



Introducere în Reinforcement Learning

Cursul #7

Ștefan Iordache, Cătălina Iordache, Ciprian Păduraru



Cuprins



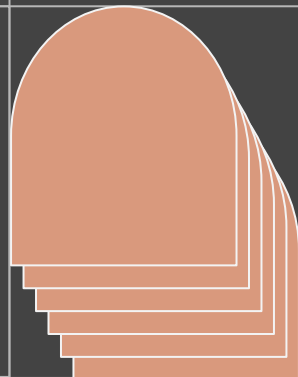
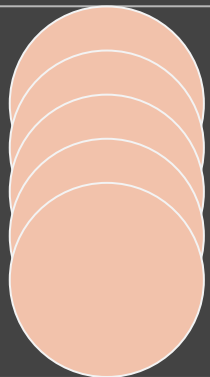
Q-Learning

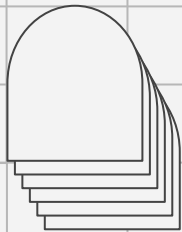


Deep Q-Networks

01

Q-Learning

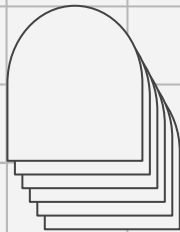




De ce să folosim Q-Learning?

- **Cel mai simplu algoritm, de înțeles și aplicat!**
- **Off-Policy & Model-Free!**
- **Idee simplă: Învățăm o politică care maximizează recompensa totală.**





De ce se numește Q-Learning?

Q = Quality (cât de utilă este o acțiune pentru rezultate viitoare?)



Formula magică!

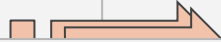
$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{current value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{current value}} \right)}_{\text{temporal difference}}$$

new value (temporal difference target)





Schelet Q-Learning

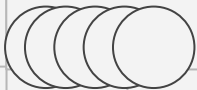


- **Inițializarea tabelului Q cu valoare zero (peste tot), urmând să ajustăm valorile în pașii de antrenare.**



```
1 import numpy as np
2
3 # state_size -> retine numarul de stari
4 # action_size -> retine numarul de actiuni posibile
5 Q = np.zeros((state_size, action_size))
```





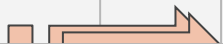
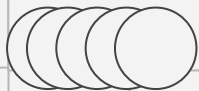
Schelet Q-Learning

- **Explorare & Exploatare**

```
1 import random
2
3 epsilon = 0.3
4
5 if random.uniform(0, 1) < epsilon:
6     """
7     EXPLORARE: selectăm o acțiune aleatorie
8     """
9 else:
10    """
11    EXPLOATARE: selectăm acțiunea cu valoare maximă
12    Valoare maximă = recompensa cea mai mare
13    """
```



Schelet Q-Learning



- **Actualizare valori Q**



```
1 current_value = Q[state, action]
2 new_value = reward + gamma * np.max(Q[new_state, :])
3 temporal_difference = new_value - current_value
4
5 Q[state, action] = (1 - lr) * current_value + lr * temporal_difference
```



Antrenare



```
1 %%time
2 """Antrenare agent"""
3
4 import random
5 from IPython.display import clear_output
6
7 # Hiperparametri
8 lr = 0.1
9 gamma = 0.6
10 epsilon = 0.1
11
12 # Colectii epoci si penalizari
13 all_epochs = []
14 all_penalties = []
15
16 for i in range(1, 100001):
17     state = env.reset()
18
19     epochs, penalties, reward, = 0, 0, 0
20     done = False
21
22     while not done:
23         if random.uniform(0, 1) < epsilon:
24             action = env.action_space.sample() # Explorare
25         else:
26             action = np.argmax(q_table[state]) # Exploatare
27
28         next_state, reward, done, info = env.step(action)
29
30         old_value = q_table[state, action]
31         next_max = np.max(q_table[next_state])
32
33         new_value = (1 - lr) * old_value + lr * (reward + gamma * next_max)
34         q_table[state, action] = new_value
35
36         if reward == -10:
37             penalties += 1
38
39         state = next_state
40         epochs += 1
41
42     if i % 100 == 0:
43         clear_output(wait=True)
44         print(f"Episodul: {i}")
45
46     print("Antrenare finalizata.\n")
```

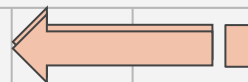
Evaluare



```
1  """Evaluarea performanței Q-learning (după antrenare)"""
2
3  total_epochs, total_penalties = 0, 0
4  episodes = 100
5
6  for _ in range(episodes):
7      state = env.reset()
8      epochs, penalties, reward = 0, 0, 0
9
10     done = False
11
12     while not done:
13         action = np.argmax(q_table[state])
14         state, reward, done, info = env.step(action)
15
16         if reward == -10:
17             penalties += 1
18
19         epochs += 1
20
21     total_penalties += penalties
22     total_epochs += epochs
23
24 print(f"Număr mediu de pași per episod: {total_epochs / episodes}")
25 print(f"Număr mediu de penalizări per episod: {total_penalties / episodes}")
```

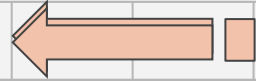


Hiperparametri - Descriere



- **Alpha (α)** – Learning rate. Trend descrescător pe parcursul învățării.
- **Gamma (γ)** – Valoare cu trend descrescător. Către finalul episoadelor preferăm să luăm în calcul tot mai mult recompensa imediat următoare în locul celor pe termen lung.
- **Epsilon (ϵ)** – Spre finalul etapei de învățare nu mai avem nevoie de explorare, ci de exploatare. Experimentarea repetată conduce către o valoare cu trend descrescător.

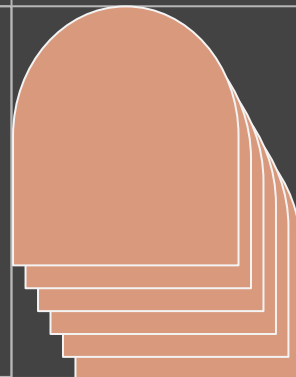
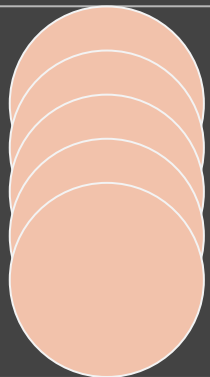
Hiperparametri - Tuning

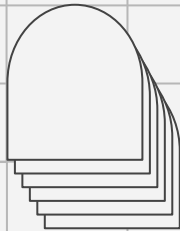


- **Grid search** – Foarte utilă la începutul experimentării cu orice algoritm de inteligență artificială.
- **Random search**
- *Ray.tune: Hyperparameter Optimization Framework*

02

Deep Q- Networks



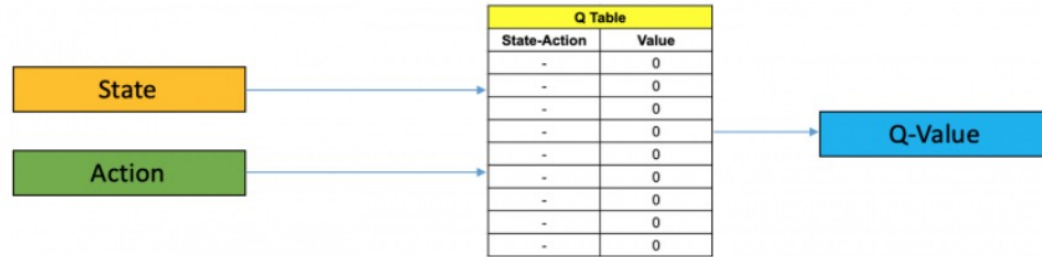


De ce avem nevoie de mai mult?

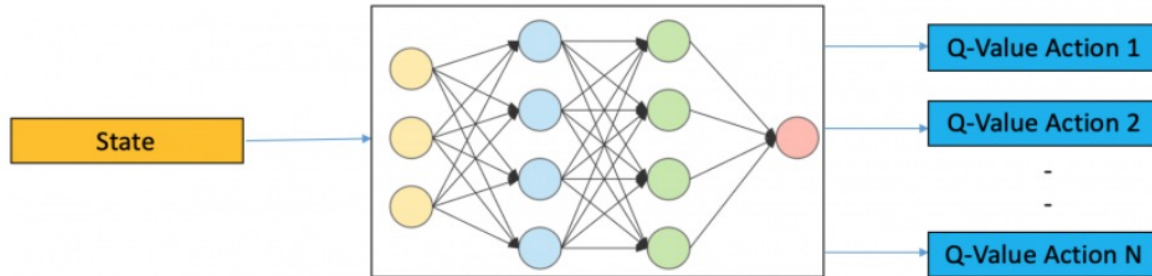
- **Mediile foarte mari sunt problematice!**
 - Necesarul de memorie pentru salvarea și actualizarea tabelului este foarte mare.
 - Timpul de explorare este nerealist.



Cum funcționează DQN/DQL?

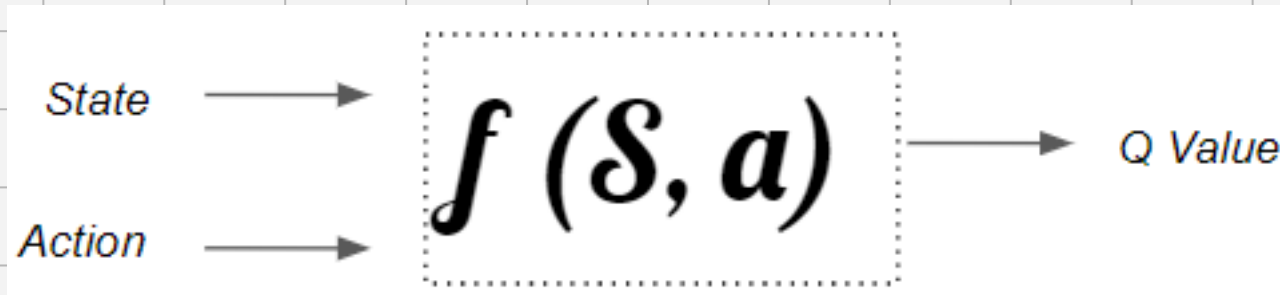
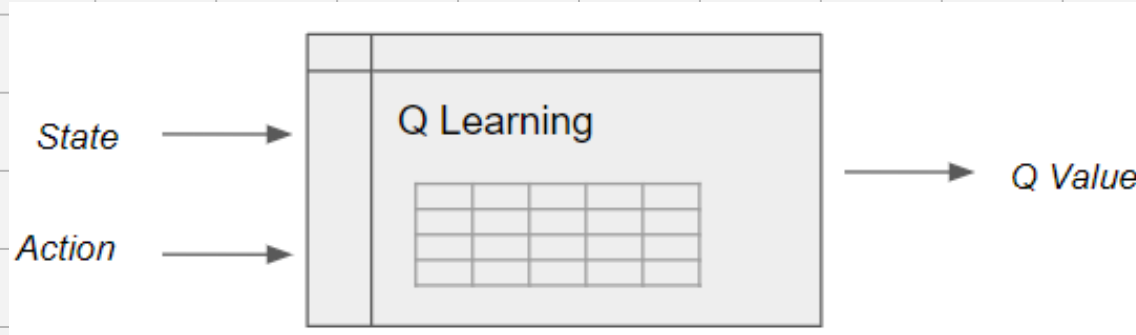


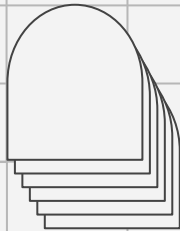
Q Learning



Deep Q Learning

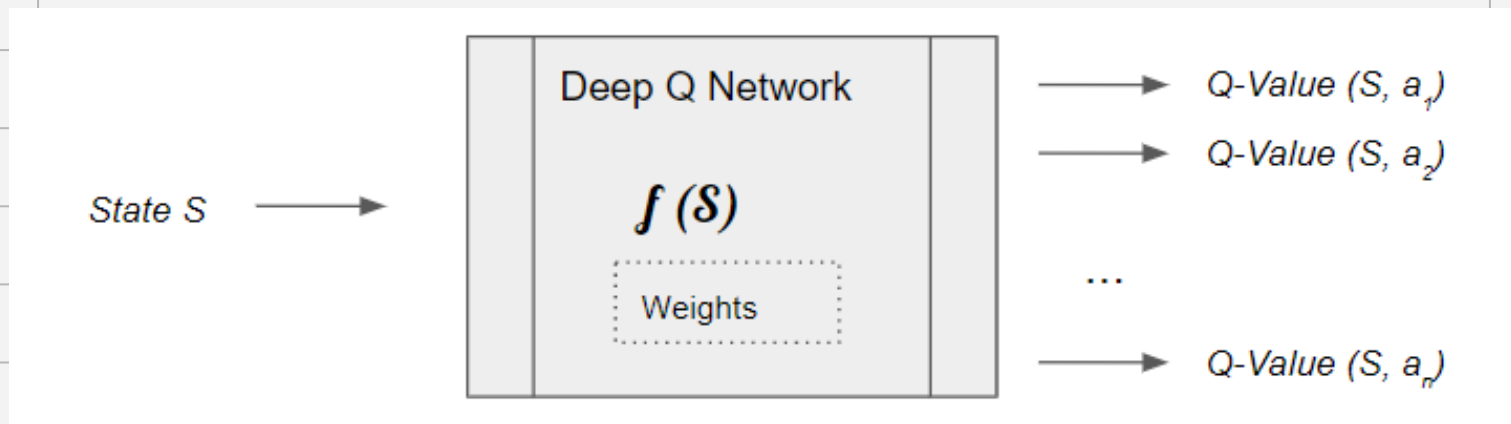
Cum funcționează DQN/DQL?



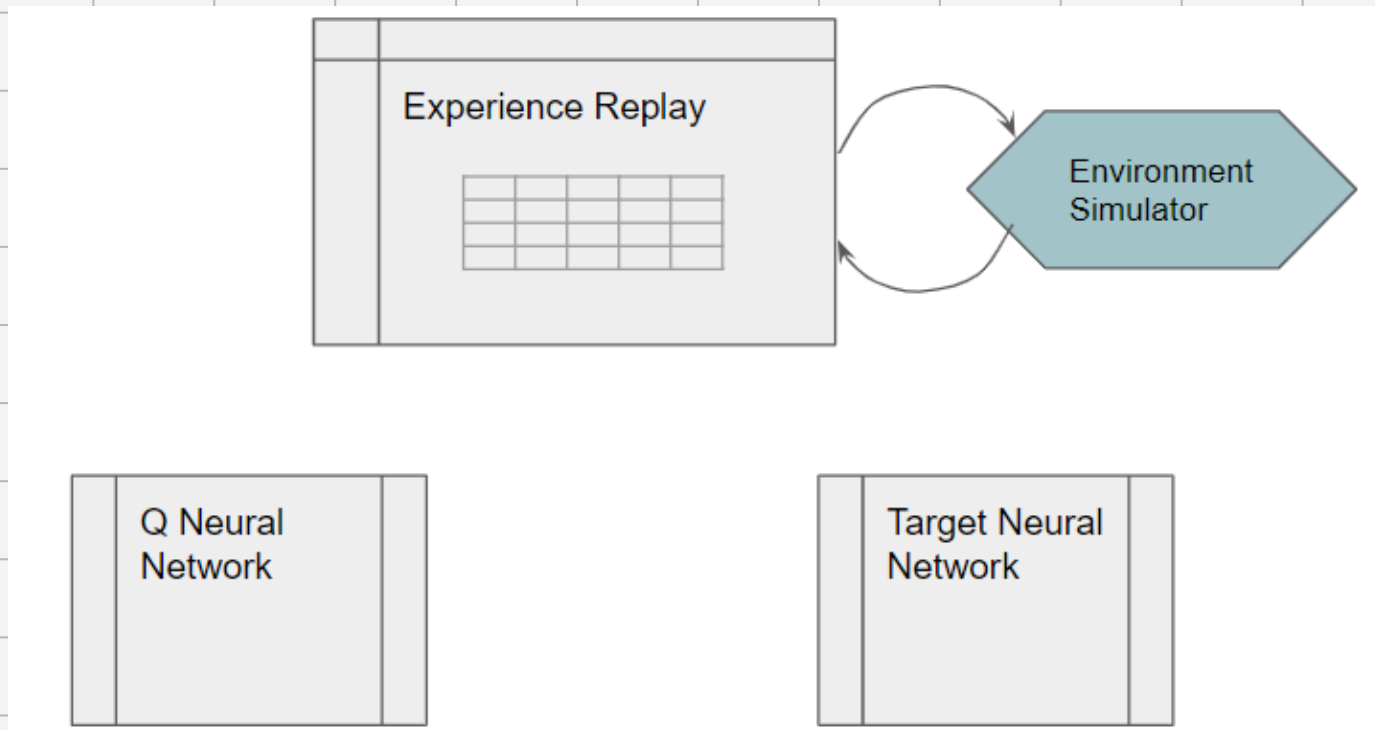


Neural Networks

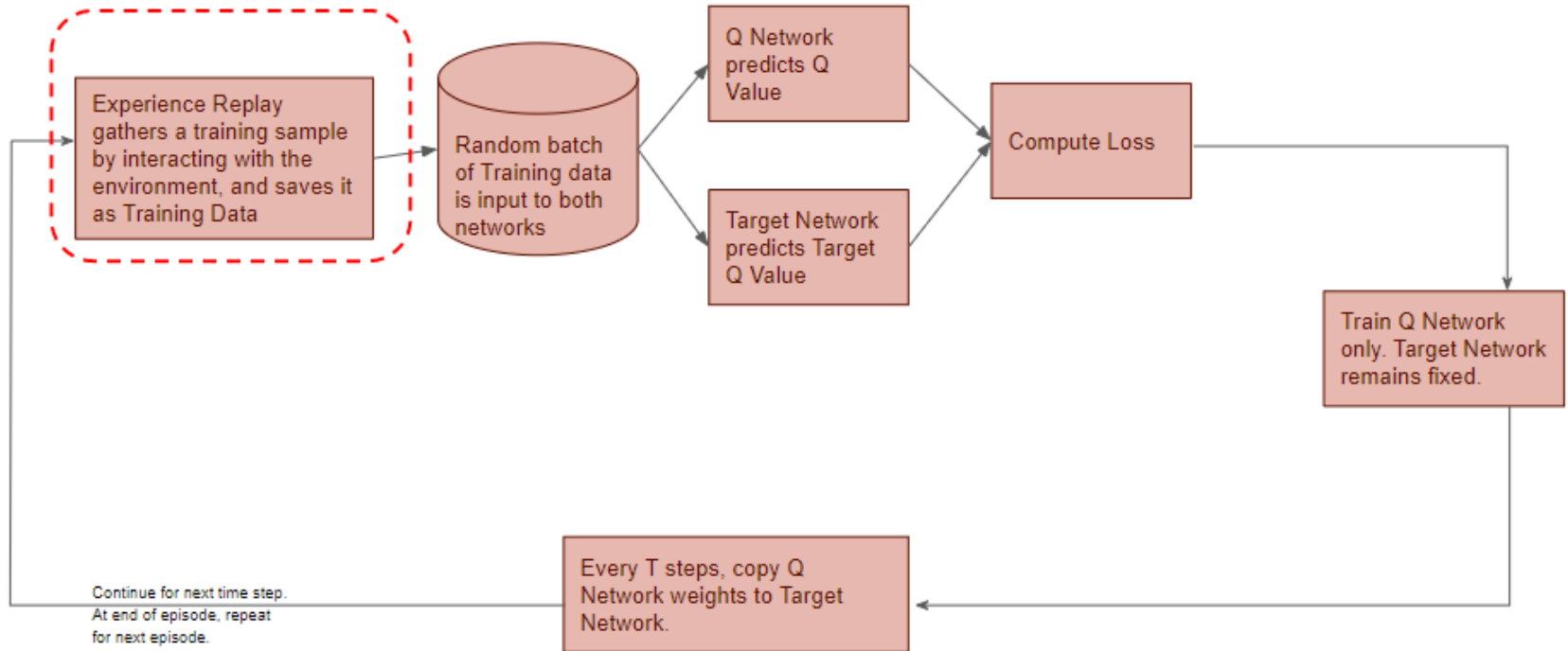
Cea mai bună metodă de aproximare a funcțiilor complexe!



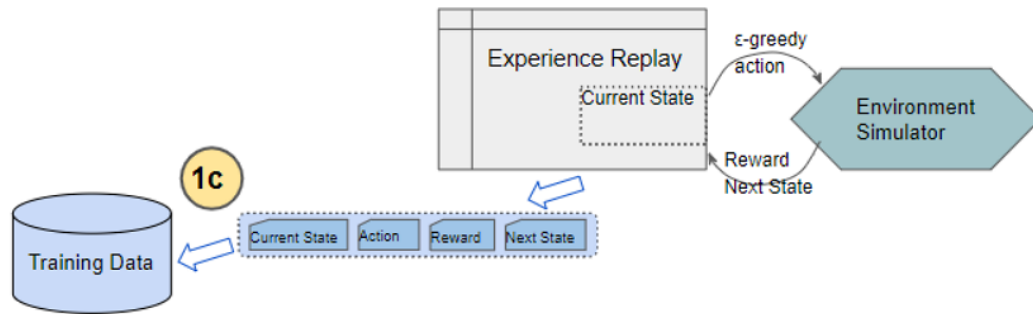
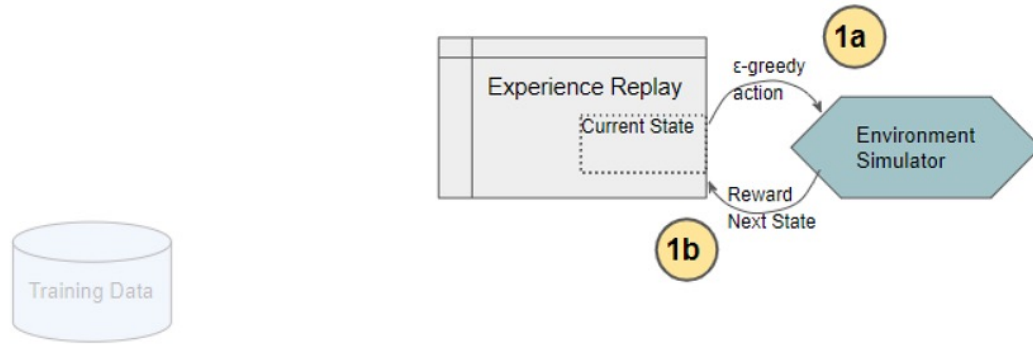
Arhitectura DQN/DQL



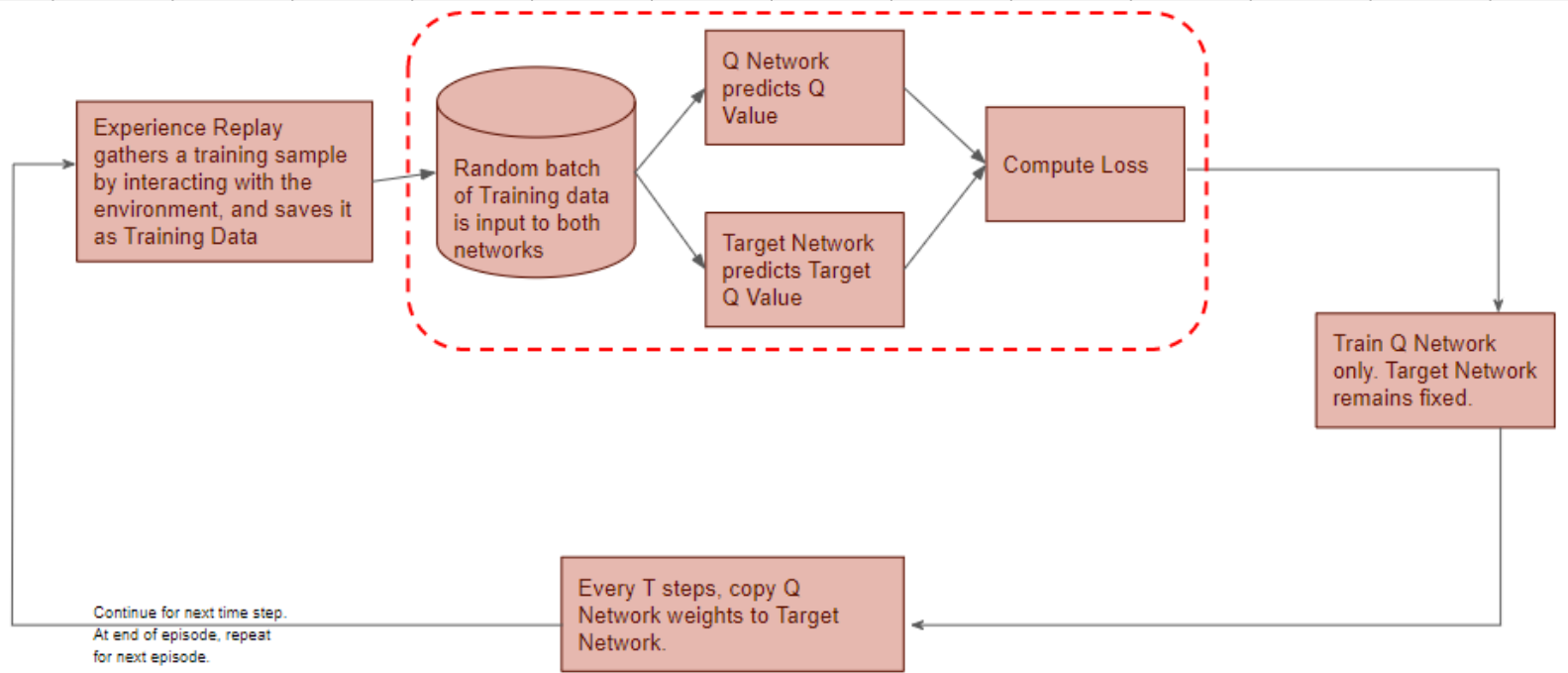
Extragerea datelor de antrenare



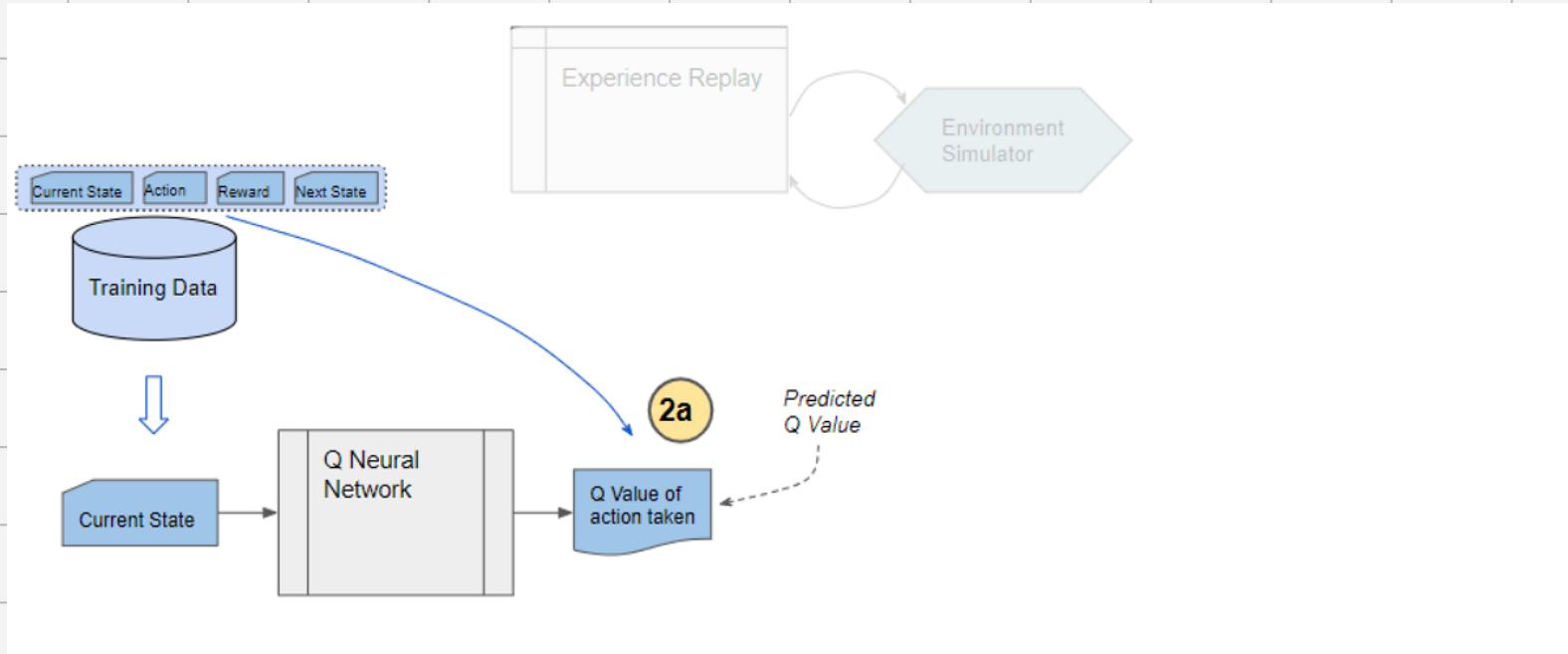
Extragerea datelor de antrenare



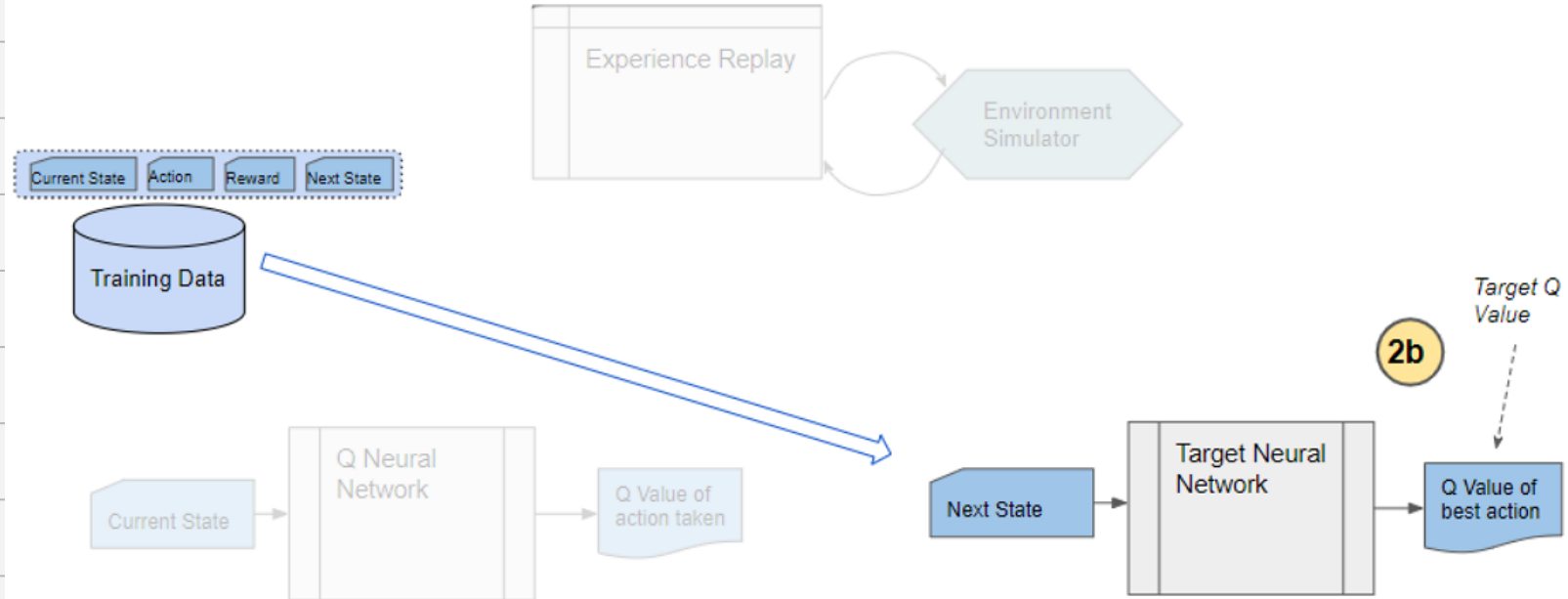
Predicții



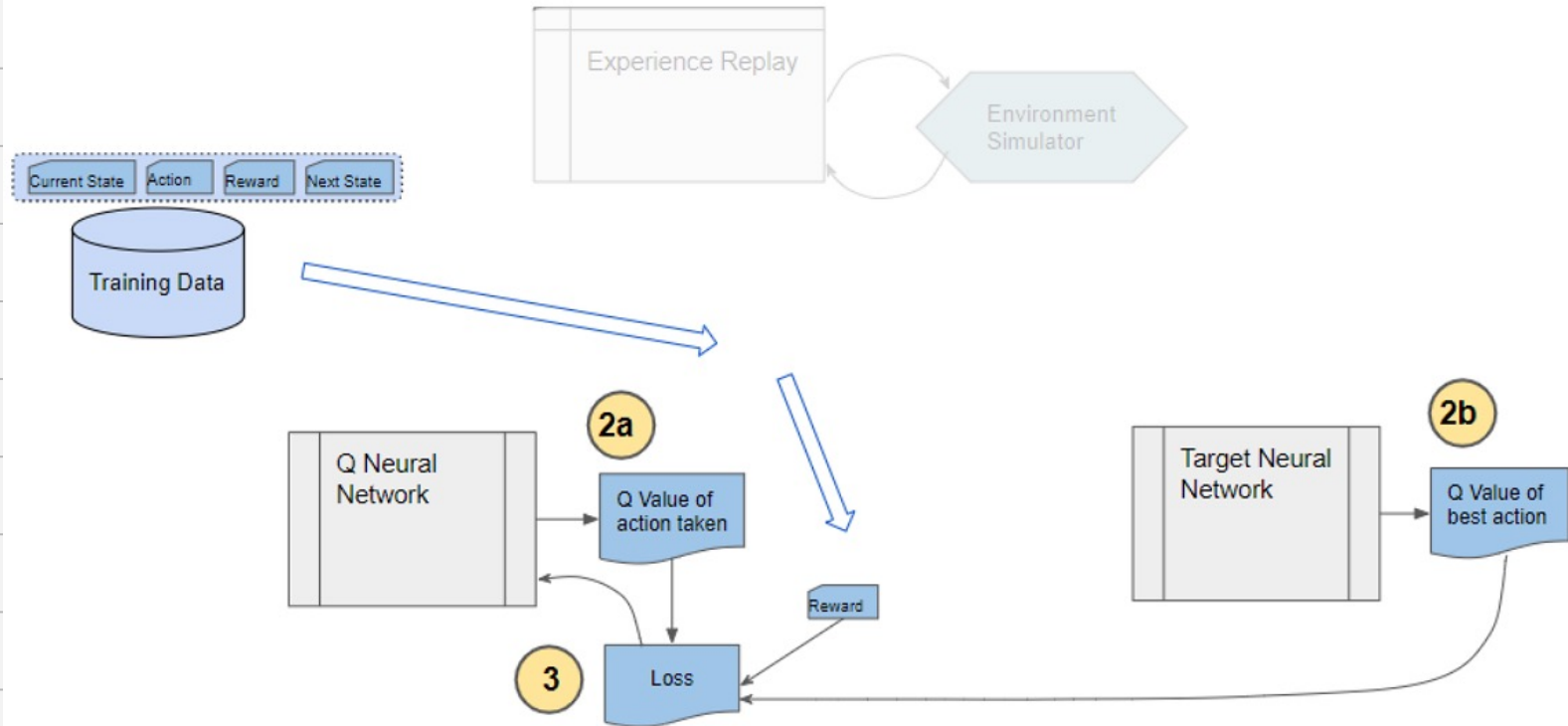
Predicții – Q-value

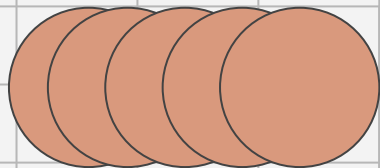


Predicții – Target Q-value



Loss & Q-Network Train



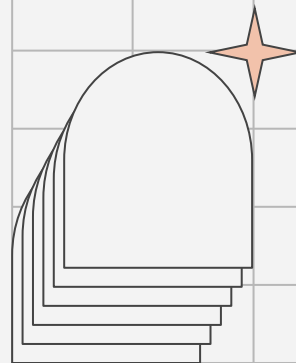


Thanks!

Este timpul pentru întrebări!!!

stefan.iordache10@s.unibuc.ro
catalina.patilea@s.unibuc.ro
ciprian.paduraru@fmi.unibuc.ro

+40 7.. ...



CREDITS: This presentation template was created by **Slidesgo**, and includes icons by **Flaticon** and infographics & images by **Freepik**

