

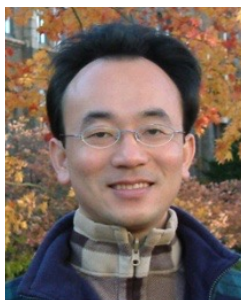
# 資料探勘 (Data Mining)

# 強化學習 (Reinforcement Learning)

1092DM11

MBA, IM, NTPU (M5026) (Spring 2021)

Tue 2, 3, 4 (9:10-12:00) (B8F40)



Min-Yuh Day

戴敏育

Associate Professor

副教授

Institute of Information Management, National Taipei University

國立臺北大學 資訊管理研究所

<https://web.ntpu.edu.tw/~myday>

2021-06-01



# 課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
1	2021/02/23	資料探勘介紹 (Introduction to data mining)
2	2021/03/02	ABC：人工智慧，大數據，雲端運算 (ABC: AI, Big Data, Cloud Computing)
3	2021/03/09	Python 資料探勘的基礎 (Foundations of Data Mining in Python)
4	2021/03/16	資料科學與資料探勘：發現，分析，可視化和呈現數據 (Data Science and Data Mining: Discovering, Analyzing, Visualizing and Presenting Data)
5	2021/03/23	非監督學習：關聯分析，購物籃分析 (Unsupervised Learning: Association Analysis, Market Basket Analysis)
6	2021/03/30	資料探勘個案研究 I (Case Study on Data Mining I)

# 課程大綱 (Syllabus)

- | 週次 (Week) | 日期 (Date)  | 內容 (Subject/Topics)   |
|-----------|------------|---|
| 7         | 2021/04/06 | 放假一天 (Day off)  |
| 8         | 2021/04/13 | 非監督學習：集群分析，行銷市場區隔<br>(Unsupervised Learning: Cluster Analysis, Market Segmentation) |
| 9         | 2021/04/20 | 期中報告 (Midterm Project Report)   |
| 10        | 2021/04/27 | 監督學習：分類和預測<br>(Supervised Learning: Classification and Prediction)                  |
| 11        | 2021/05/04 | 機器學習和深度學習<br>(Machine Learning and Deep Learning)                                   |
| 12        | 2021/05/11 | 卷積神經網絡<br>(Convolutional Neural Networks)   |

# 課程大綱 (Syllabus)

週次 (Week)	日期 (Date)	內容 (Subject/Topics)
13	2021/05/18	資料探勘個案研究 II (Case Study on Data Mining II)
14	2021/05/25	遞歸神經網絡 (Recurrent Neural Networks)
15	2021/06/01	強化學習 (Reinforcement Learning)
16	2021/06/08	社交網絡分析 (Social Network Analysis)
17	2021/06/15	期末報告 I (Final Project Report I)
18	2021/06/22	期末報告 II (Final Project Report II)

# Reinforcement Learning

# Outline

- Reinforcement Learning (RL)
  - Markov Decision Processes (MDP)
- Deep Reinforcement Learning (DRL) Algorithms
  - SARSA
  - Q-Learning
  - DQN
  - A3C
  - Rainbow

# AI, ML, DL

**Artificial Intelligence (AI)**

**Machine Learning (ML)**

**Supervised  
Learning**

**Unsupervised  
Learning**

**Deep Learning (DL)**

**CNN**

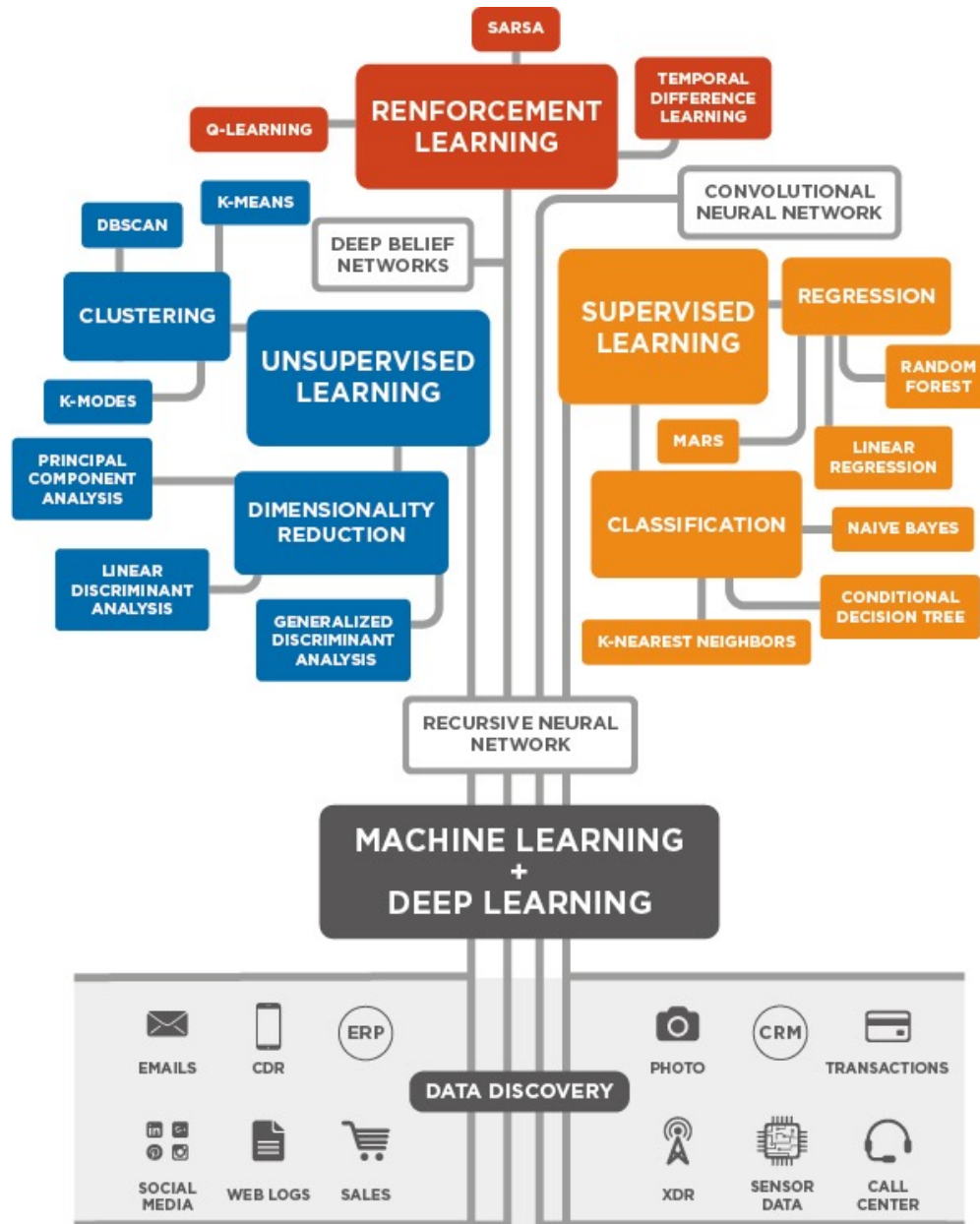
**RNN LSTM GRU**

**GAN**

**Semi-supervised  
Learning**

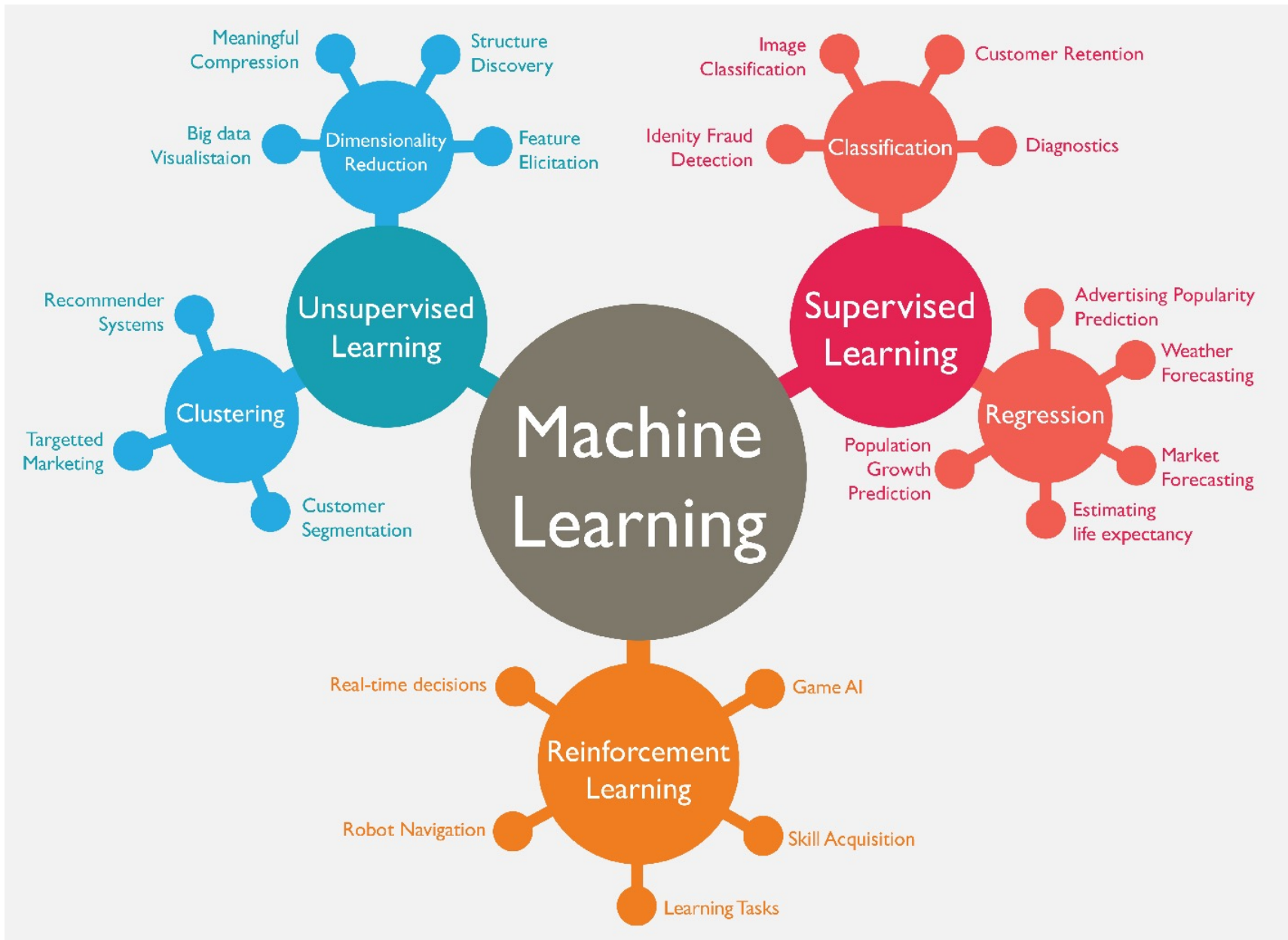
**Reinforcement  
Learning**

# 3 Machine Learning Algorithms

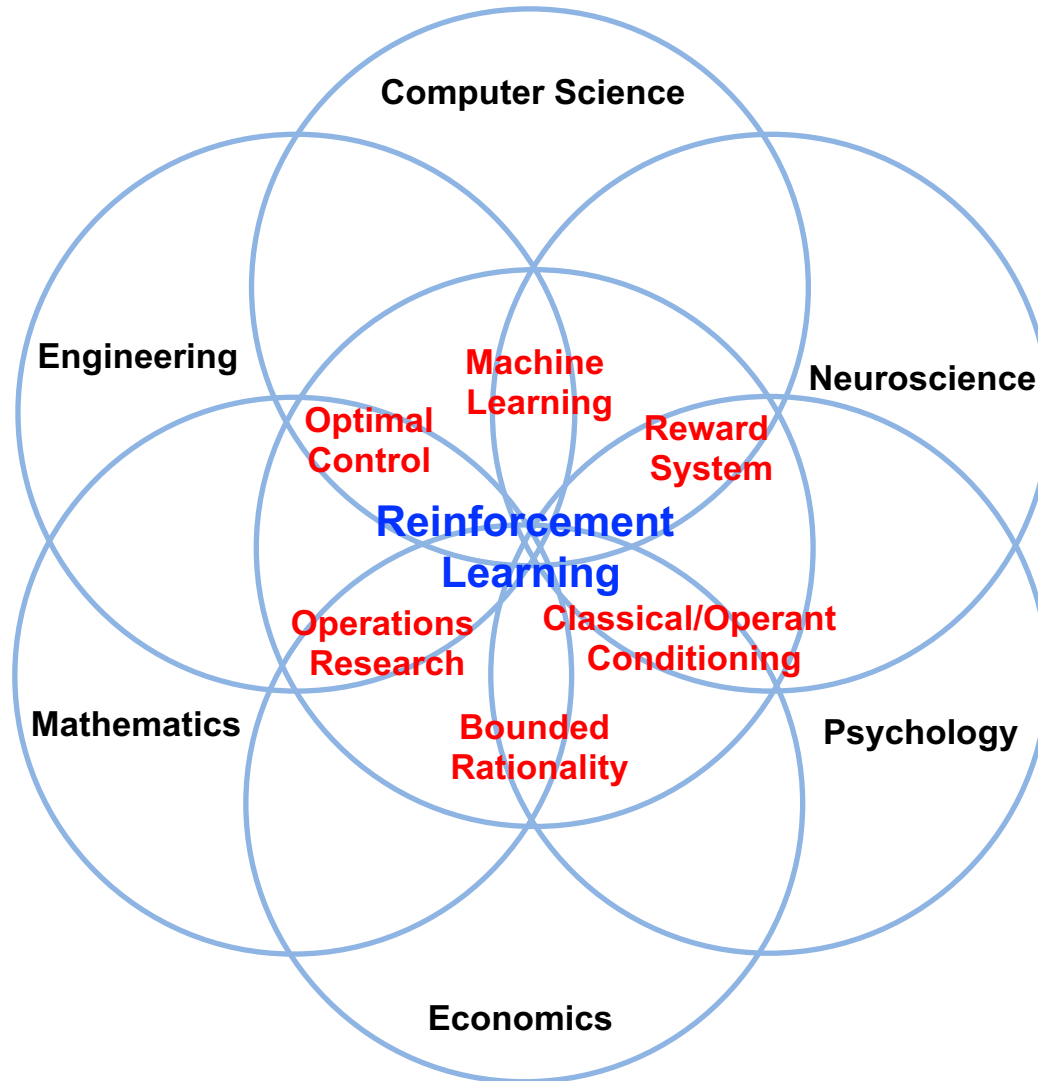




# Machine Learning (ML)



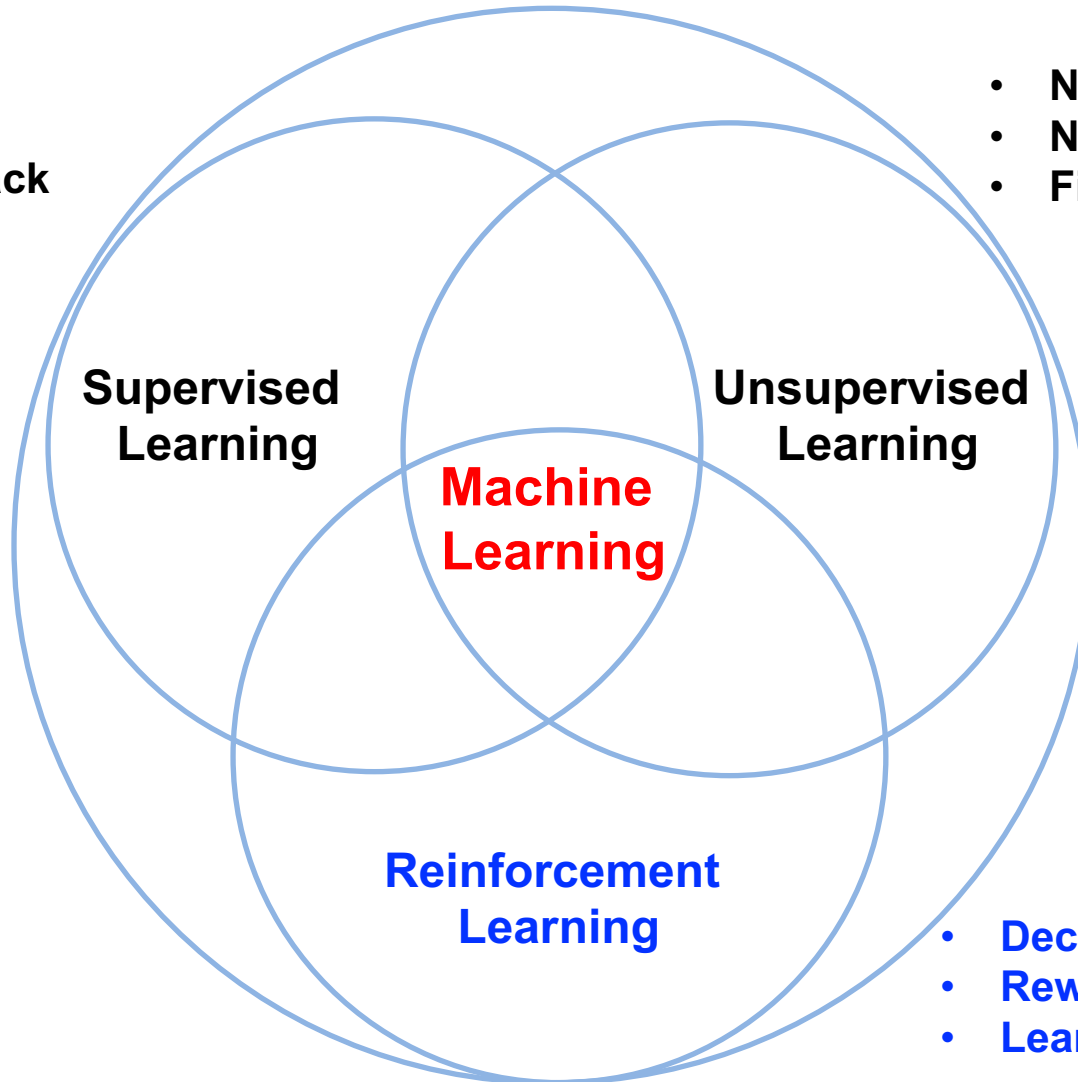
# Reinforcement Learning (RL)



# Branches of Machine Learning (ML)

## Reinforcement Learning (RL)

- Labeled data
- Direct feedback
- Predict



- No Labels
- No feedback
- Find hidden structure

- Decision process
- Reward system
- Learn series of actions

# David Silver (2015), Introduction to reinforcement learning

- **Elementary Reinforcement Learning**
  - 1: Introduction to Reinforcement Learning
  - 2: Markov Decision Processes
  - 3: Planning by Dynamic Programming
  - 4: Model-Free Prediction
  - 5: Model-Free Control
- **Reinforcement Learning in Practice**
  - 6: Value Function Approximation
  - 7: Policy Gradient Methods
  - 8: Integrating Learning and Planning
  - 9: Exploration and Exploitation
  - 10: Case Study: RL in Classic Games

# Reinforcement Learning

## AlphaZero (AZ) and AlphaGo Zero (AZ0)

- AlphaZero (Silver et al., 2018)
  - A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. (Science)
- AlphaGo Zero (Silver et al., 2017)
  - Mastering the game of Go without human knowledge (Nature)

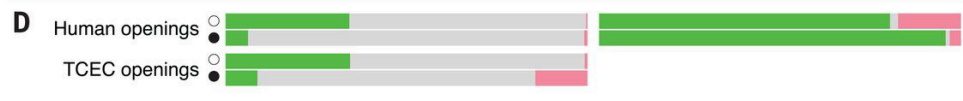
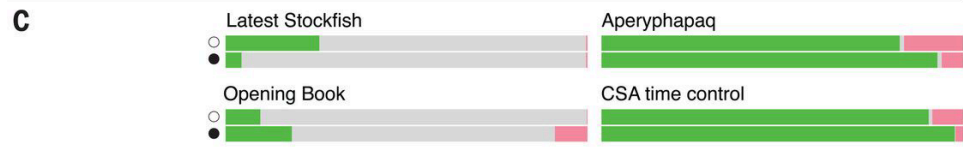
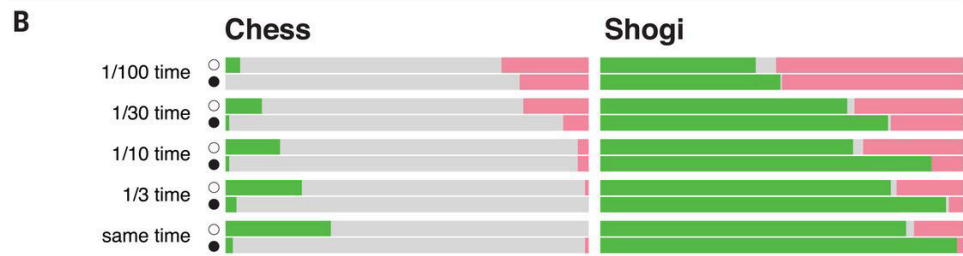
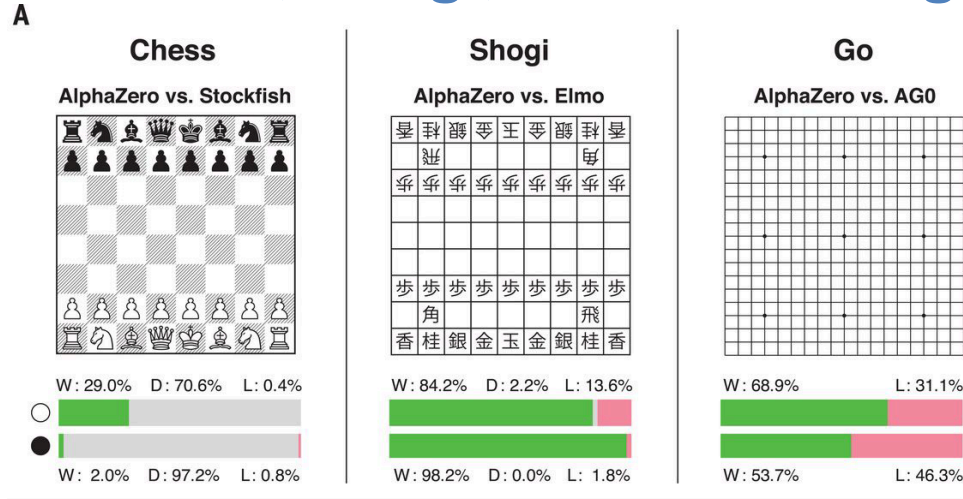
# AlphaZero: Shedding new light on the grand games of chess, shogi and Go



<https://www.youtube.com/watch?v=7L2sUGcOgh0>

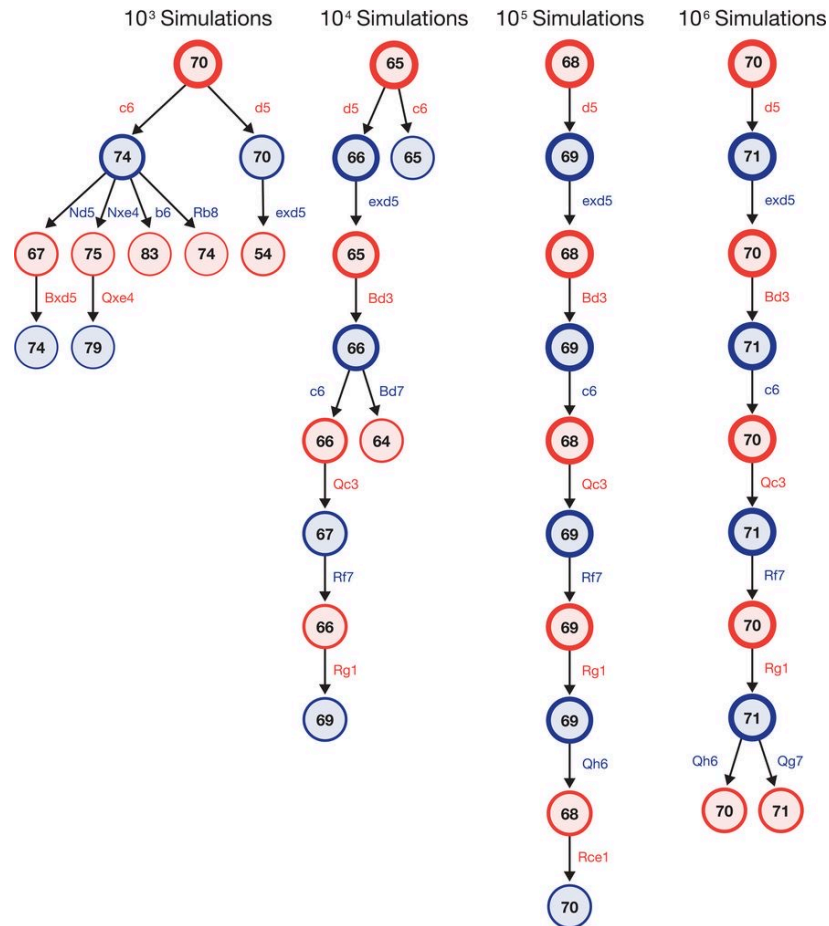
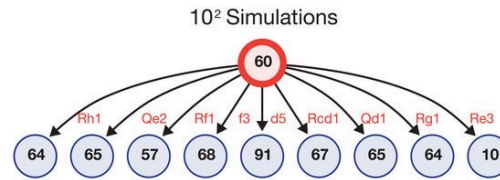
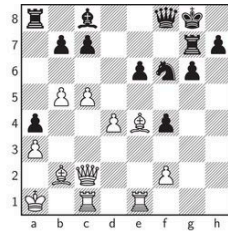
# AlphaZero

## A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play



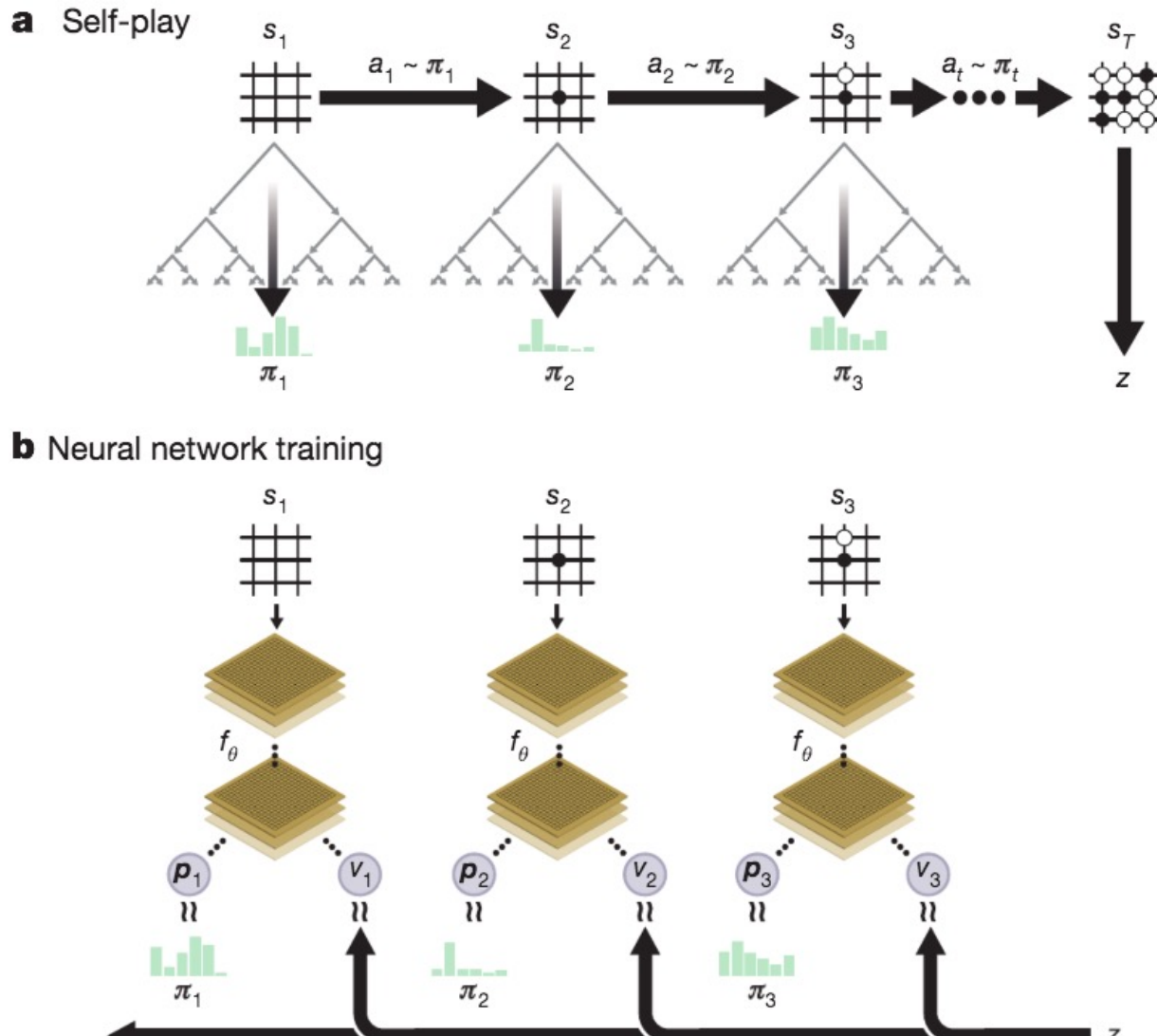
■ AlphaZero wins 
 ■ AlphaZero draws 
 ■ AlphaZero loses 
 ○ AlphaZero white 
 ● AlphaZero black

# AlphaZero's search procedure

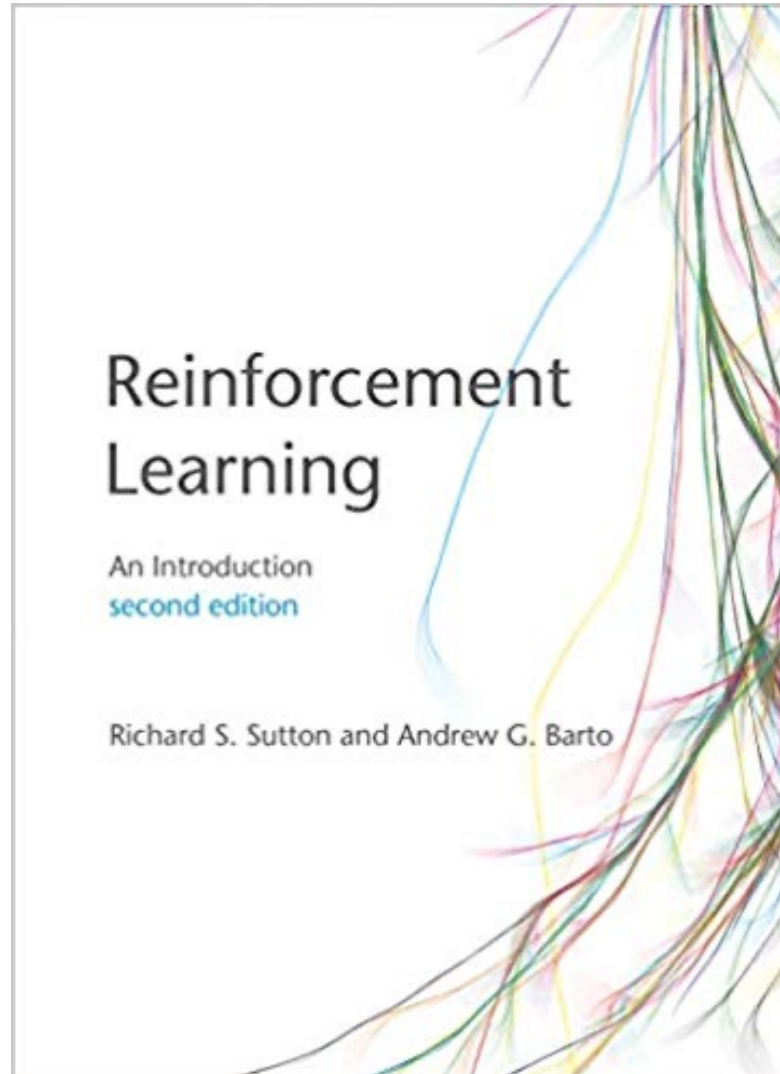




# Self-play reinforcement learning in AlphaGo Zero



Richard S. Sutton & Andrew G. Barto (2018),  
**Reinforcement Learning: An Introduction,**  
2<sup>nd</sup> Edition, A Bradford Book



Source: Richard S. Sutton & Andrew G. Barto (2018), Reinforcement Learning: An Introduction, 2nd Edition, A Bradford Book.  
<https://www.amazon.com/Reinforcement-Learning-Introduction-Adaptive-Computation/dp/0262039249>

# Reinforcement learning

- Reinforcement learning is **learning what to do**
  - how to map **situations** to **actions**
  - so as to maximize a numerical **reward** signal.

# Two most important distinguishing features of reinforcement learning

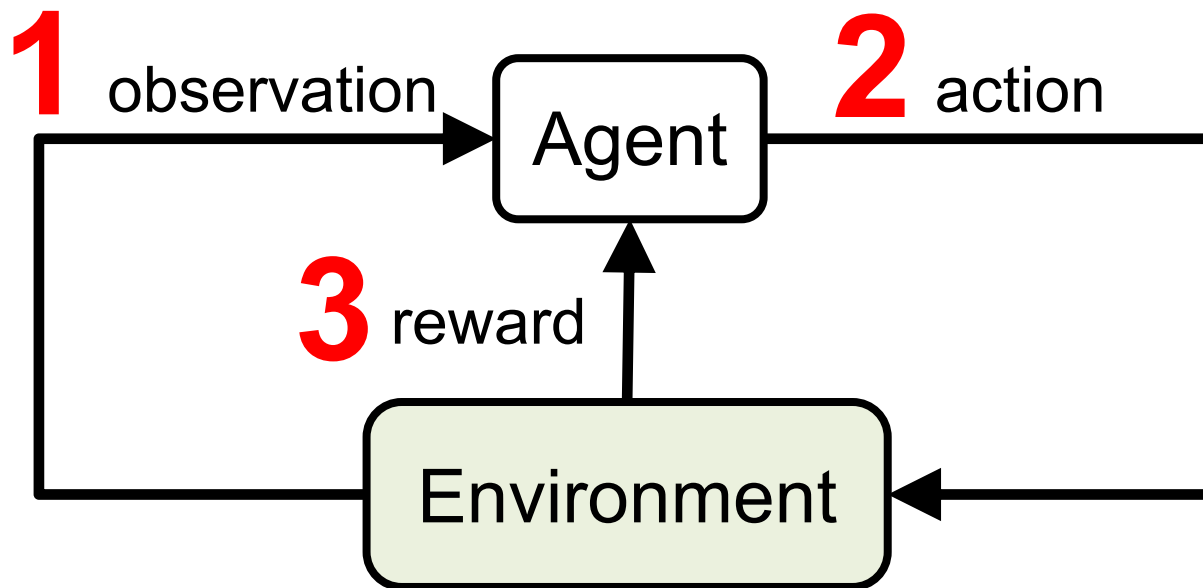
- trial-and-error search
- delayed reward

# Reinforcement Learning (DL)

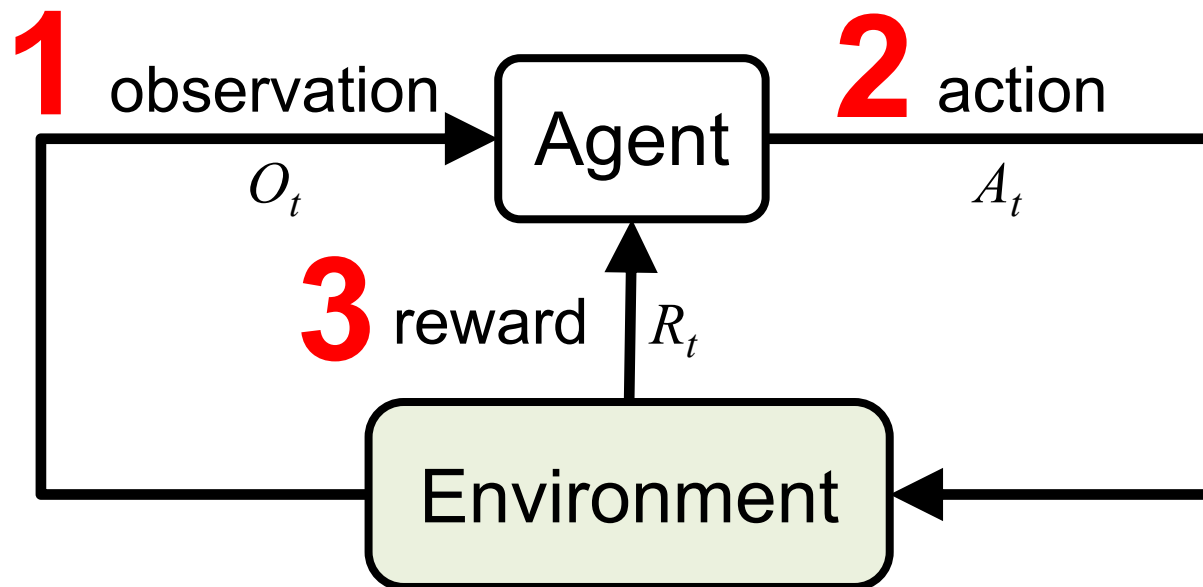
Agent

Environment

# Reinforcement Learning (DL)

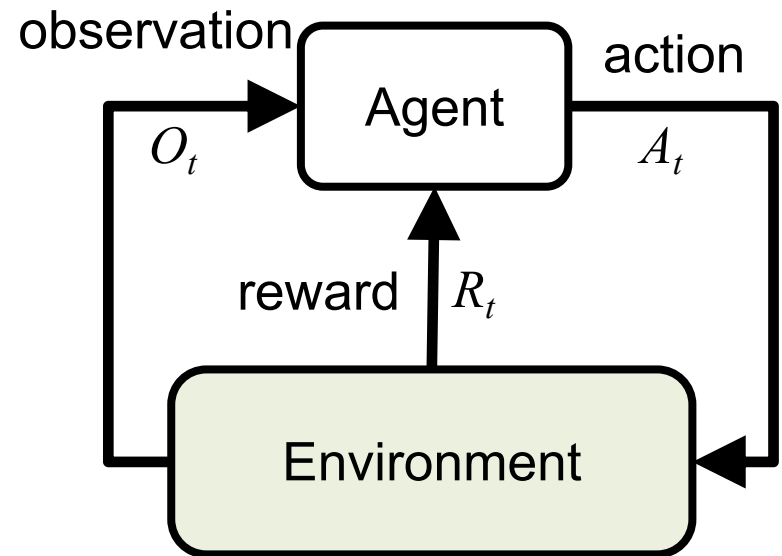


# Reinforcement Learning (DL)



# Agent and Environment

- At each step  $t$  the agent:
  - Executes **action**  $A_t$
  - Receives **observation**  $O_t$
  - Receives scalar **reward**  $R_t$
- The environment:
  - Receives action  $A_t$
  - Emits observation  $O_{t+1}$
  - Emits scalar reward  $R_{t+1}$
- $t$  increments at env. step





# History and State

- The **history** is the sequence of observations, actions, rewards

$$H_t = O_1, A_1, R_1, \dots, A_{t-1}, O_t, R_t$$

- i.e. all observable variables up to time  $t$
- i.e. the sensorimotor stream of a robot or embodied agent
- What happens next depends on the history:
  - The agent selects actions
  - The environment selects observations/rewards
- **State** is the information used to determine what happens next
- Formally, state is a function of the history:

$$S_t = f(H_t)$$

# Information State

- An **information state** (a.k.a. **Markov state**) contains all useful information from the history.

- Definition

A state  $S_t$  is **Markov** if and only if

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$$

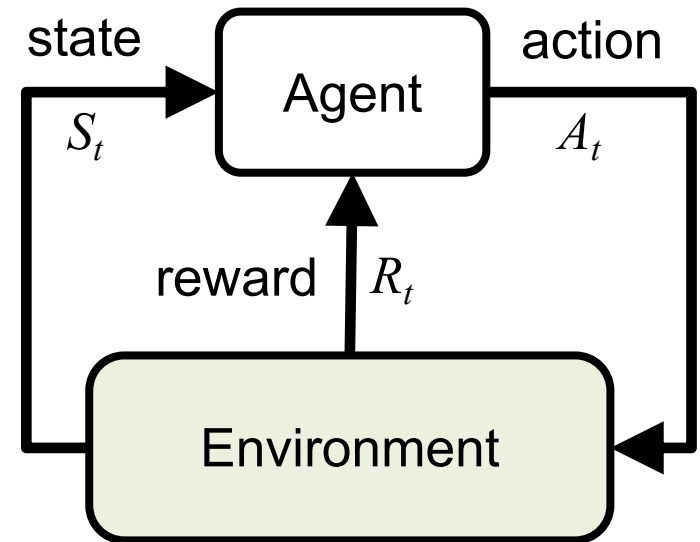
- “The future is independent of the past given the present”

$$H_{1:t} \rightarrow S_t \rightarrow H_{t+1:\infty}$$

- Once the state is known, the history may be thrown away i.e. The state is a sufficient statistic of the future
- The environment state  $S_t^e$  is Markov
- The history  $H_t$  is Markov

# Fully Observable Environments

- **Full observability:**
  - agent **directly** observes environment state
  - Agent state = environment state = information state
  - Formally, this is a **Markov decision process (MDP)**

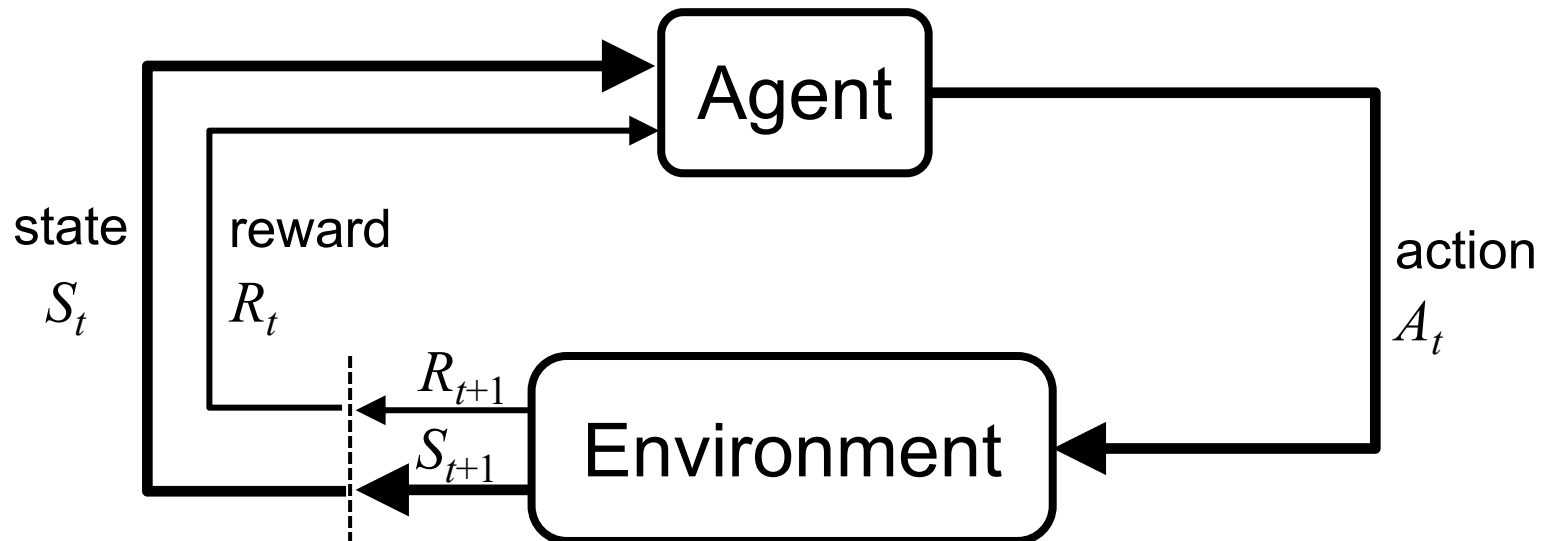


# Partially Observable Environments

- **Partial observability**: agent **indirectly** observes environment
  - A robot with camera vision isn't told its absolute location
  - A trading agent only observes current prices
  - A poker playing agent only observes public cards
- Now agent state  $\neq$  environment state
- Formally this is a **partially observable Markov decision process (POMDP)**
- Agent must construct its own state representation  $S^a_t$ , e.g.
  - Complete history:  $S^a_t = H_t$
  - **Beliefs** of environment state:  $S^a_t = (P[S^e_t = s_1], \dots, P[S^e_t = s_n])$
  - Recurrent neural network:  $S^a_t = \sigma(S^a_{t-1} W_s + O_t W_o)$

# Reinforcement Learning (DL)

The Agent-Environment Interaction  
in a Markov Decision Process (MDP)



# Characteristics of Reinforcement Learning

- No supervisor, only a **reward** signal
- Feedback is **delayed**, not instantaneous
- **Time** really matters  
(**sequential**, non i.i.d data)
- Agent's **actions** affect the subsequent data it receives

# Examples of Reinforcement Learning

- Make a humanoid robot walk
- Play many different Atari games better than humans
- Manage an investment portfolio

# Examples of Rewards

- Make a humanoid robot walk
  - +ve reward for forward motion
  - -ve reward for falling over
- Play many different Atari games better than humans
  - +/-ve reward for increasing/decreasing score
- Manage an investment portfolio
  - +ve reward for each \$ in bank



# Sequential Decision Making

- Goal: **select actions to maximize total future reward**
- **Actions** may have long term consequence
- **Reward** may be delayed
- It may be better to sacrifice immediate reward to gain more long-term reward
- Examples:
  - A financial investment (may take months to mature)
  - Blocking opponent moves (might help winning chances many moves from now)

# Elements of Reinforcement Learning

- Agent
- Environment
- Policy
- Reward signal
- Value function
- Model

# Elements of Reinforcement Learning

- Policy
  - Agent's **behavior**
  - It is a map from state to action
- Reward signal
  - The **goal** of a reinforcement learning problem
- Value function
  - How good is each state and/or action
  - A prediction of future reward
- Model
  - Agent's representation of the environment

# Major Components of an RL Agent

- 1. Policy:** agent's behaviour function
- 2. Value** function: how good is each state and/or action
- 3. Model:** agent's representation of the environment

# Policy

- A **policy** is the agent's behaviour
- It is a map from state to action, e.g.
  - Deterministic policy:  $a = \pi(s)$
  - Stochastic policy:  $\pi(a|s) = P[A_t = a | S_t = s]$

# Value Function

- **Value function** is a prediction of future reward
- Used to evaluate the goodness/badness of states
- And therefore to select between actions, e.g.

$$v_{\pi}(s) = E_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

# Model

- A **model** predicts what the environment will do next
- $P$  predicts the next state
- $R$  predicts the next (immediate) reward, e.g.

$$P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$$

$$R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

# Reinforcement Learning

- Value Based
  - No Policy (Implicit)
  - Value Function
- Policy Based
  - Policy
  - No Value Function
- Actor Critic
  - Policy
  - Value Function

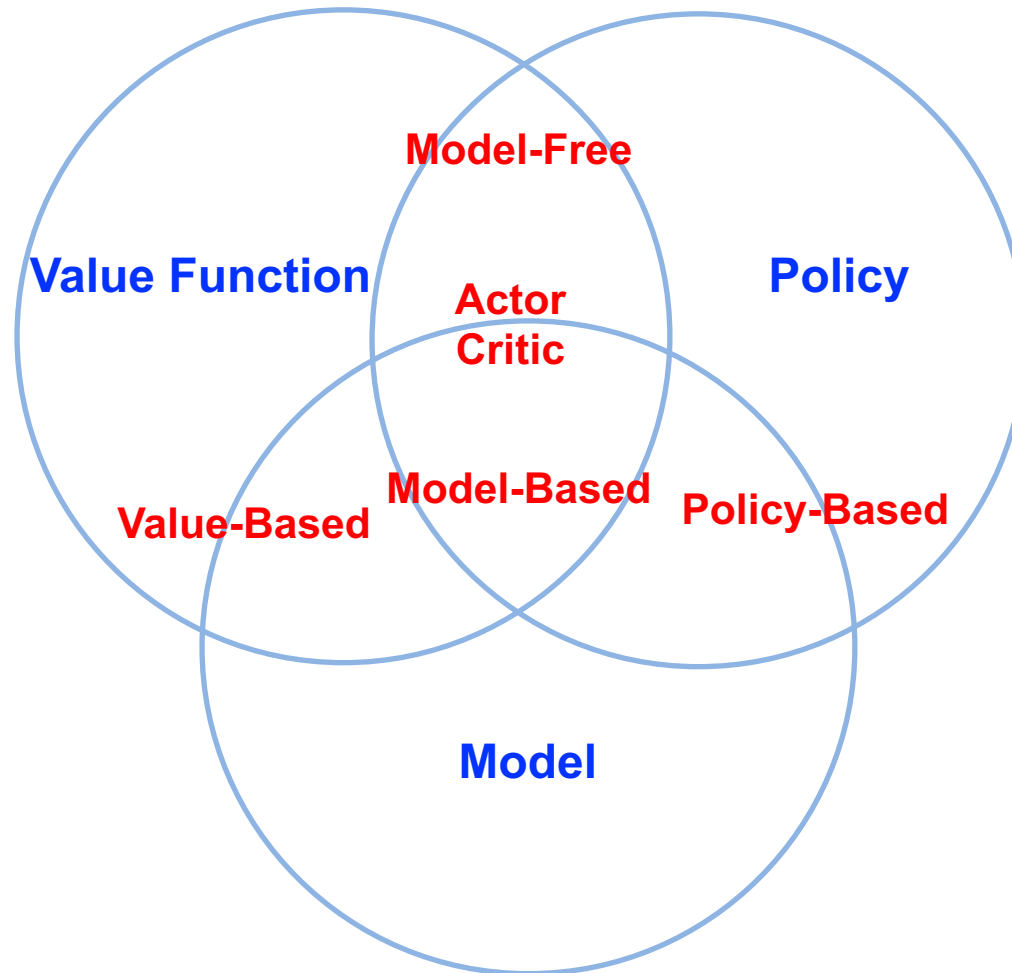


# Reinforcement Learning

- Model Free
  - Policy and/or Value Function
  - No Model
- Model Based
  - Policy and/or Value Function
  - Model

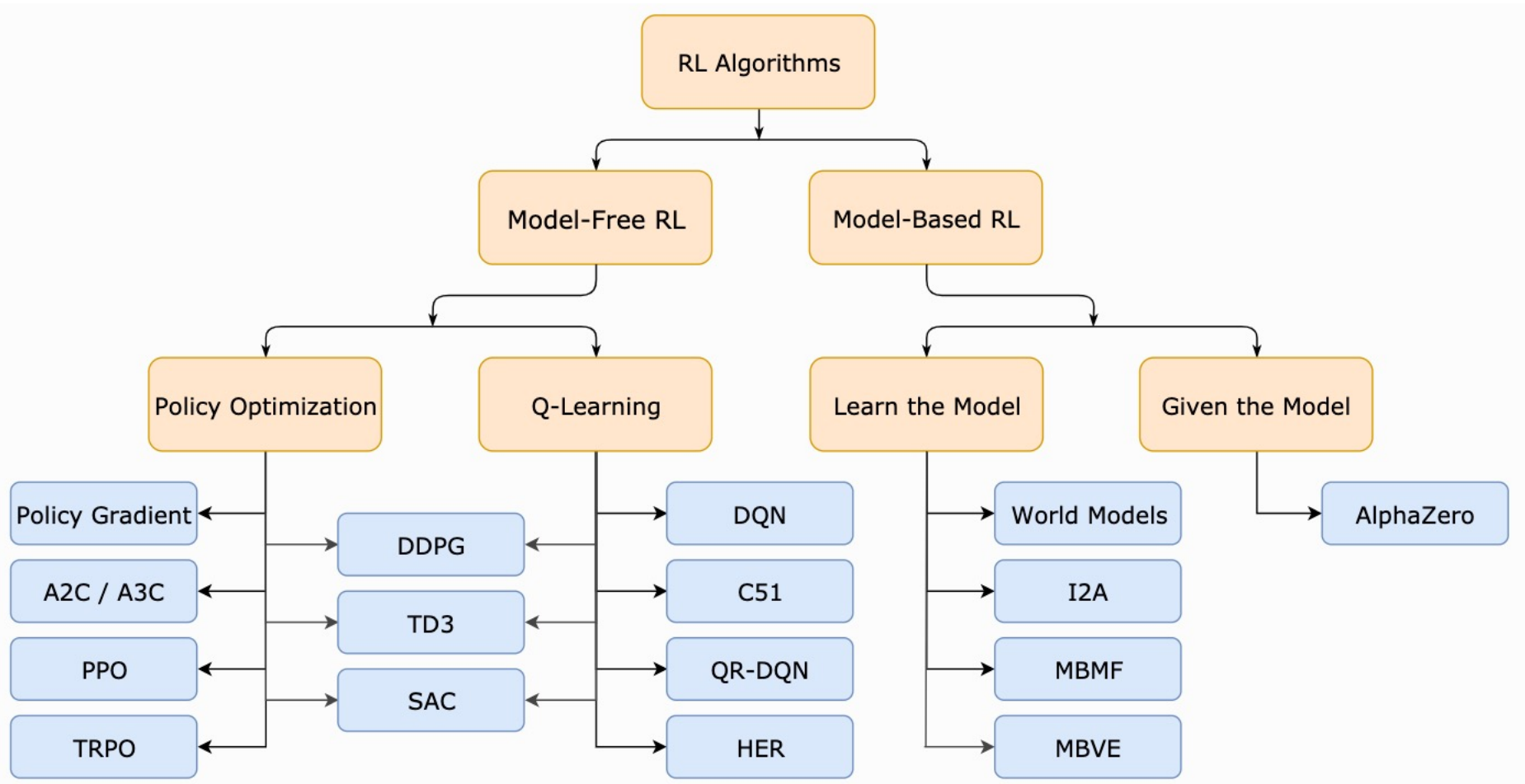
# Reinforcement Learning (RL)

## Taxonomy



# Reinforcement Learning (RL)

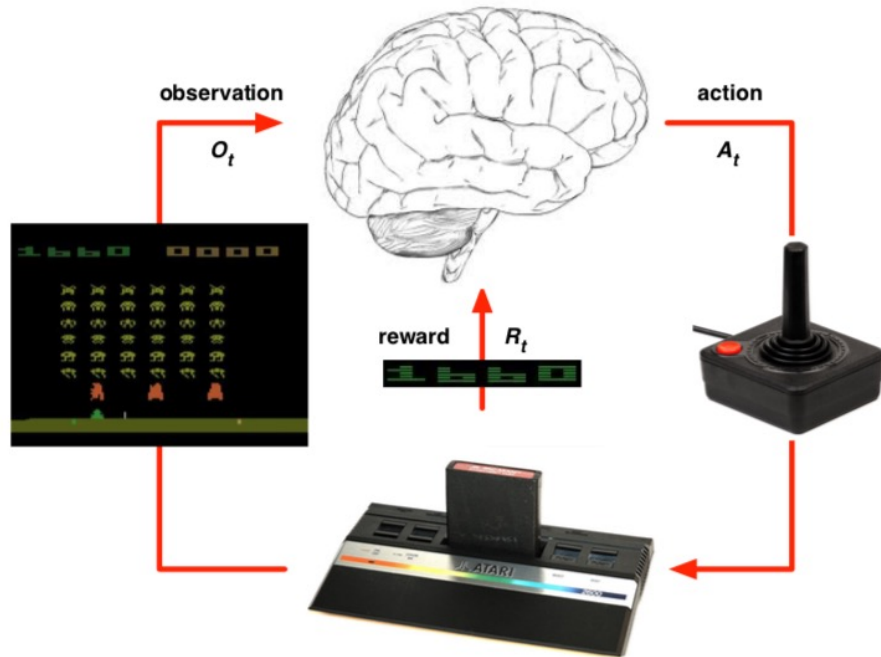
## A Taxonomy of RL Algorithms



# Learning and Planning

- Two fundamental problems in **sequential decision making**
  - **Reinforcement Learning**
    - The environment is initially unknown
    - The agent interacts with environment
    - The agent improves its policy
  - **Planning**
    - A model of the environment is known
    - The agent performs computations with its model (without any external interaction)
    - The agent improves its policy
    - a.k.a deliberation, reasoning, introspection, pondering, thought, search

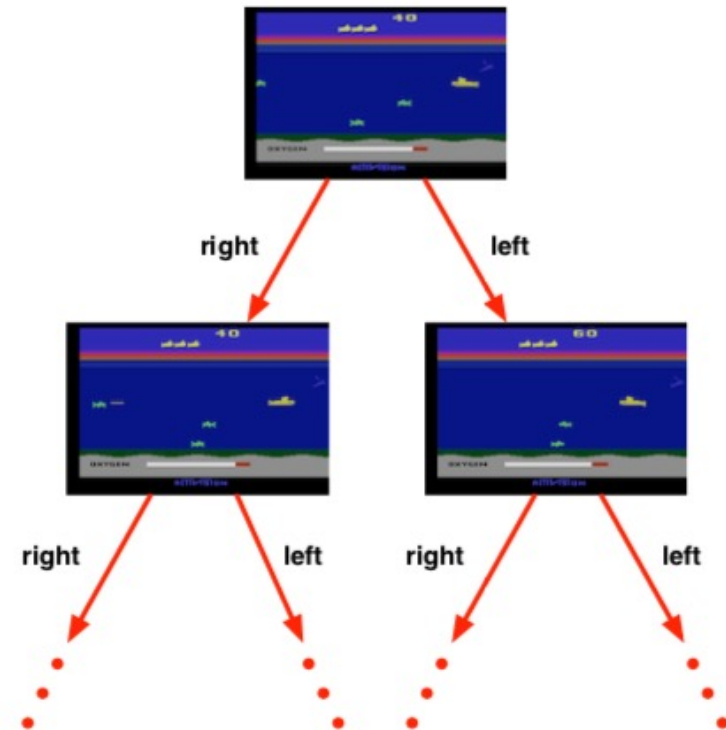
# Atari Example: Reinforcement Learning



- Rules of the game are unknown
- Learn directly from interactive game-play
- Pick actions on joystick, see pixels and scores

# Atari Example: Planning

- Rules of the game are known
- Can query emulator
  - perfect model inside agent's brain
- If I take action  $a$  from state  $s$ :
  - what would the next state be?
  - what would the score be?
- Plan ahead to find optimal policy
  - e.g. tree search



# Exploration and Exploitation

- Reinforcement learning is like **trial-and-error** learning
- The agent should discover a good **policy**
- From its **experiences** of the environment
- Without losing too much **reward** along the way
- **Exploration** finds more information about the environment
- **Exploitation** exploits known information to maximise reward
- It is usually important to explore as well as exploit

# Exploration and Exploitation

## Examples

- Restaurant Selection
  - Exploitation: Go to your favorite restaurant
  - Exploration: Try a new restaurant
- Online Banner Advertisements
  - Exploitation: Show the most successful advert
  - Exploration: Show a different advert



# Exploration and Exploitation

## Examples

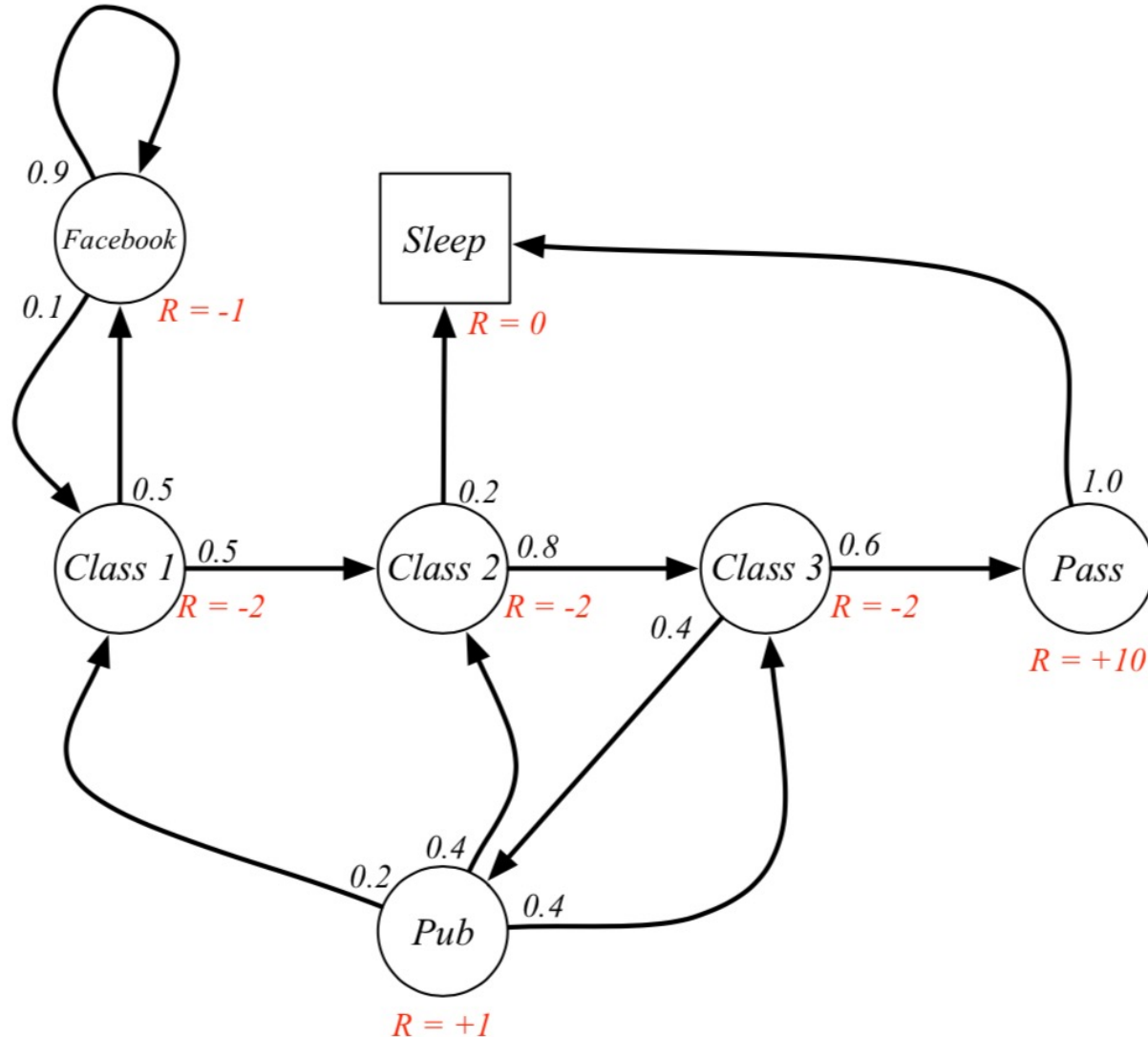
- Oil Drilling
  - Exploitation: Drill at the best known location
  - Exploration: Drill at a new location
- Game Playing
  - Exploitation: Play the move you believe is best
  - Exploration: Play an experimental move

# Prediction and Control

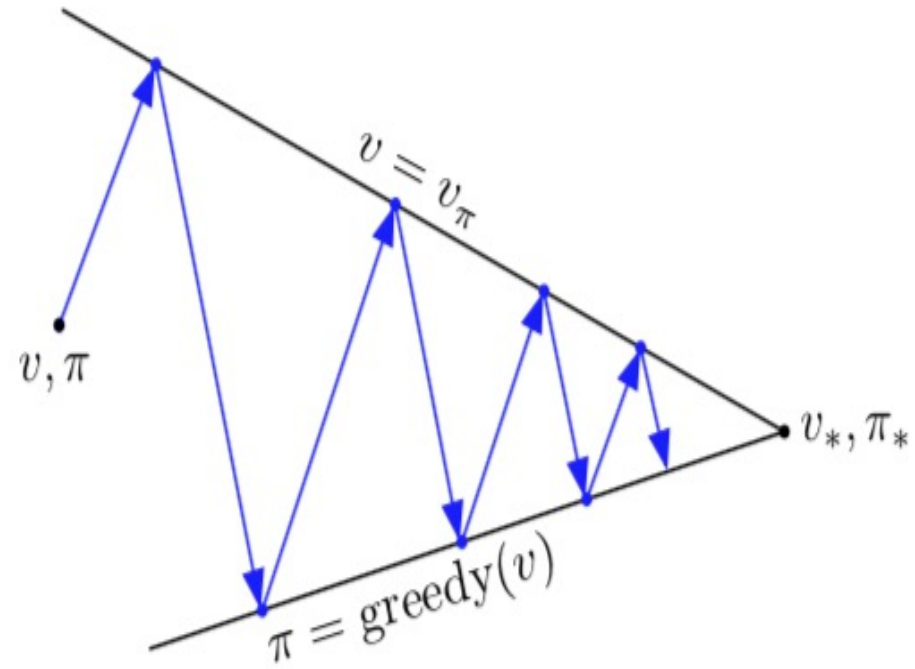
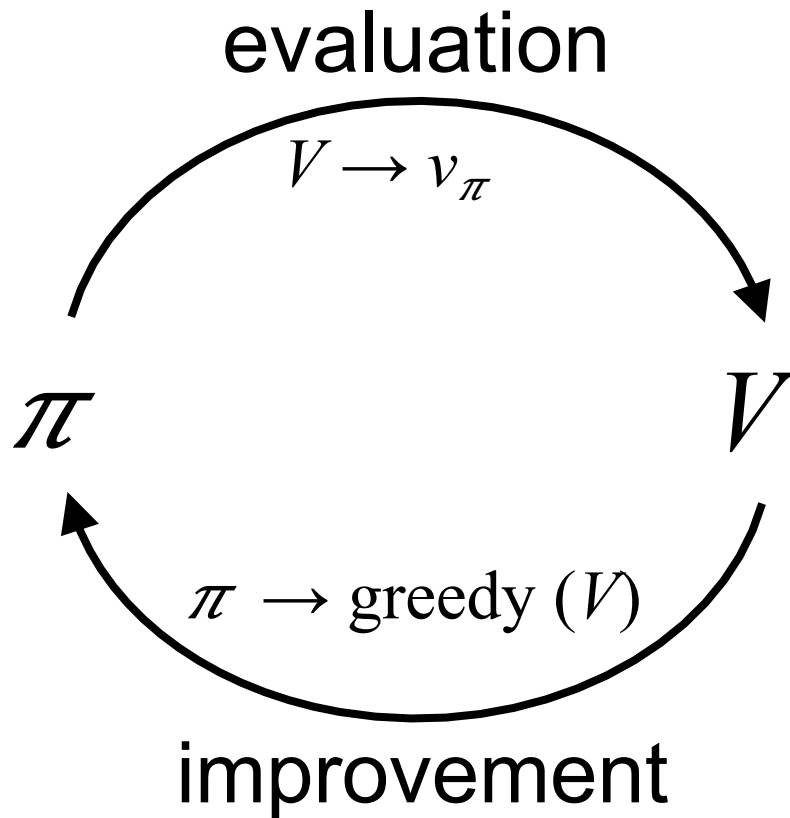
- Prediction: evaluate the future
  - Given a policy
- Control: optimize the future
  - Find the best policy

# Markov Decision Processes (MDP)

## Example: Student MDP



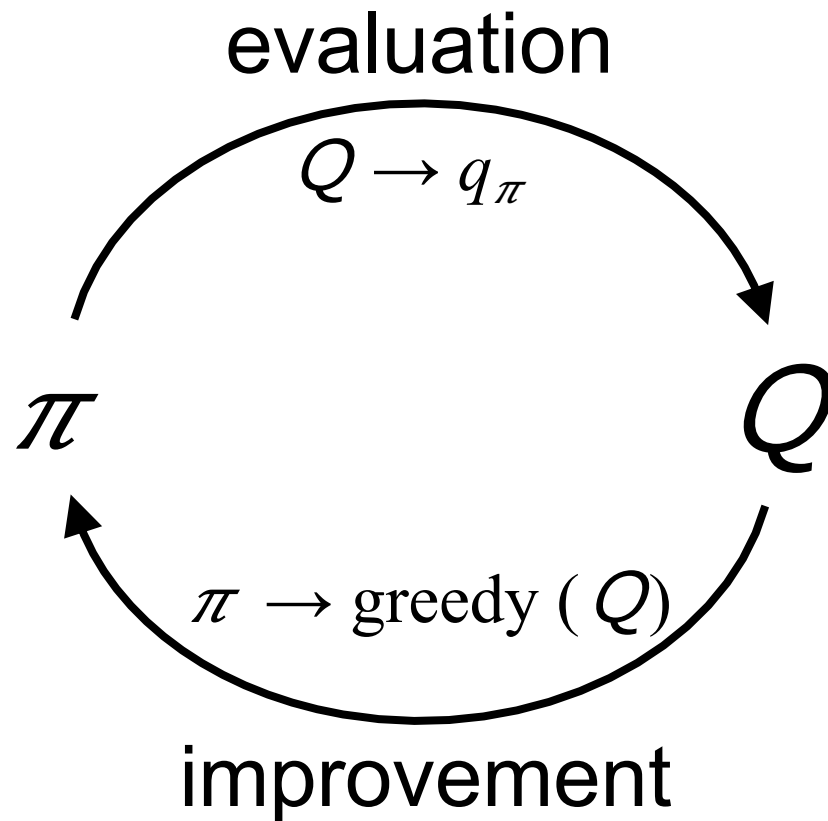
# Generalized Policy Iteration (GPI)



$$\pi_* \rightleftarrows v_*$$

# Generalized Policy Iteration (GPI)

Any iteration of **policy evaluation** and **policy improvement**, independent of their granularity.



# Temporal-Difference (TD) Learning

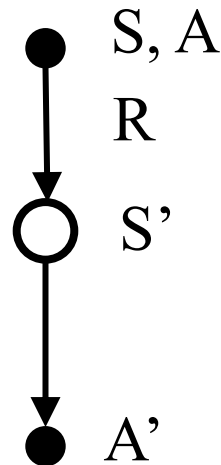
- Sarsa: On-policy TD Control
- Q-learning: Off-policy TD Control

# SARSA

(state-action-reward-state-action)

On-policy TD Control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

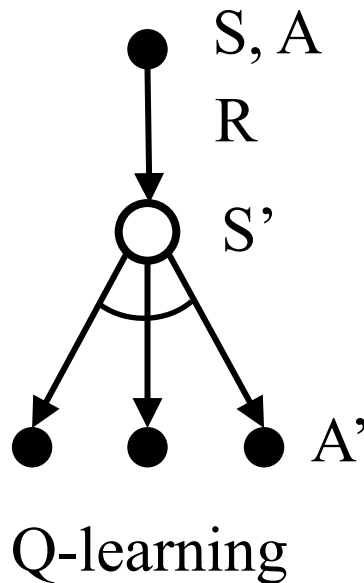


SARSA

# Q-learning (Watkins, 1989)

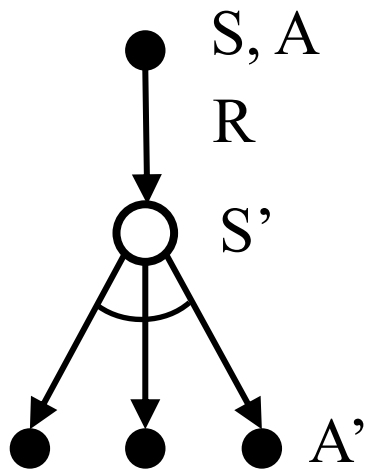
## Off-policy TD Control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

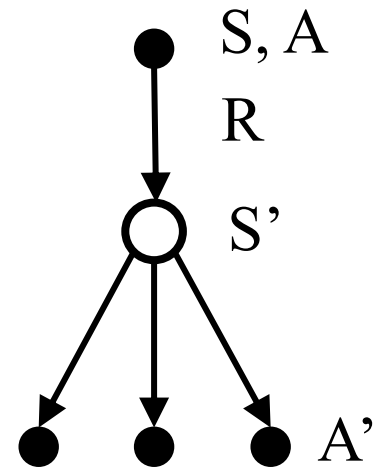




# Q-learning and Expected SARSA

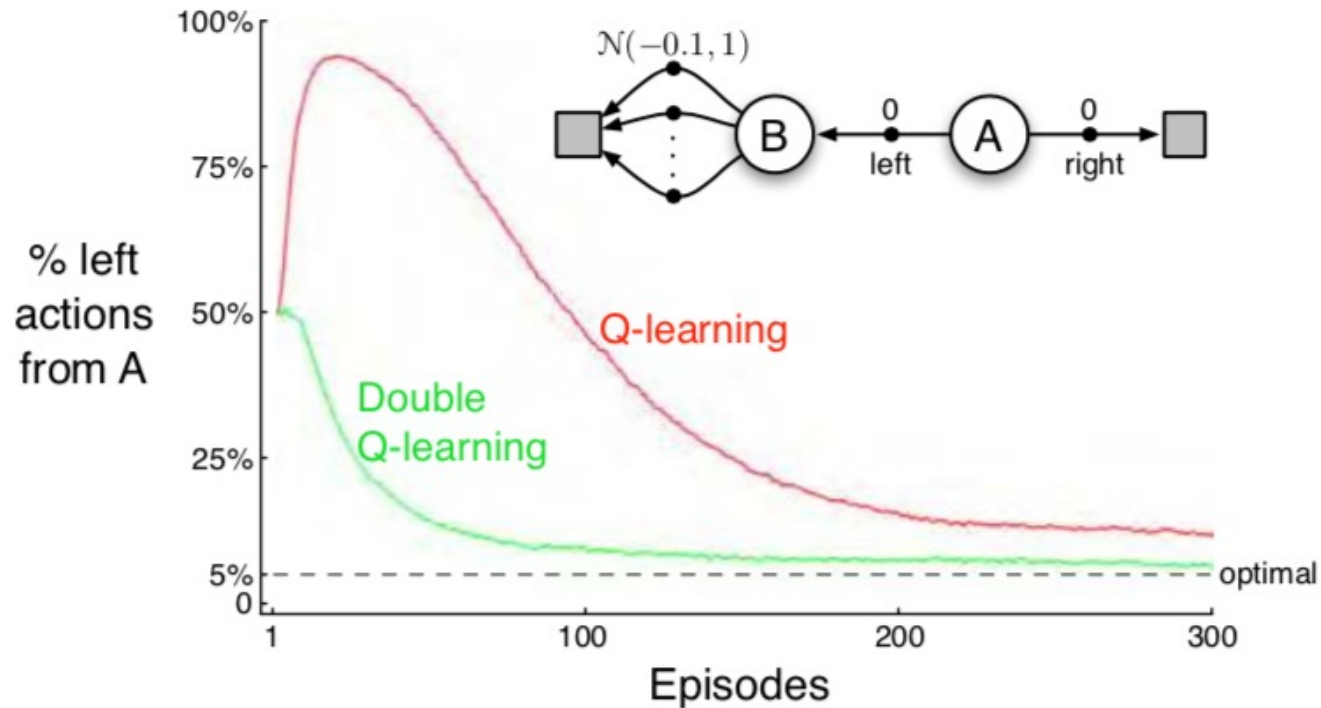


Q-learning



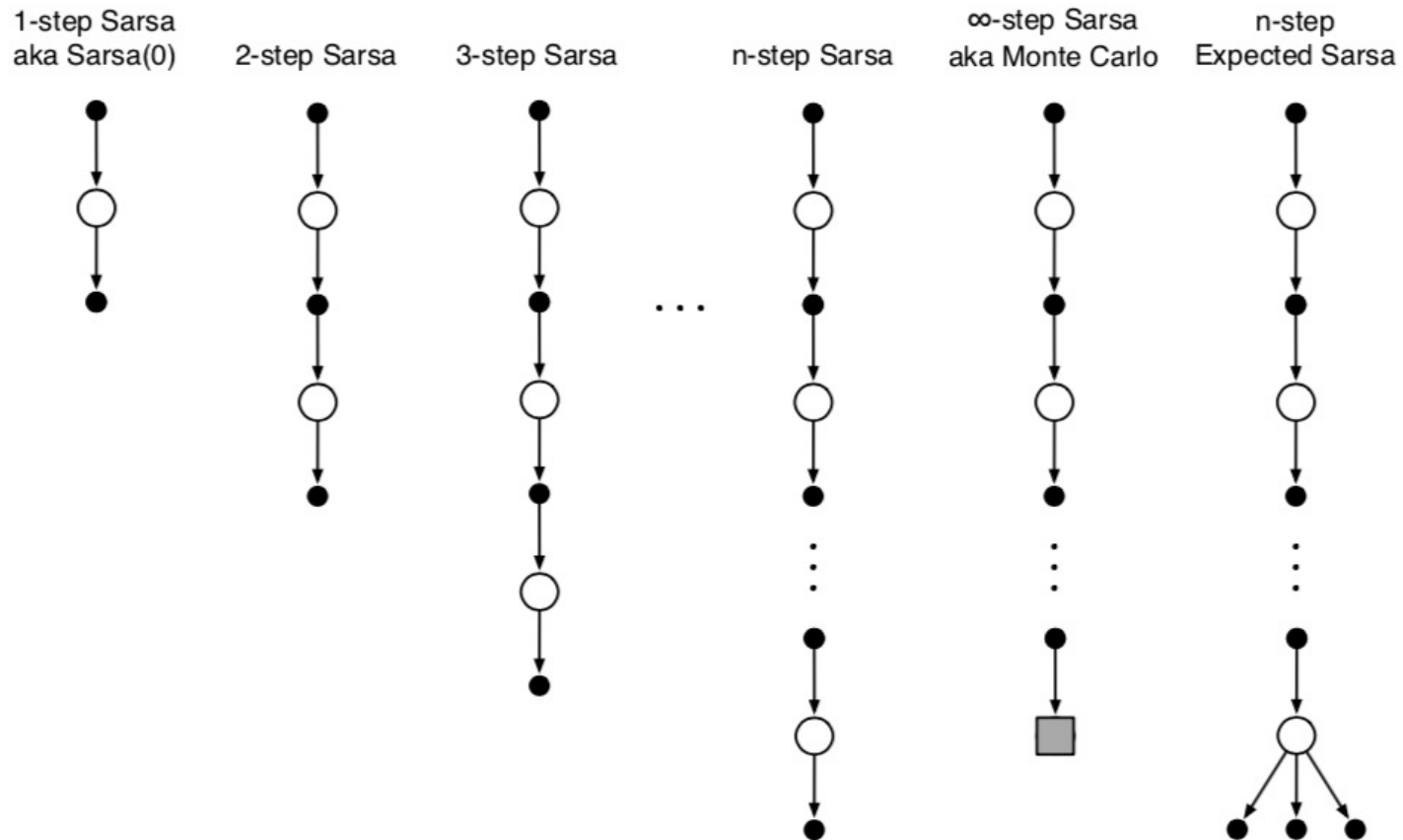
Expected SARSA

# Q-learning and Double Q-learning



**Figure 6.5:** Comparison of Q-learning and Double Q-learning on a simple episodic MDP (shown inset). Q-learning initially learns to take the left action much more often than the right action, and always takes it significantly more often than the 5% minimum probability enforced by  $\epsilon$ -greedy action selection with  $\epsilon = 0.1$ . In contrast, Double Q-learning is essentially unaffected by maximization bias. These data are averaged over 10,000 runs. The initial action-value estimates were zero. Any ties in  $\epsilon$ -greedy action selection were broken randomly.

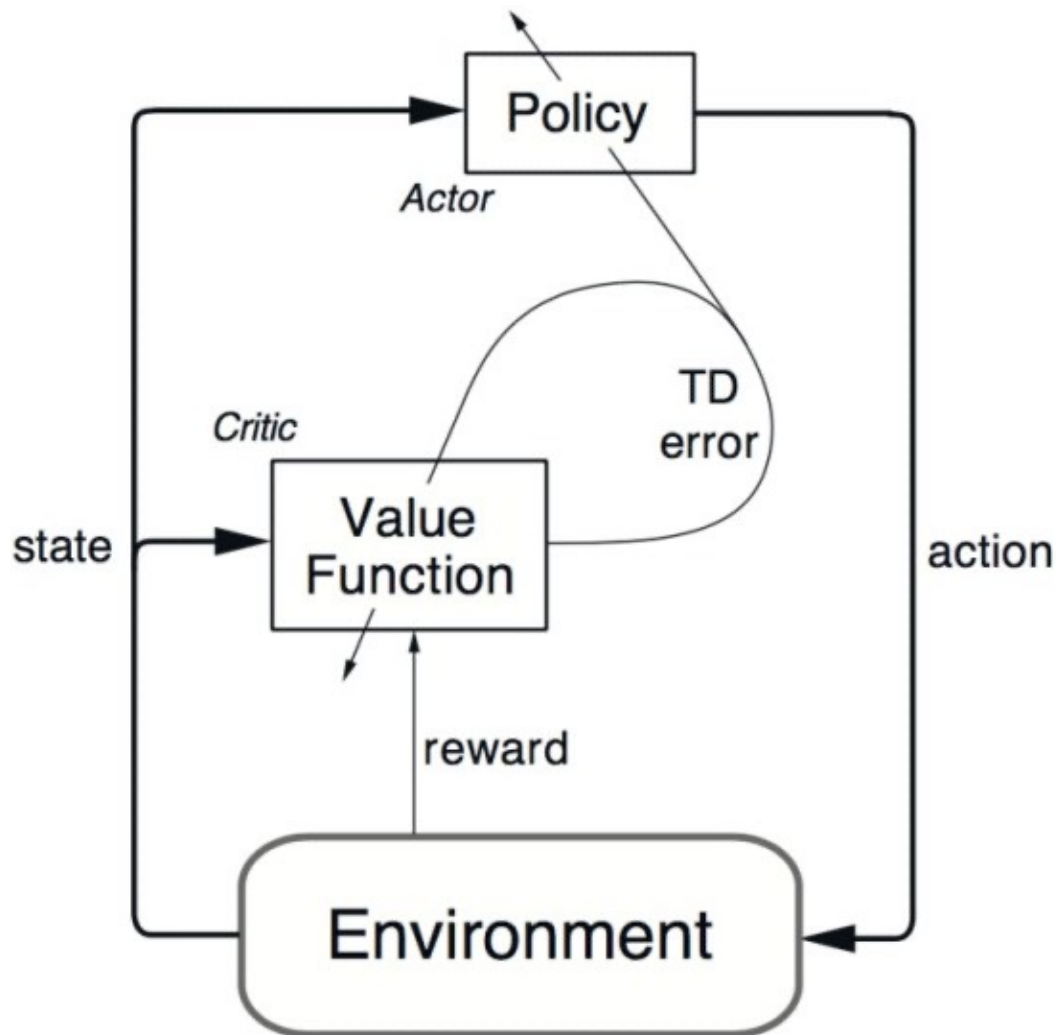
# n-step methods for state-action value



**Figure 7.3:** The backup diagrams for the spectrum of  $n$ -step methods for state-action values. They range from the one-step update of Sarsa(0) to the up-until-termination update of the Monte Carlo method. In between are the  $n$ -step updates, based on  $n$  steps of real rewards and the estimated value of the  $n$ th next state-action pair, all appropriately discounted. On the far right is the backup diagram for  $n$ -step Expected Sarsa.

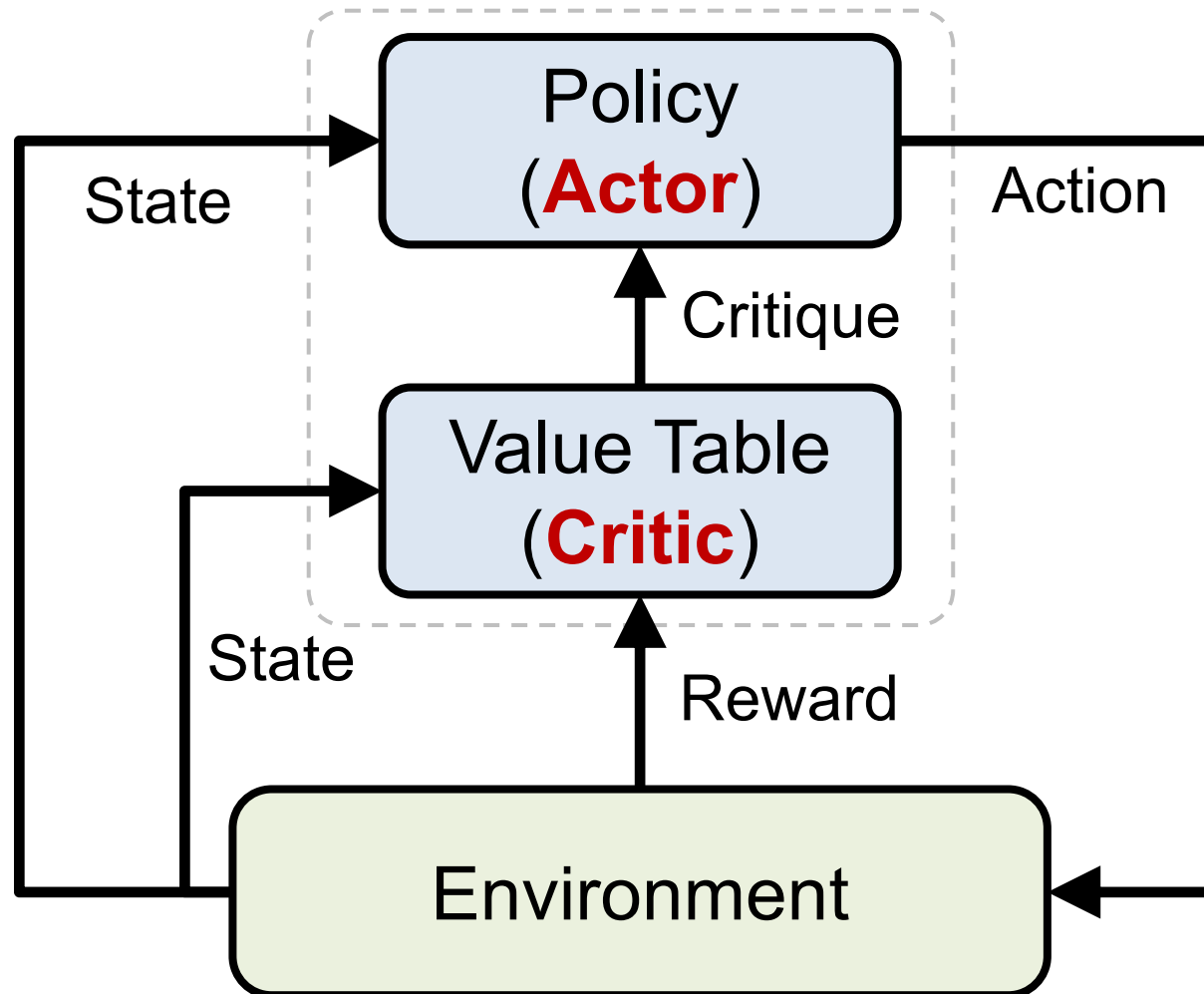
# Reinforcement Learning

## Actor-Critic (AC) Architecture

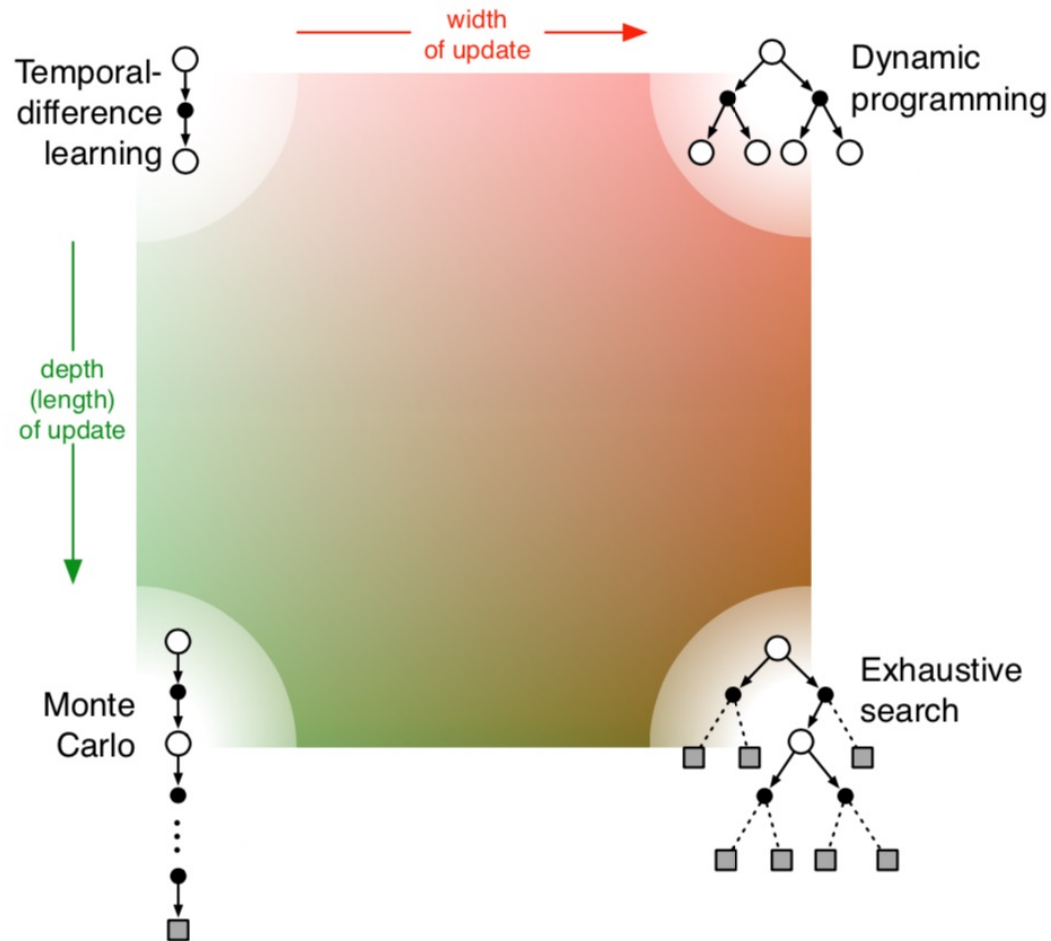


# Reinforcement Learning

## Actor-Critic (AC) Learning Methods

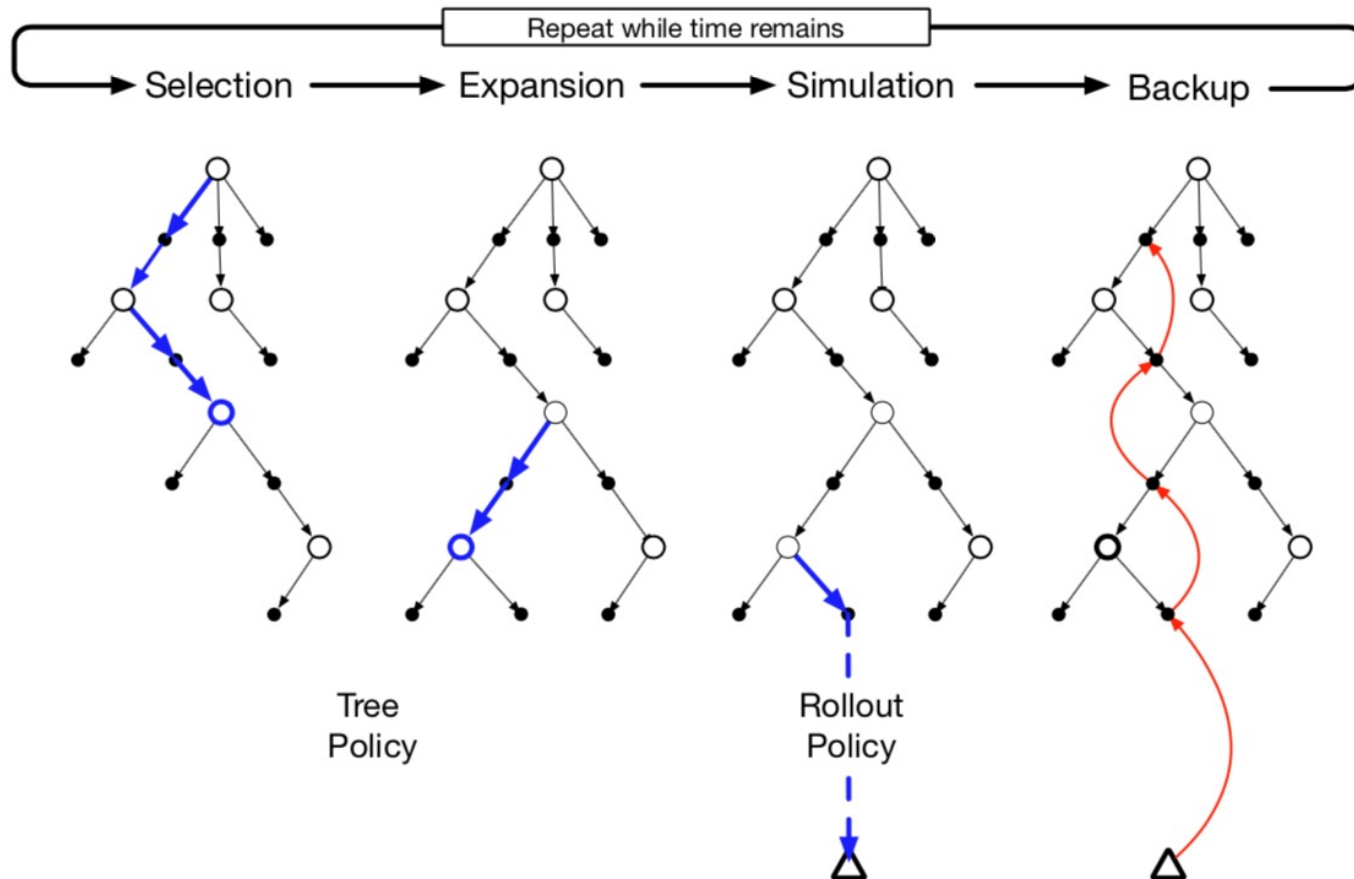


# Reinforcement Learning Methods



**Figure 8.11:** A slice through the space of reinforcement learning methods, highlighting the two of the most important dimensions explored in Part I of this book: the depth and width of the updates.

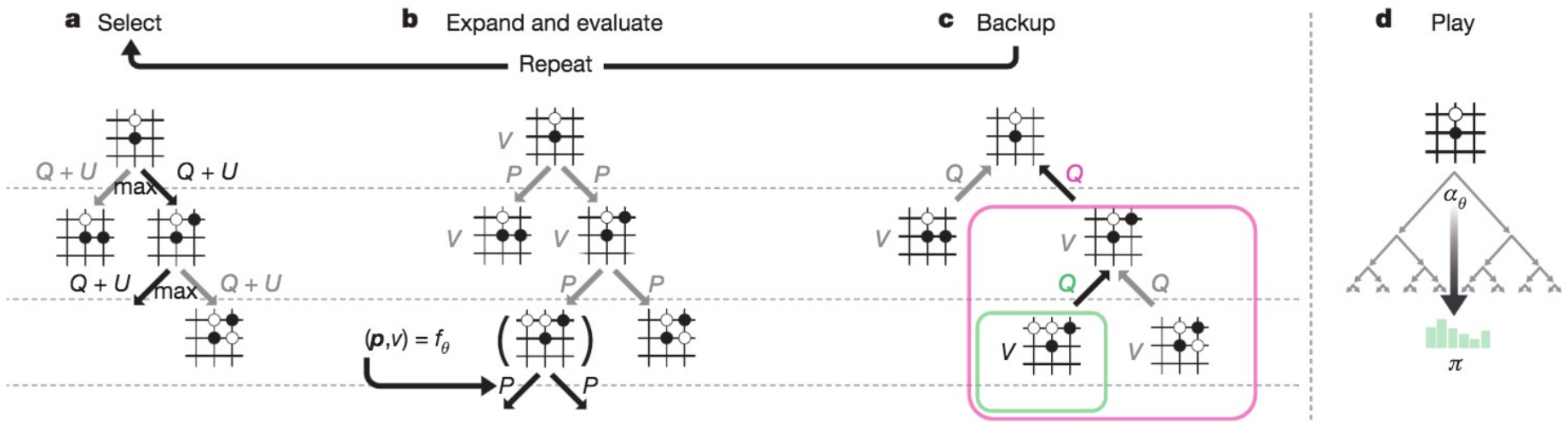
# Monte Carlo Tree Search (MCTS)



**Figure 8.10:** Monte Carlo Tree Search. When the environment changes to a new state, MCTS executes as many iterations as possible before an action needs to be selected, incrementally building a tree whose root node represents the current state. Each iteration consists of the four operations **Selection**, **Expansion** (though possibly skipped on some iterations), **Simulation**, and **Backup**, as explained in the text and illustrated by the bold arrows in the trees. Adapted from Chaslot, Bakkes, Szita, and Spronck (2008).

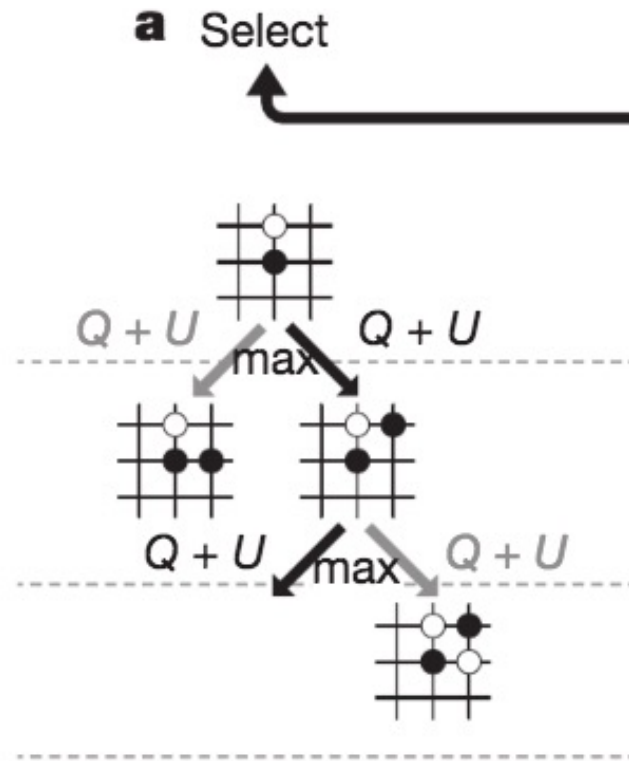
# Monte Carlo Tree Search (MCTS)

## MCTS in AlphaGo Zero





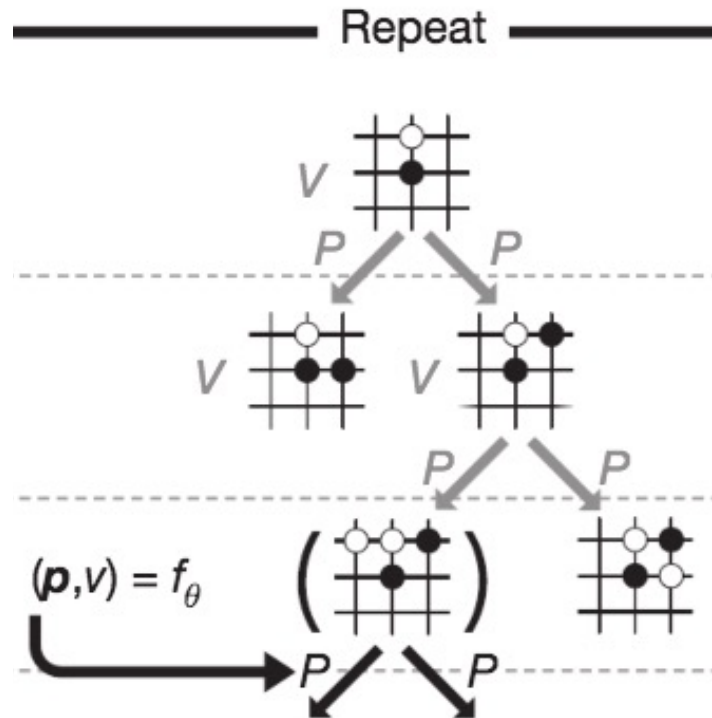
# MCTS in AlphaGo Zero



**a:** Each simulation traverses the tree by selecting the edge with maximum action value  $Q$ , plus an upper confidence bound  $U$  that depends on a stored prior probability  $P$  and visit count  $N$  for that edge (which is incremented once traversed).

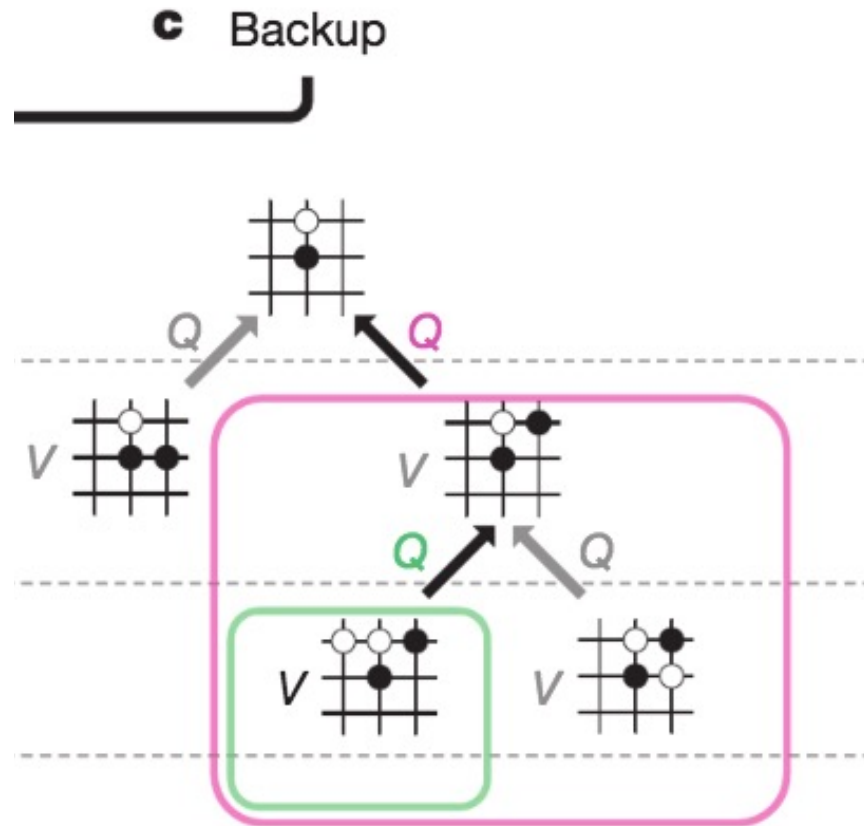
# MCTS in AlphaGo Zero

## **b** Expand and evaluate



**b:** The leaf node is expanded and the associated position  $s$  is evaluated by the neural network  $(P(s, \cdot), V(s)) = f_{\theta}(s)$ ; the vector of  $P$  values are stored in the outgoing edges from  $s$ .

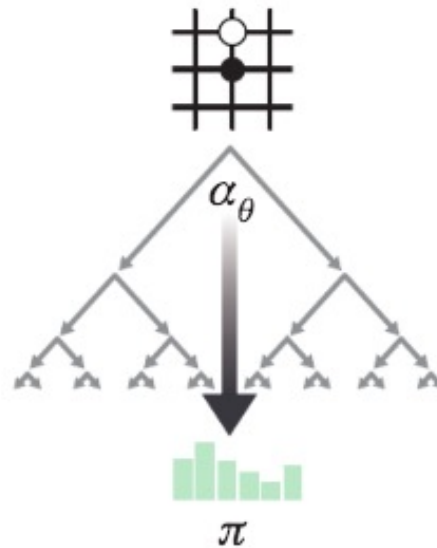
# MCTS in AlphaGo Zero



**c:** Action value  $Q$  is updated to track the mean of all evaluations  $V$  in the subtree below that action

# MCTS in AlphaGo Zero

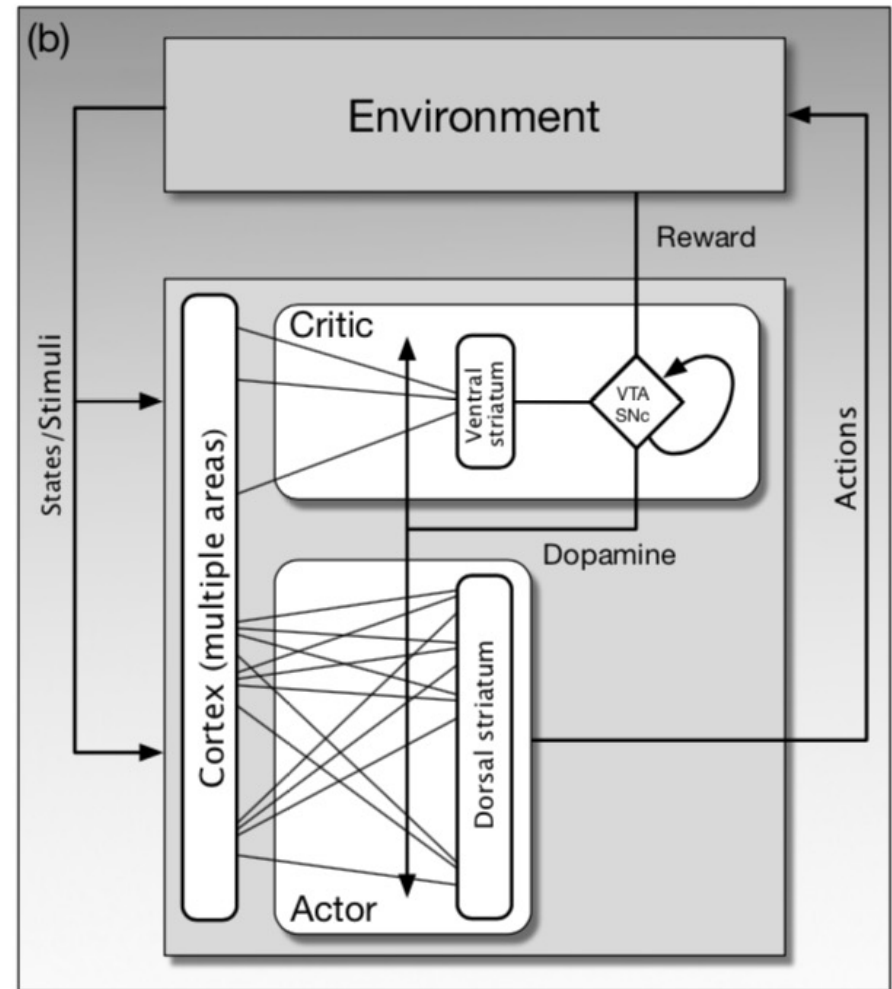
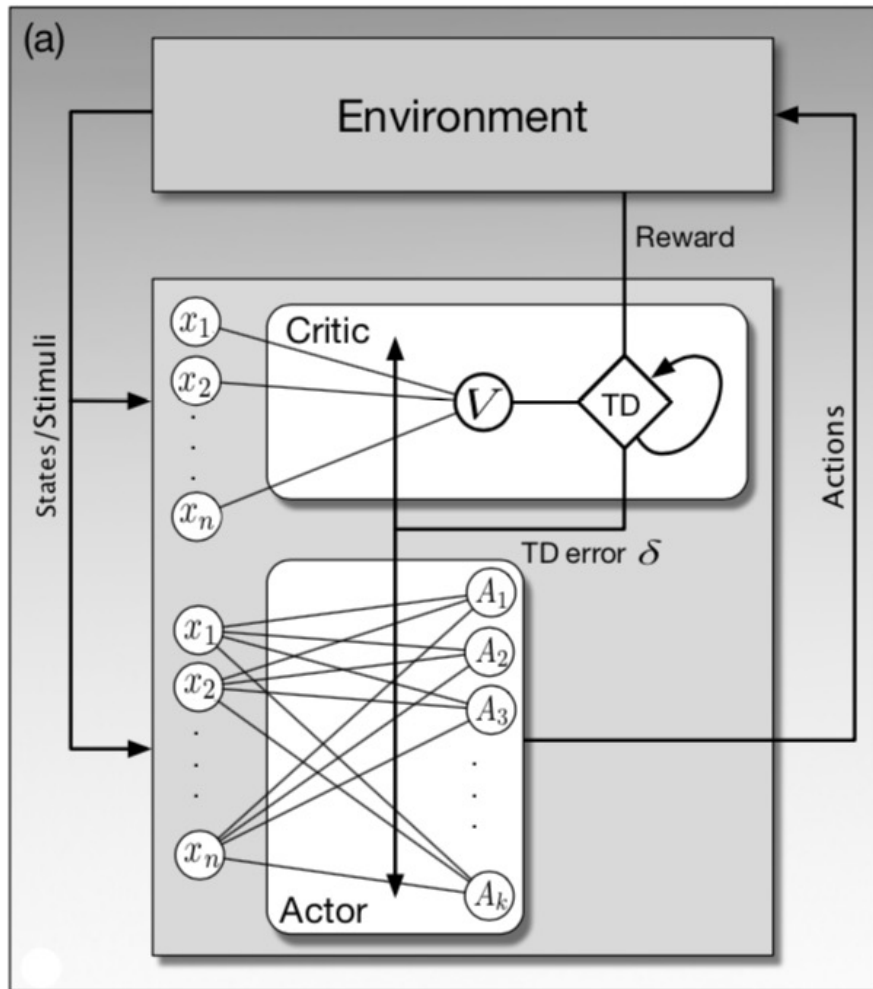
**d** Play



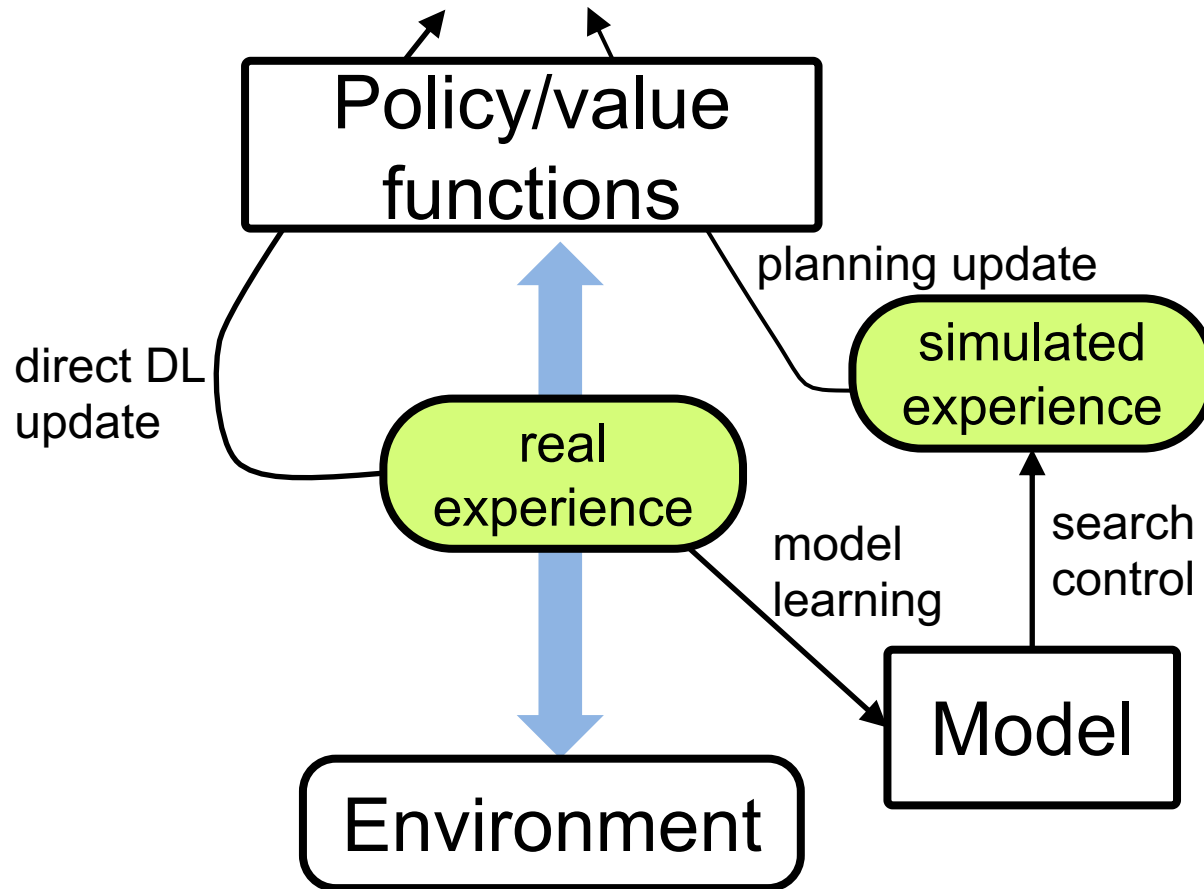
**d:** Once the search is complete, search probabilities  $\pi$  are returned, proportional to  $N^{1/\tau}$ , where  $N$  is the visit count of each move from the root state and  $\tau$  is a parameter controlling temperature.

# Reinforcement Learning

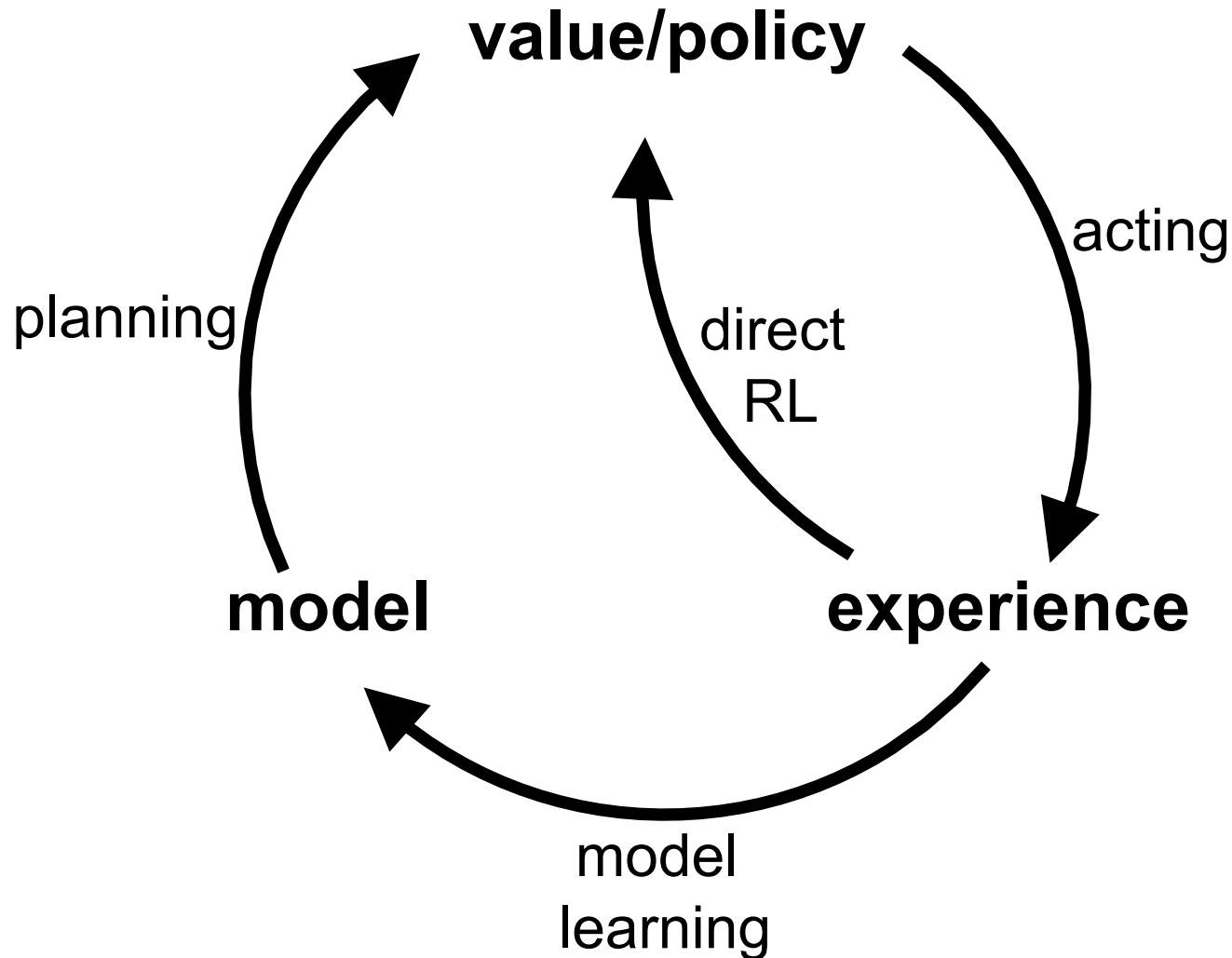
## Actor Critic ANN



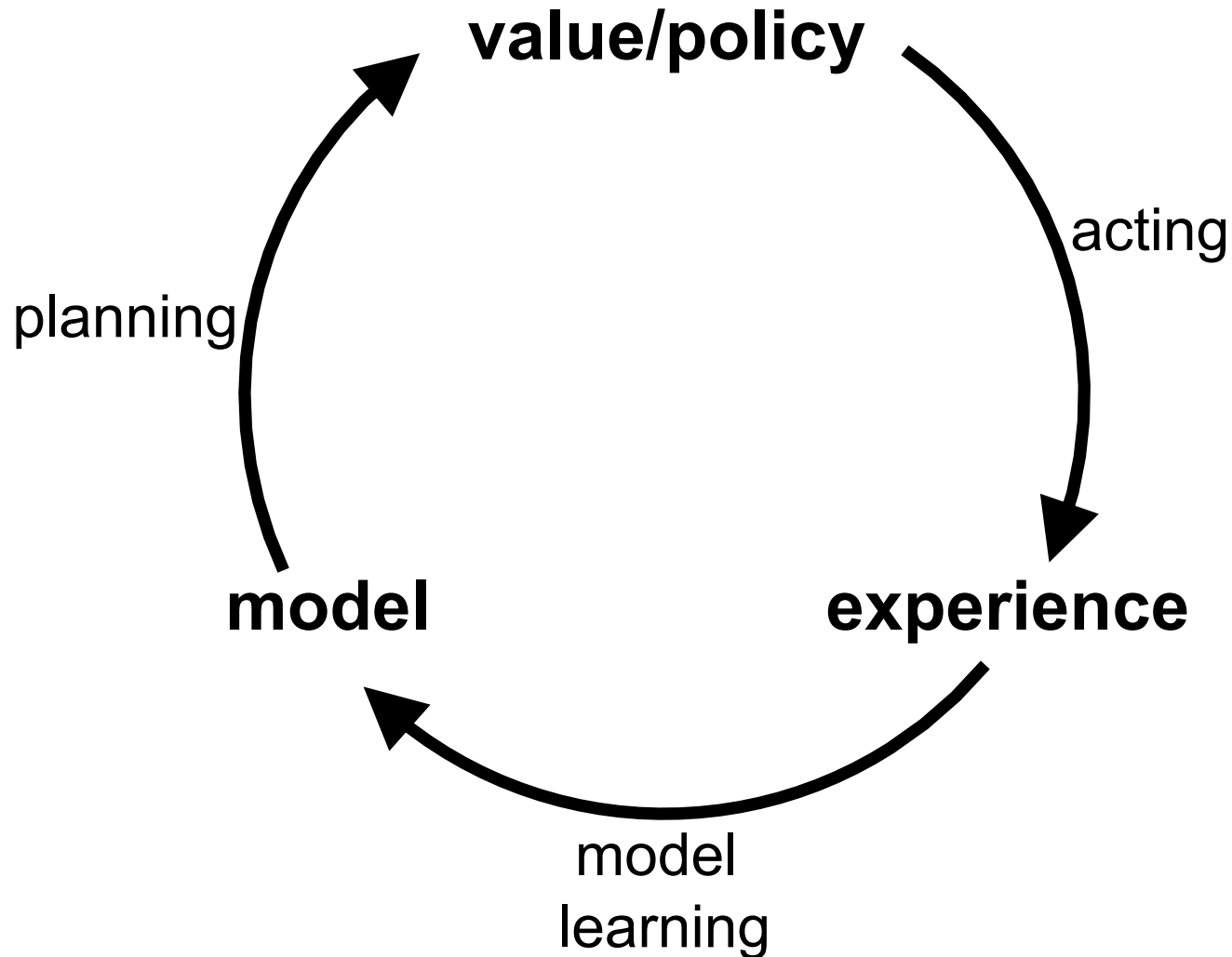
# Reinforcement Learning General Dyna Architecture



# Dyna: Integrated Planning, Acting, and Learning

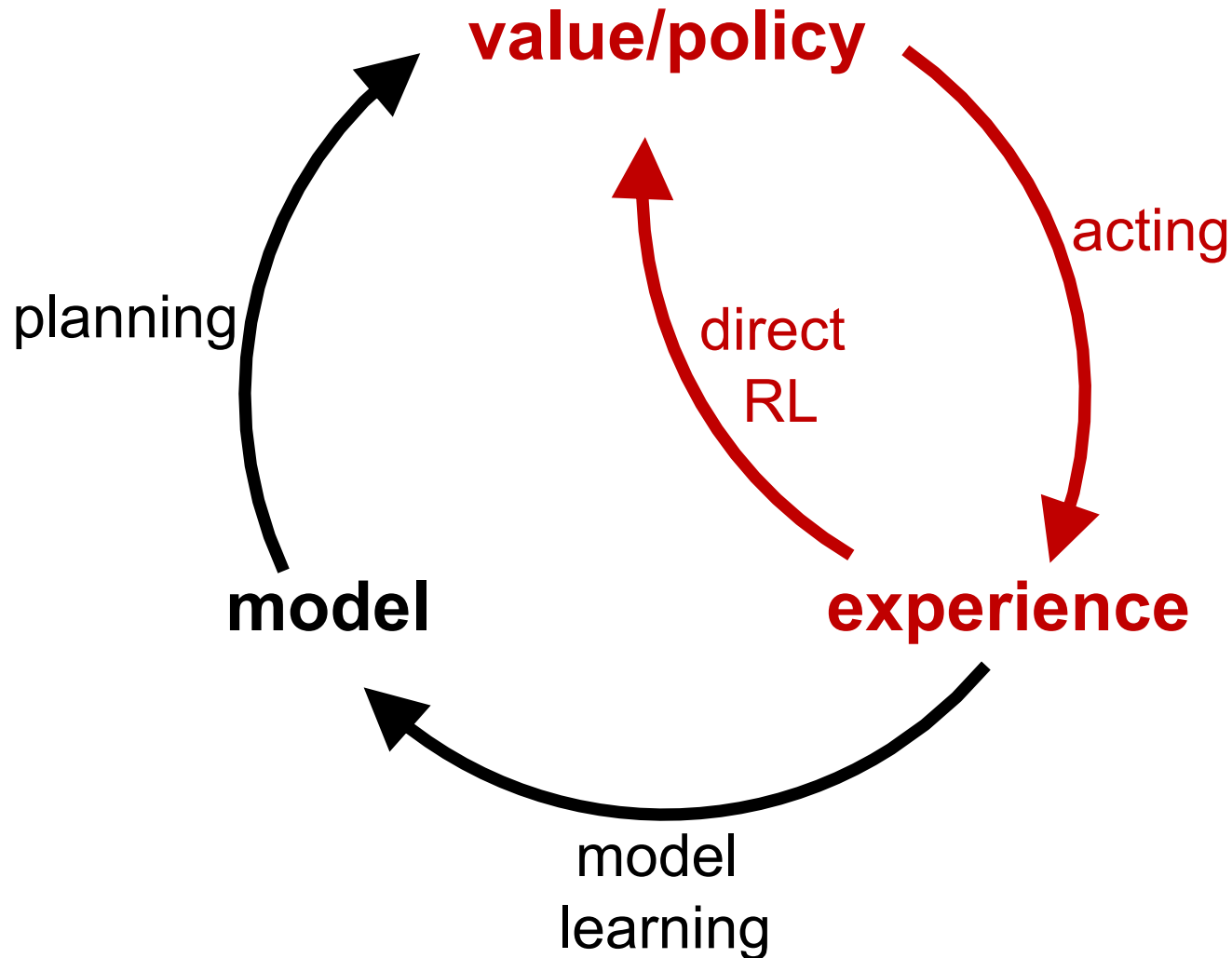


# Model-Based RL

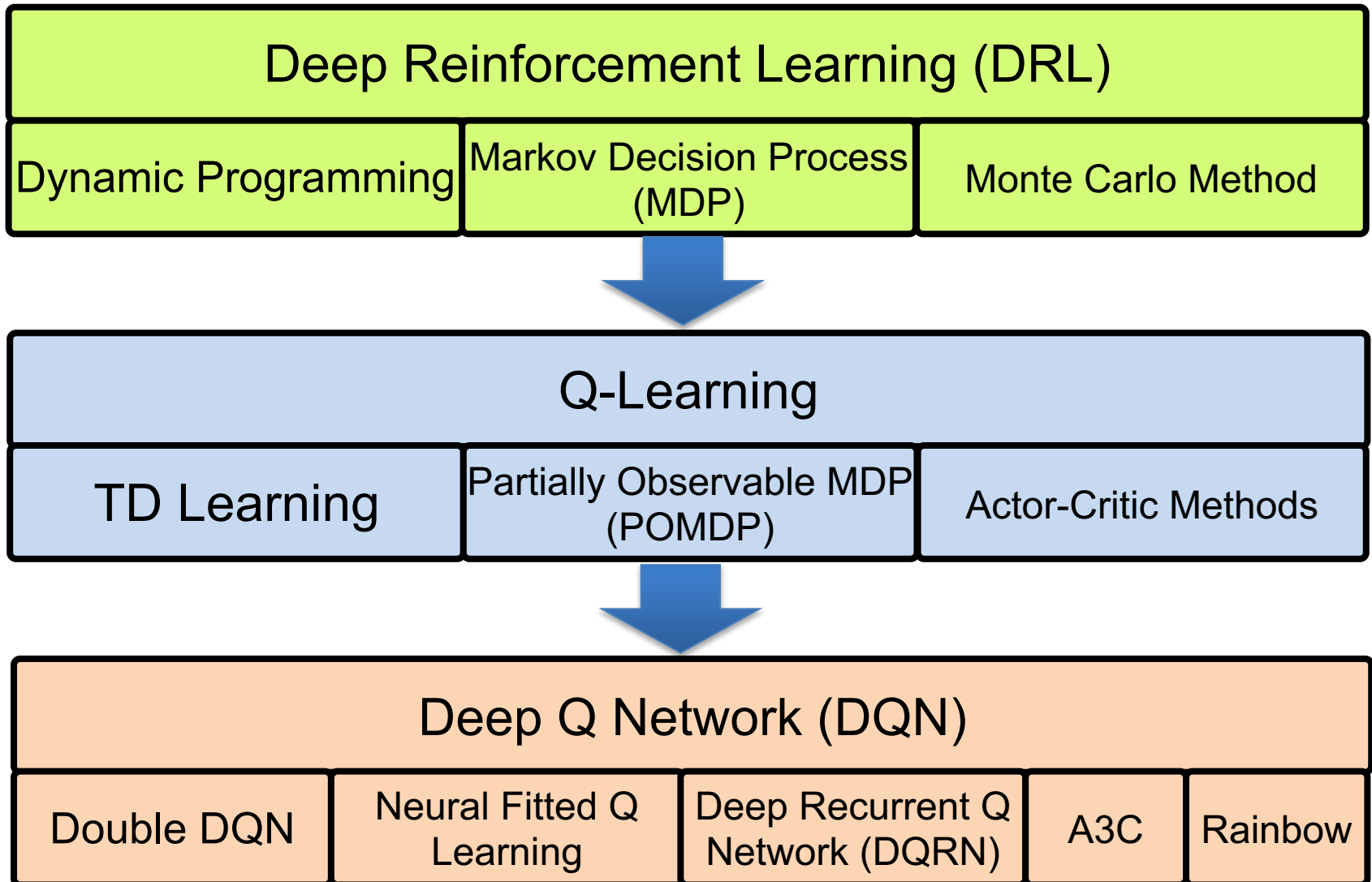




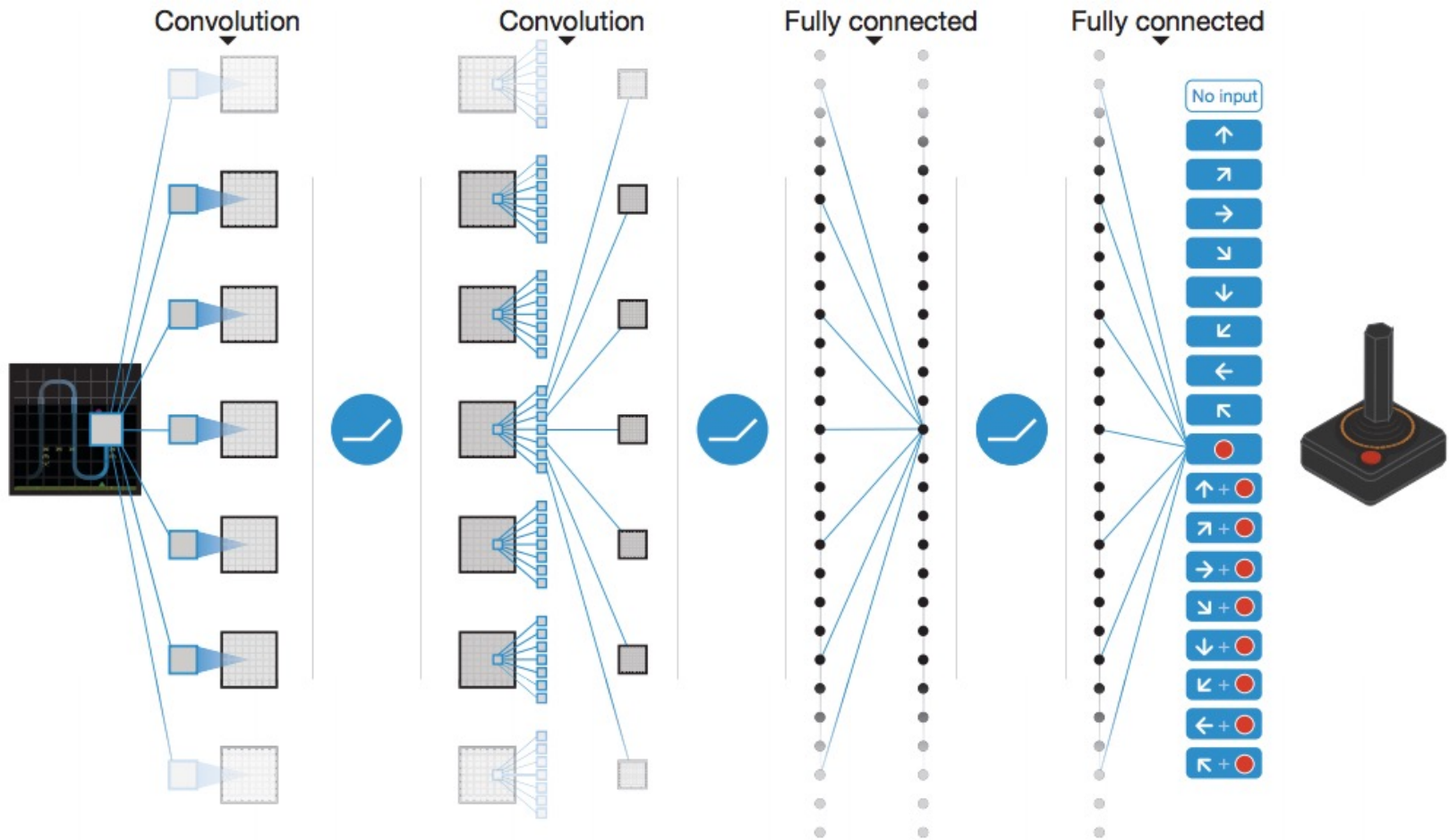
# Model-Free RL (DQN, A3C)



# Reinforcement Learning Algorithms

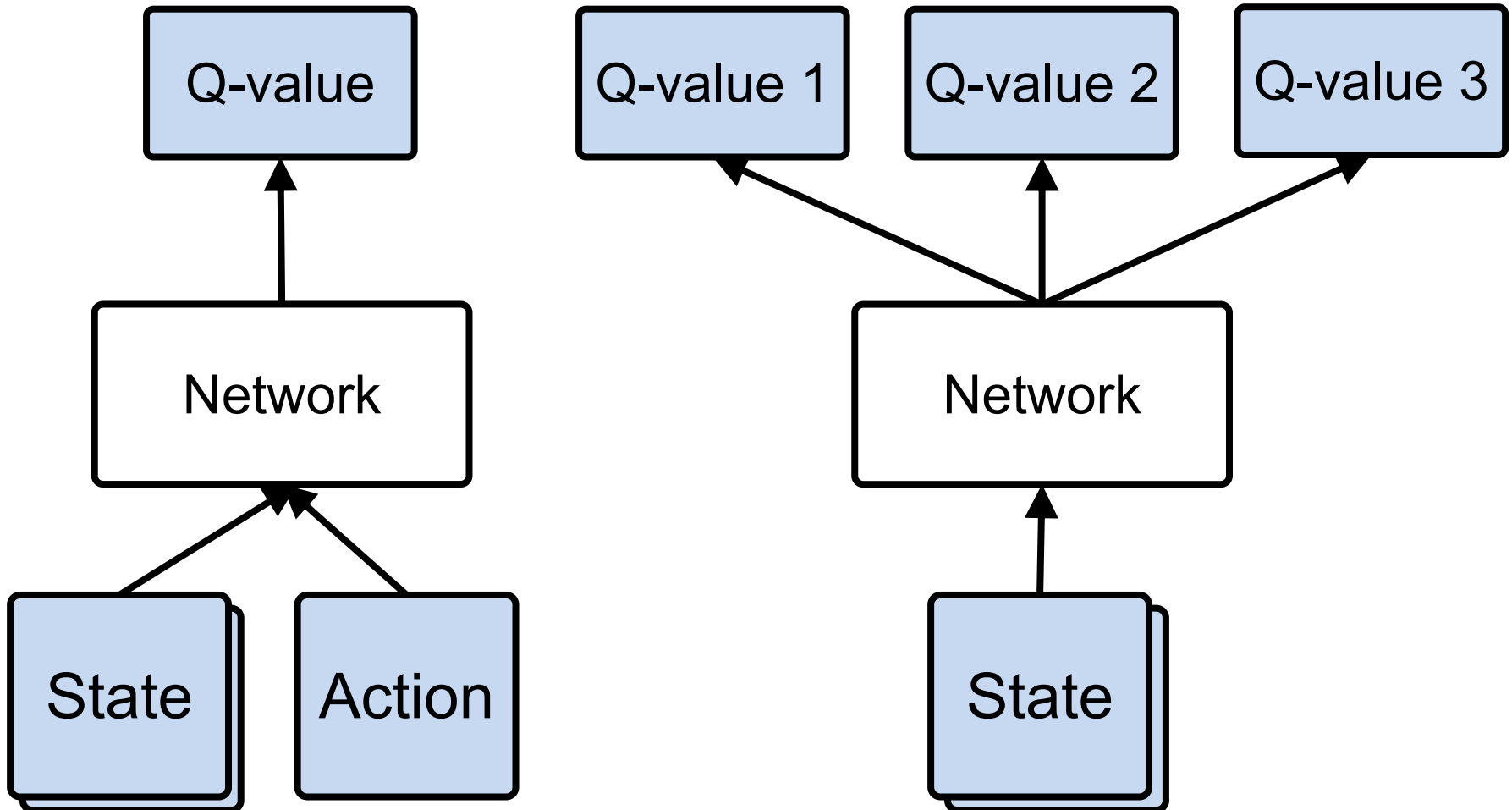


# Human-level control through deep reinforcement learning (DQN)

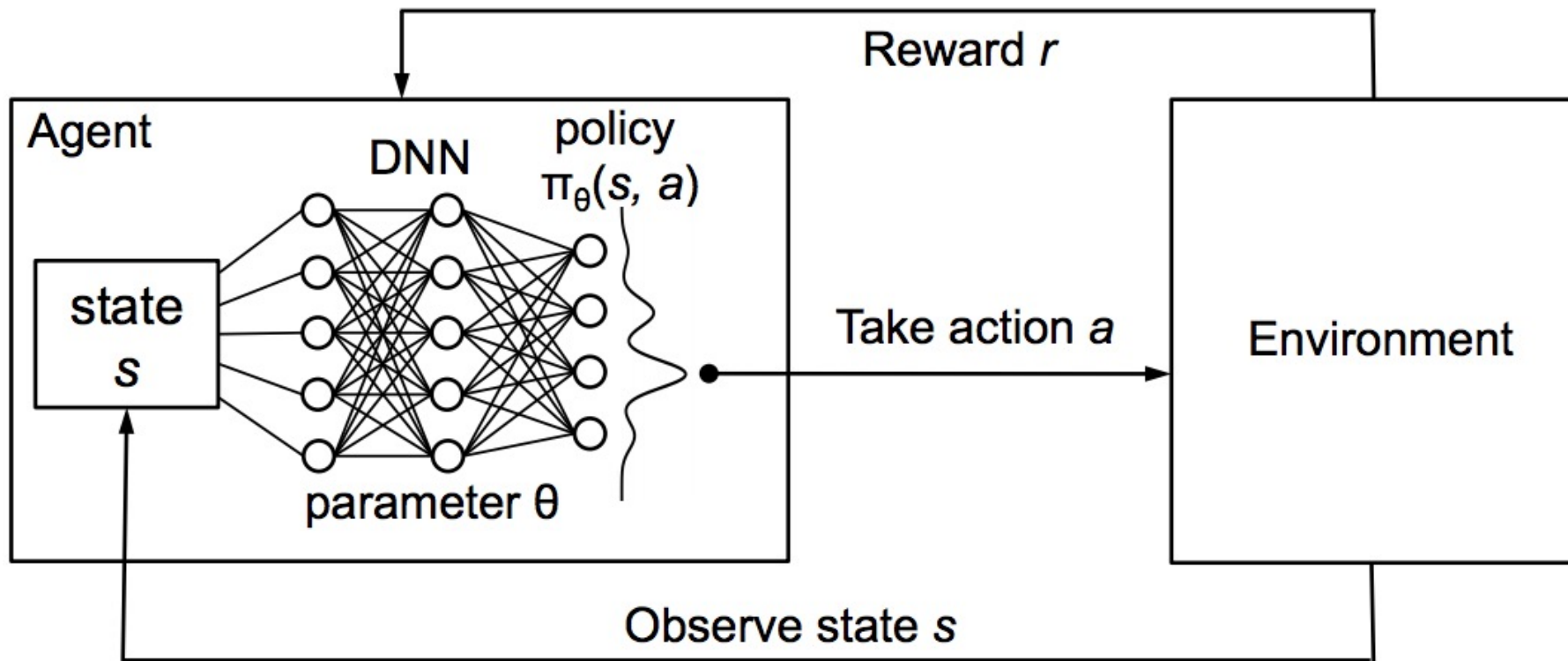


Schematic illustration of the convolutional neural network

# Deep Q-Network (DQN)

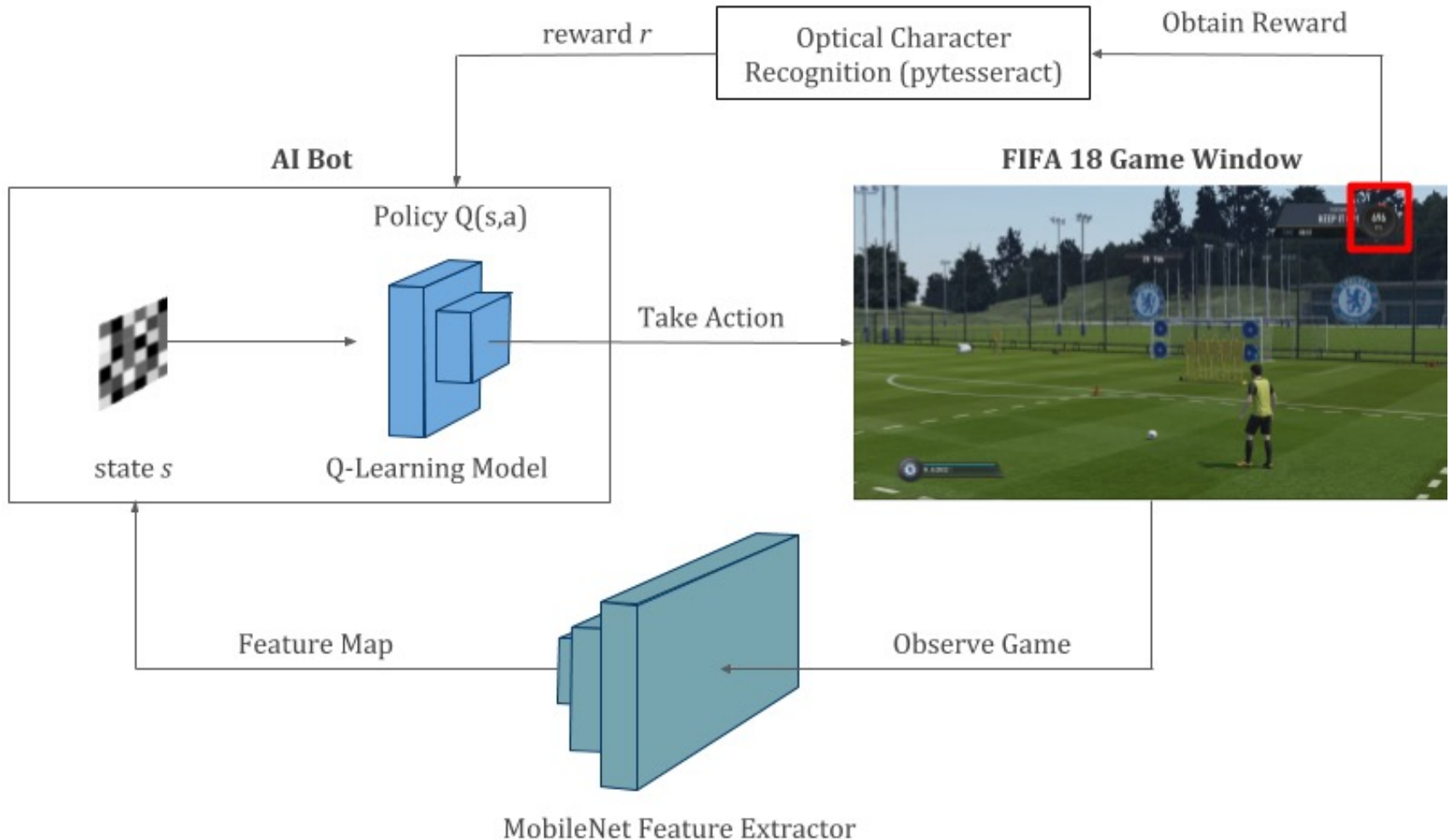


# Reinforcement Learning with policy represented via DNN

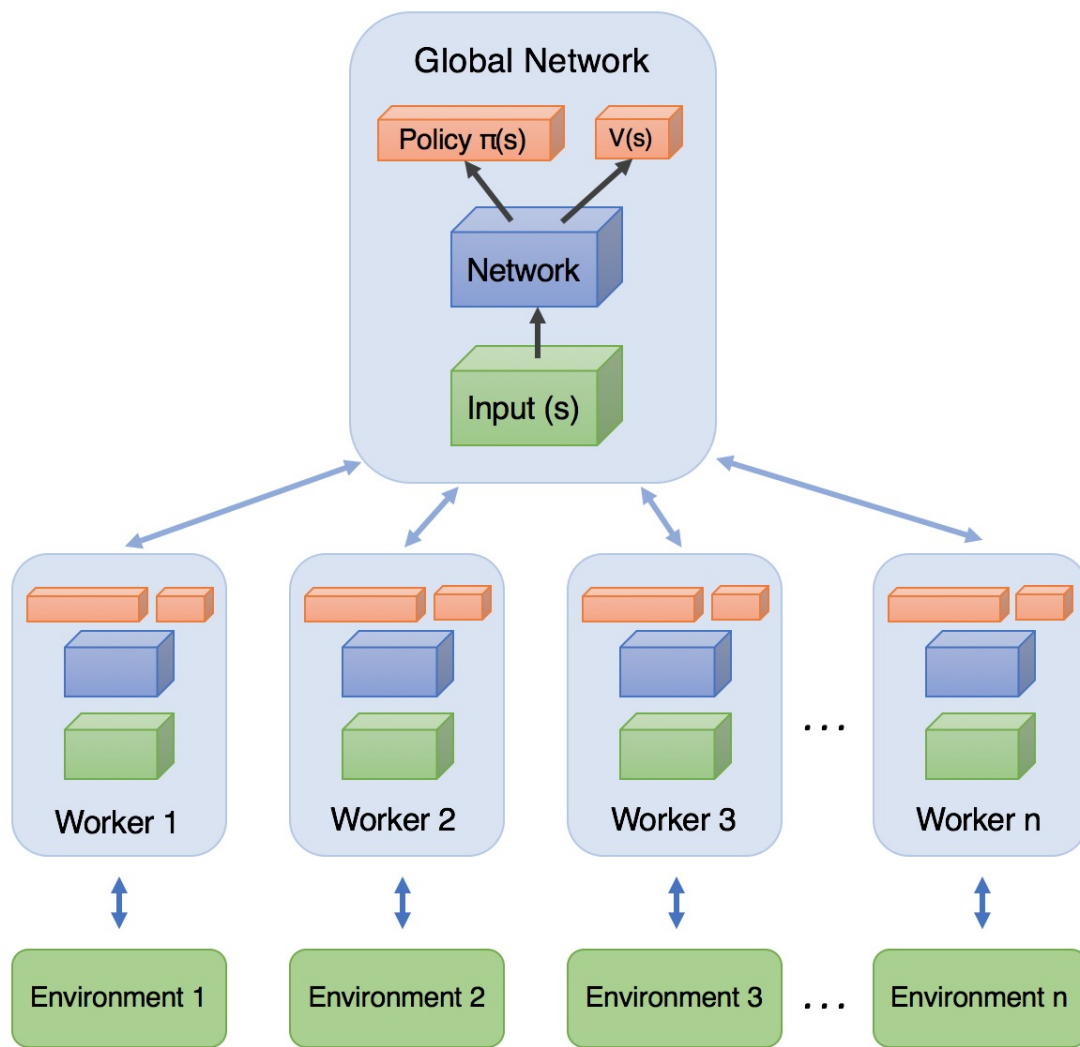


# Reinforcement Learning

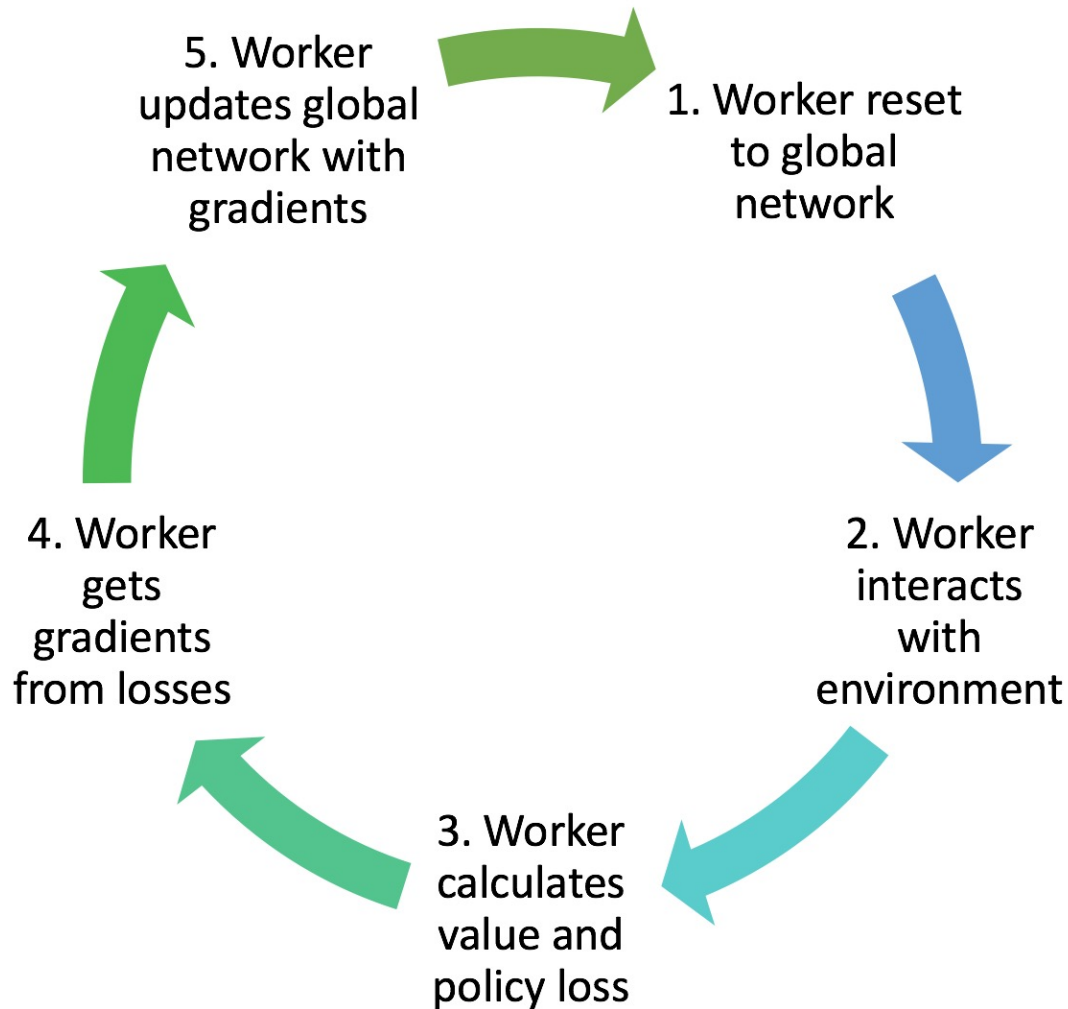
## Deep Q-Learning in FIFA 18



# Asynchronous Advantage Actor-Critic (A3C)



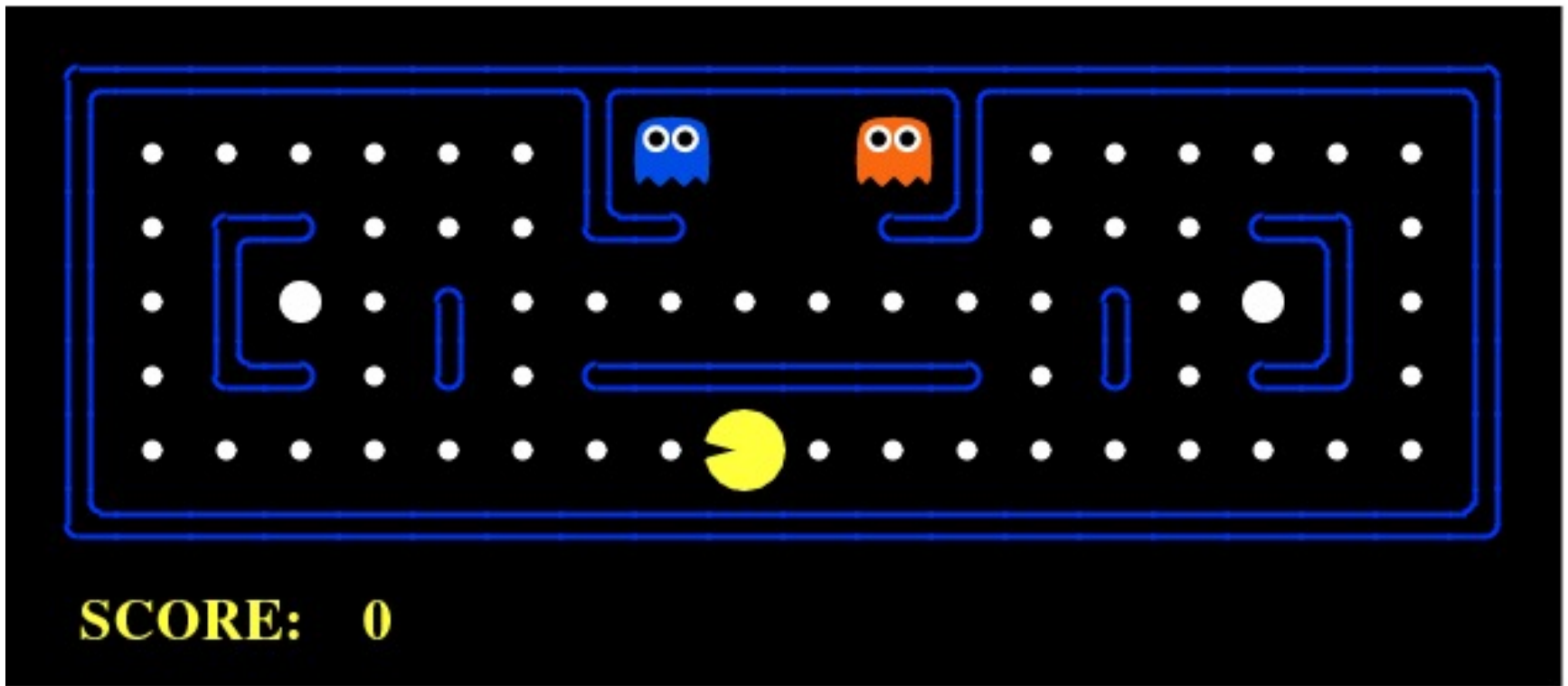
# Training workflow of each worker agent in A3C





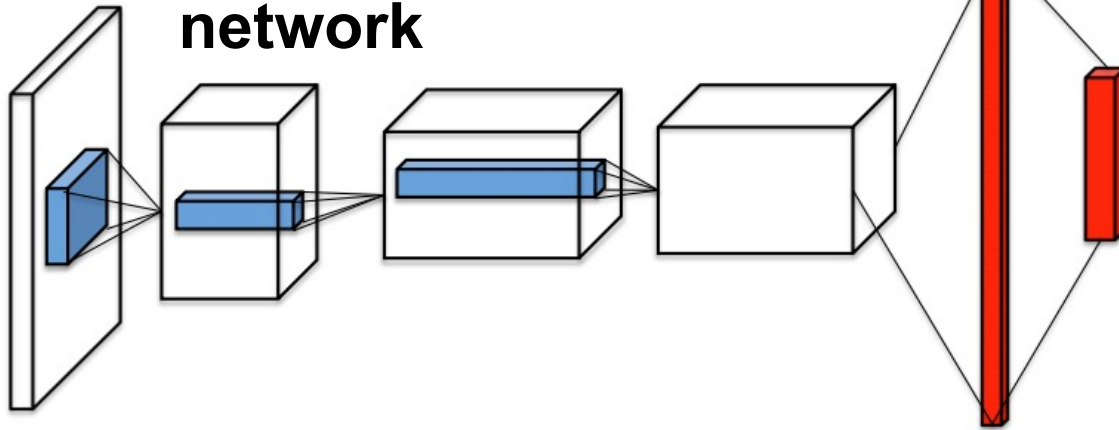
# Reinforcement Learning

## Example: PCMAN

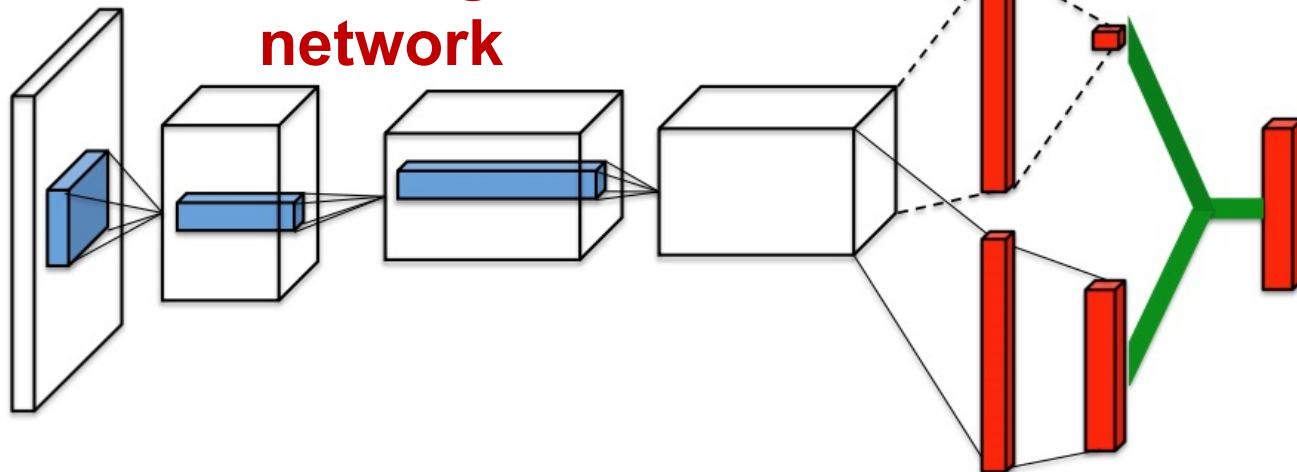


# Dueling Network Architectures for Deep Reinforcement Learning

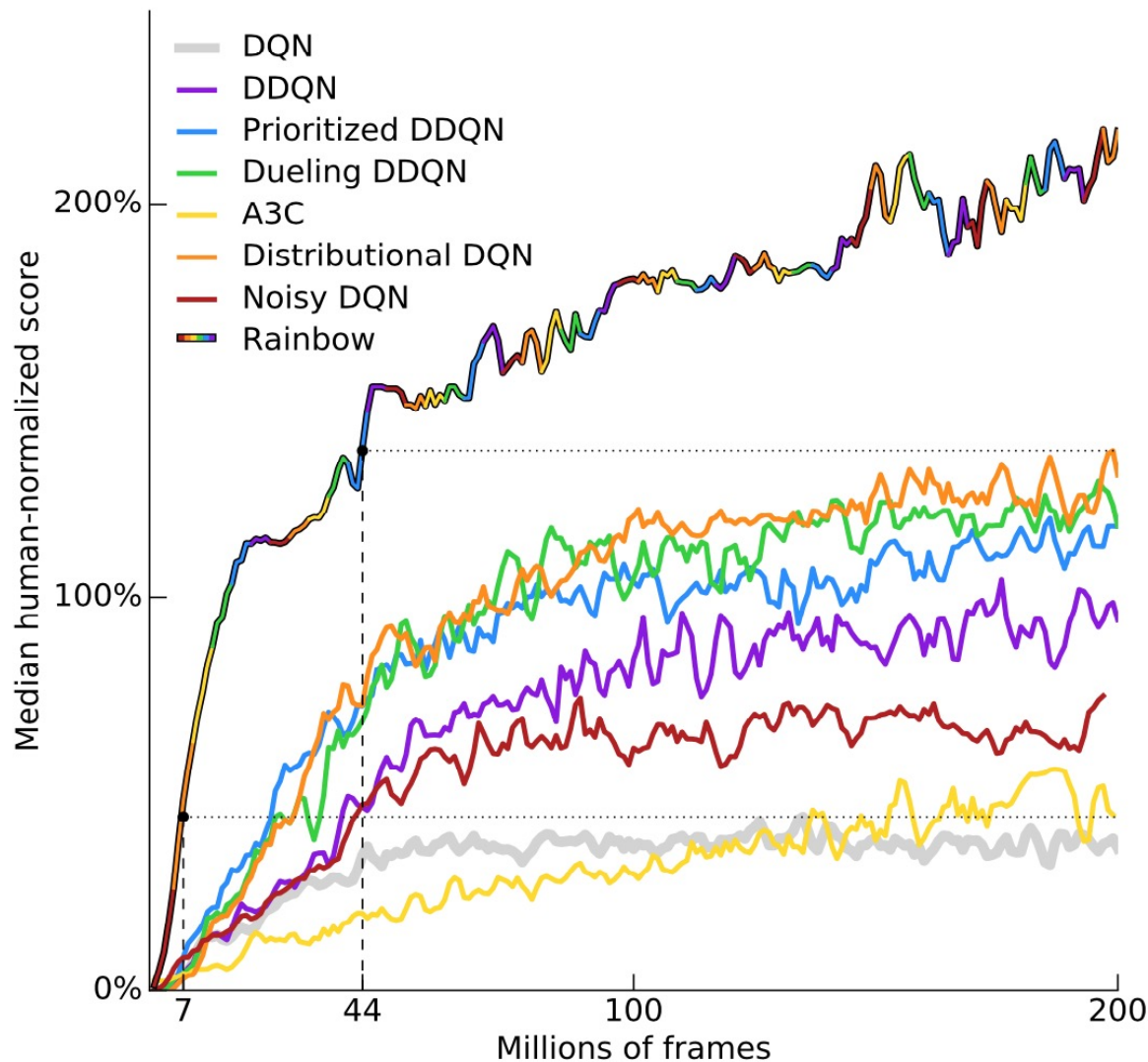
**Single stream Q-network**



**Dueling Q-network**

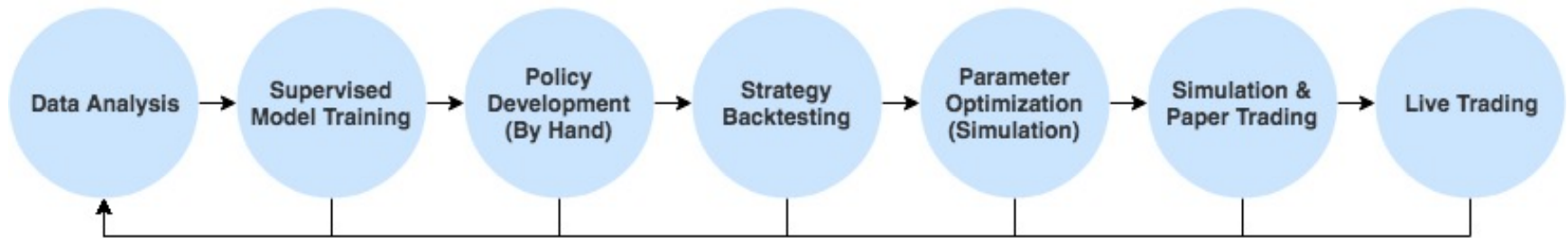


# Rainbow: Combining improvements in deep reinforcement learning

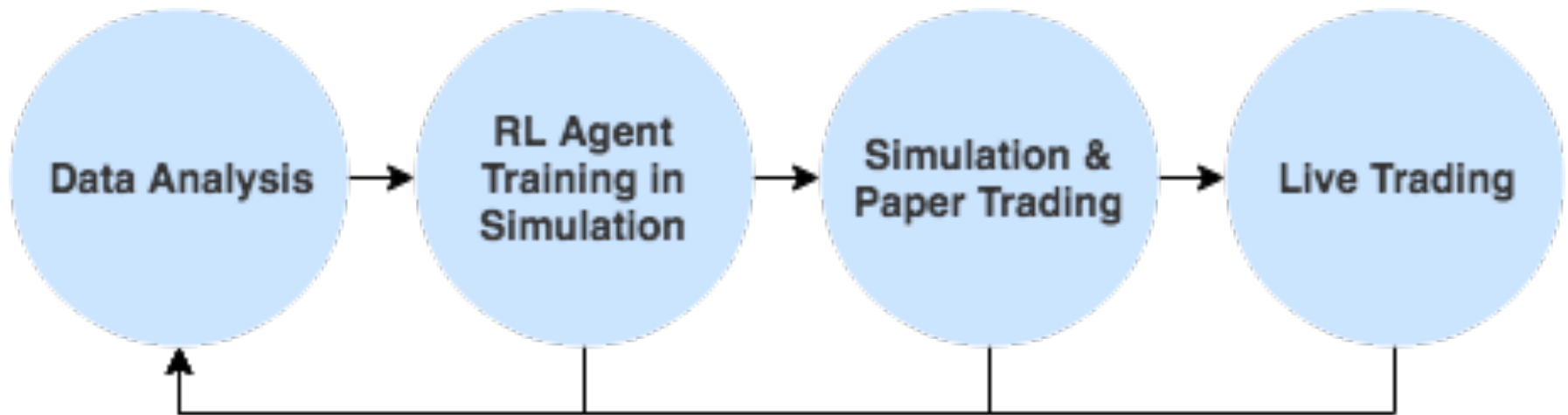


Source: Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver (2017). "Rainbow: Combining improvements in deep reinforcement learning." arXiv preprint arXiv:1710.02298 (2017).

# A Typical Strategy Development Workflow

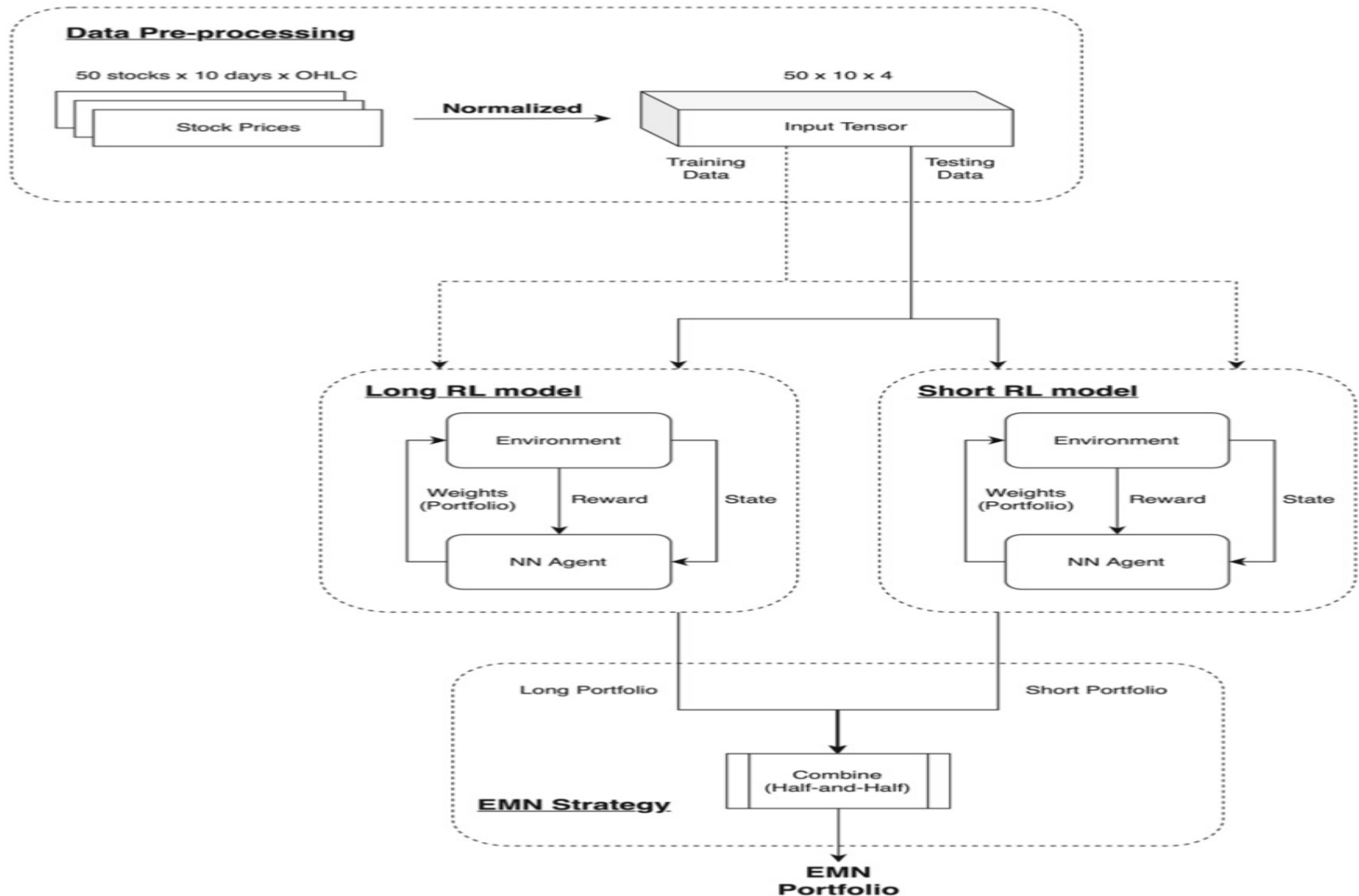


# Reinforcement Learning (RL) in Trading Strategies



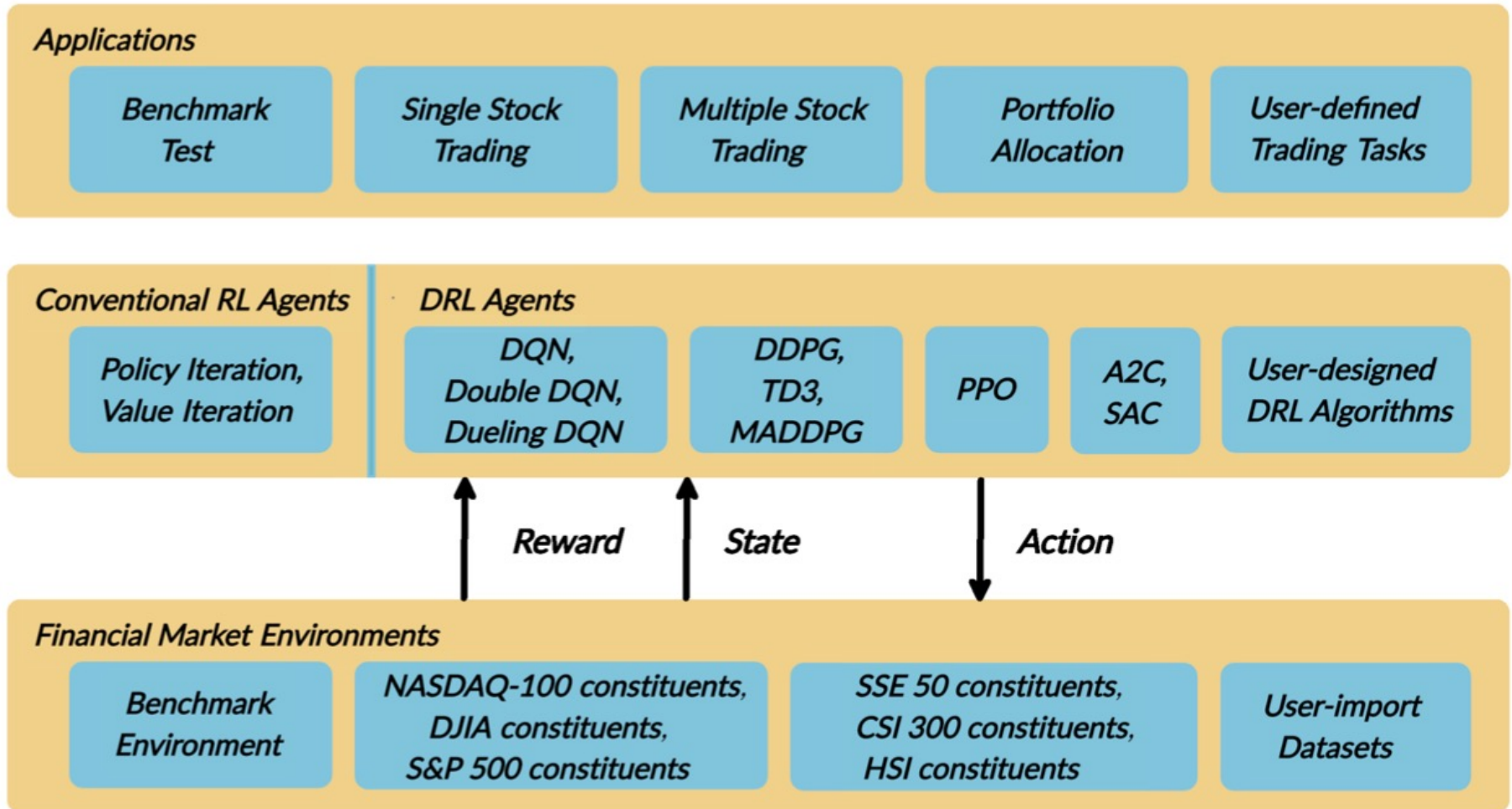
# Portfolio management system in equity market neutral using reinforcement learning

(Wu et al., 2021)



# FinRL:

## A Deep Reinforcement Learning Library for Automated Stock Trading in Quantitative Finance



# FinRL

## Deep Reinforcement Learning Algorithms

Algorithms	Input	Output	Type	State-action spaces support	Finance use cases support	Features and Improvements	Advantages
DQN	States	Q-value	Value based	Discrete only	Single stock trading	Target network, experience replay	Simple and easy to use
Double DQN	States	Q-value	Value based	Discrete only	Single stock trading	Use two identical neural network models to learn	Reduce overestimations
Dueling DQN	States	Q-value	Value based	Discrete only	Single stock trading	Add a specialized dueling Q head	Better differentiate actions, improves the learning
DDPG	State action pair	Q-value	Actor-critic based	Continuous only	Multiple stock trading, portfolio allocation	Being deep Q-learning for continuous action spaces	Better at handling high-dimensional continuous action spaces
A2C	State action pair	Q-value	Actor-critic based	Discrete and continuous	All use cases	Advantage function, parallel gradients updating	Stable, cost-effective, faster and works better with large batch sizes
PPO	State action pair	Q-value	Actor-critic based	Discrete and continuous	All use cases	Clipped surrogate objective function	Improve stability, less variance, simply to implement
SAC	State action pair	Q-value	Actor-critic based	Continuous only	Multiple stock trading, portfolio allocation	Entropy regularization, exploration-exploitation trade-off	Improve stability
TD3	State action pair	Q-value	Actor-critic based	Continuous only	Multiple stock trading, portfolio allocation	Clipped double Q-Learning, delayed policy update, target policy smoothing.	Improve DDPG performance
MADDPG	State action pair	Q-value	Actor-critic based	Continuous only	Multiple stock trading, portfolio allocation	Handle multi-agent RL problem	Improve stability and performance

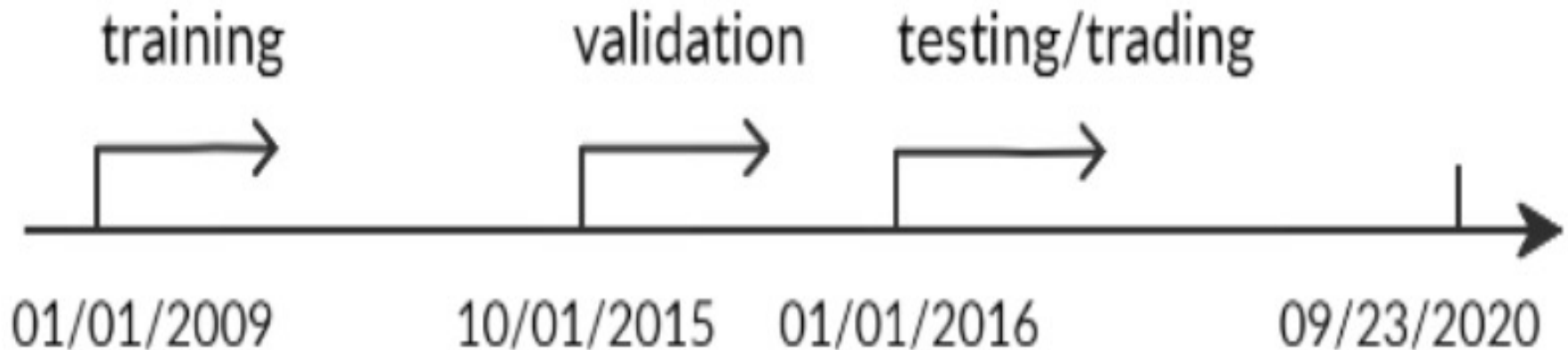


# FinRL:

A Deep Reinforcement Learning Library for  
Automated Stock Trading in Quantitative Finance

**Evaluation of Trading Performance**

**Training-Validation-Testing Flow**



# Reinforcement Learning (RL)

## FinRL

**Performance of single stock trading  
using Proximal policy optimization (PPO) in the FinRL library**



# Reinforcement Learning (RL)

## FinRL

### Performance of multiple stock trading and portfolio allocation using the FinRL library



# Reinforcement Learning (RL)

## FinRL

### Performance of single stock trading using Proximal policy optimization (PPO) in the FinRL library

2019/01/01-2020/09/23	SPY	QQQ	GOOGL	AMZN	AAPL	MSFT	S&P 500
Initial value	100,000	100,000	100,000	100,000	100,000	100,000	100,000
Final value	127,044	163,647	174,825	192,031	173,063	172,797	133,402
Annualized return	14.89%	32.33%	37.40%	44.94%	36.88%	36.49%	17.81%
Annualized Std	9.63%	27.51%	33.41%	29.62%	25.84%	33.41%	27.00%
Sharpe ratio	1.49	1.16	1.12	1.40	1.35	1.10	0.74
Max drawdown	20.93%	28.26%	27.76%	21.13%	22.47%	28.11%	33.92%

# Reinforcement Learning (RL)

## FinRL

### Performance of **multiple stock trading** and **portfolio allocation**

over the DJIA constituents stocks using the FinRL library

2019/01/01-2020/09/23	TD3	DDPG	Min-Var.	DJIA
Initial value	1,000,000	1,000,000	1,000,000	1,000,000
Final value	1,403,337; 1,381,120	1,396,607; 1,281,120	1,171,120	1,185,260
Annualized return	21.40%; 17.61%	20.34%; 15.81%	8.38%	10.61%
Annualized Std	14.60%; 17.01%	15.89%; 16.60%	26.21%	28.63%
Sharpe ratio	1.38; 1.03	1.28; 0.98	0.44	0.48
Max drawdown	11.52% 12.78%	13.72%; 13.68%	34.34%	37.01%

# Deep Reinforcement Learning Library

- OpenAI Gym
- Google Dopamine
- RLlib
- Horizon
- FinRL

# Open AI Gym

Environments Documentation

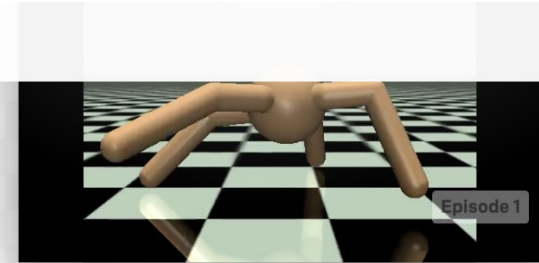


## Gym

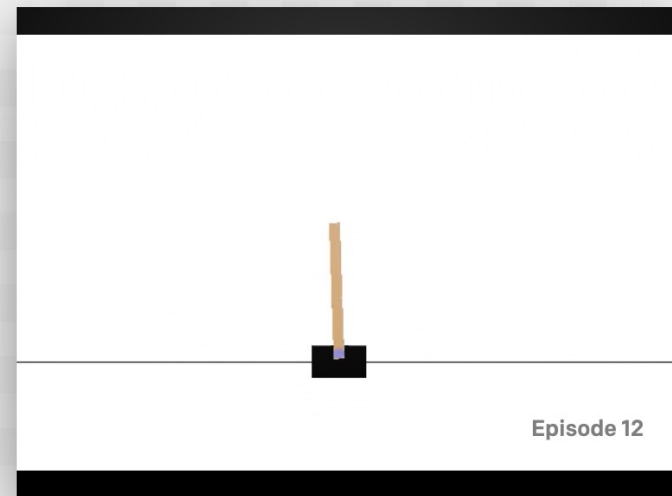
Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.

[View documentation >](#)

[View on GitHub >](#)



RandomAgent on Ant-v2



RandomAgent on CartPole-v1

# Google Dopamine



Dopamine is a research framework for fast prototyping of reinforcement learning algorithms.

<https://github.com/google/dopamine>



# Deep Reinforcement Learning

## Dopamine Colab Examples

### DQN Rainbow

agents.ipynb

File Edit View Insert Runtime Tools Help

CONNECTED EDITING

Table of contents Code snippets Files

Dopamine: How to create and train a custom agent

- Install necessary packages.
- Necessary imports and globals.
- Load baseline data

Example 1: Train a modified version of DQN

- Create an agent based on DQN, but choosing actions randomly.
- Train MyRandomDQNAgent.
- Load the training logs.
- Plot training results.

Example 2: Train an agent built from scratch.

- Create a completely new agent from scratch.
- Train StickyAgent.
- Load the training logs.

Copyright 2018 The Dopamine Authors.

Licensed under the Apache License, Version 2.0 (the "License"); you may not use this file except in compliance with the License. You may obtain a copy of the License at <https://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

### ▼ Dopamine: How to create and train a custom agent

This colab demonstrates how to create a variant of a provided agent (Example 1) and how to create a new agent from scratch (Example 2).

Run all the cells below in order.

```
[ ] Install necessary packages.
```

```
[ ] Necessary imports and globals.
```

```
    BASE_PATH: '/tmp/colab_dope_run'
```

```
    GAME: 'Asterix'
```

```
[ ] Load baseline data
```

# RLlib: Scalable Reinforcement Learning

- Examples
- Tune API Reference
- Contributing to Tune
- RLLIB**
- RLlib: Scalable Reinforcement Learning**
- RLlib Table of Contents
- RLlib Training APIs
- RLlib Environments
- RLlib Models, Preprocessors, and Action Distributions
- RLlib Algorithms
- RLlib Sample Collection and Trajectory Views
- RLlib Offline Datasets
- RLlib Concepts and Custom Algorithms
- RLlib Examples
- RLlib Package Reference
- Contributing to RLlib
- RAY SGD**
- RaySGD: Distributed Training

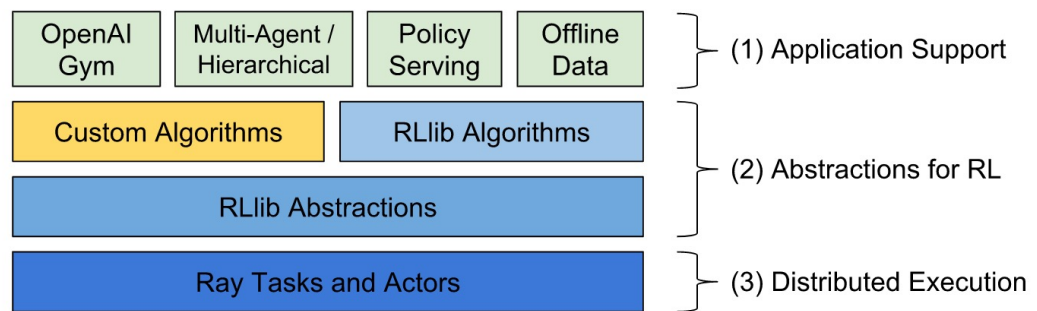


☰ Contents

## RLlib: Scalable Reinforcement Learning

RLlib is an open-source library for reinforcement learning that offers both high scalability and a unified API for a variety of applications. RLlib natively supports TensorFlow, TensorFlow Eager, and PyTorch, but most of its internals are framework agnostic.

- RLlib in 60 seconds
- Running RLlib
- Policies
- Sample Batches
- Training
- Application Support
- Customization



To get started, take a look over the [custom env example](#) and the [API documentation](#). If you're looking to develop custom algorithms with RLlib, also check out [concepts and custom algorithms](#).

## RLlib in 60 seconds

The following is a whirlwind overview of RLlib. For a more in-depth guide, see also the [full table of contents](#) and [RLlib blog posts](#). You may also want to skim the [list of built-in algorithms](#). Look out for the 🏠 and 🔄 icons to see which algorithms are [available](#) for each framework.

v: master

# Papers with Code State-of-the-Art (SOTA)



Search for papers, code and tasks



[Browse State-of-the-Art](#)

[Follow](#)

[Discuss](#)

[Trends](#)

[About](#)

[Log In/Register](#)

## Browse State-of-the-Art

1509 leaderboards • 1327 tasks • 1347 datasets • 17810 papers with code

Follow on [Twitter](#) for updates

## Computer Vision



**Semantic Segmentation**

33 leaderboards  
667 papers with code



**Image Classification**

52 leaderboards  
564 papers with code



**Object Detection**

54 leaderboards  
467 papers with code



**Image Generation**

51 leaderboards  
231 papers with code



**Pose Estimation**

40 leaderboards  
231 papers with code

[See all 707 tasks](#)

## Natural Language Processing



**Machine Translation**



**Language Modelling**



**Question Answering**



**Sentiment Analysis**

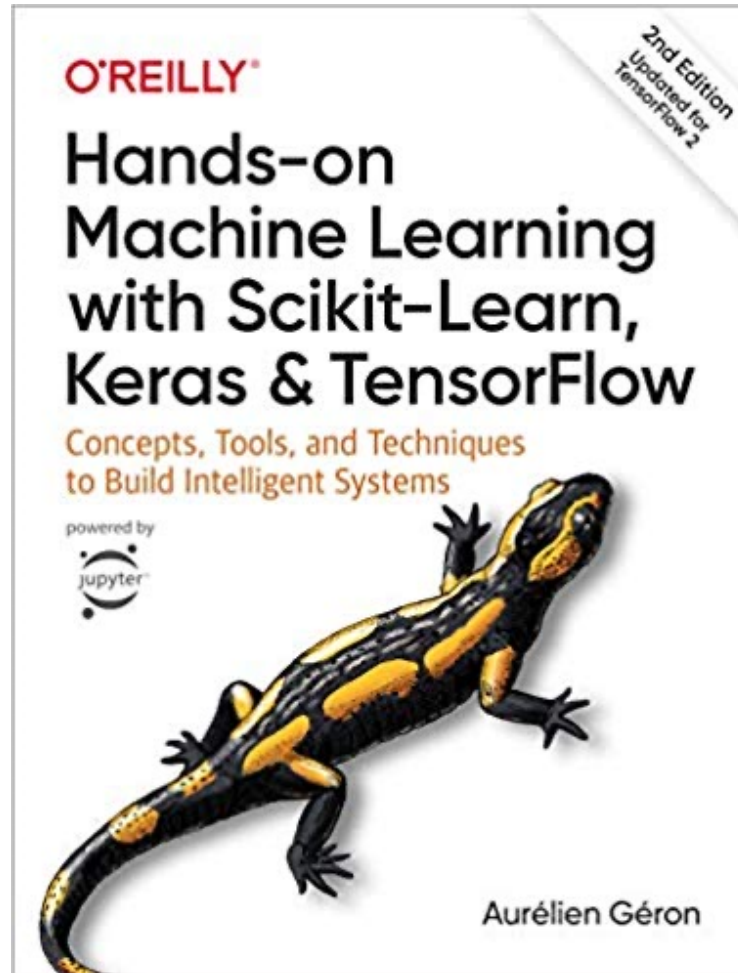


**Text Generation**

<https://paperswithcode.com/sota>

Aurélien Géron (2019),

**Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow:  
Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition**  
O'Reilly Media, 2019

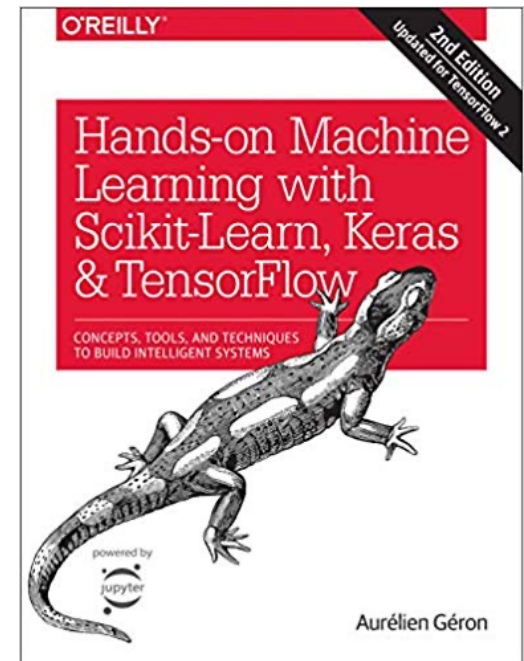


<https://github.com/ageron/handson-ml2>

# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

## Notebooks

- [1. The Machine Learning landscape](#)
- [2. End-to-end Machine Learning project](#)
- [3. Classification](#)
- [4. Training Models](#)
- [5. Support Vector Machines](#)
- [6. Decision Trees](#)
- [7. Ensemble Learning and Random Forests](#)
- [8. Dimensionality Reduction](#)
- [9. Unsupervised Learning Techniques](#)
- [10. Artificial Neural Nets with Keras](#)
- [11. Training Deep Neural Networks](#)
- [12. Custom Models and Training with TensorFlow](#)
- [13. Loading and Preprocessing Data](#)
- [14. Deep Computer Vision Using Convolutional Neural Networks](#)
- [15. Processing Sequences Using RNNs and CNNs](#)
- [16. Natural Language Processing with RNNs and Attention](#)
- [17. Representation Learning Using Autoencoders](#)
- [18. Reinforcement Learning](#)
- [19. Training and Deploying TensorFlow Models at Scale](#)



# Python in Google Colab (Python101)

<https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT>

The screenshot shows the Google Colab interface for a notebook titled 'python101.ipynb'. The top navigation bar includes 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help', with a status 'All changes saved'. On the right, there are icons for 'Comment', 'Share', and 'Settings', along with a user profile icon 'A'. Below the navigation bar, there are indicators for 'RAM' and 'Disk' usage, and a status 'Editing'.

The left sidebar contains a 'Table of contents' with the following items:

- Machine Learning with scikit-learn
  - Classification and Prediction
    - Support Vector Machine (SVM)
    - Random Forest
  - K-Means Clustering
- Deep Learning**
  - Image Classification
  - Text Classification: IMDB Movie Review
- Deep Learning for Financial Time Series Forecasting
- Portfolio Optimization and Algorithmic Trading
  - Investment Portfolio Optimisation with Python
  - Efficient Frontier Portfolio Optimisation in Python
  - Investment Portfolio Optimization
- Text Analytics and Natural Language Processing (NLP)
  - Python for Natural Language Processing
  - spaCy Chinese Model

The main content area shows a code cell with the following Python code:

```
1 import tensorflow as tf
2 mnist = tf.keras.datasets.mnist
3
4 (x_train, y_train), (x_test, y_test) = mnist.load_data()
5 x_train, x_test = x_train / 255.0, x_test / 255.0
6
7 model = tf.keras.models.Sequential([
8     tf.keras.layers.Flatten(input_shape=(28, 28)),
9     tf.keras.layers.Dense(128, activation='relu'),
10    tf.keras.layers.Dropout(0.2),
11    tf.keras.layers.Dense(10, activation='softmax')
12 ])
13
14 model.compile(optimizer='adam',
15               loss='sparse_categorical_crossentropy',
16               metrics=['accuracy'])
17
18 model.fit(x_train, y_train, epochs=5)
19 model.evaluate(x_test, y_test)
```

Below the code cell, the output shows the training progress for Epoch 1/5:

```
Epoch 1/5
1875/1875 [=====] - 4s 2ms/step - loss: 0.4790 - accuracy: 0.8606
```

<https://tinyurl.com/aintpuppython101>

# Summary

- Reinforcement Learning (RL)
  - Markov Decision Processes (MDP)
- Deep Reinforcement Learning (DRL) Algorithms
  - SARSA
  - Q-Learning
  - DQN
  - A3C
  - Rainbow

# References

- Stuart Russell and Peter Norvig (2020), *Artificial Intelligence: A Modern Approach*, 4th Edition, Pearson.
- Aurélien Géron (2019), *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*, 2nd Edition, O'Reilly Media.
- Richard S. Sutton & Andrew G. Barto (2018), *Reinforcement Learning: An Introduction*, 2nd Edition, A Bradford Book.
- David Silver (2015), Introduction to reinforcement learning, <https://www.youtube.com/playlist?list=PLqYmG7hTraZDM-OYHWgPebj2MfCFzFObQ>
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharmashan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, Demis Hassabis (2018), "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play." *Science* 362, no. 6419 (2018): 1140-1144.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis (2017), "Mastering the game of Go without human knowledge." *Nature* 550 (2017): 354–359.
- Hado Van Hasselt, Arthur Guez, and David Silver (2016). "Deep Reinforcement Learning with Double Q-Learning." In *AAAI*, vol. 2, p. 5. 2016.
- Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver (2017). "Rainbow: Combining improvements in deep reinforcement learning." *arXiv preprint arXiv:1710.02298* (2017).
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves et al. (2015) "Human-level control through deep reinforcement learning." *Nature* 518, no. 7540 (2015): 529.
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Van Hasselt, Marc Lanctot, and Nando De Freitas (2015). "Dueling network architectures for deep reinforcement learning." *arXiv preprint arXiv:1511.06581* (2015).
- Xiao-Yang Liu, Hongyang Yang, Qian Chen, Runjia Zhang, Liuqing Yang, Bowen Xiao, and Christina Dan Wang (2020). "FinRL: A Deep Reinforcement Learning Library for Automated Stock Trading in Quantitative Finance." *arXiv preprint arXiv:2011.09607* (2020).
- Mu-En Wu, Jia-Hao Syu, Jerry Chun-Wei Lin, and Jan-Ming Ho. "Portfolio management system in equity market neutral using reinforcement learning." *Applied Intelligence* (2021): 1-13.
- Min-Yuh Day (2021), Python 101, <https://tinyurl.com/aintpupython101>