



IMT Atlantique

Bretagne-Pays de la Loire

École Mines-Télécom

UE Introduction to AI

PROJECT

TEAM 2

Final restitution

Minh Triet Vo

Wenjia Fu

SUMMARY

1. Choice of strategy

2. Exploration vs Exploitation

2.1. Decayed epsilon greedy algorithm

2.2. Comparison

3. Configuration benchmarking

3.1. Model choice

3.2. Combination of RL and Combinatorial game theory



IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

RL + Game Theory !



1. Why not supervised Learning ?
2. Why not unsupervised Learning ?
3. RL only -> 53% winning rate?
4. RL + Game Theory

CHAPTER 2 : Exploration vs Exploitation

4

2.1 Epsilon-Greedy



Action at time(t)



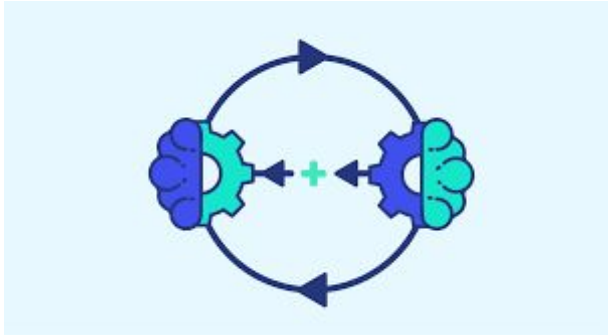
$\max Q_t(a)$

with probability $1-\epsilon$

any action (a)

with probability ϵ

2.2 Decayed Epsilon-Greedy



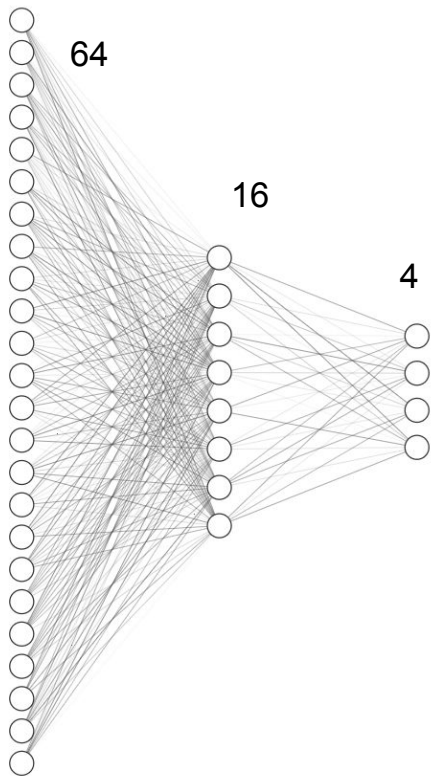
$$\epsilon = \epsilon_{end} + (\epsilon_{start} - \epsilon_{end}) * e^{-\frac{\epsilon_{step}}{\text{decay rate}}}$$

CHAPTER 2 : Decayed Epsilon-Greedy in RL

6

2.3 Reinforcement learning with vs without decayed epsilon greedy

Network configuration



Action choosing **without** epsilon greedy

```
with torch.no_grad():  
    Q_t = model(input_curr.unsqueeze(dim=0))[0]  
    action_t = torch.argmax(Q_t).item()
```

Action choosing **with** epsilon greedy

```
def eps_greedy_action(model, state, eps_start, eps_end, decay_rate, eps_step):  
    # This is the function for decayed-epsilon-greedy action-selection implementation  
    ## Formulate the decaying of epsilon  
    eps = eps_end + (eps_start - eps_end) * np.exp(-eps_step / decay_rate)  
    ## Action-selection based on epsilon value  
    if np.random.rand() > eps:  
        with torch.no_grad():  
            Q = model(state.unsqueeze(dim=0))[0]  
            ## Pick the next action that maximizes the Q value  
            action = torch.argmax(Q).item()  
    else:  
        ## Pick the next action randomly  
        action = np.random.choice(4)  
    return action  
action_t = eps_greedy_action(model, input_curr, eps_start, eps_end, decay_rate, eps_step)
```

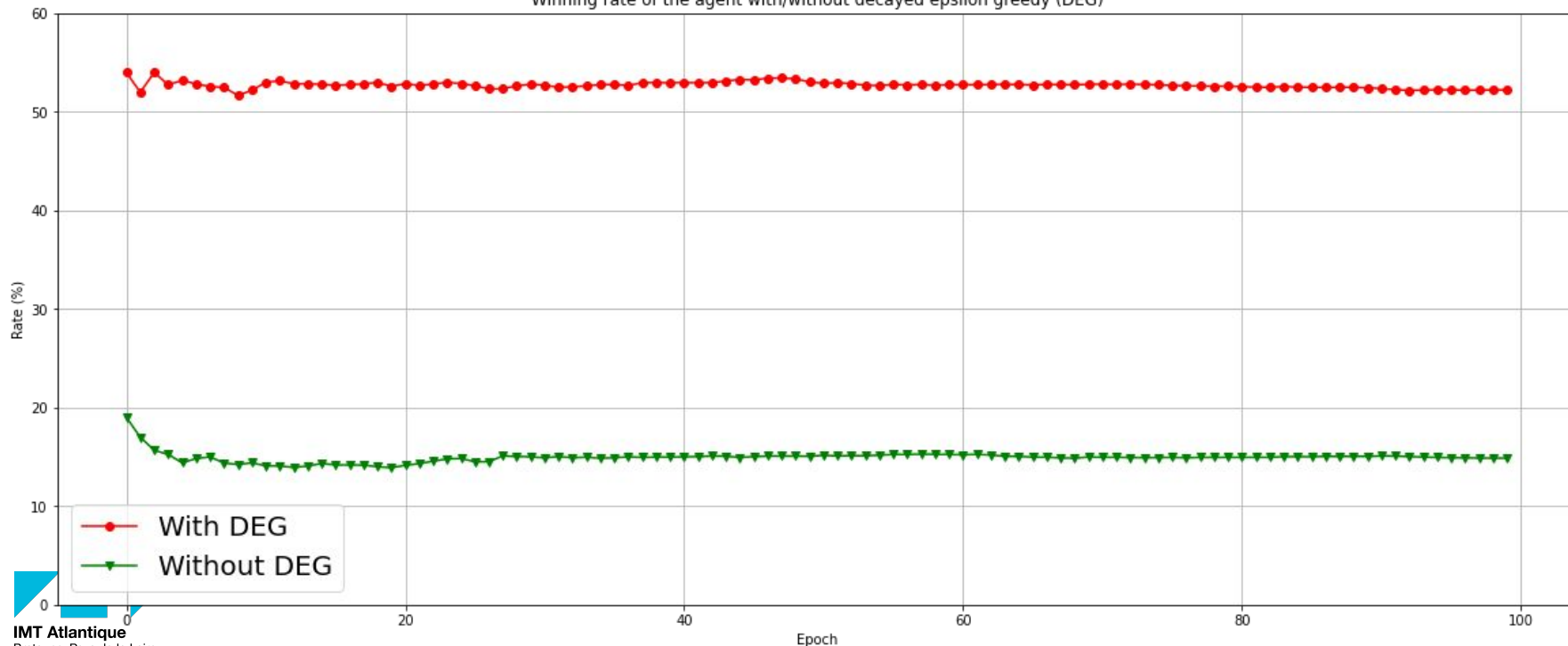
CHAPTER 2 : Decayed Epsilon-Greedy in RL

7

2.3 Reinforcement learning with vs without decayed epsilon greedy

Results

Winning rate of the agent with/without decayed epsilon greedy (DEG)

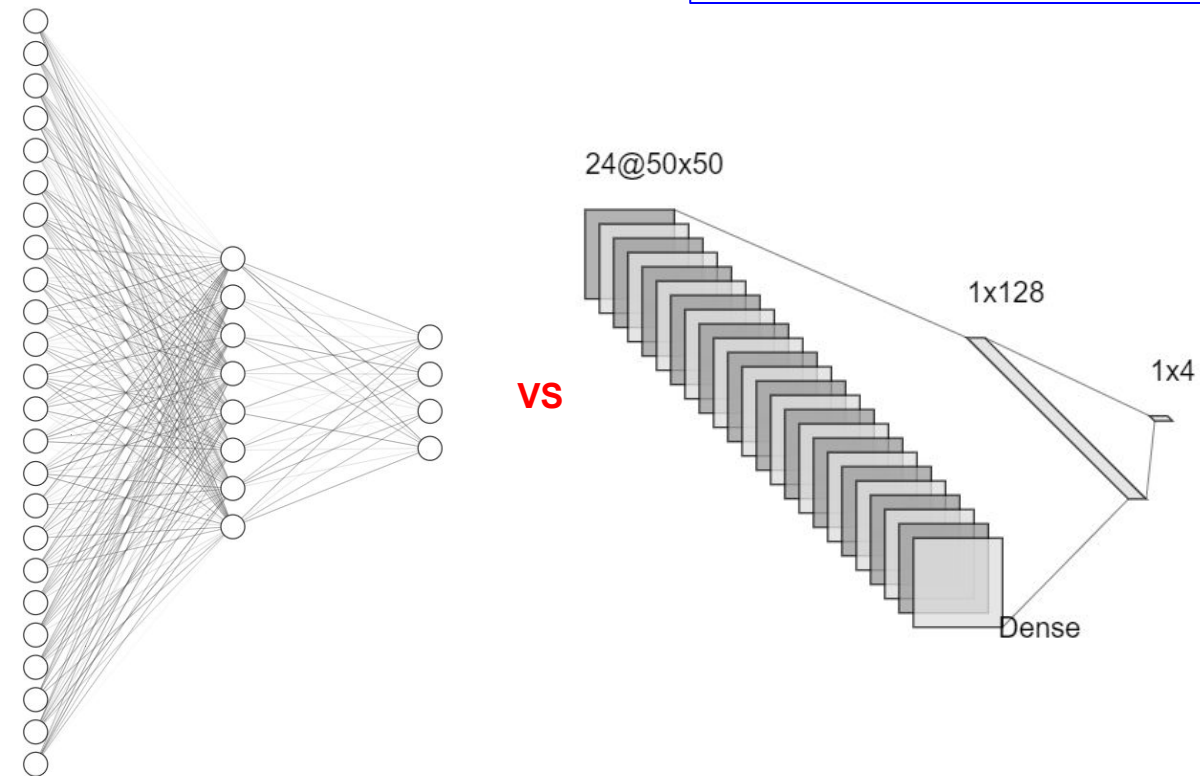


CHAPTER 3 : Configuration benchmarking

3.1 DQN network configurations

8

Network configuration

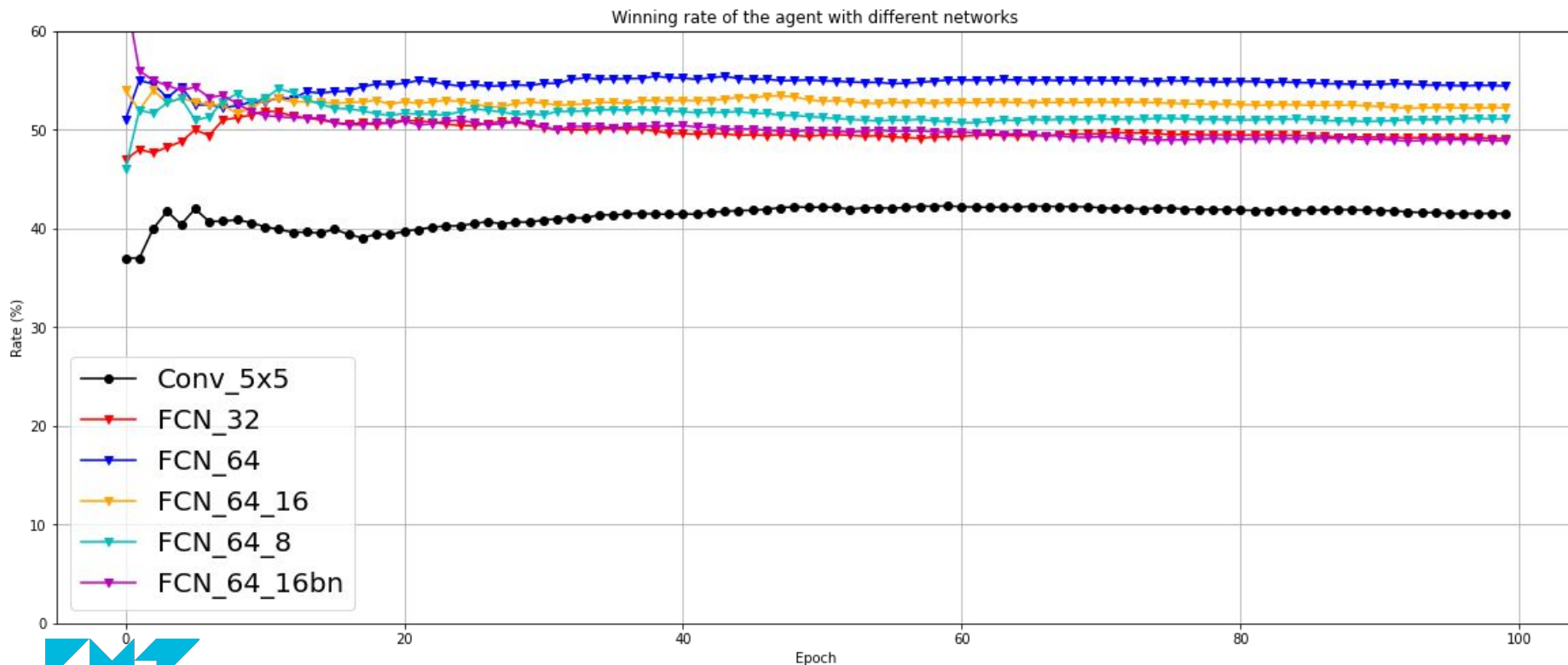


Name	Layers
Conv_5x5	2 conv (size=5, stride=1)
FCN_32	dense (32)
FCN_64	dense (64)
FCN_64_16	dense (64) → dense (16)
FCN_64_8	dense (64) → dense (8)
FCN_64_8bn	dense (64) + batchnorm → dense (8) + batchnorm

CHAPTER 3 : Configuration benchmarking

3.2 Comparison

9



CHAPTER 3 : Configuration benchmarking

10

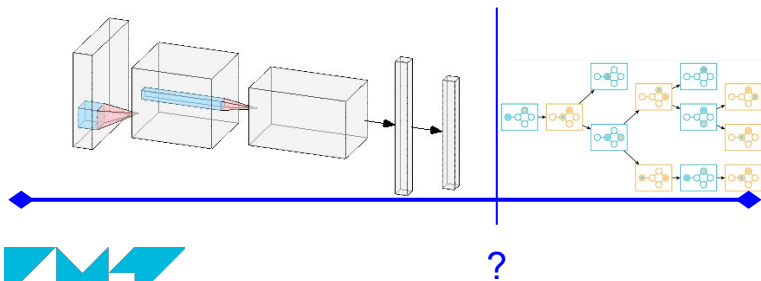
3.3 Combination of RL and game theory

How?

→ Put a threshold for the AI to choose its strategy.

Why?

- RL: Good but bounded performance.
- Game theory: unbounded performance but really slow.



Name	Winning rate (100 games)			
	7 cheese	10 cheese	12 cheese	13 cheese
Conv_5x5	53%	63%	78%	72%
FCN_32	71%	78%	73%	80%
FCN_64	57%	70%	79%	90%
FCN_64_16	63%	82%	91%	89%
FCN_64_8	74%	76%	77%	78%
FCN_64_8bn	64%	69%	78%	81%

1. <https://www.geeksforgeeks.org/epsilon-greedy-algorithm-in-reinforcement-learning/>
2. <https://aakash94.github.io/Reward-Based-Epsilon-Decay/>

Thank you!



IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom