# TDDC17 LE11 HT2024

## Machine Learning III

**Fredrik Heintz**

**Dept. of Computer Science
Linköping University**

**fredrik.heintz@liu.se**

**@FredrikHeintz**

Outline:

- **Reinforcement learning**

- **Deep reinforcement learning**

- **Multi-objective reinforcement learning**

**LiU** LINKÖPING UNIVERSITY

# Classes of Learning Problems

| Supervised Learning | Unsupervised Learning | Reinforcement Learning |
|---|---|---|
| **Data:** $(x, y)$<br>$x$ is data, $y$ is label | **Data:** $x$<br>$x$ is data, no labels! | **Data:** state-action pairs |
| **Goal:** Learn function to map $x \rightarrow y$ | **Goal:** Learn underlying structure | **Goal:** Maximize future rewards over many time steps |
| **Apple example:** | **Apple example:** | **Apple example:** |
| This thing is an apple. | This thing is like the other thing. | Eat this thing because it will keep you alive. |

# Reinforcement Learning: Key Concepts

AGENT

**Agent:** takes actions.
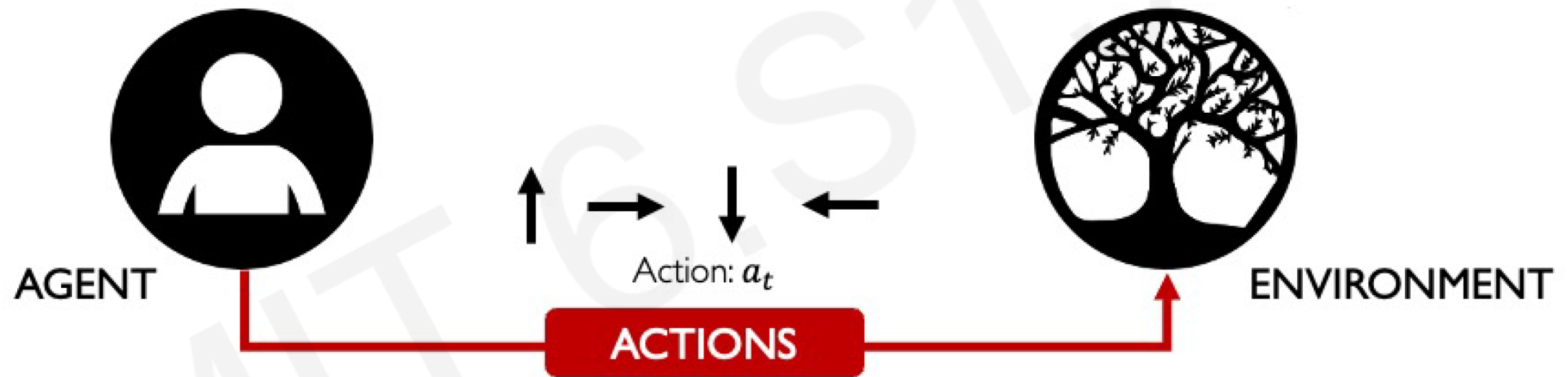
# Reinforcement Learning: Key Concepts

AGENT

ENVIRONMENT

**Environment:** the world in which the agent exists and operates.

# Reinforcement Learning: Key Concepts

**AGENT**

Action: $a_t$

**ACTIONS**

**ENVIRONMENT**

**Action**: a move the agent can make in the environment.

**Action space** $A$: the set of possible actions an agent can make in the environment

LINKÖPING
UNIVERSITY

# Reinforcement Learning: Key Concepts

**OBSERVATIONS**

**AGENT**

Action: $a_t$

**ACTIONS**

**ENVIRONMENT**

**Observations**: of the environment after taking actions.

LINKÖPING
UNIVERSITY

# Reinforcement Learning: Key Concepts

**OBSERVATIONS**

State changes: $s_{t+1}$

**AGENT**

**ENVIRONMENT**

Action: $a_t$

**ACTIONS**

**State**: a situation which the agent perceives.

**LINKÖPING UNIVERSITY**

# Reinforcement Learning: Key Concepts



**OBSERVATIONS**

State changes: $s_{t+1}$
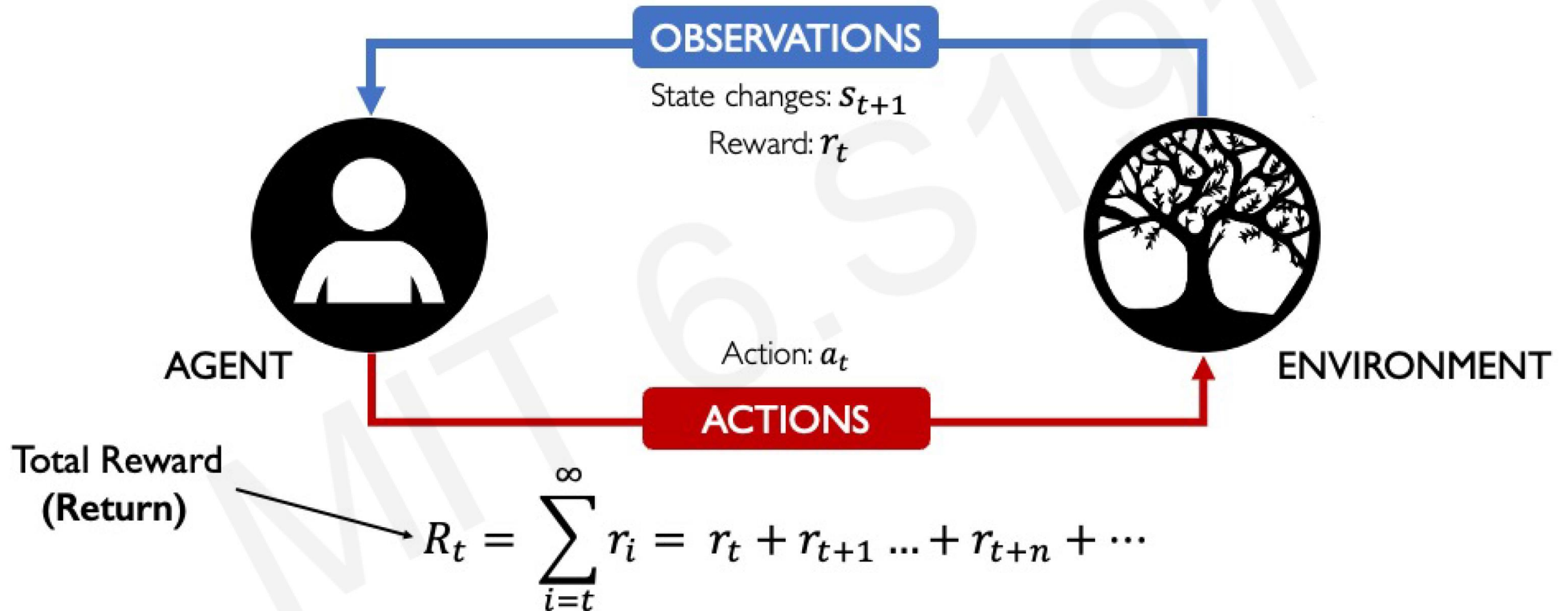
Reward: $r_t$

Action: $a_t$

**ACTIONS**

AGENT

ENVIRONMENT

**Reward**: feedback that measures the success or failure of the agent's action.

LINKÖPING UNIVERSITY

# Reinforcement Learning: Key Concepts



**OBSERVATIONS**

State changes: $s_{t+1}$
Reward: $r_t$

Action: $a_t$

**ACTIONS**

**AGENT**

**ENVIRONMENT**

Total Reward
(Return)

$$R_t = \sum_{i=t}^{\infty} r_i$$

# Reinforcement Learning: Key Concepts

**OBSERVATIONS**

State changes: $s_{t+1}$

Reward: $r_t$

AGENT

Action: $a_t$

**ACTIONS**

ENVIRONMENT

Total Reward (Return)

$$R_t = \sum_{i=t}^{\infty} r_i = r_t + r_{t+1} \ldots + r_{t+n} + \cdots$$

# Reinforcement Learning: Key Concepts

**OBSERVATIONS**

State changes: $s_{t+1}$

Reward: $r_t$

**AGENT**

Action: $a_t$

**ACTIONS**

**ENVIRONMENT**

Discounted Total Reward (Return)
$$R_t = \sum_{i=t}^{\infty} \gamma^i r_i$$

LINKÖPING UNIVERSITY

# Reinforcement Learning: Key Concepts



**OBSERVATIONS**

State changes: $s_{t+1}$
Reward: $r_t$

**AGENT**

Action: $a_t$

**ACTIONS**

**ENVIRONMENT**

Discounted
Total Reward
(Return)
$$R_t = \sum_{i=t}^{\infty} \gamma^i r_i = \gamma^t r_t + \gamma^{t+1} r_{t+1} \dots + \gamma^{t+n} r_{t+n} + \cdots$$
$\gamma$: discount factor; $0 < \gamma < 1$

# A Reinforcement Learning Problem

- The environment
- The reinforcement function $r(s,a)$
  - Pure delay reward and avoidance problems
  - Minimum time to goal
  - Games
- The value function $V(s)$
  - Policy $\pi: S \rightarrow A$
  - Value $V^\pi(s) := \Sigma_i \gamma^i r_{t+i}$
- Find the optimal policy $\pi^*$ that maximizes $V^{\pi^*}(s)$ for all states $s$.

Agent

state    reward                    action

Environment

$$s_0 \xrightarrow[r_0]{a_0} s_1 \xrightarrow[r_1]{a_1} s_2 \xrightarrow[r_2]{a_2} \dots$$

Goal: Learn to choose actions that maximize
$$r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \ , \ \text{where } 0<\gamma<1$$

LINKÖPING UNIVERSITY

# RL Value Function - Example

## A minimum time to goal world



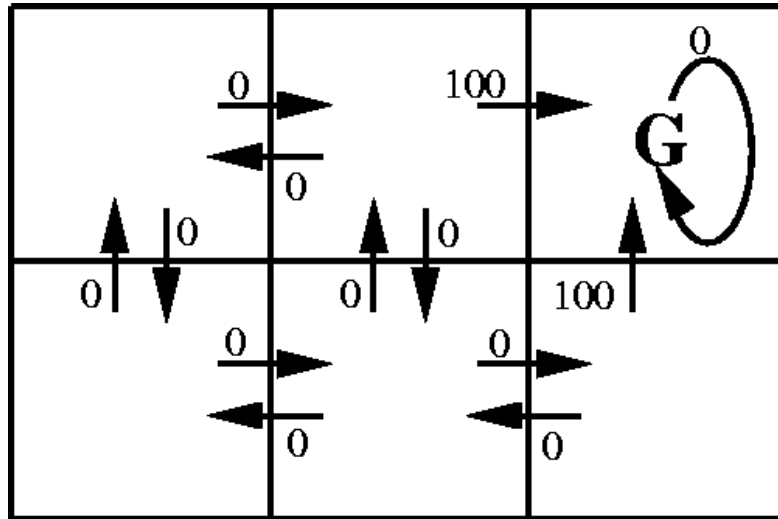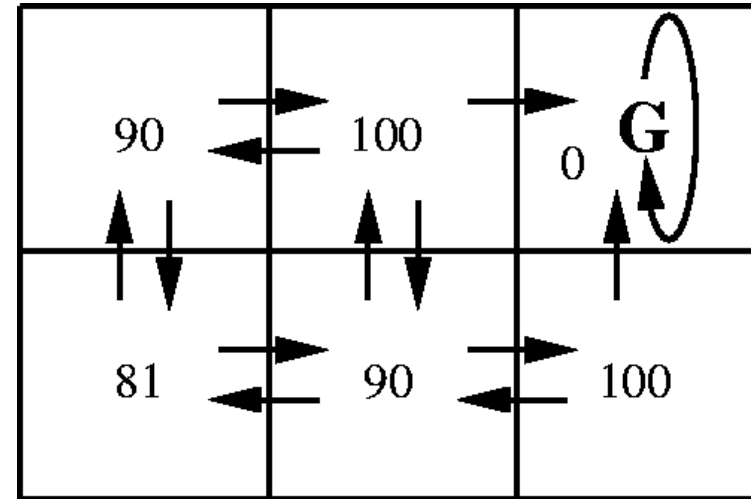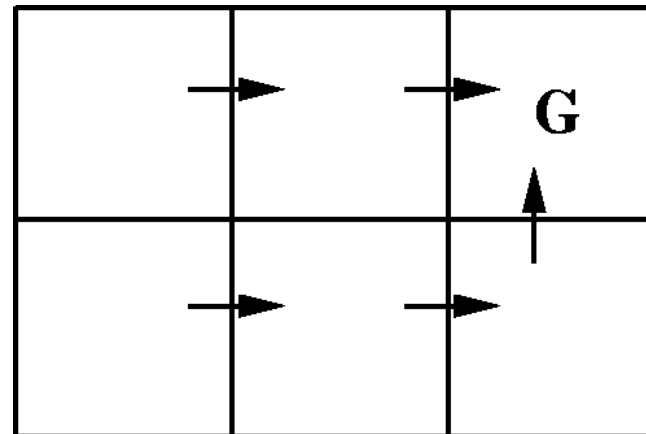| Value function for random movement | Optimal policy | Optimal value function |

# Markov Decision Processes

Assume:

- finite set of states $S$, finite set of actions $A$

- at each discrete time agent observes state $s_t \in S$ and chooses action $a_t \in A$

- then receives immediate reward $r_t$

- and state changes to $s_{t+1}$

- Markov assumption: $s_{t+1} = \delta(s_t, a_t)$ and $r_t = r(s_t, a_t)$
  - i.e. $r_t$ and $s_{t+1}$ depend only on current state and action
  - functions $\delta$ and $r$ may be non-deterministic
  - functions $\delta$ and $r$ not necessarily known to the agent

LINKÖPING
UNIVERSITY

# MDP Example



$r(s,a)$

$V*(s)$

An optimal policy

# Defining the Q-Function

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots$$

Total reward, $R_t$, is the discounted sum of all rewards obtained from time $t$

$$Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t]$$

The Q-function captures the **expected total future reward** an agent in state, $s$, can receive by executing a certain action, $a$

LINKÖPING
UNIVERSITY

# How to Take Actions Given a Q-Function

$$Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t]$$

(state, action)

Ultimately, the agent needs a **policy $\pi(s)$**, to infer the **best action to take** at its state, s

**Strategy:** the policy should choose an action that maximizes future reward
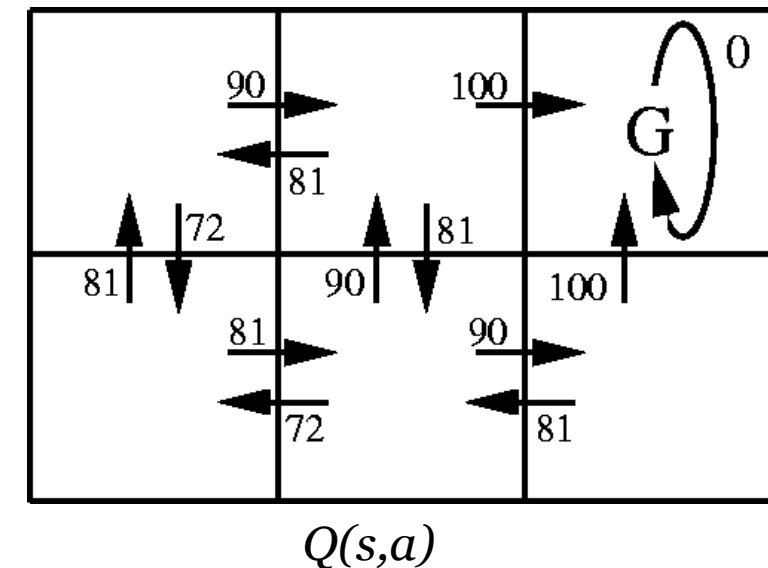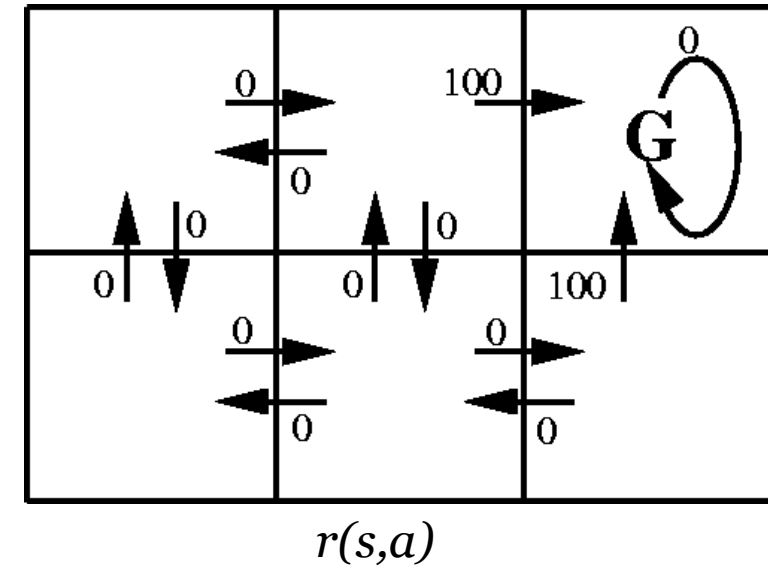
$$\pi^*(s) = \arg\max_a Q(s, a)$$

# The Q-Function

Optimal policy:

- $\pi^*(s) = \text{argmax}_a[r(s,a) + \gamma V^*(\delta(s,a))]$
- Doesn't work if we don't know $r$ and $\delta$.



$r(s,a)$

The Q-function:

- $Q(s,a) := r(s,a) + \gamma V^*(\delta(s,a))$
- $\pi^*(s) = \text{argmax}_a Q(s,a)$



$Q(s,a)$

# The Q-Function

- Note Q and V* closely related:
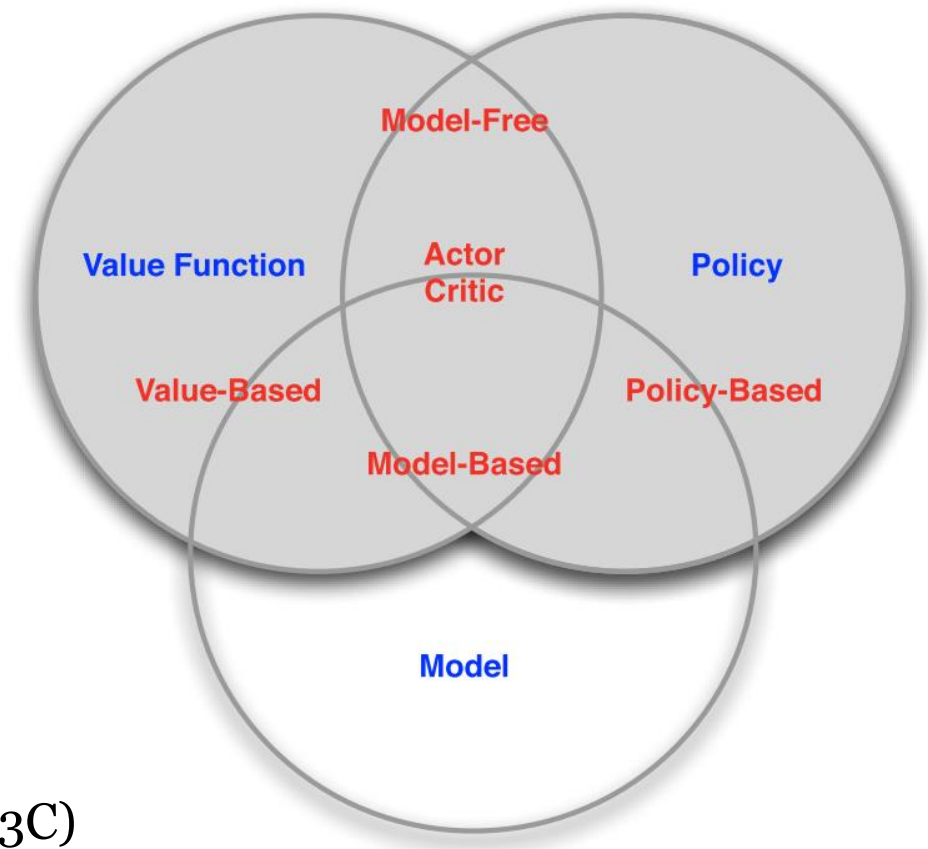  $V^*(s) = \max_{a'} Q(s,a')$

- Therefore Q can be written as:

  $Q(s_t,a_t) := r(s_t,a_t) + \gamma V^*(\delta(s_t,a_t)) =$

  $\qquad\qquad r(s_t,a_t) + \gamma \max_{a'} Q(s_{t+1},a')$

- If $Q^\wedge$ denote the current approximation of Q then it can be updated by:

  $Q^\wedge(s,a) := r + \gamma \max_{a'} Q^\wedge(s',a')$

# Reinforcement Learning Approaches

- Value-Based:
  - Learn value function
  - Implicit policy (e.g. greedy selection)
  - Example: Deep Q Networks (DQN)
- Policy-Based:
  - No value function
  - Learn explicit (stochastic) policy
  - Example: Stochastic Policy Gradients
- Model-Based:
  - Learn transition model
  - Implicit policy
  - Example: Dreamer
- Actor-Critic:
  - Learn value function
  - Learn policy using value function
  - Example: Asynchronous Advantage Actor Critic (A3C)

# Reinforcement Learning Algorithms

**Value Learning**

Find $Q(s, a)$

$a = \underset{a}{\mathrm{argmax}}\, Q(s, a)$

**Policy Learning**

Find $\pi(s)$

Sample $a \sim \pi(s)$

# Reinforcement Learning Algorithms

**Value Learning**

Find $Q(s, a)$

$$a = \underset{a}{\mathrm{argmax}}\, Q(s, a)$$

**Policy Learning**

Find $\pi(s)$

Sample $a \sim \pi(s)$
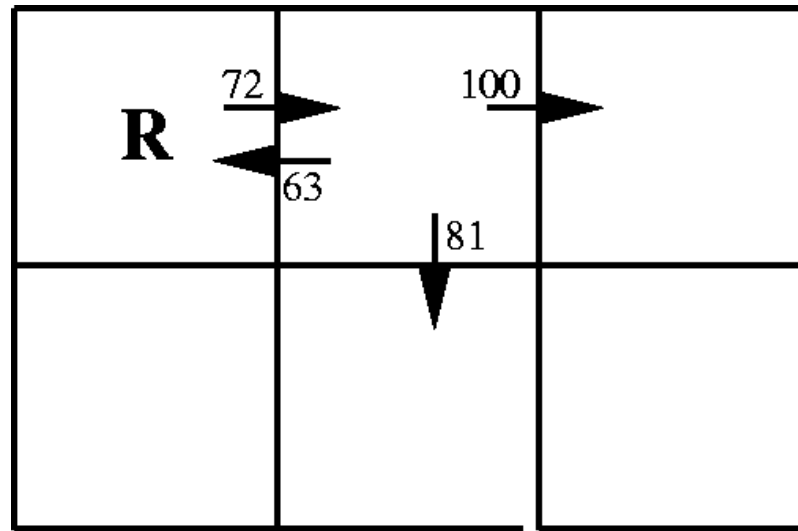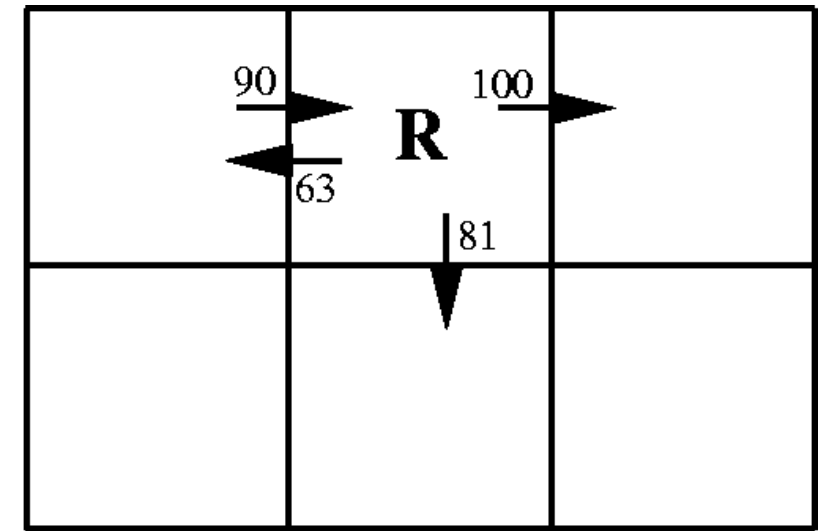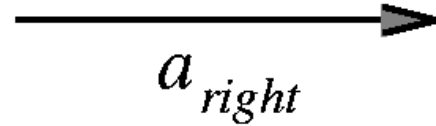
# Q-Learning for Deterministic Worlds

For each $s$, $a$ initialize table entry $Q^{\wedge}(s,a) := 0$.

Observe current state $s$.

Do forever:

1.  Select an action $a$ and execute it

2.  Receive immediate reward $r$

3.  Observe the new state $s'$

4.  Update the table entry for $Q^{\wedge}(s,a)$:
    $$Q^{\wedge}(s,a) := r + \gamma \max_{a'} Q^{\wedge}(s',a')$$

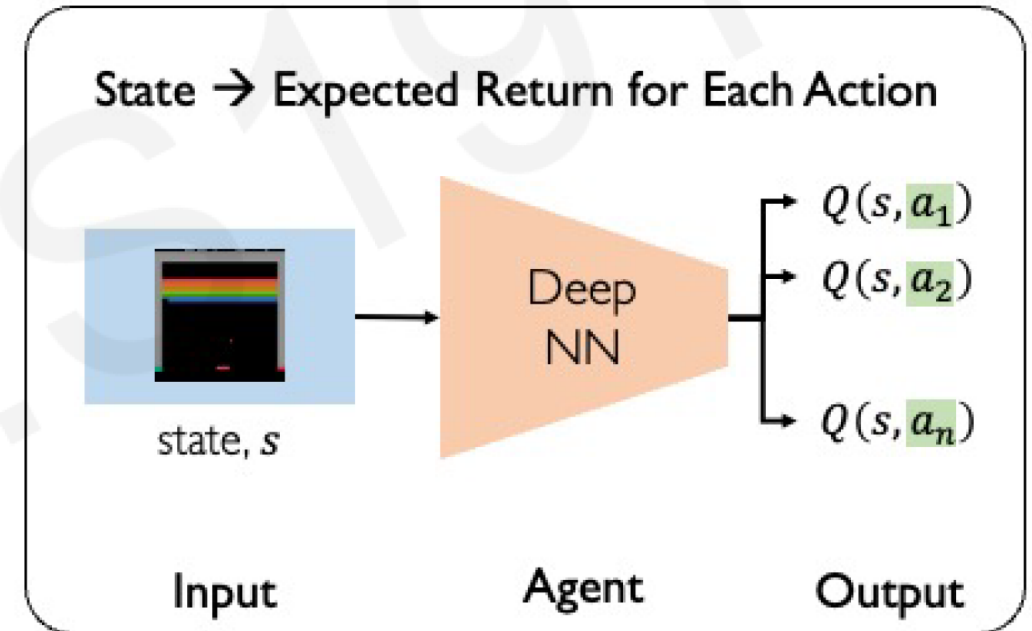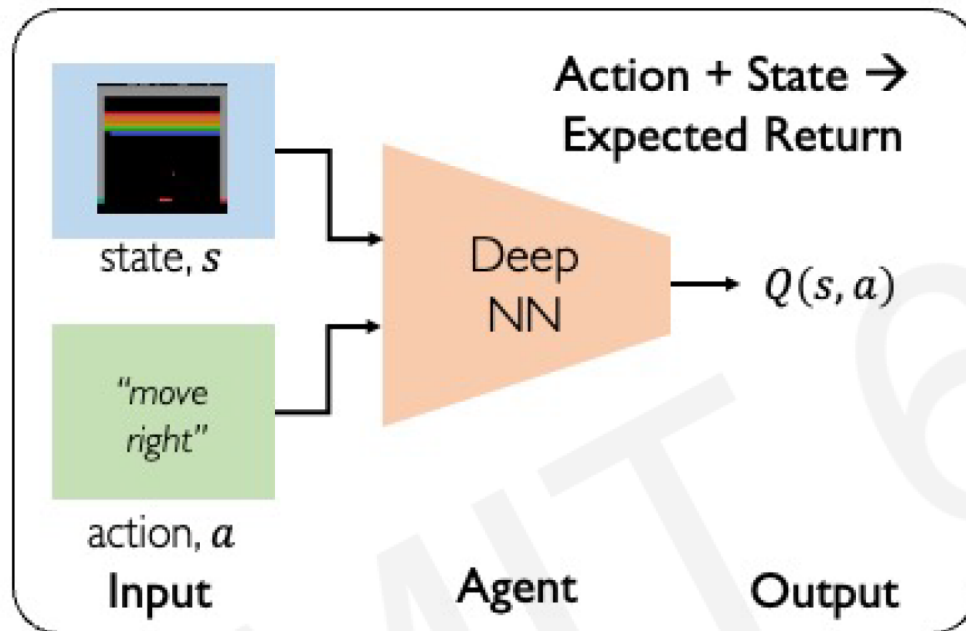5.  $s := s'$

# Q-Learning Example

Initial state: $s_1$

$a_{right}$

Next state: $s_2$

$$Q\,\hat{}\,(s_1, a_{right}) := r + \gamma \max_{a'} Q\,\hat{}\,(s_2, a')$$

$$:= 0 + 0.9 \max\{63, 81, 100\}$$

$$:= 90$$

# Q-Learning Continued

- Exploration
  - Selecting the best action
  - Probabilistic choice

- Improving convergence
  - Update sequences
  - Remember old state-action transitions and their immediate reward

- Non-deterministic MDPs
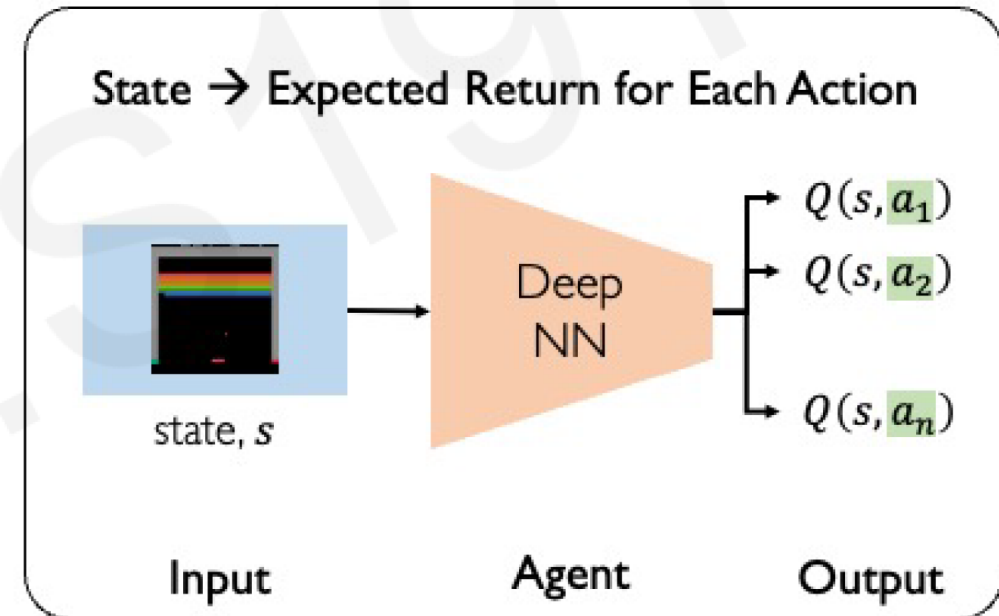
- Temporal Difference Learning

# Deep Q-Learning (DQN)

How can we use deep neural networks to model Q-functions?

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?

Action + State → Expected Return

state, $s$

"move right"

action, $a$

**Input**

Deep NN

**Agent**

$\rightarrow Q(s,a)$

**Output**

State → Expected Return for Each Action

state, $s$

**Input**

Deep NN

**Agent**

$Q(s, a_1)$

$Q(s, a_2)$

$Q(s, a_n)$

**Output**

🤔 What happens if we take all the best actions?
Maximize target return → train the agent

LINKÖPING UNIVERSITY

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?

**Action + State → Expected Return**

state, $s$

"move right"

action, $a$

Input

Deep NN

Agent

$Q(s, a)$

Output

**State → Expected Return for Each Action**

state, $s$

Deep NN

Input

Agent

$Q(s, a_1)$
$Q(s, a_2)$
$Q(s, a_n)$

Output

target

$$\left( r + \gamma \max_{a'} Q(s', a') \right)$$

💡 Take all the best actions → target return

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?

Action + State → Expected Return
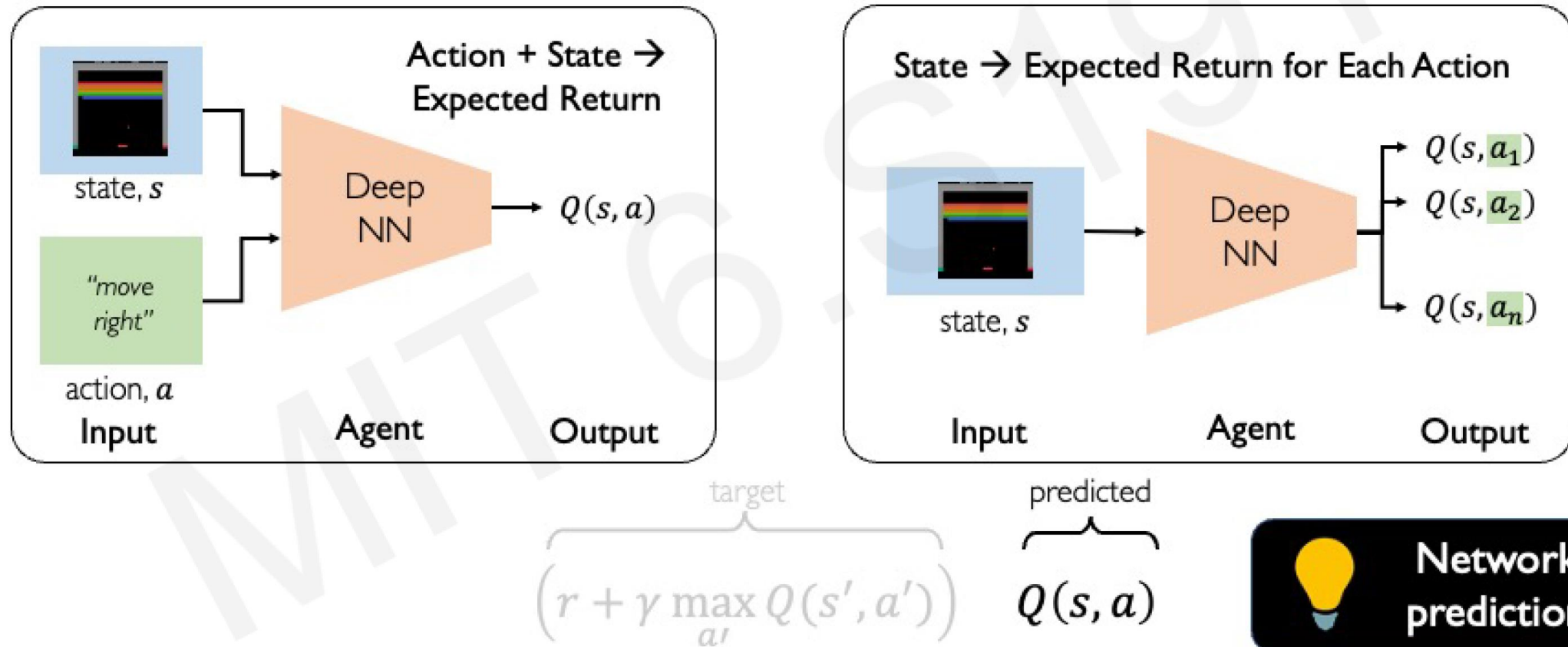
state, $s$

"move right"

action, $a$

**Input**

Deep NN

**Agent**

$Q(s, a)$

**Output**

State → Expected Return for Each Action

state, $s$

Deep NN

**Input**

**Agent**

$Q(s, a_1)$
$Q(s, a_2)$

$Q(s, a_n)$

**Output**

target

$$\left( r + \gamma \max_{a'} Q(s', a') \right)$$

predicted

$$Q(s, a)$$

💡 Network prediction

# Deep Q Networks (DQN): Training

How can we use deep neural networks to model Q-functions?

Action + State → Expected Return

state, $s$

"move right"

action, $a$

Deep NN → $Q(s, a)$

**Input        Agent        Output**

State → Expected Return for Each Action

state, $s$

Deep NN

$Q(s, a_1)$
$Q(s, a_2)$
$Q(s, a_n)$

**Input        Agent        Output**

target — predicted

$$\mathcal{L} = \mathbb{E}\left[\left\|\left(r + \gamma \max_{a'} Q(s', a')\right) - Q(s, a)\right\|^2\right]$$

**Q-Loss**

LINKÖPING UNIVERSITY

# Deep Q Network Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$

state, $s$

Deep NN

$Q(s, a_1) = 20$

$Q(s, a_2) = 3$

$Q(s, a_3) = 0$

$\pi(s) = \underset{a}{\mathrm{argmax}}\, Q(s, a)$

$= a_1$

Send action back to environment and receive next state

LINKÖPING UNIVERSITY

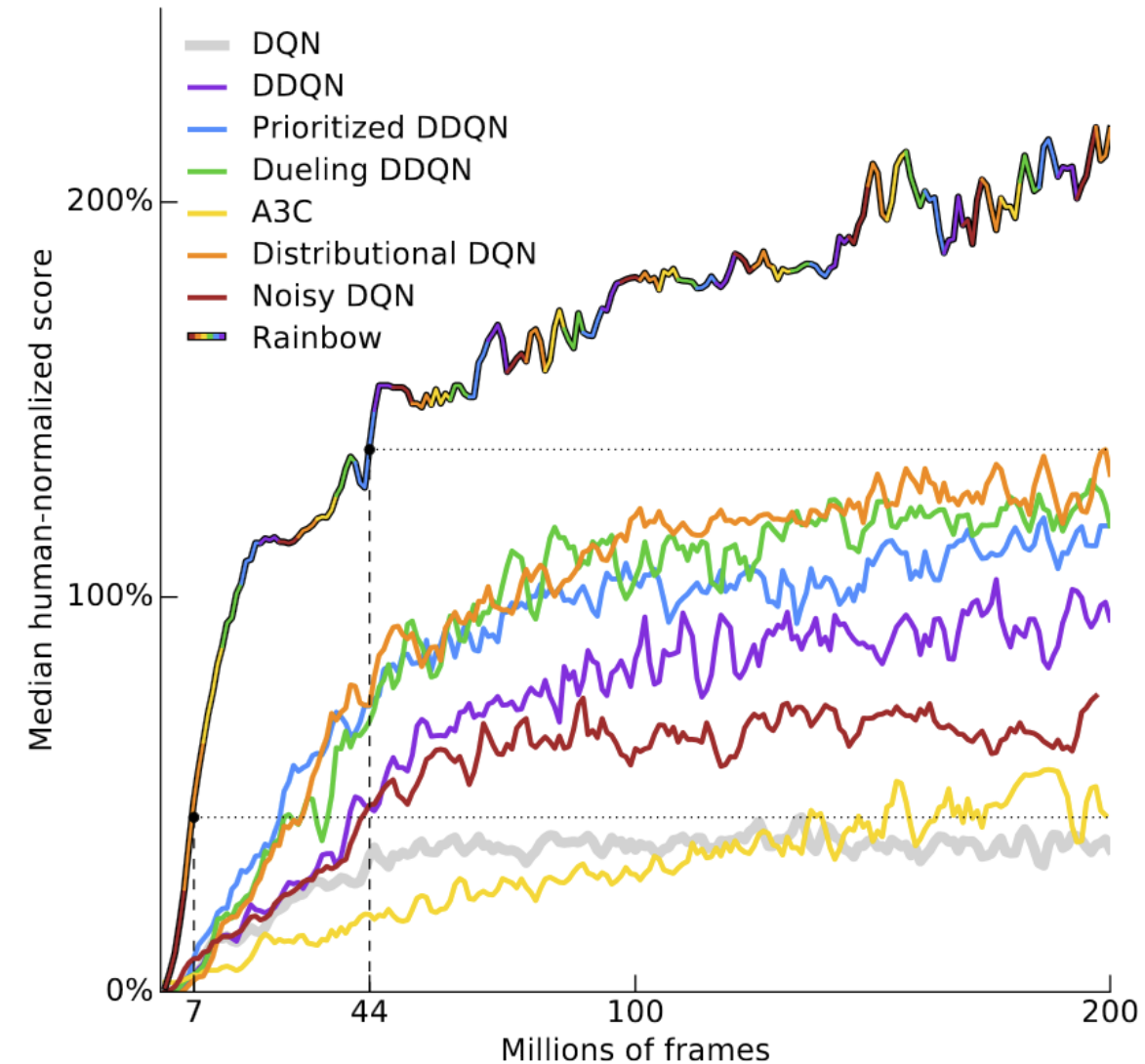# DQN Atari Results

# DQN Atari Results

# Rainbow DQN

- DQN - baseline
- Double DQN - de-overestimate values

$$(R_{t+1}+\gamma_{t+1}q_{\overline{\theta}}(S_{t+1}, \operatorname{argmax}_{a'} q_\theta(S_{t+1},a'))-q_\theta(S_t,A_t))^2$$

- Prioritized experience

$$p_t \propto \left| R_{t+1} + \gamma_{t+1} \max_{a'} q_{\overline{\theta}}(S_{t+1},a') - q_\theta(S_t,A_t) \right|^\omega$$

- Dueling networks

$$q_\theta(s,a) = v_\eta(f_\xi(s)) + a_\psi(f_\xi(s),a) - \frac{\sum_{a'} a_\psi(f_\xi(s),a')}{N_{\text{actions}}}$$

- Distributional DQN - probability distribution
- Noisy DQN - parametric noise
- -> **ADDITIVE**

# Downsides of Q-Learning

**Complexity:**
- Can model scenarios where the action space is discrete and small
- Cannot handle continuous action spaces

**Flexibility:**
- Policy is deterministically computed from the Q function by maximizing the reward → cannot learn stochastic policies

To address these, consider a new class of RL training algorithms:
Policy gradient methods

LINKÖPING
UNIVERSITY

# Reinforcement Learning Algorithms

**Value Learning**

Find $Q(s, a)$

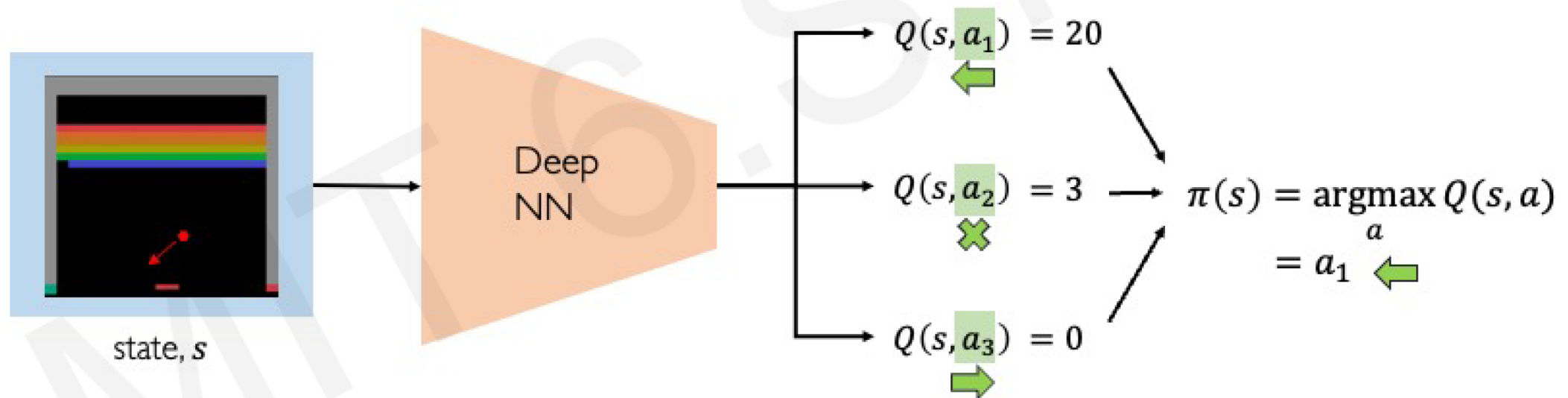$a = \underset{a}{\operatorname{argmax}} Q(s, a)$

**Policy Learning**

Find $\pi(s)$

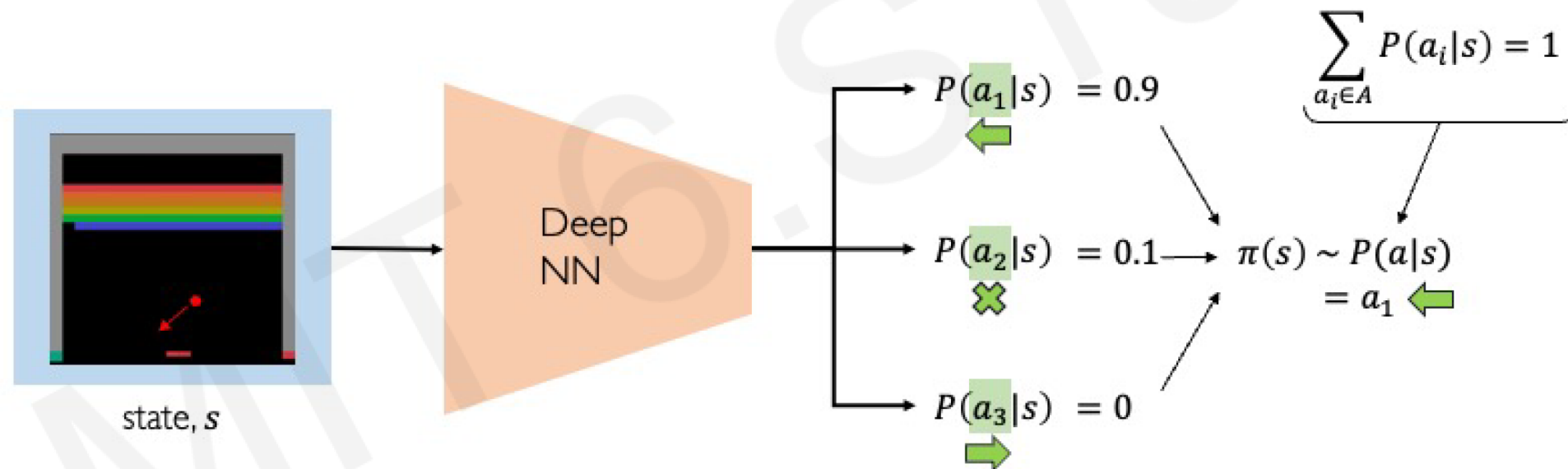Sample $a \sim \pi(s)$

# Deep Q Networks

**DQN:** Approximate Q-function and use to infer the optimal policy, $\pi(s)$

state, $s$ → Deep NN → $Q(s, a_1) = 20$
$Q(s, a_2) = 3$
$Q(s, a_3) = 0$

$\pi(s) = \underset{a}{\mathrm{argmax}}\, Q(s, a)$
$= a_1$

# Policy Gradient (PG): Key Idea

**DQN:** Approximate Q-function and use to infer the optimal policy, $\pi(s)$

**Policy Gradient:** Directly optimize the policy $\pi(s)$



$$\sum_{a_i \in A} P(a_i|s) = 1$$

$$P(a_1|s) = 0.9$$

$$P(a_2|s) = 0.1 \longrightarrow \pi(s) \sim P(a|s)$$
$$= a_1$$

$$P(a_3|s) = 0$$

state, $s$

Deep NN

🤔 What are some advantages of this formulation?

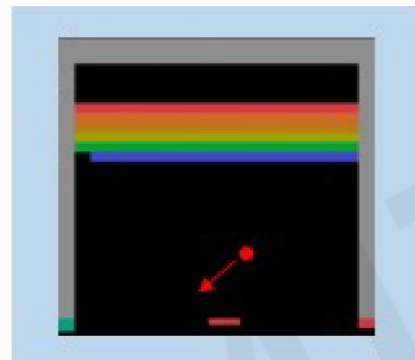# Discrete vs Continuous Action Spaces



state, $s$

# Discrete vs Continuous Action Spaces



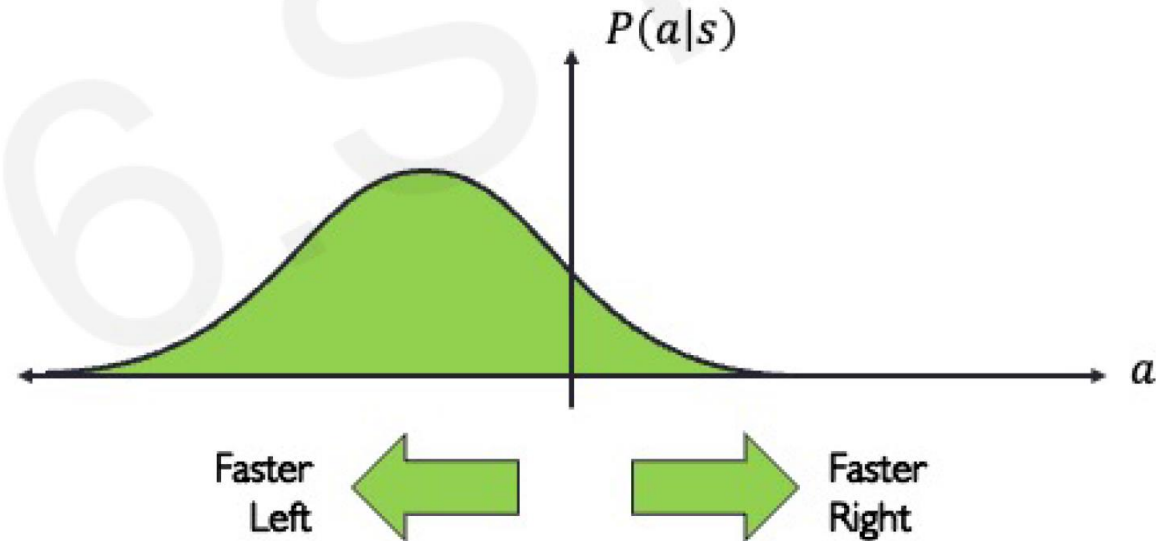**Discrete action space:** which direction should I move?

**Continuous action space:** how fast should I move?    0.7 m/s

$P(a|s)$

state, $s$

Faster Left        Faster Right

$a$

# Policy Gradient (PG): Key Idea

**Policy Gradient:** Enables modeling of continuous action space

$$\int_{a=-\infty}^{\infty} P(a|s) = 1$$

Deep NN

Mean, $\mu = -1$

Variance, $\sigma^2 = 0.5$

$P(a|s) = N(\mu, \sigma^2)$

$\pi(s) \sim P(a|s)$

$= -0.8 \, [m/s]$

state, $s$

$P(a|s) = N(\mu, \sigma^2)$

$a$

$-1$

Faster Left          Faster Right

LINKÖPING UNIVERSITY

# Training Policy Gradients: Case Study

# Training Policy Gradients

## Training Algorithm

1. Initialize the agent

2. Run a policy until termination

3. Record all states, actions, rewards

4. Decrease probability of actions that resulted in low reward

5. Increase probability of actions that resulted in high reward

# Training Policy Gradients

## Training Algorithm

1. Initialize the agent

2. Run a policy until termination

3. Record all states, actions, rewards

4. Decrease probability of actions that resulted in low reward

5. Increase probability of actions that resulted in high reward

# Training Policy Gradients

## Training Algorithm

1. Initialize the agent

2. Run a policy until termination

3. Record all states, actions, rewards

4. Decrease probability of actions that resulted in low reward

5. Increase probability of actions that resulted in high reward

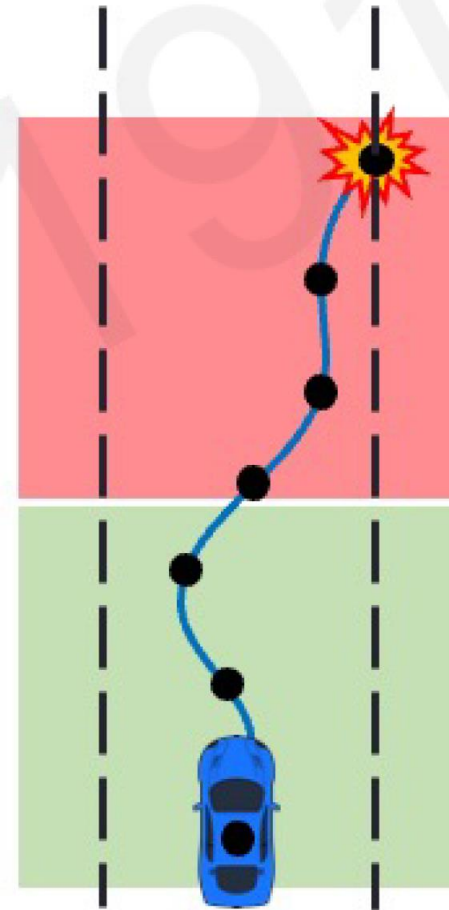# Training Policy Gradients

**Training Algorithm**

1. Initialize the agent

2. Run a policy until termination

3. Record all states, actions, rewards

4. Decrease probability of actions that resulted in low reward

5. Increase probability of actions that resulted in high reward 🤔

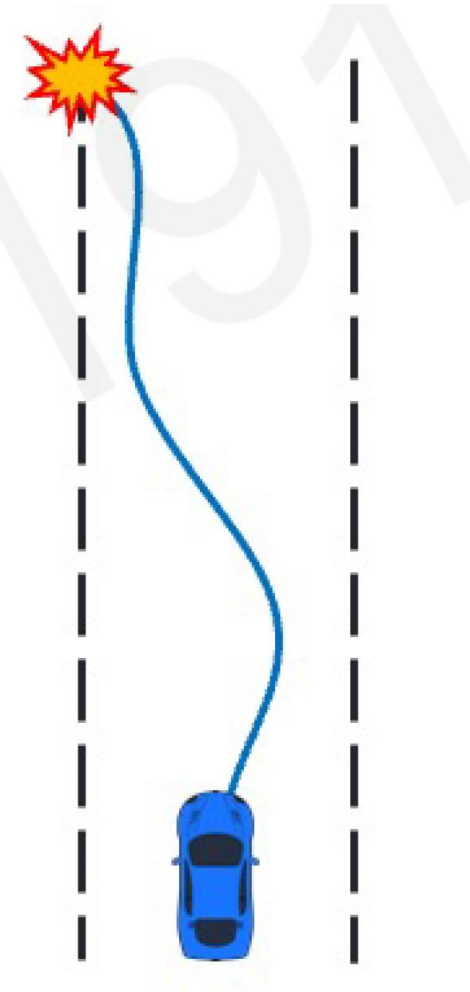# Training Policy Gradients
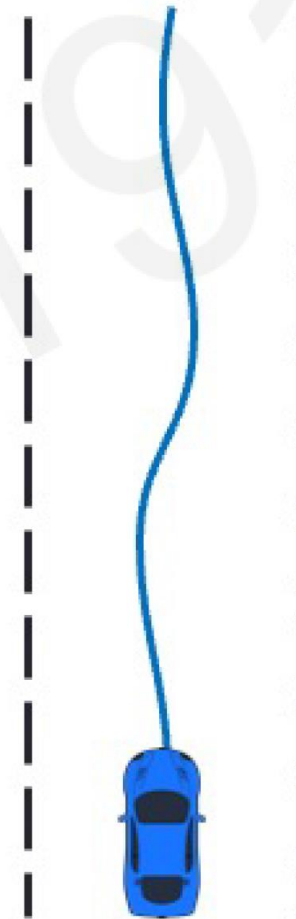
**Training Algorithm**

1. Initialize the agent

2. Run a policy until termination

3. Record all states, actions, rewards

4. Decrease probability of actions that resulted in low reward

5. Increase probability of actions that resulted in high reward

log-likelihood of action

$$\textbf{loss} = -\log \text{P}(a_t|s_t)\, R_t$$

reward

**Gradient descent update:**

$$w' = w - \nabla \textbf{loss}$$
$$w' = w + \nabla \log \text{P}(a_t|s_t)\, R_t$$

Policy gradient!

# REINFORCE

- Take parameterized policy $\pi_{\theta 0}$

- Sample an episode $\tau$ with parameters $\theta_1$

- If it is better, then push parameters in that direction

- If not, then push parameters the other way

- (aka: vanilla policy gradient)

# Policy-Gradient Theorem

$$Policy \ gradient : E_\pi[\nabla_\theta(log\pi(s, a, \theta))R(\tau)]$$

Policy function      Score function

$$Update \ rule : \quad \Delta\theta = \alpha * \nabla_\theta(log\pi(s, a, \theta))R(\tau)$$

Change in parameters        Learning rate

LINKÖPING UNIVERSITY

# REINFORCE

**function REINFORCE**
  Initialise $\theta$ arbitrarily
  **for** each episode $\{s_1, a_1, r_2, ..., s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$ **do**
    **for** $t = 1$ to $T - 1$ **do**
      $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$
    **end for**
  **end for**
  **return** $\theta$
**end function**

LINKÖPING
UNIVERSITY

# AlphaGo Beats Top Human Player (2016)

Human expert positions → Supervised Learning policy network → RL policy network → Self-play data → Value network

Classification

Self Play

Self Play

Regression

1) Initial training: human data

2) Self-play and reinforcement learning
→ super-human performance

3) "Intuition" about board state

# MuZero: Learning Dynamics for Planning (2020)



**AlphaGo** becomes the first program to master Go using neural networks and tree search (Jan 2016, Nature)

**AlphaGo Zero** learns to play completely on its own, without human knowledge (Oct 2017, Nature)

**AlphaZero** masters three perfect information games using a single algorithm for all games (Dec 2018, Science)

**MuZero** learns the rules of the game, allowing it to also master environments with unknown dynamics. (Dec 2020, Nature)

# Deep Reinforcement Learning Summary
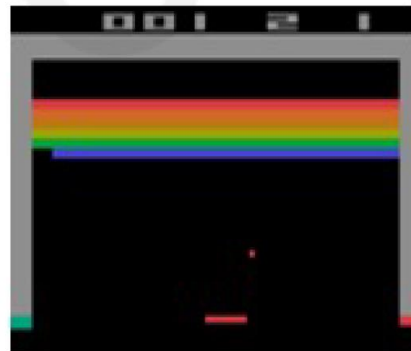
## Foundations

- Agents acting in environment
- State-action pairs → maximize future rewards
- Discounting

## Q-Learning

- Q function: expected total reward given **s, a**
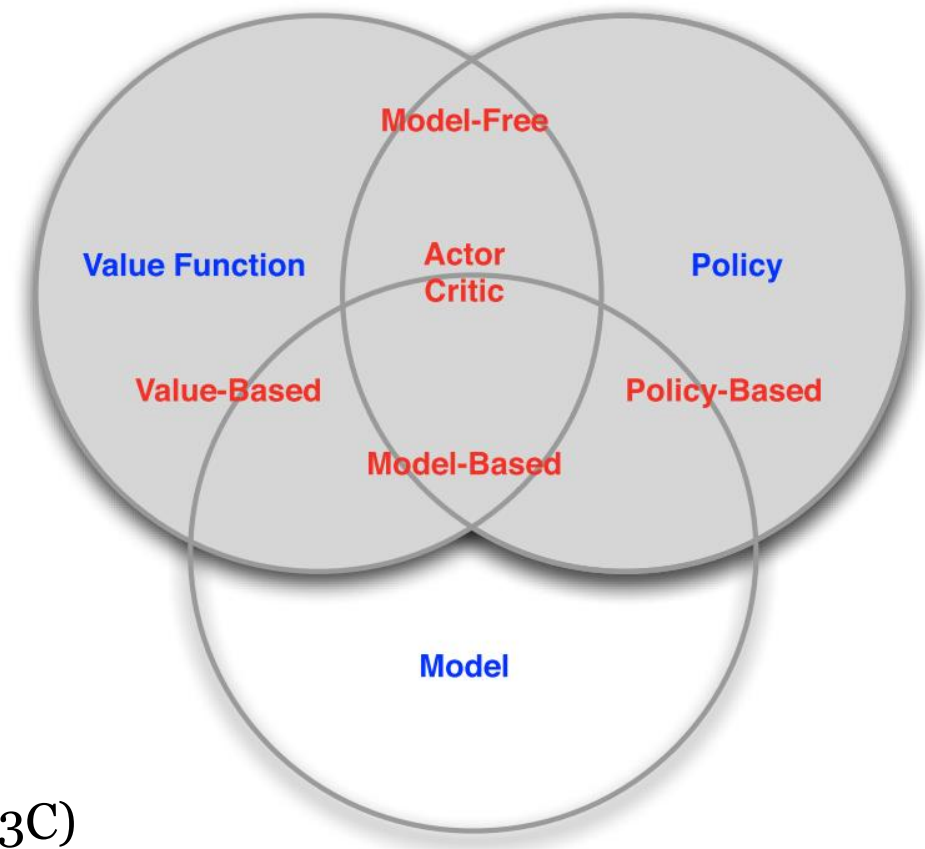- Policy determined by selecting action that maximizes Q function

## Policy Gradients

- Learn and optimize the policy directly
- Applicable to continuous action spaces

LINKÖPING UNIVERSITY

# Reinforcement Learning Approaches

- Value-Based:
  - Learn value function
  - Implicit policy (e.g. greedy selection)
  - Example: Deep Q Networks (DQN)
- Policy-Based:
  - No value function
  - Learn explicit (stochastic) policy
  - Example: Stochastic Policy Gradients
- Model-Based:
  - Learn transition model
  - Implicit policy
  - Example: Dreamer
- Actor-Critic:
  - Learn value function
  - Learn policy using value function
  - Example: Asynchronous Advantage Actor Critic (A3C)

# Model-Based vs Model-Free RL

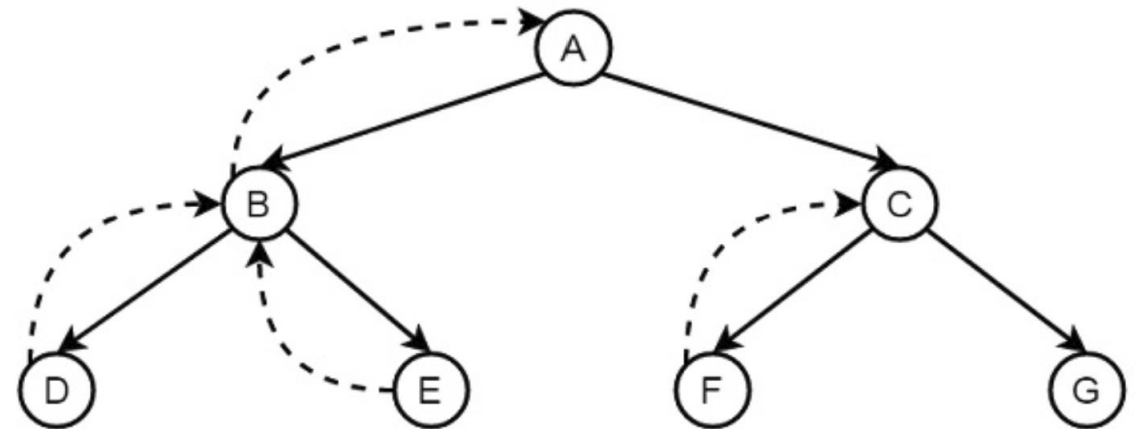

Learn *policy* direct or learn *transition* first and then policy?

# Learning Policies vs Learning Transitions

- s → a → s' → a' → s" → a" → s"' → a"' → s""

- Learning a policy s → a
  - Learning how to react in an environment

- Learning a transition s → a → s'
  - Learning how the environment reacts

# Learning vs Planning

- Learning
  - Agent changing state in the environment
  - Irreversible state change
  - Forward Path s → a → s' → a' → s'' → a'' → s''' → a''' → s''''
- Planning
  - Agent changing own local state
  - Reversible local state change
  - Backtracking Tree

# Model-Based RL



**repeat**
    Sample environment $E$ to generate data $D = (s, a, r', s')$
    Use $D$ to learn $M = T_a(s, s'), R_a(s, s')$             ▷ learning
    **for** $n = 1, \ldots, N$ **do**
        Use $M$ to update policy $\pi(s, a)$             ▷ planning
    **end for**
**until** $\pi$ converges

# Model-Based RL

Learn policy directly

Learn model
and then plan actions

Use experience to
update both model and policy

# Dyna [Sutton]

- Initialize Q-function

- Repeat

  - Initialize s; a←$\pi$(s); (s',r)←Env(s,a)   :: **Learn**

  - Q(s,a) ← Q(s,a)+$\alpha$[r+$\gamma$max$_a$Q(s',$_a$)-Q(s,a)]

  - M(s,a) ← (s',r)    :: **Model**

  - For n=1, …, N :

    - Select $\hat{s}$ and $\hat{a}$ randomly

    - (s',r) ← M($\hat{s}$,$\hat{a}$)     :: **Plan for FREE!**

    - Q($\hat{s}$,$\hat{a}$) ← Q($\hat{s}$,$\hat{a}$)+$\alpha$[r+$\gamma$max$_a$Q(s',$_a$)-Q($\hat{s}$,$\hat{a}$)]

- Until Q converges

- return Q

LINKÖPING
UNIVERSITY

# Example Model-Free RL

- Initialize Q-function

- For All Episodes:

  - Initialize s

  - For All Time Steps in this Episode:

    - Select a $\epsilon$-greedy from Q(s)

    - Perform a in Environment giving s' and r

    - Q(s,a) ← Q(s,a)+$\alpha$[r+$\gamma$max$_a$Q(s',$_a$)-Q(s,a)]

    - s ← s'

- return Q

# Example Model-Based RL

- Initialize Q-function

- Repeat

  - Initialize s; a←π(s); (s',r)←Env(s,a)   :: **Learn**

  - $Q(s,a) \leftarrow Q(s,a) + \alpha[r + \gamma \max_a Q(s',_a) - Q(s,a)]$

  - $M(s,a) \leftarrow (s',r)$   :: **Model**

  - For n=1, ..., N :

    - Select $\hat{s}$ and $\hat{a}$ randomly

    - $(s',r) \leftarrow M(\hat{s},\hat{a})$   :: **Plan for FREE!**

    - $Q(\hat{s},\hat{a}) \leftarrow Q(\hat{s},\hat{a}) + \alpha[r + \gamma \max_a Q(s',_a) - Q(\hat{s},\hat{a})]$

- Until Q converges

- return Q

# Sample Complexity
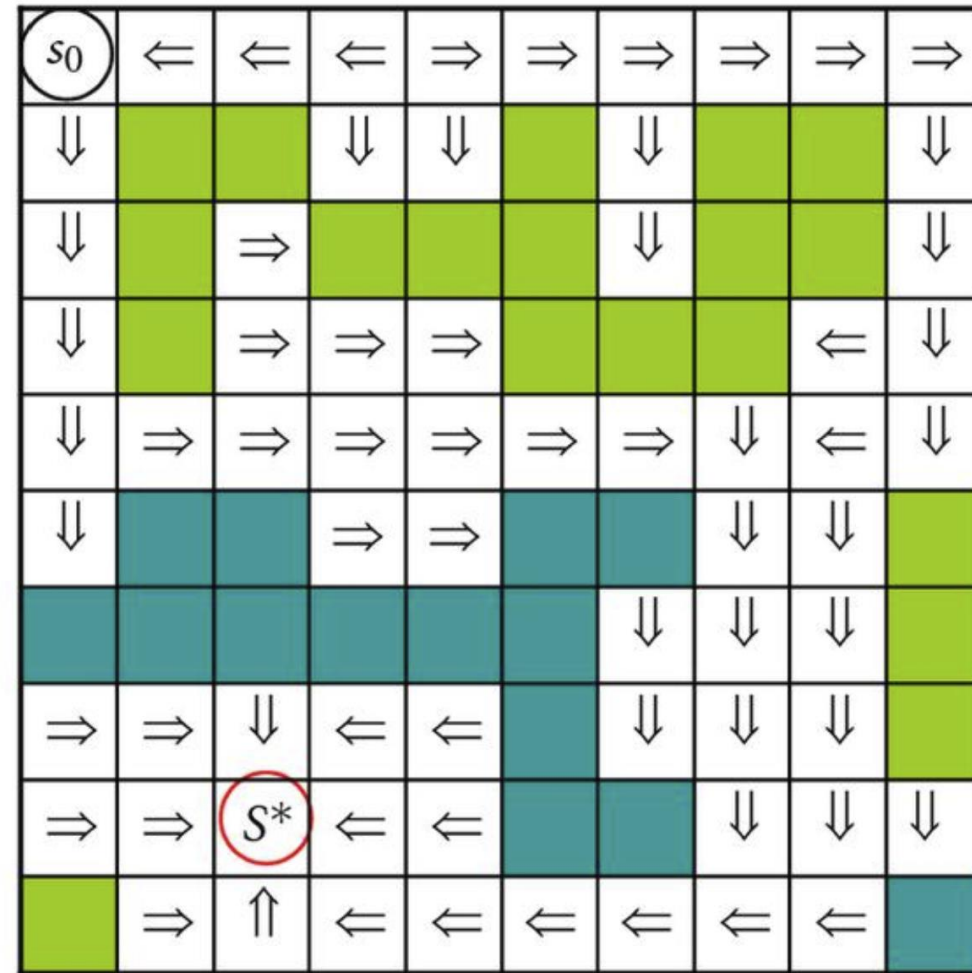
- Model-based RL reduces sample complexity.

- As soon as Model has enough transition entries, the policy can be learned from the Model, for free.

- This free learning is called planning. It does not involve environment samples, hence, "free".

# Dreamer (Latent/Traj)



World Model Learning

Learning Value and Actor Networks

Environment Interaction

The three processes of the Dreamer agent. The world model is learned from past experience. From predictions of this model, the agent then learns a value network to predict future rewards and an actor network to select actions. The actor network is used to interact with the environment.

# Reinforcement Learning Approaches

- Value-Based:
  - Learn value function
  - Implicit policy (e.g. greedy selection)
  - Example: Deep Q Networks (DQN)
- Policy-Based:
  - No value function
  - Learn explicit (stochastic) policy
  - Example: Stochastic Policy Gradients
- Model-Based:
  - Learn transition model
  - Implicit policy
  - Example: Dreamer
- Actor-Critic:
  - Learn value function
  - Learn policy using value function
  - Example: Asynchronous Advantage Actor Critic (A3C)

# Model-Based vs Model-Free RL



Learn *policy* direct or learn *transition* first and then policy?

# Actor-Critic RL

- *An Actor* that controls **how our agent behaves** (policy-based method).

- *A Critic* that measures **how good the action taken is** (value-based method).

- Two ideas to reduce variance
  - Temporal difference bootstrapping
  - Baseline subtraction

# Actor-Critic RL

# Advantage Variants

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau_0 \sim p_\theta(\tau_0)}\left[\sum_{t=0}^{n} \Psi_t \nabla_\theta \log \pi_\theta(a_t|s_t)\right]$$

- Targets:

$$\Psi_t = \hat{Q}_{MC}(s_t, a_t) = \sum_{i=t}^{\infty} \gamma^i \cdot r_i \qquad\qquad \text{Monte Carlo target}$$

$$\Psi_t = \hat{Q}_n(s_t, a_t) \quad = \sum_{i=t}^{n-1} \gamma^i \cdot r_i + \gamma^n V_\theta(s_n) \qquad \text{bootstrap (}n\text{-step target)}$$

$$\Psi_t = \hat{A}_{MC}(s_t, a_t) = \sum_{i=t}^{\infty} \gamma^i \cdot r_i - V_\theta(s_t) \qquad\qquad \text{baseline subtraction}$$

$$\Psi_t = \hat{A}_n(s_t, a_t) \quad = \sum_{i=t}^{n-1} \gamma^i \cdot r_i + \gamma^n V_\theta(s_n) - V_\theta(s_t) \qquad \text{baseline + bootstrap}$$

$$\Psi_t = Q_\phi(s_t, a_t) \qquad\qquad\qquad\qquad\qquad \text{Q-value approximation}$$

LINKÖPING UNIVERSITY

# A3C – Asynchronous Advantage Actor-Critic

- **Asynchronous**: The algorithm is an asynchronous algorithm where multiple worker agents are trained in parallel, each with their environment. This allows the algorithm to train faster as more workers are training in parallel and attain a more diverse training experience as each worker's experience is independent.

- **Advantage**: Advantage is a metric to judge how good its actions were and how they turned out. This allows the algorithm to focus on where the network's predictions were lacking. Intuitively, this will enable it to measure the advantage of taking action, following the policy $\pi$ at the given timestep.

- **Actor-Critic**: The Actor-Critic aspect of the algorithm uses an architecture that shares layers between the policy and value function.

# A3C – Asynchronous Advantage Actor-Critic

1. Fetch the global network parameters

2. Interact with the environment by following the local policy for $n$ number of steps

3. Calculate value and policy loss

4. Get gradients from losses

5. Update the global network

6. Repeat



https://pylessons.com/A3C-reinforcement-learning

# Multi-Objective Reinforcement Learning

LINKÖPING
UNIVERSITY

# Multi-Objective Reinforcement Learning (MORL)

- Many real-world tasks may present an agent with multiple, possibly conflicting objectives:
  - Time
  - Safety
  - Resource consumption
- Multi-Objective Reinforcement Learning allows an agent to learn how to prioritize among objectives at runtime
- Possible to create diverse populations of agents, or adapt agents to time-varying user needs, e.g. difficulty level or training session contents
- Training goals can also be considered by agents

J. Källström and F.Heintz, Tunable Dynamics in Agent-Based Simulation using Multi-Objective Reinforcement Learning, AAMAS Adaptive and Learning Agents Workshop 2019.

# Objectives

# Self-Driving Car Objectives



**1**

**2**

**3**

**4**

**DESTINATION**

Reach the desired destination on time

**SAFETY**

Travel in a safe manner, without any incidents of collisions

**COMFORT**

The trip should be smooth, no excessive breaking or jerking

**CONSUMPTION**

The driving style should minimise fuel comsumption

LINKÖPING UNIVERSITY

# Reward Design

- - 10?
- - 100?
- - 100 000?

How to specify the
reward for a collision?

- This dilemma of reward scale exists for each behaviour we want to encourage:
  - avoiding collisions
  - reaching the destination on time
  - staying within speed limits
  - avoiding sudden changes in speed
  - driving within the lane
  - ...

# Scalar Reward Design Process

1. Design/tweak scalar reward function

2. (Re-)Train RL agent using new/updated reward function (may take hours or days)

3. Evaluate performance (and try to figure out what went wrong!)

- Repeat until the desired agent behaviour is (finally) learned
- Wasteful and time consuming process – each trained agent must be discarded if the reward function changes
- Designers implicitly bake in tradeoffs between different behaviours
- Should AI engineers make the decisions about these tradeoffs?

LINKÖPING
UNIVERSITY

# Scalar Reward Design Process

1. Design/tweak scalar reward function

2. (Re-)Train RL agent using new/updated reward function (may take hours or days)

3. Evaluate performance (and try to figure out what went wrong!)



- GT Sophy - Super Human Racing AI Agent, Sony AI

- Objectives: high precision race car control, efficient racing tactics and maneuvers, while respecting an imprecisely defined racing etiquette

- With enough time and computation, good results can be achieved:

- Could we have done better?

# Multi-Objective Reinforcement Learning

- Vector-valued reward function
  - r = [r_objective1, r_objective2, ...]
- Length of the reward vector =
  ○ number of objectives

State
**Vectorial reward**

Action

# Multi-Objective Reinforcement Learning

- Multi-Objective MDP $< S, A, T, \gamma, \mathbf{R} >$
  - Set of states
  - Set of actions
  - A vectorial reward function $\quad \mathbf{R}: S \times A \times S \to \mathbb{R}^d$
  - $d \geq 2$ objectives
  - Transition function (dynamics of the environment)
  - Discount factor $\gamma \in [0, 1]$

LINKÖPING UNIVERSITY

# Multi-Objective Reinforcement Learning



Example: MOMDP with deterministic state transitions and stochastic rewards

$s_0$

p=0.5; r=(1,0)
p=0.5; r=(3,0)

$s_1$

A r=(1.5, 10)

B r=(1.7, 7)

C r=(2.5, 6)

D r=(3.3, 5)

E r=(3.7, 0)

Vamplew, P., Foale, C., & Dazeley, R. (2022). The impact of environmental stochasticity on value-based multiobjective reinforcement learning. Neural Computing and Applications, 1-17.

LINKÖPING UNIVERSITY

# Value Functions and Policies

- The agent behaves according to a policy:

$$\pi : S \times A \to [0, 1]$$

- The value function of a policy in a MOMDP:

$$\mathbf{V}^\pi = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k \mathbf{r}_{k+1} \mid \pi, \mu\right]$$

where $\mathbf{r}_{k+1} = \mathbf{R}(s_k, a_k, s_{k+1})$

State

**Vectorial reward**

Action

# Deterministic vs Stochastic Policies

- Episodic task, 2 objectives
- 2 deterministic policies: $\pi_i = a_i$



Vamplew, P., Dazeley, R., Barker, E., & Kelarev, A. (2009). Constructing stochastic mixture policies for episodic multiobjective reinforcement learning tasks. In AI 2009: Advances in Artificial Intelligence. Proceedings 22 (pp. 340-349). Springer Berlin Heidelberg.

# Deterministic vs Stochastic Policies

- Episodic task, 2 objectives
- 2 deterministic policies: $\pi_i = a_i$
- Always choosing action a1 will maximise the reward on objective 1, but minimise the reward for objective 2



Vamplew, P., Dazeley, R., Barker, E., & Kelarev, A. (2009). Constructing stochastic mixture policies for episodic multiobjective reinforcement learning tasks. In AI 2009: Advances in Artificial Intelligence. Proceedings 22 (pp. 340-349). Springer Berlin Heidelberg.

# Deterministic vs Stochastic Policies

- Episodic task, 2 objectives
- 2 deterministic policies: $\pi_i = a_i$
- Always choosing action a2 will maximise the reward on objective 2, but minimise the reward for objective 1



Vamplew, P., Dazeley, R., Barker, E., & Kelarev, A. (2009). Constructing stochastic mixture policies for episodic multiobjective reinforcement learning tasks. In AI 2009: Advances in Artificial Intelligence. Proceedings 22 (pp. 340-349). Springer Berlin Heidelberg.

# Deterministic vs Stochastic Policies

- Consider a stochastic policy which selects between actions a1 and a2 with probabilities p1 and (1-p1)

- The average reward received by this policy will be (p1, 1-p1)



Vamplew, P., Dazeley, R., Barker, E., & Kelarev, A. (2009). Constructing stochastic mixture policies for episodic multiobjective reinforcement learning tasks. In AI 2009: Advances in Artificial Intelligence. Proceedings 22 (pp. 340-349). Springer Berlin Heidelberg.

# Value Functions and Policies

- Vectorial value functions now supply only a partial ordering, even for a given state:

$$V_i^\pi(s) > V_i^{\pi'}(s) \text{ but } V_j^\pi(s) < V_j^{\pi'}(s)$$

- We can no longer determine which values are optimal without additional information about how to prioritize the objectives

LINKÖPING
UNIVERSITY

# Utility Functions in Multi-Object Decision Making (MODeM)

- A utility function, $u,$ is used to represent a user's preferences over objectives

  - Utility function maps a vector reward to a scalar utility:

$$u : \mathbb{R}^d \to \mathbb{R}$$

  - For MODeM, a utility function, $u$ is assumed to be monotonically increasing:

$$(\forall o, V_o^\pi \geq V\pi'_o) \implies u\left(\mathbf{V}^\pi\right) \geq u\left(\mathbf{V}^{\pi'}\right)$$

# Utility Functions

- **Linear utility function:**

$$u(\mathbf{V}^\pi) = \mathbf{w}^\top \mathbf{V}^\pi$$

- Each element $\mathbf{w}$ specifies how much one unit of value for the corresponding objective contributes to the scalarised value
- The elements of the weight vector are all positive real numbers and sum to 1

# Utility Functions

- Examples of non-linear utility functions
  - The product utility function
    - seeks to make the objective values as balanced as possible
    - [3, 1] ? [2, 2]

$$u(\mathbf{V}^\pi) = \prod_{o=1}^{d} V_o^\pi$$

# Utility Functions

- Examples of non-linear utility functions
  - The product utility function
    - seeks to make the objective values as balanced as possible
    - [3, 1] ? [2, 2]
  - The sum of squares utility function
    - tends to prioritise achieving higher values on a single objective at the expense of other objectives
    - [3, 1] ? [2, 2]

$$u(\mathbf{V}^\pi) = \prod_{o=1}^{d} V_o^\pi$$

$$u(\mathbf{V}^\pi) = \sum_{o=1}^{d} V_o^{\pi 2}$$

# Utility Functions

- Examples of non-linear utility functions
  - The product utility function
    - seeks to make the objective values as balanced as possible
    - u([3, 1]) < u([2, 2])
  - The sum of squares utility function
    - tends to prioritise achieving higher values on a single objective at the expense of other objectives
    - u([3, 1]) > u([2, 2])

$$u(\mathbf{V}^\pi) = \prod_{o=1}^{d} V_o^\pi$$

$$u(\mathbf{V}^\pi) = \sum_{o=1}^{d} V_o^{\pi 2}$$

# Solution Sets

- In single-objective RL problems, there exist a unique optimal value V, and there can be multiple optimal policies $\pi$ that all have this value

- The goal is to learn one of these optimal policies

- In multi-objective settings there can now be multiple possibly optimal value vectors **V**


- We need to reason about **sets of possibly optimal value vectors and policies** when thinking about solutions to MORL problems

# Solution Sets – Undominated Set

- The most general set of solutions: the undominated set

- The undominated set, U, is the subset of all possible policies Π and associated value vectors for which there exists a possible utility function u with a maximal scalarised value:

$$U(\Pi) = \left\{ \pi \in \Pi \; \middle| \; \exists u, \forall \pi' \in \Pi : u(\mathbf{V}^{\pi}) \geq u(\mathbf{V}^{\pi'}) \right\}$$

# Solution Sets – Coverage Set

- The undominated set may contain excess policies
- We do not need to retain all policies to retain optimal utility

- A set *C S* is a ***coverage set*** if it is a subset of *U* and if, for every *u*, it contains a policy with maximal scalarised value:

$$CS(\Pi) \subseteq U(\Pi) \wedge \left( \forall u, \exists \pi \in CS(\Pi), \forall \pi' \in \Pi : u(\mathbf{V}^\pi) \geq u(\mathbf{V}^{\pi'}) \right)$$

# Solution Sets – Pareto Front

- If the utility function u is any monotonically increasing function, then the **Pareto Front** (PF) is the undominated set:

$$PF(\Pi) = \{\pi \in \Pi \mid \nexists \pi' \in \Pi : \mathbf{V}^{\pi'} \succ_P \mathbf{V}^{\pi}\}$$

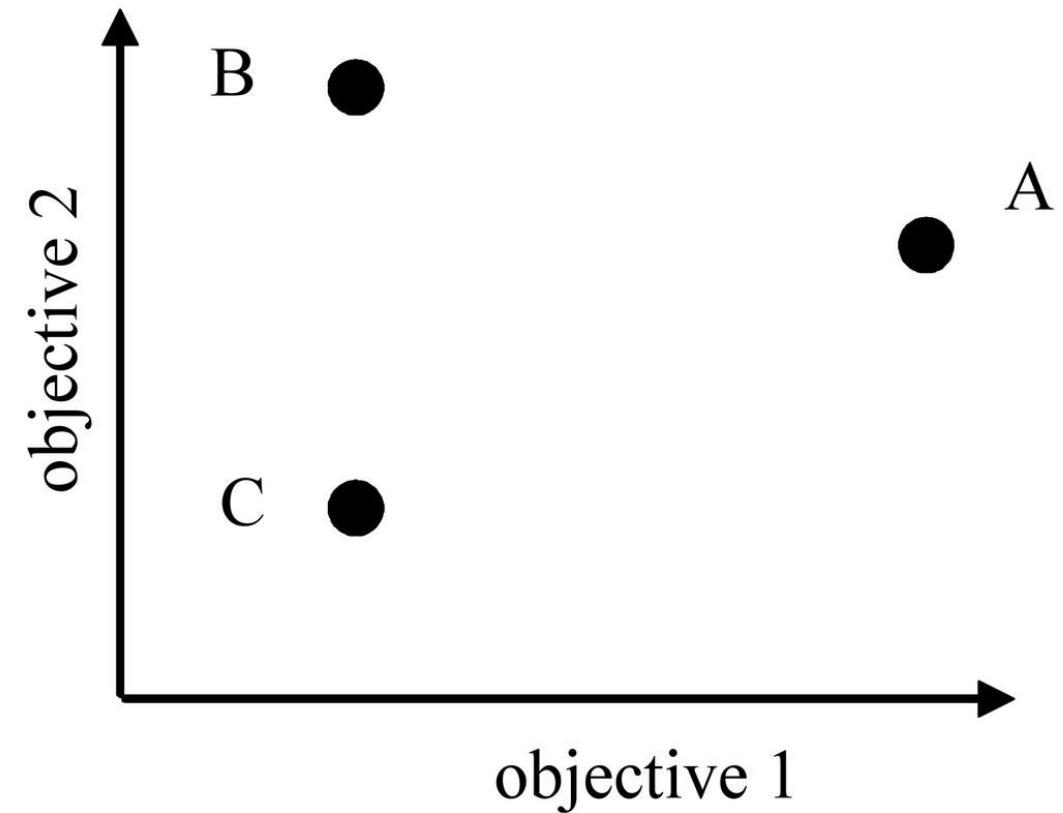where $\succ_P$ is the Pareto dominance relationship:

$$\mathbf{V}^{\pi} \succ_P \mathbf{V}^{\pi'} \iff (\forall i : \mathbf{V}^{\pi}_i \geq \mathbf{V}^{\pi'}_i) \wedge (\exists i : \mathbf{V}^{\pi}_i > \mathbf{V}^{\pi'}_i)$$

LINKÖPING
UNIVERSITY

# Solution Sets – Pareto Coverage Set

- The definition of Pareto dominance corresponds exactly to the definition of monotonically increasing value functions

- Again, we can retain one of the policies that have the same value vector

- A set of policies whose value functions correspond to the PF is called a Pareto Coverage Set (PCS)
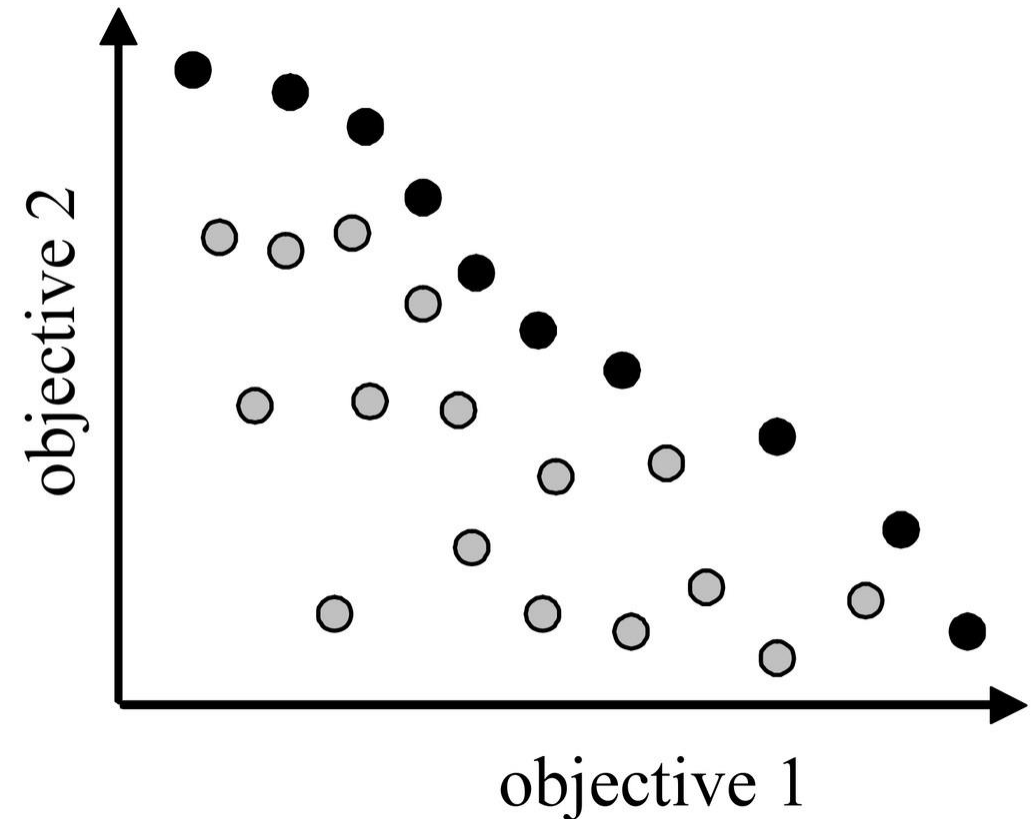
# Solution Sets – Pareto Front

- Pareto dominance illustration, maximising objectives

- Solution A strongly dominates solution C

- Solution B weakly dominates solution C

- A and B are incomparable



Vamplew, P., Dazeley, R., Berry, A., Issabekov, R., & Dekker, E. (2011). Empirical evaluation methods for multiobjective reinforcement learning algorithms. Machine learning, 84, 51-80.

# Solution Sets – Pareto Front

- Black points indicate solutions which form the Pareto front
- Grey solutions are dominated by at least one member of the Pareto front



objective 2

objective 1

Vamplew, P., Dazeley, R., Berry, A., Issabekov, R., & Dekker, E. (2011). Empirical evaluation methods for multiobjective reinforcement learning algorithms. Machine learning, 84, 51-80.
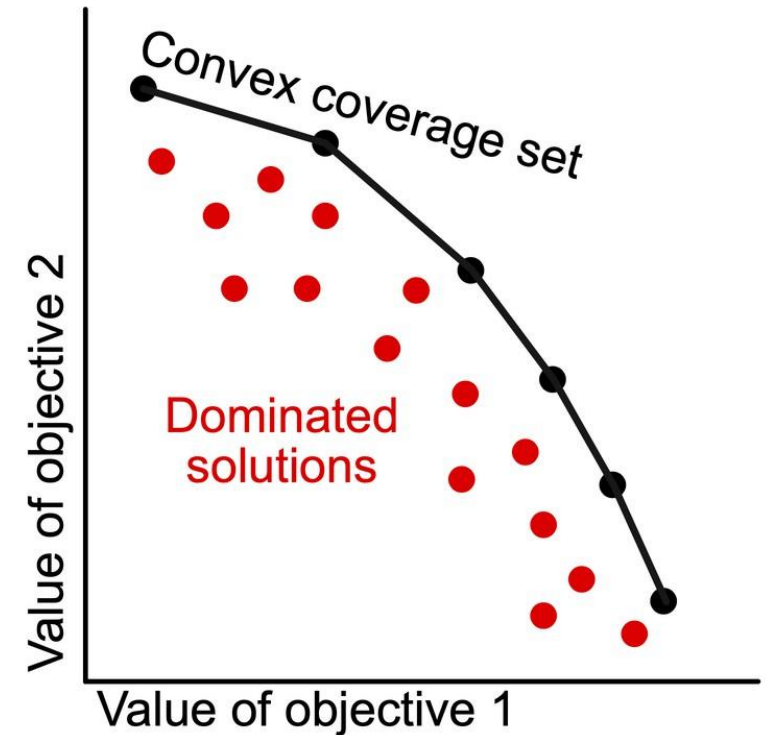
# Solution Sets – Convex Hull

- The convex hull is the undominated set for non-decreasing linear utility functions

- The ***convex hull (CH)*** is the subset of Π for which there exists a **w** (for a linear *u*) for which the linearly scalarised value is maximal:

$$CH(\Pi) = \{\pi \in \Pi \mid \exists \mathbf{w}, \forall \pi' \in \Pi : \mathbf{w}^\top \mathbf{V}^\pi \geq \mathbf{w}^\top \mathbf{V}^{\pi'}\}$$
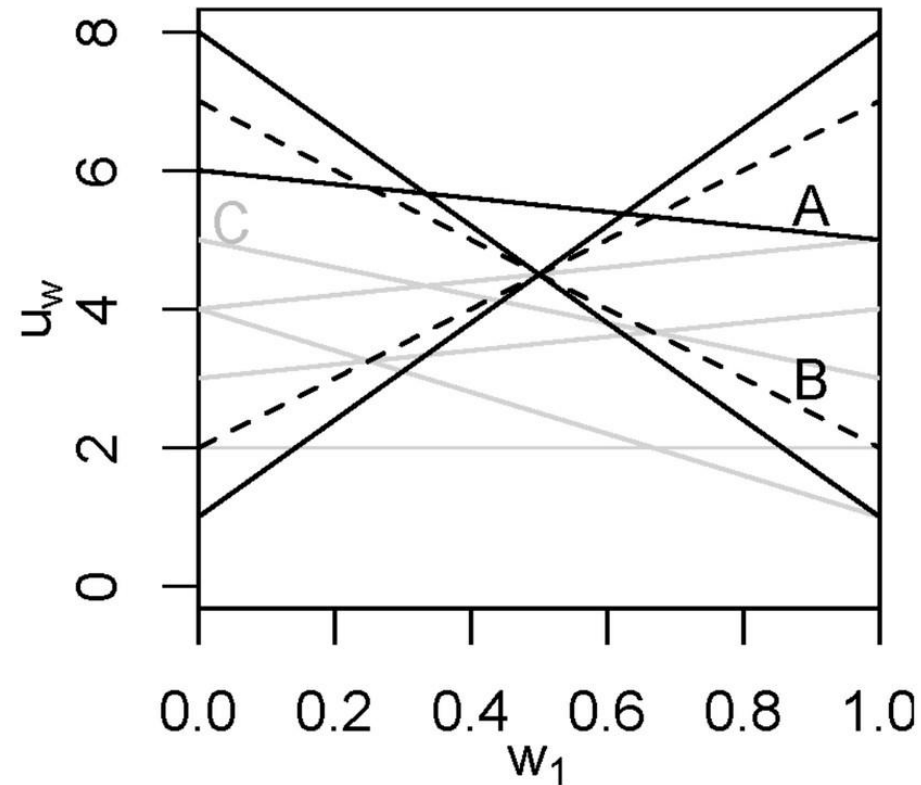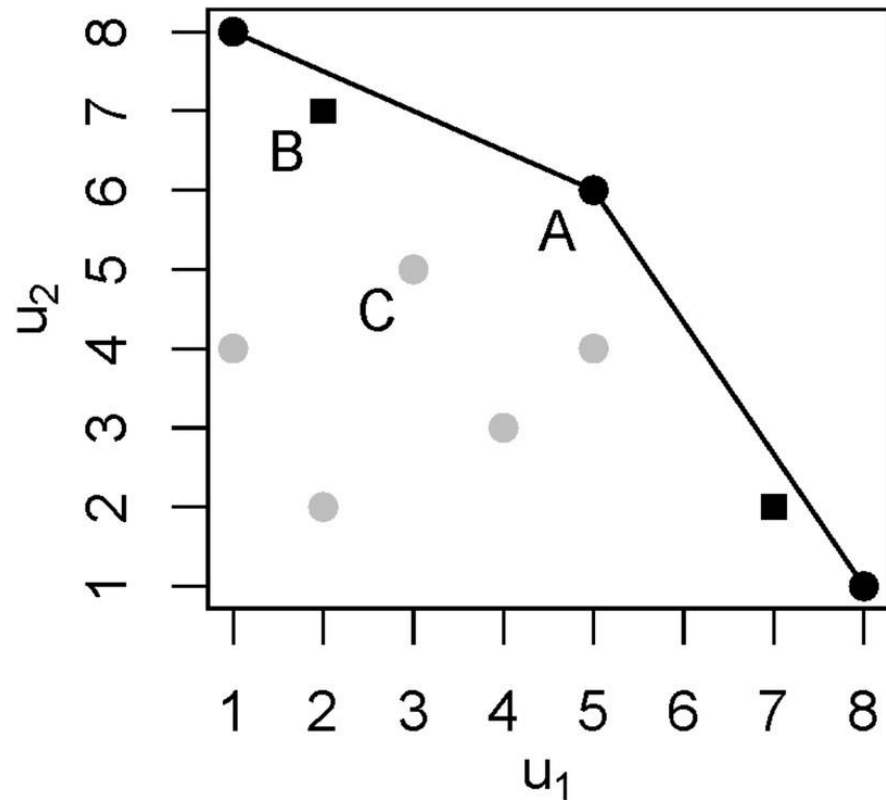
# Solution Sets – Convex Coverage Set

- A set $CCS(\Pi)$ is a **convex coverage set** if it is a subset of $CH(\Pi)$ and if for every $\mathbf{w}$ it contains a policy whose linearly scalarised value is maximal:



$$CCS(\Pi) \subseteq CH(\Pi) \land \left( \forall \mathbf{w}, \exists \pi \in CCS(\Pi), \forall \pi' \in \Pi : \mathbf{w}^\top \mathbf{V}^\pi \geq \mathbf{w}^\top \mathbf{V}^{\pi'} \right)$$
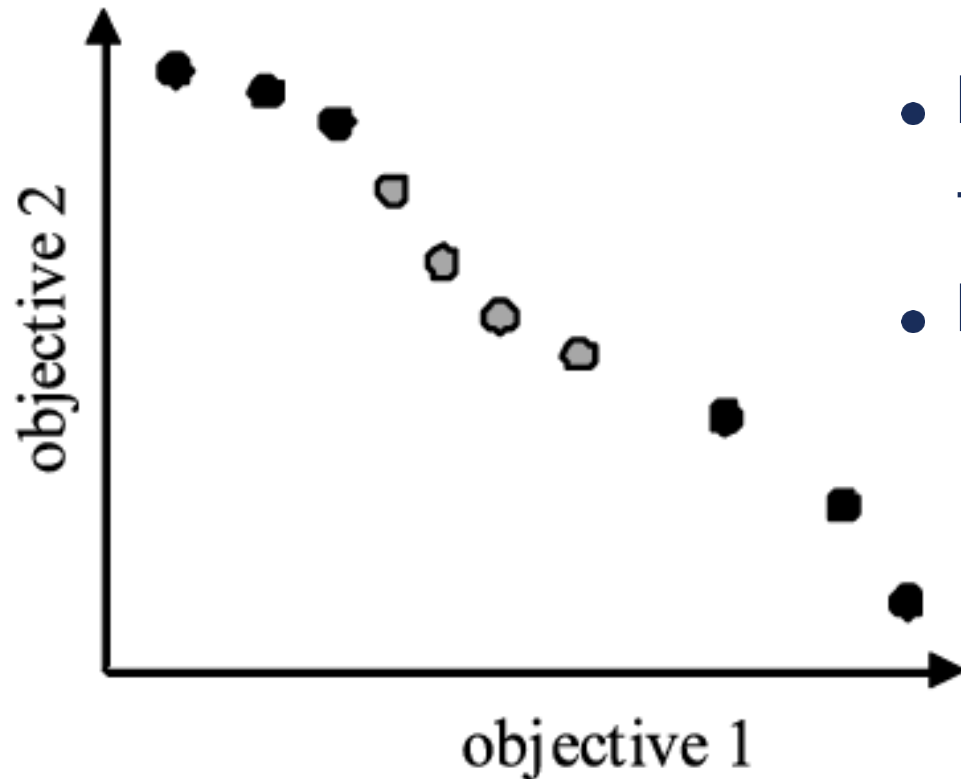
van Doorna, J., Odijkac, D., Roijersab, D. M., & de Rijkea, M. Multi-Objective Optimization for Information Retrieval.

# Solutions Sets – CCS vs PCS



● CCS

■ PCS

Roijers, D. M., Whiteson, S., & Oliehoek, F. A. (2015). Computing convex coverage sets for faster multi-objective coordination. Journal of Artificial Intelligence Research, 52, 399-443.
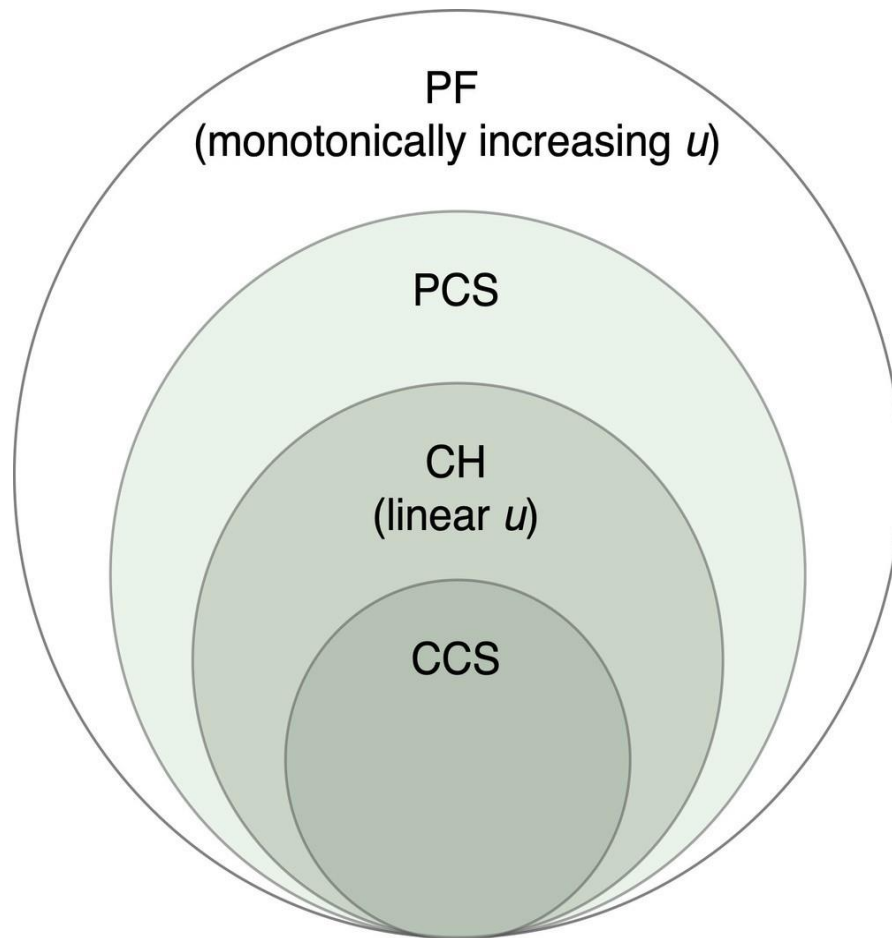
# Solution Sets – Limitation of Linear U



- Pareto front containing a concave region, indicated by the grey points
- Fundamental limitation of linear scalarisation:
  - it cannot find policies which lie in non-convex regions of the Pareto front

Vamplew, P., Dazeley, R., Berry, A., Issabekov, R., & Dekker, E. (2011). Empirical evaluation methods for multiobjective reinforcement learning algorithms. Machine learning, 84, 51-80.

# Solution Sets



- The choice of solution set is key to the efficiency of the algorithms used to solve multi- objective problems

**TDDC17 AI LE7 HT2024:**
**Reinforcement learning**
**Deep reinforcement learning**
**Multi-objective reinforcement learning**

# www.ida.liu.se/~TDDC17

LIU LINKÖPING
UNIVERSITY