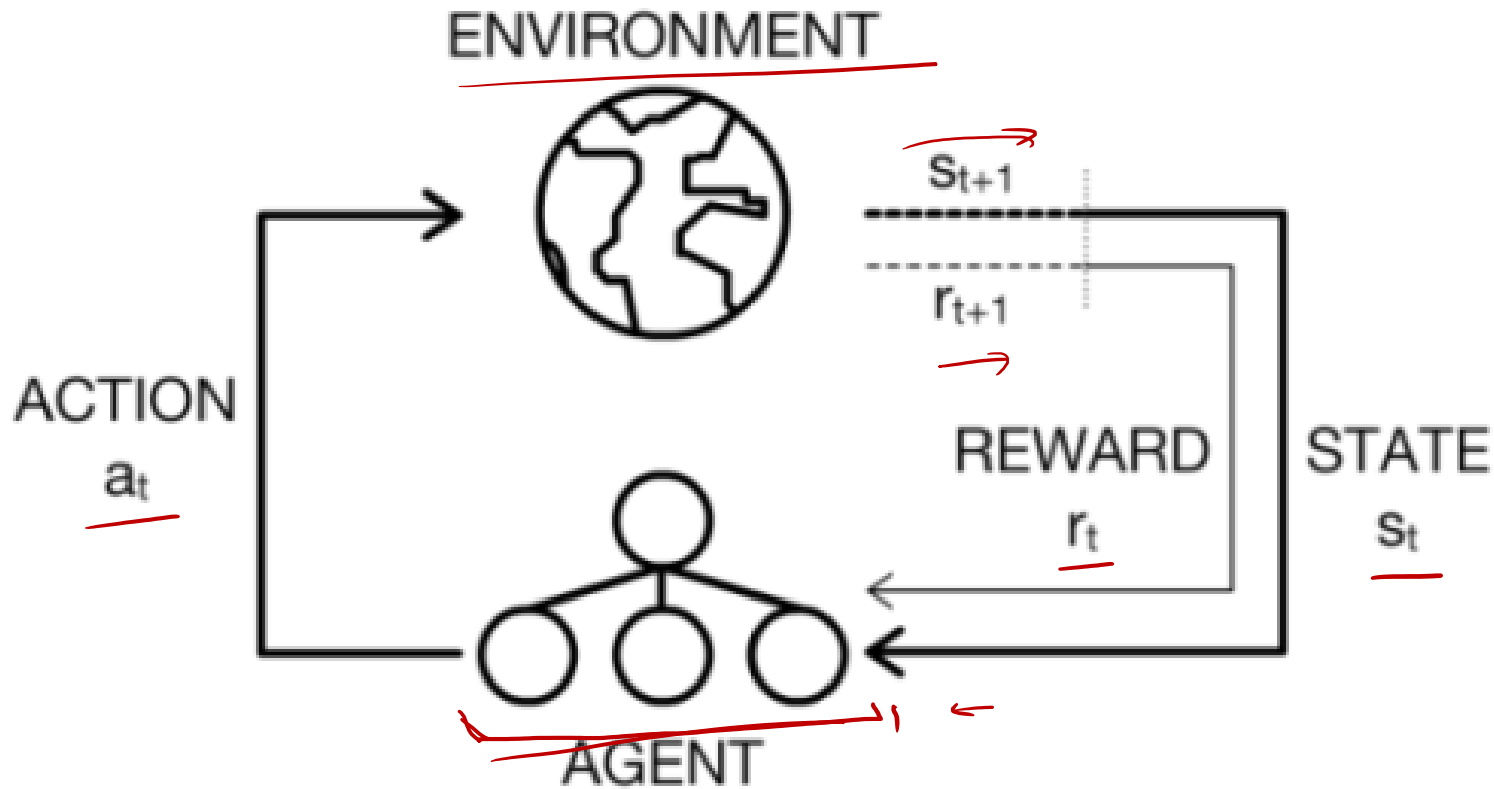


13



Обучение с подкреплением (RL)

Обучение с подкреплением Reinforcement Learning

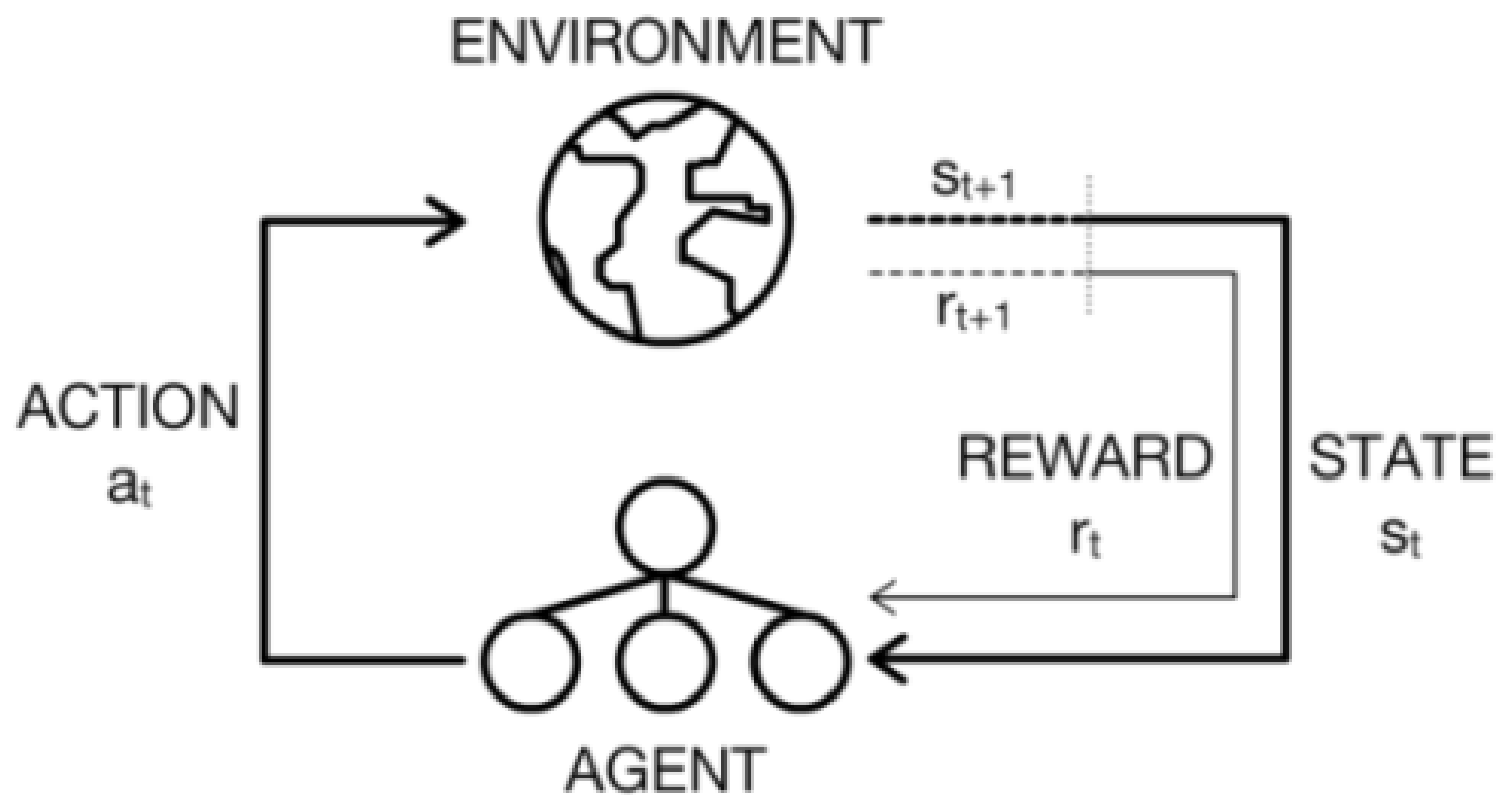


s_t – состояние среды (state) на шаге t
 a_t – действие агента (action) на шаге t
 r_t – награда (reward) на шаге t

Среда может:

- Быть недетерминированной
- Иметь внутреннее состояние

Основные определения



s_t – состояние среды (state) на шаге t
 a_t – действие агента (action) на шаге t
 r_t – награда (reward) на шаге t

$p(a|s)$ - policy function

$v(s)$ - value function

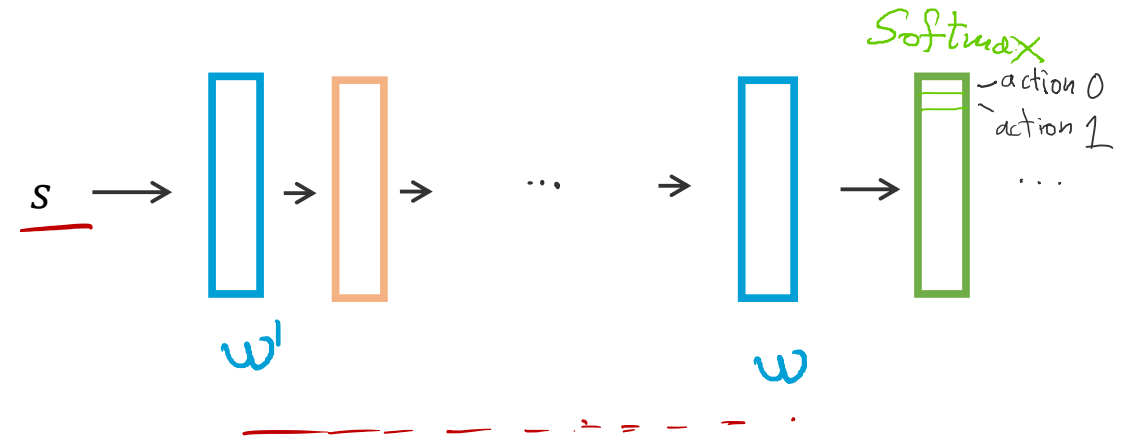
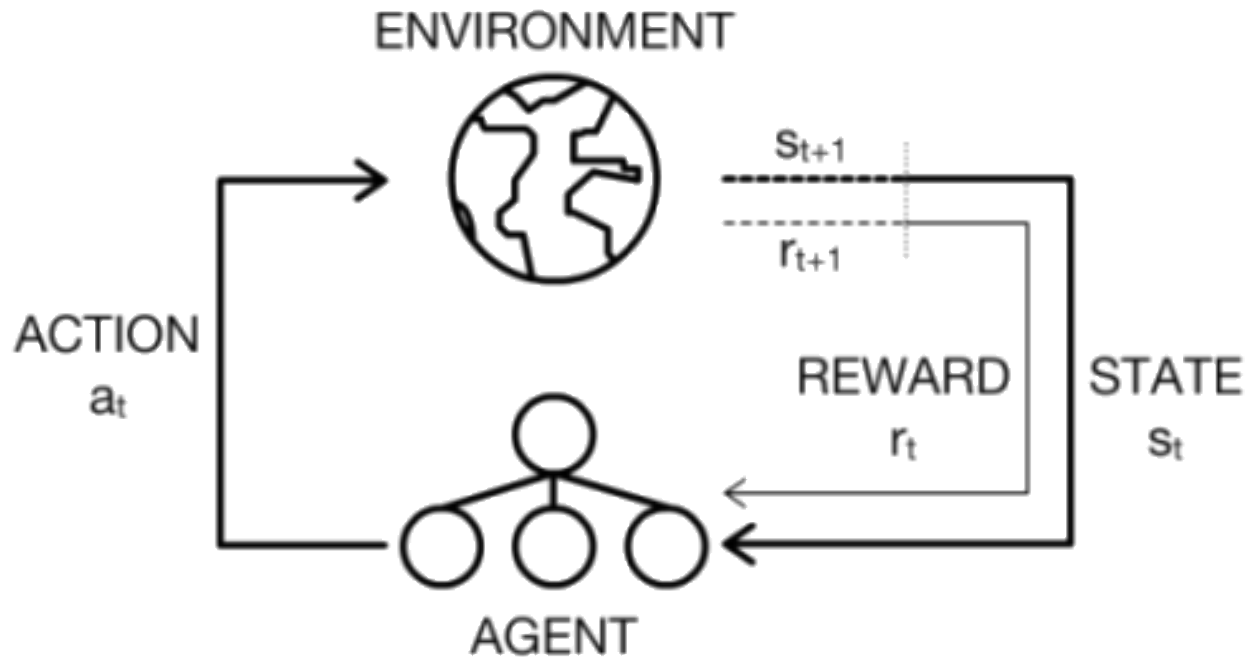
$Q(s, a)$ - Q-function

Deep Reinforcement Learning

$p(a|s)$ - policy function

$v(s)$ - value function

$Q(s, a)$ - Q-function

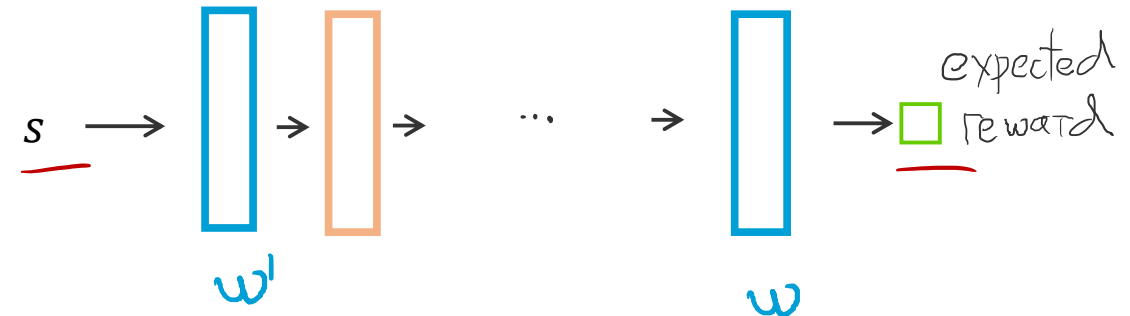
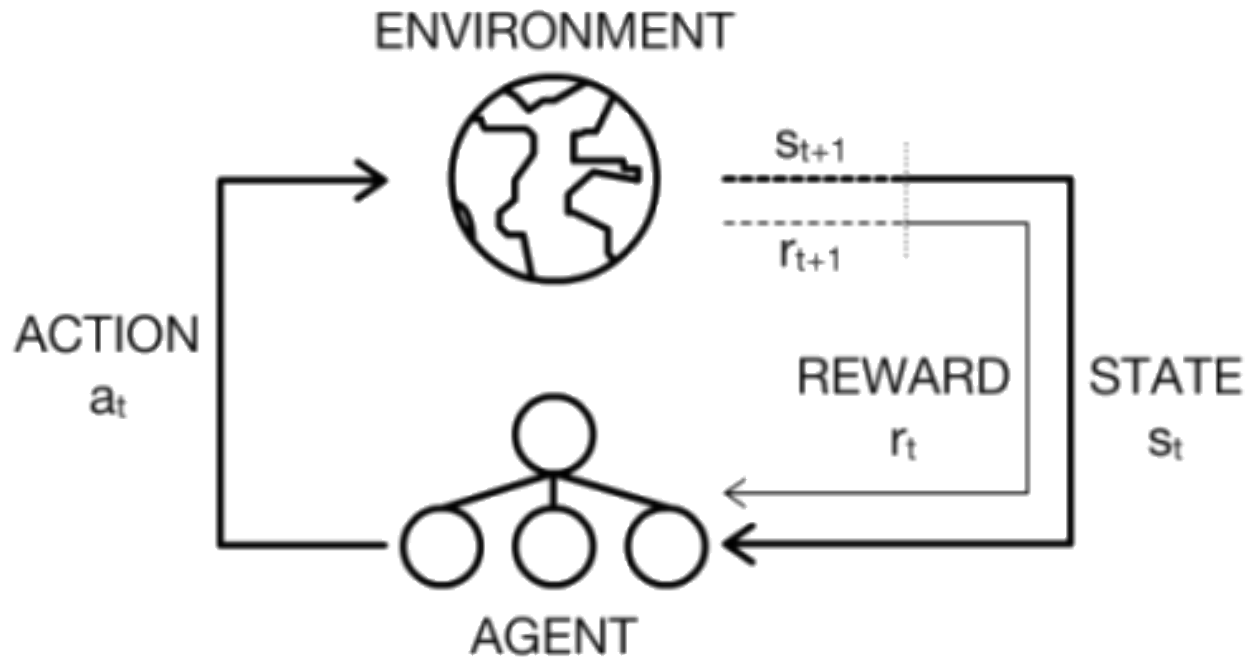


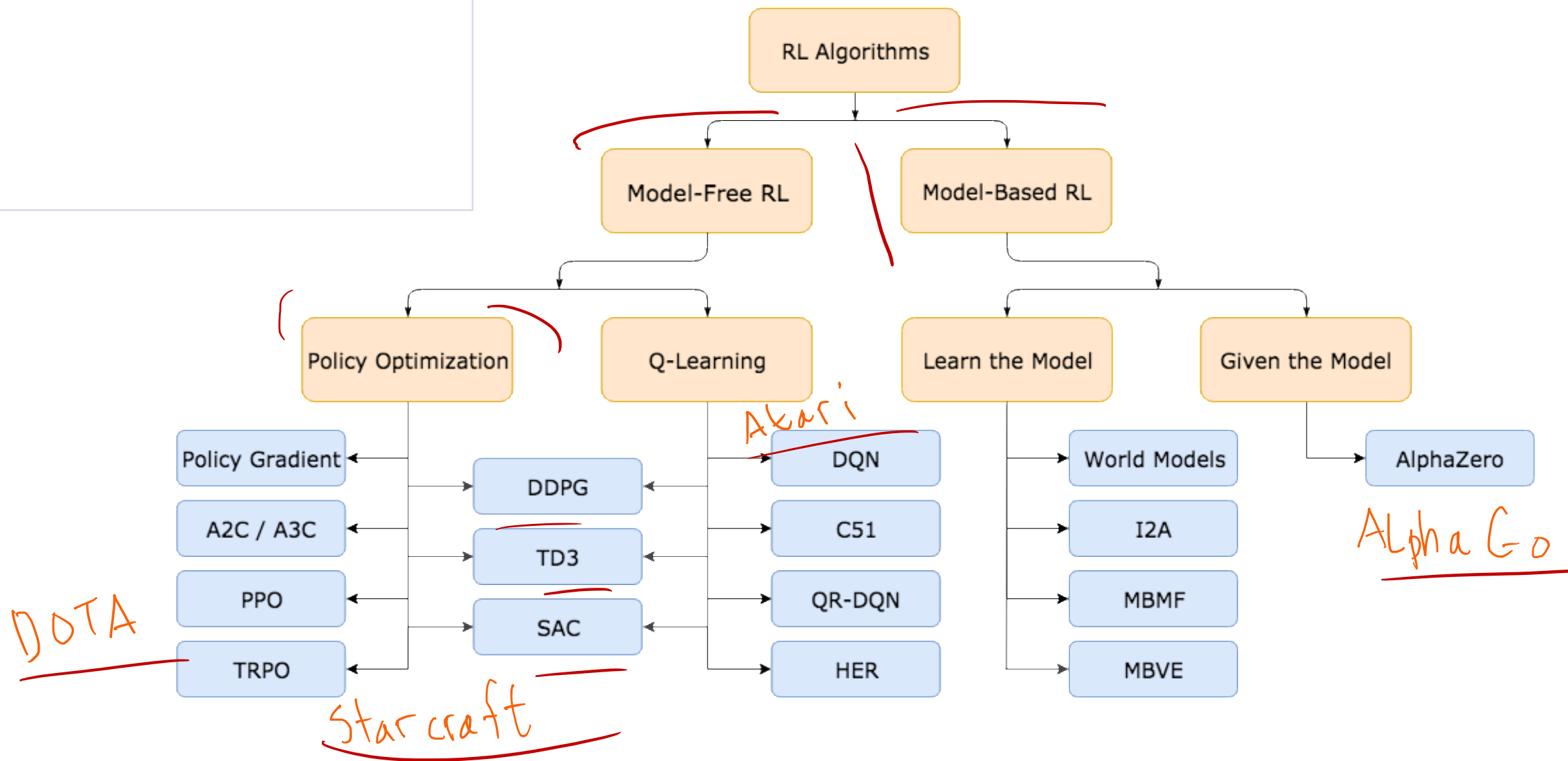
Deep Reinforcement Learning

$p(a|s)$ - policy function

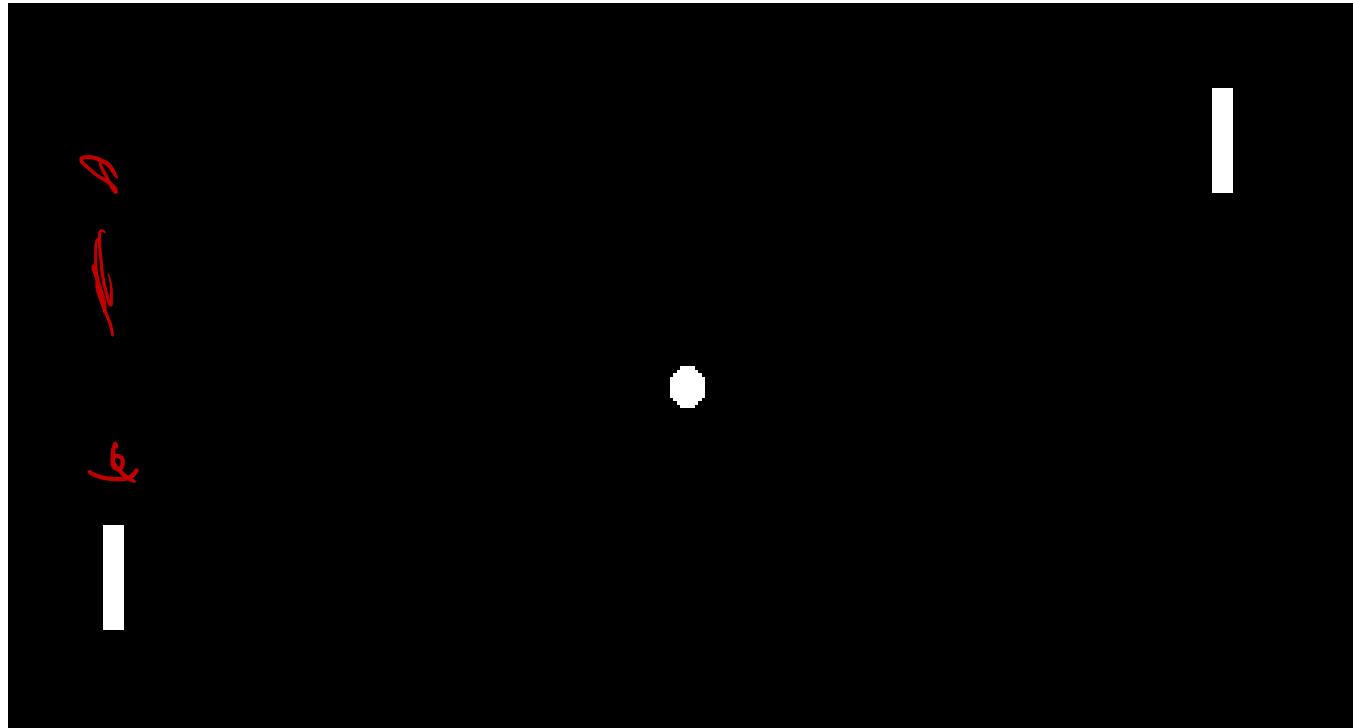
$v(s)$ - value function

$Q(s, a)$ - Q-function

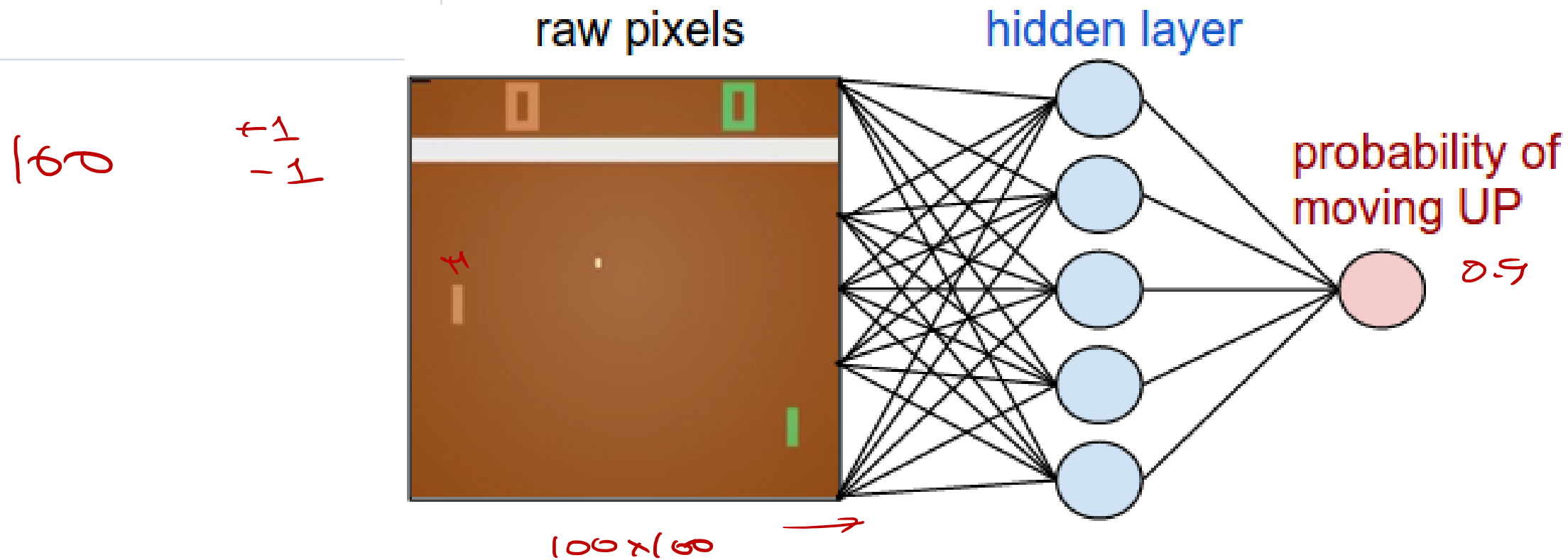




Policy Gradients на примере



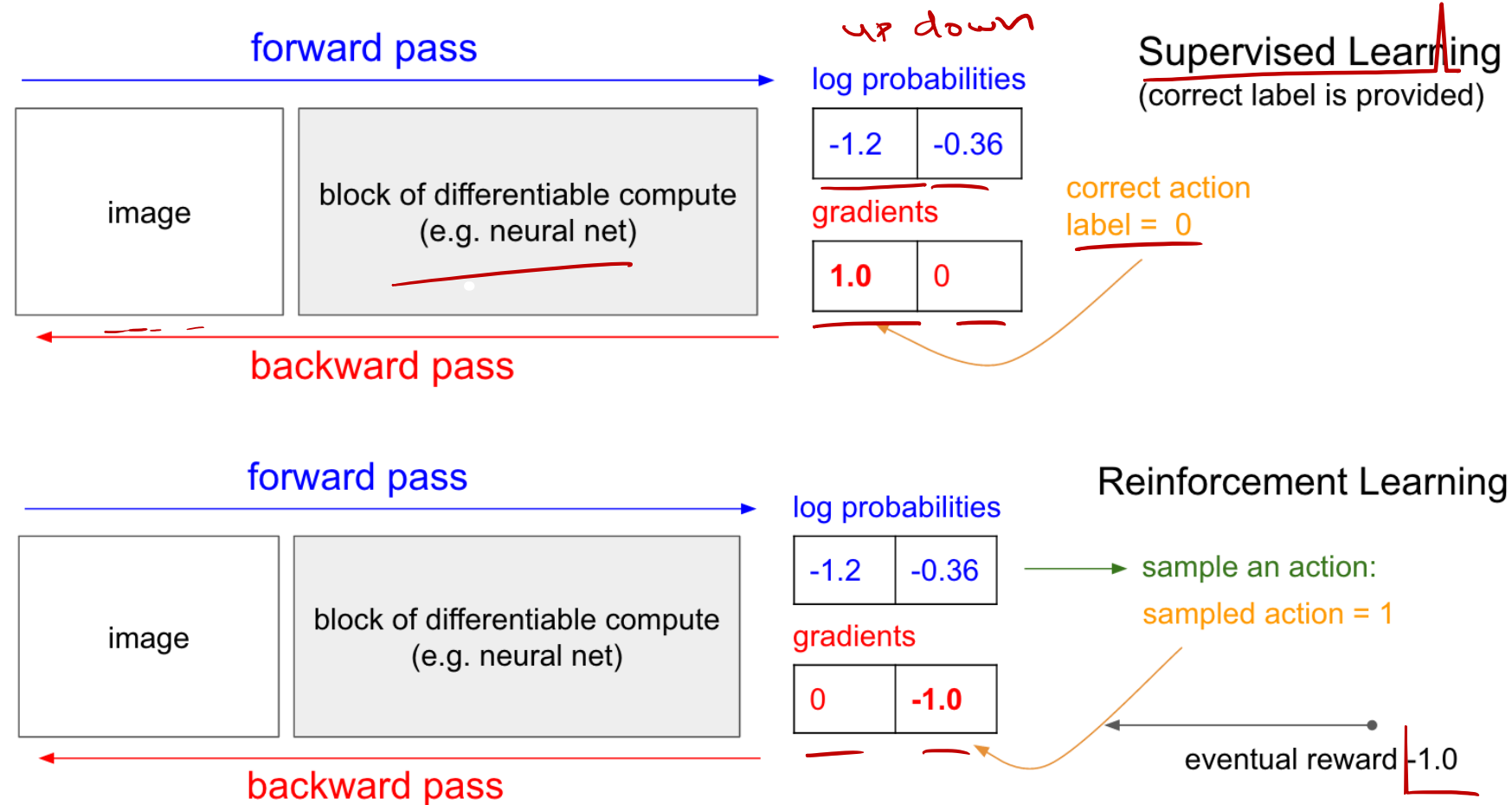
Например!



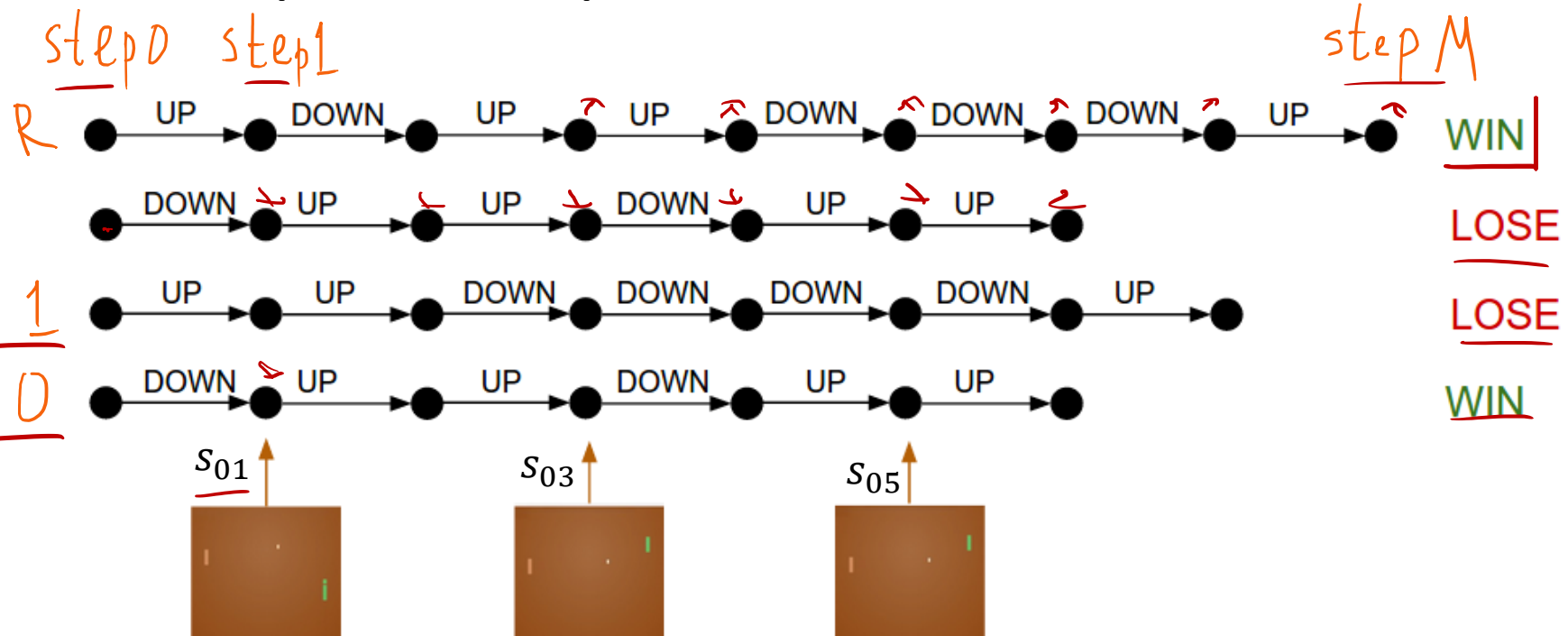
Необходимый трюк: передавать нейросети разницу между текущим и прошлым кадром

Policy Gradients

$$\sum \nabla \ln p_i$$



Набираем прогоны *rollout*



s	a	r
s_{00}	a_{00}	r_{00}
s_{01}	a_{01}	r_{01}

...

s_{RM}	a_{RM}	r_{RM}
----------	----------	----------

s_{ij} - состояние на шаге j прогона i
 a_{ij} - действие на шаге j прогона i
 r_{ij} - награда после шага j прогона i

В нашем (простом) случае награды на всех шагах прогона одинаковые $r_{ij} = r_i$

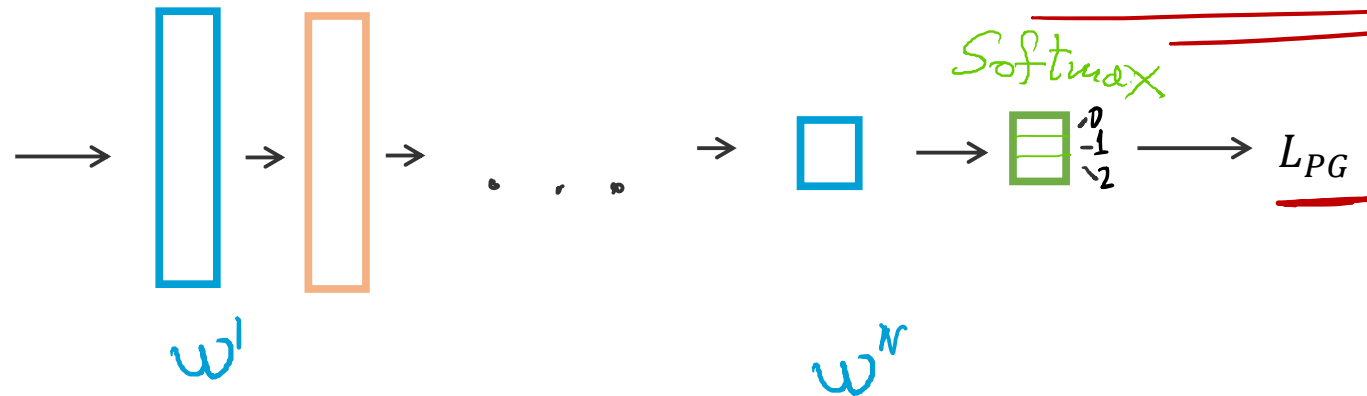
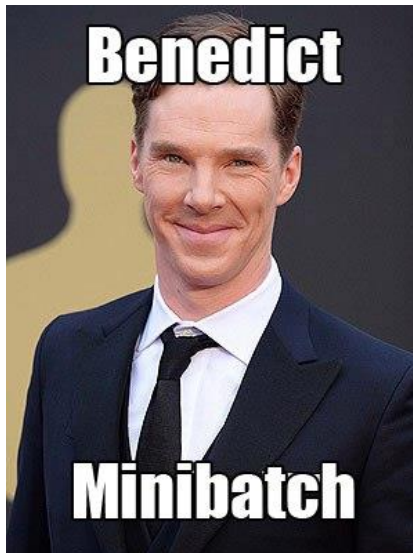
Тренируем модель

χ ψ

s_0	a_0	\bar{r}_0
s_1	a_1	<u>r_1</u>

...

s_N	a_N	r_N
-------	-------	-------



$$L_{CE} = - \sum_i \ln p(c = \underline{gt_i} | x_i)$$

$\gamma = 1$

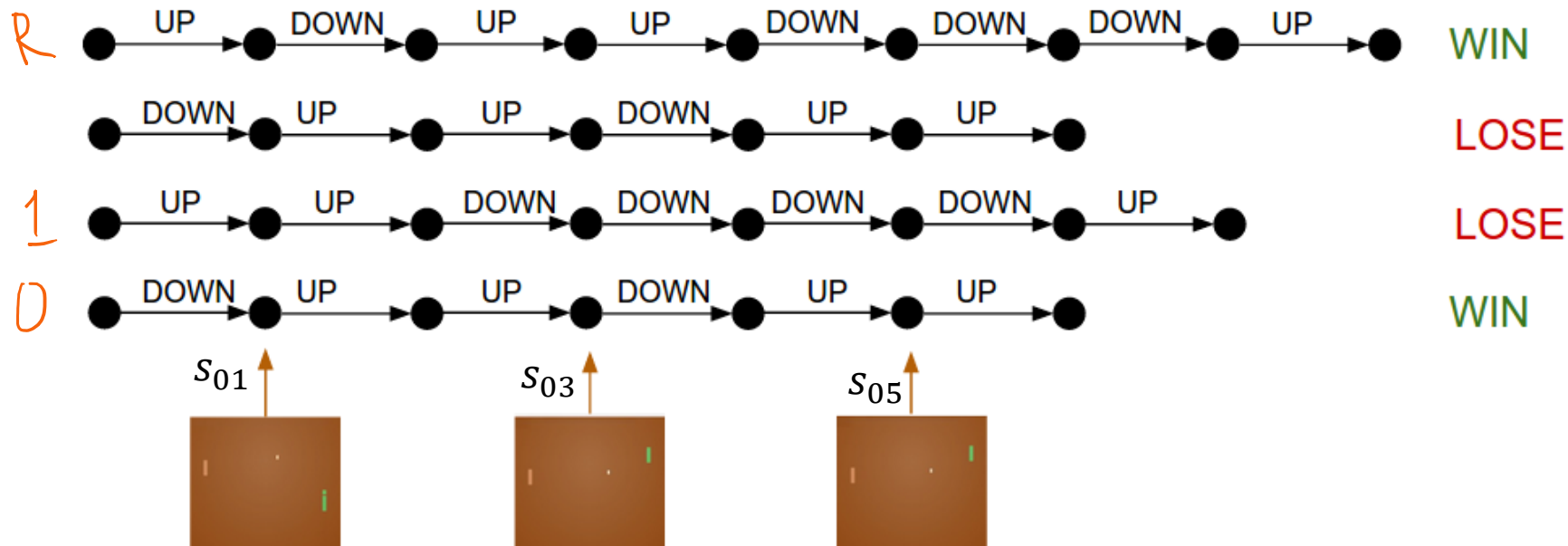
$$L_{PG} = - \sum_i \underline{r_i} \ln p(c = a_i | x_i)$$

Проходим по всем данным один раз и...

Repeat

step 0 step 1

step M

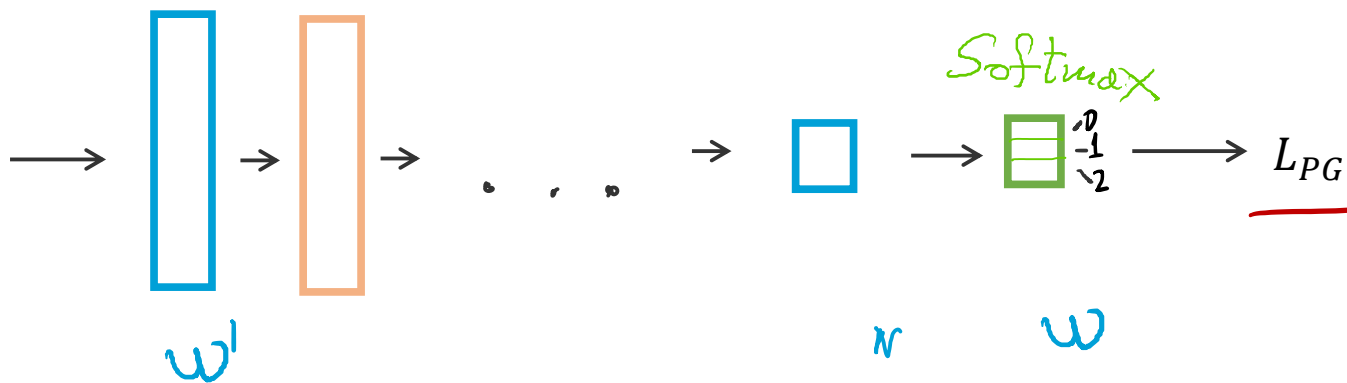


$$x \quad y$$

s_{00}	a_{00}	r_{00}
s_{01}	a_{01}	r_{01}

...

s_{RM}	a_{RM}	r_{RM}
----------	----------	----------



I KNOW

REINFORCEMENT LEARNING

memegenerator.net

Вывод на пальцах

веса сети

$p(x|s, \theta)$ - policy function

$f(x)$ - reward from environment

Хотим двигать θ так, чтобы
увеличивать $E[f(x)]$

$$\nabla_{\theta} E_x[f(x)] = \nabla_{\theta} \sum_x p(x) f(x)$$

$$= \sum_x \nabla_{\theta} p(x) f(x)$$

$$= \sum_x p(x) \left[\frac{\nabla_{\theta} p(x)}{p(x)} \right] f(x)$$

$$= \sum_x p(x) \nabla_{\theta} \log p(x) f(x)$$

$$= E_x[f(x) \nabla_{\theta} \log p(x)]$$

$$\sum_i \tau_i \log p(s_i)$$

definition of expectation

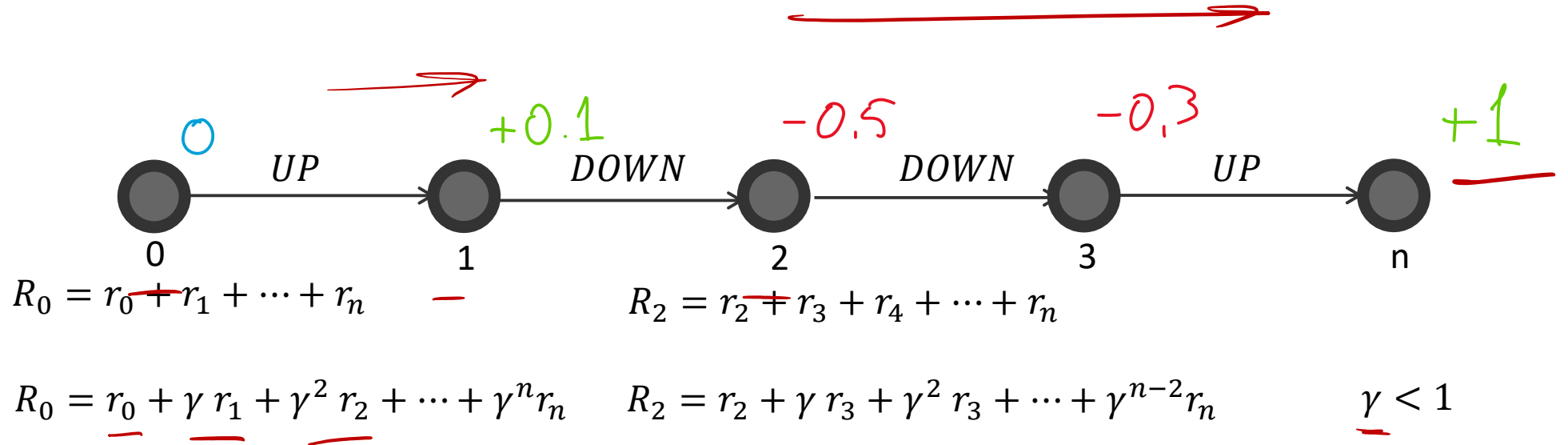
swap sum and gradient

both multiply and divide by $p(x)$

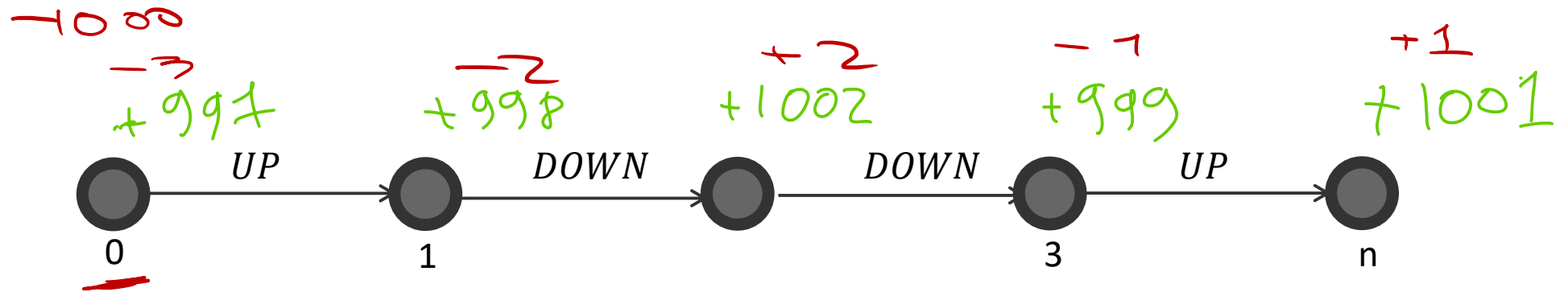
use the fact that $\nabla_{\theta} \log(z) = \frac{1}{z} \nabla_{\theta} z$

definition of expectation

Discounted rewards



Baseline



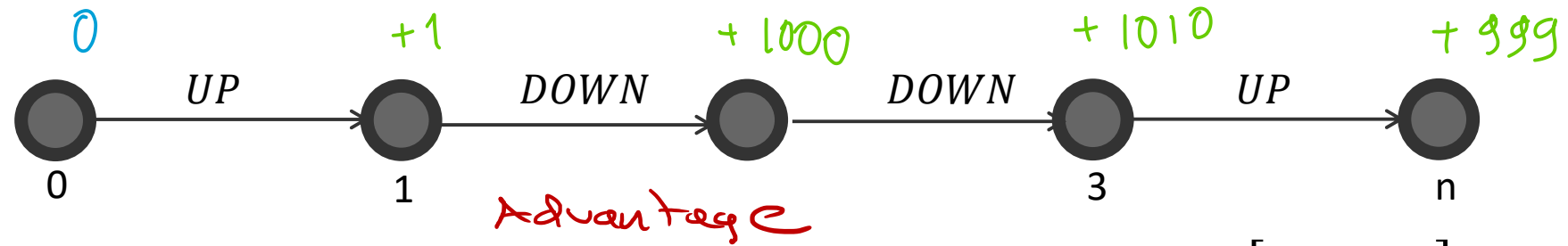
$$R_s = \sum_{i=0 \dots n-s} \gamma^i r_{s+i}$$

$$\gamma < 1$$

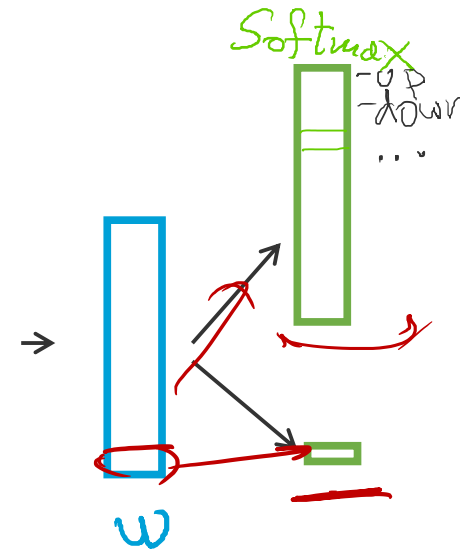
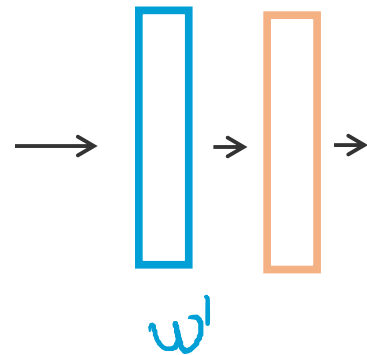
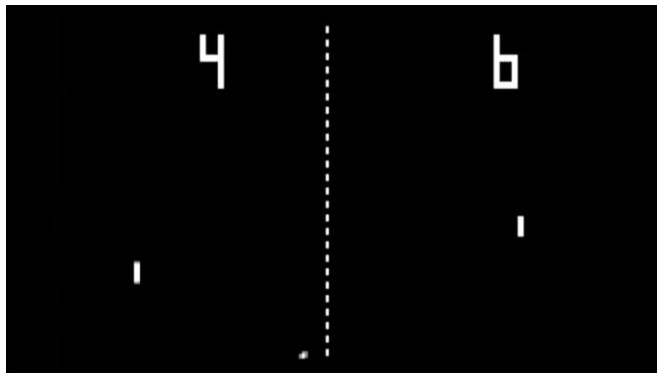
$$R_s = \sum_{i=0 \dots n-s} \gamma^i r_{s+i} - \underline{b}$$

$$b = E_{rollouts} \left[\sum_i \gamma^i r_{s+i} \right]$$

Actor-Critic



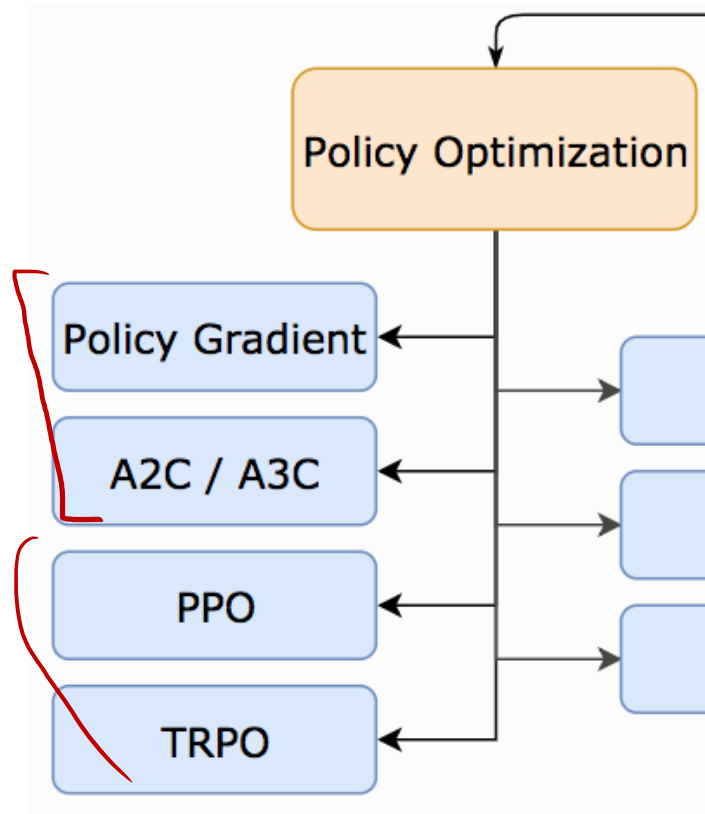
$$\underline{R_s} = \sum_{i=0 \dots n-s} \gamma^i r_{s+i} - b \quad b = \underline{E_{\pi(s)} \left[\sum_i \gamma^i r_{s+i} \right]}$$



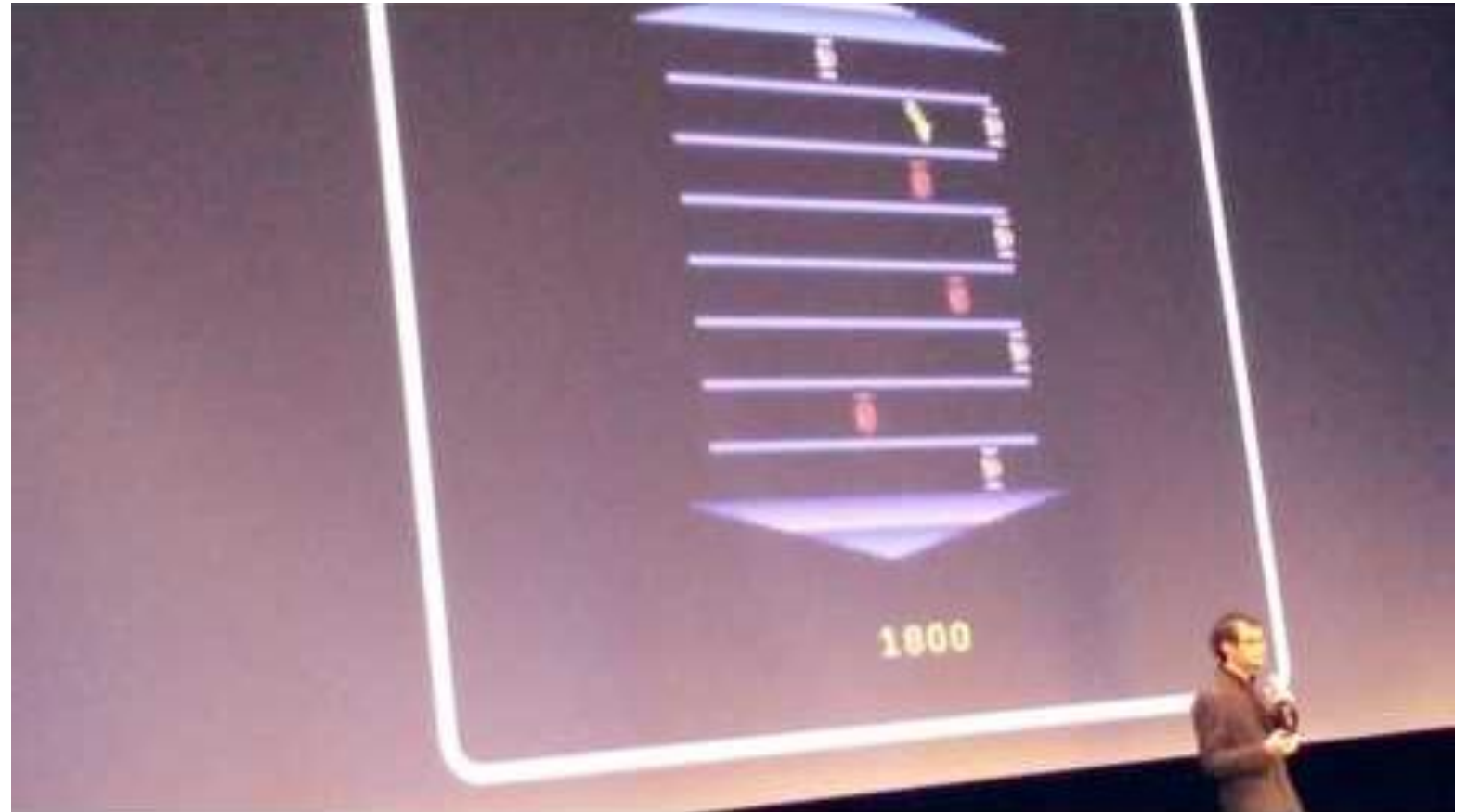
$$\log p(a) \sum_i (\gamma^i r_{s+i} - V)$$

$$\underline{L_2(V - \sum_i \gamma^i r_{s+i})}$$

Прогресс продолжается



Deep Q Learning (DQN)



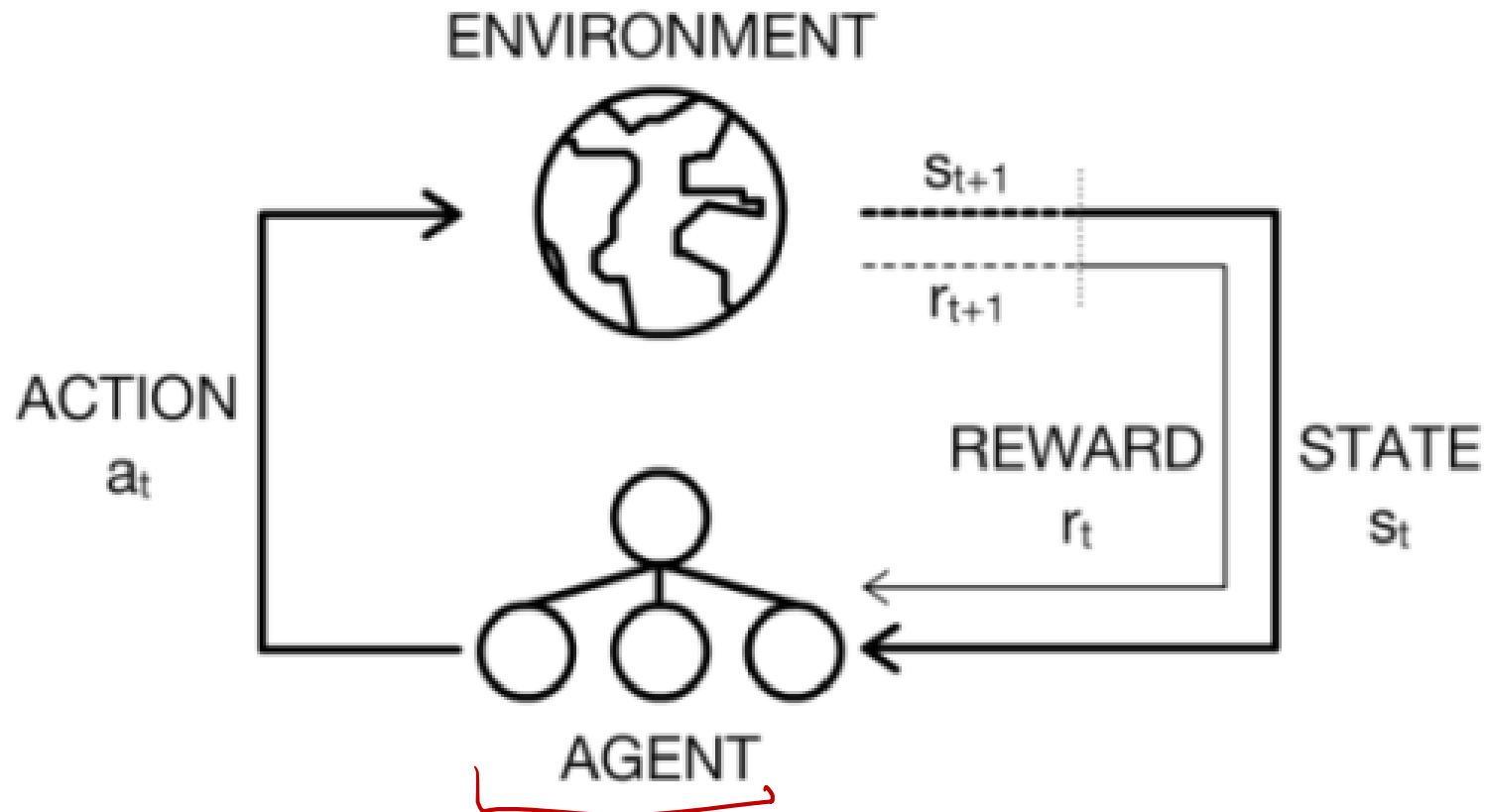
[Google Acquires Artificial Intelligence Startup DeepMind For More Than \\$500M](#)

Так. Приготовьтесь!

Сейчас будет сложно.

Q-function

$$Q(s, a) \rightarrow \tau$$



~~$p(a|s)$ - policy function~~

~~$v(s)$ - value function~~

$Q(s, a)$ - Q-function

$Q^*(s, a)$ - Q-функция оптимальной стратегии

Bellman equation

$\gamma < 1$

$$\underline{Q^*(s, a)} = \mathbb{E}_{\underline{s' \sim \mathcal{E}}} \left[\underline{r} + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

$Q^*(s, a)$ – Q-функция оптимальной стратегии

\mathcal{E} – пространство возможных следующих состояний

Итеративное приближение:

$$Q_{i+1}(s, a) = \mathbb{E} [r + \gamma \max_{a'} Q_i(s', a') \mid s, a]$$

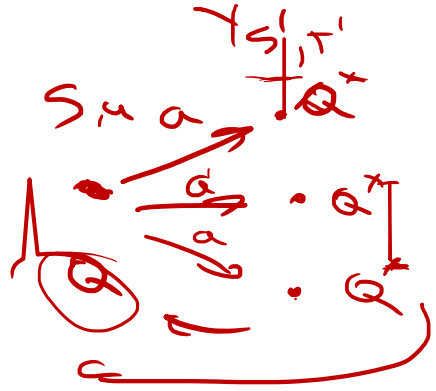
Приближаем Q^* нейросетью с параметрами θ :

$$\underline{Q(s, a; \theta)} \approx Q^*(s, a)$$

Формулируем функцию ошибки как приближение Q^* :

$$\underline{L_i(\theta_i)} = \mathbb{E}_{s, a \sim \rho(\cdot)} \left[(\underline{y_i} - \underline{Q(s, a; \theta_i)})^2 \right]$$

$$\underline{y_i} = \mathbb{E}_{s' \sim \mathcal{E}} \left[\underline{r + \gamma \max_{a'} Q(s', a'; \theta_{i-1})} \mid s, a \right]$$



Алгоритм DQN

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

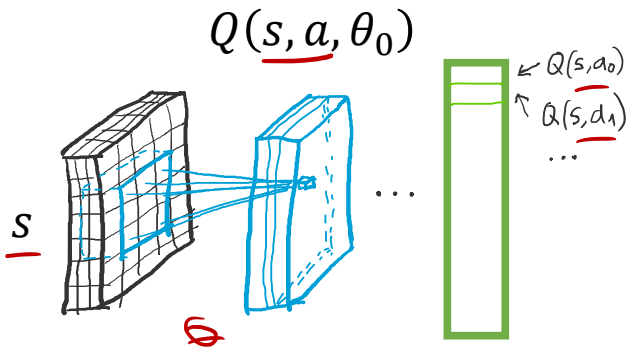
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Алгоритм DQN

D



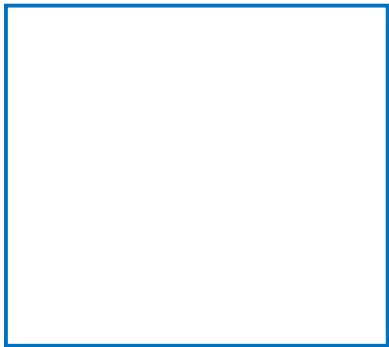
Algorithm 1 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1,  $M$  do
  Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ 
  for  $t = 1, T$  do
    With probability  $\epsilon$  select a random action  $a_t$ 
    otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ 
    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  $\mathcal{D}$ 
    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
    Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
    Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
  end for
end for
  
```

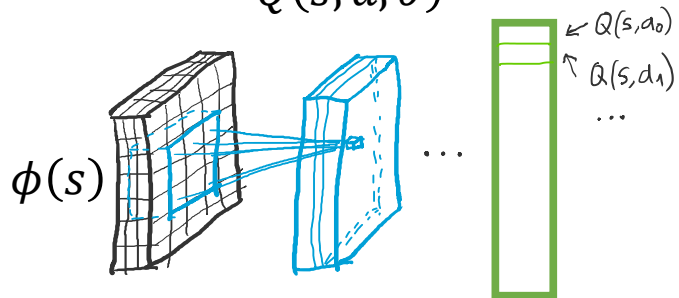
Алгоритм DQN

D



ϕ_1

$Q(s, a, \theta)$



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

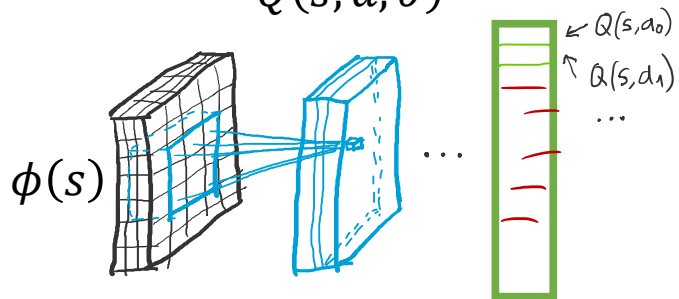
Алгоритм DQN

D



$\phi_1 a_1$

$Q(s, a, \theta)$



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t
 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

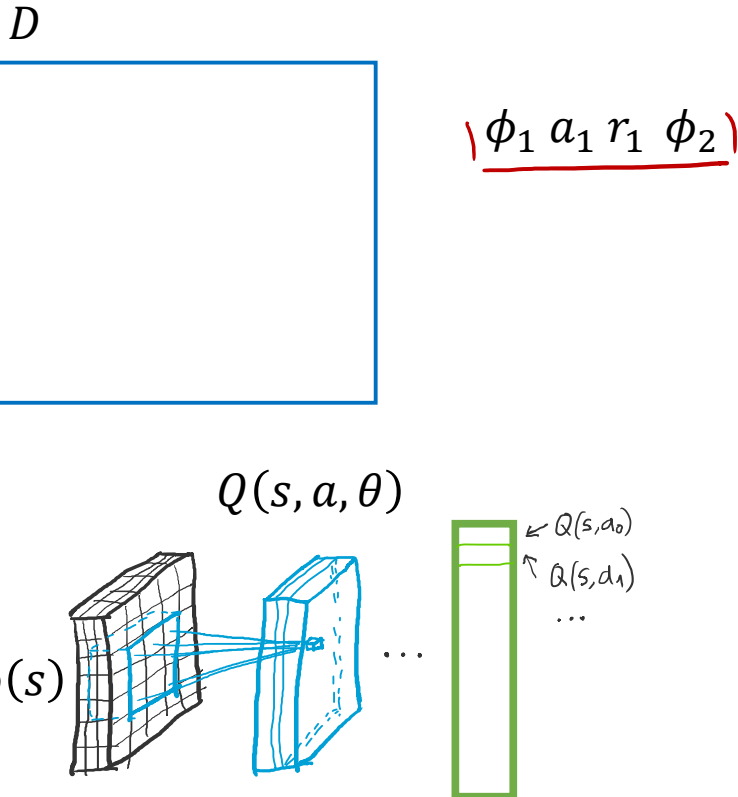
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Алгоритм DQN



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

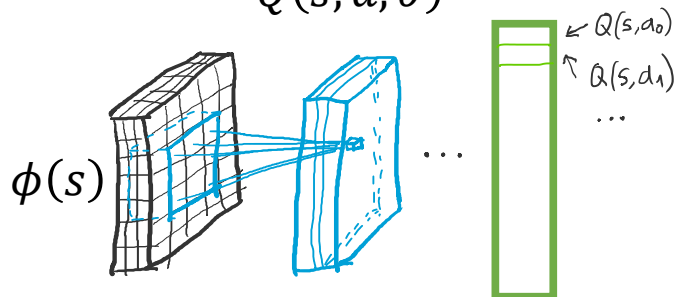
Алгоритм DQN

D

$\phi_1 \ a_1 \ r_1 \ \phi_2$
 $\phi_2 \ a_2 \ r_2 \ \phi_3$
 \dots
 $\phi_n \ a_n \ r_n \ \phi_{n+1}$

$\phi_1 \ a_1 \ r_1 \ \phi_2$

$Q(s, a, \theta)$



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

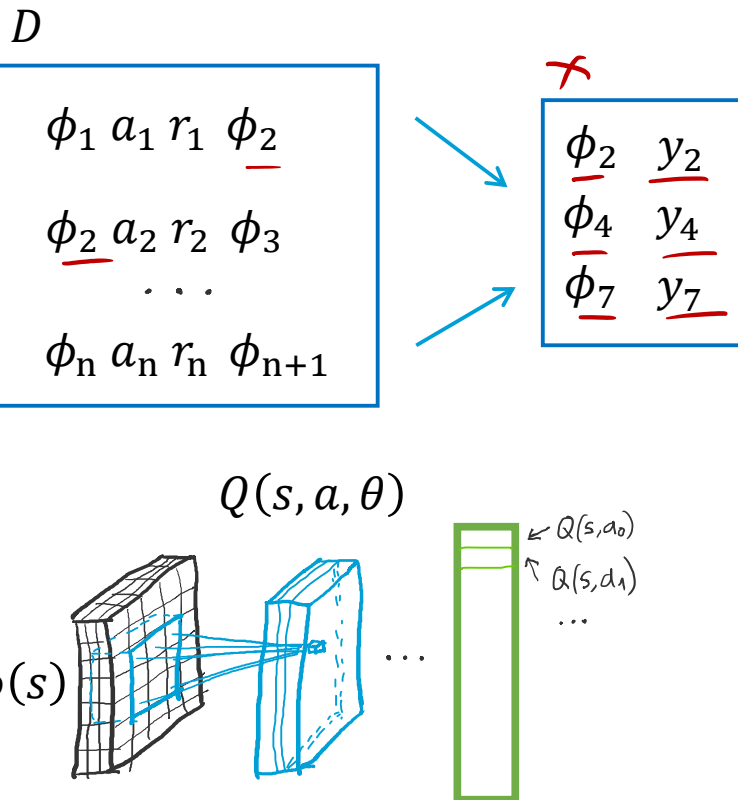
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Алгоритм DQN



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

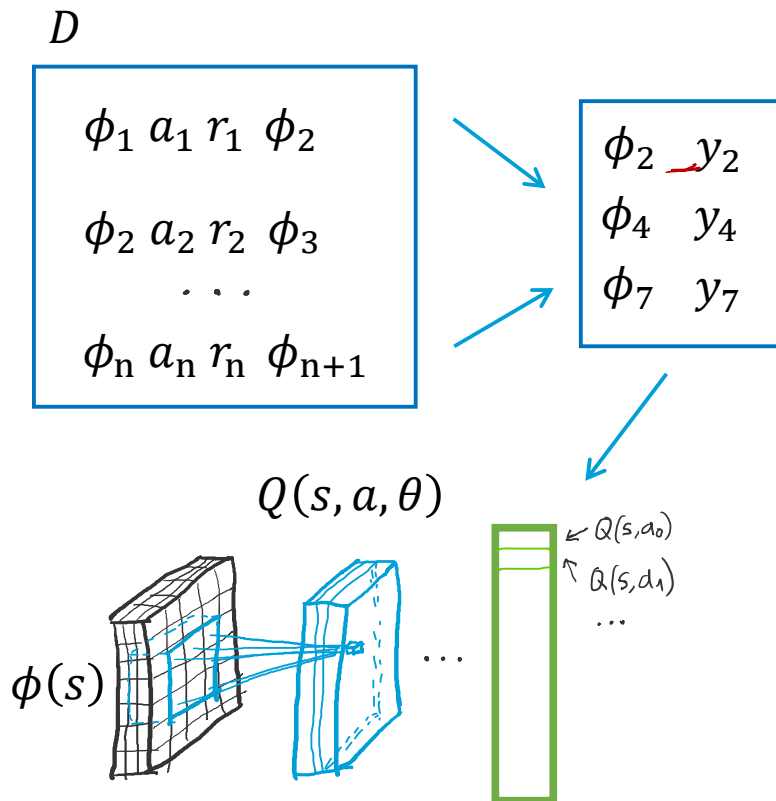
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ \frac{r_j}{r_j} + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Алгоритм DQN



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

With probability ϵ select a random action a_t

otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

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 \mathcal{D}

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

Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

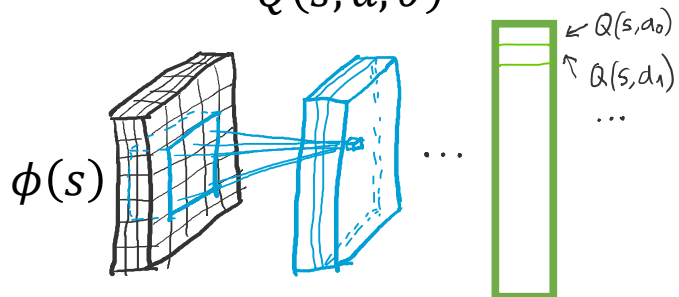
end for

Алгоритм DQN

D

ϕ_1	a_1	r_1	ϕ_2
ϕ_2	a_2	r_2	ϕ_3
...			
ϕ_n	a_n	r_n	ϕ_{n+1}

$Q(s, a, \theta)$



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

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

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

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 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 \mathcal{D}

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

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

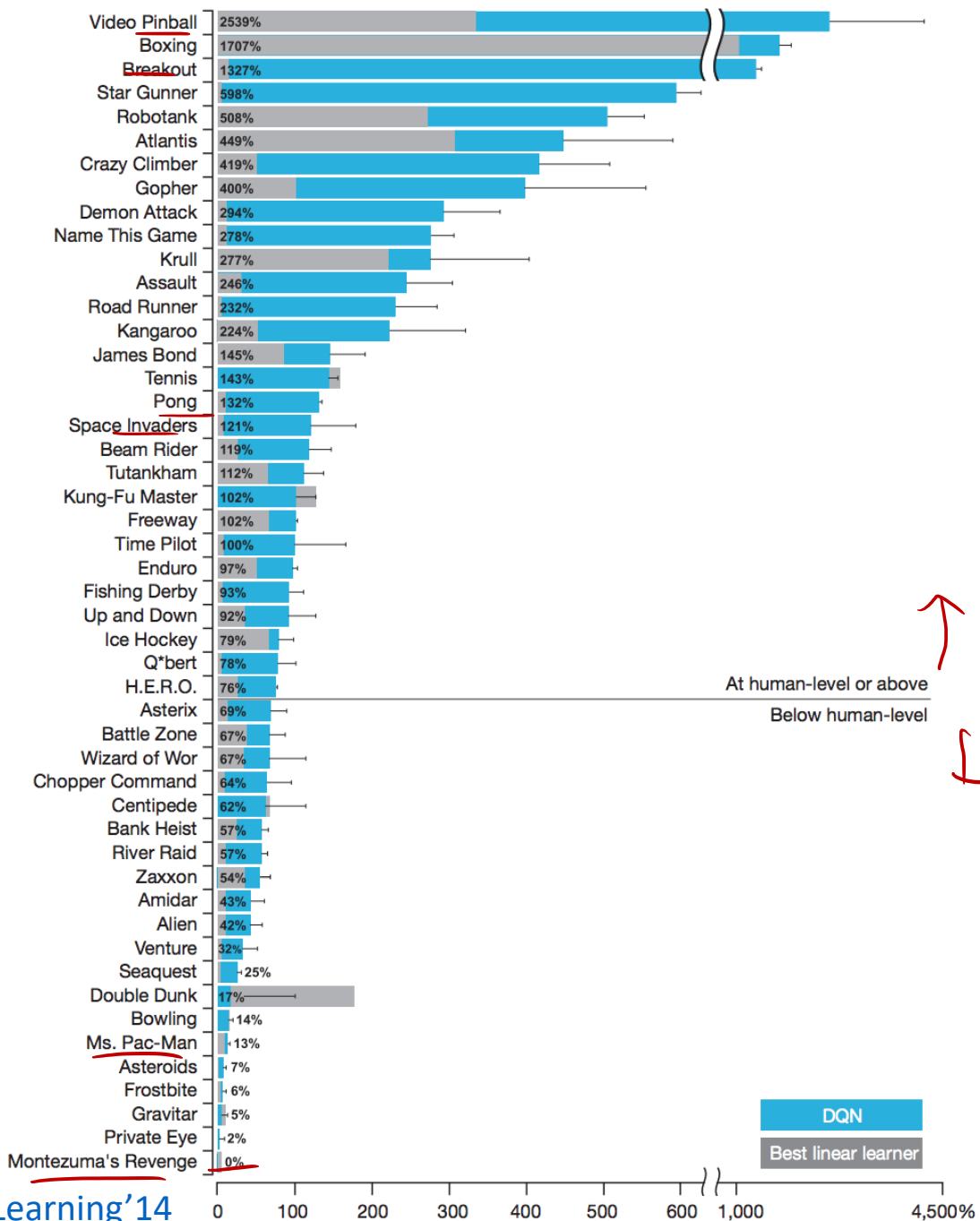
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

[Also, Habr](#)





Q-learning vs policy gradients

Policy Gradients

On-policy

Эффективнее исследует

Проще

Чаще работает

Q-learning

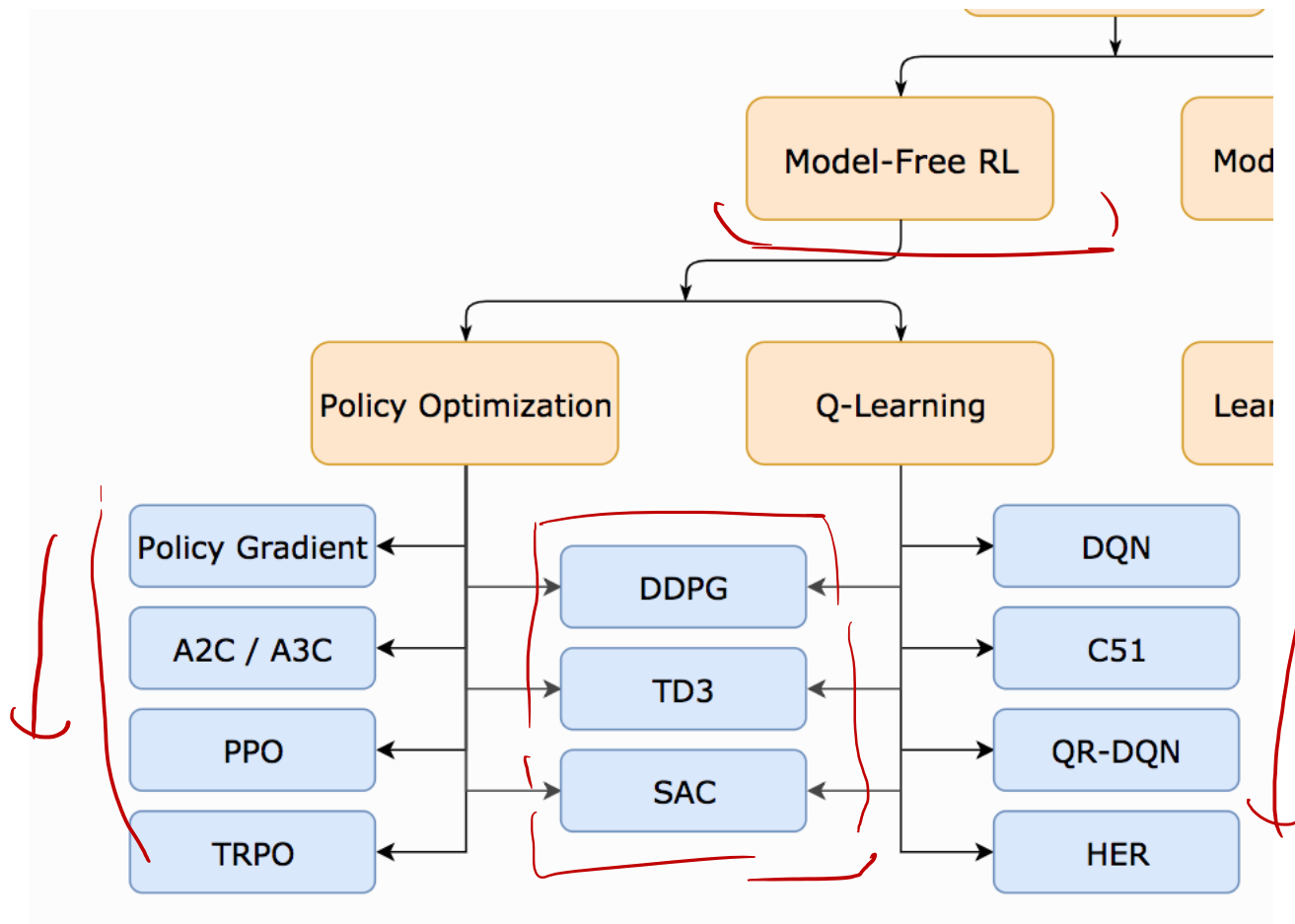
Off-policy

Эффективнее сэмплирует

Сложнее

Когда работает, сходится стабильнее

Дальнейшая работа



WHEN YOU WAKE UP



Как с этим поиграться?



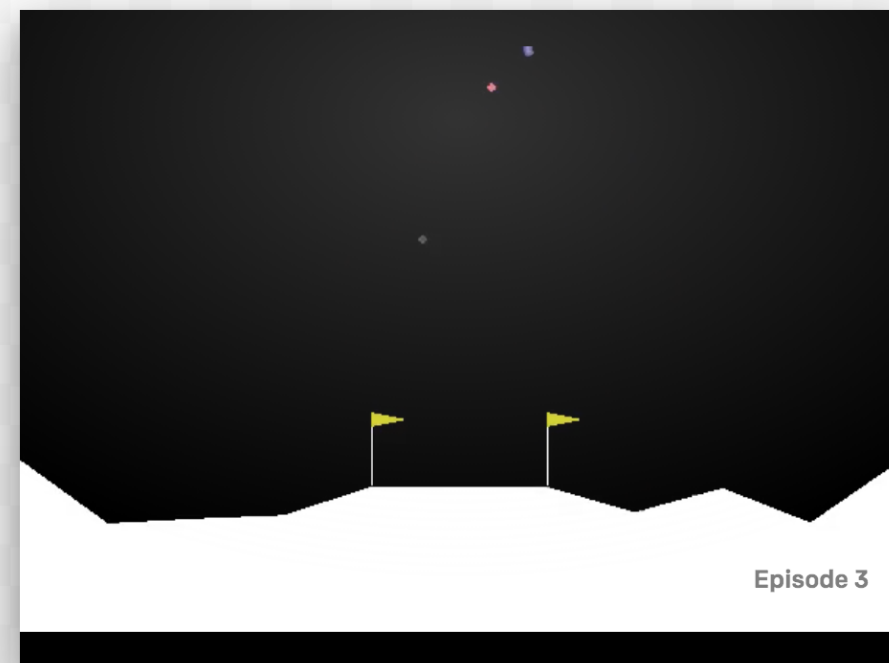
Gym

Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.

[View documentation >](#)

[View on GitHub >](#)

RandomAgent on SpaceInvaders-v0



RandomAgent on LunarLander-v2

Создаем environment

```
pip install gym  
pip install atari-py
```

```
import gym  
env = gym.make('MsPacman-v0')  
env.reset()  
print("Action space:", env.action_space)  
print("Action meanings:", env.unwrapped.get_action_meanings())  
print("Observation space:", env.observation_space)
```

Action space: Discrete(9)

Action meanings: ['NOOP', 'UP', 'RIGHT', 'LEFT', 'DOWN', 'UPRIGHT']

Observation space: Box(210, 160, 3)

Environments

Внезапный бонус
для тех кто сделал
задания 5 и 6!

Algorithms

Atari

Box2D

Classic control

MuJoCo

Robotics **NEW**

Toy text **EASY**

Box2D

Continuous control tasks in the Box2D simulator.



BipedalWalker-v2
Train a bipedal robot to walk.



BipedalWalkerHardcore-v2
Train a bipedal robot to walk over rough terrain.



CarRacing-v0
Race a car around a track.



ПУТЬ ТВОЙ



ТРУДЕН БУДЕТ

ICP
ICLR