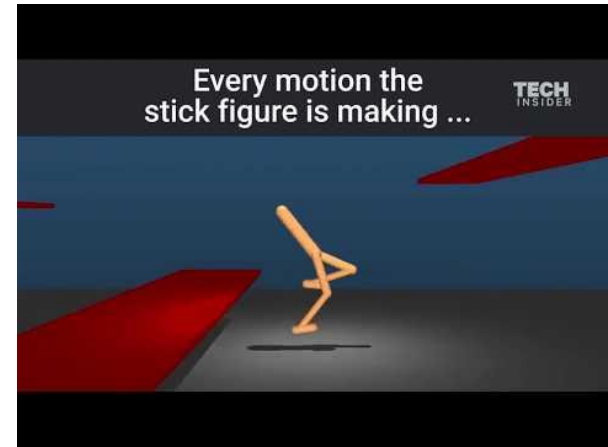
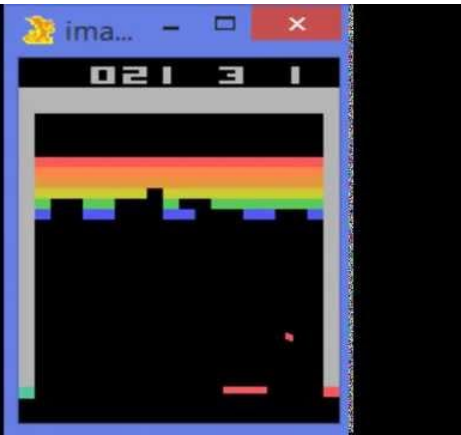


Reinforcement Learning for Robotics

Iretiayo Akinola



Reinforcement Learning for Robotics

Iretiayo Akinola

<http://www.cs.columbia.edu/~iakinola/>



COLUMBIA
UNIVERSITY

Robot Learning

Robotics: see, think, act

Direct programming can be hard for different tasks

- Degree of structure and consistency
- Perception
- Manipulation
- Deformation

Vacuum robots

Lawn mowing

Pool cleaning

Manufacturing robots

Home-cleaning robot

Cooking Robot

Laundry Robot

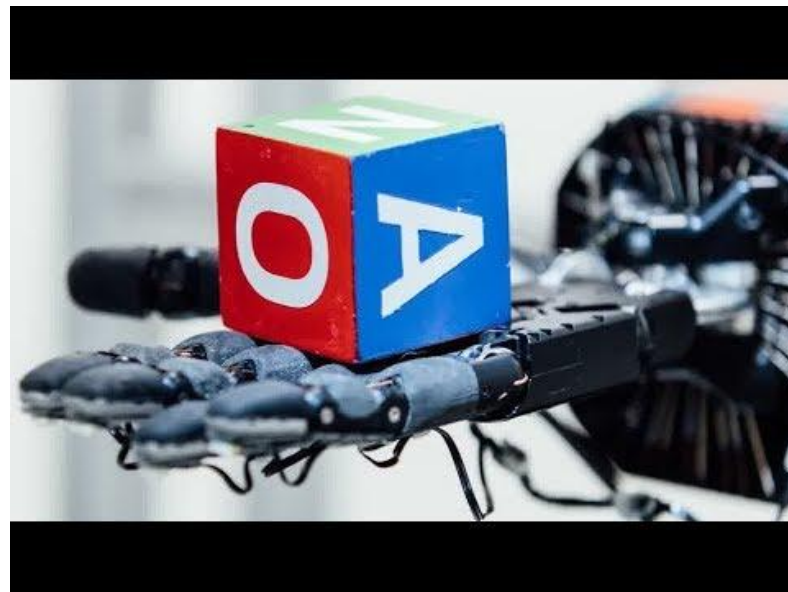
Warehouse Robot

Robot Learning

Robotics: see, think, act

Direct programming can be hard for different tasks

- Degree of structure and consistency
- Perception
- Manipulation
- Deformation



Vacuum robots

Lawn mowing

Pool cleaning

Manufacturing robots

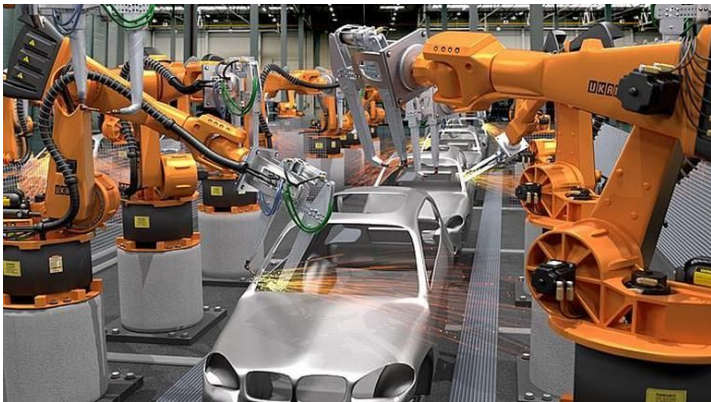
Home-cleaning robot

Cooking Robot

Laundry Robot

Warehouse Robot

Robot Learning

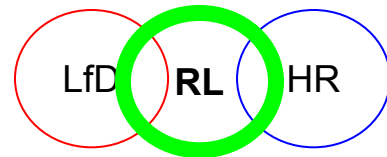


Manufacturing robots



Cooking Robot

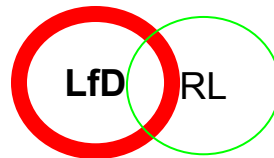
Robot Learning



An alternative to direct programming is learning-based methods:

- Learning from Demonstration
- Reinforcement Learning

Learning from Demonstration (LfD)



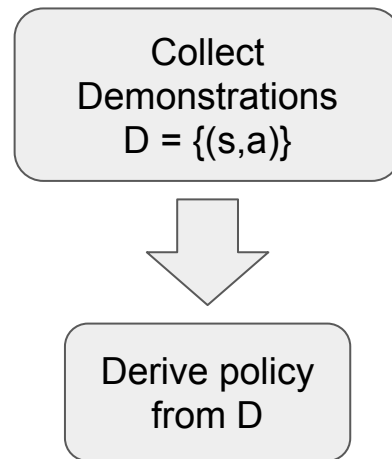
LfD Dataset (D): a set of state-action pairs

Goal: Learn $\pi(a_t|s_t)$ – policy

Assumptions: Human Teacher exists, Demonstration is possible

Key Considerations:

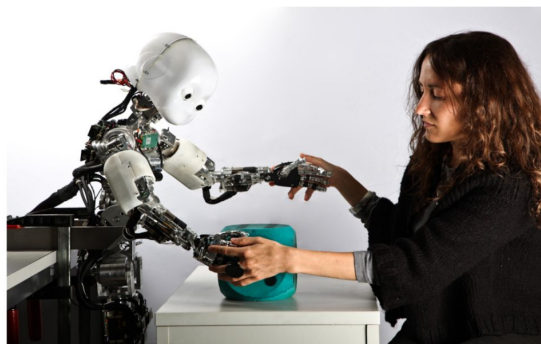
- Demonstration mode
- State representation
- Policy Derivation method
 - Supervised learning/Function approximation



Learning from Demonstration (LfD)

LfD

- Learning from Demonstration data
 - Learns a mapping from state to action
- Demonstration modes
 - Teleoperation
 - Kinesthetic teaching (e.g. in motion trajectory learning)
 - Camera recording a human teacher
 - Robotic teachers

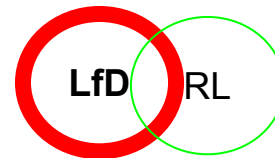


Kinesthetic Teaching



Motion Capture

LfD

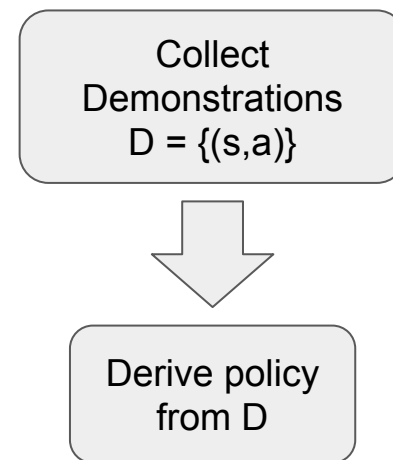


Pros

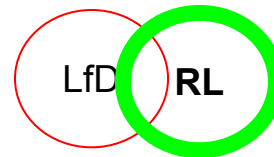
- Supervised Learning
- No need for manual reward function
- Exploration not an issue

Cons

- Might not generalize well- Covariate shift
- Volume of demonstration
- Some tasks are difficult to demonstrate
- Suboptimal Demonstrations. Limited human patience and inconsistent user input
- Performance of the robot can be limited by that of the teacher.

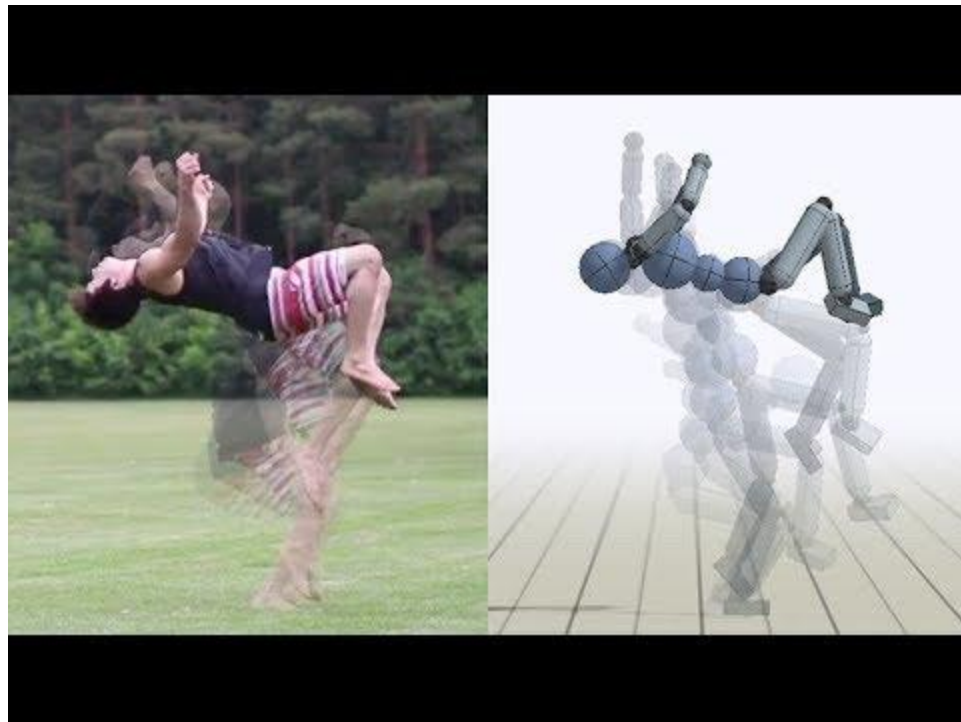


Robot Learning

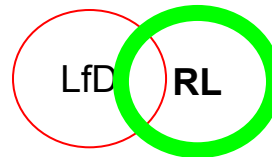


An alternative to direct programming is learning-based methods:

- Learning from Demonstration
- Reinforcement Learning
- **Hybrid**

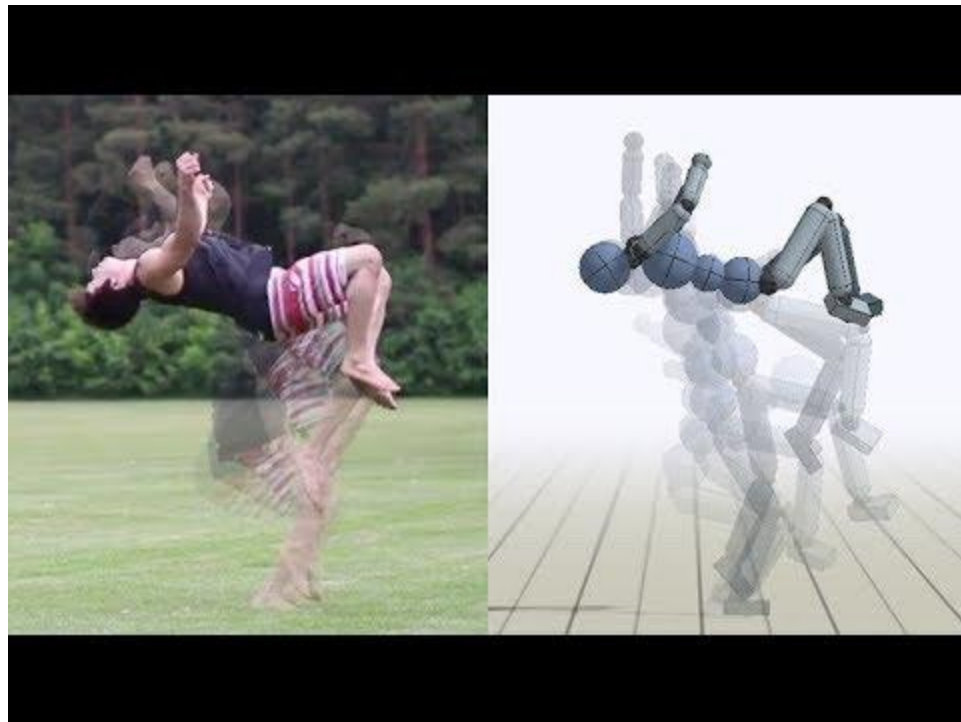


Robot Learning

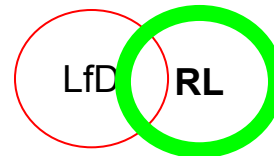


An alternative to direct programming is learning-based methods:

- Learning from Demonstration
- Reinforcement Learning
- **Hybrid**



Robot Learning



An alternative to direct programming is learning-based methods:

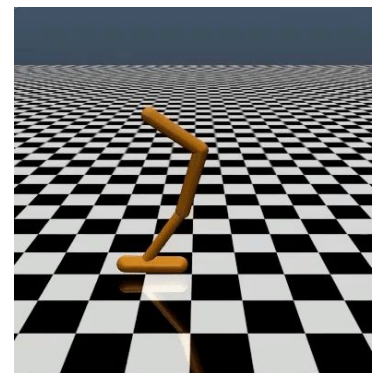
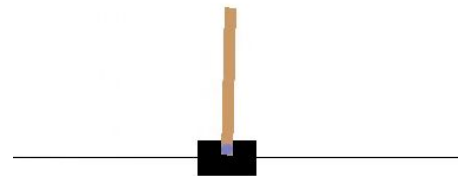
- Learning from Demonstration
- **Reinforcement Learning**

Reinforcement Learning

- Learning by trial and Error
- Maximize cumulative rewards
- Learn a policy

Reinforcement Learning Progress

- Classical Control
- Games: Atari, Go
- Robotics
 - continuous space, complex transition dynamics, complex rewards

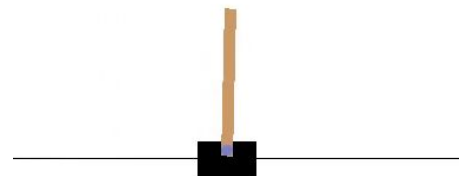


Reinforcement Learning

- Learning by trial and Error
- Maximize cumulative rewards
- Learn a policy

Reinforcement Learning Progress

- [Classical Control](#)
- Games: Atari, Go
- Robotics
 - continuous space, complex transition dynamics, complex rewards



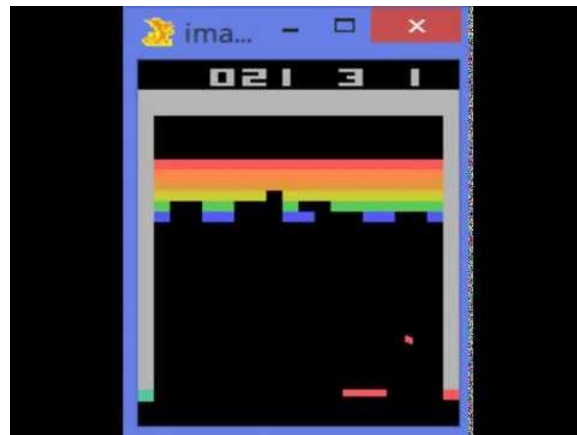
Reinforcement Learning

- Learning by trial and Error
- Maximize cumulative rewards
- Learn a policy



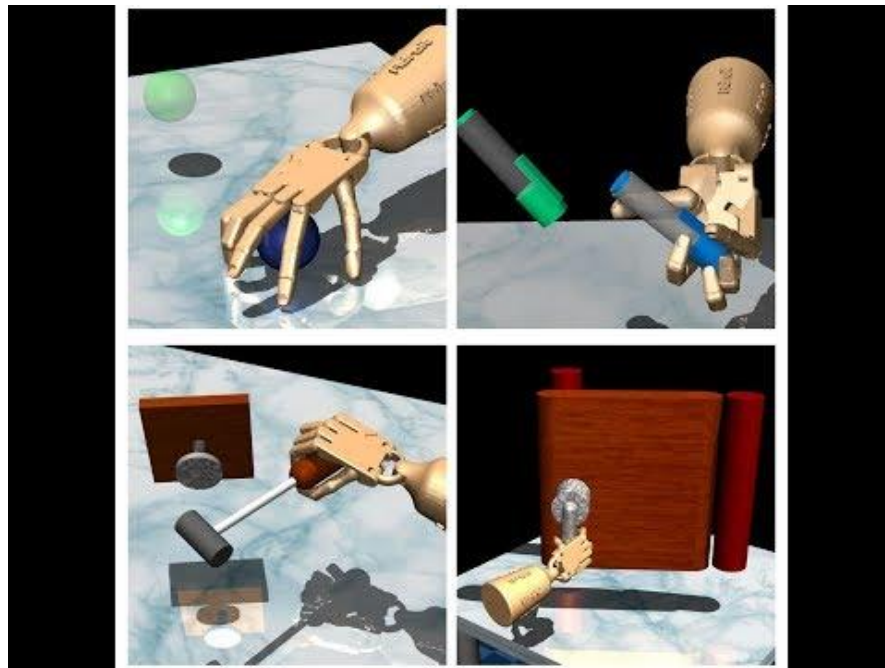
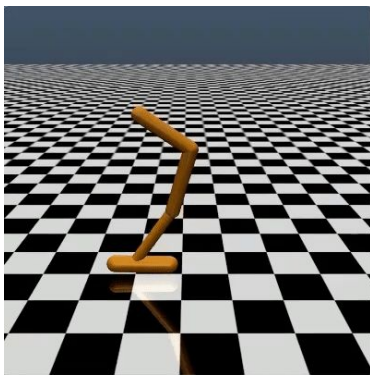
Reinforcement Learning Progress

- Classical Control
- [Games](#): Atari, Go
- Robotics
 - continuous space, complex transition dynamics, complex rewards



Reinforcement Learning

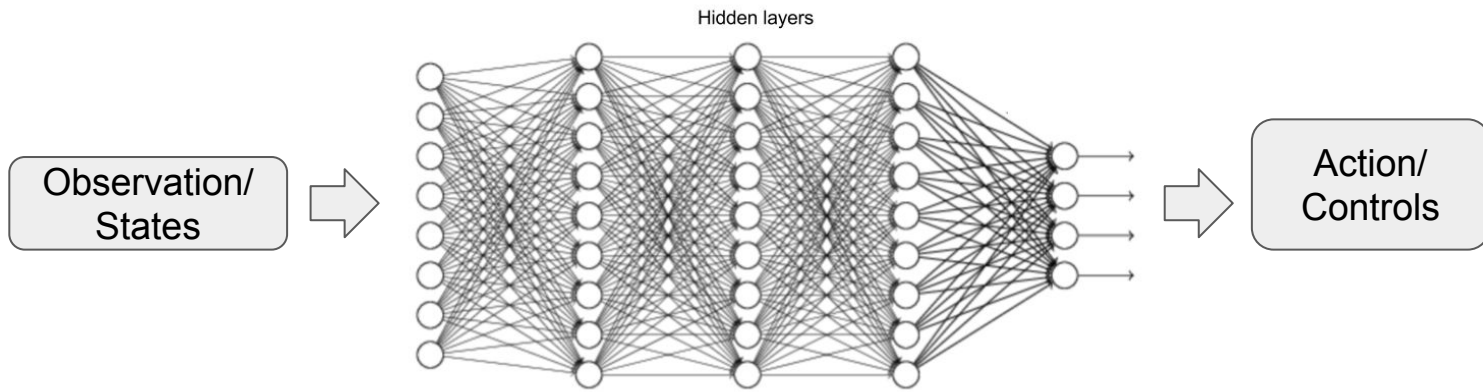
- Robotics
 - continuous space
 - complex transition dynamics
 - complex rewards



Reinforcement Learning in Robotics

Key Elements of recent Success

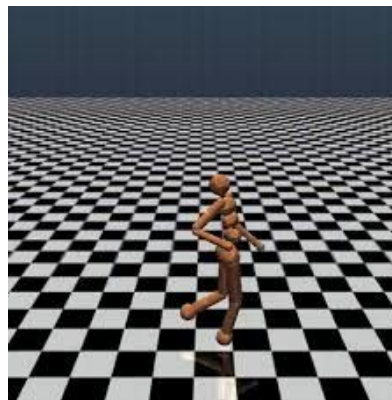
- **Deep Learning**
- Simulators (mujoco, bullet, roboschool, dart, gazebo, carsim)



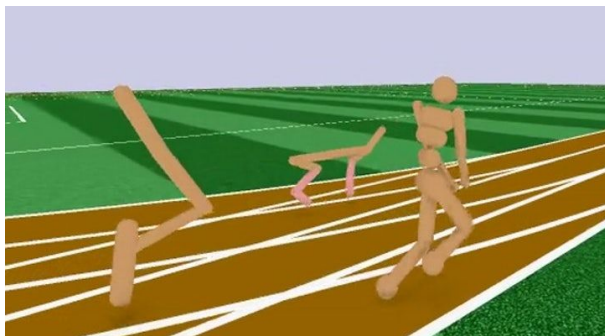
Reinforcement Learning in Robotics

Key Elements of recent Success

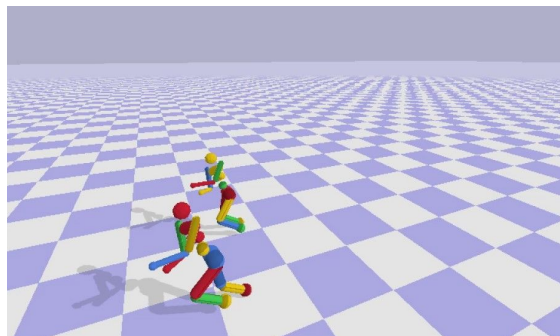
- Deep Learning
- **Simulators (mujoco, pybullet, roboschool, gazebo)**



mujoco



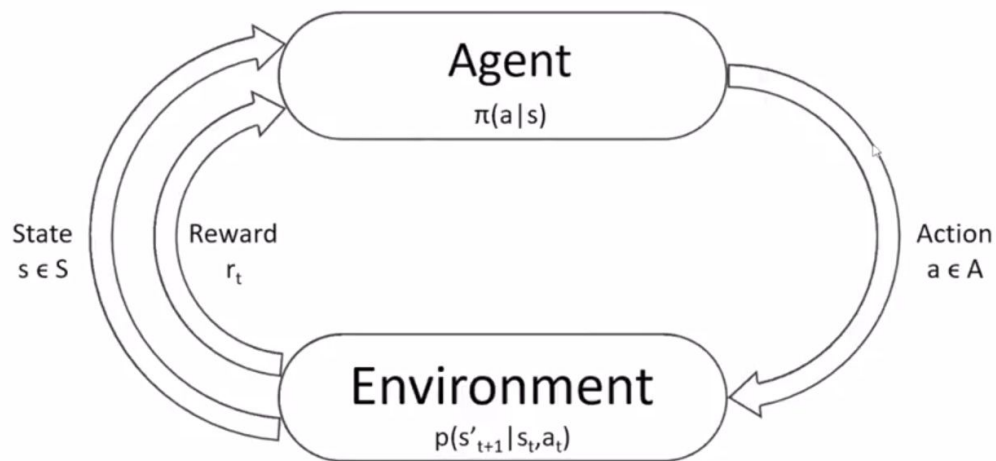
roboschool



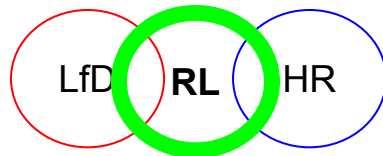
pybullet

Reinforcement Learning Formulation

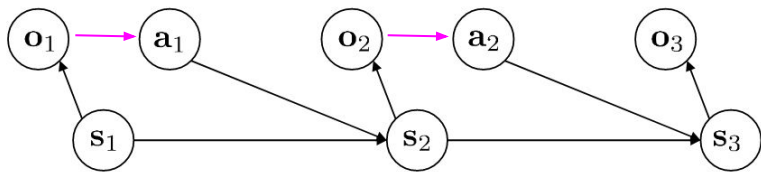
Markov Decision Process



Reinforcement Learning Formulation



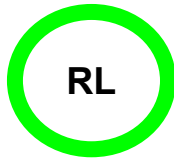
The goal of RL is to get a **policy**: $s_t \xrightarrow{\pi} a_t$



$$\underbrace{p_{\theta}(s_1, a_1, \dots, s_T, a_T)}_{\pi_{\theta}(\tau)} = p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

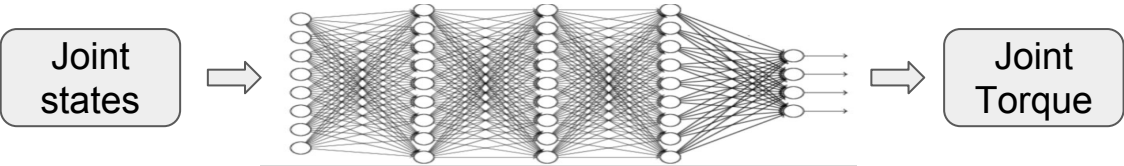
reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$



Reinforcement Learning Formulation

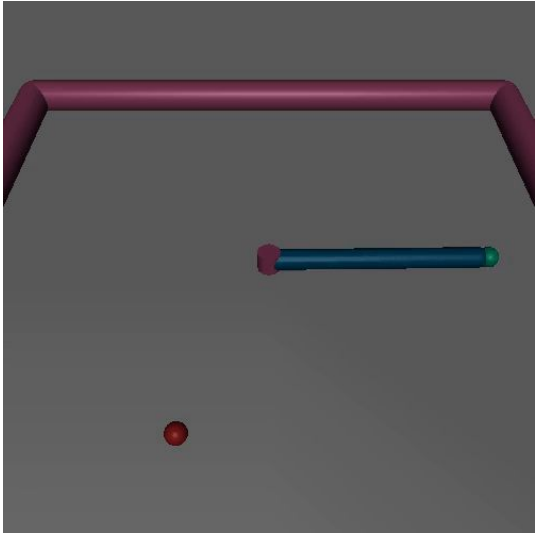
The goal of RL is to get a **policy**: $s_t \xrightarrow{\pi} a_t$



Reward function:

$$r := - \text{dist}(\text{goal}, \text{end-effector}) - \alpha \text{ magnitude}(\text{torque})$$

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$



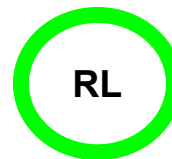
reward function

Reinforcement Learning in Robotics

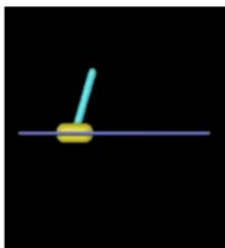
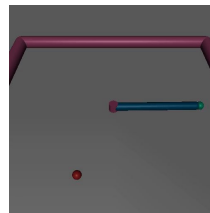
- Deterministic Policy Gradient Algorithms (David Silver et al. 2014)
 - DPG Algorithm
 - Experiments: continuous bandit, **pendulum**, mountain car, 2D puddle world and **Octopus Arm**
- Continuous Control with Deep Reinforcement Learning (Lillicrap et al. 2016)
 - DDPG



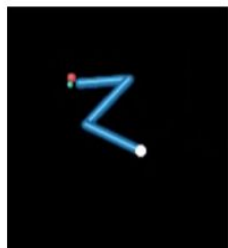
Quiz 1: Write down a reward function for each



Example: $r := -\text{dist}(\text{goal}, \text{end-effector})$



a



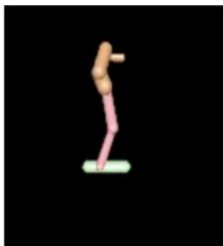
b



c



d



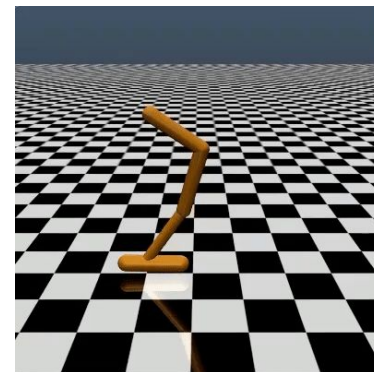
e



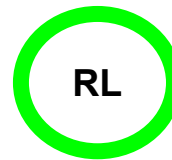
f



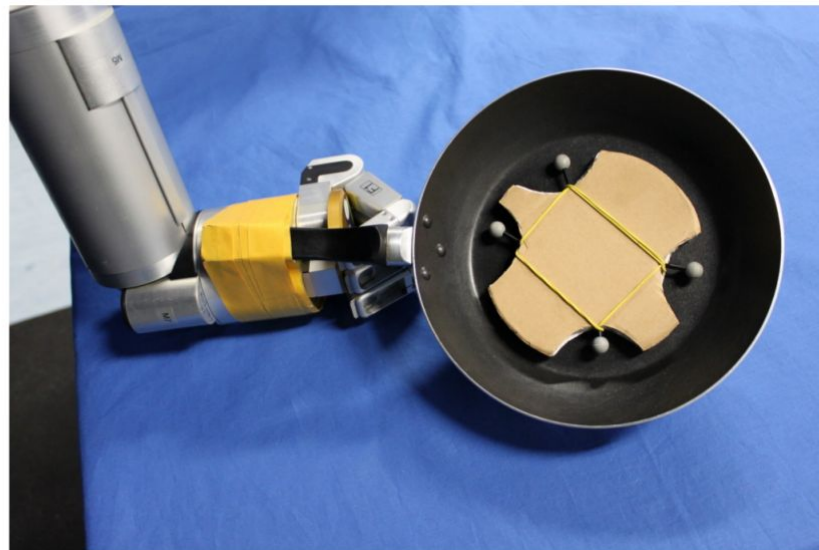
g



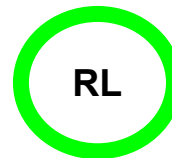
Reinforcement Learning in Robotics



Pancake Flipping Task



Reinforcement Learning in Robotics



Pancake Flipping Task

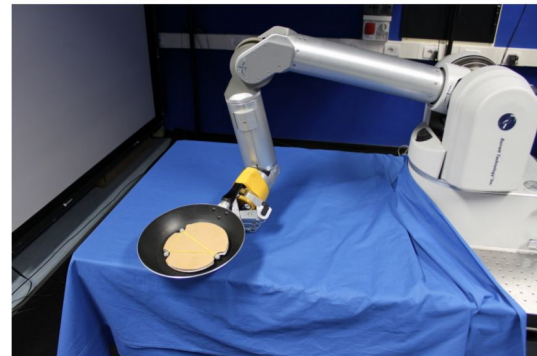


Reinforcement Learning in Robotics

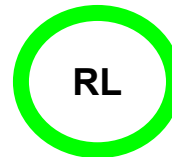
Pancake Flipping Task

$$R(\tau) = w_1 \left[\frac{\arccos(v_0 \cdot v_{t_f})}{\pi} \right] + w_2 e^{-\|x^p - x^F\|} + w_3 x_3^M$$

- Reward
 - Positional reward
 - Orientational reward



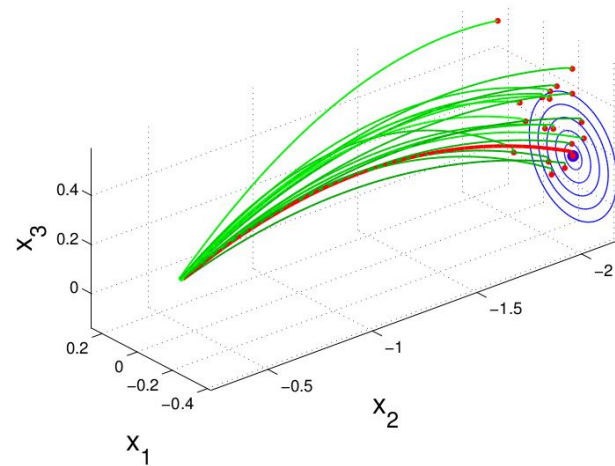
Reinforcement Learning in Robotics



Archery Task



$$R(\tau) = e^{-\|\hat{r}_T - \hat{r}_A\|}$$





Reinforcement Learning

Pros

- Fully autonomous
- Skills not explicit coded
- No demonstration needed

Cons

- Reward definition: Where does the R come from?
- Curse of Dimensionality: Exploration cost increases exponentially with dimension
- Generalization issues: Simulation to Real??
- Convergence



RL Class of Methods

- Value Iteration methods (Q-Learning, SARSA)
- Policy Gradient Methods (REINFORCE)
- Actor-Critic Methods (DDPG, TRPO, PPO, A3C)
- Model-based RL (Guided Policy Search, Dyna)

RL Class of Methods

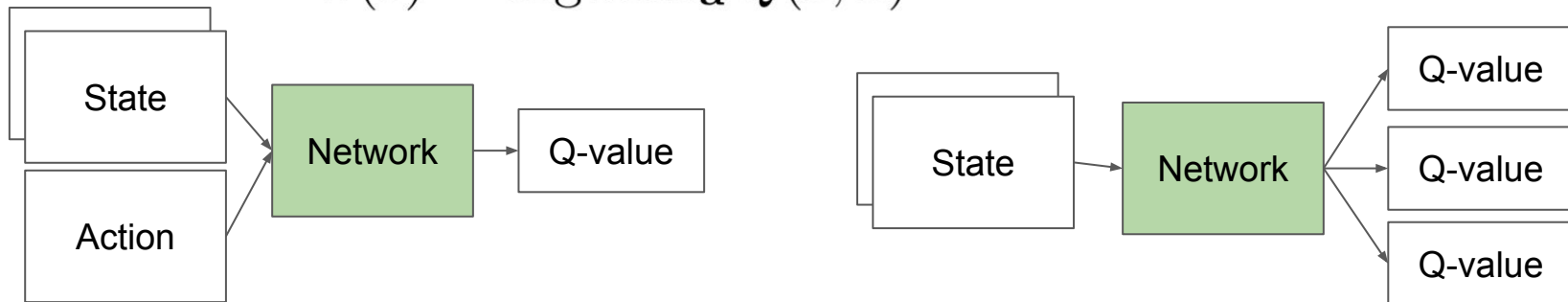
- Value Iteration methods (**Q-Learning**, SARSA)

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

$$\pi(s) = \arg \max_a (r + V(s'_s, a))$$

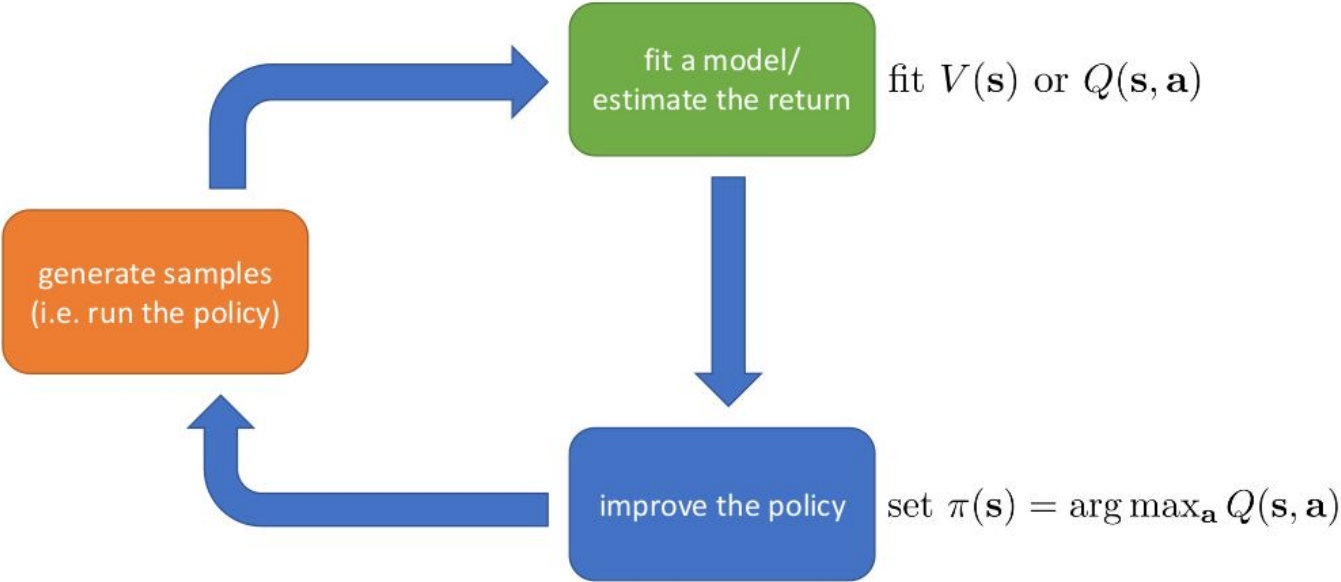
$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

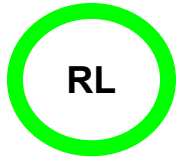
$$\pi(\mathbf{s}) = \arg \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$$



RL Class of Methods

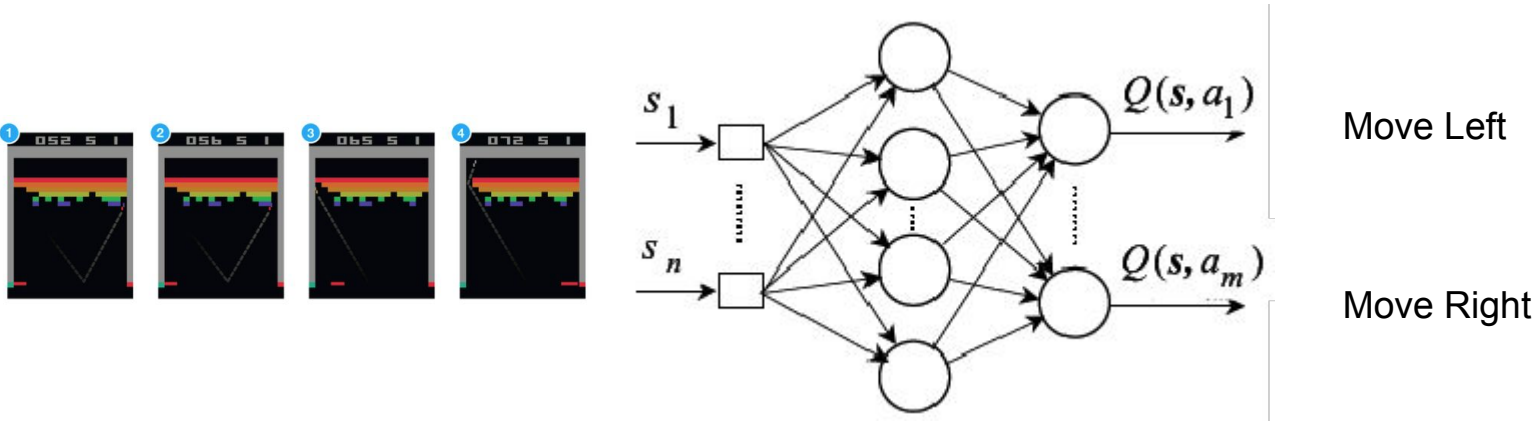
- Value Iteration methods (**Q-Learning**, SARSA)



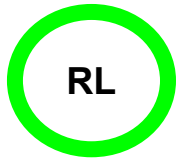


RL Class of Methods

- Value Iteration methods (**Q-Learning**, SARSA)
- DQN: Playing Atari with Deep Reinforcement Learning (Mnih et al 2015)



Works well with Games- choose between discrete actions but robotics need continuous actions



RL Class of Methods

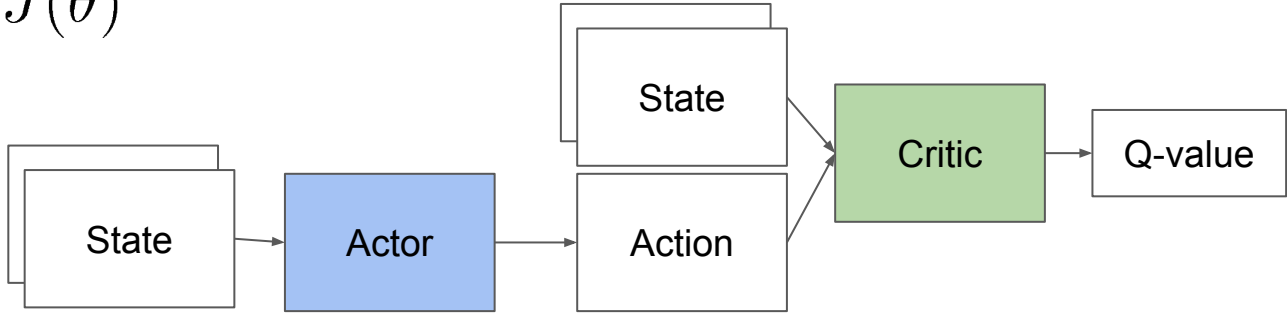
- Actor-Critic Methods (**DDPG**, A3C, TRPO, PPO)

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

$$\pi_\theta(a|s) = \mathbb{P}[a|s, \theta]$$

$$J(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a) Q(s, a)$$

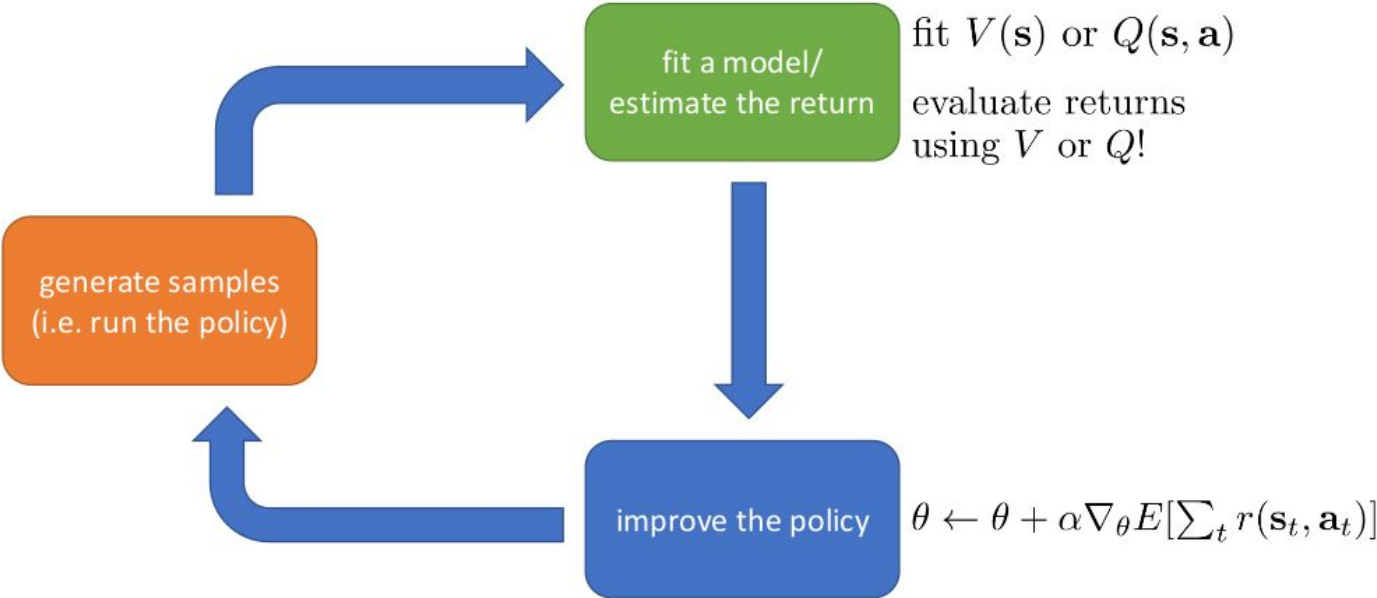
$$\theta := \overset{s}{\theta} + \alpha \nabla_{\overset{a}{\theta}} J(\theta)$$



Continuous control with deep reinforcement learning (Lillicrap et al Deepmind 2016)

RL Class of Methods

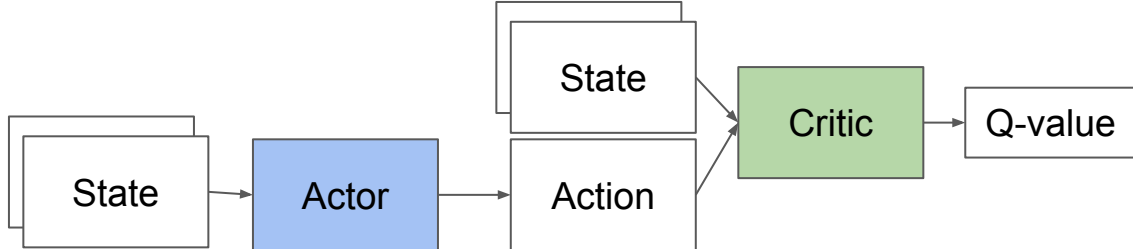
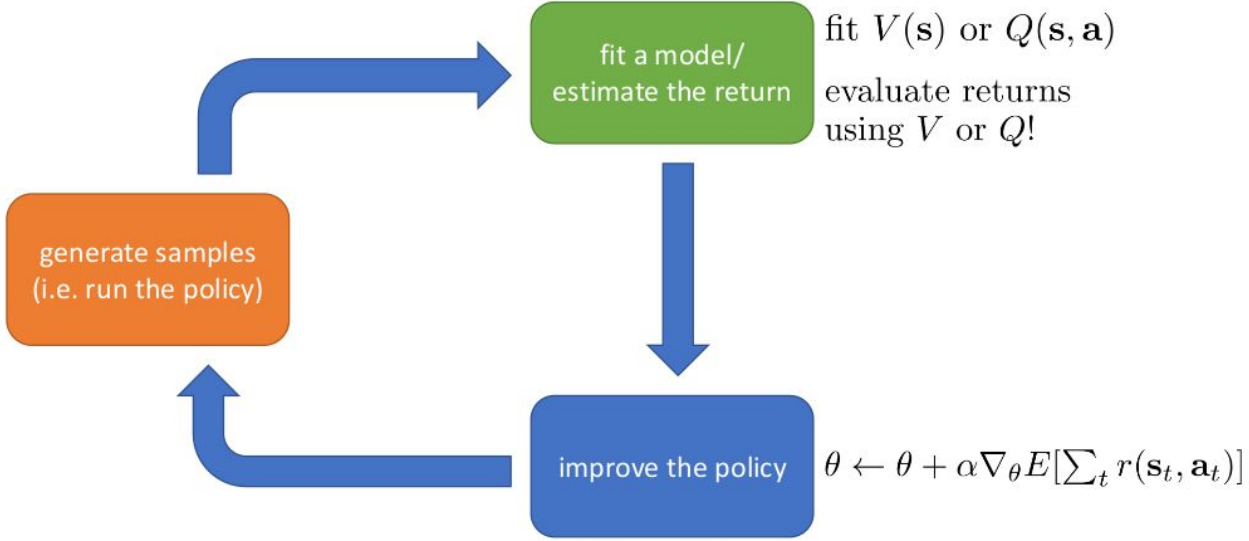
- Actor-Critic Methods (**DDPG**, TRPO, PPO, A3C)

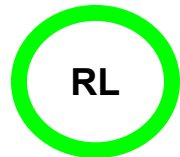




RL Class of Methods (Actor/Critic)

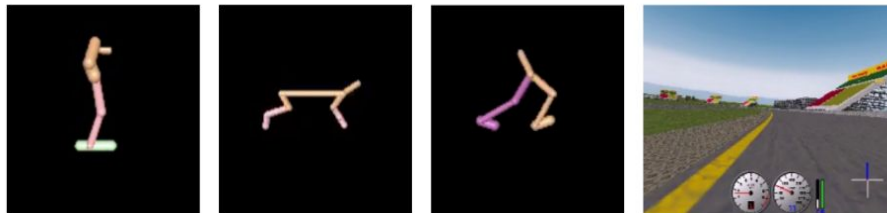
- A



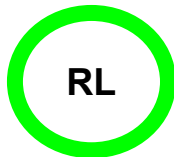


RL Class of Methods

- Actor-Critic Methods (**DDPG**)
- Deterministic Policy Gradient Algorithms (Silver et al 2014)
 - Critic: linear function, Actor: Gaussian policy
- Continuous control with deep reinforcement learning (Lillicrap et al 2016)
 - Critic: NN, Actor: NN

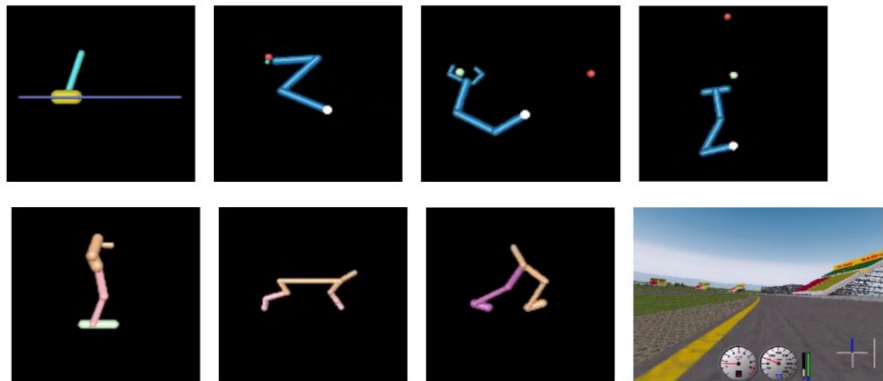


Case-Study: Actor-Critic Methods (DDPG)

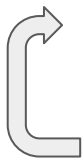


CONTINUOUS CONTROL WITH DEEP REINFORCEMENT LEARNING

Timothy P. Lillicrap*, Jonathan J. Hunt*, Alexander Pritzel, Nicolas Heess,
Tom Erez, Yuval Tassa, David Silver & Daan Wierstra
Google Deepmind
London, UK
{countzero, jjhunt, apritzel, heess,
etom, tassa, davidsilver, wierstra} @ google.com

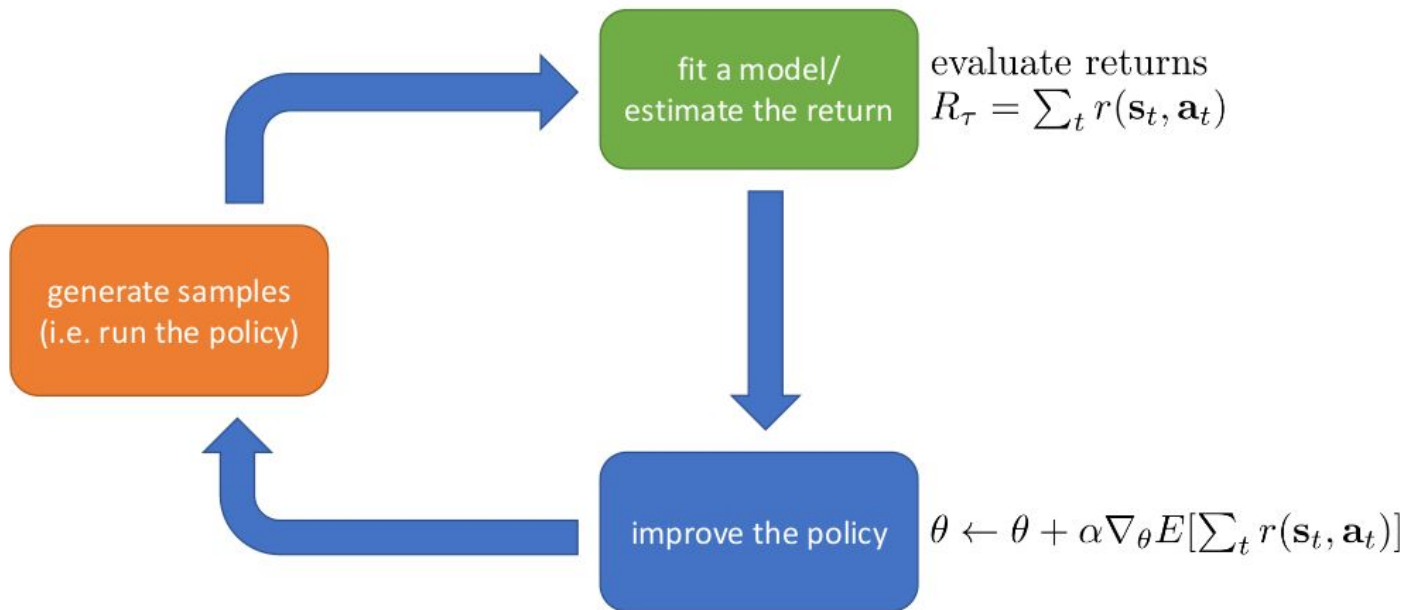


1. Initialize Actor and Critic networks
2. Generate samples from actor policy
3. Fit Critic model based on samples
4. Calculate actor gradients
5. Update actor



RL Class of Methods

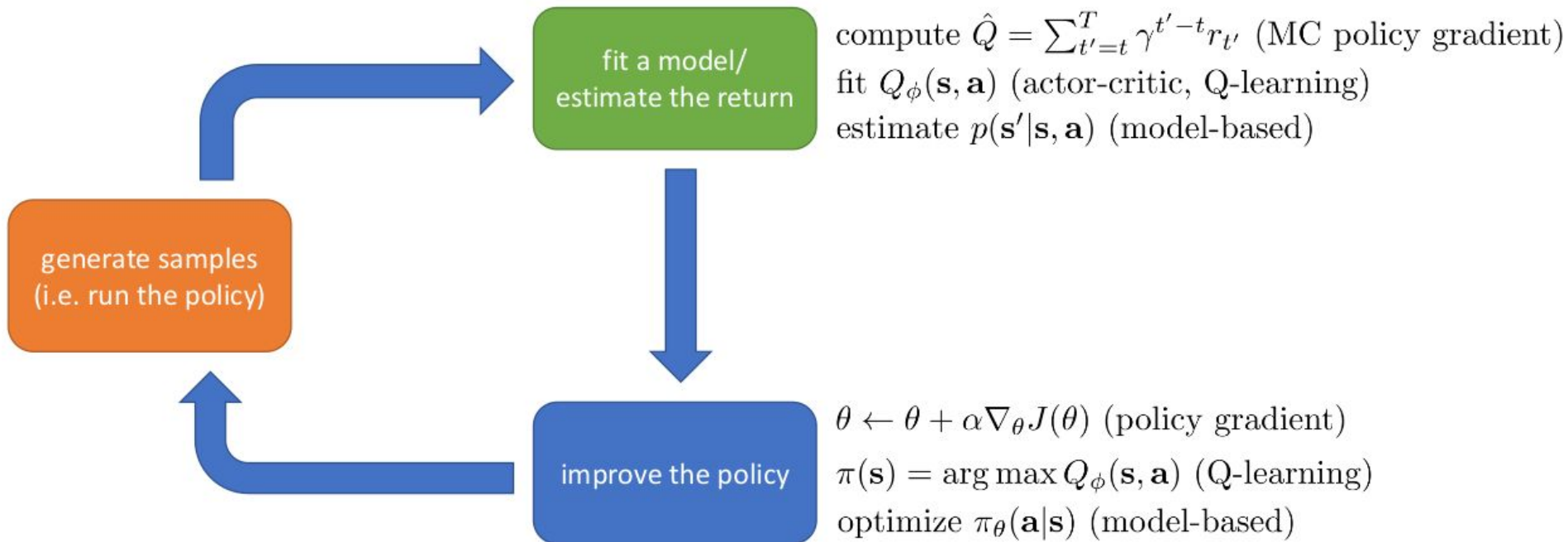
- Policy Gradient Methods (REINFORCE)



RL Class of Methods

- Guided Policy Search

RL Class of Methods



Reinforcement Learning

Pros

- Fully autonomous
- No explicit coding needed
- No demonstration needed

Cons

- Reward definition: Where does the R come from?
- Curse of Dimensionality: Exploration cost increases exponentially with dimension
- Generalization issues
- Convergence

Reinforcement Learning

Pros

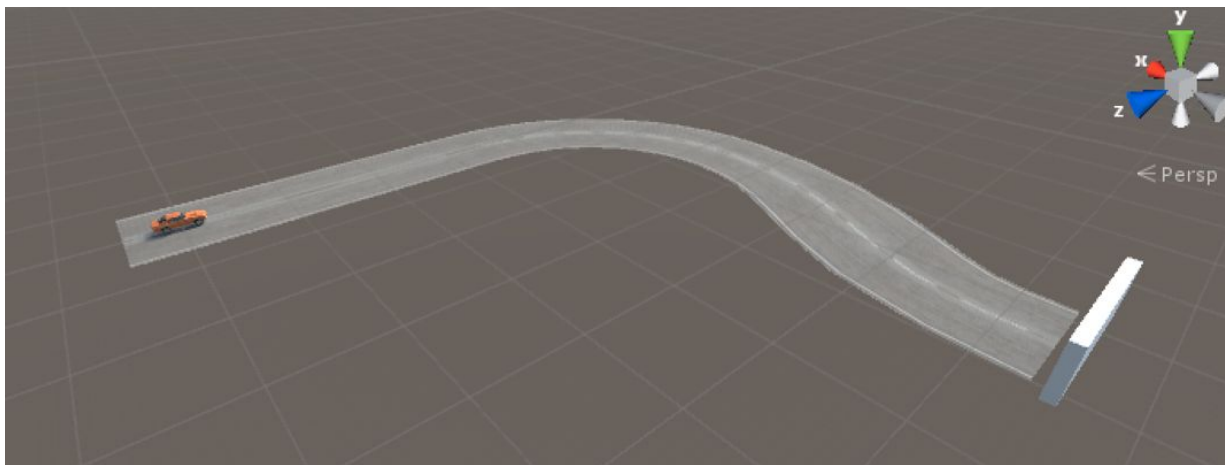
- Fully autonomous
- No explicit coding needed
- No demonstration needed

Cons

- **Reward definition: Where does R come from?**
- **Curse of Dimensionality: Exploration cost increases exponentially with dimension**
- Generalization issues
- Convergence

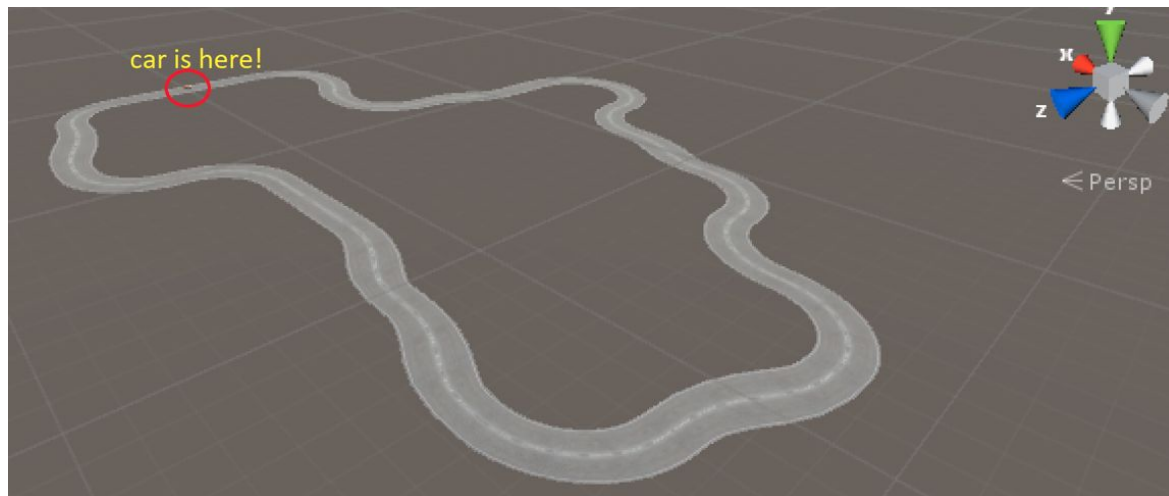
Reinforcement Learning in Robotics

- **Quiz 2: Mobile Robot Example**



Reinforcement Learning in Robotics

- Quiz 2: Mobile Robot Example



Reinforcement Learning in Robotics

- **Quiz 2: Mobile Robot Example**

- **State space:**
- **Action space:**
- **Reward function:**



Reinforcement Learning in Robotics

- **Quiz 2: Mobile Robot Example**

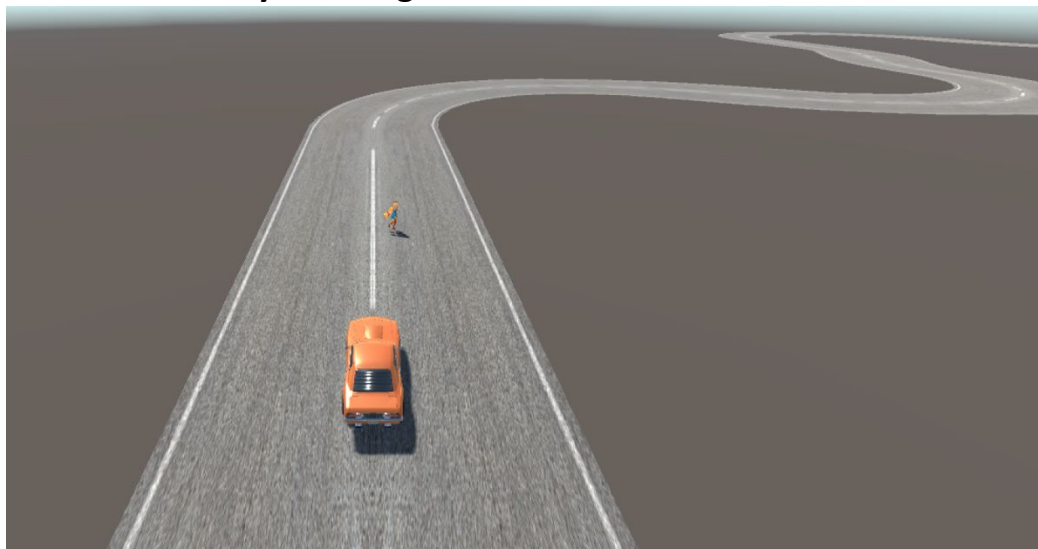
- **State space:** image + car velocity
- **Action space:** gas, wheel
- **Reward function:** forward velocity/5, and -500 if fall off



Reinforcement Learning in Robotics

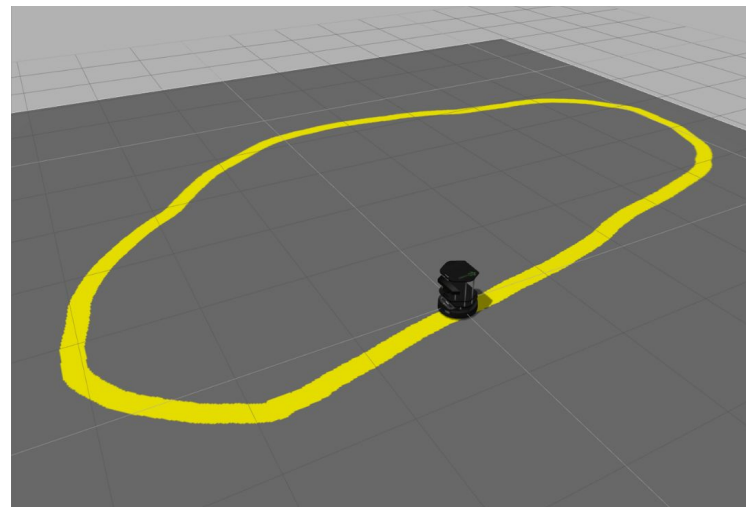
- **Quiz 2: Mobile Robot Example**

- **State space:** image + car velocity
- **Action space:** gas, wheel
- **Reward function:** forward velocity/5, and -500 if fall off
- **Pedestrians!!!** What will you change?



Reinforcement Learning in Robotics

- **Quiz 2: Mobile Robot Example**
 - Similar to Homework 5?



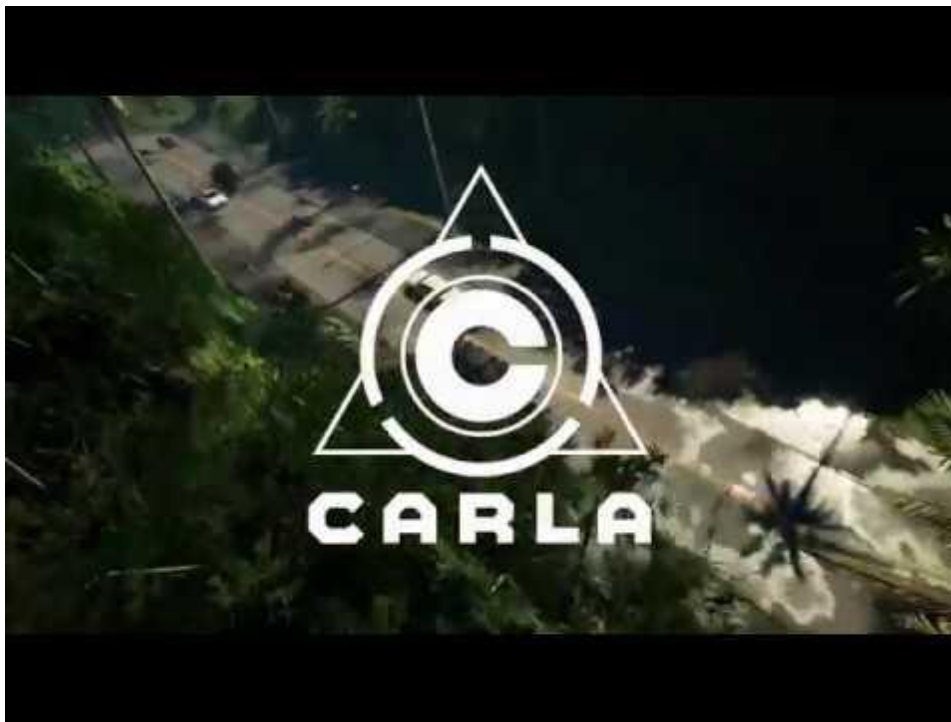
Reinforcement Learning in Robotics

Mobile Robot -> Autonomous Vehicle

- CARLA: An Open Urban Driving Simulator (Dosovitskiy et al 2017) ([GitHub](#))

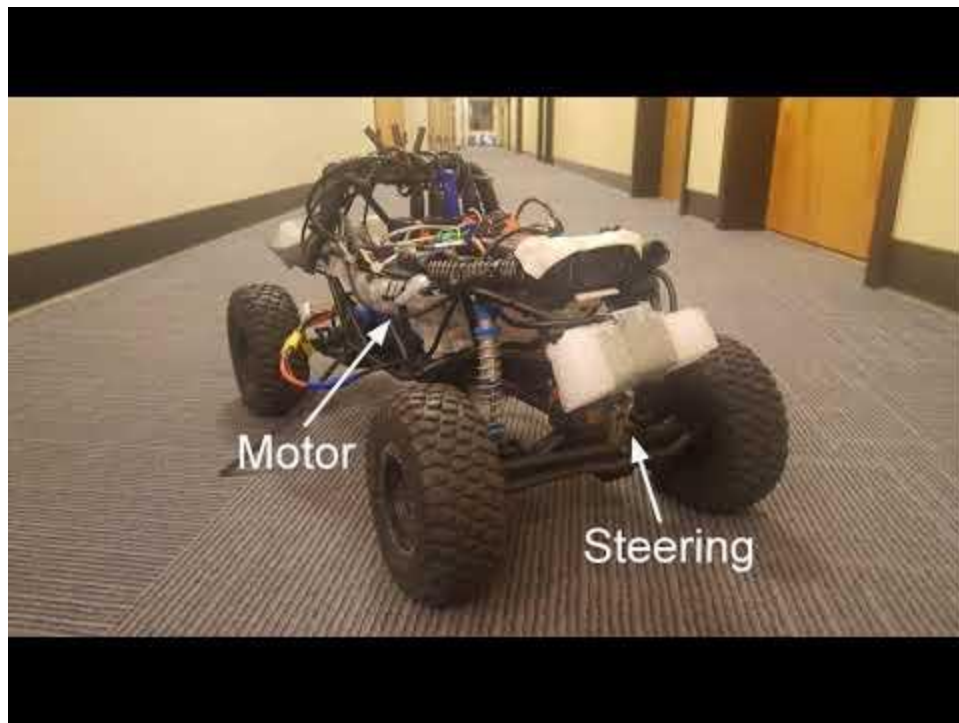
Reinforcement Learning in Robotics

CARLA: An Open Urban Driving Simulator (Dosovitskiy et al 2017) ([GitHub](#))



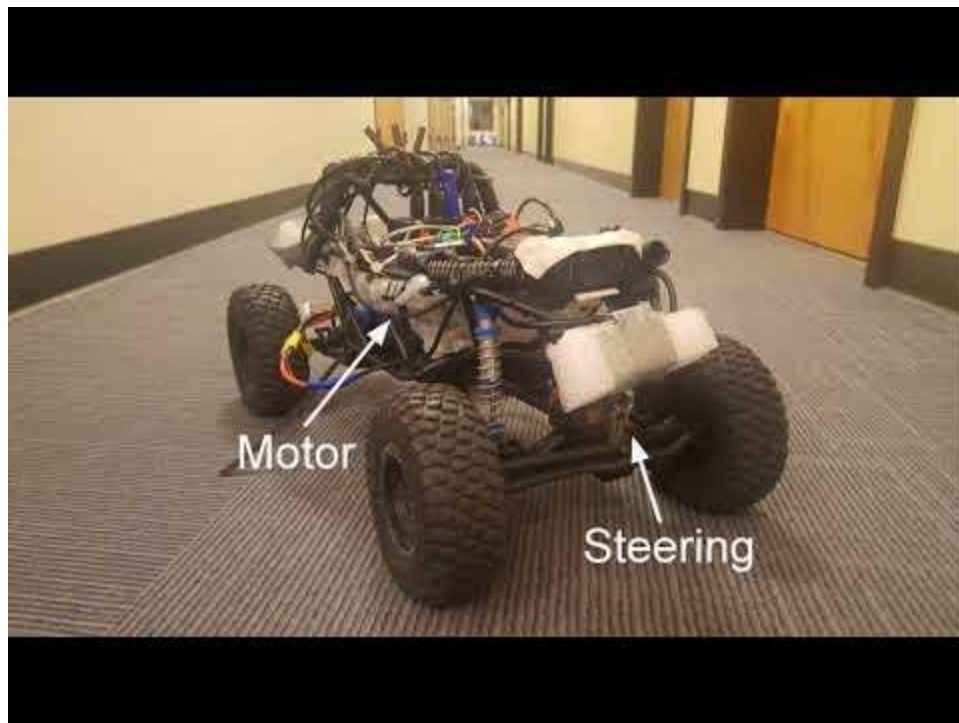
Reinforcement Learning in Robotics

Indoor Navigation



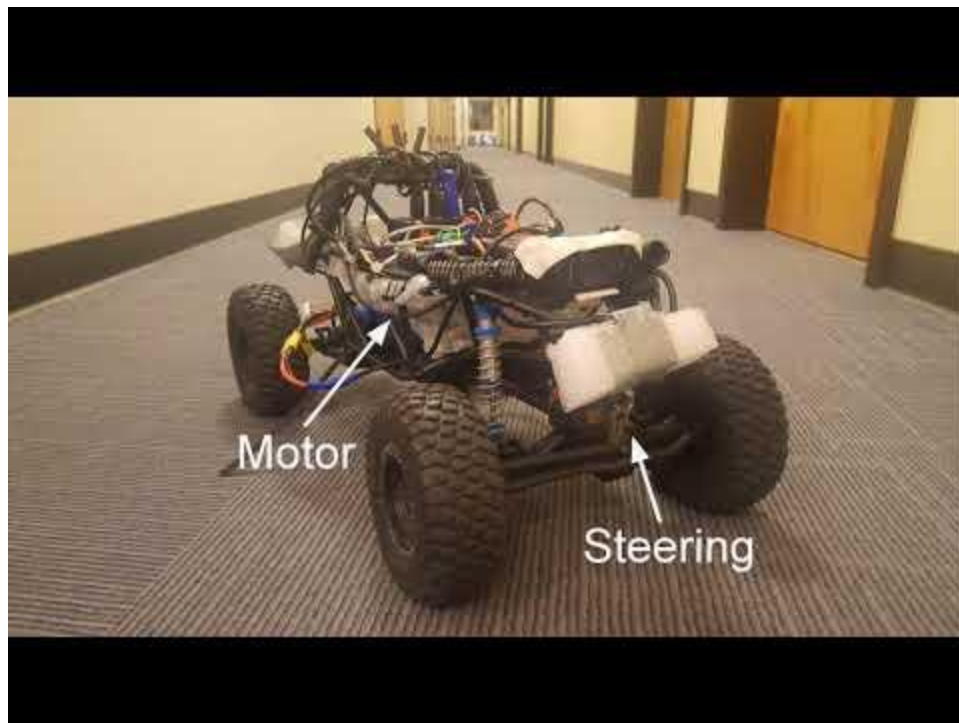
Reinforcement Learning in Robotics

Indoor Navigation



Reinforcement Learning in Robotics

Indoor Navigation



References

- “Robot Learning From Human Teachers”, Sonia Chernova and Andrea L. Thomaz (2014)
- “Deterministic Policy Gradient Algorithms” (David Silver et al. 2014)
- “Continuous control with deep reinforcement learning” (Lillicrap et al Deepmind 2016)
- "CARLA: An open urban driving simulator." Dosovitskiy, Alexey, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. (2017).
- “Self-supervised Deep Reinforcement Learning with Generalized Computation Graphs for Robot Navigation” Kahn, Gregory, Adam Villaflor, Bosen Ding, Pieter Abbeel, and Sergey Levine. (ICRA 2018)
-

Reinforcement Learning for Robotics

Iretiayo Akinola

<http://www.cs.columbia.edu/~iakinola/>

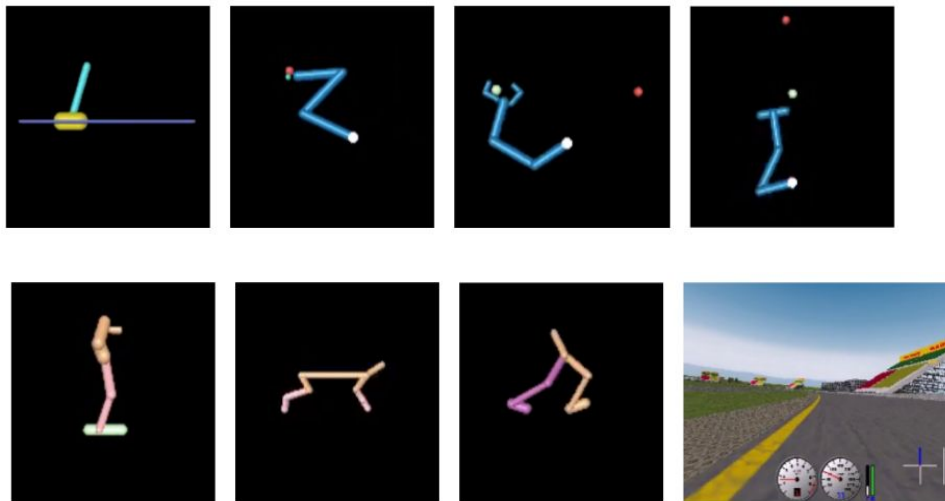


COLUMBIA
UNIVERSITY

Thanks

Reinforcement Learning in Robotics

- Deterministic Policy Gradient Algorithms (David Silver et al. 2014)
 - DPG Algorithm
 - Experiments: continuous bandit, **pendulum**, mountain car, 2D puddle world and **Octopus Arm**
- Continuous Control with Deep Reinforcement Learning (Lillicrap et al. 2016)
 - DDPG



DPG (Silver et al. 2014)

- Deterministic Policy Gradient Algorithms (David Silver et al. 2014)

Total discounted reward $r_t^\gamma = \sum_{k=t}^{\infty} \gamma^{k-t} r(s_k, a_k)$ where $0 < \gamma < 1$

Value functions: $V^\pi(s) = \mathbb{E}[r_1^\gamma | S_1 = s; \pi]$

$$Q^\pi(s, a) = \mathbb{E}[r_1^\gamma | S_1 = s, A_1 = a; \pi]$$

Performance Objective: $J(\pi) = \mathbb{E}[r_1^\gamma | \pi]$

$$\begin{aligned} J(\pi_\theta) &= \int_{\mathcal{S}} \rho^\pi(s) \int_{\mathcal{A}} \pi_\theta(s, a) r(s, a) da ds = \int_{\mathcal{S}} \int_{\mathcal{A}} \rho^\beta(s) \pi_\theta(a|s) Q^\pi(s, a) da ds \\ &= \mathbb{E}_{s \sim \rho^\pi, a \sim \pi_\theta} [r(s, a)] \end{aligned}$$

Policy Gradient: $\nabla_\theta J(\pi_\theta) \quad \theta := \theta + \alpha \nabla_\theta J(\pi_\theta)$

DPG (Silver et al. 2014)

- Deterministic Policy Gradient Algorithms (David Silver et al. 2014)

Performance Objective: $J(\pi_\theta) = \int_{\mathcal{S}} \int_{\mathcal{A}} \rho^\beta(s) \pi_\theta(a|s) Q^\pi(s, a) da ds$

Policy Gradient: $\nabla_\theta J(\pi_\theta)$

Deterministic Policy Gradient:

$$\begin{aligned} \nabla_\theta J(\mu_\theta) &= \int_{\mathcal{S}} \rho^\mu(s) \nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a)|_{a=\mu_\theta(s)} ds \\ &= \mathbb{E}_{s \sim \rho^\mu} \left[\nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a)|_{a=\mu_\theta(s)} \right] \end{aligned}$$

Actor

Critic

$$\theta := \theta + \alpha \nabla_\theta J(\mu_\theta)$$

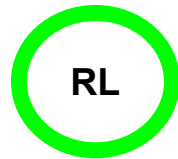
DPG (Silver et al. 2014)

- Deterministic Policy Gradient Algorithms (David Silver et al. 2014)

Actor: $\mu_{\theta}(s)$ Critic: $Q^w(s_t, a_t)$

DPG Algorithm:

$$\begin{aligned}\delta_t &= r_t + \gamma Q^w(s_{t+1}, \mu_{\theta}(s_{t+1})) - Q^w(s_t, a_t) \\ w_{t+1} &= w_t + \alpha_w \delta_t \nabla_w Q^w(s_t, a_t) \\ \theta_{t+1} &= \theta_t + \alpha_{\theta} \nabla_{\theta} \mu_{\theta}(s_t) \nabla_a Q^w(s_t, a_t)|_{a=\mu_{\theta}(s)}\end{aligned}$$

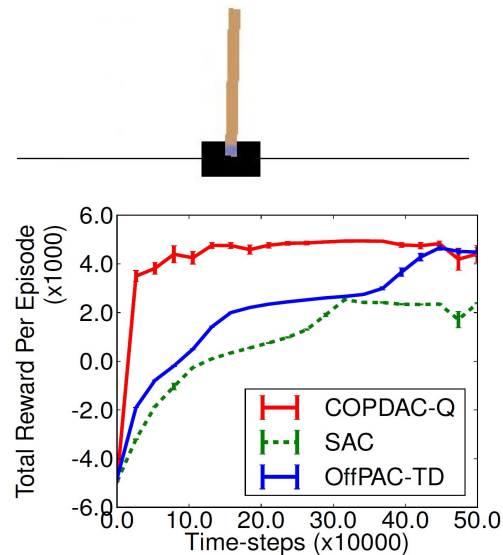


DPG (David Silver et al. 2014)

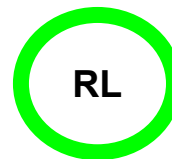
- Experiments:
 - continuous bandit, **pendulum**, mountain car, 2D puddle world, Octopus Arm

Pendulum:

- State: joint position/velocities
- Action: move left or right
- Reward: +1 for every step while rod is upright
- Critic: $V(s) = v^\top \phi(s)$
- Actor:
 - Target policy: $\mu_\theta(s) = \theta^\top \phi(s)$
 - Behavior policy: $\beta(\cdot|s) \sim \mathcal{N}(\theta^\top \phi(s), \sigma_\beta^2)$



DPG (David Silver et al. 2014)

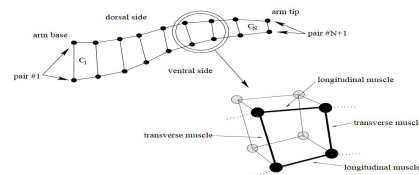
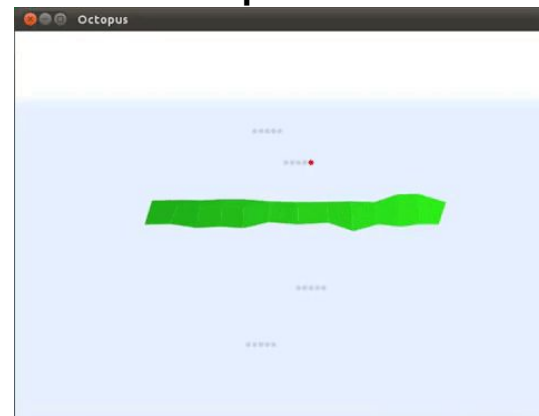


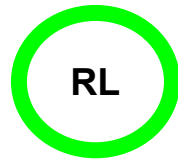
- Experiments:
 - continuous bandit, pendulum, mountain car, 2D puddle world, Octopus Arm

Octopus Arm:

- State: 50 continuous joint state variables
- Action: 20 variables to control muscles
- Reward: change in distance between arm and target
- Critic: NN Multilayer Perceptron (MLP)
- Actor: MLP (8 hidden units)

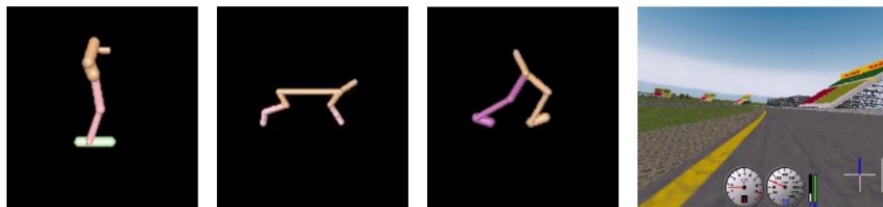
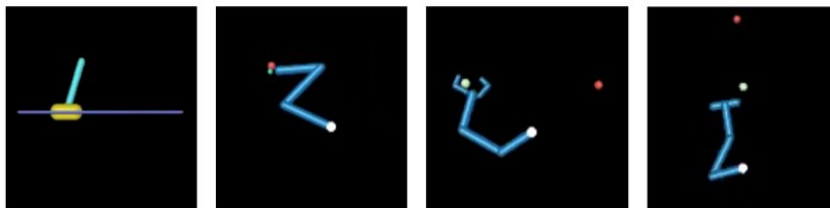
Octopus Arm



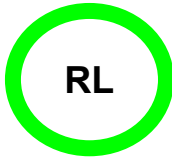


Reinforcement Learning in Robotics

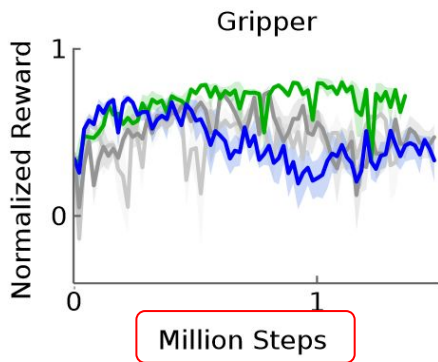
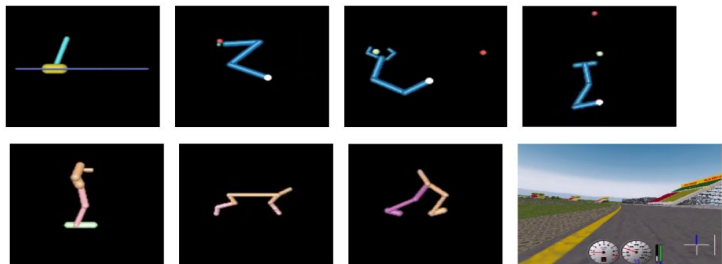
- Deterministic Policy Gradient Algorithms (Silver et al. 2014) **DPG**
 - Critic: linear function, Actor: Gaussian policy
 - Critic: NN(MLP), Actor: NN(MLP)
- Continuous control with deep reinforcement learning (Lillicrap et al. 2016) **DDPG**
 - Critic: **DNN**, Actor: **DNN**



DDPG (Lillicrap 2016)



- Batch normalization
- Target Networks



- original DDPG with batch normalization (light grey),
- with target network (dark grey),
- with target networks and batch norm (green),
- with target networks from pixel-only inputs (blue)

Reinforcement Learning in Robotics

- **Reward definition: Where does R come from?**

