



**Carnegie Mellon University**  
Language Technologies Institute

# Intro to Reinforcement Learning Part II

11-777 Multimodal Machine Learning Fall 2020

**Paul Liang**

pliang@cs.cmu.edu

 @pliang279

# Admin

## Start working on midterm assignments now!

Instructions for the **midterm presentations** are on piazza resources:

- [https://piazza.com/class\\_profile/get\\_resource/kcncr11wq24q6z7/kg742kik5kv6tg](https://piazza.com/class_profile/get_resource/kcncr11wq24q6z7/kg742kik5kv6tg)
- Deadline for pre-recorded presentation: Friday, November 13th, 2020 at 8pm ET
- 7 minutes, mostly about error analysis and updated ideas, don't try to present everything...

Instructions for **midterm report** are also online:

- [https://piazza.com/class\\_profile/get\\_resource/kcncr11wq24q6z7/kgraer741fw3n4](https://piazza.com/class_profile/get_resource/kcncr11wq24q6z7/kgraer741fw3n4)
- Deadline: Sunday, November 15th, 2020
- 8 pages for teams of 3 and 9 pages for the other teams
- Multimodal baselines, error analysis, proposed ideas

# Admin

## Reading wildcard

- Each student gets one (1) wild card to be used as a way to extend by up to 24 hours their deadline for the summary deadline (which is usually Fridays at 8pm)
- See details on piazza

 note @369   

### Wild card for Reading Assignment summaries

Hello! Bonjour!

Based on your recent feedback and internal discussions, we decided to offer all students one wild card for the Reading Assignment

Each student gets one (1) wild card to be used as a way to extend by up to 24 hours their deadline for the summary deadline (which hours). There is no need to send a note via Piazza for this wild card. We will automatically use your wild card the first time you submit

We updated the syllabus with details about both wild card types (reading assignment summaries and project assignments). We hope

Best,  
LP

P.S. If you lost some points in previous weeks because of a late submission (less than 24 hours) of your reading assignment summary

[readings](#) [logistics](#)

[edit](#) · good note | 0

# Admin

## Piazza live Q&A

The screenshot displays the Piazza web interface for a course. The browser address bar shows the URL `piazza.com/class/kcnr11wq24q6z7?cid=43`. The page title is "Piazza" and the user is identified as "Louis-Philippe Morency". The interface includes a navigation bar with "LIVE Q&A" highlighted in a red box. Below the navigation bar, there are tabs for "Unread", "Updated", "Unresolved", and "Following". A "New Post" button is also highlighted in a red box. The main content area shows a list of posts, including a question titled "Question" with the text "When is the lecture starting?" and an answer titled "the instructors' answer, where instructors collectively construct a single answer" with the text "At 3:20pm EST". Both the question and answer have an "edit" button and a "good question/answer" indicator with a count of 0. The interface also includes a search bar and a "Ban User Console" section.

Please share your questions and comments on Piazza Live Q&A



Live responses by your TAs and follow-up by the instructor after the main lecture

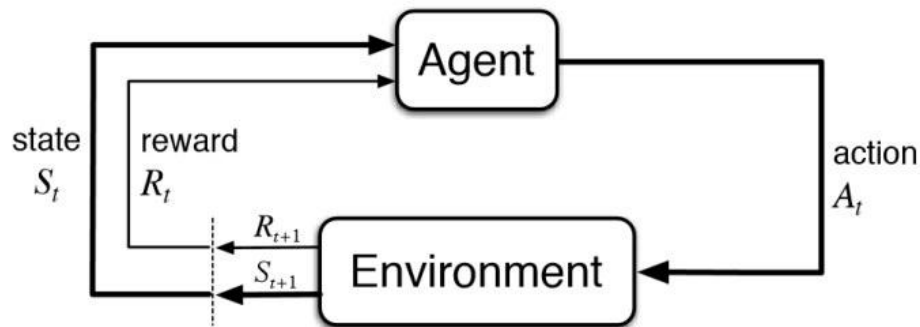
## Used Materials

Acknowledgement: Some of the material and slides for this lecture were borrowed from Pieter Abbeel, Yan Duan, Xi Chen, and Andrej Karpathy's Deep RL Bootcamp at UC Berkeley, as well as Katerina Fragkiadaki and Ruslan Salakhutdinov's 10-703 course at CMU, who in turn borrowed much from Rich Sutton's class and David Silver's class on Reinforcement Learning.

# Recap: Markov Decision Process (MDPs)

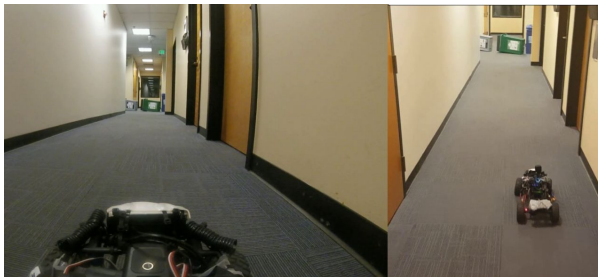
An MDP is defined by:

- Set of states  $S$
- Set of actions  $A$
- Transition function  $P(s' | s, a)$
- Reward function  $R(s, a, s')$
- Start state  $s_0$
- Discount factor  $\gamma$
- Horizon  $H$



**Trajectory**

$$S_0, a_0, r_0, S_1, a_1, r_1, S_2, a_2, r_2, \dots$$



# Return

We aim to maximize *total discounted reward*:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

**Discount  
factor**

$\gamma$  close to 0 leads to "myopic" evaluation

$\gamma$  close to 1 leads to "far-sighted" evaluation

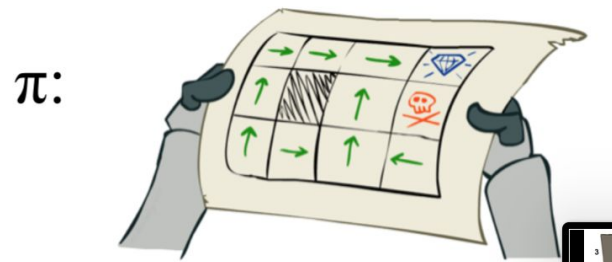
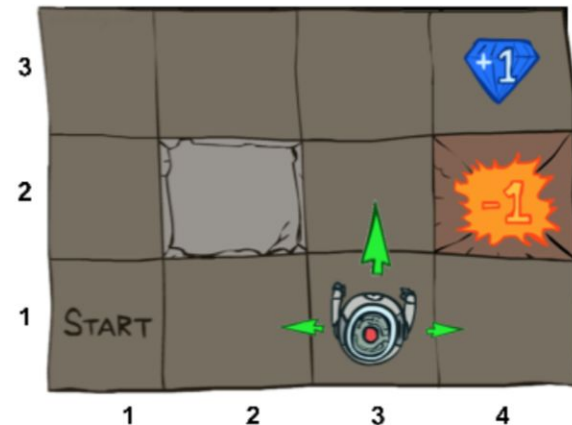
# Policy

**Definition:** A policy is a distribution over actions given states

$$\pi(a | s) = \mathbf{Pr}(A_t = a | S_t = s), \forall t$$

- A policy fully defines the behavior of an agent
- The policy is stationary (time-independent)
- During learning, the agent changes its policy as a result of experience

Special case: deterministic policies

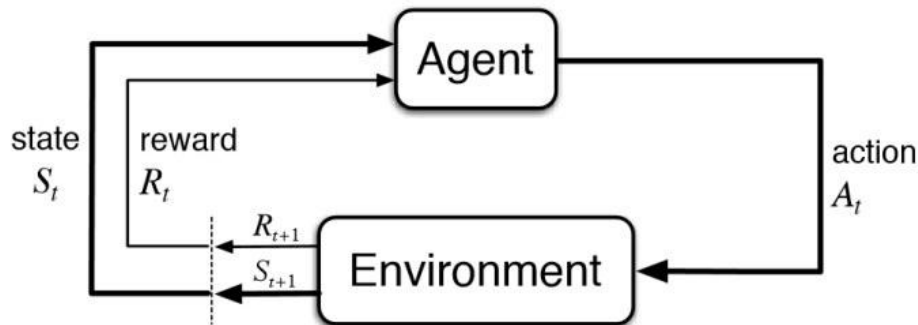




# Recap: MDPs, Returns, Policies

An MDP is defined by:

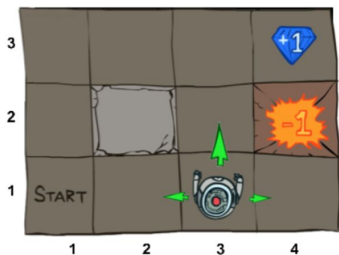
- Set of states  $S$
- Set of actions  $A$
- Transition function  $P(s' | s, a)$
- Reward function  $R(s, a, s')$
- Start state  $s_0$
- Discount factor  $\gamma$
- Horizon  $H$



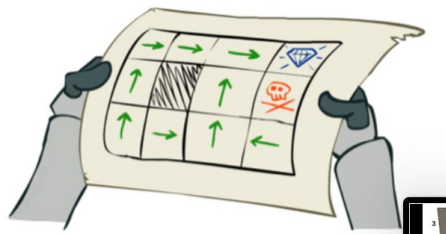
**Return:**

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

**Goal:**  $\arg \max_{\pi} \mathbb{E} \left[ \sum_{t=0}^H \gamma^t R_t | \pi \right]$



$\pi$ :



# Reinforcement Learning vs Supervised Learning

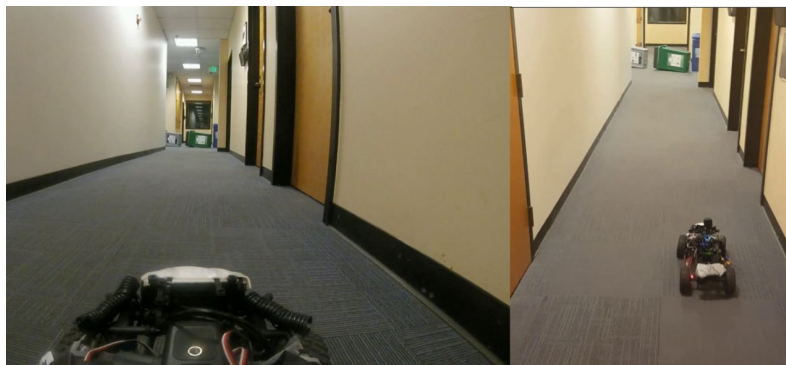
## Reinforcement Learning

- Sequential decision making
- Maximize cumulative reward
- Sparse rewards
- Environment maybe unknown



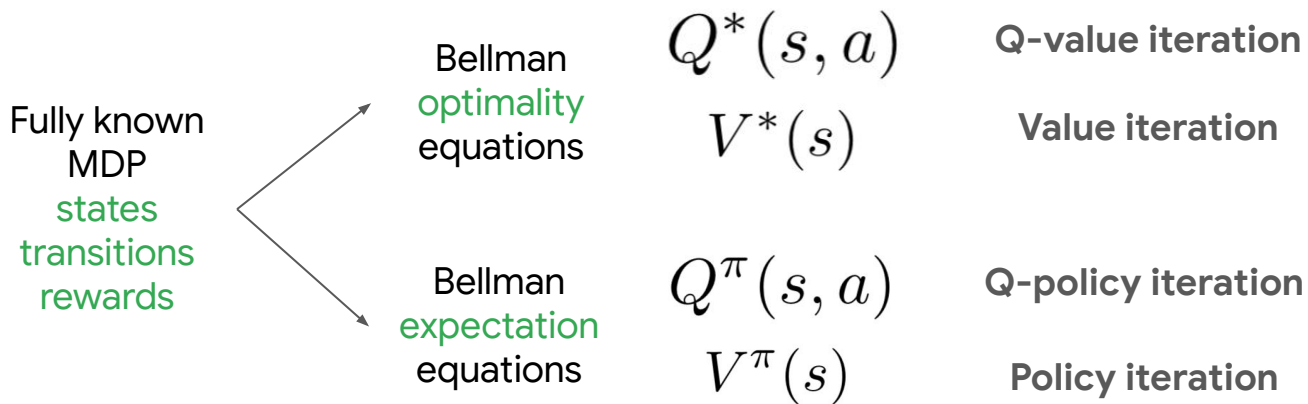
## Supervised Learning

- One-step decision making
- Maximize immediate reward
- Dense supervision
- Environment always known



# Recap: Exact methods

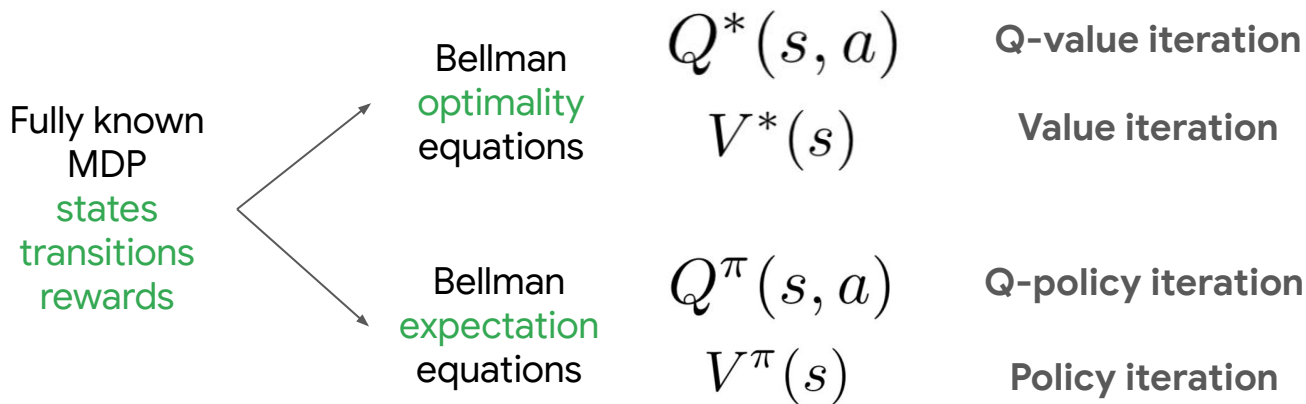
$$Q^*(s, a) = \mathbb{E}_{s'} \left[ r(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$



Repeat until policy converges. Guaranteed to converge to optimal policy.

## Recap: Exact methods

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



Repeat until policy converges. Guaranteed to converge to optimal policy.

**Iterate over and storage for all states and actions**  
**Requires small, discrete state and action space**  
**Update equations require fully observable MDP and known transitions**

# Recap: Tabular Q-learning

MDP  
with  
unknown  
transitions



Bellman  
optimality  
equations



Replace **true**  
expectation over  
transitions with  
estimates

**Tabular Q-learning**

$s' \sim P(s'|s, a)$  simulation and exploration, epsilon greedy is important!

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

# Recap: Tabular Q-learning

MDP  
with  
unknown  
transitions



Bellman  
optimality  
equations



Replace **true**  
expectation over  
transitions with  
estimates

**Tabular Q-learning**

$s' \sim P(s'|s, a)$  simulation and exploration, epsilon greedy is important!

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

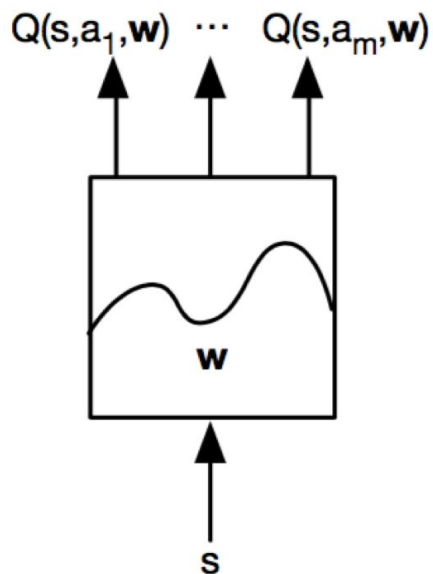
old estimate

target

$$Q_{k+1}(s, a) \leftarrow Q_k(s, a) + \alpha \left( r(s, a, s') + \gamma \max_{a'} Q_k(s', a') - Q_k(s, a) \right)$$

**Tabular: keep a  $|S| \times |A|$  table of  $Q(s,a)$**   
**Still requires small and discrete state and action space**  
**How can we generalize to unseen states?**

# Recap: Deep Q-learning



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

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s, a, r, s' \sim \mathcal{D}_i} \left[ \left( \underbrace{r + \gamma \max_{a'} Q(s', a'; w_i^-)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right)^2 \right]$$

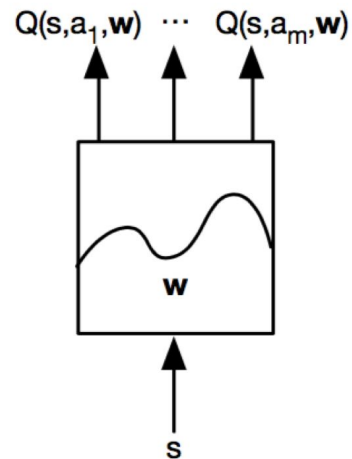
**Correlated samples + non-stationary targets**

# Recap: Deep Q-learning

- Sample random mini-batch of transitions  $(s, a, r, s')$  from  $D$
- Compute Q-learning targets w.r.t. old fixed parameters  $\mathbf{w}$
- Optimize MSE between Q-network and Q-learning targets

$s_1, a_1, r_2, s_2$
$s_2, a_2, r_3, s_3$
$s_3, a_3, r_4, s_4$
...
$s_t, a_t, r_{t+1}, s_{t+1}$

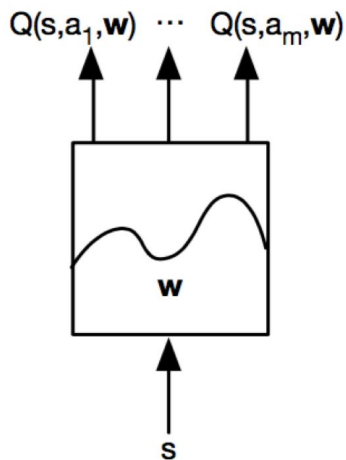
$$\mathcal{L}_i(\mathbf{w}_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \underbrace{\left( r + \gamma \max_{a'} Q(s', a'; \mathbf{w}_i^-) \right)}_{\text{Q-learning target}} - \underbrace{Q(s, a; \mathbf{w}_i)}_{\text{Q-network}} \right]^2$$



- Use stochastic gradient descent
- Update  $\mathbf{w}$  with updated  $\mathbf{w}$  every  $\sim 1000$  iterations



# Recap: Deep Q-learning



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

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s, a, r, s' \sim \mathcal{D}_i} \left[ \left( \underbrace{r + \gamma \max_{a'} Q(s', a'; w_i^-)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right)^2 \right]$$

Stochastic gradient descent + Exploration + Experience replay + Fixed Q-targets

Works for high-dimensional state and action spaces  
Generalizes to unseen states

## Recap: Obtaining the optimal policy

Optimal policy can be found by maximizing over  $Q^*(s,a)$

$$\pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a Q^*(s, a) \\ \epsilon, & \text{else} \end{cases}$$

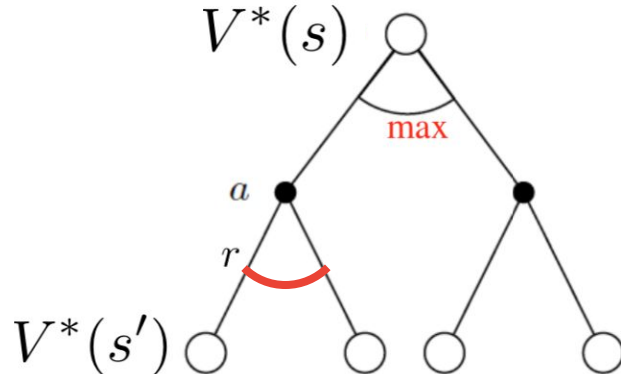
# Recap: Obtaining the optimal policy

Optimal policy can be found by maximizing over  $Q^*(s,a)$

$$\pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a Q^*(s, a) \\ \epsilon, & \text{else} \end{cases}$$

Optimal policy can also be found by maximizing over  $V^*(s')$  with **one-step look ahead**

$$\pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a \mathbb{E}_{s'} [r(s, a, s') + \gamma V^*(s')] \\ \epsilon, & \text{else} \end{cases}$$



# Contents

- Policy gradient methods
- Actor-critic
- Applications: Language and RL
- Applications: RL for language (e.g. text generation)

# Value-based and Policy-based RL

## ▶ Value Based

- Learned Value Function
- Implicit policy (e.g.  $\epsilon$ -greedy)

### State value functions

$$V^\pi(s)$$

$$V^*(s)$$

### Action value functions

$$Q^\pi(s, a)$$

$$Q^*(s, a)$$

$$\pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a \mathbb{E}_{s'} [r(s, a, s') + \gamma V^*(s')] \\ \epsilon, & \text{else} \end{cases} \quad \pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a Q^*(s, a) \\ \epsilon, & \text{else} \end{cases}$$

# Value-based and Policy-based RL

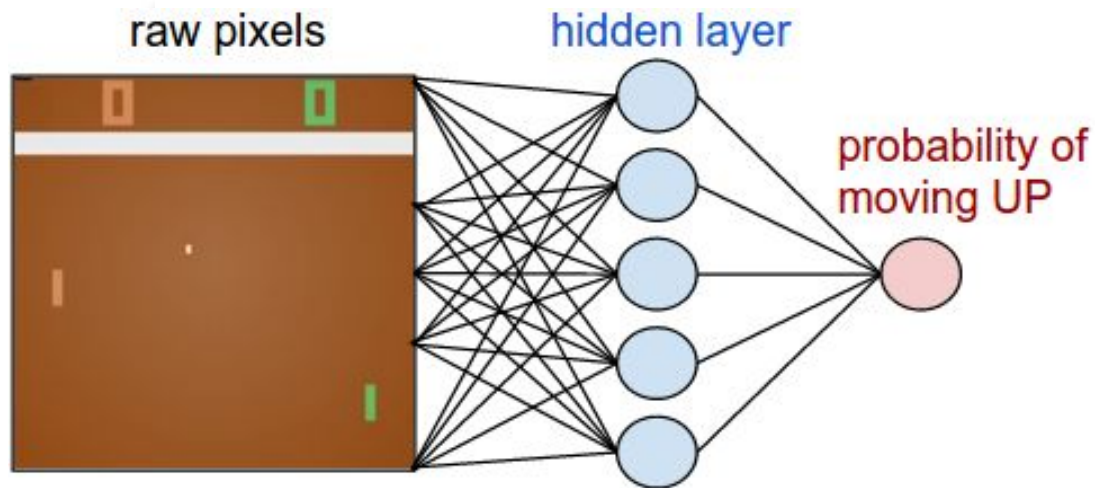
- ▶ Value Based

- Learned Value Function
- Implicit policy (e.g.  $\epsilon$ -greedy)

- ▶ Policy Based

- No Value Function
- Learned Policy

$$\pi_{\theta}(s, a) = \mathbb{P}[a \mid s, \theta]$$



# Directly learning the policy

- Often  $\pi$  can be simpler than Q or V

- E.g., robotic grasp

Q(s,a) and V(s) very high-dimensional  
But policy could be just 'open/close hand'

# Directly learning the policy

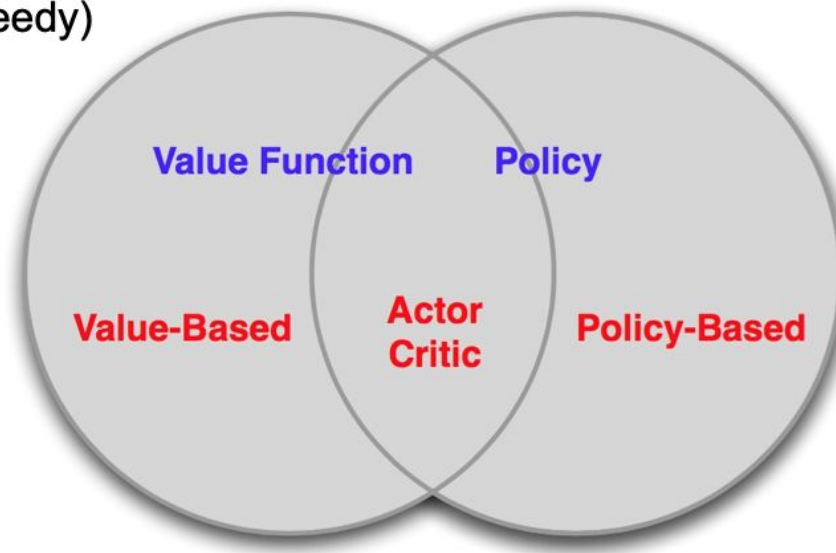
- Often  $\pi$  can be simpler than Q or V  $Q(s,a)$  and  $V(s)$  very high-dimensional  
But policy could be just 'open/close hand'
  - E.g., robotic grasp
- V: doesn't prescribe actions
  - Would need dynamics model (+ compute 1 Bellman back-up)
- Q: need to be able to efficiently solve  $\arg \max_u Q_\theta(s, u)$ 
  - Challenge for continuous / high-dimensional action spaces\*

$$\pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a \mathbb{E}_{s'} [r(s, a, s') + \gamma V^*(s')] \\ \epsilon, & \text{else} \end{cases} \quad \pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_a Q^*(s, a) \\ \epsilon, & \text{else} \end{cases}$$



# Value-based and Policy-based RL

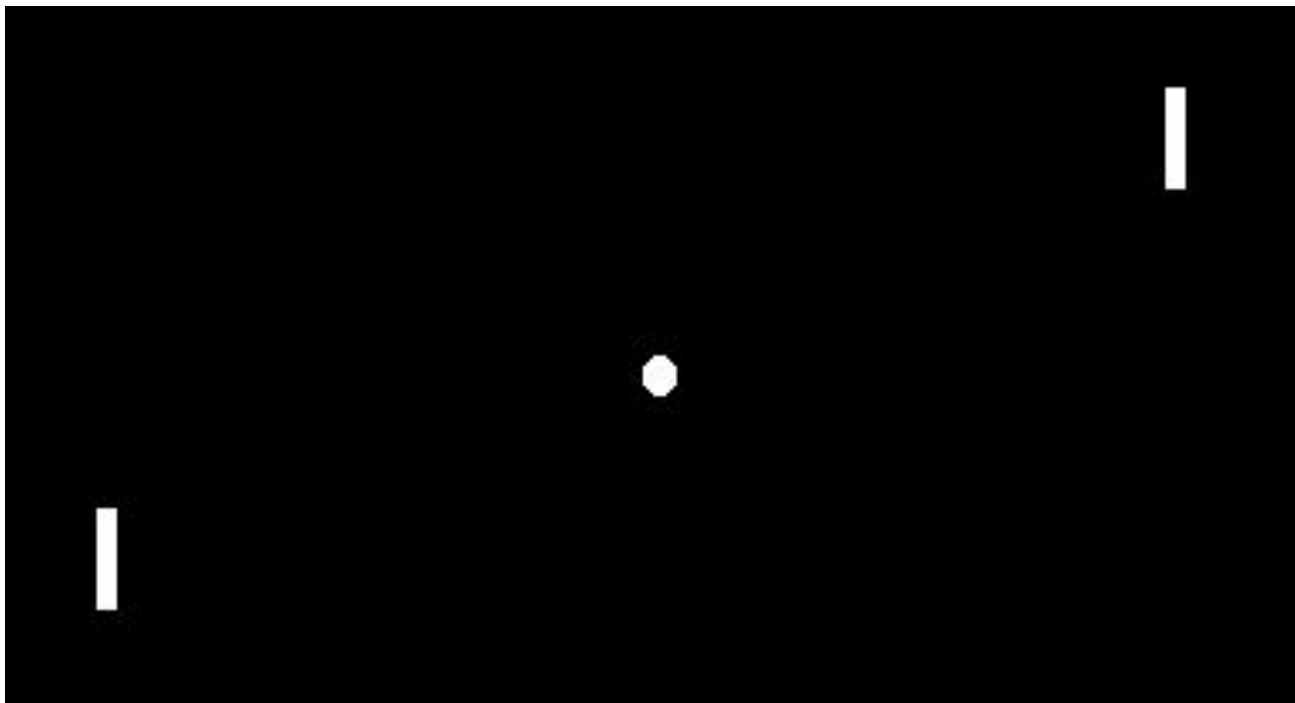
- ▶ **Value Based**
  - Learned Value Function
  - Implicit policy (e.g.  $\epsilon$ -greedy)
- ▶ **Policy Based**
  - No Value Function
  - Learned Policy
- ▶ **Actor-Critic**
  - Learned Value Function
  - Learned Policy



# Value-based and Policy-based RL

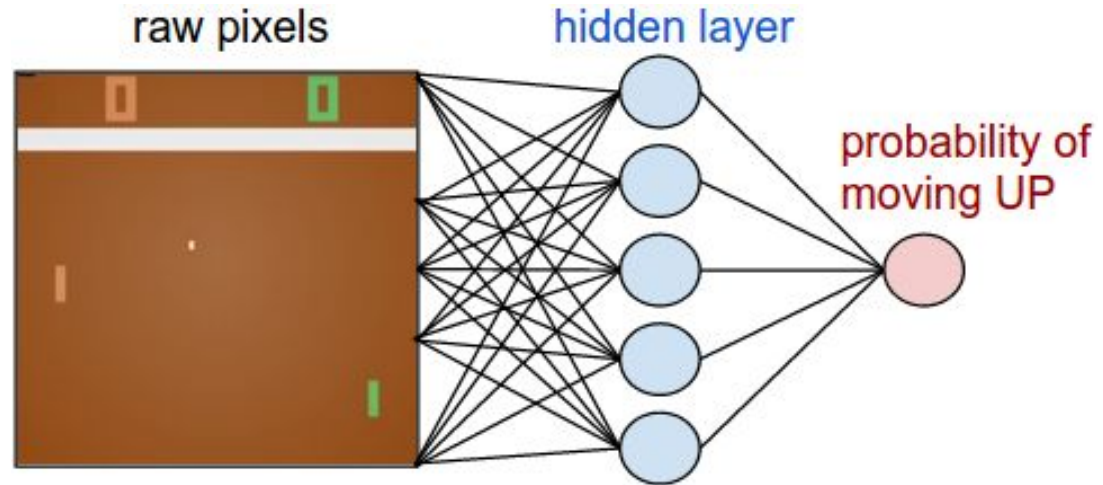
	Policy-based	Value-based
■ Conceptually:	Optimize what you care about	Indirect, exploit the problem structure, self-consistency
■ Empirically:	More compatible with rich architectures (including recurrence)  More versatile  More compatible with auxiliary objectives	More compatible with exploration and off-policy learning  More sample-efficient when they work

# Pong from pixels



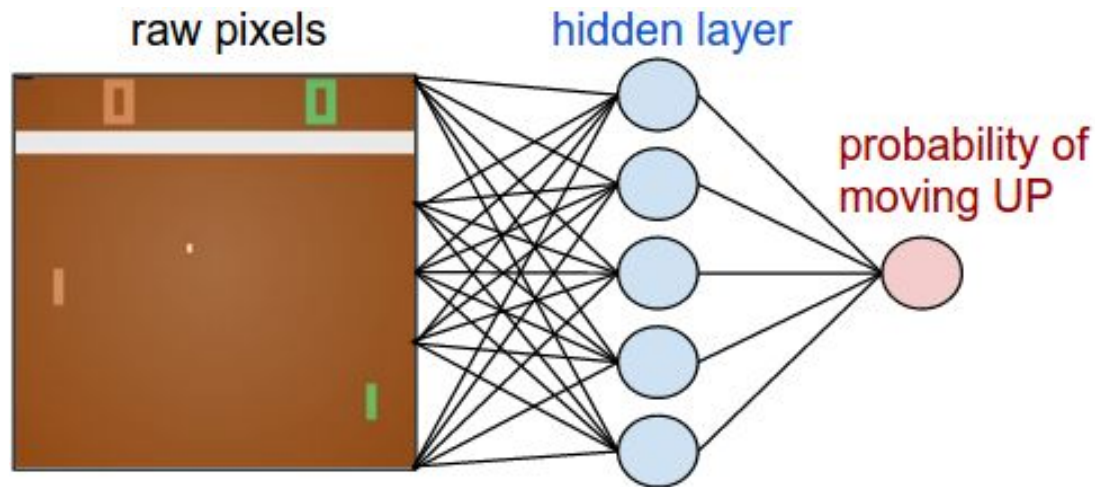
# Pong from pixels

e.g.,  
height width  
[80 x 80]  
array of



# Pong from pixels

e.g.,  
height width  
[80 x 80]  
array of



**Network sees +1 if it scored a point, and -1 if it was scored against.  
How do we learn these parameters?**

# Pong from pixels

Suppose we had the training labels...  
(we know what to do in any state)

```
(x1,UP)  
(x2,DOWN)  
(x3,UP)  
...
```

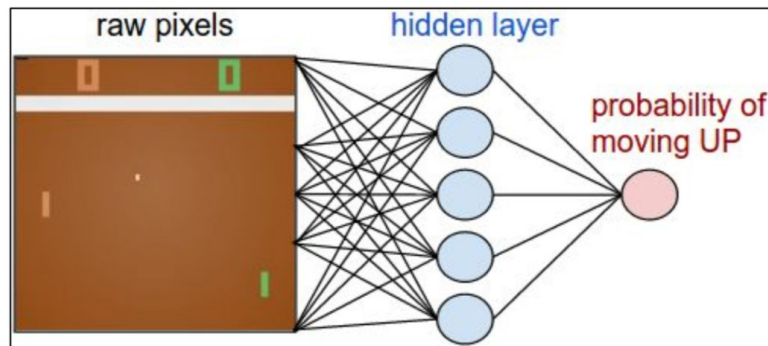
# Pong from pixels

Suppose we had the training labels...  
(we know what to do in any state)

(x1,UP)  
(x2,DOWN)  
(x3,UP)  
...

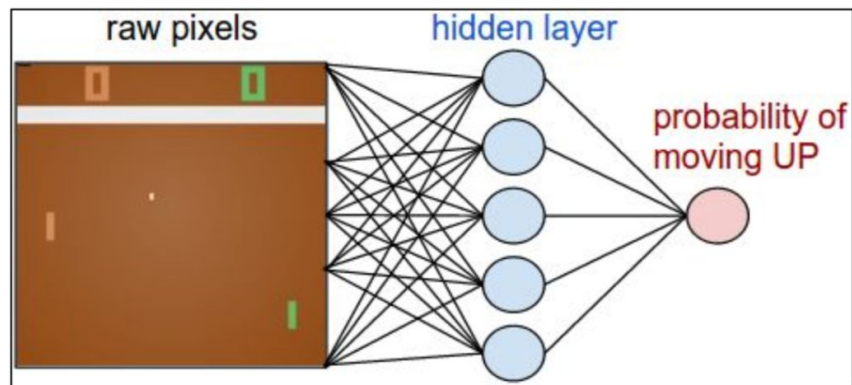
maximize:

$$\sum_i \log p(y_i | x_i)$$



# Pong from pixels

Except, we don't have labels...

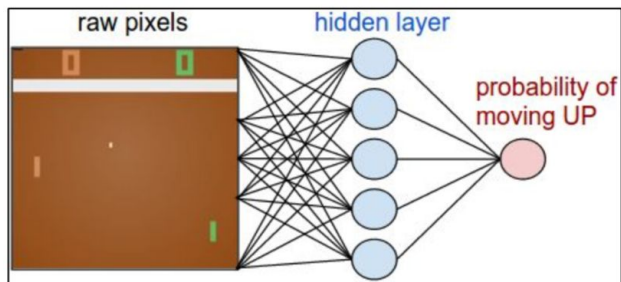


Should we go UP or DOWN?

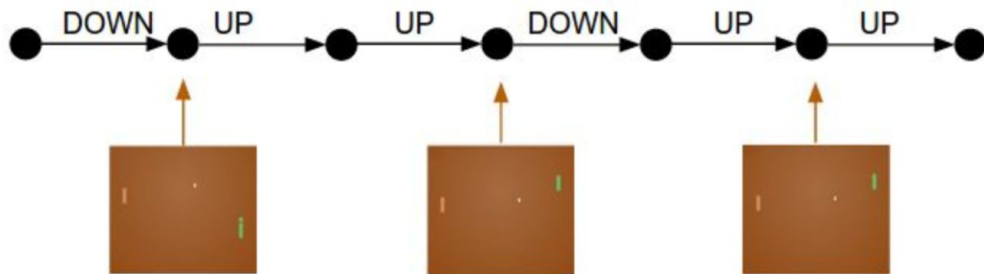


# Pong from pixels

Let's just act according to our current policy...



Rollout the policy and collect an episode

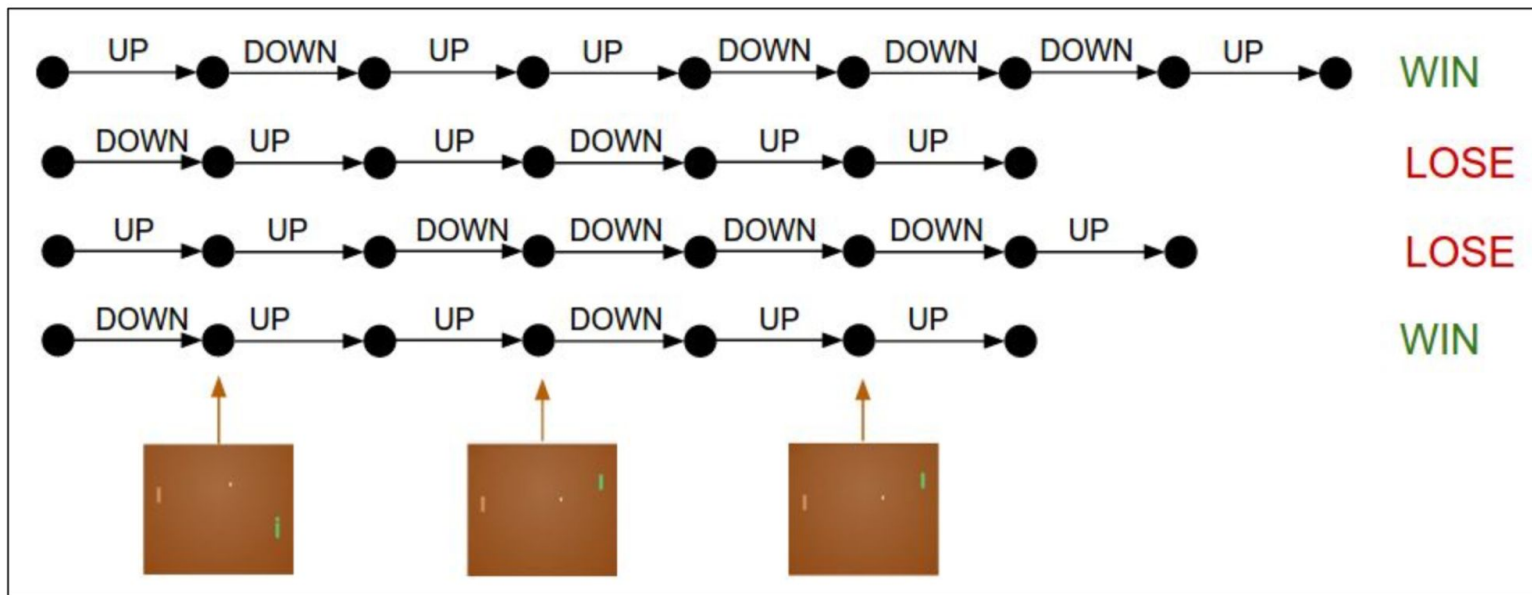


WIN

# Pong from pixels

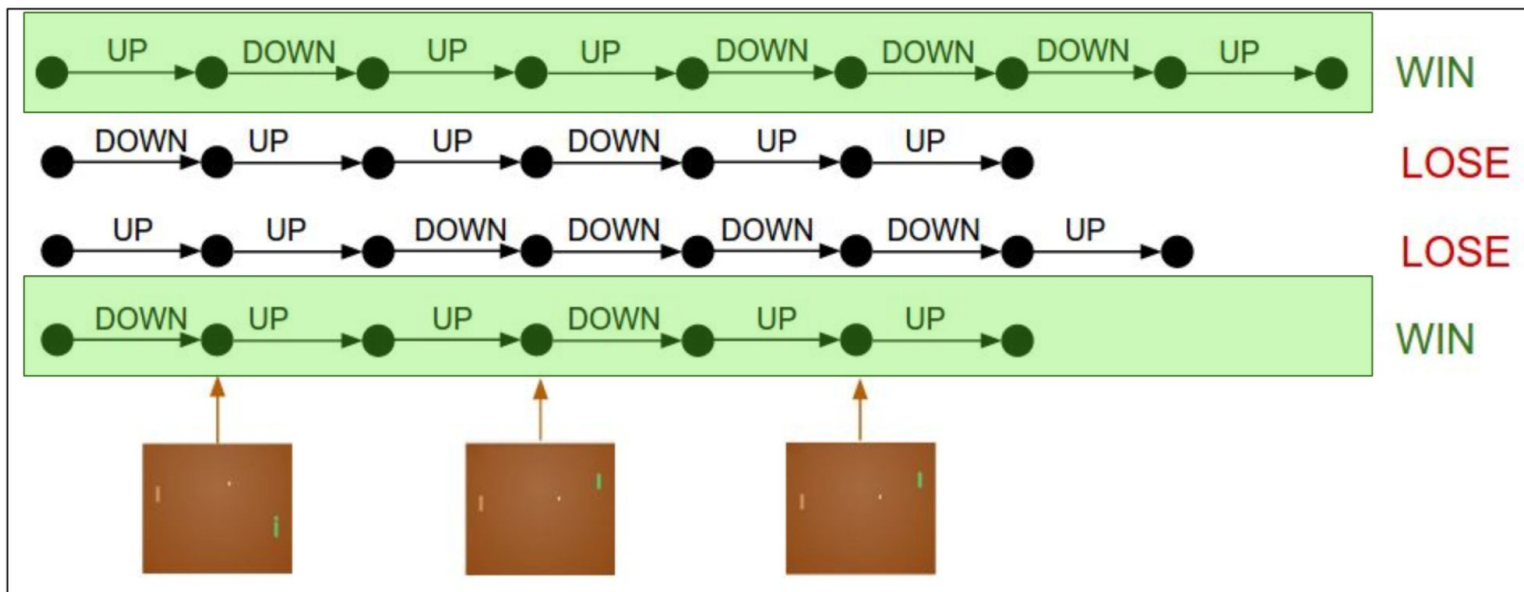
Collect many rollouts...

**4 rollouts:**



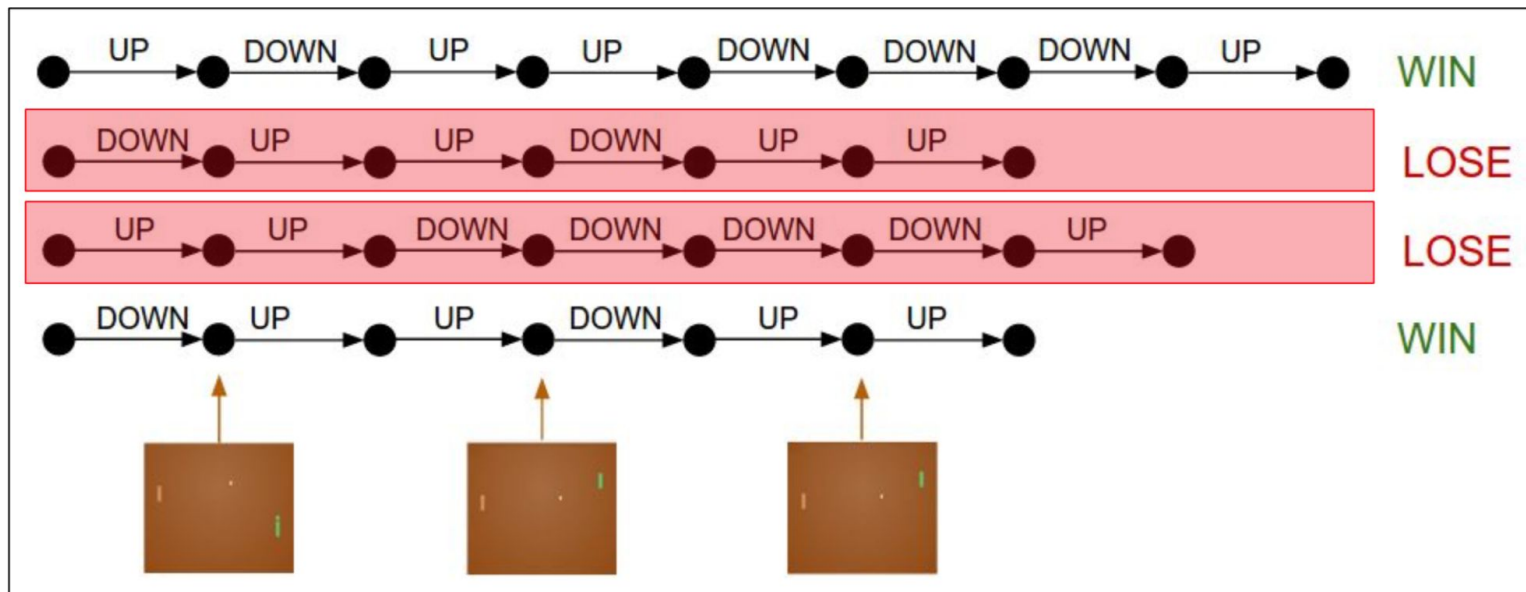
# Pong from pixels

Not sure whatever we did here, but apparently it was good.



# Pong from pixels

Not sure whatever we did here, but it was bad.



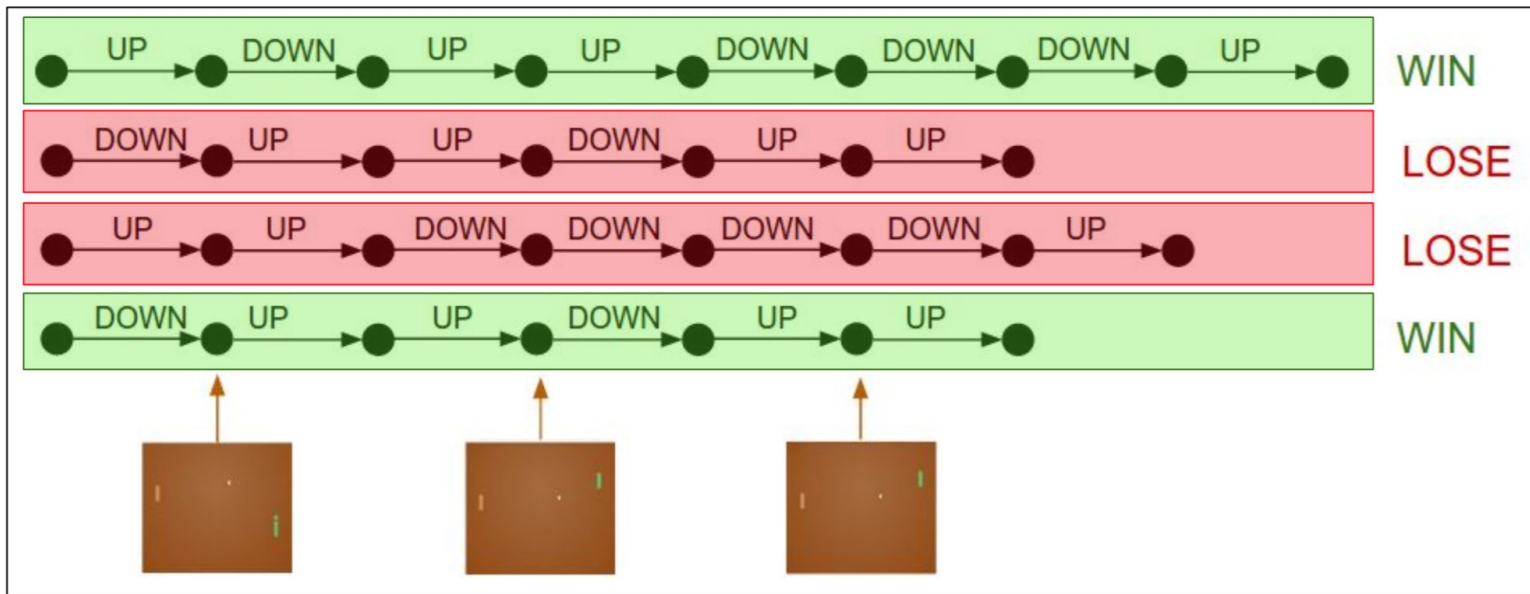
# Pong from pixels

Pretend every action we took here was the correct label.

maximize:  $\log p(y_i | x_i)$

Pretend every action we took here was the wrong label.

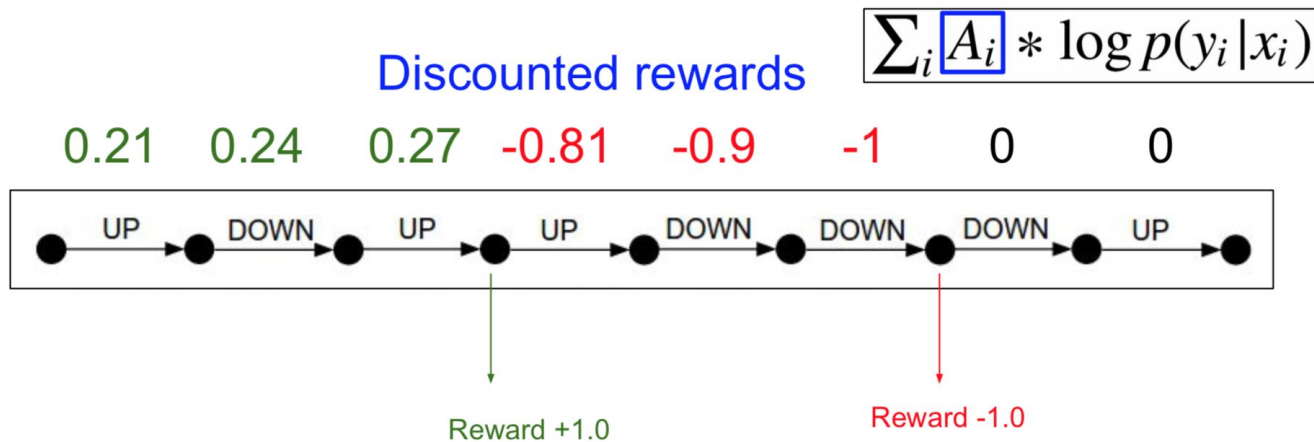
maximize:  $(-1) * \log p(y_i | x_i)$



# Pong from pixels

## Discounting

Blame each action assuming that its effects have exponentially decaying impact into the future.

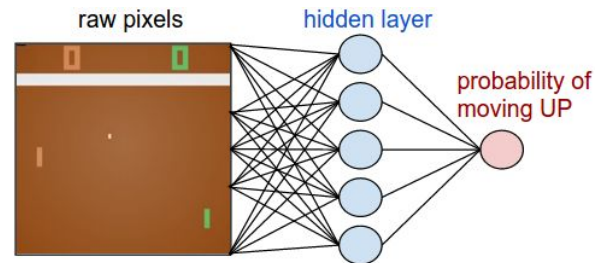


$\gamma = 0.9$

# Pong from pixels

1. Initialize a policy network at random

$$\pi(a | s)$$



# Pong from pixels

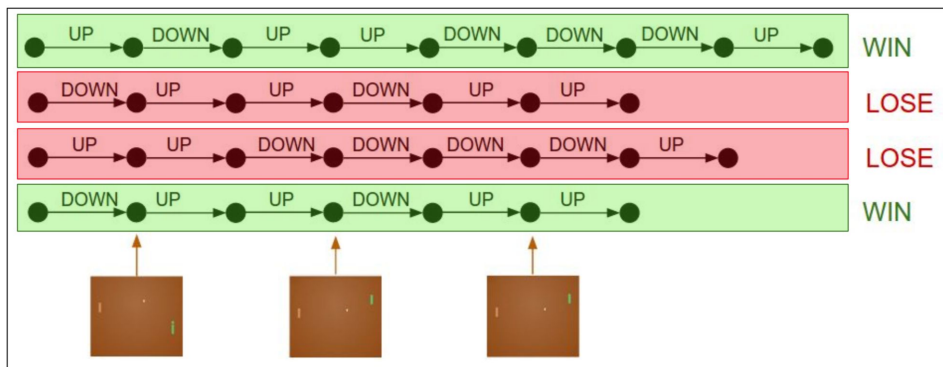
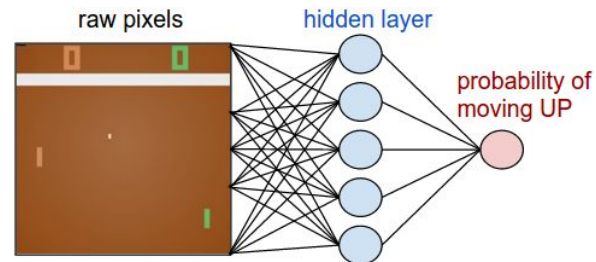
1. Initialize a policy network at random

2. **Repeat Forever:**

3. Collect a bunch of rollouts with the policy

**epsilon greedy!**

$$\pi(a | s)$$

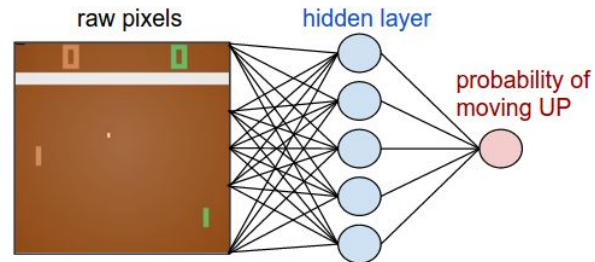




# Pong from pixels

1. Initialize a policy network at random
2. **Repeat Forever:**
3. Collect a bunch of rollouts with the policy **epsilon greedy!**
4. Increase the probability of actions that worked well

$$\pi(a | s)$$



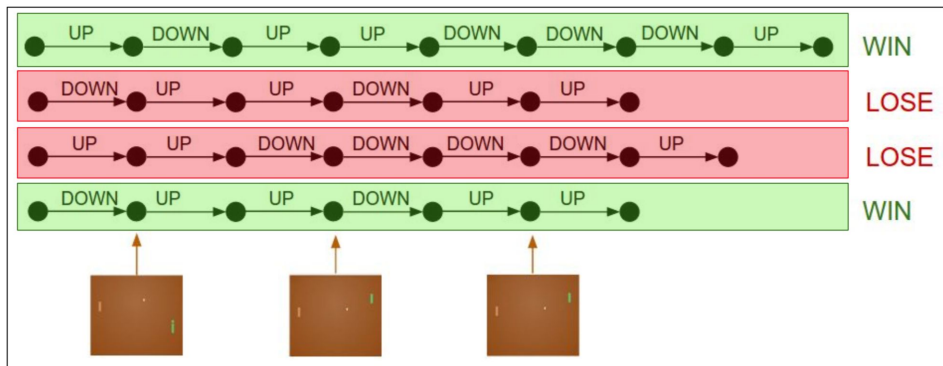
Pretend every action we took here was the correct label.

$$\text{maximize: } \log p(y_i | x_i)$$

Pretend every action we took here was the wrong label.

$$\text{maximize: } (-1) * \log p(y_i | x_i)$$

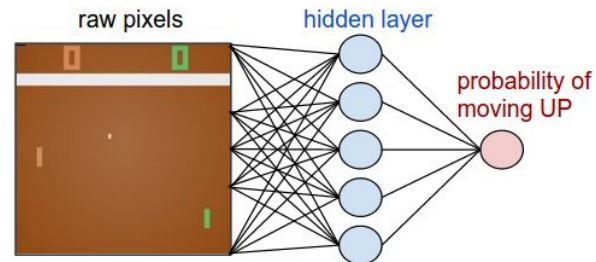
$$\sum_i A_i * \log p(y_i | x_i)$$



# Pong from pixels

1. Initialize a policy network at random
2. **Repeat Forever:**
3. Collect a bunch of rollouts with the policy **epsilon greedy!**
4. Increase the probability of actions that worked well

$$\pi(a | s)$$

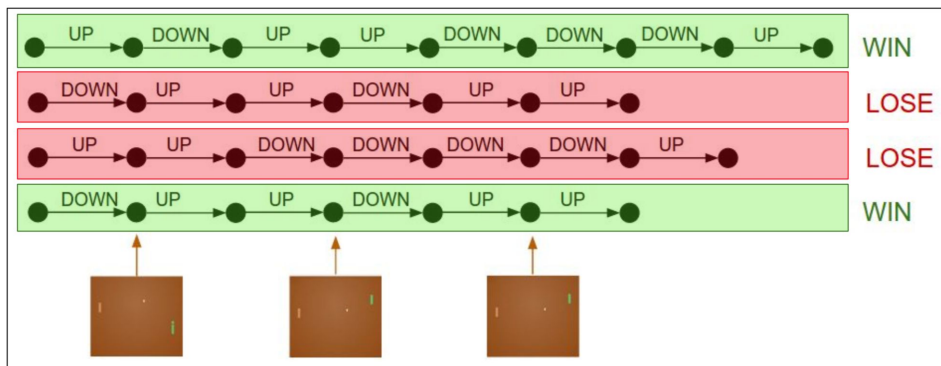


Pretend every action we took here was the correct label.

$$\text{maximize: } \log p(y_i | x_i)$$

Pretend every action we took here was the wrong label.

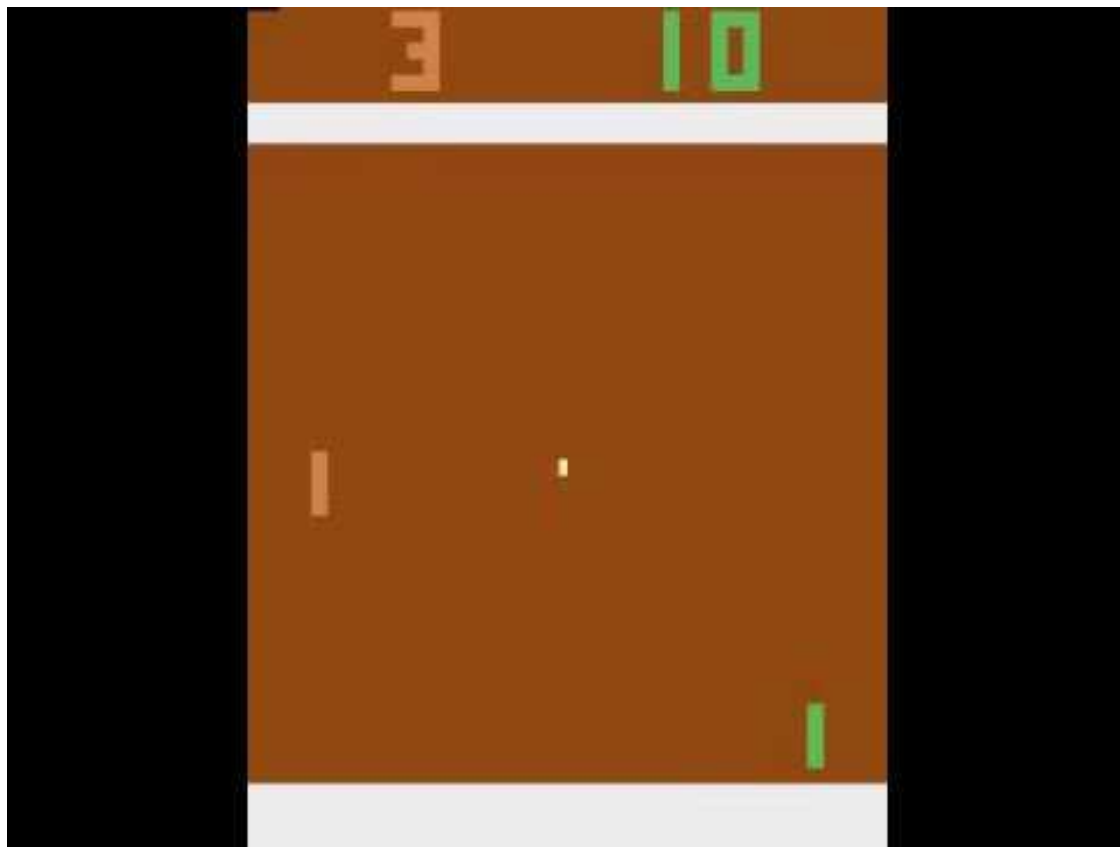
$$\text{maximize: } (-1) * \log p(y_i | x_i)$$



$$\sum_i A_i * \log p(y_i | x_i)$$

Does not require transition probabilities  
Does not estimate Q(), V()  
Predicts policy directly

# Pong from pixels



# Policy gradients

## Why does this work?

1. Initialize a policy network at random
2. **Repeat Forever:**
3.       Collect a bunch of rollouts with the policy
4.       Increase the probability of actions that worked well

$$\sum_i A_i * \log p(y_i | x_i)$$

# Policy gradients

Formally, let's define a class of parameterized policies  $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid \pi_\theta \right]$$

# Policy gradients

Writing in terms of trajectories  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$

Probability of a trajectory

Reward of a trajectory

$$\begin{aligned} p(\tau; \theta) &= \pi_{\theta}(a_0|s_0)p(s_1|s_0, a_0) \\ &\times \pi_{\theta}(a_1|s_1)p(s_2|s_1, a_1) \\ &\times \pi_{\theta}(a_2|s_2)p(s_3|s_2, a_2) \\ &\times \dots \\ &= \prod_{t \geq 0} p(s_{t+1}|s_t, a_t)\pi_{\theta}(a_t|s_t) \end{aligned}$$

$$r(\tau) = \sum_{t \geq 0} \gamma^t r_t$$

# Policy gradients

Writing in terms of trajectories  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$

Probability of a trajectory

Reward of a trajectory

$$\begin{aligned} p(\tau; \theta) &= \pi_{\theta}(a_0 | s_0) p(s_1 | s_0, a_0) \\ &\times \pi_{\theta}(a_1 | s_1) p(s_2 | s_1, a_1) \\ &\times \pi_{\theta}(a_2 | s_2) p(s_3 | s_2, a_2) \\ &\times \dots \\ &= \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t) \end{aligned}$$

$$r(\tau) = \sum_{t \geq 0} \gamma^t r_t$$

$$J(\theta) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | \pi_{\theta} \right] = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$$

# Policy gradients

Formally, let's define a class of parameterized policies  $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t \mid \pi_\theta \right]$$

We want to find the optimal policy

How can we do this?

$$\theta^* = \arg \max_{\theta} J(\theta)$$

**Gradient ascent on policy parameters**



# REINFORCE algorithm

Expected reward:  $J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$

$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

# REINFORCE algorithm

Expected reward:  $J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$

$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

Now let's differentiate this:  $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$

Intractable! Gradient of an expectation is problematic when p depends on  $\theta$

# REINFORCE algorithm

Expected reward:  $J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$

$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

Now let's differentiate this:  $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$

Intractable! Gradient of an expectation is problematic when  $p$  depends on  $\theta$

However, we can use a nice trick:  $\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$

# REINFORCE algorithm

Expected reward:  $J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$

$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

Now let's differentiate this:  $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$

Intractable! Gradient of an expectation is problematic when p depends on  $\theta$

However, we can use a nice trick:  $\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$

If we inject this back:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau \\ &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)] \end{aligned}$$

**Tractable :-)**

# REINFORCE algorithm

Can we compute these without knowing the transition probabilities?

We have:  $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$

# REINFORCE algorithm

Can we compute these without knowing the transition probabilities?

We have: 
$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

Thus: 
$$\log p(\tau; \theta) = \sum_{t \geq 0} (\log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t))$$

# REINFORCE algorithm

Can we compute these without knowing the transition probabilities?

We have: 
$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

Thus: 
$$\log p(\tau; \theta) = \sum_{t \geq 0} (\log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t))$$

And when differentiating: 
$$\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Doesn't depend on  
transition probabilities

# REINFORCE algorithm

Can we compute these without knowing the transition probabilities?

We have: 
$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

Thus: 
$$\log p(\tau; \theta) = \sum_{t \geq 0} (\log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t))$$

And when differentiating: 
$$\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Doesn't depend on transition probabilities

Therefore when sampling a trajectory, we can estimate gradients:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)] \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$



# Intuition

Gradient estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If **r(trajecory)** is high, push up the probabilities of the actions seen
- If **r(trajecory)** is low, push down the probabilities of the actions seen

# Intuition

Gradient estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

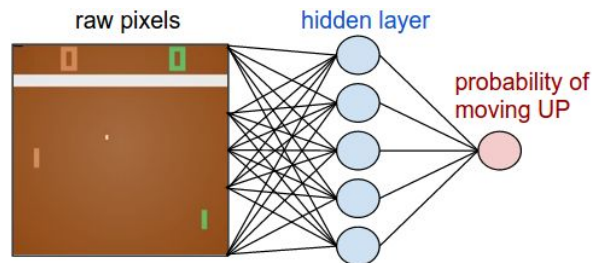
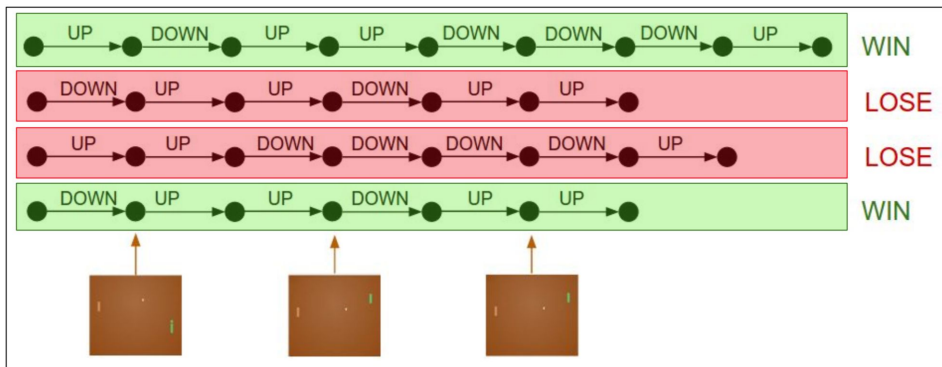
- If **r(trajectory)** is high, push up the probabilities of the actions seen
- If **r(trajectory)** is low, push down the probabilities of the actions seen

Pretend every action we took here was the correct label.

maximize:  $\log p(y_i | x_i)$

Pretend every action we took here was the wrong label.

maximize:  $(-1) * \log p(y_i | x_i)$



$$\sum_i A_i * \log p(y_i | x_i)$$

# Intuition

Gradient estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If **r(trajecory)** is high, push up the probabilities of the actions seen
- If **r(trajecory)** is low, push down the probabilities of the actions seen

## REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \theta), \forall a \in \mathcal{A}, s \in \mathcal{S}, \theta \in \mathbb{R}^n$

Initialize policy weights  $\theta$

Repeat forever:

Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

For each step of the episode  $t = 0, \dots, T - 1$ :

$G_t \leftarrow$  return from step  $t$

$\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi(A_t | S_t, \theta)$

# Intuition

Gradient estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If **r(trajecory)** is high, push up the probabilities of the actions seen
- If **r(trajecory)** is low, push down the probabilities of the actions seen

## REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization  $\pi(a|s, \theta), \forall a \in \mathcal{A}, s \in \mathcal{S}, \theta \in \mathbb{R}^n$

Initialize policy weights  $\theta$

Repeat forever:

Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$  following  $\pi(\cdot | \cdot, \theta)$

For each step of the episode  $t = 0, \dots, T - 1$ :

$G_t \leftarrow$  return from step  $t$

$\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi(A_t | S_t, \theta)$

**epsilon greedy**

# Intuition

Gradient estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If **r(trajectory)** is high, push up the probabilities of the actions seen
- If **r(trajectory)** is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. **But in expectation, it averages out!**

However, this also suffers from high variance because credit assignment is really hard - can we help this estimator?

# Variance reduction with a baseline

**Problem:** The raw reward of a trajectory isn't necessarily meaningful. E.g. if all rewards are positive, you keep pushing up probabilities of all actions.

**What is important then?** Whether a reward is higher or lower than what you expect to get.

# Variance reduction with a baseline

**Problem:** The raw reward of a trajectory isn't necessarily meaningful. E.g. if all rewards are positive, you keep pushing up probabilities of all actions.

**What is important then?** Whether a reward is higher or lower than what you expect to get.

**Idea:** Introduce a baseline function dependent on the state, which gives us an estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} (r(\tau) - b(s_t)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

e.g. exponential moving average of the rewards. Provably **reduces variance** while **remaining unbiased**.

# Actor-critic methods

A better baseline: want to push the probability of an action from a state, if this action was better than the expected value of what we should get from that state

Recall: **Q** and **V** - **action value** and **state value functions!**



# Actor-critic methods

A better baseline: want to push the probability of an action from a state, if this action was better than the expected value of what we should get from that state

Recall: **Q** and **V** - **action value** and **state value functions!**

We are happy with an action **a** in a state **s** if **Q(s,a) - V(s)** is large.  
Otherwise we are unhappy with an action if it's small.

Using this, we get the estimator:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

# Actor-critic methods

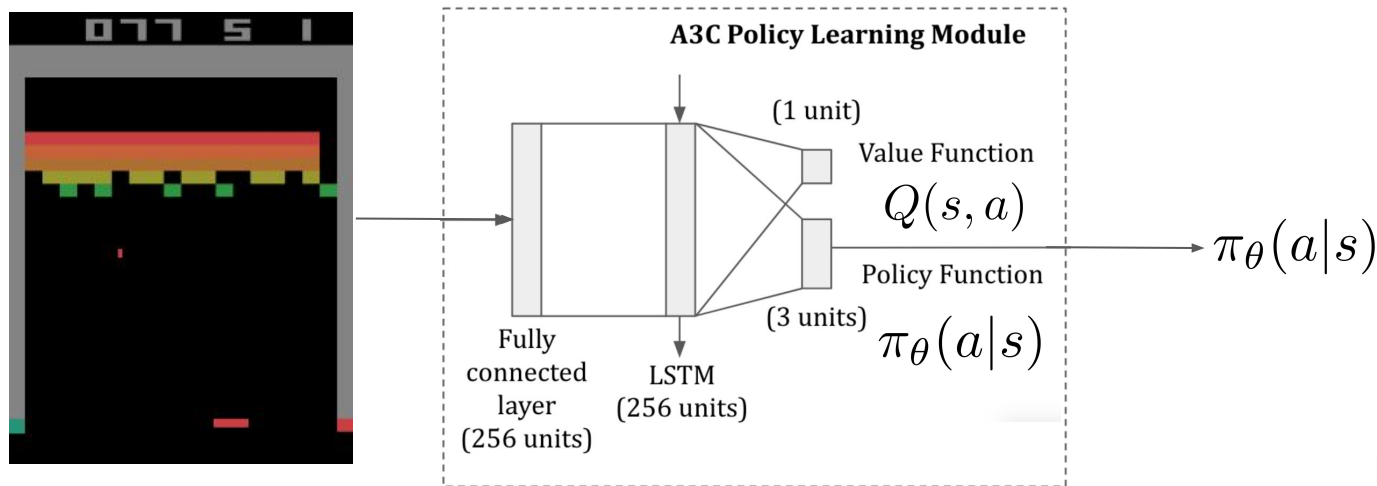
**Problem:** we don't know  $Q$  and  $V$  - can we learn them?

**Yes**, using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the  $Q$  function)

# Actor-critic methods

**Problem:** we don't know  $Q$  and  $V$  - can we learn them?

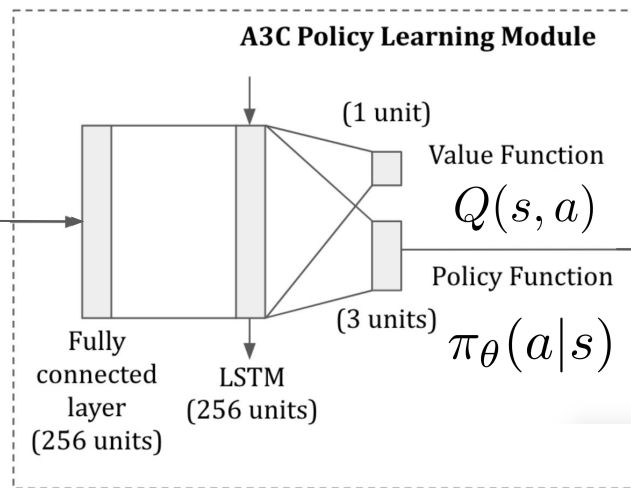
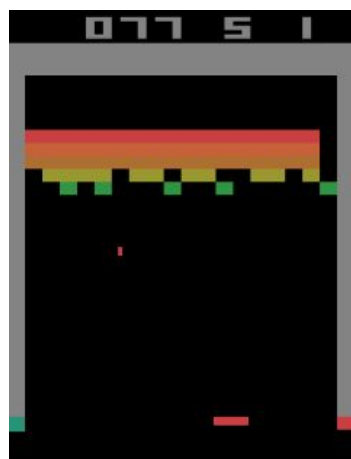
**Yes,** using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q function)



# Actor-critic methods

**Problem:** we don't know  $Q$  and  $V$  - can we learn them?

**Yes,** using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q function)



**Critic: evaluates how good the action is**

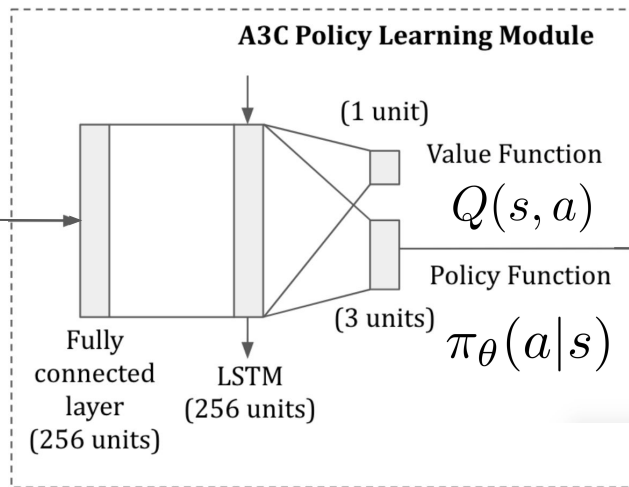
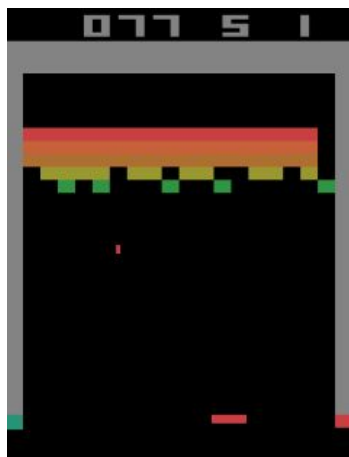
$$\pi_{\theta}(a|s)$$

**Actor: decides what actions to take**

# Actor-critic methods

**Problem:** we don't know Q and V - can we learn them?

**Yes,** using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q function)



**Critic: evaluates how good the action is**

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \underbrace{\left( r + \gamma \max_{a'} Q(s', a'; w_i^-) \right)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right]^2$$

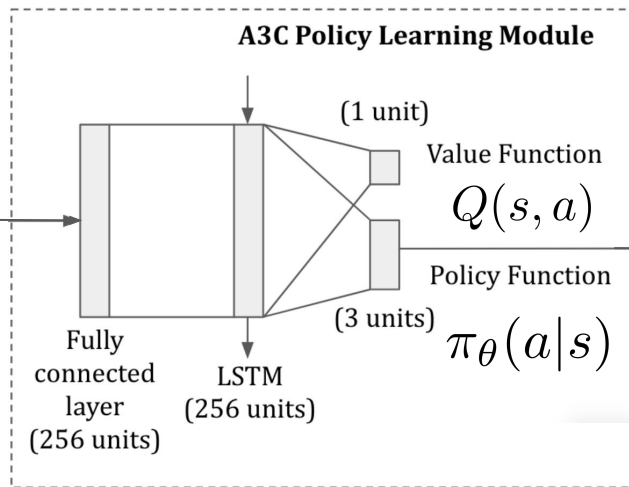
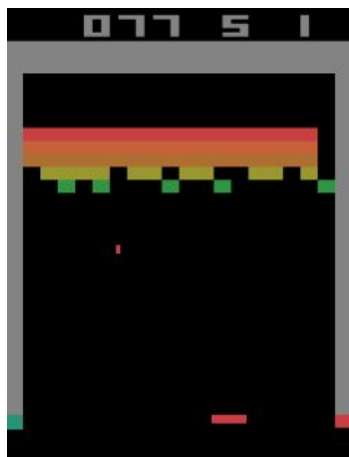
$$\pi_{\theta}(a|s)$$

**Actor: decides what actions to take**

# Actor-critic methods

**Problem:** we don't know Q and V - can we learn them?

**Yes,** using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q function)



**Critic: evaluates how good the action is**

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \underbrace{\left( r + \gamma \max_{a'} Q(s', a'; w_i^-) \right)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right]^2$$

$$\pi_{\theta}(a|s)$$

**Actor: decides what actions to take**

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)) \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$$

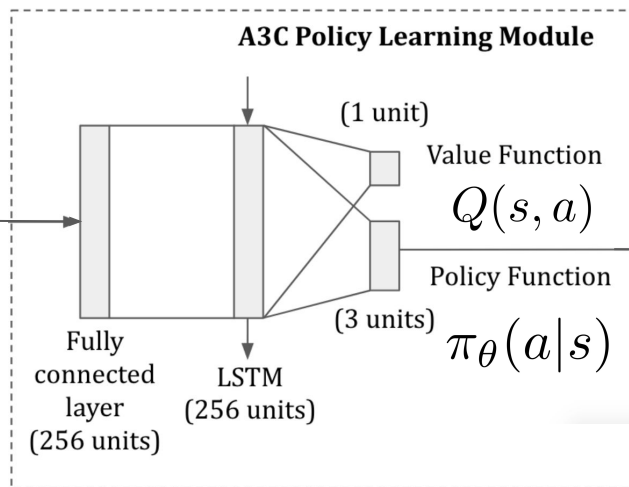
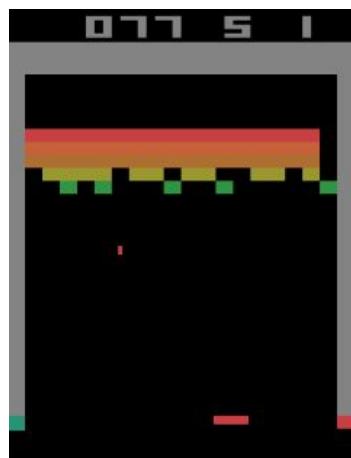
Minh et. al., ICML 2016

# Actor-critic methods

Exploration + experience replay  
Decorrelate samples  
Fixed targets

**Problem:** we don't know Q and V - can we learn them?

**Yes,** using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q function)



**Critic: evaluates how good the action is**

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \underbrace{\left( r + \gamma \max_{a'} Q(s', a'; w_i^-) \right)}_{\text{Q-learning target}} - \underbrace{Q(s, a; w_i)}_{\text{Q-network}} \right]^2$$

$\pi_{\theta}(a|s)$

**Actor: decides what actions to take**

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Minh et. al., ICML 2016

# Summary of RL methods

## Value Based

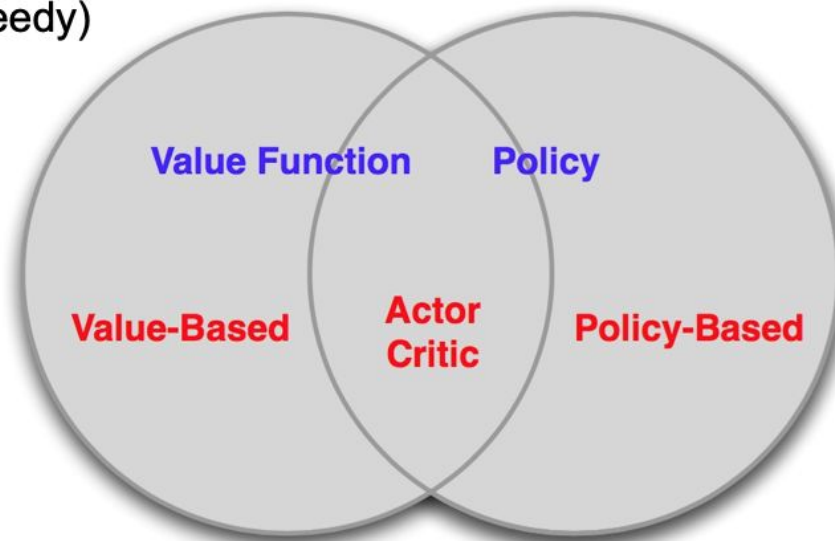
- Value iteration
- Policy iteration
- (Deep) Q-learning
- Learned Value Function
- Implicit policy (e.g.  $\epsilon$ -greedy)

## Policy Based

- Policy gradients
- No Value Function
- Learned Policy

## Actor-Critic

- Actor (policy)
- Critic (Q-values)
- Learned Value Function
- Learned Policy





Applications: Stochastic optimization

# Stochastic Optimization

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

# Stochastic Optimization

$$\max_{\phi} E_{q_{\phi}(z)}[f(z)]$$

## VAEs

$$\max_{\theta, \phi} \mathcal{L}(x; \theta, \phi) \quad \text{Evidence lower bound}$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)] - D_{KL}(q_{\phi}(z|x) || p(z))$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)]$$

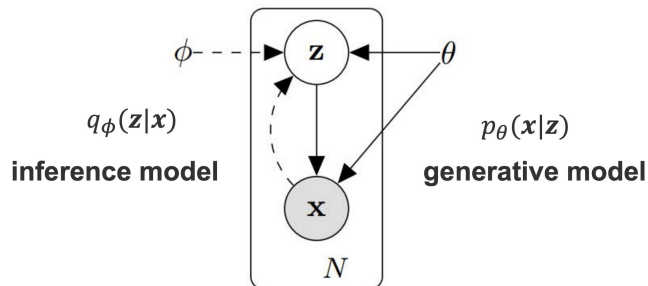


Figure courtesy: Kingma & Welling, 2014

# Stochastic Optimization

$$\max_{\phi} E_{q_{\phi}(z)}[f(z)]$$

## VAEs

$$\max_{\theta, \phi} \mathcal{L}(x; \theta, \phi) \quad \text{Evidence lower bound}$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)] - D_{KL}(q_{\phi}(z|x) || p(z))$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)] \quad \text{Solve by reparameterization!}$$

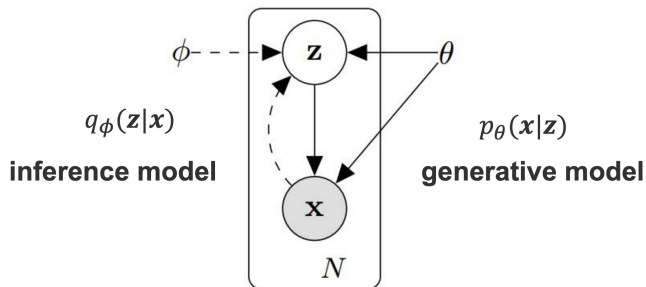
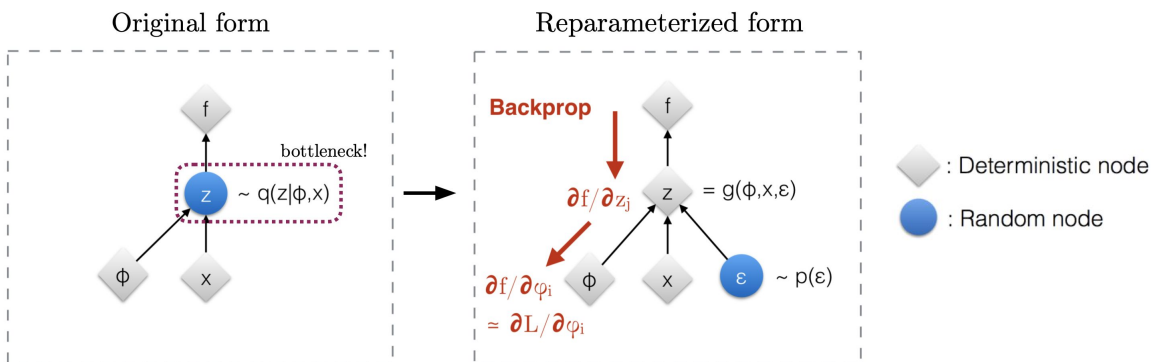


Figure courtesy: Kingma & Welling, 2014



# Stochastic Optimization

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

## VAEs

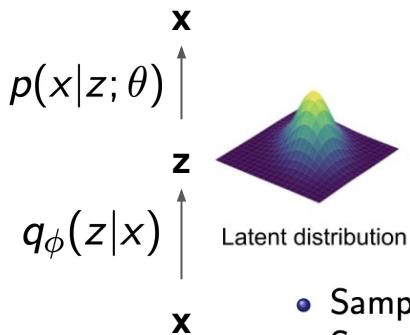
$$\max_{\theta, \phi} \mathcal{L}(x; \theta, \phi) \quad \text{Evidence lower bound}$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)]$$

Solve by reparameterization!

We require that:

- $z$  is continuous
- $q(z)$  is reparameterizable
- $f(z)$  is differentiable wrt  $\phi$



Latent distribution

- Sample  $\mathbf{z} \sim q_{\phi}(\mathbf{z})$
- Sample  $\epsilon \sim \mathcal{N}(0, I)$ ,  $\mathbf{z} = \mu + \sigma\epsilon$

# Stochastic Optimization

## VAEs

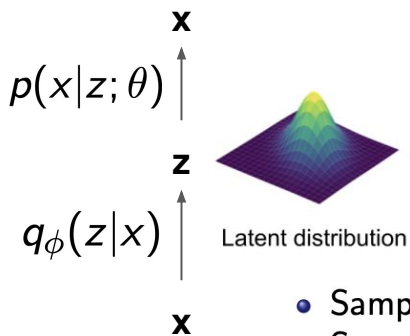
$$\max_{\theta, \phi} \mathcal{L}(x; \theta, \phi) \quad \text{Evidence lower bound}$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)]$$

Solve by reparameterization!

We require that:

- $z$  is continuous
- $q(z)$  is reparameterizable
- $f(z)$  is differentiable wrt  $\phi$



- Sample  $\mathbf{z} \sim q_{\phi}(\mathbf{z})$
- Sample  $\epsilon \sim \mathcal{N}(0, I)$ ,  $\mathbf{z} = \mu + \sigma\epsilon$

$$\max_{\phi} E_{q_{\phi}(z)}[f(z)]$$

## RL

$$\max_{\phi} J(\phi) \quad \text{Reward}$$

$$\max_{\phi} E_{\tau \sim p(\tau; \phi)}[r(\tau)]$$

# Stochastic Optimization

## VAEs

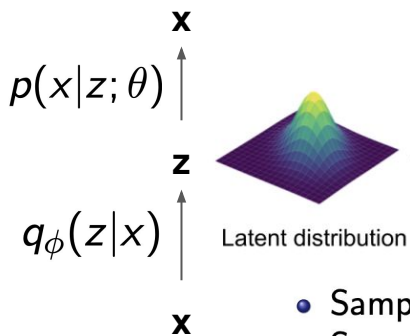
$$\max_{\theta, \phi} \mathcal{L}(x; \theta, \phi) \quad \text{Evidence lower bound}$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)]$$

Solve by reparameterization!

We require that:

- $z$  is continuous
- $q(z)$  is reparameterizable
- $f(z)$  is differentiable wrt  $\phi$



- Sample  $\mathbf{z} \sim q_{\phi}(\mathbf{z})$
- Sample  $\epsilon \sim \mathcal{N}(0, I)$ ,  $\mathbf{z} = \mu + \sigma\epsilon$

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

## RL

$$\max_{\phi} J(\phi) \quad \text{Reward}$$

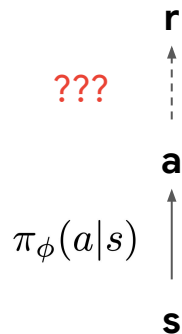
$$\max_{\phi} E_{\tau \sim p(\tau; \phi)}[r(\tau)]$$

Reparameterization???

In RL (at least for discrete actions):

- $T$  is a sequence of discrete actions
- $p(T; \phi)$  is not reparameterizable
- $r(T)$  is a black box function

i.e. the environment



# Stochastic Optimization

## VAEs

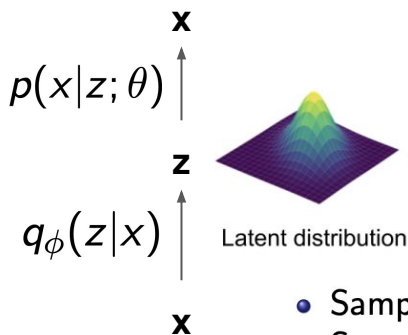
$$\max_{\theta, \phi} \mathcal{L}(x; \theta, \phi) \quad \text{Evidence lower bound}$$

$$\max_{\theta, \phi} E_{q_{\phi}(z|x)}[\log p(x|z; \theta)]$$

Solve by reparameterization!

We require that:

- $z$  is continuous
- $q(z)$  is reparameterizable
- $f(z)$  is differentiable wrt  $\phi$



- Sample  $z \sim q_{\phi}(z)$
- Sample  $\epsilon \sim \mathcal{N}(0, I)$ ,  $z = \mu + \sigma\epsilon$

$$\max_{\phi} E_{q_{\phi}(z)}[f(z)]$$

## RL

$$\max_{\phi} J(\phi) \quad \text{Reward}$$

$$\max_{\phi} E_{\tau \sim p(\tau; \phi)}[r(\tau)]$$

Reparameterization???

In RL (at least for discrete actions):

- $T$  is a sequence of discrete actions
- $p(T; \phi)$  is not reparameterizable
- $r(T)$  is a black box function

i.e. the environment

???

$\pi_{\phi}(a|s)$

$r$   
↑  
 $a$   
↑  
 $s$

**REINFORCE is a general-purpose solution!**



# Revisiting REINFORCE

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

We want to take gradients wrt  $\phi$  of the term:

$$\nabla_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})] = E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})\nabla_{\phi} \log q_{\phi}(\mathbf{z})]$$

# Revisiting REINFORCE

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

We want to take gradients wrt  $\phi$  of the term:

$$\nabla_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})] = E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z})]$$

We can now compute a Monte Carlo estimate:

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_{\phi}(\mathbf{z})$  and estimate

$$\nabla_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_{\phi} \log q_{\phi}(\mathbf{z}^k)$$

# Revisiting REINFORCE

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

We want to take gradients wrt  $\phi$  of the term:

$$\nabla_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})] = E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})\nabla_{\phi} \log q_{\phi}(\mathbf{z})]$$

We can now compute a Monte Carlo estimate:

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_{\phi}(\mathbf{z})$  and estimate

$$\nabla_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_{\phi} \log q_{\phi}(\mathbf{z}^k)$$

what we derived. sample trajectories and compute.

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

# Revisiting REINFORCE

$$\max_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})]$$

We can now compute a Monte Carlo estimate:

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_{\phi}(\mathbf{z})$  and estimate

$$\nabla_{\phi} E_{q_{\phi}(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_{\phi} \log q_{\phi}(\mathbf{z}^k)$$

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

We just need the distribution  $q()$  to allow for easy sampling

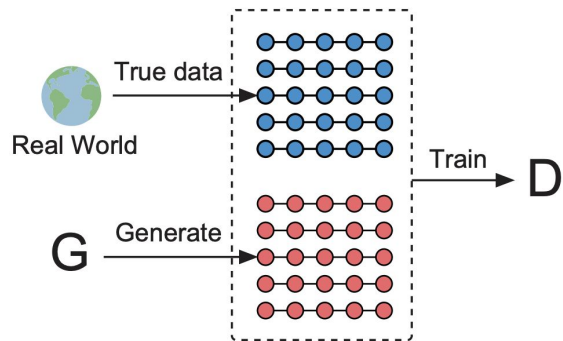
- $\mathbf{z}$  can be discrete or continuous!
- $q(\mathbf{z})$  can be a discrete and continuous distribution! (but must be differentiable wrt  $\phi$ )
- $f(\mathbf{z})$  can be a black box!

# Applications: Text generation

GANs for text generation

1. Text data is discrete

- Discriminator gradient does not exist for samples from categorical distribution
- Gradient sparse due to large dictionary size



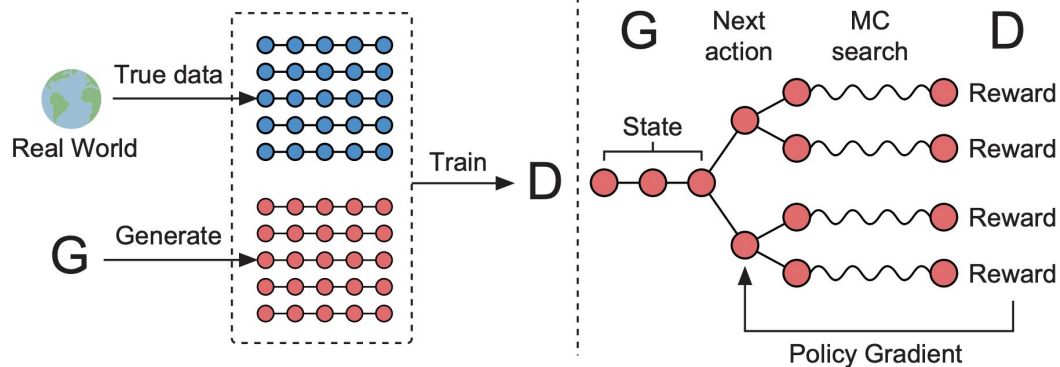
# Applications: Text generation

GANs for text generation

1. Text data is discrete

- Discriminator gradient does not exist for samples from categorical distribution
- Gradient sparse due to large dictionary size

More efficient search strategy for most likely sentence



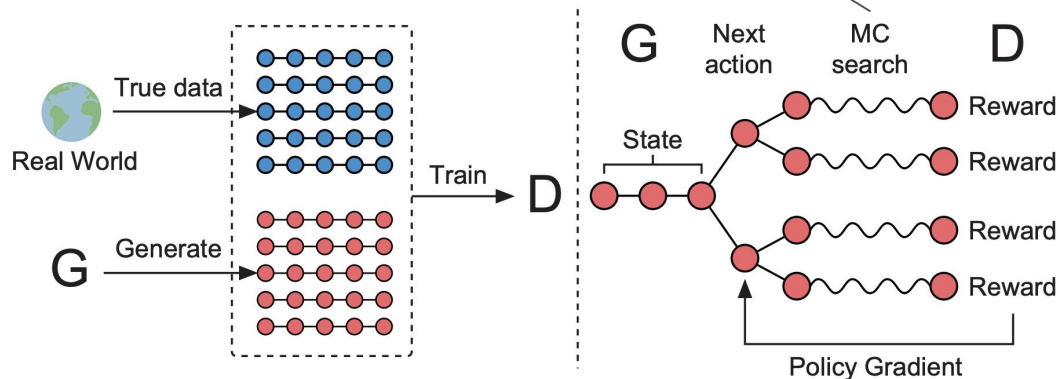
# Applications: Text generation

GANs for text generation

1. Text data is discrete

- Discriminator gradient does not exist for samples from categorical distribution
- Gradient sparse due to large dictionary size

More efficient search strategy for most likely sentence



**Training generator:**

After sampling all words using Monte Carlo search, compute reward for generator based on discriminator feedback

- if similar to real text, high reward
- if different from real text, low reward

complete sentences LM

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_\phi(\mathbf{z})$  and estimate

$$\nabla_\phi E_{q_\phi(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_\phi \log q_\phi(\mathbf{z}^k)$$

disc reward

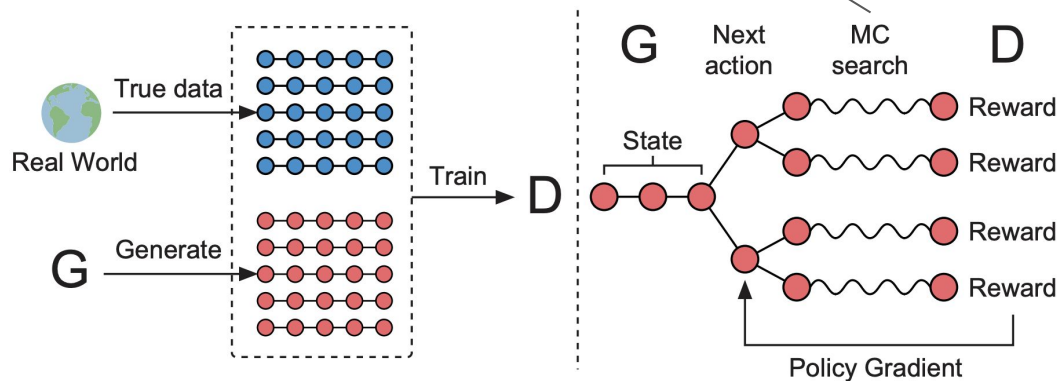
# Applications: Text generation

GANs for text generation

1. Text data is discrete

- Discriminator gradient does not exist for samples from categorical distribution
- Gradient sparse due to large dictionary size

More efficient search strategy for most likely sentence



**Training generator:**

After sampling all words using Monte Carlo search, compute reward for generator based on discriminator feedback

- if similar to real text, high reward
- if different from real text, low reward

complete sentences LM

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_\phi(\mathbf{z})$  and estimate

$$\nabla_\phi E_{q_\phi(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_\phi \log q_\phi(\mathbf{z}^k)$$

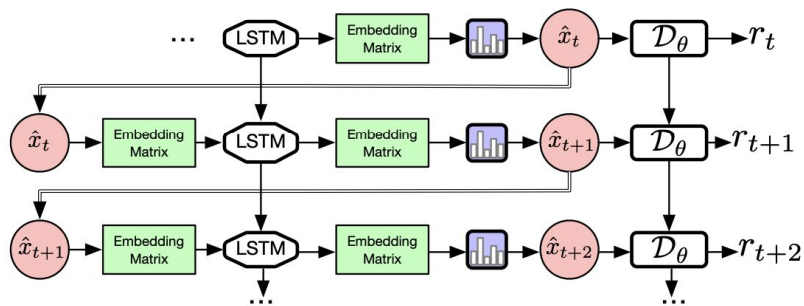
disc reward



# Applications: Text generation

## GANs for text generation

1. Text is sensitive to noise (small disturbances easily alters the meaning of text)
2. Sparse discriminator feedback (feedback only makes sense on full sentences)



+ Instead of discriminator reward after whole sentence, use a recurrent discriminator which compares generated vs real prefixes to give dense rewards at all time steps  
+ large variance from REINFORCE: use large batch sizes and subtract baseline (moving average of rewards)  
+ other tricks, see paper

prefixes                      LM

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_\phi(\mathbf{z})$  and estimate

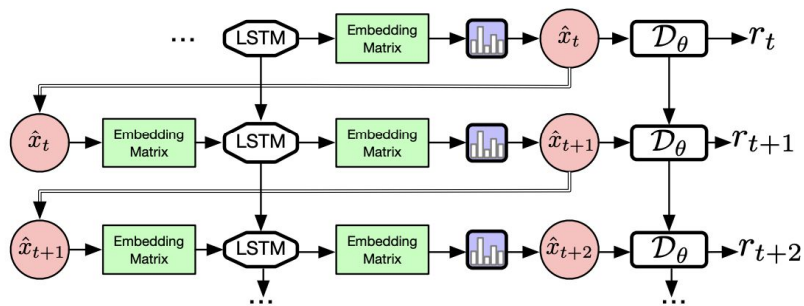
$$\nabla_\phi E_{q_\phi(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_\phi \log q_\phi(\mathbf{z}^k)$$

disc reward

# Applications: Text generation

## GANs for text generation

1. Text is sensitive to noise (small disturbances easily alters the meaning of text)
2. Sparse discriminator feedback (feedback only makes sense on full sentences)



Other approaches as well without using policy gradients

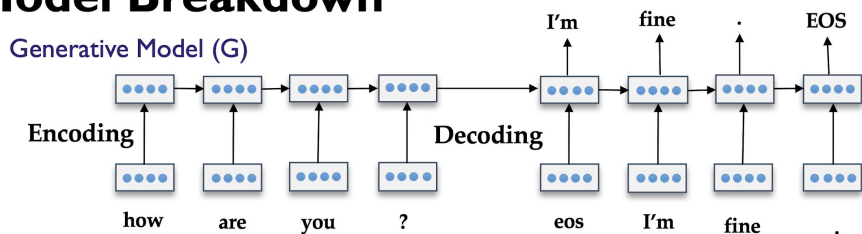
- discriminator directly comparing in logit space
- use Gumbel softmax (Jang et al. 2016)
- see more in

<https://www.cl.uni-heidelberg.de/statnlpgroup/blog/rl4nmt/>  
and  
[https://deepgenerativemodels.github.io/assets/slides/cs236\\_lecture15.pdf](https://deepgenerativemodels.github.io/assets/slides/cs236_lecture15.pdf)

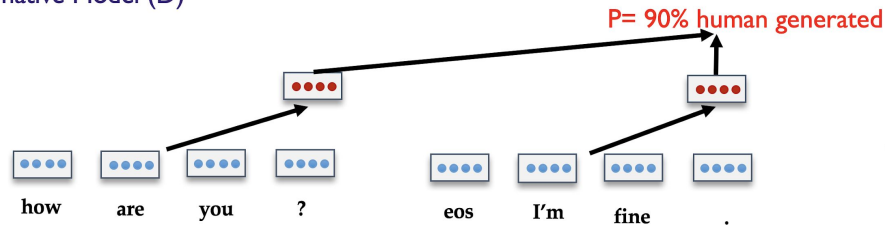
# Applications: Dialog generation

GANs for dialog generation and evaluation

## Model Breakdown



**Discriminative Model (D)**



**Sample:**  
Input message  
Response 1  
Response 2  
...  
Response K

**Define rewards:**  
1. Ease of answering  
2. Information flow  
3. Meaningfulness  
4. Discriminator wrt human dialog

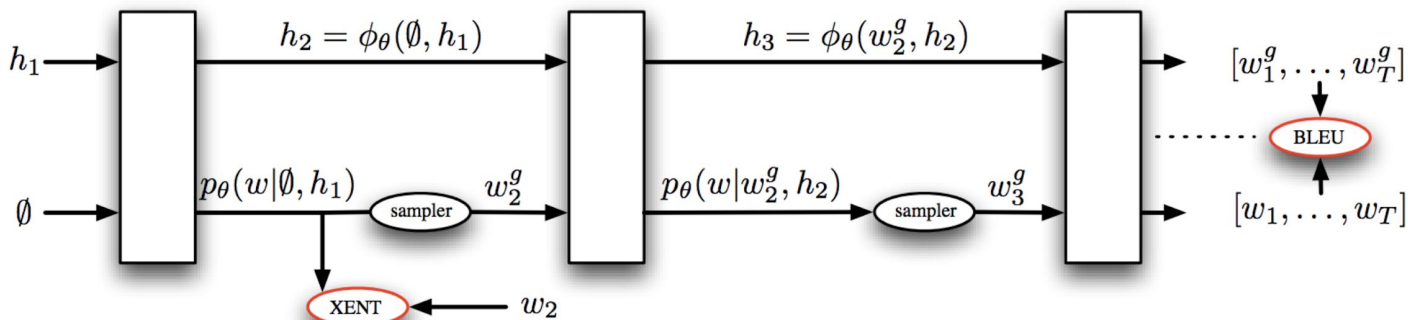
Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_\phi(\mathbf{z})$  and estimate

$$\nabla_\phi E_{q_\phi(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_\phi \log q_\phi(\mathbf{z}^k)$$

# Applications: Optimizing general rewards

Instead of optimizing for cross-entropy (not final evaluation metric), optimize directly for the evaluation metric e.g. BLEU score

- BLEU score only defined on raw text after sampling from softmax
- Not differentiable through standard gradient methods.



Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_\phi(\mathbf{z})$  and estimate

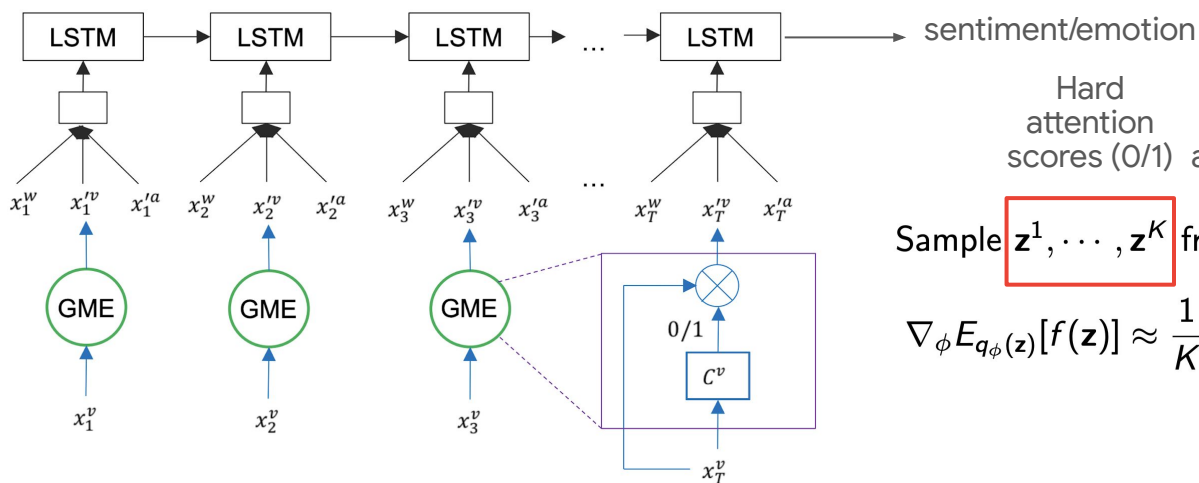
$$\nabla_\phi E_{q_\phi(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_\phi \log q_\phi(\mathbf{z}^k)$$

BLEU score

# Applications: Hard attention

Hard attention 'gates' (0/1) rather than soft attention (softmax between 0-1)

- Can be seen as discrete layers in between differentiable neural net layers



Hard attention scores (0/1)      Controller: input to hard attention scores

Sample  $\mathbf{z}^1, \dots, \mathbf{z}^K$  from  $q_\phi(\mathbf{z})$  and estimate

$$\nabla_\phi E_{q_\phi(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_k f(\mathbf{z}^k) \nabla_\phi \log q_\phi(\mathbf{z}^k)$$

classification accuracy



[Xu et. al., ICML 2015]

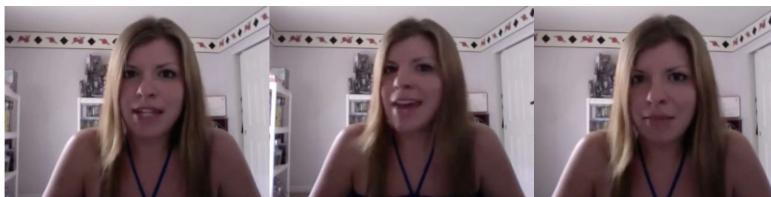
[Chen et al., ICMI 2017]

# Applications: Hard attention

Hard attention ‘gates’ (0/1) rather than soft attention (softmax between 0-1)

- Can be seen as discrete layers in between differentiable neural net layers

Sentiment analysis,  
emotion recognition



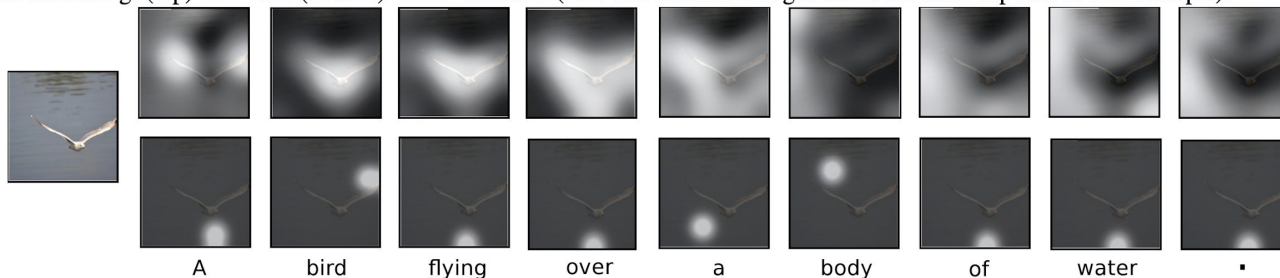
Reject

Pass

Reject

Figure 3. Visualization of the attention for each generated word. The rough visualizations obtained by upsampling the attention weights and smoothing. (top) “soft” and (bottom) “hard” attention (note that both models generated the same captions in this example).

Image captioning



[Xu et. al., ICML 2015]

[Chen et al., ICMI 2017]

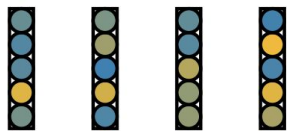
# Applications: RL and Language

# RL and Language

## Task-independent

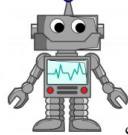
[...] having the correct **key** can open the lock [...]  
[...] known lock and **key** device was discovered [...]  
[...] unless the correct **key** is inserted [...]

Pre-training



$V_{\text{key}}$   $V_{\text{skull}}$   $V_{\text{ladder}}$   $V_{\text{rope}}$

Pre-trained



Agent

Action

State, Reward



Environment

## Task-dependent

### **Language-assisted**

**Key** Opens a door of the same color as the key.

**Skull** They come in two varieties, rolling skulls and bouncing skulls ... you must jump over rolling skulls and walk under bouncing skulls.

### **Language-conditional**

Go down the ladder and walk right immediately to avoid falling off the conveyor belt, jump to the yellow rope and again to the platform on the right.

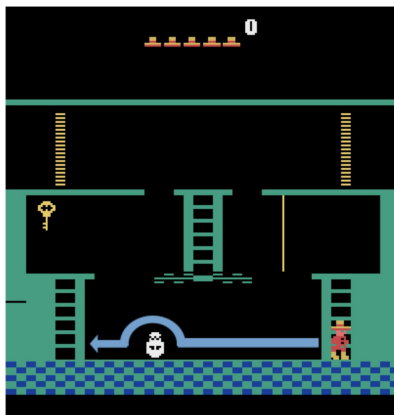


# Language-conditional RL

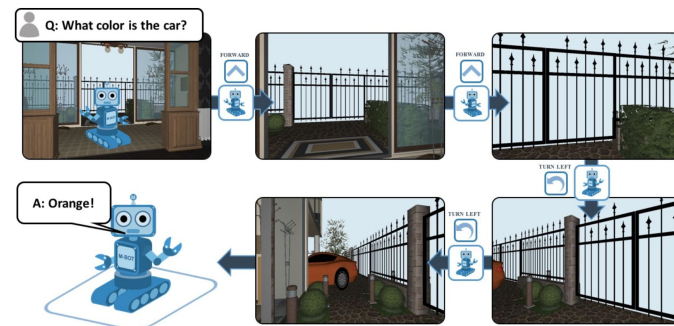
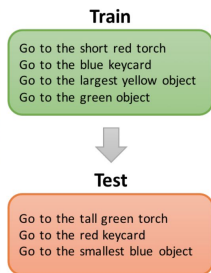
Instruction following



Rewards from instructions



Language in S and A



# Language-assisted RL

- Language for communicating domain knowledge
- Language for structuring policies

# Language-assisted RL: Domain knowledge

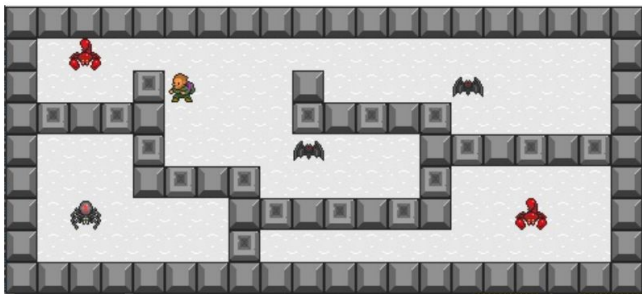
- Properties of entities in the environment are annotated by language



is an enemy who chases you



is a stationary collectible



is a randomly moving enemy

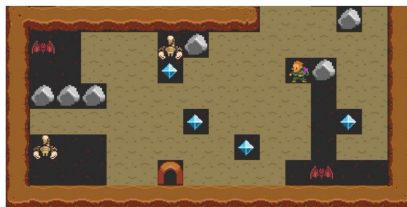




is a stationary immovable wall

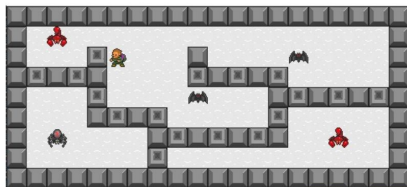
from Amazon Mturk :-  
(asked annotators to play  
the game and describe  
entities)



# Language-assisted RL: Domain knowledge

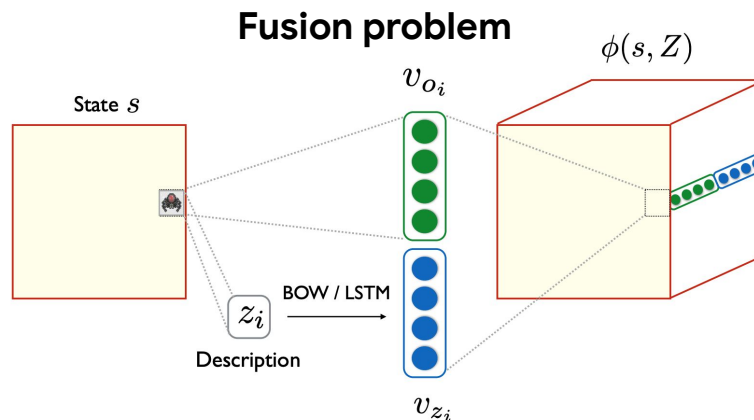
- Properties of entities in the environment are annotated by language



-  is an enemy who chases you
-  is a stationary collectible





-  is a randomly moving enemy
-  is a stationary immovable wall

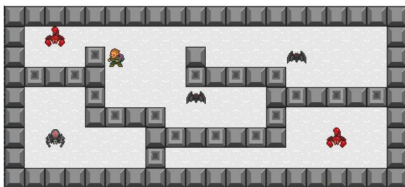




# Language-assisted RL: Domain knowledge

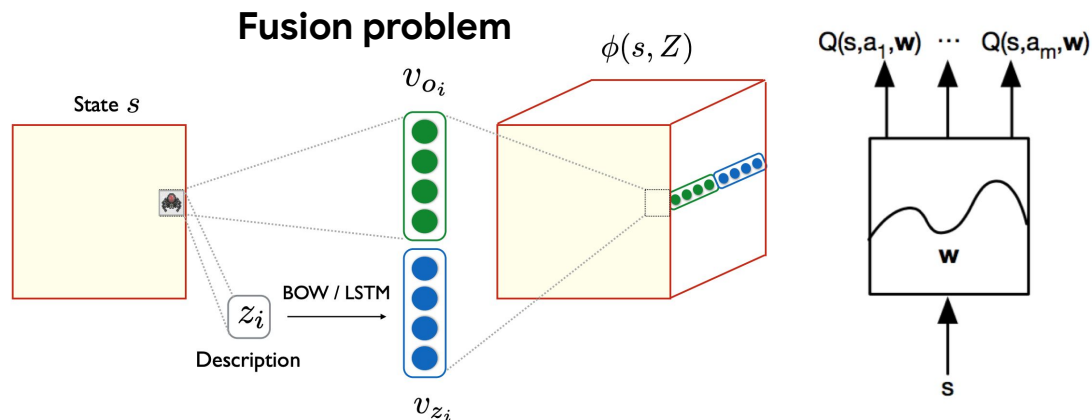
- Properties of entities in the environment are annotated by language



-  is an enemy who chases you
-  is a stationary collectible

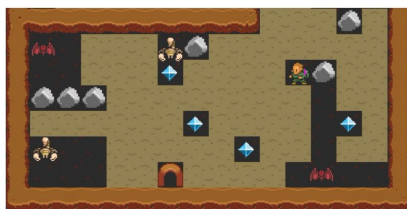




-  is a randomly moving enemy
-  is a stationary immovable wall

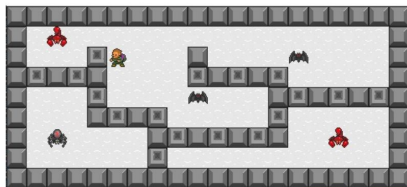




# Language-assisted RL: Domain knowledge

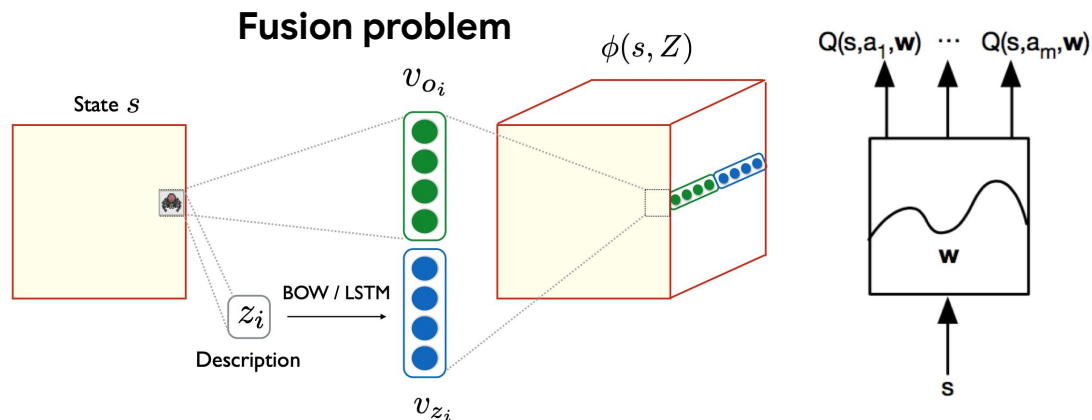
- Properties of entities in the environment are annotated by language



-  is an enemy who chases you
-  is a stationary collectible



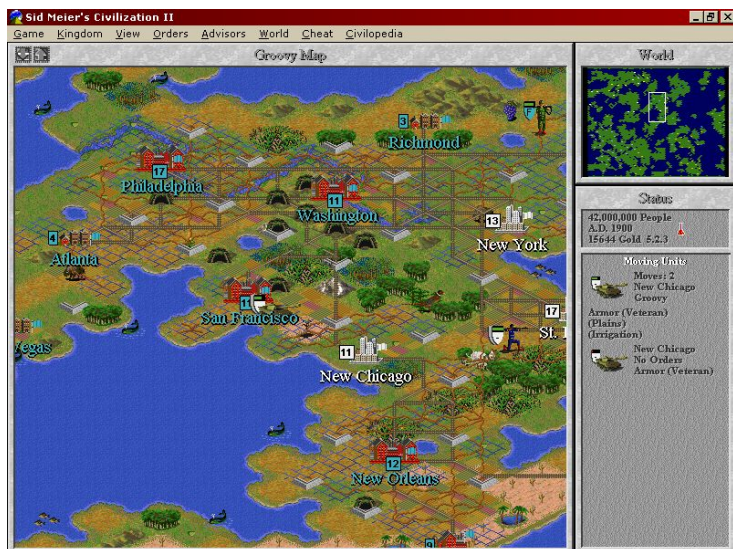
-  is a randomly moving enemy
-  is a stationary immovable wall



**Grounded language learning**  
Helps to ground the meaning of text to the dynamics, transitions, and rewards  
Language helps in multi-task learning and transfer learning

# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals

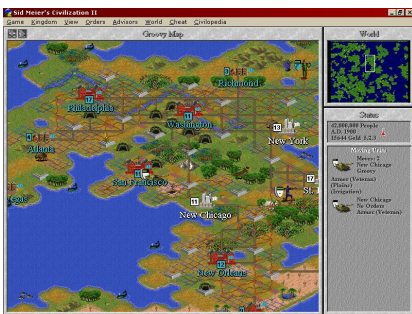


*The natural resources available where a population settles affects its ability to produce food and goods. Build your city on a plains or grassland square with a river running through it if possible.*

Figure 1: An excerpt from the user manual of the game Civilization II.

# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals



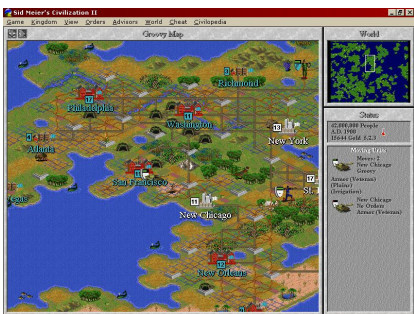
*The natural resources available where a population settles affects its ability to produce food and goods. Build your city on a plains or grassland square with a river running through it if possible.*

1. Choose **relevant** sentences
2. Label words into **action-description**, **state-description**, or **background**



# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals



*The natural resources available where a population settles affects its ability to produce food and goods. Build your city on a plains or grassland square with a river running through it if possible.*

1. Choose **relevant** sentences
2. Label words into **action-description**, **state-description**, or **background**

## Map tile attributes:

- Terrain type (e.g. grassland, mountain, etc)
- Tile resources (e.g. wheat, coal, wildlife, etc)

## City attributes:

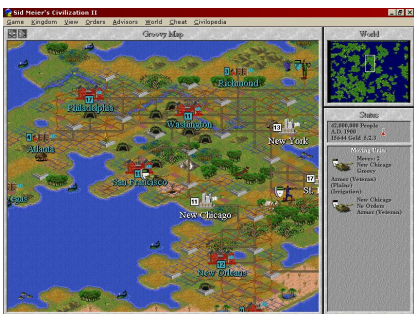
- City population
- Amount of food produced

## Unit attributes:

- Unit type (e.g., worker, explorer, archer, etc)
- Is unit in a city ?

# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals



*The natural resources available where a population settles affects its ability to produce food and goods. Build your city on a plains or **grassland** square with a river running through it if possible.*

1. Choose **relevant** sentences
2. Label words into **action-description, state-description, or background**

## Map tile attributes:

- Terrain type (e.g. **grassland**, mountain, etc)
- Tile resources (e.g. wheat, coal, wildlife, etc)

## City attributes:

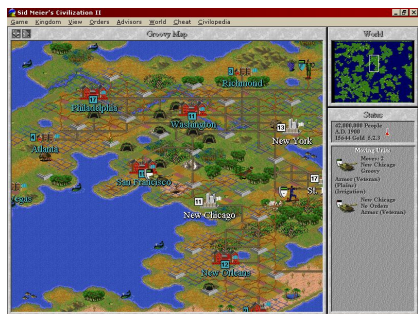
- City population
- Amount of food produced

## Unit attributes:

- Unit type (e.g., worker, explorer, archer, etc)
- Is unit in a city ?

# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals



*The natural resources available where a population settles affects its ability to produce food and goods. Build your city on a plains or grassland square with a river running through it if possible.*

## Map tile attributes:

- Terrain type (e.g. grassland, mountain, etc)
- Tile resources (e.g. wheat, coal, wildlife, etc)

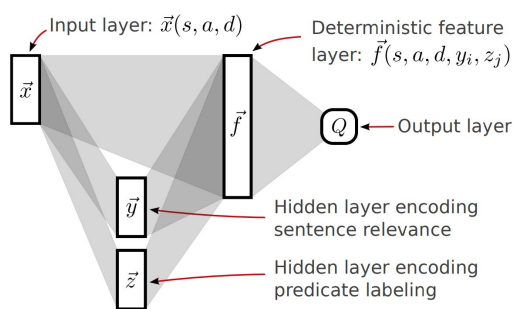
## City attributes:

- City population
- Amount of food produced

## Unit attributes:

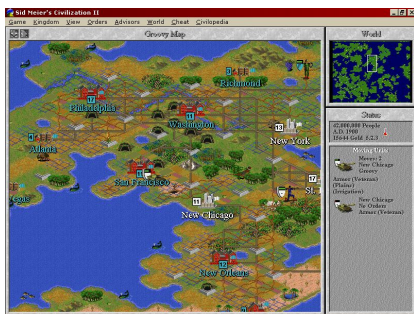
- Unit type (e.g., worker, explorer, archer, etc)
- Is unit in a city ?

1. Choose **relevant** sentences
2. Label words into **action-description**, **state-description**, or **background**



# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals



- Phalanxes are twice as effective at defending cities as warriors. ✓
- Build the city on plains or grassland with a river running through it. ✓
- You can rename the city if you like, but we'll refer to it as Washington.
- There are many different strategies dictating the order in which advances are researched

Relevant sentences

- After the road is built, use the settlers to start improving the terrain.  
S S S A A A A A
- When the settlers becomes active, chose build road.  
S S S A A A
- Use settlers or engineers to improve a terrain square within the city radius.  
A S X A A S A X S S S S

A: action-description  
S: state-description

# Language-assisted RL: Domain knowledge

- Learning to read instruction manuals

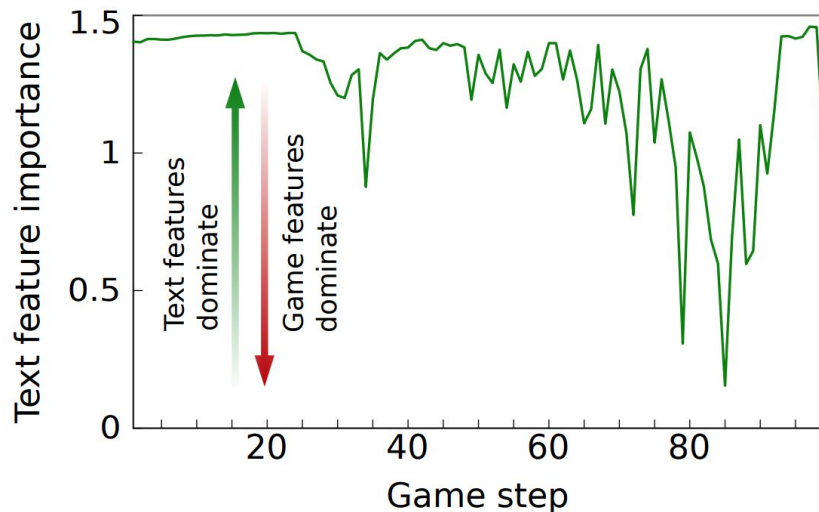
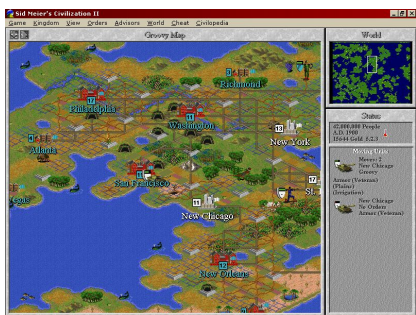


Method	% Win	% Loss	Std. Err.
Random	0	100	—
Built-in AI	0	0	—
Game only	17.3	5.3	$\pm 2.7$
Sentence relevance	46.7	2.8	$\pm 3.5$
<b>Full model</b>	<b>53.7</b>	5.9	$\pm 3.5$
Random text	40.3	4.3	$\pm 3.4$
Latent variable	26.1	3.7	$\pm 3.1$

Grounded language learning  
Ground the meaning of text to the dynamics, transitions, and rewards  
Language helps in learning

# Language-assisted RL: Domain knowledge

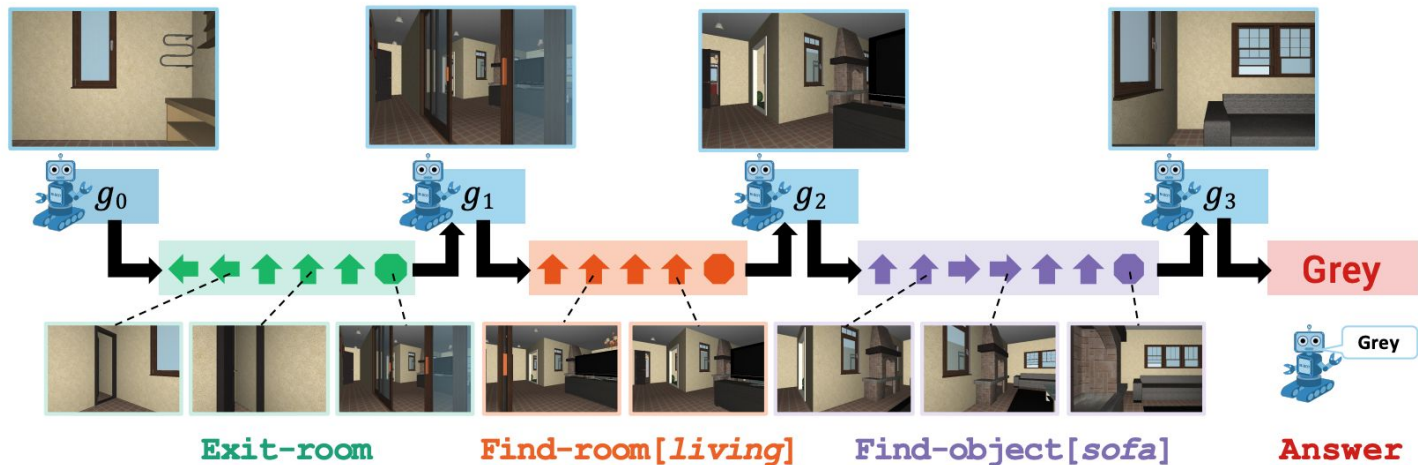
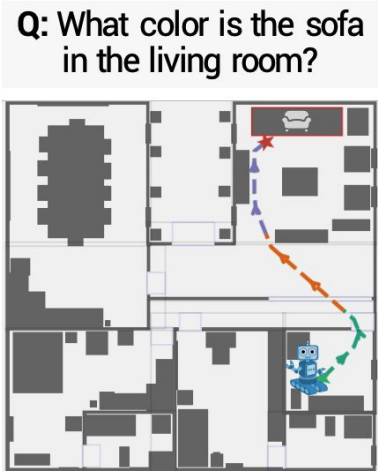
- Learning to read instruction manuals



Language is most important at the start when you don't have a good policy  
Afterwards, the model relies on game features

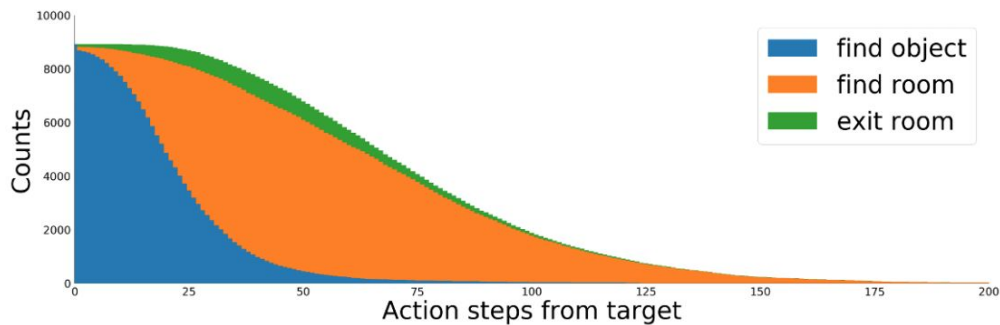
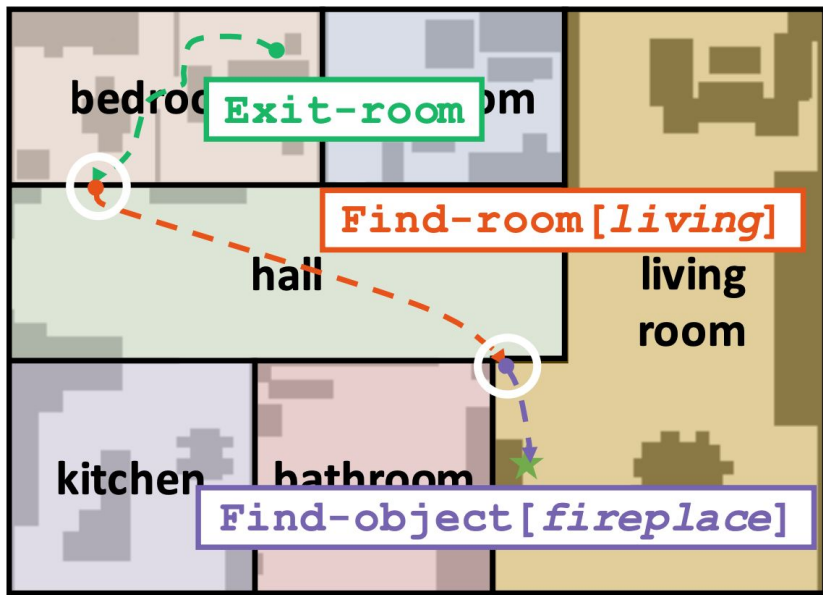
# Language for structuring policies

- Composing modules for Embodied QA



# Language for structuring policies

- Composing modules for Embodied QA





# Summary of RL methods

## Value Based

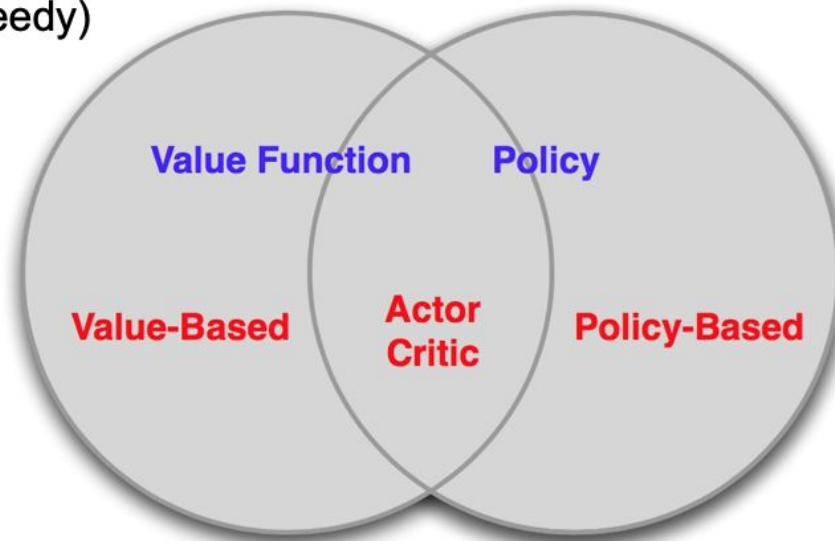
- Value iteration
- Policy iteration
- (Deep) Q-learning
- Learned Value Function
- Implicit policy (e.g.  $\epsilon$ -greedy)

## Policy Based

- Policy gradients
- No Value Function
- Learned Policy

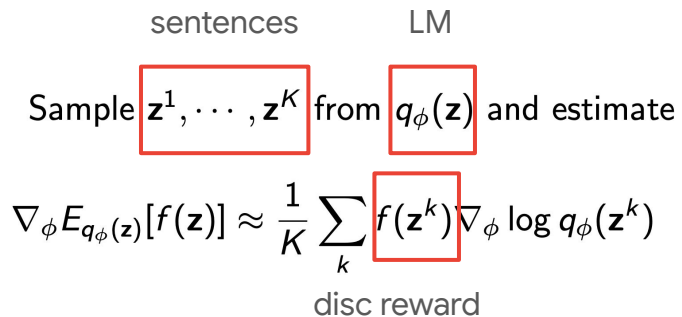
## Actor-Critic

- Actor (policy)
- Critic (Q-values)
- Learned Value Function
- Learned Policy

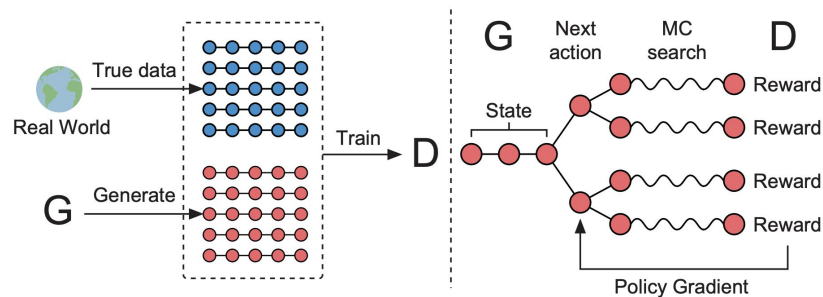


# Summary of applications

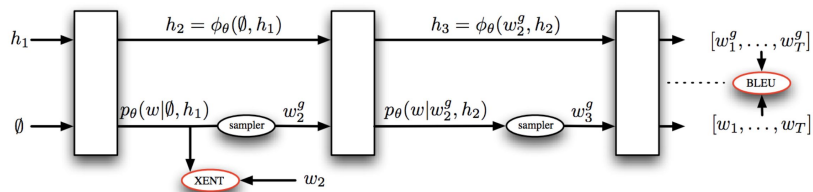
## Stochastic optimization



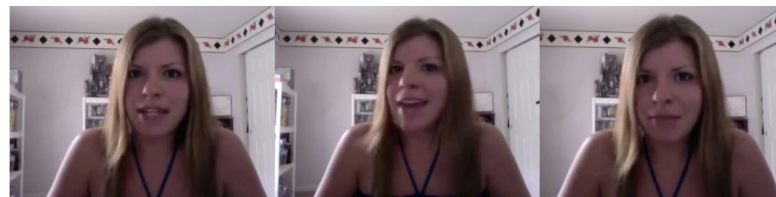
## Text generation



## General reward functions



## Discrete layers



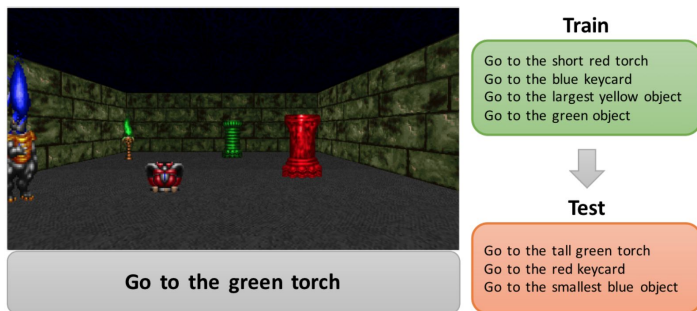
Reject

Pass

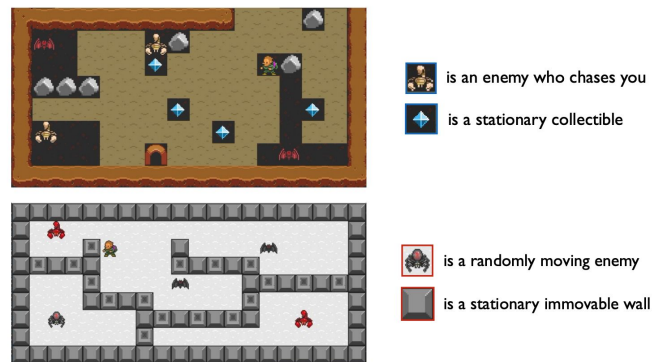
Reject

# Summary of applications

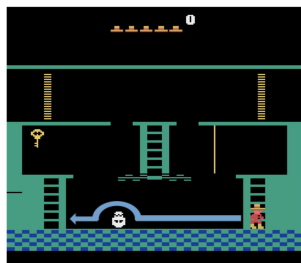
## Instruction following



## Language as domain knowledge



## Language for rewards



"Jump over the skull while going to the left"

## Language to structure policies



Q: What color is the sofa in the living room?

