



Partially Observable RL

CS885 Reinforcement Learning

Module 4: October 4, 2021

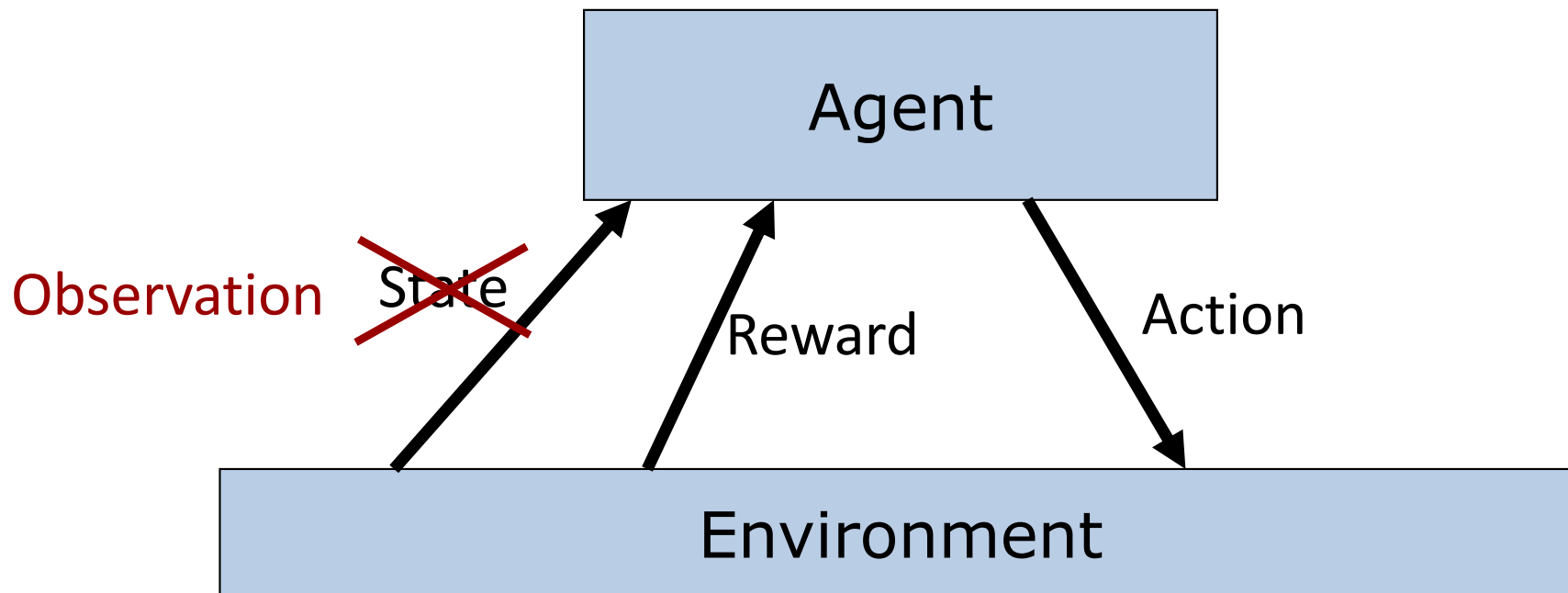
Hausknecht, M., & Stone, P. (2015). Deep recurrent q-learning for partially observable MDPs. In 2015 AAAI fall symposium series.

Outline

- Partially Observable Markov Decision Processes
- Hidden Markov Models
- Recurrent neural networks
 - Long short term memory (LSTM) networks
- Deep recurrent Q-networks



Reinforcement Learning Problem



Goal: Learn to choose actions that maximize rewards

Markov Process

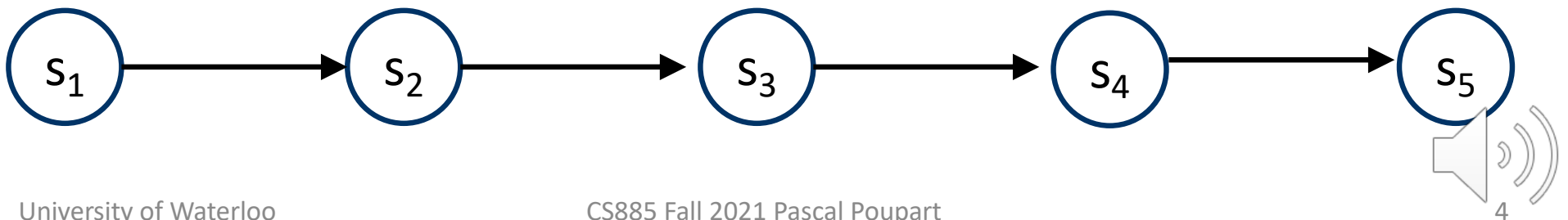
- Assumptions:

- (first-order) Markovian:

$$\Pr(s_t | s_{t-1}, \dots, s_0) = \Pr(s_t | s_{t-1})$$

- Stationary:

$$\Pr(s_t | s_{t-1}) = \Pr(s_{t+1} | s_t) \forall t$$



Hidden Markov Model

- Assumptions:

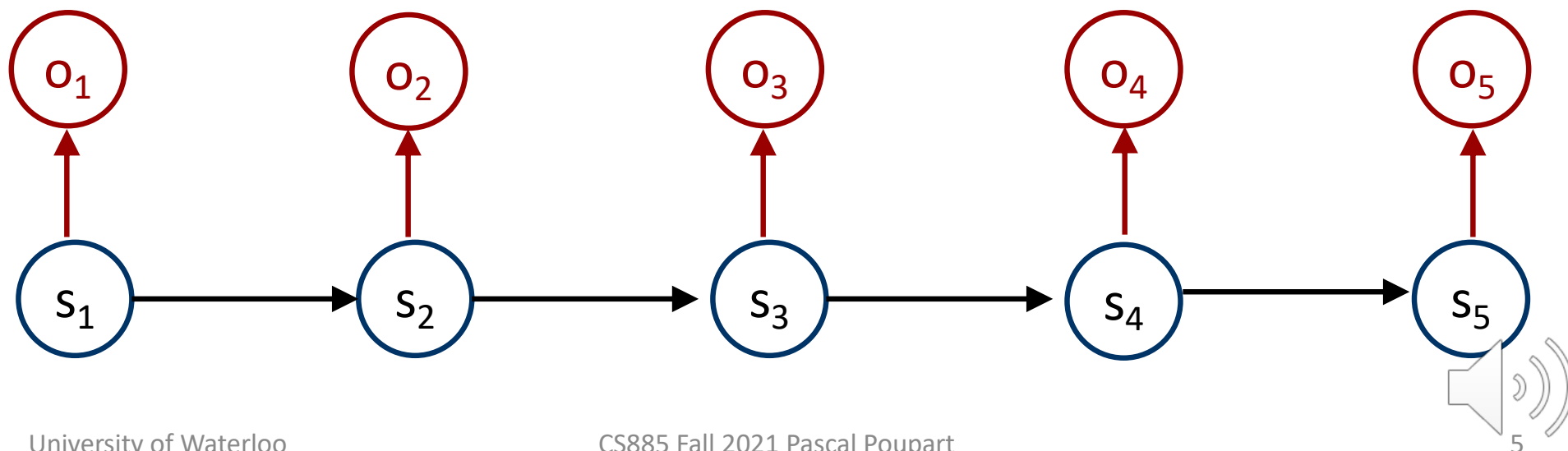
- (first-order) Markovian:

$$\Pr(s_t | s_{t-1}, \dots, s_0) = \Pr(s_t | s_{t-1})$$

- Stationary:

$$\Pr(s_t | s_{t-1}) = \Pr(s_{t+1} | s_t) \forall t$$

$$\Pr(o_t | s_t) = \Pr(o_{t+1} | s_{t+1}) \forall t$$

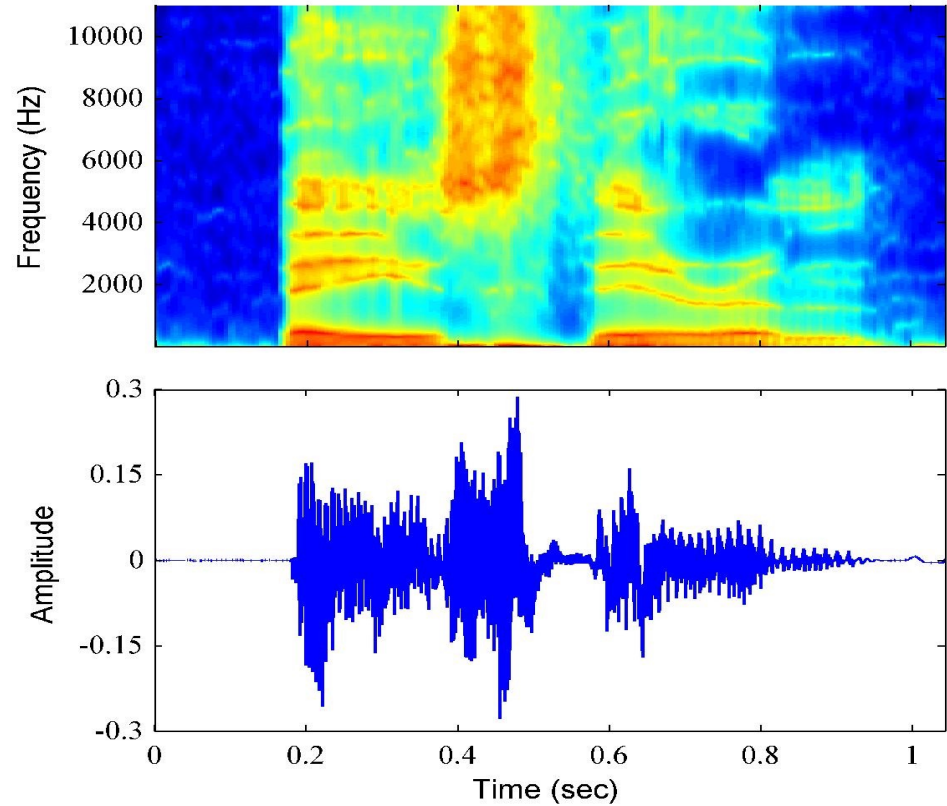


Speech Recognition

Compute:

$$\Pr(\textit{word}|\textit{speech})$$

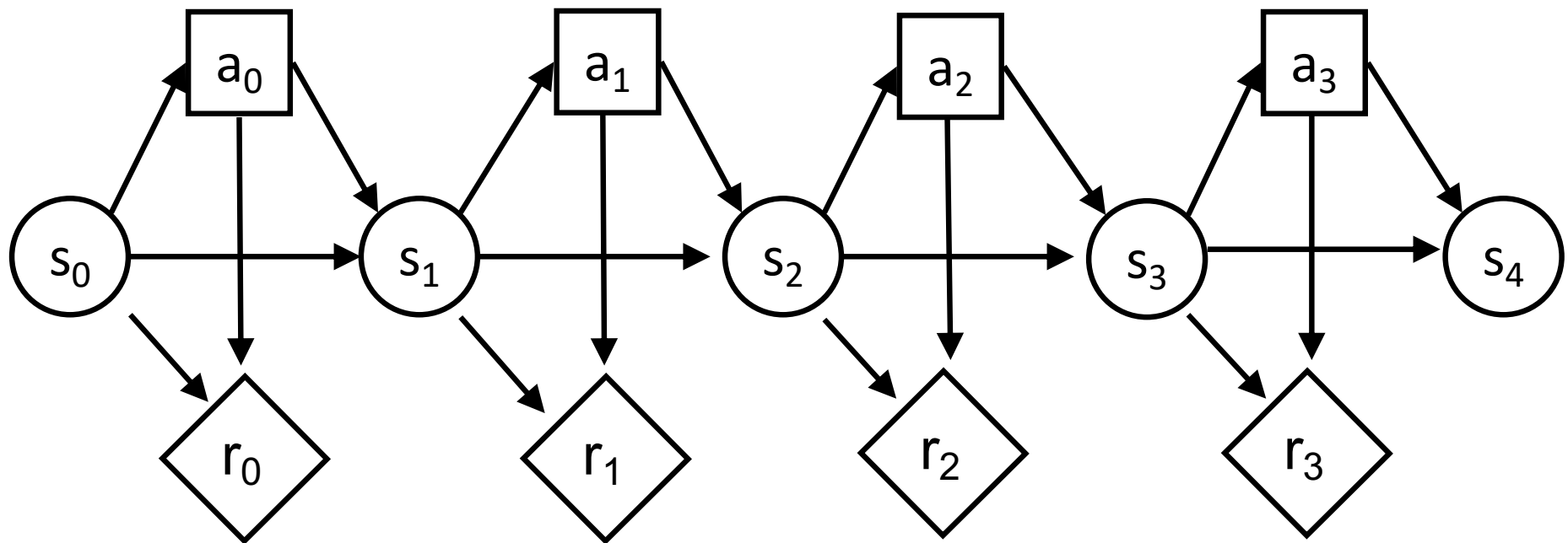
$$\Pr(s_t|o_t, o_{t-1}, \dots, o_1)$$



| b | ey | z | th | ih | er | em |
| Bayes' | Theorem |

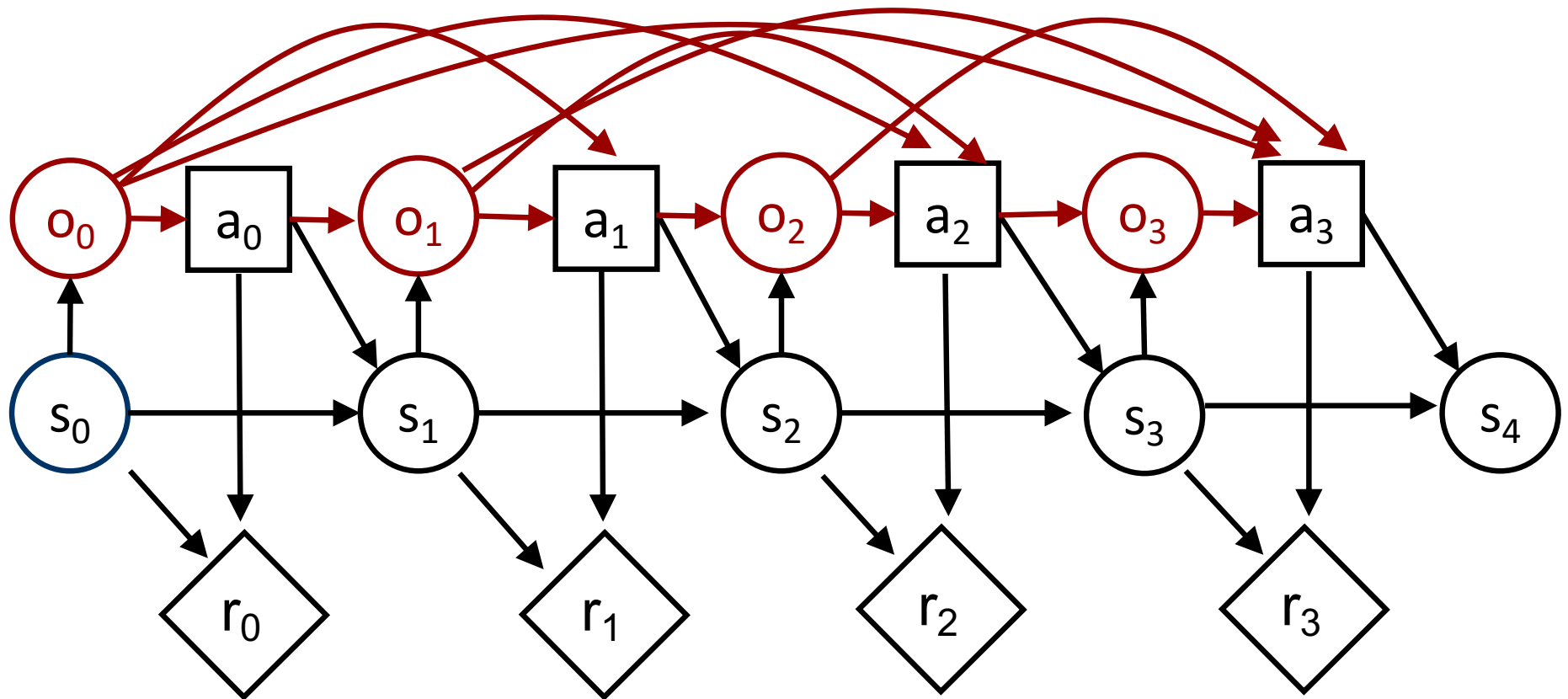


(Fully Observable) Markov Decision Process (MDP)



Partially Observable Markov Decision Process (POMDP)

- MDP augmented with observations



Partially Observable RL

- Definition

- States: $s \in S$
- Observations: $o \in O$
- Actions: $a \in A$
- Rewards: $r \in \mathbb{R}$
- Transition model: $\Pr(s_t | s_{t-1}, a_{t-1})$
- Observation model: $\Pr(o_t | a_{t-1}, s_t)$
- Reward model: $\Pr(r_t | s_t, a_t)$
- Discount factor: $0 \leq \gamma \leq 1$
 - discounted: $\gamma < 1$ undiscounted: $\gamma = 1$
- Horizon (i.e., # of time steps): h
 - Finite horizon: $h \in \mathbb{N}$ infinite horizon: $h = \infty$

} unknown model

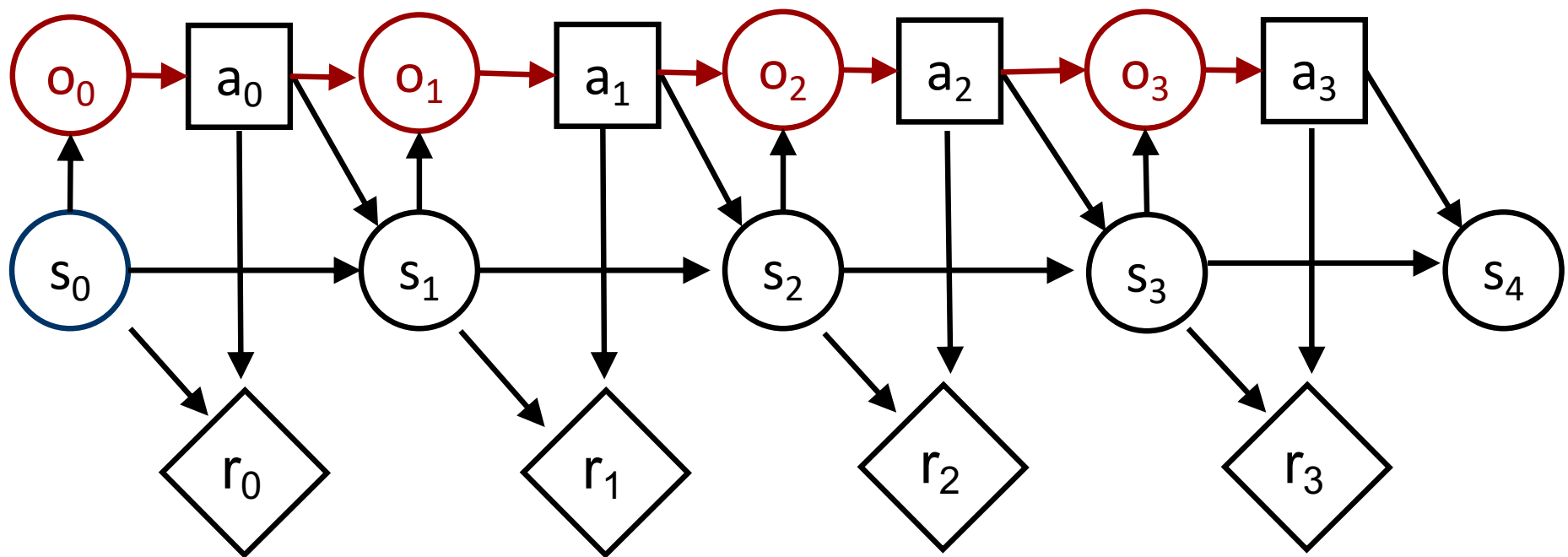
- Goal: find optimal policy π^* such that

$$\pi^* = \operatorname{argmax}_{\pi} \sum_{t=0}^h \gamma^t E_{\pi}[r_t]$$



Simple Heuristic

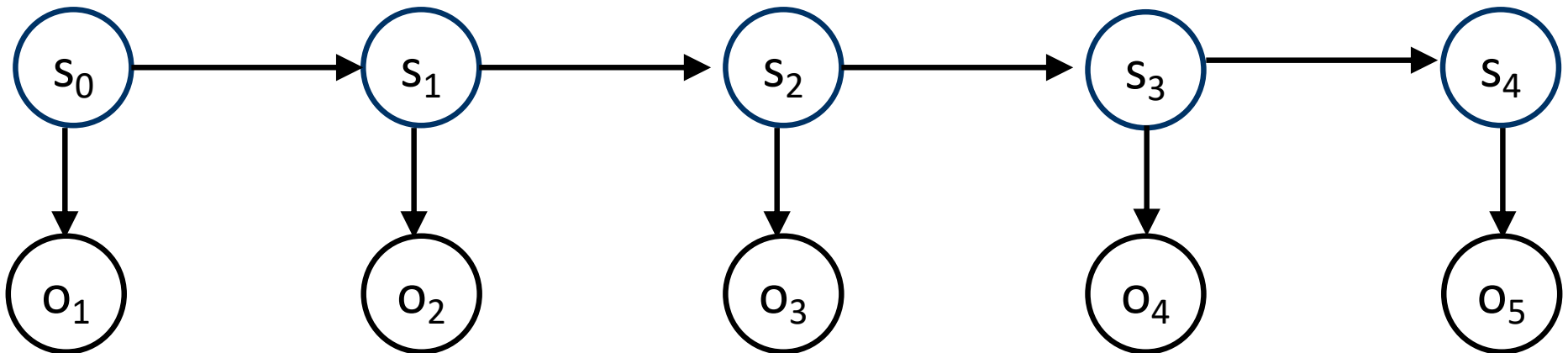
- Approximate s_t by o_t (or finite window of previous observations: $o_{t-k}, o_{t-k+1}, \dots, o_t$)
- Use favorite RL algo on observations instead of states



Belief Monitoring

- Hidden Markov model
 - Initial state distribution: $\Pr(s_0)$
 - Transition probabilities: $\Pr(s_{t+1}|s_t)$
 - Observation probabilities: $\Pr(o_t|s_t)$
- Belief monitoring

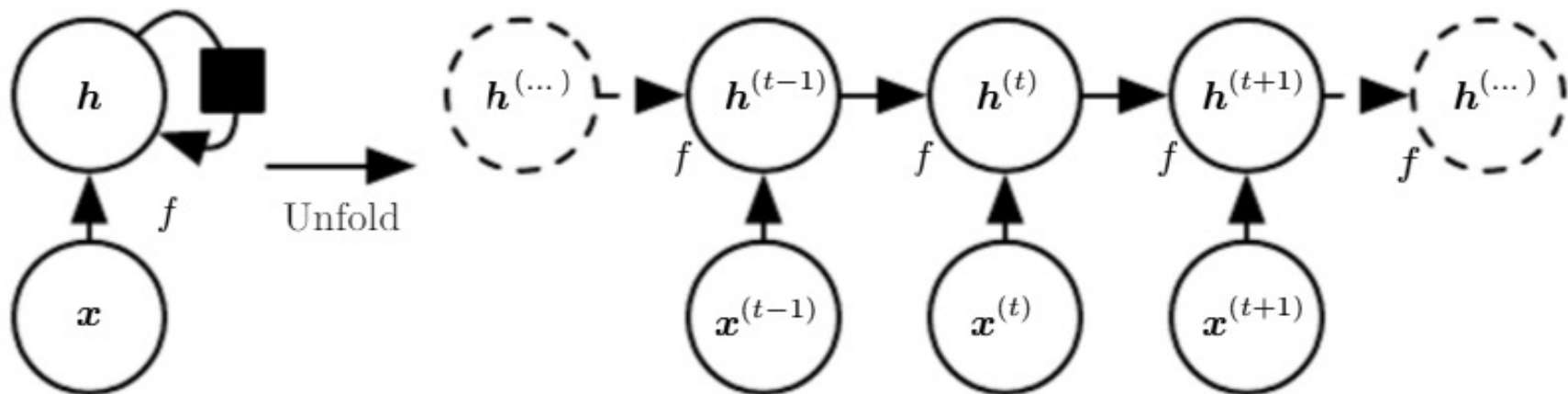
$$\Pr(s_t|o_{1..t}) \propto \Pr(o_t|s_t) \sum_{s_{t-1}} \Pr(s_t|s_{t-1}) \Pr(s_{t-1}|o_{1..t-1})$$



Recurrent Neural Network (RNN)

- In RNNs, outputs can be fed back to the network as inputs, creating a recurrent structure
- HMMs can be simulated and generalized by RNNs
- RNNs can be used for belief monitoring

\mathbf{x}_t : vector of observations \mathbf{h}_t : belief state



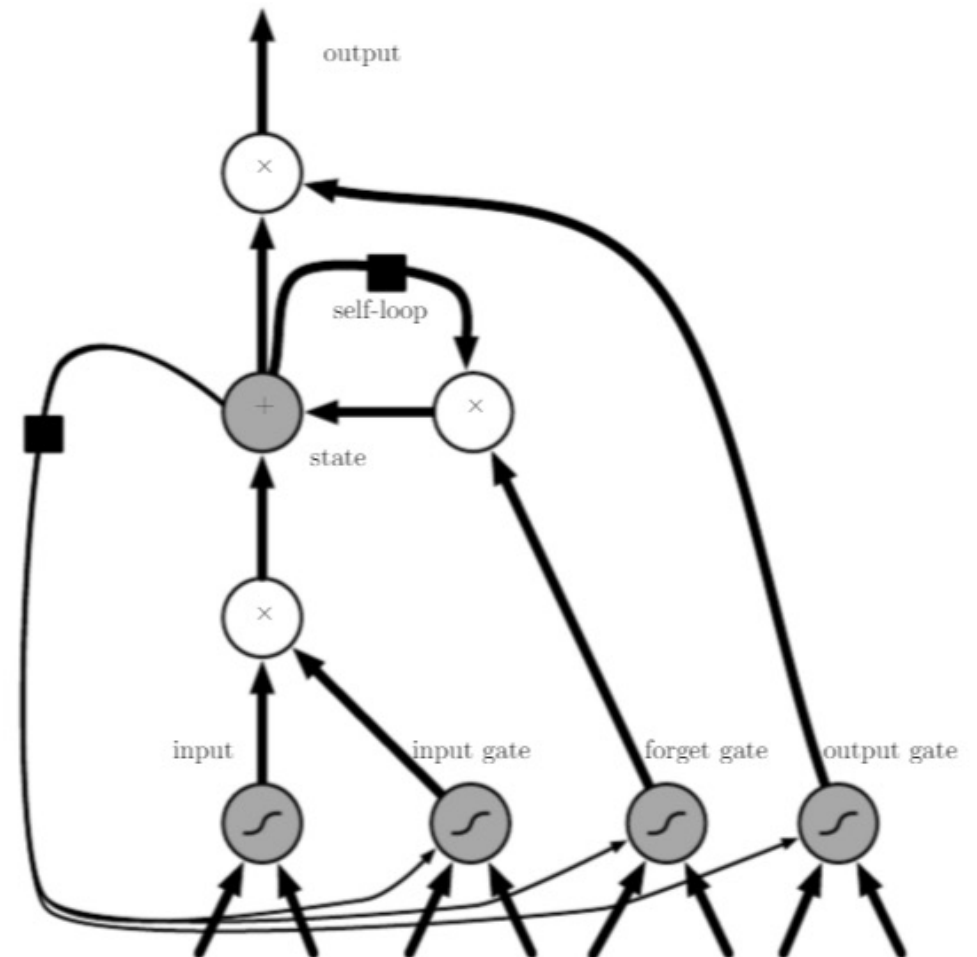
Training

- Recurrent neural networks are trained by backpropagation on the unrolled network
 - E.g., backpropagation through time
- Weight sharing:
 - Combine gradients of shared weights into a single gradient
- Challenges:
 - Gradient vanishing (and explosion)
 - Long range memory
 - Prediction drift

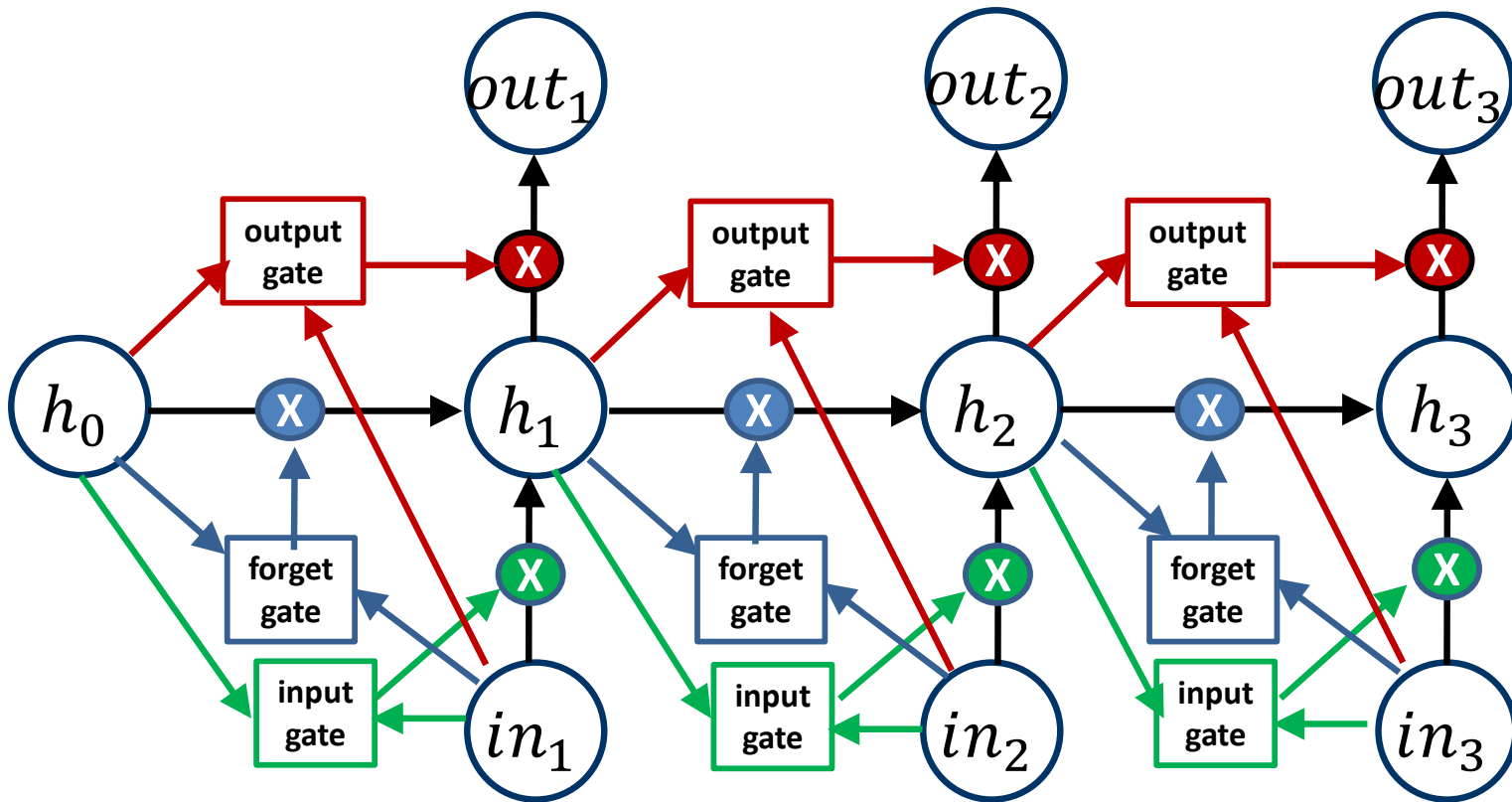


Long Short Term Memory (LSTM)

- Special gated structure to control memorization and forgetting in RNNs
- Mitigate gradient vanishing
- Facilitate long term memory

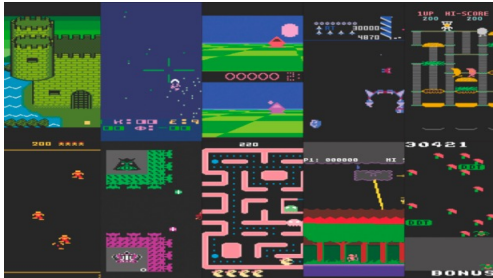


Unrolled long short term memory

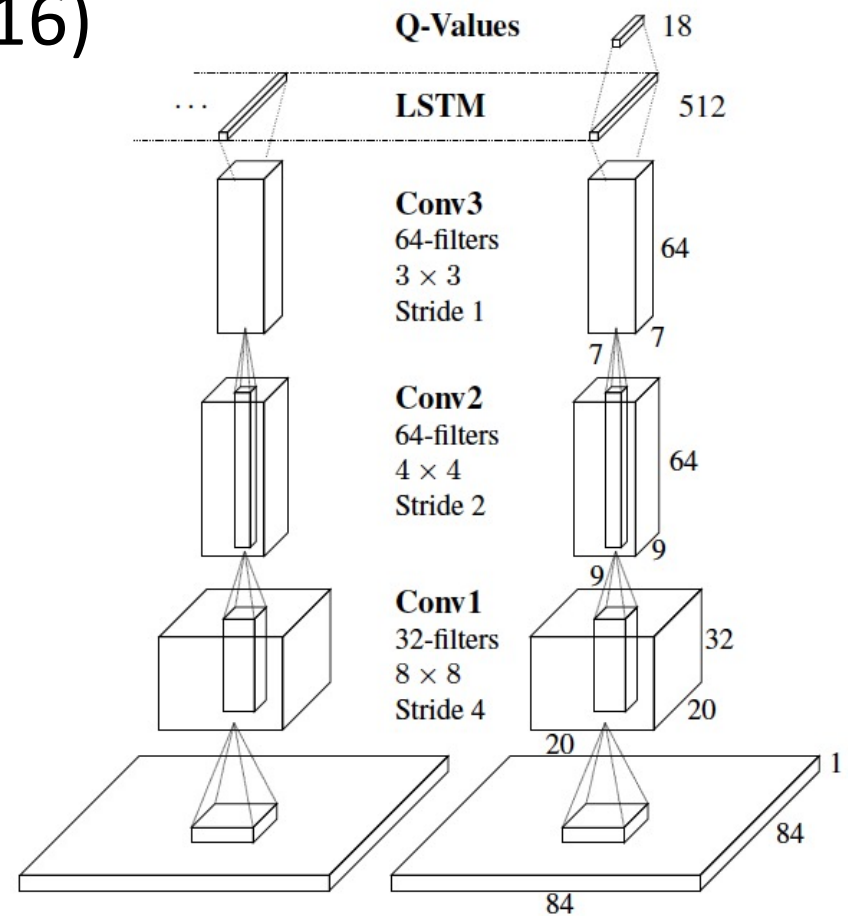


Deep Recurrent Q-Network

- Hausknecht and Stone (2016)
 - Atari games



- Transition model
 - LSTM network
- Observation model
 - Convolutional network



Deep Recurrent Q-Network

Initialize weights \mathbf{w} and $\bar{\mathbf{w}}$ at random

Observe current state s

Loop

Execute policy for entire episode

Add episode $(o_1, a_1, r_1, o_2, a_2, r_2, \dots, o_T, a_T, r_T)$ to experience buffer

Sample episode from buffer

Initialize h_0

For $t = 1$ till the end of the episode do

$$\frac{\partial \text{Err}}{\partial \mathbf{w}} = \left[Q_{\mathbf{w}}(\text{RNN}_{\mathbf{w}}(\hat{o}_{1..t}), \hat{a}_t) - \hat{r} - \gamma \max_{\hat{a}_{t+1}} Q_{\bar{\mathbf{w}}}(\text{RNN}_{\bar{\mathbf{w}}}(\hat{o}_{1..t+1}), \hat{a}_{t+1}) \right] \frac{\partial Q_{\mathbf{w}}(\text{RNN}_{\mathbf{w}}(\hat{o}_{1..t}), \hat{a}_t)}{\partial \mathbf{w}}$$

Update weights: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \text{Err}}{\partial \mathbf{w}}$

Every c steps, update target: $\bar{\mathbf{w}} \leftarrow \mathbf{w}$



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)

Table 1: On standard Atari games, DRQN performance parallels DQN, excelling in the games of Frostbite and Double Dunk, but struggling on Beam Rider. Bolded font indicates statistical significance between DRQN and our DQN.⁵

Flickering games
(missing observations)

