# Reinforcement Learning

# Credit Assignment Problem

Previous action/rewards

| |
|---|
| |
| |
| |
| a-2-r-2 |
| a-1-r-1 |
| a0-r0 |
| a1-r1 |
| a2-r2 |
| a3-r3 |
| a4-r4 |

Who is responsible for reward r1. This action (a1) or previous actions

$\gamma$    action return

$\gamma$    action return

*Sum of discounted rewards that come after a reward at a step helps evaluate action impact*
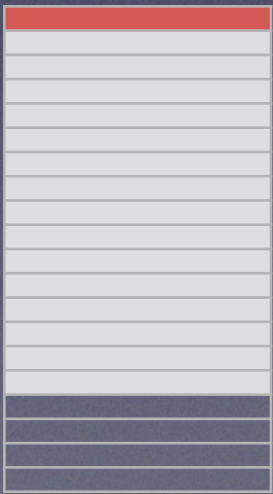
$\gamma$ **action return** $= \Sigma \; Discounted \; Rewards$

$\gamma \approx 0$ future rewards have little impact; how does this action affect immediate rewards
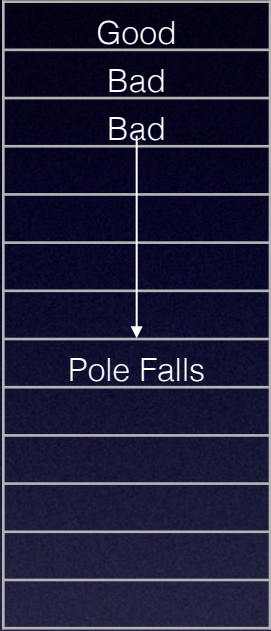
Cart-Pole actions have short term effects.

Why?

Bad or good actions may undone by immediate future actions

$\gamma \approx 1$ future rewards have large effect; how does this action affect future rewards

# Credit Assignment Problem

* Good actor in bad movie
* Top runner losing race
* Secret drive route flooded with traffic

| Good |
| --- |
| Bad |
| Bad |
| |
| |
| |
| |
| Pole Falls |
| |
| |
| |
| |
| |

A **good** action may followed by **bad** actions, but this will be a low occurrence if many episodes are run.

Good actions are more likely to be followed by good actions.

Run many action episodes -> solve action returns for each step in each episode step -> Standardization on each action return is called action advantage

(-) Action Advantage (action will result in bad future)
(+) Action Advantage (action will result in good future)

Multiply gradient by action advantage

(+ Action Advantage Mult) The gradients will reduce loss. Makes action move likely in future.
(- Action Advantage Mult)  The gradients will not reduce loss. Makes action less likely in future.

Note:
Negative action advantages shorten the number of steps in an episode( falling pole ).
Thus, making negative action advantages smaller than positive action advantages,
This produces larger action-advantage and gradient products ( training will learn to fit positive actions )

# Policy Gradients

Episode

step

| A | B | C | D | E | F | G | H | I | J |

$Discount_a = Reward_a + \gamma * Discount_b$

$Discount_b = Reward_b + \gamma * Discount_c$
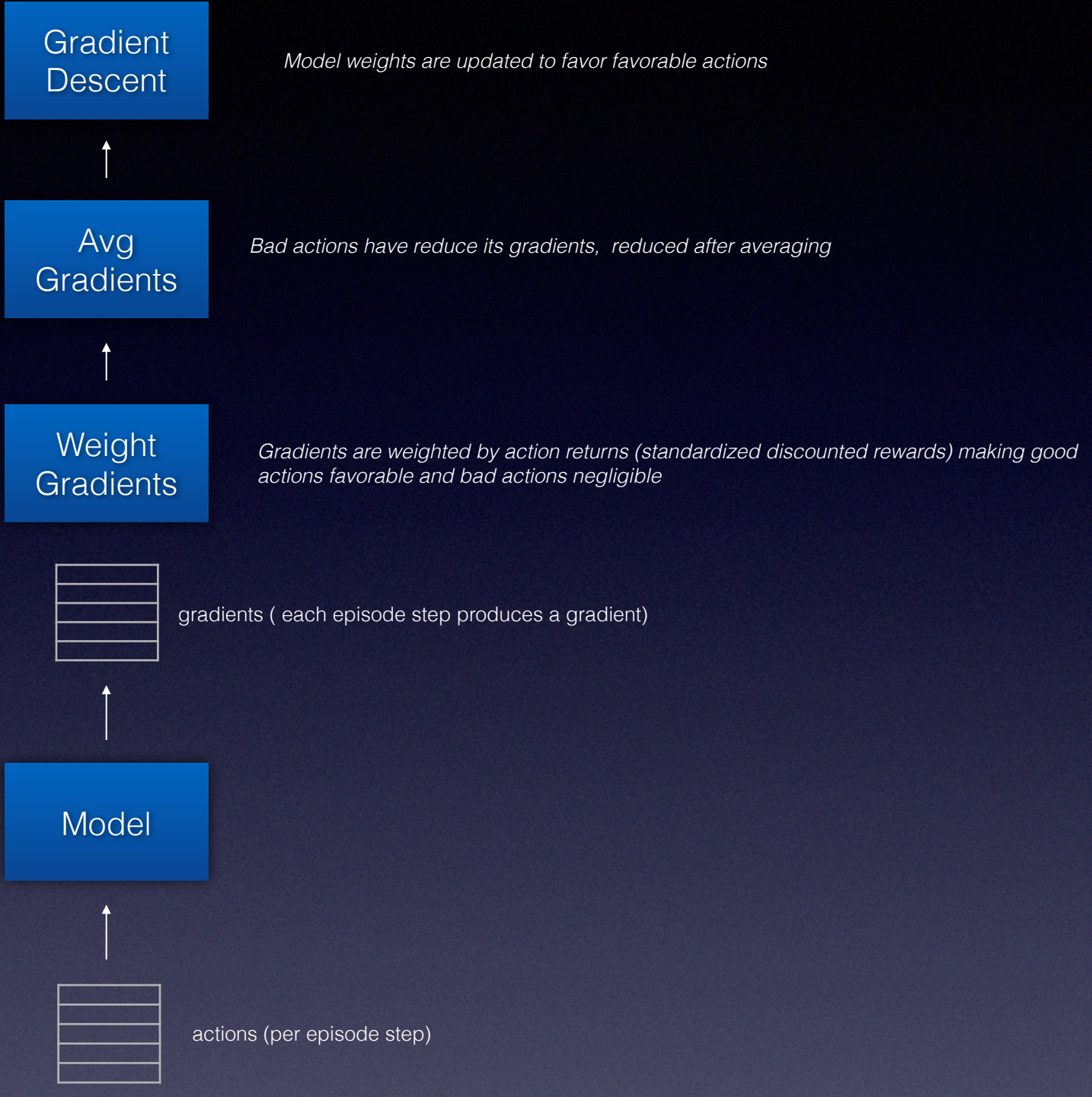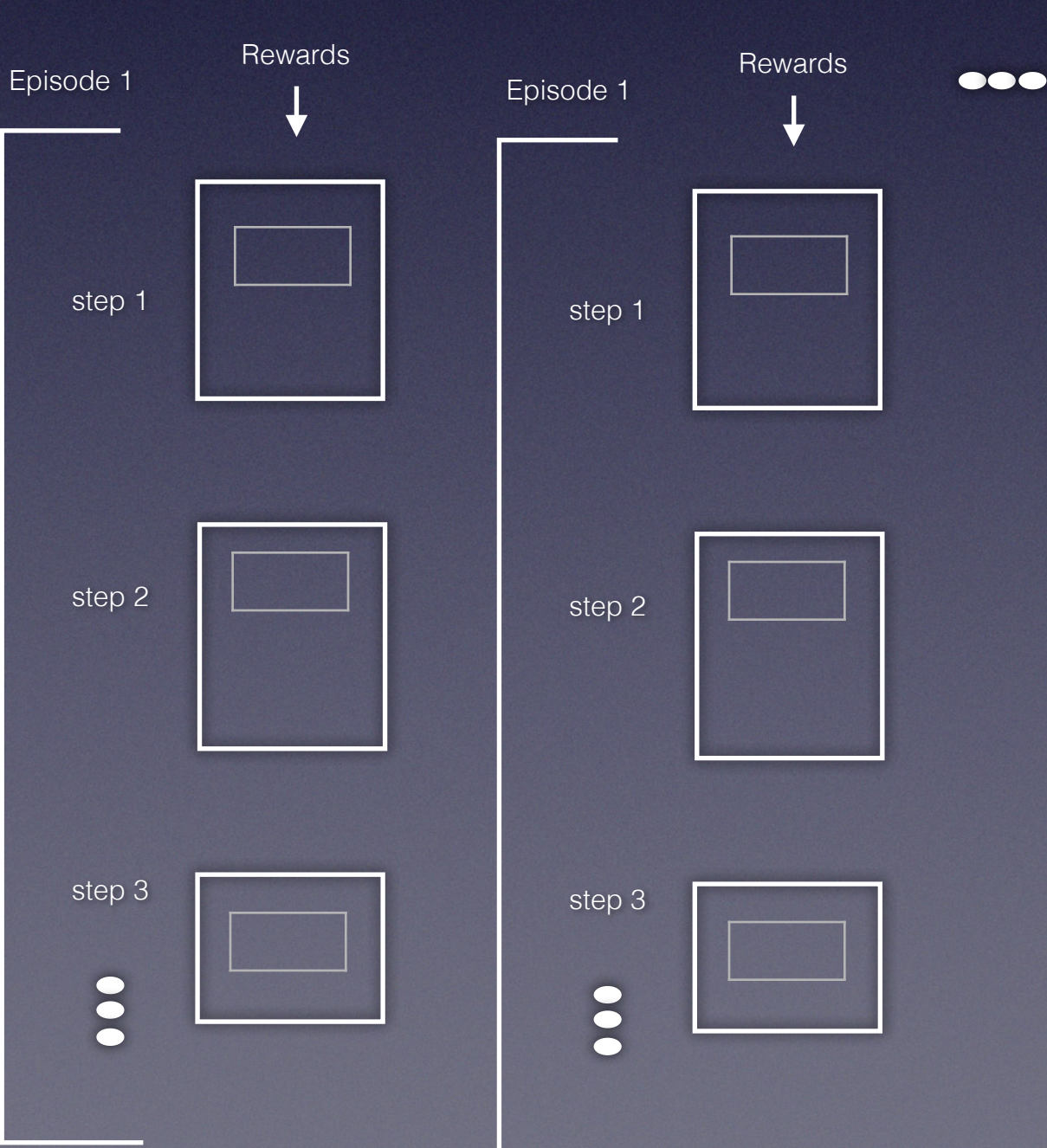
Used to normalize each discounted reward to produce
**how much better or worse an action is compared to other actions on average**

Episode

step

| A | B | C | D | E | F | G | H |

$Discount_a = Reward_a + \gamma * Discount_b$

$Discount_b = Reward_b + \gamma * Discount_c$

**Mean**

**Std**

Episode

step

| A | B | C | D | E | F | G | H | I | J | K | L | M |

$Discount_a = Reward_a + \gamma * Discount_b$

$Discount_b = Reward_b + \gamma * Discount_c$

**Gradient Descent**

Model weights are updated to favor favorable actions

**Avg Gradients**

Bad actions have reduce its gradients, reduced after averaging

**Weight Gradients**

Gradients are weighted by action returns (standardized discounted rewards) making good actions favorable and bad actions negligible

gradients ( each episode step produces a gradient)

**Model**

actions (per episode step)

# Cart Pole Gradient

Episode 1

Episode 2

Episode 1

Action advantage

Episode 2

Action advantage ●●●

① Multiply each gradient by final reward

●●●

**Episode 1**

step 1

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 1

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

**Episode 2**

step 1

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 2

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 2

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 2

step 3

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 3

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 3

| Kernel (1) |
| Kernel (5,1) |
| Kernel (5) |
| Kernel (4, 5) |

step 3

Note: In this policy a large final reward evaluates into a longer running episode(i..e game)

Avg gradients used for gradient descent

| Avg Kernel (1) |
| Avg Kernel (5,1) |
| Avg Kernel (5) |
| Avg Kernel (4, 5) |

# Cart Pole

Some episode

$\gamma = 0.8$

$\rightarrow$

reward

$\rightarrow$

reward

$\rightarrow$

reward

10

0

-50

$10 + \gamma \times 0 + \gamma^2 \times (-50)$

RETURN ACTION
(ADVANTAGE)

$0 + \gamma \times (-50)$

$-50$

# Cart Pole Policy

Random Sample

Estimate:
Probability of moving left (action 0)?

Note:
If left prob small —> model will estimate to move right
If left prob large —> model will estimate to move left

$x_1$

$Kernel = (1)$

$\sigma$
$\Sigma$

$Kernel = (5,1)$

$x_1$          $x_2$          $x_3$          $x_4$          $x_5$          $Kernel = (5)$

$\sigma$        $\sigma$        $\sigma$        $\sigma$
$\Sigma$        $\Sigma$        $\Sigma$        $\Sigma$

$Kernel = (4,5)$

Observation          $x_1$          $x_2$          $x_3$          $x_4$

# Cart Pole Policy

PDF

Probability of action = 0

Obs

The policy probability distribution is changed during training to increase rewards (keeping cart standing)

# *Discounted Reward*

Discount allows us to evaluate actions. Good actions will get higher returns than bad ones on average.

Discount Factor : $\gamma = 0.8$

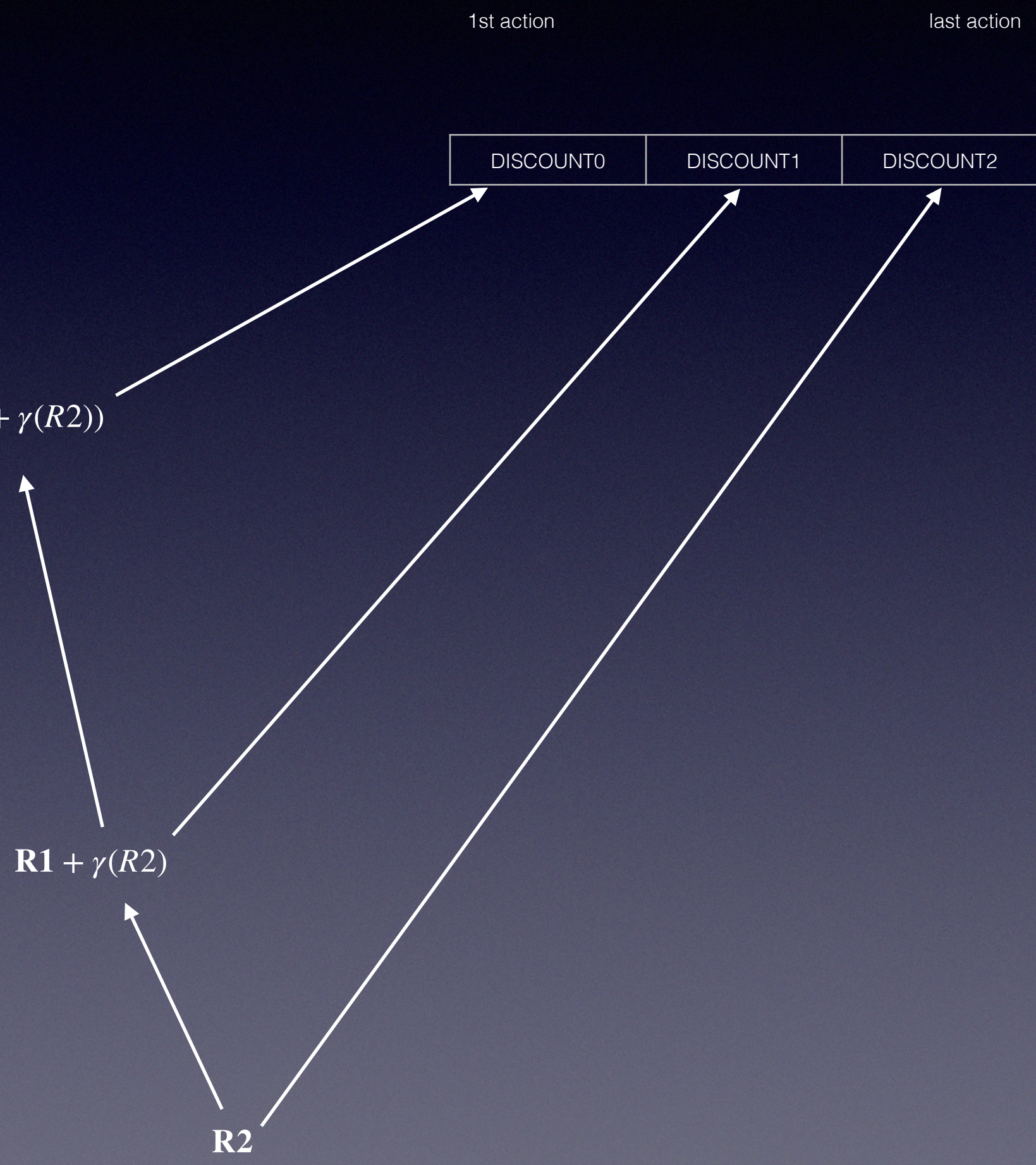| R0 | R1 | R2 |
|----|----|----|

1st action                     last action

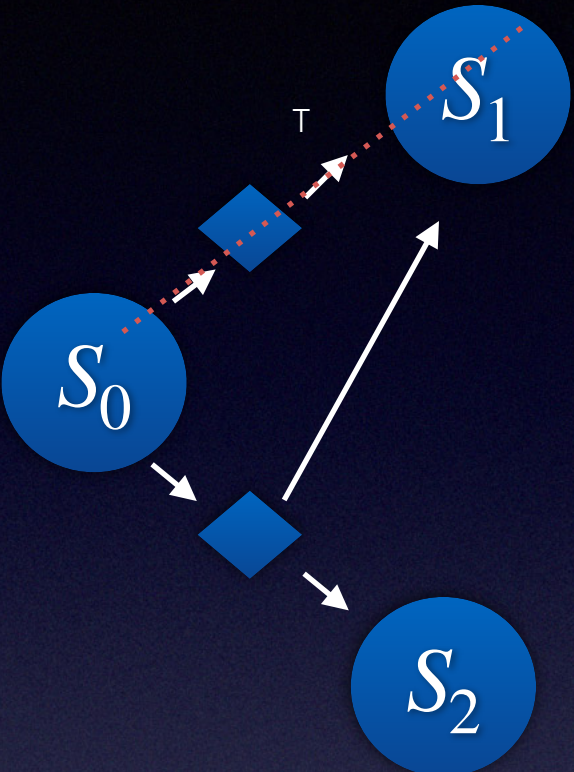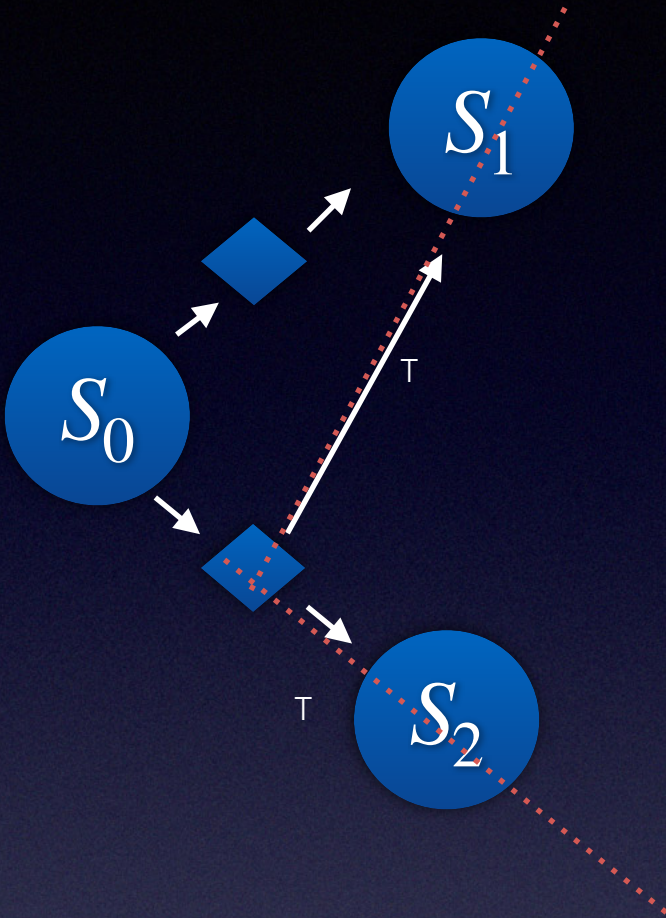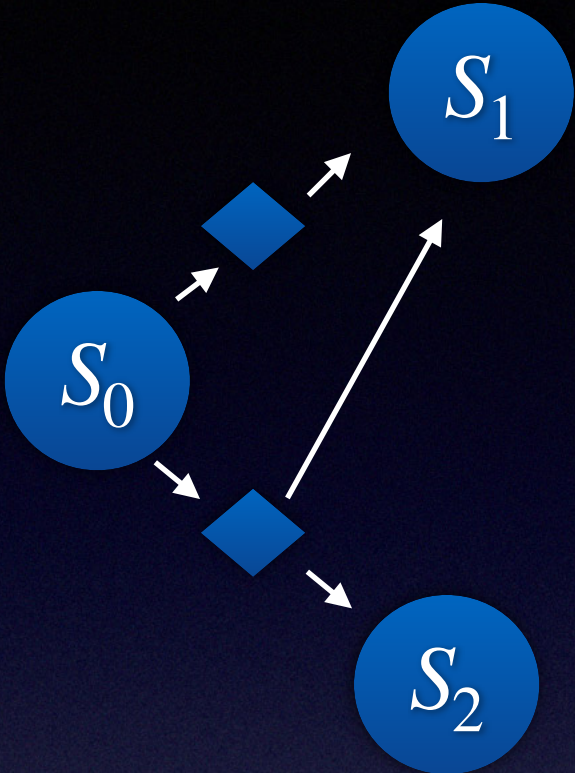| DISCOUNT0 | DISCOUNT1 | DISCOUNT2 |
|-----------|-----------|-----------|

$R0 + R1 \times \gamma^1 + R2 \times \gamma^2 = R0 + \gamma(R1 + \gamma(R2))$

$\textbf{R1} + \gamma(R2)$

$\textbf{R2}$

# Bellman Optimal



Optimal State Value (Expectation)

$$V_k*(s_0) \quad \boxed{T} \quad \boxed{R\_01} \quad \gamma \quad \boxed{V_k*(s_1)}$$
$$\boxed{T} \quad \boxed{R\_02} \quad \gamma \quad \boxed{V_k*(s_2)} \quad +$$

$$\boxed{T} \quad \boxed{R\_01} \quad \gamma \quad \boxed{V_k*(s_1)}$$
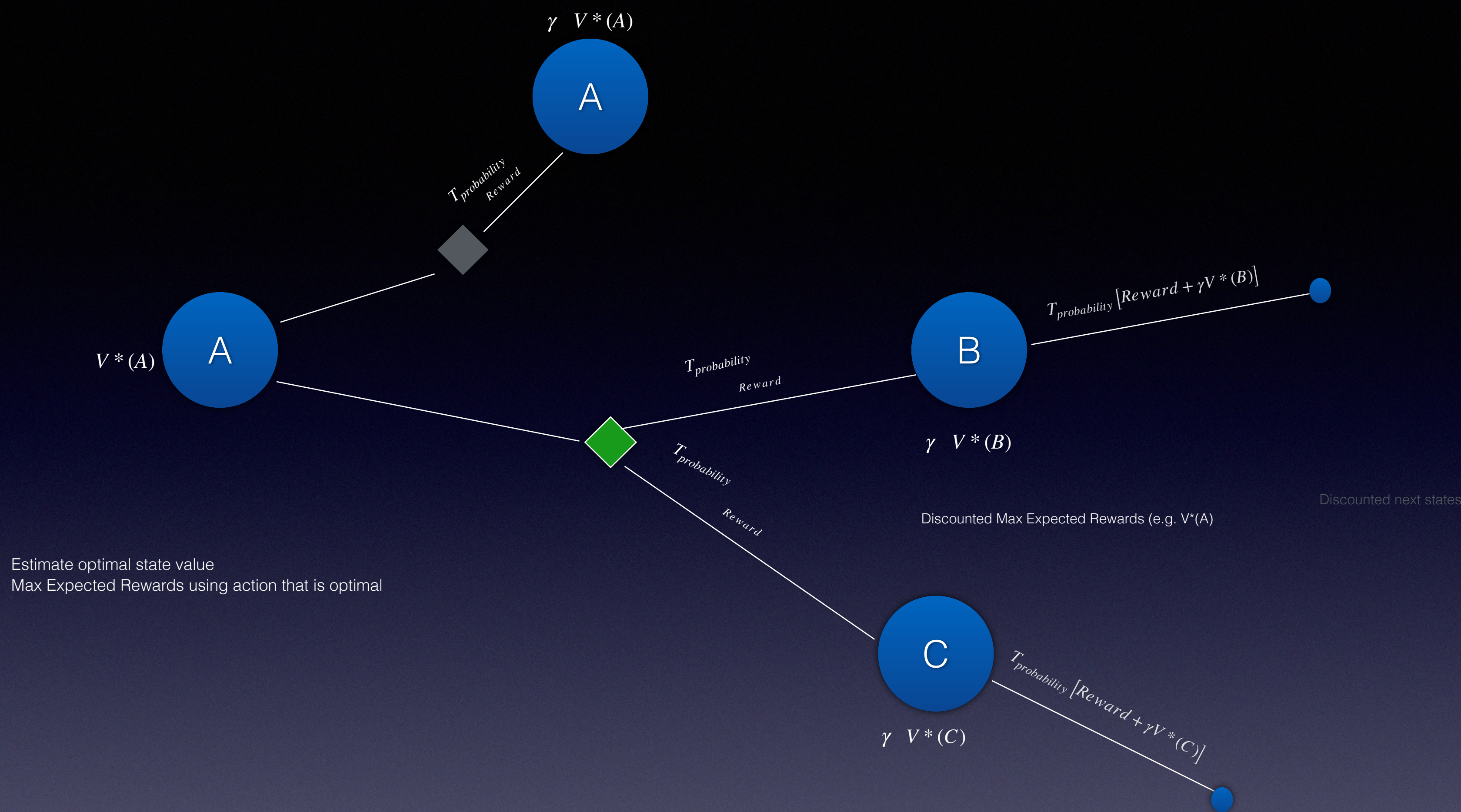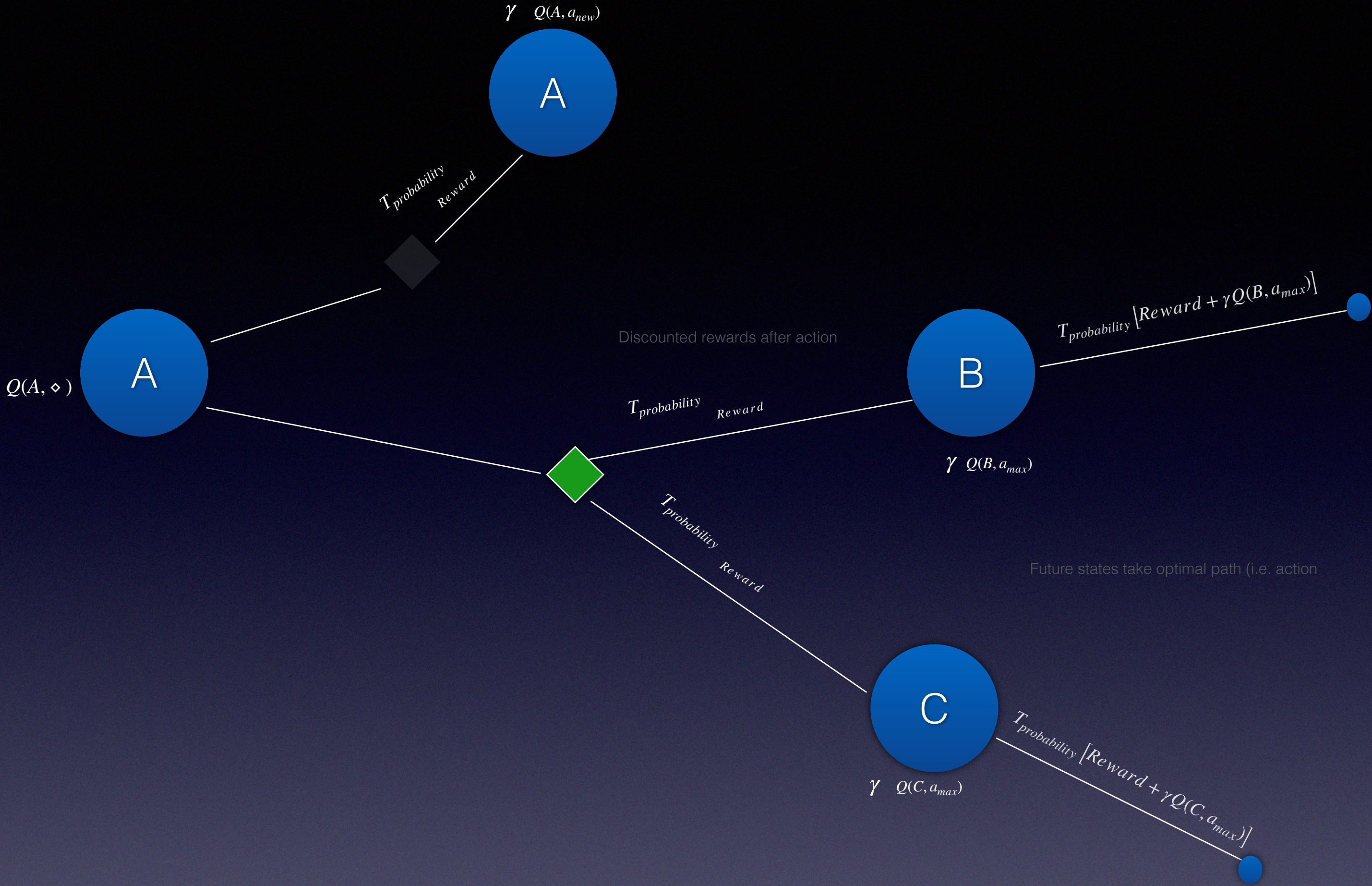
Max

Note:
T= Transition Probability
Statistical Average of an agent moving to a state after an action

# Bellman Optimal



$\gamma \ V*(A)$

A

$T_{probability}$
$Reward$

$V*(A)$

A

$T_{probability}$
$Reward$

$T_{probability}$
$Reward$

$T_{probability}\left[Reward + \gamma V*(B)\right]$

B

$\gamma \ V*(B)$

Discounted next states

Discounted Max Expected Rewards (e.g. V*(A))

Estimate optimal state value
Max Expected Rewards using action that is optimal

C

$\gamma \ V*(C)$

$T_{probability}\left[Reward + \gamma V*(C)\right]$

# Quality Value



$\gamma \quad Q(A, a_{new})$

A

$T_{probability}$  $Reward$

A

$Q(A, \diamond)$

Discounted rewards after action

$T_{probability}$  $Reward$

B

$T_{probability} \left[ Reward + \gamma Q(B, a_{max}) \right]$

$\gamma \quad Q(B, a_{max})$

$T_{probability}$  $Reward$

Future states take optimal path (i.e. action

C

$\gamma \quad Q(C, a_{max})$

$T_{probability} \left[ Reward + \gamma Q(C, a_{max}) \right]$

# Quality Value



Optimal Policy for agent (Expectation)

$$Q(s_0, b) \quad\boxed{\text{T}}\quad \boxed{\text{R\_01}} \quad \gamma \quad \boxed{Q(s_1, max_{action})}$$

$$\boxed{\text{T}}\quad \boxed{\text{R\_02}} \quad \gamma \quad \boxed{Q(s_2, max_{action})} \quad +$$

$$Q(s_0, a) \quad\boxed{\text{T}}\quad \boxed{\text{R\_01}} \quad \gamma \quad \boxed{Q(s_1, max_{action})}$$

$Q_{init}(s, a)$

|       | $a$ | $b$ |
|-------|-----|-----|
| $s_0$ | 0   | 0   |
| $s_1$ | 0   | 0   |
| $s_2$ | 0   | 0   |

$Q_{next}(s, a)$

|       | $a$ | $b$ |
|-------|-----|-----|
| $s_0$ | –   | –   |
| $s_1$ |     |     |
| $s_2$ |     |     |

Random policy explores network, updates  state values (i.e. estimations) based on rewards and transitions explored.

# Algorithm

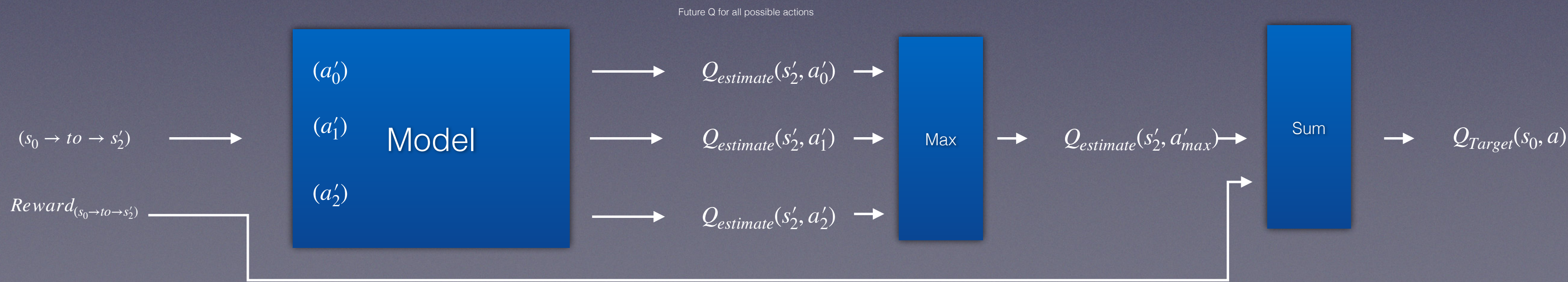sum of future discounted rewards

$$Q_{k+1}(s, a) = (1 - \alpha)Q_k(s, a) + \alpha[r + \gamma Q_k(s', a'_{max})]$$

Target Q Value (not scalable)

Approximate $Q_k(s', a'_{max}) = Q_{k\approx}(s', a'_{max})$

# Approx Algorithm

Estimate sum of future discounted rewards

$$Q_{k+1}(s, a) = (1 - \alpha)Q_k(s, a) + \alpha[r + \gamma Q_{\approx}(s'_2, a'_{max})]$$

Error ( goal is to reduce error)

$$Q_{k+1}(s, a) = Q_k(s, a) + \alpha([r + \gamma Q_{\approx}(s', a'_{max})] - Q_k(s, a))$$

learn_rate

Target Q Value

Future Q for all possible actions

$(s_0 \to to \to s'_2)$

$Reward_{(s_0 \to to \to s'_2)}$

$(a'_0)$

$(a'_1)$  Model

$(a'_2)$

$Q_{estimate}(s'_2, a'_0)$

$Q_{estimate}(s'_2, a'_1)$ → Max

$Q_{estimate}(s'_2, a'_2)$
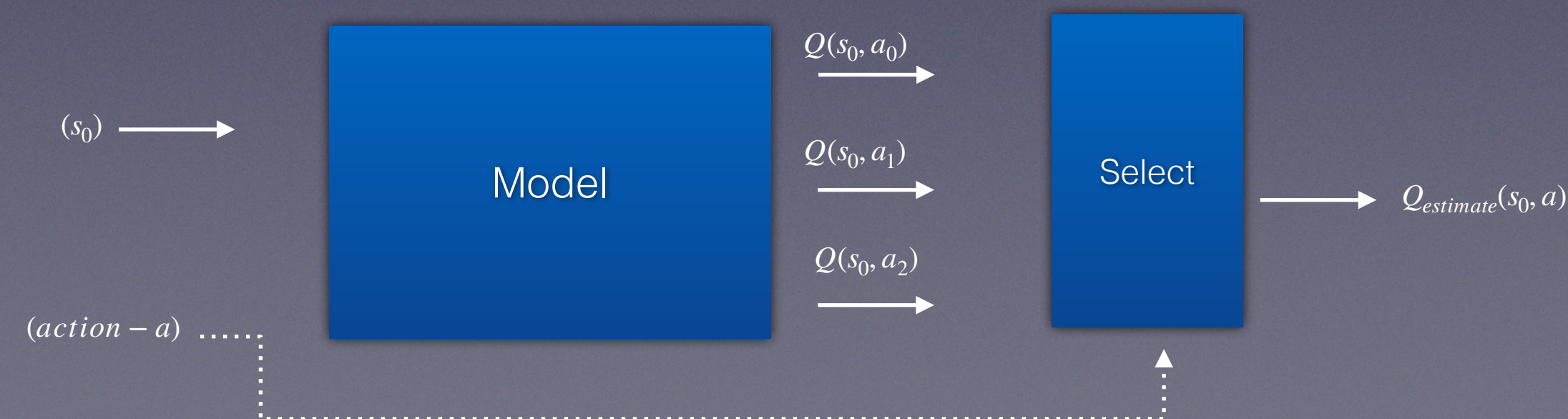
$Q_{estimate}(s'_2, a'_{max})$ → Sum

$Q_{Target}(s_0, a)$

## Algorithm

$$Q_{k+1}(s, a) = (1 - \alpha)Q_k(s, a) + \alpha[r + \gamma Q_k(s', a'_{max})]$$

Target Q Value (not scalable)

Approximate $Q_k(s', a'_{max}) = Q_{k\approx}(s'_2, a'_{max})$

$$Q_{k+1}(s, a) = (1 - \alpha)Q_k(s, a) + \alpha[r + \gamma Q_{\approx}(s'_2, a'_{max})]$$

Error ( algorithm will eventually converge when there is negligible error between Q(s,a) and QTarget(s,a). So model trains to lower error for all state-action combinations explored in network .

$$Q_{k+1}(s, a) = Q_k(s, a) + \alpha([r + \gamma Q_{\approx}(s'_2, a'_{max})] - Q_k(s, a))$$

Target Q Value

Estimated using model



$(s_0) \longrightarrow$

Model

$Q(s_0, a_0)$

$Q(s_0, a_1)$

$Q(s_0, a_2)$

Select

$Q_{estimate}(s_0, a)$

$(action - a)$

# Lunar Lander v3

# Cart Pole

state

Policy

$left\_proba$

$sample_{uniform} > left\_proba$

$action$

1 - action

$Target$

$estimate$

Loss

System actions:
1 - move cart right
0 - move cart left

# Lunar Lander v3

state

Policy

*Action Probabilities n = 4*

ArgMax → One-Hot

Target

Estimate

Categorical
CrossEntropy
Loss

Softmax(4)

Dense(32)

Dense(16)

Dense(8)

Input*8)