# Reinforcement Learning

## Lecture 9: Model-based methods

Chris G. Willcocks

Durham University
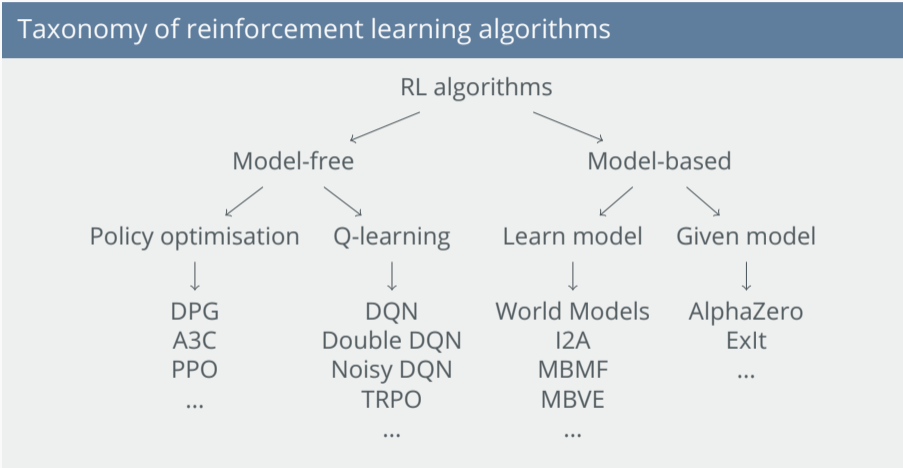
# Lecture overview

Lecture covers chapter 8 in Sutton & Barto [1] and examples from David Silver [2]

**1** **Model-based reinforcement learning**

- taxonomy
- overview
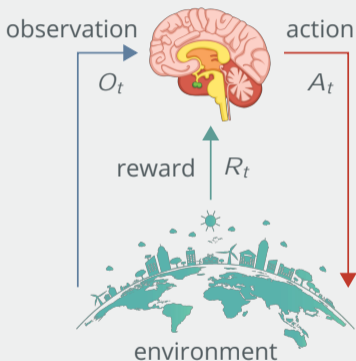- the simulation cycle
- characteristics

**2** **Integrated learning and planning**

- Dyna-Q
- characteristics
- Monte Carlo tree search
- simulated policy learning

## Taxonomy of reinforcement learning algorithms

RL algorithms

Model-free — Model-based

Policy optimisation — Q-learning — Learn model — Given model

Policy optimisation:
DPG
A3C
PPO
...

Q-learning:
DQN
Double DQN
Noisy DQN
TRPO
...

Learn model:
World Models
I2A
MBMF
MBVE
...

Given model:
AlphaZero
ExIt
...

This figure does not capture overlap, for example between policy optimsiation and Q-learning algorithms
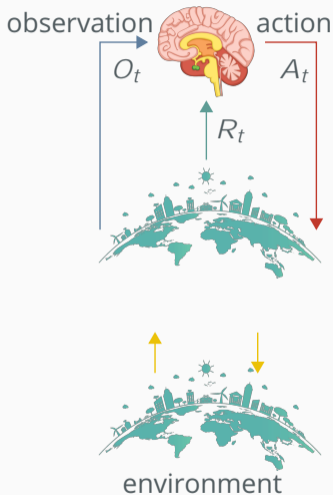
## RL Agents

In **model-free** RL:

No model

**Learn** the value function $q(s, a)$ and/or the policy $(a|s)$ from experience

In **model-based** RL:

Learn the model from experience

**Plan** the value function and/or the policy from the model

Figure based on [2, 1]

observation          action

$O_t$          $A_t$

$R_t$

environment

**Model-based** RL cycle:

The agent experiences the real environment

We learn a model to predict what the real environment does (when you take an action)

We then use this simulated model to plan

This allows us to estimate the value function and/or policy without directly interacting with the real environment

But we use this policy to take real actions again

**Model-based RL advantages**:

> The model can sometimes be a simpler and more useful representation of the environment than you can otherwise access by experience
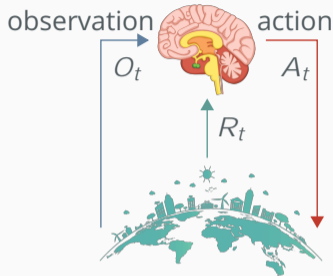
> Can be learnt by supervised learning

> Can reason about model uncertainty

**Model-based RL disadvantages**:

> This is another component which introduces some approximation error
>> Value function and/or policy approximation and now model approximation

> We can only be as good as our model

observation $O_t$   action $A_t$

$R_t$

environment

## Definition: model

A model $\mathcal{M} = \langle P, R \rangle$ is a parameterised representation of an MDP: $\langle S, A, P, R \rangle$. It approximates state transitions $P \approx P$ and rewards $R \approx R$, learning a distribution over the next states and rewards:

$$S_{t+1} \sim P(S_{t+1} | S_t, A_t)$$
$$R_{t+1} = R(R_{t+1} | S_t, A_t),$$

which typically are conditionally independent of each other:

$$P(S_{t+1}, R_{t+1} | S_t, A_t) = P(S_{t+1} | S_t, A_t) P(R_{t+1} | S_t, A_t)$$

## Example: environment model

### Learning the model

We learn the model $\mathcal{M}$ from experience $\{S_1, A_1, R_2, ..., S_T\}$ using **supervised learning**.

We receive a stream of actual experiences
This gives us a dataset:

$$S_1, A_1 \rightarrow R_2, S_2$$
$$S_2, A_2 \rightarrow R_3, S_3$$
$$...$$

$s, a \rightarrow r$ is a regression problem
$s, a \rightarrow s'$ is a density estimation problem

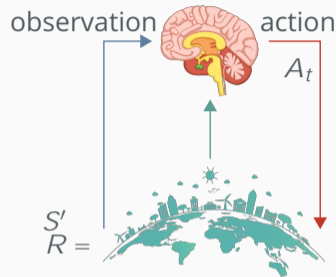## Experience can be simulated and real

**Simulated experience** sampled from $\mathcal{M}$

$$S' \sim \mathcal{P}(S'|S, A)$$
$$R = \mathcal{R}(R|S, A)$$

**Real experience** sampled from the true MDP

$$S' \sim P^a_{s, s'}$$
$$R = R^a_s$$

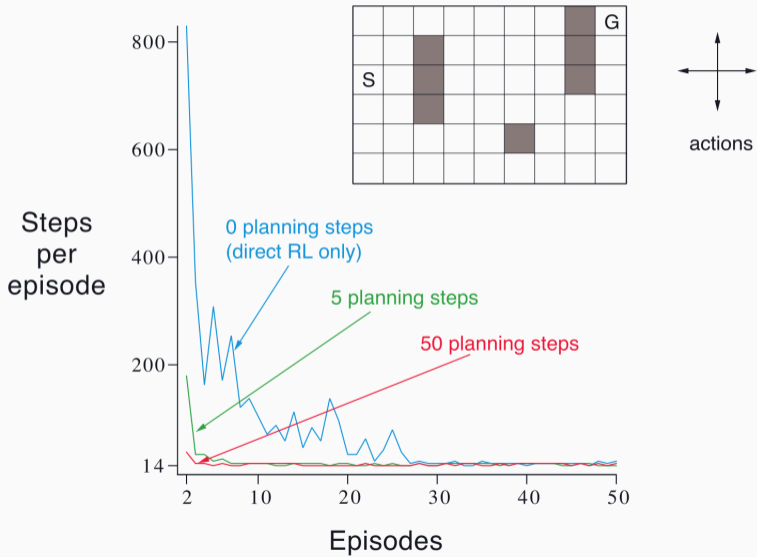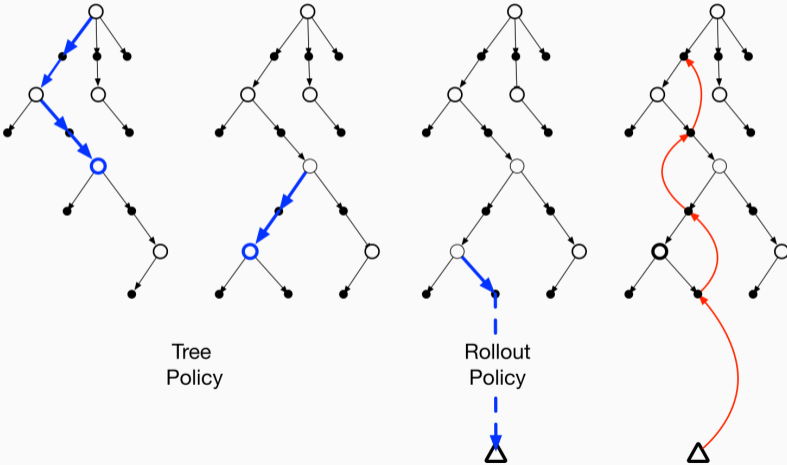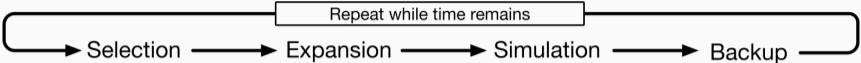observation    action

$A_t$

$S'$
$R =$

$S'$
$R =$

environment

**Algorithm:** Dyna-Q [3, 4]

```
initialise Q(s; a) and model M(s; a) for all s ∈ S and a ∈ A(s)
while True:
  s ← current (nonterminal) state
  a ← greedy(s; Q)
  r; s' ← env.step(s; a)
  Q(s; a) ← Q(s; a) + α(r + γ max_a Q(s'; a) − Q(s; a))
  M(s; a) ← r; s' (assuming deterministic environment)
  for i in range(n):
    s ← random previously observed state
    a ← random action previously taken in s
    r; s' ← M(s; a)
    Q(s; a) ← Q(s; a) + α(r + γ max_a Q(s'; a) − Q(s; a))
```

Repeat while time remains

Selection ⟶ Expansion ⟶ Simulation ⟶ Backup

Tree
Policy

Rollout
Policy

"Model-Based Reinforcement Learning for Atari" [5]

## Summary

In summary, model-based methods:

are easy to train with supervised learning
allow for planning ahead
can be very data efficient
can be used to imagine situations without experiencing them
but the value and policy learnt can only be as good as the model
they can be combined with model-free methods

[1] Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction (second edition). Available online 📥. MIT press, 2018.

[2] David Silver. Reinforcement Learning lectures. https://www.davidsilver.uk/teaching/. 2015.

[3] Richard S Sutton. "Integrated architectures for learning, planning, and reacting based on approximating dynamic programming". In: Machine learning proceedings 1990. Elsevier, 1990, pp. 216–224.

[4] Baolin Peng et al. "Deep Dyna-Q: Integrating planning for task-completion dialogue policy learning". In: arXiv preprint arXiv:1801.06176 (2018).

[5] Lukasz Kaiser et al. "Model-based reinforcement learning for atari". In: arXiv preprint arXiv:1903.00374 (2019).