

Historical Best Q-Networks for Deep Reinforcement Learning

基于历史最优Q网络的强化学习方法

Wenwu Yu, Rui Wang, Ruiying Li, Jing Gao, Xiaohui Hu

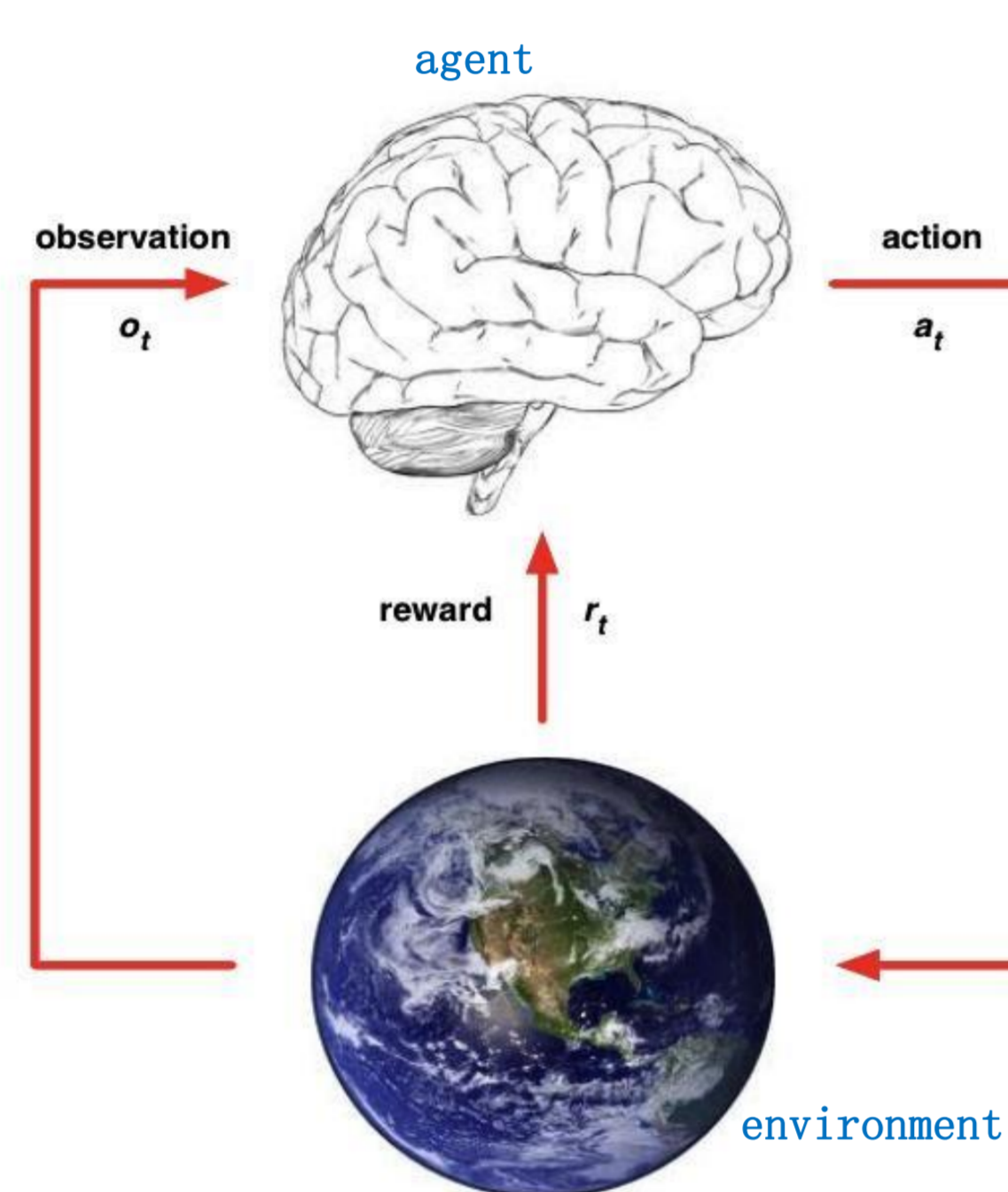
2018 IEEE 30th International Conference on Tools with Artificial Intelligence (ICTAI). IEEE, 2018: 6-11.

联系人: 俞文武

联系方式: 15652326966 (wenwu2016@iscas.ac.cn)

Background

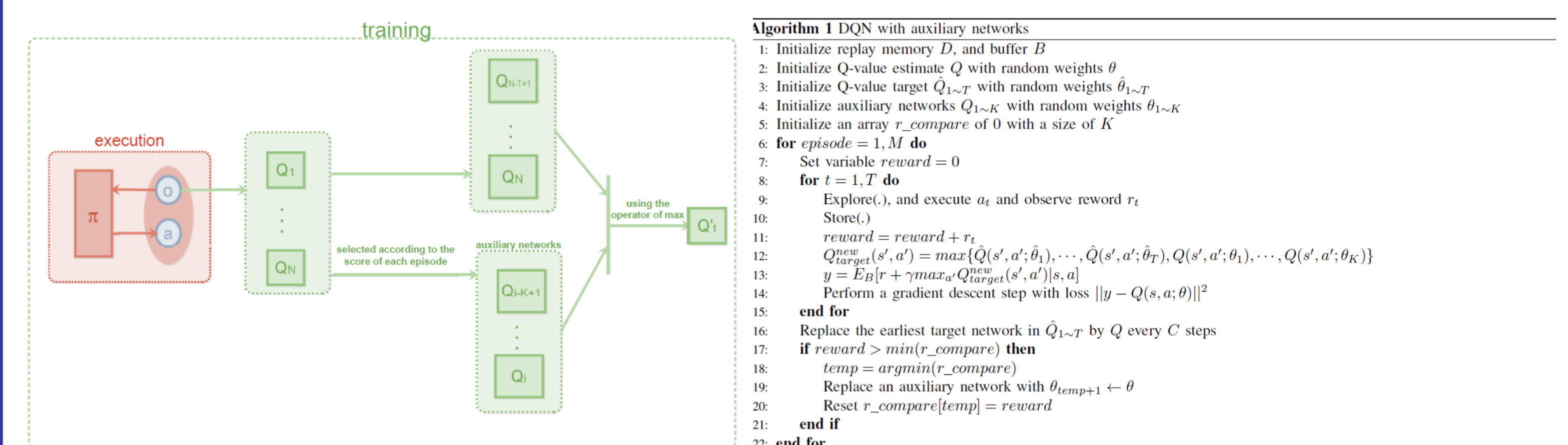
• In Reinforcement Learning (RL) an agent seeks the optimal policies, for sequential decision problems. And at each time step t , the agent observes state S_t , takes an action a_t and receive a scalar reward r .



However, the criteria used to determine which network is better as an auxiliary network is indeed a problem. To overcome this issue, naturally, we adopt the following measures:

- according to the score
- using the operator of max

The whole algorithm, which we call DQN with auxiliary networks, is presented in Algorithm (right).



```

Algorithm 1 DQN with auxiliary networks
1: Initialize replay memory  $D$  and buffer  $B$ 
2: Initialize Q-value estimate  $Q$  with random weights  $\theta$ 
3: Initialize Q-value target  $Q_{1-T}$  with random weights  $\theta_{1-T}$ 
4: Initialize auxiliary networks  $Q_{1-K}$  with random weights  $\theta_{1-K}$ 
5: Initialize an array  $r\_compare$  of 0 with a size of  $K$ 
6: for episode = 1,  $M$  do
7:   Set variable  $reward = 0$ 
8:   for  $t = 1, T$  do
9:     Explore(), and execute  $a_t$  and observe reward  $r_t$ 
10:    Store()
11:     $reward = reward + r_t$ 
12:     $Q_{new}(s', a') = \max\{Q(s', a'; \theta_1), \dots, Q(s', a'; \theta_T), Q(s', a'; \theta_K), \dots, Q(s', a'; \theta_K)\}$ 
13:     $y = E[y] + \gamma \max_{a'} Q_{new}(s', a'; \theta)$ 
14:    Perform a gradient descent step with loss  $\|y - Q(s, a; \theta)\|^2$ 
15:  end for
16: Replace the earliest target network in  $Q_{1-T}$  by  $Q$  every  $C$  steps
17: if  $reward > \min(r\_compare)$  then
18:    $temp = \arg\min(r\_compare)$ 
19:   Replace an auxiliary network with  $\theta_{temp+1} \leftarrow \theta$ 
20:   Reset  $r\_compare[temp] = reward$ 
21: end if
22: end for
  
```

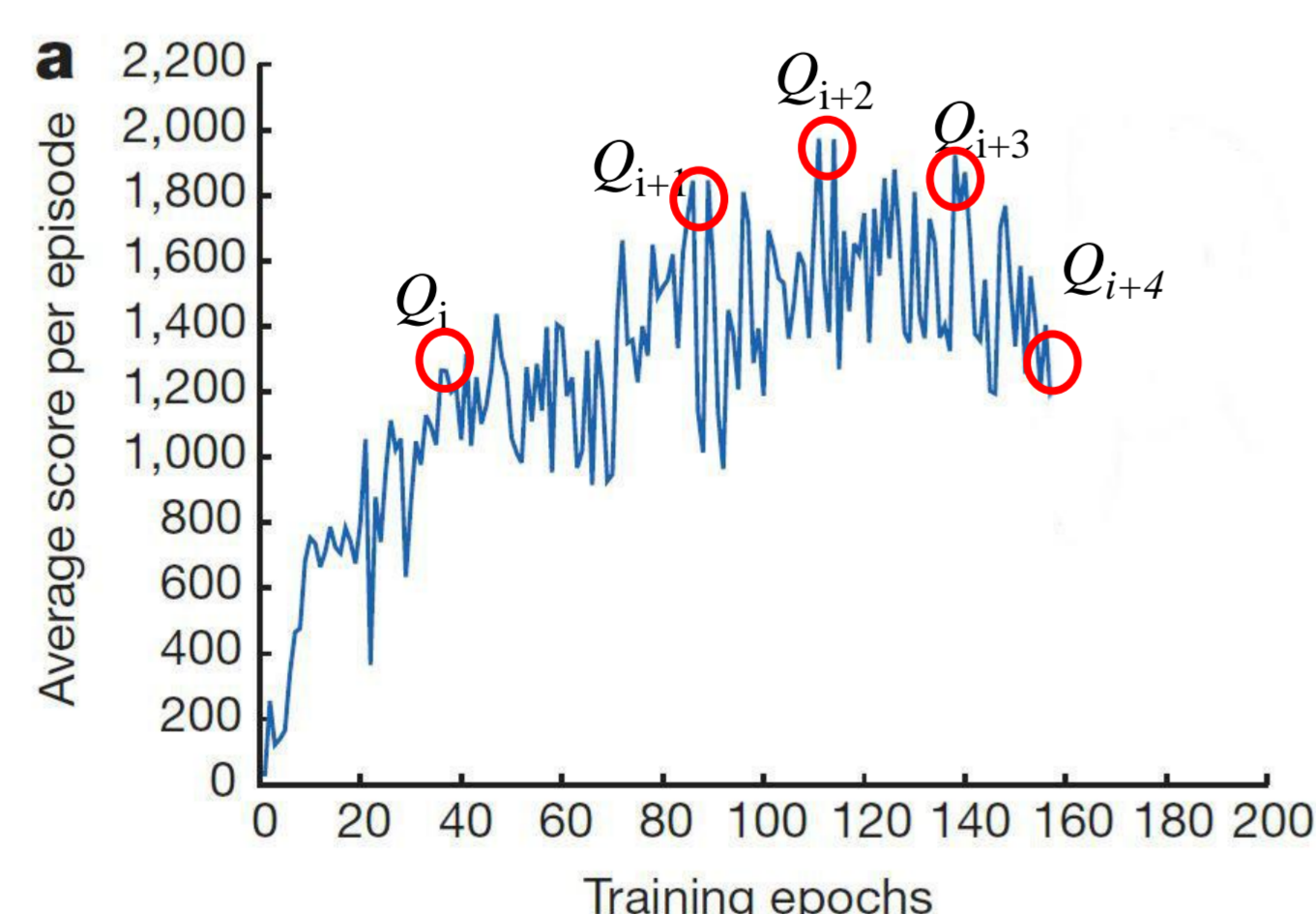
Difficulties

Some features in DRL (Deep Reinforcement Learning):

- have overestimation phenomena
- can not use the samples sufficiently
- have one target network: updated by the latest learned Q-value estimat

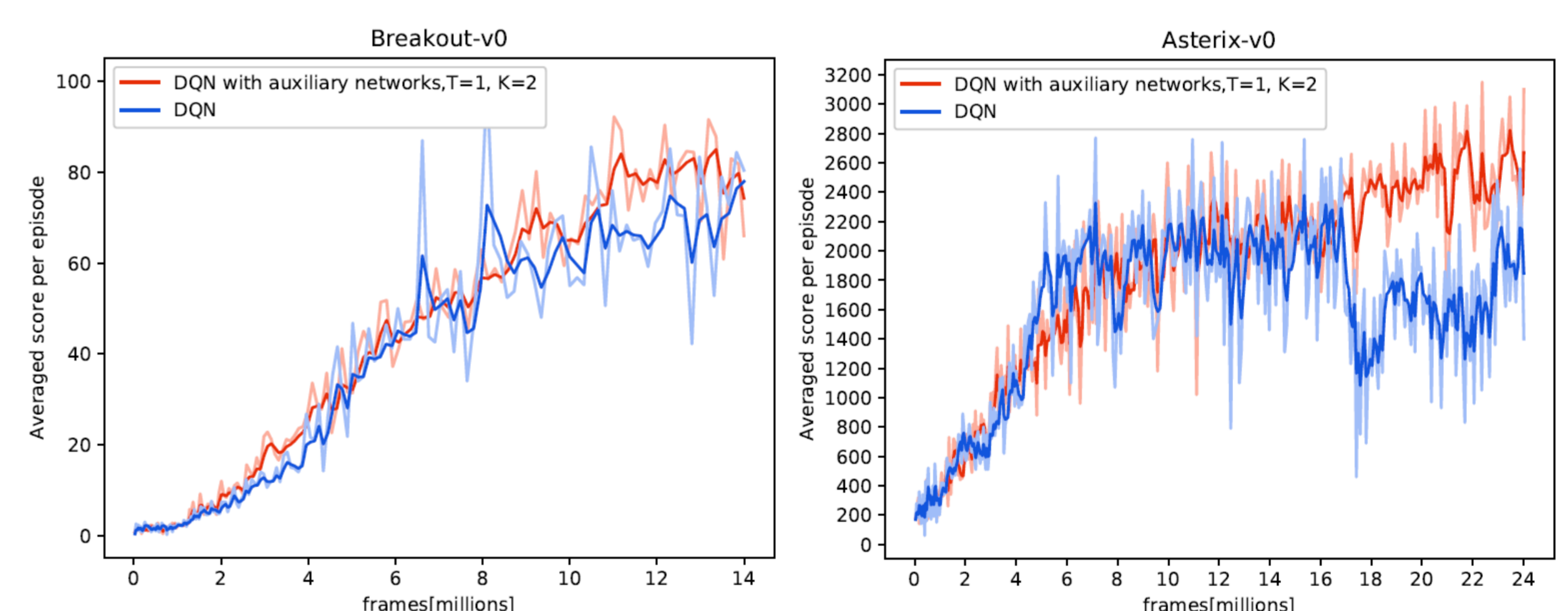
Can we use the historical Q-networks to generate a new target?

using the formula: $Q_{i+4} \leftarrow \{Q_i, Q_{i+1}, Q_{i+2}, Q_{i+3}\}$

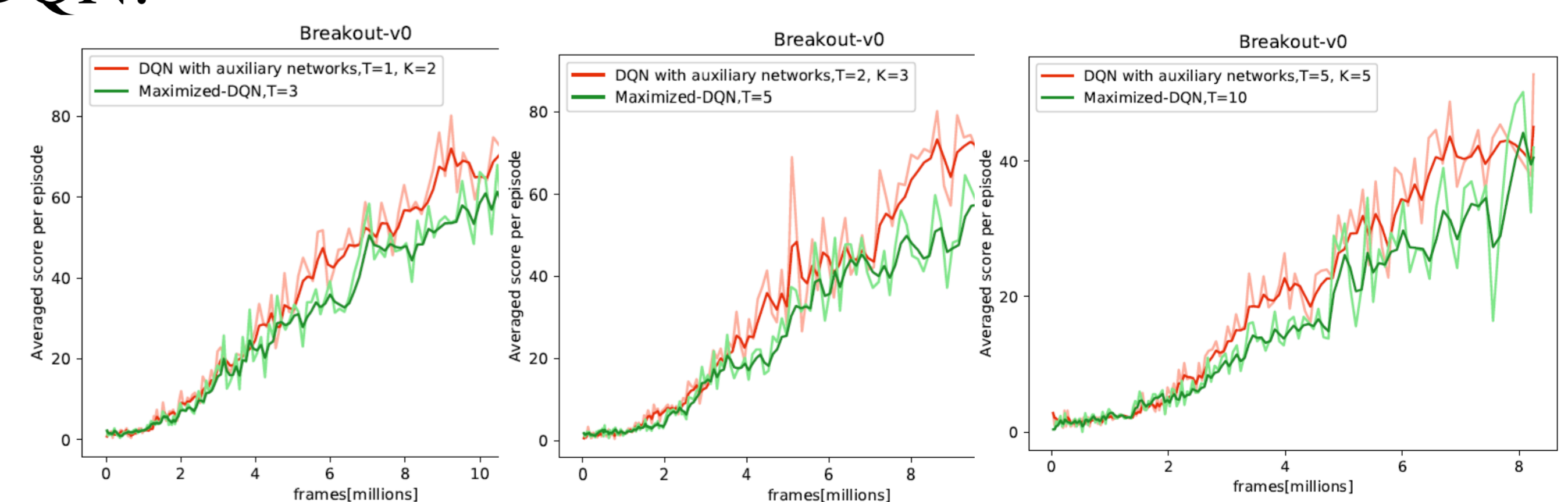


Experiment

• DQN with auxiliary networks compare with DQN:



• DQN with auxiliary networks compare with Maximized-DQN:



Method

This is the overview (left) of our auxiliary networks for deep learning approach. Our method, named DQN with auxiliary networks, has these networks:

- multiple target networks
- T latest previous target networks
- K auxiliary networks

Conclusion

choose several historical best networks as our auxiliary networks
use the score of each episode as the criteria
demonstrate that the auxiliary networks play an important role, not the operation of maximizing