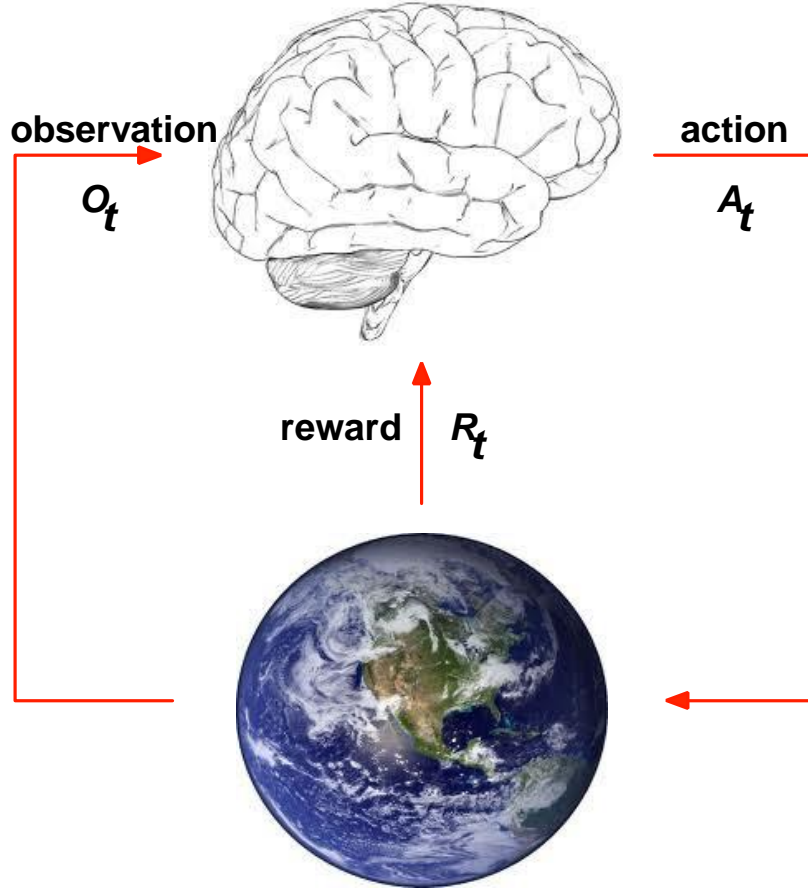


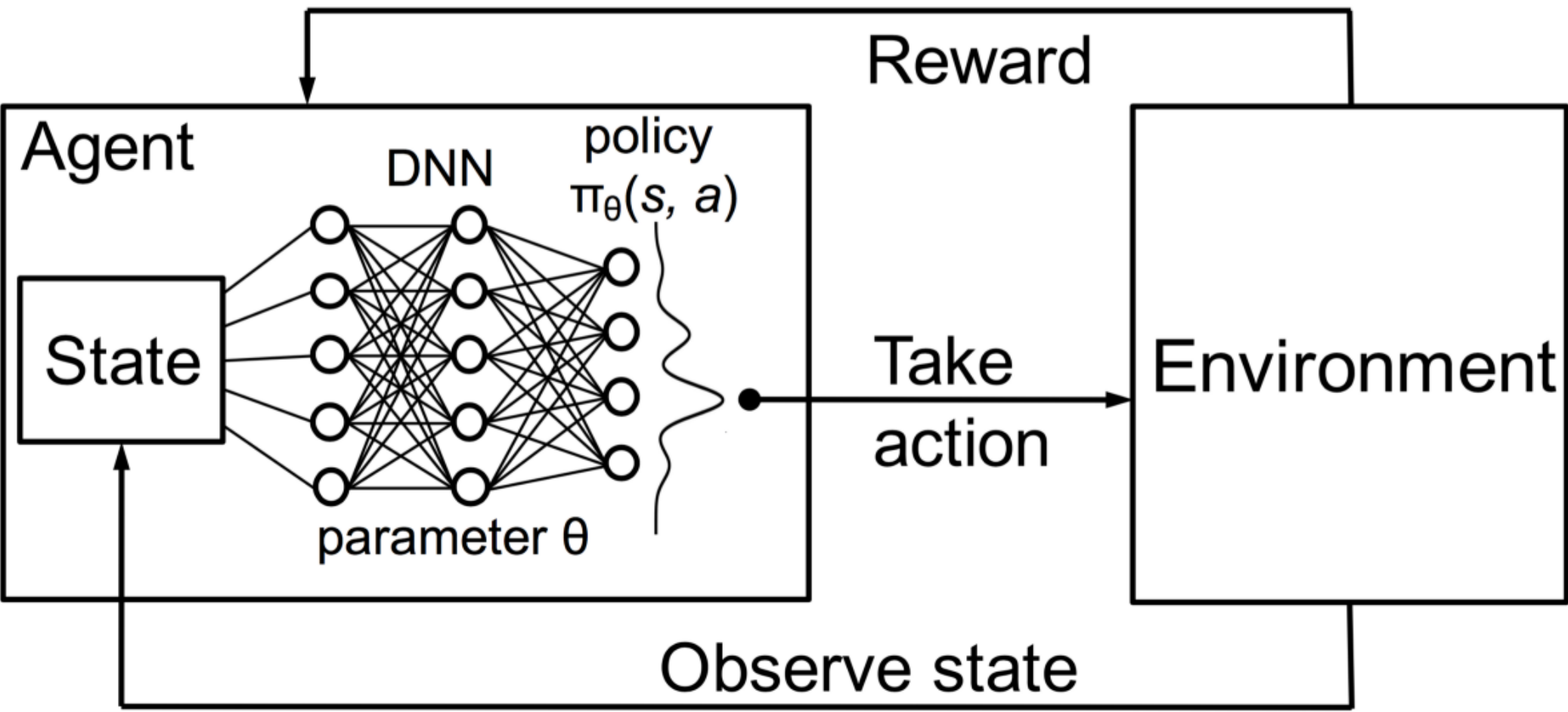
RL: DEEP RL

# План

- ☐ Deep Q learning
- ☐ Experience replay
- ☐ Double Q-learning
- ☐ Dueling Network

# Взаимодействие среды и агента. MDP





# Базовые определения

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) = r_{t+1} + \gamma G_{t+1}$$

$$V^\pi(s) = E_\pi[G_t | s_t = s]$$

$$Q^\pi(s) = E_\pi[G_t | s_t = s, a_t = a]$$

Рекуррентные формулы

$$Q^\pi(s) = E_{s_{t+1}} [r_t + \gamma V^\pi(s_{t+1})]$$

$$Q^\pi(s) = E_{s_{t+1}, a_{t+1} \sim \pi} [r_t + Q^\pi(s_{t+1}, a_{t+1})]$$

## Оптимальная политика

Для всех  $\pi, s, a$ :  $Q^{\pi^*}(s, a) \geq Q^{\pi}(s, a)$

$$\pi^*(s) = \operatorname{argmax}_a Q^{\pi^*}(s, a)$$

Уравнение Беллмана

$$Q^*(s_t, a) = E_{s_{t+1}} [r_t + \max_{a'} Q^*(s_{t+1}, a')]$$

## Q - learning

Шаг обучения

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left( r_t + \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)$$

Q-learning как минимизация MSE

$$L = \left( r_t + \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)^2$$

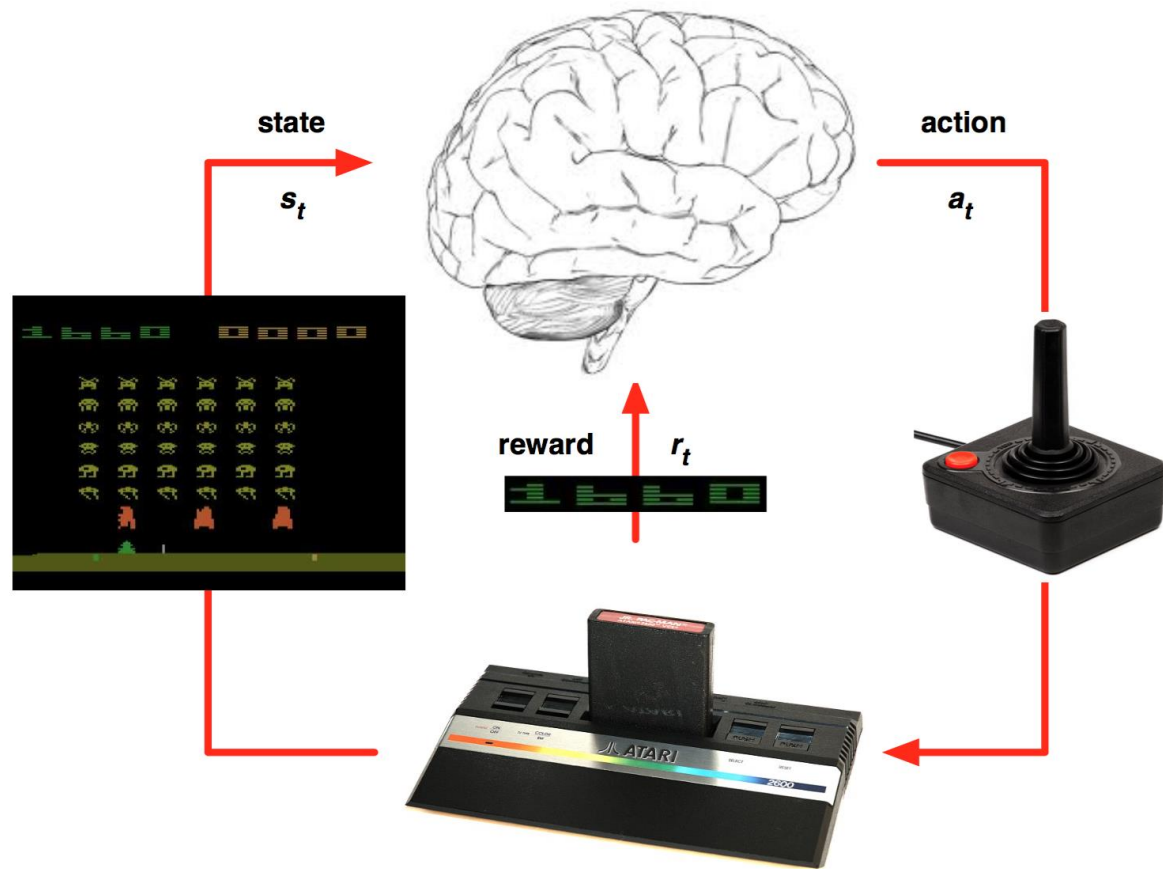
$$\nabla L = 2 \left( r_t + \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)$$

Какая проблема?

# Реальный мир

Сколько приблизительно состояний?

$$S = 2^{210 \cdot 160 \cdot 8 \cdot 3}$$





# Проблема

- ❑ Пространство состояний может быть очень большим или непрерывным
  
- ❑ Два решения:
  - ❑ Бинаризация пространства состояний
  - ❑ Функциональная аппроксимация агента
  
- ❑ Что из этого лучше подойдет для Atari?

## От таблицы к аппроксимации

□ Для таблицы:

- Для всех состояний и действий мы должны запомнить  $Q(s, a)$

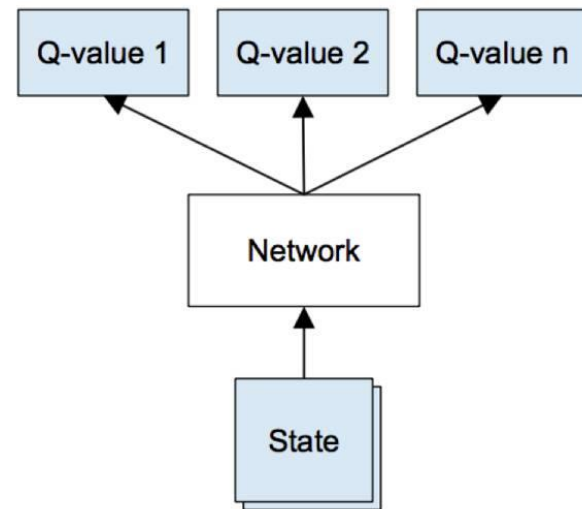
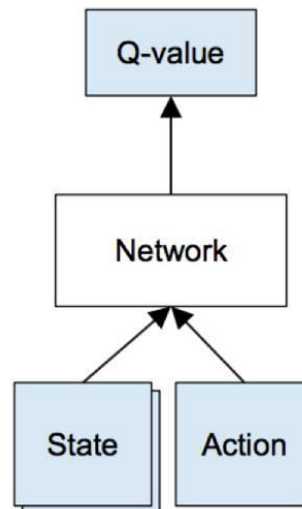
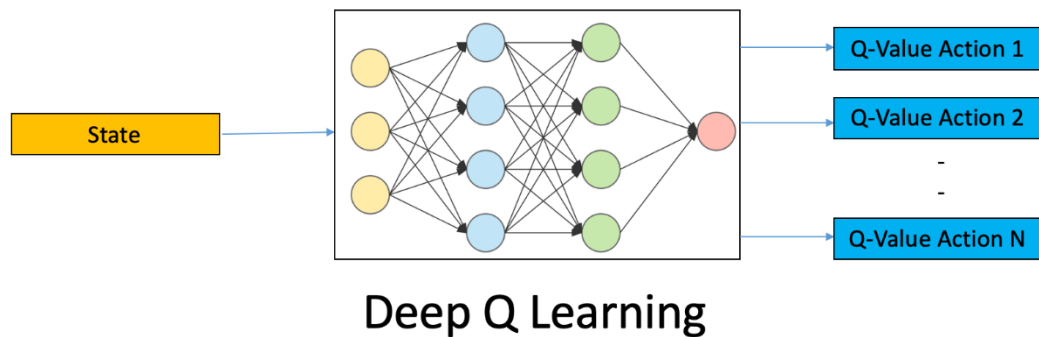
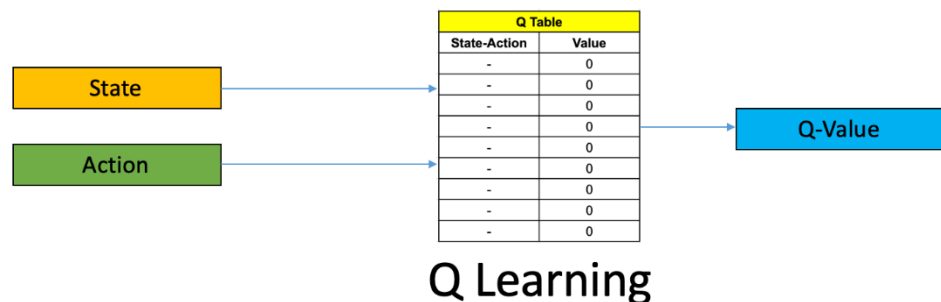
□ Аппроксимация:

- Аппроксимация  $Q(s, a)$  некоторой функцией
- Например линейной функцией

$$\operatorname{argmin}_{w,b} \left( r_t + \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)^2$$

Вопрос: что мы используем, классификация или регрессию?

# Возможные архитектуры



# Аппроксимация Q-learning

## □ Q-values

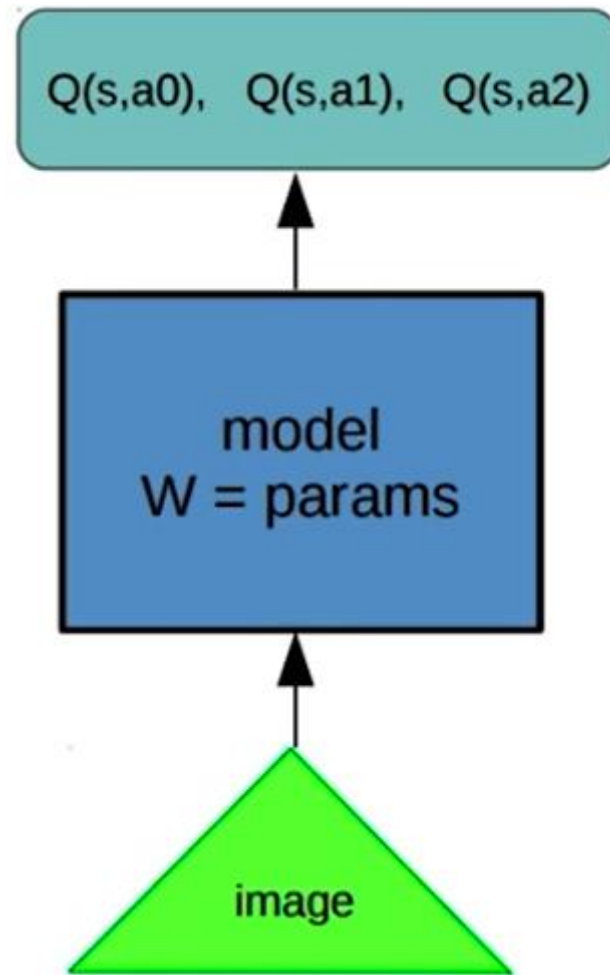
- $\hat{Q}(s_t, a_t) = r_t + \max_{a'} Q(s_{t+1}, a')$

## □ Целевая функция:

- $L = \left( Q(s_t, a_t) - [r_t + \max_{a'} Q(s_{t+1}, a')] \right)^2$

## □ Шаг градиента:

$$w_{t+1} = w_t - \alpha \frac{\partial L}{\partial w}$$



# Аппроксимация Q-learning

❑ Целевая функция:

$$\blacksquare L = \left( Q(s_t, a_t) - \hat{Q}(s_t, a_t) \right)^2$$

❑ Q-learning

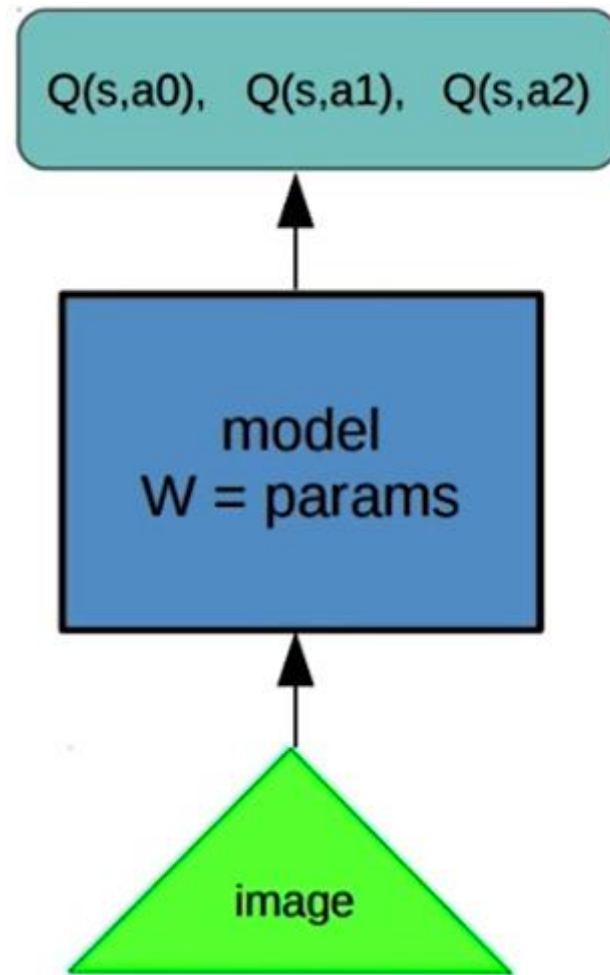
$$\blacksquare \hat{Q}(s_t, a_t) = r_t + \gamma \max_{a'} Q(s_{t+1}, a')$$

❑ SARSA

$$\blacksquare \hat{Q}(s_t, a_t) = r_t + \gamma Q(s_{t+1}, a_{t+1})$$

❑ Expected Value SARSA:

$$\hat{Q}(s_t, a_t) = r_t + \gamma E_{\pi} Q(s_{t+1}, a_{t+1})$$

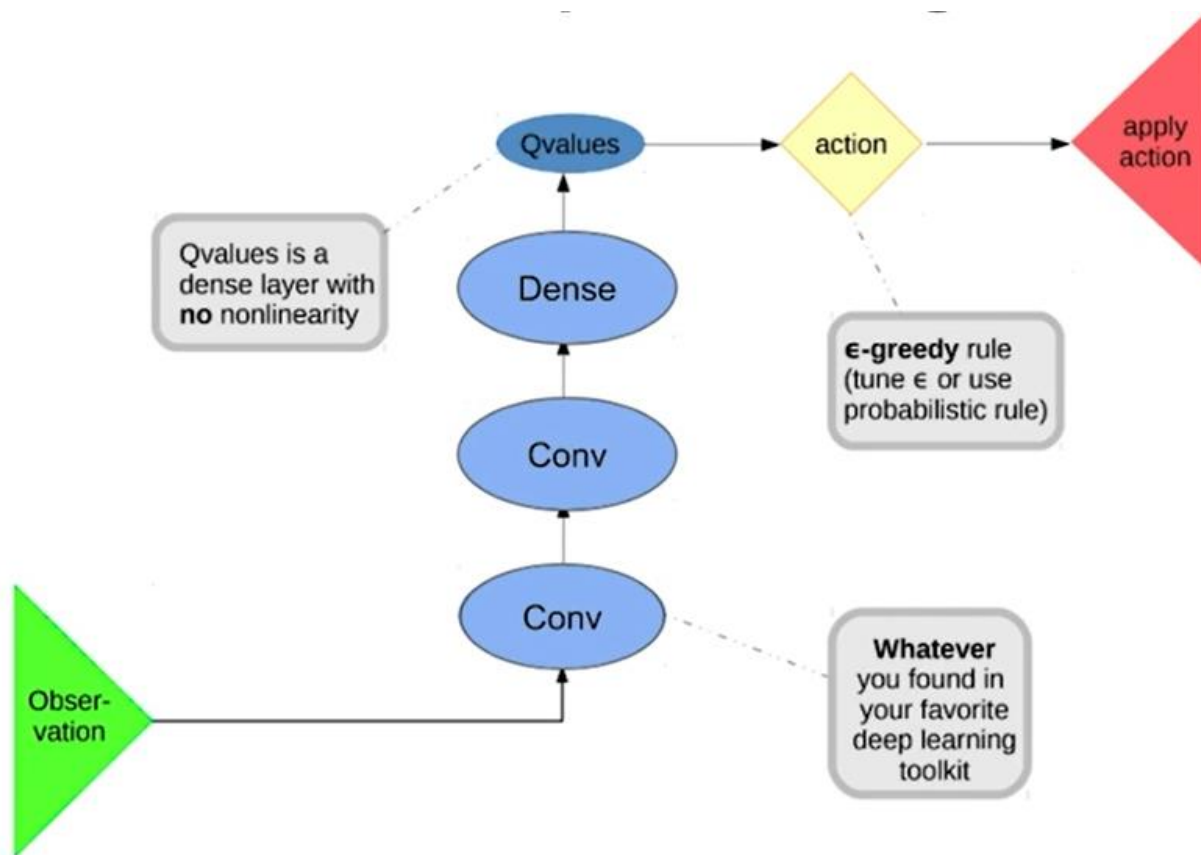


# Deep Q-learning



Какую нейросеть применить для этой задачи?

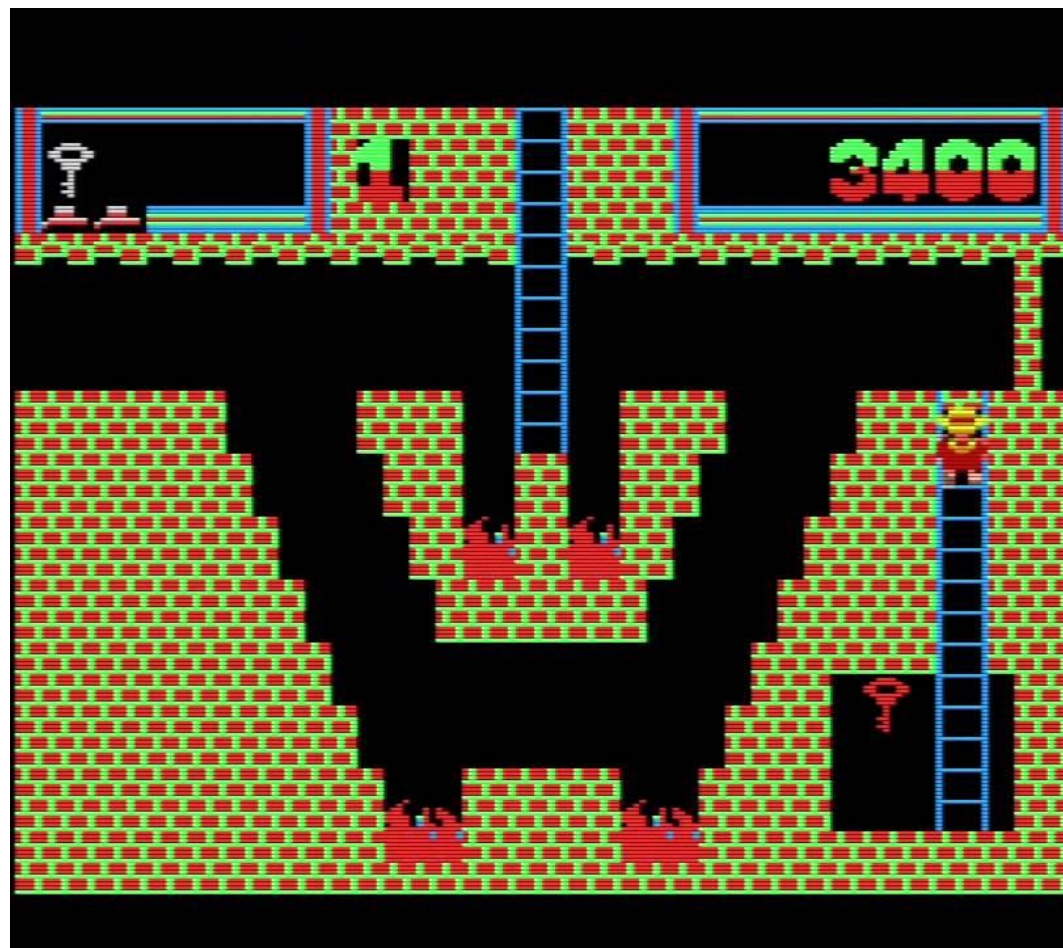
# Deep Q-learning



**Вы уже играете в лигу 8 лет, я так полагаю, что вы хороший игрок**







# Проблема

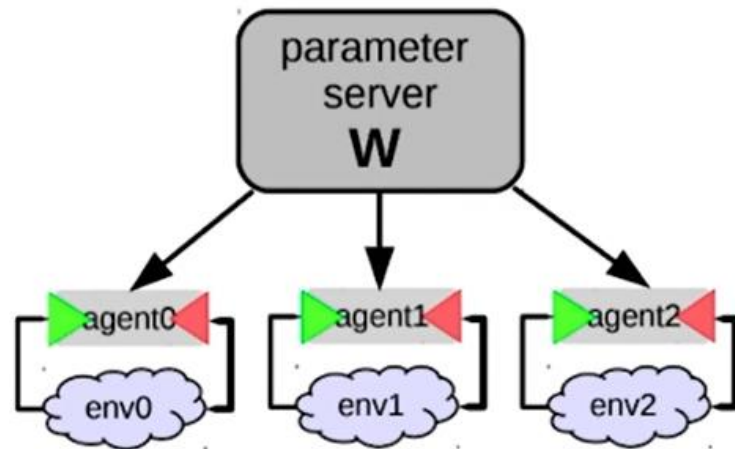
- ❑ Обучаемая выборка не независима, т.е. изображения не являются независимыми одинаково распределенными величинами
- ❑ Модели забывают о частях среды, которые они не посещали в течение определенного времени
- ❑ Снижение темпов обучения
- ❑ Идеи?

## Много агентов

- ❑ Идея: Использовать несколько агентов
- ❑ скорее всего, они будут исследовать разные части среды
- ❑ Более стабильное обучение
- ❑ Требуется множество взаимодействий

Вопрос: агент автономный транспорт.

Какие проблемы?



# Experience replay

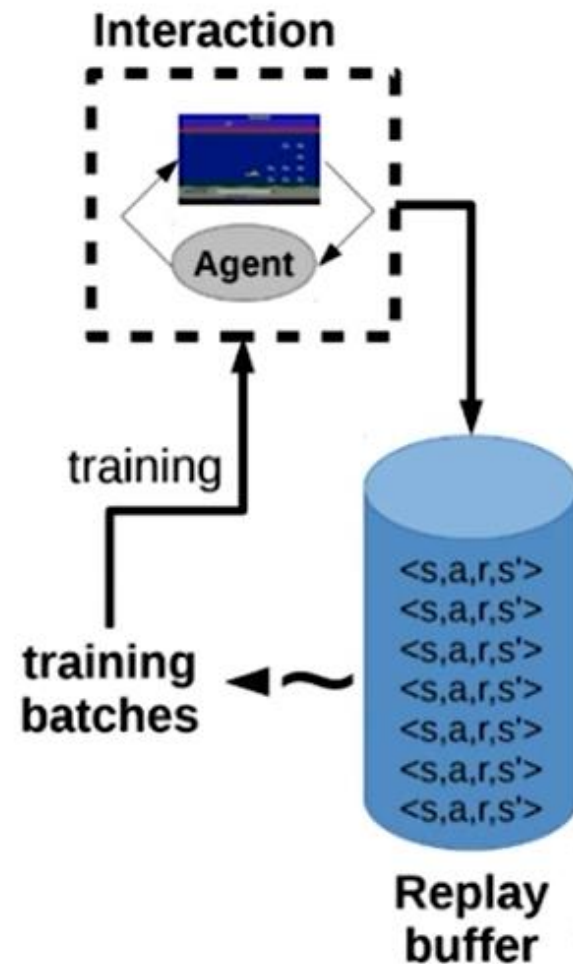
- ❑ Идея: сохранить несколько прошлых итераций

$\langle s, a, r, s' \rangle$

*Обучение проходит на случайных подвыборках*

- ❑ Подвыборки ближе к независимым

- ❑ более старые взаимодействия были получены при более слабой политике



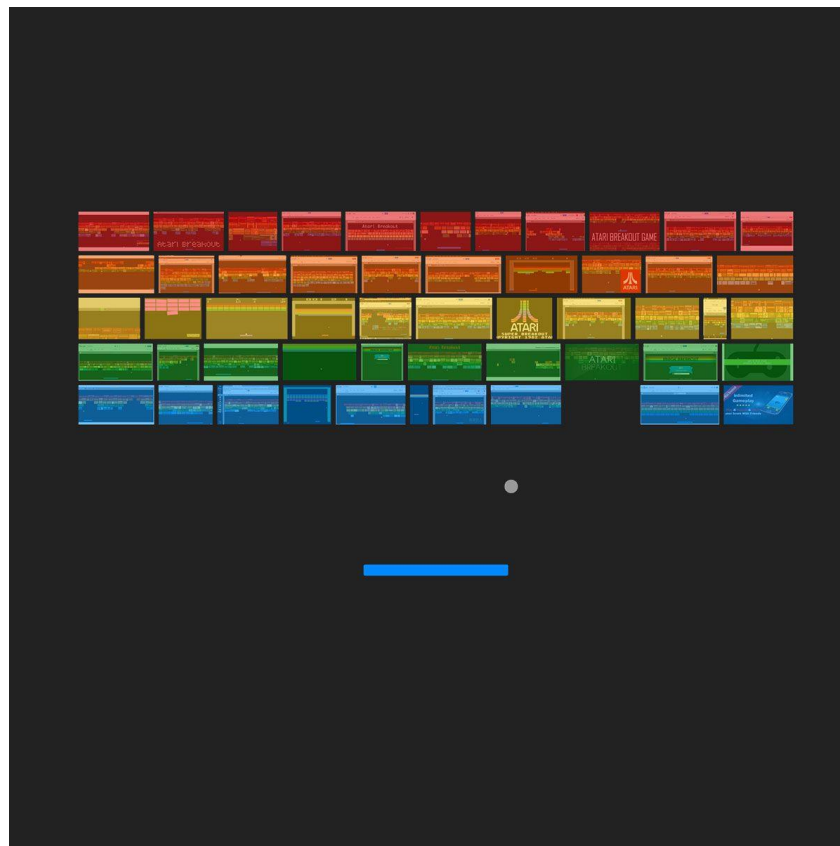
# Experience replay

- ❑ Мы аппроксимируем  $Q(s,a)$  нейронной сетью
- ❑ Используем *experience replay* для обучения
  
- ❑ *Вопрос: какие из перечисленных алгоритмов можем использовать?*
  - Q-learning
  - SARSA
  - CEM
  - Expected Value SARSA

## При обучении on-policy алгоритмов

- Не использовать (или использовать не большой) experience replay
- Можно компенсировать сессиями параллельных игр

# Лево или право?

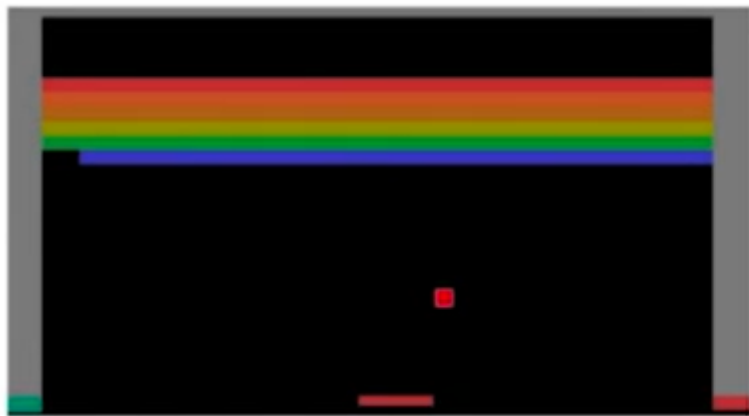


# N-gram trick

Идея:

$$s_t \neq o(s_t)$$

$$s_t \approx (o(s_{t-n}), a_{t-n}, \dots, o(s_{t-1}), a_{t-1}, o(s_t))$$



Одна игра



Несколько игр



# N-gram trick

- ❑ Предполагается марковость  $n$ -го порядка:
- ❑ Работает для скоростей/ускорений
- ❑ Потерпим неудачу, если информация выходит за рамки  $N$  кадров
- ❑ Непрактично для больших  $N$

# Автокорреляция

# Target network

- Идея: использовать сеть с замороженными весами для расчета таргета.

$$L(\theta) = E_{s \sim S, a \sim A} \left( Q(s_t, a_t, \theta) - \left[ r_t + \gamma \max_{a'} Q(s_{t+1}, a', \bar{\theta}) \right] \right)^2$$

где  $\bar{\theta}$  - замороженные веса

- *Hard target network:*

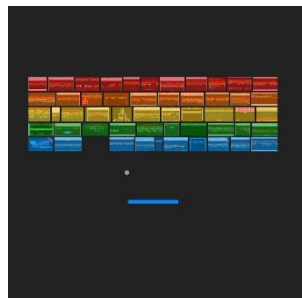
Обновлять  $\bar{\theta}$  каждые  $n$  шагов и устанавливать значения  $\theta$

- *Soft target network:*

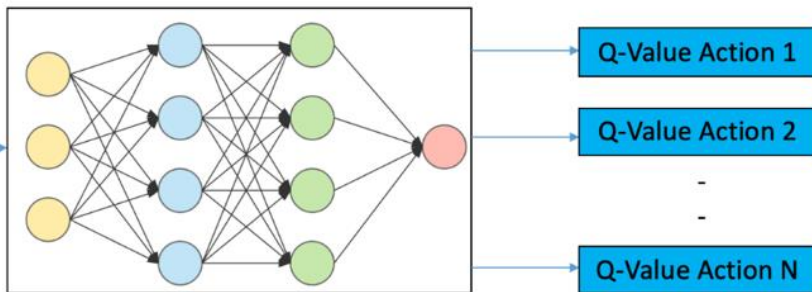
Обновлять  $\bar{\theta}$  каждый шаг:

$$\bar{\theta} = (1 - \alpha)\bar{\theta} + \alpha\theta$$

# Playing Atari with Deep Reinforcement Learning (2013, DeepMind)



4 последних кадра



□ Обновление весов:

$$L(\theta) = E_{s \sim S, a \sim A} \left( Q(s_t, a_t, \theta) - \left[ r_t + \gamma \max_{a'} Q(s_{t+1}, a', \bar{\theta}) \right] \right)^2$$

Обновление  $\bar{\theta}$  каждые 5000 шагов обучения

Experience replay



$10^6$  последних переходов

# Проблема завышенной оценки

- Мы используем оператор  $\max$  для расчета таргета

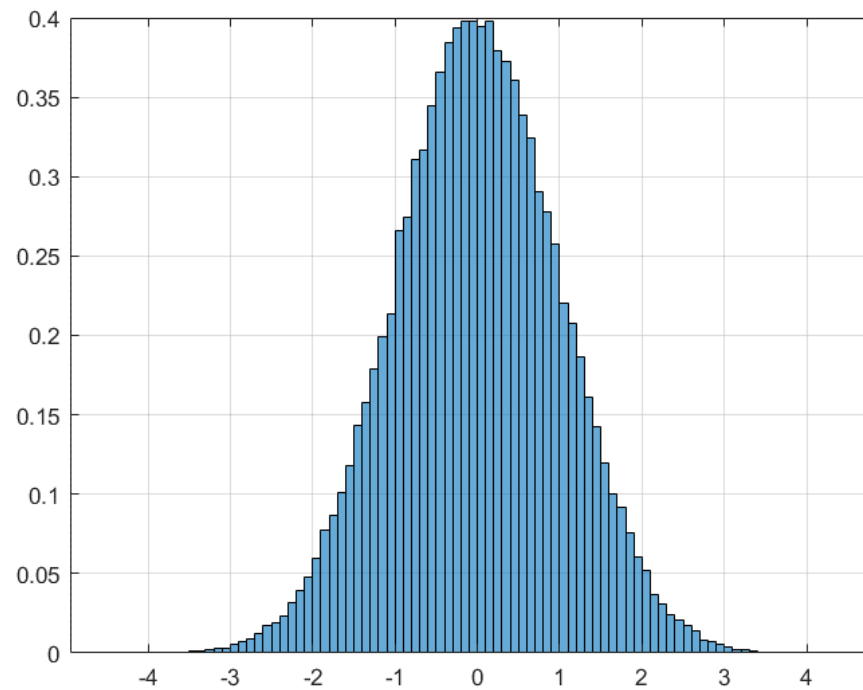
$$L(s, a) = \left( Q(s, a) - \left[ r_t + \gamma \max_{a'} Q(s', a') \right] \right)^2$$

Проблема

(мы хотим чтобы  $E_{s \sim S, a \sim A}[L(s, a)] = 0$ )

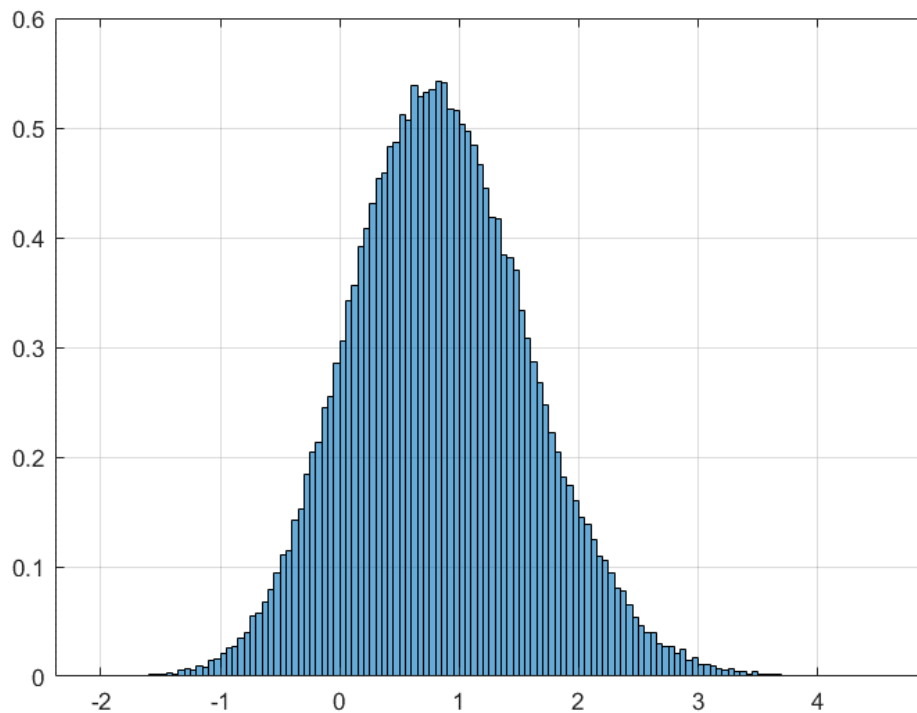
# Проблема завышенной оценки

- $Mean(data) \sim 0.001$

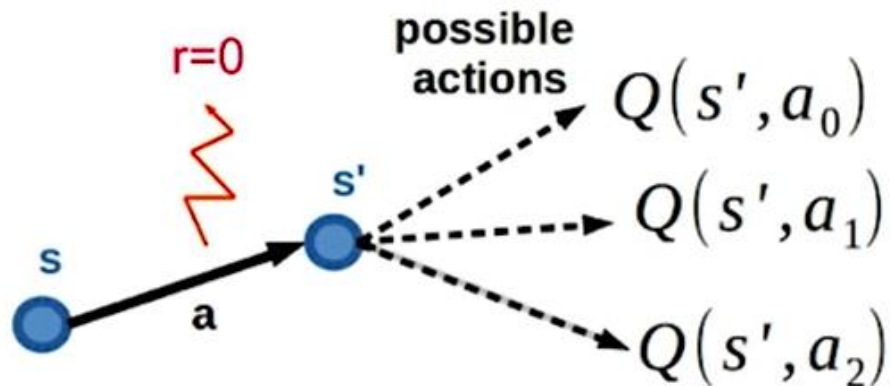


# Проблема завышенной оценки

- $Mean(data) \sim 0.8482$



## Проблема завышенной оценки



- Если мы обновляем  $Q(s,a)$  в соответствии с  $r_t + \gamma \max_{a'} Q(s', a')$ , мы будем иметь завышенную оценку

$$E[\max_{a'} Q(s', a')] \geq \max_{a'} E[Q(s', a')]$$



# Double Q-learning (NIPS 2010)

$$y = r_t + \gamma \max_{a'} Q(s', a')$$

Q-learning target

$$y = r_t + \gamma \max_{a'} Q(s', \operatorname{argmax}_{a'} Q(s', a'))$$

Q-learning target

*Идея: использовать 2 оценки q-values:  $Q^A, Q^B$*

*Q-learning Они должны компенсировать несоответствия друг друга, т.к. они независимы.*

$$y = r_t + \gamma \max_{a'} Q^A(s', \operatorname{argmax}_{a'} Q^B(s', a'))$$

Double Q-learning

# Double Q-learning (NIPS 2010)

---

**Algorithm 1** Double Q-learning

---

```
1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end
```

---

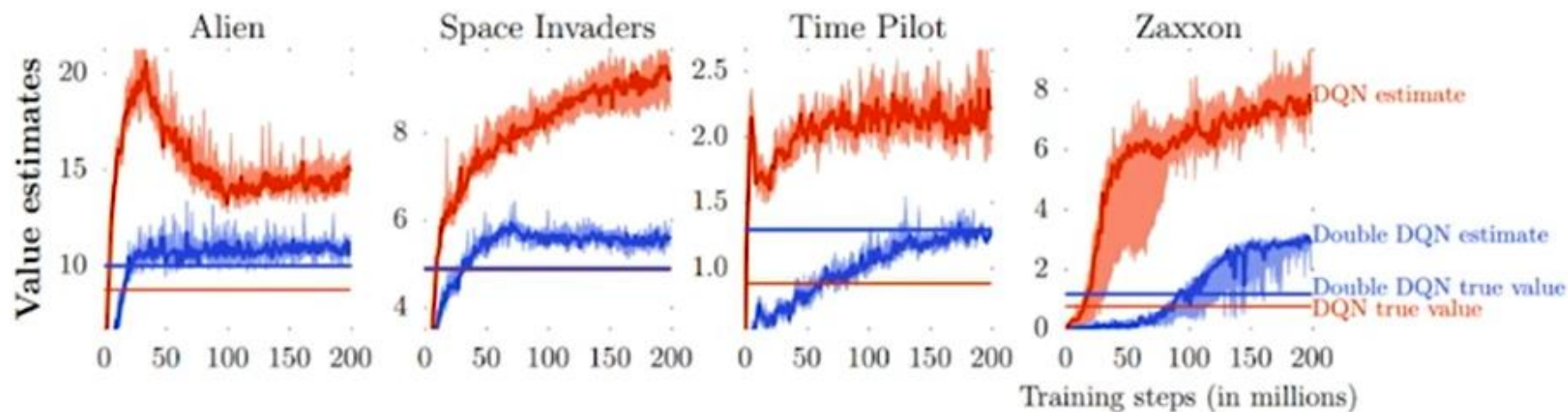
*Сможем ли мы скомбинировать этот алгоритм с DQN?*

# Deep Reinforcement Learning with Double Q-learning (Deepmind, 2015)

Идея: использовать основную сеть для выбора действия

$$y_{dq\eta} = r_t + \gamma \max_{a'} Q(s', a', \bar{\theta})$$

$$y_{ddq\eta} = r_t + \gamma Q(s', \operatorname{argmax}_{a'} Q(s', a', \theta), \bar{\theta})$$



	DQN	Double DQN	Double DQN (tuned)
Median	47.5%	88.4%	116.7%
Mean	122.0%	273.1%	475.2%

## Experience replay

State	Action	Reward	Next state
s_0	a_0	0	s_1
s_1	a_1	0	s_2
...	...	...	...
s_(n-1)	a_(n-1)	0	s_n
<b>s_n</b>	<b>a_n</b>	<b>100</b>	<b>s_(n+1)</b>
s_(n+1)	a_(n+1)	0	s_(n+2)
...	...	...	...

## Experience replay: приоритезация (DeerMind, 2016)

Идея: семплировать переходы из буфера более грамотно

Мы хотим назначить вероятности для каждого перехода. Можем использовать абсолютное значение TD ошибки перехода как вероятности

$$\text{TD-error } \delta = Q(s, a) - (r + \gamma Q(s', \argmax_{a'} Q(s', a'), \theta), \bar{\theta})$$

$$p = |\delta|$$

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha} \text{ здесь } \alpha \text{ параметр приоритета (0 – равномерный выбор)}$$

Какую видите проблему?

## Experience replay: приоритезация (DeepMind, 2016)

Решение: корректно обучать используя importance-sampling weights

$$w_i = \left( \frac{1}{N} \frac{1}{P(i)} \right)^\beta \quad \text{где } \beta \text{ параметр}$$

Таким образом мы семплируем  $P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$  и умножая ошибку на  $w_i$

# Experience replay: приоритезация (DeerMind, 2016)

Дополнительные детали

Нормировка весов на  $1 / \max_i w_i$

При записи переходов в experience replay, принимаем maximal priority  $p_t = \max_{i < t} p_i$

*<https://youtu.be/UXurvvdY93o>*



# Dueling Network Architectures for Deep Reinforcement Learning (DeepMind, 2016)

Идея: изменить архитектуру НС

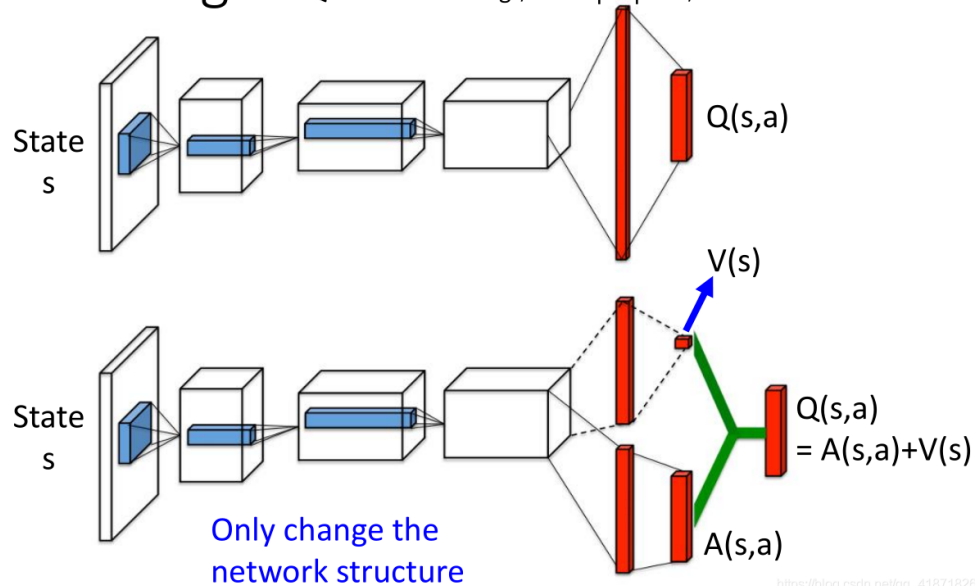
Advantage function  $A(s,a) = Q(s,a) - V(s)$

Таким образом  $Q(s,a) = A(s,a) + V(s)$

Проблема?

## Dueling DQN

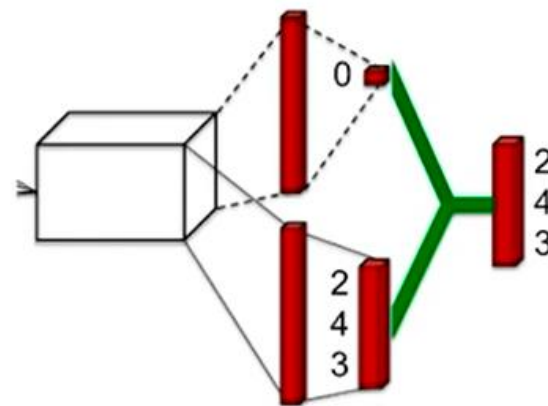
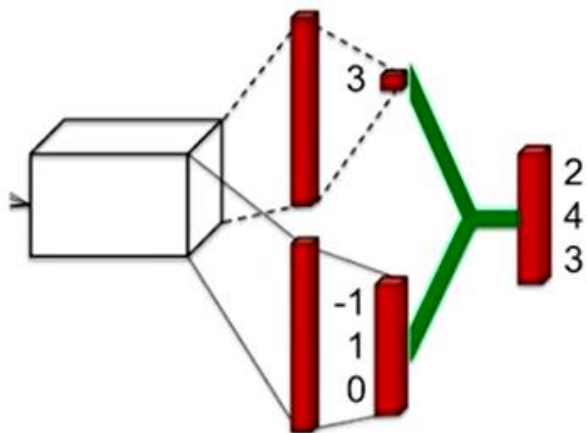
Ziyu Wang, Tom Schaul, Matteo Hessel, Hado van Hasselt, Marc Lanctot, Nando de Freitas, "Dueling Network Architectures for Deep Reinforcement Learning", arXiv preprint, 2015



[https://blog.csdn.net/qz\\_41871826](https://blog.csdn.net/qz_41871826)

# Dueling Network Architectures for Deep Reinforcement Learning (DeepMind, 2016)

Имеем дополнительную степень свободы



*Какой вариант лучше?*

# Dueling Network Architectures for Deep Reinforcement Learning (DeepMind, 2016)

Решение: требуем  $\max_{a' \in |A|} A(s, a'; \theta, \alpha)$  равнялось нулю

Таким образом Q функция рассчитывается как:

$$Q(s, a, \alpha, \beta) = V(s; \theta, \beta) + \left( A(s, a; \theta, \alpha) - \max_{a' \in |A|} A(s, a'; \theta, \alpha) \right)$$

Авторы статьи также предложили следующий способ расчета Q-values:

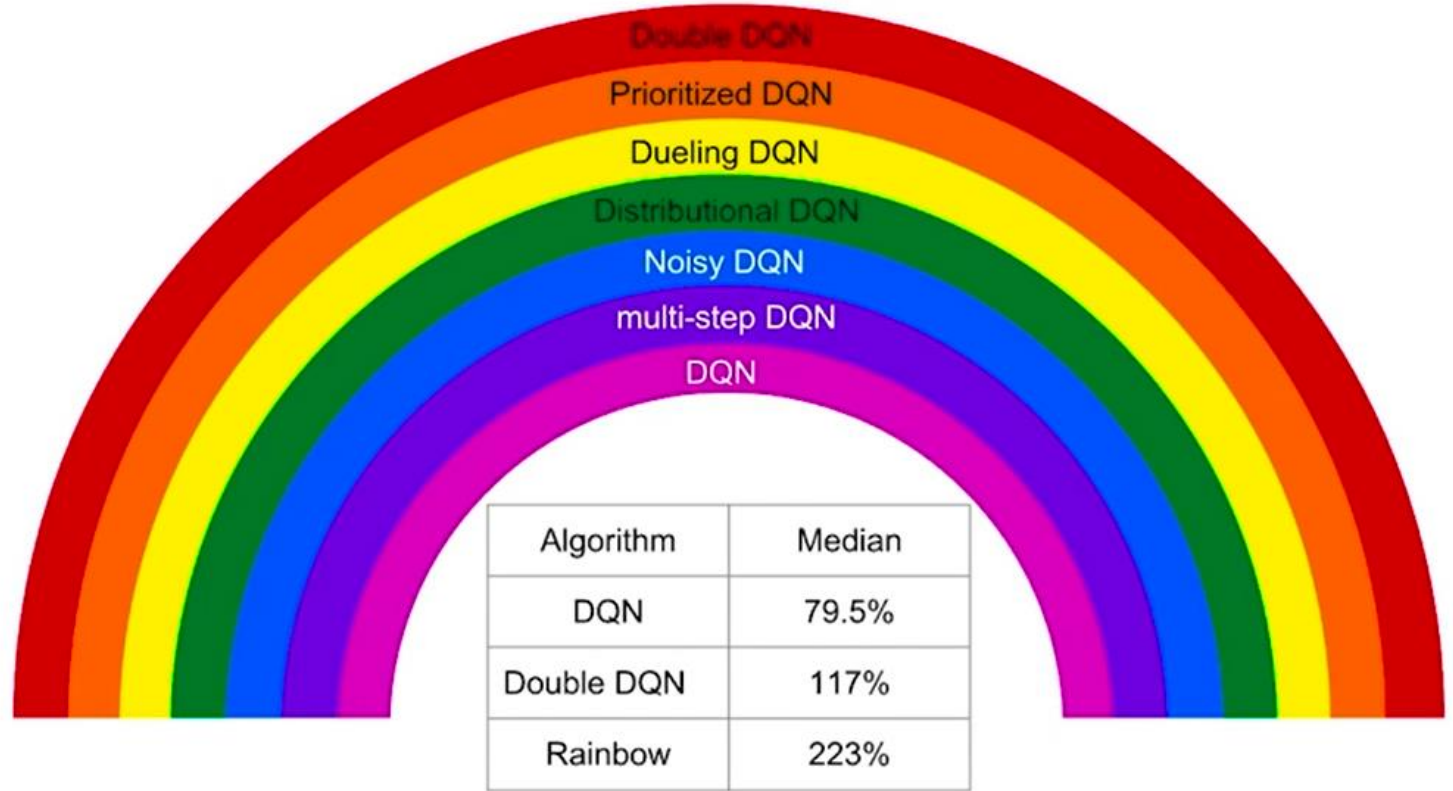
$$Q(s, a, \alpha, \beta) = V(s; \theta, \beta) + \left( A(s, a; \theta, \alpha) - \frac{1}{|A|} \sum_{a'} A(s, a'; \theta, \alpha) \right)$$

Этот вариант повышает стабильность оптимизации

# Асинхронное обучение для DLR (DeerMind, 2016)

# Rainbow (DeepMind, 2017)

Rainbow (2017, Deepmind)



## R2D2: (Deepmind, 2018)

- LSTM
- Distributed Prioritized Experience replay
- N-step DQN
- Reword re-scaling
- Double DQN
- Dueling DQN

Медианная производительность: 1920% от человеческой