

RL: Lecture 18

Harvey Mudd College

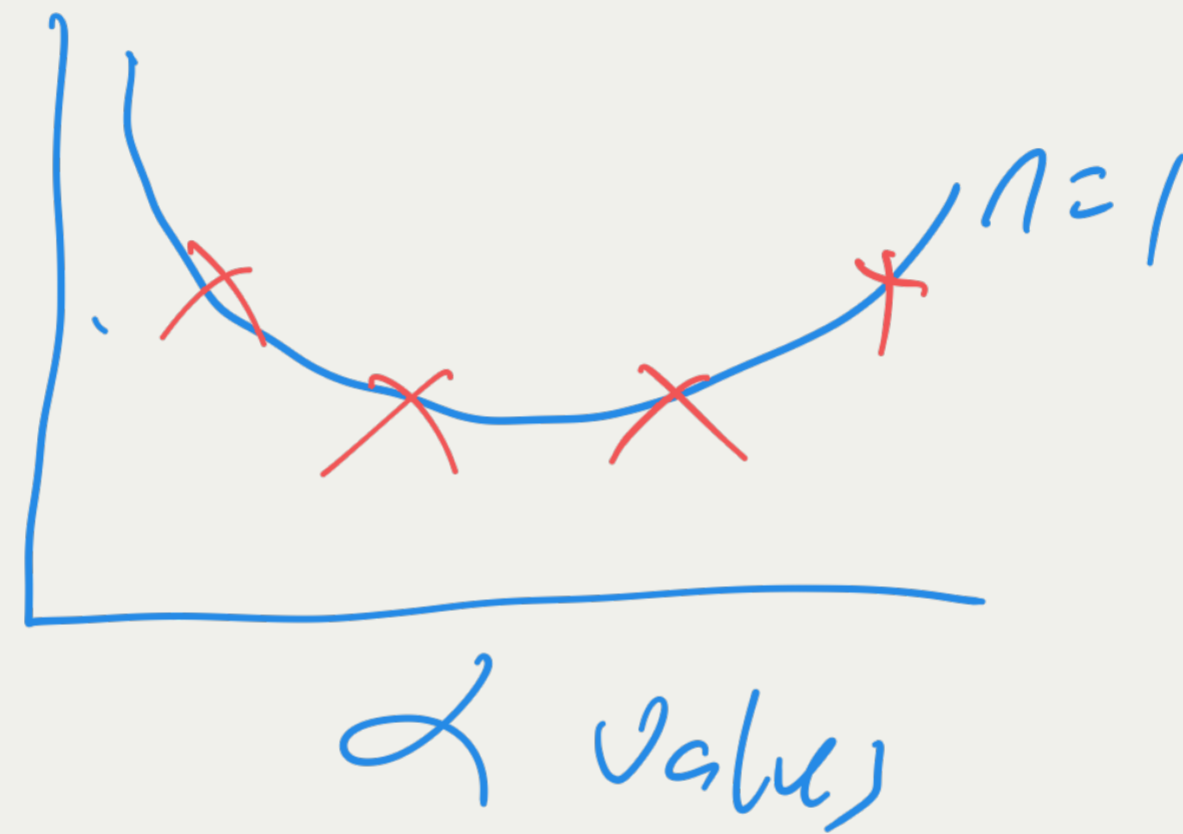
April 6, 2020

Neil Rhodes

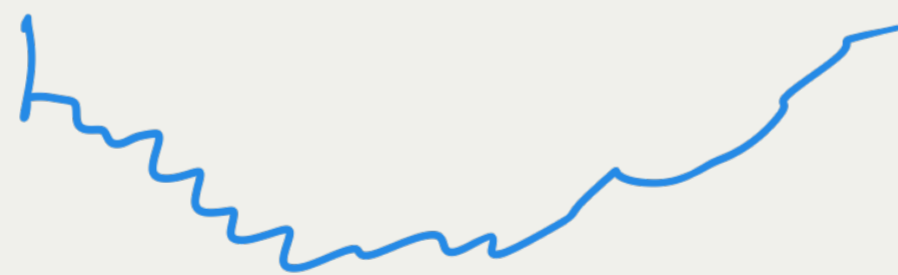
PA 5

Quiz 8 available
today
from 1-11:59?
~~PE~~ PDT?

PRNG
pseudo random
number generator

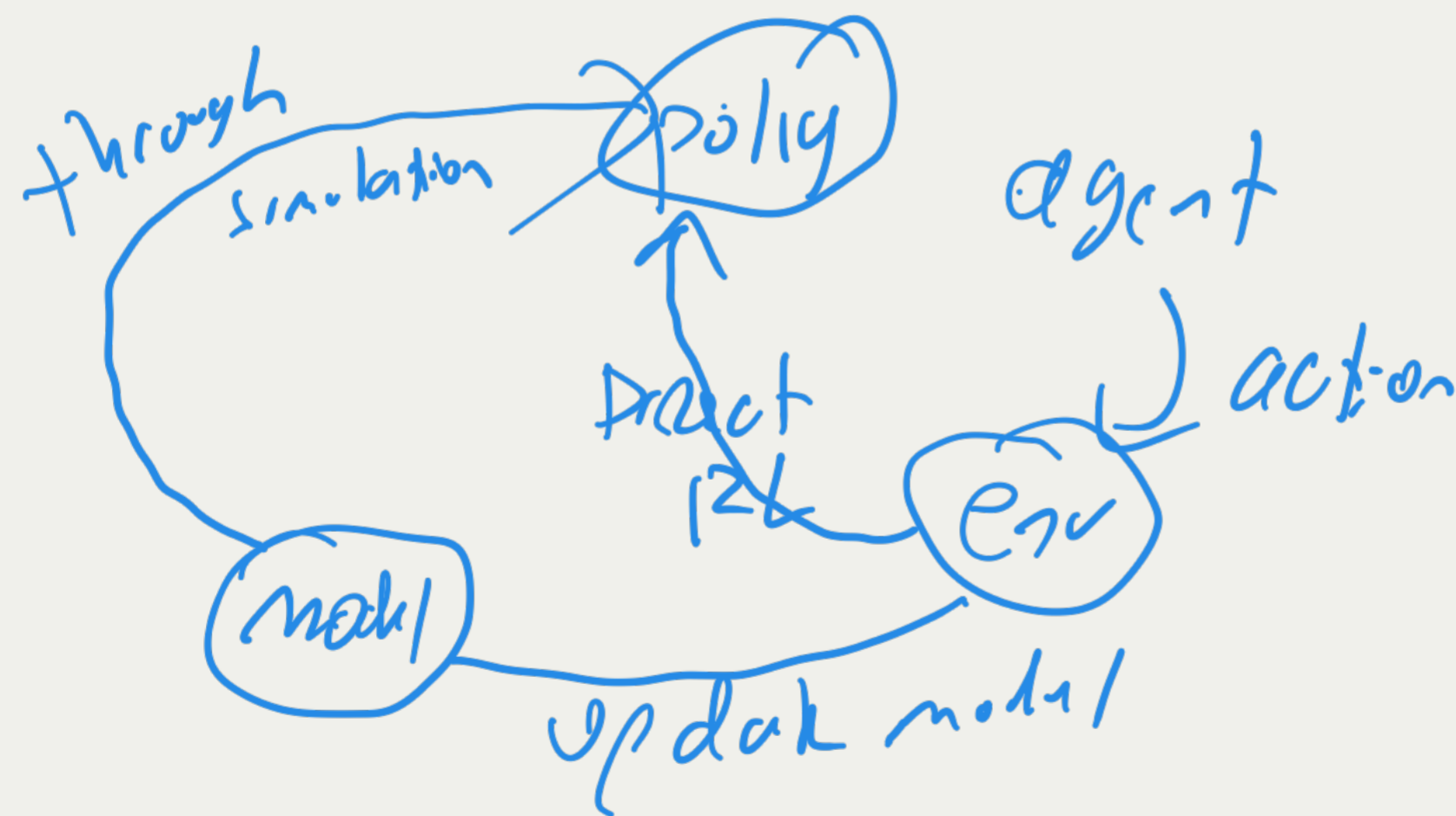


same
episode



What is planning?

Given a model, updating policy w/o interacting w/ env



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. $A \leftarrow \epsilon - \text{greedy}(S, Q)$

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

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

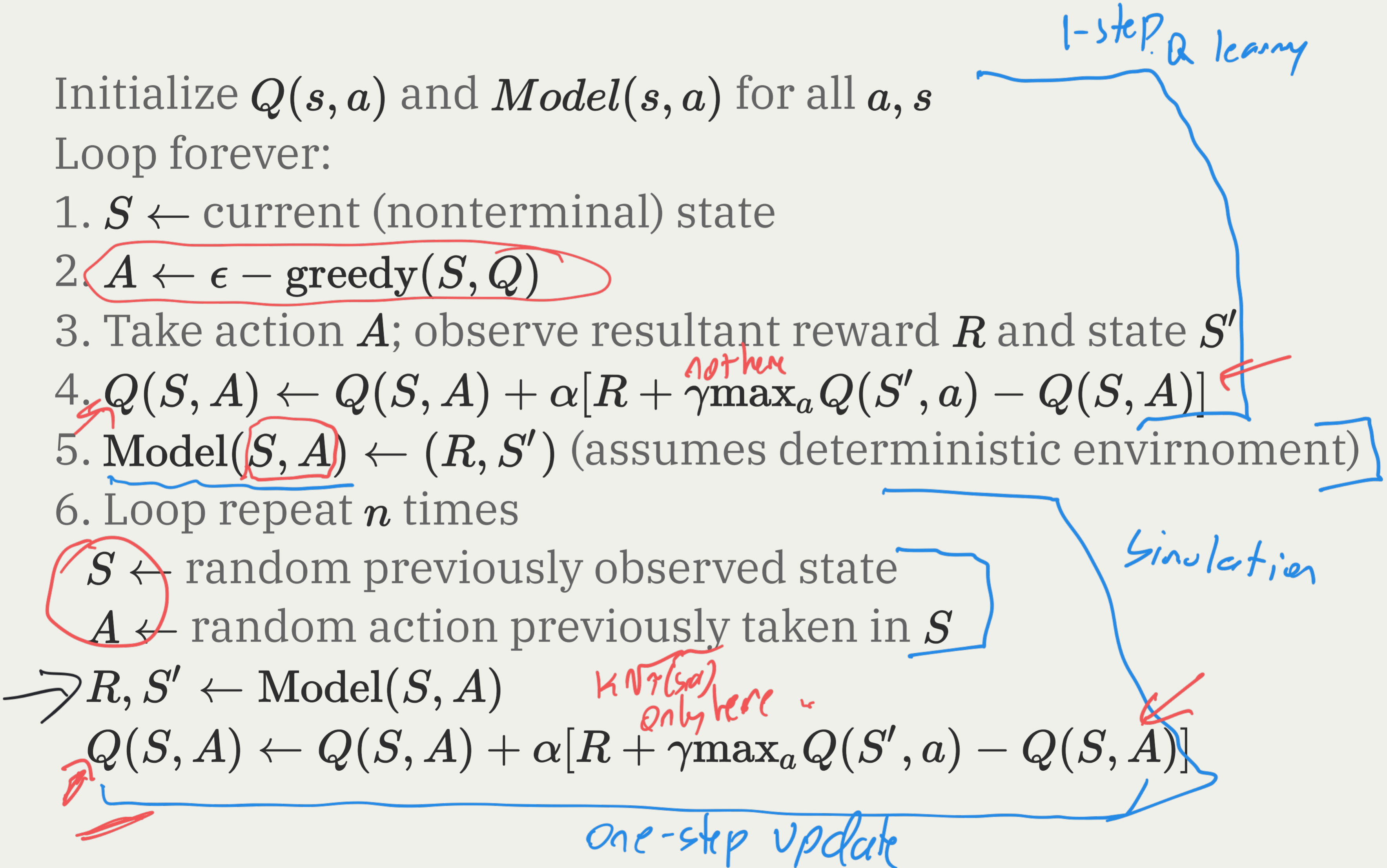
$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

1-step Q learning

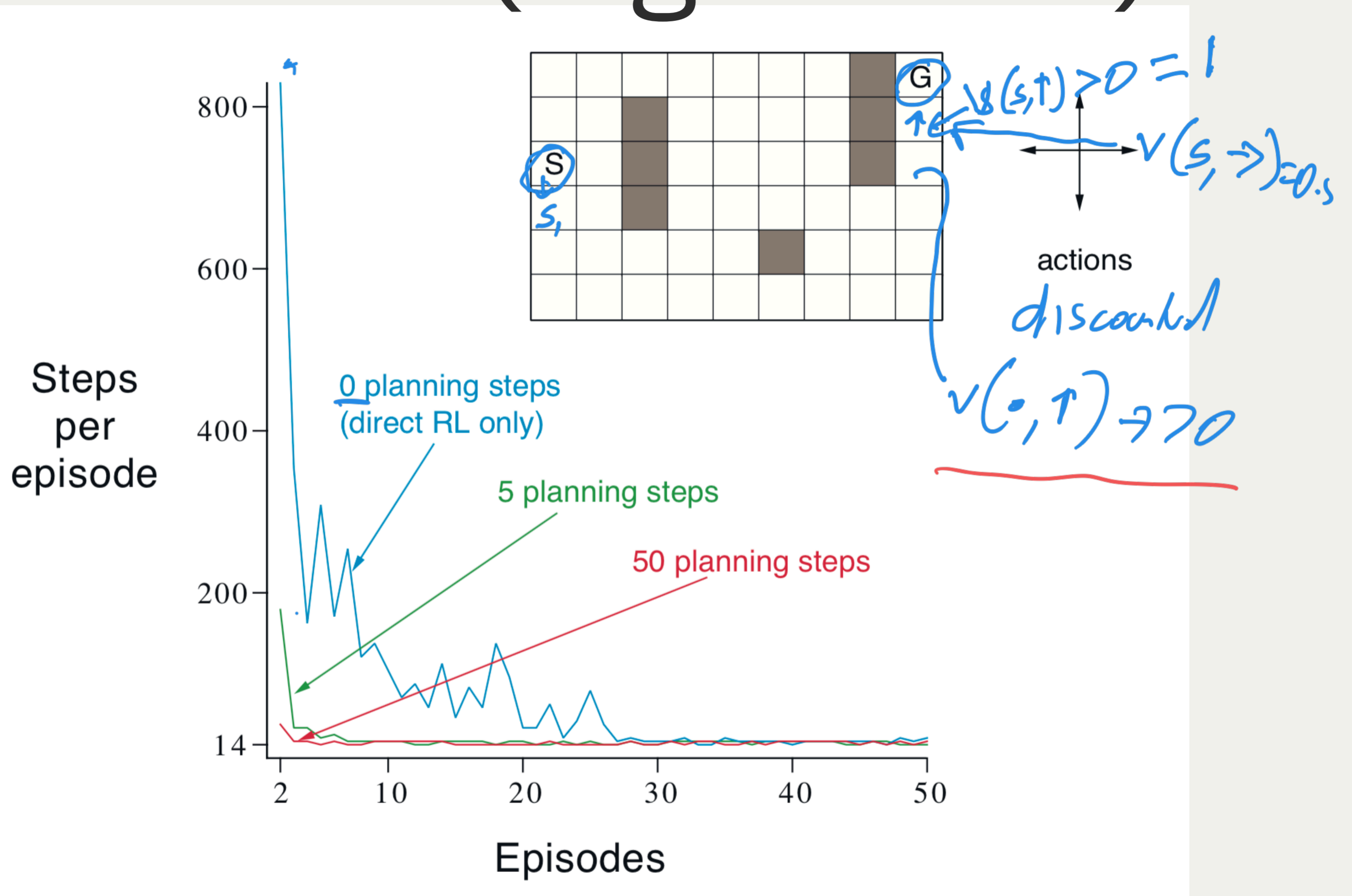
Simulation

known $\gamma(s,a)$ only here

one-step update

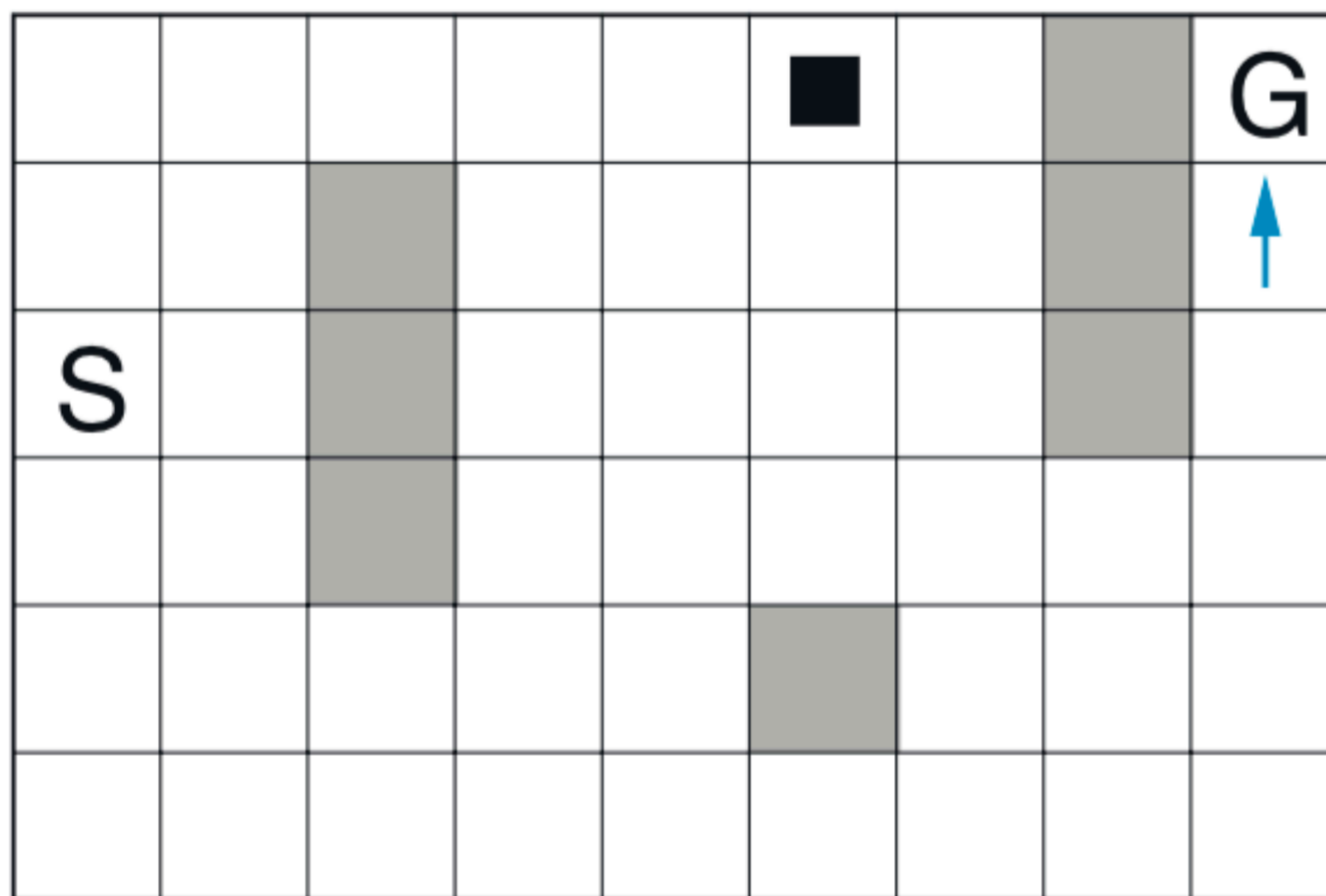


Dyna Maze (Figure 8.2)



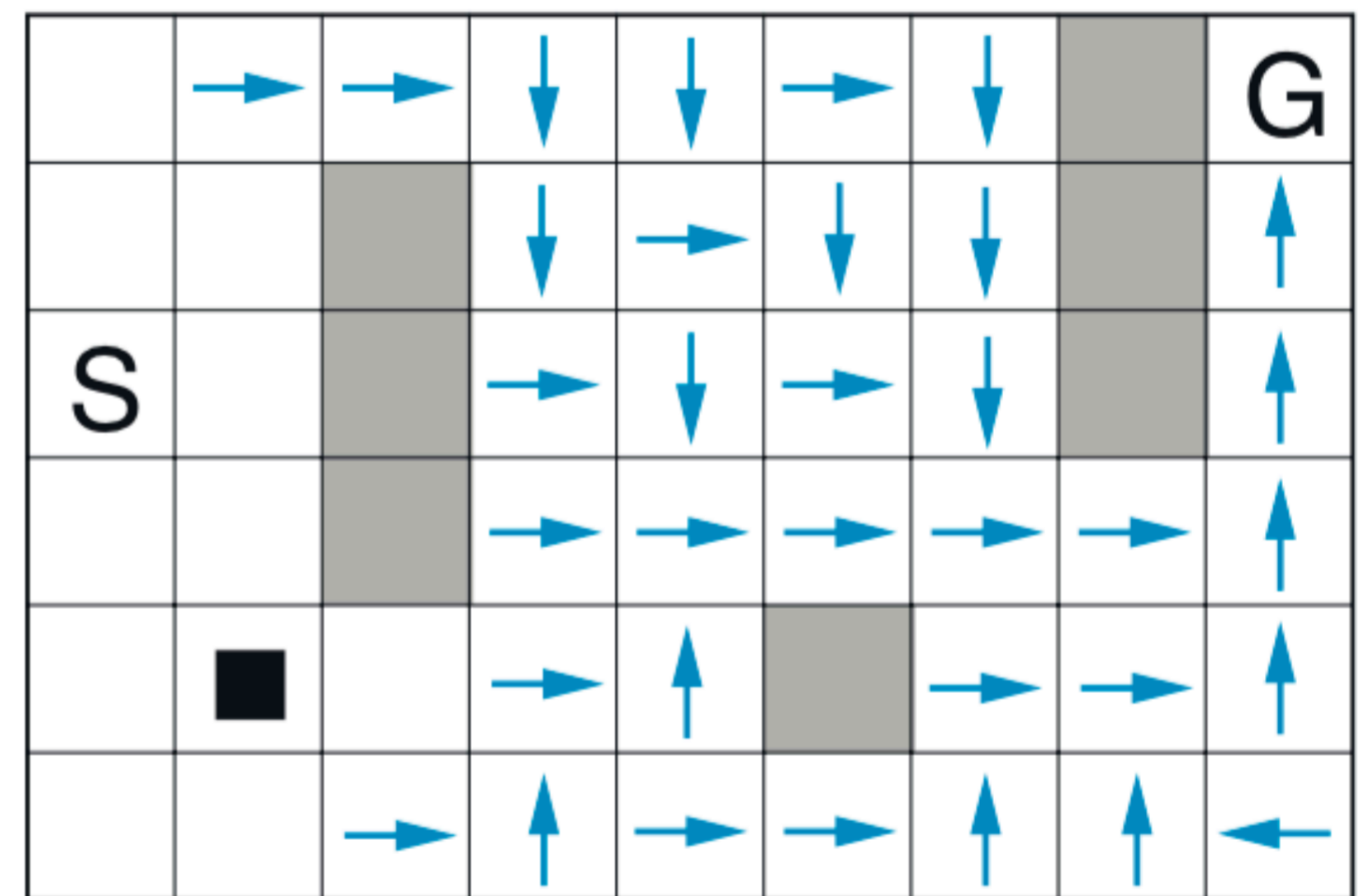
Dyna Maze (Figure 8.3)

WITHOUT PLANNING ($n=0$)



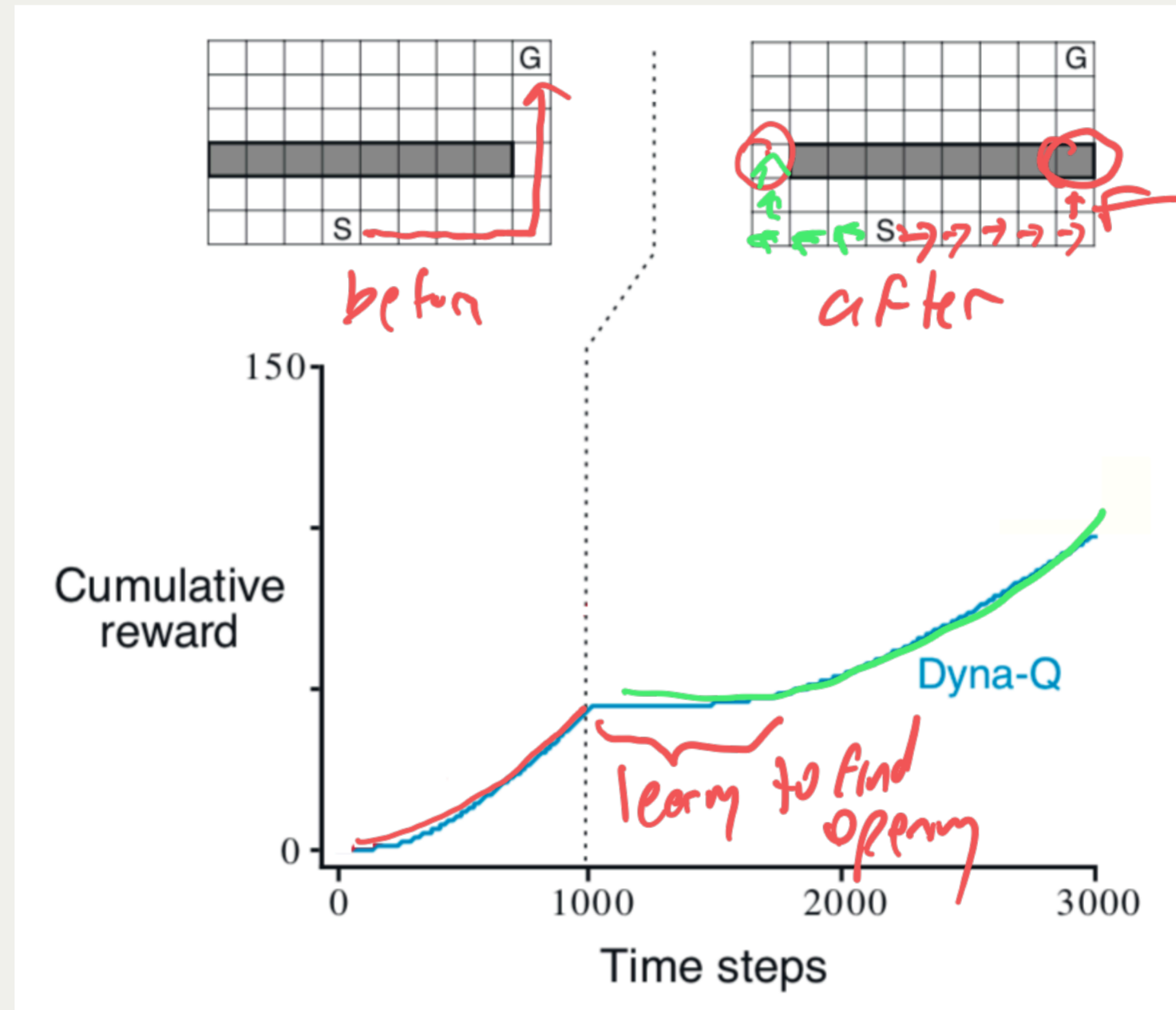
after ≈ 400 steps

WITH PLANNING ($n=50$)



about
after ≈ 17 steps

When the Model is Wrong: Optimistic Model



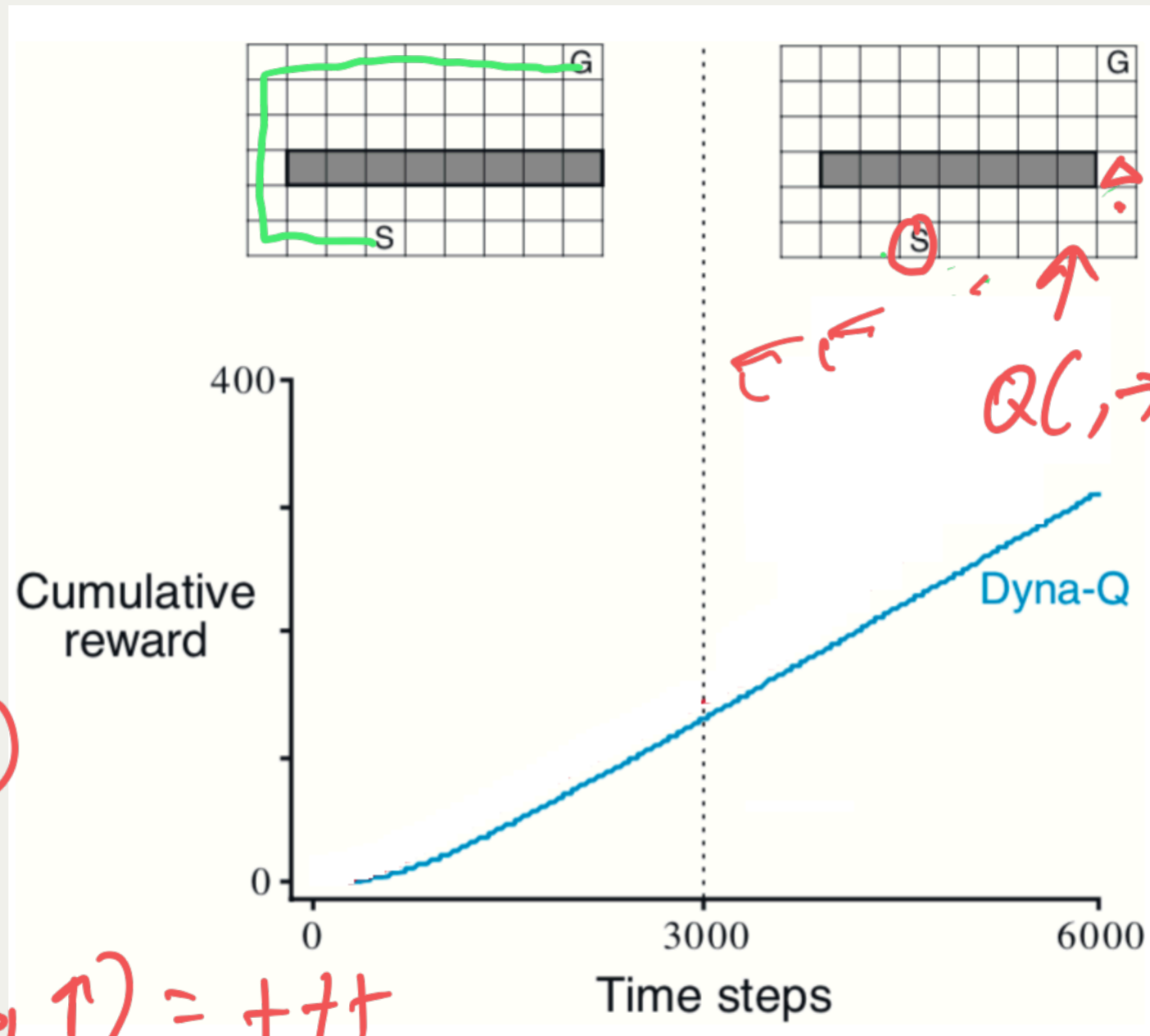
changing environment

$Q(\cdot, \pi)$ reduce fris
 $\text{model}(\cdot, \pi) = (\cdot, 0)$

planning

$$\begin{aligned}
 Q(\cdot, \pi) &+ = \alpha \left[0 + \gamma \max_a Q(\cdot, a) \right] \\
 &+ = \alpha \left[0 + \gamma Q(\cdot, \pi) - Q(\cdot, \pi) \right] \\
 &= \alpha \left[(\gamma - 1) Q(\cdot, \pi) \right]
 \end{aligned}$$

When the Model is Wrong: Pessimistic Model



$Q(S, \downarrow) \rightarrow -$
model remain (S, \downarrow)

$Q(\cdot, \uparrow) \rightarrow$
model change to (A, \uparrow)
 ϵ^7

$Q(A, \uparrow) = + + +$

(A, a)

In real world, only way to do that is a bunch of exploratory steps

$Q(\cdot, \uparrow) \rightarrow +$

$Q(S, \uparrow) \rightarrow +$

has $Q(\cdot, \uparrow) \rightarrow +$

model $Q(S, \downarrow) \rightarrow +$

$\epsilon \epsilon \epsilon$
 $Q(\cdot, \uparrow)$

[model $(\cdot, \uparrow) =$
 (\cdot, \downarrow)

in order to change

is to try

(\cdot, \uparrow) in real-world

encourage explore stale parts of the model

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

- Provide an implicit reward to exploring stale transitions

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

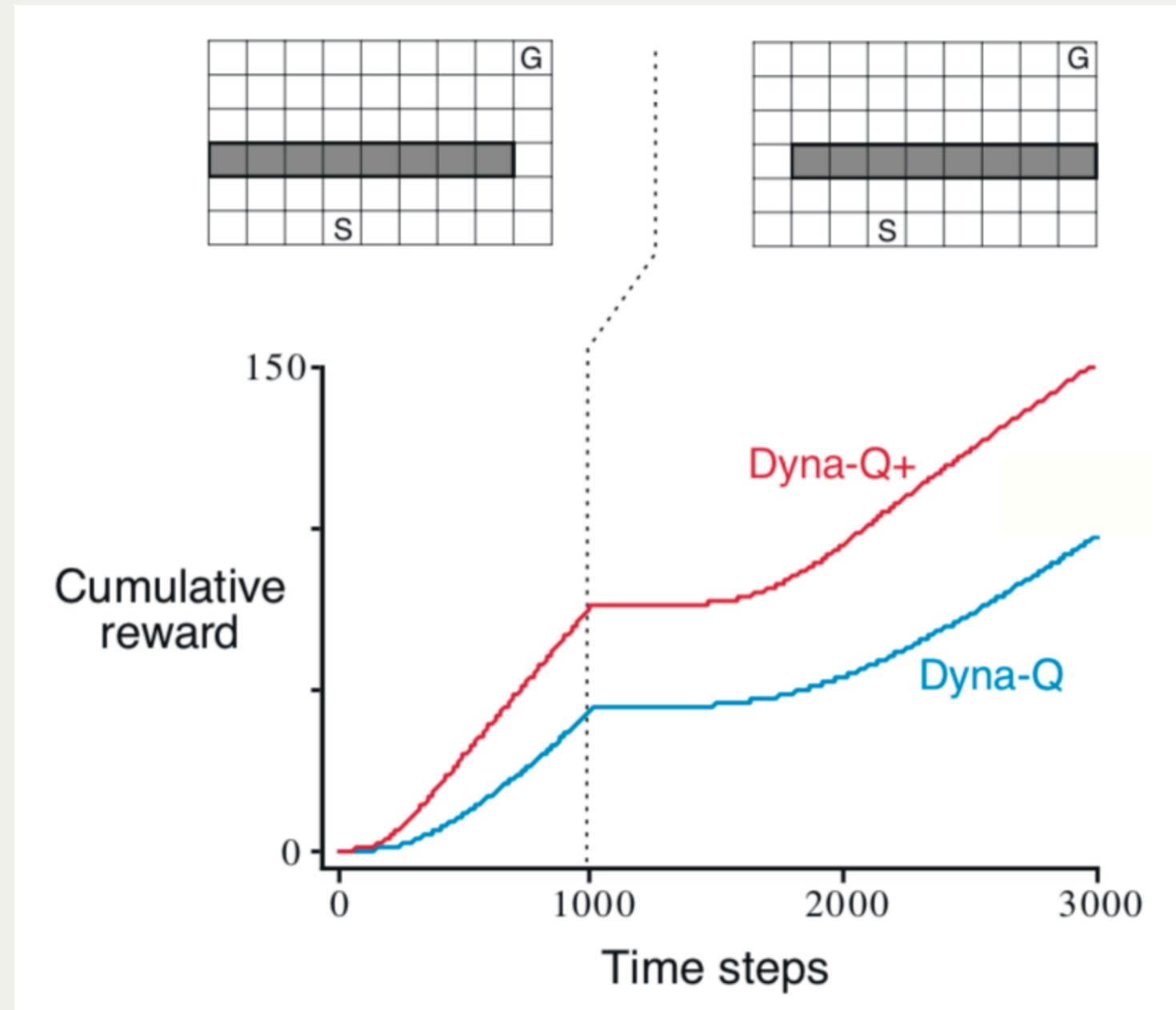
only in planning
based on model

- Allows actions that have never been tried from a state to be considered in planning (initial model: such an action leads back to the same state with a reward of 0)

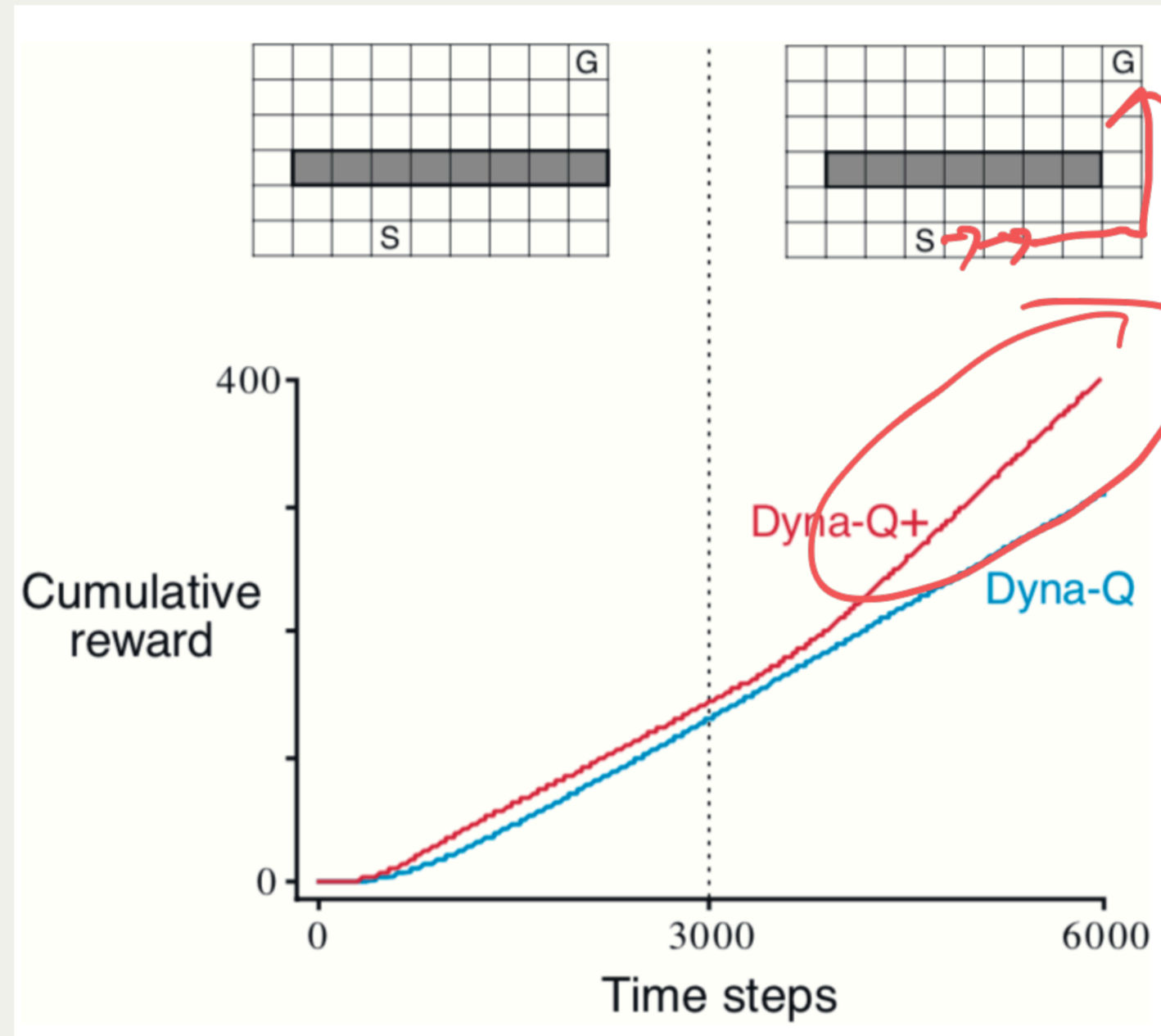
$\kappa = \text{constant}$

$\tau(S, A)$ measures how long since we tried S, A in real world

When the Model is Wrong: Optimistic Model



When the Model is Wrong: Pessimistic Model



Dyna-Q tries all state-action pairs uniformly.

Is there a better way?

$$\text{model}(\bullet, \tau) = (\bullet, 0)$$

@ beginning of second episode

how many non-zero

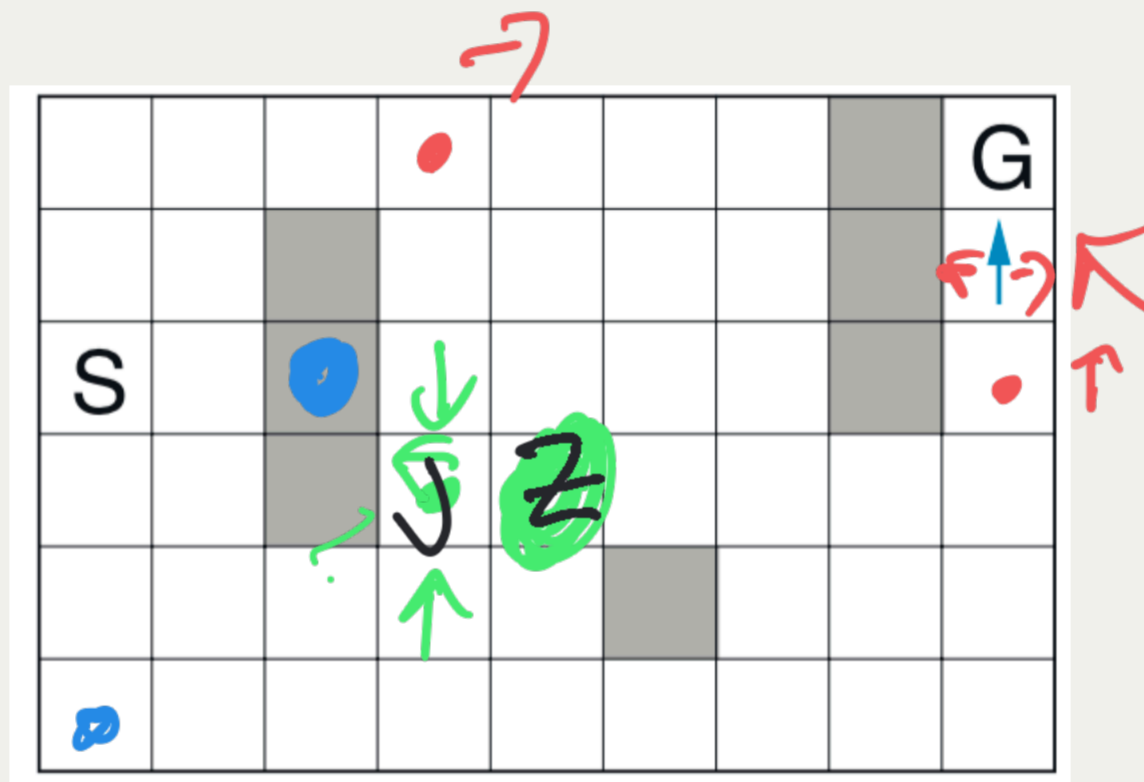
$$Q(\text{state}, \text{action})$$

one

↓ in our model

$$(z, \uparrow) \rightarrow (\cdot, \emptyset)$$

$Q(s', A')$ changes from 0.5 to 1.5 large



$Q(s, A)$ changes from 0.5623 to 0.5624

based on how changed small

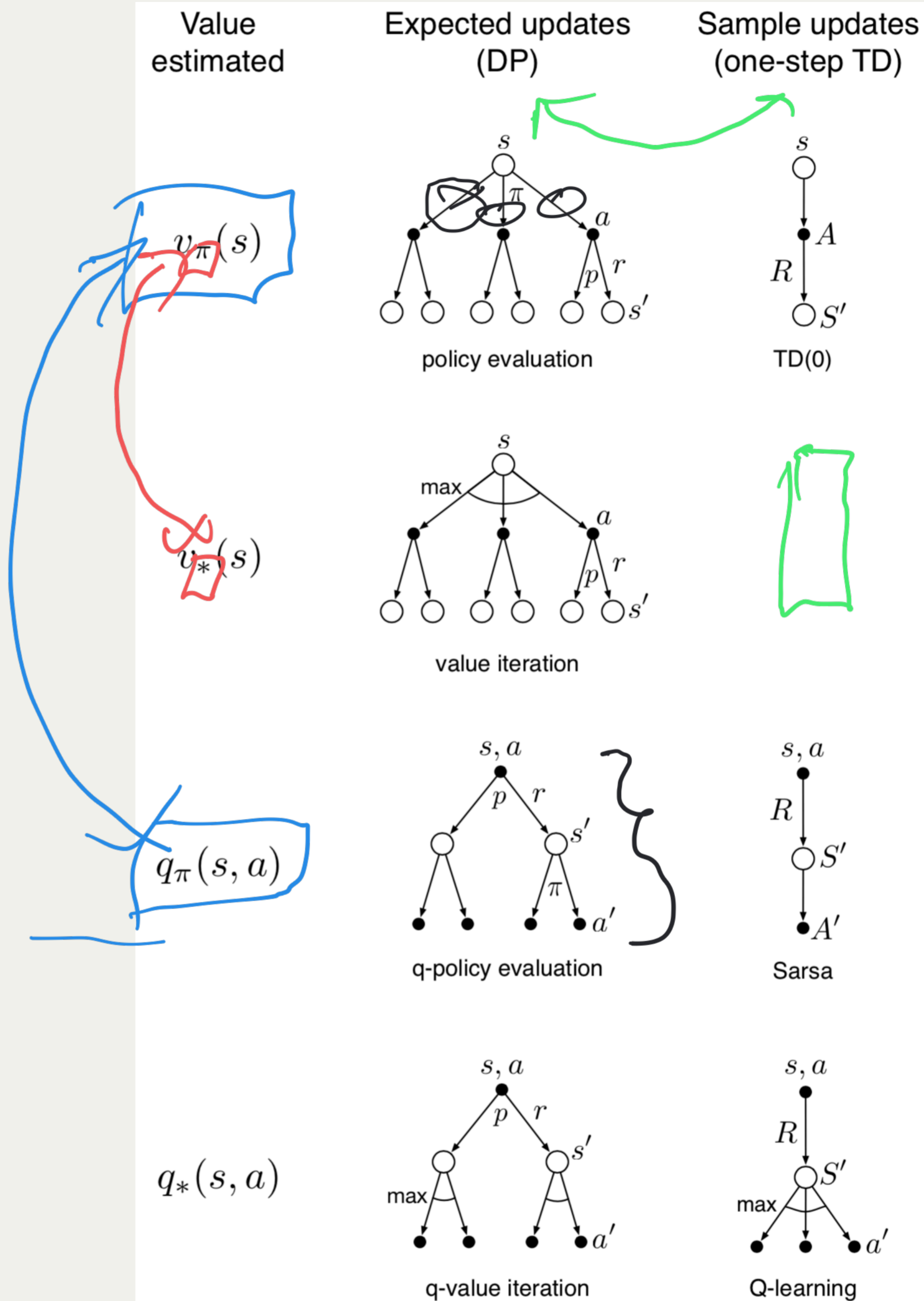
Prioritized Sweeping (Det. Env.)

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

1. $S \leftarrow$ current (nonterminal) state; 2. $A \leftarrow$ 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. $P = R + \gamma \max_a Q(S', a) - Q(S, A)$ Δ in change of Q
6. If $P > \Theta$, insert (S, A) into $PQueue$ with priority P $\Theta = \text{min}$
viable
change
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 *From the model*
 - ii. $P = \bar{R} + \gamma \max_a Q(S', a) - Q(\bar{S}, \bar{A})$ ΔQ would be
 - iii. If $P > \Theta$, insert (\bar{S}, \bar{A}) into $PQueue$ with priority P

Dimensions

- Update state/action values
- Optimal vs. arbitrary policy
- Expected vs. sample updates



Expected updates > sample updates?

Advantage of expected updates

low variance
 need to know policy probabilities

higher α
 could converge quicker

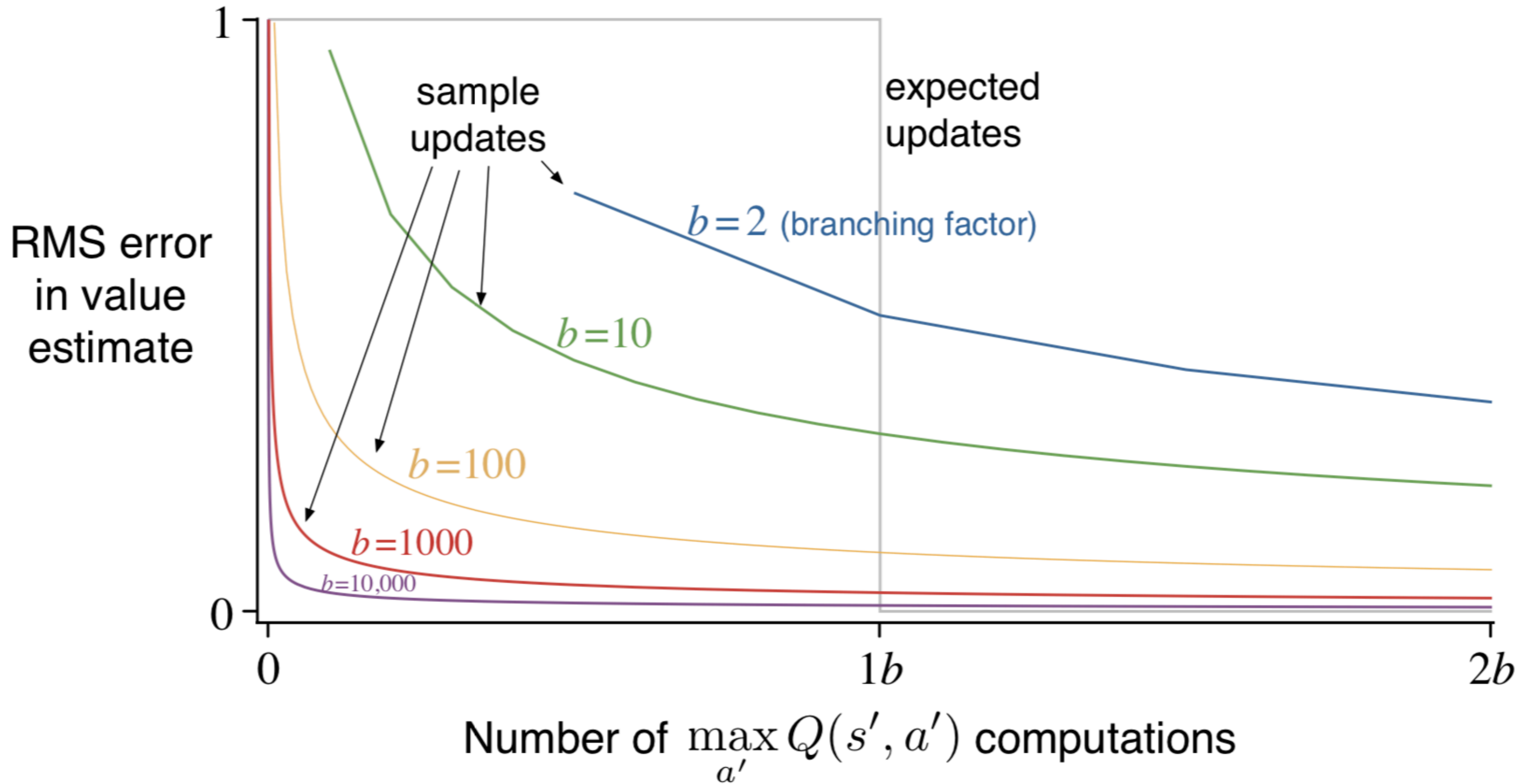


Sample Once

$$\max_a \underbrace{Q(s, a)}_{\text{expected}}$$

only pick action

Expected updates > sample updates?



Trajectory Sampling

Sampling:

- Uniform
- According to on-policy distribution

Trajectory Sampling

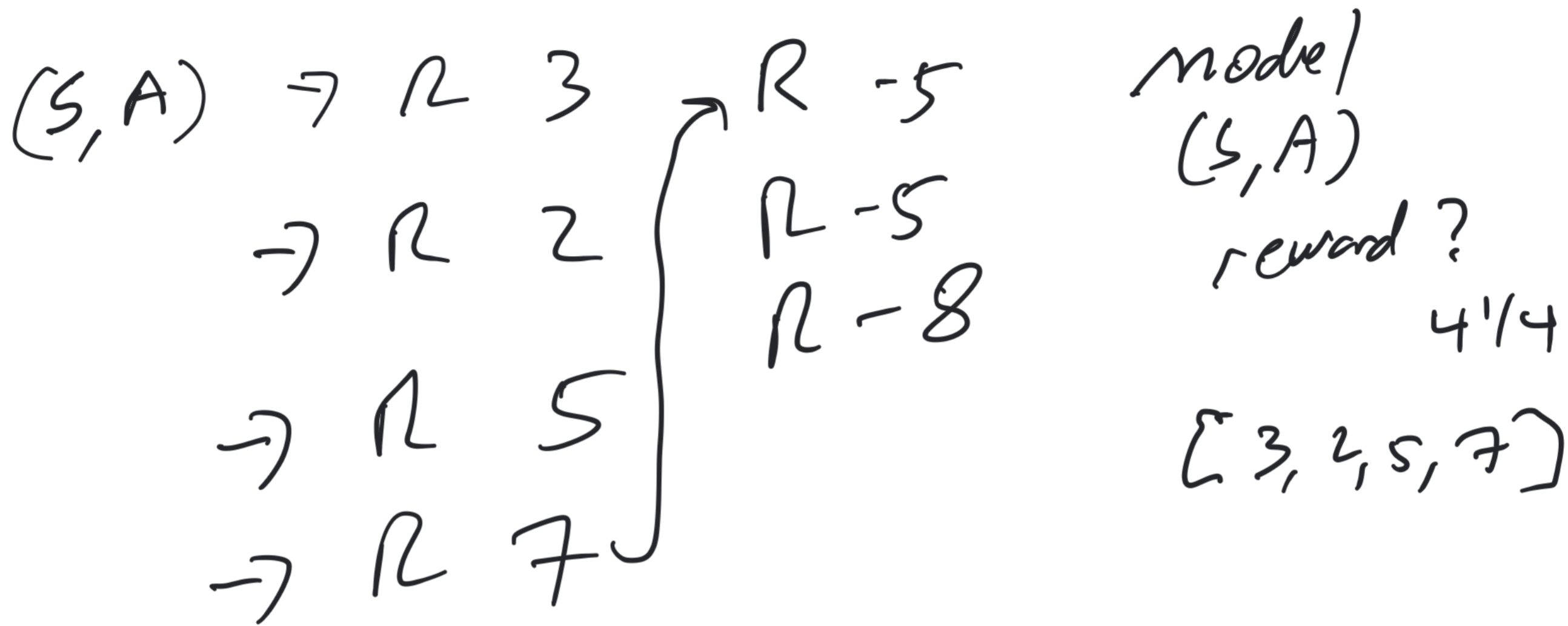
Sampling:

- Uniform
- According to on-policy distribution

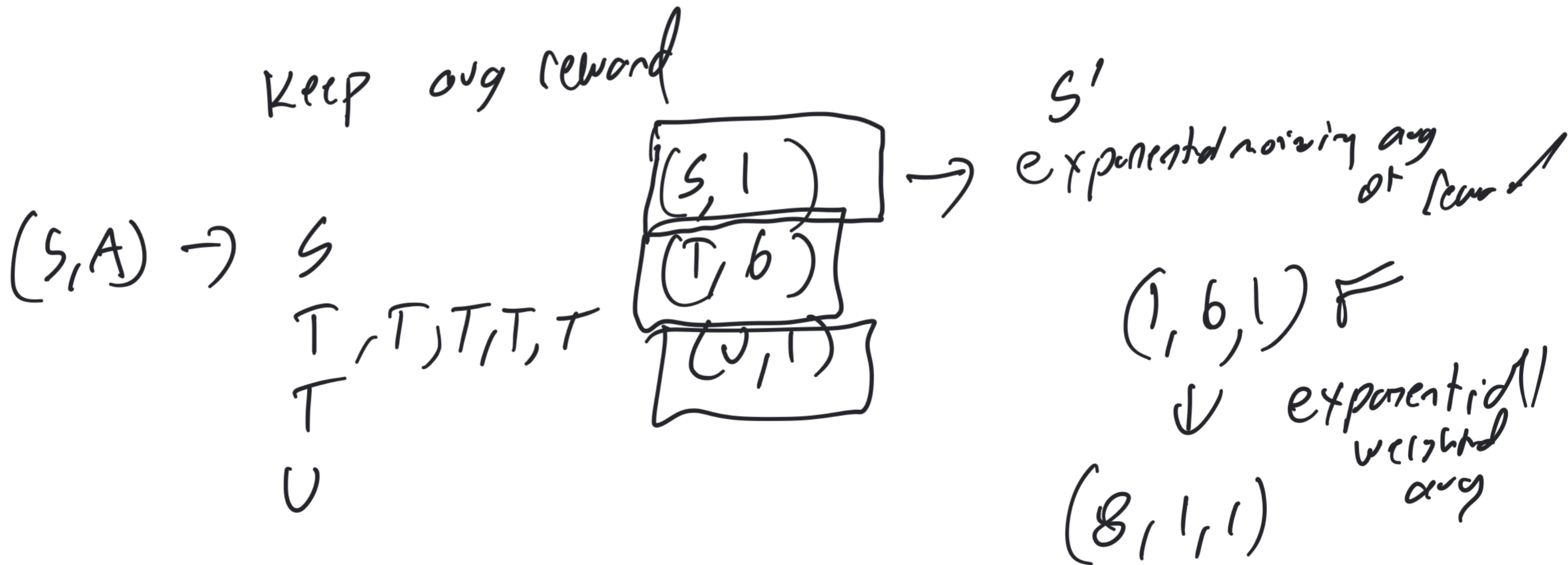
When to plan

- In the background
- At decision time

Assum s' & R are independent
 \uparrow



Keep avg reward



~~x_1~~
 ~~x_2~~
 ~~x_3~~

$t-1$
 $\left[\begin{array}{c} R_1 \\ R_2 \\ \vdots \\ i \end{array} \right]$
 t

$t-1+1$
 $\left[\begin{array}{c} R_1 \\ R_2 \\ R_3 \\ \vdots \\ i \end{array} \right]$
 $t+1$