

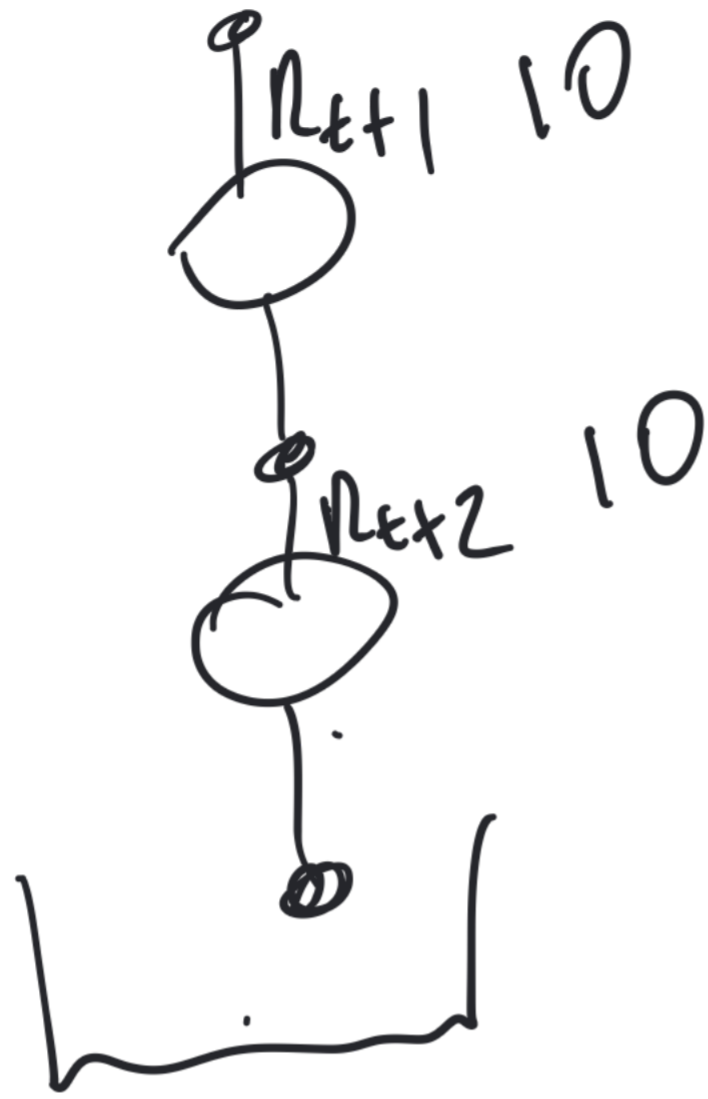
RL: Lecture 17

Harvey Mudd College

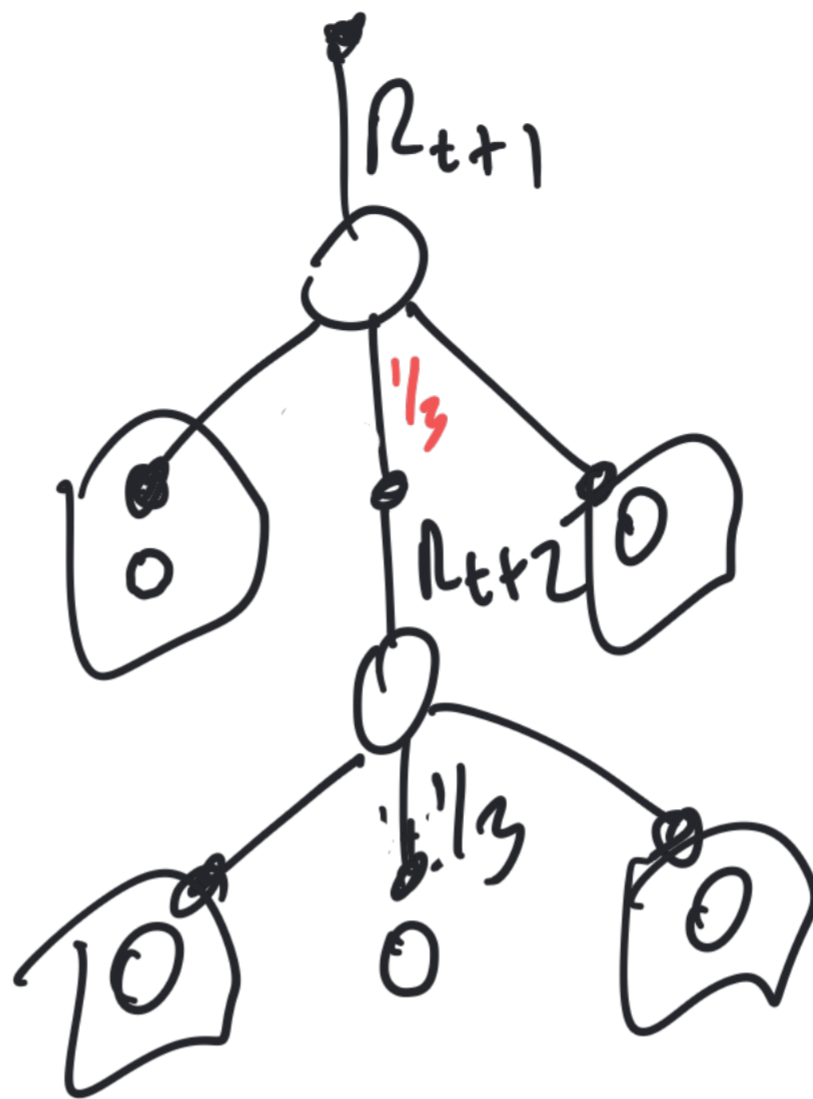
March 30, 2020

Neil Rhodes

on-policy



tree backup



tree backup

on-policy update

assum

initially

$$Q(s, a) = 0 \forall s, a$$

equiprobable policy

$$Q(s_t, A_t) \neq \alpha [R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(s_{t+2}, A_{t+2}) - Q(s_t, A_t)]$$

$$Q(s_t, A_t) = \alpha [R_{t+1}]$$

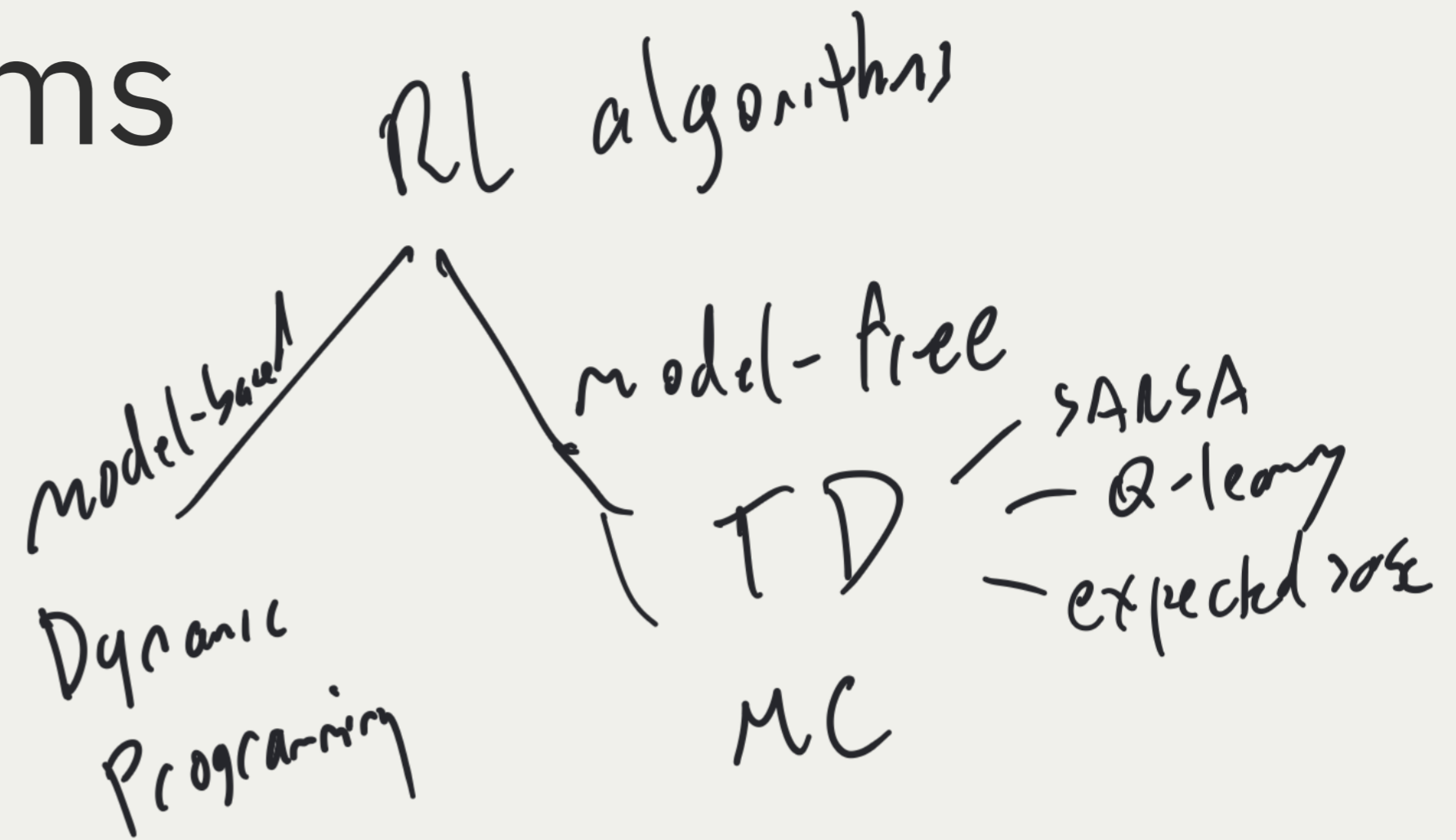
new tree backup

$$Q = R_t \quad \text{AND} \quad 2.20$$

$$Q = \alpha [R_{t+1} + \frac{1}{3} R_{t+2}]$$

Types of RL algorithms

- Model-based
- Model-free



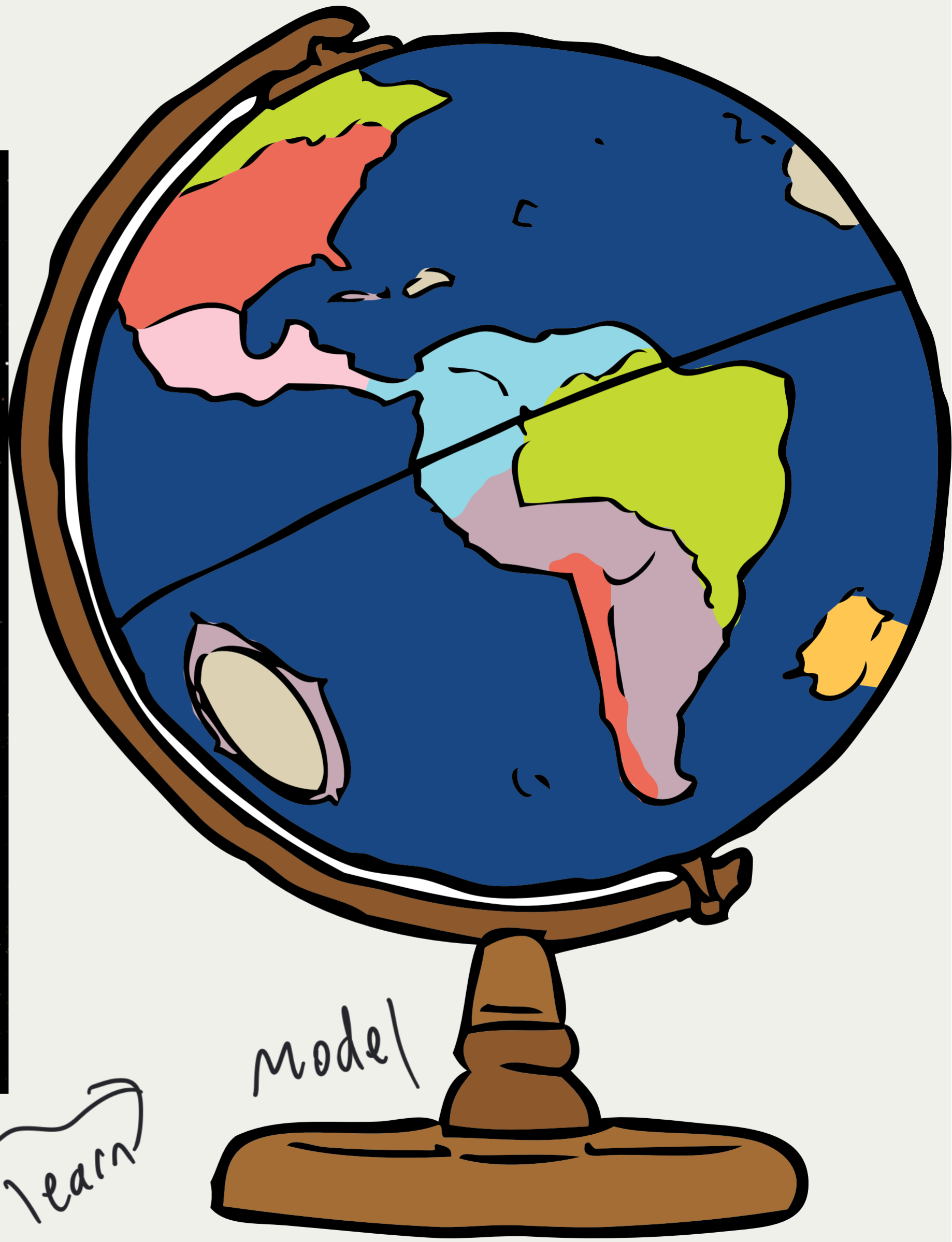
$$\sum P(s, r | s', a')$$

MDP environment



↑
environment

learn



model

Types of Models

- Distributional model
- Sample model

— actually have all the given probs $S, A \rightarrow$

prob $\left\{ \begin{array}{l} (S', A) \\ \vdots \end{array} \right.$

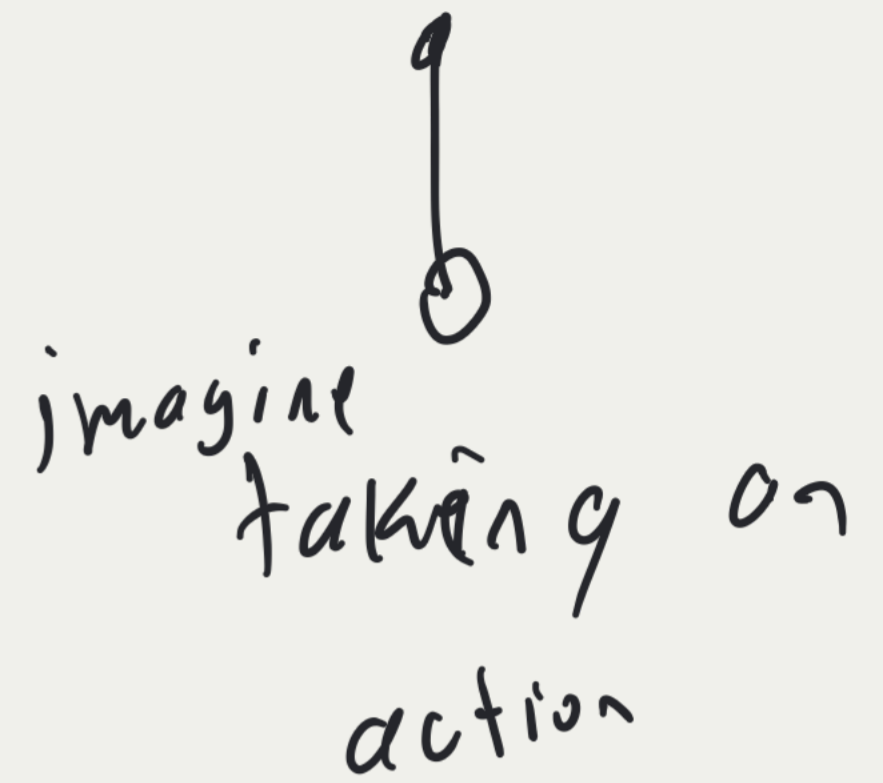
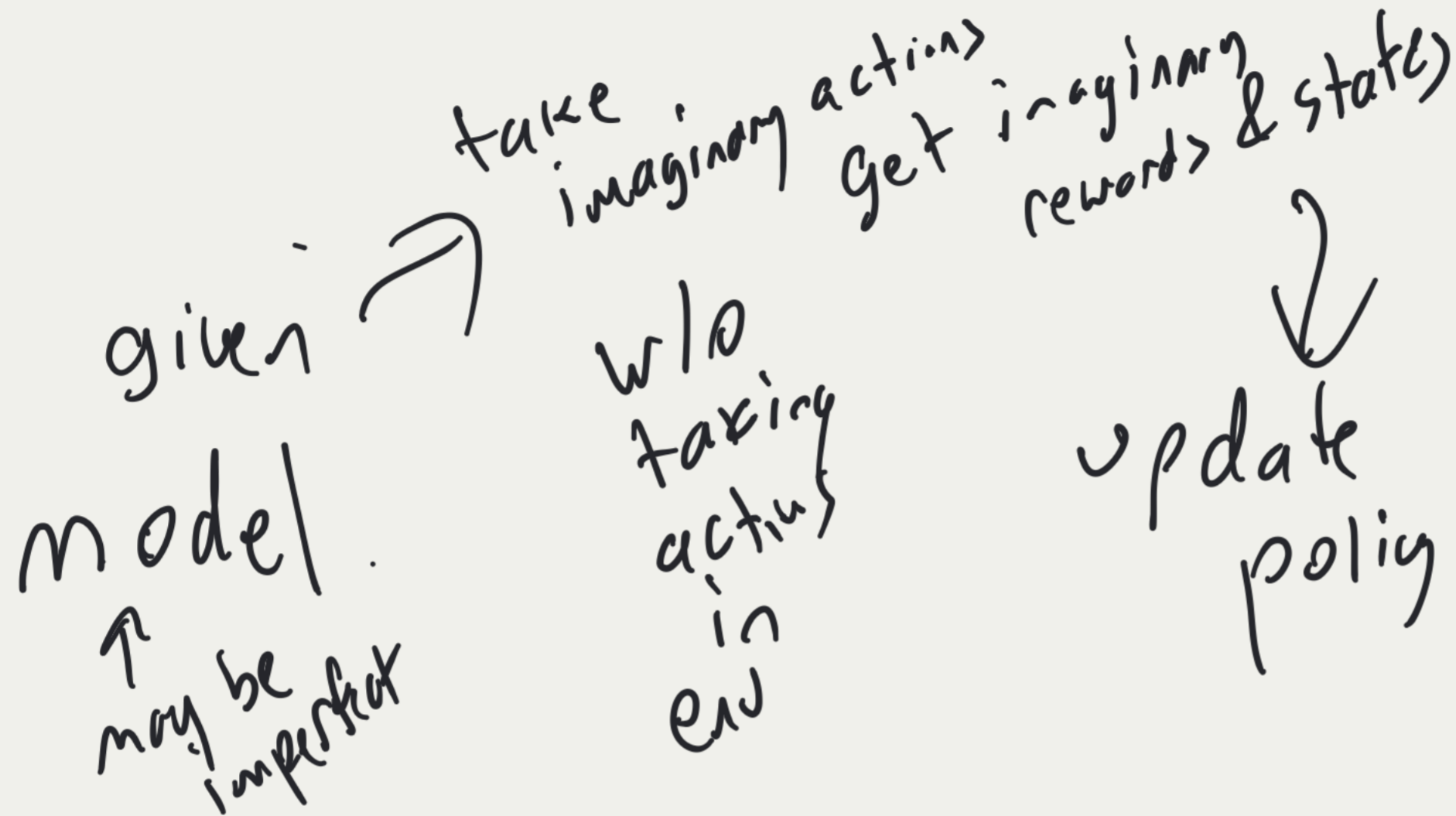
more constrained model

given an (S, A) it returns a sample S', A

Given a distribution model,

easy to come up w/ samples

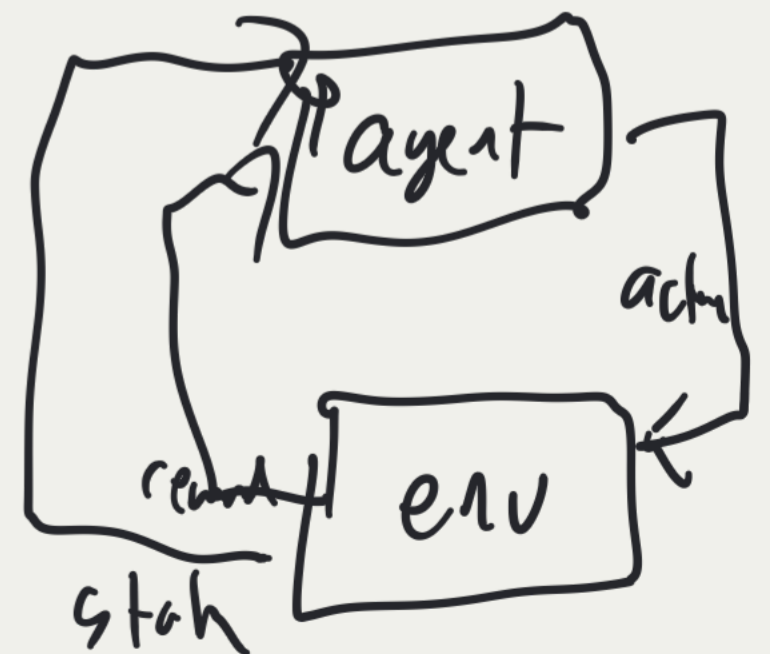
What is planning?



w/o model

Direct RL

use estimates of state values
take actions in env.



Random-sample one-step tabular Q-planning

Not learning planning ← uses model, not env

Loop forever:

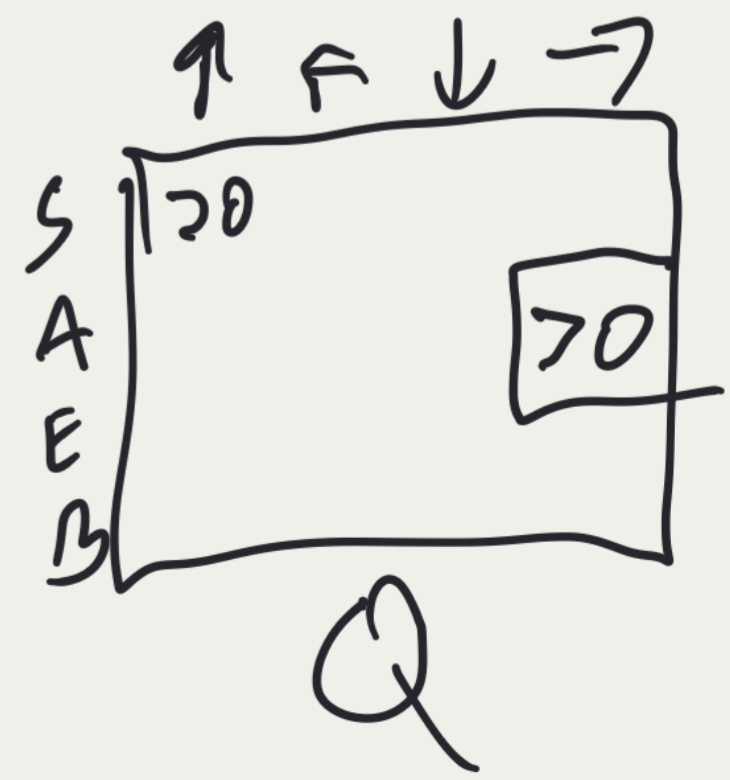
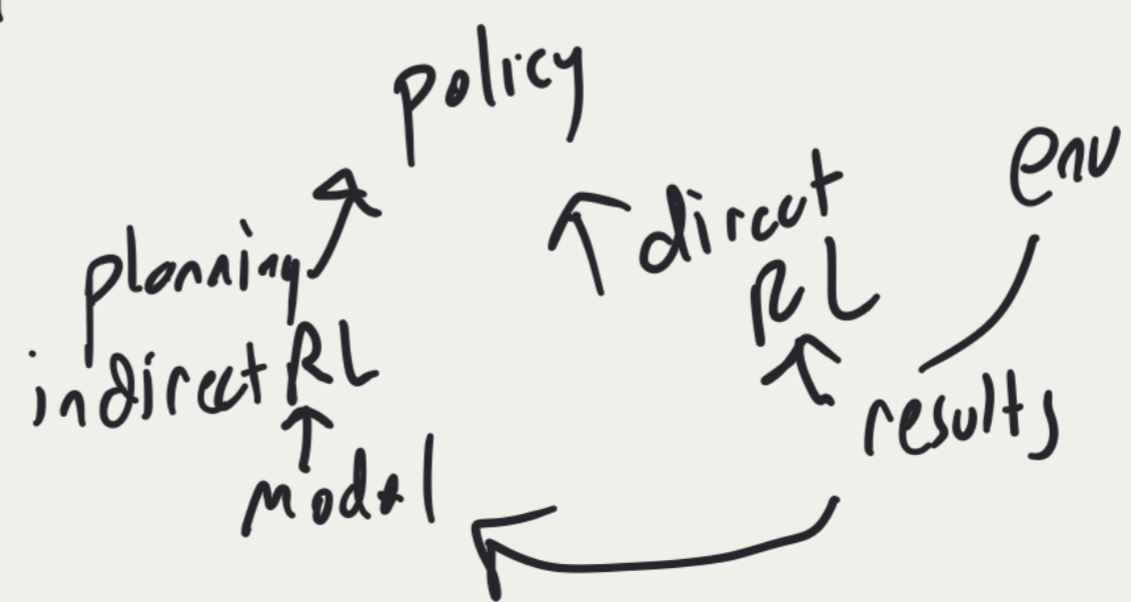
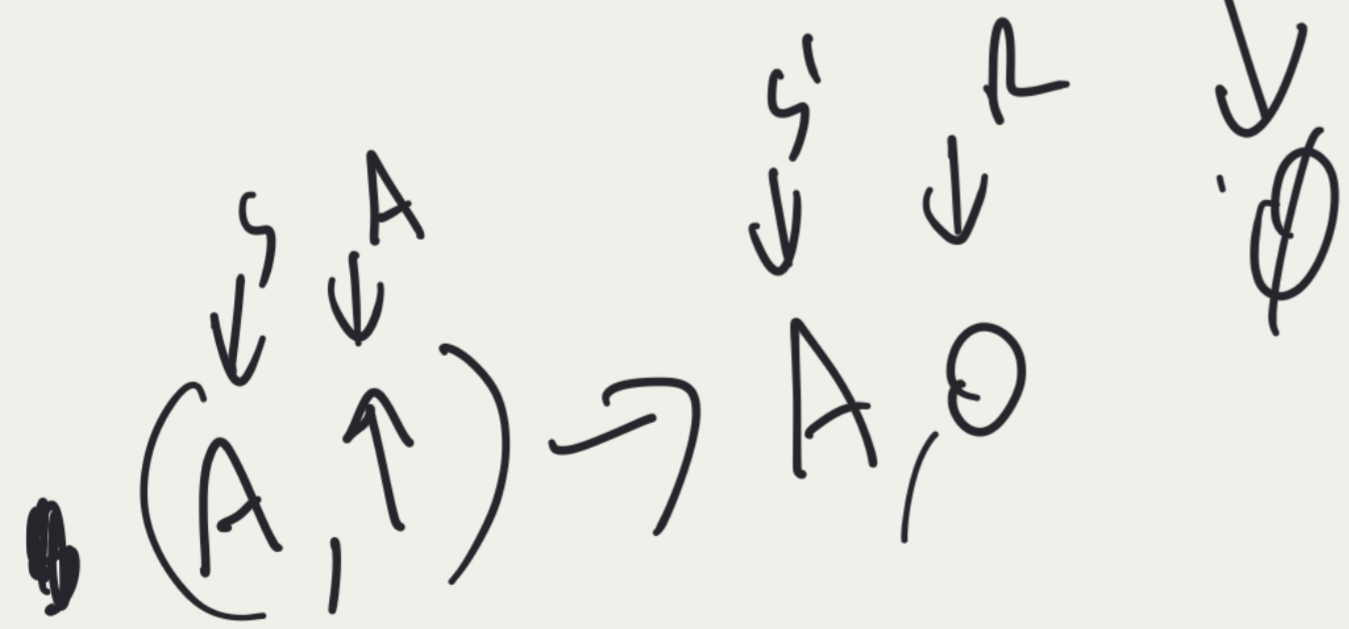
1. Select a state, S , and an action A at random
2. Send S, A to a sample model and obtain a sample next reward R , and a sample next state, S'
3. Apply one-step tabular Q-learning to S, A, R, S'

not the environment

imagined modeled S', R

updating actual policy

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$



Real episodes

stochastic

- 1) A 0 B 0
- 2) B 1
- 3) B 1
- 4) B 0
- 5) B 1
- 6) B 1
- 7) B 1

deterministic

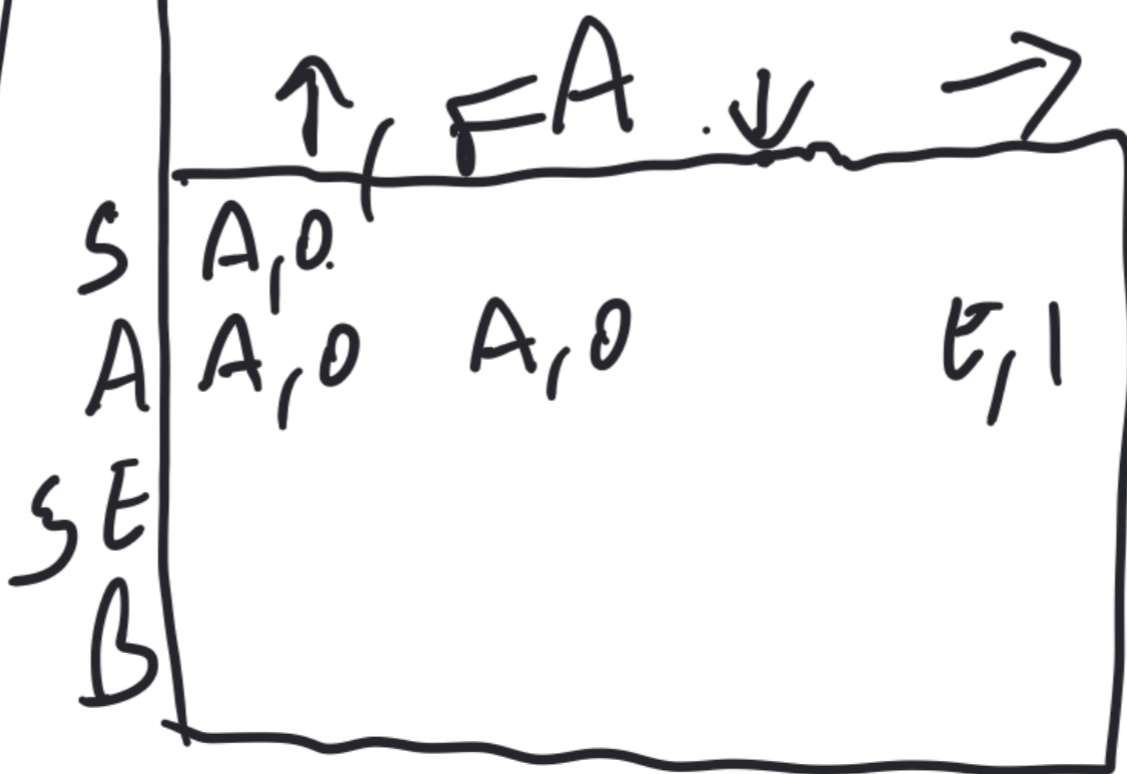
maze



$(S, \uparrow) \rightarrow A$

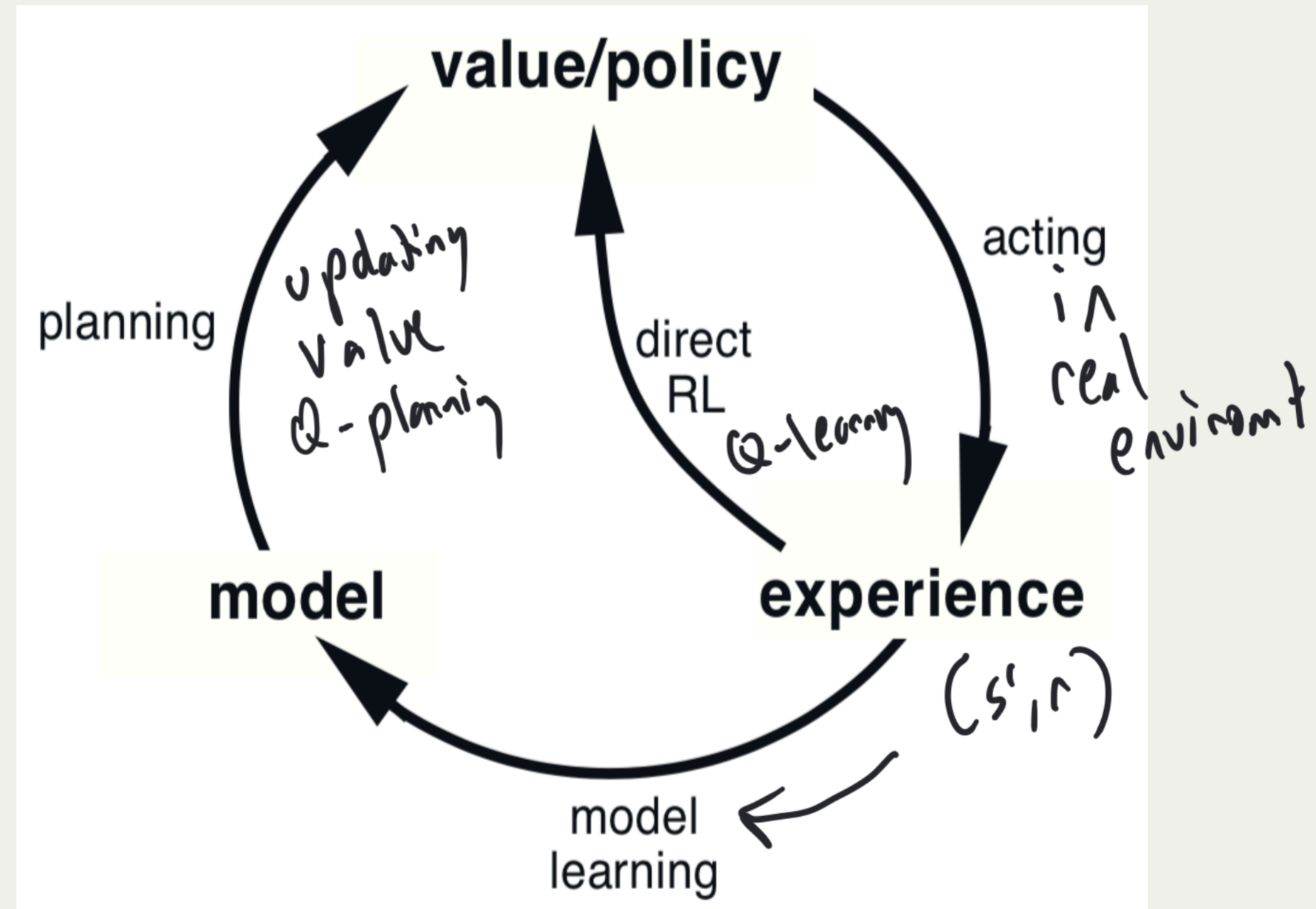
$(S, A) \rightarrow (S', R)$
 $(A, \uparrow) \rightarrow (A, 0)$

model
 $S \times A \rightarrow S \times R$

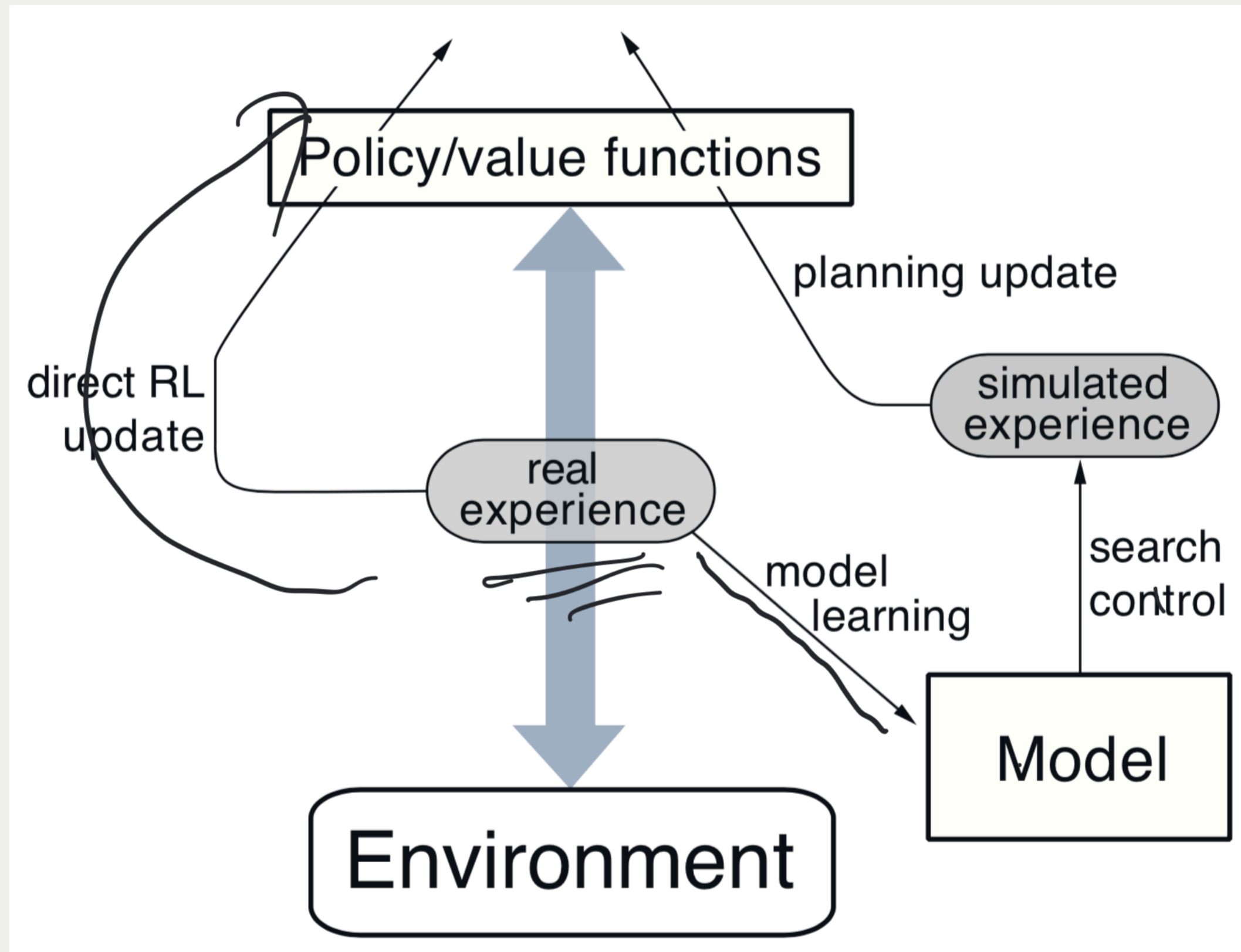


model of maze env

Online planning



Dyna Architecture



Tabular Dyna-Q Overview

- Direct RL method: one-step tabular Q-learning
- Model-learning method:
 - Assumes environment is deterministic
 - Table-based
 - Given $A_t, S_t \rightarrow R_{t+1}, S_{t+1}$, stores $\text{model}[(S_t, A_t)] = (R_{t+1}, S_{t+1})$

Tabular Dyna-Q Algorithm

Initialize $Q(s, a)$ and $Model(s, a)$ for all a, s

Loop forever:

1. $S \leftarrow$ current (nonterminal) state

2. $A \leftarrow \epsilon$ -greedy(S, Q)

3. Take action A ; observe resultant reward R and state S'

4. $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

5. $Model(S, A) \leftarrow (R, S')$ (assumes deterministic environment)

6. Loop repeat n times

$S \leftarrow$ random previously observed state

$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

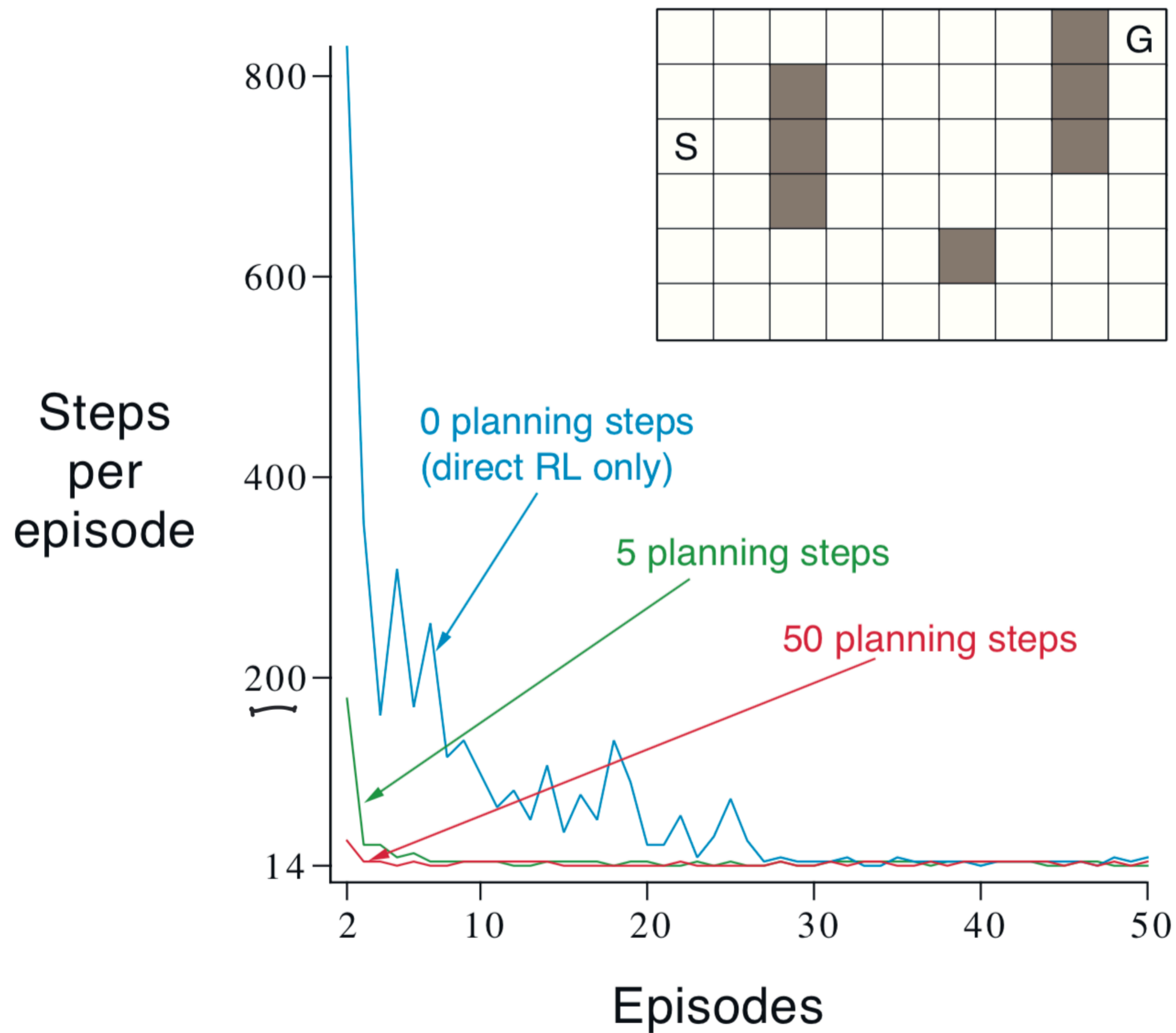
$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

direct RL

update model

planning

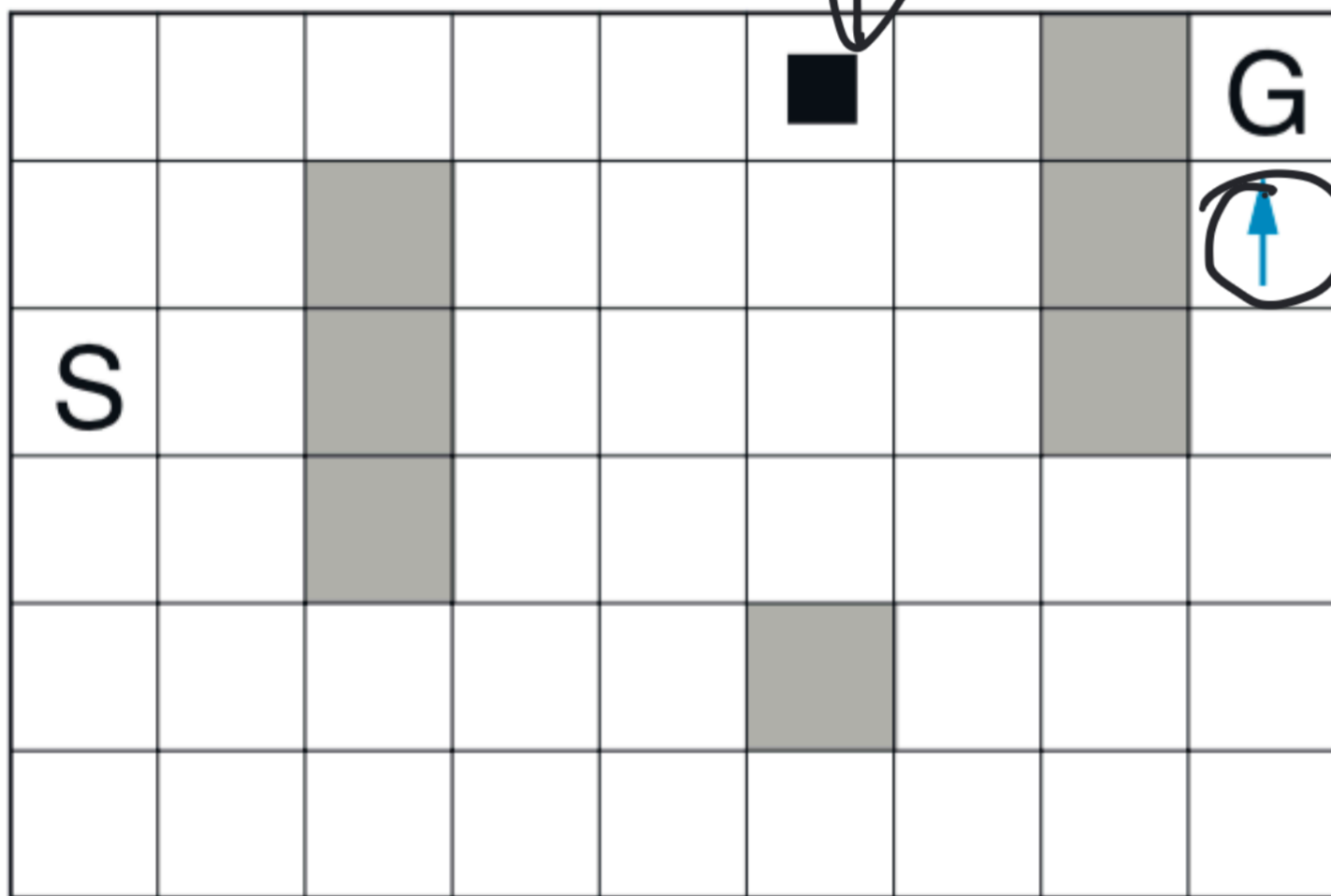
Dyna Maze (Figure 8.2)



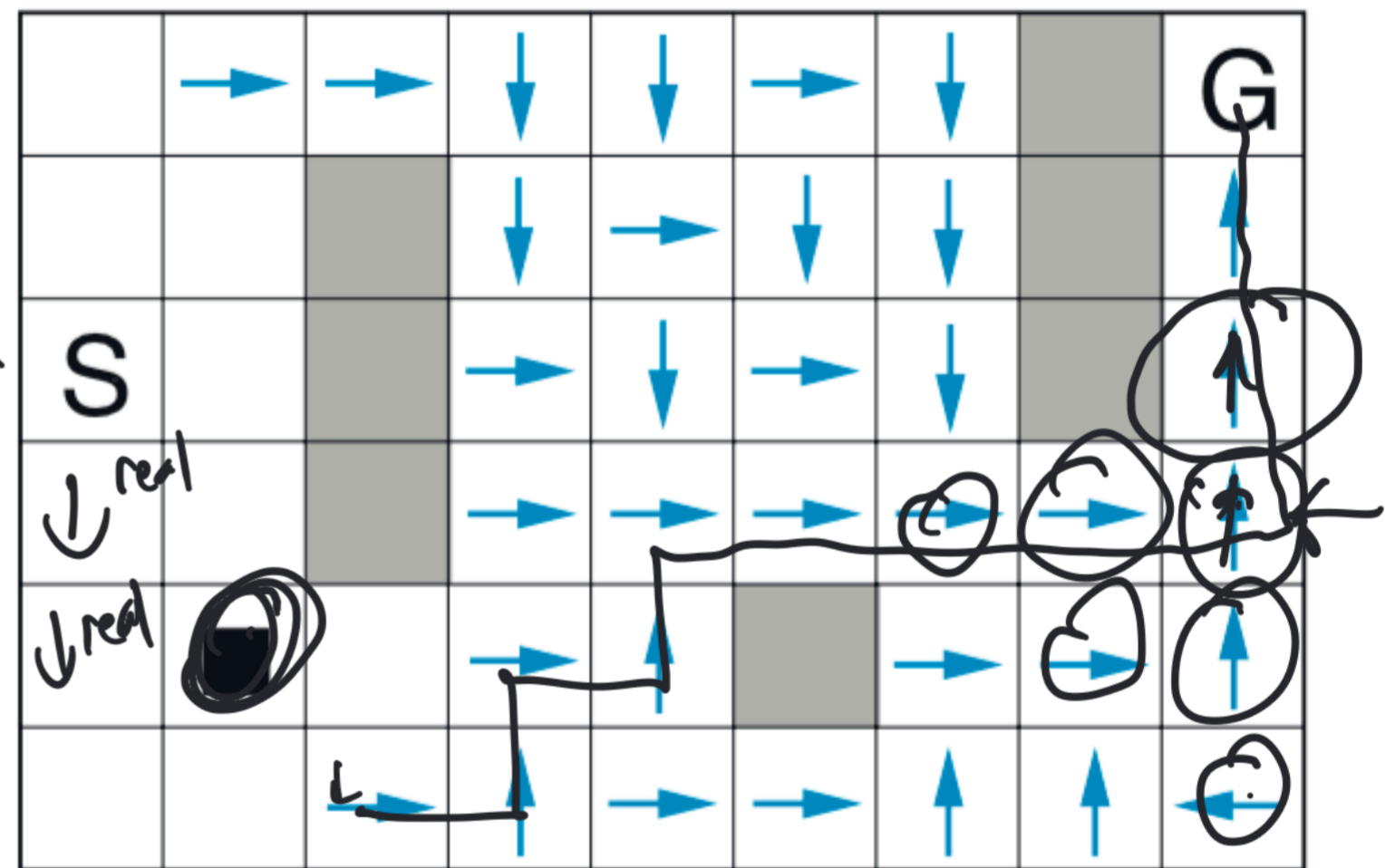
reward = 0
except @ goal
discount = .9

Dyna Maze (Figure 8.3)

WITHOUT PLANNING ($n=0$)



WITH PLANNING ($n=50$)

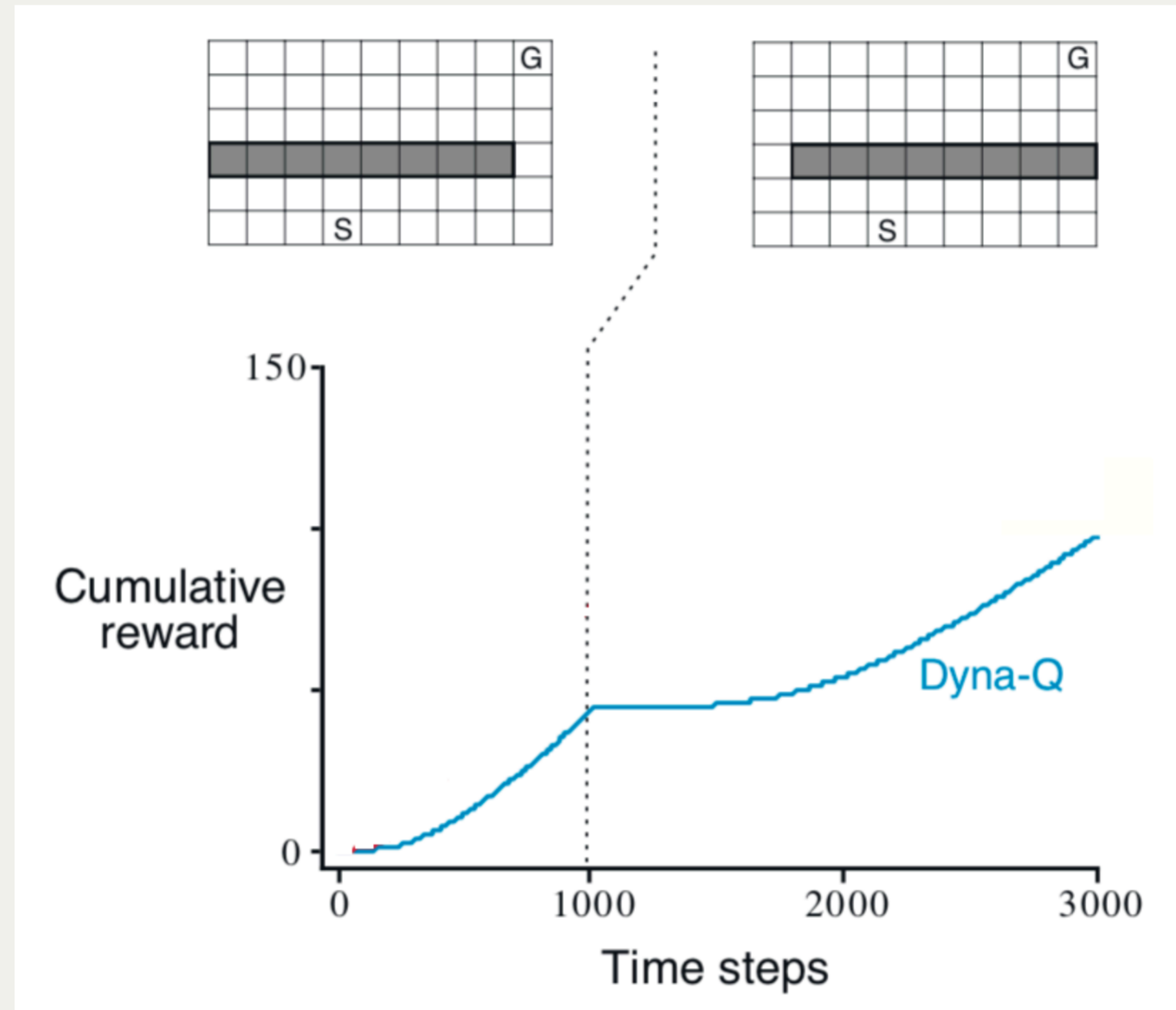


1/2 way through 2nd episode ↑

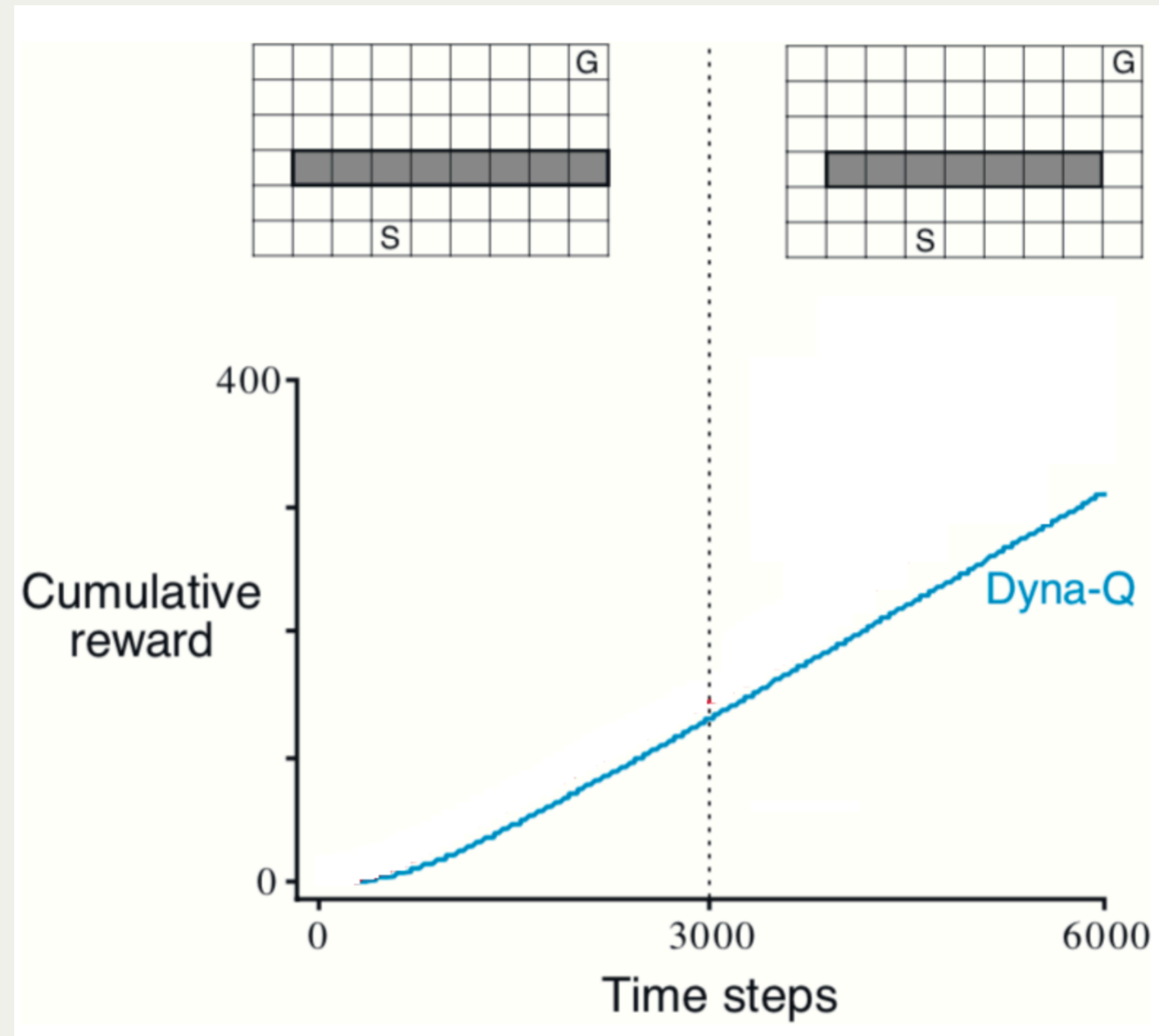
model-free + source of uncertainty: Q-values

model-based + additional source of uncertainty: model

When the Model is Wrong: Optimistic Model



When the Model is Wrong: Pessimistic Model



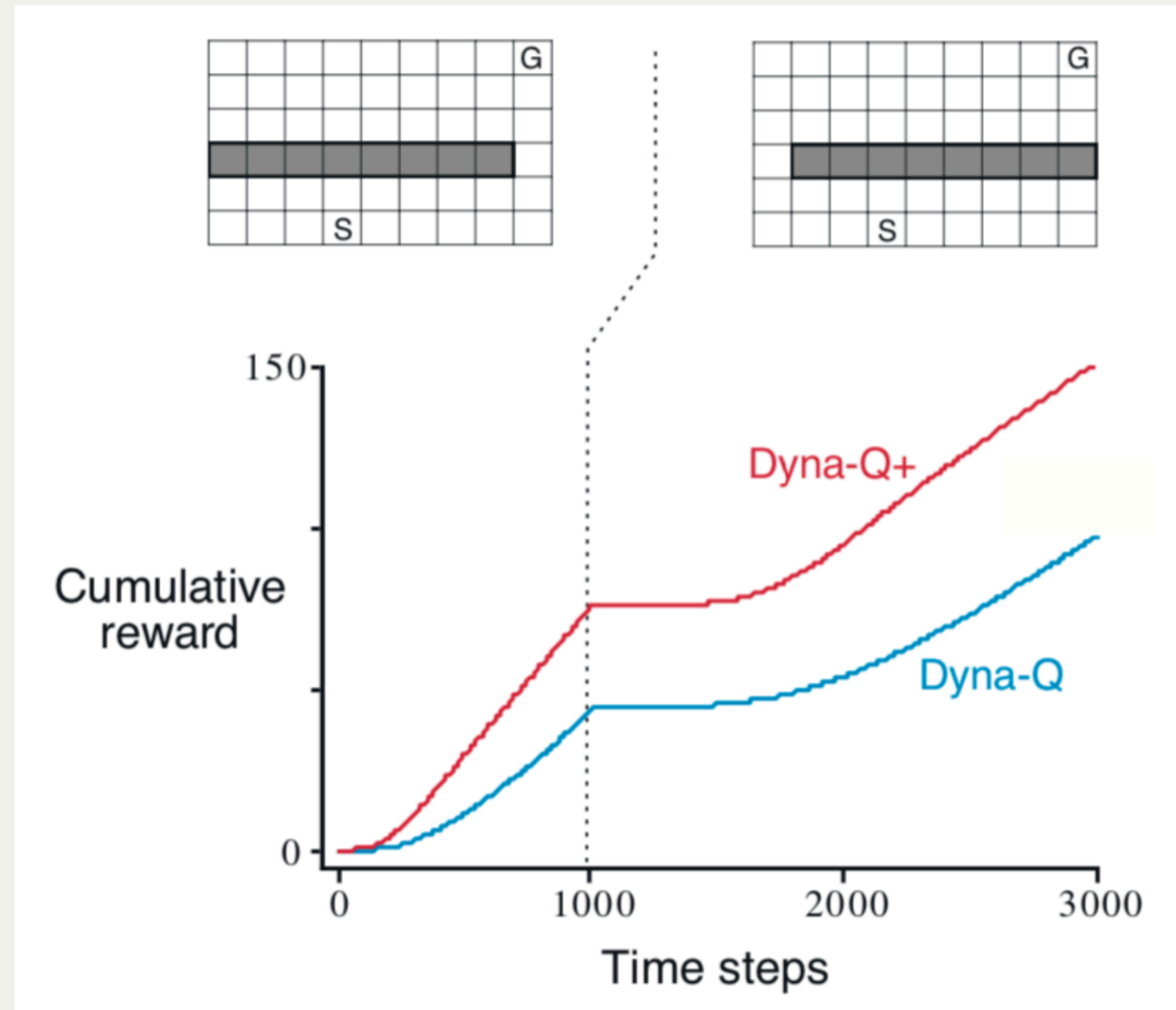
Dyna-Q+: Dyna-Q + heuristics for encouraging

- Provide an implicit reward to exploring stale transitions

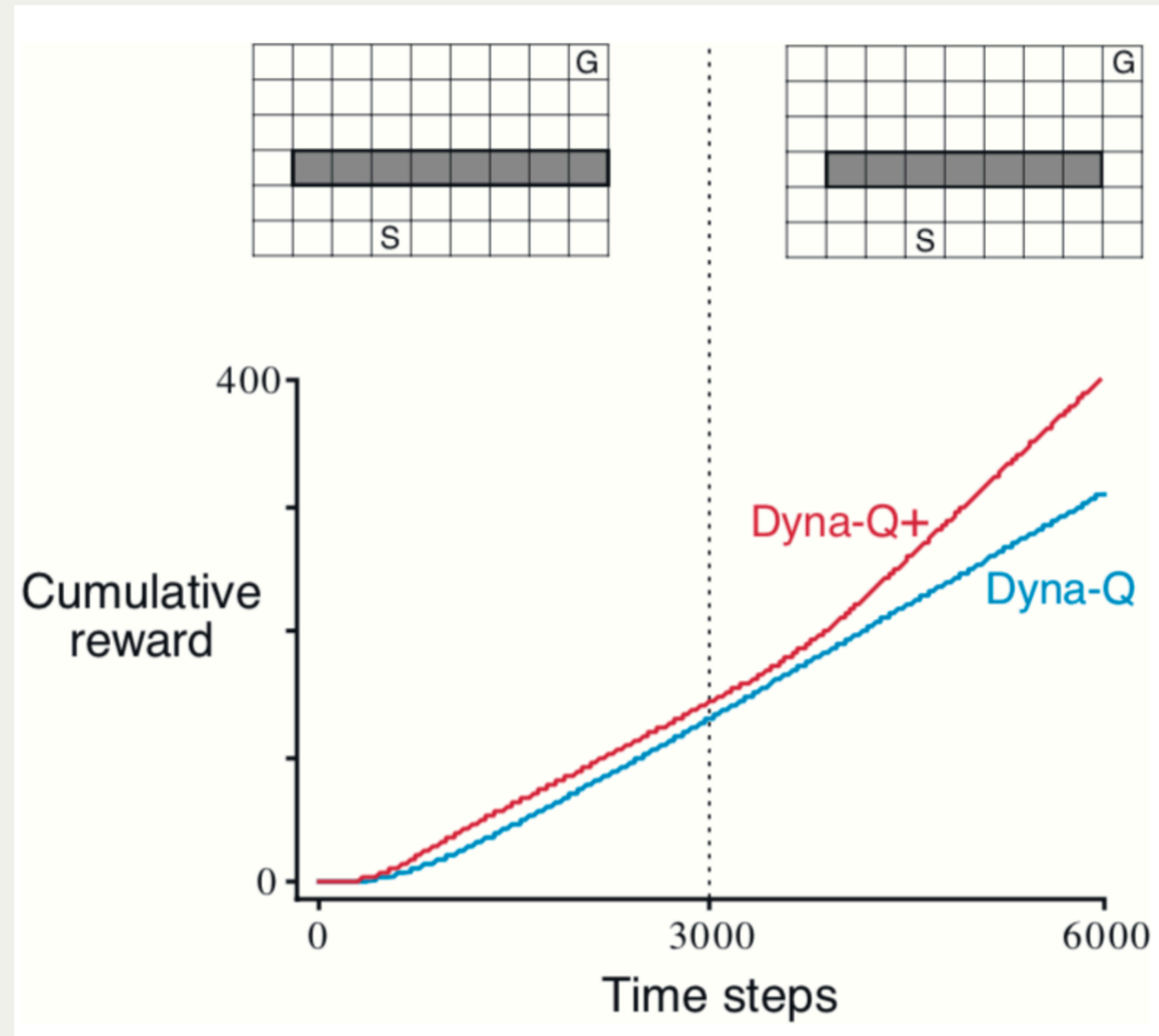
$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + k\sqrt{\tau(S, A)} + \gamma\max_a Q(S', a) - Q(S, A)]$$

- Allow actions that had never been tried from a state to be considered in planning (initial model was that such an action led back to the same state with a reward of 0)

When the Model is Wrong: Optimistic Model



When the Model is Wrong: Pessimistic Model



Prioritized Sweeping (Det. Env.)

Initialize $Q(s, a)$ and $Model(s, a)$ for all a, s , $PQueue$ to empty
Loop forever:

1. $S \leftarrow$ current (nonterminal) state; 2. $A \leftarrow \text{policy}(S, Q)$
3. Take action A ; observe resultant reward R and state S'
4. $Model(S, A) \leftarrow (R, S')$ (assumes deterministic environment)
5. $P = R + \gamma \max_a Q(S', a) - Q(S, A)$
6. If $P > \Theta$, insert (S, A) into $PQueue$ with priority P
7. Loop repeat n times while $PQueue$ is not empty
 - a. $S, A = \text{first}(PQueue); R, S' \leftarrow Model(S, A)$
 - b. $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
 - c. Loop for all \bar{S}, \bar{A} predicted to lead to S :
 - i. $\bar{R} = \text{pred. reward for } \bar{S}, \bar{A}, S$
 - ii. $P = \hat{R} + \gamma \max_a Q(S', a) - Q(\bar{S}, \bar{A})$
 - iii. If $P > \Theta$, insert (\bar{S}, \bar{A}) into $PQueue$ with priority P