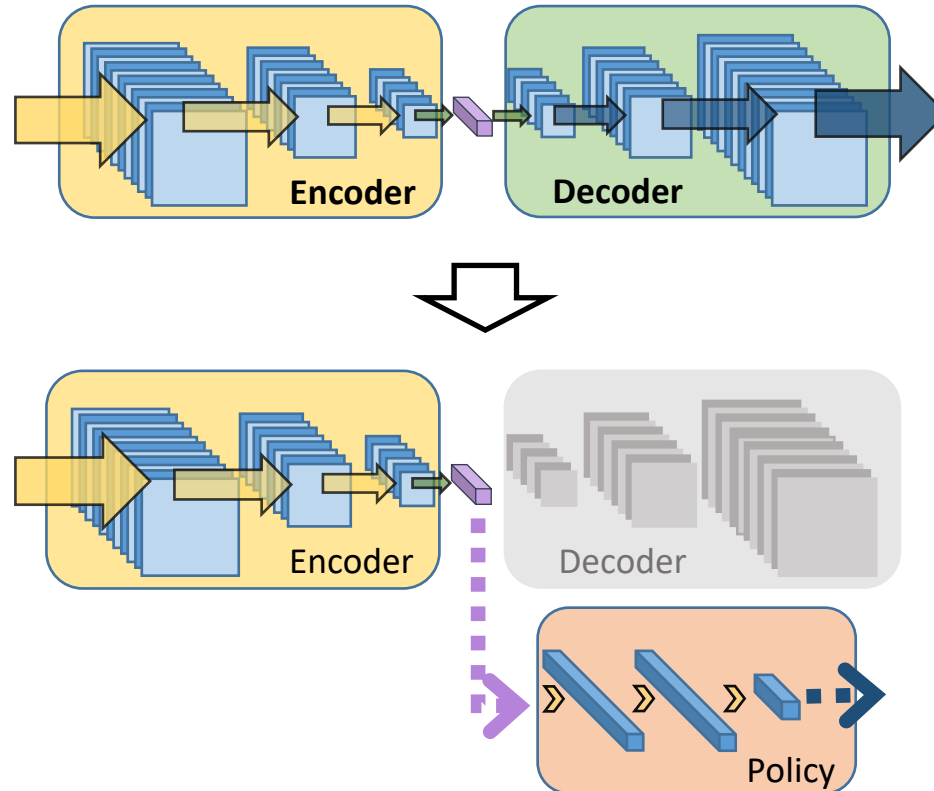


Rodrigo Pérez

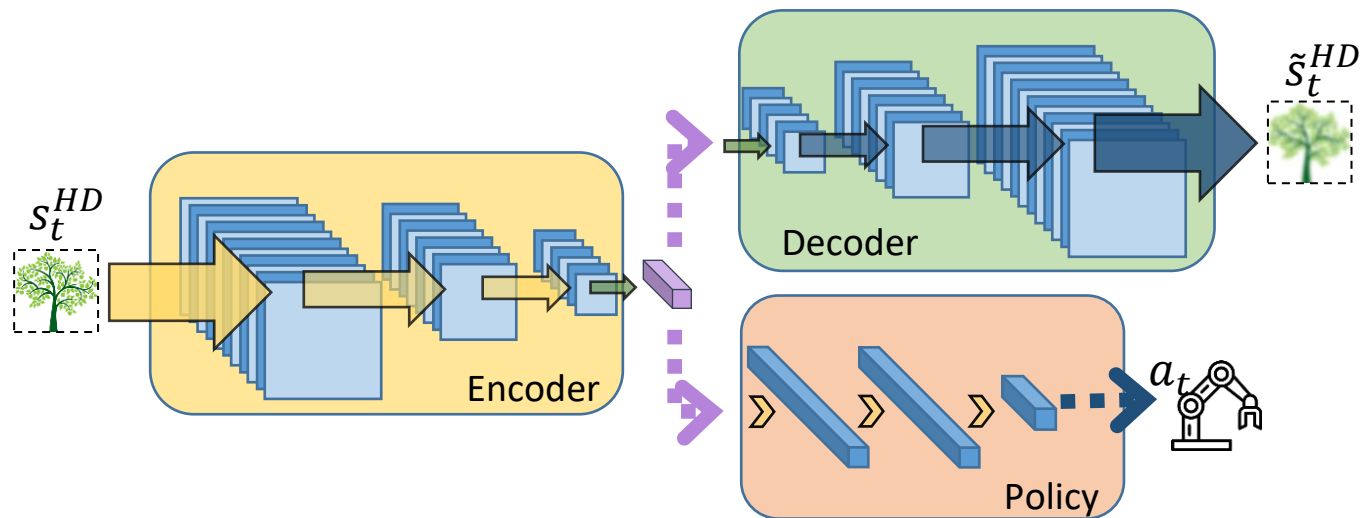
Carlos Celemin

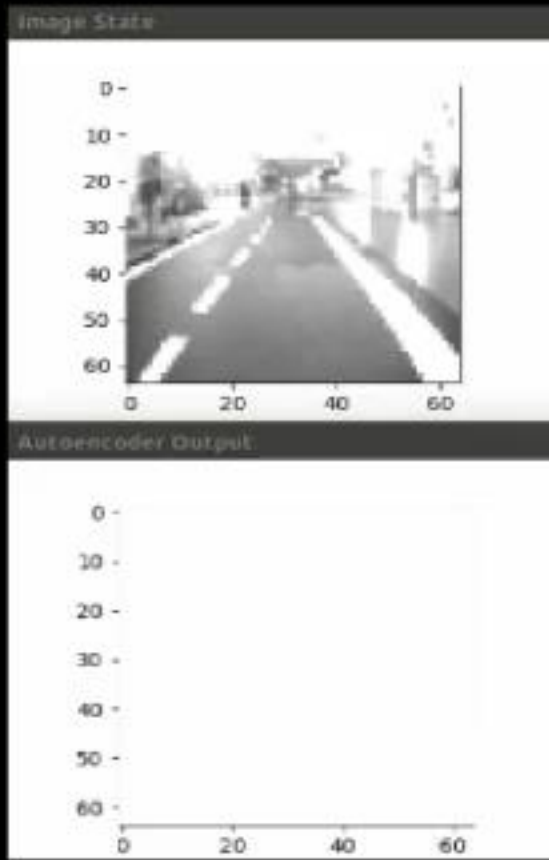


Deep COACH - Offline



Deep COACH - Online







Rodrigo Pérez

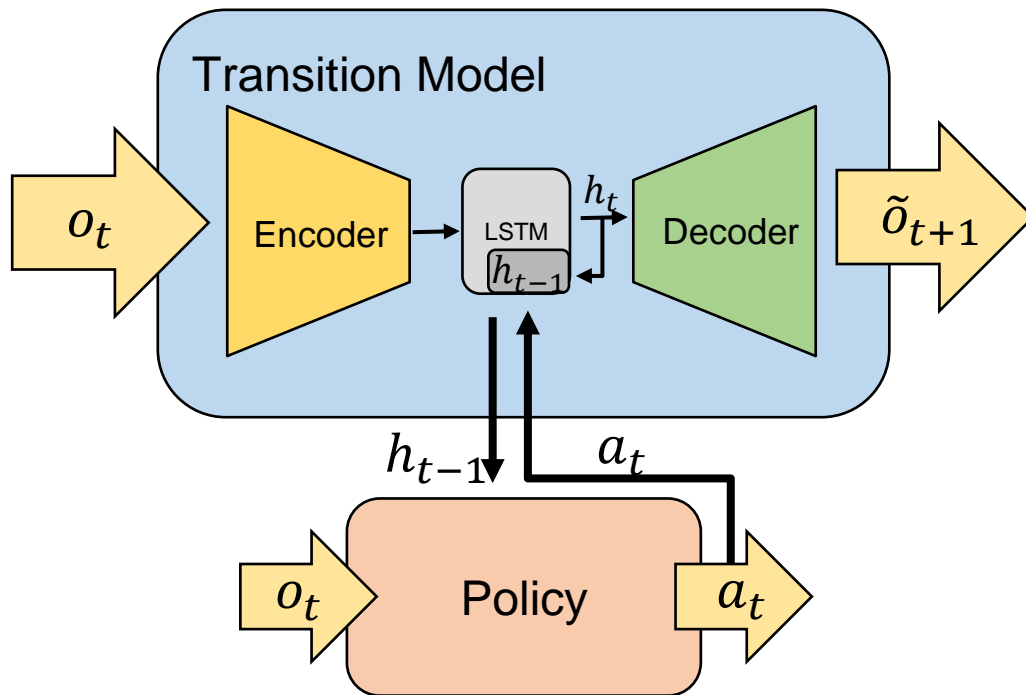


Giovanni Franzese



Carlos Celemin

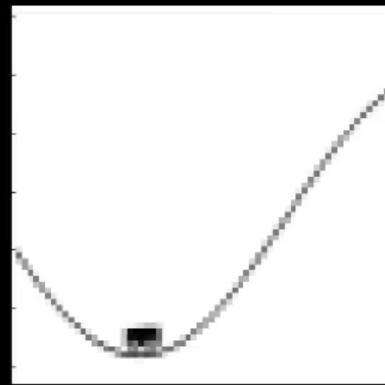
Deep COACH - Memory



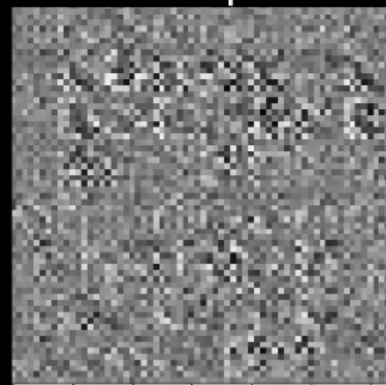
x4

Mountain Car

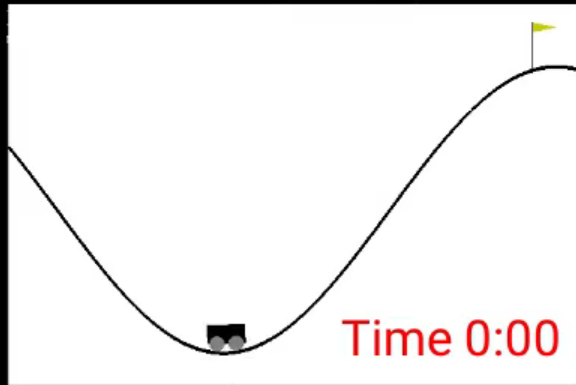
Network Input



Observation prediction



Environment

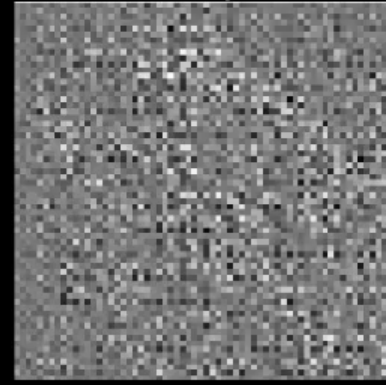


Swing-up Pendulum

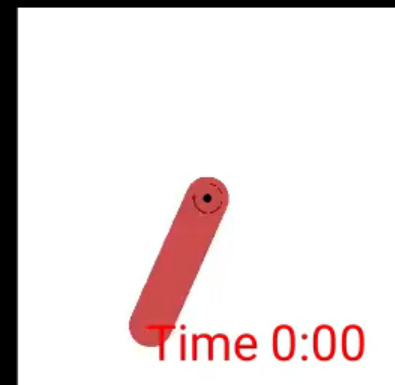
Network Input



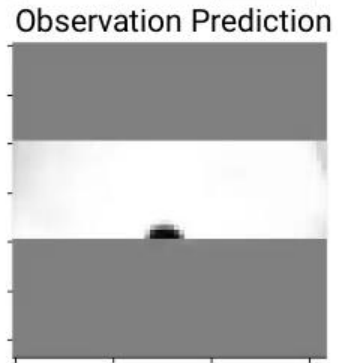
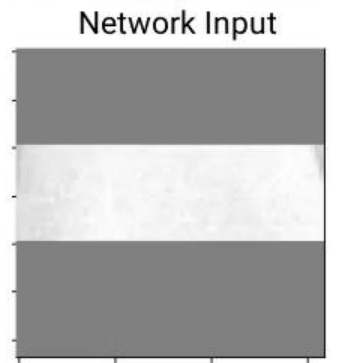
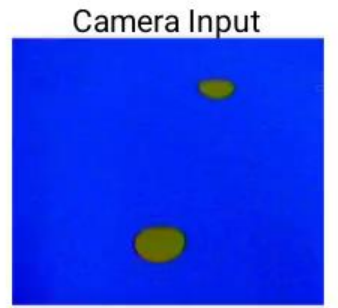
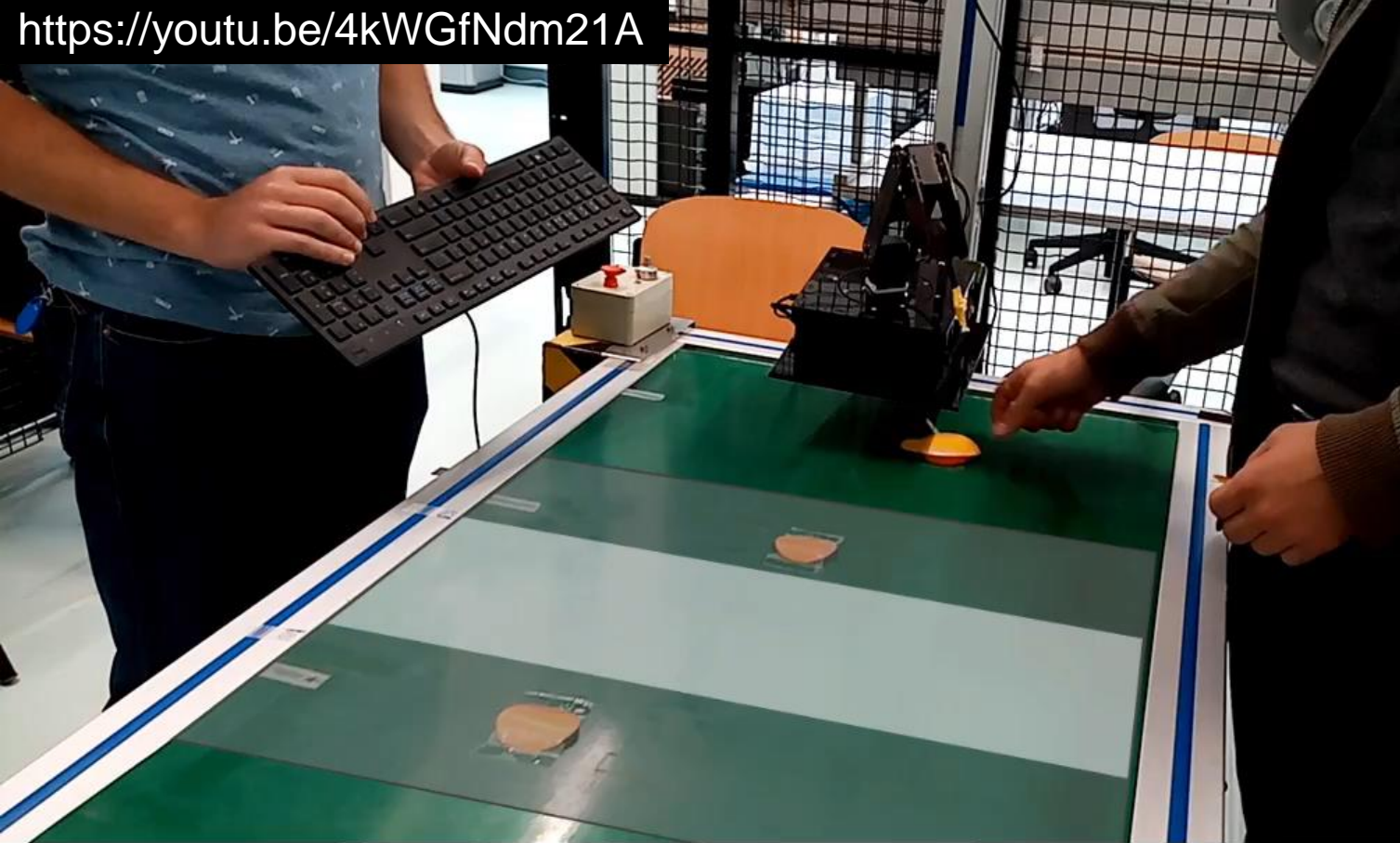
Observation prediction



Environment



<https://youtu.be/4kWGfNdm21A>



Start Training

x4

Pérez, Celemin Paez, Franzese, Ruiz-del-Solar, & Kober, RAM 2020

<https://youtu.be/apxeeW-hm6I>

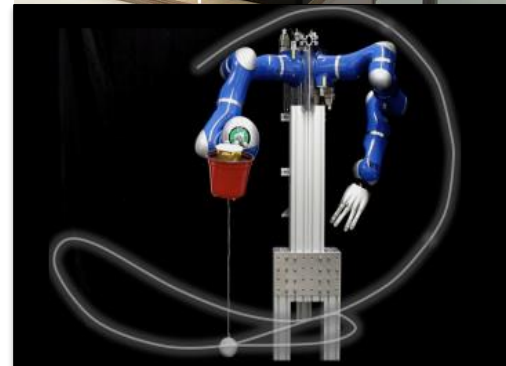


Walt Disney Studios – Real Steel

Conclusion

Summary

- Robot learning constraints
 - “Small” data
 - Safety
- Interactive Learning
 - Efficient & intuitive
- Challenges
 - Uncertainty
 - Variations



Questions?

j.kober@tudelft.nl

www.jenskober.de

