

CS4400

DEEP REINFORCEMENT LEARNING

Lecture 5: Advanced DQN

Wendelin Böhmer

`<j.w.bohmer@tudelft.nl>`



30th of November 2023

Content of this lecture



5.1 Q-learning extensions

5.2 Partial observability

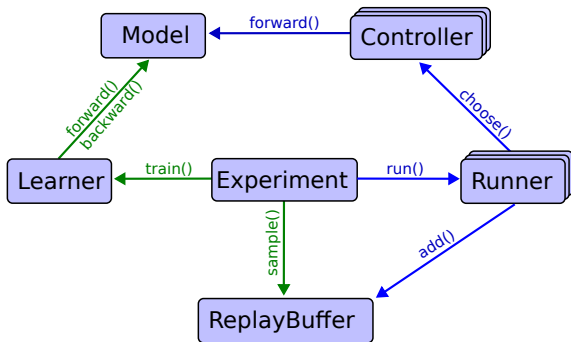
5.1

Advanced DQN Q-learning extensions

5.1 Synchronous software architecture



- Alternate `run()` and `train()`
 - Learner on GPU, Runner on CPU/GPU
 - large `ReplayBuffer` on CPU
 - many parallel `Runner` possible



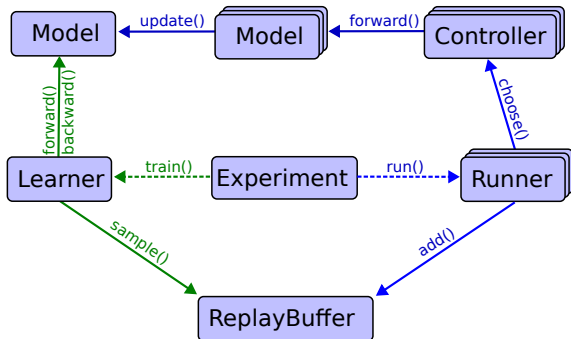
assignment sheet 2

see [Daley and Amato \(2021\)](#) for state-of-the-art DQN architectures

5.1 Asynchronous software architecture



- Endless `run()` and `train()` in parallel threads
 - Controller regularly update copy of Model
 - Runner can run on different servers
 - multiple Learner regularly average Model



see [Daley and Amato \(2021\)](#) for state-of-the-art DQN architectures

5.1 Double Q-learning



- Uncertain bootstrapping *overestimates* values
 - amplifies over multiple time-steps

$$\mathbb{E}\left[\max_a Q(s, a)\right] \geq \max_a \mathbb{E}[Q(s, a)]$$



5.1 Double Q-learning



- Uncertain bootstrapping *overestimates* values
 - amplifies over multiple time-steps

$$\mathbb{E}\left[\max_a Q(s, a)\right] \geq \max_a \mathbb{E}[Q(s, a)]$$

- Double Q-learning: use independent estimate Q'
 - Q and Q' are alternately trained
 - requires twice as many transitions

$$\mathbb{E}\left[Q\left(s, \arg \max_a Q'(s, a)\right)\right] \approx \max_a \mathbb{E}[Q(s, a)]$$



- Uncertain bootstrapping *overestimates* values
 - amplifies over multiple time-steps

$$\mathbb{E}\left[\max_a Q(s, a)\right] \geq \max_a \mathbb{E}[Q(s, a)]$$

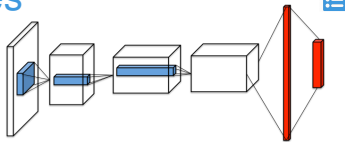
- Double Q-learning: use independent estimate Q'
 - Q and Q' are alternately trained
 - requires twice as many transitions

$$\mathbb{E}\left[Q\left(s, \arg \max_a Q'(s, a)\right)\right] \approx \max_a \mathbb{E}[Q(s, a)]$$

- Deep Double Q-learning (DDQN)
 - target network is Q and current network is Q'
 - not uncorrelated, but works in practice



5.1 Dueling network architectures



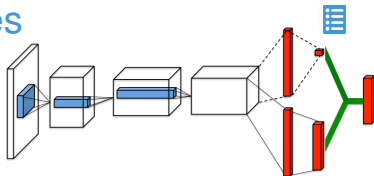
- Estimating $Q(s, a)$ for many actions
 - only one head at a time updated
 - update all actions with same successor state s' ?
- $Q^\pi(s, a) =: V^\pi(s) + A^\pi(s, a)$
 - advantage of action $A^\pi(s, a)$
 - $Q^*(s, a^*) = V^*(s) \Rightarrow A^*(s, a^*) = 0$
 - estimating $A^* + V^*$ easier than Q^* ?

images and dueling architectures from [Wang et al. \(2016\)](#)

5.1 Dueling network architectures

- Estimating $Q(s, a)$ for many actions

- only one head at a time updated
- update all actions with same successor state s' ?

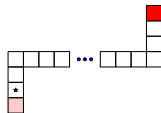


- $Q^\pi(s, a) =: V^\pi(s) + A^\pi(s, a)$
 - advantage of action $A^\pi(s, a)$
 - $Q^*(s, a^*) = V^*(s) \Rightarrow A^*(s, a^*) = 0$
 - estimating $A^* + V^*$ easier than Q^* ?

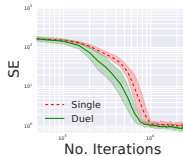
- Dueling architectures: $Q(s, a) :=$

- $V(s) + \left(A(s, a) - \max_{a \in \mathcal{A}} A(s, a) \right)$
- $V(s) + \left(A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} A(s, a) \right)$

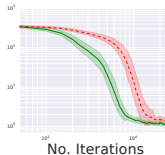
CORRIDOR ENV



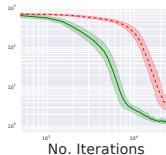
5 ACTIONS



10 ACTIONS



20 ACTIONS



images and dueling architectures from [Wang et al. \(2016\)](#)

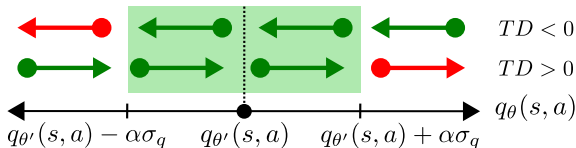
5.1 Trust-region online bootstrapping



- Target networks θ' stabilize, but delay value propagation
- Bootstrapping with (detached) *online network* θ must be stabilized
 - mask out TD losses when online and TD error very different, when:

$$\begin{aligned}
 |q_\theta(s_t, a_t) - q_{\theta'}(s_t, a_t)| &> \alpha \sqrt{\mathbb{V}[q_\theta(s, a) | (s, a) \sim \mathcal{D}]}^{\sigma_q} \\
 \text{sgn}\left(q_\theta(s_t, a_t) - q_{\theta'}(s_t, a_t)\right) &\neq \text{sgn}\left(\underbrace{q_\theta(s_t, a_t) - \left(r_t + \gamma \max_a q_\theta(s_{t+1}, a)\right)}_{TD_t}\right)
 \end{aligned}$$

- only gradient steps towards trust-region allowed



method and figure replicated from [Badia et al. \(2020\)](#)

5.1 Distributional Q-learning

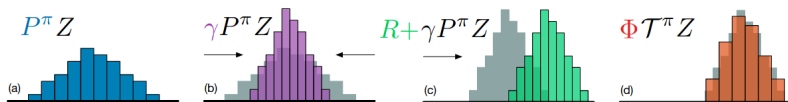


- Expectations over returns not always enough
 - e.g. for risk-sensitive or constrained RL

5.1 Distributional Q-learning



- Expectations over returns not always enough
 - e.g. for risk-sensitive or constrained RL
- Estimate distribution over possible returns $Z_t = \sum_{i=0}^{\infty} \gamma^i R_{t+i} \in \mathbb{R}$
 - $\mathbb{P}(Z_t | s_t, a_t) \stackrel{D}{=} \mathbb{P}(R_t | s_t, a_t) + \gamma \mathbb{P}(Z_{t+1} | s_{t+1}, a_{t+1})$

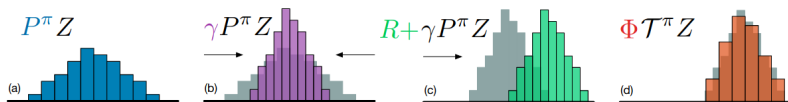


originally by [Bellemare et al. \(2017\)](#), see [Lyle et al. \(2019\)](#) for a recent comparison with DQN

5.1 Distributional Q-learning



- Expectations over returns not always enough
 - e.g. for risk-sensitive or constrained RL
- Estimate distribution over possible returns $Z_t = \sum_{i=0}^{\infty} \gamma^i R_{t+i} \in \mathbb{R}$
 - $\mathbb{P}(Z_t | s_t, a_t) \stackrel{D}{=} \mathbb{P}(R_t | s_t, a_t) + \gamma \mathbb{P}(Z_{t+1} | s_{t+1}, a_{t+1})$



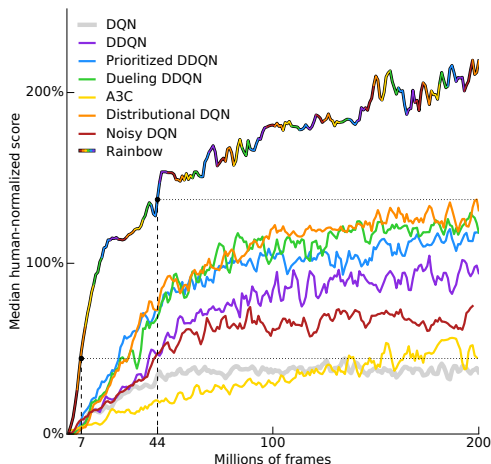
- Many ways to model the return-distribution
 - categorical distribution (Bellemare et al., 2017)
 - quantile regression (Dabney et al., 2018)
 - moment matching (Nguyen-Tang et al., 2021)

originally by [Bellemare et al. \(2017\)](#), see [Lyle et al. \(2019\)](#) for a recent comparison with DQN

5.1 Modern DQN versions



- Combine DQN extensions
- e.g. RAINBOW (2018)
 - double Q-learning
 - prioritized replay
 - dueling architectures
 - distributional DQN
- e.g. Agent57 (2020)
 - trust-region
 - short-term memory
 - exploration
 - episodic memory
 - meta-controllers



see [Deepmind's blog](#) about Agent57 for a good overview

RAINBOW and image by [Hessel et al. \(2018\)](#), Agent57 by [Badia et al. \(2020\)](#), newest results in [Kapturowski et al. \(2023\)](#)

- Standard architecture can be easily parallelized
- Double Q-learning reduces overestimation
- Dueling networks for many similar actions
- Trust-regions for faster value propagation
- Distributional Q-learning to estimate risk
- Modern DQN versions combine tricks to learn superhuman AI

Learning Objectives

LO5.1: Explain when the discussed extensions are beneficial

LO5.2: Implement and test double Q-learning

5.2

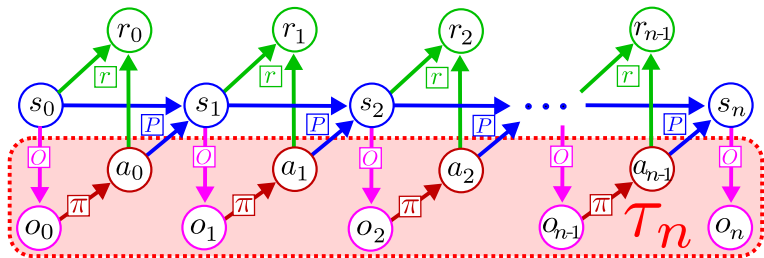
Advanced DQN

Partial observability

5.2 Partially observable MDP



- POMDP: MDP + observations $o_t \sim O(\cdot|s_t) \in \mathcal{O}$
 - true state not observed by agent
 - histories $\tau_t := (o_0, a_0, \dots, o_t) \in (\mathcal{O} \times \mathcal{A})^t \times \mathcal{O}$

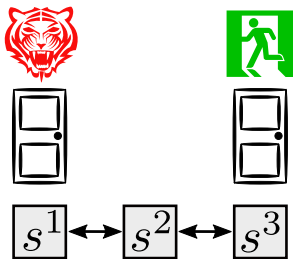


see e.g. [Spaan \(2012\)](#) for an introduction to POMDP

5.2 Partially observable MDP



- POMDP: MDP + observations $o_t \sim O(\cdot|s_t) \in \mathcal{O}$
 - true state not observed by agent
 - histories $\tau_t := (o_0, a_0, \dots, o_t) \in (\mathcal{O} \times \mathcal{A})^t \times \mathcal{O}$
 - example: escape tiger by only observing doors

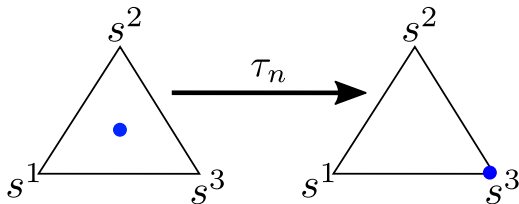
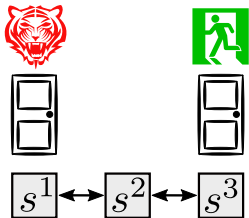


see e.g. [Spaan \(2012\)](#) for an introduction to POMDP

5.2 Partially observable MDP



- POMDP: MDP + observations $o_t \sim O(\cdot|s_t) \in \mathcal{O}$
 - true state not observed by agent
 - histories $\tau_t := (o_0, a_0, \dots, o_t) \in (\mathcal{O} \times \mathcal{A})^t \times \mathcal{O}$
 - example: escape tiger by only observing doors
- Equivalent deterministic MDP over belief-distributions b
 - $b(s|\tau_t) := \mathbb{P}(s_t = s|b_0, o_0, a_0, \dots, o_t)$
 - $b(s|\tau_{t+1}) \propto O(o_{t+1}|s) \int P(s|s', a_t) b(s'|\tau_t) ds'$



see e.g. [Spaan \(2012\)](#) for an introduction to POMDP

- POMDP: MDP + observations $o_t \sim O(\cdot|s_t) \in \mathcal{O}$
 - true state not observed by agent
 - histories $\tau_t := (o_0, a_0, \dots, o_t) \in (\mathcal{O} \times \mathcal{A})^t \times \mathcal{O}$
 - example: escape tiger by only observing doors
- Equivalent deterministic MDP over belief-distributions b
 - $b(s|\tau_t) := \mathbb{P}(s_t = s | b_0, o_0, a_0, \dots, o_t)$
 - $b(s|\tau_{t+1}) \propto O(o_{t+1}|s) \int P(s|s', a_t) b(s'|\tau_t) ds'$
- History τ_t is *sufficient statistic* of belief $b(s|\tau_t)$
 - $r(\tau_t, a) = \int r(s, a) b(s|\tau_t) ds$
 - $P(\tau_{t+1}|\tau_t, a_t) = \iint O(o_{t+1}|s') P(s'|s, a_t) b(s|\tau_t) ds ds'$
 - $Q^*(\tau_t, a) = r(\tau_t, a) + \gamma \int P(\tau_{t+1}|\tau_t, a) \max_{a' \in \mathcal{A}} Q^*(\tau_{t+1}, a') d\tau_{t+1}$

see e.g. [Spaan \(2012\)](#) for an introduction to POMDP



- Use LSTM to encode the history $\tau_t = (o_0, a_0, \dots, o_t)$
 - $\mathbf{h}_t = \text{LSTM}(\mathbf{h}_{t-1}, \underbrace{a_{t-1}, o_t}_{\text{input}}) \in \mathbb{R}^m$, $\mathbf{h}_0 = \text{LSTM}(\underbrace{\mathbf{0}, 0}_{\text{init}}, o_0)$
- End-to-end Q-learning of histories

$$\mathcal{L}[\theta] = \mathbb{E} \left[\sum_{t=0}^n \left(r_t + \gamma \max_{a \in \mathcal{A}} q_{\theta'}(\mathbf{h}_{t+1}, a) - q_{\theta}(\mathbf{h}_t, a_t) \right)^2 \mid \tau_n \sim \mathcal{D} \right]$$

- How does this change the replay buffer?



- Use LSTM to encode the history $\tau_t = (o_0, a_0, \dots, o_t)$
 - $\mathbf{h}_t = \text{LSTM}(\mathbf{h}_{t-1}, \underbrace{a_{t-1}, o_t}_{\text{input}}) \in \mathbb{R}^m$, $\mathbf{h}_0 = \text{LSTM}(\underbrace{\mathbf{0}, 0}_{\text{init}}, o_0)$
- End-to-end Q-learning of histories

$$\mathcal{L}[\theta] = \mathbb{E} \left[\sum_{t=0}^n \left(r_t + \gamma \max_{a \in \mathcal{A}} q_{\theta'}(\mathbf{h}_{t+1}, a) - q_{\theta}(\mathbf{h}_t, a_t) \right)^2 \mid \tau_n \sim \mathcal{D} \right]$$

- How does this change the replay buffer?
 - sample mini-batch of episodes, not transitions
 - contiguous tensor stores episodes of varying length
 - zero-out non-existent losses afterwards

5.2 DRQN implementation (1)



- We start with a normal DQN gradient descent

```
1 import torch as th
2 batch = self.replay_buffer.sample()
3
4
5
6
7
8
9
10 targets = ???
11
12 current_v = ???
13 loss = mse_loss(current_v, targets.detach())
14 # gradient descent step on supervised regression loss
15 optimizer.zero_grad()
16 (loss).backward()
17 optimizer.step()
```


5.2 DRQN implementation (2)



- Replay buffer returns batch of episodes (e.g. time-dim=0)
- Tensor slice `[:-1]` and `[1:]` for 'current' and 'next'

```
1 import torch as th
2 batch = self.replay_buffer.sample()
3 acts = batch['actions']
4
5
6
7 values = ???
8 # derive the regression targets from the values
9 max_v = ~batch['terminals'] * values.max(dim=-1, keepdim=True)[0]
10 targets = (batch['rewards'][:-1] + gamma * max_v[1:])
11 # derive the MSE loss to be optimized
12 current_v = values[:-1].gather(dim=-1, index=acts[:-1])
13 loss = mse_loss(current_v, targets.detach())
14 # gradient descent step on supervised regression loss
15 optimizer.zero_grad()
16 (loss).backward()
17 optimizer.step()
```

5.2 DRQN implementation (3)



- LSTM Q-value function q on observations o_t and actions a_{t-1}
- $h_t = \text{LSTM}(h_{t-1}, a_{t-1}, o_t)$ starts with o_0 and $a_{-1} = 0$

```
1 import torch as th
2 batch = self.replay_buffer.sample()
3 acts = batch['actions']
4 # compute all recurrent Q-values with function q
5 first_act = th.zeros(1, *acts.shape[1:]) # a_{-1} = 0
6 delayed_acts = th.cat([first_act, acts[:-1]], dim=0)
7 values = q(th.cat([delayed_acts, batch['observations']], dim=-1))[0]
8 # derive the regression targets from the values
9 max_v = ~batch['terminals'] * values.max(dim=-1, keepdim=True)[0]
10 targets = (batch['rewards'][:-1] + gamma * max_v[1:])
11 # derive the MSE loss to be optimized
12 current_v = values[:-1].gather(dim=-1, index=acts[:-1])
13 loss = mse_loss(current_v, targets.detach())
14 # gradient descent step on supervised regression loss
15 optimizer.zero_grad()
16 (loss).backward()
17 optimizer.step()
```

5.2 DRQN implementation (4)



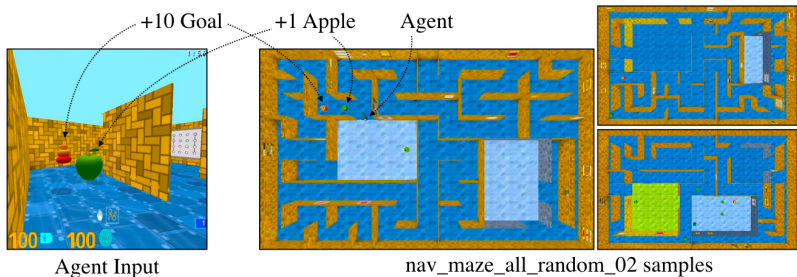
- Multiplying with a `mask` tensor removes 'invalid' entries
- Normalize loss by number of 'valid' TD-errors
- Can be implemented with any DQN extension

```
1 import torch as th
2 batch = self.replay_buffer.sample()
3 acts, mask = batch['actions'], batch['mask'][:-1]
4 # compute all recurrent Q-values with function q
5 first_act = th.zeros(1, *acts.shape[1:]) #  $a_{-1} = 0$ 
6 delayed_acts = th.cat([first_act, acts[:-1]], dim=0)
7 values = q(th.cat([delayed_acts, batch['observations']], dim=-1))[0]
8 # derive the regression targets from the values
9 max_v = ~batch['terminals'] * values.max(dim=-1, keepdim=True)[0]
10 targets = (batch['rewards'][:-1] + gamma * max_v[1:]) * mask
11 # derive the MSE loss to be optimized
12 current_v = values[:-1].gather(dim=-1, index=acts[:-1]) * mask
13 loss = mse_loss(current_v, targets.detach(), reduction='sum')
14 # gradient descent step on supervised regression loss
15 optimizer.zero_grad()
16 (loss / mask.sum()).backward()
17 optimizer.step()
```



Question: Design a neural architecture

- Design a neural network architecture for the shown task
 - agent observes (ambiguous) first-person images
 - actions $a \in \mathbb{R}^2$: move forward by a_1 , rotate by a_2
 - collecting apples and reaching goal rewarded
- Choose a fitting neural architecture for the PPO algorithm
 - sketch outputs and which modules feed into each other





Question: Design a neural architecture

- Design a neural network architecture for the shown task
 - agent observes (ambiguous) first-person images
 - actions $\mathbf{a} \in \mathbb{R}^2$: move forward by a_1 , rotate by a_2
 - collecting apples and reaching goal rewarded
- Choose a fitting neural architecture for the PPO algorithm
 - sketch outputs and which modules feed into each other

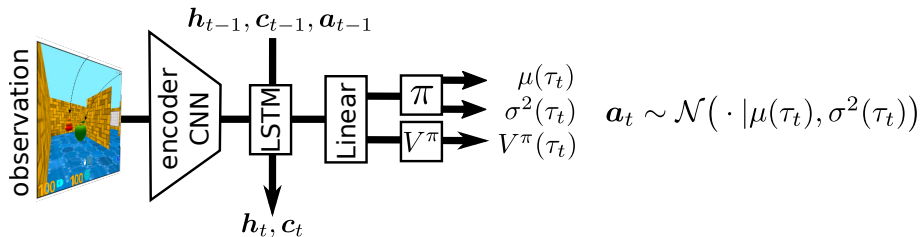


image from the Deepmind Control Suite ([Jaderberg et al., 2017](#))

- POMDP induce deterministic Belief-MDP
- Histories are a sufficient statistic for the belief
- DRQN encodes histories with LSTMs
- LSTM, CNN and other modules can be combined

Learning Objectives

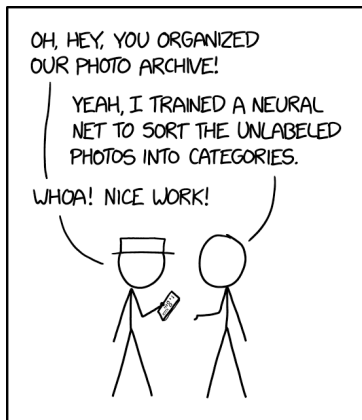
LO5.3: Explain why deep recurrent Q-networks solve POMDP

LO5.4: Design neural architectures for POMDP algorithms

5.2 Next lecture



- Next lecture:
continuous actions!
- Don't forget assignment sheet 2
- Questions? Ask them here:
answers.ewi.tudelft.nl



ENGINEERING TIP:
WHEN YOU DO A TASK BY HAND,
YOU CAN TECHNICALLY SAY YOU
TRAINED A NEURAL NET TO DO IT.

image source: xkcd.com

- Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human benchmark. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020*, volume 119 of *Proceedings of Machine Learning Research*, pages 507–517, 2020. URL <https://arxiv.org/abs/2003.13350>.
- Marc G. Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 449–458, 2017. URL <https://proceedings.mlr.press/v70/bellemare17a.html>.
- Will Dabney, Georg Ostrovski, David Silver, and Remi Munos. Implicit quantile networks for distributional reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1096–1105. PMLR, 2018. URL <https://proceedings.mlr.press/v80/dabney18a.html>.
- Brett Daley and Christopher Amato. Human-level control without server-grade hardware, 2021. URL <https://arxiv.org/abs/2111.01264>.
- Matthew J. Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable mdps. In *2015 AAAI Fall Symposia*, pages 29–37, 2015. URL <http://www.aaai.org/ocs/index.php/FSS/FSS15/paper/view/11673>.
- Matteo Hessel, Joseph Modayil, Hado van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Gheshlaghi Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI*, pages 3215–3222, 2018. URL <https://arxiv.org/abs/1710.02298>.
- Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z. Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. In *5th International Conference on Learning Representations, ICLR*. OpenReview.net, 2017. URL <https://arxiv.org/abs/1611.05397>.
- Steven Kapturowski, Víctor Campos, Ray Jiang, Nemanja Rakicevic, Hado van Hasselt, Charles Blundell, and Adria Puigdomenech Badia. Human-level atari 200x faster. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://arxiv.org/abs/2209.07550>.



- Clare Lyle, Marc G. Bellemare, and Pablo Samuel Castro. A comparative analysis of expected and distributional reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):4504–4511, Jul. 2019. doi: 10.1609/aaai.v33i01.33014504. URL <https://ojs.aaai.org/index.php/AAAI/article/view/4365>.
- Thanh Nguyen-Tang, Sunil Gupta, and Svetha Venkatesh. Distributional reinforcement learning via moment matching. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(10):9144–9152, 2021. URL <https://ojs.aaai.org/index.php/AAAI/article/view/17104>.
- M. T. J. Spaan. Partially observable markov decision processes. In M. Wiering and M. van Otterlo, editors, *Reinforcement Learning: State of the Art*, pages 387–414. Springer Verlag, 2012. URL <https://www.st.ewi.tudelft.nl/mtjspaان/pub/Spaan12pomdp.pdf>.
- Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the 13th AAAI Conference on Artificial Intelligence*, pages 2094–2100, 2016. URL <https://arxiv.org/abs/1509.06461>.
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *the 33rd International Conference on Machine Learning (ICML)*, volume 48 of *Proceedings of Machine Learning Research*, pages 1995–2003, 2016. URL <http://proceedings.mlr.press/v48/wangf16.html>.