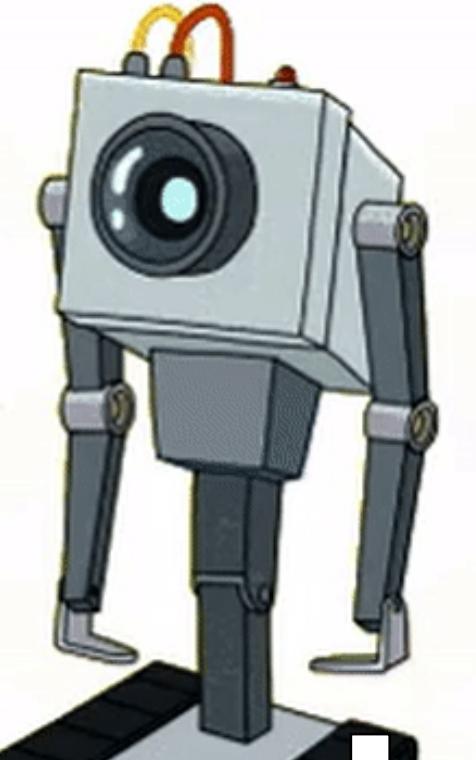


Обучение с подкреплением

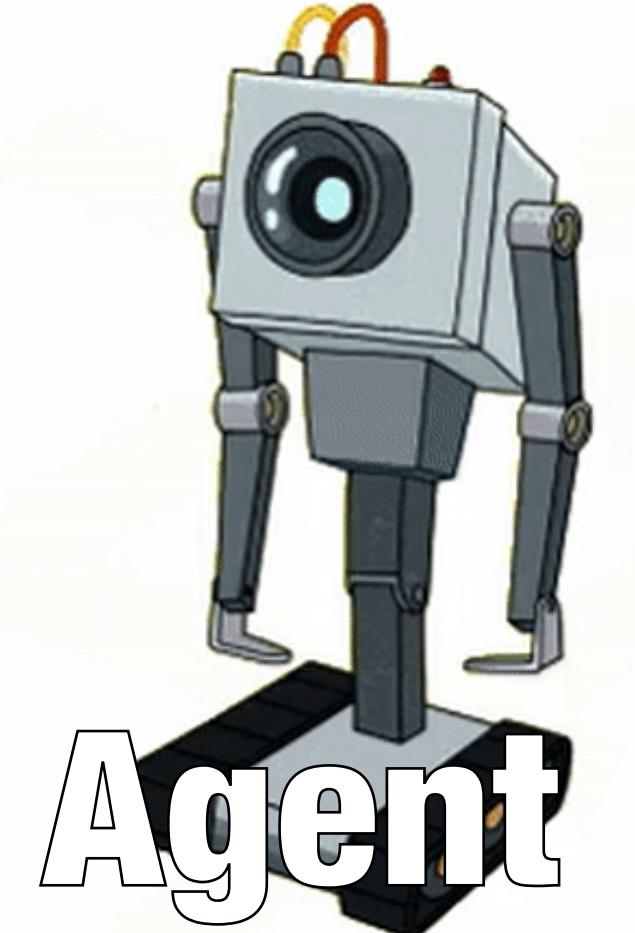
Lecture 2 - Deep Reinforcement Learning

Власов Кирилл Вячеславович



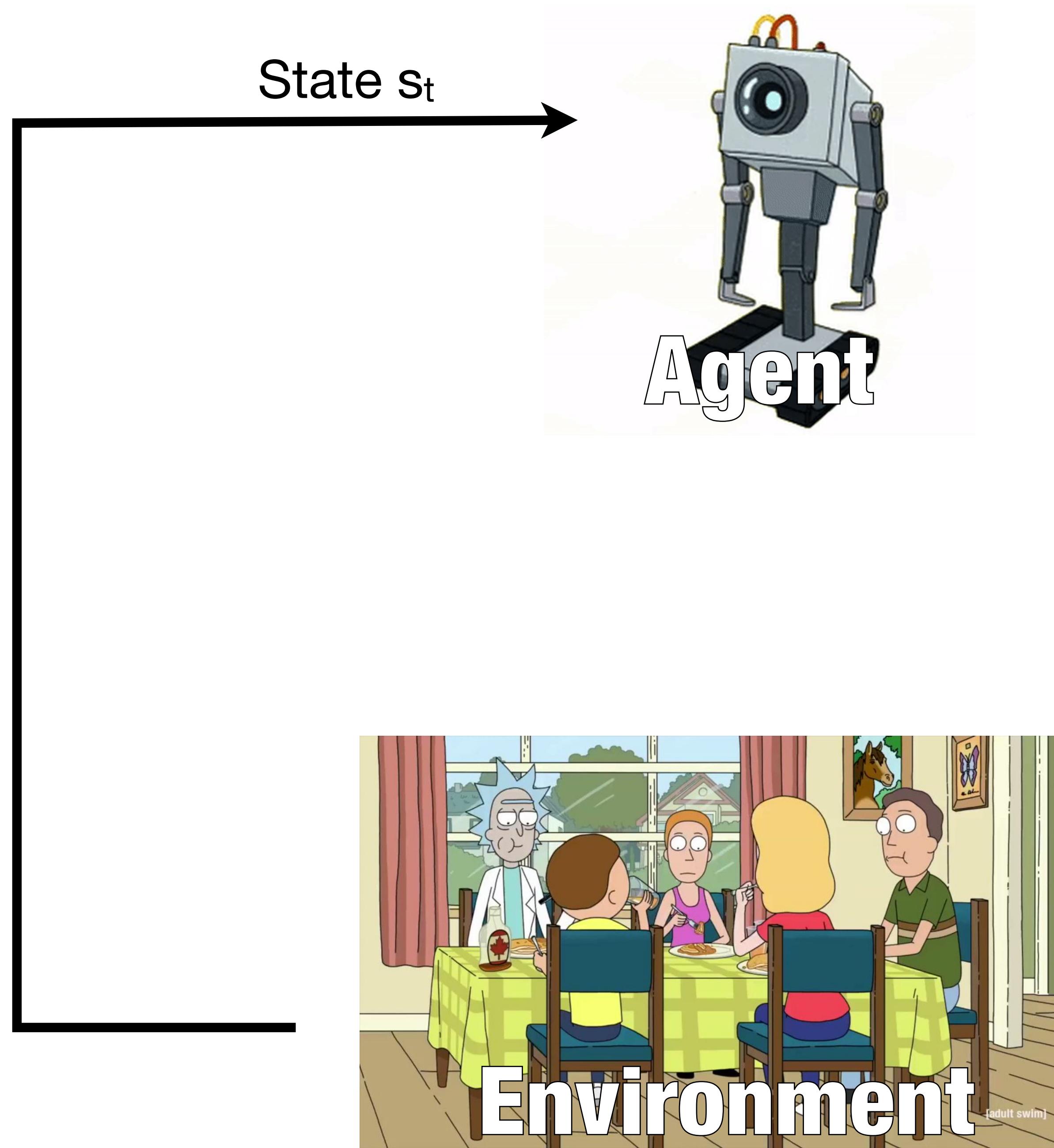


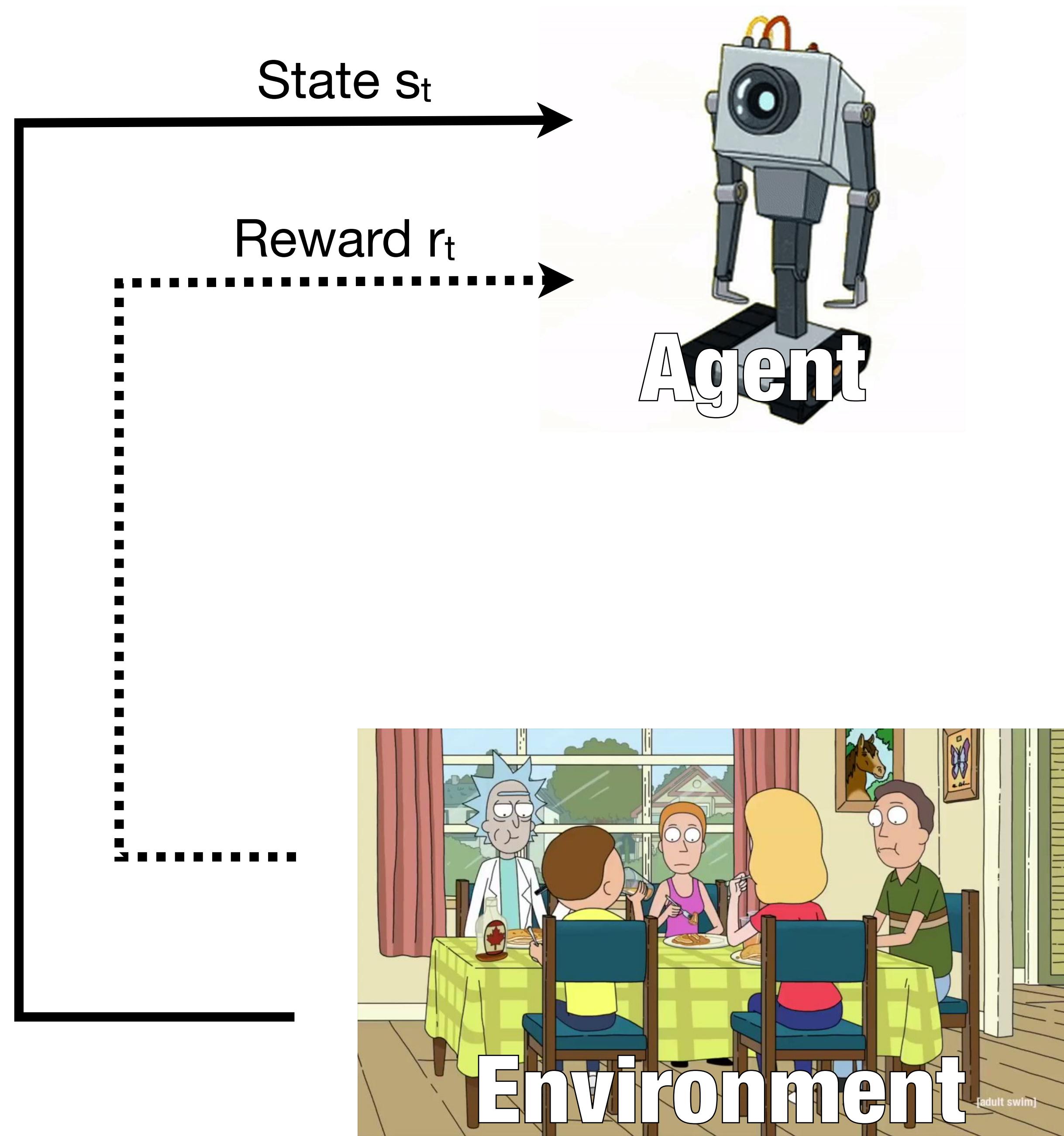
Agent

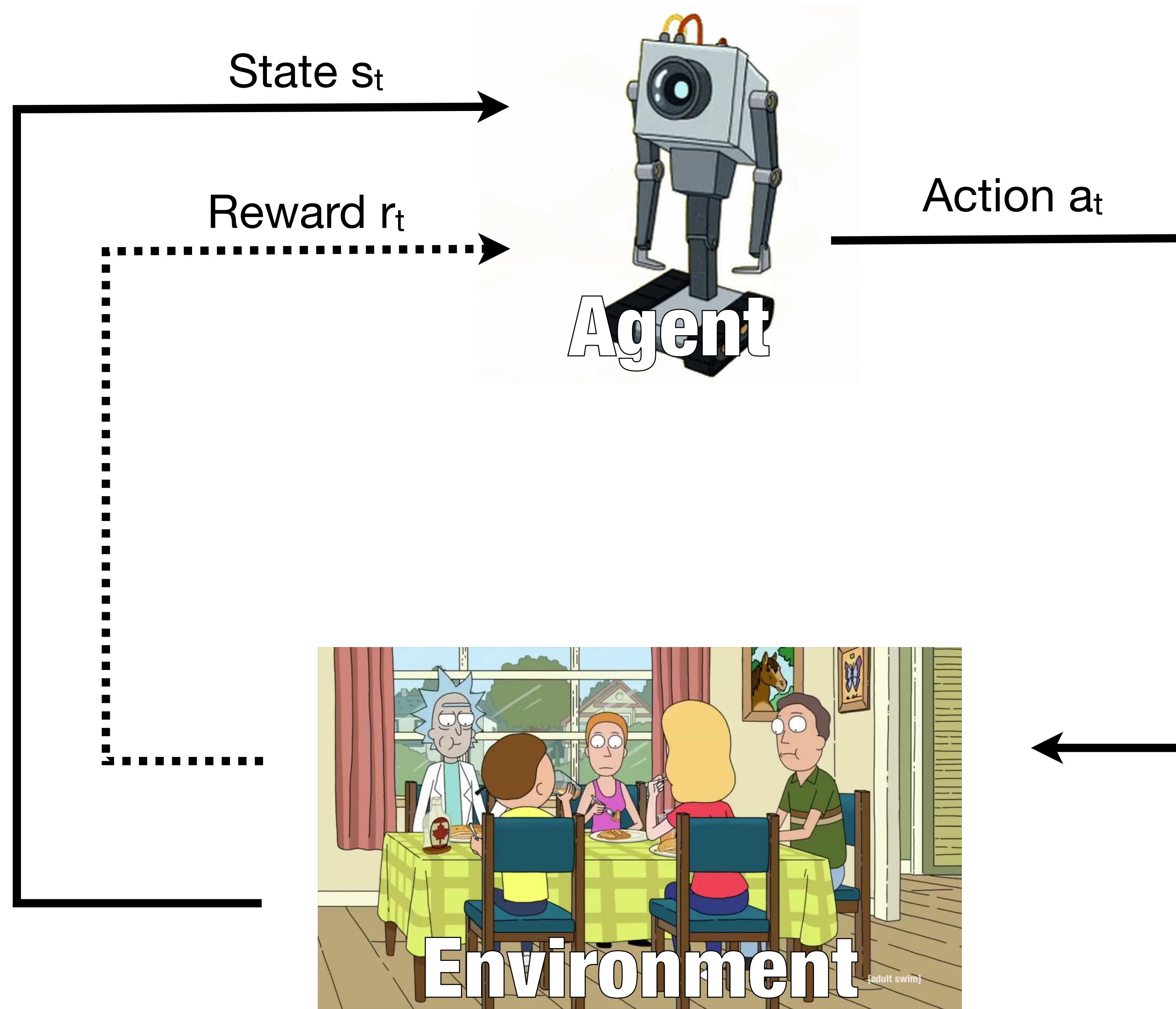


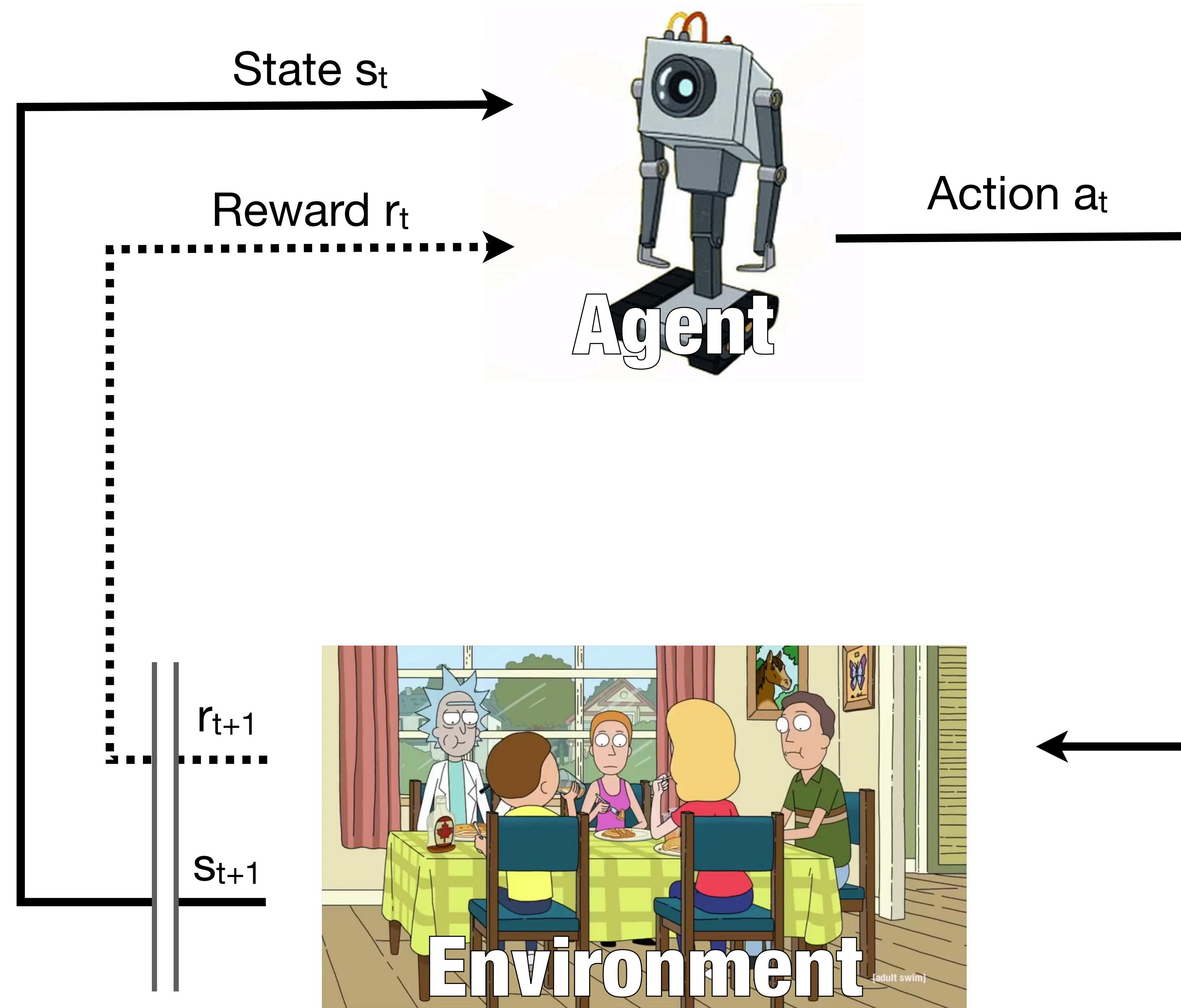
Environment

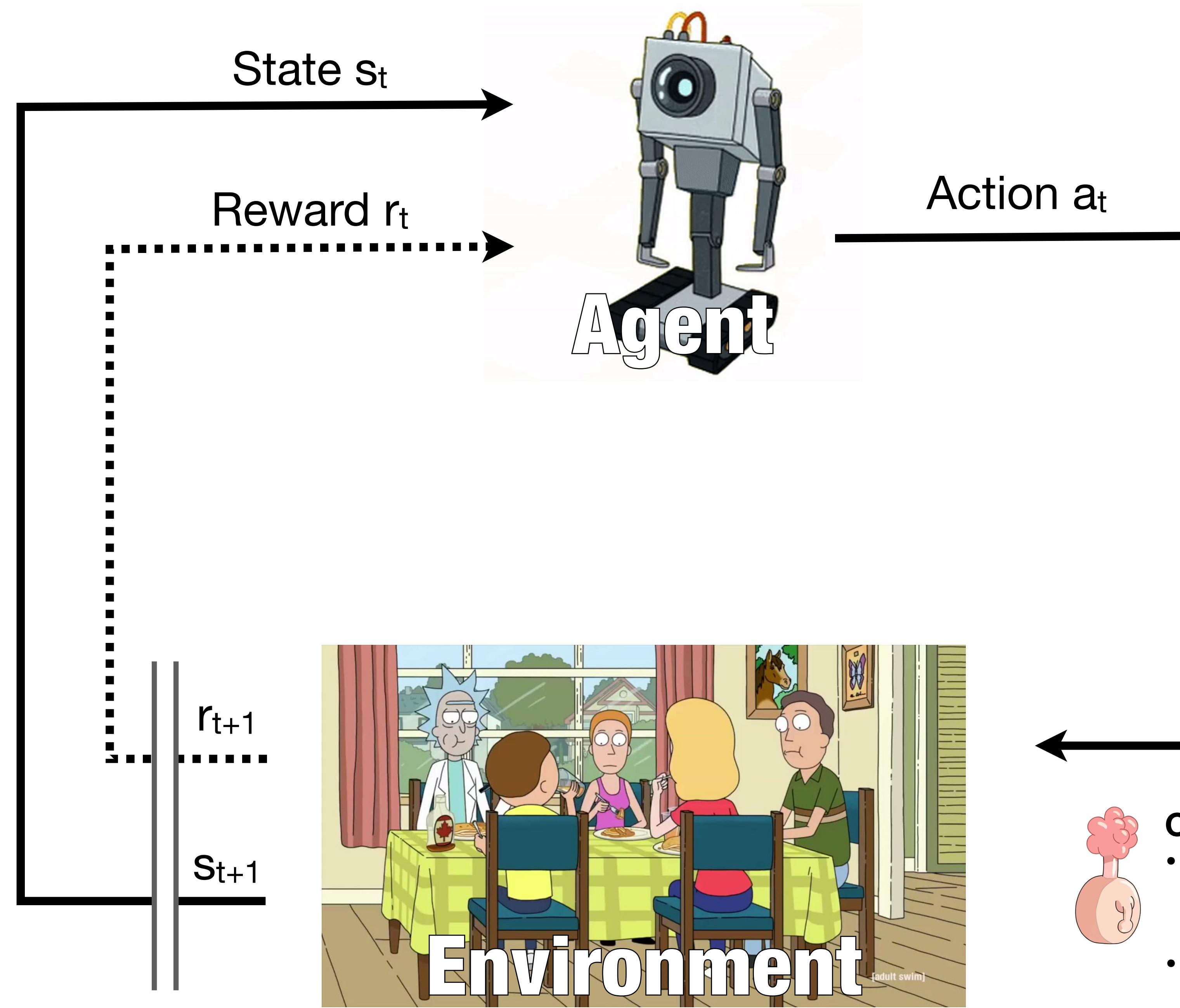
State s_t











Среда может:

- быть не детерминированной (выдавать каждый раз разный стейт)
- внутреннее состояние (которое явно можно не наблюдать в стейте)

$$\sum_{t=1}^n r_t \longrightarrow max$$



$$\sum_{t=1}^n r_t \longrightarrow \max$$

$p(a | s)$ policy function

$v(s)$ value function

$Q(s, a)$ Q-function



$$\sum_{t=1}^n r_t \longrightarrow \max$$



$p(a | s)$

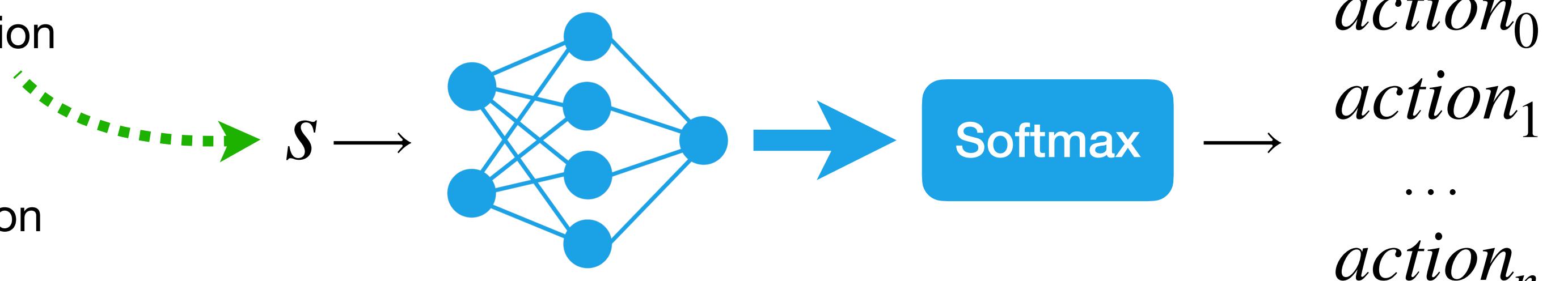
$v(s)$

$Q(s, a)$

policy function

value function

Q-function



$$\sum_{t=1}^n r_t \longrightarrow \max$$



$p(a | s)$

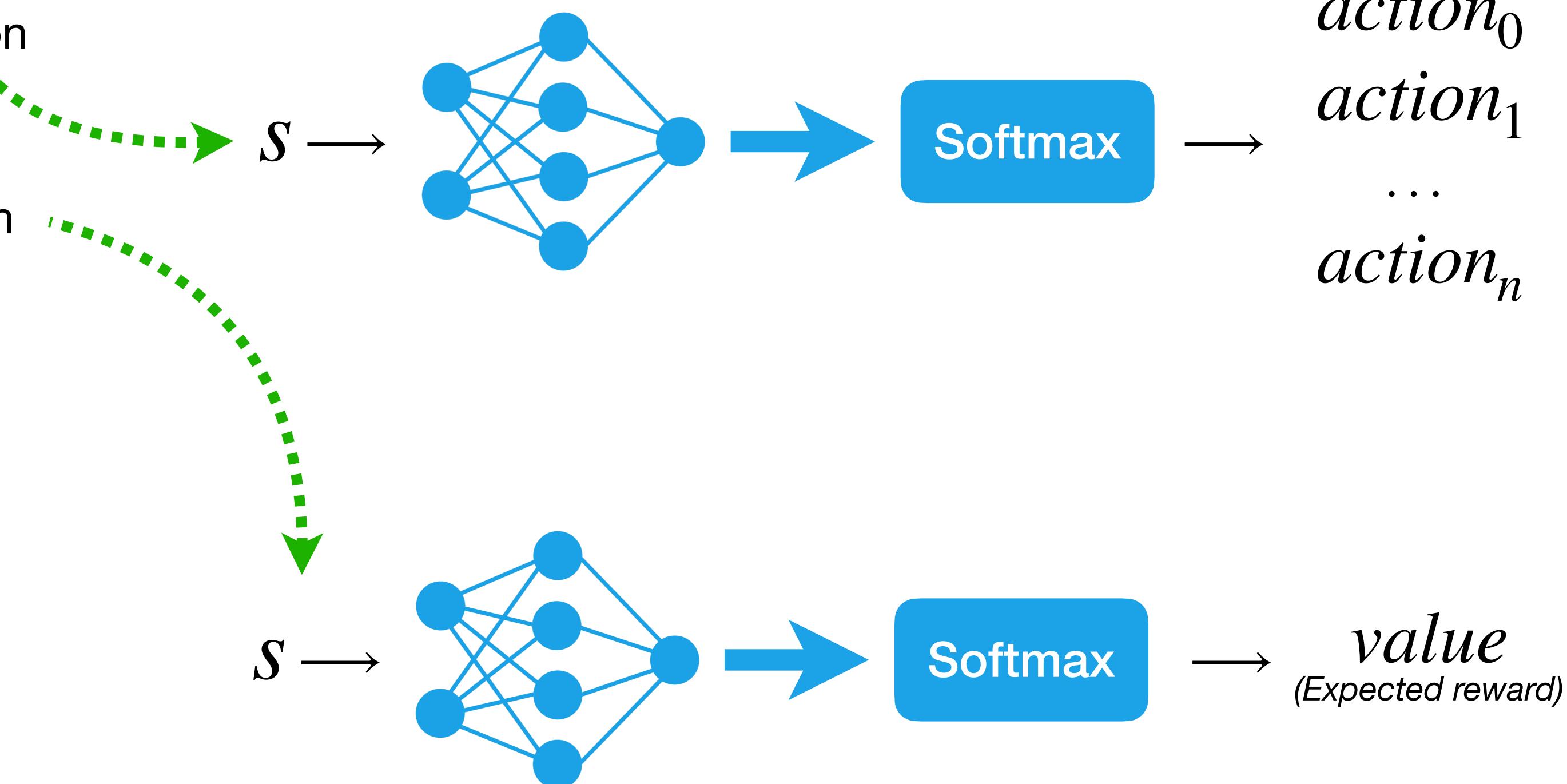
$v(s)$

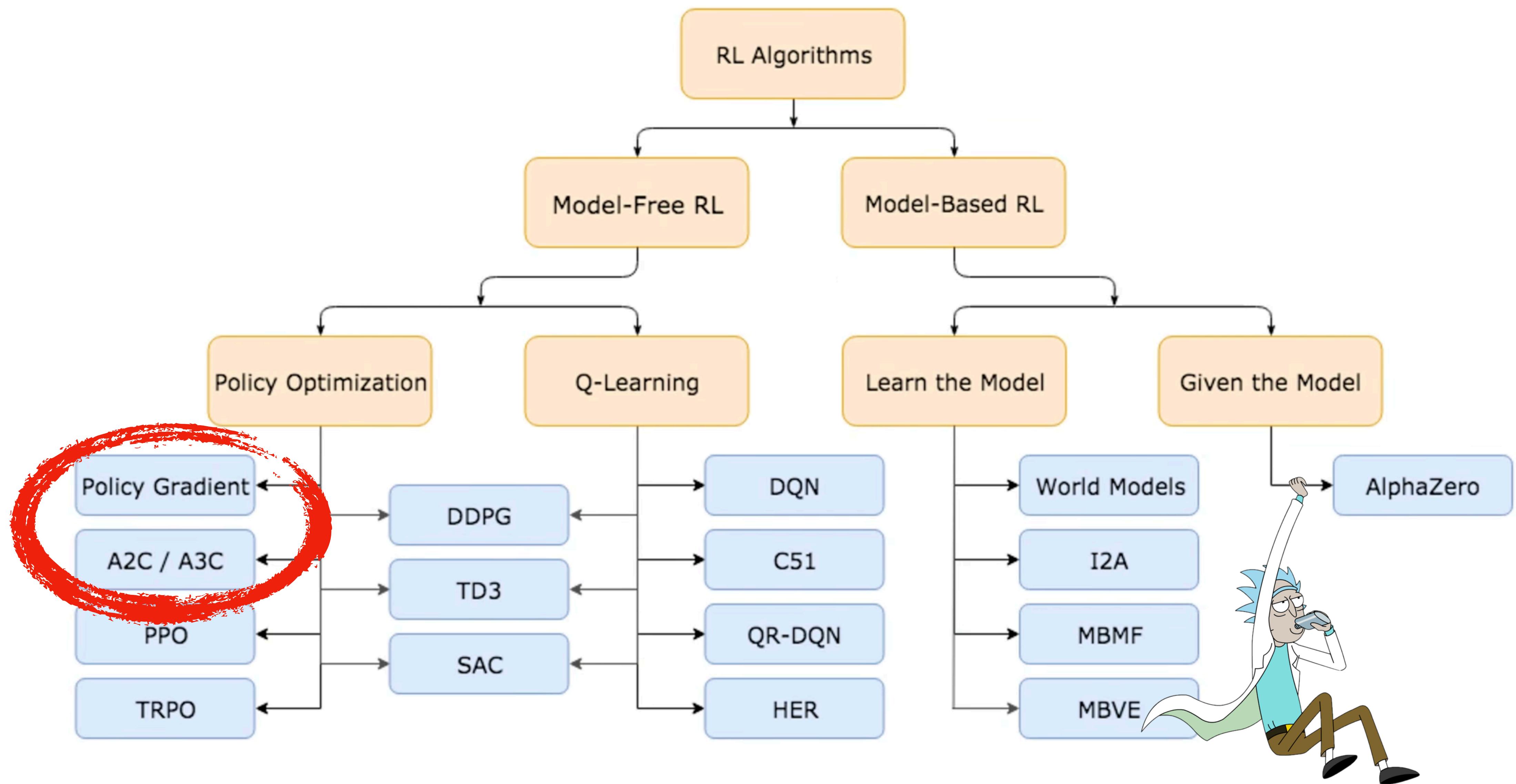
$Q(s, a)$

policy function

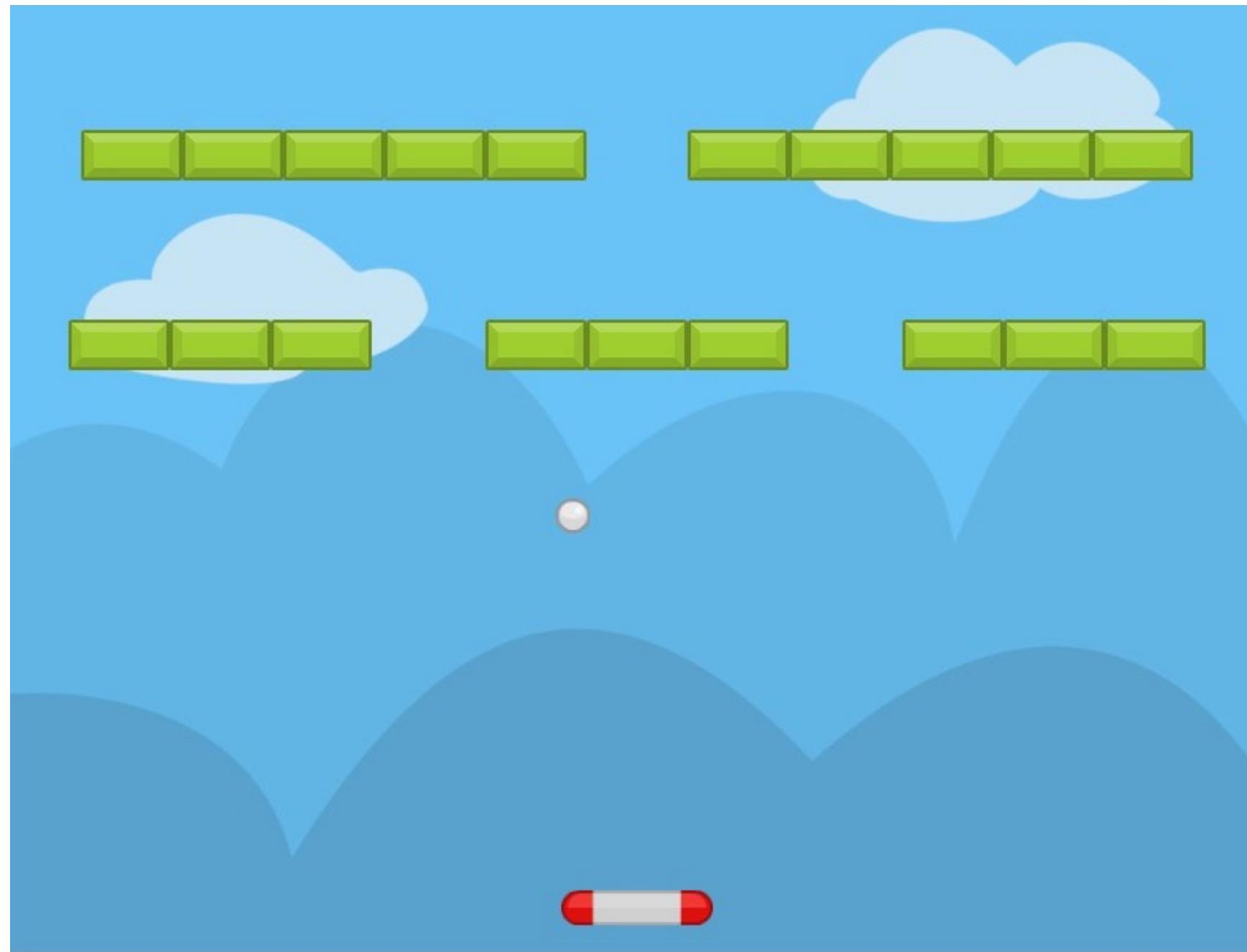
value function

Q-function

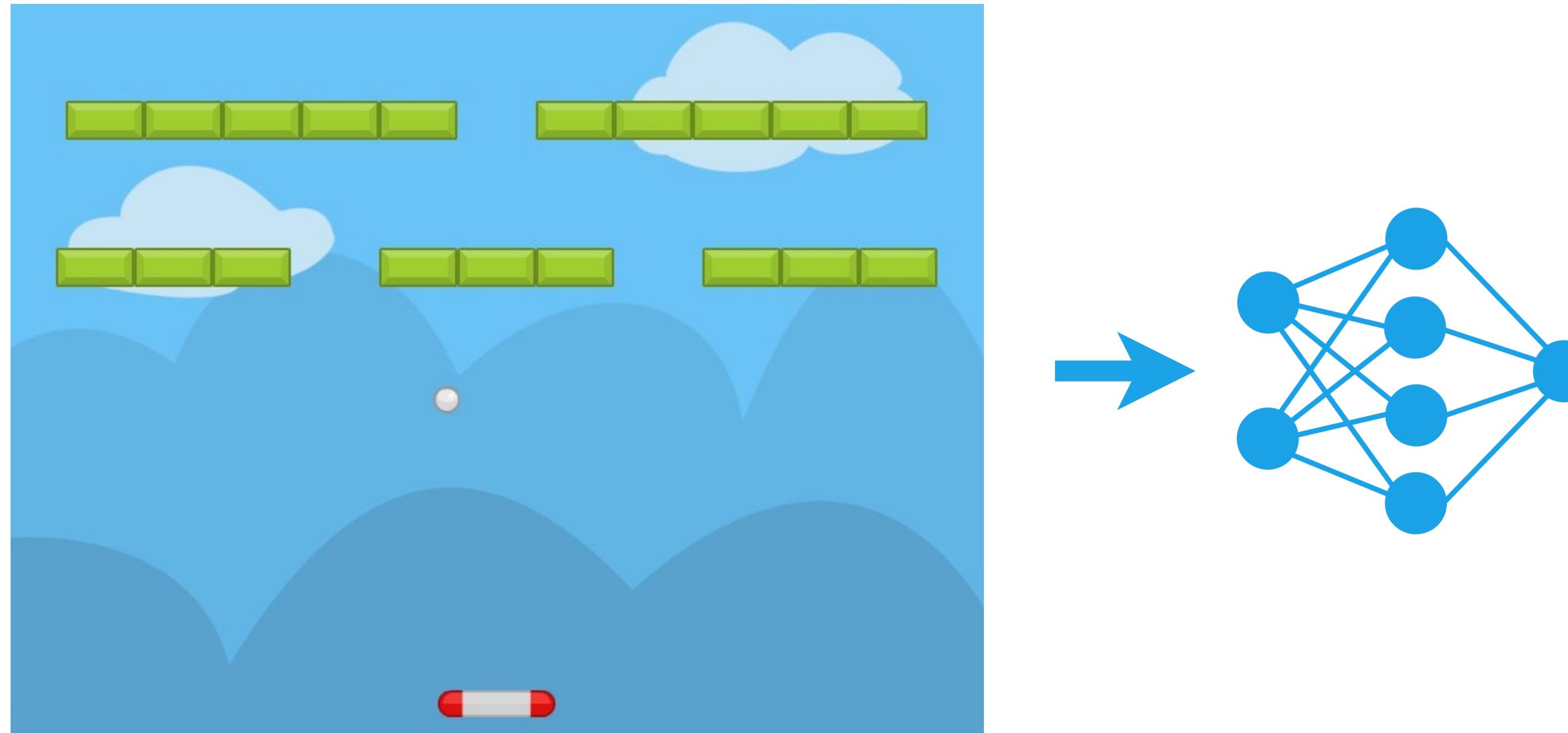




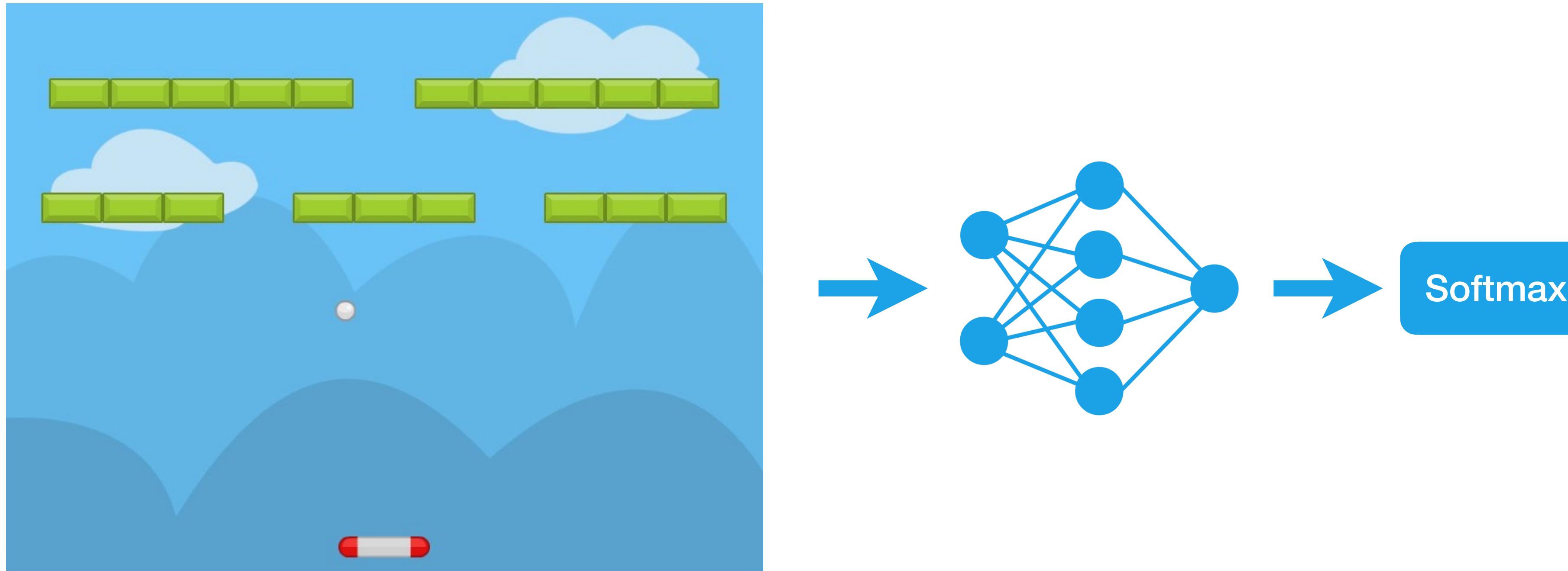
Policy Gradients



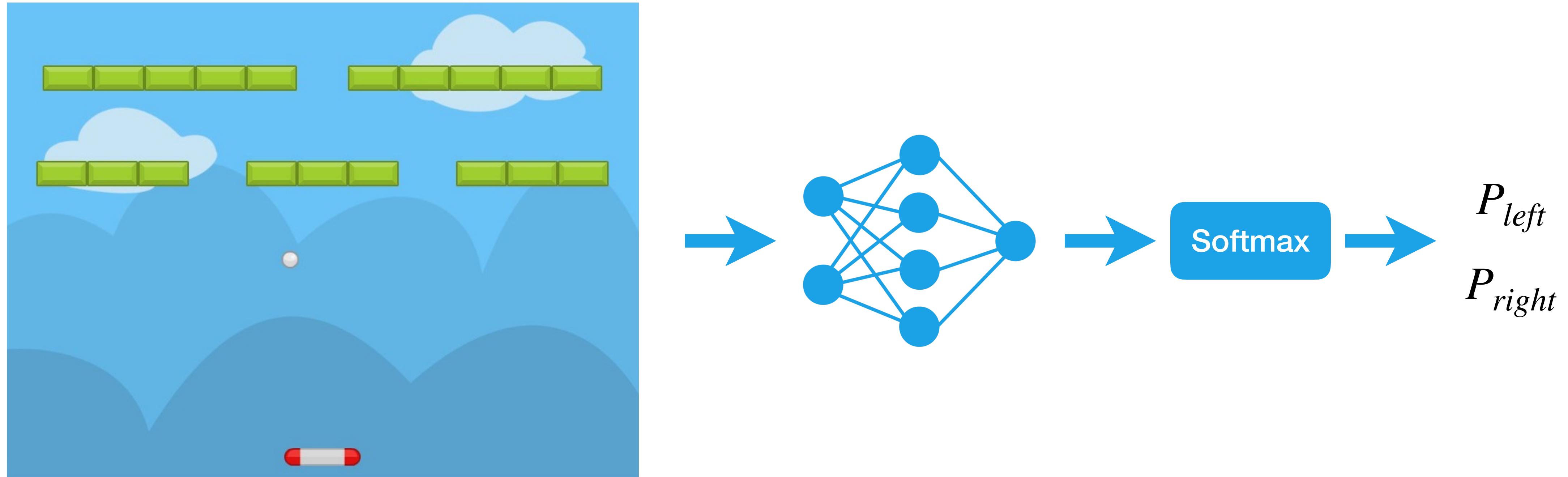
Policy Gradients



Policy Gradients



Policy Gradients



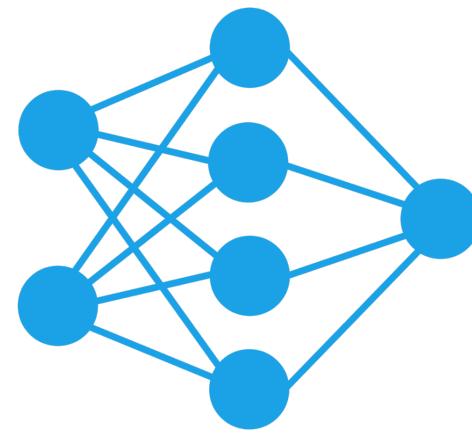
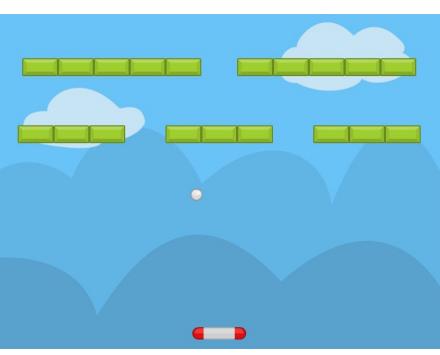
Отметим:

- Несколько эпизодов игры
- Каждый эпизод содержит разное количество ходов
- Мы можем оценить реворд (выиграл или проиграл) только после конца эпизода

Policy Gradients

*Supervised
learning*

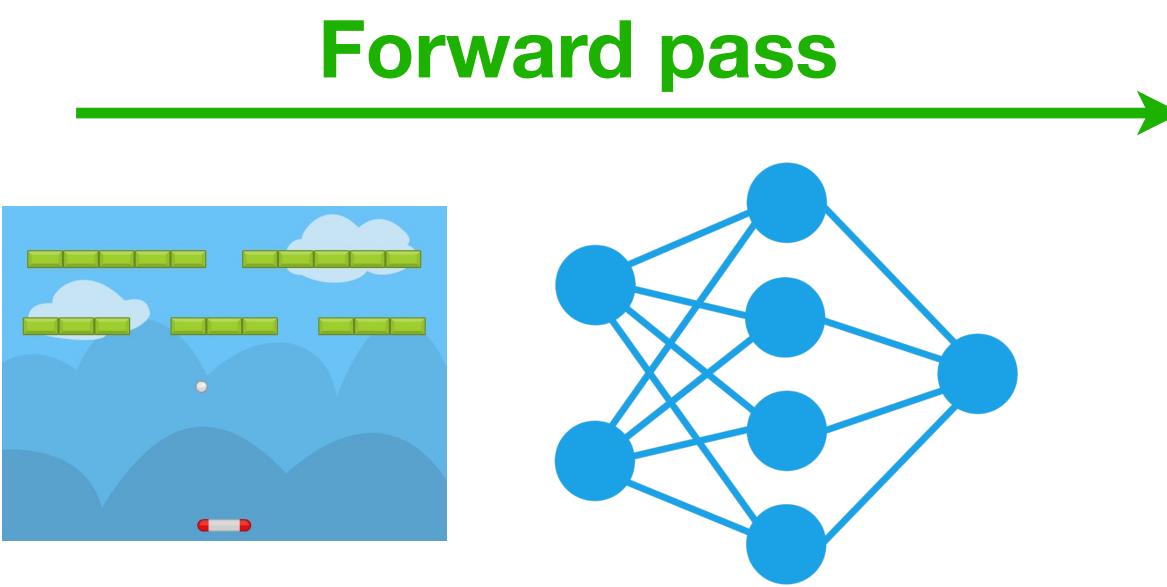
$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Policy Gradients

*Supervised
learning*

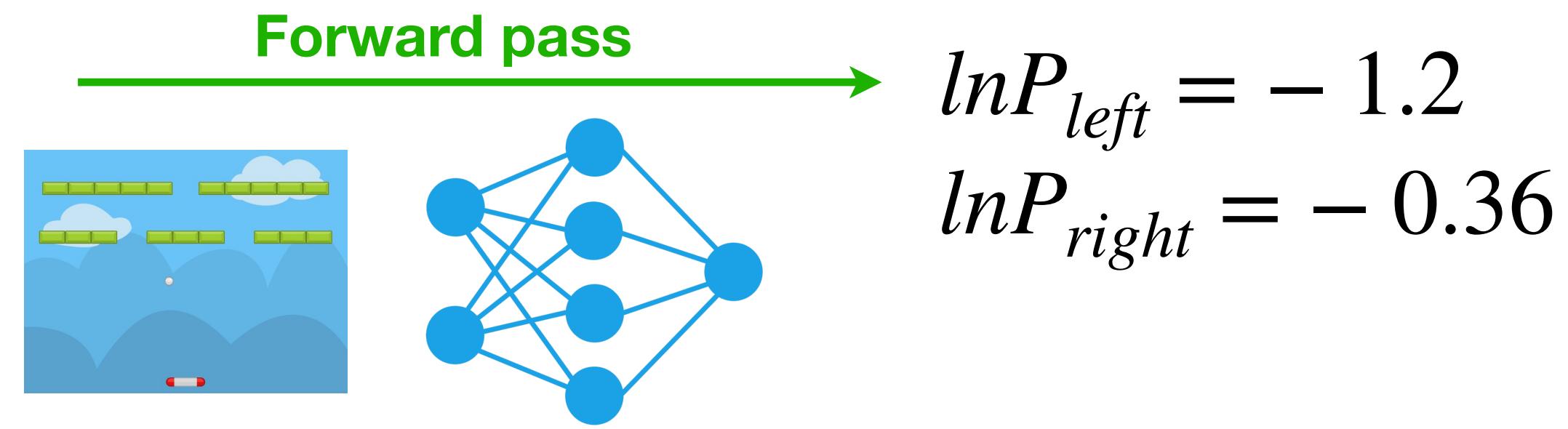
$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Policy Gradients

*Supervised
learning*

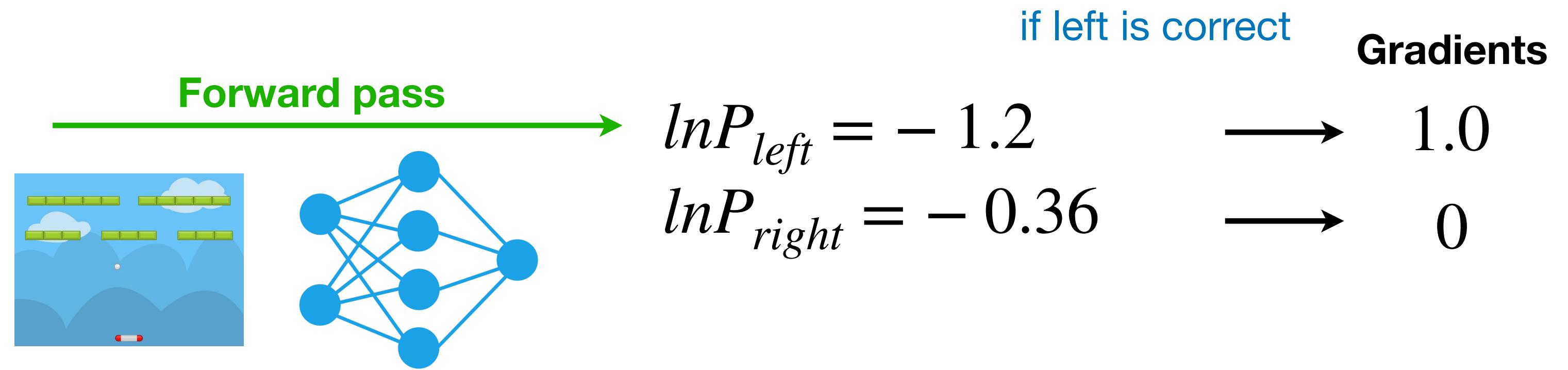
$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Policy Gradients

**Supervised
learning**

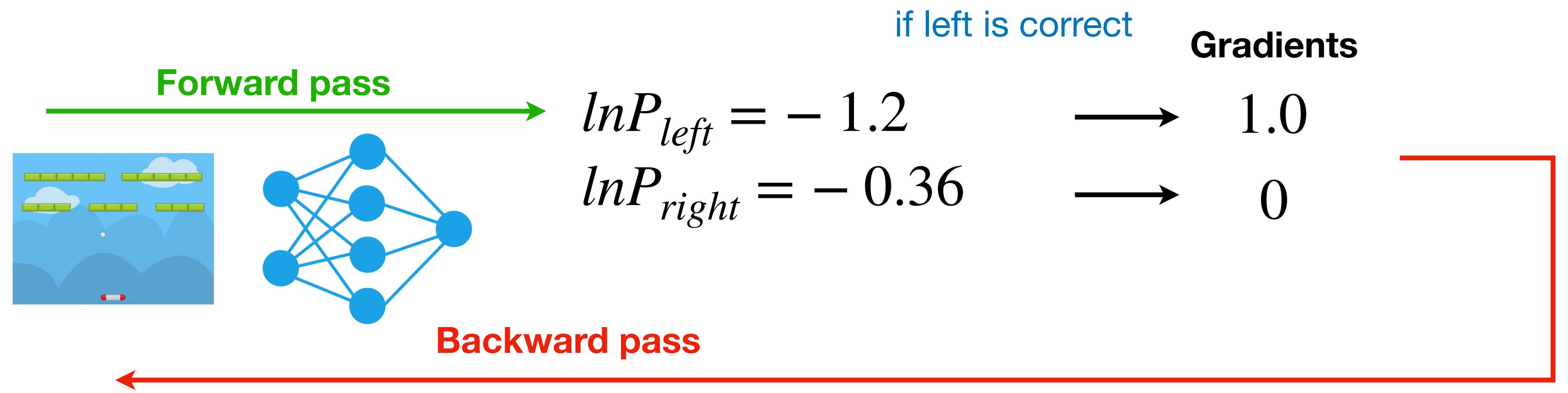
$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Policy Gradients

Supervised learning

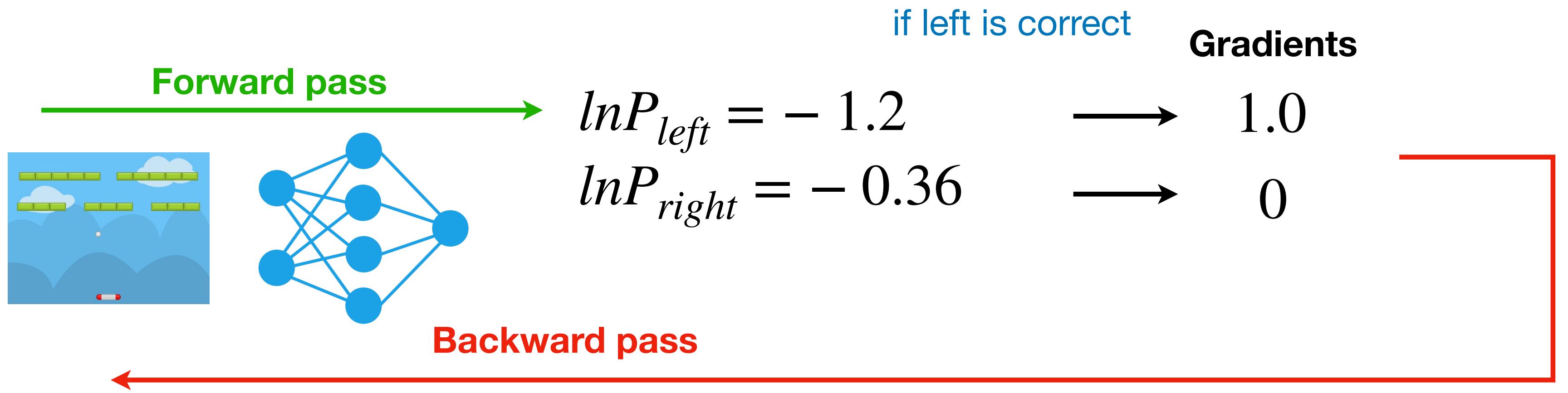
$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Проблема:

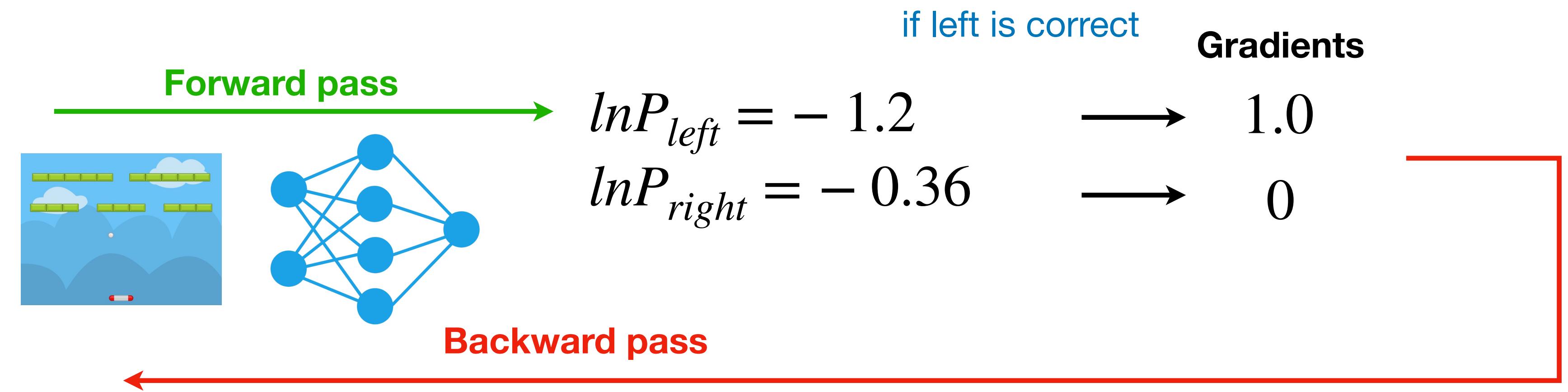
Мы получаем оценку, что действие **правильное**,
только в конце эпизода, а не после каждого шага



Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

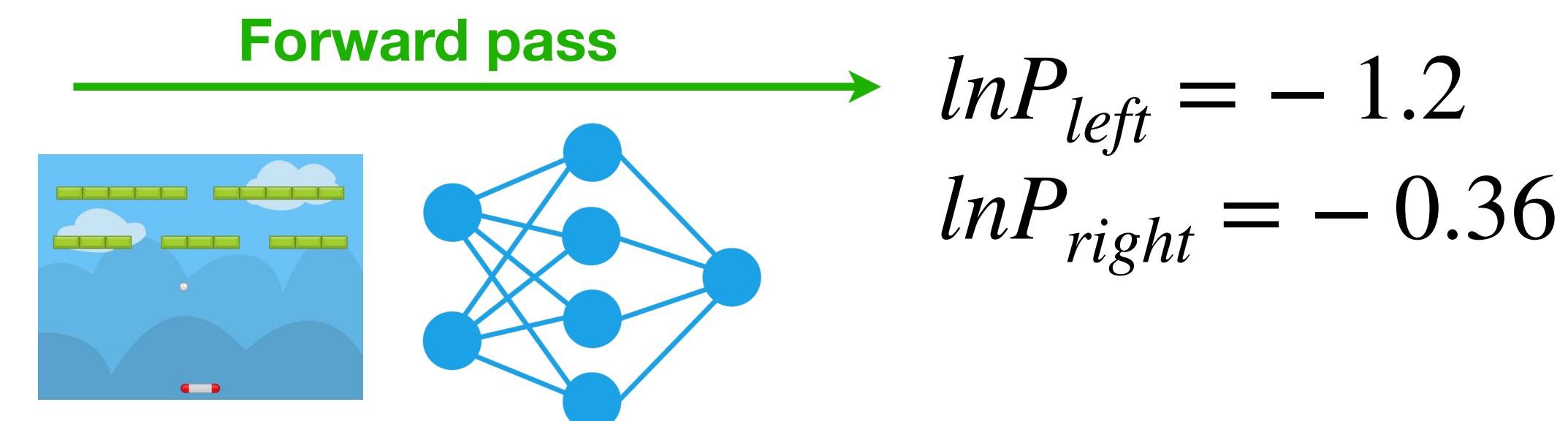
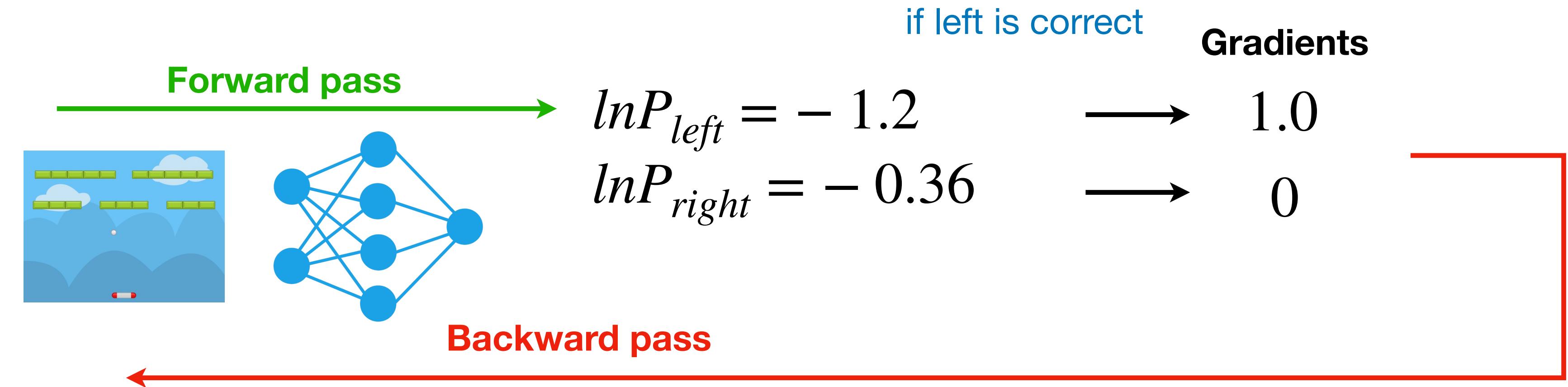


Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

Reinforcement learning

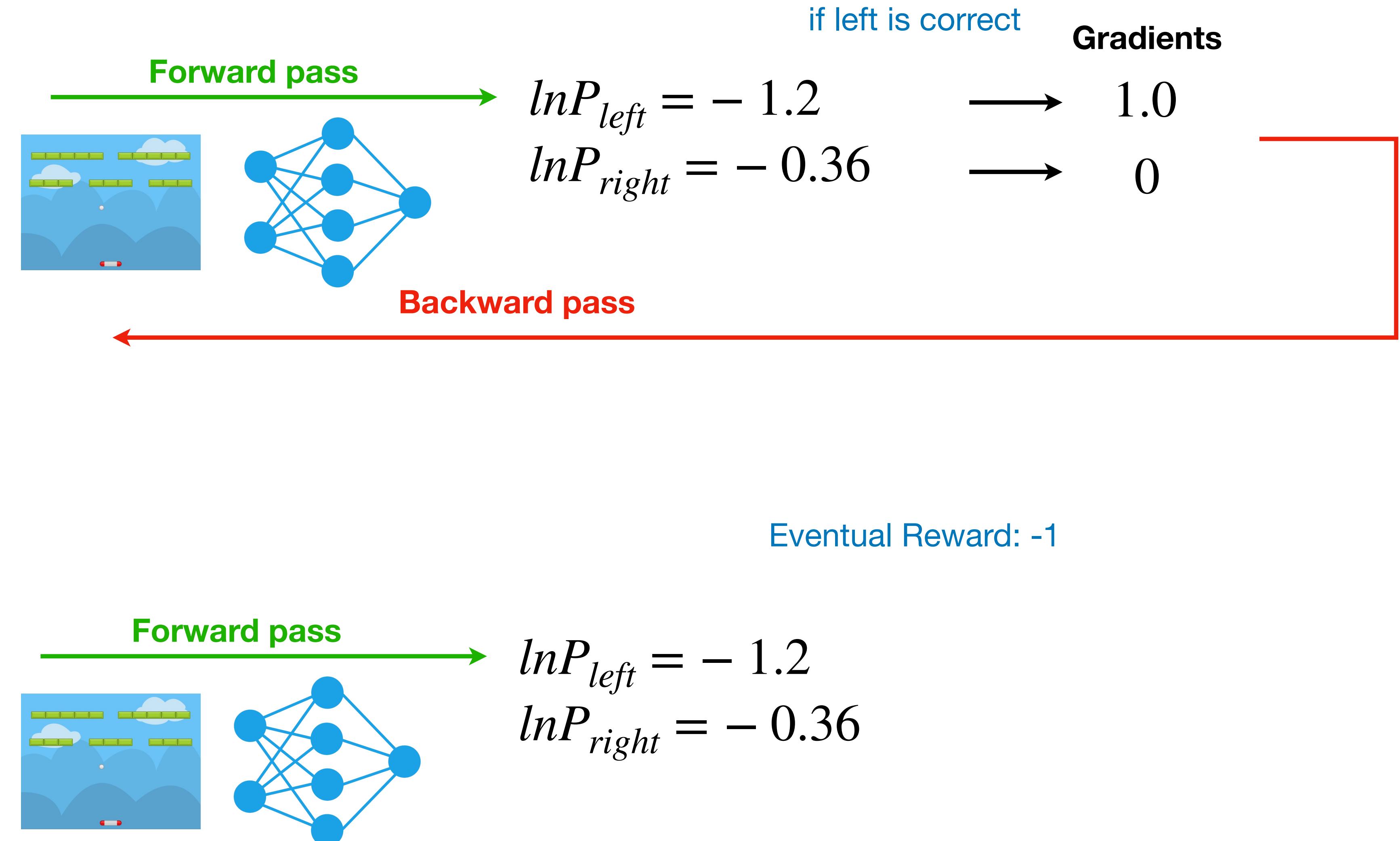


Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

Reinforcement learning



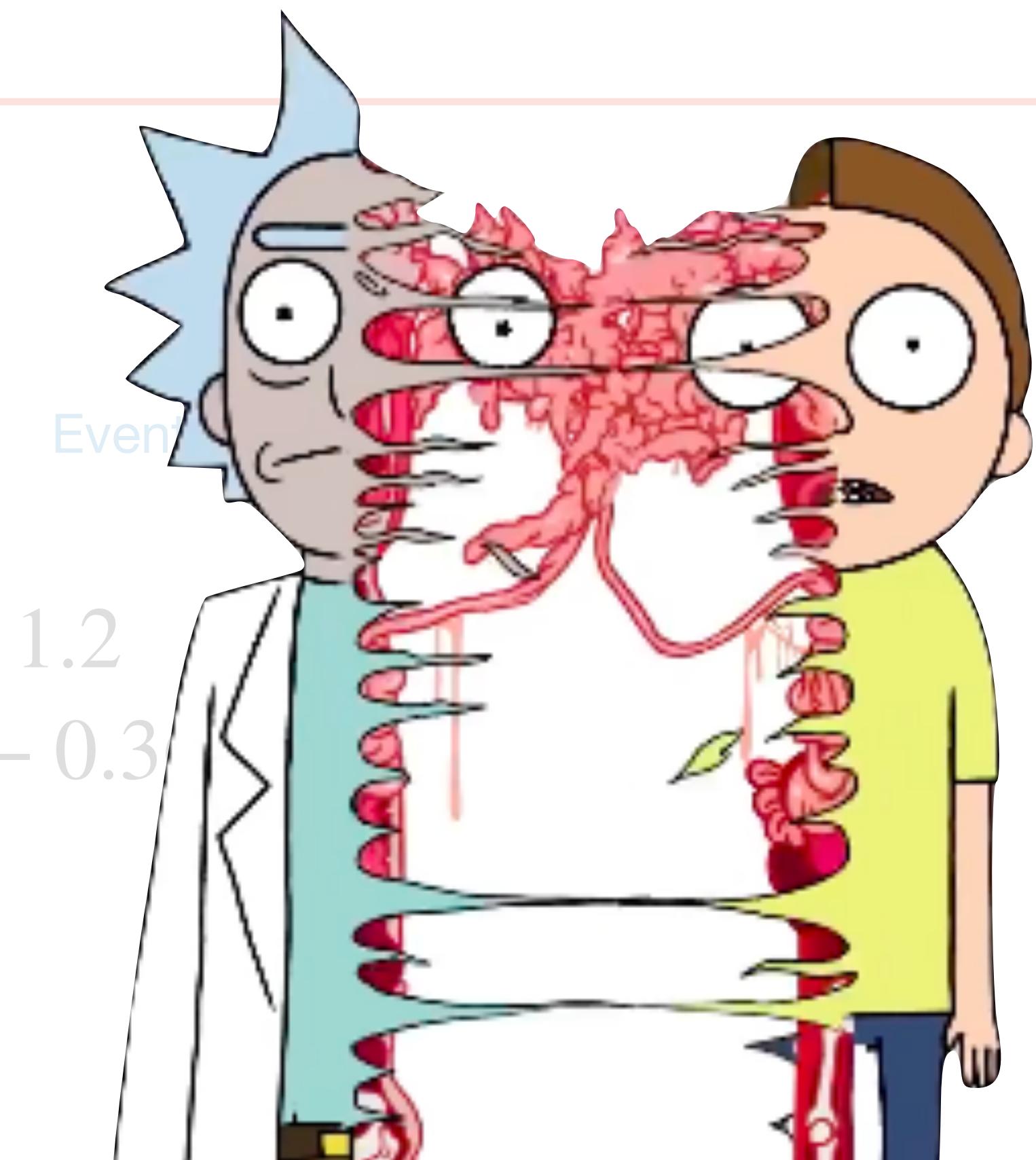
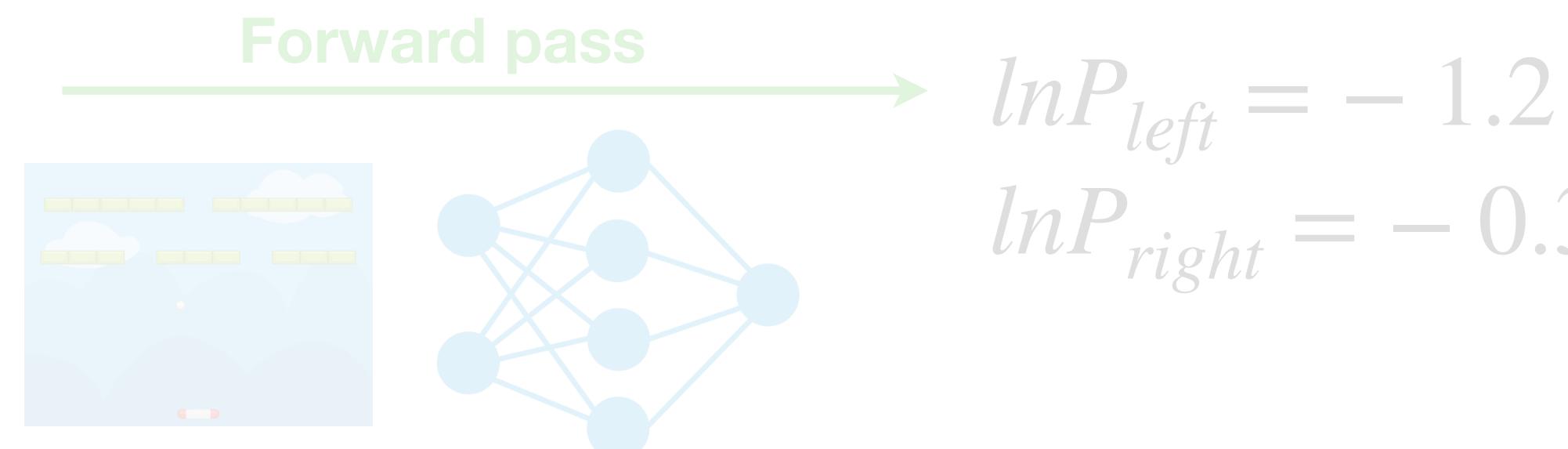
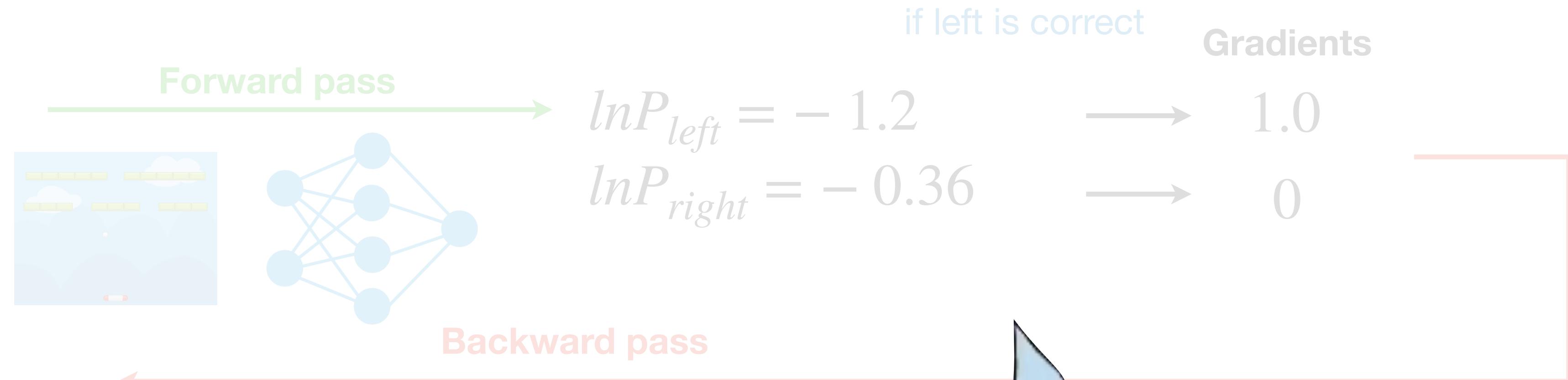
Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$



Reinforcement learning



Policy Gradients

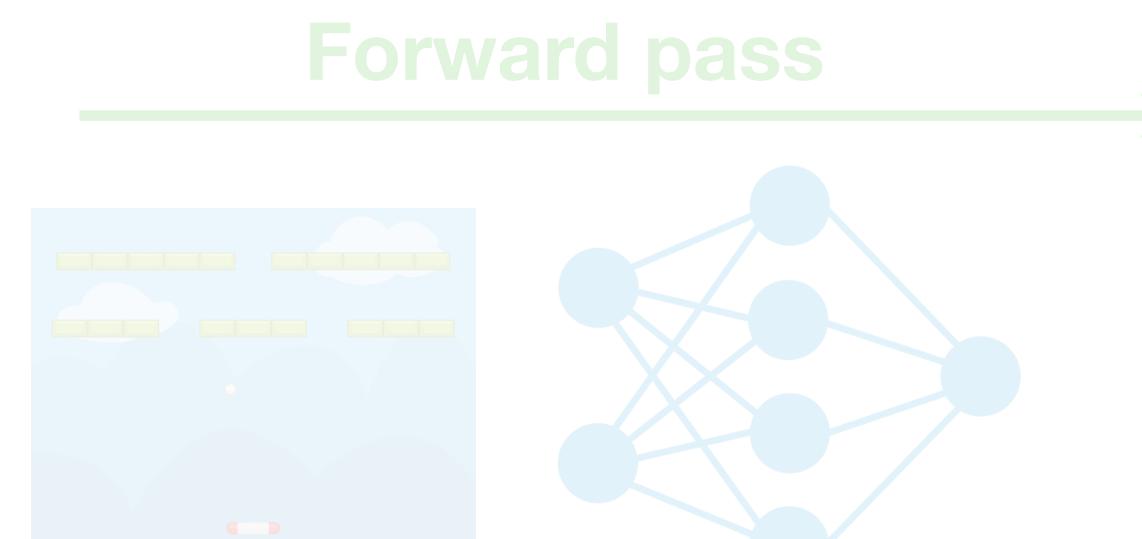
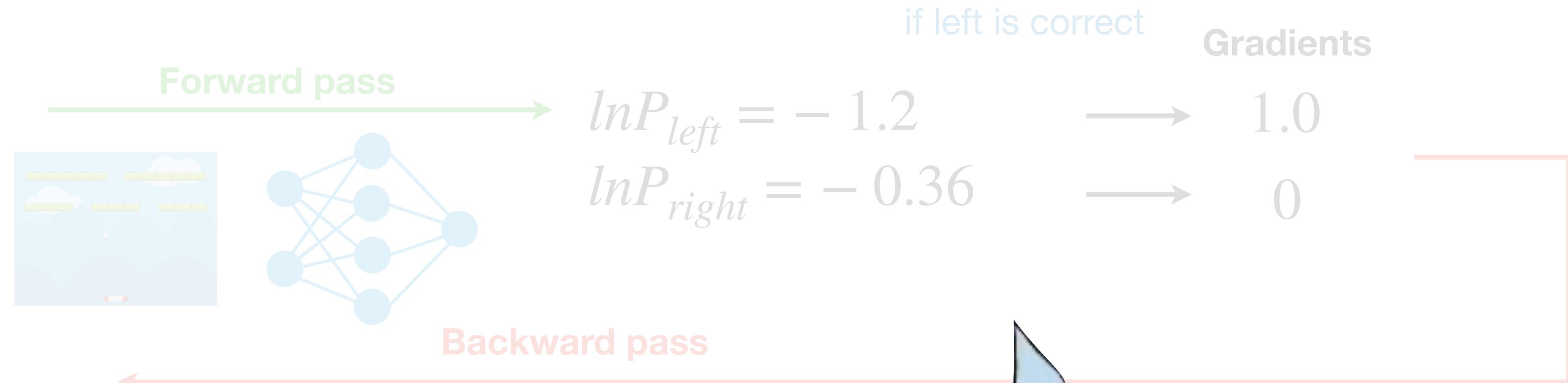
Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

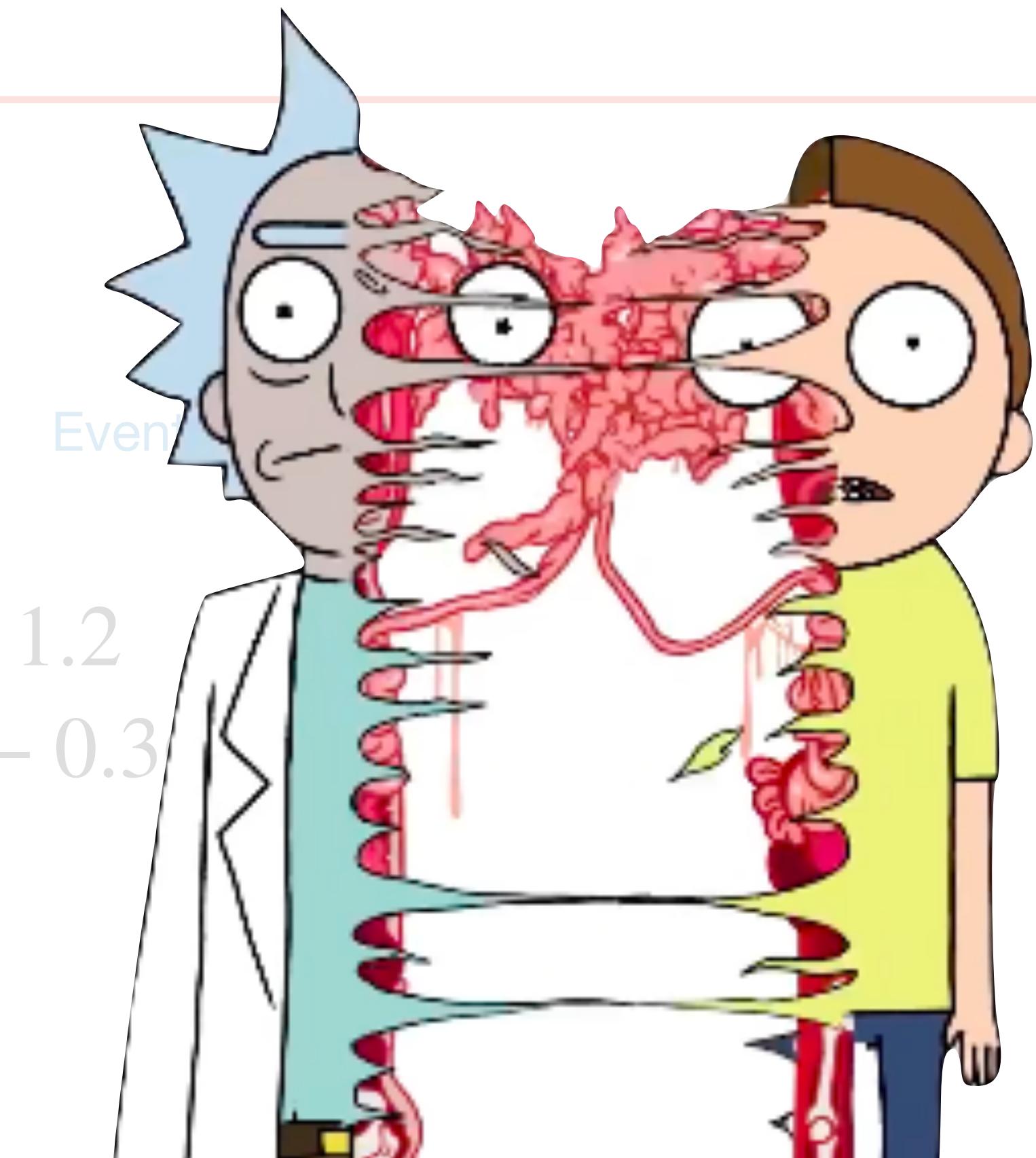


Reinforcement learning

$$\mathcal{L}_{PG} = - \sum_i r_i \ln p(c = a_i | x_i)$$



$$\begin{aligned}\ln P_{left} &= -1.2 \\ \ln P_{right} &= -0.3\end{aligned}$$



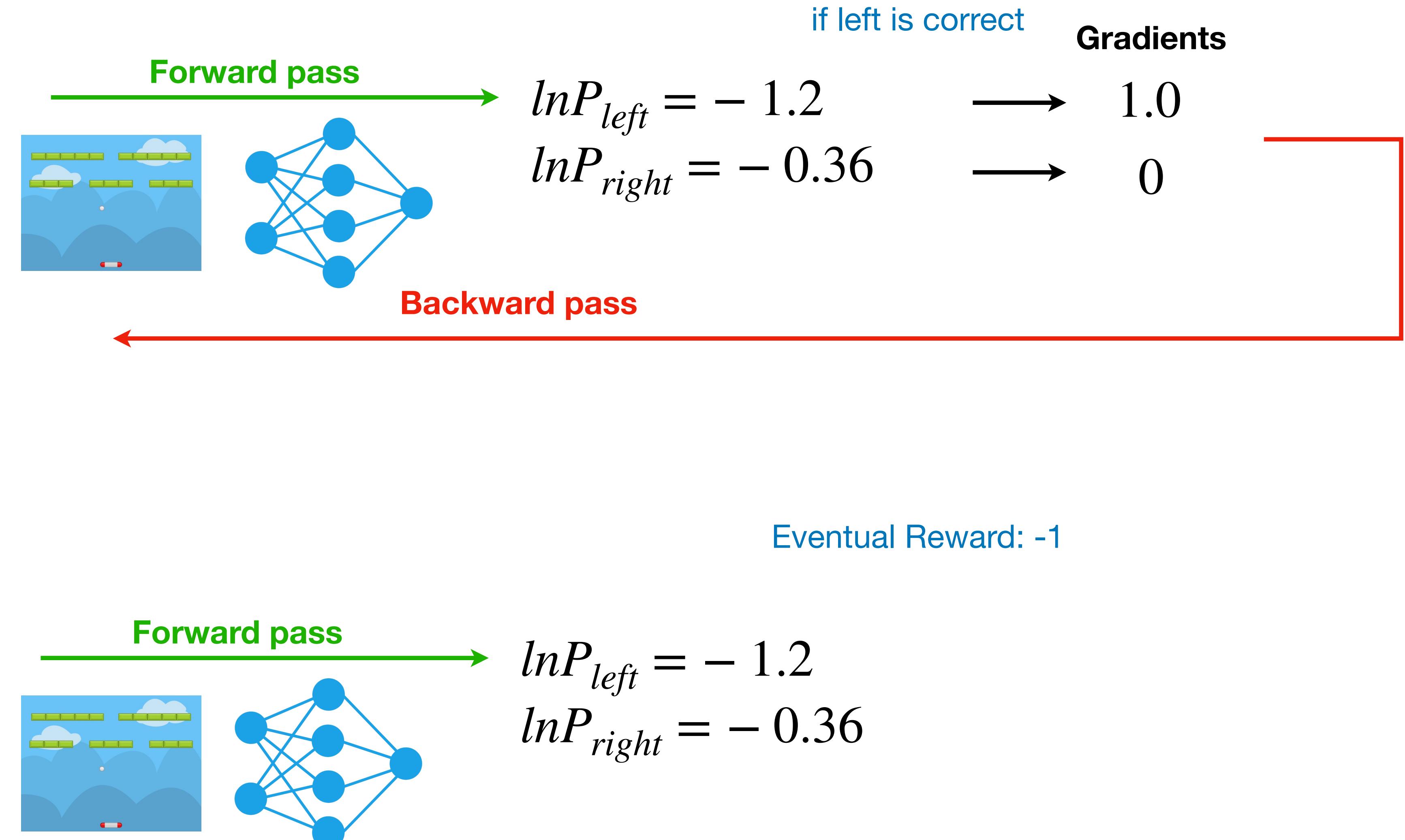
Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

Reinforcement learning

$$\mathcal{L}_{PG} = - \sum_i r_i \ln p(c = a_i | x_i)$$



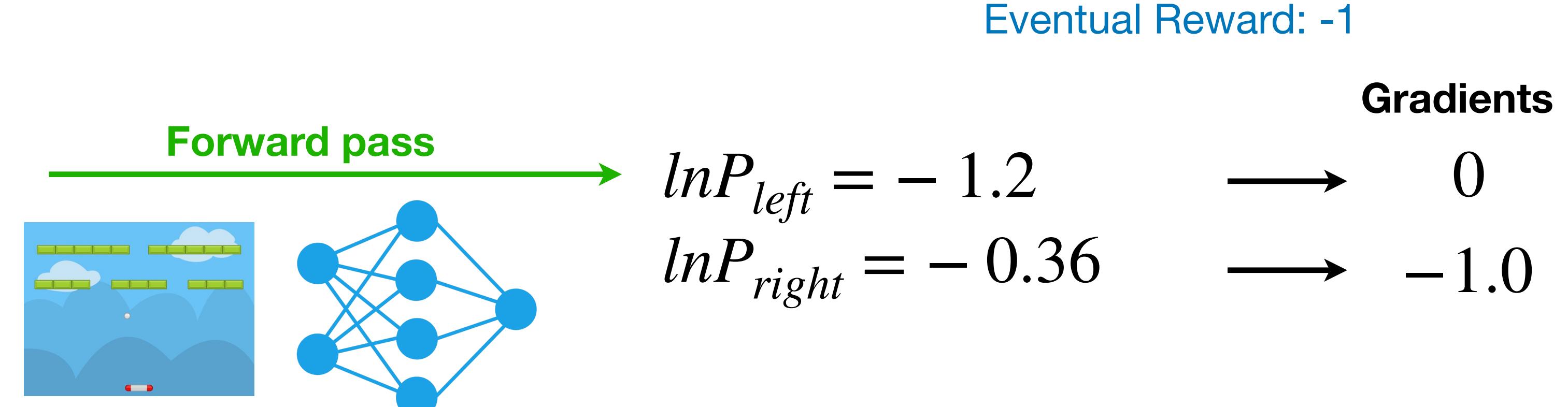
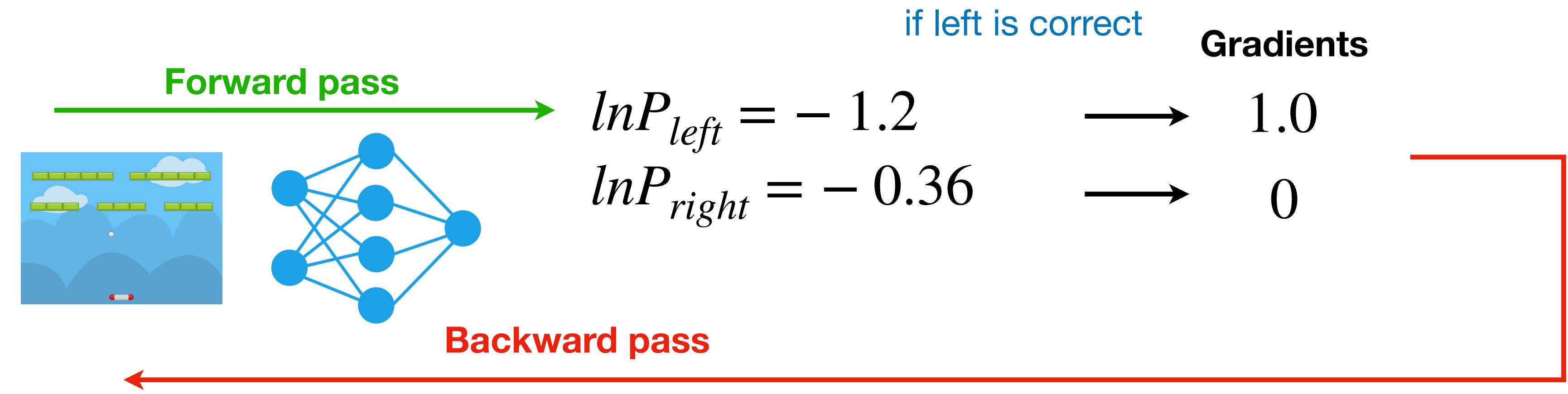
Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

Reinforcement learning

$$\mathcal{L}_{PG} = - \sum_i r_i \ln p(c = a_i | x_i)$$



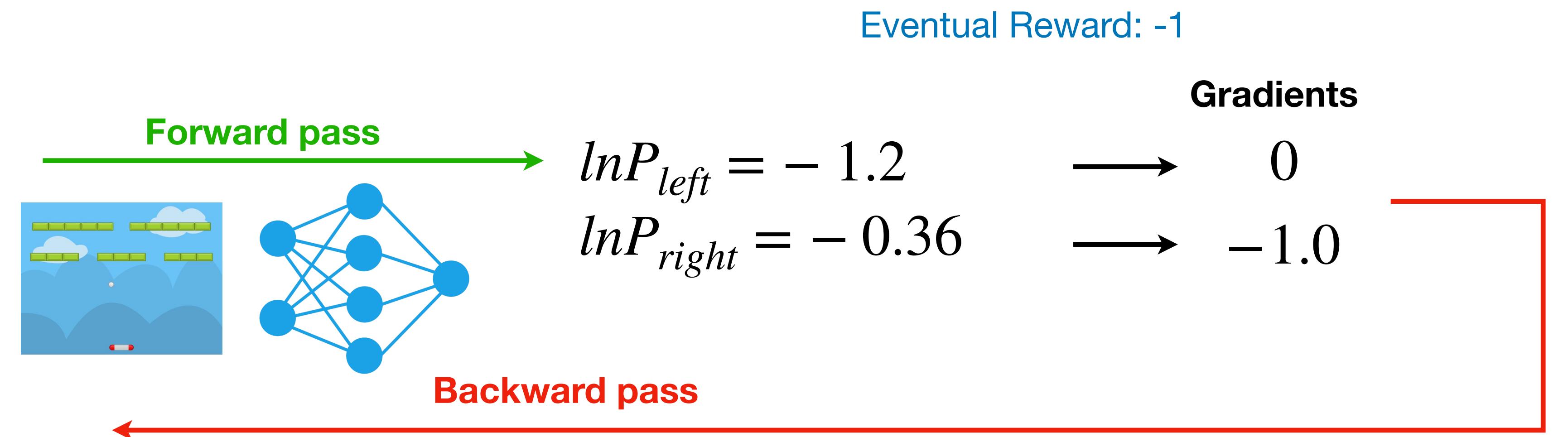
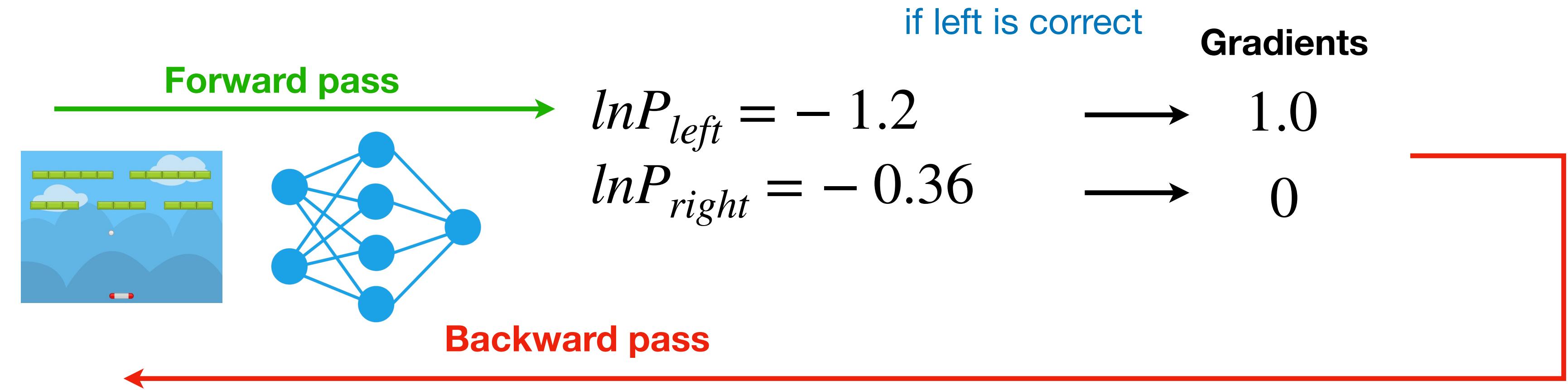
Policy Gradients

Supervised learning

$$\mathcal{L}_{CE} = - \sum_i y_i \ln p(c = t_i | x_i)$$

Reinforcement learning

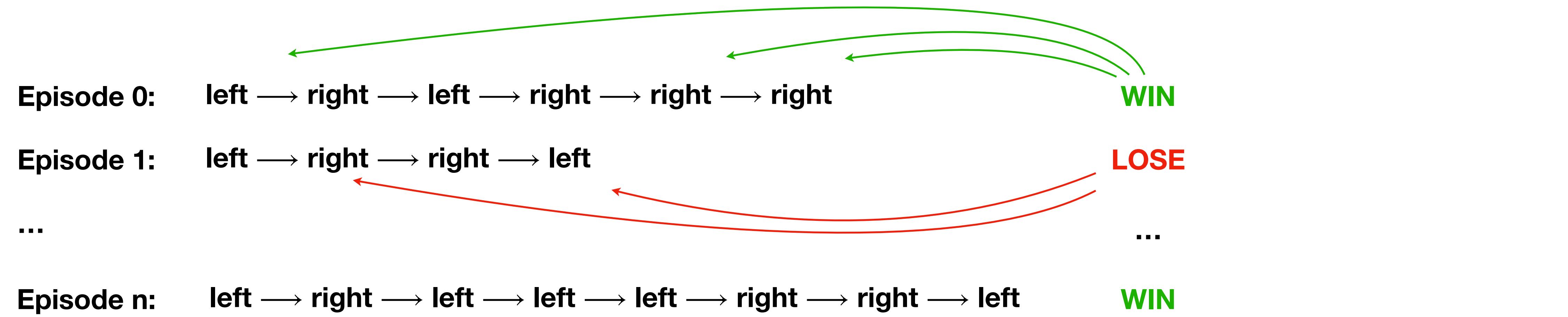
$$\mathcal{L}_{PG} = - \sum_i r_i \ln p(c = a_i | x_i)$$



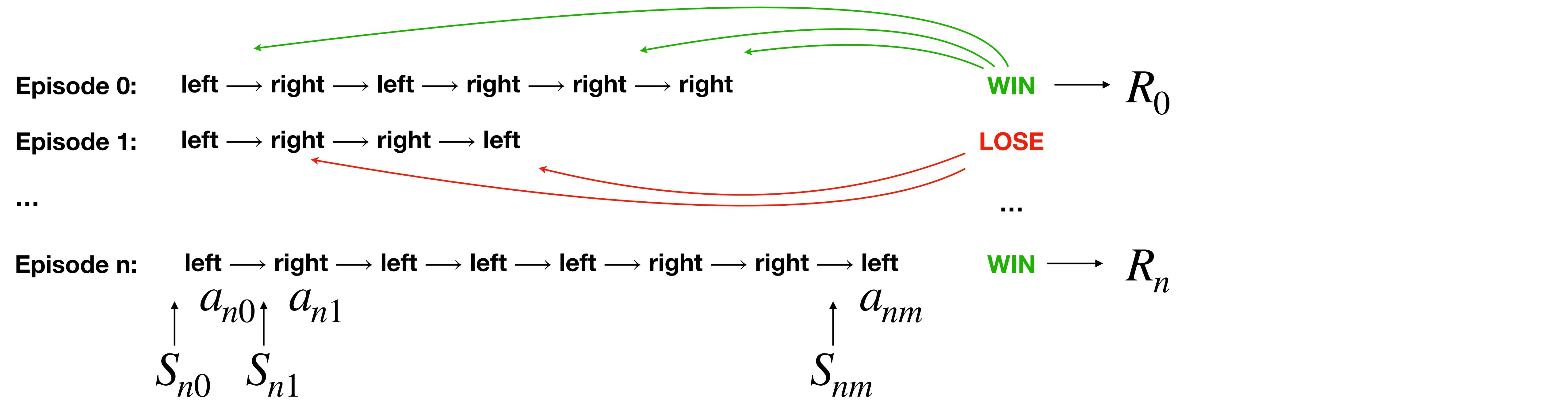
Policy Gradients

Episode 0:	left → right → left → right → right → right	WIN
Episode 1:	left → right → right → left	LOSE
...		...
Episode n:	left → right → left → left → left → right → right → left	WIN

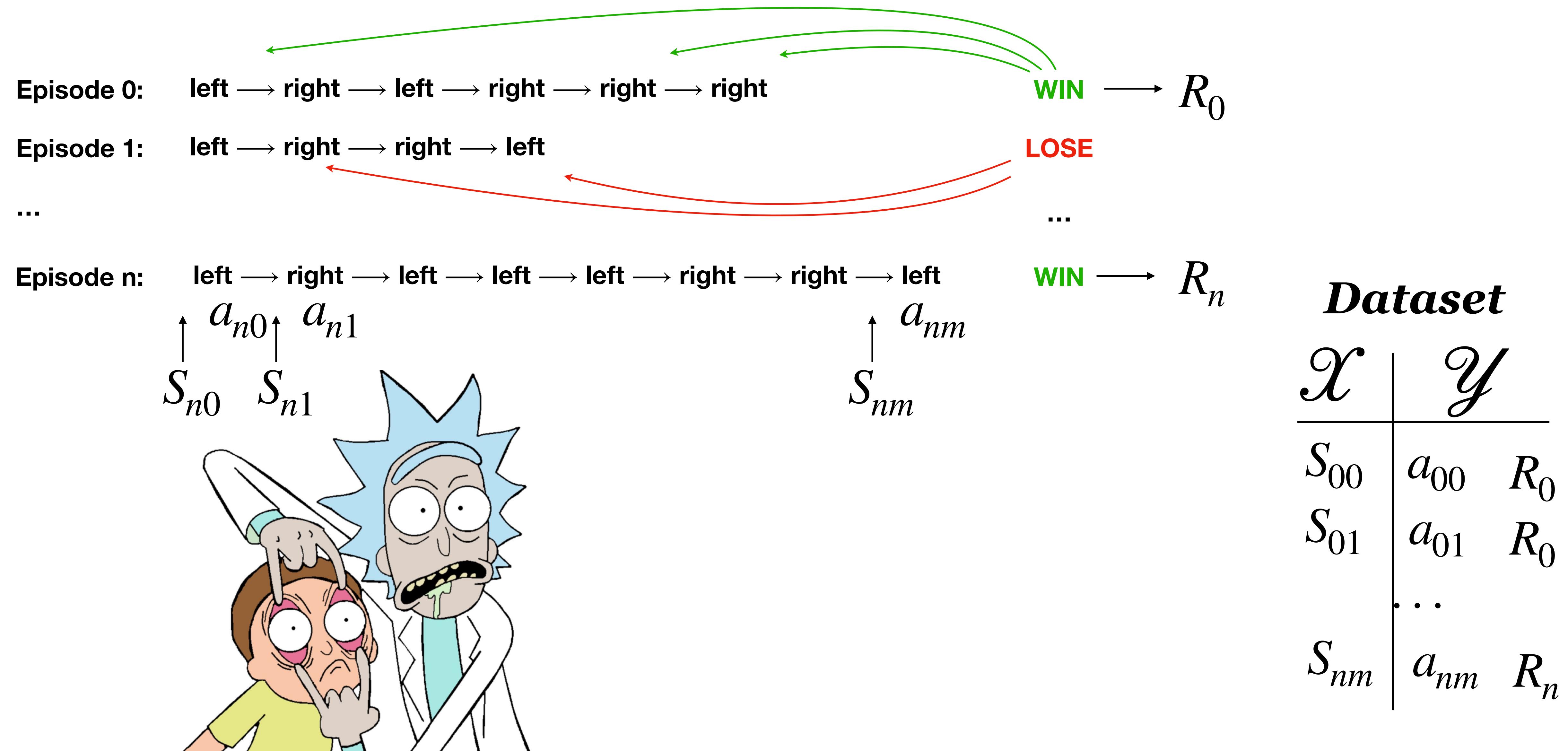
Policy Gradients



Policy Gradients



Policy Gradients



Policy Gradients

Policy Gradients

1. Инициализировали сеть со случайными весами

Policy Gradients

1. Инициализировали сеть со случайными весами
2. Прогнали n эпизодов

Policy Gradients

1. Инициализировали сеть со случайными весами
2. Прогнали n эпизодов
3. Обучили на 1 эпохе

Policy Gradients

1. Инициализировали сеть со случайными весами
2. Прогнали n эпизодов
3. Обучили на 1 эпохе
4. Прогнали n эпизодов

Policy Gradients

1. Инициализировали сеть со случайными весами
2. Прогнали n эпизодов
3. Обучили на 1 эпохе
4. Прогнали n эпизодов
5. Перешли к п. 3

Policy Gradients

1. Инициализировали сеть со случайными весами
2. Прогнали n эпизодов
3. Обучили на 1 эпохе
4. Прогнали n эпизодов*
5. Перешли к п. 3

*и выкинули данные прошлого прогона



Policy Gradients

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

$f(a)$ - reward from environment

Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \longrightarrow \max$

Policy Gradients

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

$f(a)$ - reward from environment

$$E[f(a)] \longrightarrow \max$$



Policy Gradients

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

$f(a)$ - reward from environment

$$E[f(a)] = \sum_a p(a)f(a)$$

$$E[f(a)] \longrightarrow \max$$



Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$

$$E[f(a)] = \sum_a p(a)f(a)$$

$$\nabla_{\theta} E[f(a)] = \nabla_{\theta} \sum_a p(a)f(a)$$



Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$

$$E[f(a)] = \sum_a p(a)f(a)$$

$$\begin{aligned}\nabla_{\theta} E[f(a)] &= \nabla_{\theta} \sum_a p(a)f(a) \\ &= \sum_a \nabla_{\theta} p(a)f(a)\end{aligned}$$



Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$



$$E[f(a)] = \sum_a p(a)f(a)$$

$$\nabla_{\theta} E[f(a)] = \nabla_{\theta} \sum_a p(a)f(a)$$

$$= \sum_a \nabla_{\theta} p(a)f(a)$$

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

Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$



$$E[f(a)] = \sum_a p(a)f(a)$$

$$\nabla_{\theta} E[f(a)] = \nabla_{\theta} \sum_a p(a)f(a)$$

$$= \sum_a \nabla_{\theta} p(a)f(a)$$

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

$$\frac{1}{z} \nabla_{\theta} z = \nabla_{\theta} \log(z)$$

Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$



$$E[f(a)] = \sum_a p(a)f(a)$$

$$\nabla_{\theta} E[f(a)] = \nabla_{\theta} \sum_a p(a)f(a)$$

$$= \sum_a \nabla_{\theta} p(a)f(a)$$

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

$$= \sum_a p(a) [\nabla_{\theta} \log(p(a))] f(a)$$

$$\frac{1}{z} \nabla_{\theta} z = \nabla_{\theta} \log(z)$$

Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$



$$E[f(a)] = \sum_a p(a)f(a)$$

$$\nabla_{\theta} E[f(a)] = \nabla_{\theta} \sum_a p(a)f(a)$$

$$= \sum_a \nabla_{\theta} p(a)f(a)$$

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

$$= \sum_a p(a) [\nabla_{\theta} \log(p(a))] f(a)$$

$$= E [f(a) \nabla_{\theta} \log(p(a))]$$

$$\frac{1}{z} \nabla_{\theta} z = \nabla_{\theta} \log(z)$$

Policy Gradients

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

$f(a)$ - reward from environment

$E[f(a)] \rightarrow \max$



$$E[f(a)] = \sum_a p(a)f(a)$$

$$\nabla_{\theta} E[f(a)] = \nabla_{\theta} \sum_a p(a)f(a)$$

$$= \sum_a \nabla_{\theta} p(a)f(a)$$

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

$$= \sum_a p(a) [\nabla_{\theta} \log(p(a))] f(a)$$

$$= E [f(a) \nabla_{\theta} \log(p(a))]$$

$$\frac{1}{z} \nabla_{\theta} z = \nabla_{\theta} \log(z)$$

$$\mathcal{L}_{PG} = - \sum_i r_i \ln p(c = a_i | x_i)$$

Tips and Tricks



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
\dots	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

1

Discounted rewards



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

2

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

2 Baseline

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$

2

Baseline

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

$$r_0 = -1$$

$$r_1 = 2$$

$$r_2 = -4$$

$$r_3 = 1$$



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
\dots	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

2

Baseline

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

$$r_0 = -1$$

$$r_1 = 2$$

$$r_2 = -4$$

$$r_3 = 1$$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$

3

Actor-Critic

$$r_0 = 1$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

2

Baseline

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

$$r_0 = -1$$

$$r_1 = 2$$

$$r_2 = -4$$

$$r_3 = 1$$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$

3

Actor-Critic

$$r_0 = 1$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

2

Baseline

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

$$r_0 = -1$$

$$r_1 = 2$$

$$r_2 = -4$$

$$r_3 = 1$$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$

3

Actor-Critic

$$r_0 = 1$$

$$r_1 = 10002$$

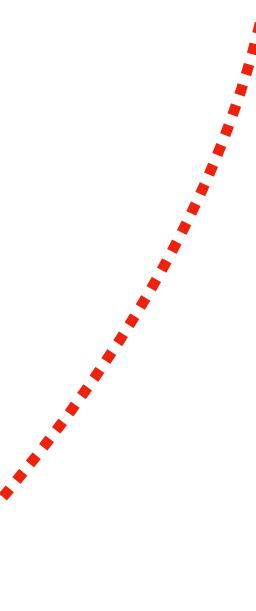
$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

Хотим, чтобы он
динамически
изменялся



Tips and Tricks

Dataset

\mathcal{X}	\mathcal{Y}
S_{00}	$a_{00} \quad R_0$
S_{01}	$a_{01} \quad R_0$
...	
S_{nm}	$a_{nm} \quad R_n$

After action 0: $R_0 = r_0 + r_1 + \dots + r_n$

After action 1: $R_0 = r_1 + r_2 + \dots + r_n$

After action 2: $R_0 = r_2 + r_3 + \dots + r_n$

2

Baseline

$$r_0 = 9999$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right]$$

$$r_0 = -1$$

$$r_1 = 2$$

$$r_2 = -4$$

$$r_3 = 1$$

1

Discounted rewards

$$R_s = \sum_{i=0..n-s} \gamma^i r_{s+i} \quad 0 < \gamma < 1$$

3

Actor-Critic

$$r_0 = 1$$

$$r_1 = 10002$$

$$r_2 = 9996$$

$$r_3 = 10001$$

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

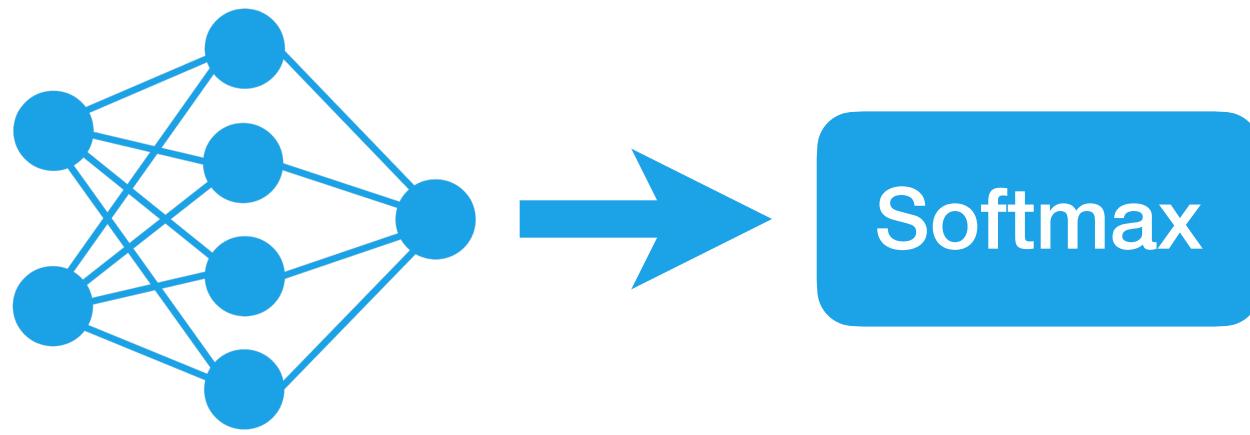
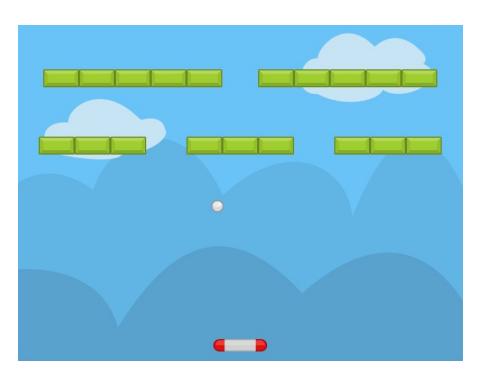
$$b = E \left[\cancel{\sum_{i=0..n-s}} \gamma^i r_{s+i} \right] \rightarrow b = v(s)$$

Actor-Critic

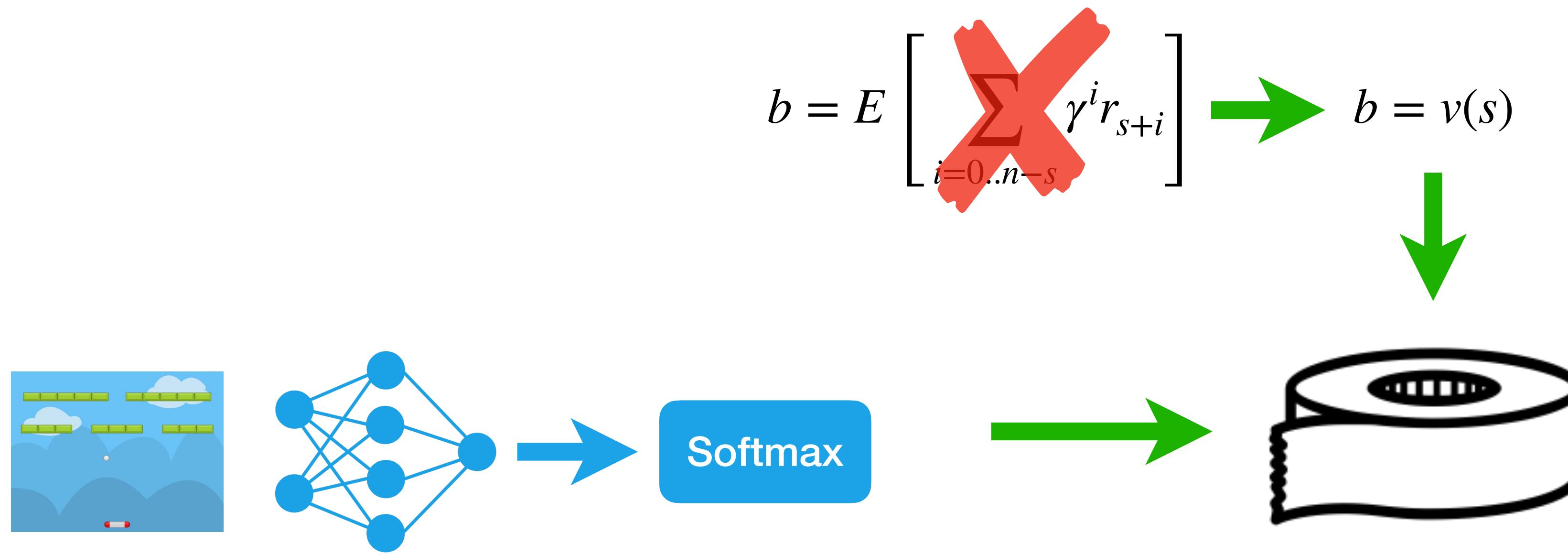
$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right] \xrightarrow{\text{red X}} b = v(s)$$

Actor-Critic

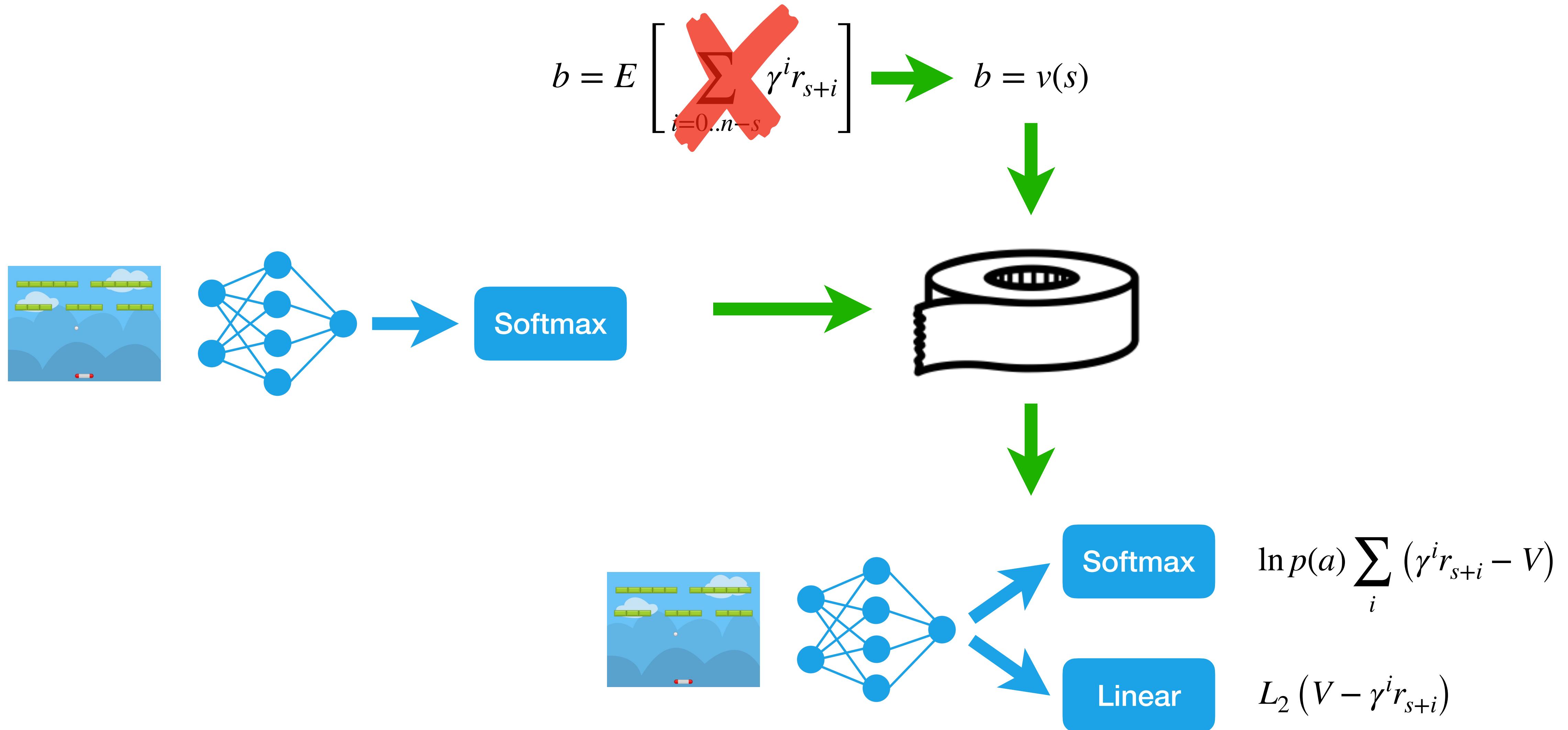
$$b = E \left[\sum_{i=0..n-s} \gamma^i r_{s+i} \right] \xrightarrow{\text{Red X}} b = v(s)$$

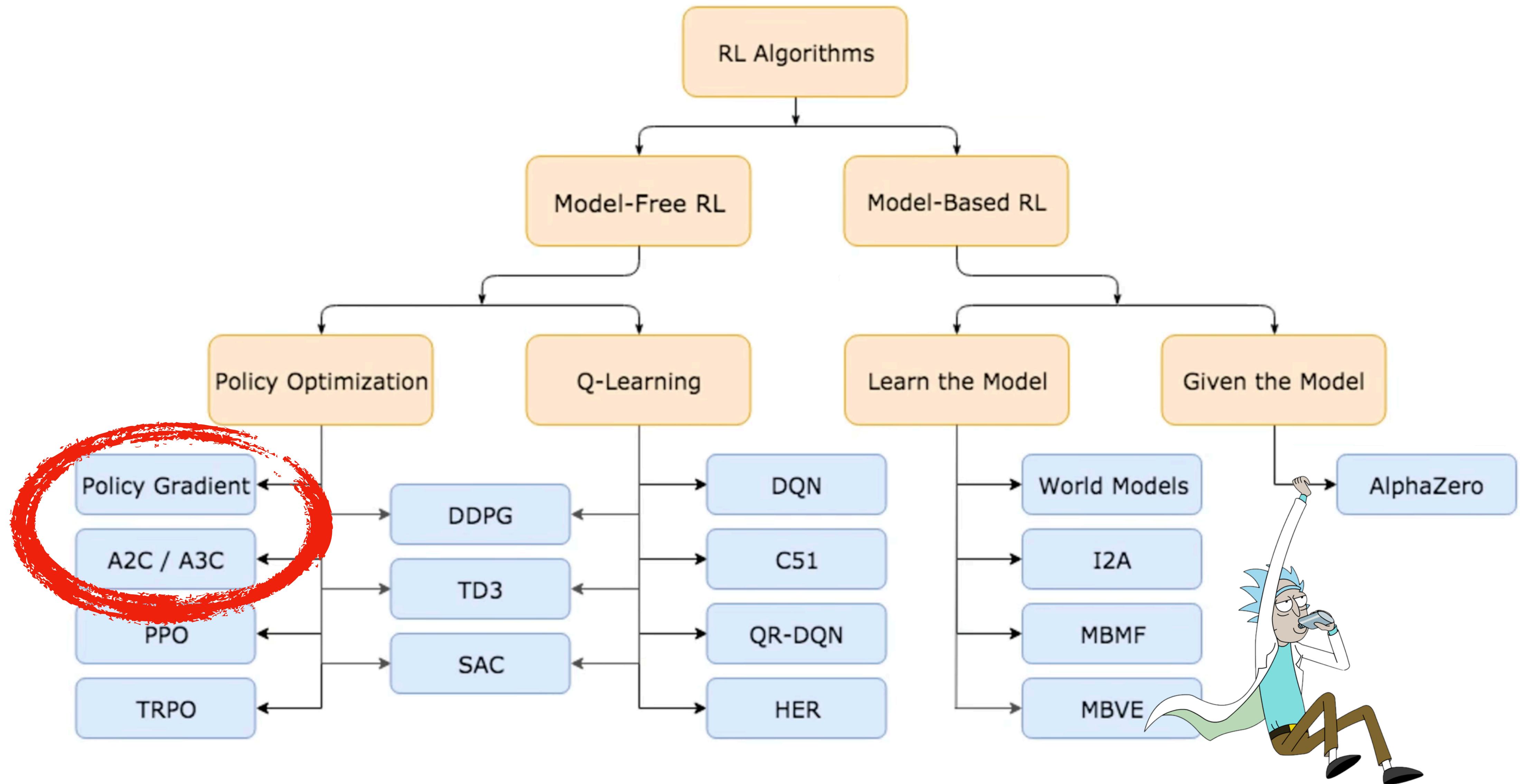


Actor-Critic



Actor-Critic

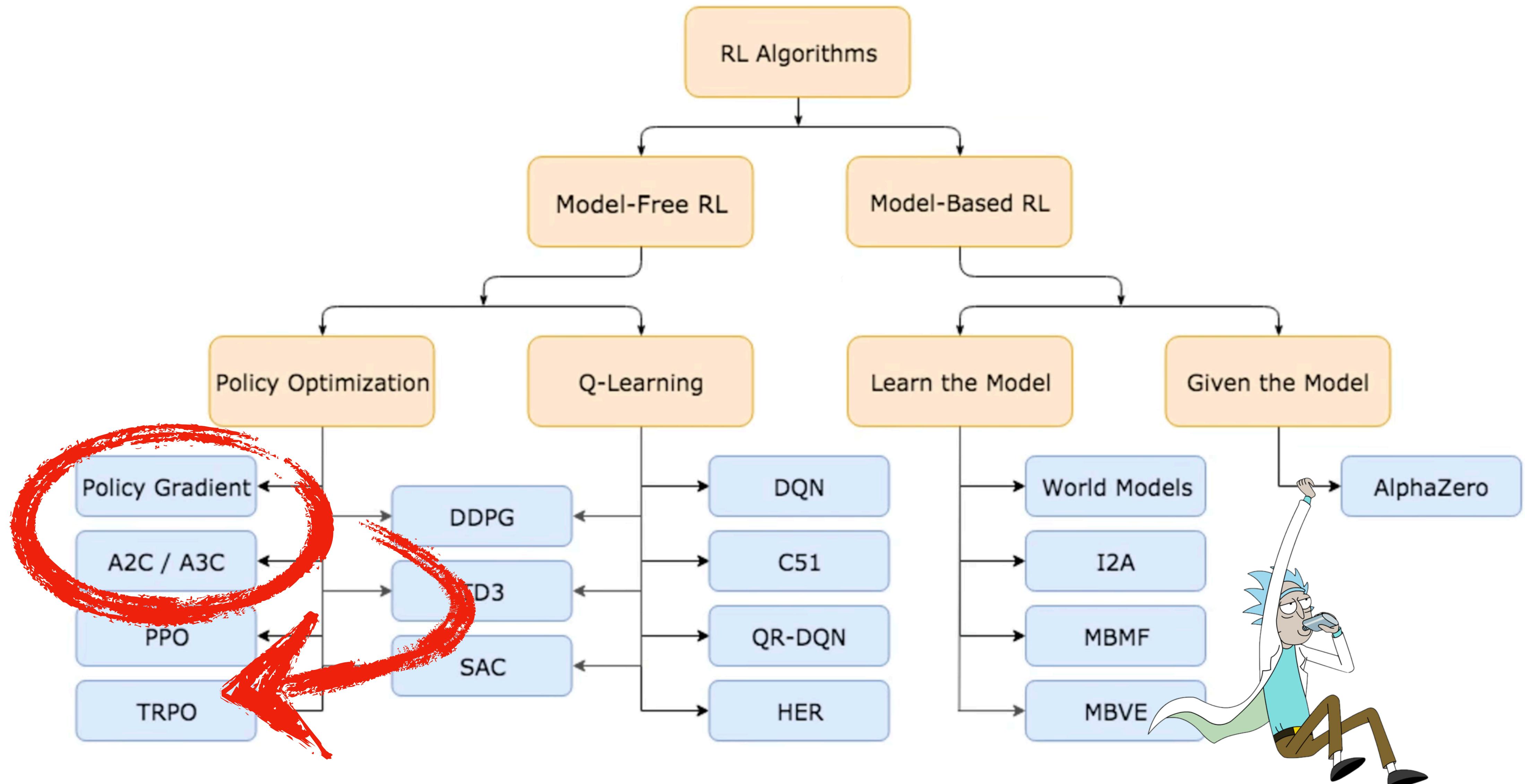




Проблема:

Если вдруг позволить агенту менять политику слишком сильно, то он может сделать катастрофически плохой шаг.

⇒ не сможем получать позитивных примеров



Trust Region Policy Optimization (TRPO)

<https://spinningup.openai.com/en/latest/algorithms/trpo.html>

<https://arxiv.org/pdf/1502.05477.pdf>

Чтобы гарантировать, что политика не изменяется слишком сильно и не попадала в катастрофические точки, мы добавим ограничение к задаче оптимизации. Ограничение представим в виде KL-дивергенции между новой и старой политикой (которую можно представить как расстояние между двумя распределениями вероятностей) должно быть меньше дельты (δ)

$$L^{PG}(\theta) = \hat{\mathbb{E}}_t \left[\log \pi_\theta(a_t | s_t) \hat{A}_t \right].$$

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right]$$

$$\text{subject to} \quad \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)]] \leq \delta.$$

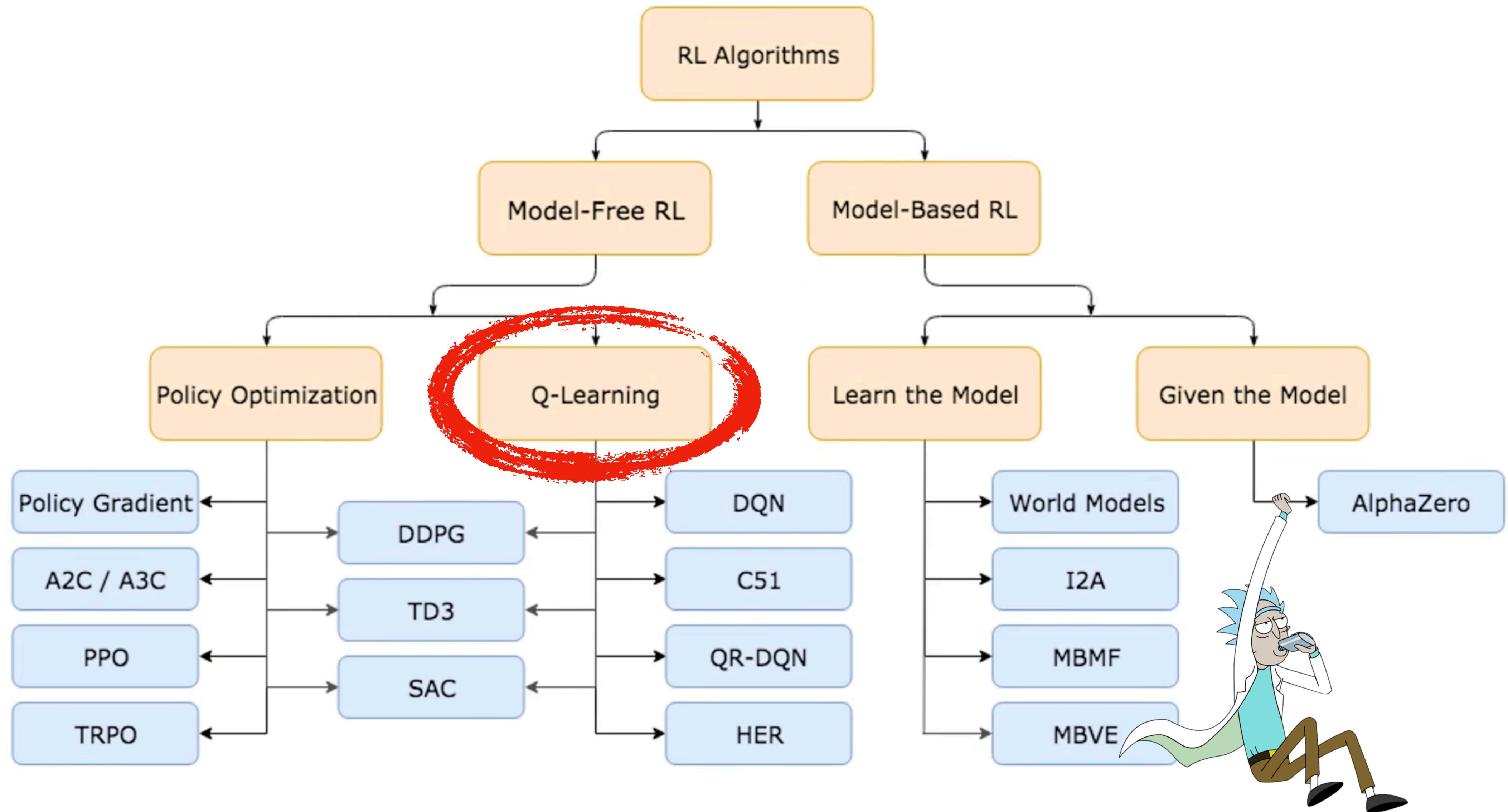
Proximal Policy Optimization (PPO)

<https://openai.com/blog/openai-baselines-ppo/>

<https://arxiv.org/pdf/1707.06347.pdf>

вместо того, чтобы добавлять ограничение отдельно, мы включаем его в целевую функцию в качестве штрафа (мы вычитаем из функции расхождение KL, умноженное на константу C).

$$\mathcal{L}_{PG} = - \sum_i r_i \ln \pi_\theta + C \times \text{KL}_{\pi_{\theta_{\text{old}}}}(\pi_\theta)$$

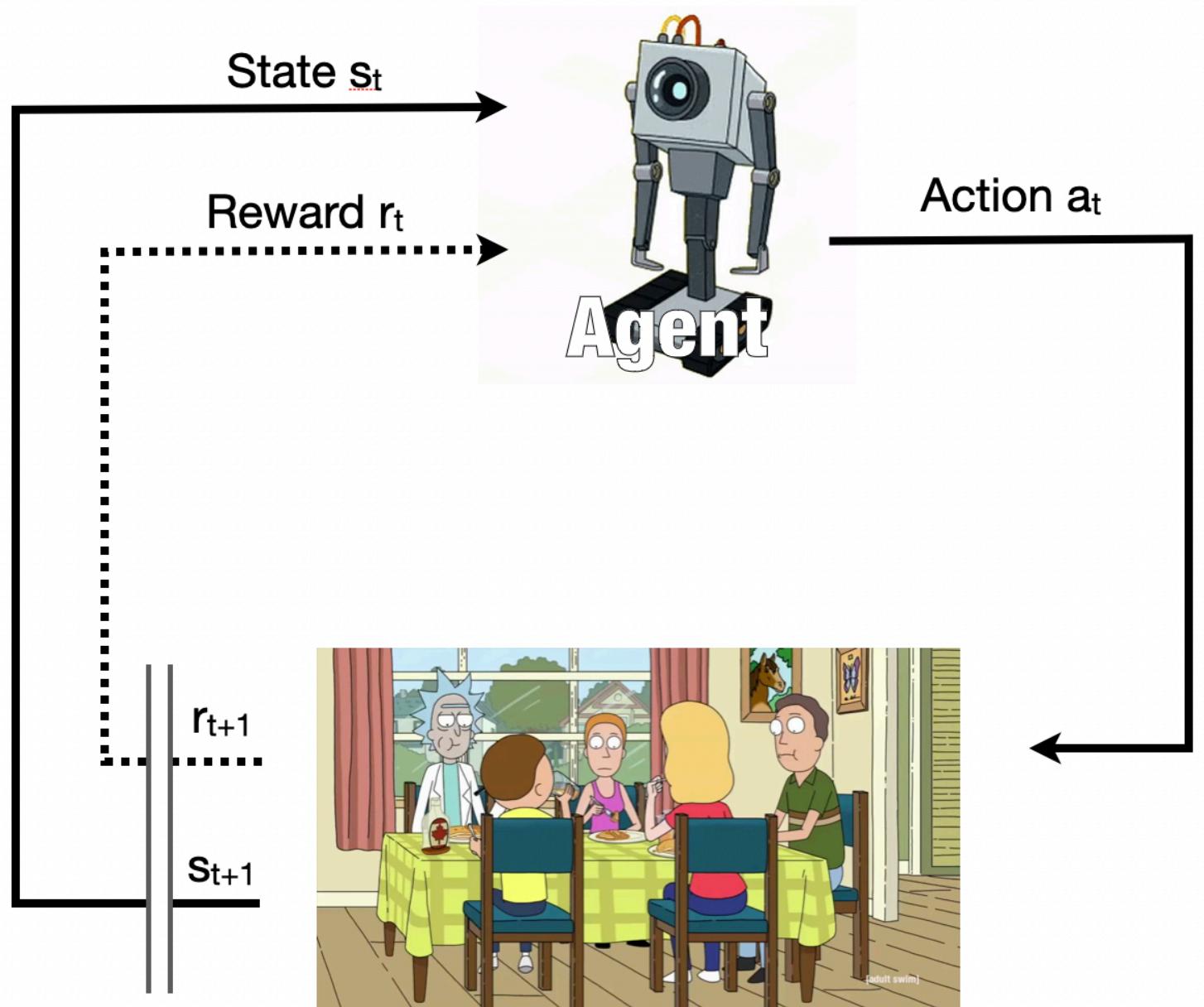


$$\sum_{t=1}^n r_t \longrightarrow \max$$

$$p(a | s)$$

$$v(s)$$

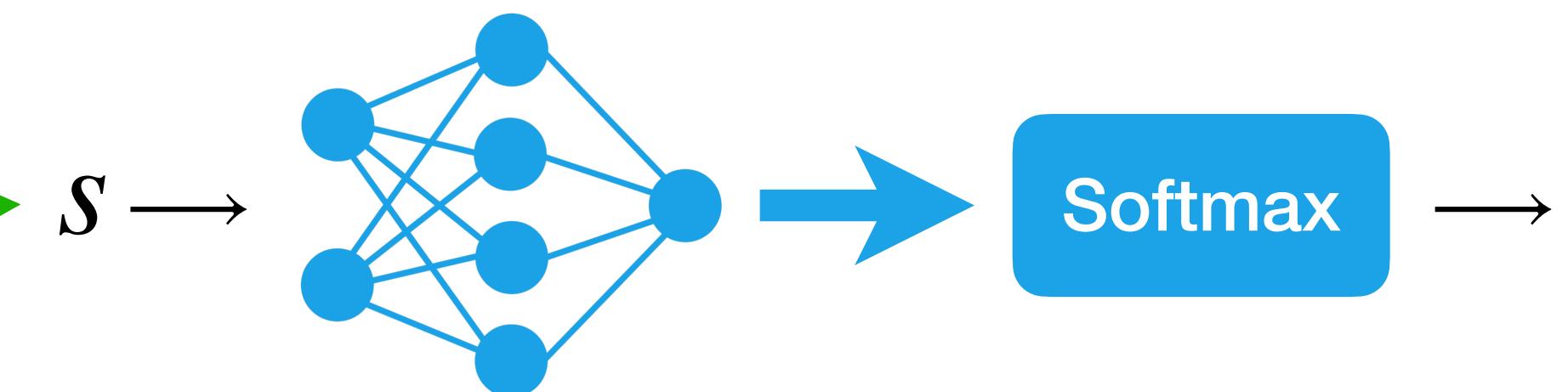
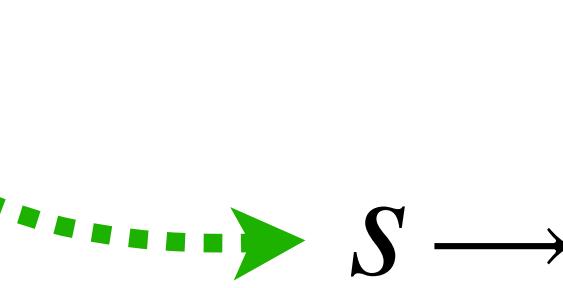
$$Q(s, a)$$



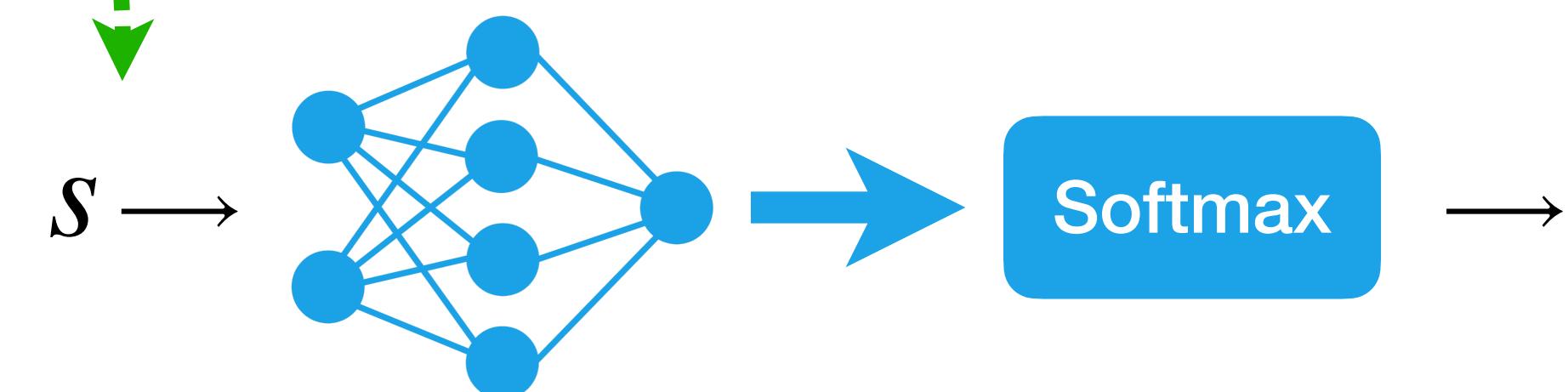
policy function

value function

Q-function

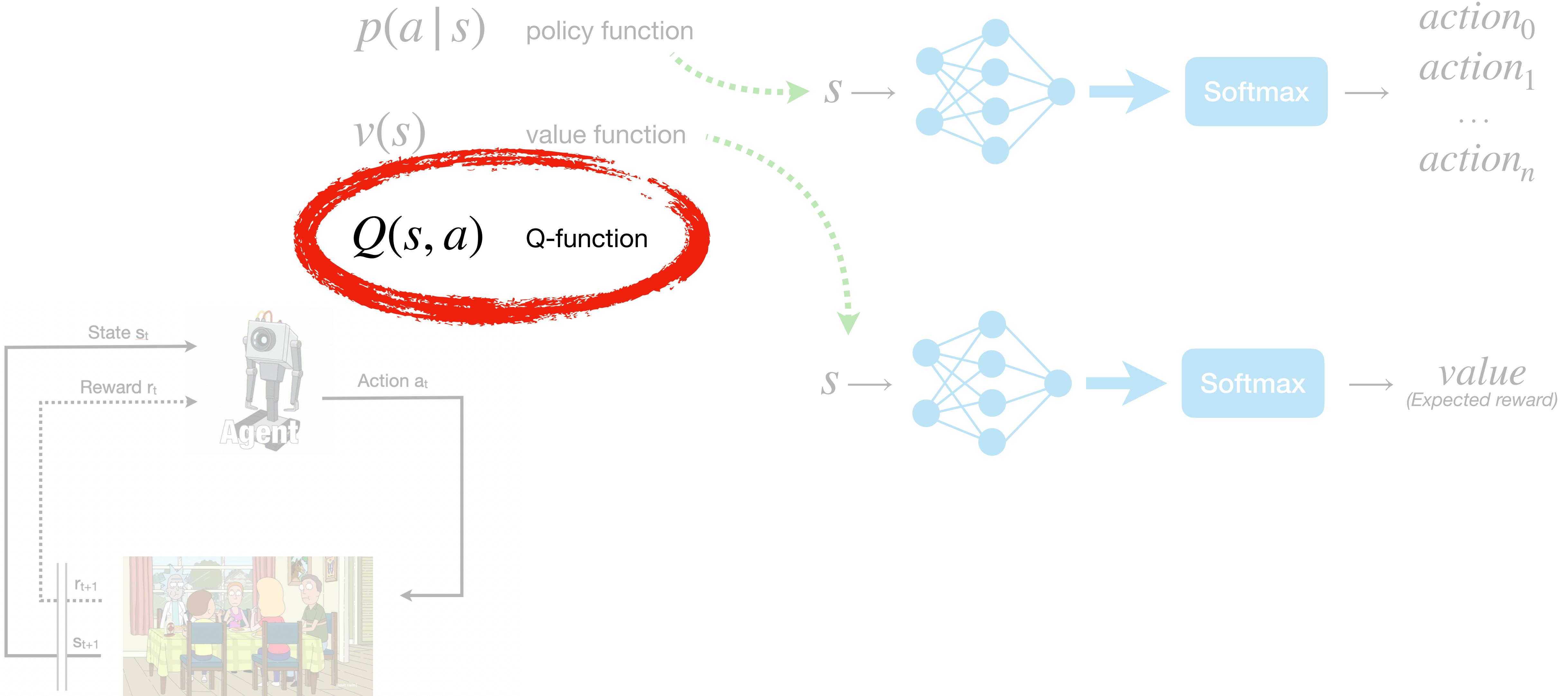


action₀
action₁
...
action_n



value
(Expected reward)

$$\sum_{t=1}^n r_t \longrightarrow \max$$



Bellman equation

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

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

\mathcal{E} – пространство возможных следующих состояний

Приближаем Q^* нейросетью с параметрами θ :

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

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

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

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right].$$

1. Initialize Table Q(S,A) with random values.

2. Take a action (A) with epsilon — greedy policy and move to next state S'

3. Update the Q value of a previous state by following the update equation:

$$Q(s, a) = \frac{\text{Old value}}{\text{learning rate}} + \alpha \left(\frac{\text{reward}}{\text{discount}} + \frac{\text{optimal future value}}{\text{Old value}} - Q(s, a) \right)$$

Step-by-step

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 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



Playing Atari with Deep Reinforcement Learning

<https://arxiv.org/pdf/1312.5602.pdf>



Что почитать?

**Deep Reinforcement Learning:
Pong from Pixels**

[http://karpathy.github.io/
2016/05/31/rl/](http://karpathy.github.io/2016/05/31/rl/)

**Deep Reinforcement Learning
(или за что купили DeepMind)**

<https://habr.com/ru/post/279729/>

**Курс Deep Learning на пальцах
от Семена Козлова**

- 13 - Reinforcement Learning
- 14 - Еще RL

[https://youtu.be/ x0ASf9jV9U](https://youtu.be/x0ASf9jV9U)
https://youtu.be/aOIK1i1xt_M

**Introduction to Reinforcement
Learning with David Silver**

[https://deepmind.com/learning-
resources/-introduction-
reinforcement-learning-david-silver](https://deepmind.com/learning-resources/-introduction-reinforcement-learning-david-silver)
<https://www.davidsilver.uk/teaching/>

**Reinforcement Learning: An
Introduction. Richard S. Sutton
and Andrew G. Barto**

[https://web.stanford.edu/class/
psych209/Readings/
SuttonBartoIPRLBook2ndEd.pdf](https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf)