

Map

Map

Model
Based



Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy



Off Policy

ML MB TRL
(learn T, R)

Q-learning

Map

Model
Based



Model
Free



learn Q
SARSA

learn π
Policy
Gradient

On Policy



Off Policy

ML MB TRL
(learn T, R)

Q-learning

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Map

Model
Based



Model
Free



learn Q
SARSA

learn π
Policy
Gradient

On Policy



Off Policy

ML MB TRL
(learn T, R)

Q-learning

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

Map

Model
Based

Model
Free

learn Q
SARSA

learn π
Policy
Gradient

On Policy

Off Policy

ML MB TRL
(learn T, R)

Q-learning

Part I

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

Map

Model
Based



Model
Free



learn Q
SARSA

learn π
Policy Gradient
Part 2

On Policy

Off Policy

ML MB TRL
(learn T, R)

Q-learning

Part 1

Challenges:

1. Exploration vs Exploitation
2. Credit Assignment
3. Generalization

Last Time: Neural Networks

Part 1

DQN

Discussion

Discussion

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Discussion

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Discussion

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate Q with Q_{θ}

Discussion

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate Q with Q_{θ}
- What should (x, y) be?

Discussion

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x, y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate Q with Q_{θ}
- What should (x, y) be?
- What should l be?

(s, a, r, s') *experience tuples*

$$l(s, a, r, s') = (r + \gamma \max_{a'} Q(s', a') - Q(s, a))^2$$

Discussion

Q-Learning:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Neural Networks

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \in \mathcal{D}} l(f_{\theta}(x), y)$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} l(f_{\theta}(x), y)$$

Deep Q learning:

- Approximate Q with Q_{θ}
- What should (x, y) be?
- What should l be?

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\underline{\theta} \leftarrow \underline{\theta} - \underline{\alpha} \nabla_{\theta} \| r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a) \|_2$$

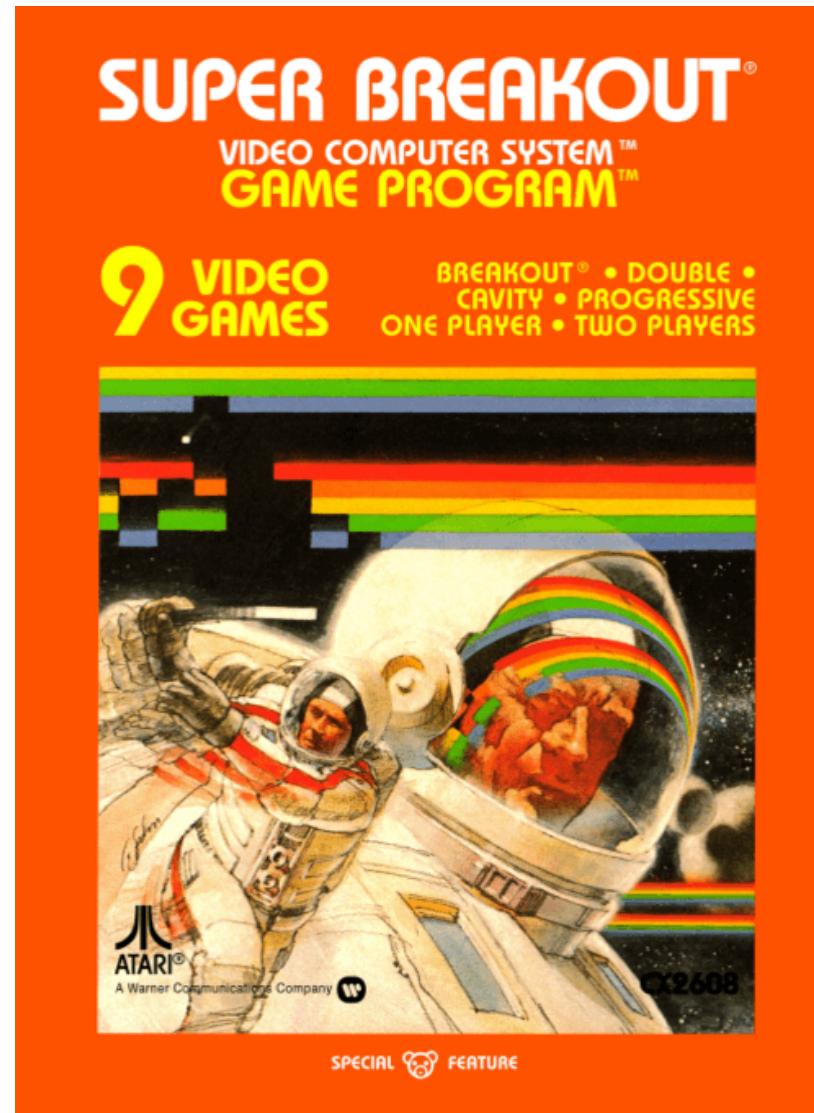
$$s \leftarrow s'$$

DQN: The Atari Benchmark

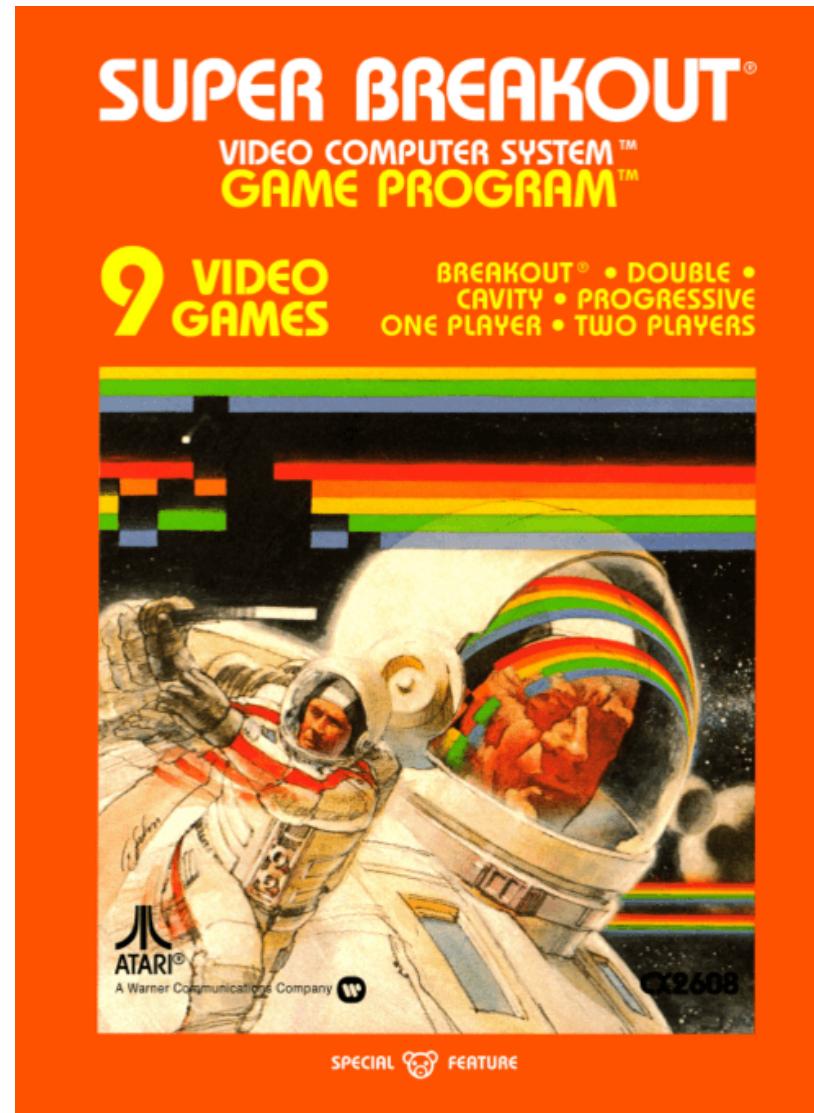
DQN: The Atari Benchmark



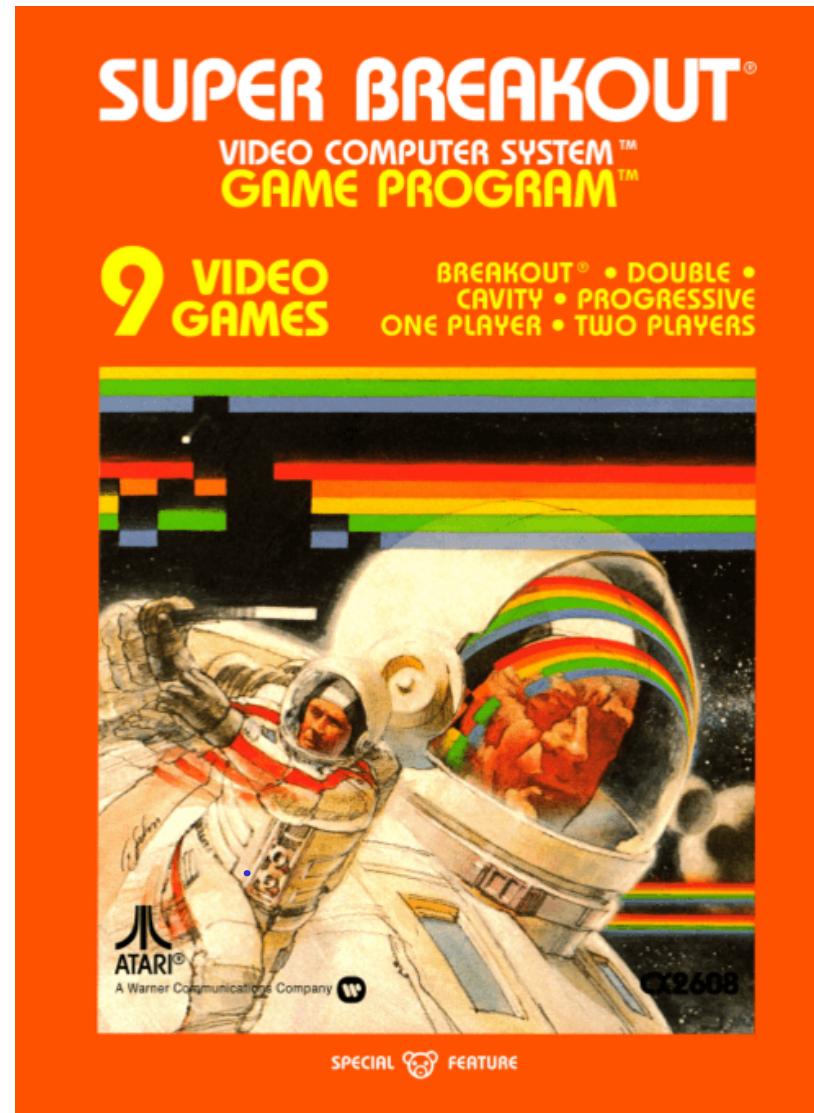
DQN: The Atari Benchmark



DQN: The Atari Benchmark



DQN: The Atari Benchmark



DQN: Problems with Naive Approach

Problems:

- } data buffer/experience replay
- } periodically freeze target

DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \| r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a) \|_2$$

$$s \leftarrow s'$$

Problems:

} data buffer/experience replay
} periodically freeze target

DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \| r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a) \|_2$$

$$s \leftarrow s'$$

Problems:

1. Samples Highly Correlated

} data buffer/experience replay
} periodically freeze target

DQN: Problems with Naive Approach

Candidate Algorithm:

loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \| r + \gamma \max_{a'} Q_{\theta}(s', a') - Q_{\theta}(s, a) \|_2$$

$$s \leftarrow s'$$

Problems:

1. Samples Highly Correlated
2. Size-1 batches

} data buffer/experience replay
} periodically freeze target

DQN: Problems with Naive Approach

Candidate Algorithm:

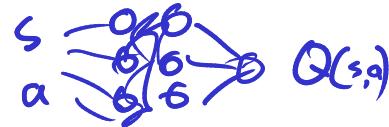
loop

$$a \leftarrow \operatorname{argmax} Q(s, a) \text{ w.p. } 1 - \epsilon, \quad \operatorname{rand}(A) \text{ o.w.}$$

$$r \leftarrow \operatorname{act}!(\text{env}, a)$$

$$s' \leftarrow \operatorname{observe}(\text{env})$$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \| r + \gamma \max_{a'} \underbrace{Q_{\theta}(s', a') - Q_{\theta}(s, a)}_{\text{target network}} \|_2$$



Problems:

1. Samples Highly Correlated
2. Size-1 batches
3. Moving target

} data buffer/experience replay
} periodically freeze target

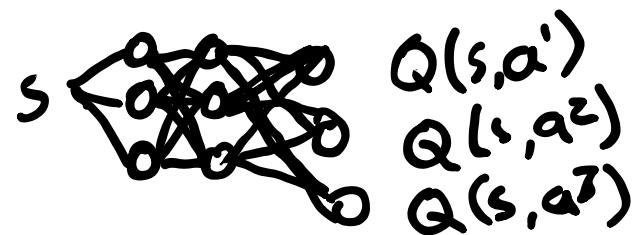
DQN

DQN

Q Network Structure:

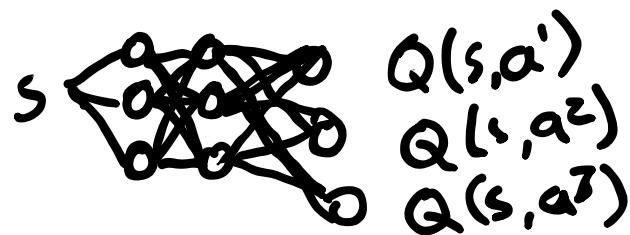
DQN

Q Network Structure:



DQN

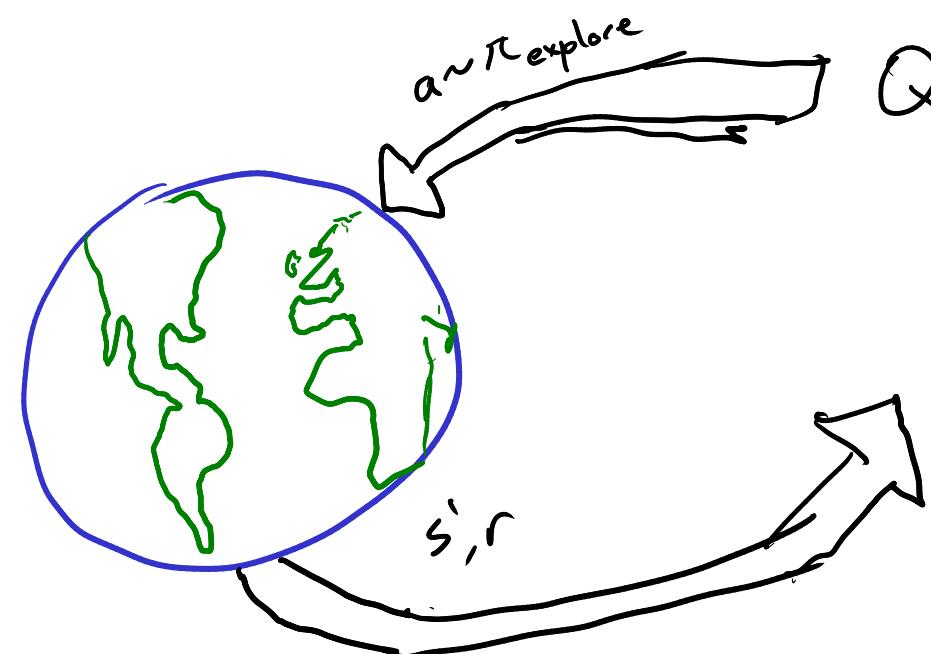
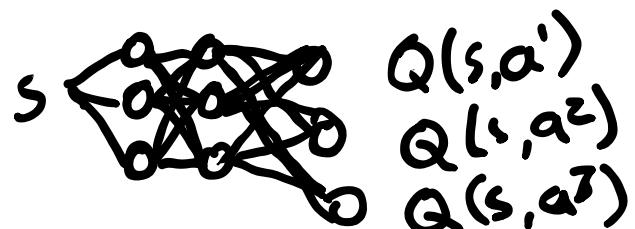
Q Network Structure:



Experience Tuple: (s, a, r, s')

DQN

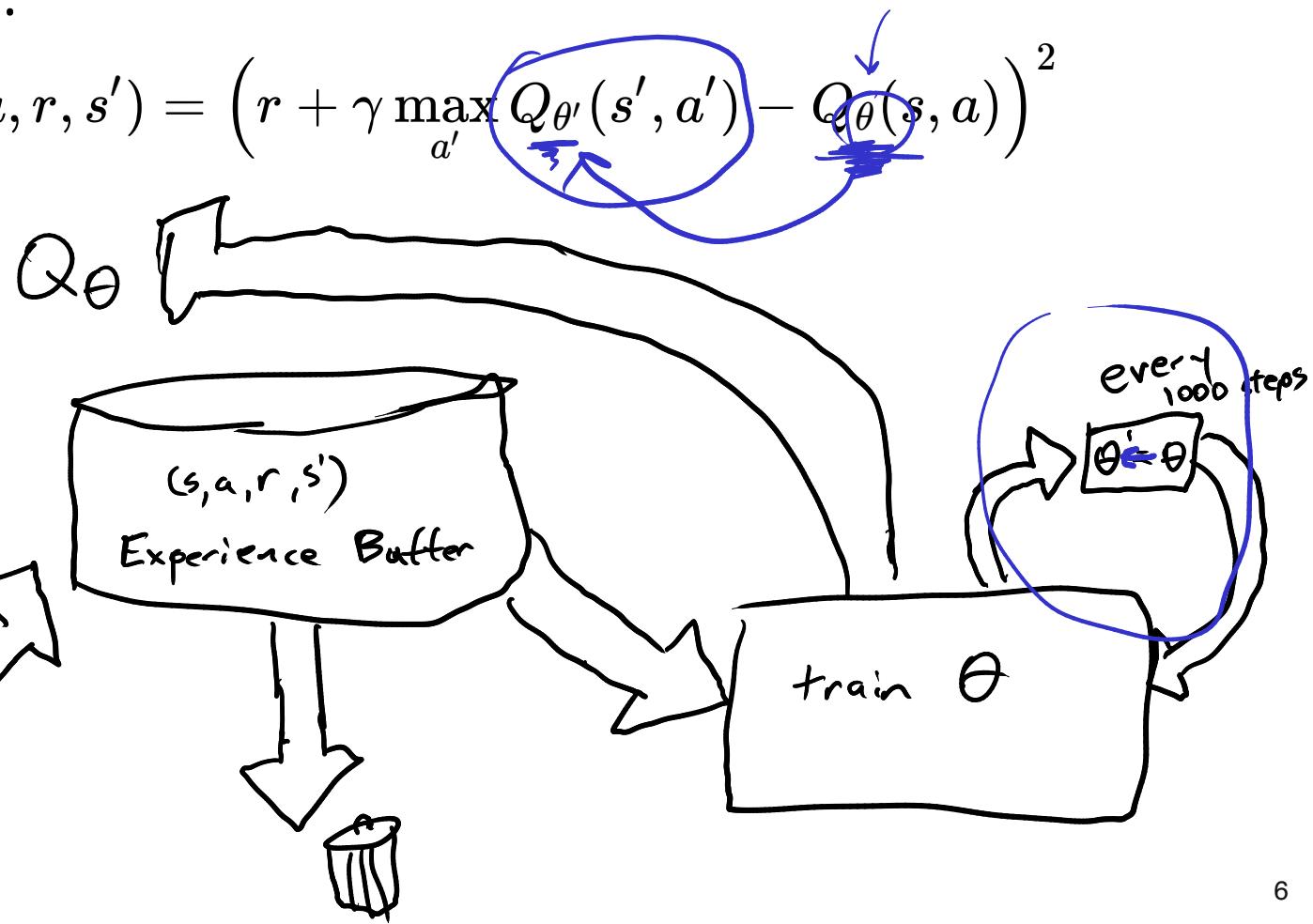
Q Network Structure:



Experience Tuple: (s, a, r, s')

Loss:

$$l(s, a, r, s') = \left(r + \gamma \max_{a'} Q_{\theta'}(s', a') - Q_\theta(s, a) \right)^2$$



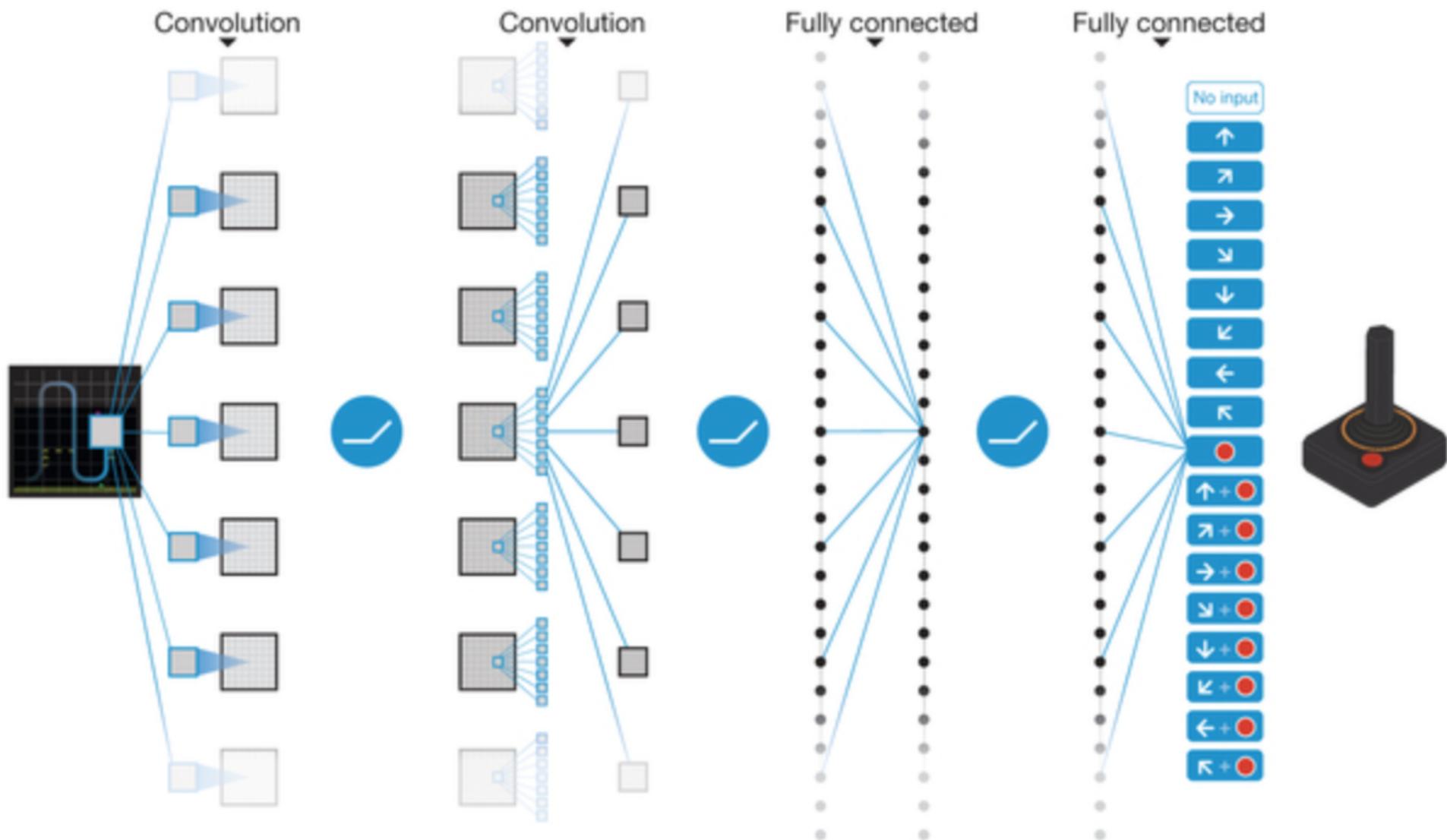
Playing Atari with Deep Reinforcement Learning

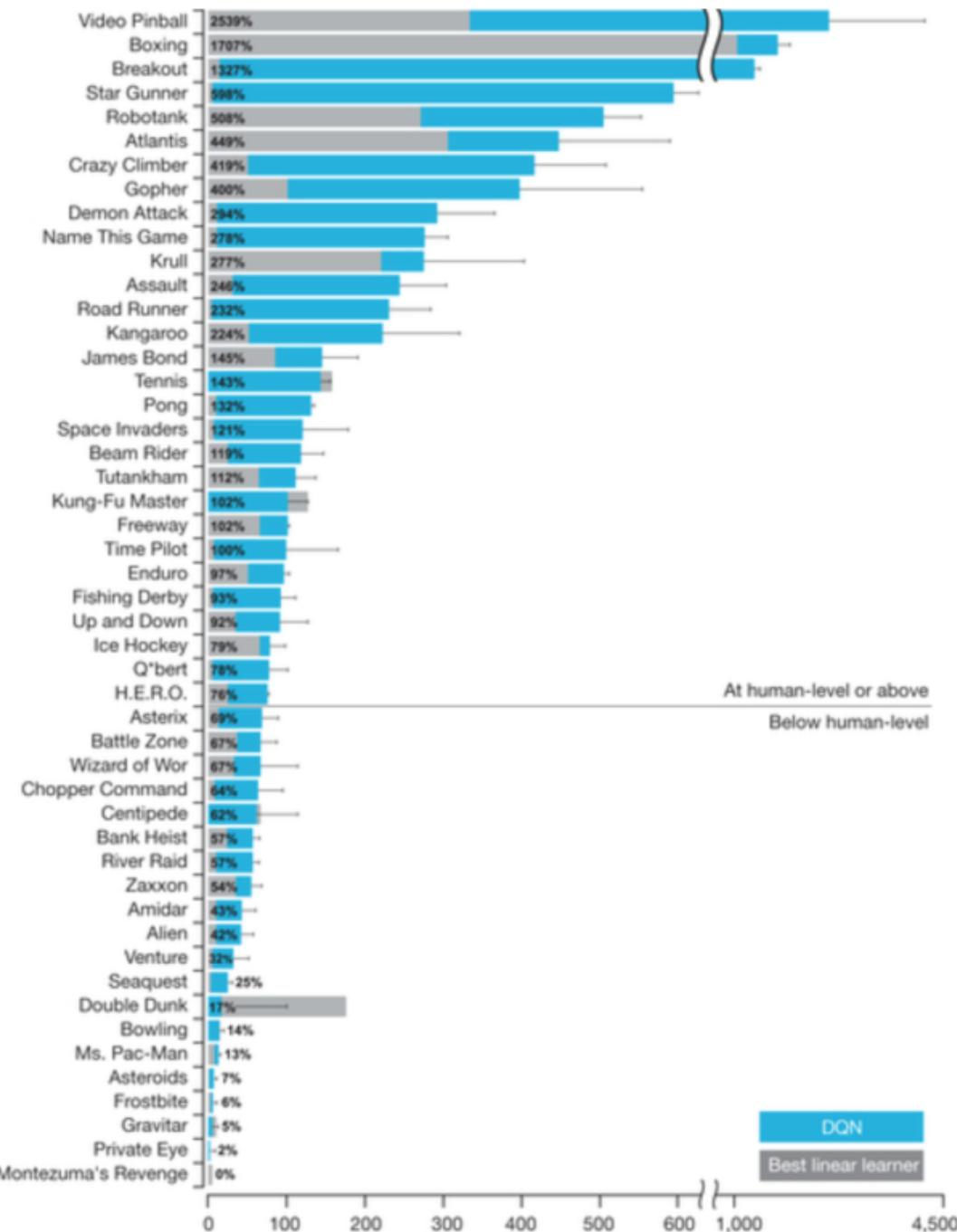
Volodymyr Mnih Koray Kavukcuoglu David Silver Alex Graves Ioannis Antonoglou

Daan Wierstra Martin Riedmiller

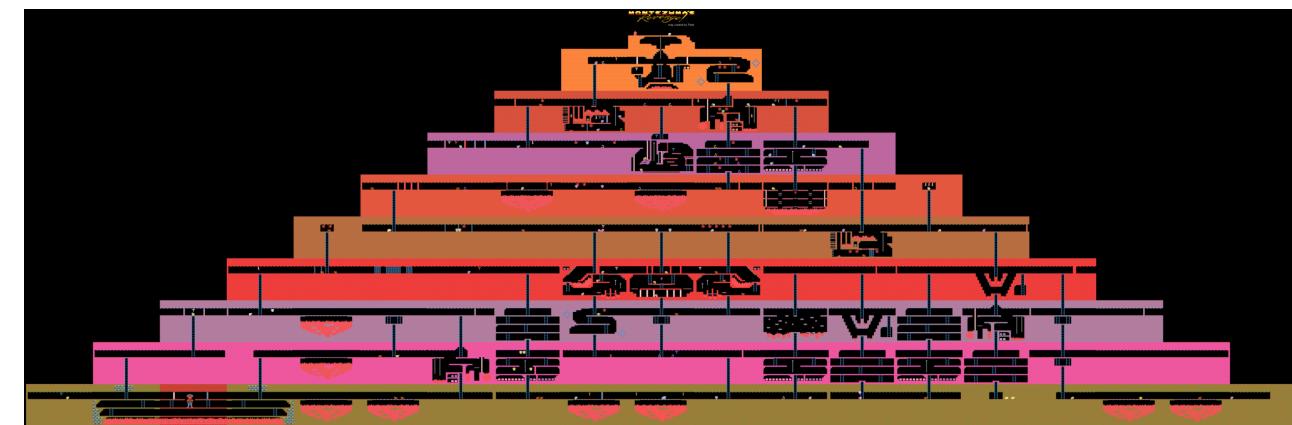
DeepMind Technologies







[https://www.youtube.com/watch?
v=SuZVyOlgVek](https://www.youtube.com/watch?v=SuZVyOlgVek)



Rainbow

Rainbow

- Double Q Learning

Rainbow

- Double Q Learning
- Prioritized Replay
 - (priority proportional to last TD error)