

# Reinforcement Learning workshop

Florian Goebels Sep 2019

September 19th 2019



# Agenda

- Day 1: RL Basics (theory)
- Day 2: Q-learning (coding)
- Day 3: Policy Gradient methods (coding)

# Today's Agenda

- Day 2: Q-Learning
  - Q-Learning
  - Deep Q-Learning
  - Rainbow

# Q-Learning

- Combines ideas from DP and MC methods
  - Like MC: learns directly from raw experiences
  - Like DP: iteratively updates estimates

# TD $V(S)$ estimate

**MC**

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

**DP**

$$v_{\pi}(s) = \mathbf{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

**TD**

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

# TD $V(S)$ estimate

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$



Better estimate of  $V(S)$

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

TD ERROR

Current estimate of  $V(S)$

# TD advantages

- Does not require model of environment (unlike DP)
- MC needs to wait until episode finish, TP can online update
  - MC it's hard to estimate value of action-state pair
- Both TD and MC converge to optimal policy?

# TD demo

- <https://cs.stanford.edu/people/karpathy/reinforcejs/>



# Q-learning

Q-learning (off-policy TD control) for estimating  $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

    until  $S$  is terminal

# Deep Q-Learning

# Deep Q-Learning

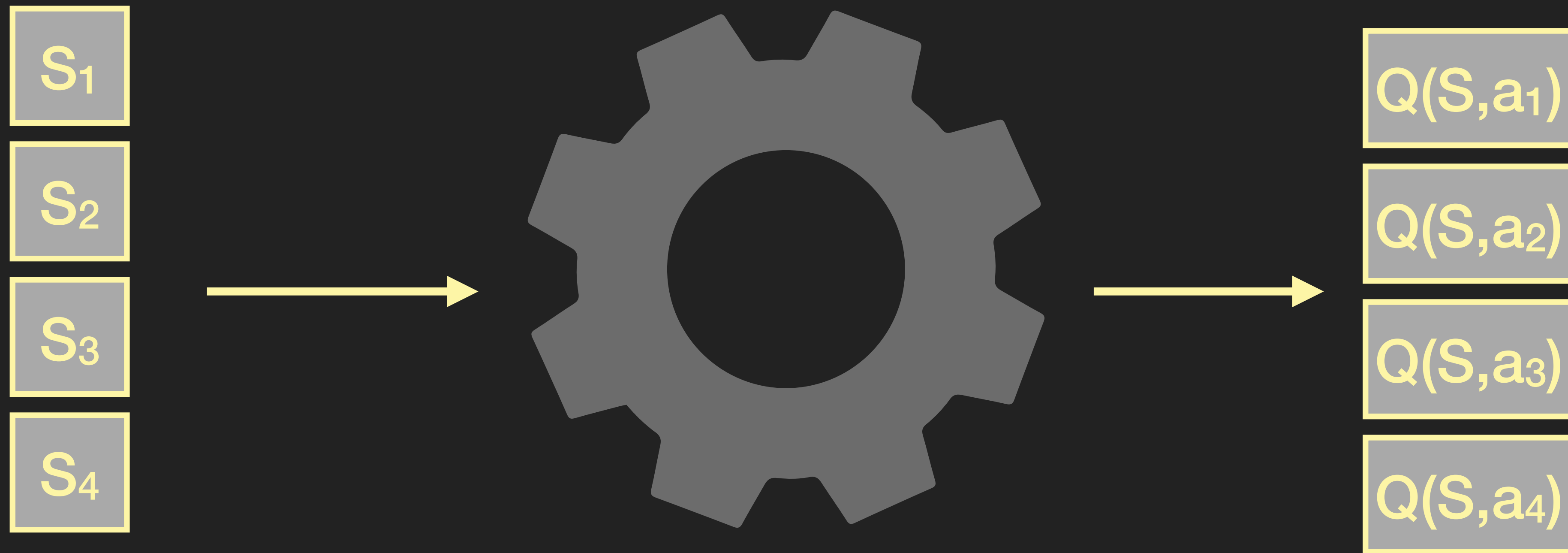
- In Q learning we iteratively update  $Q(A,S)$  via policy iteration
- Neural networks are general function approximators  $NN(S) \rightarrow V(A)$
- Cost function:

Better estimate of  $V(S)$

$$\left( \text{NN}(S_t, a) - r + \gamma \max_a \text{NN}(S_{t+1}, a) \right)^2$$

NN estimation of  $V(S)$

# Deep Q-Learning in a nutshell



# Deep Q-Learning

**For** episode = 1,  $M$  **do**

Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$

**For**  $t = 1, T$  **do**

With probability  $\varepsilon$  select a random action  $a_t$

otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$

Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$

Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$

Every  $C$  steps reset  $\hat{Q} = Q$

**End For**

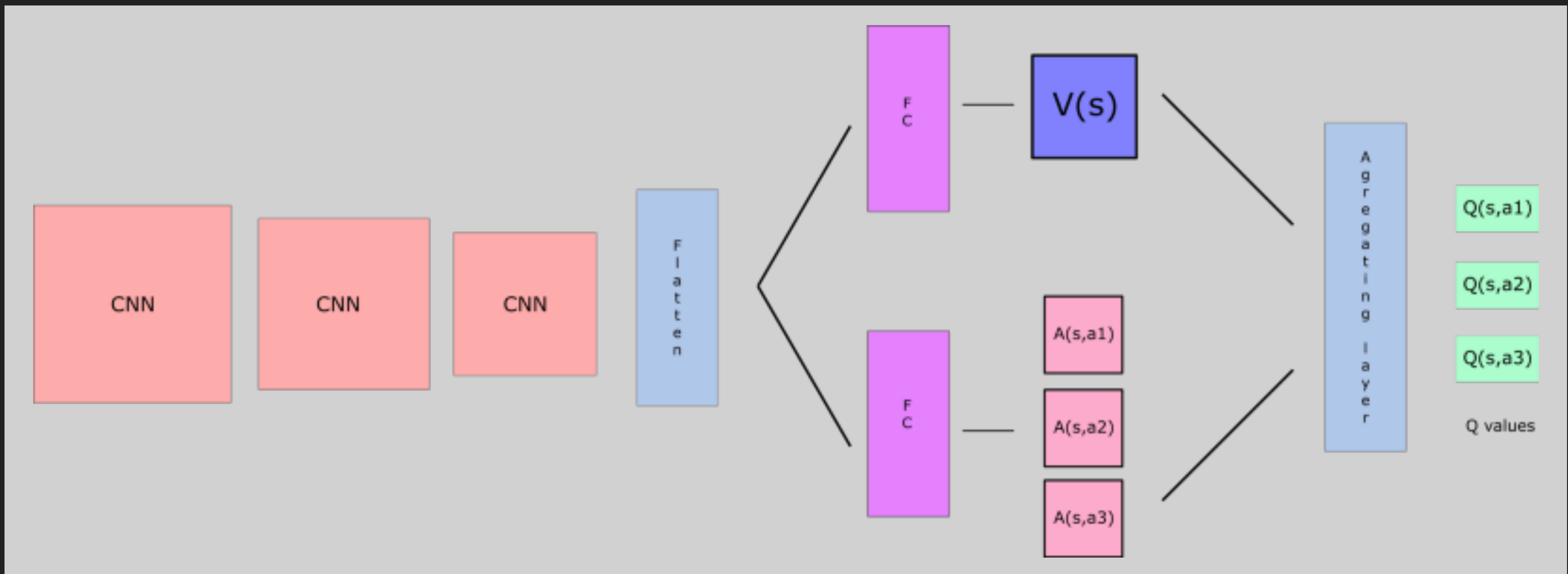
**End For**

# Double deep Q-Learning

$$Y_t^{\text{DQN}} \equiv R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \boldsymbol{\theta}_t^-)$$

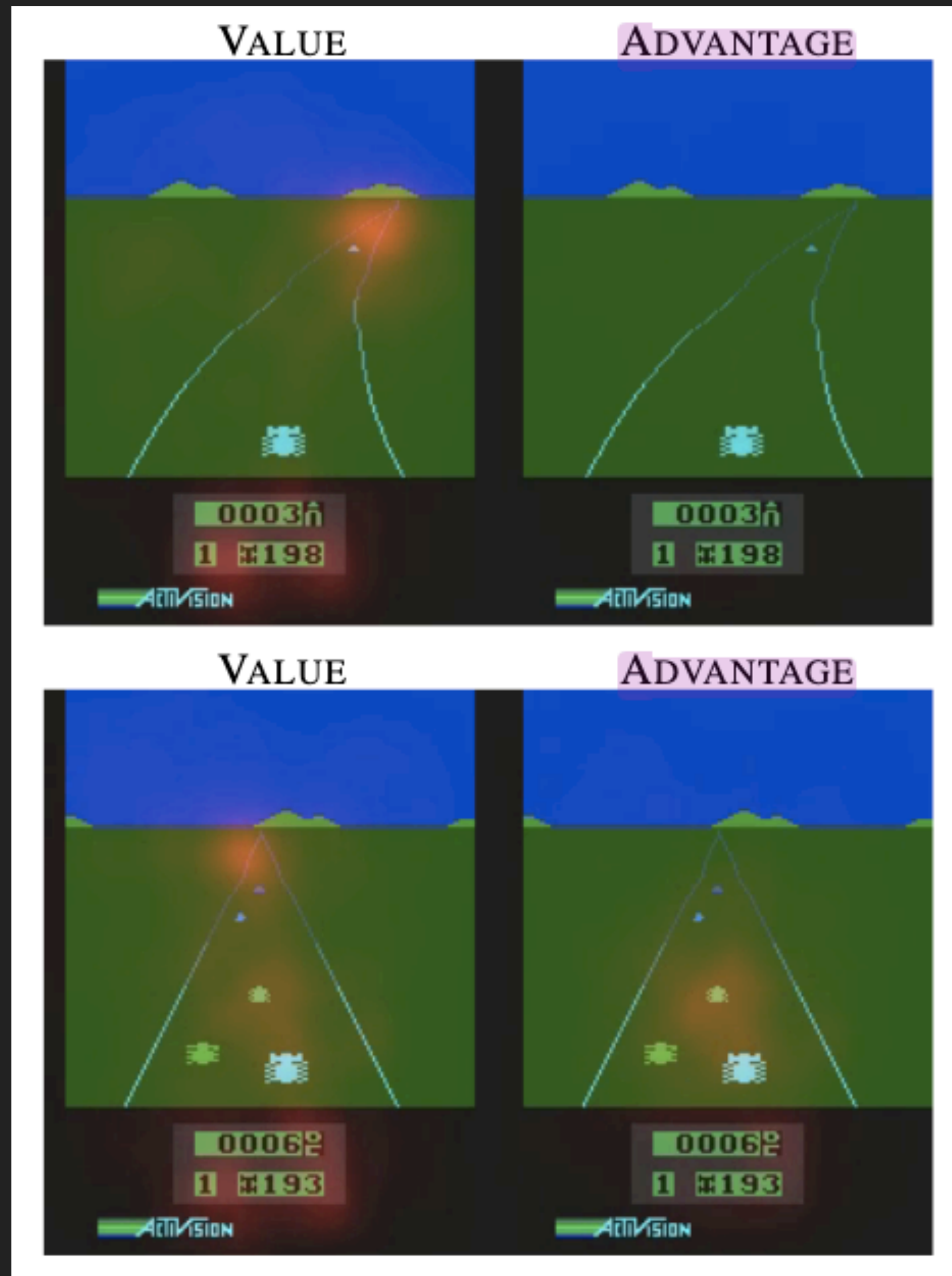
$$Y_t^{\text{DoubleQ}} \equiv R_{t+1} + \gamma Q(S_{t+1}, \underset{a}{\operatorname{argmax}} Q(S_{t+1}, a; \boldsymbol{\theta}_t); \boldsymbol{\theta}_t')$$

# Duelling deep Q-Learning





# Advantage of Action



- Value stream pays attention to road
- Advantage stream pays to care in front to avoid collision
- Final output  $V(s) + (A - \text{mean}(A))$
- How much better is action than other in average



# Noisy Nets



# Multi-step learning

$$\left( \mathbf{NN}(S_t, a) - r + \gamma \max_a \mathbf{NN}(S_{t+1}, a) \right)^2$$

$$\left( \mathbf{NN}(S_t, a) - r + \gamma^n \max_a \mathbf{NN}(S_{t+n}, a) \right)^2$$

# Prioritized Experience Replay

---

**Algorithm 1** Double DQN with proportional prioritization

---

- 1: **Input:** minibatch  $k$ , step-size  $\eta$ , replay period  $K$  and size  $N$ , exponents  $\alpha$  and  $\beta$ , budget  $T$ .
  - 2: Initialize replay memory  $\mathcal{H} = \emptyset$ ,  $\Delta = 0$ ,  $p_1 = 1$
  - 3: Observe  $S_0$  and choose  $A_0 \sim \pi_\theta(S_0)$
  - 4: **for**  $t = 1$  **to**  $T$  **do**
  - 5:   Observe  $S_t, R_t, \gamma_t$
  - 6:   Store transition  $(S_{t-1}, A_{t-1}, R_t, \gamma_t, S_t)$  in  $\mathcal{H}$  with maximal priority  $p_t = \max_{i < t} p_i$
  - 7:   **if**  $t \equiv 0 \pmod K$  **then**
  - 8:     **for**  $j = 1$  **to**  $k$  **do**
  - 9:       Sample transition  $j \sim P(j) = p_j^\alpha / \sum_i p_i^\alpha$
  - 10:       Compute importance-sampling weight  $w_j = (N \cdot P(j))^{-\beta} / \max_i w_i$
  - 11:       Compute TD-error  $\delta_j = R_j + \gamma_j Q_{\text{target}}(S_j, \arg \max_a Q(S_j, a)) - Q(S_{j-1}, A_{j-1})$
  - 12:       Update transition priority  $p_j \leftarrow |\delta_j|$
  - 13:       Accumulate weight-change  $\Delta \leftarrow \Delta + w_j \cdot \delta_j \cdot \nabla_\theta Q(S_{j-1}, A_{j-1})$
  - 14:     **end for**
  - 15:     Update weights  $\theta \leftarrow \theta + \eta \cdot \Delta$ , reset  $\Delta = 0$
  - 16:     From time to time copy weights into target network  $\theta_{\text{target}} \leftarrow \theta$
  - 17:   **end if**
  - 18:   Choose action  $A_t \sim \pi_\theta(S_t)$
  - 19: **end for**
-

# Distributional Reinforcement Learning

# License

All course material- including this slide deck, hands-on workbooks, and any other guided exercise for example assignments- are part of the workshop "RL workshop" and owned by Aggregate Intellect Inc. (<https://ai.science>), and released under "Creative Commons Attribution-NonCommercial-ShareAlike CC BY-NC-SA" license.

This material can be altered and distributed for non-commercial use with reference to Aggregate Intellect Inc. as the original owner, and any material generated from it must be released under similar terms.

For more details see: <https://creativecommons.org/licenses/by-nc-sa/4.0/>