

Introduction to Deep Reinforcement Learning

Yen-Chen Wu

2015/12/11

Outline

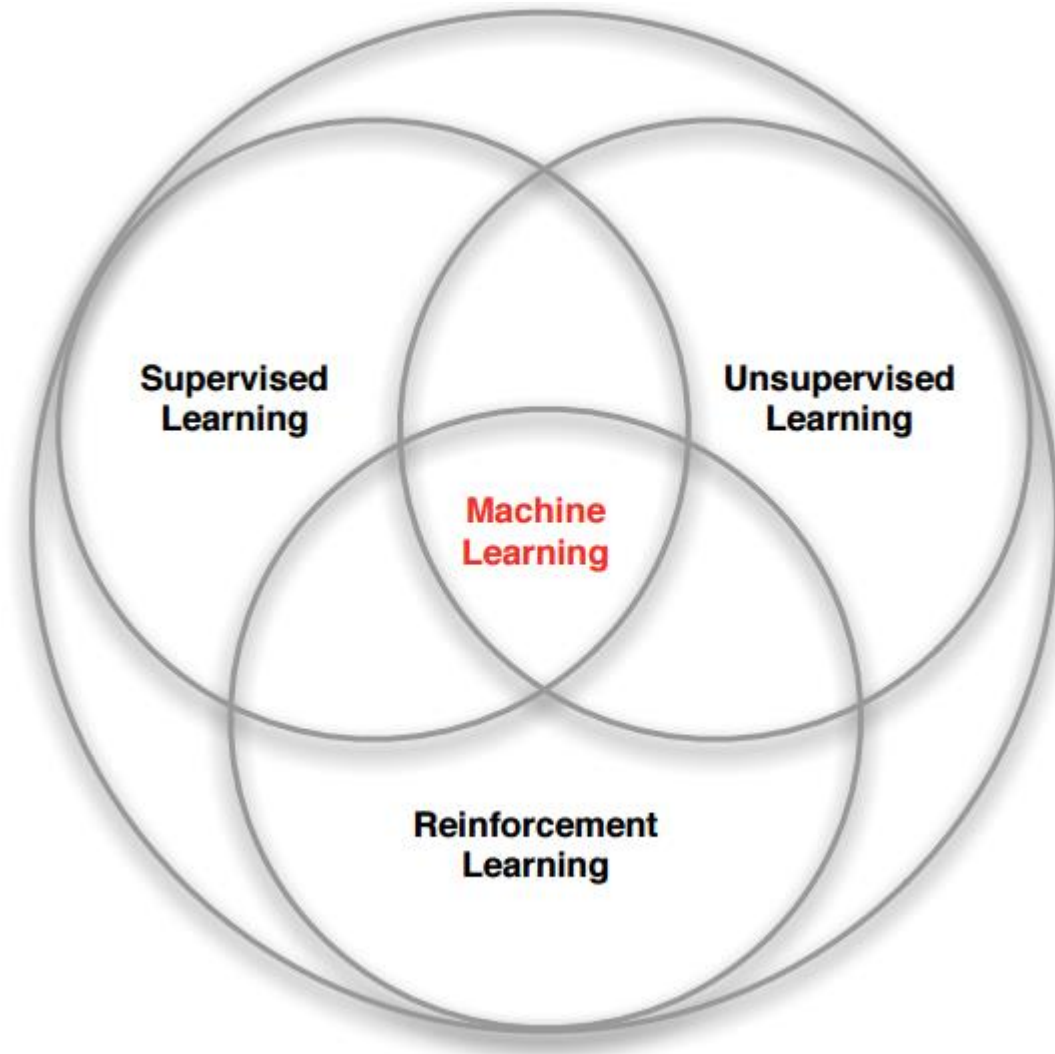
- Reinforcement Learning
- Markov Decision Process
- How to Solve MDPs
 - DP
 - MC
 - TD
 - Q-learning (DQN)
- Paper Review





REINFORCEMENT LEARNING

Branches of Machine Learning



What makes different?

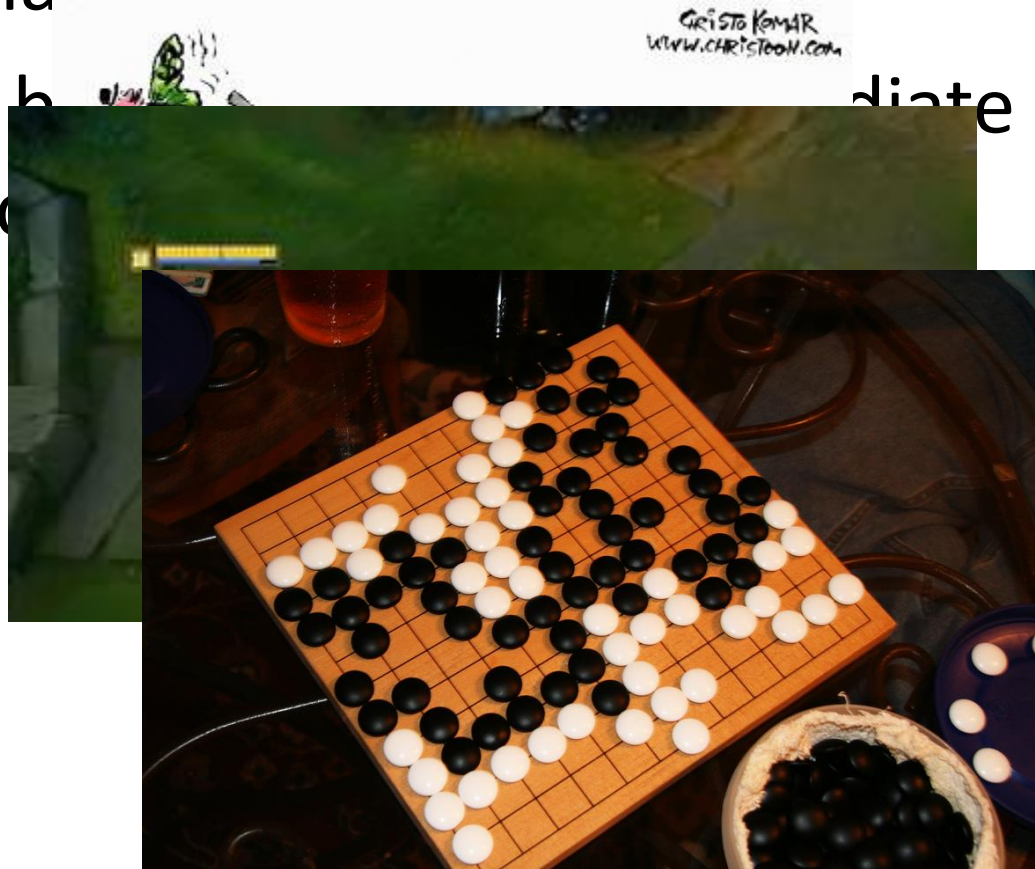
- There is no supervisor, only a **reward** signal
- Feedback is delayed, not instantaneous
- Time really matters (sequential, non i.i.d data)
- Agent's actions affect the subsequent data it receives



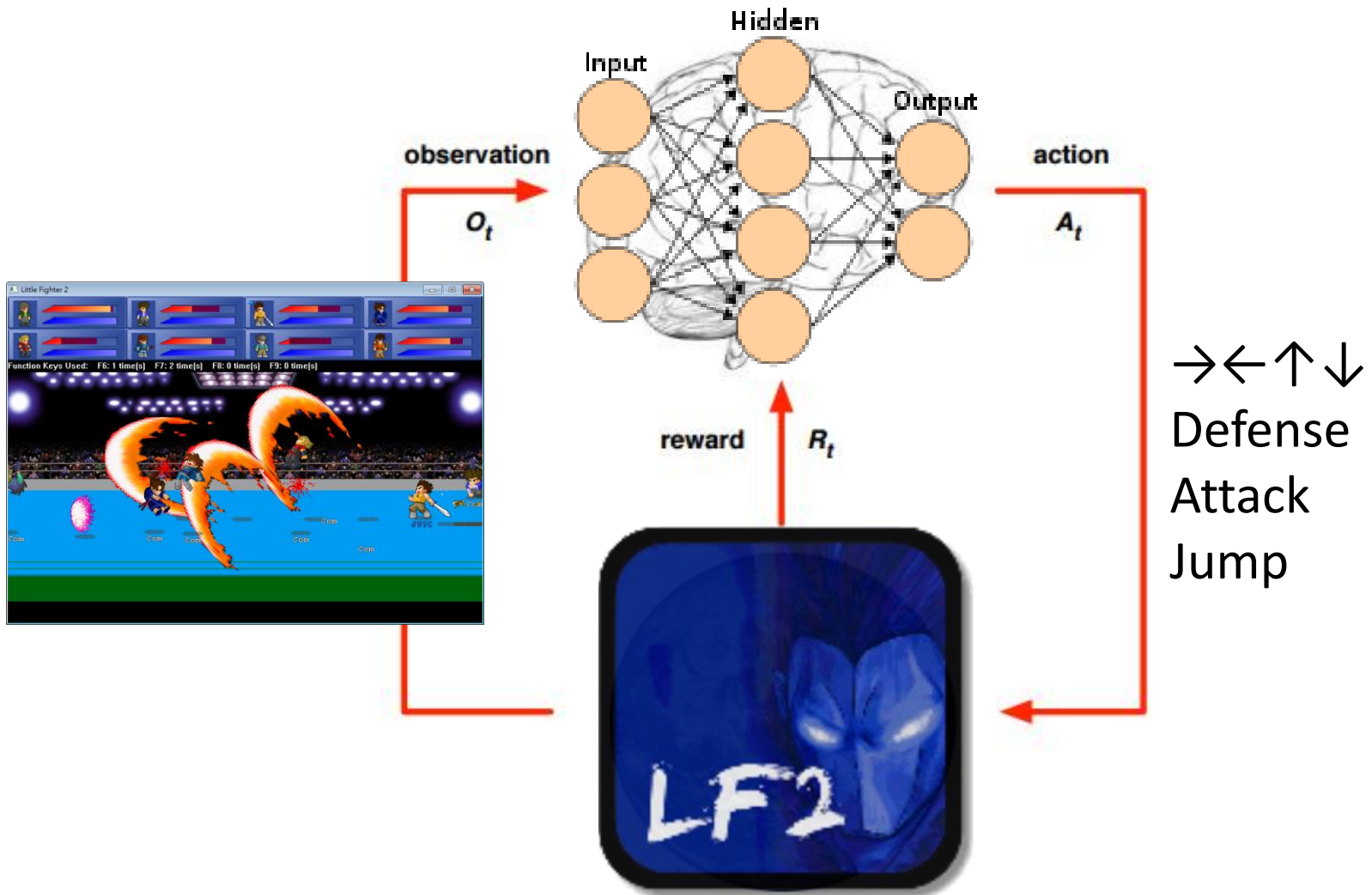
Goal:

Maximize Cumulative Reward

- Actions may have long term consequences
- Reward may be delayed
- It may be better to sacrifice immediate reward to gain more



Agent & Environment



Markov Processes

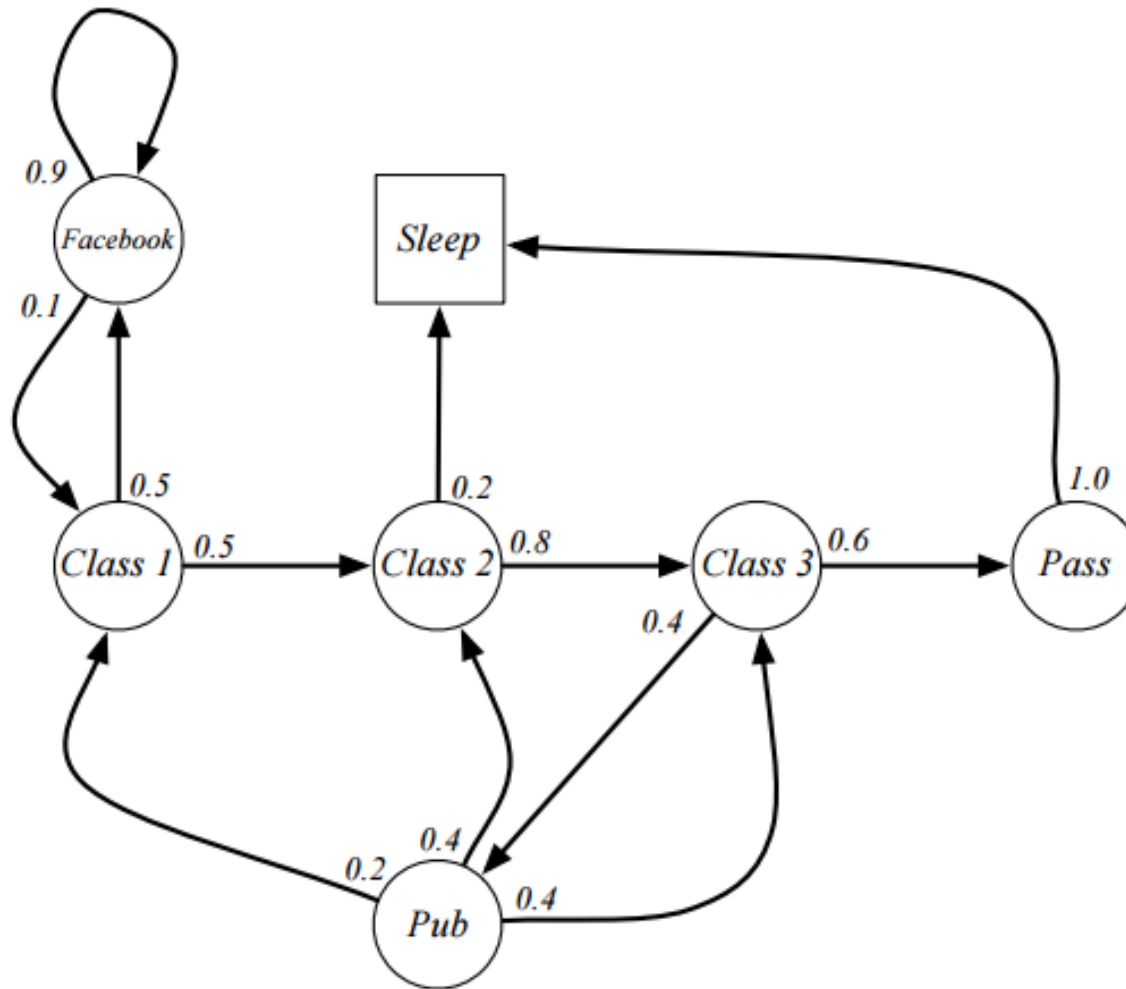
Markov Reward Processes

Markov Decision Processes

MARKOV DECISION PROCESS

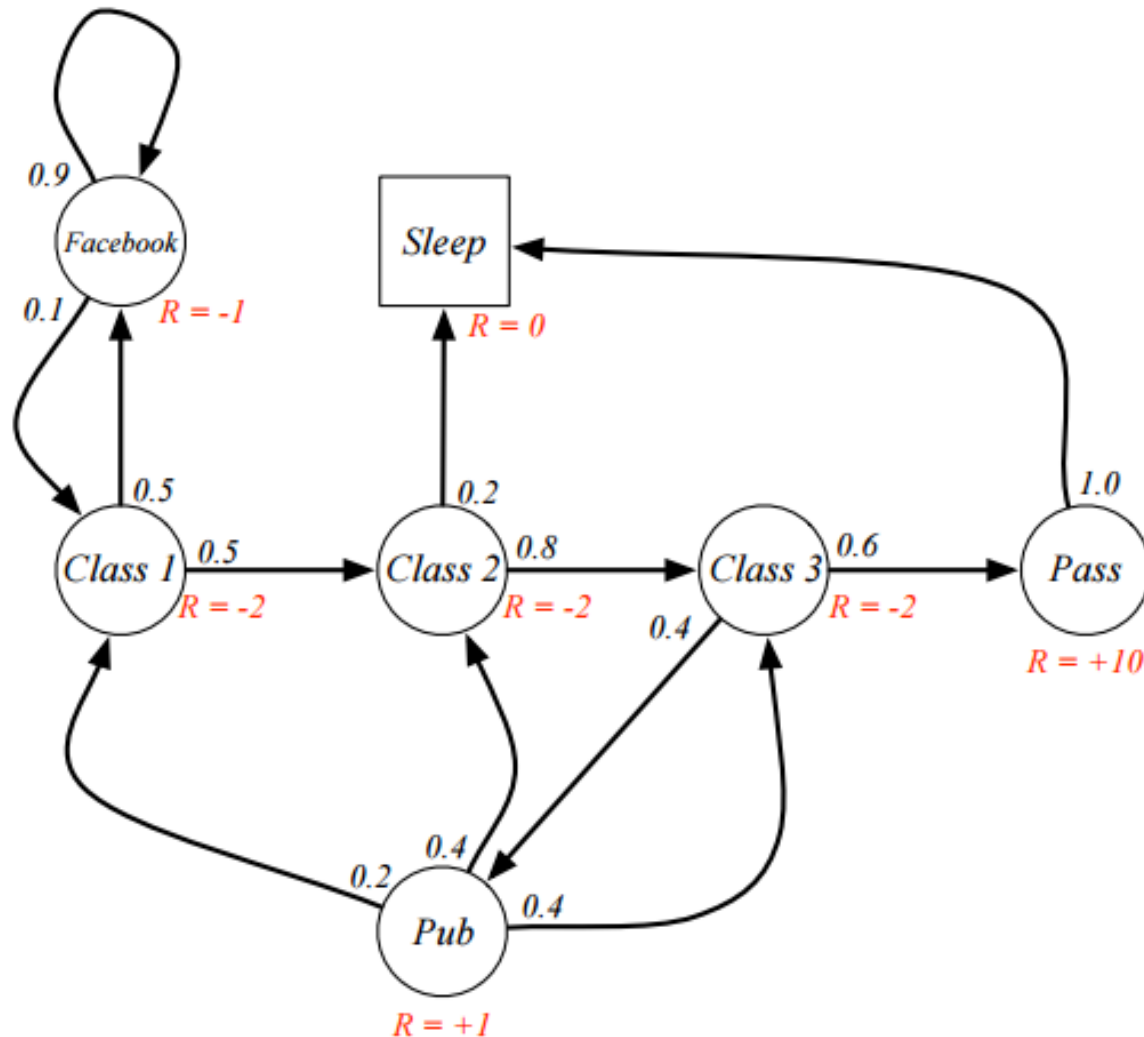
Markov Process

Example: Student Markov Chain



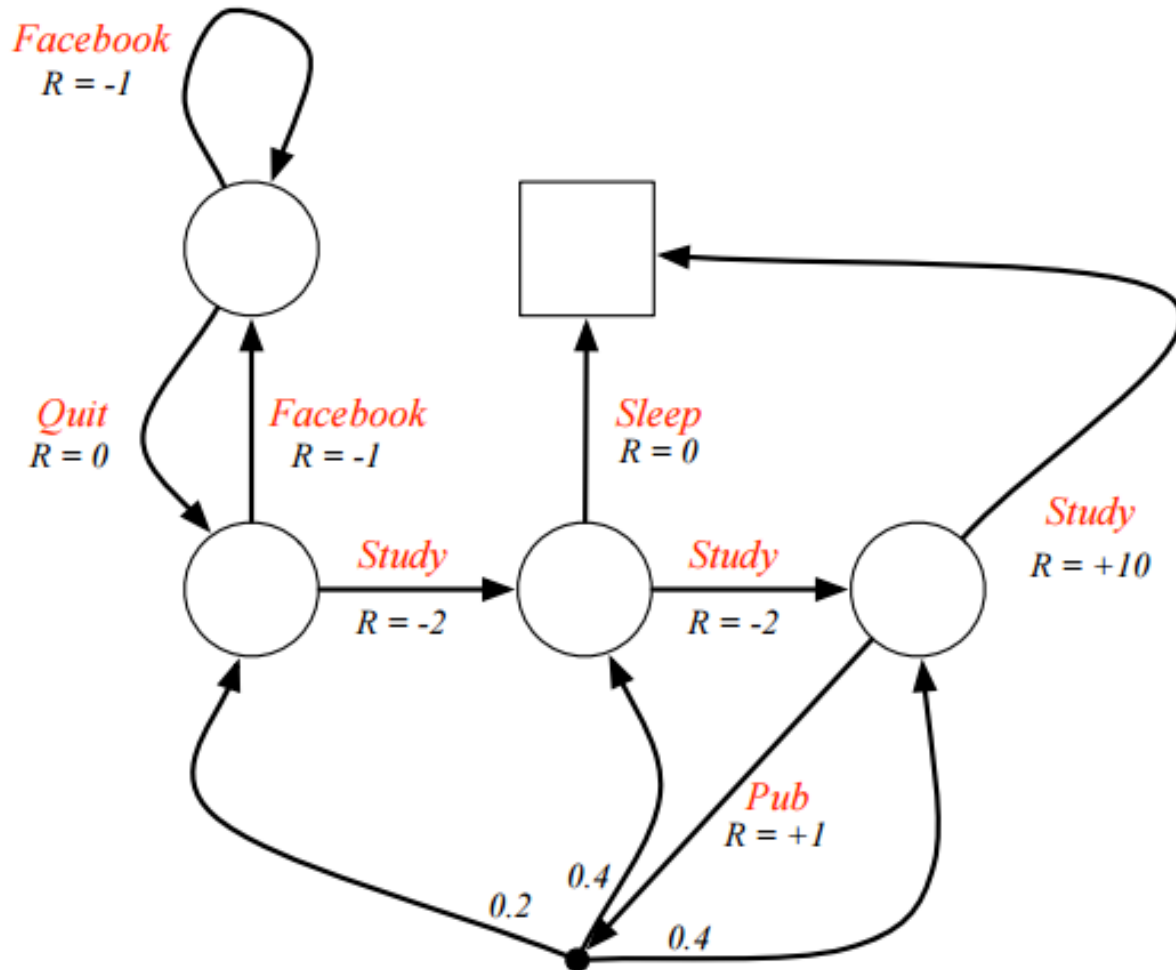
Markov Reward Processes

Example: Student MRP



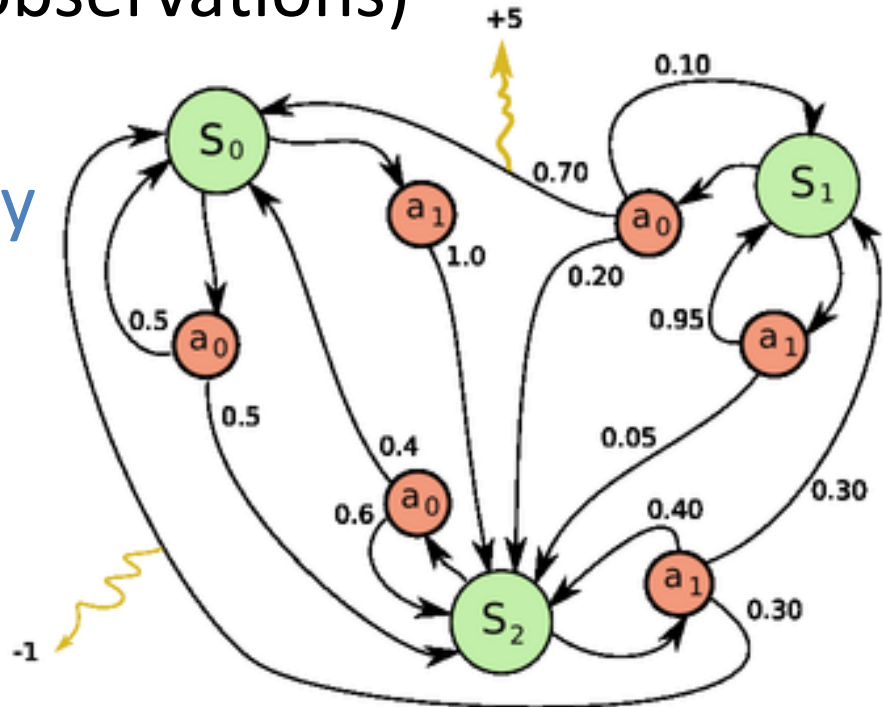
Markov Decision Process

Example: Student MDP



Markov Decision Process(MDP)

- S : finite set of **states** (observations)
- A : finite set of **actions**
- P : transition **probability**
- R : immediate **reward**
- γ : discount factor



- Goal :
 - Choose **policy** π
 - Maximize expected **return** : $R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$

Dynamic Programming

Monte-Carlo

Temporal-Difference

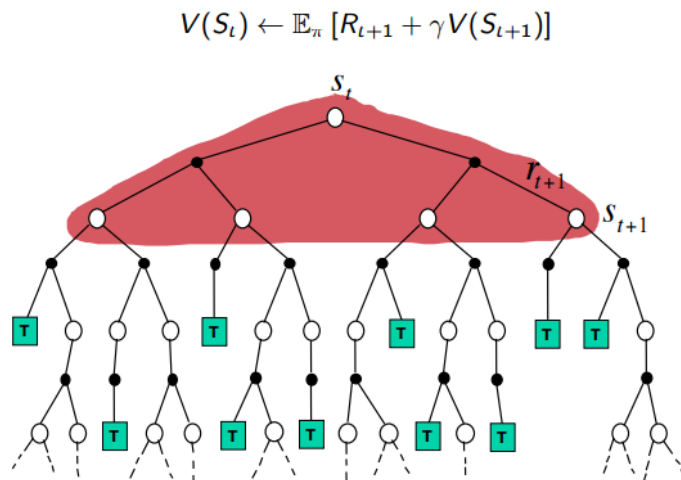
Q-Learning

HOW TO SOLVE MDP

Model-based

- Dynamic Programming
 - Evaluate policy
 - Update policy

Dynamic Programming Backup



0.22 ↖	0.25 ↖	0.27 ↖	0.30 ↖	0.34 ↖	0.38 ↓	0.34 ↖	0.30 ↖	0.34 ↖	0.38 ↓
0.25 →	0.27 →	0.30 →	0.34 →	0.38 →	0.42 ↓	0.38 ←	0.34 ←	0.38 →	0.42 ↓
0.27 ↑					0.46 ↓				0.46 ↓
0.20 ↘	0.22 ↖	0.25 ↓	-0.78 ↖ R=-1.0		0.52 →	0.57 →	0.64 ↓	0.57 ↖	0.52 ↖
0.22 ↖	0.25 ↖	0.27 ↓	0.25 ↖		0.08 ↓ R=-1.0	-0.36 → R=-1.0	0.71 ↓	0.64 ←	0.57 ←
0.25 ↖	0.27 ↖	0.30 ↓	0.27 ↖		1.20 ↑ R=1.0	0.08 ← R=-1.0	0.79 ↓	-0.29 ← R=-1.0	0.52 ↓
0.27 ↖	0.30 ↖	0.34 ↓	0.30 ←		1.08 ↑	0.97 ←	0.87 ←	-0.21 ← R=-1.0	0.57 ↓
0.31 ↖	0.34 ↖	0.38 ↓	-0.58 ↓ R=-1.0		-0.08 ↑ R=-1.0	-0.18 ↑ R=-1.0	0.78 ↑	0.71 ←	0.64 ←
0.34 →	0.38 →	0.42 →	0.46 →	0.52 →	0.57 →	0.64 →	0.71 ↑	0.64 ↖	0.52 ↖
0.31 ↖	0.34 ↖	0.38 ↖	0.42 ↖	0.46 ↖	0.52 ↖	0.57 ↖	0.64 ↑	0.52 ↖	0.52 ↖

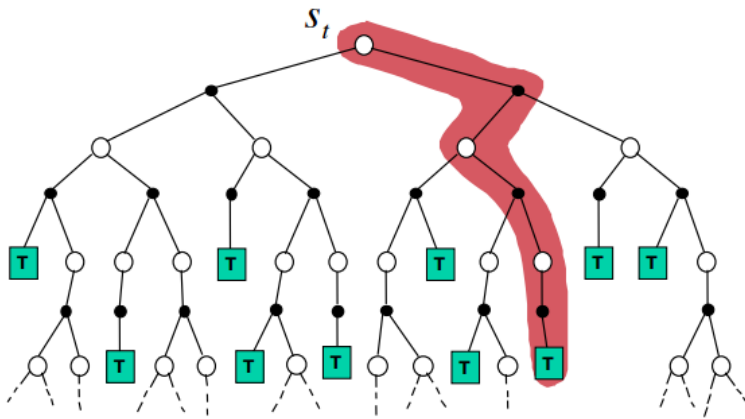
Right: A simple Gridworld solved with a Dynamic Programming. Very exciting. Head over to the [GridWorld: DP demo](#) to play with the GridWorld environment and policy iteration.

Model Free

- Unknown Transition Probability & Reward
- MC vs TD

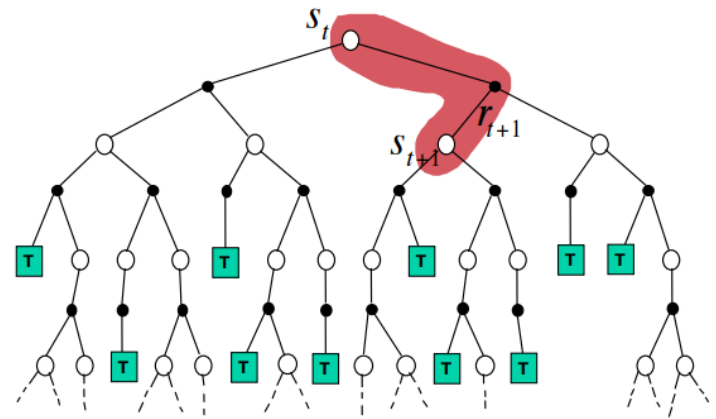
Monte-Carlo Backup

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



Temporal-Difference Backup

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



Model Free: Q-learning

- Instead of tabular
- optimal action-value **function** (Q-learning)
 - $Q^*(s,a) = \max_{\pi} \mathbb{E}[R_t | s_t = s, a_t = a, \pi]$
- Bellman equation

$$Q^*(s,a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s',a') | s,a \right]$$

- Basic idea : iterative update (lack of generalization)
- In practical : function approximator $Q(s,a; \theta) \approx Q^*(s,a)$
 - Linear ?
 - Using DNN !

LETTER

doi:10.1038/nature14236

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

DEEP Q-NETWORK (DQN)

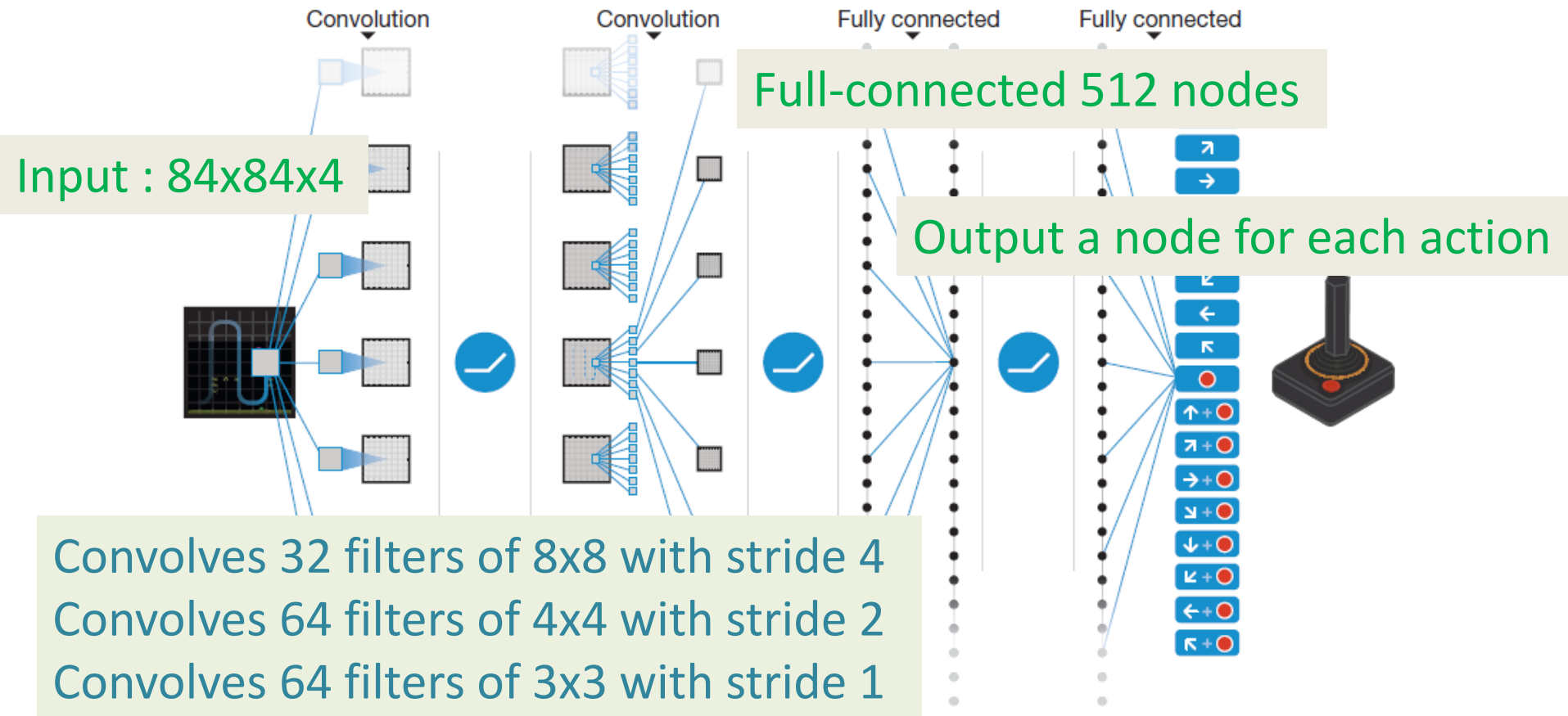
Video

- <https://www.youtube.com/watch?v=LJ4oCb6u7kk>



Deep Q-Network

- compute Q-values for all actions



Update DQN

- Loss function

$$\begin{aligned} L_i(\theta_i) &= \mathbb{E}_{s,a,r} \left[\left(\mathbb{E}_{s'} [y|s,a] - Q(s,a; \theta_i) \right)^2 \right] \\ &= \mathbb{E}_{s,a,r,s'} \left[(y - Q(s,a; \theta_i))^2 \right] + \mathbb{E}_{s,a,r} [\mathbb{V}_{s'} [y]] \end{aligned}$$

- Gradient

$$\nabla_{\theta_i} L(\theta_i) = \mathbb{E}_{s,a,r,s'} \left[\left(r + \gamma \max_{a'} Q(s',a'; \theta_i^-) - Q(s,a; \theta_i) \right) \nabla_{\theta_i} Q(s,a; \theta_i) \right]$$

Two Technique

- Experience Replay
 - Experience
 - Pooled Mer $e_t = (s_t, a_t, r_t, s_{t+1})$
 - Data efficiency ($D_t = \{e_1, \dots, e_t\}$)
 - Avoid correlation between samples (variance between batches)
 - Off-policy is suitable for Q-learning
 - Random sampled mini-batch

	Example	Learn the value of...	Pros & Cons
• On-policy	SARSA	policy being carried out by the agent	Fast but weak
Off-policy	DQN	optimal policy independently of the agent's actions	Slow but robust

DEMO



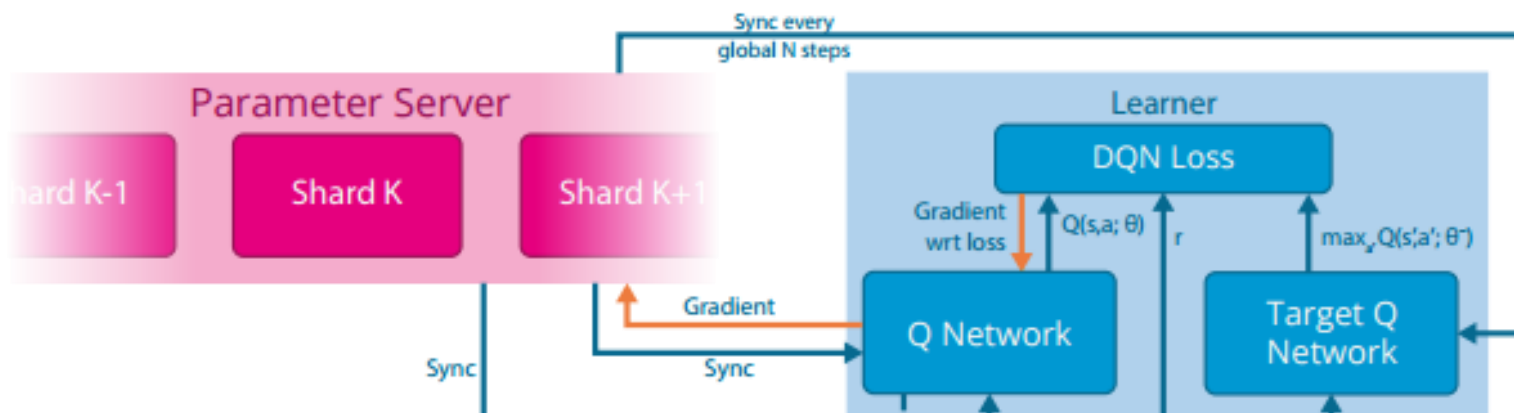


PAPER REVIEW

Paper list

- **Massively Parallel Methods for Deep Reinforcement Learning**
- **Continuous control with deep reinforcement learning**
- **Deep Reinforcement Learning with Double Q-learning**
- **Policy Distillation**
- **Dueling Network Architectures for Deep Reinforcement Learning**
- **Multiagent Cooperation and Competition with Deep Reinforcement Learning**

Gorila (GOogle ReInforcement Learning Architecture)



Massively Parallel Methods for Deep Reinforcement Learning

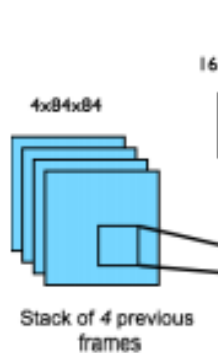
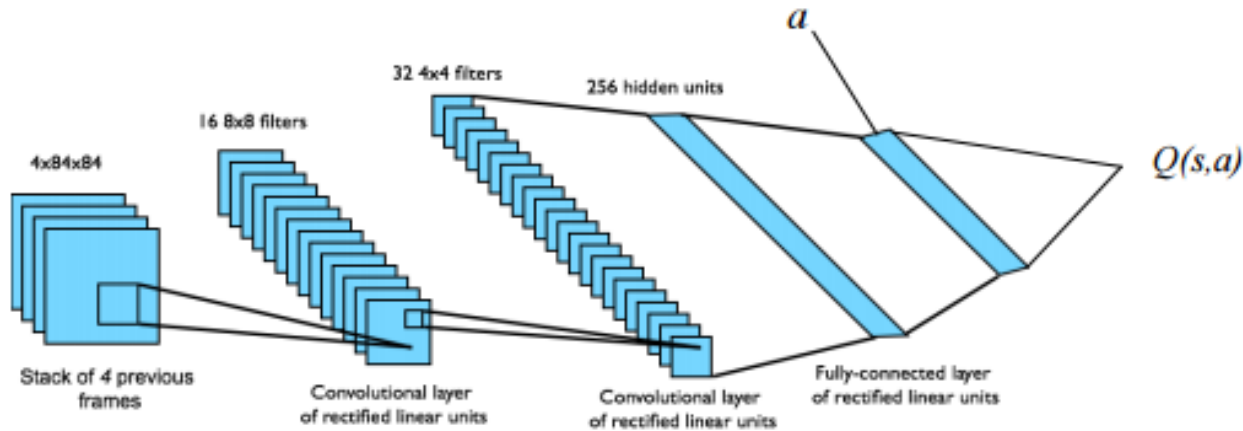
Arun Nair

arXiv:1507.04296

- ▶ **Parallel acting:** generate new interactions
- ▶ **Distributed replay memory:** save interactions
- ▶ **Parallel learning:** compute gradients from replayed interactions
- ▶ **Distributed neural network:** update network from gradients

DDPG (Deterministic Policy Gradient)

- DDAC (Deep Deterministic Actor-Critic)



Continuous control with deep reinforcement learning

Timothy P. Lillicrap

arXiv:1509.02971

<https://goo.gl/J4PIAz>

of rectified linear units

of rectified linear units

Double Q-learning

$$Y_t^{\text{DQN}} \equiv R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t^-).$$

$$Y_t^{\text{DoubleQ}} \equiv R_{t+1} + \gamma Q(S_{t+1}, \operatorname{argmax}_a Q(S_{t+1}, a; \theta_t); \theta'_t).$$



Policy Distillation

- Soft target

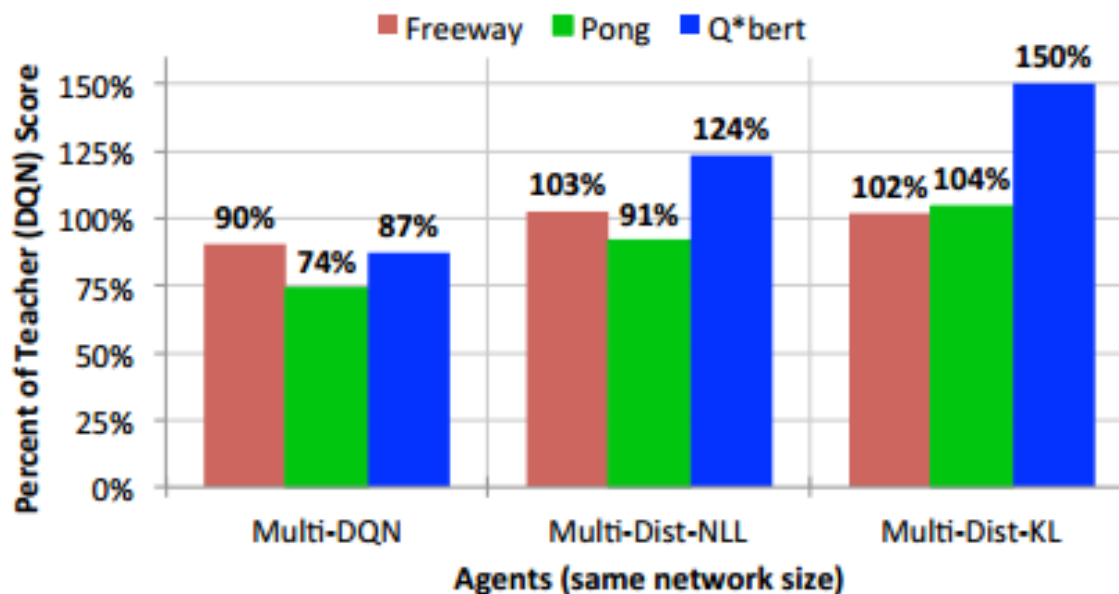


Figure 3: Performance of multi-task agents with identical network architecture and size, relative to respective single-task DQN teachers. A detailed results table is given in Appendix B

Dueling Network

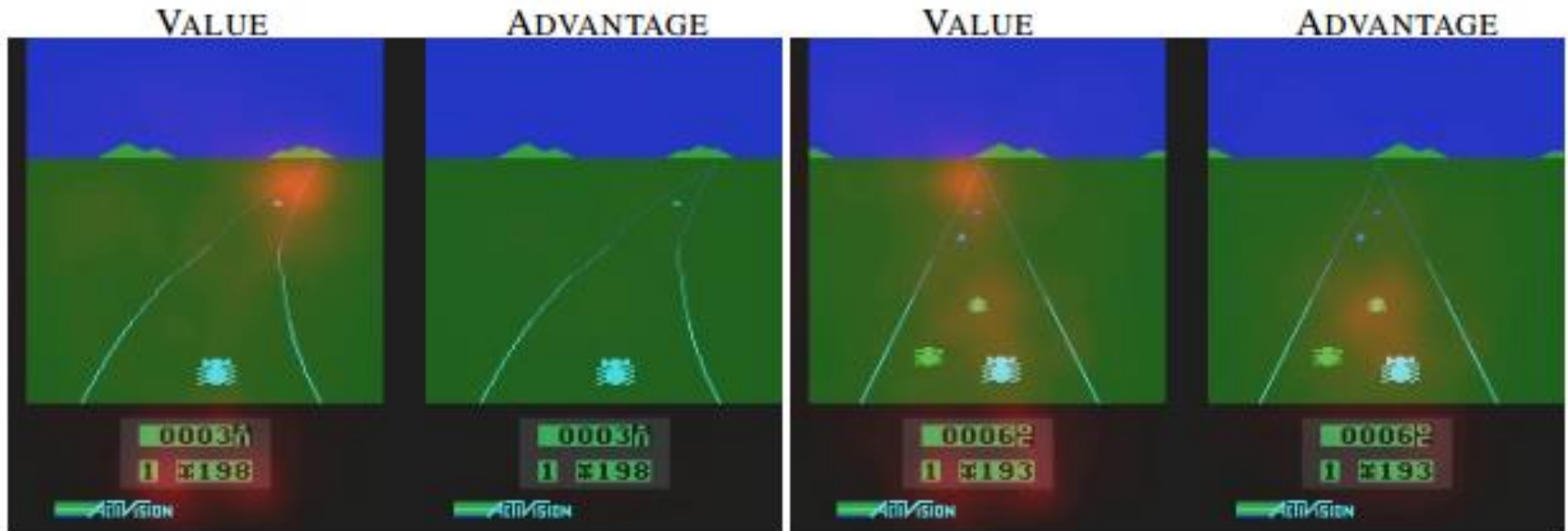
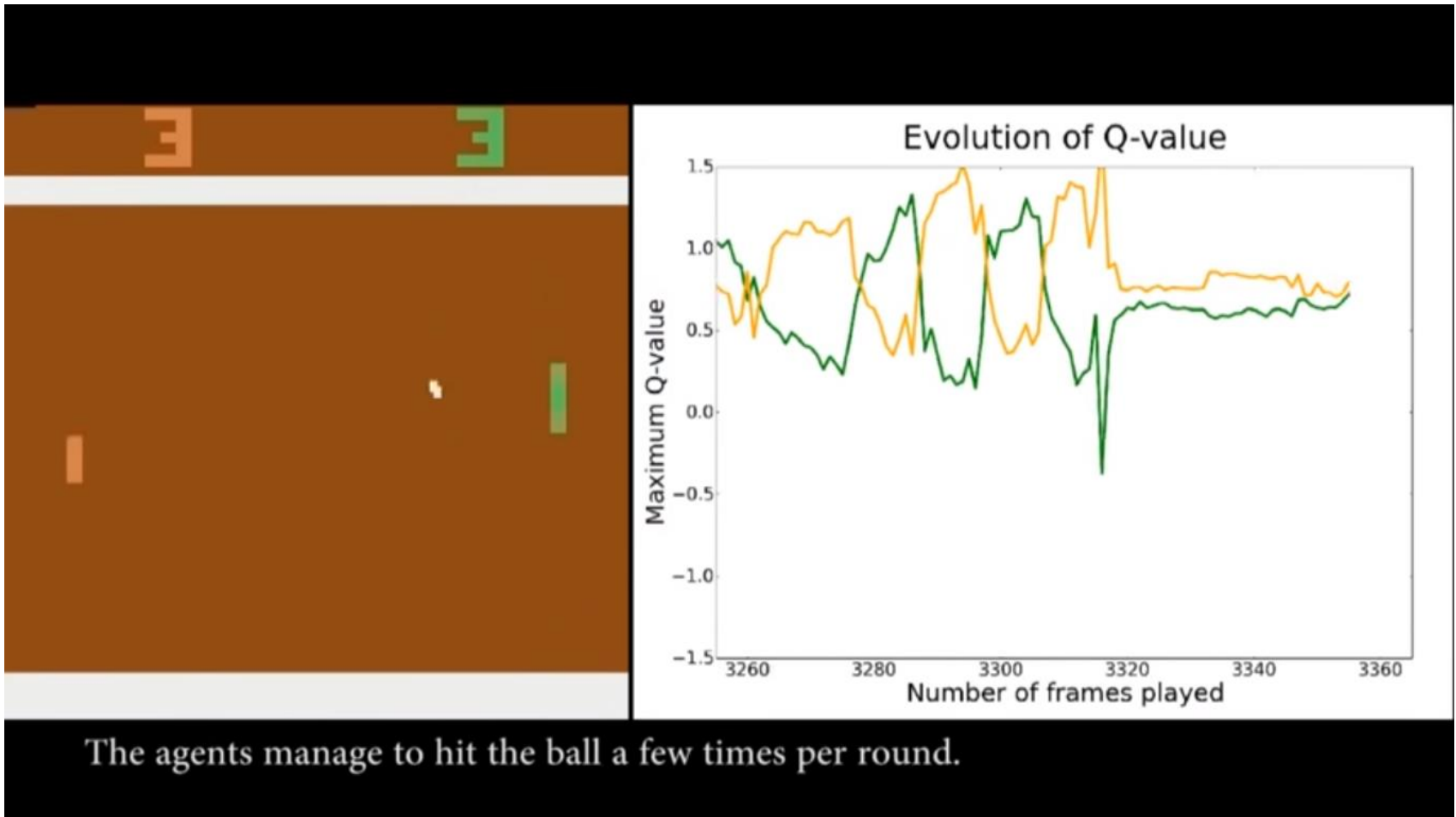


Figure 2: See, attend and drive: Value and advantage saliency maps on the Enduro game for a trained dueling architecture. The value stream learns to pay attention to the road. The advantage stream learns to pay attention only when there are cars immediately in front, so as to avoid collisions.

Multiagent



The agents manage to hit the ball a few times per round.