

# 用領域知識提升DQN學習效率之研究

Use Domain Knowledge to Make DQN  
Learn Faster for Tic-Tac-Toe Game

何佩蓁

台北市私立東山高中

楊宗憲 博士

賽微科技股份有限公司

# Objectives

- Use **Deep Q Networks** (DQN) to achieve unbeatable Tic-Tac-Toe.
- Use **domain knowledges** to make DQN utilize training data more efficiently.

O		
O	X	
X		X

# Deep Q Networks

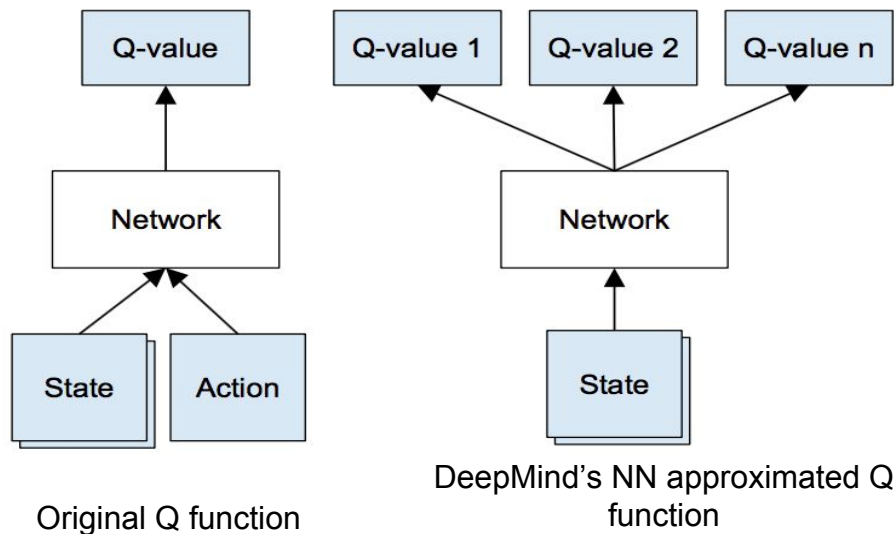
- **DQN = Q-Learning + NN (DeepMind 2013, 2015)**

- $Q(s,a)$  approximated by Neural Networks
- $Q(s) \rightarrow$  a list of Q-values for all actions, for speedup
- Experience Replay: random sampling from recently collected training data.



Atari 2600 games

<http://neuro.cs.ut.ee/demystifying-deep-reinforcement-learning/>

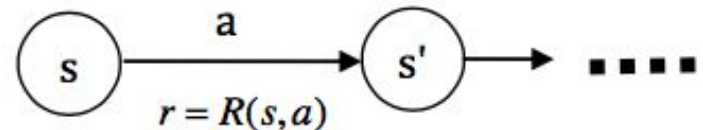


- **Bellman Estimates target  $Q^t$  from  $(s, a, r, s')$**

$$Q(s; \theta) \rightarrow Q = [q_1, q_2, \dots, q_n]$$

$$Q(s'; \theta^-) \rightarrow Q' = [q'_1, q'_2, \dots, q'_n]$$

$$q_a^t = \begin{cases} r, & \text{if } s' \in S_T \\ r + \gamma \times \max_i q'_i, & \text{else} \end{cases}$$



$$Q^t \leftarrow Q$$

$$Q^t[a] \leftarrow q_a^t \Rightarrow Q^t = [q_1, q_2, \dots, q_a^t, \dots, q_n]$$

# DQN Implementing Tic-Tac-Toe

- **Markov Decision Process**

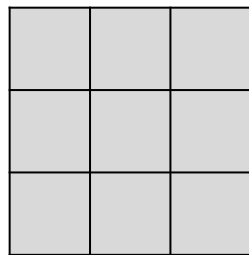
- Stochastic environment: opponent may move randomly

- $s$  : current board state

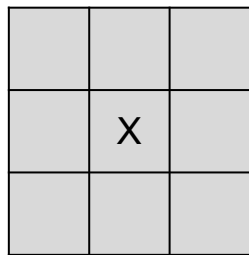
- $a$  : 9 board locations to put pieces

- $r = R(s') = \begin{cases} 1 & \text{if DQN wins} \\ -1 & \text{if DQN loses} \\ 0.1 & \text{tie} \\ 0 & \text{game not finished} \end{cases}$

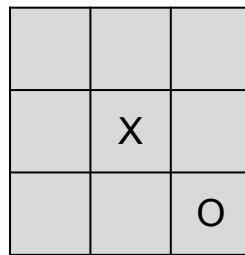
- $s'$  : the board state after opponent puts the piece



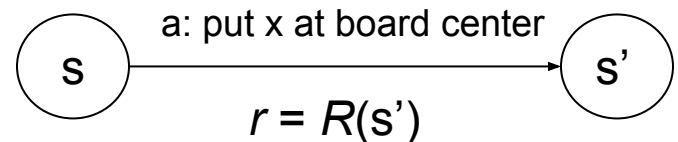
$s$



$temp$

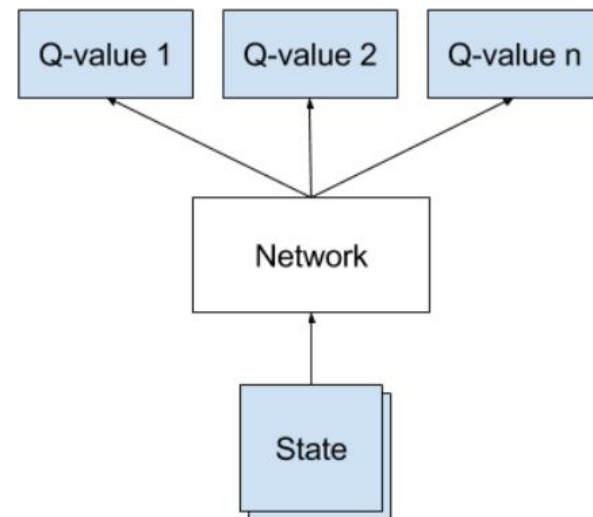
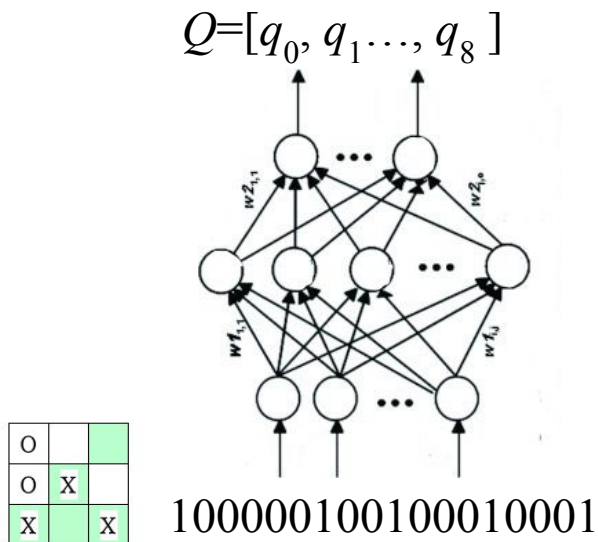


$s'$



# MLP Q Function

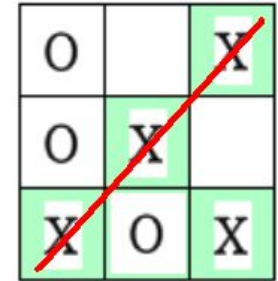
- Input: 2 bits for each square on board, totally  $9 \times 2 = 18$  neurons.
- Hidden: one layer, 36 neurons,  $\tanh(x)$
- Output: 9 neurons,  $f(x) = x$ , for Q values of 9 squares on board.



# Make DQN Learn Faster

- **Extra copies of determinative transitions into Experience Replay**

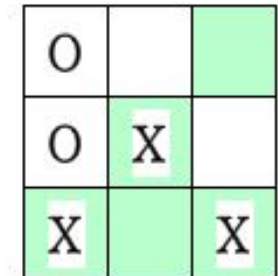
- For  $(s, a, r, s')$ ,  $s'$  is terminal state.
- Higher probability for determinative transitions to be selected as training samples



O		X
O	X	
X	O	X

- **Early rewarding for 2-way winning board**

- DQN learns to deliberately move toward 2-way winning board states to win the game.



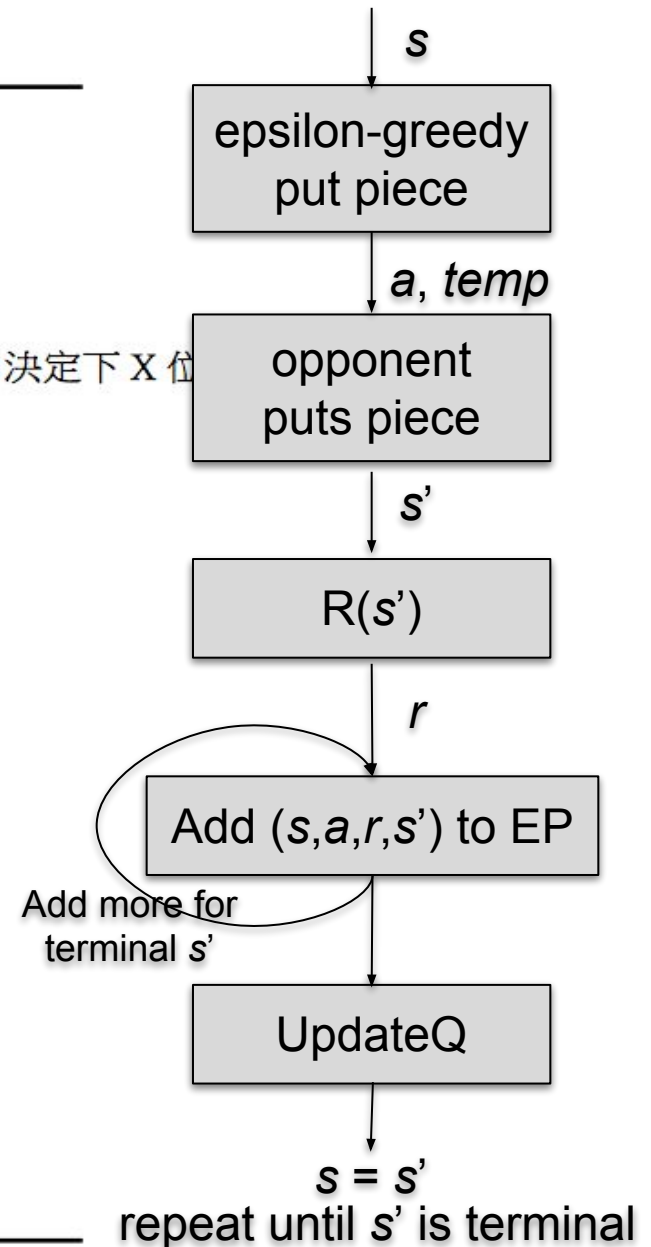
O		
O	X	
X		X

# 演算法 1：用 DQN 訓練井字遊戲先手下 X

```

1: 以隨機值初始化 MLP 參數  $\theta$ 
2:  $\theta^- \leftarrow$  複製  $\theta$ 
3: for  $episode = 1, 2, \dots, T$  do
4:    $s \leftarrow$  初始化空盤面
5:   while  $s$  非終局盤面 do
6:      $a \leftarrow$  以  $\epsilon$ -greedy 方式隨機或  $\arg \max MLP(s; \theta)$  決定下 X 位
7:      $temp \leftarrow$  先手於盤面  $s$  的  $a$  位置下 X
8:     if  $temp$  為終局盤面 then
9:        $s' \leftarrow temp$ 
10:    else
11:       $s' \leftarrow$  對手在盤面  $temp$  隨機下 O
12:    end if
13:     $r = R(s')$ 
14:    將  $(s, a, r, s')$  加入經驗回放
15:    if  $s'$  為終局盤面 then
16:      增加  $(s, a, r, s')$  加入經驗回放的數量
17:    end if
18:     $\theta \leftarrow UpdateQ(\theta, \theta^-)$ 
19:     $s \leftarrow s'$ 
20:  end while
21: 每隔若干次  $episode$  更新  $\theta^- \leftarrow$  複製  $\theta$ 
22: 每隔若干次  $episode$  衰減  $\epsilon \leftarrow 0.9 \times \epsilon$ 
23: end for
  
```

move one piece



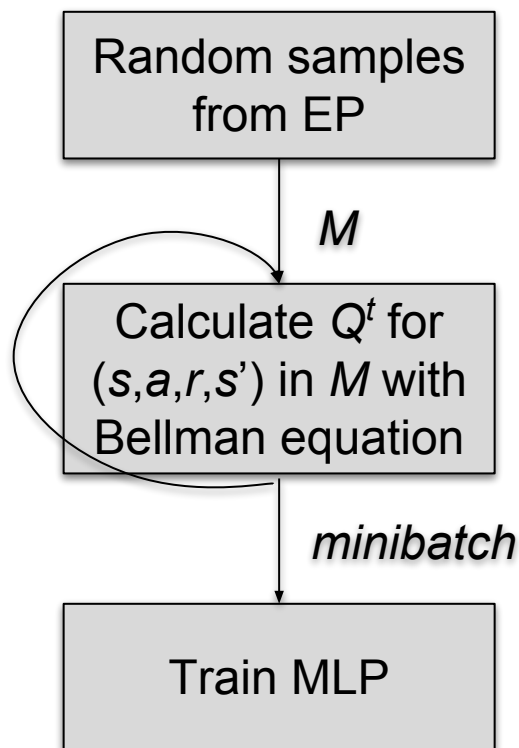


## 演算法 2：更新 MLP 的參數 $\theta$

---

```
1: function UpdateQ( $\theta, \theta^-$ )
2:   從經驗回放隨機挑選若干筆資料放入  $M$ ，包含最新加入的該筆資料
3:    $minibatch \leftarrow$  空串列
4:   for each  $(s, a, r, s')$  in  $M$  do
5:      $MLP(s; \theta) \rightarrow Q = [q_0, q_1, \dots, q_8]$ 
6:     if  $s'$  為終局盤面 then
7:        $q_a^t = r$ 
8:     else
9:        $MLP(s'; \theta^-) \rightarrow Q' = [q'_0, q'_1, \dots, q'_8]$ 
10:       $q_a^t = r + \gamma \max_{a'} (q'_{a'}, s'[a'] \notin \{O, X\})$ 
11:    end if
12:     $Q^t \leftarrow Q$ 
13:     $Q^t[a] = q_a^t$ 
14:    將  $(s, Q^t)$  加入  $minibatch$ 
15:  end for
16:   $\theta \leftarrow TrainMLP(\theta, minibatch)$ 
17:  return  $\theta$ 
```

---



# Experiments

- **Experimental environment**

- Training episodes: 20k for 1<sup>st</sup> mover, 100k for 2<sup>nd</sup>
- MLP Learning rate  $\alpha$  :  $\alpha = 0.01$ ,  $\alpha \leftarrow \alpha \times 0.9$ , decay 50 times
- $\epsilon$ -greedy :  $\epsilon = 1$ ,  $\epsilon \leftarrow \epsilon \times 0.9$ , decay 50 times
- Experience Replay capacity : 2048 , minibatch size: 32
- 10 identical experiments for each setting, randomly initialize  $\theta$
- Store  $\theta$  for every 100 training episodes

- **Test Criterion**

- Unbeatable, playing Tic Tac Toe against randomly moving program for 10,000 rounds without losing any game.

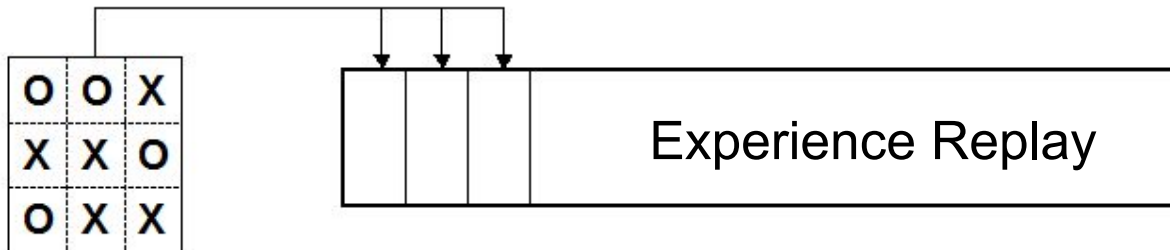
# Experiment 1

- Adding extra copies of determinative transitions to Experience Replay

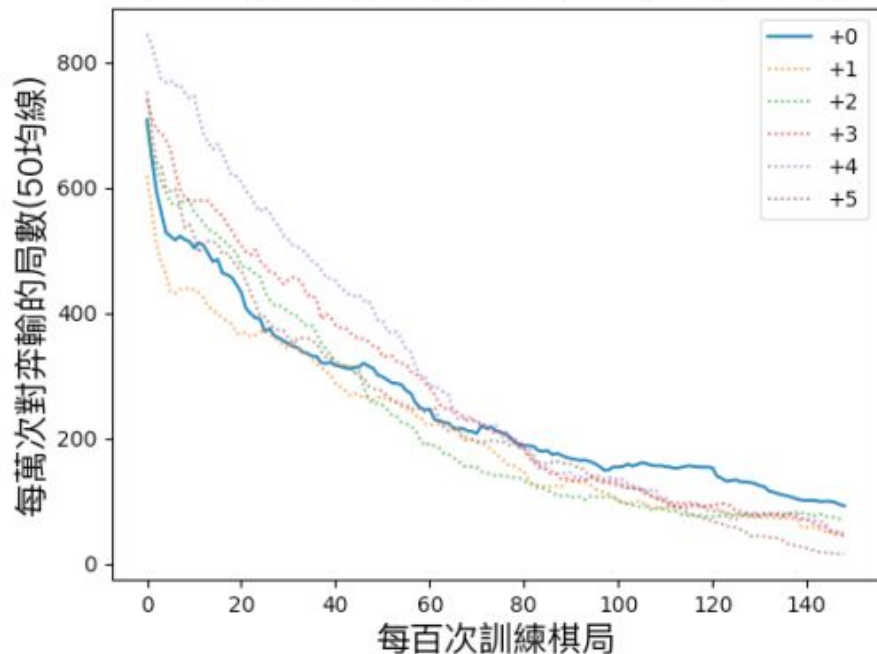
Average number of episodes for 1<sup>st</sup> and 2<sup>nd</sup> movers to reach unbeatable

extra copies	0	1	2	3	4	5
1 <sup>st</sup> mover	9.7k	8.9k	11.0k	10.5k	11.3k	9.6k
2 <sup>nd</sup> mover	82.5k	68.9k	72.0k	62.0k	58.0k	65.6k

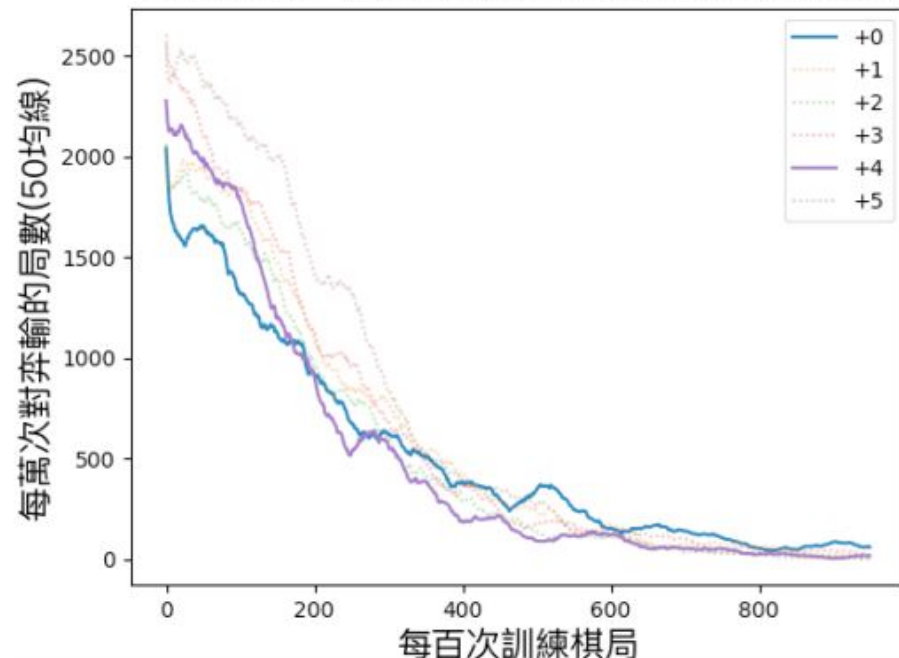
30% reduction



先手訓練時，在經驗回放中複製有價值盤面的數量



後手訓練時，在經驗回放中複製有價值盤面的數量



For both DQN trained player 1 and 2, adding extra copies of determinative transitions to Experience Replay can make learning faster during the training progress, as compared with the baselines.

# Experiment 2

- **Early rewarding for 2-way winning for 1<sup>st</sup> mover**

Numbers of episodes for 1<sup>st</sup> mover to reach unbeatable after adding reward

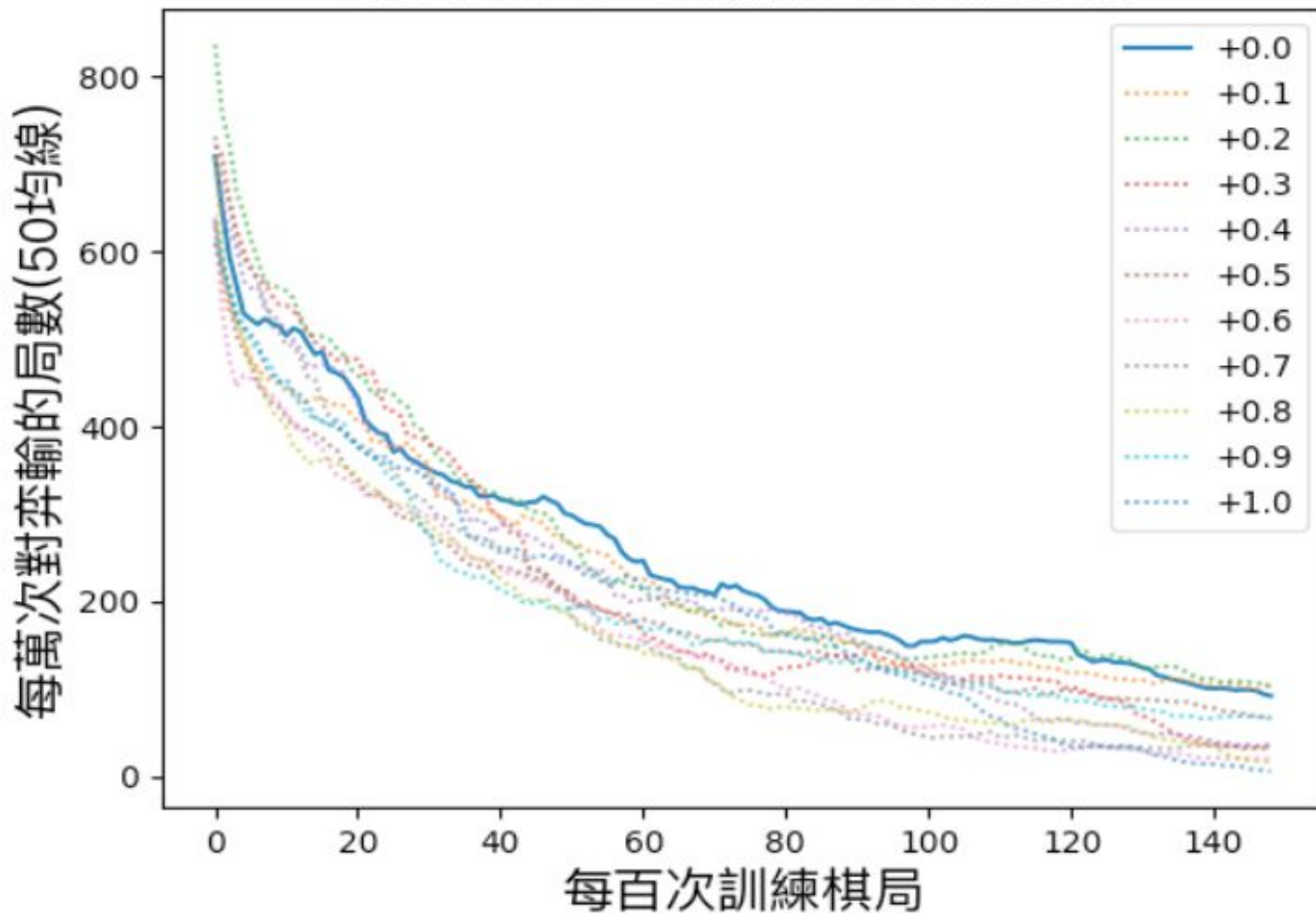
Reward	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
Episodes	9.7k	10.2k	8.8k	9.4k	9.2k	8.9k	7.9k	7.8k	8.5k	7.0k	8.5k

28% reduction

- When the reward is higher than 0.5, it starts to improve significantly in terms of reaching unbeatable criterion.

O		
O	X	
X		X

先手訓練時，雙活路盤面獎勵分數



In terms of learning speed, offering extra rewards for two-way winning board outperforms the baseline noticeably no matter how much the reward is.

# Experiment 3

- **Cutting training data in half when training 2<sup>nd</sup> mover**
  - Reducing to 50k training episodes
  - Due to the second mover's disadvantage, DQN training may fail to reach unbeatable criterion without sufficient training episodes.
  - Result: improving training reliability by nearly 100%.

The effect on learning stability of DQN when adding extra copies

extra copies	0	1	2	3	4	5
reach unbeatable (%)	40	80	70	70	80	80

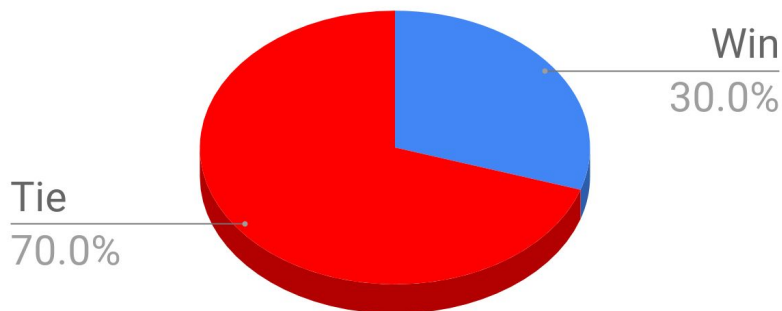
improve nearly 100%

# Experiment 4

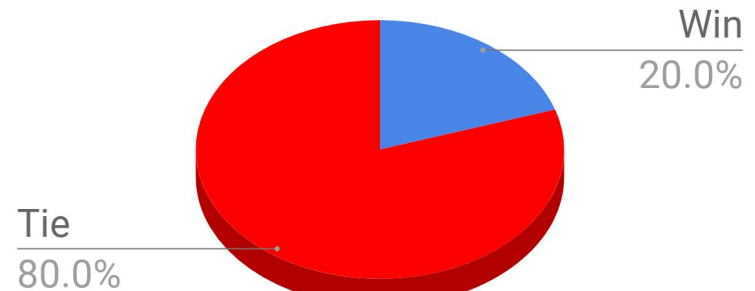
- **Playing against human players**

- To ensure that DQN trained players can really reach the goal of being unbeatable (without losing any game).
- 10 testers of ages above 15, each plays 2 rounds.
- DQN moves first: two-way winning reward **0.9**
- DQN moves second: **4** extra copies of determinative transitions

DQN moves first



DQN moves second





# Experiment 5: DQN Self Playing

<table><tr><td></td><td></td><td></td></tr><tr><td></td><td>X</td><td></td></tr><tr><td></td><td></td><td></td></tr></table>					X					<table><tr><td>O</td><td></td><td></td></tr><tr><td></td><td>X</td><td></td></tr><tr><td></td><td></td><td></td></tr></table>	O				X					<table><tr><td>O</td><td></td><td></td></tr><tr><td></td><td>X</td><td></td></tr><tr><td>X</td><td></td><td></td></tr></table>	O				X		X			<table><tr><td>O</td><td></td><td>O</td></tr><tr><td></td><td>X</td><td></td></tr><tr><td>X</td><td></td><td></td></tr></table>	O		O		X		X		
	X																																						
O																																							
	X																																						
O																																							
	X																																						
X																																							
O		O																																					
	X																																						
X																																							
<table><tr><td>O</td><td>X</td><td>O</td></tr><tr><td></td><td>X</td><td></td></tr><tr><td>X</td><td></td><td></td></tr></table>	O	X	O		X		X			<table><tr><td>O</td><td>X</td><td>O</td></tr><tr><td></td><td>X</td><td></td></tr><tr><td>X</td><td>O</td><td></td></tr></table>	O	X	O		X		X	O		<table><tr><td>O</td><td>X</td><td>O</td></tr><tr><td></td><td>X</td><td></td></tr><tr><td>X</td><td>O</td><td>X</td></tr></table>	O	X	O		X		X	O	X	<table><tr><td>O</td><td>X</td><td>O</td></tr><tr><td></td><td>X</td><td>O</td></tr><tr><td>X</td><td>O</td><td>X</td></tr></table>	O	X	O		X	O	X	O	X
O	X	O																																					
	X																																						
X																																							
O	X	O																																					
	X																																						
X	O																																						
O	X	O																																					
	X																																						
X	O	X																																					
O	X	O																																					
	X	O																																					
X	O	X																																					
<table><tr><td>O</td><td>X</td><td>O</td></tr><tr><td>X</td><td>X</td><td>O</td></tr><tr><td>X</td><td>O</td><td>X</td></tr></table>	O	X	O	X	X	O	X	O	X	<p>X: two-way winning reward <b>0.9</b></p> <p>O: <b>4</b> extra copies of determinative transit</p>																													
O	X	O																																					
X	X	O																																					
X	O	X																																					

# Conclusion

- **Summary**

- A DQN approach using MLP approximated Q function is proposed to train a computer to play Tic-Tac-Toe invincibly.
- Domain knowledges can make DQN utilize training data more efficiently, including
  - Adding extra copies of determinative training samples into Experience Replay
  - Early rewarding for some specific game states that can lead to winning

- **More Domain Knowledges**

- While training 2<sup>nd</sup> mover, give penalty when the 1<sup>st</sup> mover gets to the 2-way winning board.