

# 10703 Deep Reinforcement Learning

Tom Mitchell

Machine Learning Department

September 17, 2018

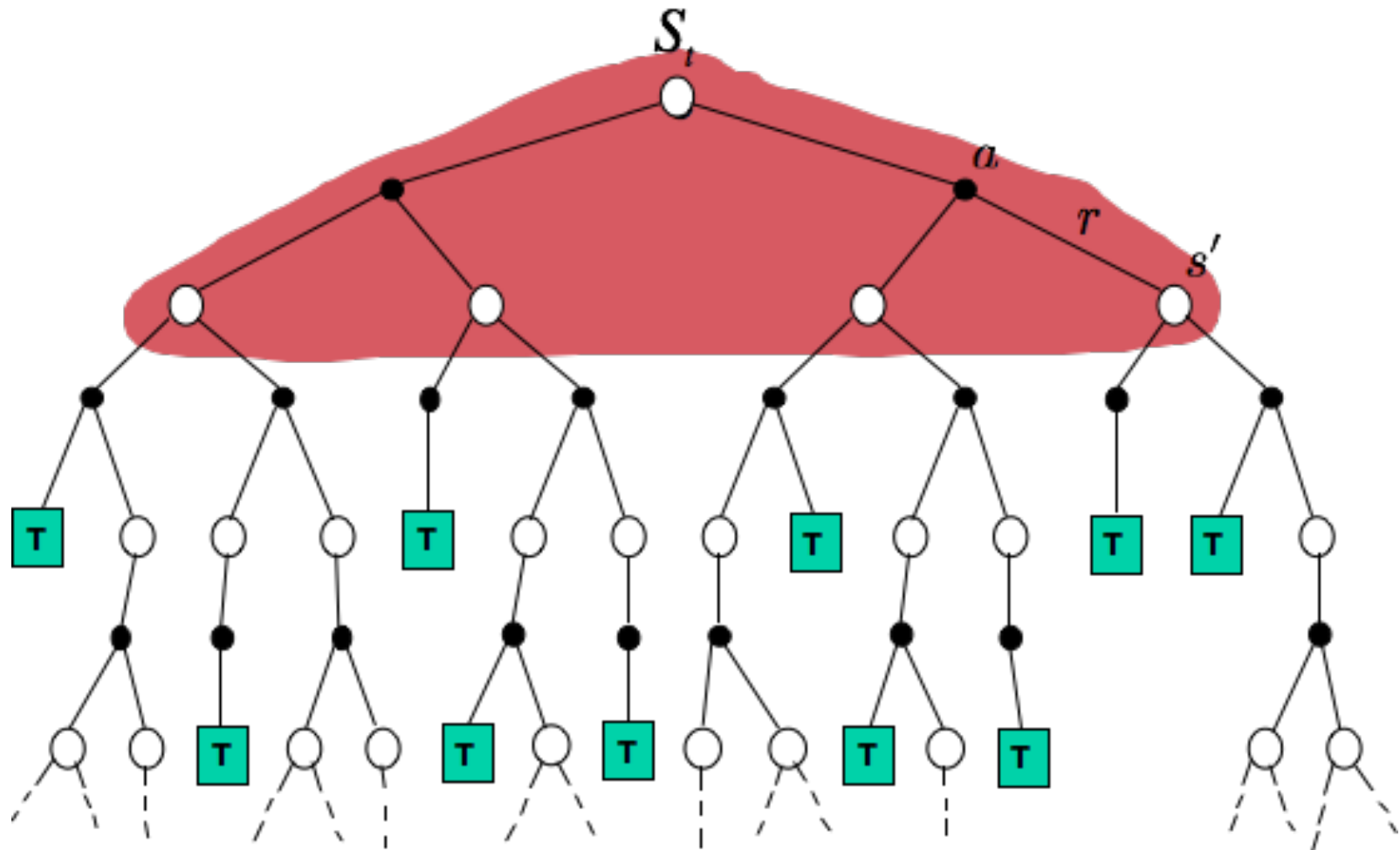
**Temporal Difference Methods**

# Used Materials

- **Acknowledgement:** Much of the material and slides for this lecture were borrowed from Ruslan Salakhutdinov, who in turn borrowed much from Rich Sutton's class and David Silver's class on Reinforcement Learning.

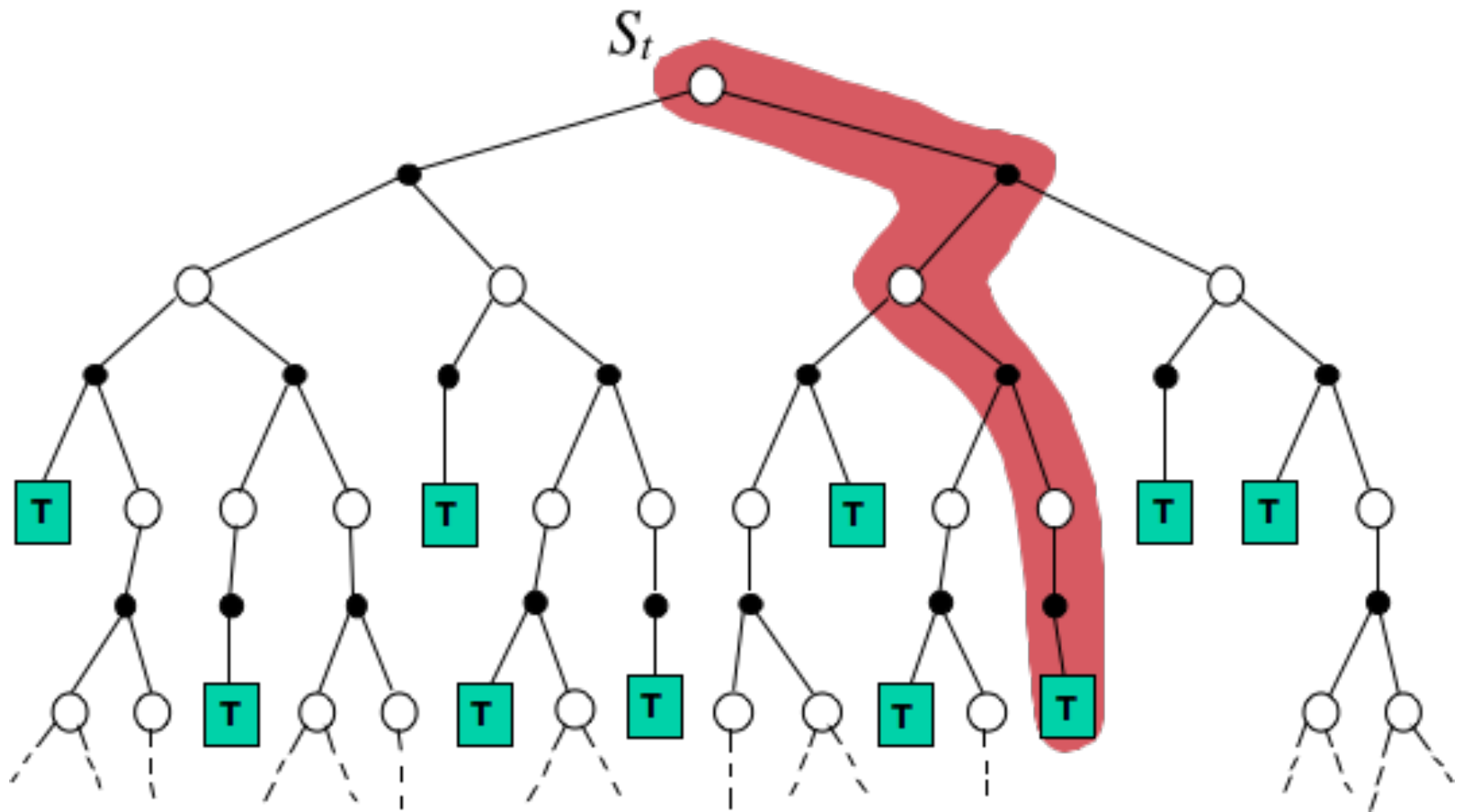
# Dynamic Programming

$$V(S_t) \leftarrow E_{\pi} [R_{t+1} + \gamma V(S_{t+1})] = \sum_a \pi(a|S_t) \sum_{s', r} p(s', r|S_t, a) [r + \gamma V(s')]$$



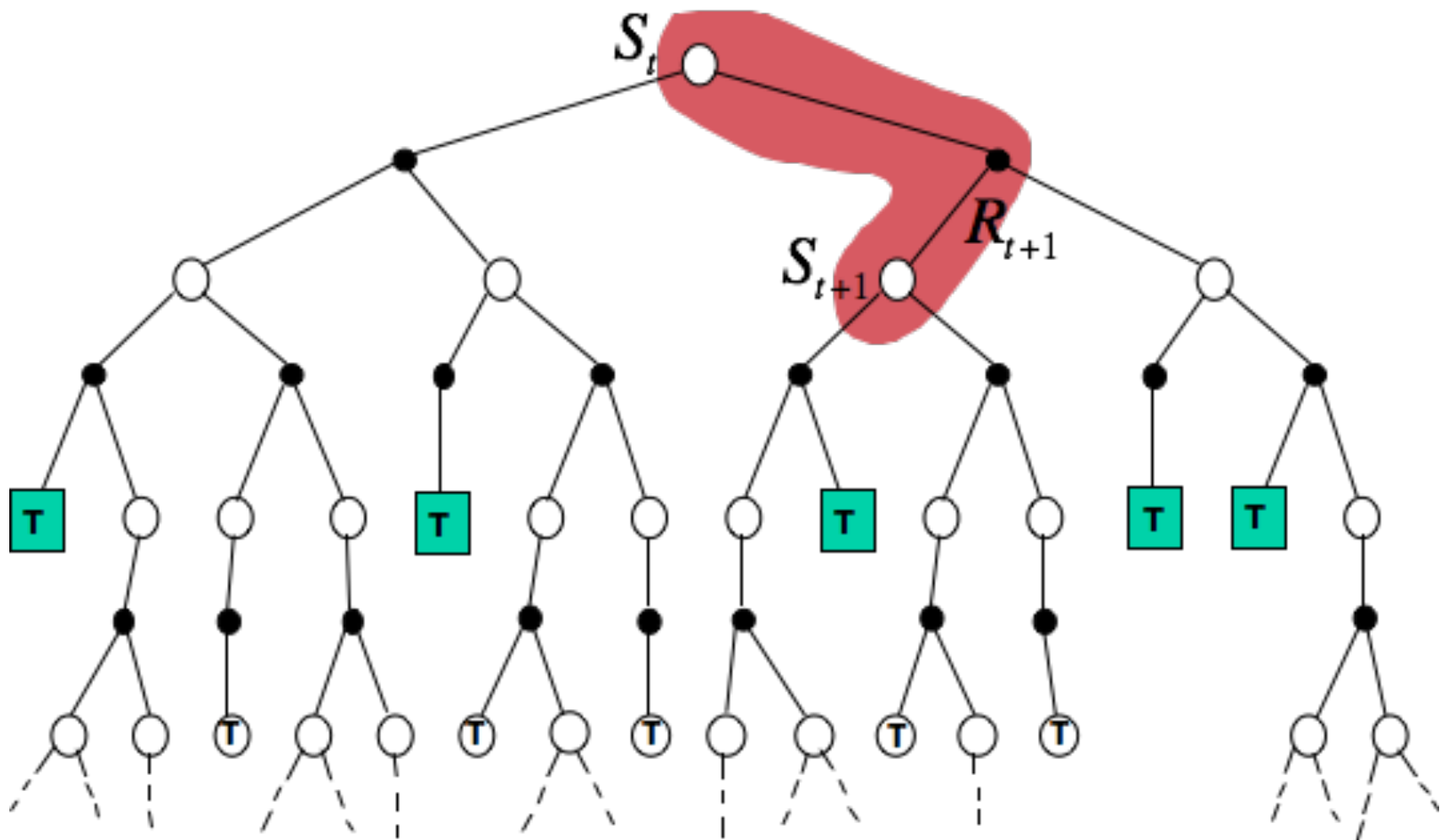
# Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



# Simplest TD(0) Method

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



# TD Prediction

- ▶ **Policy Evaluation** (the prediction problem):
  - for a given policy  $\pi$ , compute the state-value function  $v_\pi$

- ▶ **Remember:** Simple every-visit **Monte Carlo method**:

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ G_t - V(S_t) \right]$$

**target:** the actual return after time  $t$

- ▶ The simplest **Temporal-Difference** method TD(0):

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ \underbrace{R_{t+1} + \gamma V(S_{t+1})}_{\text{target}} - V(S_t) \right]$$

**target:** an estimate of the return

# TD(0) Prediction

- ▶ **Policy Evaluation** (the prediction problem):
  - for a given policy  $\pi$ , compute the state-value function  $v_\pi$
- ▶ The simplest **Temporal-Difference** method TD(0):

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ \underbrace{R_{t+1} + \gamma V(S_{t+1})}_{\text{target}} - V(S_t) \right]$$

- ▶ Or estimate  $Q(S_t, A_t)$  with TD(0)

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ \underbrace{R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a)}_{\text{target}} - Q(S_t, A_t) \right]$$

**target**: an estimate of the return

# Q-Learning: Off-Policy TD(0) Control

- ▶ One-step Q-learning:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

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

$S \leftarrow S'$ ;

until  $S$  is terminal

- ▶ Converges, if every  $(s,a)$  pair is visited infinitely often



Many ways to blend sampled  $R_{t+k}$  and model-based estimates  $Q(s_{t+k}, a)$

Possible target values for training  $Q(s_t, a_t)$ :

$$[R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a)]$$

$$[R_{t+1} + \gamma R_{t+2} + \gamma^2 \max_{a \in A} Q(S_{t+2}, a)]$$

$$[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 \max_{a \in A} Q(S_{t+3}, a)]$$

...

Many ways to blend sampled  $R_{t+k}$  and current estimates of  $Q(s_{t+k}, a)$

Possible target values for training  $Q(s_t, a_t)$ :

$$[R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a)]$$

$$[R_{t+1} + \gamma R_{t+2} + \gamma^2 \max_{a \in A} Q(S_{t+2}, a)]$$

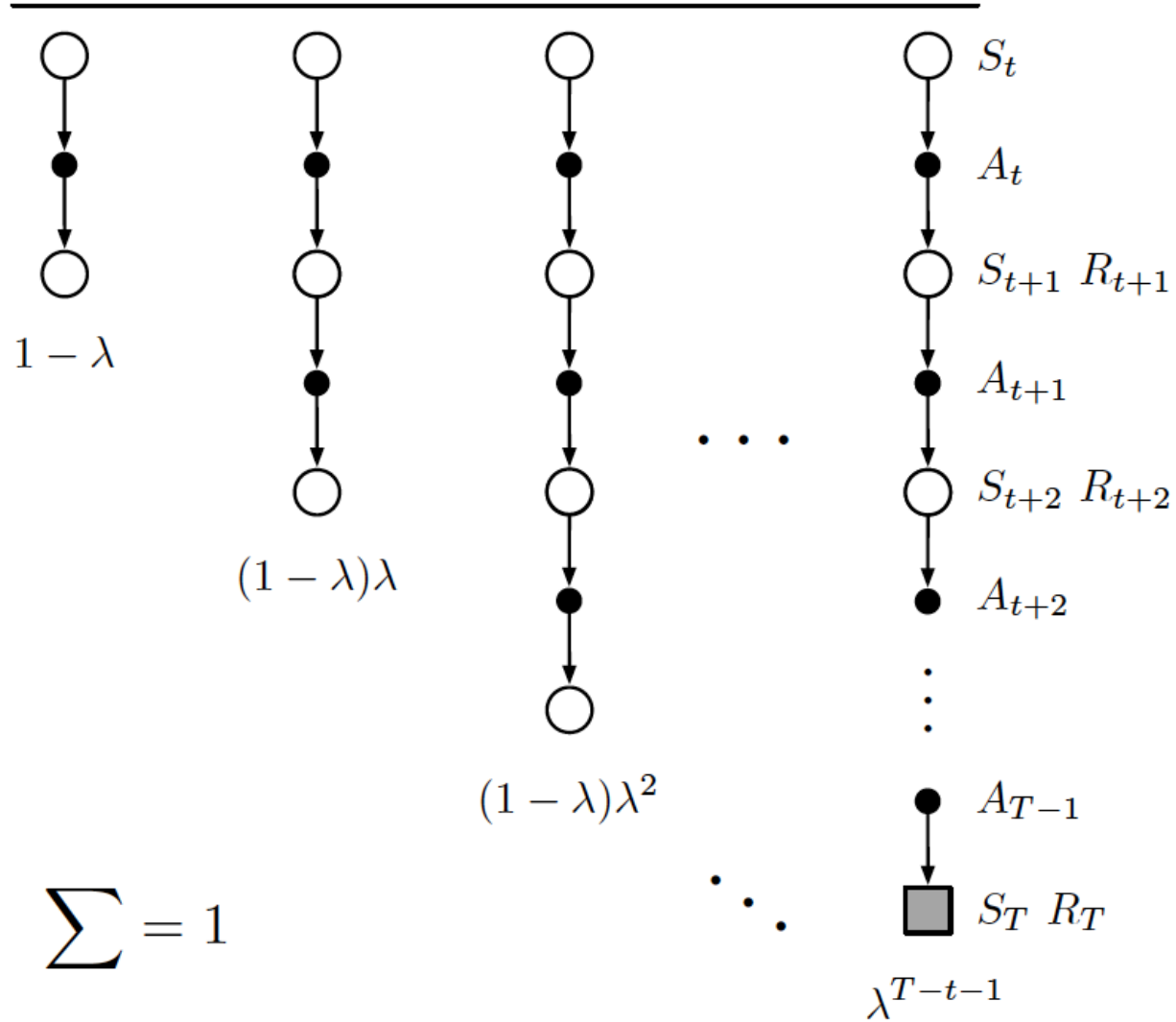
$$[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 \max_{a \in A} Q(S_{t+3}, a)]$$

...

**TD( $\lambda$ ) – Blend all of these**

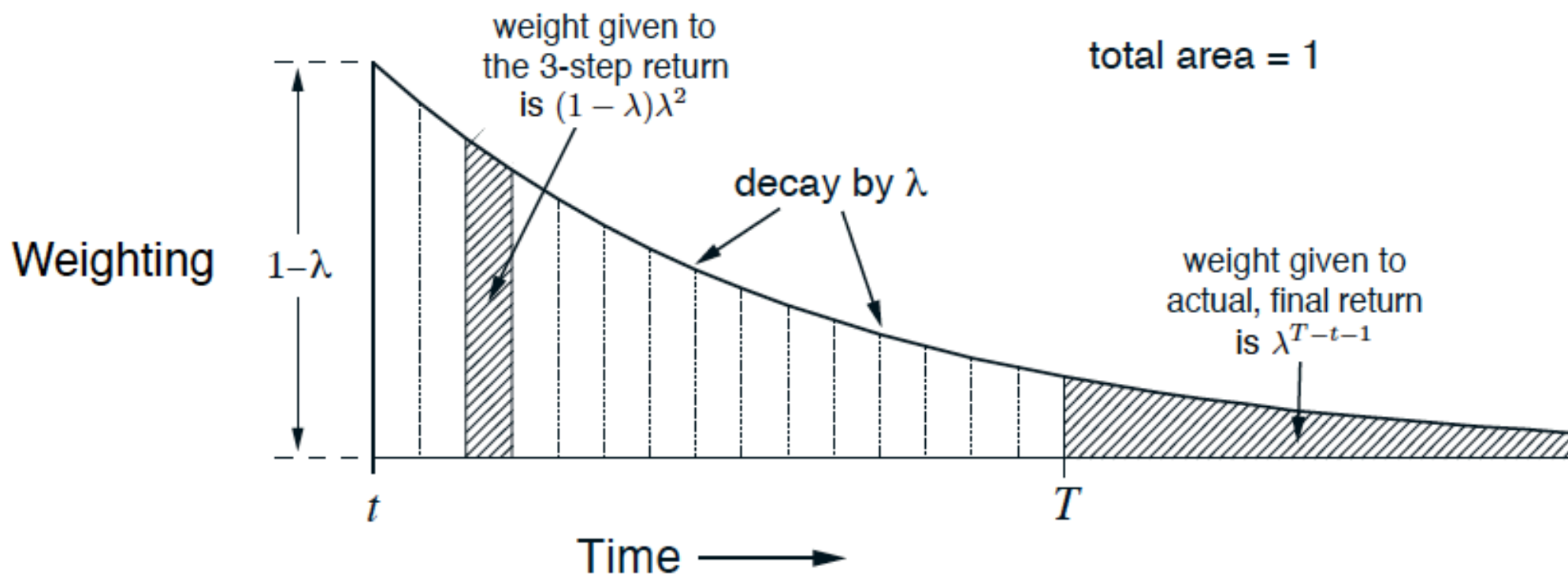
# TD( $\lambda$ )

$$0 \leq \lambda \leq 1$$



**Figure 12.1:** The backup digram for TD( $\lambda$ ). If  $\lambda = 0$ , then the overall update reduces to its first component, the one-step TD update, whereas if  $\lambda = 1$ , then the overall update reduces to its last component, the Monte Carlo update.

$$\begin{aligned}
& (1 - \lambda) [R_{t+1} + \gamma \max_{a \in A} Q(S_{t+1}, a)] \\
& + (1 - \lambda)\lambda [R_{t+1} + \gamma R_{t+2} + \gamma^2 \max_{a \in A} Q(S_{t+2}, a)] \\
& + (1 - \lambda)\lambda^2 [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 \max_{a \in A} Q(S_{t+3}, a)] \\
& + \dots
\end{aligned}$$



**Figure 12.2:** Weighting given in the  $\lambda$ -return to each of the  $n$ -step returns.

# Bias-Variance Trade-Off

- ▶ MC has high variance, zero bias
  - Good convergence properties
  - Even with function approximation
  - Not very sensitive to initial value
  - Very simple to understand and use
- ▶ TD has low variance, some bias
  - Usually more efficient than MC
  - More sensitive to initial values of  $Q(s,a)$  and  $V(s)$