

CS 4803 / 7643: Deep Learning

Topics:

- (Deep) Reinforcement Learning
- Closing time

Dhruv Batra
Georgia Tech

Administrativa

- Last class today
- Project submission
 - Due: 12/04, 11:55pm
 - Last deliverable in the class
 - Can't use late days
 - <https://piazza.com/class/jkujs03pgu75cd?cid=225>

Recap from last time

Types of Learning

- Supervised learning

- Learning from a “teacher”
- Training data includes desired outputs

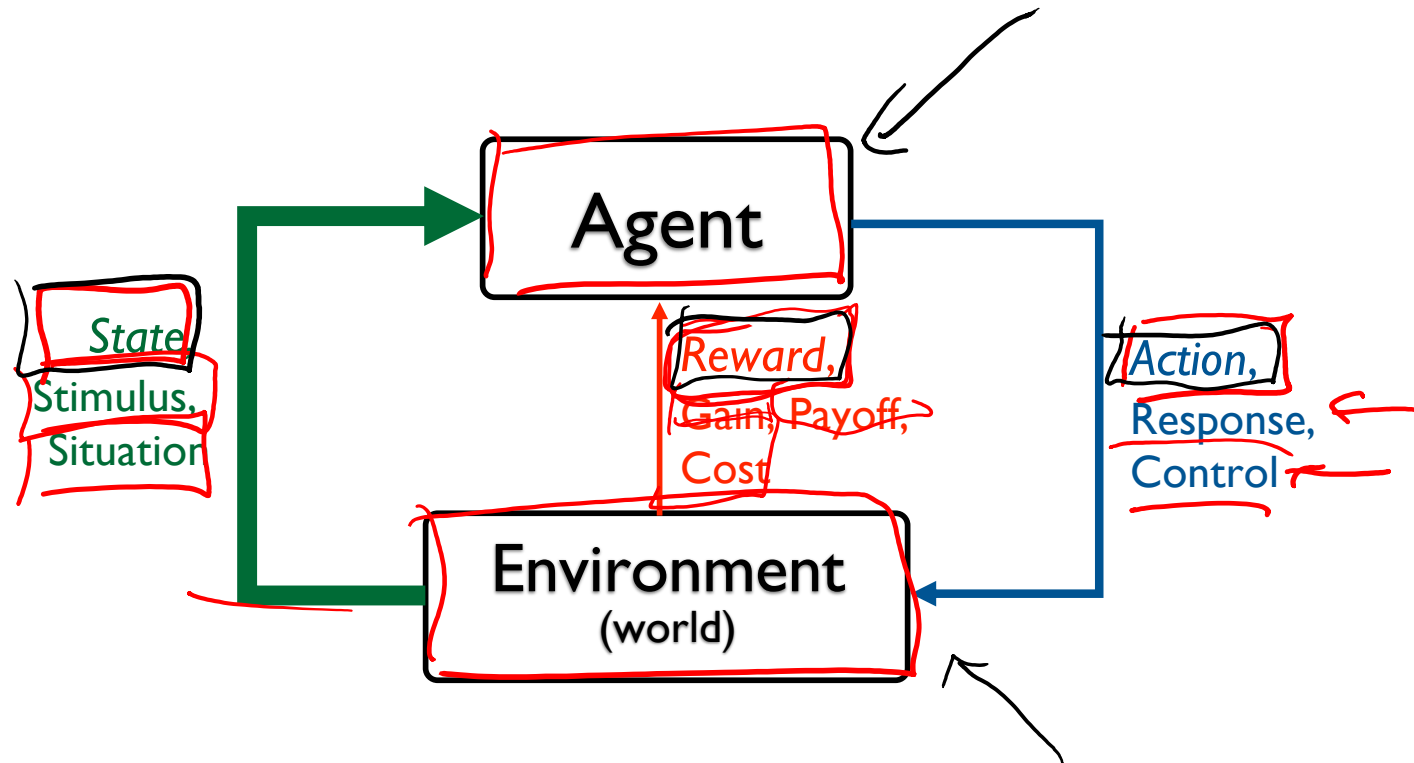
- Unsupervised learning

- Discover structure in data
- Training data does not include desired outputs

- Reinforcement learning

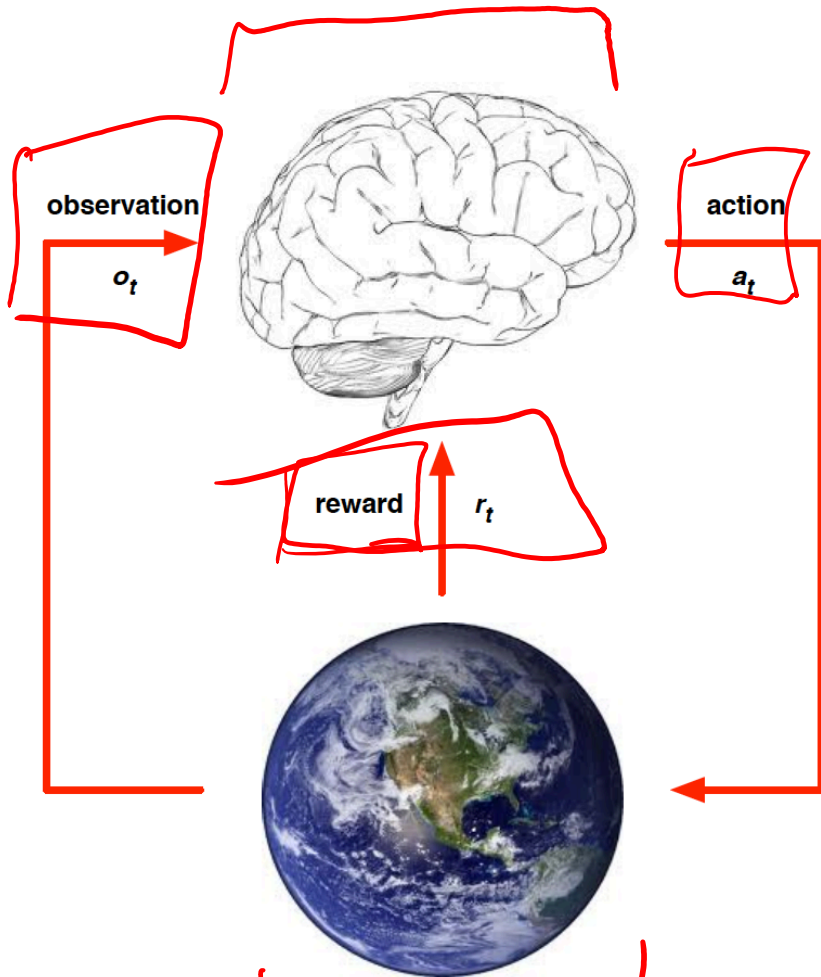
- Learning to act under evaluative feedback (rewards)

RL API



- Environment may be unknown, nonlinear, stochastic and complex
- Agent learns a policy mapping states to actions
 - Seeking to maximize its cumulative reward in the long run

RL API



- ▶ At each step t the agent:
 - ▶ Executes action a_t
 - ▶ Receives observation o_t
 - ▶ Receives scalar reward r_t

- ▶ The environment:
 - ▶ Receives action a_t
 - ▶ Emits observation o_{t+1}
 - ▶ Emits scalar reward r_{t+1}

RL

$$r_t(s_t, a_t)$$

RL	SL
s_t	x_t
a_t	y_t

SL

$$l(y^*, \hat{y}(z))$$

$$o_t = f(s_t)$$

Signature challenges of RL

↓ *How*

s_0, a_1, s_1, \dots, i

- Evaluative feedback (reward)
- Sequentiality, delayed consequences
- Need for trial and error, to explore as well as exploit

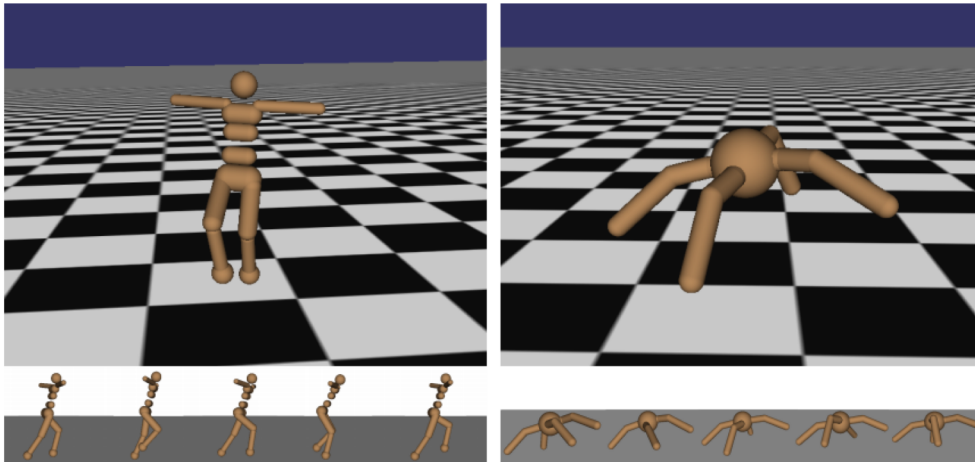
- Non-stationarity

$$x, y \sim P_{data}()$$

$$S \sim \left[\begin{array}{c} (s_{t+1} | s_t) \\ \boxed{q_t} \end{array} \right]$$

- The fleeting nature of time and online data

Robot Locomotion



Objective: Make the robot move forward

State: Angle and position of the joints!

Action: Torques applied on joints

Reward: 1 at each time step upright + forward movement

Figures copyright John Schulman et al., 2016. Reproduced with permission.

Atari Games



Objective: Complete the game with the highest score

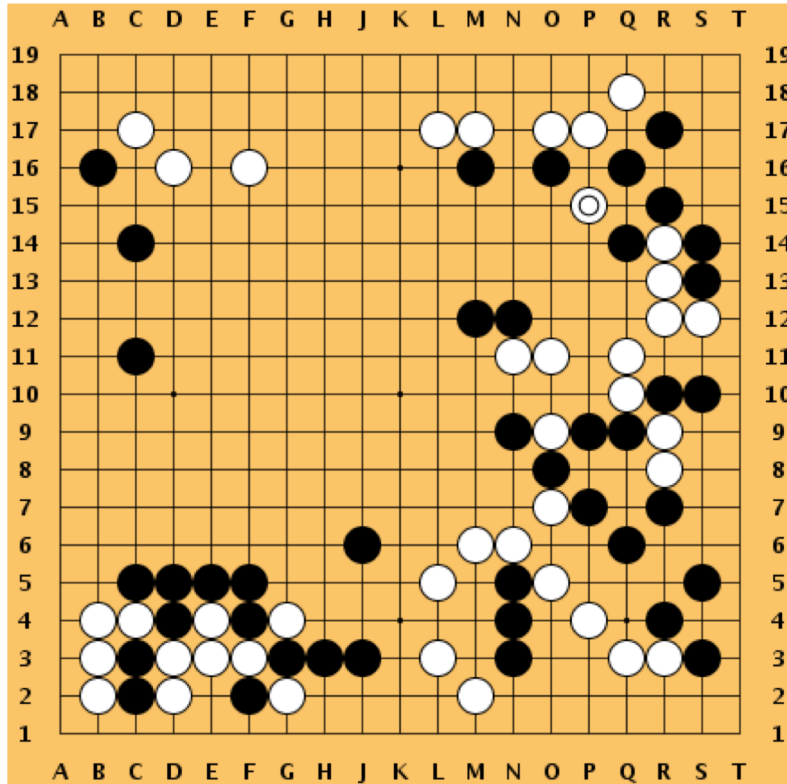
State: Raw pixel inputs of the game state

Action: Game controls e.g. Left, Right, Up, Down

Reward: Score increase/decrease at each time step

Figures copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Go



Objective: Win the game!

State: Position of all pieces

Action: Where to put the next piece down

Reward: 1 if win at the end of the game, 0 otherwise

[This image is CC0 public domain](#)

Markov Decision Process

- Mathematical formulation of the RL problem

Defined by: (S, A, R, P, γ)

π

S : set of possible states

A : set of possible actions

R : distribution of reward given (state, action) pair

P : transition probability i.e. distribution over next state given (state, action) pair

γ : discount factor

$$P(\underline{S}_{t+1} \mid \underline{S}_t, \underline{a}_t)$$

Markov Decision Process

- Mathematical formulation of the RL problem

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

\mathcal{S} : set of possible states

\mathcal{A} : set of possible actions

\mathcal{R} : distribution of reward given (state, action) pair

\mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair

γ : discount factor

- Life is trajectory: $\dots [S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots]$

- Markov property: Current state completely characterizes the state of the world

$$p(\underline{r}, \underline{s}' | s, a) = \text{Prob} \left[\underline{R_{t+1} = r}, \underline{S_{t+1} = s'} \mid \underline{S_t = s}, \underline{A_t = a} \right]$$

Components of an RL Agent

- Policy
 - How does an agent behave?
- Value function
 - How good is each state and/or state-action pair?
- Model
 - Agent's representation of the environment

Policy

- A policy is how the agent acts
- Formally, map from states to actions

Deterministic policy: $a = \pi(s)$

Stochastic policy: $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$

$\pi(\cdot)$

e.g.

State	Action
A	2
B	1
⋮	

The optimal policy π^*

What's a good policy?

Maximizes current reward? Sum of all future reward?

Discounted future rewards!

$$\sum_{t=0}^{\infty} \gamma^t r_t + \underbrace{\gamma^2 r_{t+2}} + \underbrace{\gamma^3 r_{t+3}} \dots$$

The image shows a handwritten equation representing discounted future rewards. The first term is $\sum_{t=0}^{\infty} r_t$, with a box around the index t . The second term is $\gamma^2 r_{t+2}$, where the γ^2 is written above the term and the term is underlined. The third term is $\gamma^3 r_{t+3}$, where the γ^3 is circled and the term is underlined. A red γ is written above the first underlined term. The equation is followed by an ellipsis.

The optimal policy π^*

What's a good policy?

Maximizes current reward? Sum of all future reward?

Discounted future rewards!

Formally:

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi \right]$$

with

$$s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$$

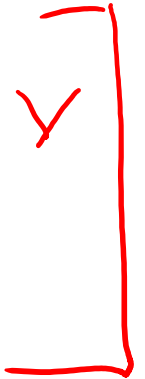
$$a_t = \pi(s_t) \quad s_{t+1} = \mathcal{P}(s_t, a_t)$$

Components of an RL Agent

- Policy
 - How does an agent behave?
- Value function
 - How good is each state and/or state-action pair?
- Model
 - Agent's representation of the environment

Value Function

- A value function is a prediction of future reward
- “State Value Function” or simply “Value Function”
 - How good is a state?
 - Am I screwed? Am I winning this game?
- “Action Value Function” or Q-function
 - How good is a state action-pair?
 - Should I do this now?



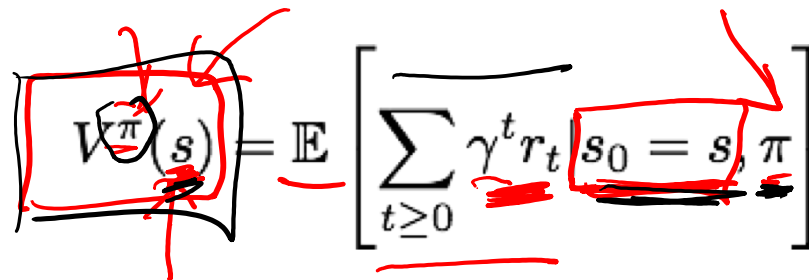
Q

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

How good is a state?

The **value function** at state s , is the expected cumulative reward from state s (and following the policy thereafter):



The diagram shows the equation $V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$ with several red annotations. A red box encloses the left side of the equation, $V^\pi(s)$. Another red box encloses the right side, $\mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$. A red arrow points from the s in the left box to the $s_0 = s$ in the right box. There are also red scribbles and lines around the summation symbol and the expectation operator.

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$$

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

How good is a state?

The **value function** at state s , is the expected cumulative reward from state s (and following the policy thereafter):

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid \underline{s_0 = s, \pi} \right]$$

How good is a state-action pair?

The **Q-value function** at state s and action a , is the expected cumulative reward from taking action a in state s (and following the policy thereafter):

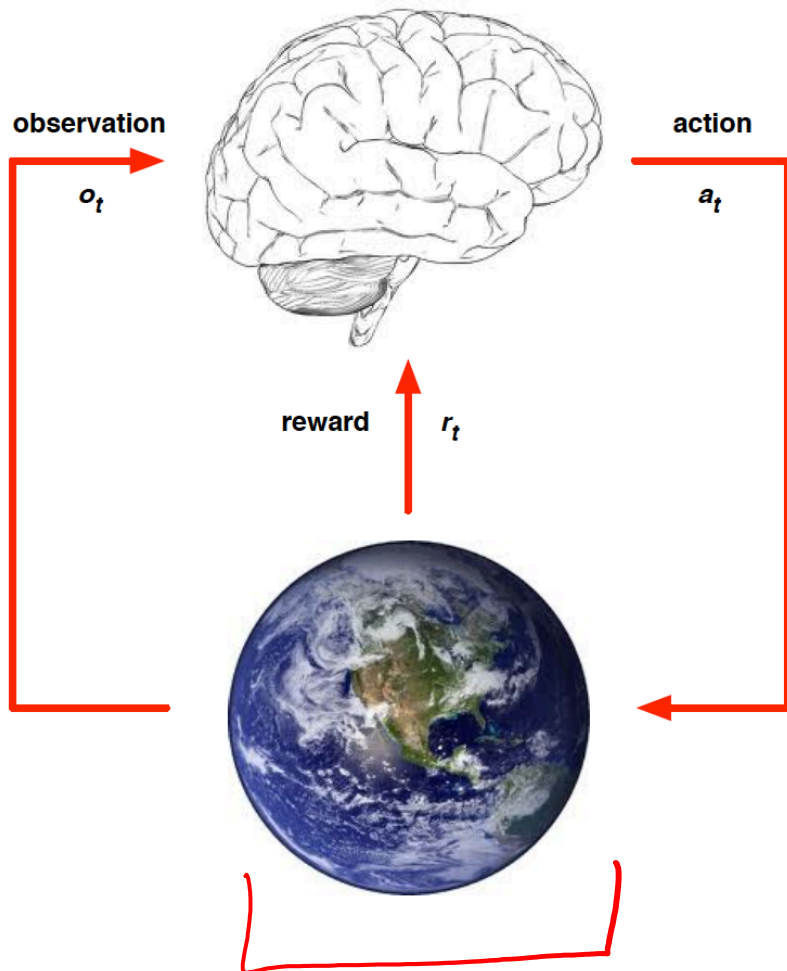
\max_a

$$\underline{Q^\pi(s, a)} = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid \underline{s_0 = s, a_0 = a, \pi} \right]$$

Components of an RL Agent

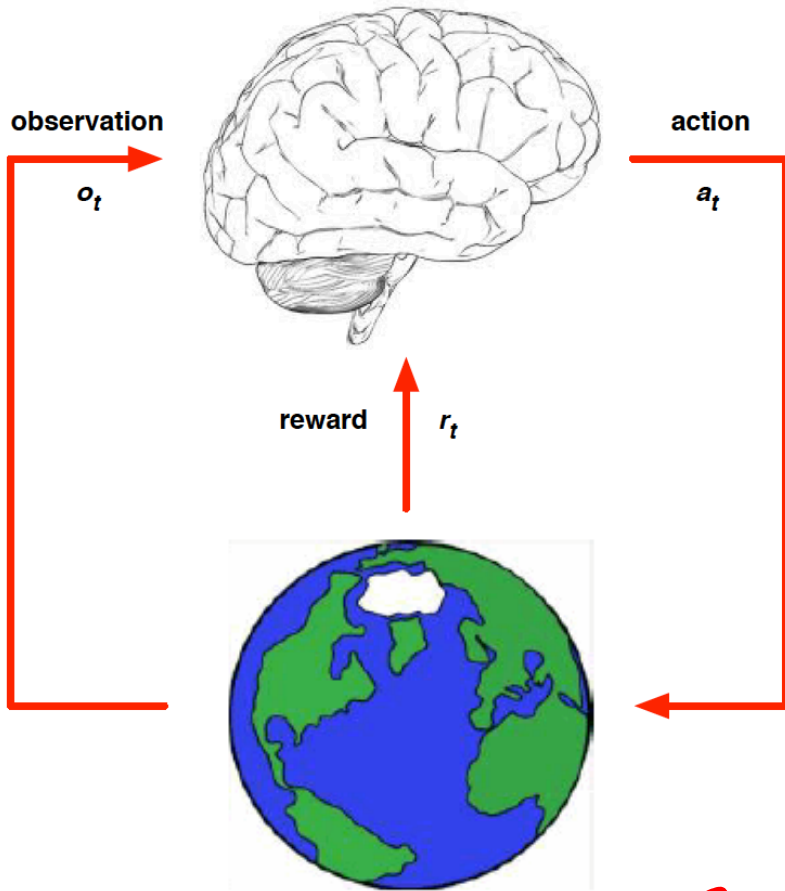
- Policy
 - How does an agent behave?
- Value function
 - How good is each state and/or state-action pair?
-)• Model
 - Agent's representation of the environment

Model

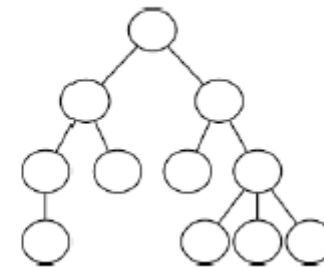


Model

- Model predicts what the world will do next



Model is learnt from experience
Acts as proxy for environment
Planner interacts with model
e.g. using lookahead search



$$s_t, a_t \rightarrow s_{t+1}$$

$$P(s_{t+1} | s_t, a_t)$$

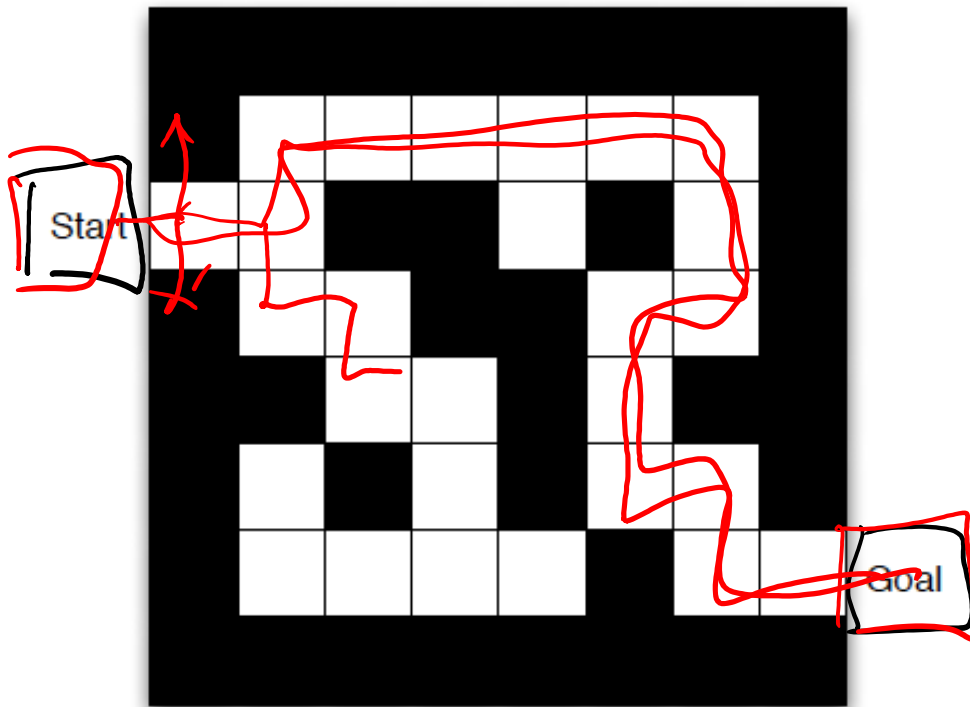
Plan for Today

- (Deep) Reinforcement Learning
 - Policy gradients
- Closing the loop

Components of an RL Agent

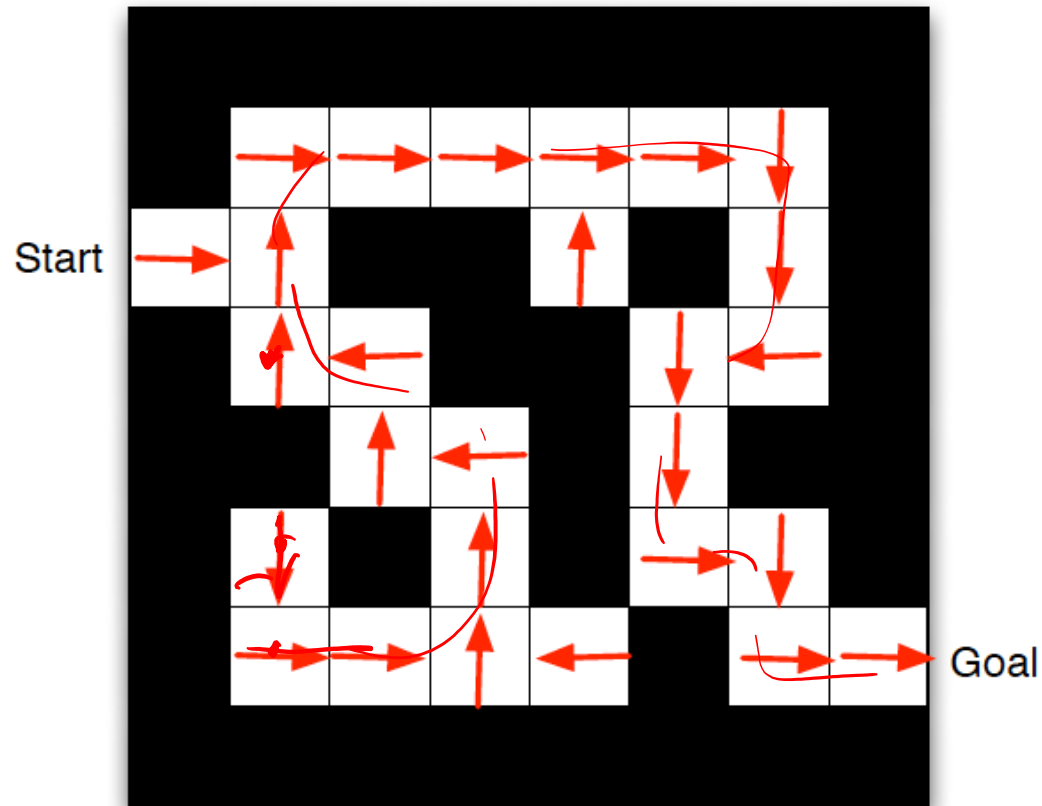
- Policy
 - How does an agent behave?
- Value function
 - How good is each state and/or state-action pair?
- Model
 - Agent's representation of the environment

Maze Example



- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: Agent's location

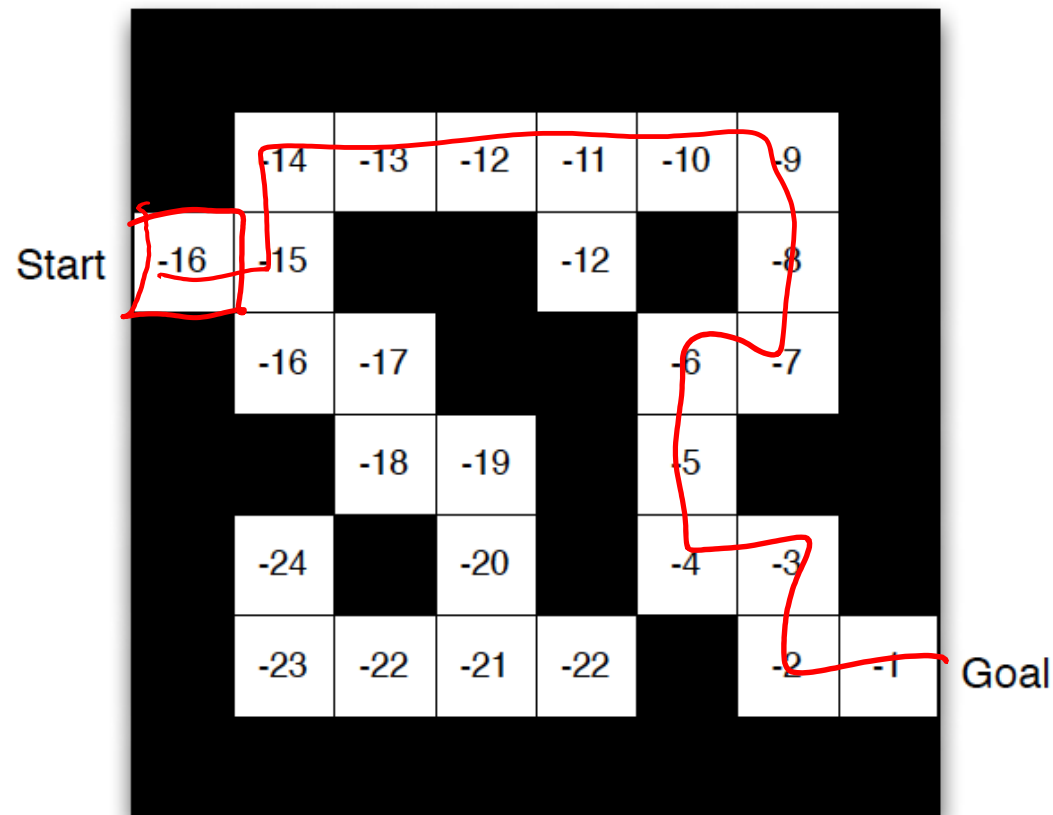
Maze Example: Policy



- Arrows represent policy $\pi(s)$ for each state s

Maze Example: Value

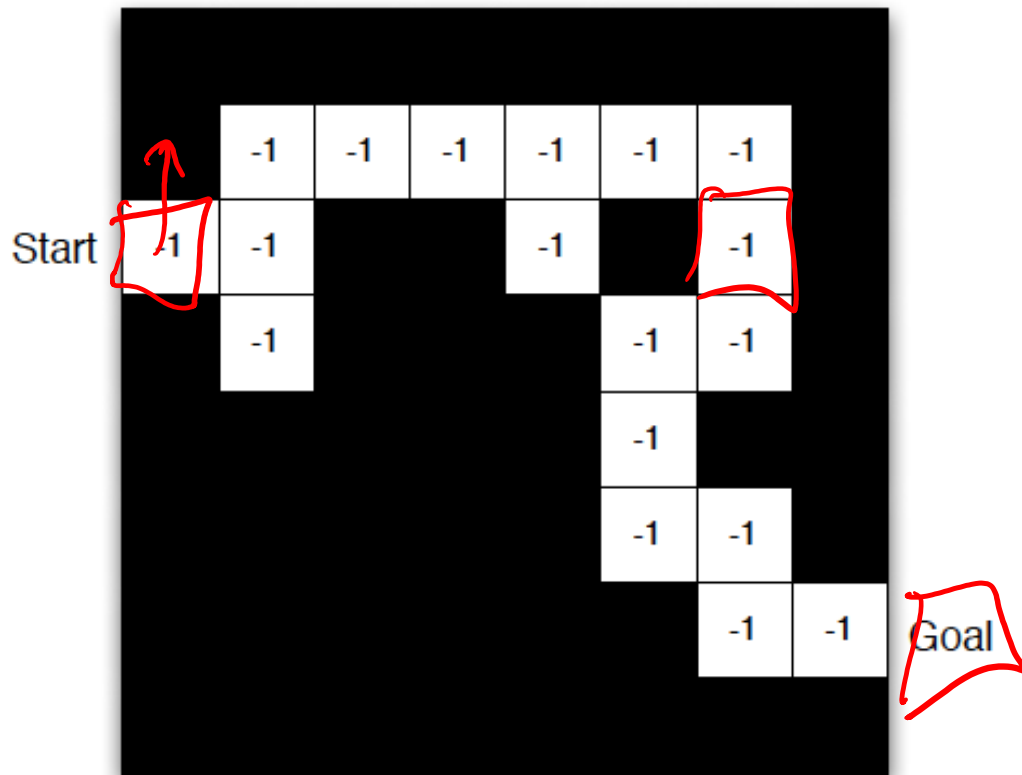
$V(s)$



- Numbers represent value $V_{\pi}(s)$ of each state s

V_{π}^*

Maze Example: Model



- Agent may have an internal model of the environment
- Dynamics: how actions change the state
- Rewards: how much reward from each state
- The model may be imperfect

- Grid layout represents transition model $\mathcal{P}_{ss'}^a$
- Numbers represent immediate reward \mathcal{R}_s^a from each state s (same for all a)

Components of an RL Agent

- Value function
 - How good is each state and/or state-action pair?
- Policy
 - How does an agent behave?
- Model
 - Agent's representation of the environment

Approaches to RL

- Value-based RL

- Estimate the optimal action-value function $Q^*(s, a)$

- Policy-based RL

- Search directly for the optimal policy

π^*

- Model-based RL

- Build a model of the world
 - State transition, reward probabilities
- Plan (e.g. by look-ahead) using model

Deep RL

- Value-based RL

- Use neural nets to represent Q function

$$Q(\bar{s}, \bar{a}; \theta)$$

$$Q(s, a; \theta^*) \approx Q^*(s, a)$$

- Policy-based RL

- Use neural nets to represent policy

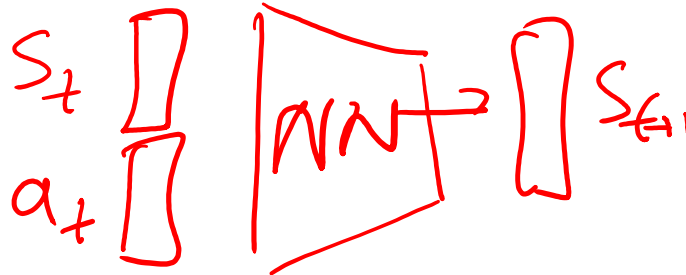
$$\pi_{\theta}$$

$$\pi_{\theta^*} \approx \pi^*$$



- Model

- Use neural nets to represent and learn the model



Deep RL

- Value-based RL

- Use neural nets to represent Q function $Q(s, a; \theta)$

$$Q(s, a; \theta^*) \approx Q^*(s, a)$$

- **Policy-based RL**

- Use neural nets to represent policy π_θ

$$\pi_{\theta^*} \approx \pi^*$$

- Model

- Use neural nets to represent and learn the model



Policy Gradients

$$s_t \rightarrow a_t$$

$$s_t \xrightarrow{\theta} a_t$$

$$\pi_{\theta}(a_t | s_t)$$

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

For each policy, define its value:

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

$$\max_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [F(z)]$$

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:

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

We want to find the optimal policy $\underline{\theta^*} = \underline{\arg \max_{\theta} J(\theta)}$

How can we do this?

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:

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

We want to find the optimal policy $\theta^* = \arg \max_{\theta} J(\theta)$

How can we do this?

Gradient ascent on policy parameters!

REINFORCE algorithm

Mathematically, we can write:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$$
$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

Where $r(\tau)$ is the reward of a trajectory

$$\tau = (s_0, a_0, r_0, s_1, a_1, \dots, s_T)$$

Trajectory

$$P(\tau; \theta) = P(s_0, a_0, r_0, \dots)$$
$$P(s_0) \times \prod_{t=1}^T P(s_t, a_t, r_t \mid s_{t-1}, a_{t-1})$$

REINFORCE algorithm

Mathematically, we can write:

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

Where $r(\tau)$ is the reward of a trajectory $\tau = (s_0, a_0, r_0, s_1, \dots)$

world *policy*

$$\begin{aligned} p(\tau; \theta) &= \prod_{t \geq 0} \underbrace{p(s_{t+1} | s_t, a_t)}_{\text{world}} \underbrace{\pi_{\theta}(a_t | s_t)}_{\text{policy}} \\ \log p(\tau; \theta) &= \sum_{t \geq 0} \log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t) \end{aligned}$$

REINFORCE algorithm

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

$$\nabla J(\theta) = \mathbb{E} \left[\underbrace{r(\tau)} \cdot \underbrace{\nabla_{\theta} \log p(\tau; \theta)} \right]$$

REINFORCE algorithm

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

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

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

REINFORCE algorithm

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

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

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

Intractable! Expectation of gradient is problematic when p depends on θ

REINFORCE algorithm

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

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

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

Intractable! Expectation of gradient is problematic when p depends on θ

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

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

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

Intractable! Expectation of gradient is problematic when p depends on θ

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:

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

Can estimate with
Monte Carlo sampling

REINFORCE algorithm

Can we compute those quantities without knowing the transition probabilities?

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

Handwritten annotations:
- A red underline under $p(\tau; \theta)$.
- A red bracket under the product term $\prod_{t \geq 0} p(s_{t+1} | s_t, a_t)$.
- A red bracket under the policy term $\tau_{\theta}(a_t | s_t)$, with the word "policy" written in red above it.

REINFORCE algorithm

Can we compute those quantities without knowing the transition probabilities?

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

$\frac{\partial}{\partial \theta}$ 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 those quantities without knowing the transition probabilities?

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

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

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

Doesn't depend on transition probabilities!

REINFORCE algorithm

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

Can we compute those quantities 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 τ , we can estimate $J(\theta)$ with

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

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 $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ 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 $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ 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!**

Pong from pixels

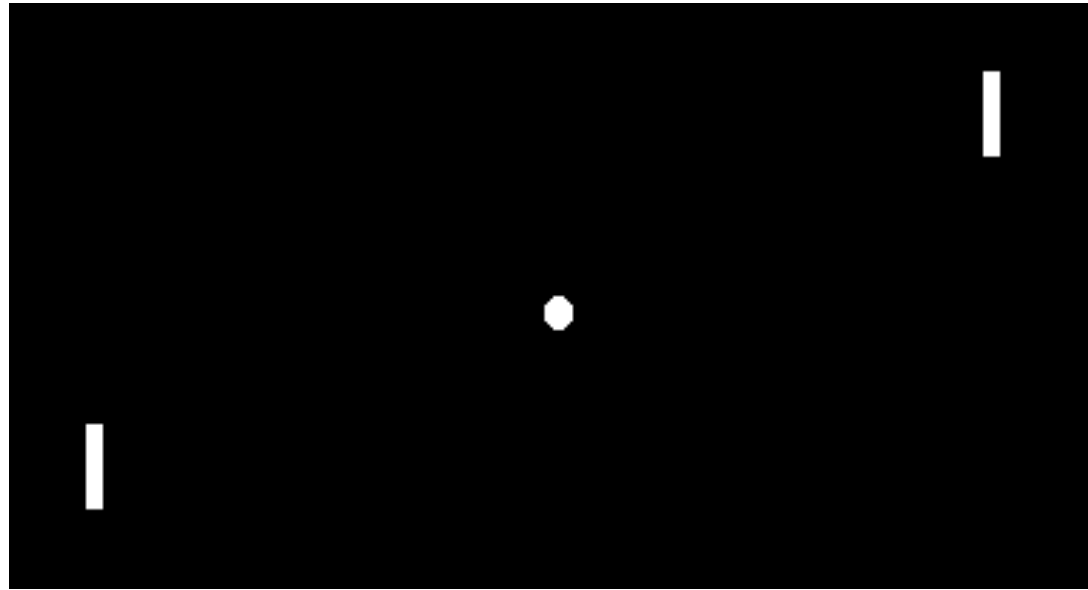
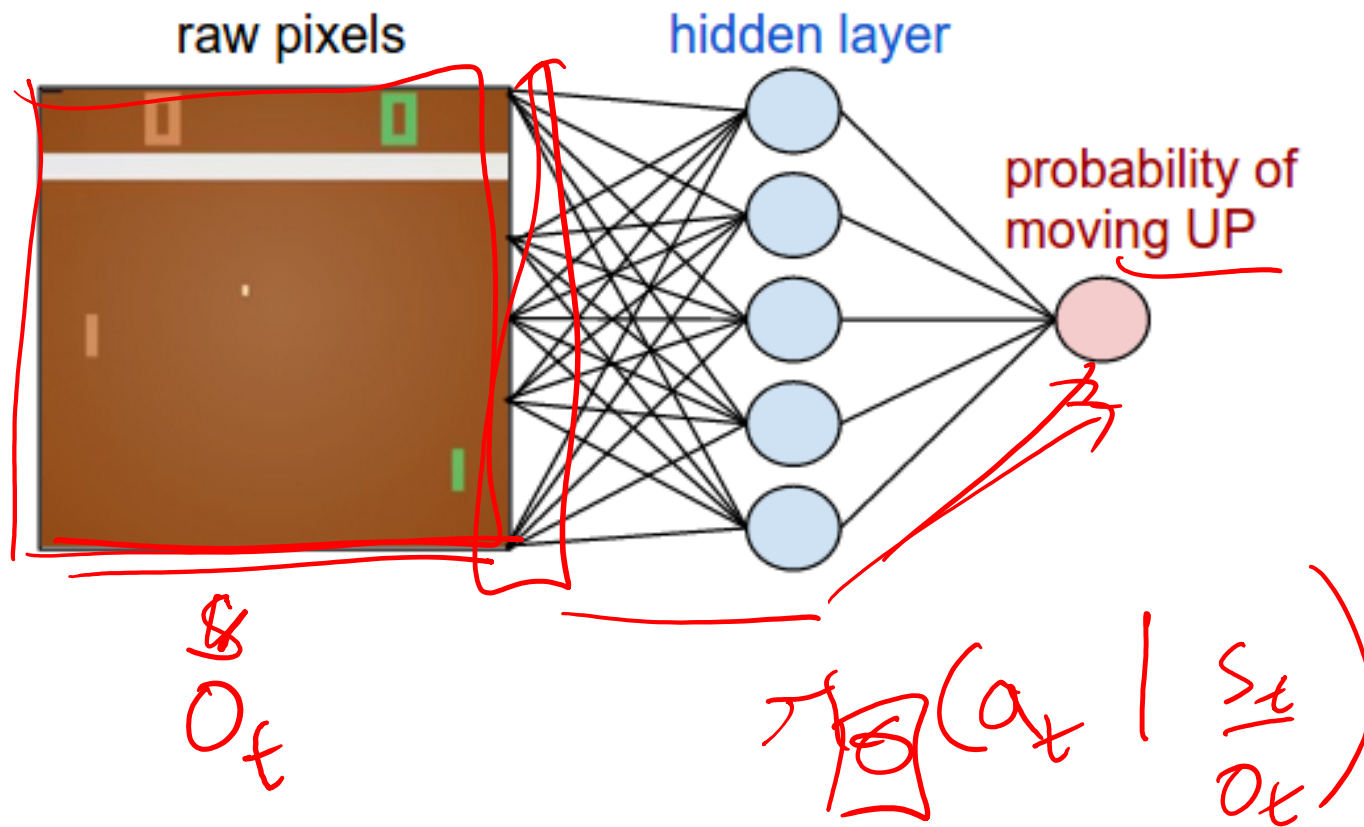
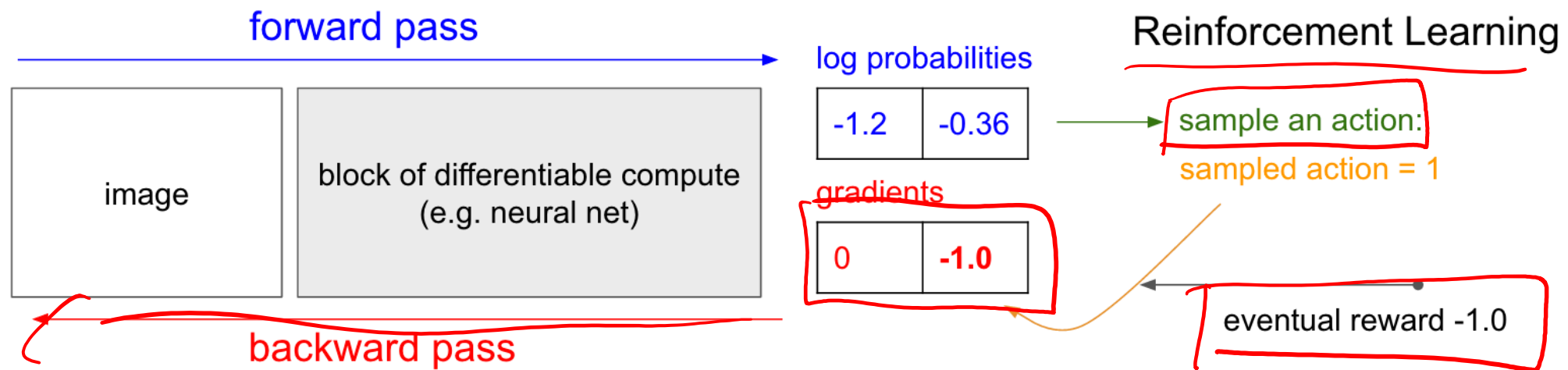
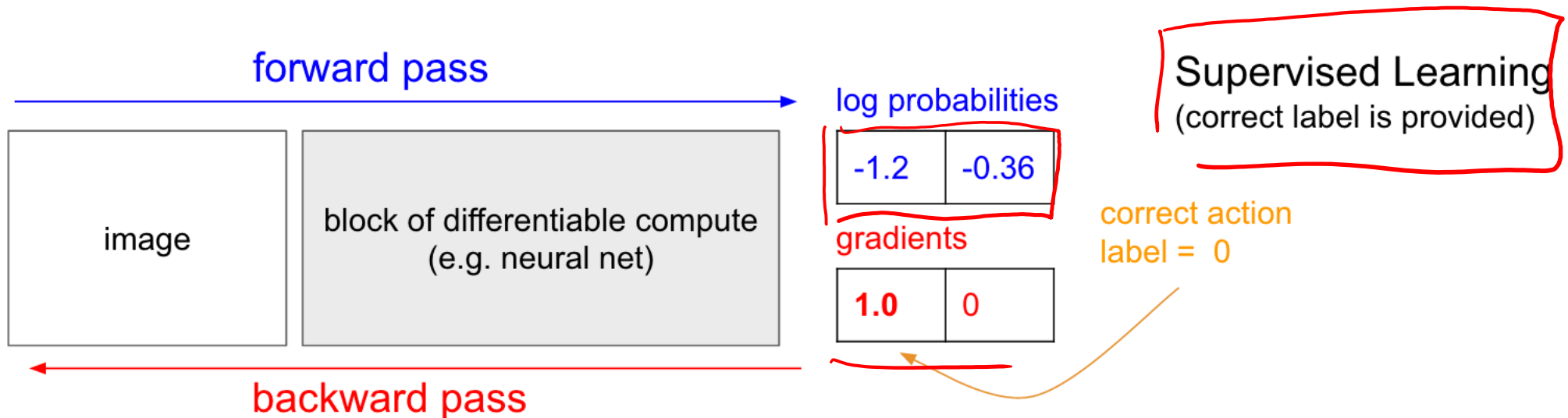


Image Credit: <http://karpathy.github.io/2016/05/31/rl/>

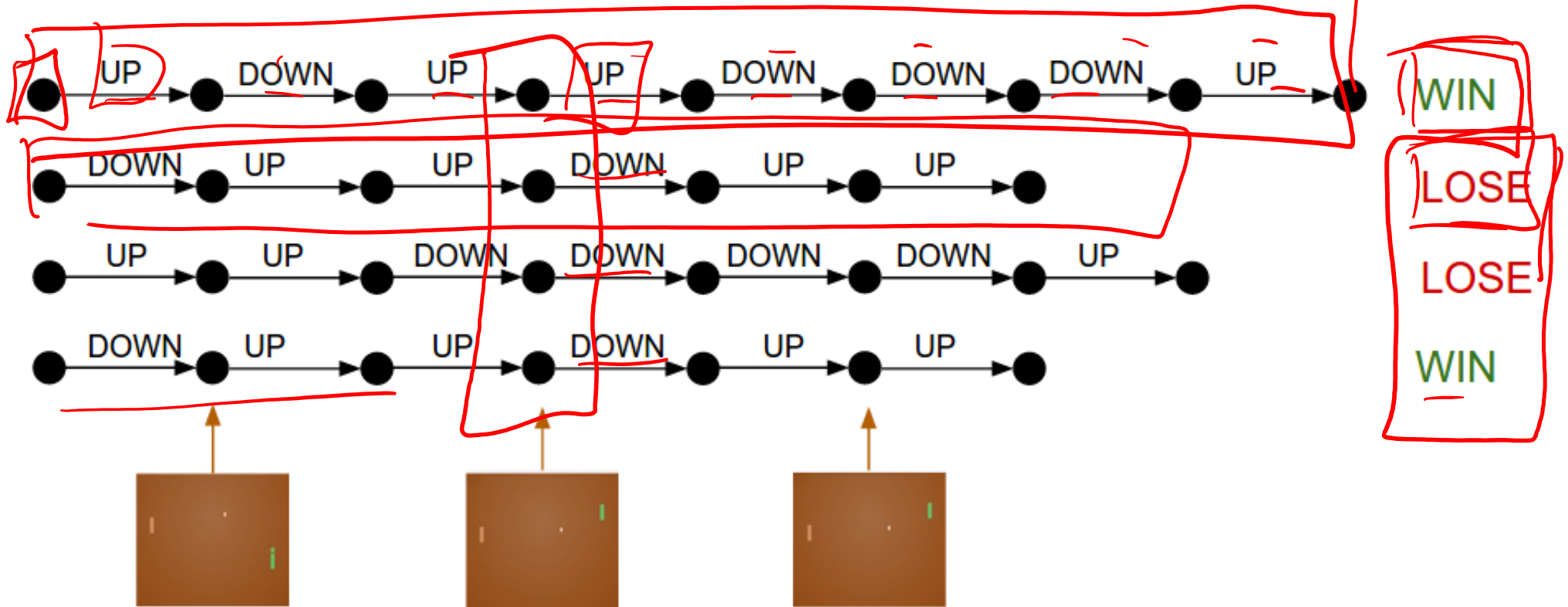
Pong from pixels



Pong from pixels



Intuition



REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

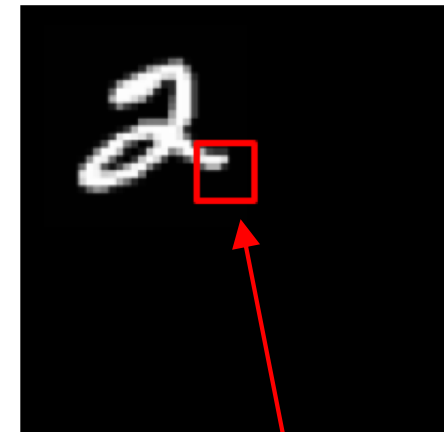
Take a sequence of “glimpses” selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

Reward: 1 at the final timestep if image correctly classified, 0 otherwise



glimpse

[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

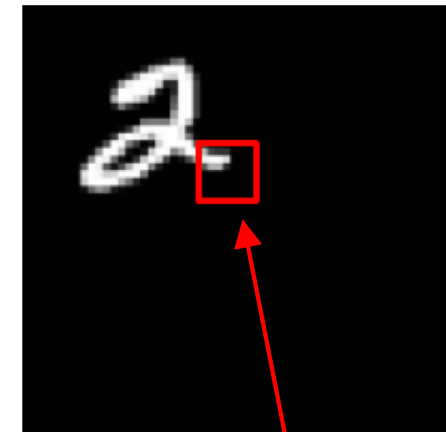
Take a sequence of “glimpses” selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

Reward: 1 at the final timestep if image correctly classified, 0 otherwise

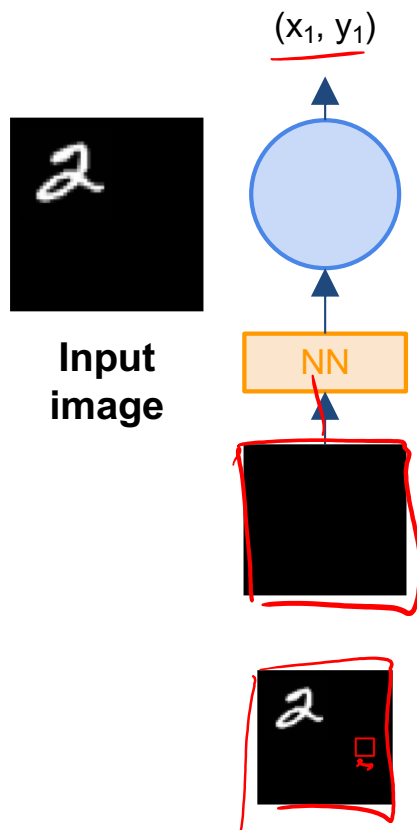


glimpse

Glimpsing is a non-differentiable operation => learn policy for how to take glimpse actions using REINFORCE
Given state of glimpses seen so far, use RNN to model the state and output next action

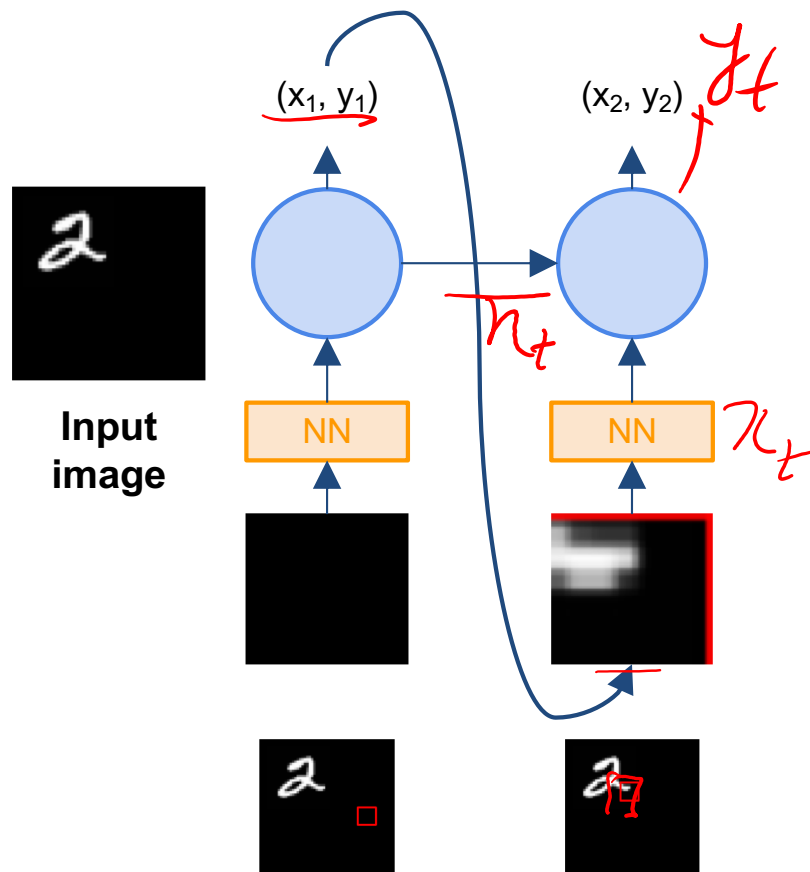
[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



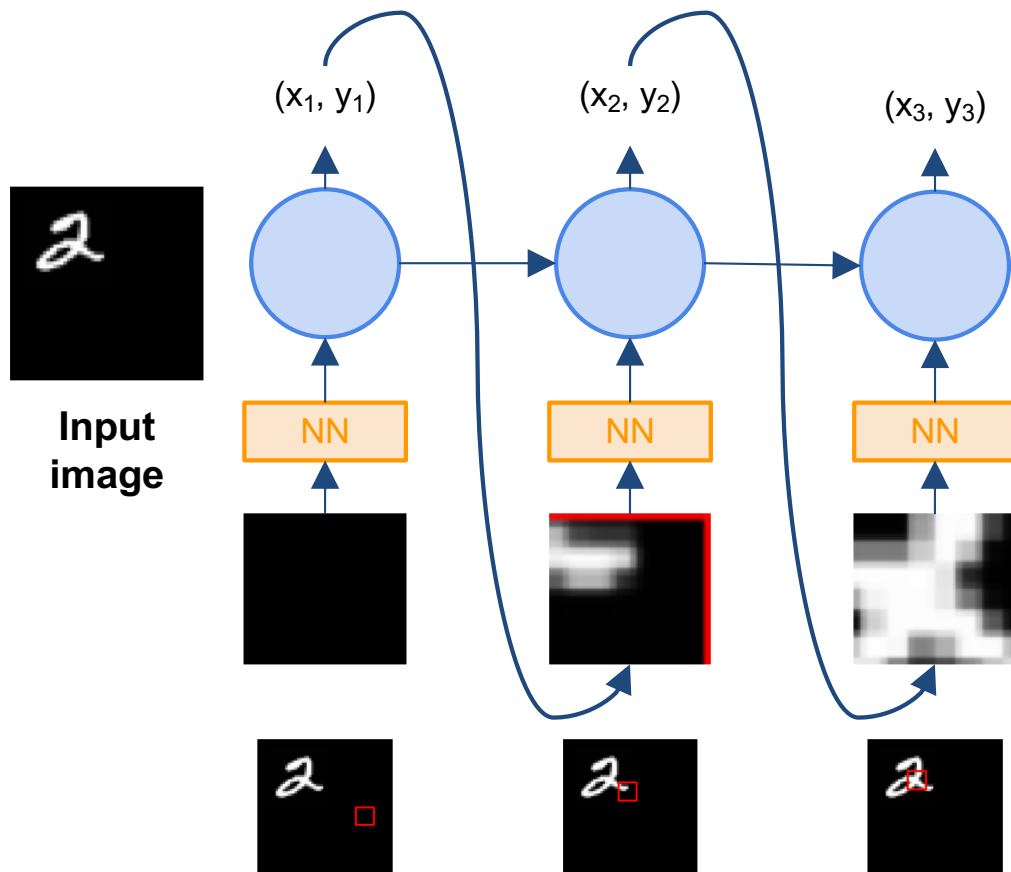
[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



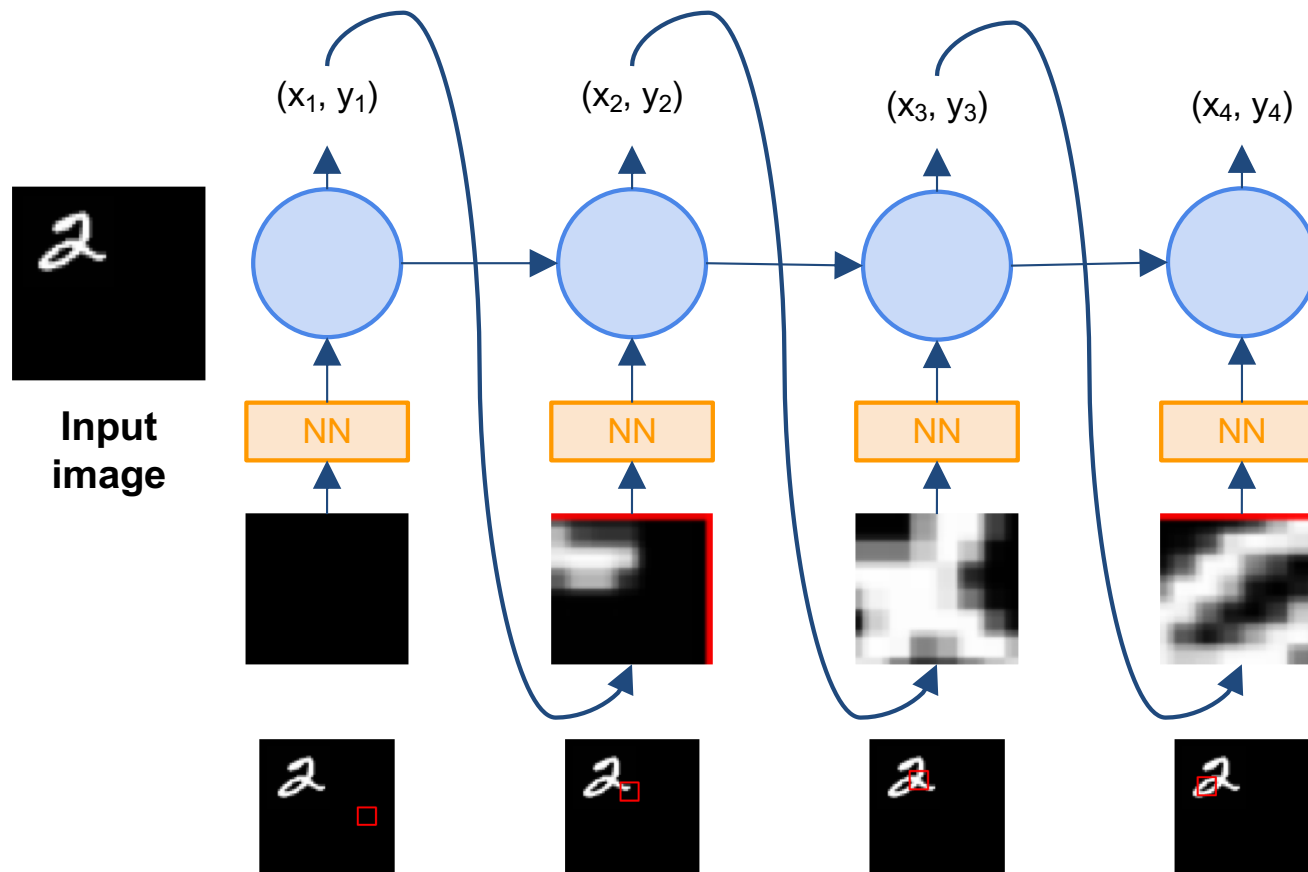
[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



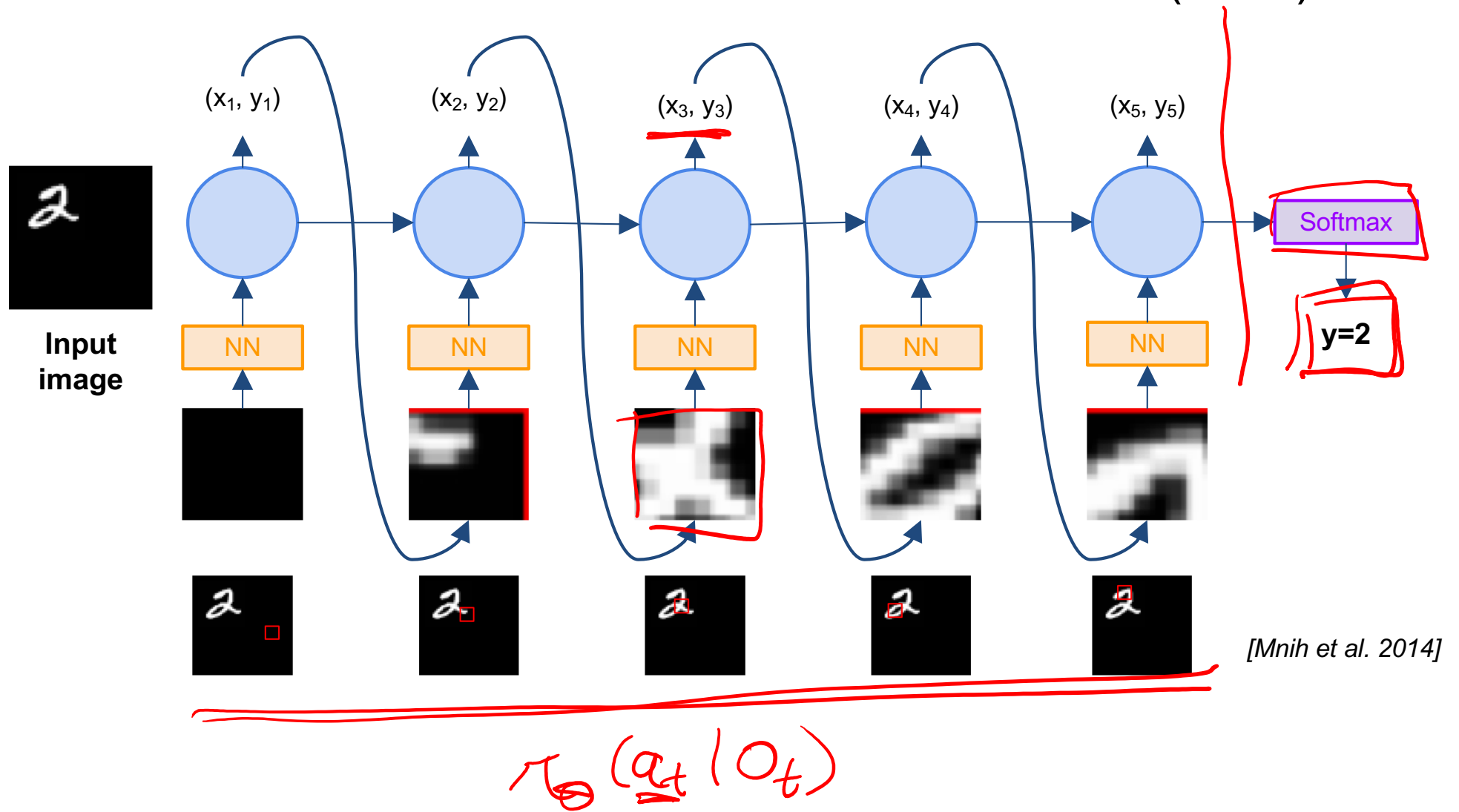
[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



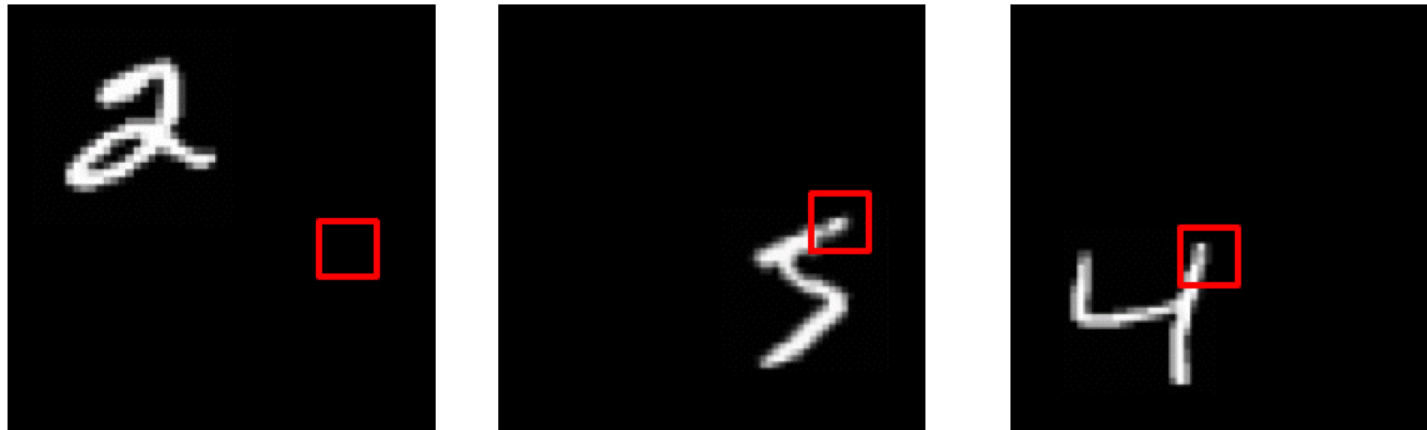
[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)



Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

Figures copyright Daniel Levy, 2017. Reproduced with permission.

[Mnih et al. 2014]

Visual Dialog



A cat drinking water out of a coffee mug.

What color is the mug?

White and red

Are there any pictures on it?

No, something is there can't tell what it is

Is the mug and cat on a table?

Yes, they are

Are there other items on the table?

Yes, magazines, books, toaster and basket, and a plate



Start typing question here ...



Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning

[ICCV '17]



Abhishek Das*
(Georgia Tech)



Satwik Kottur*
(CMU)



José Moura
(CMU)



Stefan Lee
(Virginia Tech)



Dhruv Batra
(Georgia Tech)

Image Guessing Game

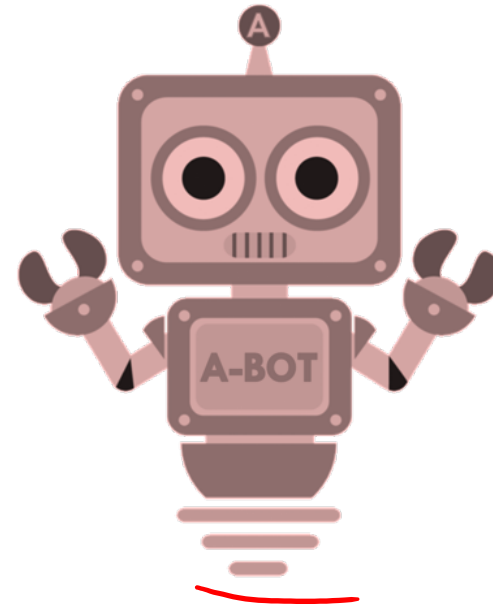
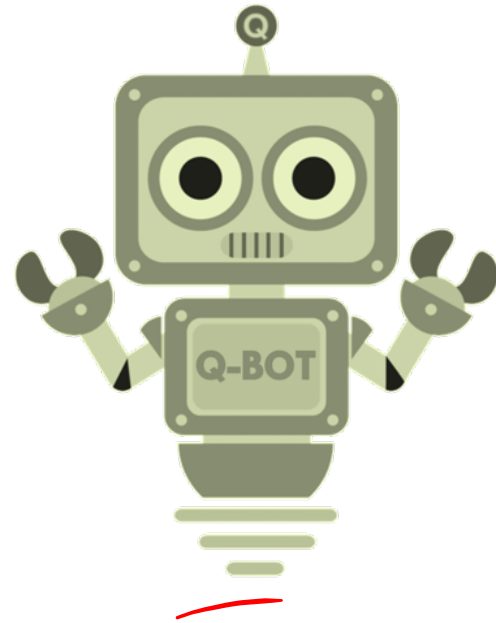
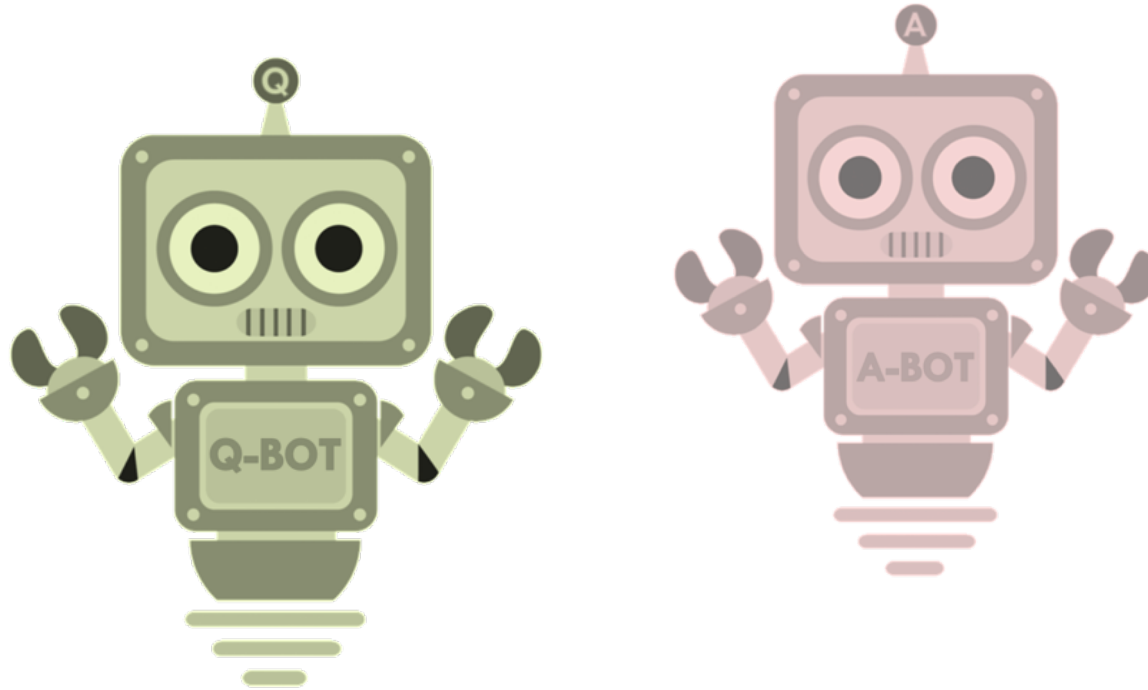
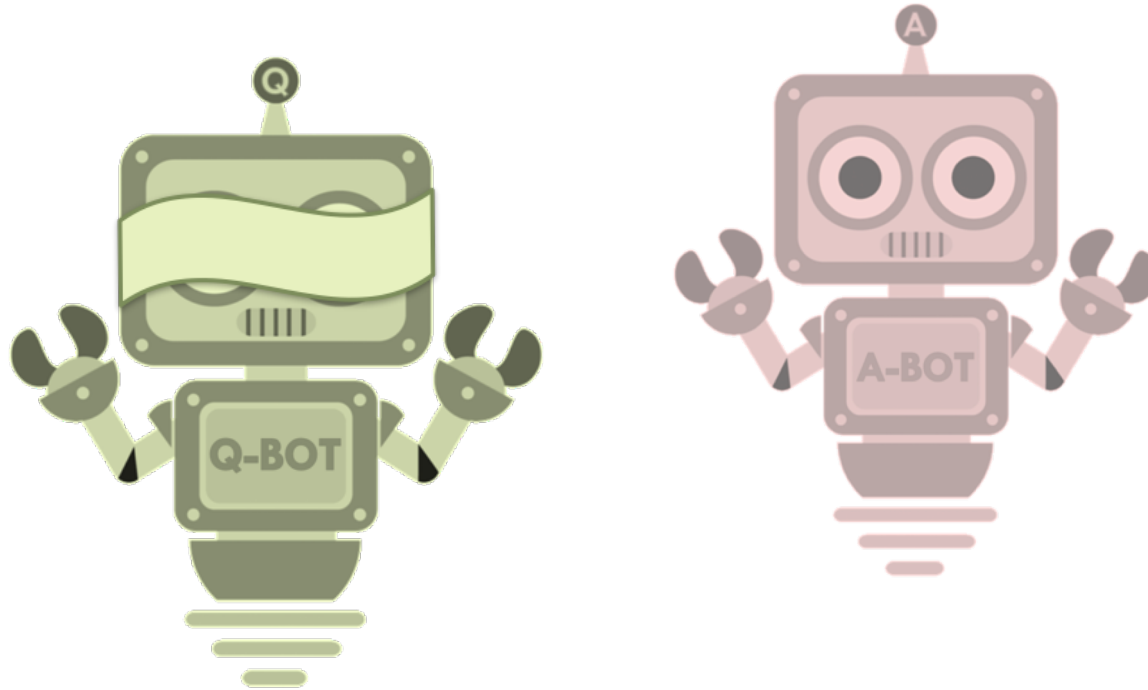


Image Guessing Game



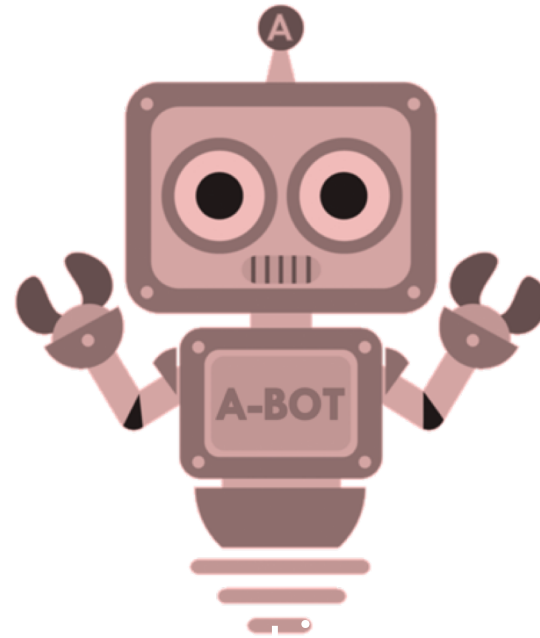
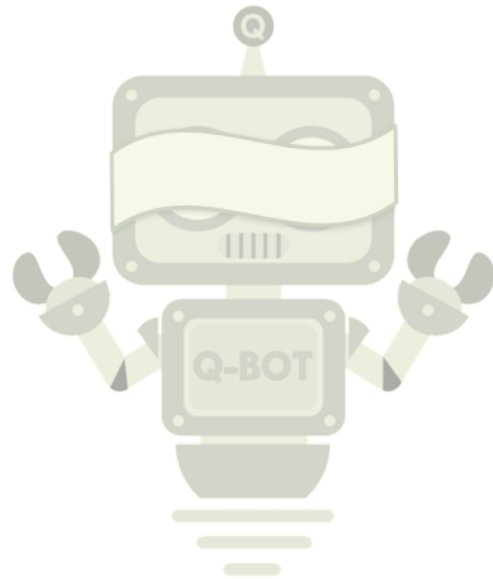
Q-Bot asks questions

Image Guessing Game



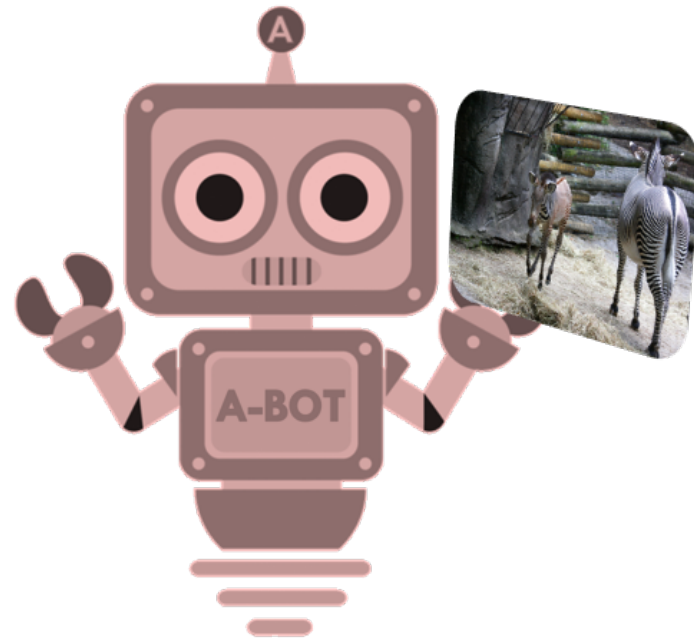
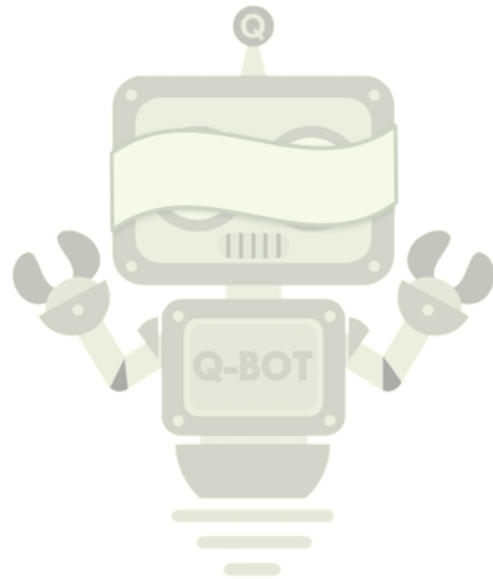
Q-Bot is blindfolded

Image Guessing Game



A-Bot answers questions

Image Guessing Game



A-Bot sees an image

Image Guessing Game

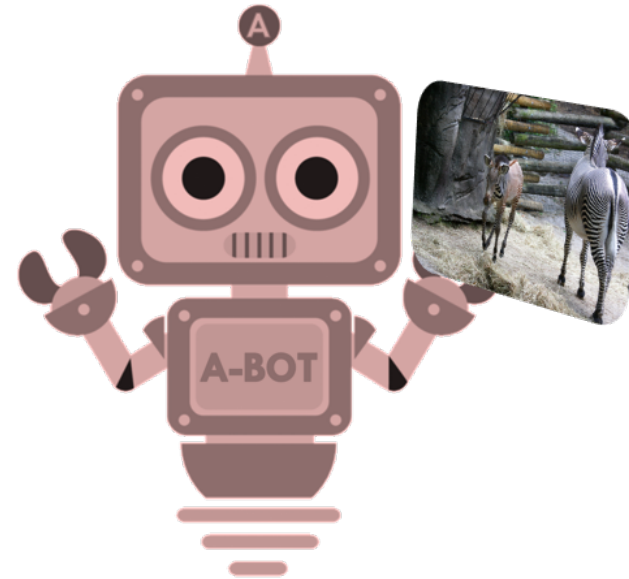
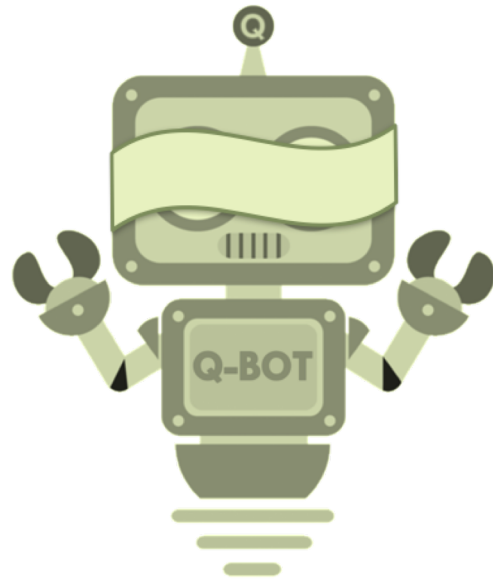


Image Guessing Game

Q Two zebra are walking around their pen at the zoo. A

Q1: Any people in the shot?

A1: No, there aren't any.

Q2: Any other animal?

A2: No, just zebras.

Q3: Are they facing each other?

A3: They aren't.

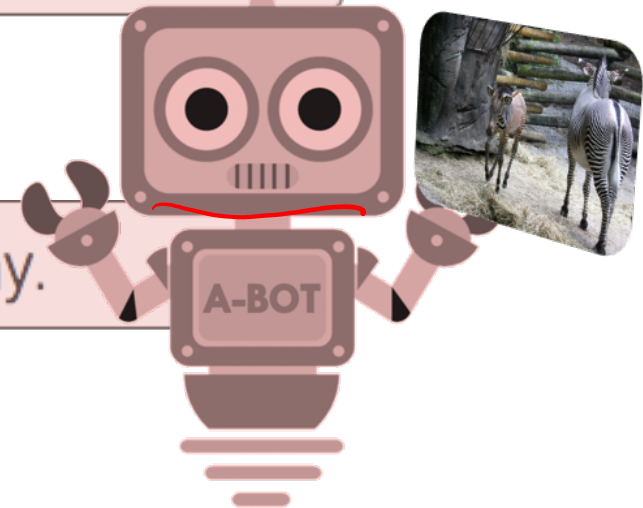
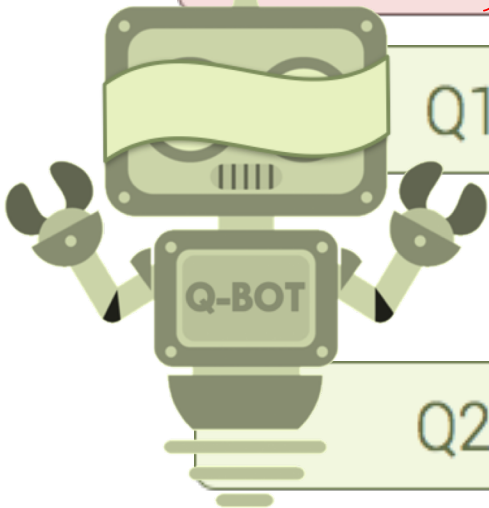
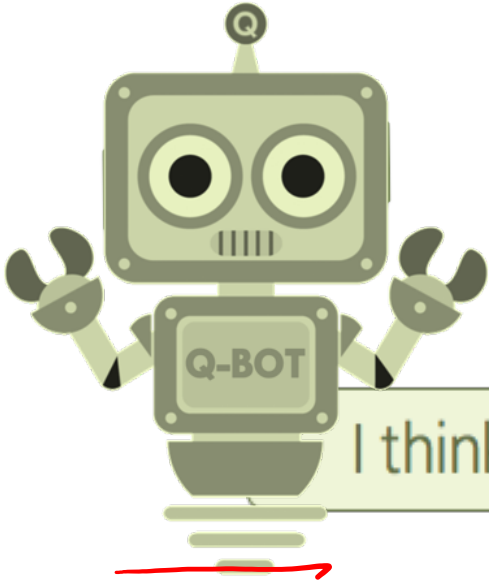
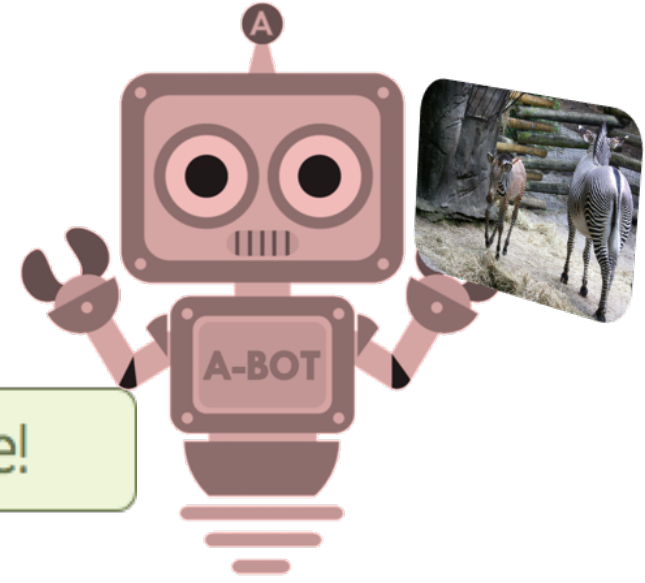


Image Guessing Game

A3: They aren't.



I think we were talking about this image!

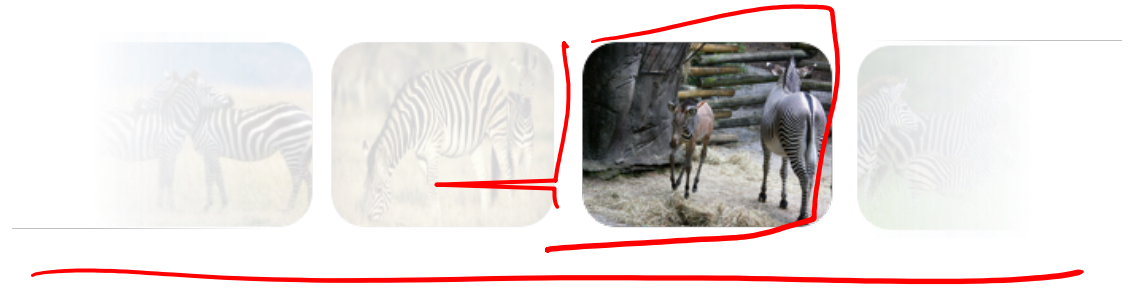


RL for Cooperative Dialog Agents

- Agents: (Q-bot, A-bot)



- Environment: Image



- Action:

- Q-bot: question (symbol sequence)
- A-bot: answer (symbol sequence)
- Q-bot: image regression

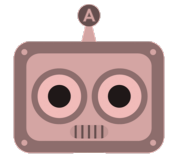
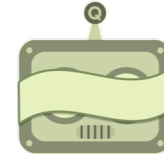
q_t Any people in the shot?
 a_t No, there aren't any.

$$\hat{y}_t \in \mathbb{R}^{4096}$$

- State

- Q-bot: $s_t^Q = [c, q_1, a_1, \dots, q_{t-1}, a_{t-1}]$
- A-bot: $s_t^A = [I, c, q_1, a_1, \dots, q_{t-1}, a_{t-1}, q_t]$

RL for Cooperative Dialog Agents



- Action:

- Q-bot: question (symbol sequence) q_t *Any people in the shot?*
- A-bot: answer (symbol sequence) a_t *No, there aren't any.*
- Q-bot: image regression $\hat{y}_t \in \mathbb{R}^{4096}$

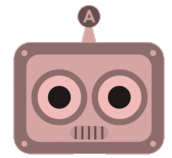
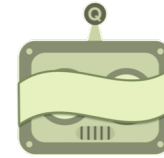
- State

- Q-bot: $s_t^Q = [c, q_1, a_1, \dots, q_{t-1}, a_{t-1}]$
- A-bot: $s_t^A = [I, c, q_1, a_1, \dots, q_{t-1}, a_{t-1}, q_t]$

RL for Cooperative Dialog Agents

- Action:

- Q-bot: question (symbol sequence)
- A-bot: answer (symbol sequence)
- Q-bot: image regression



q_t Any people in the shot?

a_t No, there aren't any.

$\hat{y}_t \in \mathbb{R}^{4096}$

- State

- Q-bot: $s_t^Q = [c, q_1, a_1, \dots, q_{t-1}, a_{t-1}]$

- A-bot: $s_t^A = [I, c, q_1, a_1, \dots, q_{t-1}, a_{t-1}, q_t]$

- Policy

Q-bot
A-bot

$\pi_Q(q_t | s_{t-1}^Q)$
 $\pi_A(a_t | s_{t-1}^A)$

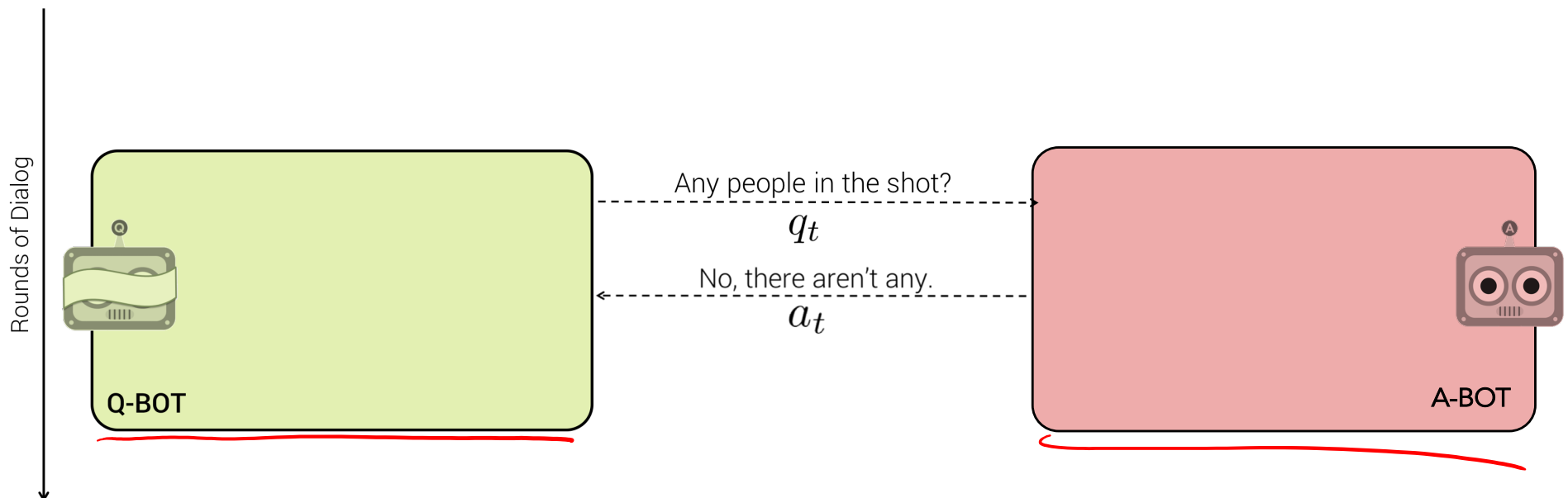
- Reward

$$\underbrace{r_t}_{\text{reward}} \left(\underbrace{s_t^Q}_{\text{state}} \underbrace{(q_t, a_t, y_t)}_{\text{action}} \right) = \underbrace{\ell(\hat{y}_{t-1}, y^{gt})}_{\text{distance at } t-1} - \underbrace{\ell(\hat{y}_t, y^{gt})}_{\text{distance at } t}$$

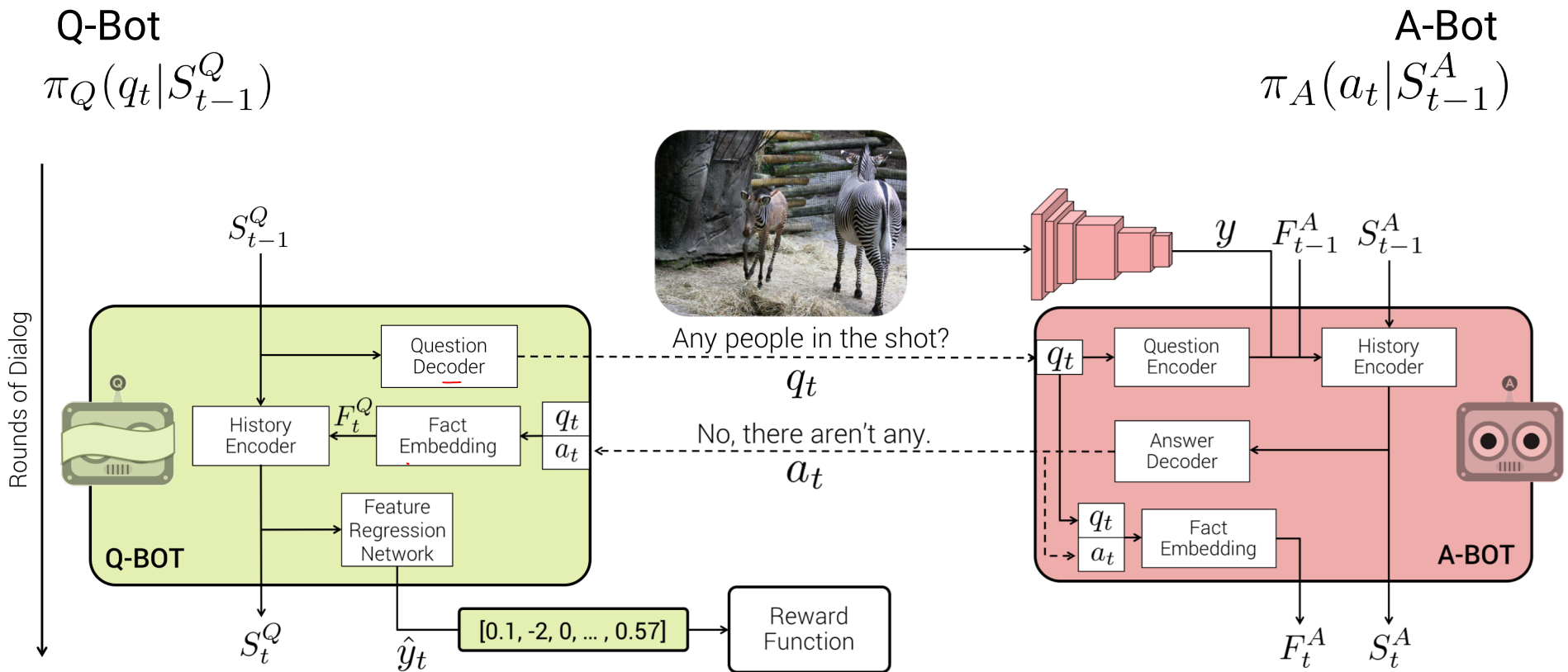
Policy Networks

Q-Bot
 $\pi_Q(q_t | S_{t-1}^Q)$

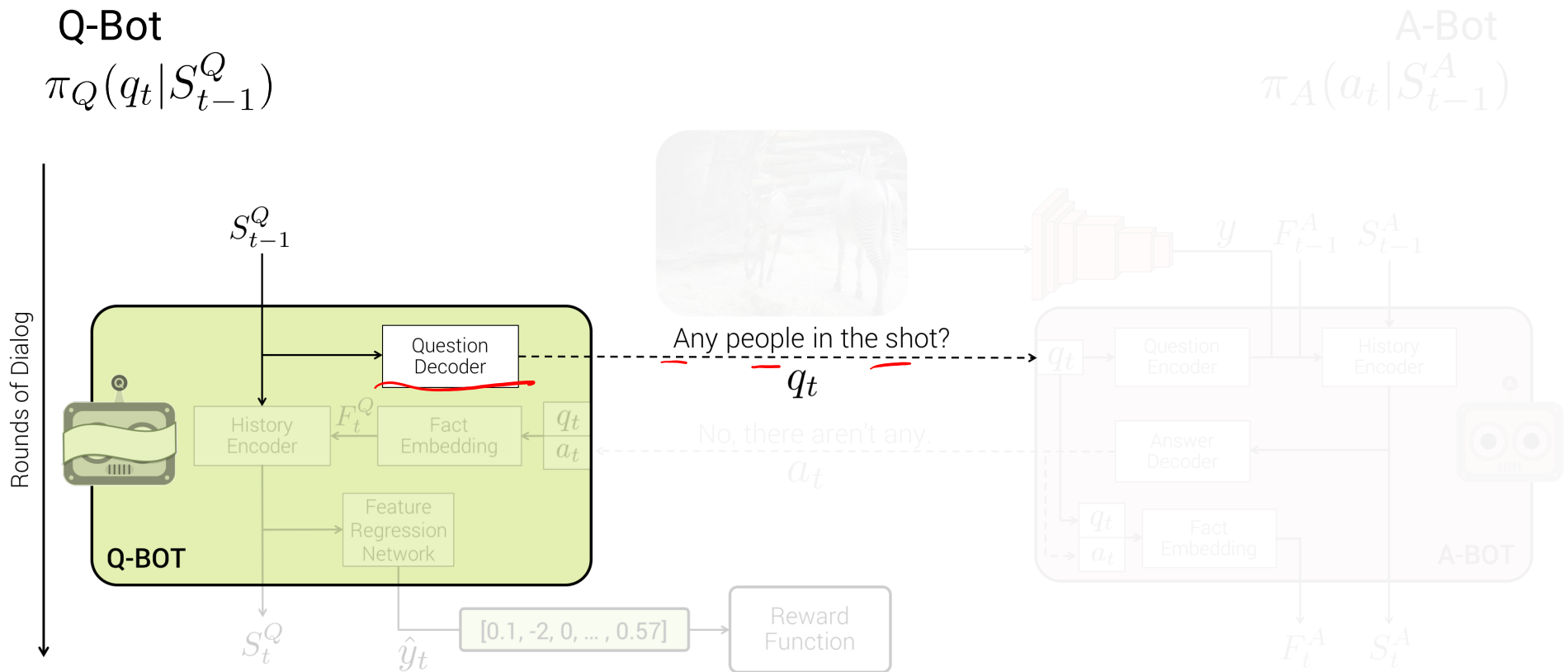
A-Bot
 $\pi_A(a_t | S_{t-1}^A)$



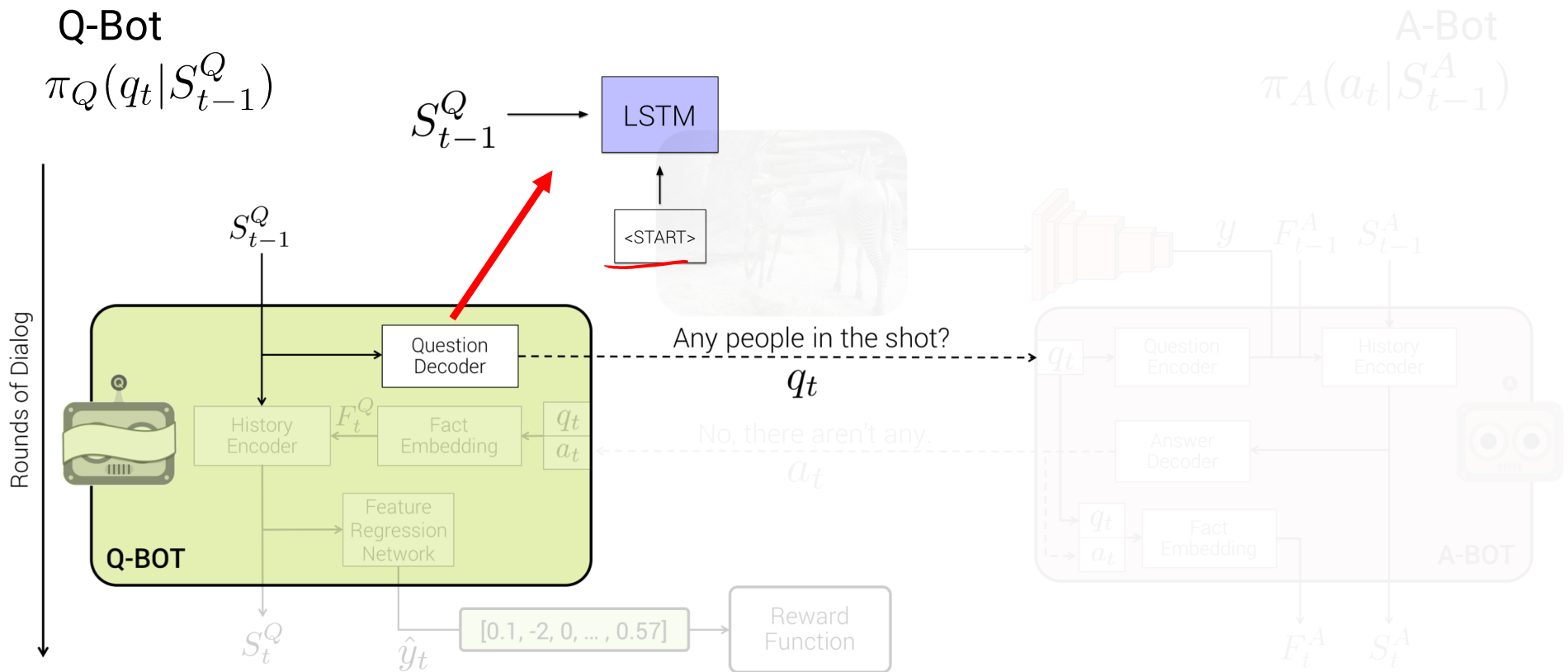
Policy Networks



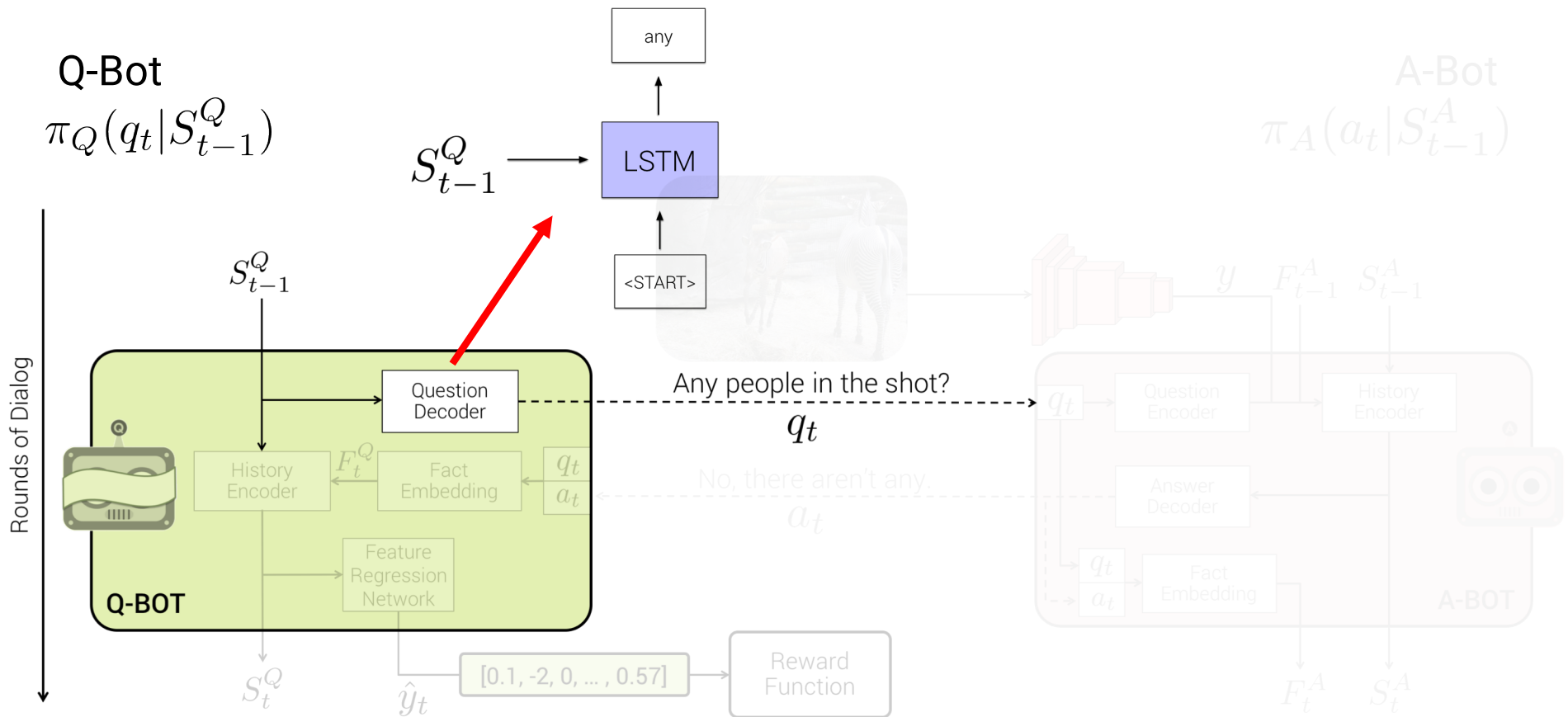
Policy Networks



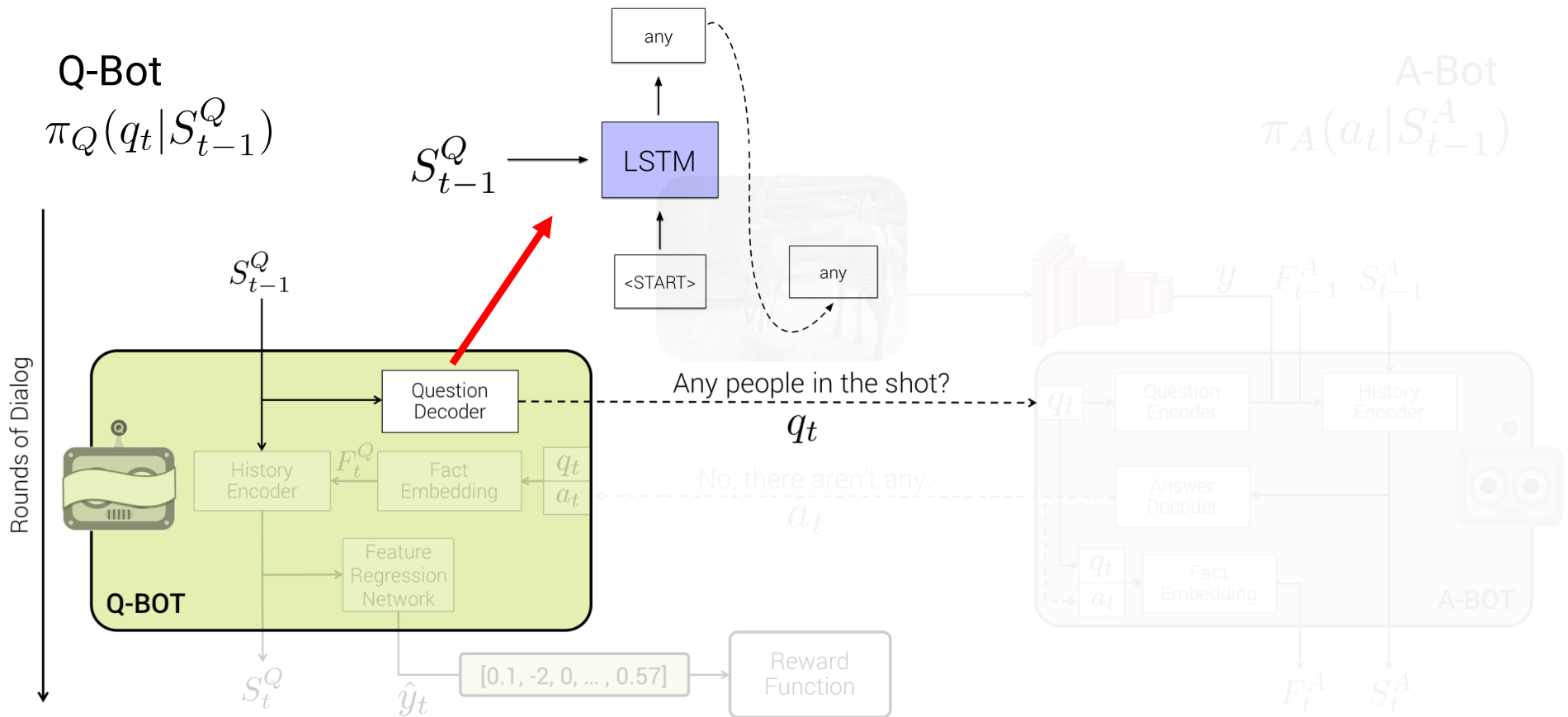
Policy Networks



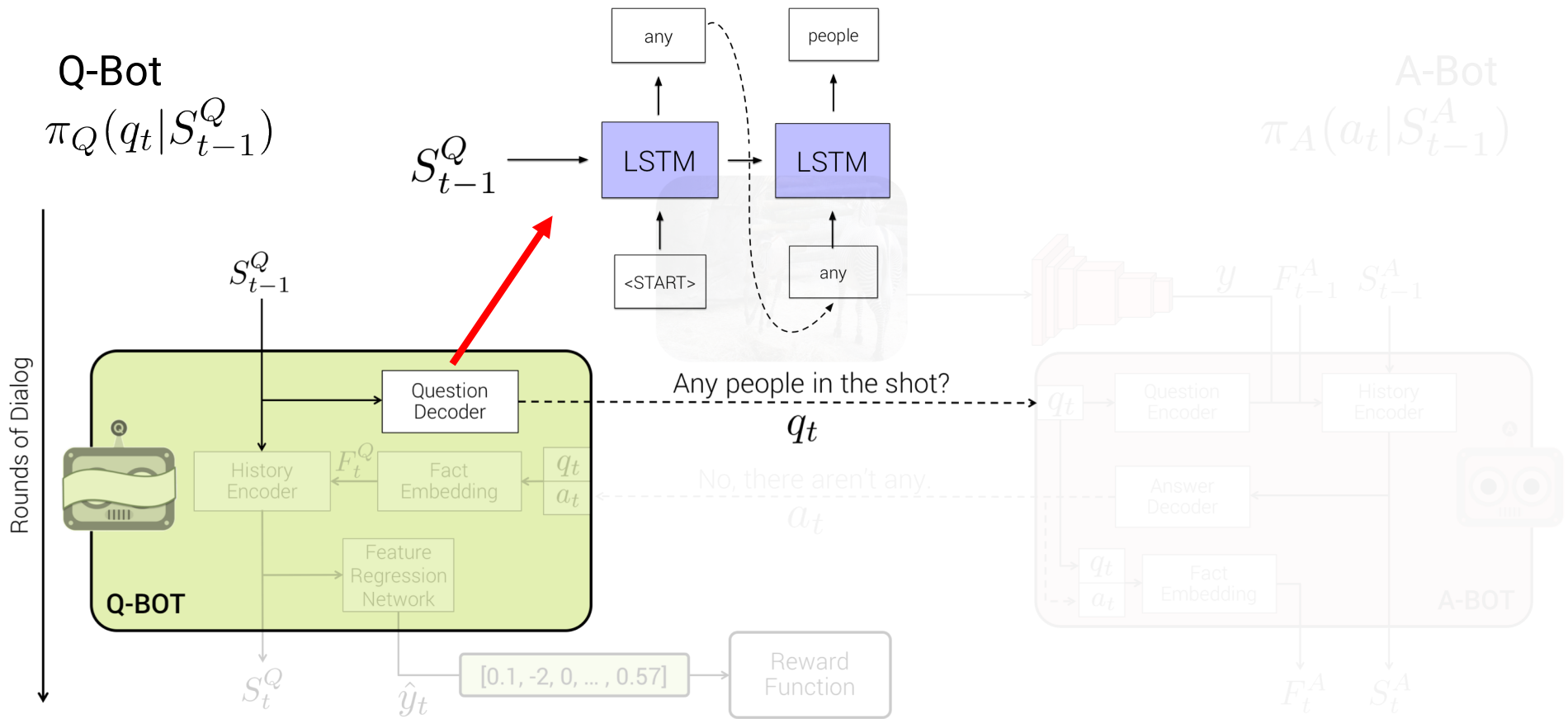
Policy Networks



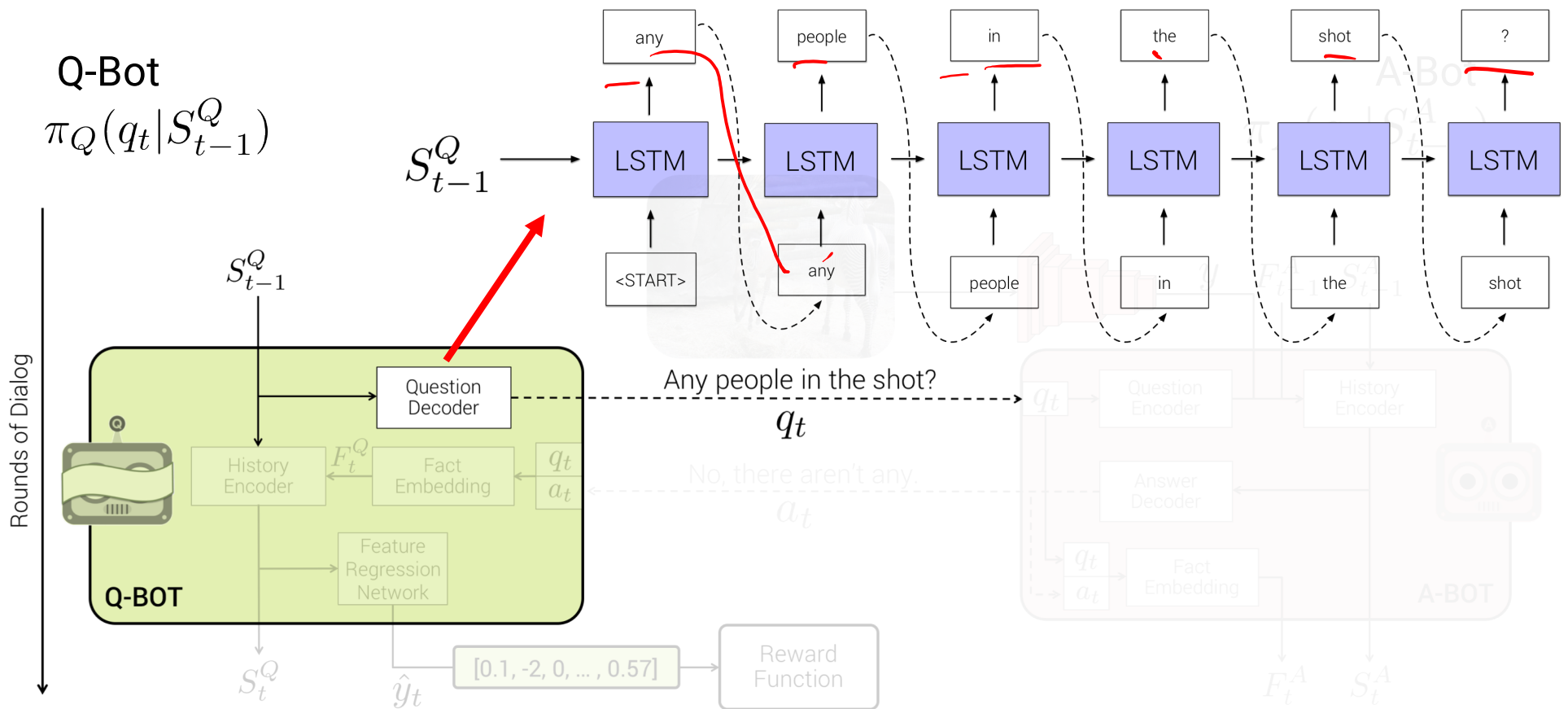
Policy Networks



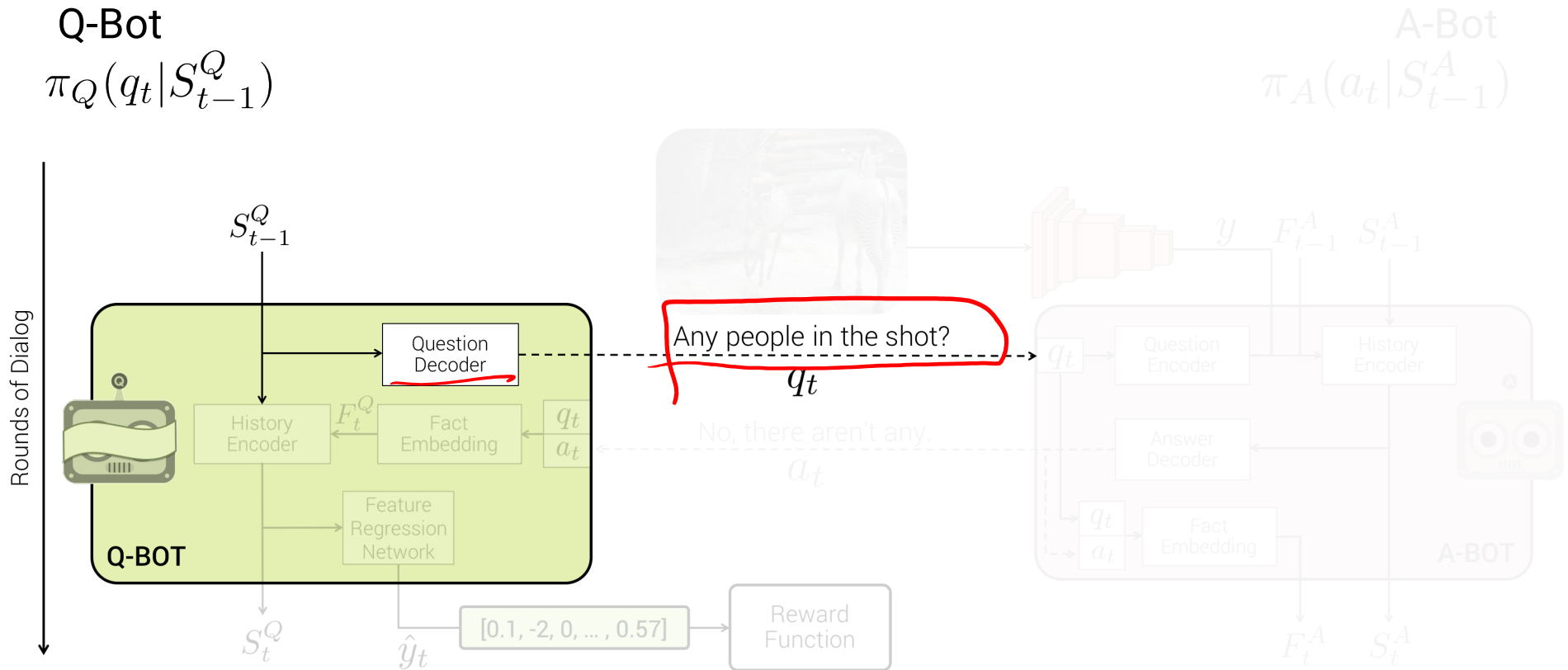
Policy Networks



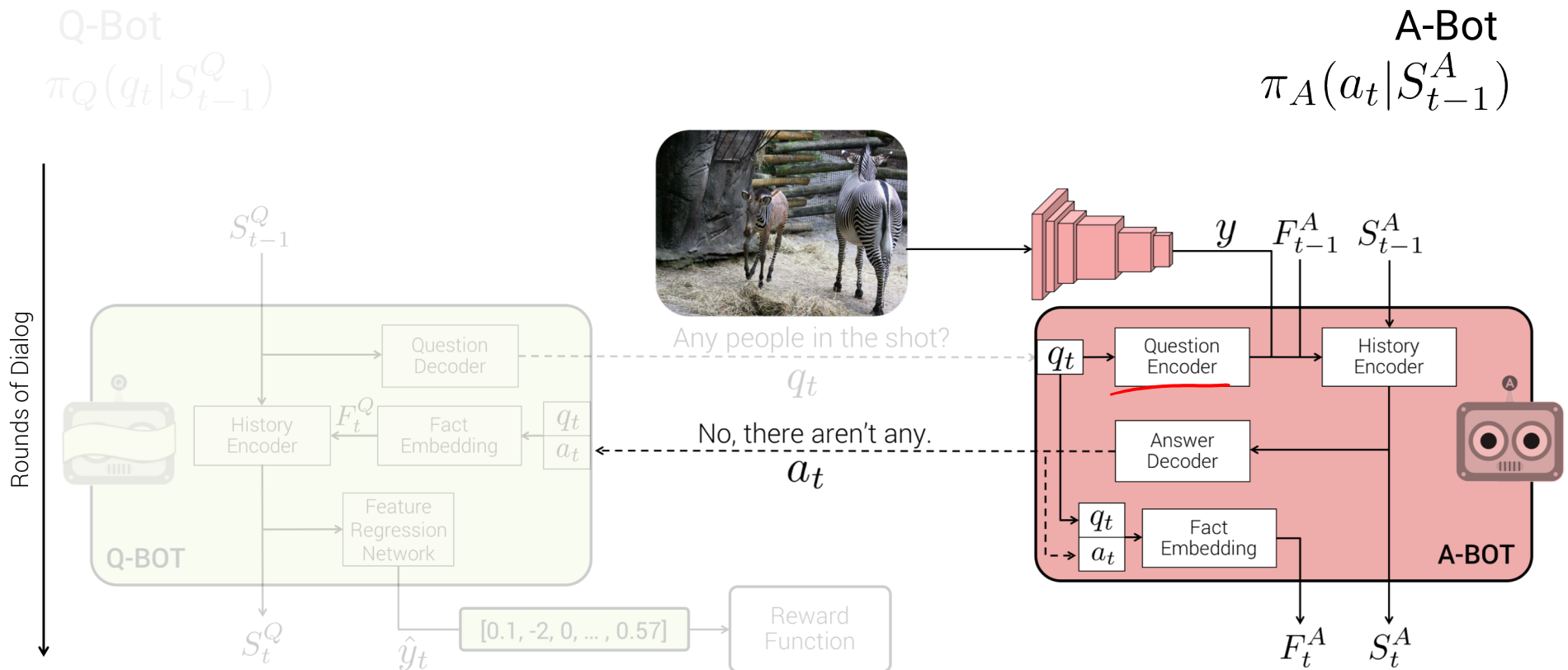
Policy Networks



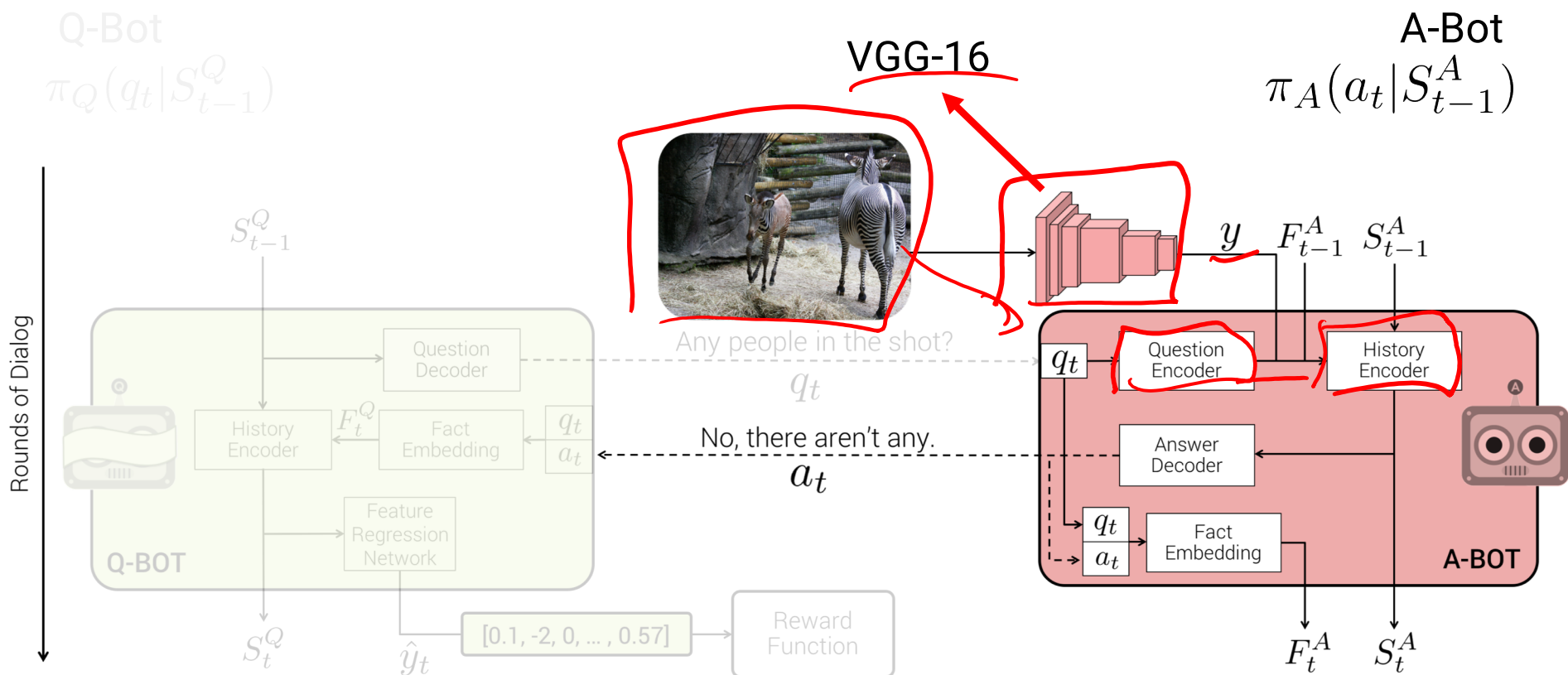
Policy Networks



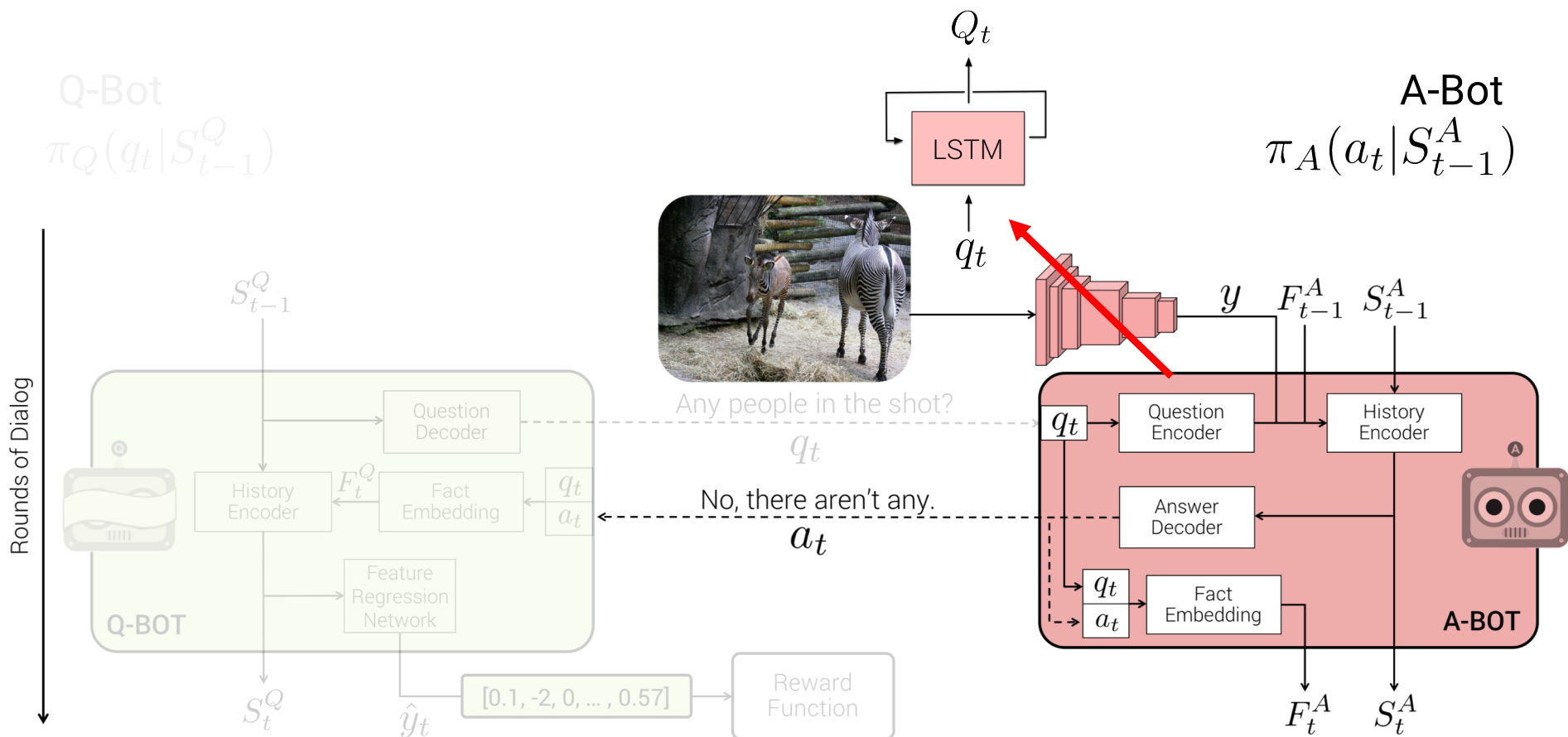
Policy Networks



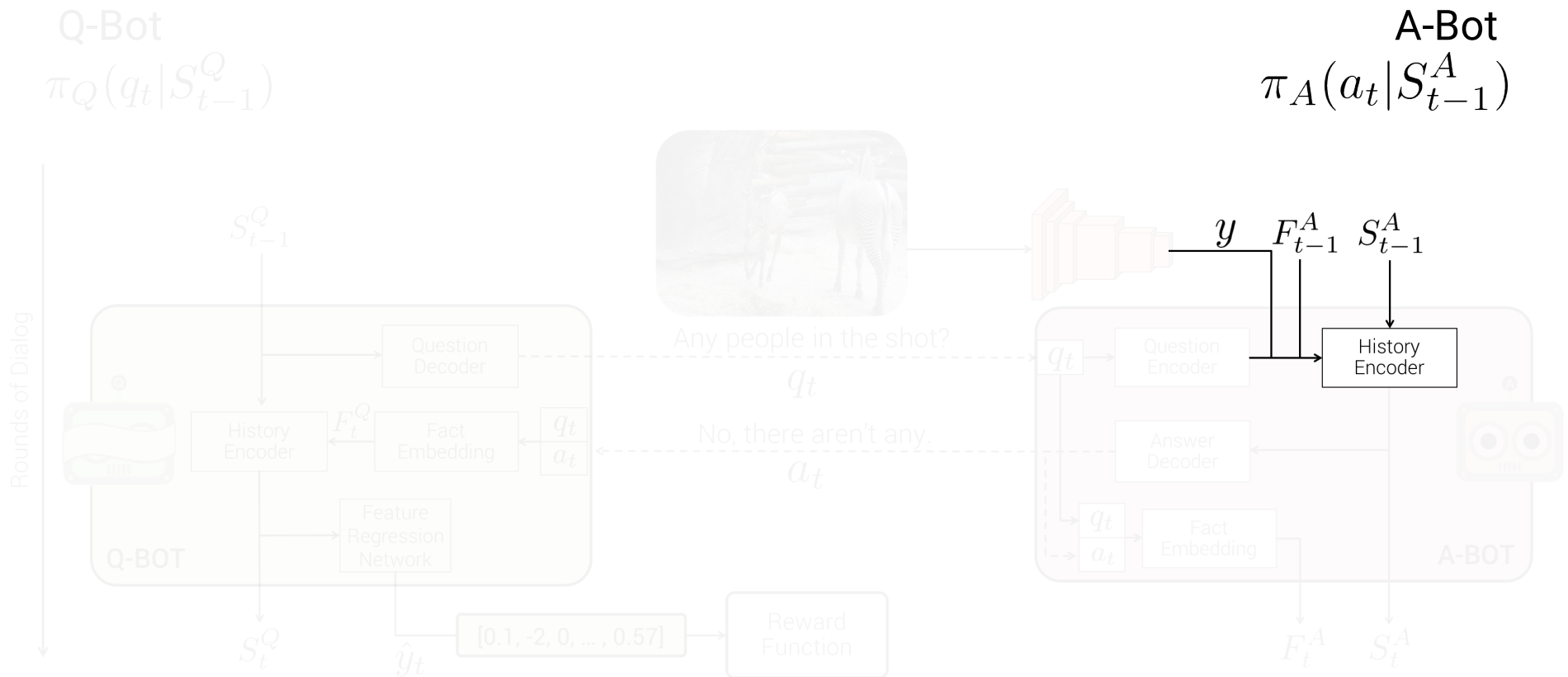
Policy Networks



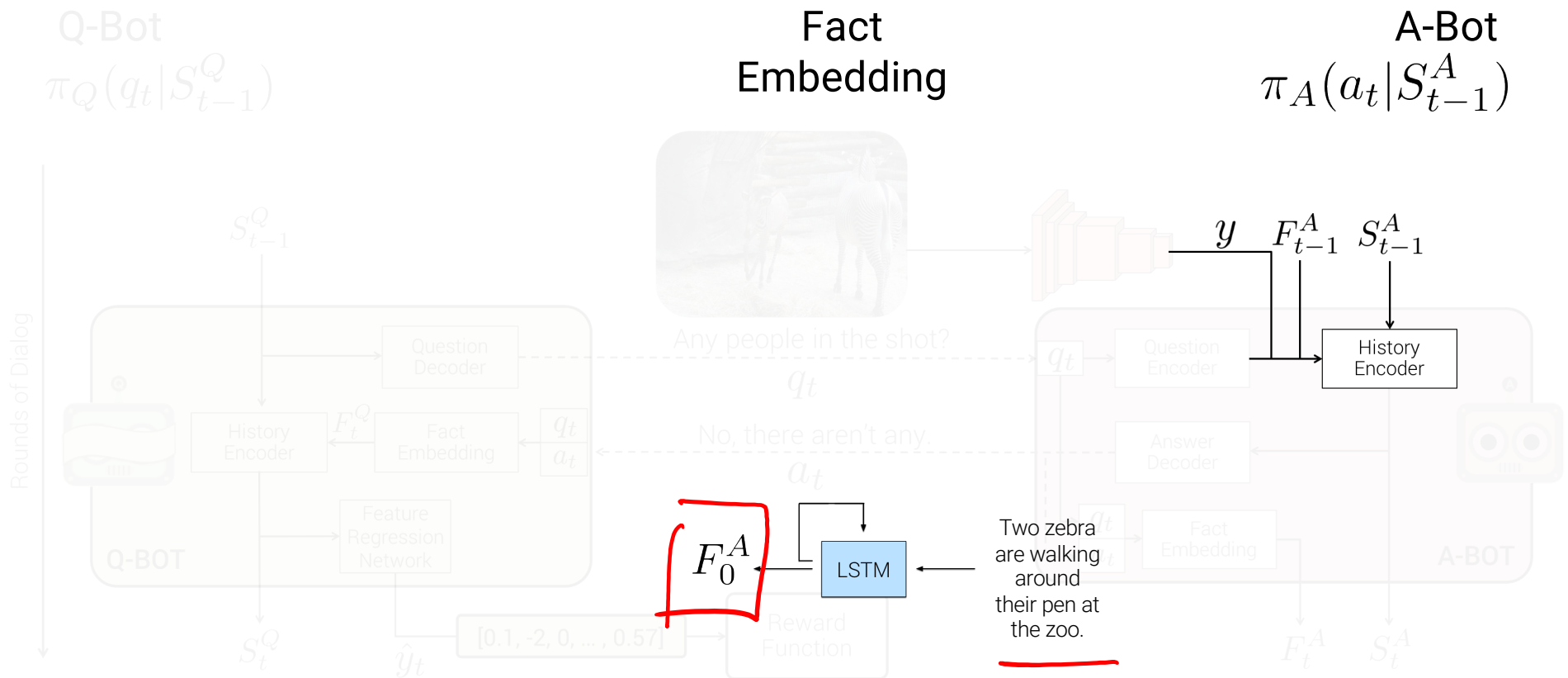
Policy Networks



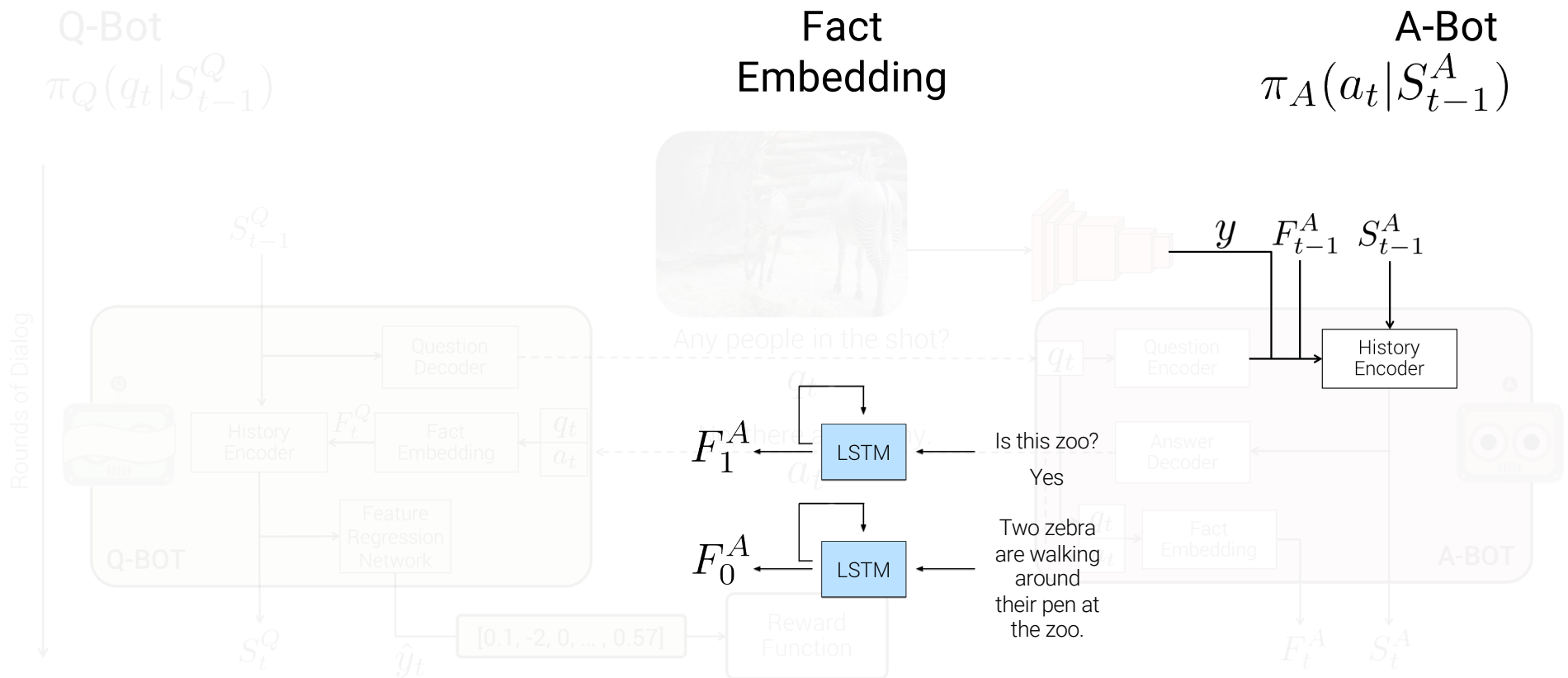
Policy Networks



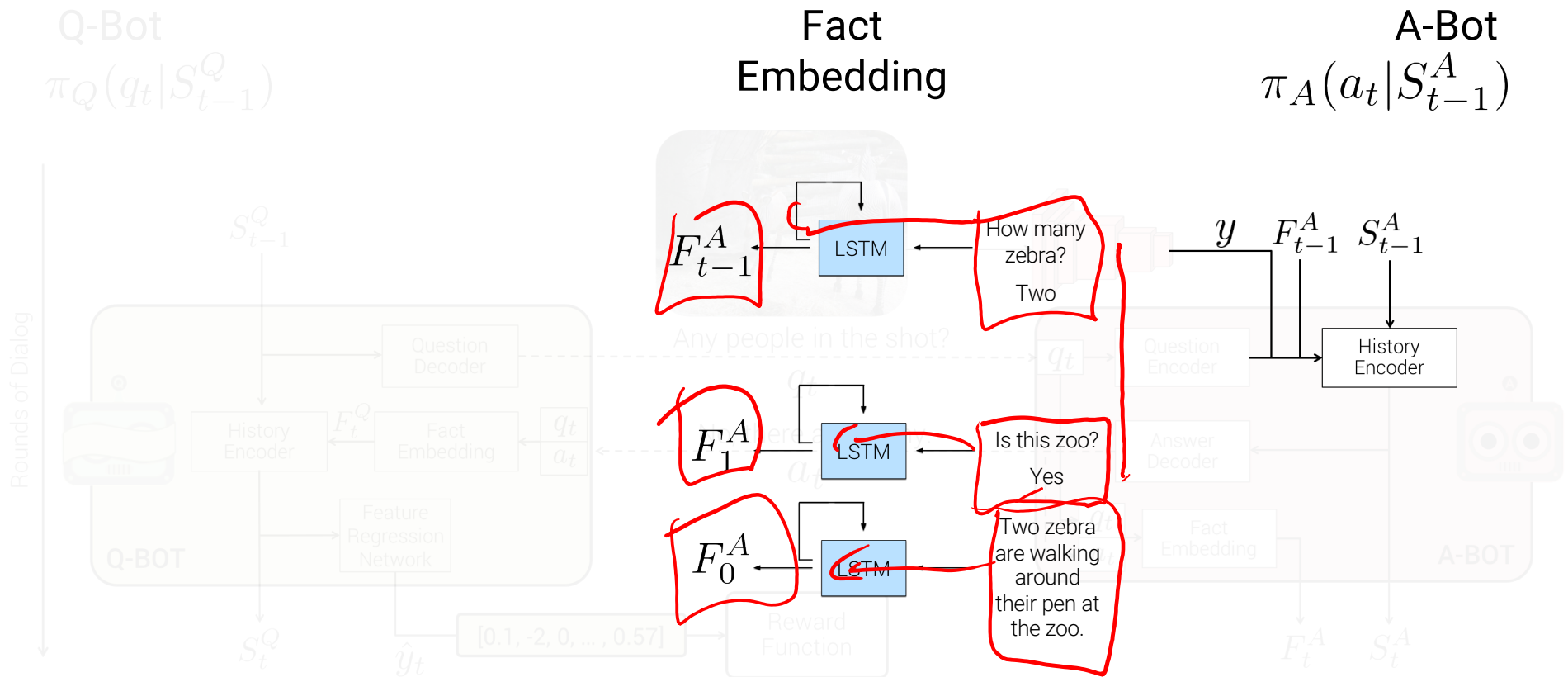
Policy Networks



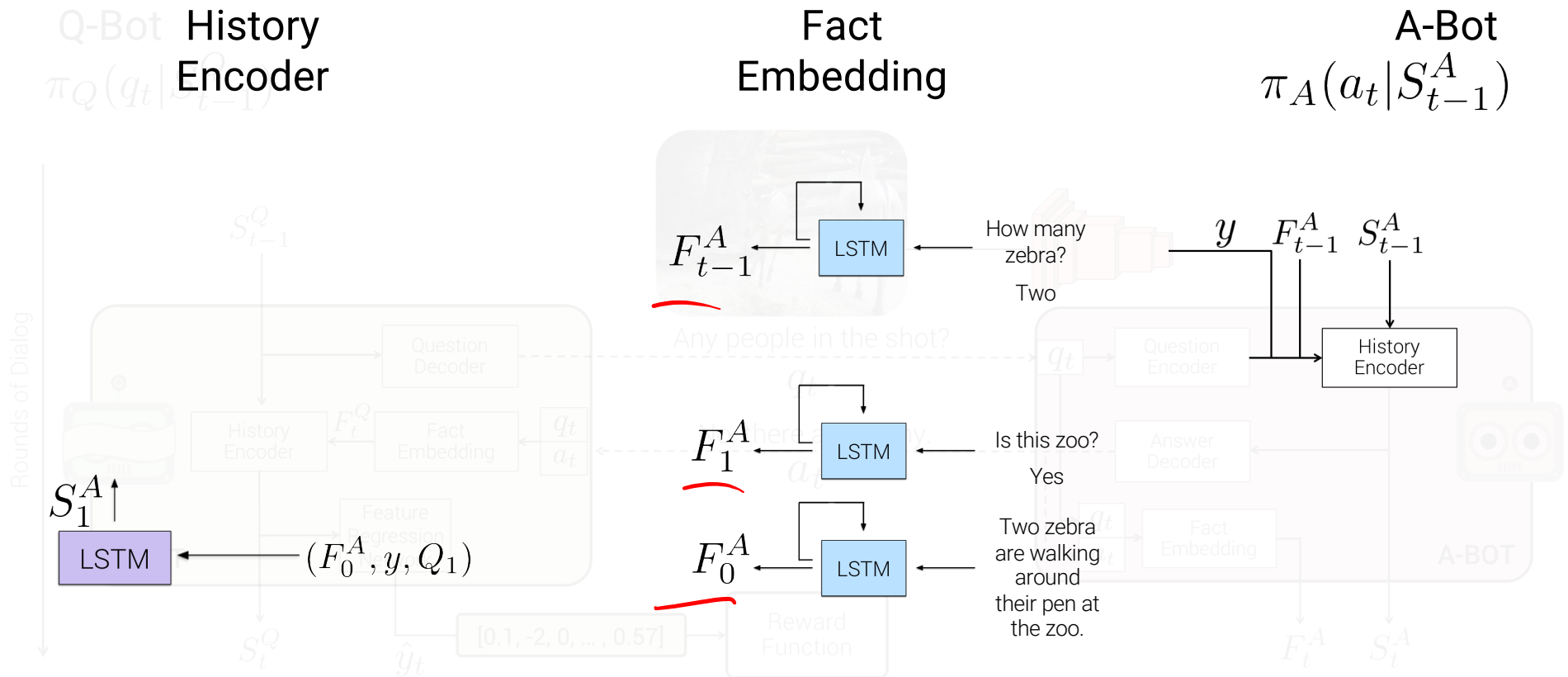
Policy Networks



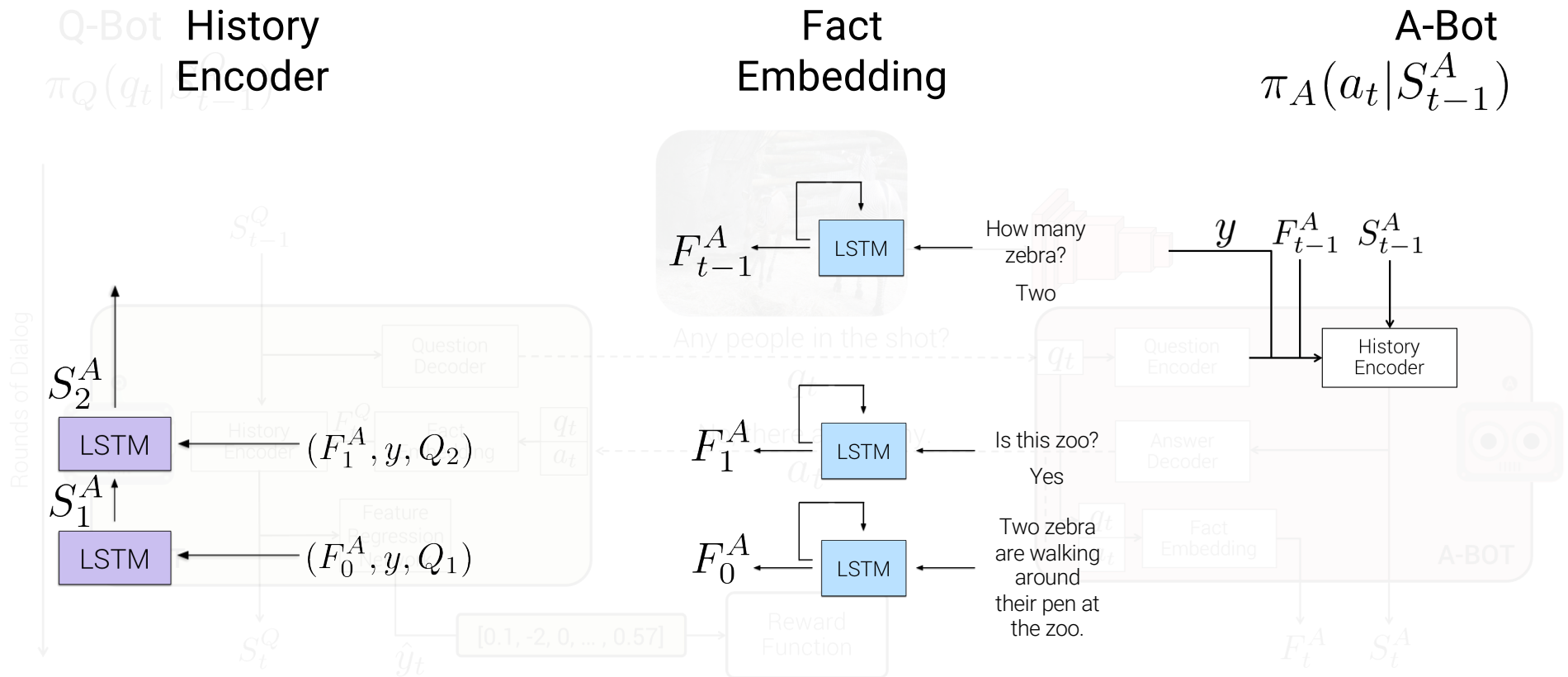
Policy Networks



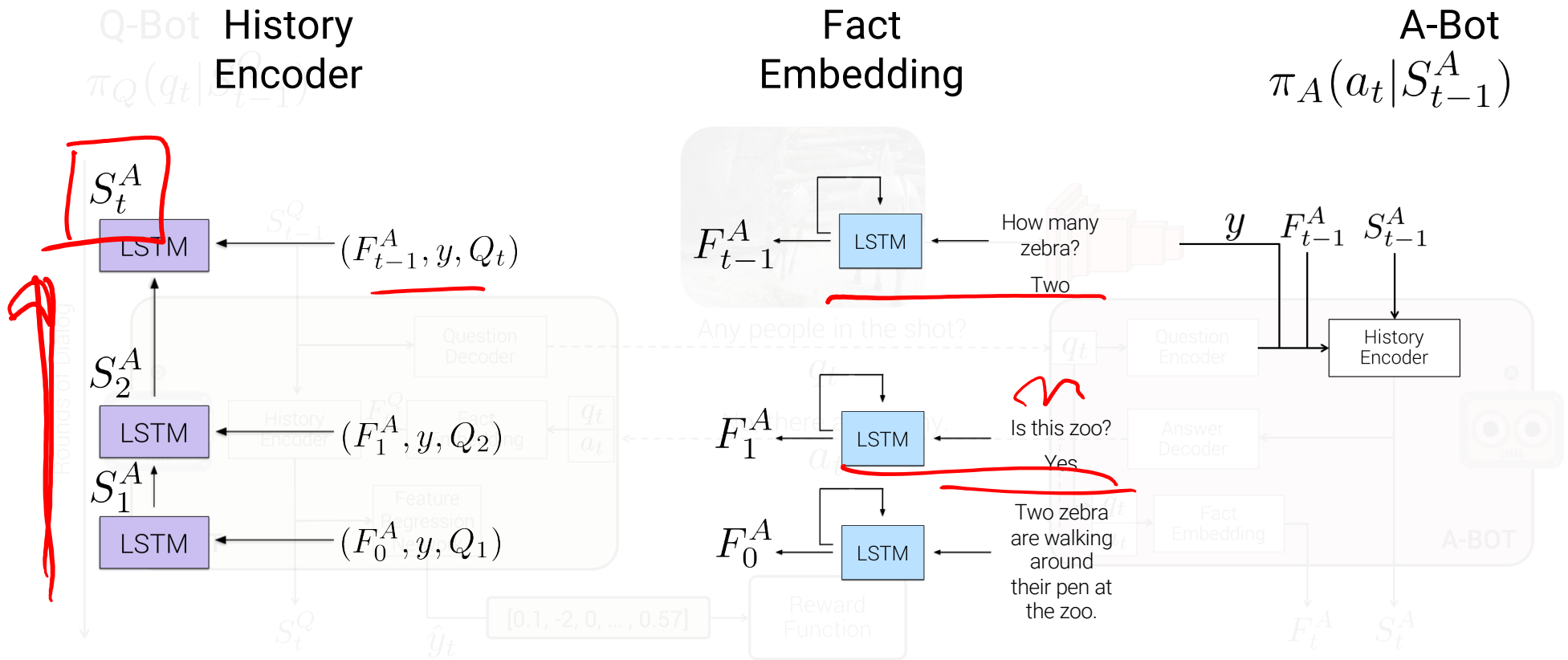
Policy Networks



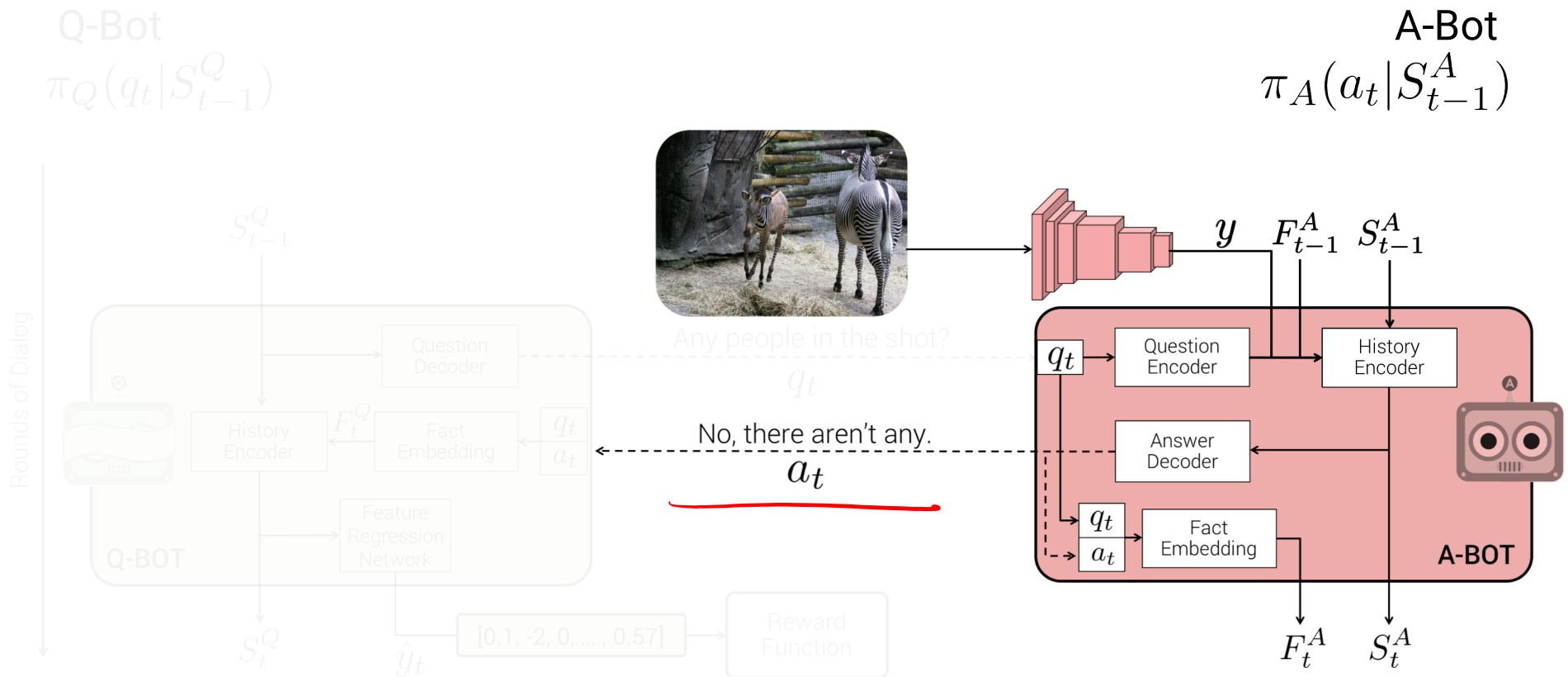
Policy Networks



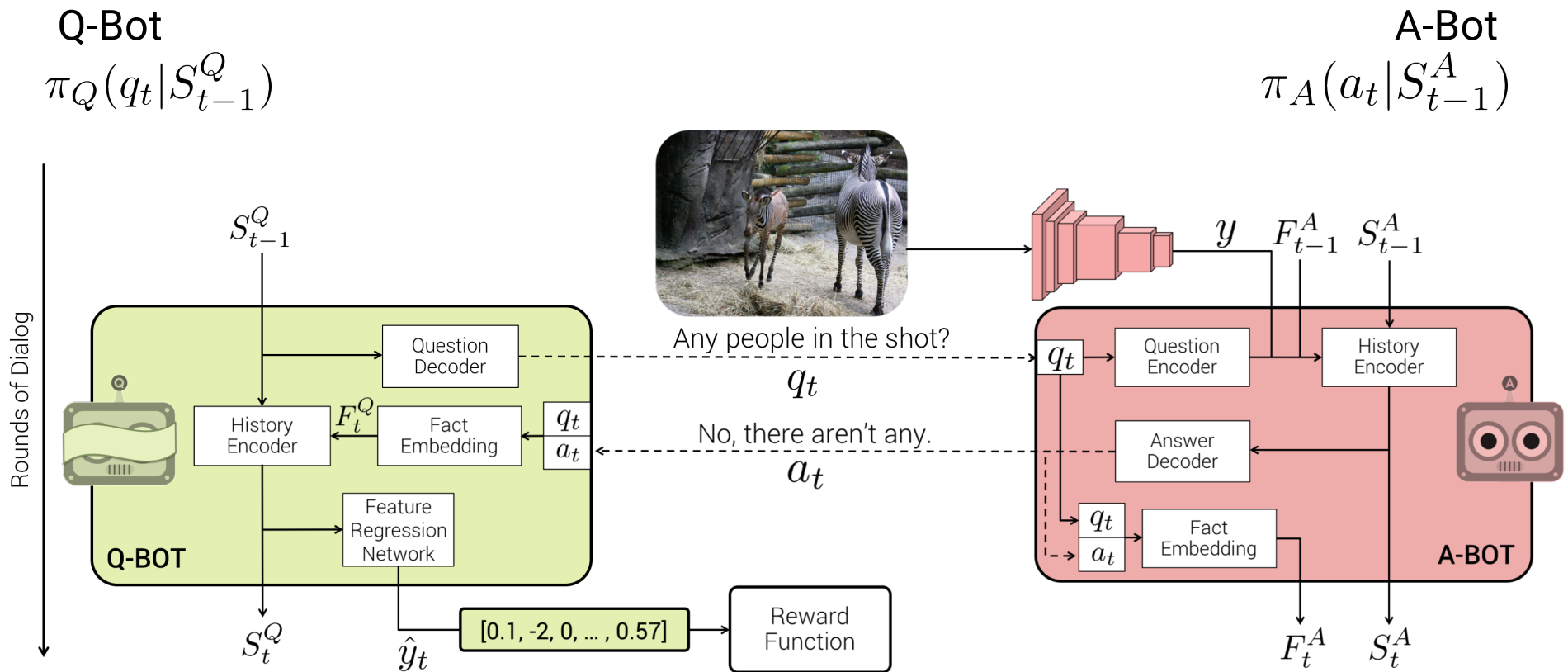
Policy Networks



Policy Networks



Policy Networks



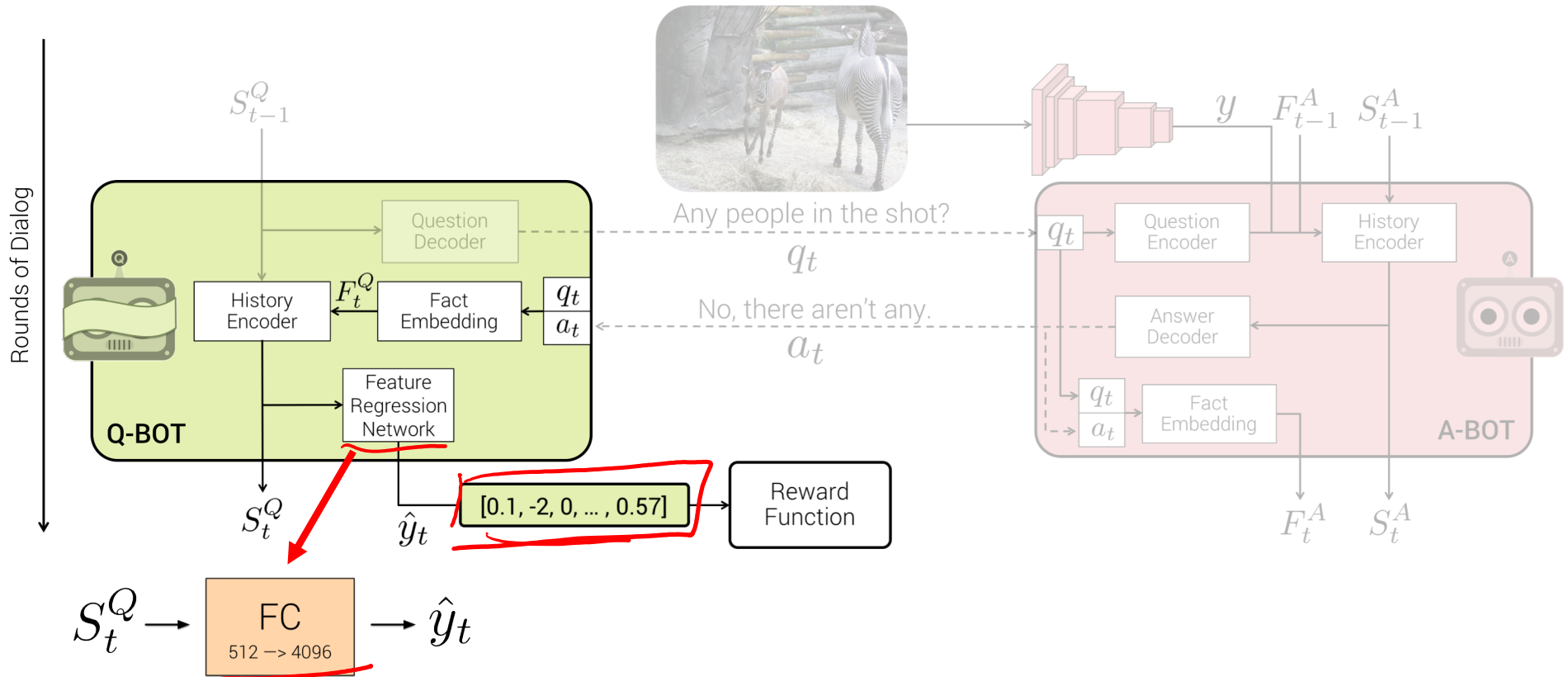
Policy Networks

Q-Bot

$$\pi_Q(q_t | S_{t-1}^Q)$$

A-Bot

$$\pi_A(a_t | S_{t-1}^A)$$



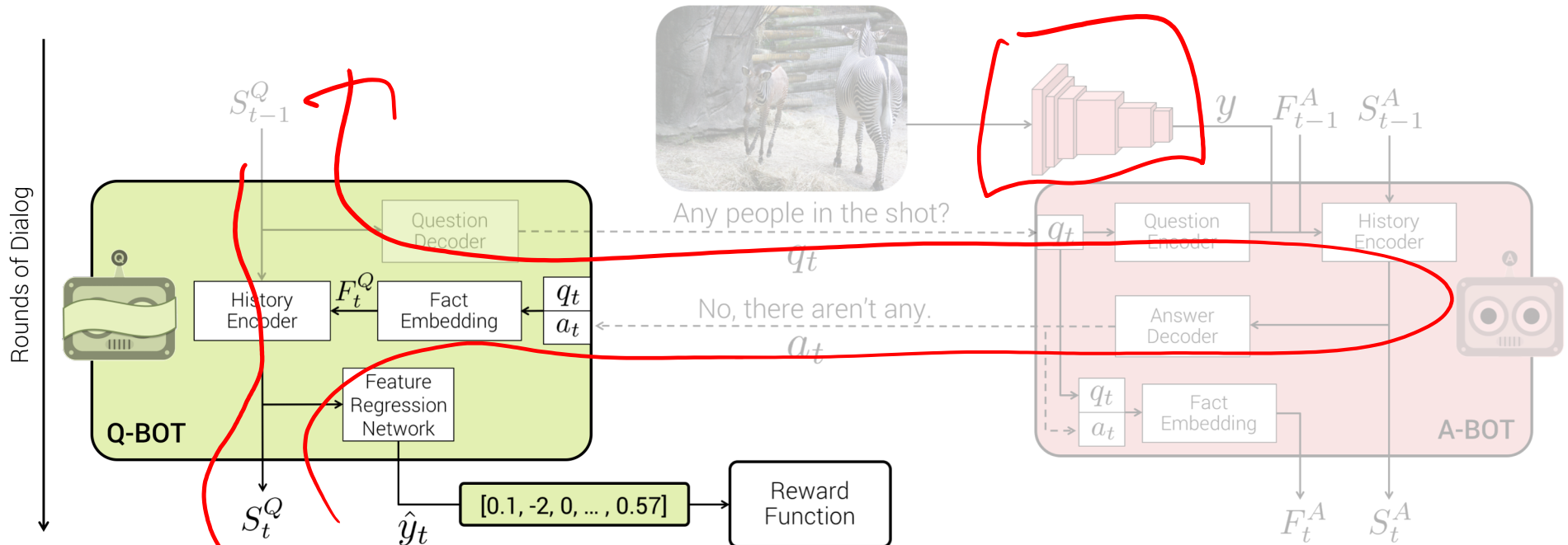
Policy Networks

Q-Bot

$$\pi_Q(q_t | S_{t-1}^Q)$$

A-Bot

$$\pi_A(a_t | S_{t-1}^A)$$



$$r_t \left(\underbrace{s_t^Q}_{\text{state}}, \underbrace{(q_t, a_t, y_t)}_{\text{action}} \right) = \underbrace{\ell(\hat{y}_{t-1}, y^{gt})}_{\text{distance at } t-1} - \underbrace{\ell(\hat{y}_t, y^{gt})}_{\text{distance at } t}$$

Policy Gradients

$$J(\theta_A, \theta_Q) = \mathbb{E}_{\pi_Q, \pi_A} \left[r_t(s_t^Q, (q_t, a_t, y_t)) \right]$$

REINFORCE Gradients

$$\begin{aligned} \nabla_{\theta_Q} J &= \nabla_{\theta_Q} \left[\mathbb{E}_{\pi_Q, \pi_A} [r_t(\cdot)] \right] \\ &= \sum_{q_t, a_t} \pi_Q(q_t | s_{t-1}^Q) \nabla_{\theta_Q} \log \pi_Q(q_t | s_{t-1}^Q) \pi_A(a_t | s_t^A) r_t(\cdot) \\ &= \mathbb{E}_{\pi_Q, \pi_A} \left[r_t(\cdot) \nabla_{\theta_Q} \log \pi_Q(q_t | s_{t-1}^Q) \right] \end{aligned}$$

Turing Test



How old is this man?	He looks like late teens
Does he have short hair or long hair?	Medium
What color is his hair?	Brown
Is he wearing glasses?	No
Does he have facial hair?	No
How is he dressed?	He's white dress shirt with red tie ... school uniform
Is he larger man or is he skinny?	He's skinny
Is this picture taken outside?	Inside
Can you tell if it's daytime?	I'm not sure
Can you tell what room he is standing in?	He's in hallway

A man making the live long and prosper sign from star trek.

Where is man located?	Looks like classroom of some sort
<u>What is man wearing?</u>	<u>Black t shirt and jeans</u>
Are there any other people in picture?	Yes, quite few people in background
What color is table?	Picture is black and white, but it's wood table
Is there anything else on table?	Not that I can see
What are other people doing?	They are standing in front of him, but I don't ...
Are there any windows?	Not that I can see
What color are walls?	I can't see walls
What time of day do you think it is?	I can't tell
<u>What is man wearing?</u>	<u>Black t shirt and jeans</u>



How old is man?
 What is he wearing?
 What color is his shirt?
 what color are his pants?
 what color is his shirt?
 What color are his shoes?
 What color are his shoes?
 What color are his shoes?
 What color are his shoes?
 What color are his shoes?

I can't see his face, I don't know
 shorts and t shirt
 white
 black
 blue
 Black
 Black
 Black
 Black
 Black

SL Agents

What color is skateboard?
 Is he wearing helmet?
 Is he wearing knee pads?
 Is it day or night?
 Is it sunny?
 Can you see any other people?
 Is it day or night?
 Is it sunny?
 What color is skateboard?
 Is he wearing helmet?

RL Agents

It is hard to tell but I think it's black
 Yes
 No, no knee pads
 It is daytime
 It is bright out, but I see no shadows, so I can not
 There are 2 other people in background
 It is daytime
 Yes, it is sunny
 It is hard to tell but I think it's black
 Yes, he is wearing helmet

Plan for Today

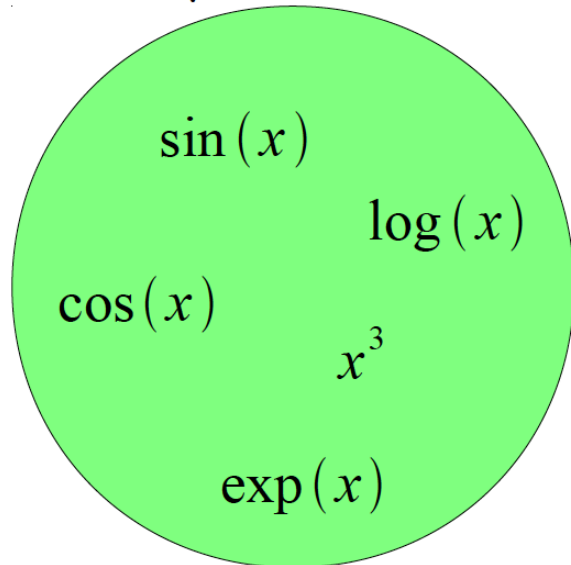
- (Deep) Reinforcement Learning
 - Policy gradients
- Closing the loop

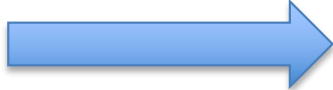
So what *is* Deep (Machine) Learning?

- A few different ideas:
- (Hierarchical) Compositionality
 - Cascade of non-linear transformations
 - Multiple layers of representations
- End-to-End Learning
 - Learning (goal-driven) representations
 - Learning to feature extraction
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons work together

Building A Complicated Function

Given a library of simple functions

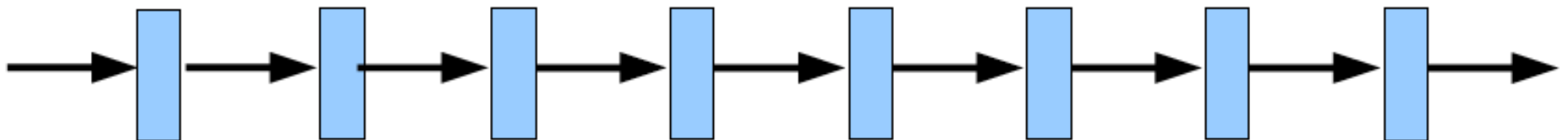


Compose into a

complicate function

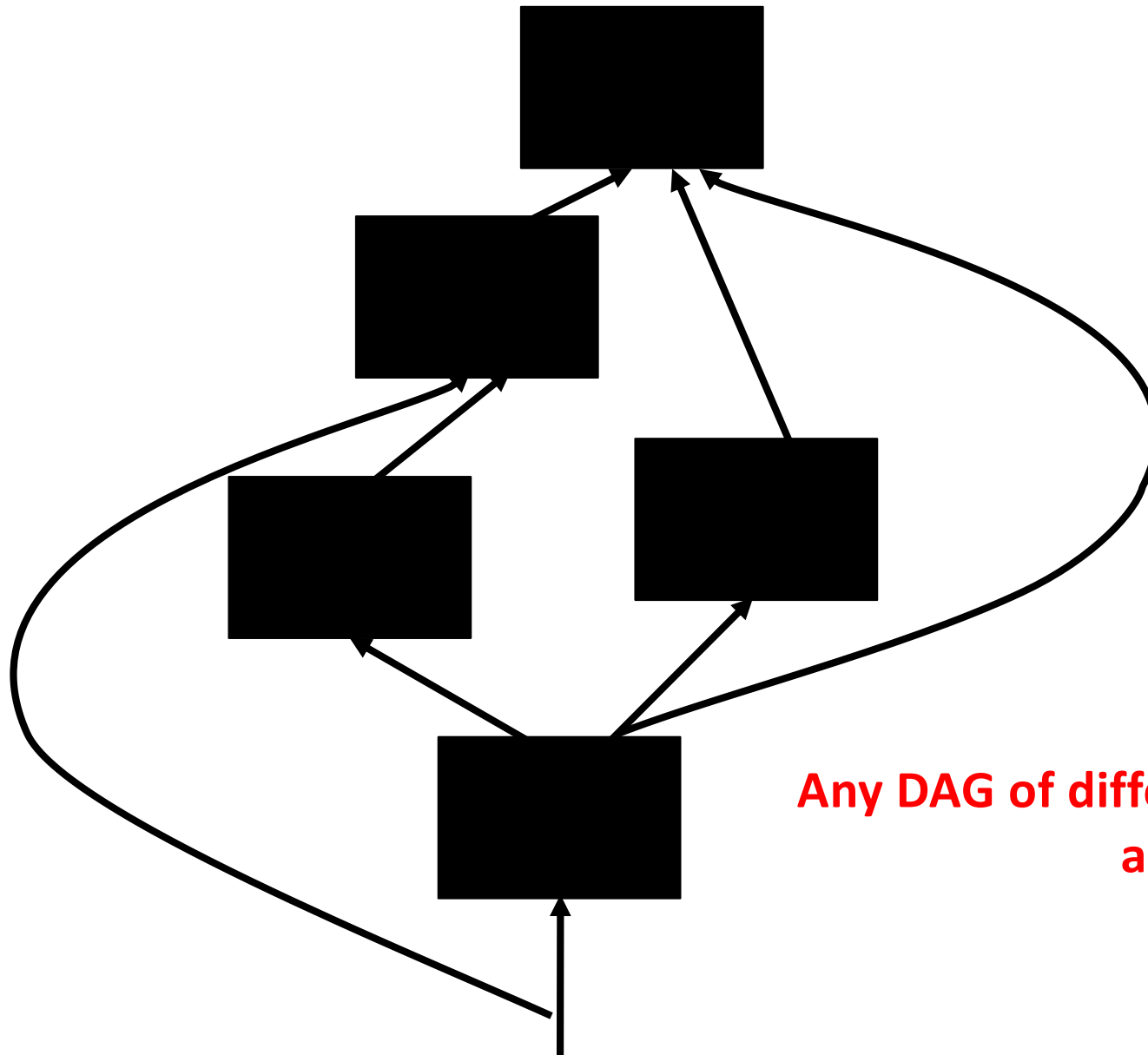
Idea 2: Compositions

- Deep Learning
- Grammar models
- Scattering transforms...

$$f(x) = g_1(g_2(\dots(g_n(x)\dots)))$$



Differentiable Computation Graph



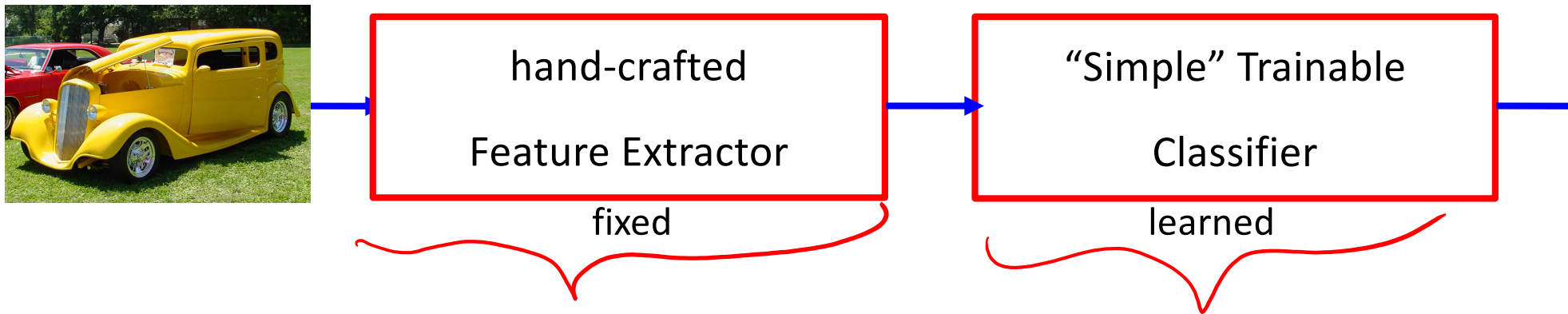
Any DAG of differentiable modules is allowed!

So what *is* Deep (Machine) Learning?

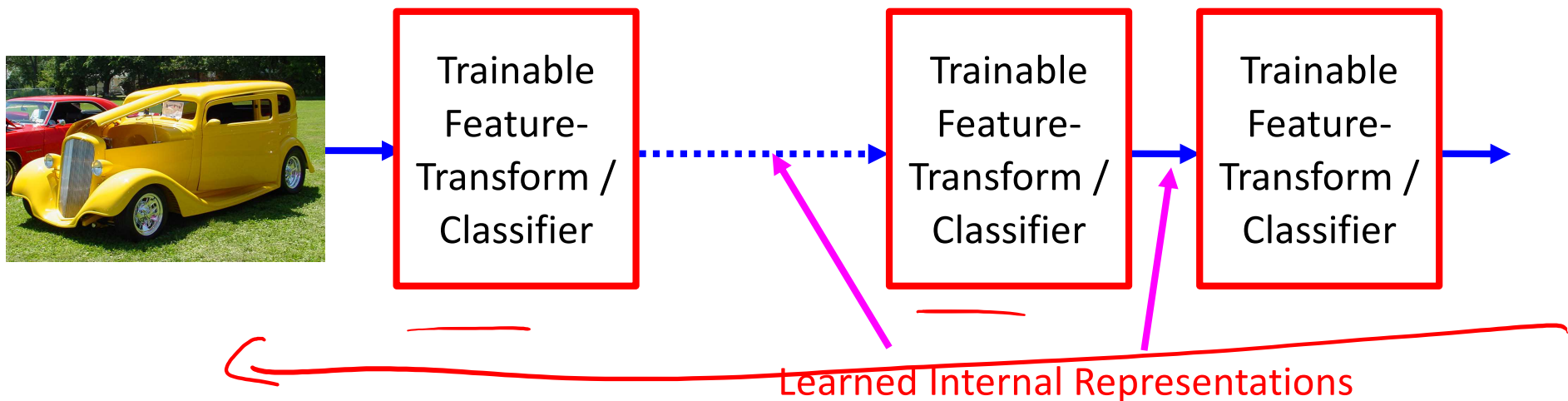
- A few different ideas:
 - (Hierarchical) Compositionality
 - Cascade of non-linear transformations
 - Multiple layers of representations
 - End-to-End Learning
 - Learning (goal-driven) representations
 - Learning to feature extraction
 - Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons work together

“Shallow” vs Deep Learning

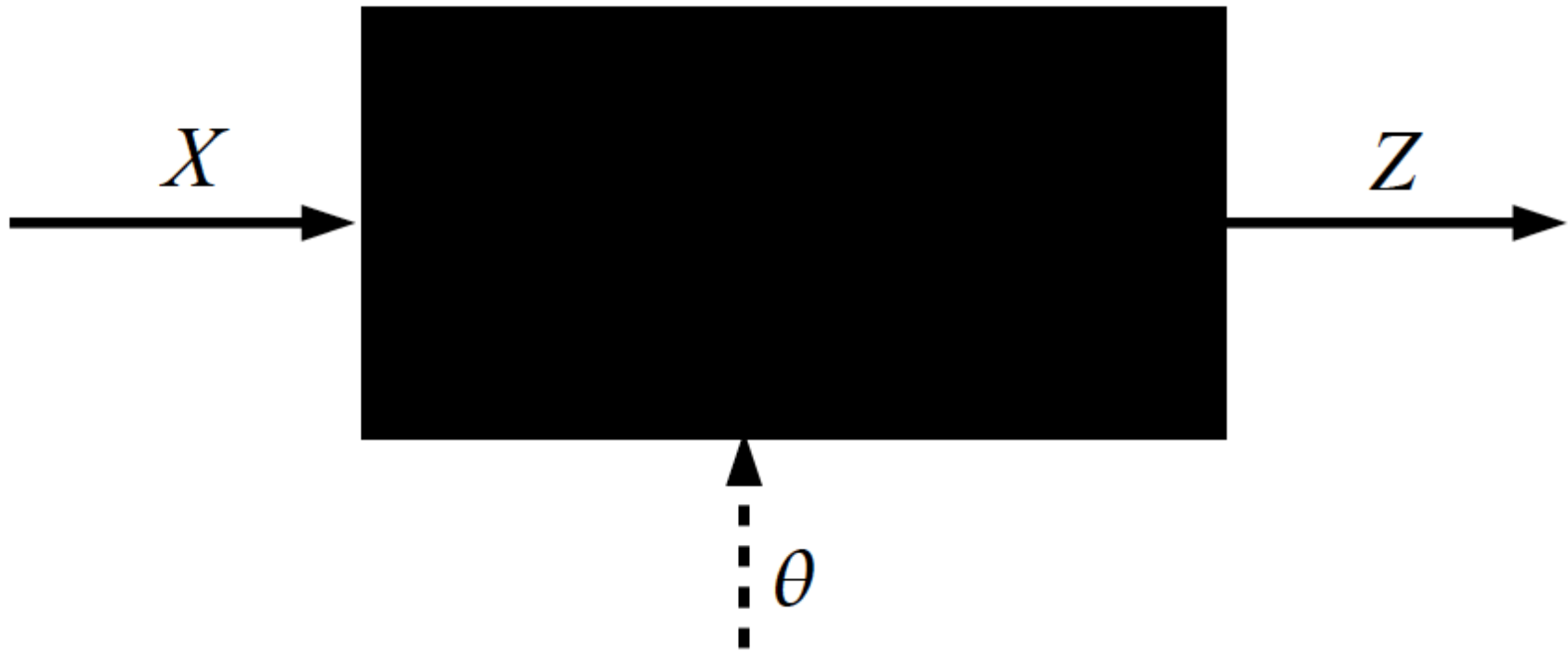
- “Shallow” models



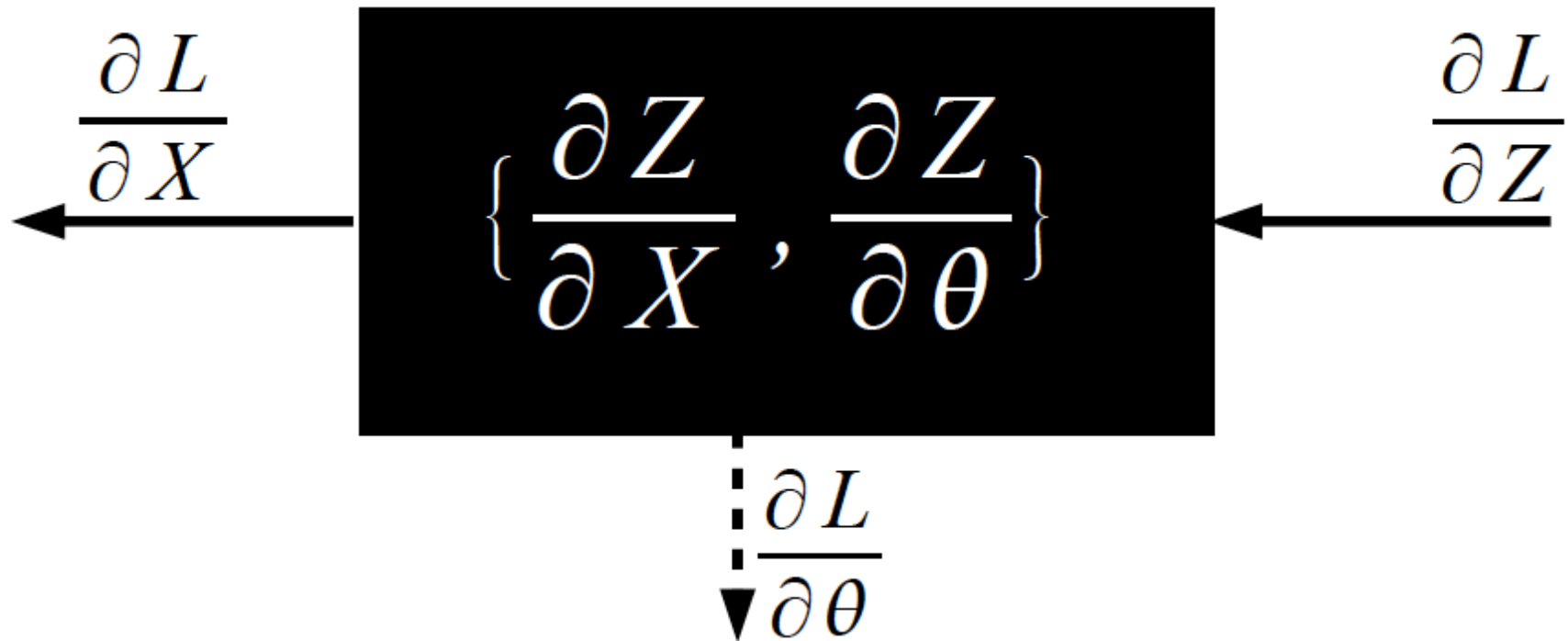
- Deep models



Key Computation: Forward-Prop



Key Computation: Back-Prop

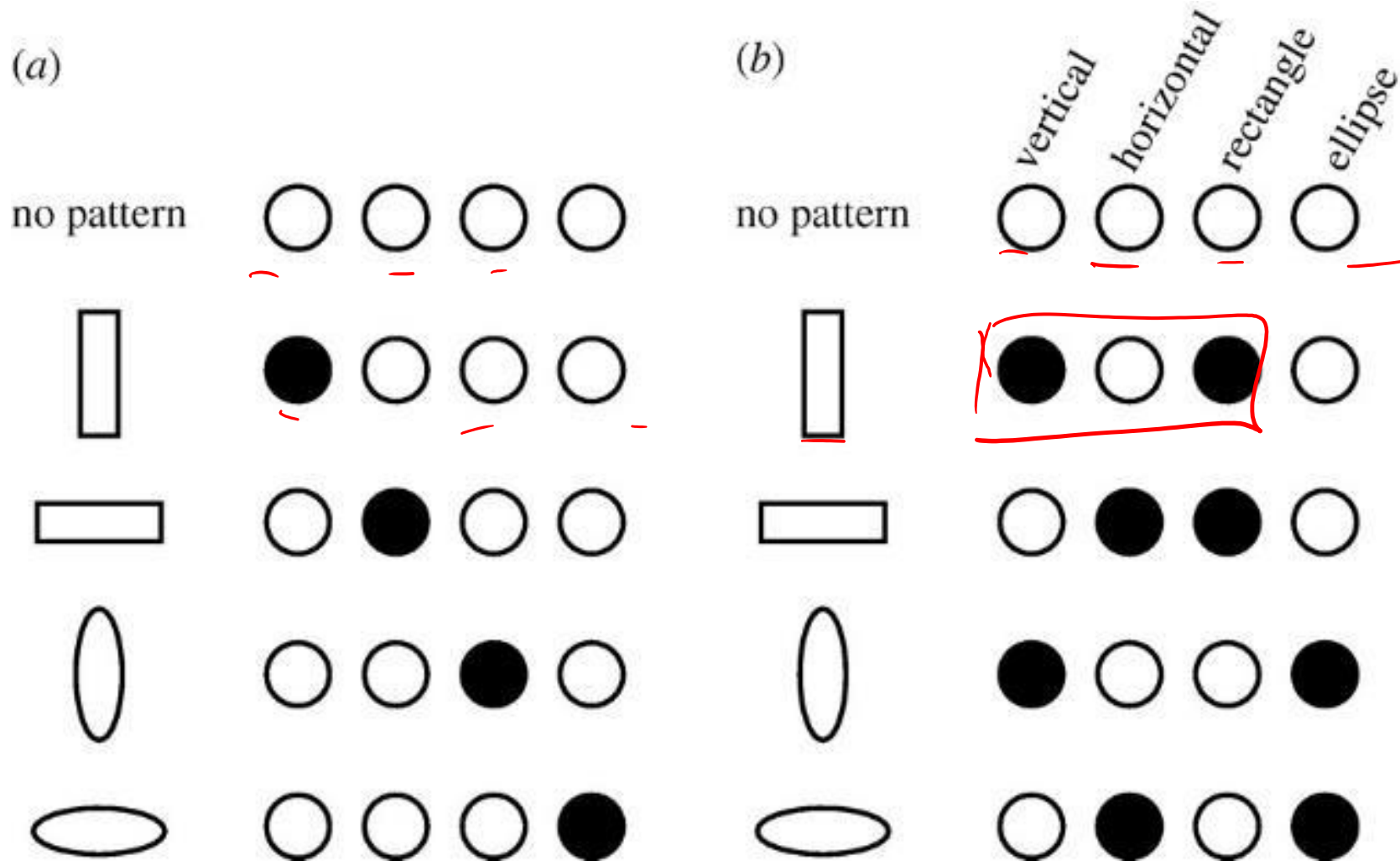


So what *is* Deep (Machine) Learning?

- A few different ideas:
- (Hierarchical) Compositionality
 - Cascade of non-linear transformations
 - Multiple layers of representations
- End-to-End Learning
 - Learning (goal-driven) representations
 - Learning to feature extraction
- Distributed Representations
 - No single neuron “encodes” everything
 - Groups of neurons work together

Distributed Representations Toy Example

- Can we interpret each dimension?



Power of distributed representations!

Local

$$\underbrace{\bullet \bullet} \quad \circ \quad \underbrace{\bullet} = VR + HR + HE = ?$$

Distributed

$$\underbrace{\bullet \bullet} \quad \circ \quad \underbrace{\bullet} = V + H + E \approx \bigcirc$$

What is this class about?

What is this class about?

- Introduction to Deep Learning
- Goal:
 - After finishing this class, you should be ready to get started on your first DL research project.
 - CNNs
 - RNNs
 - Deep Reinforcement Learning
 - Generative Models (VAEs, GANs)
- Target Audience:
 - Senior undergrads, MS-ML, and new PhD students

What did we learn?

- **Background & Basics**
 - Neural Networks, Backprop, Optimization (SGD)
- **Module 1: Convolutional Neural Networks (CNNs)**
 - Architectures, Pre-training, Fine-tuning
 - Visualizations, Fooling CNNs, Adversarial examples
 - Different tasks: detection CNNs, segmentation CNNs
- **Module 2: Recurrent Neural Networks (RNNs)**
 - Difficulty of learning; “Vanilla” RNNs, LSTMs, GRU
 - RNNs for Sequence-to-Sequence (machine translation & image captioning, VQA, Visual Dialog)
- **Module 3: Deep Reinforcement Learning**
 - Overview, policy gradients
 - Optimizing Neural Sequence Models for goal-driven rewards
- **Module 4: Deep Structured Prediction**
 - Crash course on Bayes Nets, Variational Inference
 - Variational Auto Encoders (VAEs)
- **Module 5: Advanced Topics**
 - GANs, Adversarial Learning

Arxiv Fire Hose

PhD Student

Deep Learning papers



Feedback

CIOS Help

GT Course Instructor Opinion Survey (CIOS) now open

To: Batra, Dhruv,

Reply-To: cioshelp@gatech.edu

Inbox - GT November 26, 2018 at 1:37 AM

CH

Siri found new contact info in this email: Help Cios evaluations@smartevals.com

[add to Contacts...](#)

Dear Dhruv,

Good morning. The Course/Instructor Opinion Survey (CIOS) is now available for the following courses. Your courses, their survey start and end dates, and your current response rate are shown in the table below.

Eval	Course Prefix	Course Number	Sec	Type	Name	Begin	End	Not Resp.	Resp.	Tot.
Preview	CS	4803	DL	A	Special Topics	11-26	12-16	20	0	20
Preview	CS	7643	A	A	Deep Learning	11-26	12-16	84	0	84

Students have received an announcement indicating that surveys have begun, and they will continue to receive periodic reminder emails with all of the necessary information to complete the survey. However, you can ALSO set up additional reminders within the system that would come from you. Simply login at the link below, click the "Not Set" button near the left of the table, and follow the directions to set up auto-email reminders for your all of your courses.

Reports with your results will be available 5 days after full semester grades are due and you will receive an email with report access information at that time.

If you would like to view your response rates at any time, you can log in with your GT account here: <http://b.gatech.edu/cios>

If you have any problems with the survey system, please email cioshelp@gatech.edu.



Thanks!