**LECTURER: TAI LE QUY**

# INTRODUCTION TO

# REINFORCEMENT LEARNING

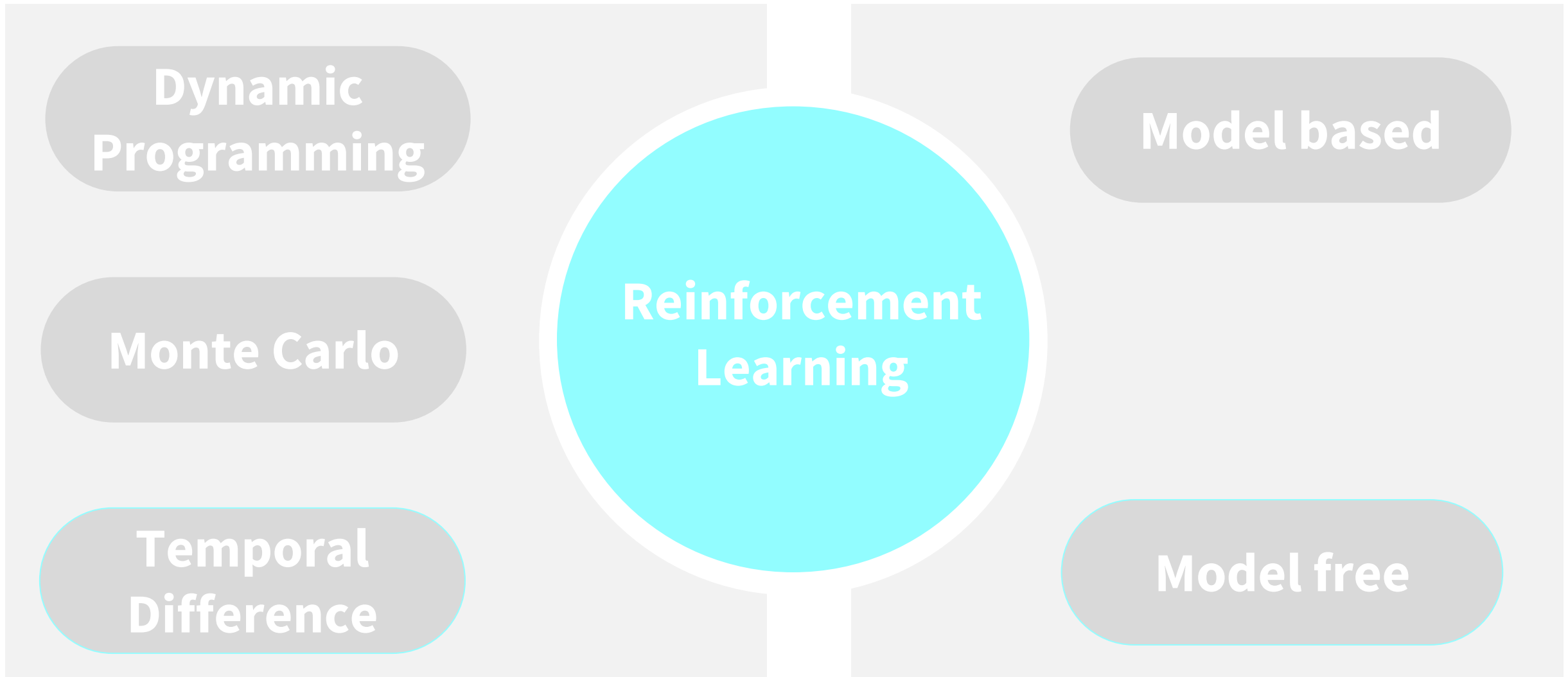# REINFORCEMENT LEARNING

# ALGORITHMS AND THEIR PROPERTIES

— Describe the mechanisms of temporal difference algorithms and their importance in reinforcement leaning (RL)

— Compare and contrast exploration and exploitation strategies used by RL agents to balance their behaviors

— Understand how SARSA and Q-learning algorithms behave in different RL scenarios

1. Explain how RL algorithms achieve optimal policies through mapping between states and actions

2. Describe the differences between model-based and model-free RL algorithms

3. Compare and contrast SARSA and Q-learning algorithms discussing their strengths and weaknesses

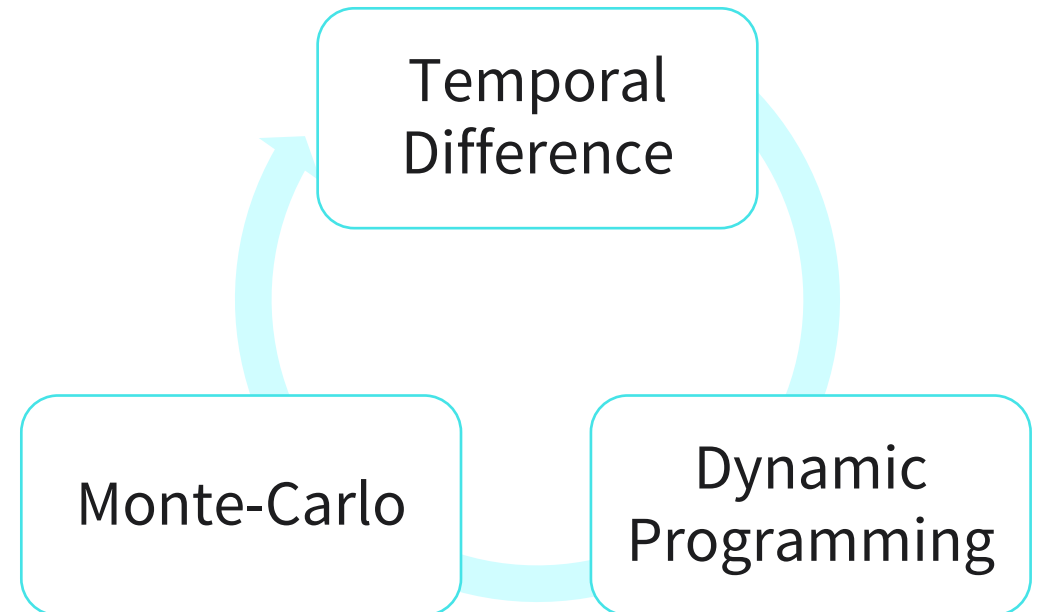# CATEGORIZATIONS OF REINFORCEMENT LEARNING

– Categorized based on whether the algorithm has **access** to the internal dynamics of the environment that represents the problem being solved or not.

– Model-based learn optimal policies using **transition functions** (could be true models of the environment, if available, or they could be learned approximations)

  – Dynamic programming (DP)

– Model-free algorithms learn optimal policies directly from the raw experiences they gather through their interactions with the environment

  – Monte-Carlo (MC) method

  – Temporal difference (TD) learning (build on DP and MC)

    – Q-learning

    – State-action-reward-state-action (SARSA)

## Predict a quantity based on future values of a signal

– Dynamic programming assumes perfect world model

– Monte-Carlo methods are model-free and learn from raw experiences

– Temporal difference methods combine the best of both worlds

Temporal Difference

Monte-Carlo

Dynamic Programming

- State value function $\quad V_\pi(s_t) = E[G_t \mid s_0 = s_t]$

- Action value function $\quad Q_\pi(s_t, a_t) = E[G_t \mid s_0 = s_t, a_0 = a_t]$

- Discounted rewards $\quad G_t = \sum_{t=0}^{T} \gamma^t \cdot r_t$

- Bellman equations

$$V_\pi(s_t) = E_{s_{t+1} \sim P}[r_t + \gamma \cdot V_\pi(s_{t+1})]$$

$$Q_\pi(s_t, a_t) = E_{s_{t+1} \sim P}\left[r_t + \gamma \cdot E_{a_{t+1} \sim \pi}[Q_\pi(s_{t+1}, a_{t+1})]\right]$$

- The next state $s_{t+1}$ is sampled from the transition function P and the next action $a_{t+1}$ is sampled from the policy $\pi$

– Value functions in DP are updated as per the Bellman equation

$$V_{k+1}(s_t) = E\big[r_t + \gamma \,.\, V_k(s_{t+1})\big]$$

– The estimate of the value function is updated to $V_{k+1}$ based on the current estimate $V_k$. (**bootstrapping** – "taking a guess from a guess")

– Assume access to a perfect model of the world

- MC methods learn **directly** from the raw experiences that are gathered through the agent-environment interactions (model-free), **sampling**.

$$V_{k+1}(s_t) = V_k(s_t) + \alpha \cdot [G_t - V_k(s_t)]$$

- $G_t$ : the true return from time step $t$, $[G_t - V_k(s_t)]$: error term identifying how far off the V -function's estimate is from the true return
- The hyperparameter α: step size or the learning rate.
  - Controls the rate at which new information, If α is large, then the error term has a large influence on the factor by which the estimate is updated
  - Algorithms must wait until the end of the **episode** when $G_t$ is available to determine the update that would be applied to $V_k$.
- MC methods only work with episodic tasks, i.e., tasks, such as games like Chess and Go that have a clear end as marked by a terminal state.

– TD is model-free and learns by sampling, **do not wait** till the end of the episode to perform an update step.
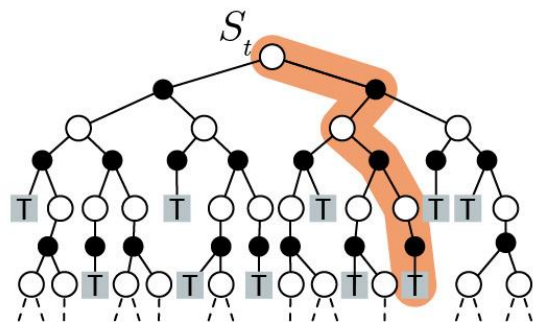
$$V_{k+1}(s_t) = V_k(s_t) + \alpha \cdot \left[ r_t + \gamma \cdot V_k(s_{t+1}) - V_k(s_t) \right]$$

- The true return $G_t$ has been replaced by an **estimated** value
- $V_k(s_{t+1})$ is called TD target
- The difference between the target and the current value is an error term called TD error, denoted as $\delta_t$

## TEMPORAL DIFFERENCE LEARNING AND Q-FACTORS

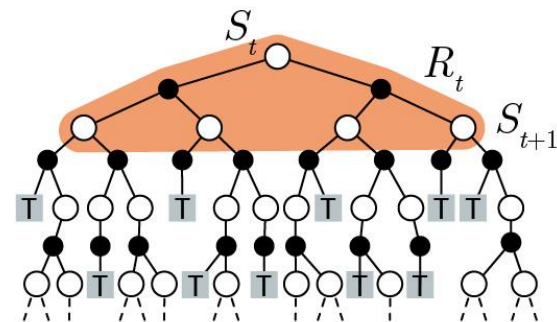### Monte Carlo
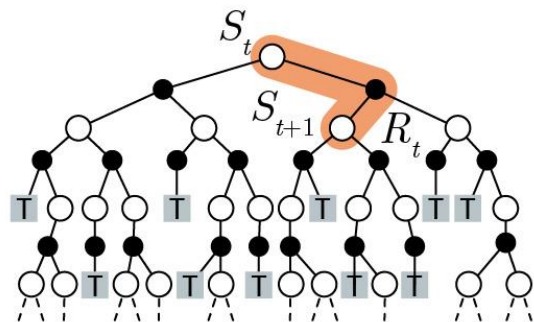$$V(S_t) \mapsto V(S_t) + \alpha \cdot (G_t - V(S_t))$$



### Dynamic Programming
$$V(S_t) \mapsto E_\pi [R_t + \gamma \cdot V(S_{t+1})]$$
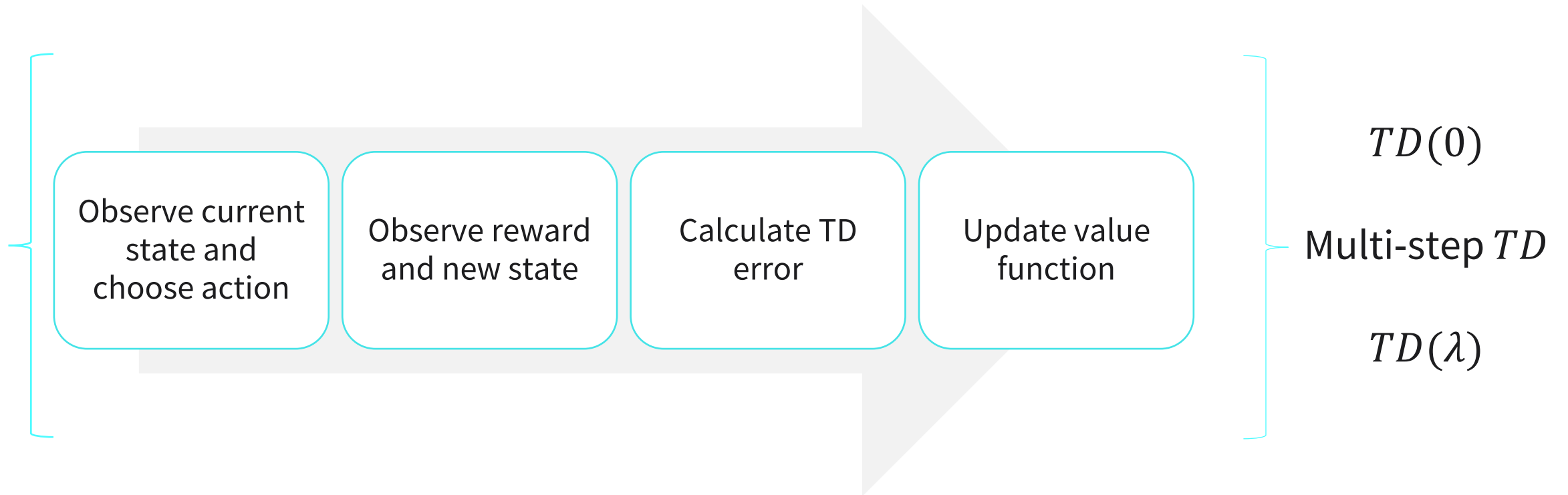


### Temporal Difference
$$V(S_t) \mapsto V(S_t) + \alpha(R_t + \gamma \cdot V(S_{t+1}) - V(S_t))$$

# TD algorithm: efficient learning from raw experiences

| Observe current state and choose action | Observe reward and new state | Calculate TD error | Update value function |
|---|---|---|---|

$TD(0)$

Multi-step $TD$

$TD(\lambda)$

- TD(0): just one step into the future, performing updates after every single transition

- 2-step looks ahead

$$V_{k+1}(s_t) = V_k(s_t) + \alpha . \left[ r_t + \gamma . r_{t+1} + \gamma^2 . V_k(s_{t+2}) - V_k(s_t) \right]$$

- Updated equation

$$V_{k+1}(s_t) = V_k(s_t) + \alpha . [G_t - V_k(s_t)]$$

- N-step TD

$$G_{t:t+n} = r_t + \gamma . r_{t+1} + \gamma^2 . r_{t+2} + ... + \gamma^n . V_k(s_{t+n})$$

- n = ∞, the TD update rule becomes equivalent to the MC update rule

- if t + n = T, where T is the length of the entire trajectory, we will get the true return $G_t$. → choosing n is difficult

- The λ-return is computed as geometrical weighted sum of the **n-step** returns

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} . G_{t:t+n}$$

  - The number of steps being considered ranging from 1, all the way to ∞.
  - λ ∈ [0, 1] . When λ = 0, we get the 1-step formulation of TD and increasing the value of λ increases the number of rewards to be sampled to estimate the return

**BENEFITS AND SHORT-COMINGS OF TD LEARNING**

BENEFITS

- Online learning

- Learn from incomplete sequences

- Episodic and non-episodic tasks

- Efficient

- Faster convergence

- Sensitive to start state

- Biased compared to MC

- Overestimates with function approximation

- Requires knowledge of rewards

- Limited to Markov process

SHORT COMINGS

## Exploration vs Exploitation Dilemma

"To obtain reward, an agent must prefer effective actions that it has tried in the past. But to discover such actions, it has to try actions that it has not selected before"

### Exploitation

— Pick action based on policy

— Gain large reward

### Exploration

— Pick novel action

— Gain environment knowledge

Source of the text: Sutton and Barto, 2018.

## Balancing exploration vs exploitation in RL

— Epsilon ($\epsilon$) controls exploration and exploitation balance

— Ranges from completely exploitative to completely exploratory

$$a_t = \begin{cases} \pi(s_t) & \textit{with probability } (1 - \epsilon) \\ \\ \textit{random action with probability } \epsilon \end{cases}$$

— Epsilon value decays over time

## Iteratively estimate Q-value using state-action pairs

— TD method to learn Q-values for state-action pairs

— Policy update by action selection

— The update rule for SARSA

$$Q_{k+1}(s_t,\ a_t) = Q_k(s_t,\ a_t) + \alpha \cdot [r_t + \gamma \cdot Q_k(s_{t+1}, a_{t+1}) - Q_k(s_t, a_t)]$$

— SARSA updates its Q-value estimate based on the Q-value of the next state, $s_{t+1}$, and the action, $a_{t+1}$, suggested by the current policy

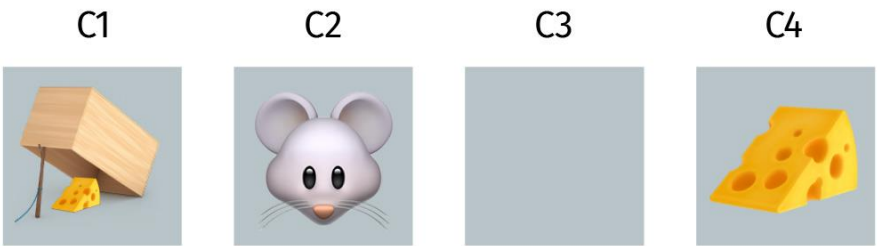— SARSA updates its Q-value based on the policy it currently follows → on-policy

Figure 33: Pseudocode for the SARSA Algorithm

1: Initialize $Q(s,a)$ arbitrarily

2: **for** episode $= 1,2,..., M$ **do**

3:   Initialize state $s_1$

4:   Select action $a_1$ per policy derived from Q (e.g., $\epsilon$-greedy)

5:   **for** $t = 1,2,...,T$ **do**

6:     Perfom $a_t$ and observe $r_t$, $s_{t+1}$

7:     Select action $a_{t+1}$ per policy derived from Q(e.g., $\epsilon$-greedy)

8:     Update $Q(s_t, a_t)$ as $Q(s_t, a_t) + \alpha[r_t + \gamma.Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$

9:     $s_t = s_{t+1}$, $a_t = a_{t+1}$

10:   **end for**

11: **end for**

# ON POLICY LEARNING: SARSA, EXAMPLE

| C1 | C2 | C3 | C4 |
|----|----|----|----|



The Mouse Gridworld Example

| Cell | Left | Right |
|------|------|-------|
| C1 | 0 | 0 |
| C2 | -0.9 | 0.7 |
| C3 | 0.4 | 0.9 |
| C4 | 0 | 0 |

Q-Table Representing Q-Values for the Mouse Gridworld Example

Estimate optimal Q-value using highest expected return

— Q-learning update: reward
+ maximum future Q-value

$$Q_{k+1}(s_t, a_t)$$
$$= Q_k(s_t, a_t) + \alpha \cdot \left[ r_t + \gamma \cdot \max_{a_{t+1}} Q_k(s_{t+1}, a_{t+1}) - Q_k(s_t, a_t) \right]$$
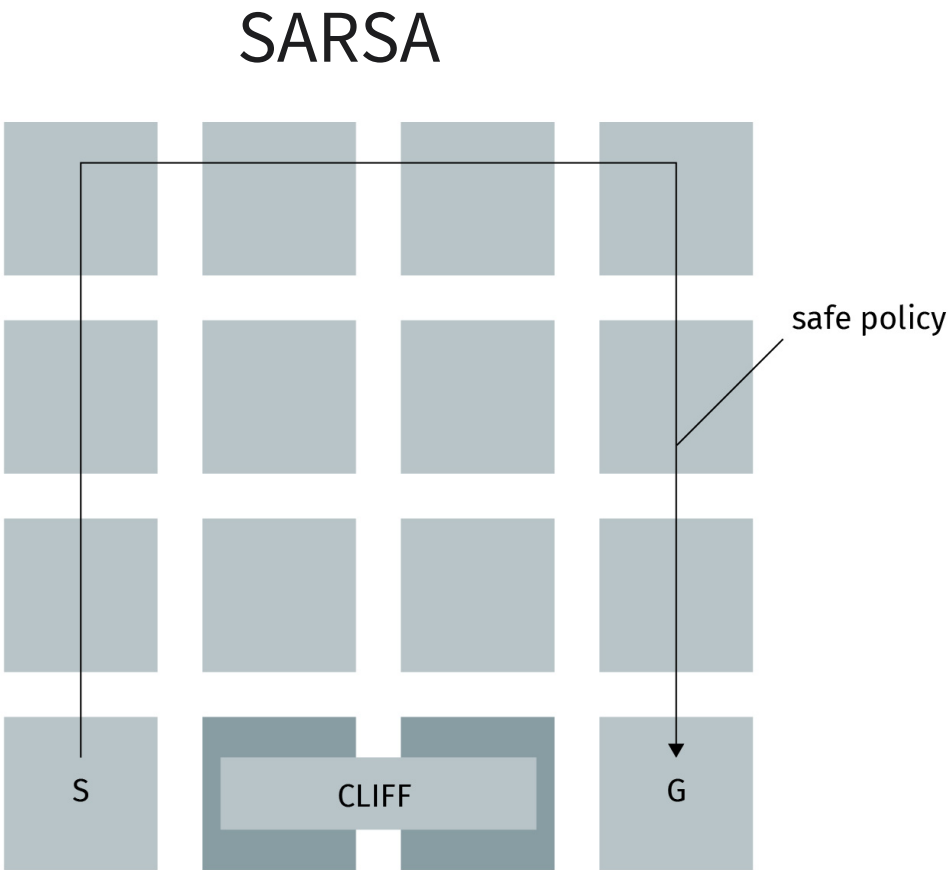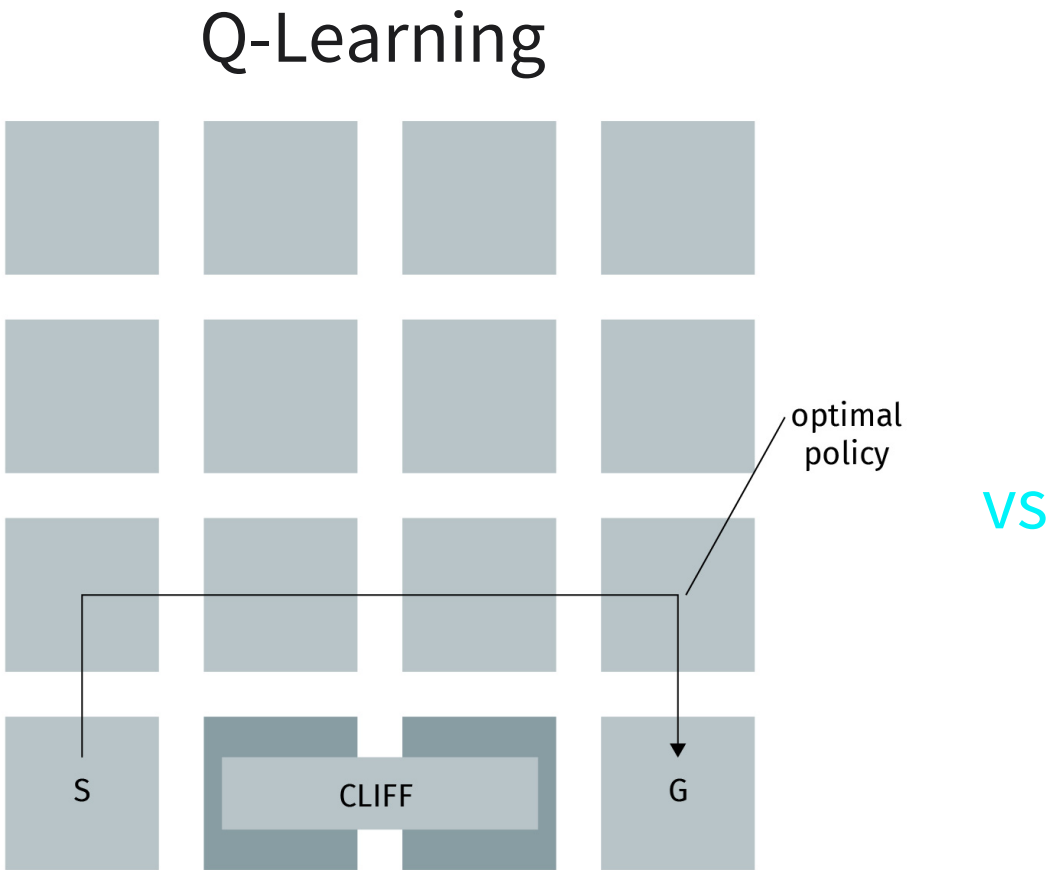
— Future Q-value computed
assuming optimal policy

**Figure 35: Q-Learning**

1: Initialize $Q(s,a)$ arbitrarily
2: **for** episode $= 1,2,..., M$ **do**
3:     Initialize state $s_1$
4:     **for** $t = 1,2,...,T$ **do**
5:         Select $a_t$ per policy derived from Q (e.g., $\epsilon$-greedy)
6:         Perfom $a_t$ and observe $r_t$, $s_{t+1}$
7:         Update $Q(s_t, a_t)$ as $Q(s_t, a_t) + \alpha[r_t + \gamma.\max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$
8:         $s_t = s_{t+1}$
9:     **end for**
10: **end for**

**EXAMPLE: Q-LEARNING VERSUS SARSA IN THE CLIFF WALKING GRIDWORLD**



Source of the image: Nair, 2023.

— Describe the mechanisms of temporal difference algorithms and their importance in reinforcement leaning (RL)

— Compare and contrast exploration and exploitation strategies used by RL agents to balance their behaviors

— Understand how SARSA and Q-learning algorithms behave in different RL scenarios

# TRANSFER TASK

# Case study

You are tasked with developing an AI-agent to play Tic-Tac-Toe. Your goal is to develop an optimal policy that maximizes the chances of an agent to win against a human opponent.

# Task

What type of reinforcement learning algorithm would be the most appropriate for this problem? What are the key steps of implementing this algorithm for game playing? Consider factors such as exploration vs exploitation, reward design, state space and action space representation

# Please present your results.

# The results will be discussed in plenary.

1. Temporal difference learning uses the benefits of which methods?
   a) DP and unsupervised learning
   b) MC and genetic algorithms
   c) MC and ML-models
   d) DP and MC

2. Which a popular approach for the agent to explore the environment
   a) Gradient descent
   b) Dynamic programming
   c) Epsilon greedy
   d) TD-learning

3. Using the Q maximizing action in Q-learning makes it a _____ method
   a) Off policy
   b) On policy
   c) Model-based
   d) Model-free

# LIST OF SOURCES

## Text

Sutton, R. S., & Barto A. G. (2018). *Reinforcement learning: an Introduction.* The MIT Press. https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf

Silver, D. (2020). Model-Free Prediction [lecture notes]. https://www.davidsilver.uk/wp-content/uploads/2020/03/MC-TD.pdf

## Images

Silver, 2020.

Nair, 2023.