# Reinforcement Learning
## Computer Engineering Department
## Sharif University of Technology

**Mohammad Hossein Rohban, Ph.D.**

Spring 2025

# Motivation (cont.) ChatGPT; Why RL?!

# Large-scale Reasoning-Oriented Reinforcement Learning

DeepSeek-v3-Base

Training step 1

DeepSeek-R1-Zero

Solution score (reward)

**Training prompt**

Write python code that takes a list of numbers, returns them in a sorted order, but also adds 42 at the start.

Model checkpoint under training

Generate 4 possible solutions

here's a joke about frogs — Low

echo 42 — Low

def sort(a) ... — Low

def sort_and_prepend(a) ... — High

Update the model so its less likely to output low score solutions like these and more likely to output high-score solutions in response to such a prompt

Courtesy: J. Alammar, The Illustrated DeepSeek R1

## SFT Memorizes, RL Generalizes:
## A Comparative Study of Foundation Model Post-training

# History

| 2013 | Atari (DQN) [Deepmind] |



Pong



Enduro



Beamrider



Q*bert

# A Few Deep RL Highlights

2013 | Atari (DQN)
[Deepmind]

**2014** | **2D locomotion (TRPO)**
**[Berkeley]**



Our algorithm was tested on
three locomotion problems
in a physics simulator

The following gaits were obtained

**Play 0:06 – 0:25**

# History

2013    Atari (DQN)
        [Deepmind]

2014    2D locomotion (TRPO)
        [Berkeley]

2015    **AlphaGo
        [Deepmind]**



**AlphaZero** Silver et al, 2017
Tian et al, 2016; Maddison et al, 2014; Clark et al, 2015

# A Few Deep RL Highlights

| | |
|---|---|
| 2013 | Atari (DQN) [Deepmind] |
| 2014 | 2D locomotion (TRPO) [Berkeley] |
| 2015 | AlphaGo [Deepmind] |
| **2016** | **3D locomotion (TRPO+GAE) [Berkeley]** |



Iteration 0

[Schulman, Moritz, Levine, Jordan, Abbeel, ICLR 2016]

# A Few Deep RL Highlights

| 2013 | Atari (DQN) [Deepmind] |
| 2014 | 2D locomotion (TRPO) [Berkeley] |
| 2015 | AlphaGo [Deepmind] |
| 2016 | 3D locomotion (TRPO+GAE) [Berkeley] |
| **2016** | **Real Robot Manipulation (GPS) [Berkeley]** |

[Levine*, Finn*, Darrell, Abbeel, JMLR 2016]

# History

| Year | Event |
|------|-------|
| 2013 | Atari (DQN) [Deepmind] |
| 2014 | 2D locomotion (TRPO) [Berkeley] |
| 2015 | AlphaGo [Deepmind] |
| 2016 | 3D locomotion (TRPO+GAE) [Berkeley] |
| 2016 | Real Robot Manipulation (GPS) [Berkeley, Google] |
| **2017** | **Dota2 (PPO) [OpenAI]** |



OpenAI Dota Bot beat best humans 1:1 (Aug 2018)

# A Few Deep RL Highlights

| Year | Highlight |
|------|-----------|
| 2013 | Atari (DQN) [Deepmind] |
| 2014 | 2D locomotion (TRPO) [Berkeley] |
| 2015 | AlphaGo [Deepmind] |
| 2016 | 3D locomotion (TRPO+GAE) [Berkeley] |
| 2016 | Real Robot Manipulation (GPS) [Berkeley, Google] |
| 2017 | Dota2 (PPO) [OpenAI] |
| **2018** | **DeepMimic [Berkeley]** |



[Peng, Abbeel, Levine, van de Panne, 2018]

# History

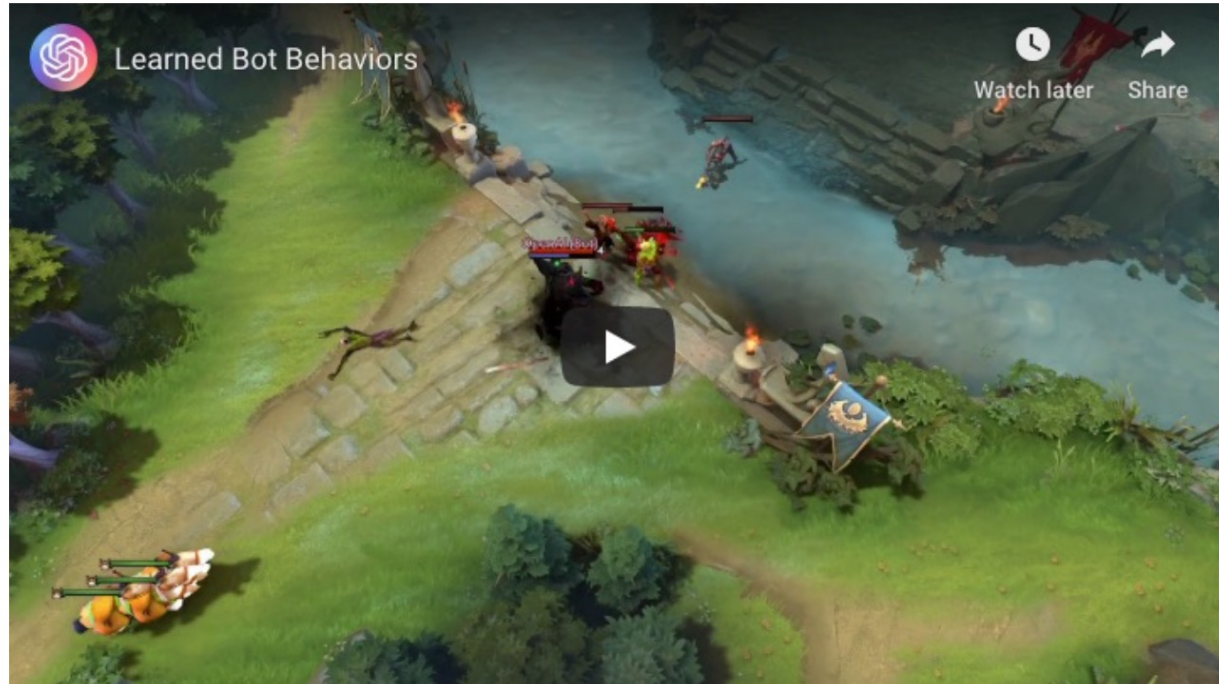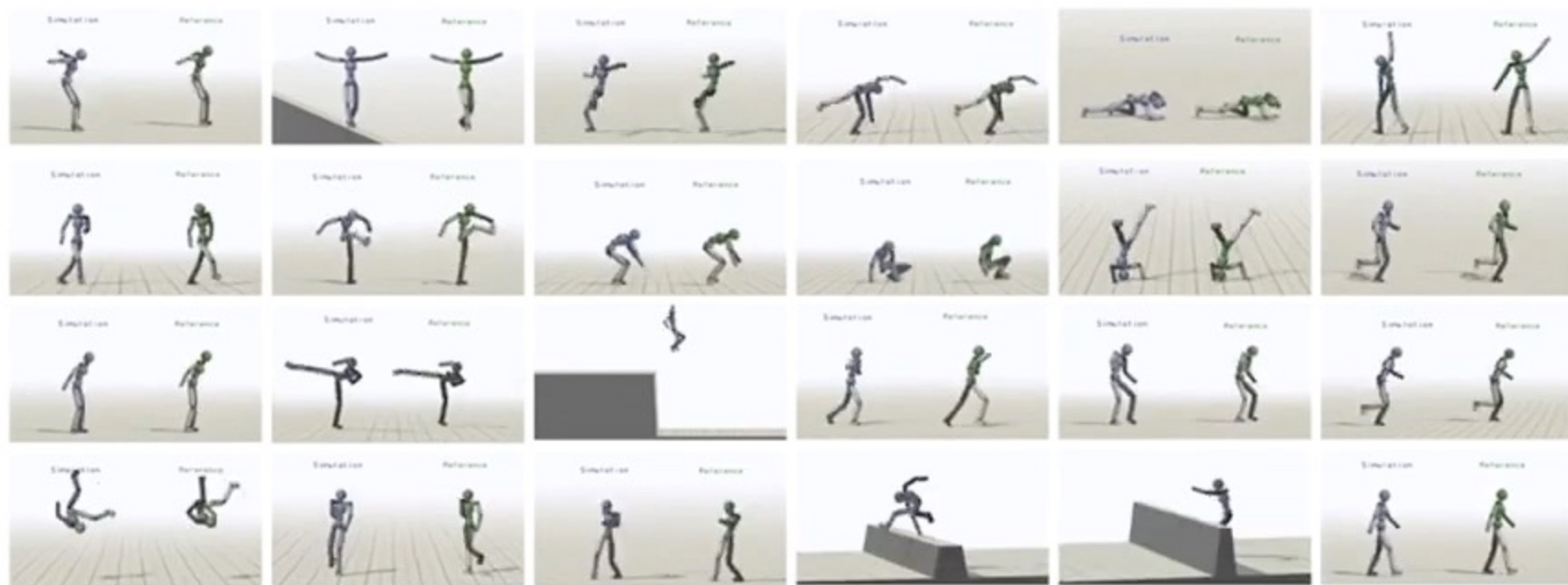| | |
|---|---|
| 2013 | Atari (DQN) [Deepmind] |
| 2014 | 2D locomotion (TRPO) [Berkeley] |
| 2015 | AlphaGo [Deepmind] |
| 2016 | 3D locomotion (TRPO+GAE) [Berkeley] |
| 2016 | Real Robot Manipulation (GPS) [Berkeley, Google] |
| 2017 | Dota2 (PPO) [OpenAI] |
| 2018 | DeepMimic [Berkeley] |
| **2019** | **AlphaStar [Deepmind]** |

# A Few Deep RL Highlights

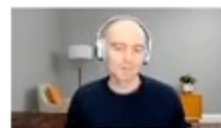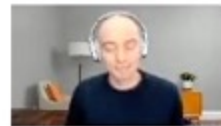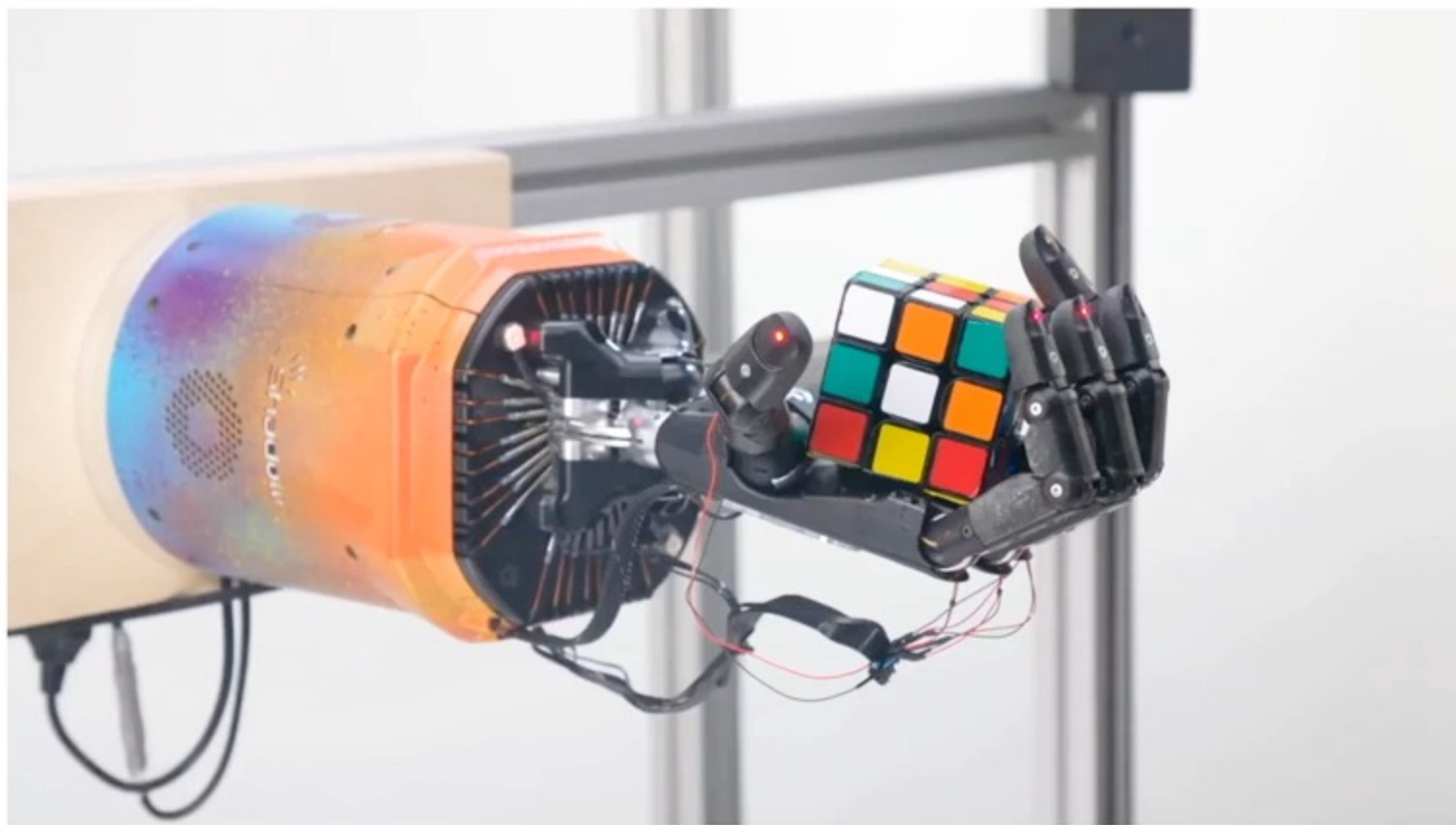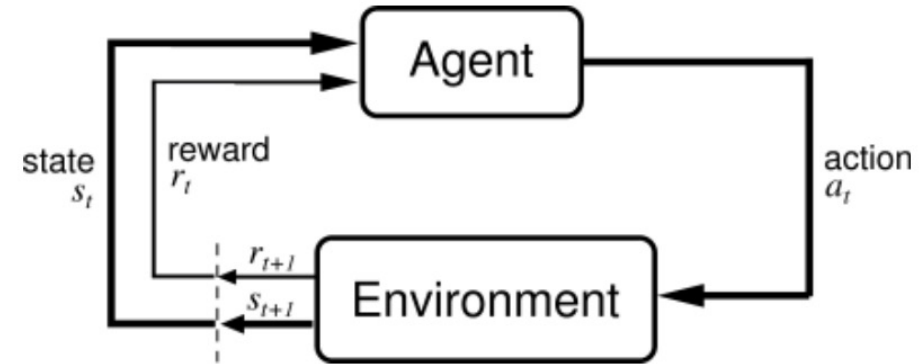| | |
|---|---|
| 2013 | Atari (DQN) [Deepmind] |
| 2014 | 2D locomotion (TRPO) [Berkeley] |
| 2015 | AlphaGo [Deepmind] |
| 2016 | 3D locomotion (TRPO+GAE) [Berkeley] |
| 2016 | Real Robot Manipulation (GPS) [Berkeley, Google] |
| 2017 | Dota2 (PPO) [OpenAI] |
| 2018 | DeepMimic [Berkeley] |
| 2019 | AlphaStar [Deepmind] |
| **2019** | **Rubik's Cube (PPO+DR) [OpenAI]** |

# Let's Begin: Markov Decision Processes (MDPs)

An MDP is defined by:

- Set of states $S$

- Set of actions $A$

- Transition function $P(s' \mid s, a)$

- Reward function $R(s, a, s')$

- Start state $s_0$

- Discount factor $\gamma$

- Horizon $H$

# The Goal

- The policy is $\pi_\theta : S \to A$ for infinite horizon or

  $\pi_\theta : S \times \{0, \ldots, H\} \to A$ for finite horizon MDP.

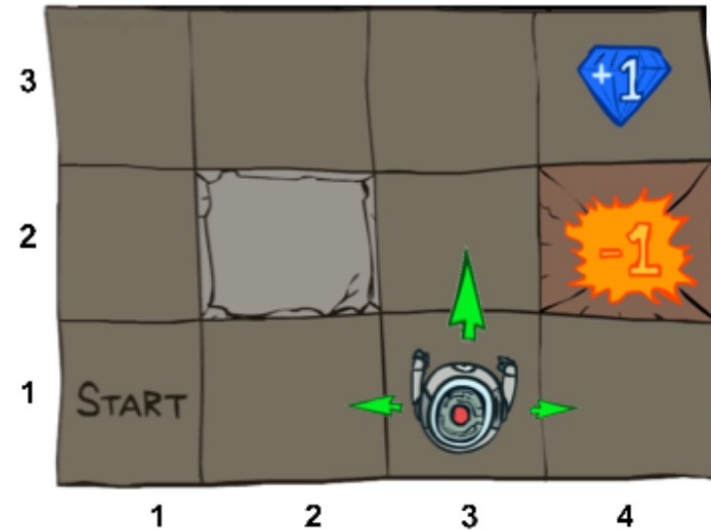MDP (S, A, T, R, γ, H),          goal:     $max_\pi \mathrm{E}[\sum_{t=0}^{H} \gamma^t R(S_t, A_t, S_{t+1}) | \pi]$

Sometimes the policy could be stochastic: $\pi : S \times A \to [0,1]$, which is

$\pi(a|s) = \Pr(A_t = a | S_t = s)$.
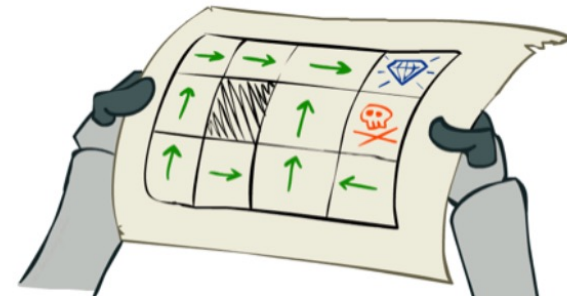
# Example: Grid World

An MDP is defined by:

- Set of states $S$

- Set of actions $A$

- Transition function $P(s' \mid s, a)$

- Reward function $R(s, a, s')$

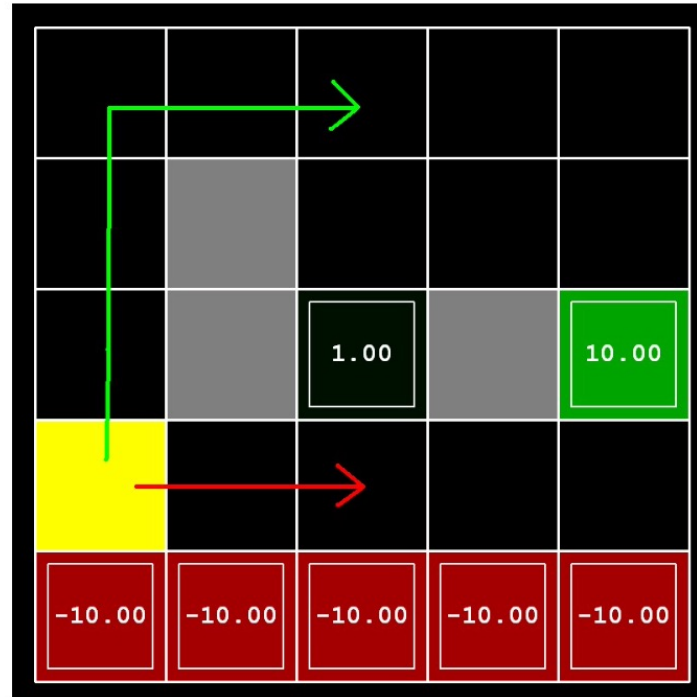- Start state $s_0$

- Discount factor $\gamma$

- Horizon $H$

Goal: $$max_\pi \mathrm{E}\left[\sum_{t=0}^{H} \gamma^t R(S_t, A_t, S_{t+1}) \mid \pi\right]$$

$\pi$:

# Exercise



(a) Prefer the close exit (+1), risking the cliff (-10)

(b) Prefer the close exit (+1), but avoiding the cliff (-10)

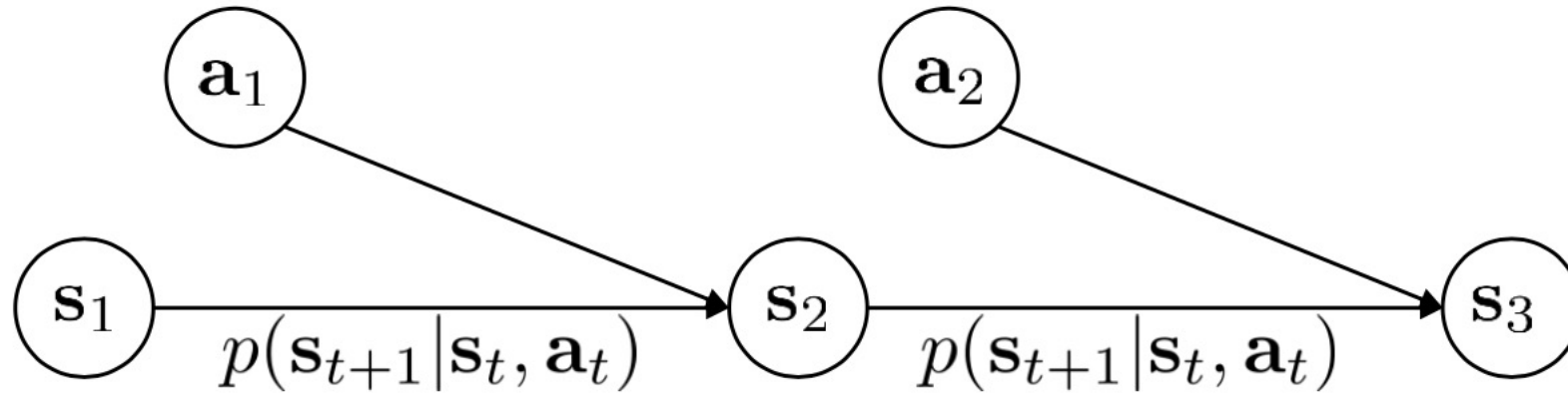(c) Prefer the distant exit (+10), risking the cliff (-10)

(d) Prefer the distant exit (+10), avoiding the cliff (-10)

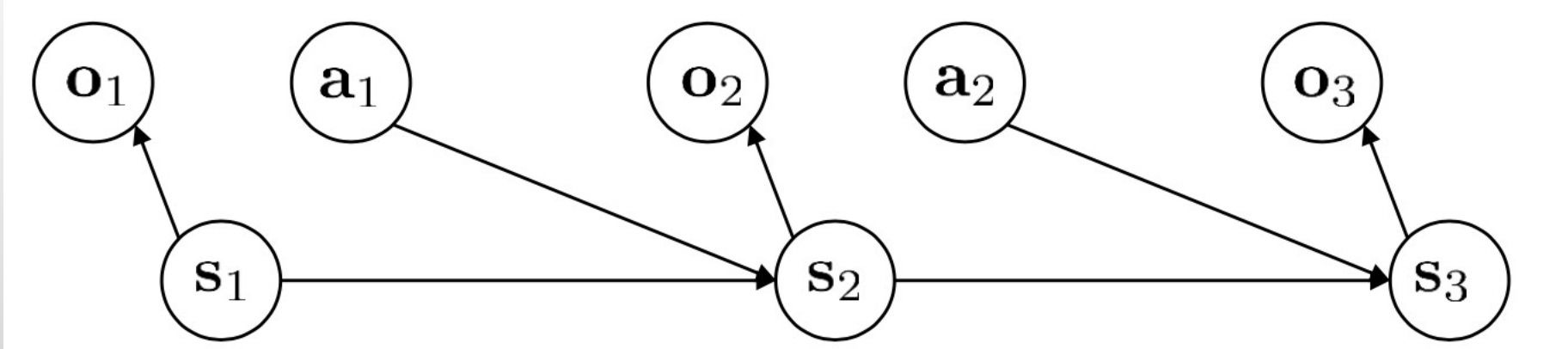(1) γ = 0.1, noise = 0.5

(2) γ = 0.99, noise = 0

(3) γ = 0.99, noise = 0.5

(4) γ = 0.1, noise = 0

# Graphical Model of MDPs

# Partially Observable MDPs (POMDPs)
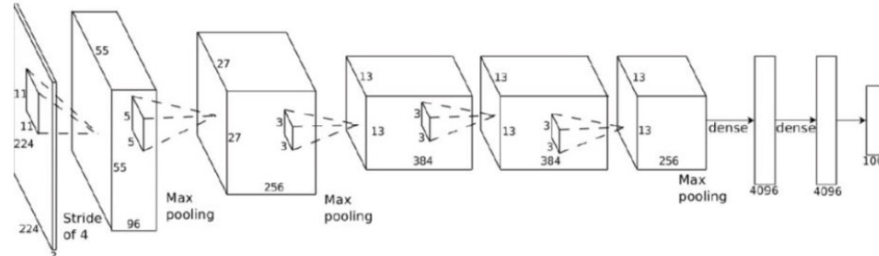
- Often times the state $S_t$ is hidden from the agent,

  and only noisy or incomplete measurement of it is available $O_t$.

$$\mathbf{o}_t \qquad \pi_\theta(\mathbf{a}_t|\mathbf{o}_t) \qquad \mathbf{a}_t$$

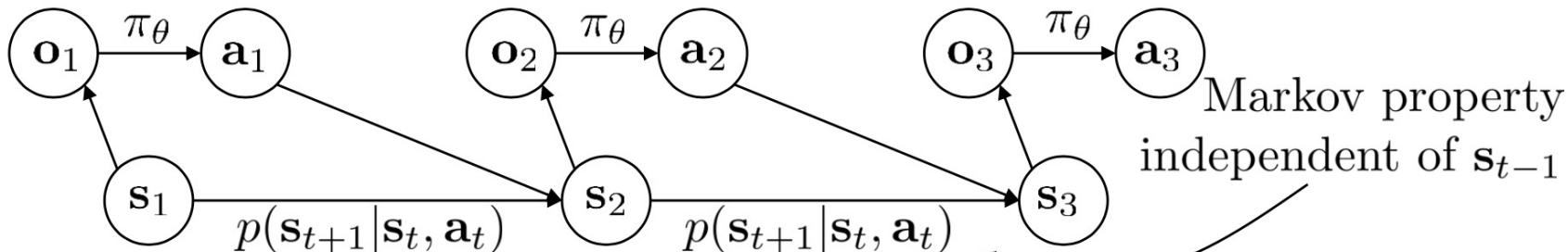$\mathbf{s}_t$ – state

$\mathbf{o}_t$ – observation

$\mathbf{a}_t$ – action

$\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ – policy

$\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ – policy (fully observed)



Markov property independent of $\mathbf{s}_{t-1}$

# Optimal Value Function

MDP (S, A, T, R, γ, H),  goal:  $max_\pi \mathrm{E}[\sum_{t=0}^{H} \gamma^t R(S_t, A_t, S_{t+1}) | \pi]$

$$V^*(s) = \max_\pi \mathbb{E}\left[\sum_{t=0}^{H} \gamma^t R(s_t, a_t, s_{t+1}) \mid \pi, s_0 = s\right]$$

= sum of discounted rewards when starting from state s and acting optimally

# Optimal Value Function

$$V^*(s) = \max_{\pi} \mathbb{E} \left[ \sum_{t=0}^{H} \gamma^t R(s_t, a_t, s_{t+1}) \mid \pi, s_0 = s \right]$$

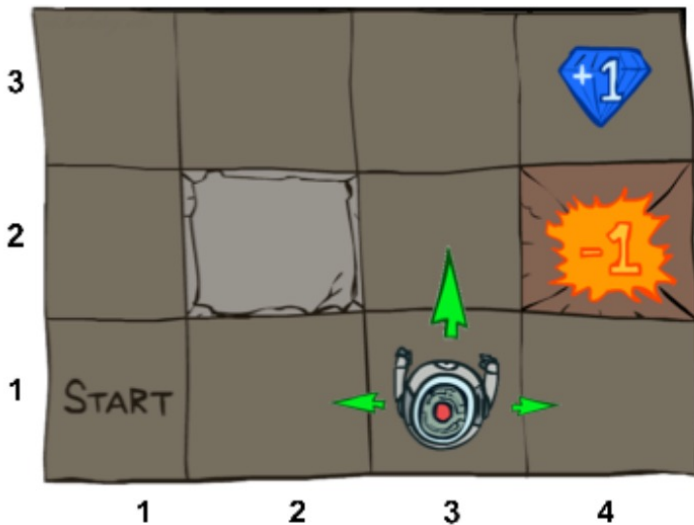= sum of discounted rewards when starting from state s and acting optimally



Let's assume:
actions deterministically successful,  gamma = 1, H = 100

V*(4,3) =

V*(3,3) =

V*(2,3) =

V*(1,1) =

V*(4,2) =

# Optimal Value Function

$$V^*(s) = \max_{\pi} \mathbb{E}\left[\sum_{t=0}^{H} \gamma^t R(s_t, a_t, s_{t+1}) \mid \pi, s_0 = s\right]$$

= sum of discounted rewards when starting from state s and acting optimally



Let's assume:
actions deterministically successful,  gamma = 0.9, H = 100
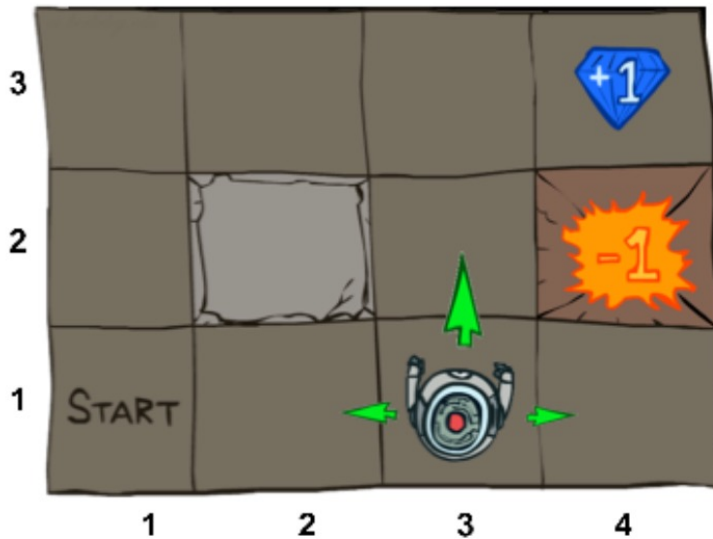
V*(4,3) =

V*(3,3) =
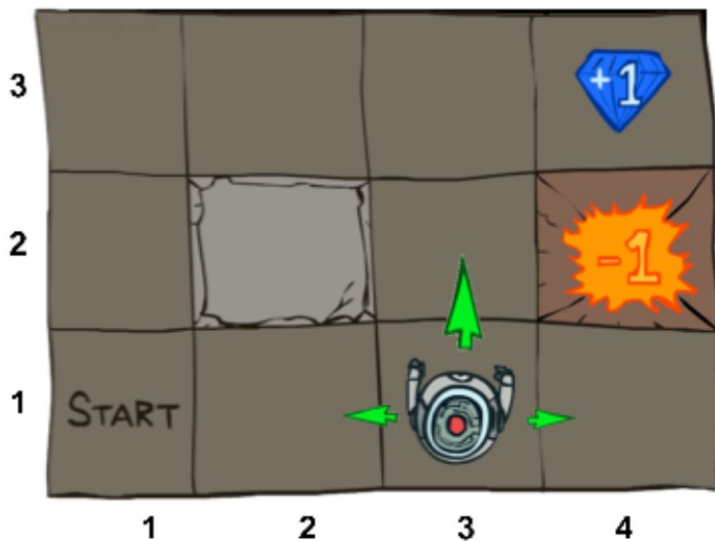
V*(2,3) =

V*(1,1) =

V*(4,2) =

# Optimal Value Function

$$V^*(s) = \max_{\pi} \mathbb{E} \left[ \sum_{t=0}^{H} \gamma^t R(s_t, a_t, s_{t+1}) \mid \pi, s_0 = s \right]$$

= sum of discounted rewards when starting from state s and acting optimally

Let's assume:
actions successful w/probability 0.8,  gamma = 0.9, H = 100

V*(4,3) =

V*(3,3) =

V*(2,3) =

V*(1,1) =

V*(4,2) =