

# A Tutorial on RL

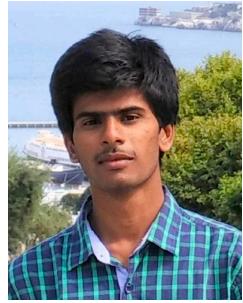
Benjamin Van Roy  
Stanford University and DeepMind



Maria Dimakopolou  
Netflix



Shi Dong  
Stanford



Vikranth Dwaracherla  
Stanford



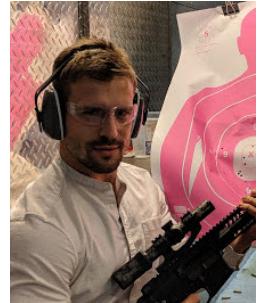
Morteza Ibrahimi  
DeepMind



Abbas Kazerouni  
Stanford



Xiuyuan Lu  
Stanford



Ian Osband  
DeepMind



Dan Russo  
Columbia



Zheng Wen  
DeepMind



Bill Zhu  
Facebook

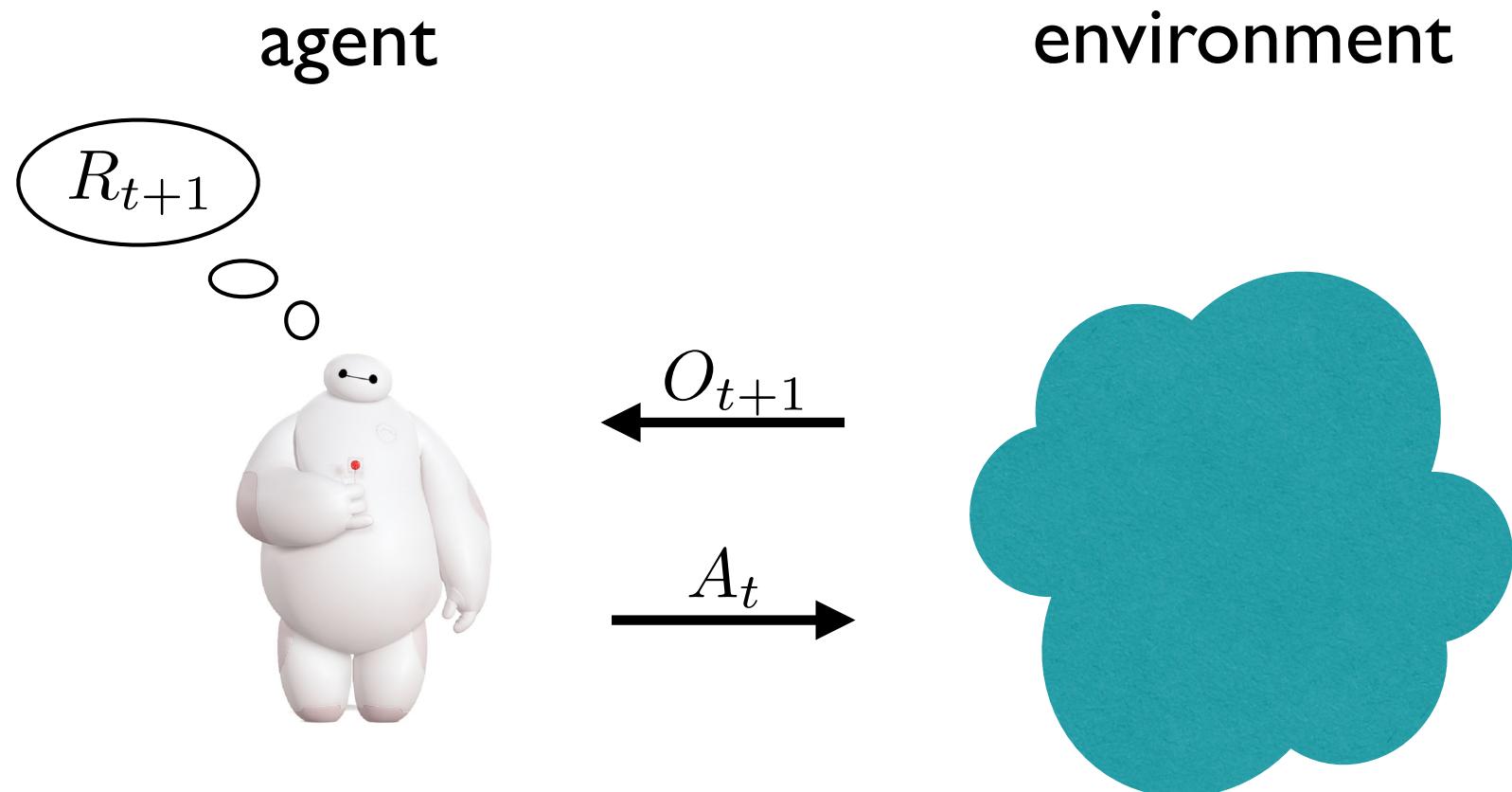
# Objectives

- Introduce you to reinforcement learning
  - Problem, formalism, and examples
  - Algorithms
  - Some current research directions
- Prerequisites
  - Probability / stochastic processes
  - Supervised learning
    - e.g., neural networks, SGD

# What is reinforcement learning?

- A problem
- A community who work on the problem
- Methods produced by that community
  - Some of which are also useful for other problems

# The Reinforcement Learning Problem



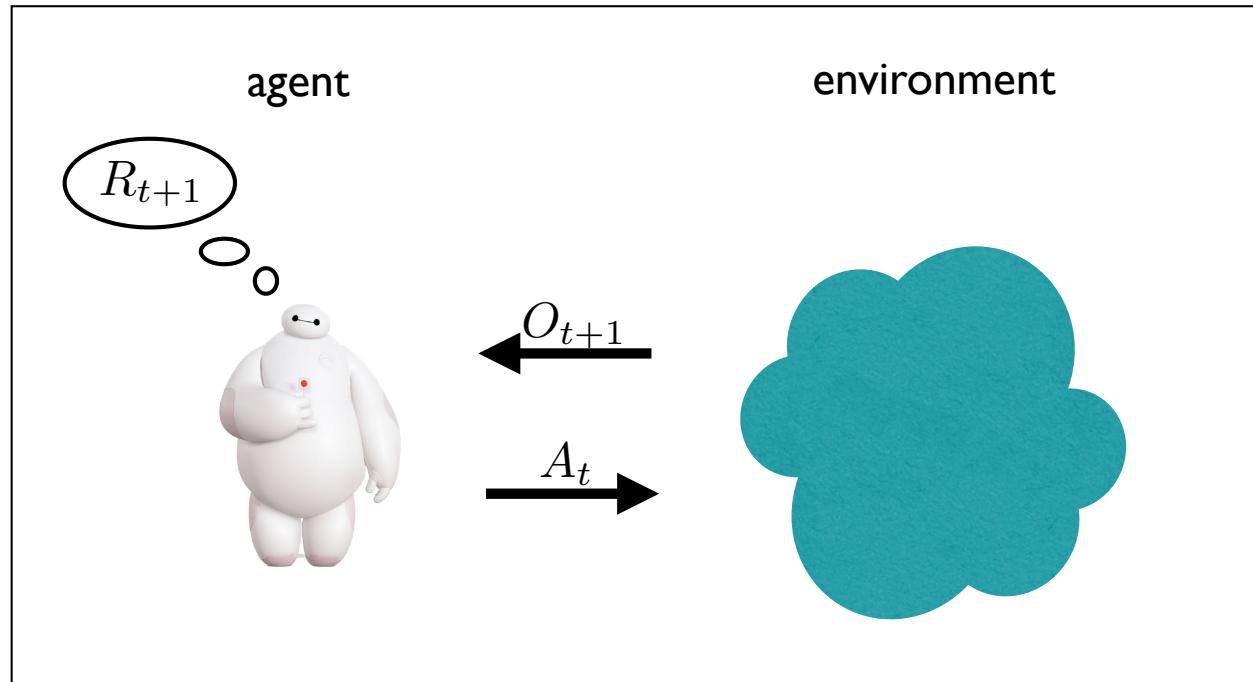
“other learning paradigms are about **minimization**,  
reinforcement learning is about **maximization**”

- Klopf

# Artificial General Intelligence

- Develop a system that can achieve human-level competence across *all* intellectual tasks
- Can RL achieve AGI?

# Episodic Version of the RL Problem



data from episodes 1 through  $\ell$

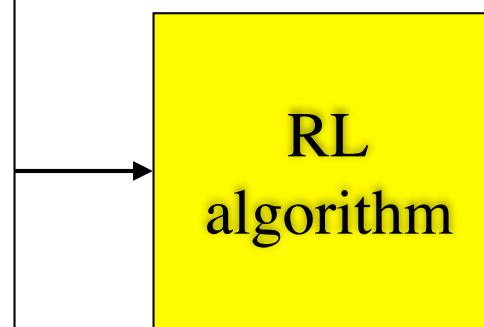
$$O_0^1, A_0^1, O_1^1, \dots, A_{\tau_1-1}^1, O_{\tau_1}^1$$

•

•

•

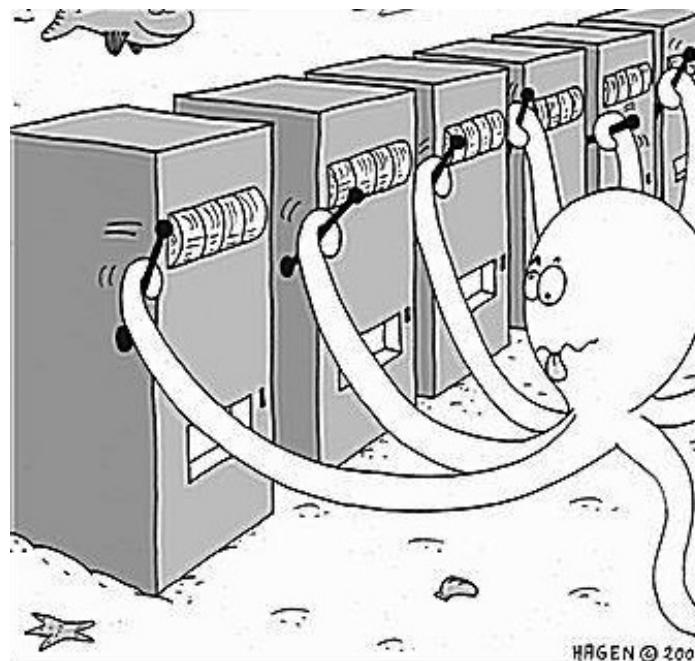
$$O_0^\ell, A_0^\ell, O_1^\ell, \dots, A_{\tau_\ell-1}^\ell, O_{\tau_\ell}^\ell$$



policy for  
next episode

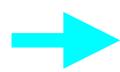
$$\pi_\ell$$

# RL versus Bandits

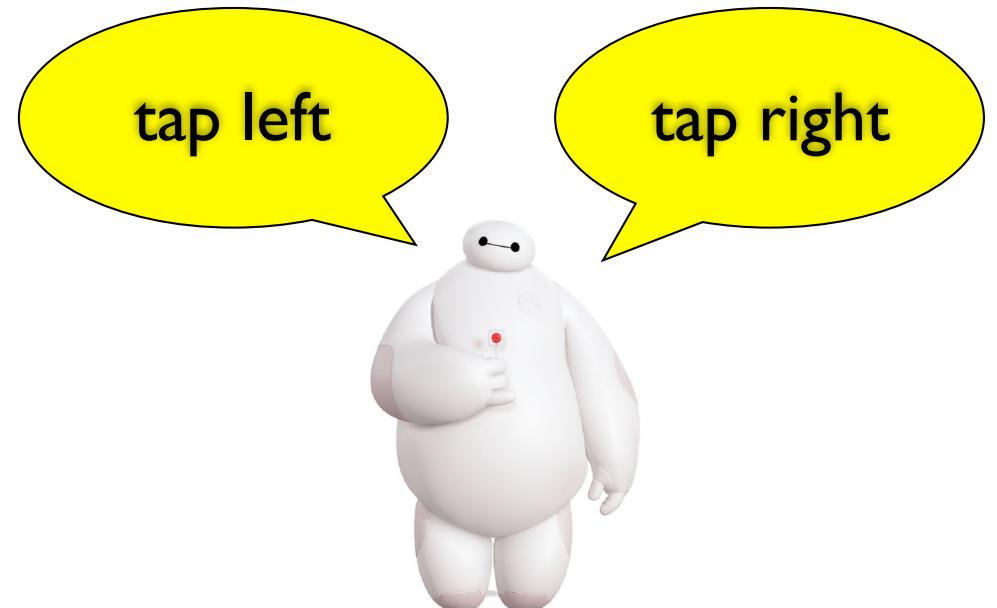


- Bandit problem = RL with episode of duration 1
- Bandit algorithms are RL algorithms applied to bandits
- Bandits much better understood
- RL formulation can be overkill for specific applications
- RL can be critical for delayed consequences
- RL is critical for AGI

# Outline

- 
- Examples
  - Mathematical formalism
  - Common algorithms
  - Challenges of data efficiency

# Example: Cart-Pole Swingup



20 seconds per episode

# Example: Robotics

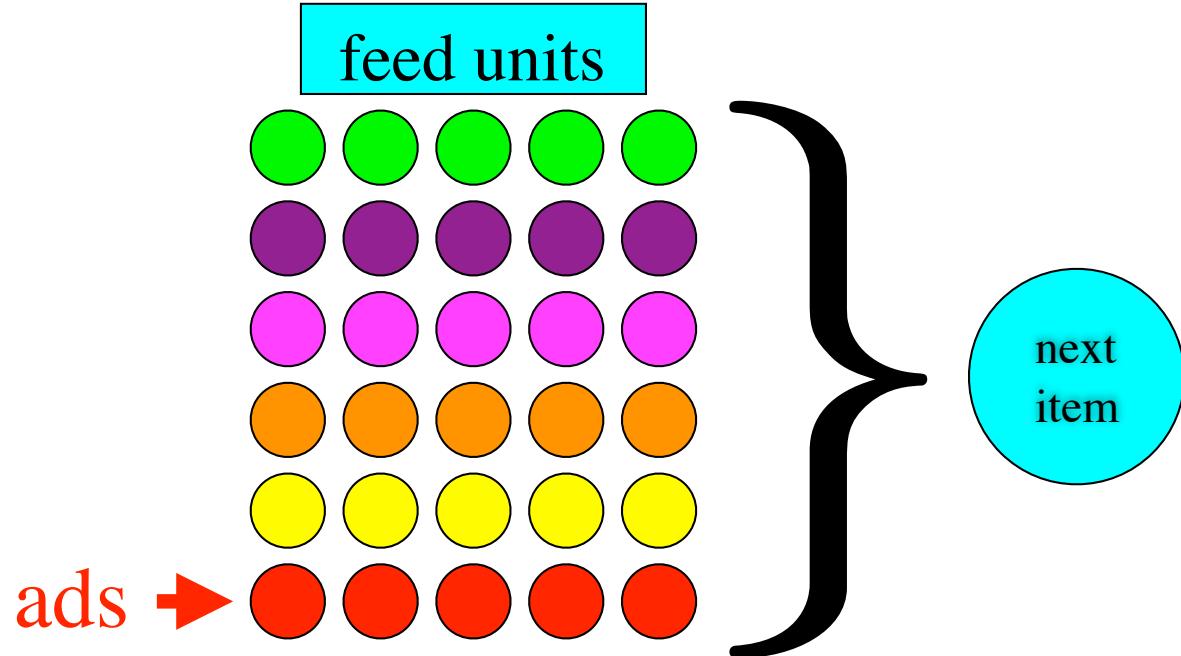


# Example: Arcade Games



action: joystick actuation  
observation: screen shot

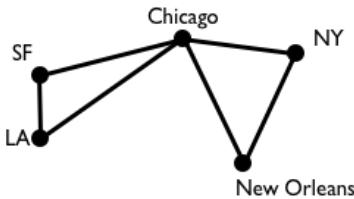
# Example: News Feed Engagement



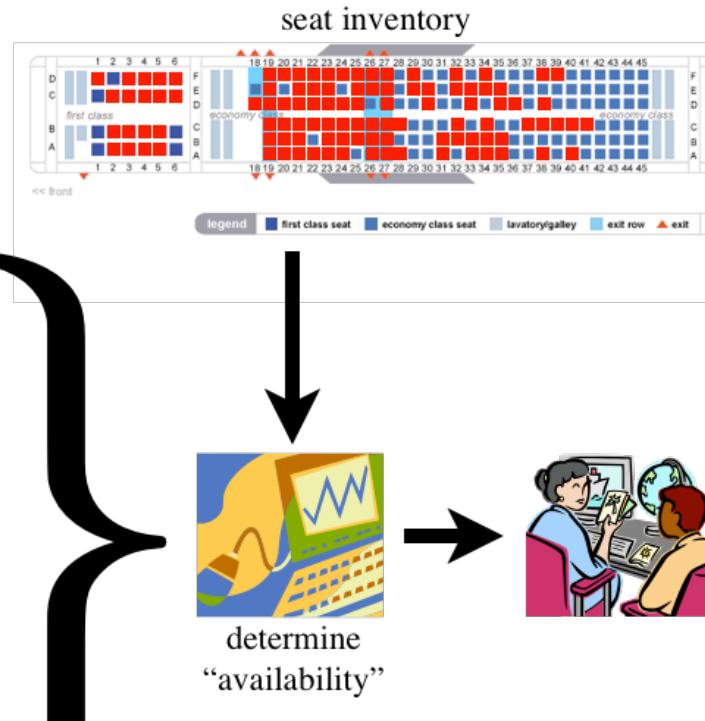
action: next item

observation: click, like, share, comment, scroll, terminate

# Example: Revenue Management



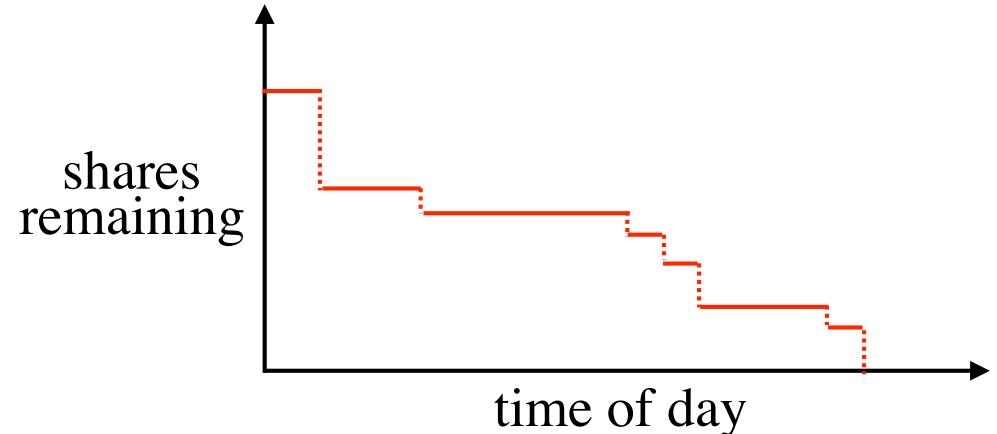
Fare Product	Specification	Price
1	Economy SF-LA 8AM 12/11/10 LA-SF 1PM 15/11/10 Nonrefundable ...	\$120
2	First SF-Ch 8AM 12/24/10 Ch-NY 4PM 12/24/10 Refundable ...	\$2500
...	...	...



action: accept/reject  
observation: requests, cancellations, forecasts, etc.

# Example: Trade Execution

sell 1000 shares by end of day



action: orders to exchange  
observation: order book, other market data

# Outline

- Examples
- Mathematical formalism
- Common algorithms
- Challenges of data efficiency



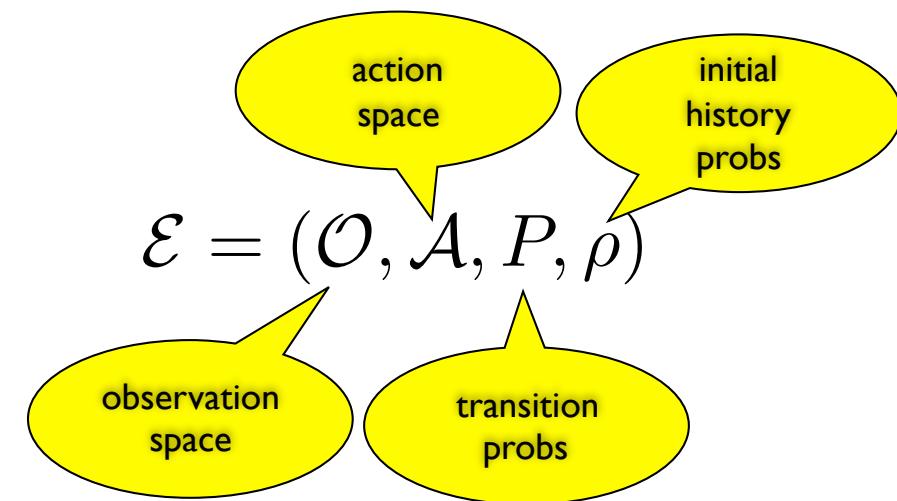
# The Role of Formalism

- Clarify the class of problems under consideration
- Provide notation for making rigorous statements
- Establish properties that apply across the class
- Design algorithms that apply across the class
- Guide modeling
  - Epistemic versus aleatoric uncertainty

# Formalism: Environment

- Environment

$$\mathcal{E} = (\mathcal{O}, \mathcal{A}, P, \rho)$$



- Assume finite
  - observation space
  - action space

- News feed model

$$\mathcal{E} = (\mathcal{O}, \mathcal{A}, P, \rho)$$

$$\mathcal{A} \in \{\text{feed unit}\}$$

$$\mathcal{O} \in \left\{ \begin{array}{l} (\text{inventory, user profile}) \\ (\text{like/swipe/quit}) \end{array} \right\}$$

$$\rho \rightarrow (\text{inventory, user profile})$$

$$P \rightarrow (\text{like/swipe/quit})$$



# Formalism: Dynamics

- History

$$H_0 = (O_o)$$

$$H_t = (O_0, A_0, \dots, O_t)$$

$$H_t \in \mathcal{H} = \mathcal{H}^+ \cup \mathcal{H}^-$$



- Dynamics

$$\rho_h = \mathbb{P}(H_0 = h | \mathcal{E})$$

$$H_0 = ((\text{inventory}, \text{user profile}))$$

$$H_t = (\dots, \text{feed unit } \#22, (\text{like}, \text{swipe}))$$

$$\mathcal{H}^- = \{(\dots, (\dots, \text{quit}))\}$$



$$P_{ahh'} = \mathbb{P}(H_{t+1} = h' | \mathcal{E}, H_t = h, A_t = a)$$

# Formalism: Policies

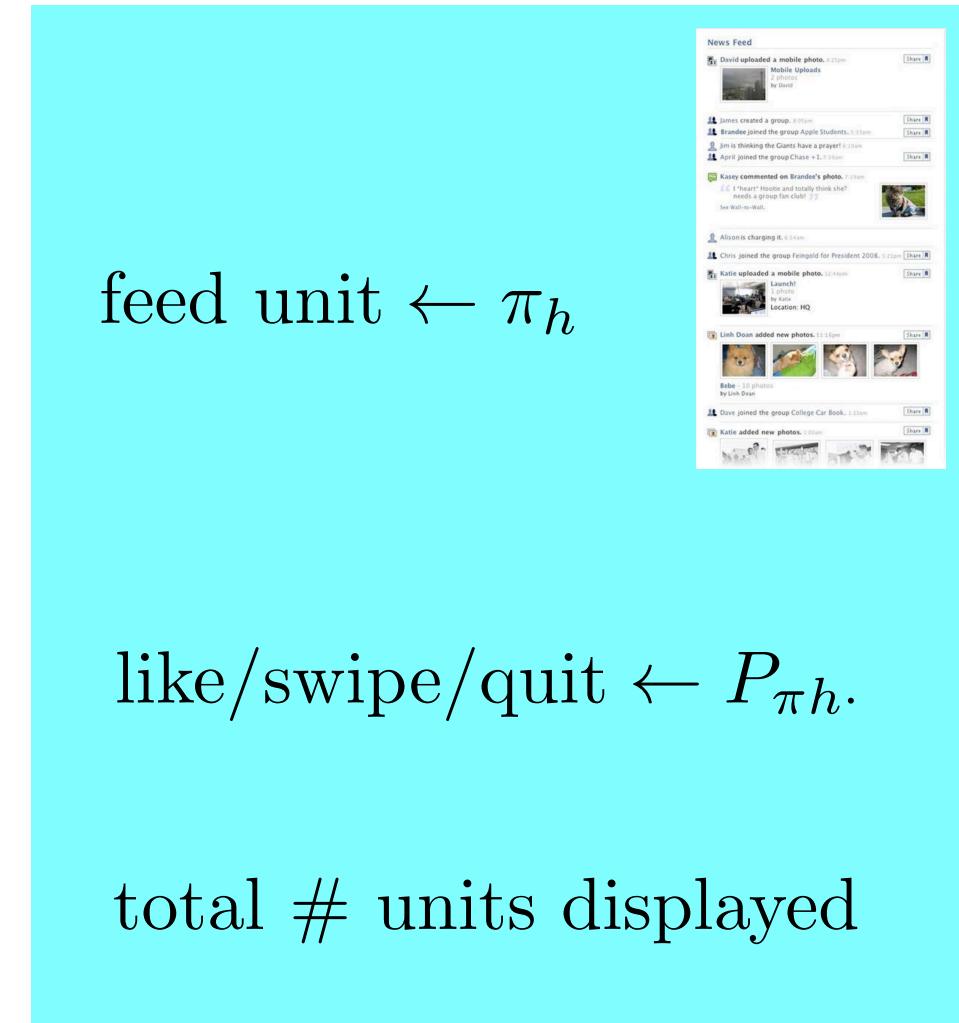
- Policy

$$\pi_h(a) = \mathbb{P}(A_t = a | \mathcal{E}, H_t = h)$$

- Transition probs

$$P_{\pi h h'} = \sum_{a \in \mathcal{A}} P_{a h h'}$$

- Finite episode duration  $\tau$



feed unit  $\leftarrow \pi_h$

like/swipe/quit  $\leftarrow P_{\pi h}.$

total # units displayed

# Formalism: Rewards

- Period reward

$$R_{t+1} = r(H_{t+1})$$

$$r(H_{t+1}) = \begin{cases} 1 & \text{if } A_t \in \text{ads} \\ 0 & \text{otherwise} \end{cases}$$

- Episode reward

$$\sum_{t=0}^{\tau-1} R_{t+1}$$

total # ads displayed

How will an optimal policy behave?

# Value Functions

- Value functions  $V^\pi : \mathcal{H}^+ \rightarrow \mathbb{R}$

$$V^\pi(h) = \mathbb{E} \left[ \sum_{k=t}^{\tau-1} R_{k+1} \middle| \mathcal{E}, H_t = h, \pi \right]$$

$$V^*(h) = \sup_{\pi} V^\pi(h)$$

remaining # ads under  $\pi$

remaining # ads under  $\pi^*$

- Value per episode  $\bar{V}^\pi \in \mathbb{R}$

$$\bar{V}^\pi = \sum_{h \in \mathcal{H}^+} \rho_h V^\pi(h)$$

total # ads under  $\pi$

$$\bar{V}^* = \sum_{h \in \mathcal{H}^+} \rho_h V^*(h)$$

total # ads under  $\pi^*$

# Action Value Functions

- Value conditioned on immediate action  $Q^\pi : \mathcal{H}^+ \times \mathcal{A} \rightarrow \mathbb{R}$

$$Q^\pi(h, a) = \mathbb{E} \left[ \sum_{k=t}^{\tau-1} R_{k+1} \middle| \mathcal{E}, H_t = h, A_t = a, \pi \right]$$

$$Q^\pi(h, a) = \sum_{h' \in \mathcal{H}} P_{ahh'} r(h') + \sum_{h' \in \mathcal{H}^+} P_{ahh'} V^\pi(h')$$

- Optimal action value function

$$Q^*(h, a) = \sup_{\pi} Q^\pi(h, a)$$

# From Value to Policy

- Greedy policy with respect to  $Q$

$$\pi_h \sim \text{unif} \left( \arg \max_{a \in \mathcal{A}} Q(h, a) \right)$$

- Optimal policy is greedy with respect to  $Q^*$

$$\pi_h^* \sim \text{unif} \left( \arg \max_{a \in \mathcal{A}} Q^*(h, a) \right)$$

- Actually, any distribution over the argmax will work...

# Bellman's Equation

- Dynamic programming operators

$$(TV)(h) = \max_{a \in \mathcal{A}} \left( \sum_{h' \in \mathcal{H}} P_{ahh'} r(h') + \sum_{h' \in \mathcal{H}^+} P_{ahh'} V(h') \right)$$

$$(FQ)(h, a) = \sum_{h' \in \mathcal{H}} P_{ahh'} r(h') + \sum_{h' \in \mathcal{H}^+} P_{ahh'} \max_{a' \in \mathcal{A}} Q(h', a')$$

- Bellman's equation

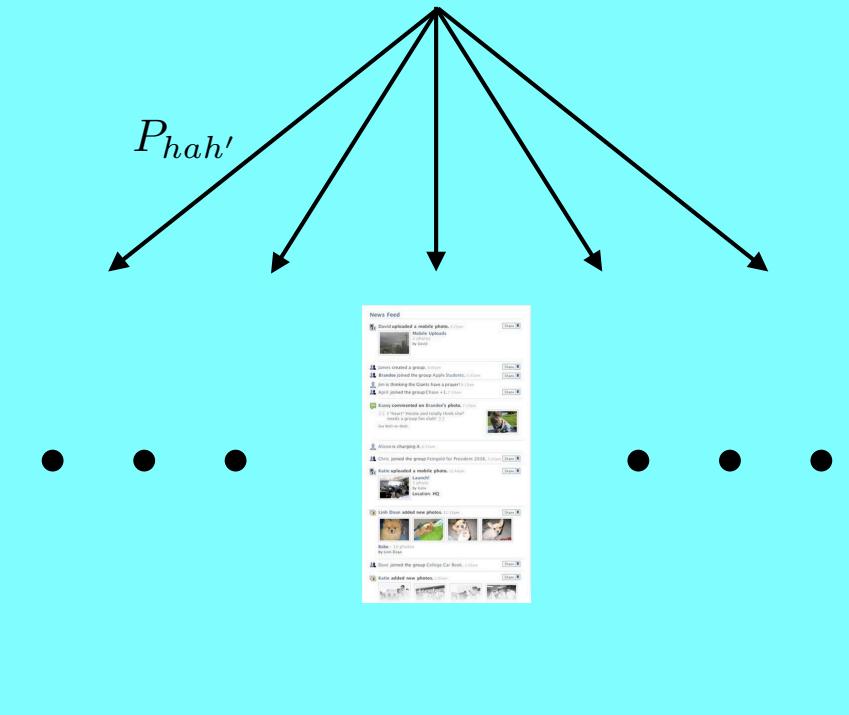
$$TV^* = V^*$$

$$FQ^* = Q^*$$

# Optimization via Value Iteration

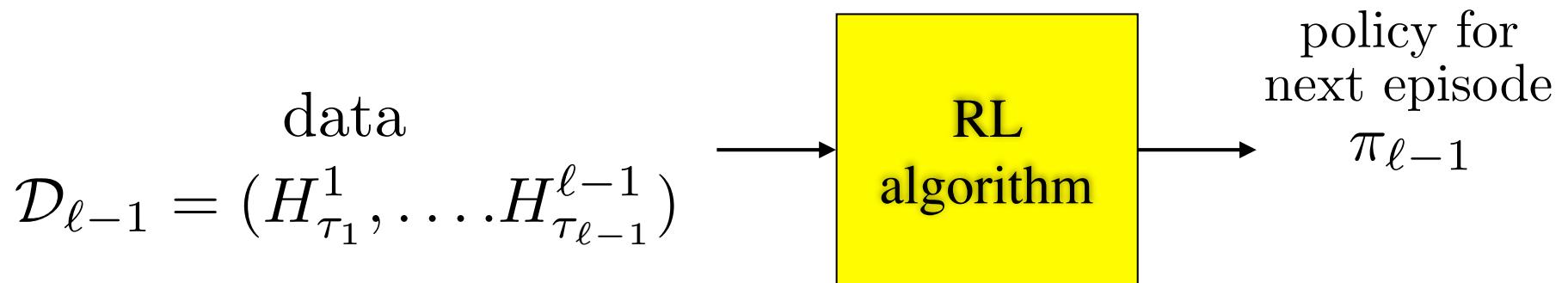
- Computing  $V^*$ 
  - Repeat  $V \leftarrow TV$  forever
- Computing  $Q^*$ 
  - Repeat  $Q \leftarrow FQ$  forever
- VI converges, but
  - Computationally intractable
  - Requires knowledge of  $P$
- RL algorithms address these issues

value  $Q(h, a)$   
history  $h$   
action  $a$



# Formalism: Episodic RL

- Environment  $\mathcal{E}$  is a random variable
  - In particular,  $\rho$  and  $P$  are random variables
  - *Epistemic* versus *aleatoric* uncertainty
- Episodes  $\ell = 1, 2, \dots, L$



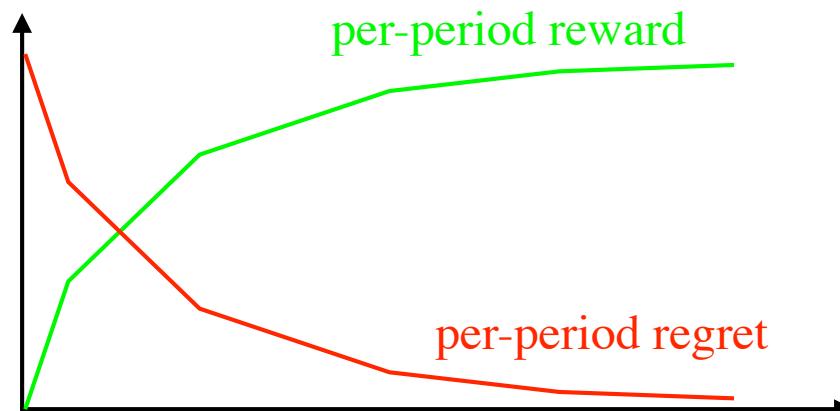
- RL algorithms address
  - Computational issues
  - Statistical issues, like exploration versus exploitation

# Formalism: Objective

- Expected reward       $\text{Reward}(\text{alg}) = \mathbb{E} \left[ \sum_{\ell=1}^L \bar{V}^{\pi_{\ell-1}} \right]$
- “Bayesian decision problem”
  - RL algorithms can be viewed as heuristics for this problem

- Regret

$$\text{Regret}(\text{alg}) = \mathbb{E} \left[ \sum_{\ell=1}^L \left( \bar{V}^* - \bar{V}^{\pi_{\ell-1}} \right) \right]$$



# Atari Arcade Environment



Atari 2600 supports 49 different games  
develop an RL algorithm that learns to play  
**all games of this flavor**

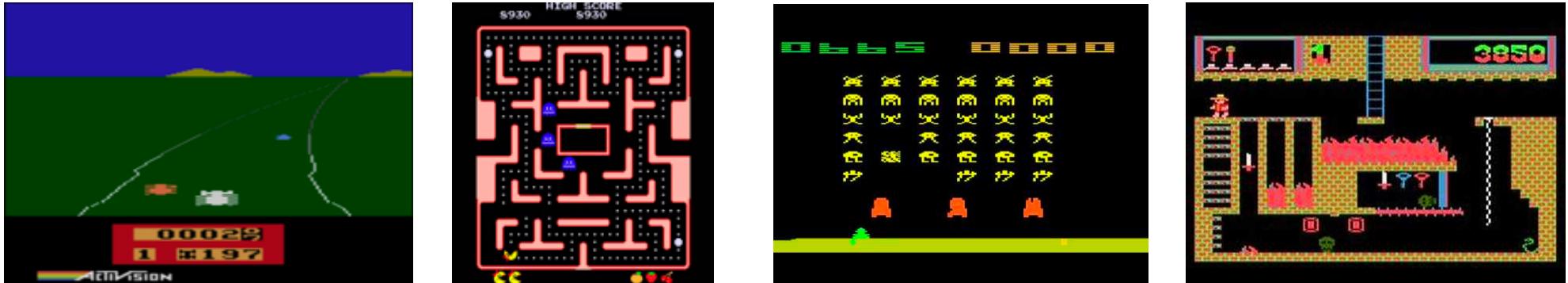
# Arcade Games: Formulation

- Each game is an environment  $\mathcal{E} = (\mathcal{O}, \mathcal{A}, P, \rho)$ 
  - Time increment: 1/60 of a second
- Observations  $O_t \in \mathcal{O}$ 
  - Video frame
- Actions  $A_t \in \mathcal{A}$ 
  - 18 joystick modes
- Probabilities  $\rho, P$ 
  - As produced by game
- Reward  $R_{t+1} = r(H_{t+1})$ 
  - -1 if score decreasing
  - +1 if score increasing
  - 0 otherwise

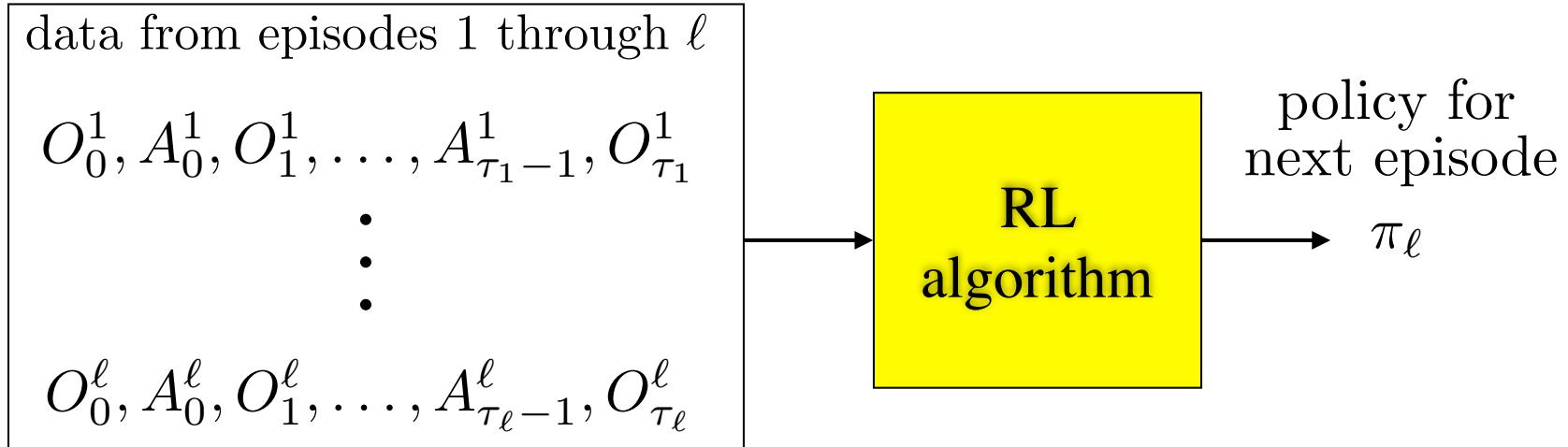


# Arcade Games: RL

- Epistemic uncertainty
  - $\rho$  and  $P$  are random variables



- RL Algorithm applied over episodes of the same game



# Outline

- Examples
  - Mathematical formalism
  - Common algorithms
  - Challenges of data efficiency
- 

# Knowledge Representation

- Model learning: learn  $P$ 
  - Left with a computationally intractable planning problem
  - Exponential explosion of roll-forward errors
- Policy learning: learn  $\pi^*$
- Value function learning: learn  $Q^*$

# Generic RL Loop

---

**Algorithm 1** live

---

**Input:** agent  
          environment

```
1: for  $\ell$  in  $(1, 2, \dots)$  do
2:   agent.learn_from_buffer()
3:   environment.reset()
4:   while environment.history is not terminal do
5:     action  $\leftarrow$  agent.act(environment.history)
6:     environment.step(action)
7:   end while
8:   agent.update_buffer(environment.history)
9: end for
```

---

<b>object</b>	environment
<b>attributes</b>	history
<b>methods</b>	reset step

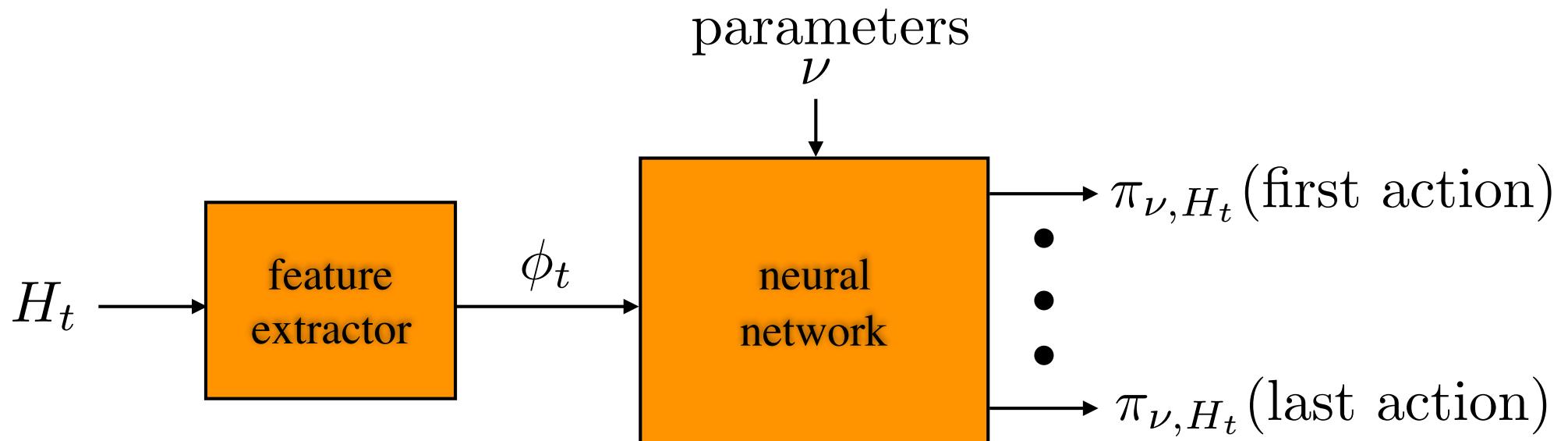
<b>object</b>	agent
<b>methods</b>	learn_from_buffer act update_buffer

# Agent Methods

- act
  - policy learning: use stored policy
  - value function learning: greedy
  - possibly dither, e.g.,  $\epsilon$ -greedy
- update\_buffer
  - append history
  - possibly subsample or drop very old samples
  - possibly compress to what is needed for learning
- learn\_from\_buffer
  - update parameters of policy and/or value function
  - sophisticated algorithms
  - we will focus on this part

# Policy Representation

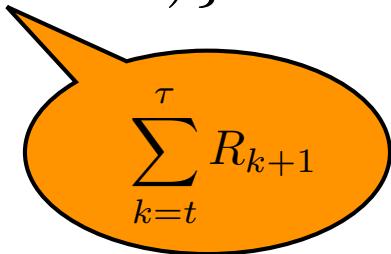
- Policy parameterized by vector  $\nu$
- Example
  - neural network with multinomial logistic output



# The *Reinforce* Algorithm

```
policy =  $\tilde{\pi}_{\text{parameters}=\cdot, \text{features}=\cdot}(\text{action} = \cdot)$ 
```

```
buffer = {(features, action, subsequent reward)}
```



---

## Algorithm 2 learn\_reinforce

---

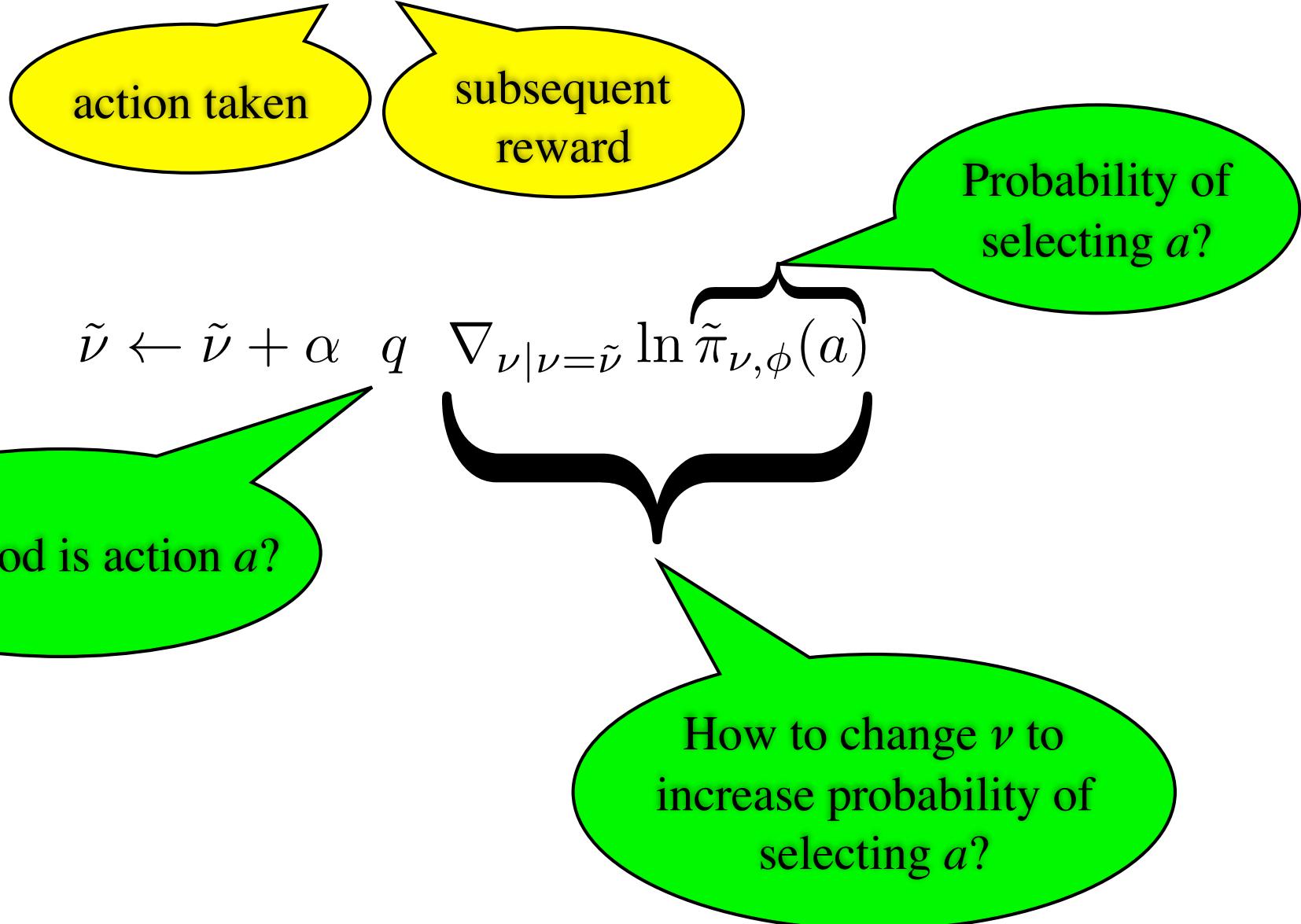
**Input:** agent

```
1:  $\tilde{\pi} \leftarrow \text{agent.policy}$ 
2:  $\tilde{\nu} \leftarrow \text{agent.parameters}$ 
3: buffer  $\leftarrow \text{agent.buffer}$ 
4:  $N \leftarrow \text{agent.num_iters}$ 
5:  $\alpha \leftarrow \text{agent.learning_rate}$ 
6: for  $n$  in  $(1, 2, \dots, N)$  do
7:    $\tilde{\mathcal{D}} \leftarrow \text{buffer.sample_minibatch}()$ 
8:    $\tilde{\nu} \leftarrow \tilde{\nu} + \alpha \sum_{(\phi, a, q) \in \tilde{\mathcal{D}}} q \nabla_{\nu | \nu = \tilde{\nu}} \ln \tilde{\pi}_{\nu, \phi}(a)$ 
9: end for
10:  $\text{agent.parameters} = \tilde{\nu}$ 
```

---

# Intuition about the *Reinforce* Parameter Update

- Consider case with minibatch of size 1
  - Data point  $(\phi, a, q)$



# *Reinforce* is a *Policy Gradient* Algorithm

- Reinforce is a *policy gradient* method
  - The mean field follows the gradient of  $\bar{V}^{\pi_\nu}$ 
    - Under certain technical conditions...
  - Agent performance improves over episodes
  - Popular due to reliable improvement
- An *on-policy* algorithm
  - Estimates gradient at buffered policy
  - To be effectiveness, only retain recent data
  - Lack of *counterfactual analysis*
    - As opposed to *off-policy* algorithms

# Policy Gradient Methods: Obstacles

- SGD can be slow and converge to local minima
  - Use adaptive step size rules as with supervised learning
- Statistically inefficient
  - On-policy learning
  - Naive exploration
  - High variance of gradient estimates
    - Variance reduction techniques
    - Baselining

# Variance Reduction Through Baselining

- Thrashing if  $q$  is large
  - Reducing step size makes learning too slow when  $q$  is small

$$\tilde{\nu} \leftarrow \tilde{\nu} + \alpha \quad q \quad \nabla_{\nu|_{\nu=\tilde{\nu}}} \ln \tilde{\pi}_{\nu,\phi}(a)$$

- Baselined update with any  $B(\phi)$  retains policy gradient

$$\tilde{\nu} \leftarrow \tilde{\nu} + \alpha \quad (q - B(\phi)) \quad \nabla_{\nu|_{\nu=\tilde{\nu}}} \ln \tilde{\pi}_{\nu,\phi}(a)$$

- A good baseline is given by  $V^{\tilde{\pi}_{\tilde{\nu}}}$ 
  - Use approximation  $\tilde{V}_{\theta}(\phi) \approx V^{\tilde{\pi}_{\tilde{\nu}}}$

# Learning a Baseline

---

## Algorithm 3 learn\_reinforce\_baseline

---

**Input:** agent

```
1:  $\tilde{\pi} \leftarrow \text{agent.policy}$ 
2:  $\tilde{V} \leftarrow \text{agent.value\_function}$ 
3:  $(\tilde{\nu}, \tilde{\theta}) \leftarrow \text{agent.parameters}$ 
4: buffer  $\leftarrow \text{agent.buffer}$ 
5:  $N \leftarrow \text{agent.num\_iters}$ 
6:  $\alpha \leftarrow \text{agent.learning\_rate}$ 
7: for  $n$  in  $(1, 2, \dots, N)$  do
8:    $\tilde{\mathcal{D}} \leftarrow \text{buffer.sample\_minibatch()}$ 
9:    $\tilde{\nu} \leftarrow \tilde{\nu} + \alpha \sum_{(\phi, a, q) \in \tilde{\mathcal{D}}} (q - \tilde{V}_{\tilde{\theta}}(\phi)) \nabla_{\nu|_{\nu=\tilde{\nu}}} \ln \tilde{\pi}_{\nu, \phi}(a)$ 
10:   $\tilde{\theta} \leftarrow \tilde{\theta} + \alpha \sum_{(\phi, a, q) \in \tilde{\mathcal{D}}} \left( \nabla_{\theta|\theta=\tilde{\theta}} \tilde{V}_{\theta}(\phi) \right) \left( q - \tilde{V}_{\tilde{\theta}}(\phi) \right)$ 
11: end for
12: agent.parameters =  $(\tilde{\nu}, \tilde{\theta})$ 
```



adapt baseline  
to approximate  
 $V^{\tilde{\pi}_{\tilde{\nu}}}$

- *Actor*: parameterized policy
- *Critic*: parameterized value function
- Broader class of *actor-critic methods*...

# Actor-Critic Algorithm

- Actor can learn directly from the action-value critic rather than observed rewards

---

**Algorithm 4** learn\_actor\_critic

---

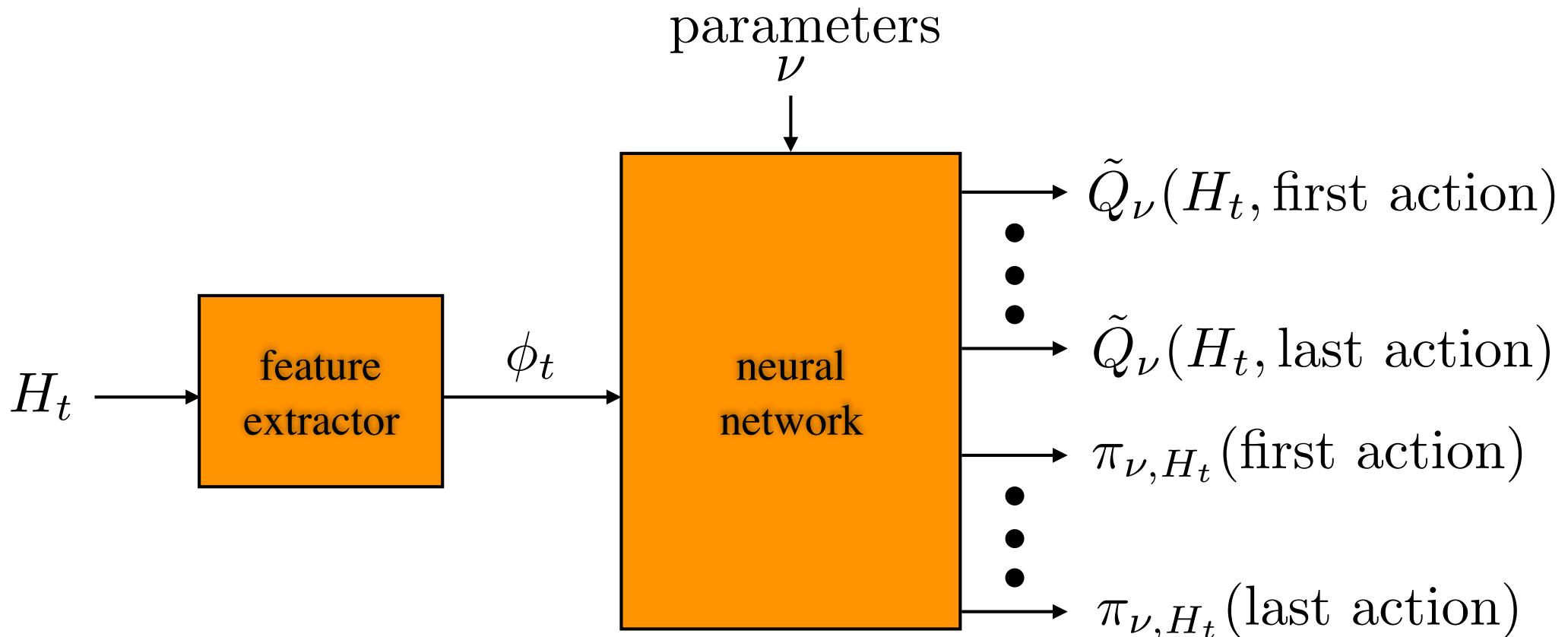
**Input:** agent

```
1:  $\tilde{\pi} \leftarrow \text{agent.policy}$ 
2:  $\tilde{Q} \leftarrow \text{agent.value\_function}$ 
3:  $(\tilde{\nu}, \tilde{\theta}) \leftarrow \text{agent.parameters}$ 
4: buffer  $\leftarrow \text{agent.buffer}$ 
5:  $N \leftarrow \text{agent.num\_iters}$ 
6:  $\alpha \leftarrow \text{agent.learning\_rate}$ 
7: for  $n$  in  $(1, 2, \dots, N)$  do
8:    $\tilde{\mathcal{D}} \leftarrow \text{buffer.sample\_minibatch}()$ 
9:    $\tilde{\nu} \leftarrow \tilde{\nu} + \alpha \sum_{(\phi, a, q) \in \tilde{\mathcal{D}}} \left( \tilde{Q}_{\tilde{\theta}}(\phi, a) - \sum_{a' \in \mathcal{A}} \tilde{\pi}_{\tilde{\nu}, \phi}(a') \tilde{Q}_{\tilde{\theta}}(a') \right) \nabla_{\nu | \nu = \tilde{\nu}} \ln \tilde{\pi}_{\nu, \phi}(a)$ 
10:   $\tilde{\theta} \leftarrow \tilde{\theta} + \alpha \nabla_{\theta | \theta = \tilde{\theta}} \sum_{(\phi, a, q) \in \tilde{\mathcal{D}}} \tilde{Q}_{\theta}(\phi, a) \left( q - \tilde{Q}_{\tilde{\theta}}(\phi, a) \right)$ 
11: end for
12:  $\text{agent.parameters} = (\tilde{\nu}, \tilde{\theta})$ 
```

---

# Combined Actor-Critic Architecture

- Single neural network
  - Shared features
    - Policy-gradient theorem
  - Shared hidden layers

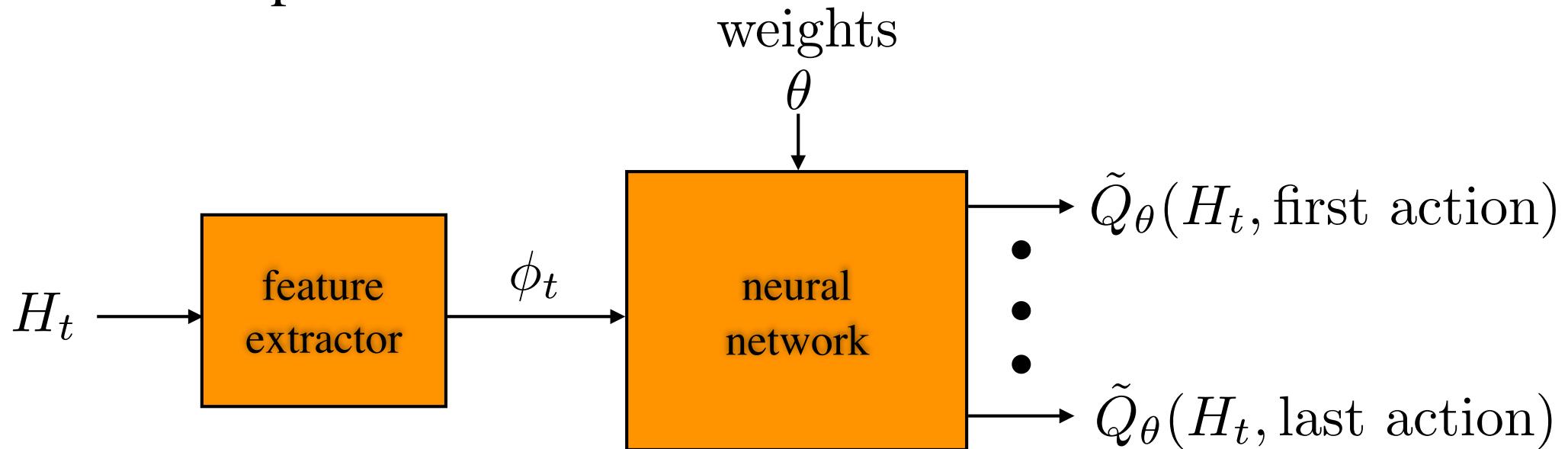


# Value Function Learning

- Represent knowledge as parameterized value function
  - Possibly the optimal value function
  - Enables off-policy learning

$$\tilde{Q}_\theta(h, a) \approx Q^*(h, a)$$

- Example



# Incremental TD

---

**Algorithm 5** learn\_incremental\_td

---

**Input:** agent

```
1:  $\tilde{Q} \leftarrow \text{agent.value\_function}$ 
2:  $\tilde{\theta} \leftarrow \text{agent.parameters}$ 
3: buffer  $\leftarrow \text{agent.buffer}$ 
4:  $N \leftarrow \text{agent.num\_iters}$ 
5:  $\alpha \leftarrow \text{agent.learning\_rate}$ 
6: for  $n$  in  $(1, 2, \dots, N)$  do
7:    $\tilde{\mathcal{D}} \leftarrow \text{buffer.sample\_minibatch()}$ 
8:    $\delta \leftarrow \text{buffer.minibatch\_size} / \text{buffer.size}$ 
9:    $\tilde{\theta} \leftarrow \tilde{\theta} - \alpha \nabla_{\theta|\theta=\tilde{\theta}} (\mathcal{L}(\theta, \tilde{\theta}, \tilde{\mathcal{D}}) + \delta \mathcal{R}(\theta))$  ▷ TD step
10: end for
```

---

TD loss 
$$\mathcal{L}(\theta, \theta^-, \mathcal{D}) = \frac{1}{\sigma_{\text{outcome}}^2} \sum_{(\phi, a, r, \phi') \in \mathcal{D}} \left( r + \max_{a' \in \mathcal{A}} \tilde{Q}_{\theta^-}(\phi', a') - \tilde{Q}_\theta(\phi, a) \right)$$

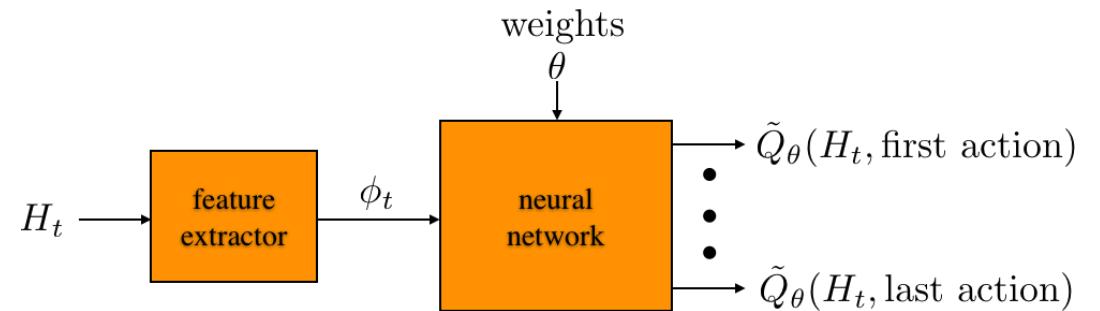
regularizer 
$$\mathcal{R}(\theta) = \frac{1}{\sigma_{\text{prior}}^2} \|\theta\|_2^2$$

# Off-Policy Learning

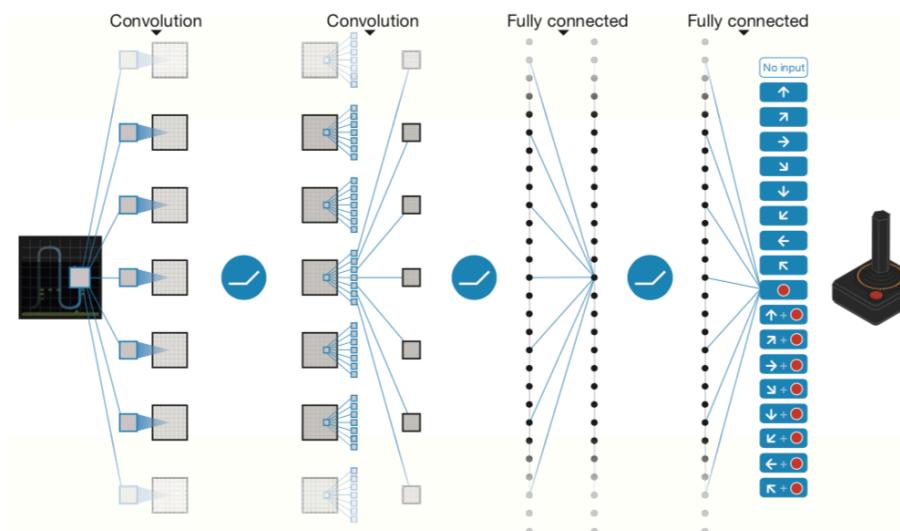
- TD error aligns prediction with max at next step
  - Not total reward under current policy

TD loss       $\mathcal{L}(\theta, \theta^-, \mathcal{D}) = \frac{1}{\sigma_{\text{outcome}}^2} \sum_{(\phi, a, r, \phi') \in \mathcal{D}} \left( r + \max_{a' \in \mathcal{A}} \tilde{Q}_{\theta^-}(\phi', a') - \tilde{Q}_\theta(\phi, a) \right)$

# Representation for Arcade RL Algorithm

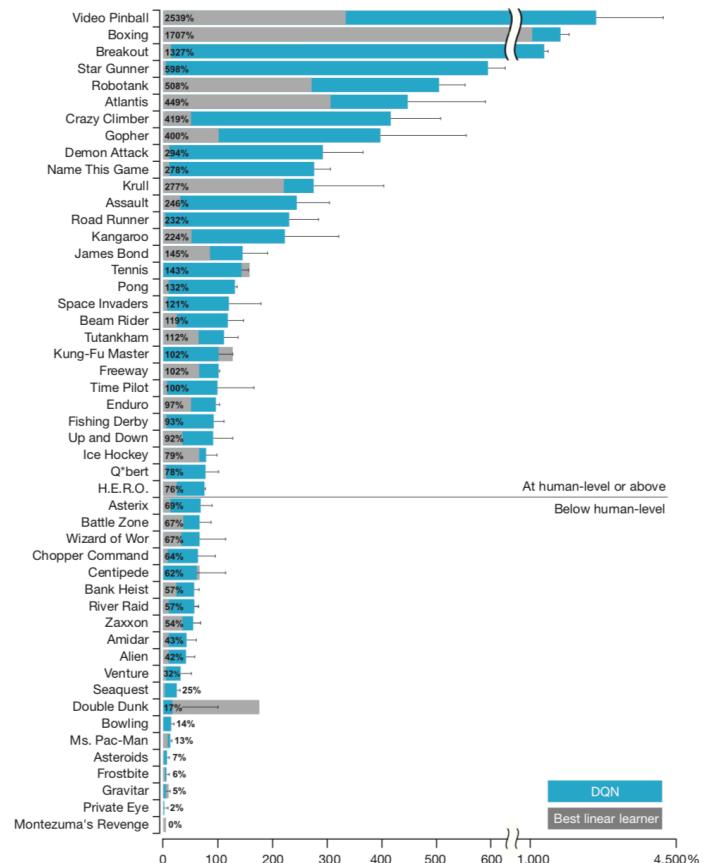


- Features
  - Last four frames, with some preprocessing
- Neural network architecture
  - Convolutional, RELUs, one linear output node per action



# Performance in Atari Arcade Environment

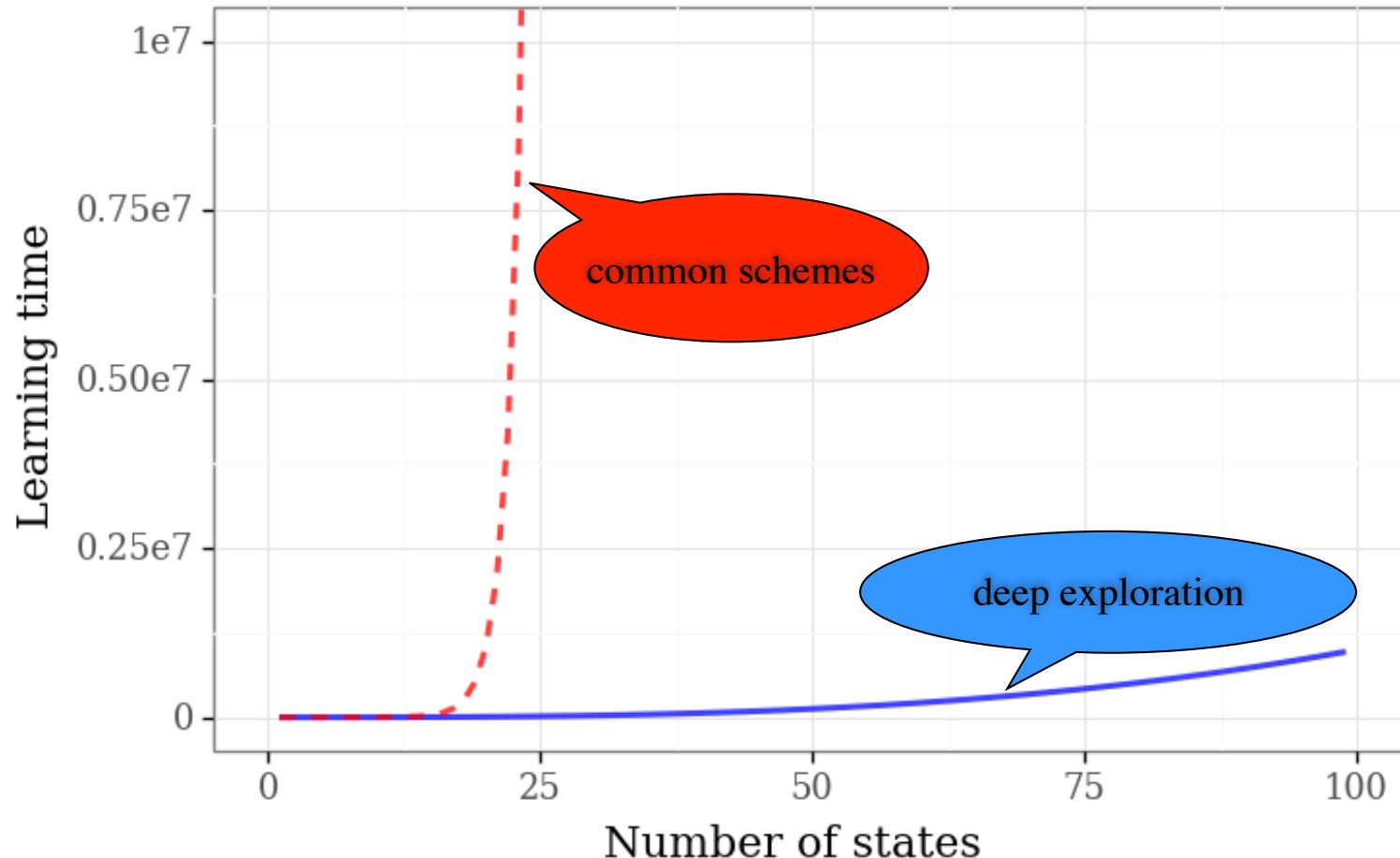
- DQN (incremental TD)
  - outperforms humans in most games
- Policy gradients approaches
  - perform even better
  - require much more data



# Outline

- Examples
- Mathematical formalism
- Common algorithms
- ● Challenges of data efficiency

# Big Data vs Smart Data

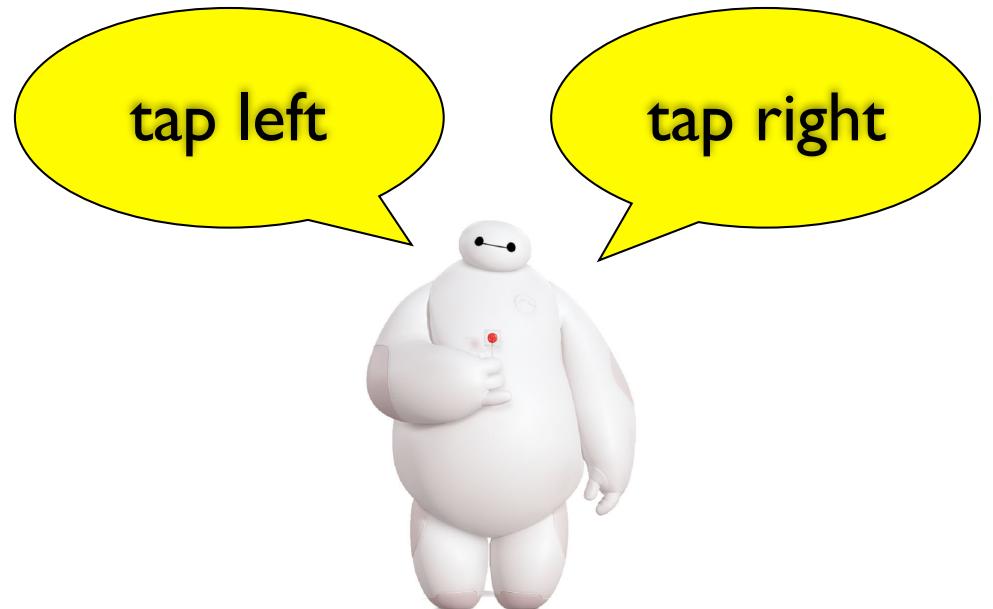
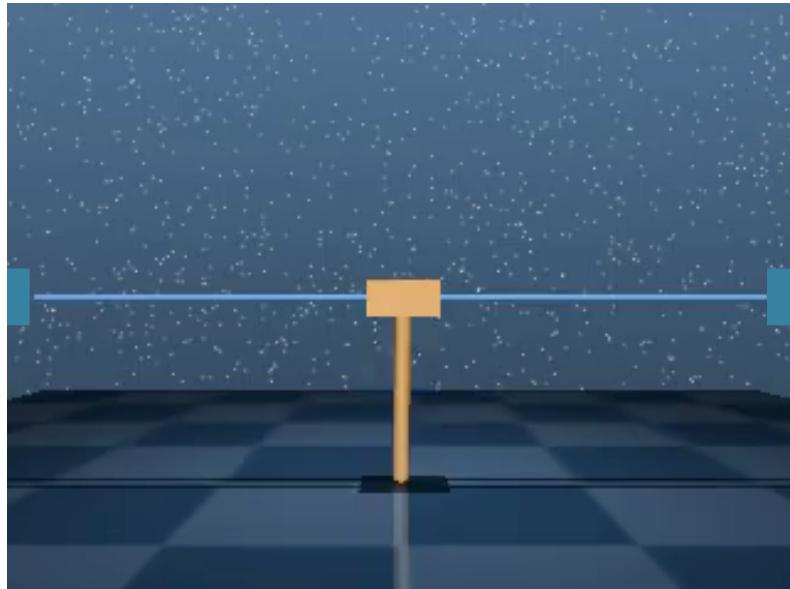


common exploration schemes can require exponentially large data sets even for simple reinforcement learning tasks

# Representing and Resolving Uncertainty

- Many common RL methods learn point estimates
- Efficient exploration requires uncertainty representation
- How to learn neural networks that do this

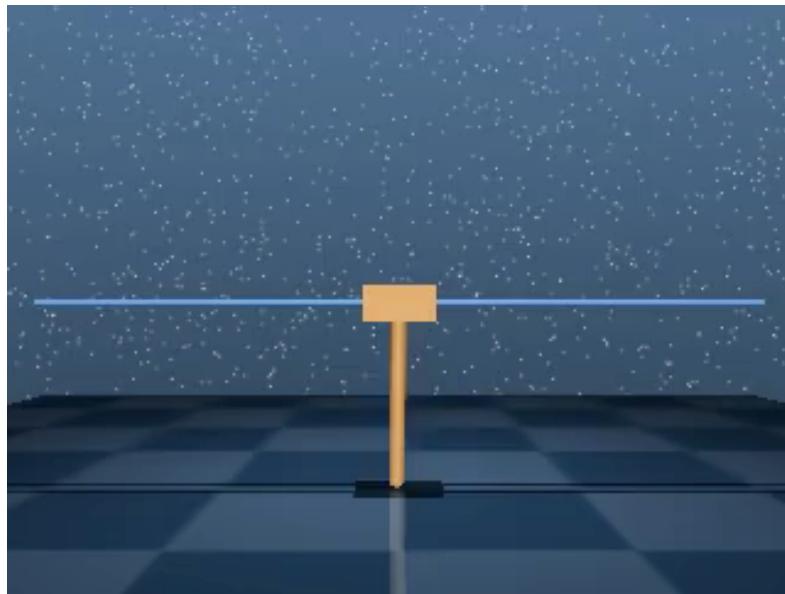
# Cart-Pole Problem



20 seconds per episode

# Cart-Pole Swing-Up

$\epsilon$ -greedy



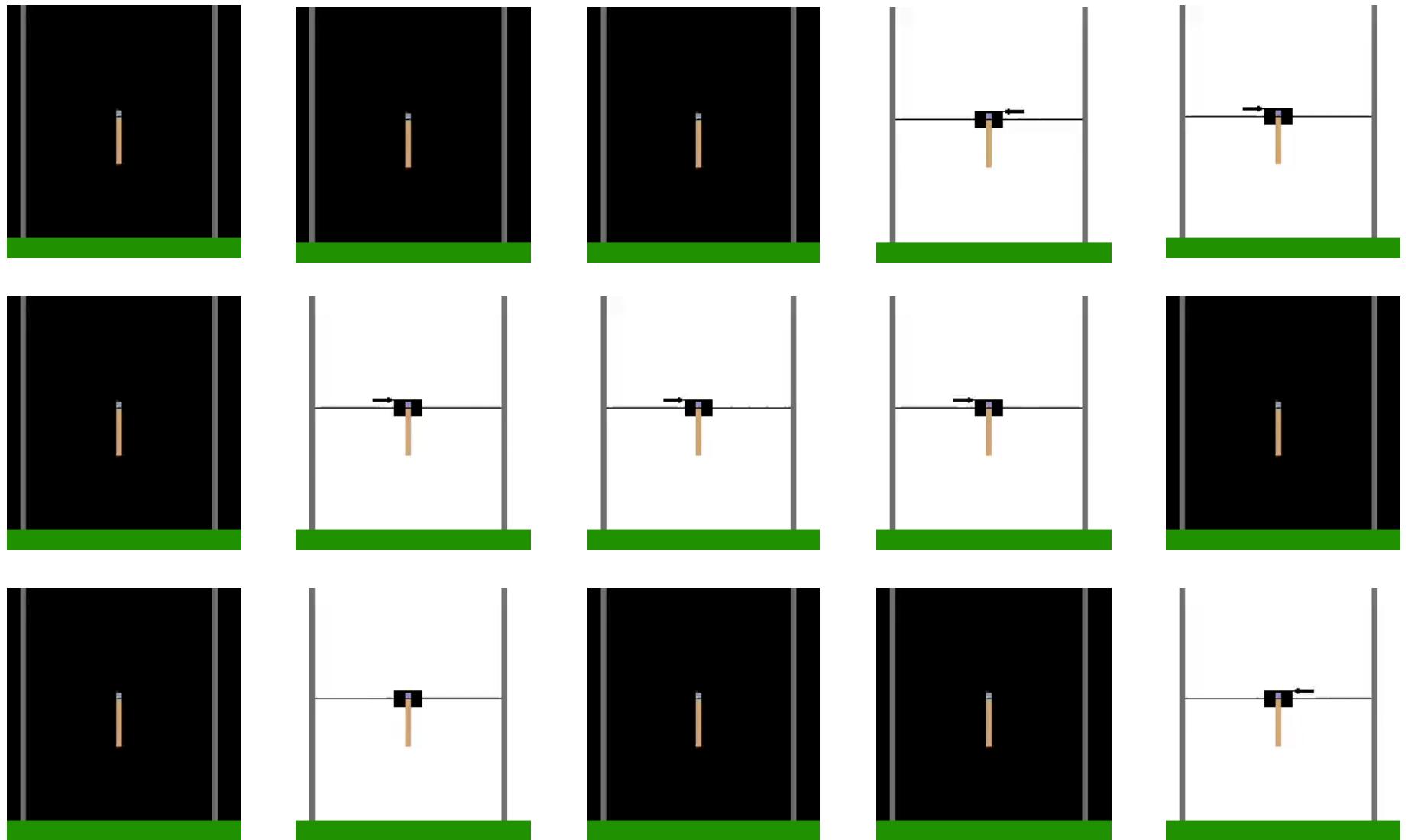
deep exploration



after 1000 episodes

after 1000 episodes

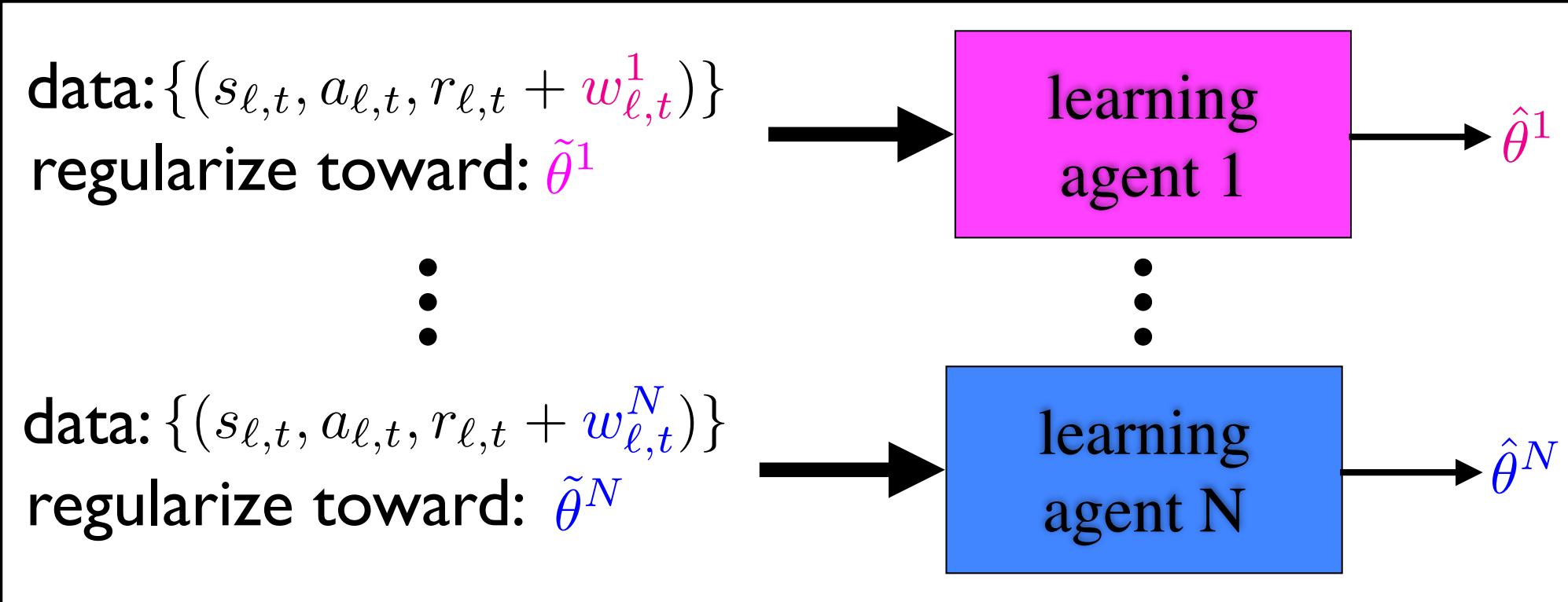
# Information Sharing and Coordination



100 agents total

# Ensemble Sampling

- Learn ensemble of  $N$  value functions, streaming data



- Before each episode, sample one uniformly, then argmax
- Other incremental approaches...

# Hierarchical Representations

- Specific model classes, like linear, are too restrictive
- Can we do better than “flat” neural networks?
- CNNs induce useful inductive biases for vision
- What inductive biases are right for RL?
- HRL: effective policies decompose hierarchically

# Leveraging Rich Observations

- Common methods often focus on predicting and maximizing future reward
- Richer observations can dramatically accelerate learning
- Don't learn to optimize GDP based only on GDP
- Robots should learn from outcomes, not only rewards

