

Deep Recurrent Q-Learning for Partially Observable MDPs

Matthew Hausknecht and Peter Stone

AAAI Fall Symposium on Sequential Decision Making for Intelligent Agents (AAAI-SDMIA15)

Why Do We Need Recurrency?

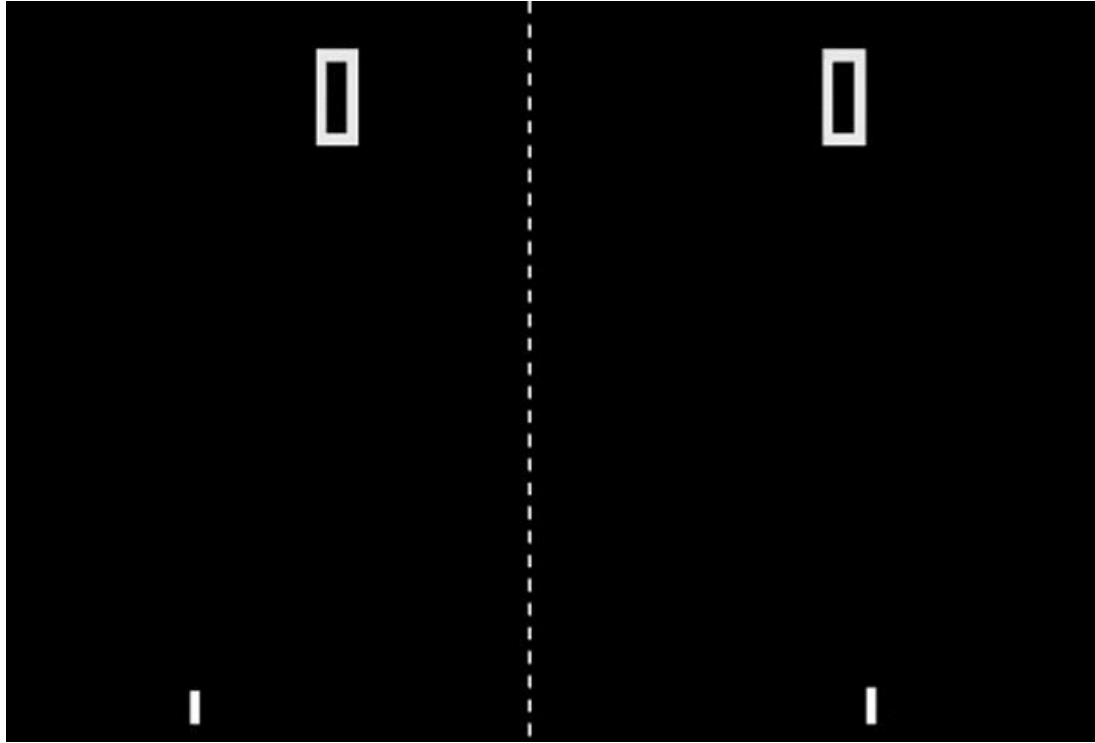


Image retrieved from: [link](#)

Why Do We Need Recurrency?

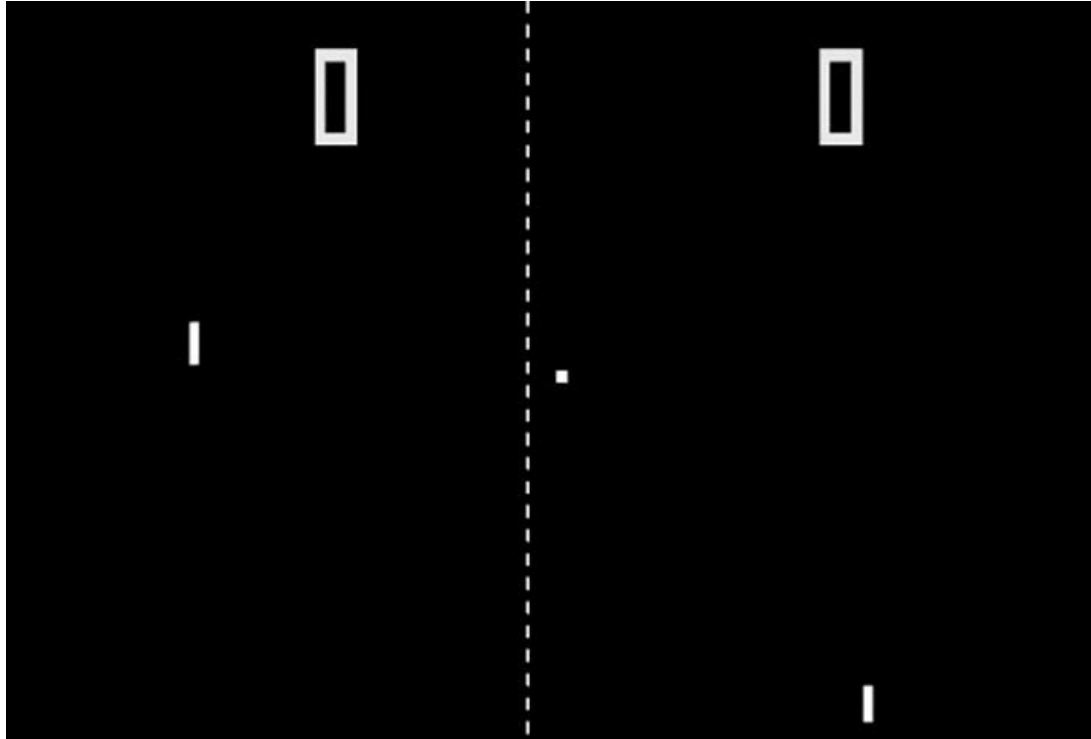


Image retrieved from: [link](#)

Why Do We Need Recurrency?



 alamy stock photo

R98MNK
www.alamy.com

Image retrieved from: [link](#)

Outline

- ~~1. Why use Recurrency?~~
2. MDPs and POMDPs
3. Deep Q-Learning
4. Putting the R in DRQN
5. Atari Experiments
6. Experiment Review
7. Conclusion

MDPs and POMDPs

Markov Decision Process

- Agent receives state \mathbf{s}
- \mathbf{s} is the true system state

Partially Observable Markov Decision Process

- Agent receives observation \mathbf{o}
- \mathbf{o} only partial description of system state

Recurrency helps to “narrow the gap” between $Q(\mathbf{s},\mathbf{a})$ and $Q(\mathbf{o},\mathbf{a})$

Q-Learning

Learn

$$Q(s, a)$$

Using

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

Deep Q-Learning

Learn

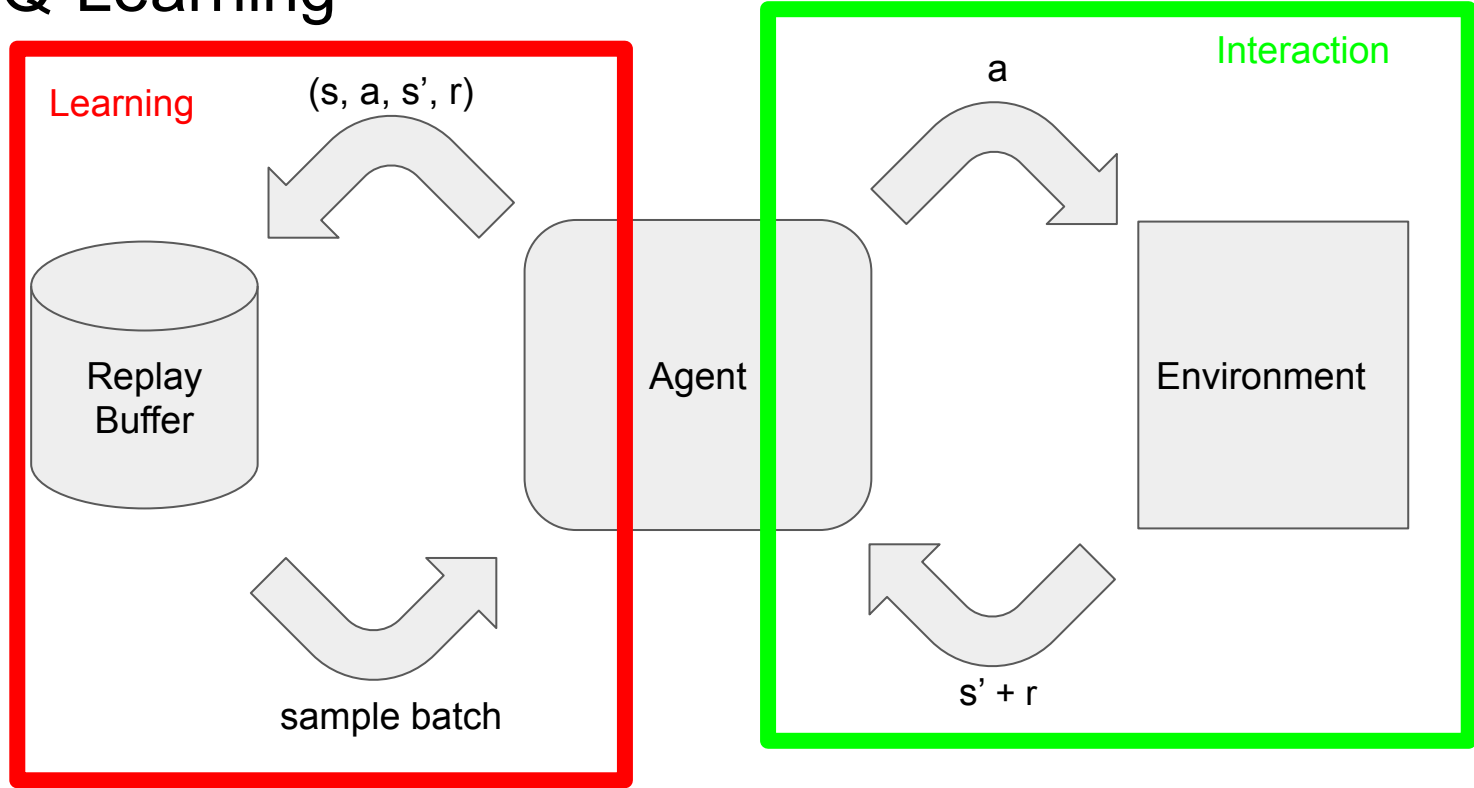
$$Q(s, a|\theta_i)$$

Using

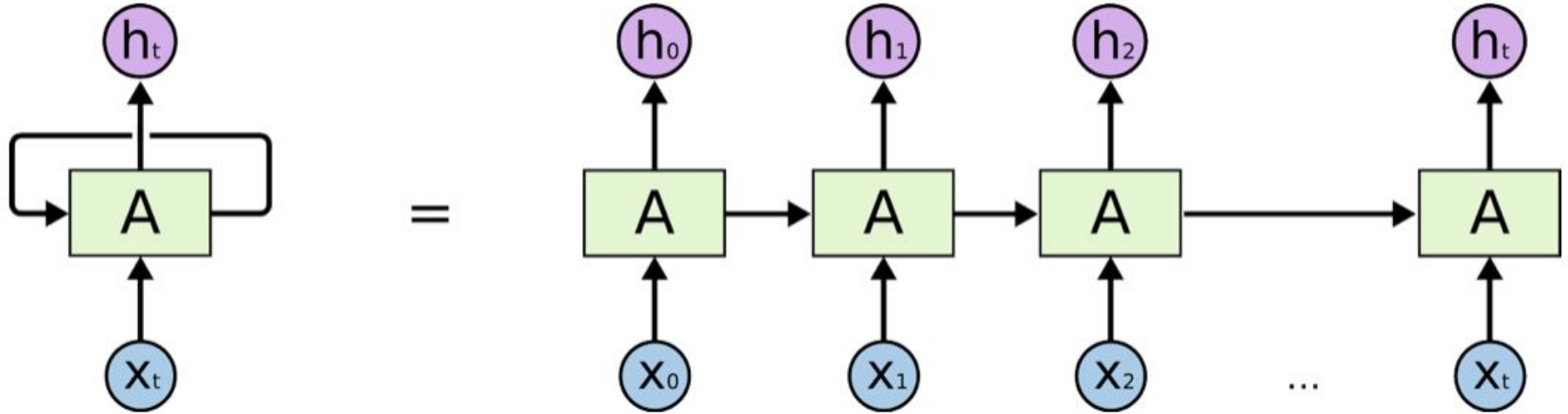
$$L(s, a|\theta_i) = \left(r + \gamma \max_{a'} Q(s', a'|\theta_i) - Q(s, a|\theta_i) \right)^2$$

$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta} L(\theta_i)$$

Deep Q-Learning



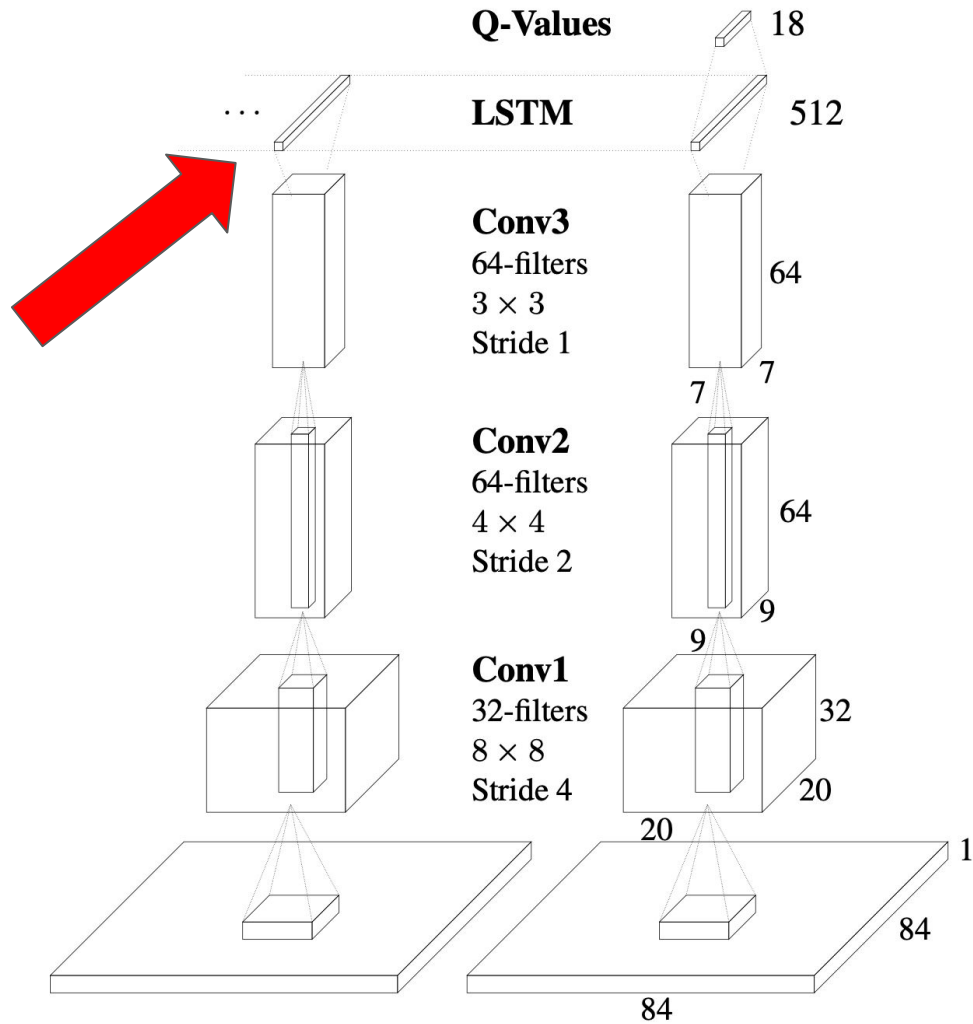
Putting the R in DRQN



An unrolled recurrent neural network.

Putting the R in DRQN

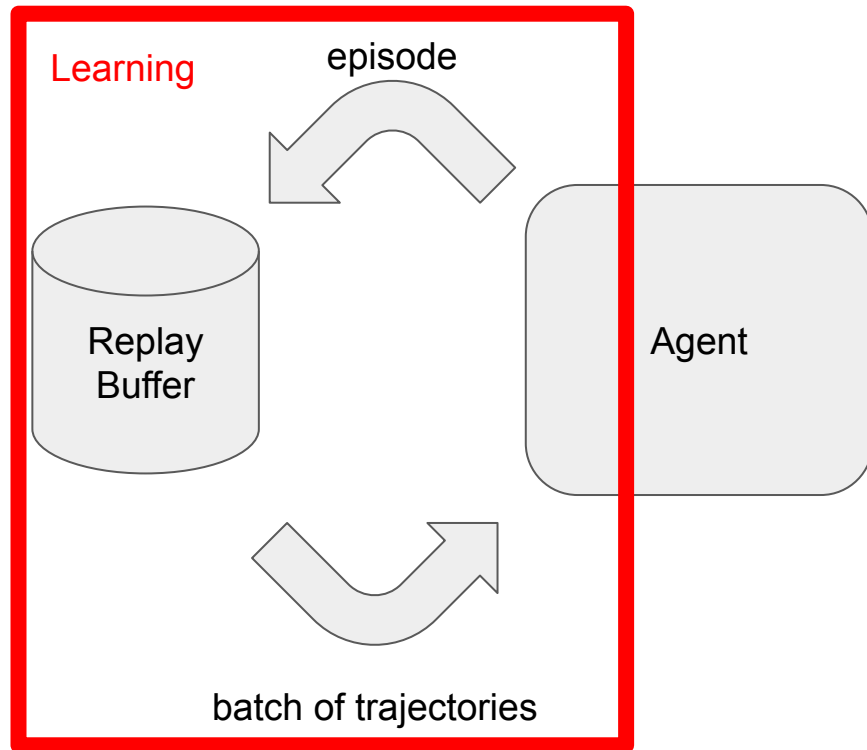
first fully connected layer is replaced with LSTM



Putting the R in DRQN

Bootstrapped Random Updates

- Sample sequences of length k
- Hidden state zeroed at start of update

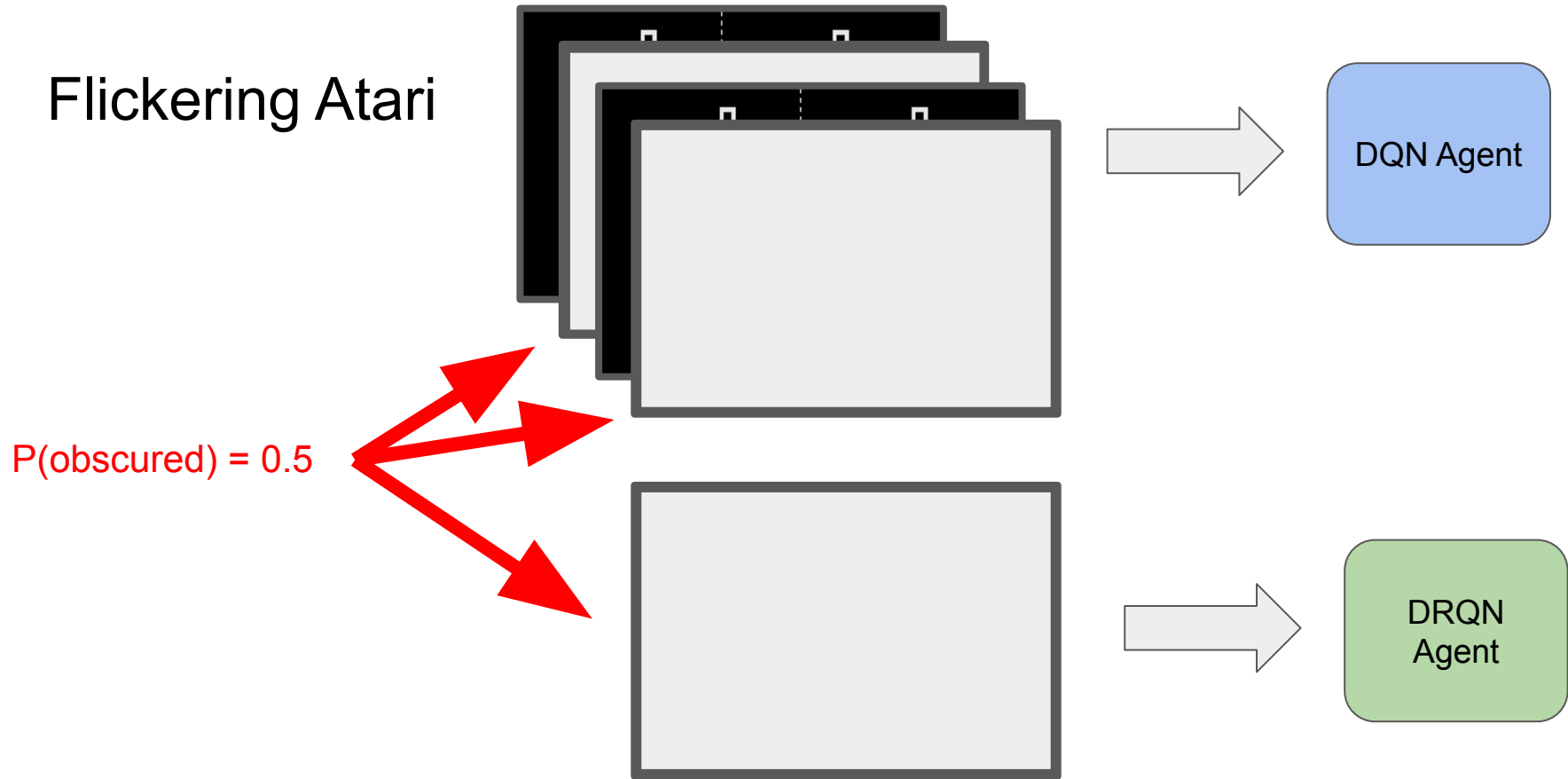


Outline

- ~~1. Why use RNNs?~~
- ~~2. MDPs and POMDPs~~
- ~~3. Deep Q Learning~~
- ~~4. Putting the R in DRQN~~
5. Atari Experiments
6. Experiment Review
7. Conclusion

Questions?

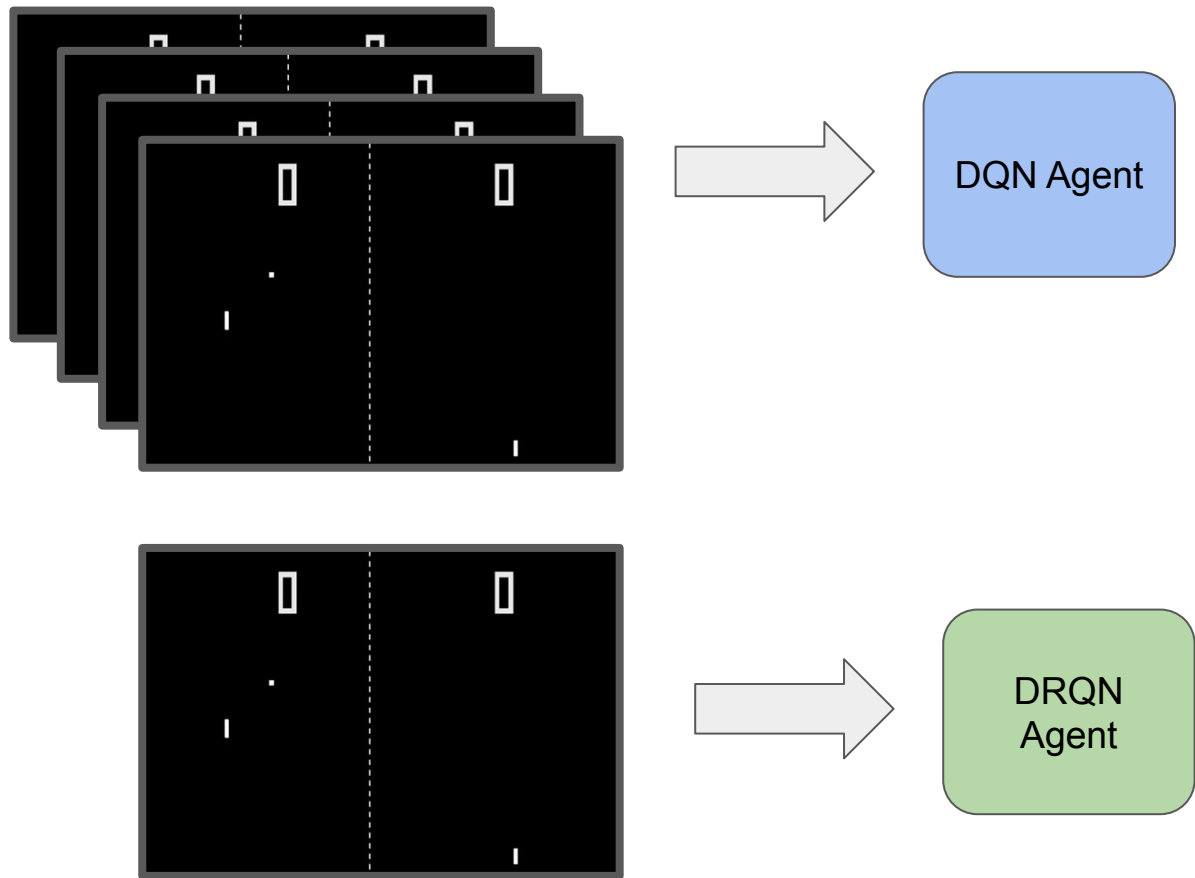
Flickering Atari



Results

Flickering	DRQN $\pm std$	DQN $\pm std$
Asteroids	1032 (± 410)	1010 (± 535)
Beam Rider	618 (± 115)	1685.6 (± 875)
Bowling	65.5 (± 13)	57.3 (± 8)
Centipede	4319.2 (± 4378)	5268.1 (± 2052)
Chopper Cmd	1330 (± 294)	1450 (± 787.8)
Double Dunk	-14 (± 2.5)	-16.2 (± 2.6)
Frostbite	414 (± 494)	436 (± 462.5)
Ice Hockey	-5.4 (± 2.7)	-4.2 (± 1.5)
Ms. Pacman	1739 (± 942)	1824 (± 490)
Pong	12.1 (± 2.2)	-9.9 (± 3.3)

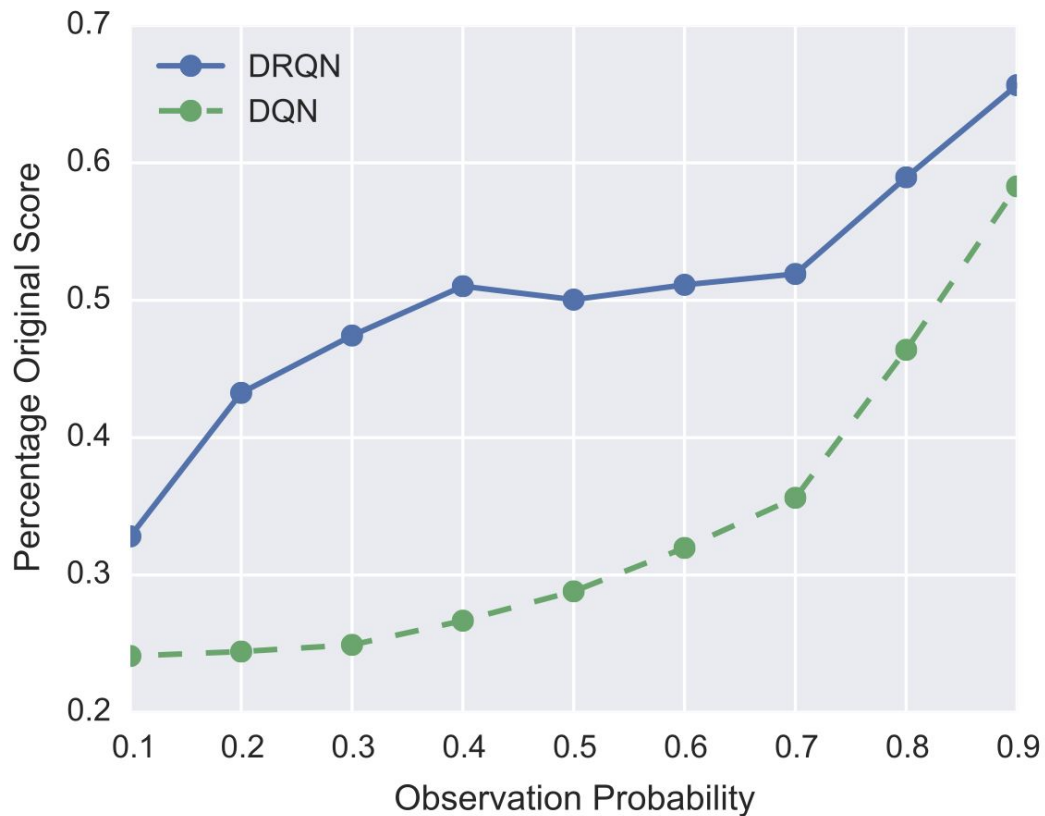
Standard Atari



Standard Atari Results

Game	DRQN $\pm std$	DQN $\pm std$	
		Ours	Mnih et al.
Asteroids	1020 (± 312)	1070 (± 345)	1629 (± 542)
Beam Rider	3269 (± 1167)	6923 (± 1027)	6846 (± 1619)
Bowling	62 (± 5.9)	72 (± 11)	42 (± 88)
Centipede	3534 (± 1601)	3653 (± 1903)	8309 (± 5237)
Chopper Cmd	2070 (± 875)	1460 (± 976)	6687 (± 2916)
Double Dunk	-2 (± 7.8)	-10 (± 3.5)	-18.1 (± 2.6)
Frostbite	2875 (± 535)	519 (± 363)	328.3 (± 250.5)
Ice Hockey	-4.4 (± 1.6)	-3.5 (± 3.5)	-1.6 (± 2.5)
Ms. Pacman	2048 (± 653)	2363 (± 735)	2311 (± 525)

MDP to POMDP Generalization



Benjamini-Hochberg Procedure

Problem: Multiple hypothesis testing can result in false positives

Solution: Adjust p-values so that the % of false positives is controlled

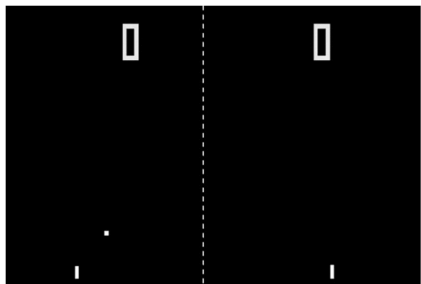
B-H with threshold of 0.05 means that 5% of **significant** results are false positives

Experiment Review

Good	Needs Improvement
<ul style="list-style-type: none">● Agents evaluated offline at intervals as recommended in [1]● Mostly small adjustments from DQN paper● BH procedure to control false positive rate	<ul style="list-style-type: none">● Individual t-tests assume normality● Compare against original DQN results but make modifications● Hyperparameter selection not explained● Ambiguous experiment details● Missing error bars● Some conclusions seem unjustified

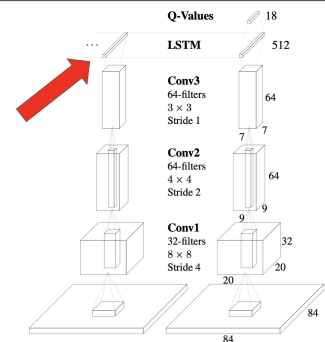
Conclusion

Why Do We Need Recurrency?



Putting the R in DRQN

first fully connected layer is replaced with LSTM



Standard Atari

