

#DLUPC

Applications 2

# Reinforcement Learning



Xavier Giro-i-Nieto

[xavier.giro@upc.edu](mailto:xavier.giro@upc.edu)

Associate Professor

Universitat Politècnica de Catalunya  
Barcelona Supercomputing Center



# Acknowledgements



Víctor Campos

victor.campos@bsc.es

PhD Candidate

Barcelona Supercomputing Center



Míriam Bellver

miriam.bellver@bsc.edu

PhD Candidate

Barcelona Supercomputing Center



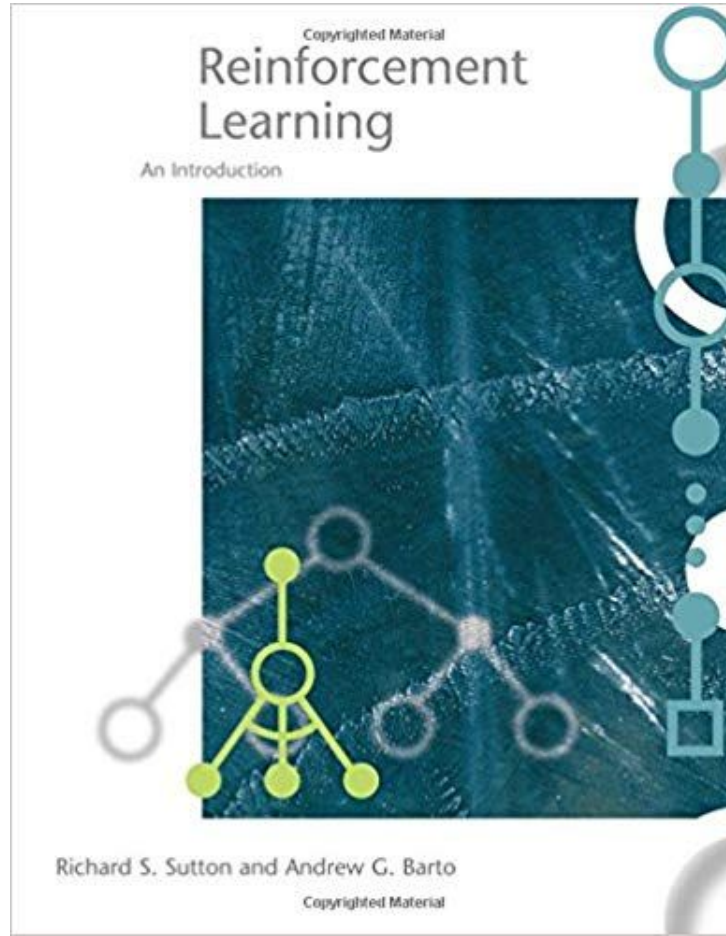
# Acknowledgements



# OpenAI

Lilian Weng, [“A \(Long\) Peek into Reinforcement Learning”](#) (2018)

# Bibliography





# Temporal-Difference Learning

*Rich Sutton*

Reinforcement Learning & Artificial Intelligence Laboratory  
Alberta Machine Intelligence Institute  
Dept. of Computing Science, University of Alberta  
Canada



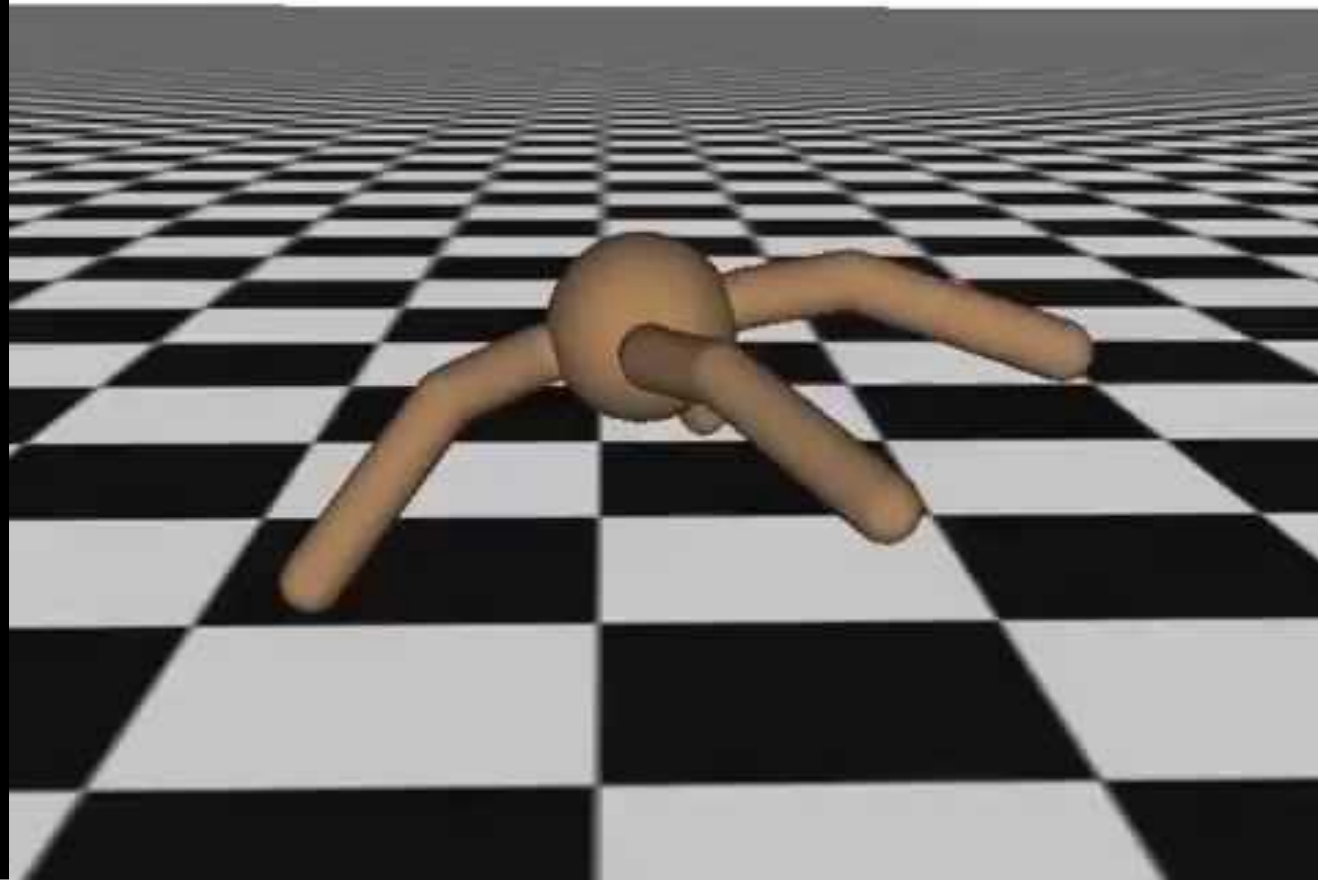
# Outline

1. **Control tasks with reinforcement learning**
2. Markov Decision Processes
3. Learning action-value functions  $Q$
4. Learning policies



Mnih, Volodymyr, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. "Playing atari with deep reinforcement learning." NIPS Deep Learning Workshop (2013).

Iteration 20







Wayve, "The first example of deep reinforcement learning on-board an autonomous car" (2018)

# Types of machine learning

Yann Lecun's Black Forest cake

## ■ "Pure" Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

## ■ Supervised Learning (icing)

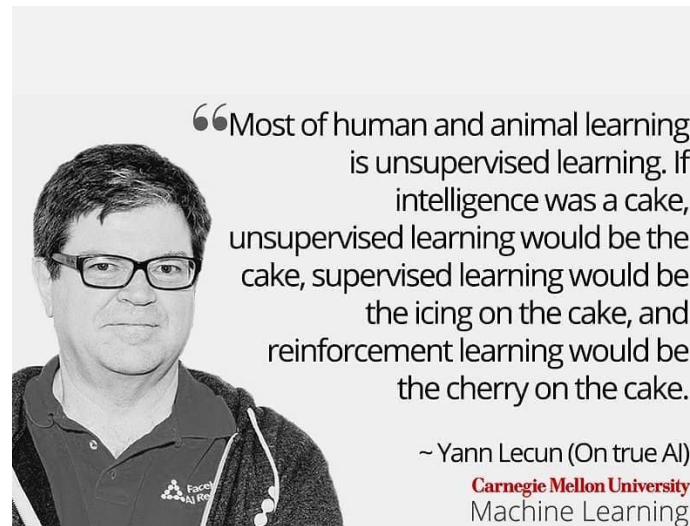
- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

## ■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)





Yoshua Bengio, Geoffrey Hinton and Yann LeCun, the fathers of [#DeepLearning](#), receive the 2018 [#ACMTuringAward](#) for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing today.

[bit.ly/2HVJtdV](https://bit.ly/2HVJtdV)

Tradueix el tuit





## Le Cake



# Our vision

## Le Cake

♥ A Pieter Abbeel i 7 més els agrada



**Yann LeCun**  
@ylecun



Pierre Pieter Abbeel gave a seminar at NYU today.

He showed up at the Turing reception and was given the RL piece of "Le Cake" (and lots of cherries).

[Tradueix el tuit](#)

1:10 · 13/4/19 · [Facebook](#)



# A broader picture of types of learning...

	...with a teacher	...without a teacher
Active agent...	Reinforcement learning (with extrinsic reward)	Intrinsic motivation / Exploration.
Passive agent...	Supervised learning	Unsupervised learning



Slide inspired by Alex Graves (Deepmind) at  
[“Unsupervised Learning Tutorial”](#) @ NeurIPS 2018.

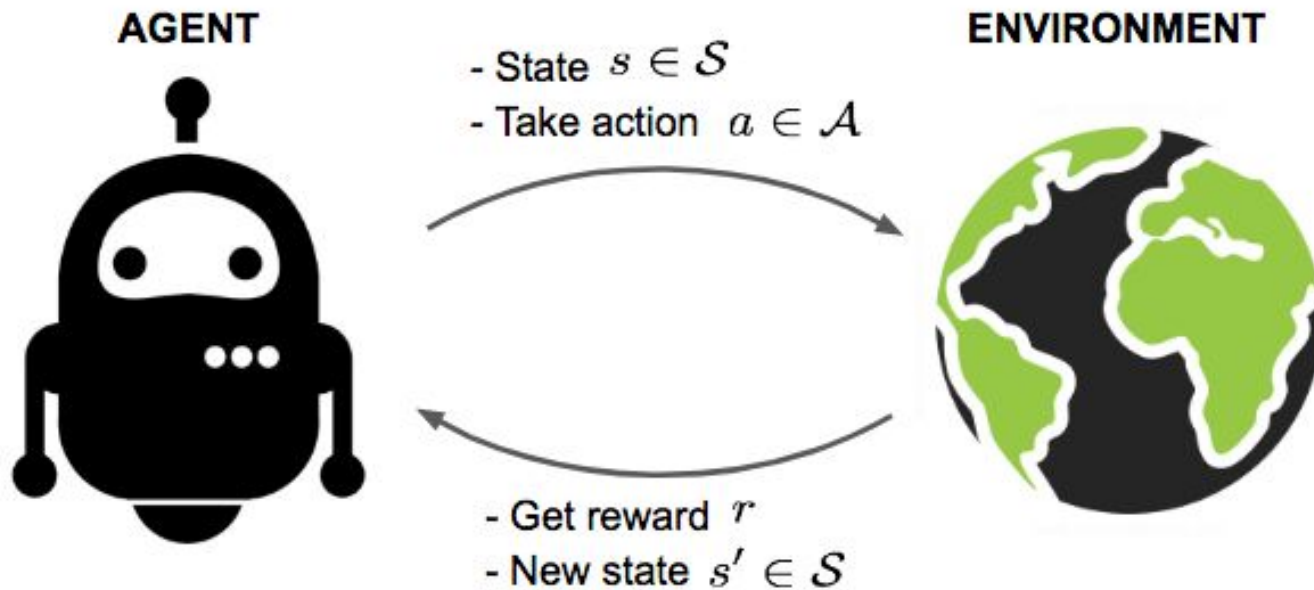
# A broader picture of types of learning...

	...with a teacher	...without a teacher
Active agent...	Reinforcement learning (with extrinsic reward)	Intrinsic motivation / Exploration.
Passive agent...	Supervised learning	Unsupervised learning



Slide inspired by Alex Graves (Deepmind) at  
[“Unsupervised Learning Tutorial”](#) @ NeurIPS 2018.

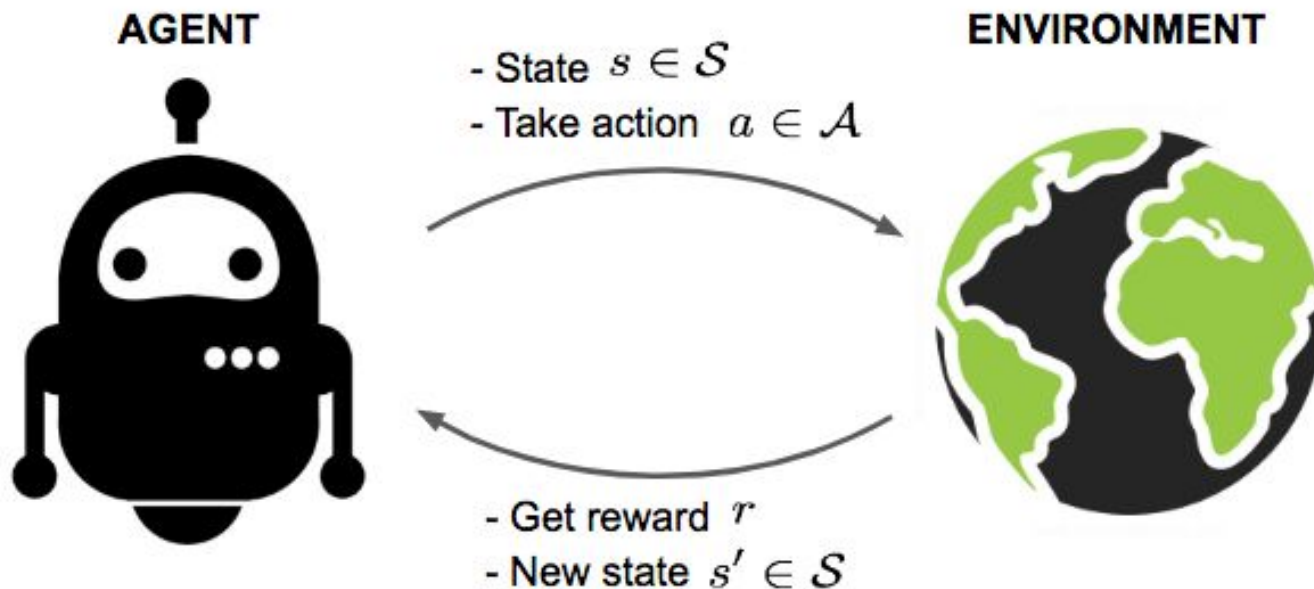
# Reinforcement Learning (RL)





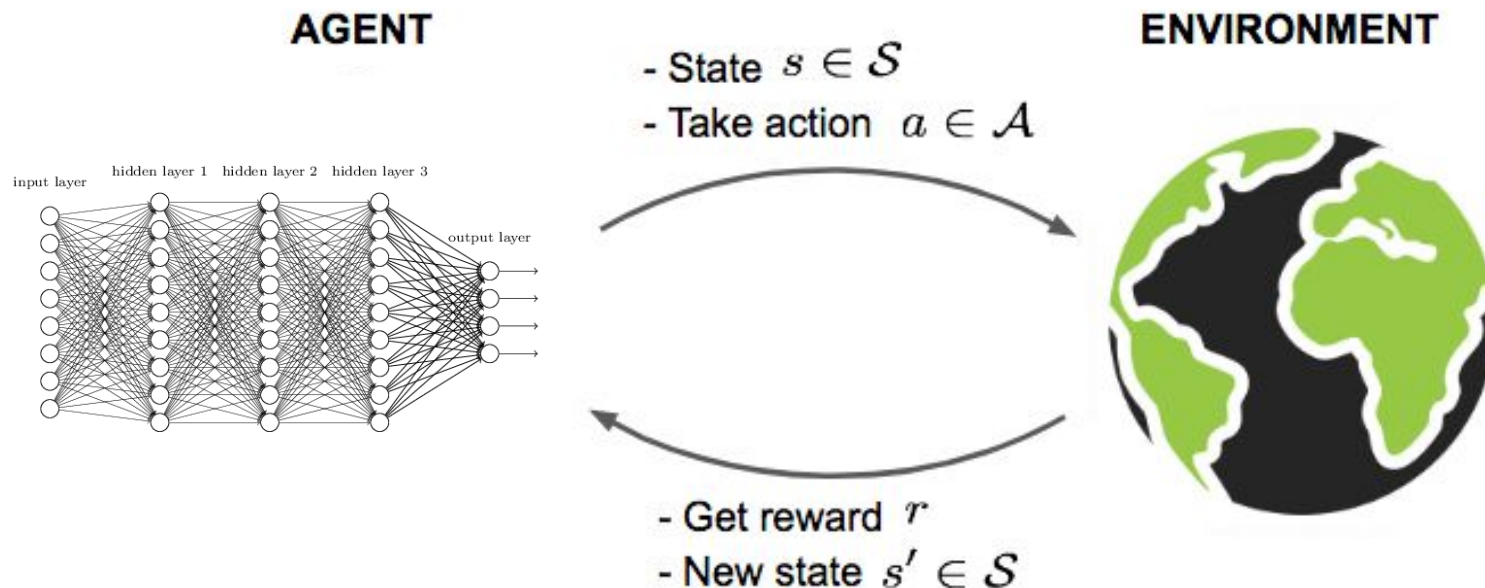
# Reinforcement Learning (RL)

The Policy  $\pi$  is a function  $S \rightarrow A$  that specifies which action the agent will take given a state.



# Deep Reinforcement Learning (RL)

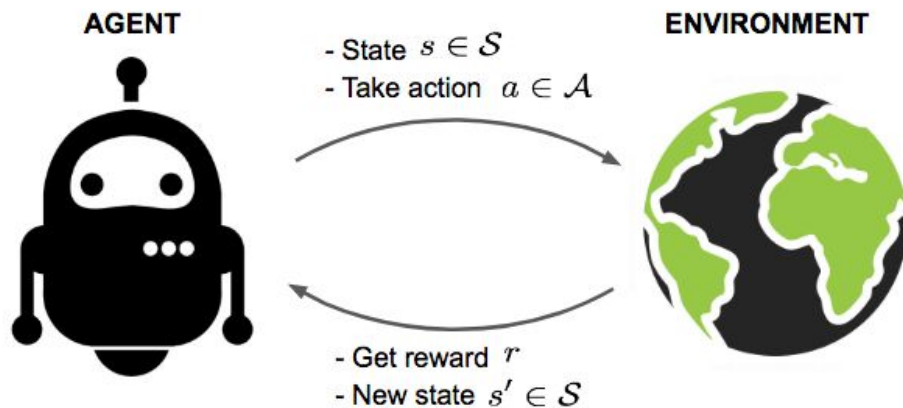
The Policy  $\pi$  can be governed through a deep neural network.



# Training data by computation

In RL, training data is obtained by computing interaction sequences of:

## State , Action , Reward



A complete episode consists of  $T$  interactions:

$$S_1, A_1, R_2, S_2, A_2, \dots, S_T$$

One experience corresponds to a single tuple within an episode:

$$(s, a, r, s')$$

# Outline

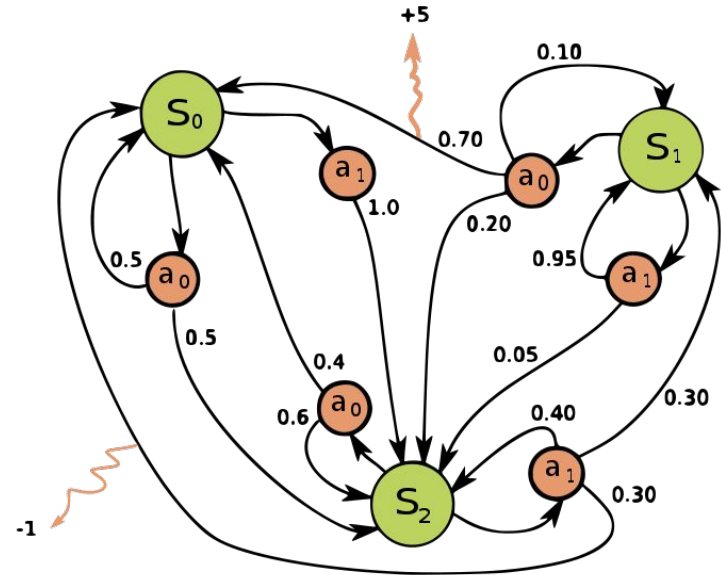
1. Control tasks with reinforcement learning
- 2. Markov Decision Processes**
3. Learning action-value functions  $Q$
4. Learning policies

# Markov Decision Processes (MDP)

Markov Decision Processes provide a formalism for reinforcement learning problems.

## Markov property:

Current state completely characterises the state of the world.



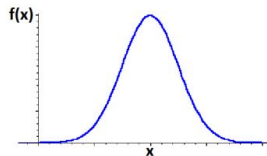
# Markov Decision Process (MDP)



S



A



R



P



Y



Environment  
samples initial  
state  $s_0 \sim p(s_0)$

state (s)



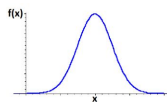
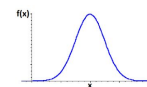
Agent  
selects  
action a

action (a)



Environment samples  
reward  $r \sim R(\cdot | s, a)$

reward  
(r)



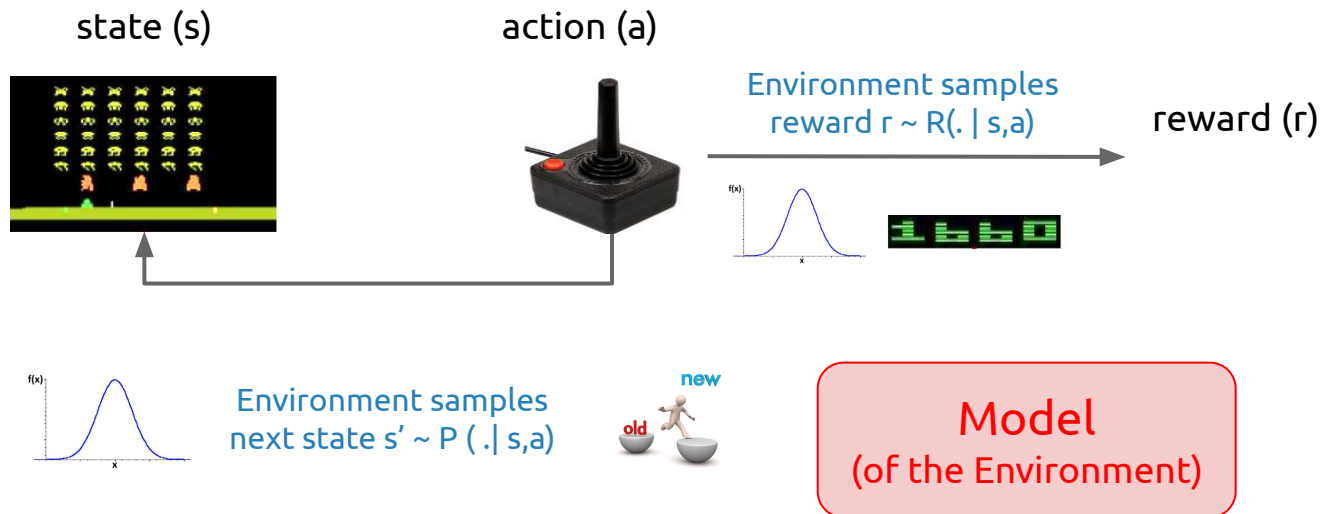
Environment samples  
next state  $s' \sim P(\cdot | s, a)$



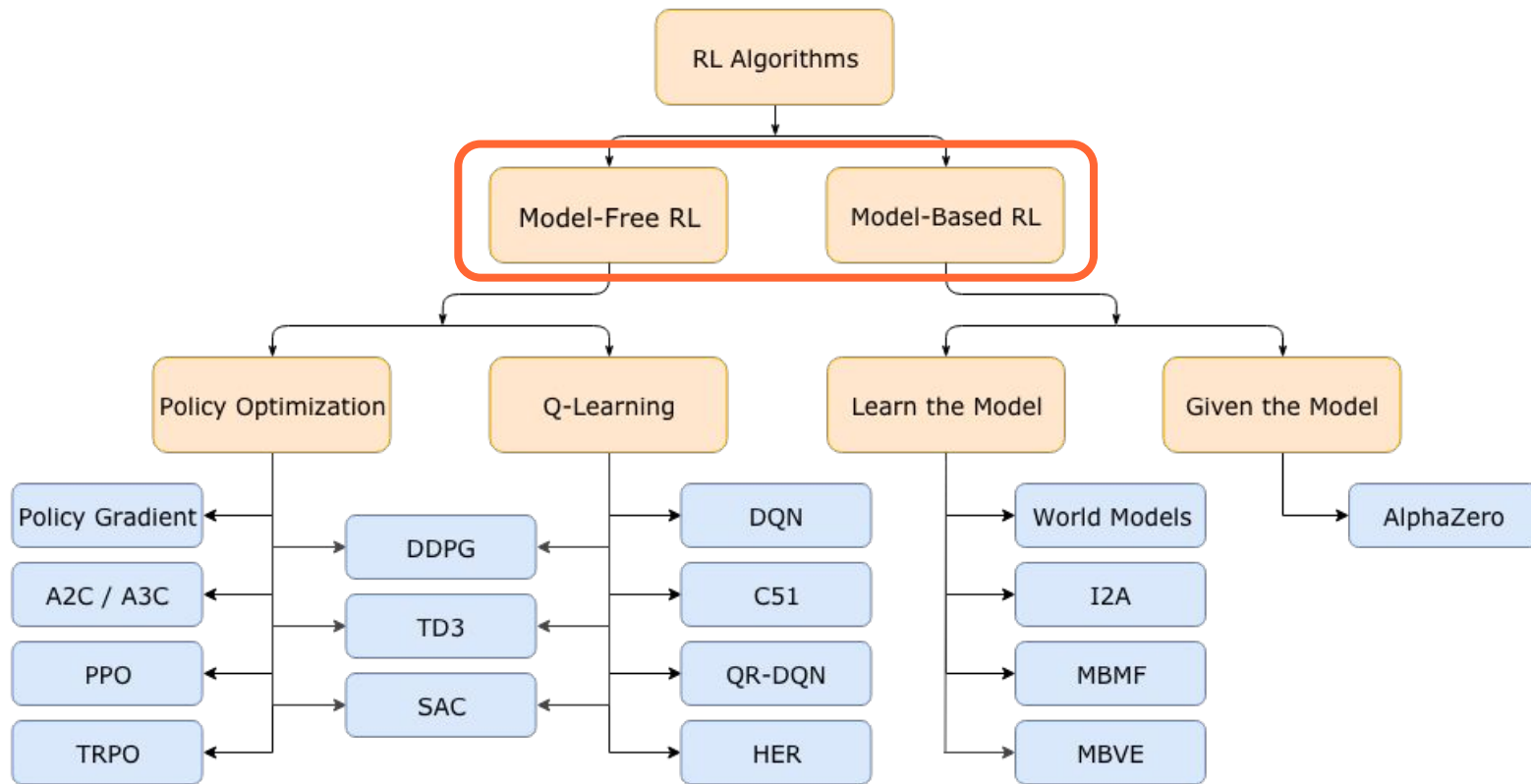
# MDP: Model (of the Environment)

The Model (of the Environment) is defined by:

- $P(s'|s,a)$ : State-transition function
- $R(s,a)$ : Reward function.



# Model (of the Environment)





# Model (of the Environment)

- **Model-based RL:**  $R(\cdot | s, a)$  and  $P(\cdot | s, a)$  are known, so an optimal solution can be found with [Dynamic Programming](#).

More details at [Sergey Levine's slides](#) @ Berkeley.

- **Model-free RL:** there is no prior knowledge of the world. The agent needs to learn all dynamics from scratch, resulting in poor data efficiency which limits its application to real-world agents.

Requires [sim2real](#) transfer for real deployment.



Wayve, "Sim2Real: Learning to Drive from Simulation without Real World Labels" (2018)

# Model (of the Environment)

Example of sim2real: Dexterity (OpenAI 2018)

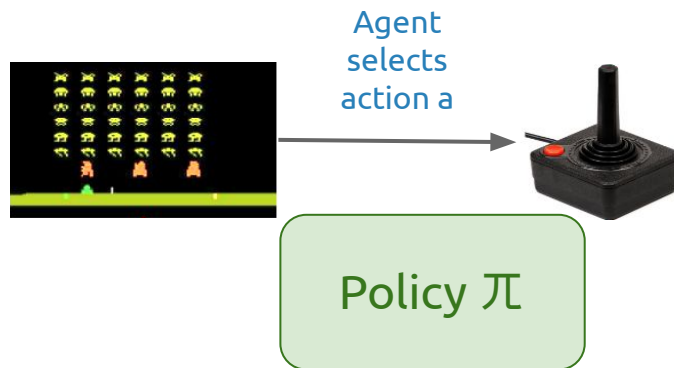


OpenAI

# On- vs Off-Policy Learning

The **Policy**  $\pi$  is a function  $S \rightarrow A$  that specifies which action to take given a state.

- **On-policy learning:** the agent is trained with a sequence of interactions of the *target policy*.
- **Off-policy learning:** the agent is trained with a sequence of interactions obtained from a *behaviour policy* different from the *target policy*.



# Policy: Deterministic vs Stochastic

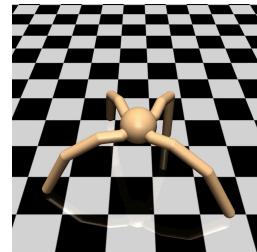
The **Policy**  $\pi$  is a function  $S \rightarrow A$  that specifies which action to take in each state.

- **Deterministic:**  $\pi(s)=a$
- **Stochastic:**  $\pi(a|s)=\mathbb{P}_{\pi}[A=a|S=s]$

# Policy: Discrete vs Continuous

The **Policy**  $\pi$  is a function  $S \rightarrow A$  that specifies which action to take in each state.

- **Discrete actions:** categorical distribution over actions.
  - In the deterministic case: take the argmax action.
  - In the stochastic case: sample from the categorical distribution.
- **Continuous actions:** Gaussian distribution over actions; i.e. the policy generally outputs the mean and std per dimension.
  - In the deterministic case: take the mean.
  - In the stochastic case: sample from a normal distribution.



# Return $G_t$

The future reward, also known as return  $G_t$ , is a total sum of discounted rewards from  $t$  and onwards in time...

Return  
 $G_t$

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

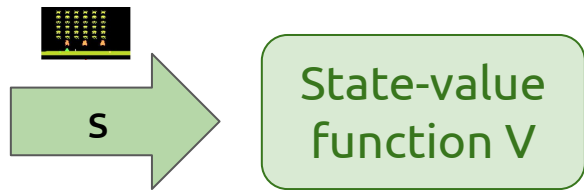


...where  $\gamma$  is the discount factor between  $[0,1]$ .

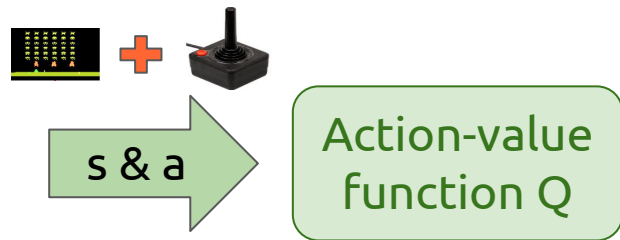


# Value Functions $V_{\pi}(s)$ & $Q_{\pi}(s,a)$

The value function of a policy  $\pi$  is the expected return  $G_t$  given the state (V) or the state & action (Q):



$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$$



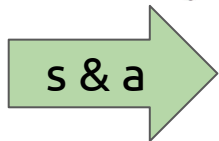
$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$$

“Expected return...”  
...by following policy  $\pi$  starting from state  $s$  & action  $a$ ”



# Advantage Function $A_{\pi}(s,a)$

The Advantage function (A-value) of a policy  $\pi$  is the defined as:



Advantage  
function A

$$A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(s)$$

$A_{\pi}(s,a)$  expresses how good/bad an action is with respect to the expectation of returns over all possible actions for a given state.

# Optimal Value Functions & Optimal Policy (\*)

The optimal value functions produce the maximum return...

Optimal  
value  
functions

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

The optimal policy achieves optimal value functions  $V_*(s)$  and  $Q_*(s, a)$

Optimal  
policy  $\pi_*$

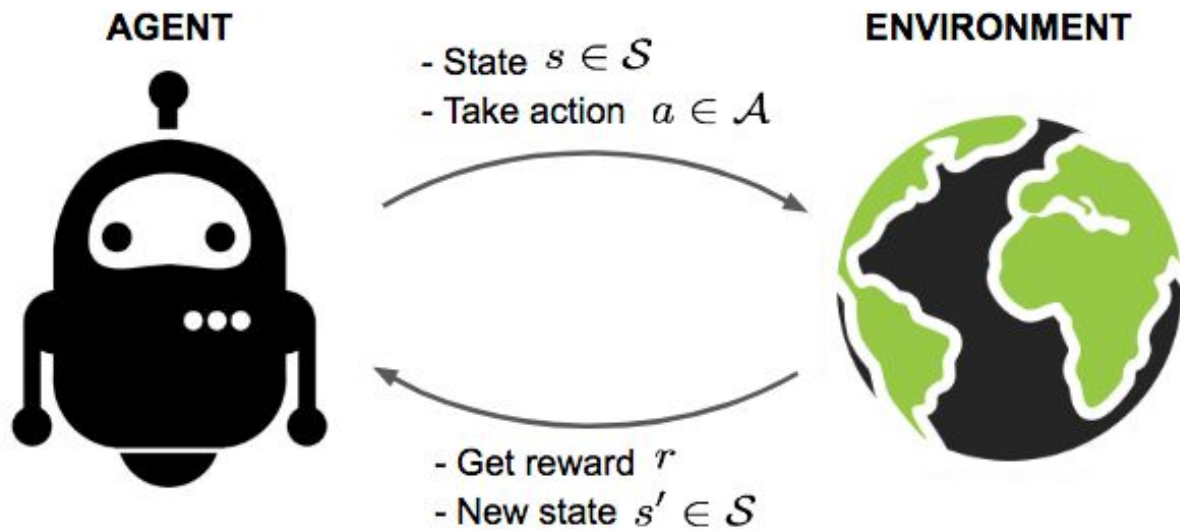
$$\pi_* = \arg \max_{\pi} V_{\pi}(s)$$

$$\pi_* = \arg \max_{\pi} Q_{\pi}(s, a)$$

# Outline

1. Control tasks with reinforcement learning
2. Markov Decision Processes
- 3. Learning action-value function  $Q$**
4. Learning policies

# Value vs Policy vs Model

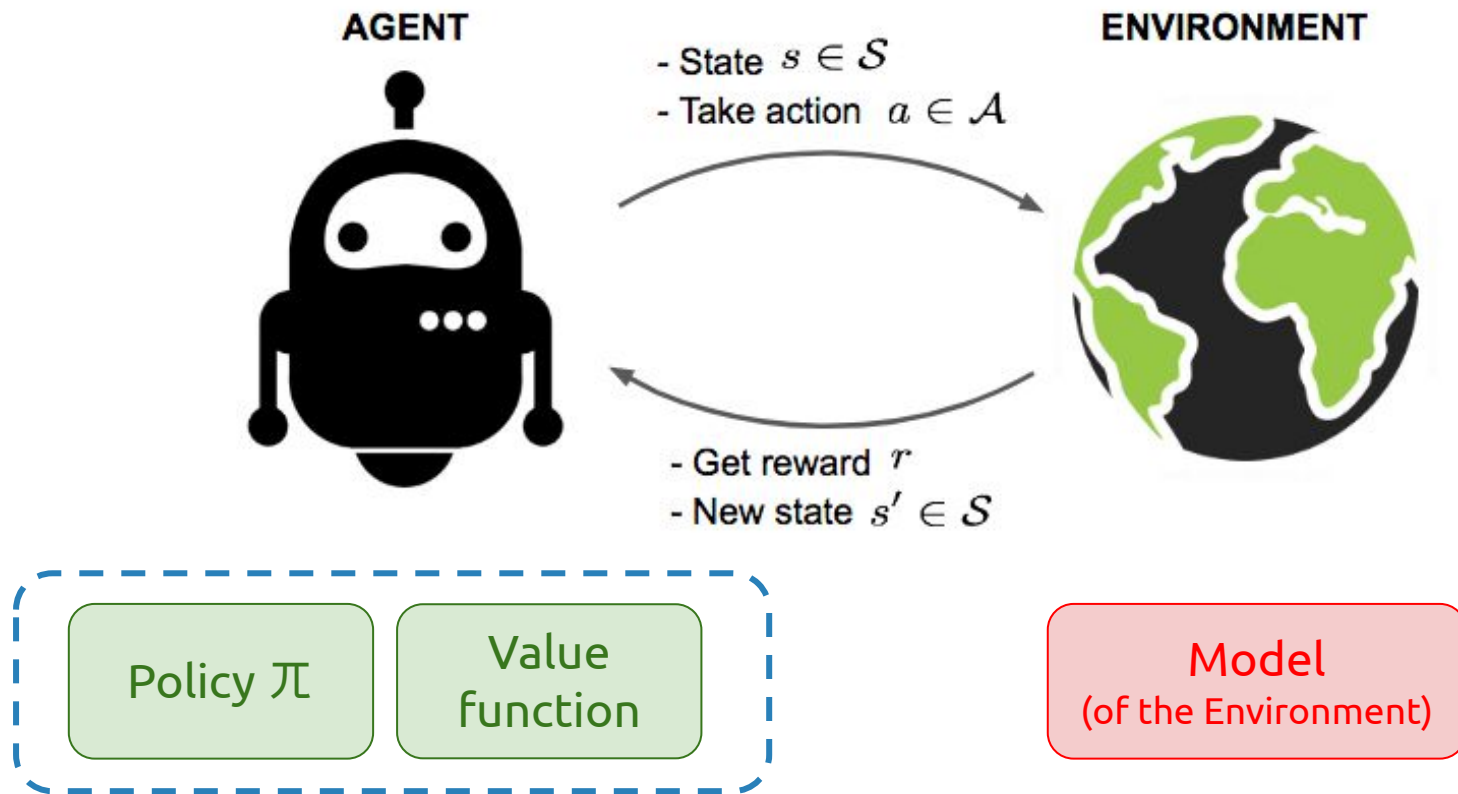


Policy  $\pi$

Value  
function

Model  
(of the Environment)

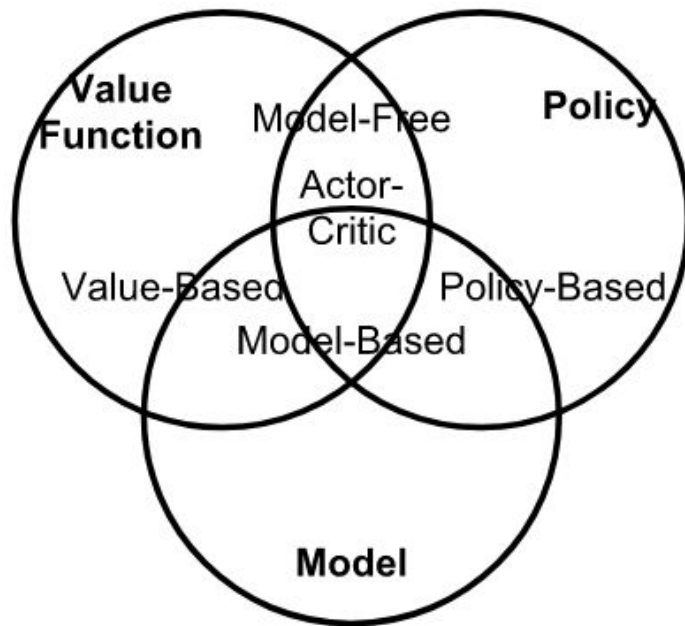
# Value vs Policy vs Model



Goals of Reinforcement Learning

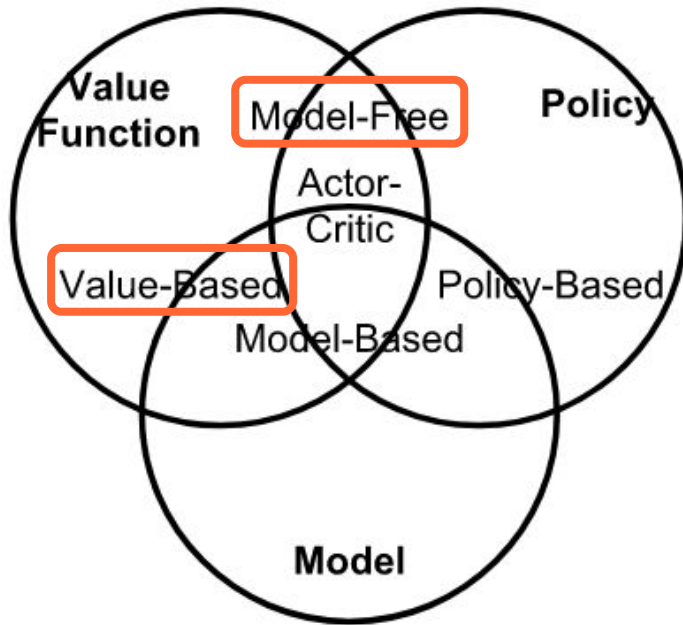
# Value vs Policy vs Model

Summary of approaches in RL based on whether we want to learn the value, policy, or the model of the environment.

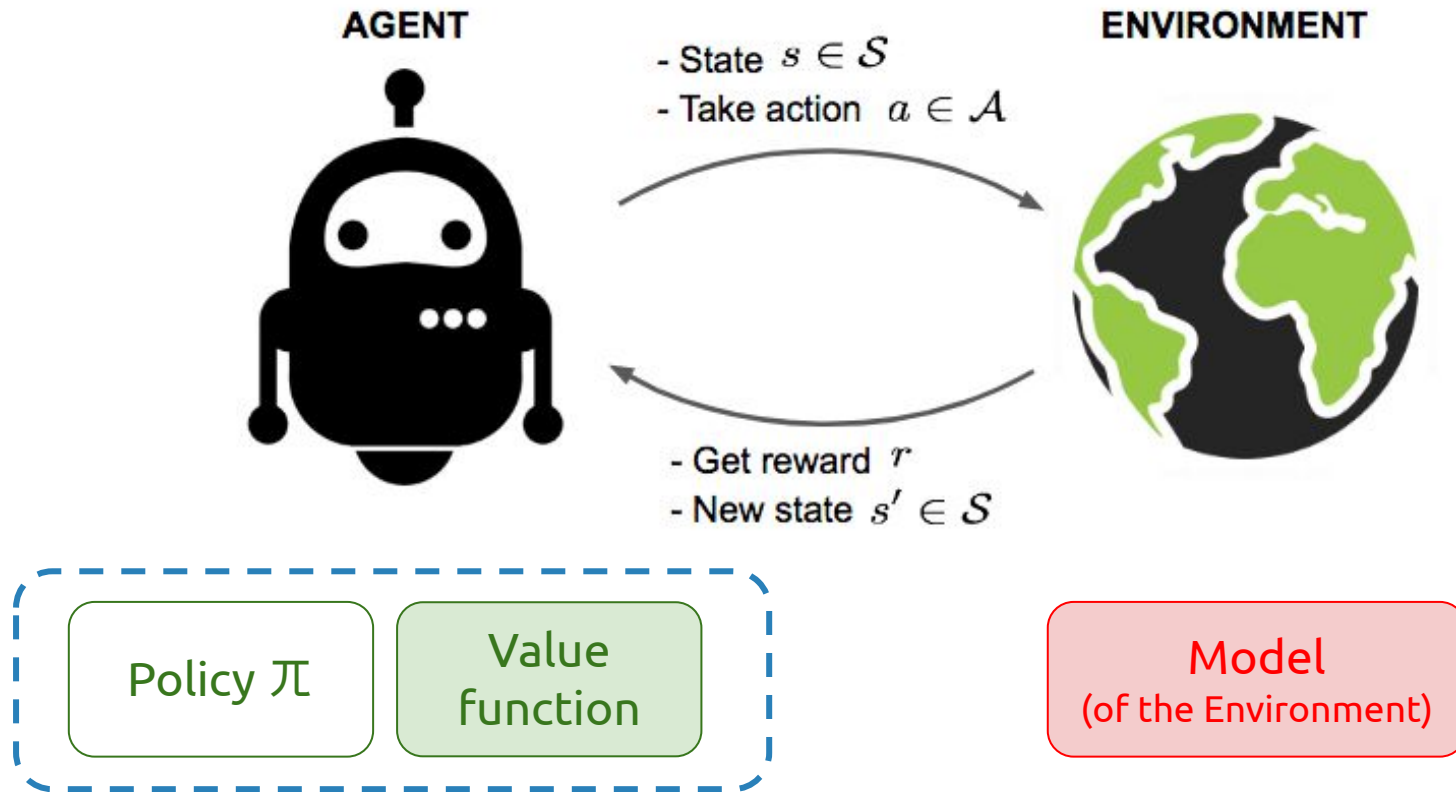


# Value vs Policy vs Model

Summary of approaches in RL based on whether we want to learn the value, policy, or the model of the environment.



# Value vs Policy vs Model

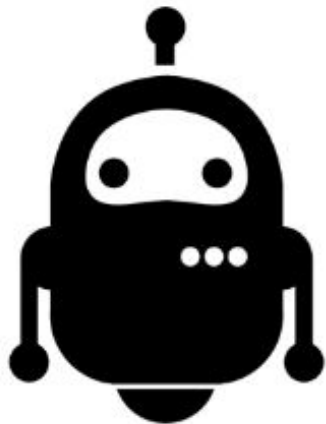


Goals of Reinforcement Learning



# Inferring a Policy from the Value

AGENT



Given a value function, a policy can be easily defined. For example:

Greedy policy

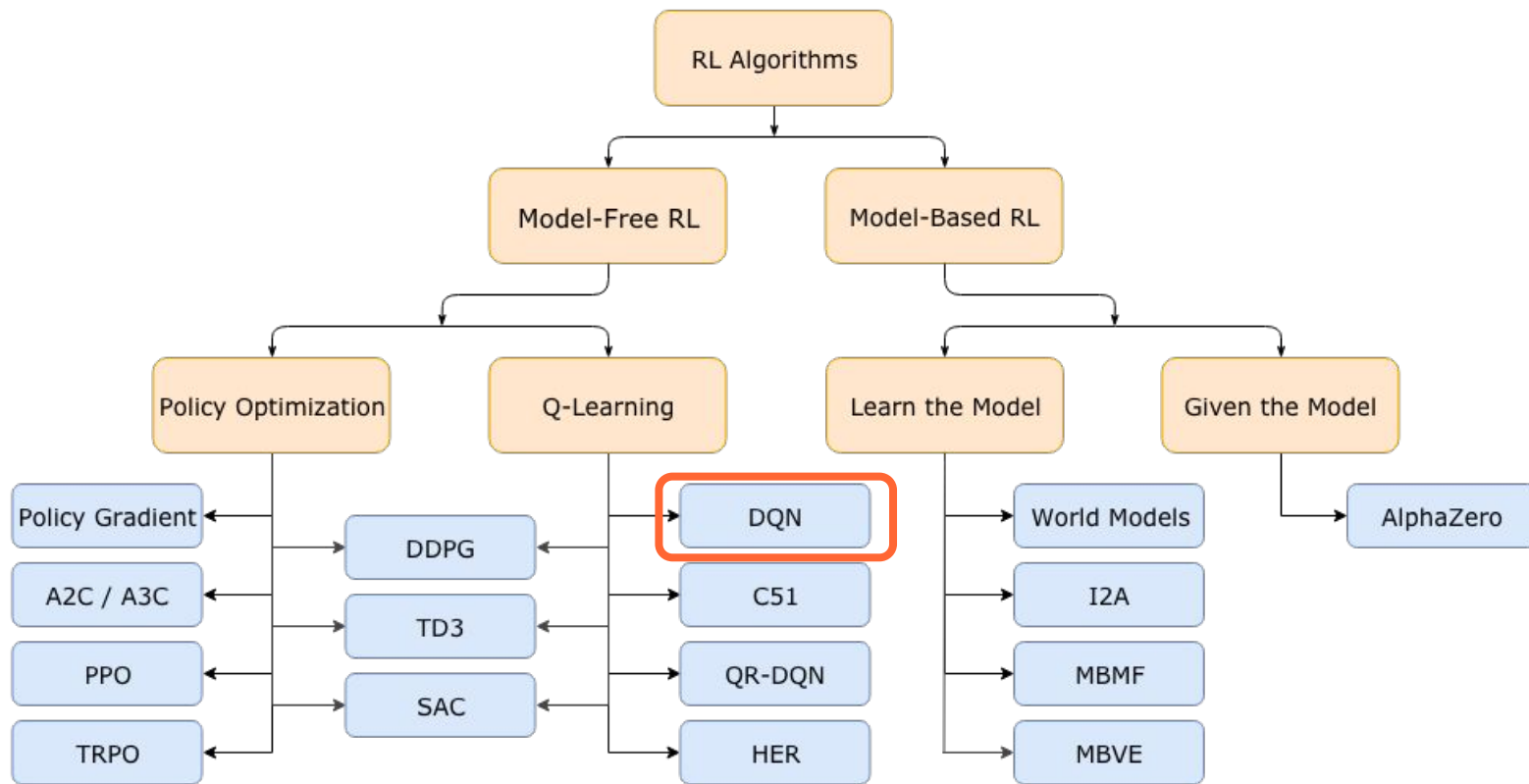
$$\pi(s) = \arg \max_{a \in A} Q(s, a)$$

Policy  $\pi$

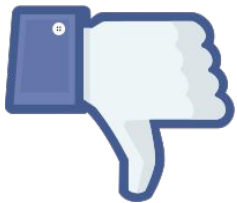
Value  
function



# Deep Q-learning (DQN)



# Deep Q-learning (DQN)



**Problem:** Estimating the optimal action-value function  $Q_*(s,a)$  might be feasible with exhaustive search if few state-action pairs are possible, but impossible for large spaces.



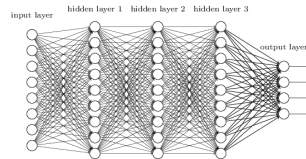
**Solution:** Learn an approximation of the optimal action-value function  $Q(s,a, \theta)$  with a NN:

$$Q(s,a, \theta)$$



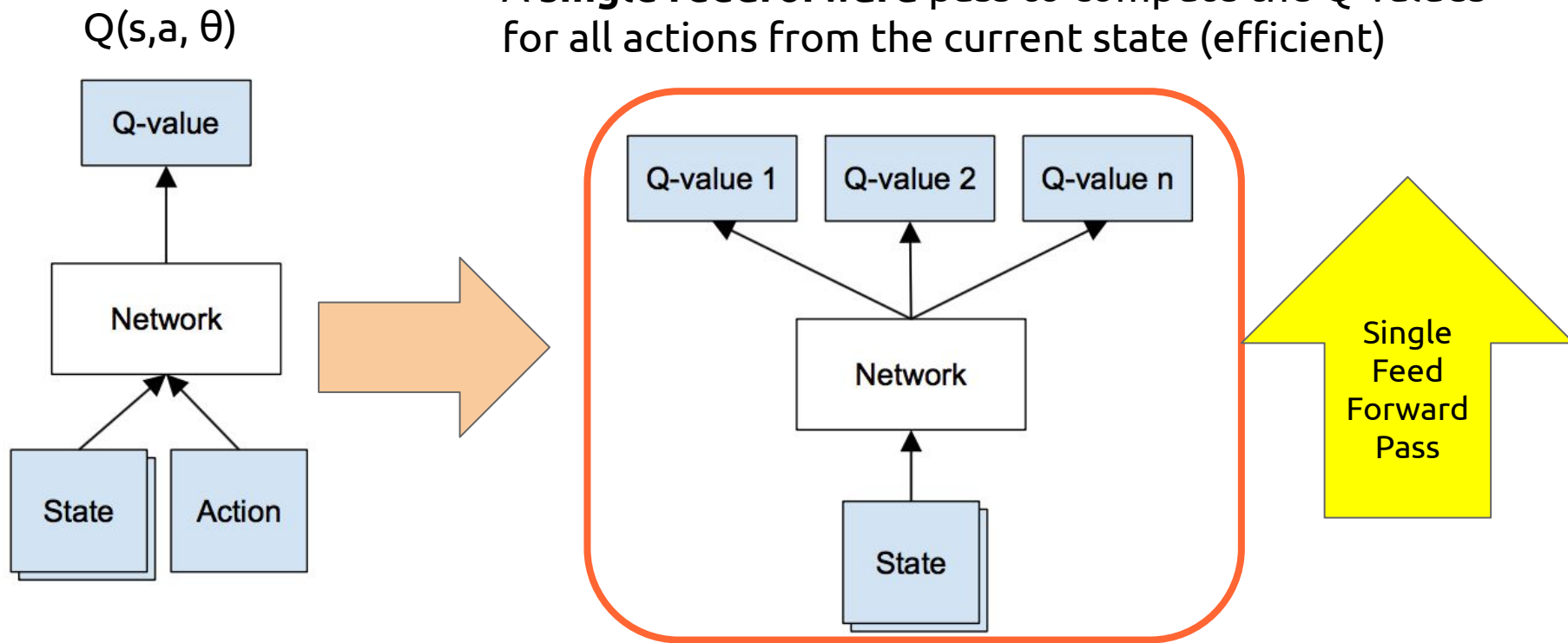
$$Q_*(s) = \arg \max_{a \in A} Q(s,a, \theta)$$

Eg. Neural Network  
parameters

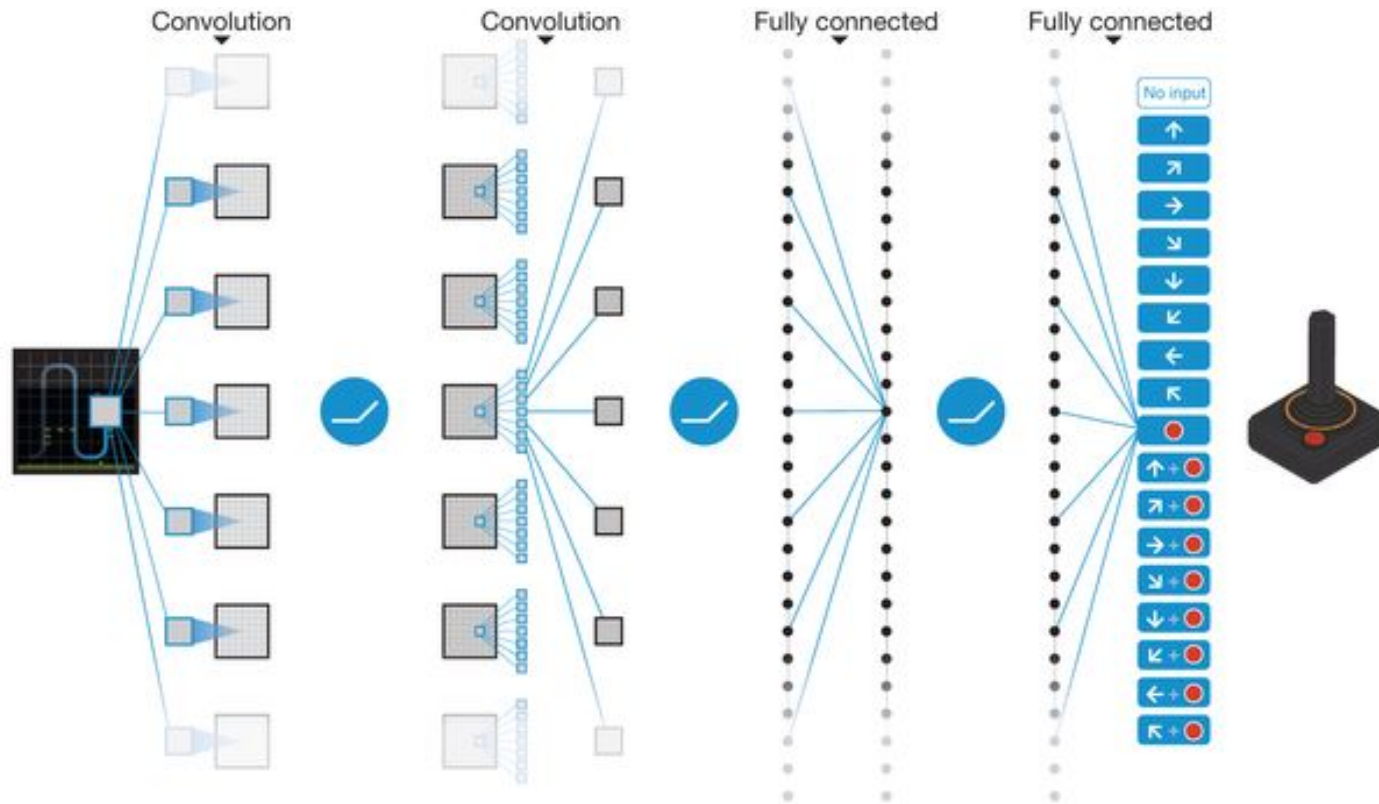


# Deep Q-learning (DQN)

A **single feedforward** pass to compute the Q-values for all actions from the current state (efficient)



# Deep Q-learning (DQN)

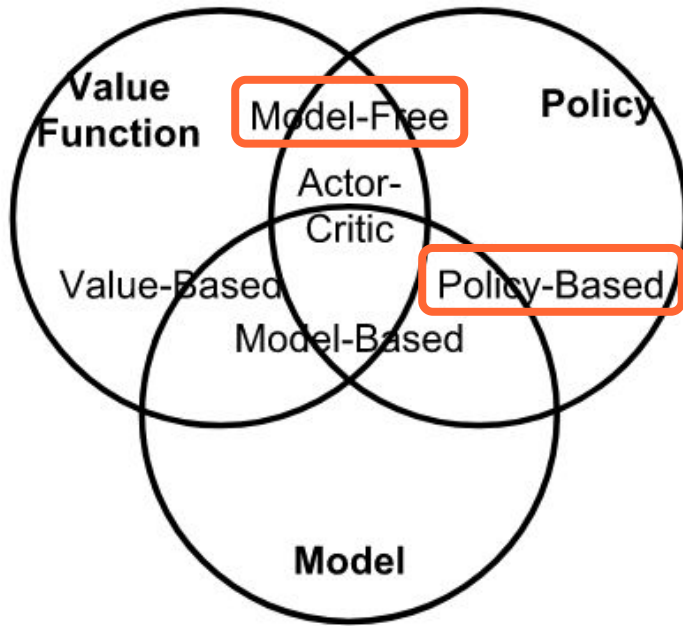


# Outline

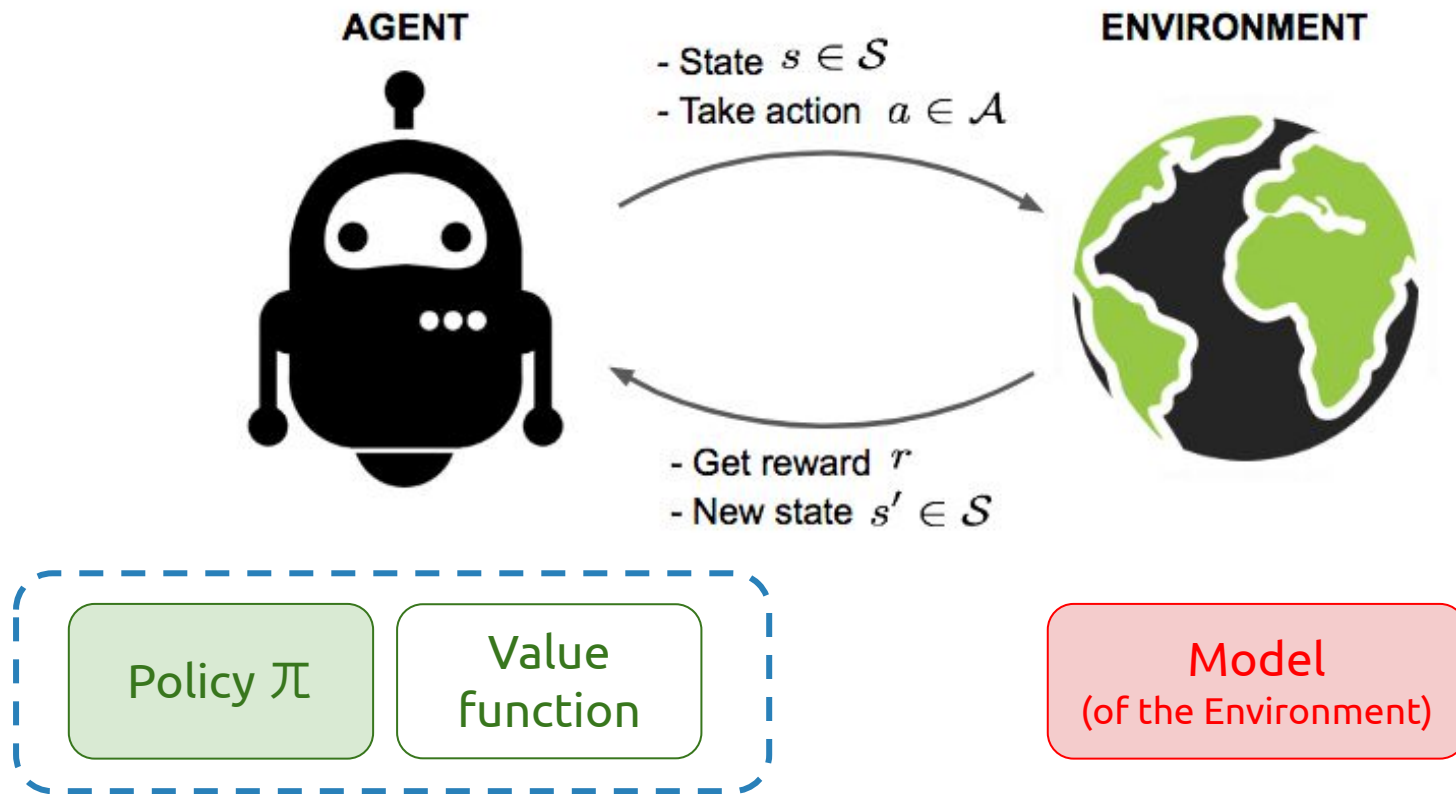
1. Control tasks with reinforcement learning
2. Markov Decision Processes
3. Learning action-value function  $Q$
4. **Learning policies**

# Value vs Policy vs Model

Summary of approaches in RL based on whether we want to learn the value, policy, or the model of the environment.



# Value vs Policy vs Model

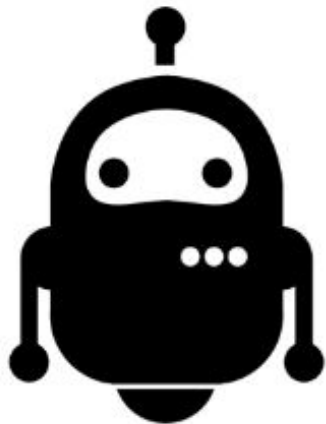


Goals of Reinforcement Learning



# Value vs Policy vs Model

AGENT



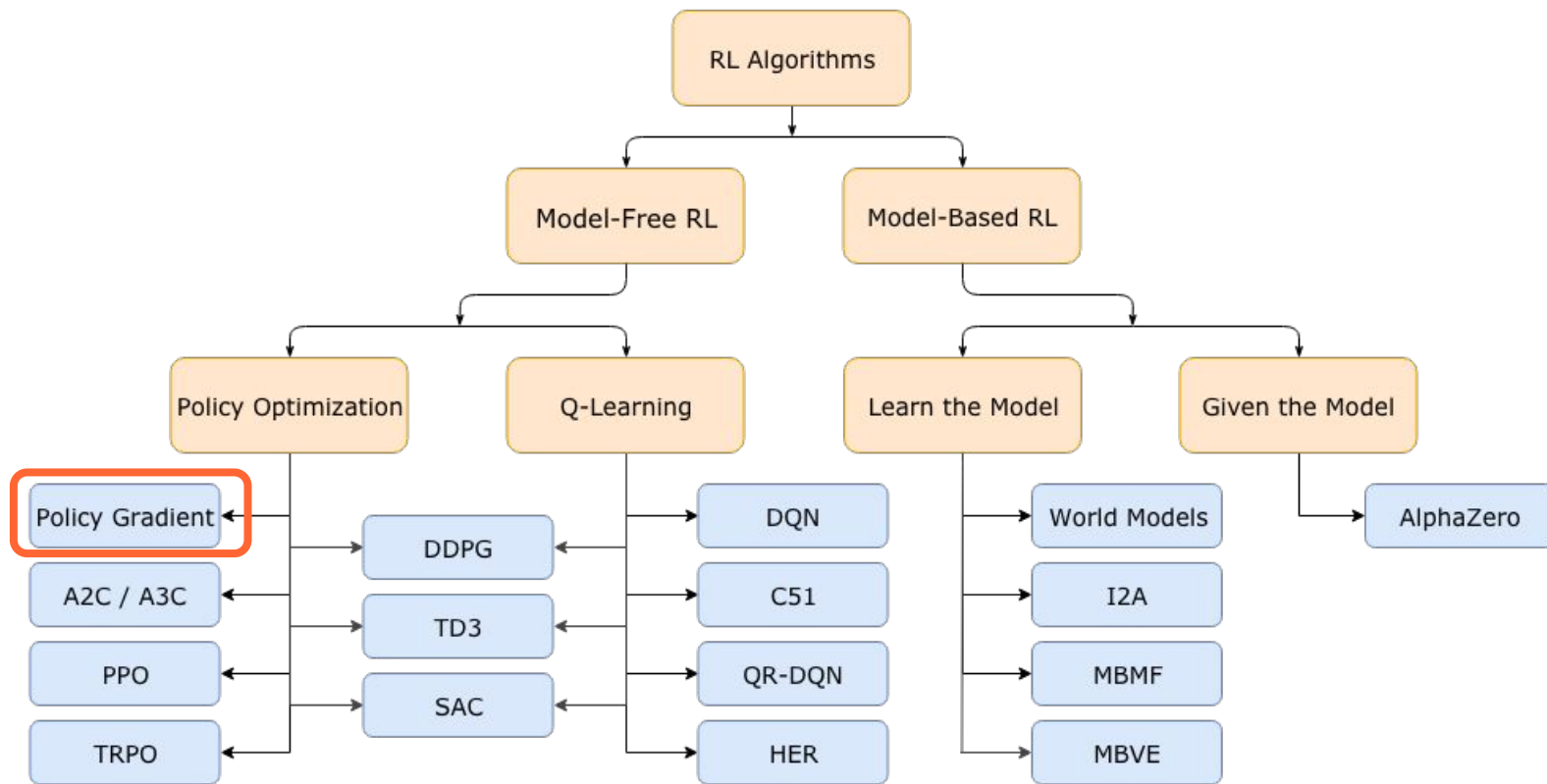
Policy  $\pi$

Value  
function

Directly learn the policy by estimating the parameters  $\theta$  of a stochastic policy function:

$$\pi(a|s;\theta)$$

# REINFORCE (Vanilla Policy Gradients - VPN)



# REINFORCE (Vanilla Policy Gradients - VPN)



## REINFORCE (Vanilla Policy Gradients - VPN)

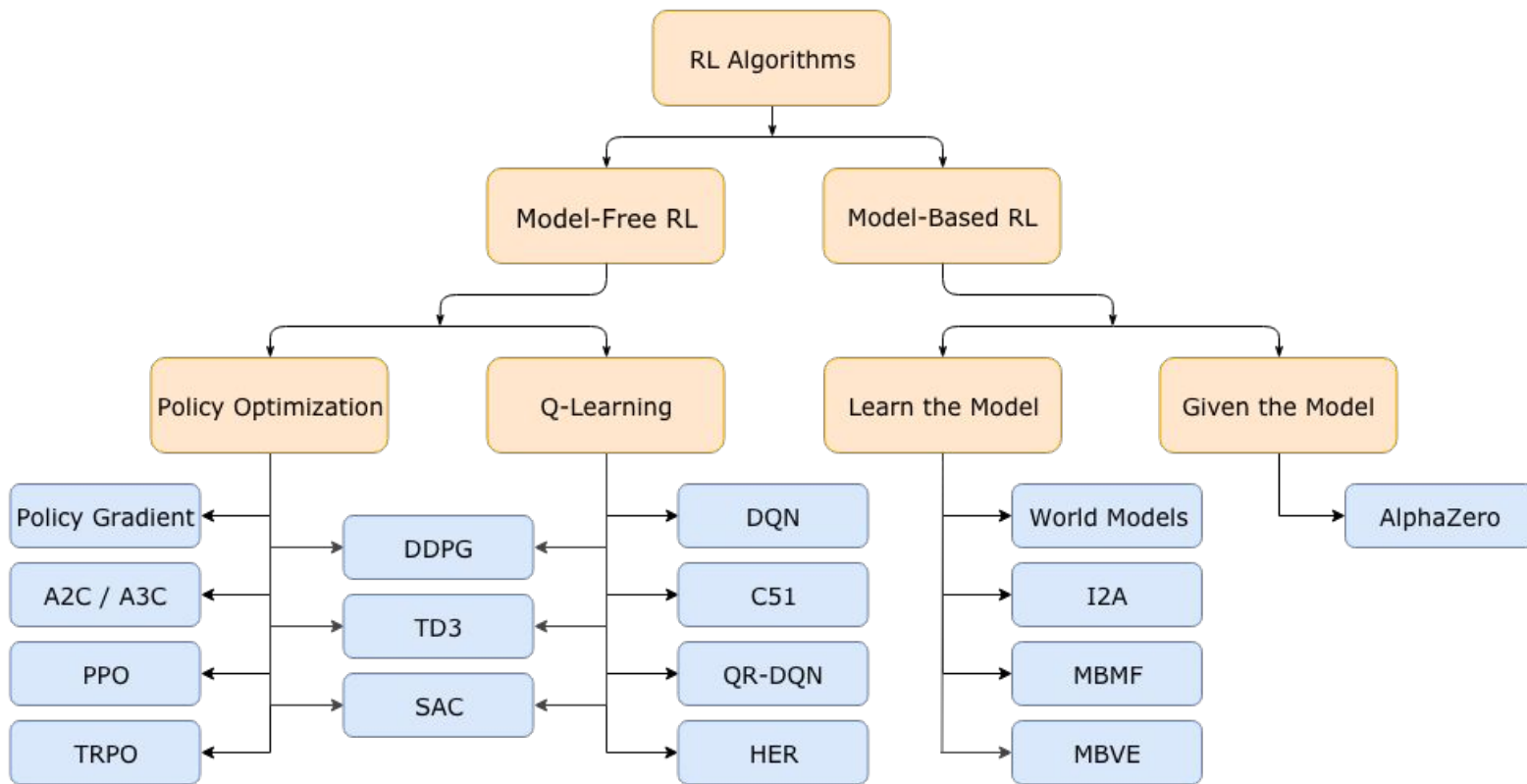
the mathematical formulation of 'trial-and-error': try an action, and make it more likely if it resulted in positive reward; otherwise, make it less likely.

Opposite goals in:

- **Supervised learning:** minimize a loss function by **gradient descent**.
- **Reinforcement learning:** maximize  $J(\theta)$ , the expected return of the policy, by **gradient ascent**.

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T R_{i,t} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t}, \mathbf{s}_{i,t}) \right]$$

# A taxonomy of RL Algorithms



# Outline

1. Control tasks with reinforcement learning
2. Markov Decision Processes
3. Learning action-value functions  $Q$
4. Learning policies

# Learn more

## UCL & Deepmind: [“Advanced Deep Learning & Reinforcement Learning”](#) (2018) [[slides](#)]









### Advanced Deep Learning & Reinforcement Learning

18 videos • 13,816 views • Updated 3 days ago



This course, taught originally at UCL and recorded for online access, has two interleaved parts that converge towards the end of the course. One part is on machine learning with deep neural networks, the other part is about prediction and control using reinforcement learning. The two strands come together when we discuss deep reinforcement learning, where deep neural networks are trained as function approximators in a reinforcement learning setting.

The deep learning stream of the course will cover a short

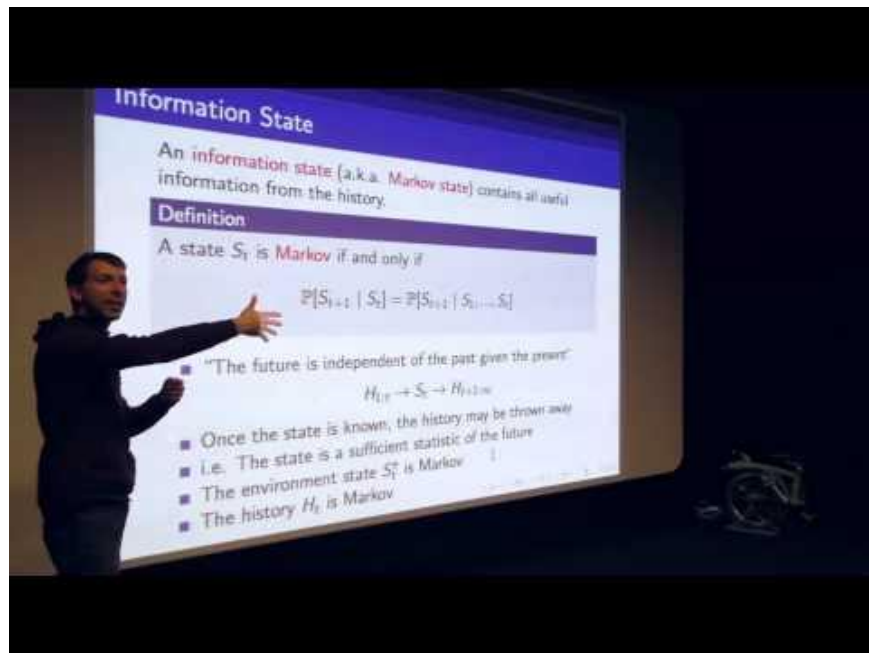
- 1  **Deep Learning 1: Introduction to Machine Learning Based AI**  
DeepMind  
1:43:08
- 2  **Deep Learning 2: Introduction to TensorFlow**  
DeepMind  
1:46:51
- 3  **Deep Learning 3: Neural Networks Foundations**  
DeepMind  
1:44:36
- 4  **Reinforcement Learning 1: Introduction to Reinforcement Learning**  
DeepMind  
1:43:17
- 5  **Reinforcement Learning 2: Exploration and Exploitation**  
DeepMind  
1:48:24
- 6  **Reinforcement Learning 3: Markov Decision Processes and Dynamic Programming**  
DeepMind  
1:44:24

Reinforcement Learning 3: Markov Decision Processes and Dynamic Programming

# Learn more



David Silver, UCL COMP050, [Reinforcement Learning](#)



# Learn more

Pieter Abbeel and John Schulman, [CS 294-112 Deep Reinforcement Learning](#), Berkeley.

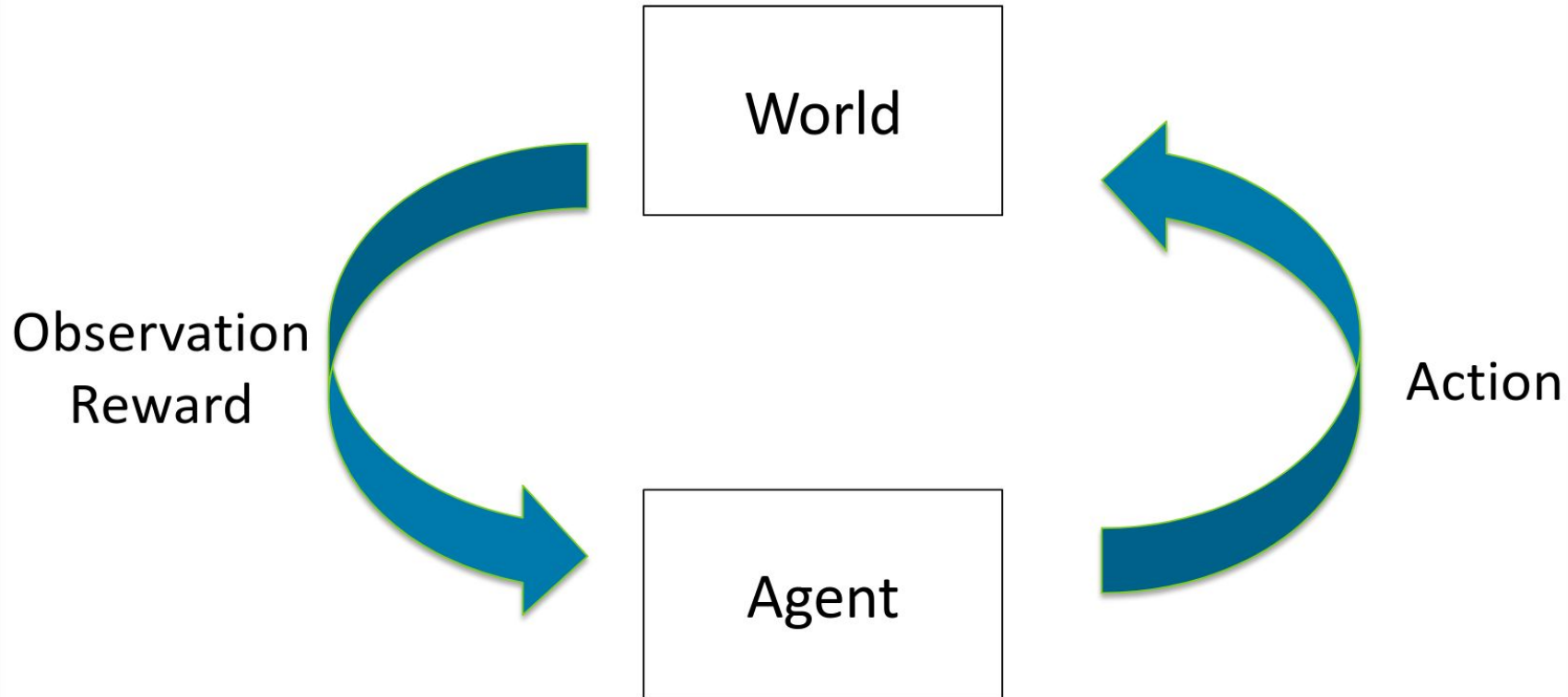
Slides: [“Reinforcement Learning - Policy Optimization”](#) OpenAI / UC Berkeley (2017)





# Learn more

Emma Brunskill, [“CS234: Reinforcement Learning”](#). Stanford University.



# Learn more

[OpenAI Spinning Up in Deep RL](#)



# OpenAI



# Final Questions

## Undergradese

What undergrads ask vs. what they're REALLY asking

"Is it going to be an open book exam?"

Translation: "I don't have to actually memorize anything, do I?"

"Hmm, what do you mean by that?"

Translation: "What's the answer so we can all go home."

"Are you going to have office hours today?"

Translation: "Can I do my homework in your office?"

"Can i get an extension?"

Translation: "Can you re-arrange your life around mine?"

"Is this going to be on the test?"

Translation: "Tell us what's going to be on the test."

"Is grading going to be curved?"

Translation: "Can I do a mediocre job and still get an A?"

JORGE CHAM © 2008



WWW.PHDCOMICS.COM