

# Tree-backup and $Q(\sigma)$

Sutton & Barto - 7.5, 7.6, 12.10



KRIS DE ASIS

# Overview

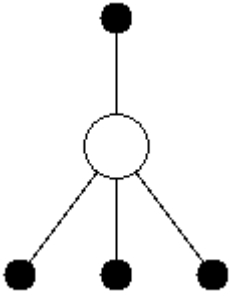
- n-step Tree-backup
- n-step  $Q(\sigma)$
- n-step to Traces
- $TB(\lambda)$  and  $Q(\sigma, \lambda)$

# One-step Expected Sarsa

Computes expectation under target policy directly, can do off-policy learning **without** importance sampling

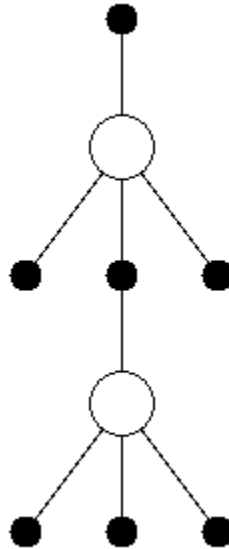
TD Target:

$$\hat{G}_t = R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a)$$

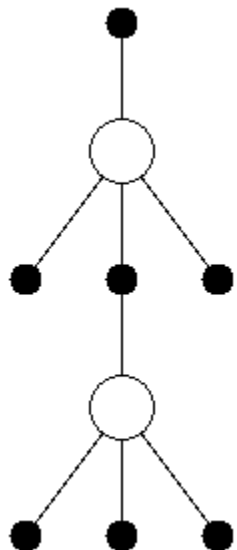


# n-step Tree-backup

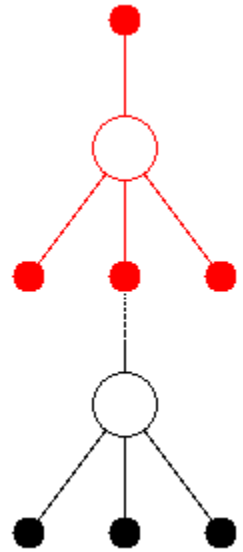
Tree-backup is a multi-step generalization of Expected Sarsa which computes this expectation under the target policy at **every step**



# n-step Tree-backup's Target



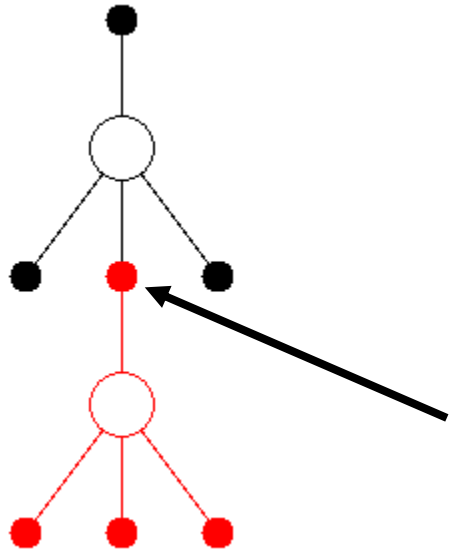
# n-step Tree-backup's Target



$$\hat{G}_t = R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a)$$

Expected Sarsa

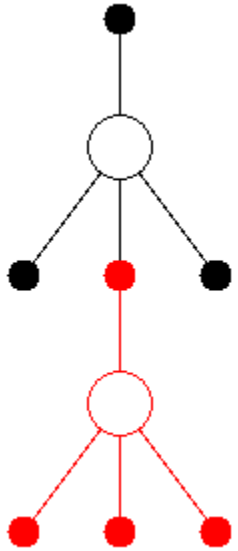
# n-step Tree-backup's Target



$$\hat{G}_t = R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a)$$

Except we know the outcome of the action we took!

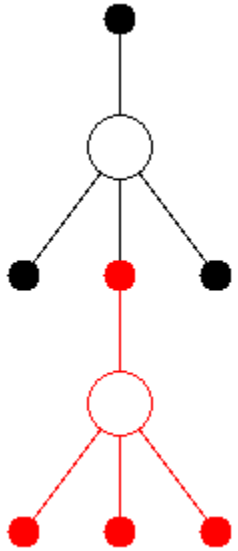
# n-step Tree-backup's Target



$$\hat{G}_t = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1})Q(S_{t+1}, A_{t+1})$$



# n-step Tree-backup's Target

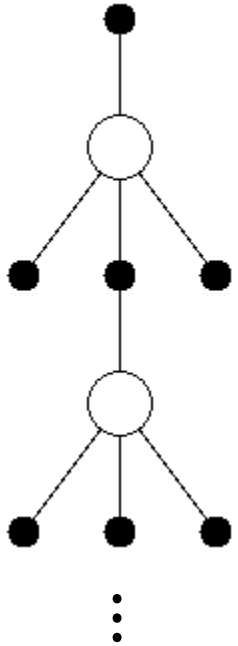


$$\hat{G}_t = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1})\hat{Q}(S_{t+1}, A_{t+1})$$

$$\hat{Q}(S_{t+1}, A_{t+1}) = R_{t+2} + \gamma \sum_{a'} \pi(a'|S_{t+2})Q(S_{t+2}, a')$$

2-step Tree-backup

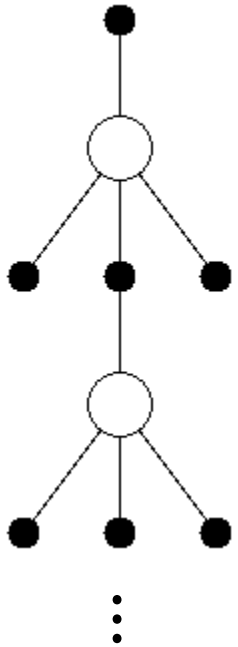
# n-step Tree-backup's Target



$$\hat{G}_{t:h} = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q(S_{t+1}, a) \\ + \gamma \pi(A_{t+1}|S_{t+1})\hat{G}_{t+1:h}$$

$$\hat{G}_{h:h} = Q(S_h, A_h)$$

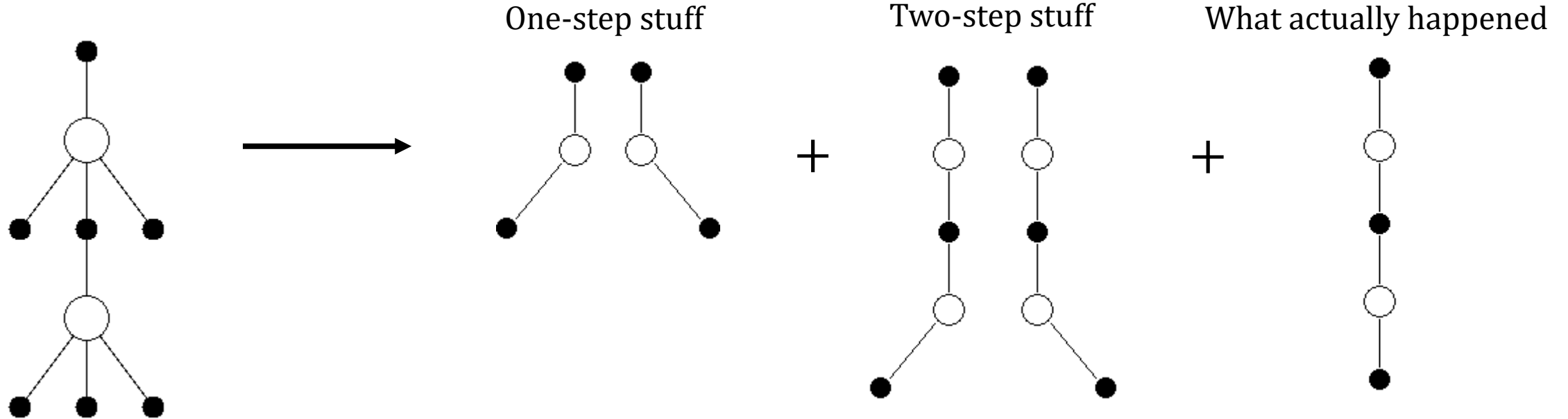
# n-step Tree-backup



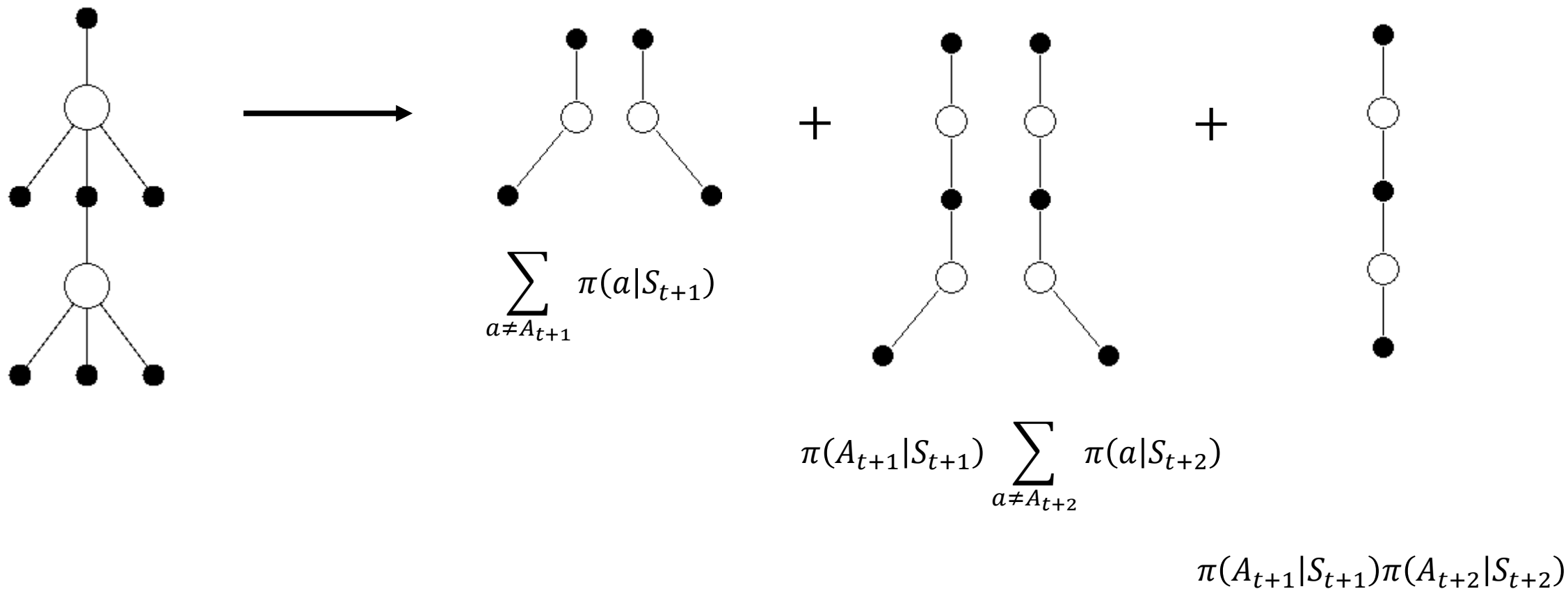
Allows for **multi-step** off-policy learning **without** importance sampling

Can be seen as bootstrapping off of what **didn't happen** to reduce the variance of sampled information

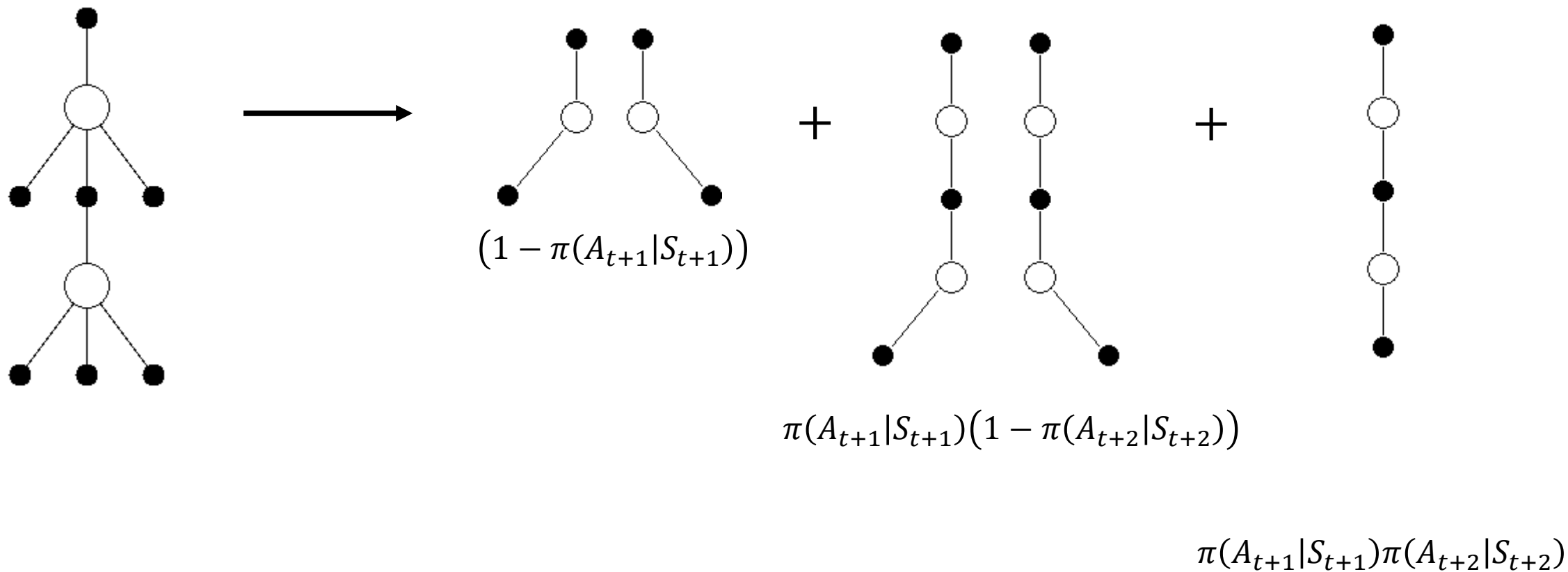
# n-step Tree-backup



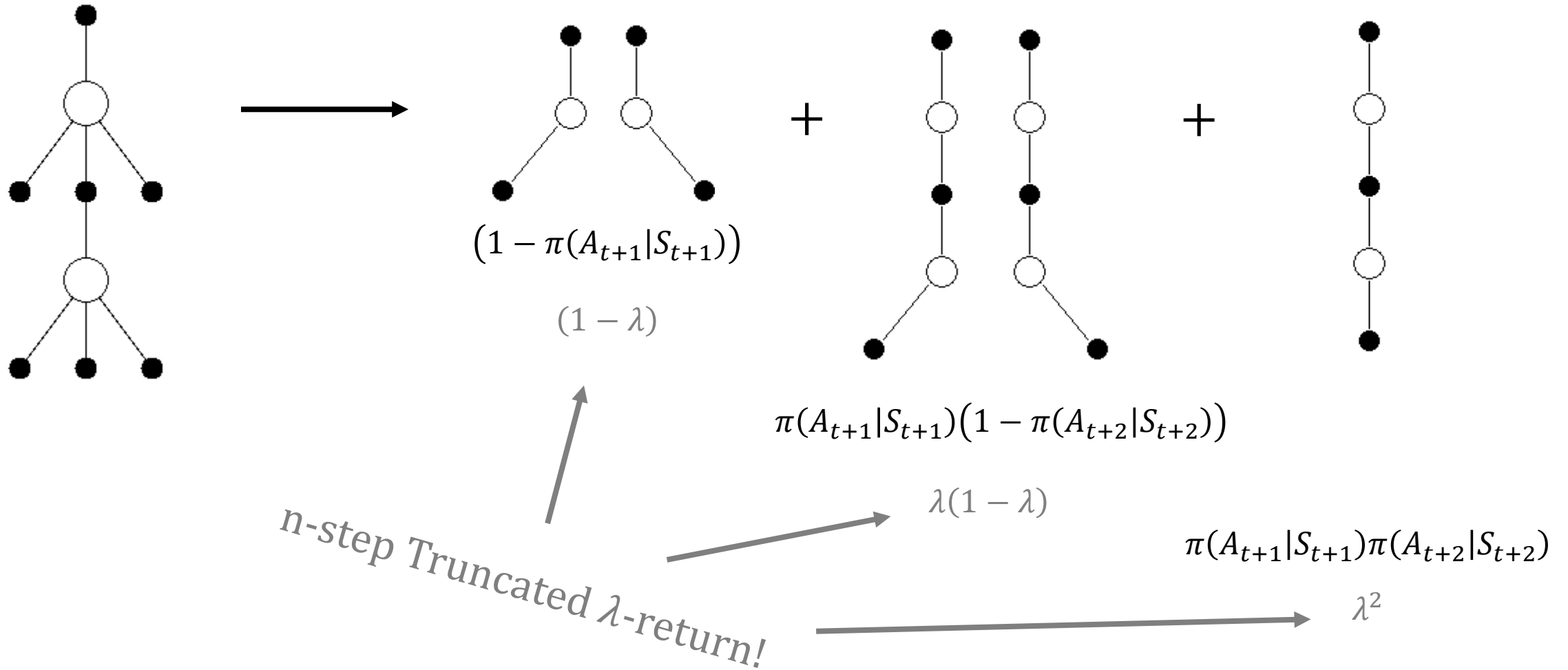
# n-step Tree-backup



# n-step Tree-backup



# n-step Tree-backup



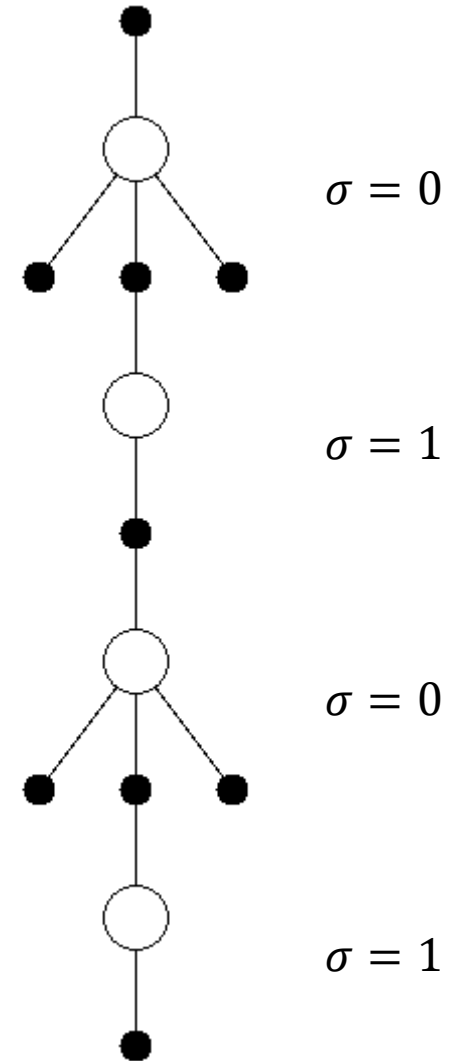
# n-step $Q(\sigma)$

Decide on a per-decision basis whether to:

- Use sampled information ( $\sigma = 1$ )
- Compute an expectation ( $\sigma = 0$ )
- A mix of the two ( $0 < \sigma < 1$ )

Unifies several existing action-value multi-step TD methods:

- n-step Sarsa
- n-step Expected Sarsa
- n-step Tree-backup





# n-step $Q(\sigma)$

Constant values of  $\sigma$ :

$$\sigma \left[ \begin{array}{c} \bullet \\ \circ \\ \bullet \\ \circ \\ \bullet \\ \circ \\ \bullet \end{array} \right] + (1 - \sigma) \left[ \begin{array}{c} \bullet \\ \circ \\ \bullet \quad \bullet \quad \bullet \\ \circ \\ \bullet \quad \bullet \quad \bullet \\ \circ \\ \bullet \quad \bullet \quad \bullet \\ \circ \\ \bullet \quad \bullet \quad \bullet \end{array} \right]$$


But  $\sigma$  doesn't have to be constant!

# n-step $Q(\sigma)$ 's Target

$$\hat{G}_{t:h} = R_{t+1} + \gamma(\sigma_{t+1}(\text{Sarsa}) + (1 - \sigma_{t+1})(\text{TB}))$$

$\rho_{t+1}G_{t+1:h}$



$$\pi(A_{t+1}|S_{t+1})G_{t+1:h} + \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q(S_{t+1}, a)$$


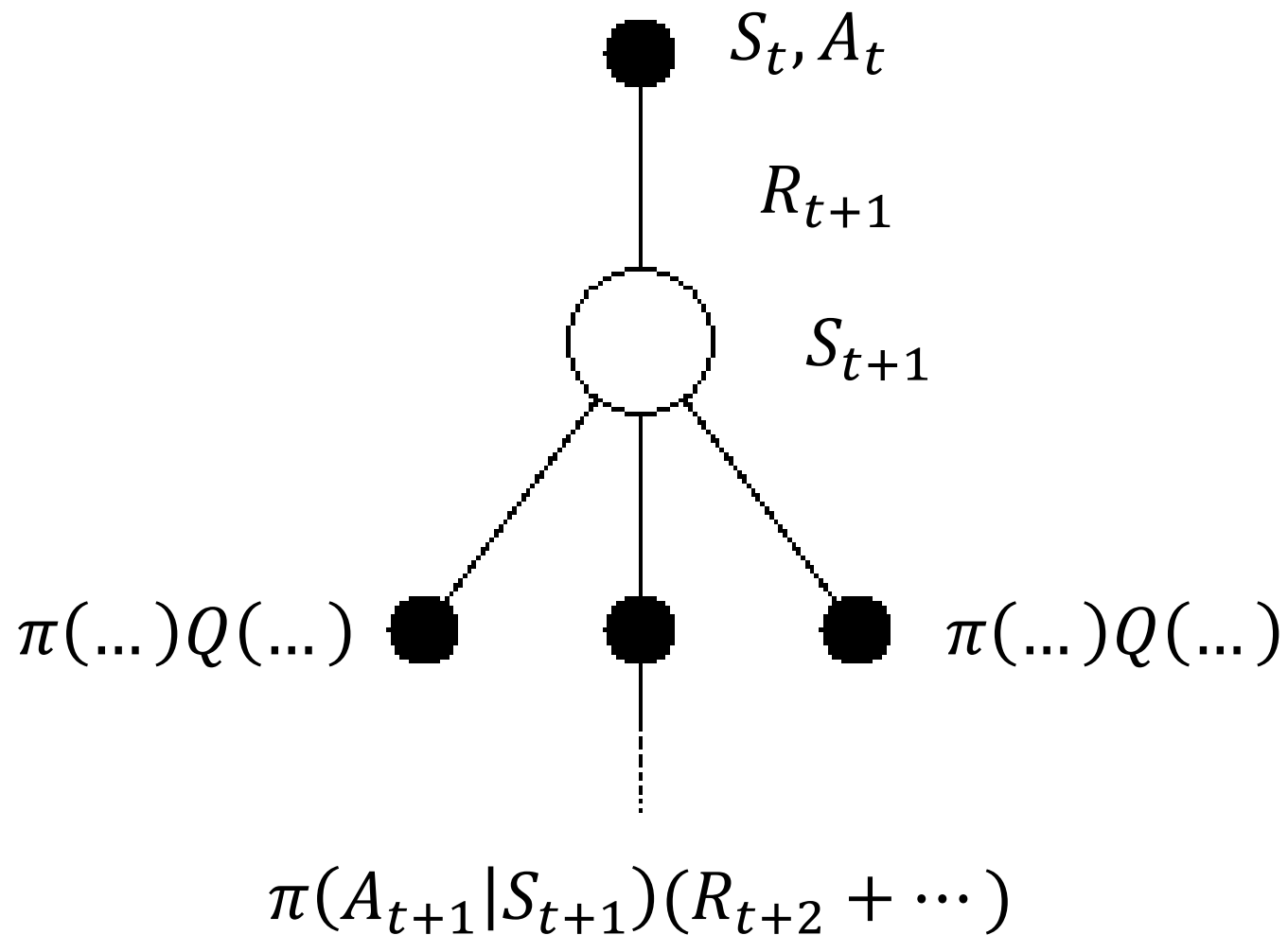
$$\hat{G}_{h:h} = Q(S_h, A_h)$$

# Why unify?

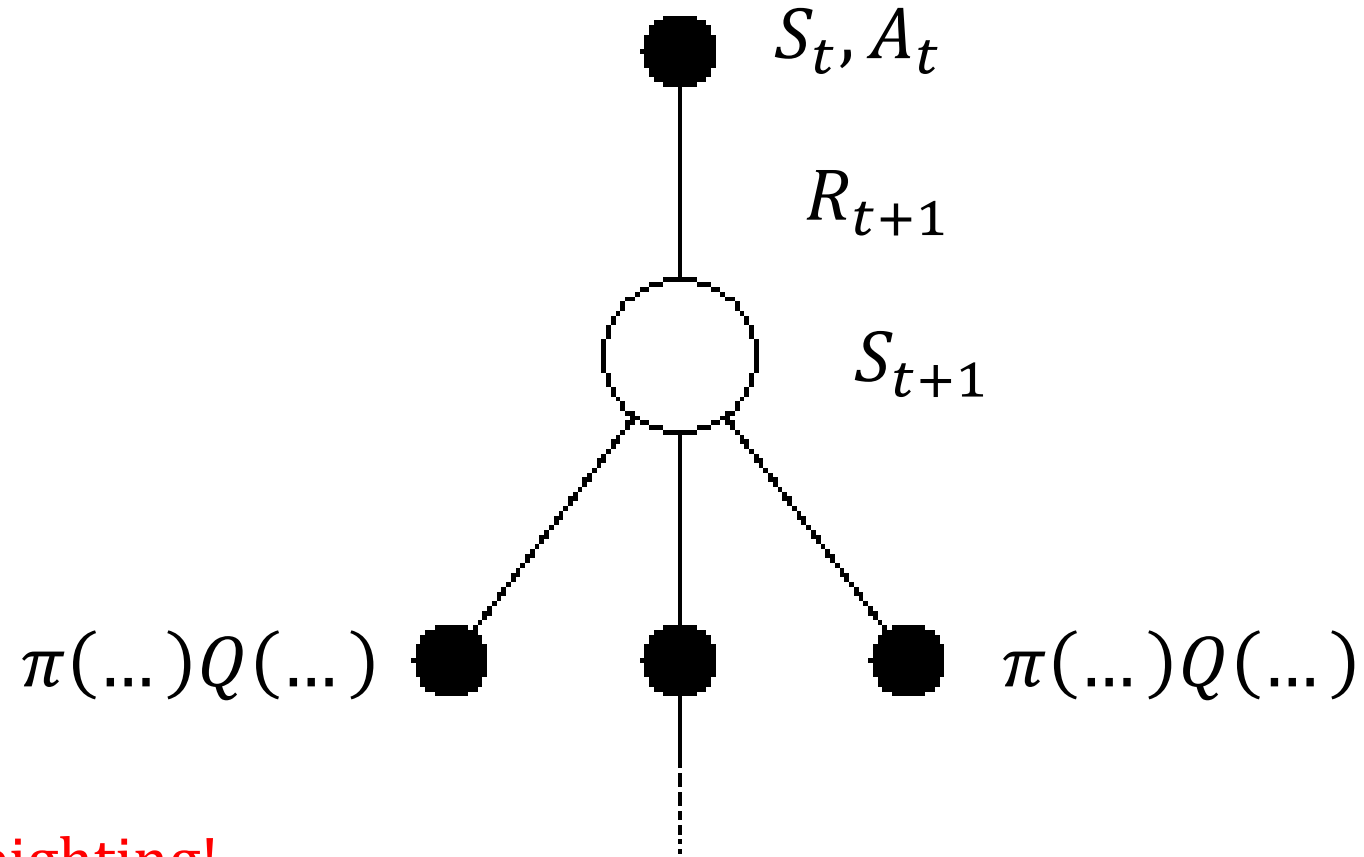
In the one-step case, Expected Sarsa generally outperforms Sarsa

How does n-step Tree-backup compare to n-step Sarsa?

# Why unify?



# Why unify?



$R_{t+2}$  has a lower weighting!

$\rightarrow \pi(A_{t+1}|S_{t+1})(R_{t+2} + \dots)$

# Why unify?

Reward

Weight

$$R_{t+1}$$

$$1$$

$$R_{t+2}$$

$$\pi(A_{t+1}|S_{t+1})$$

$$R_{t+3}$$

$$\pi(A_{t+1}|S_{t+1})\pi(A_{t+2}|S_{t+2})$$

$$R_{t+4}$$

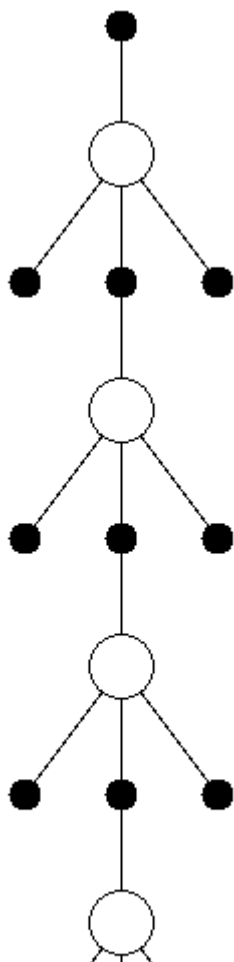
$$\pi(A_{t+1}|S_{t+1})\pi(A_{t+2}|S_{t+2})\pi(A_{t+3}|S_{t+3})$$

$\vdots$

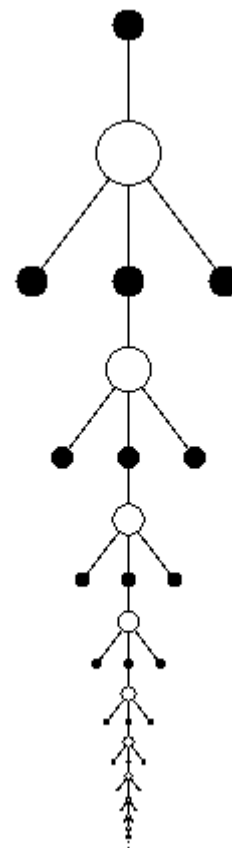
$\vdots$

# Why unify?

# Expectations



# Reality



# Why unify?

$n$ -step Tree-backup has considerably low variance, but is more biased than  $n$ -step Sarsa (for the same  $n$ )

It gives less weight to the sampled rewards, and compensates by giving additional weight to bootstrapping (in the direction of **actions not taken**)

$n$ -step  $Q(\sigma)$  can be viewed as a way to address this bias by adjusting the weighting given to sampled information



# Dynamic $\sigma$

When value estimates are poor, it's better to use more sampled information

As estimates get better, can leverage bootstrapping to get targets with lower variance

Simple heuristic- start with  $\sigma = 1$  (full sampling) and decay towards  $\sigma = 0$  (no sampling- increased bootstrapping through computing expectations)

# 19-state Random Walk

A  $1 \times 19$  grid world with deterministic 2-directional movement

Terminal transitions on each end with rewards -1 and +1, all other transitions have a reward of 0.

On-policy policy evaluation with an equiprobable random policy

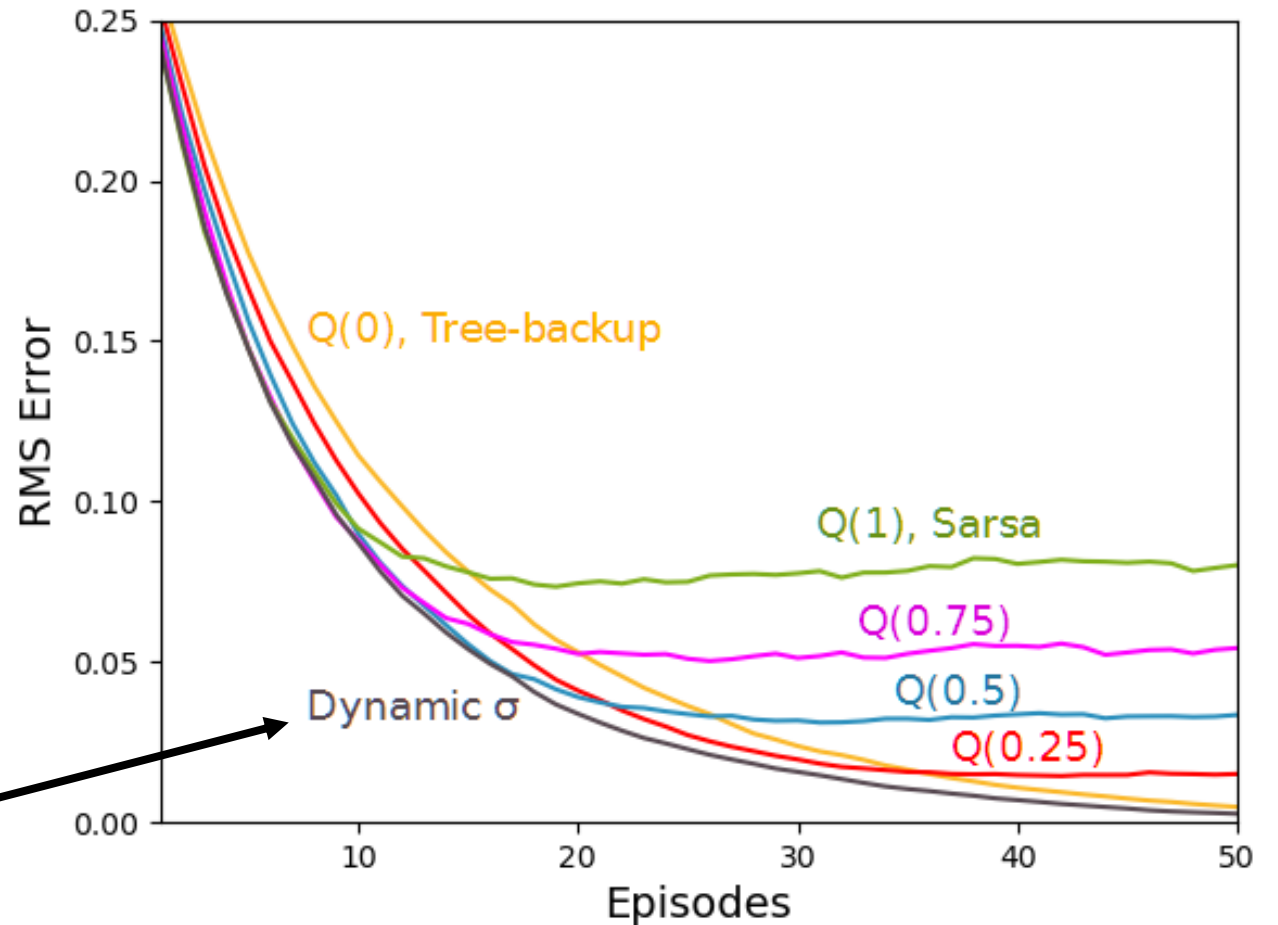


# 19-state Random Walk

3-step  $Q(\sigma)$

Showing best constant step size in terms of RMS error after 50 episodes

Decaying  $\sigma$  after each episode from 1 to 0



# n-step $Q(\sigma)$

$\sigma$  controls a trade-off between initial and final performance- more sampling results in quicker learning but larger variance, and vice versa

$n$  (or  $\lambda$ ) tunes the degree of bootstrapping in the direction of **actions taken**,  $\sigma$  tunes the degree of bootstrapping in the direction of **actions not taken**

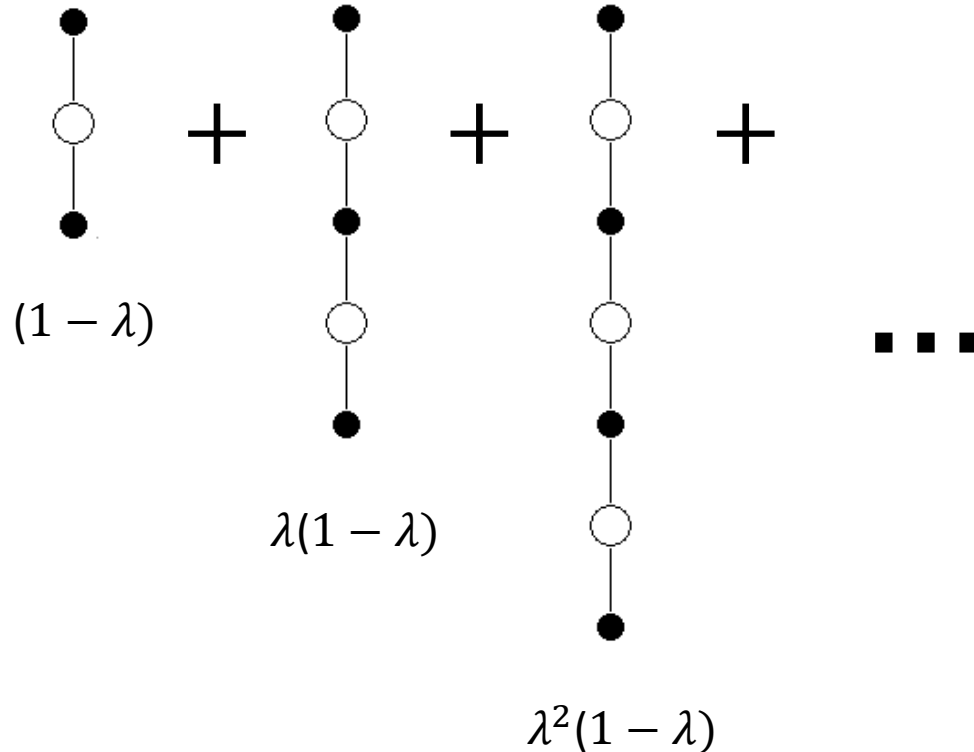
# From n-step to traces

$\lambda$  is often referred to as a **degree of bootstrapping**, but can also be viewed as the **probability of (not) bootstrapping** in a stochastic process

With probability  $1 - \lambda$ , we bootstrap off of current estimates. With probability  $\lambda$ , we continue to the next reward in the return

# From n-step to traces

On-policy Sarsa( $\lambda$ ):



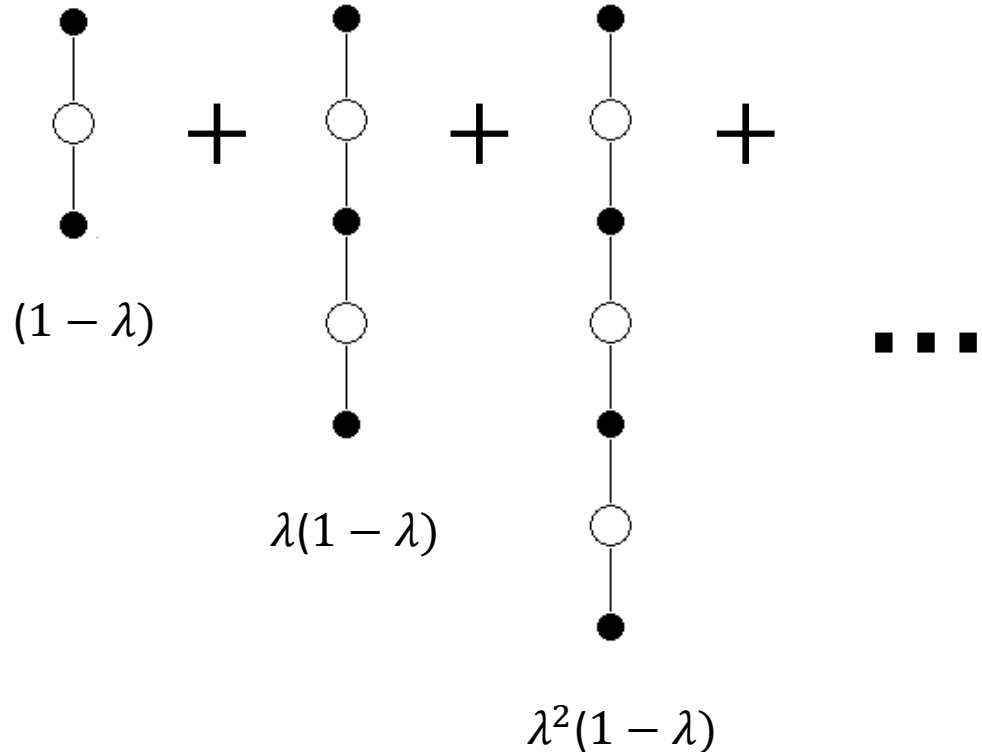
$$\hat{G}_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{G}_{t:t+n}$$



$$\hat{G}_t^\lambda = (1 - \lambda)(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})) + \lambda(R_{t+1} + \gamma \hat{G}_{t+1}^\lambda)$$

# From n-step to traces

On-policy Sarsa( $\lambda$ ):



$$\hat{G}_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{G}_{t:t+n}$$



$$\hat{G}_t^\lambda = R_{t+1} + \gamma(\dots)$$

$$\underbrace{(1 - \lambda)Q(S_{t+1}, A_{t+1})}_{R_{t+1}} + \lambda G_{t+1}^\lambda$$

# Tree-backup( $\lambda$ )

$$\hat{G}_{t:h} = R_{t+1} + \gamma \left( \pi(A_{t+1}|S_{t+1})\hat{G}_{t+1:h} + \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q(S_{t+1}, a) \right)$$
$$\hat{G}_{h:h} = Q(S_h, A_h)$$



# Tree-backup( $\lambda$ )

$$\begin{aligned}\hat{G}_{t:h} &= R_{t+1} + \gamma \left( \pi(A_{t+1}|S_{t+1})\hat{G}_{t+1:h} + \mathbb{E}_{\pi}[Q(S_{t+1}, \cdot)] - \pi(A_{t+1}|S_{t+1})Q(S_{t+1}, A_{t+1}) \right) \\ \hat{G}_{h:h} &= Q(S_h, A_h)\end{aligned}$$

# Tree-backup( $\lambda$ )

$$\begin{aligned}\hat{G}_{t:h} &= R_{t+1} + \gamma \left( \pi_{t+1} \hat{G}_{t+1:h} + \underbrace{\mathbb{E}_{\pi}[Q_{t+1}] - \pi_{t+1} Q_{t+1}}_{\text{Expectation correction-like a control variate!}} \right) \\ \hat{G}_{h:h} &= Q_h\end{aligned}$$

$$\begin{aligned}Q_t &= Q(S_t, A_t) \\ \pi_t &= \pi(A_t | S_t)\end{aligned}$$

$\hat{G}_{t:h}$  with  $\hat{G}_{t+1:h}$  outlines the continuation condition in the recursive  $\lambda$ -return

$\hat{G}_{t:h}$  with  $\hat{G}_{h:h}$  substituted into  $\hat{G}_{t+1:h}$  outlines the bootstrapping condition

# Tree-backup( $\lambda$ )

$$\begin{aligned}\hat{G}_{t:h} &= R_{t+1} + \gamma(\pi_{t+1}\hat{G}_{t+1:h} + \mathbb{E}_{\pi}[Q_{t+1}] - \pi_{t+1}Q_{t+1}) \\ \hat{G}_{h:h} &= Q_h\end{aligned}$$

$$\hat{G}_t^\lambda =$$

$$\begin{aligned}(1 - \lambda) &\left(R_{t+1} + \gamma(\pi_{t+1}Q_{t+1} + \mathbb{E}[Q_{t+1}] - \pi_{t+1}Q_{t+1})\right) + \\ &\lambda \left(R_{t+1} + \gamma(\pi_{t+1}\hat{G}_{t+1}^\lambda + \mathbb{E}[Q_{t+1}] - \pi_{t+1}Q_{t+1})\right)\end{aligned}$$

Easier to work with than:

$$\hat{G}_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{G}_{t:t+n}$$

# Tree-backup( $\lambda$ )

$$\hat{G}_t^\lambda = R_{t+1} + \gamma \left( (1 - \lambda) \mathbb{E}_\pi[Q_{t+1}] + \lambda (\pi_{t+1} G_{t+1}^\lambda + \mathbb{E}_\pi[Q_{t+1}] - \pi_{t+1} Q_{t+1}) \right)$$

$$\hat{G}_t^\lambda = \sum_{k=t}^{\infty} \underbrace{(R_{k+1} + \gamma \mathbb{E}_\pi[Q_{k+1}] - Q_k)}_{\text{One-step TD error}} \prod_{i=t+1}^k \underbrace{\gamma \lambda \pi_i}_{\text{Trace-decay rate}}$$

# Tree-backup( $\lambda$ )

$$\hat{G}_t^\lambda = R_{t+1} + \gamma \left( (1 - \lambda) \mathbb{E}_\pi[Q_{t+1}] + \lambda (\pi_{t+1} G_{t+1}^\lambda + \mathbb{E}_\pi[Q_{t+1}] - \pi_{t+1} Q_{t+1}) \right)$$

$$\hat{G}_t^\lambda = \sum_{k=t}^{\infty} (R_{k+1} + \gamma \mathbb{E}_\pi[Q_{k+1}] - Q_k) \prod_{i=t+1}^k \gamma \lambda \pi_i$$

$$\pi_t \leq \rho_t = \frac{\pi_t}{\mu_t}$$

For the same  $\lambda$ , Tree-backup's traces decay quicker- performs "shorter" backups

$$\hat{G}_t^\lambda = \sum_{k=t}^{\infty} (R_{k+1} + \gamma \rho_{k+1} Q_{k+1} - Q_k) \prod_{i=t+1}^k \gamma \lambda \rho_i$$

Sarsa( $\lambda$ )

# $Q(\sigma, \lambda)$

$$\hat{G}_t^\lambda = R_{t+1} + \gamma(\sigma_{t+1}(\text{Sarsa}) + (1 - \sigma_{t+1})(\text{TB}))$$

$$(1 - \lambda)\rho_{t+1}Q_{t+1} + \lambda\rho_{t+1}G_{t+1}^\lambda$$

$$(1 - \lambda)\mathbb{E}_\pi[Q_{t+1}] + \lambda(\pi_{t+1}G_{t+1}^\lambda + \mathbb{E}_\pi[Q_{t+1}] - \pi_{t+1}Q_{t+1})$$

$$\hat{G}_t^\lambda = \sum_{k=t}^{\infty} \underbrace{(R_{k+1} + \gamma(\sigma_{k+1}\rho_{k+1}Q_{k+1} + (1 - \sigma_{k+1})\mathbb{E}_\pi[Q_{k+1}]) - Q_k)}_{\text{One-step TD error}} \prod_{i=t+1}^k \gamma\lambda(\sigma_i\rho_i + (1 - \sigma_i)\pi_i)$$

One-step TD error

Trace-decay rate

# $Q(\sigma, \lambda)$

$$\hat{G}_t^\lambda = R_{t+1} + \gamma(\sigma_{t+1}(\text{Sarsa}) + (1 - \sigma_{t+1})(\text{TB}))$$

$$(1 - \lambda)\rho_{t+1}Q_{t+1} + \lambda\rho_{t+1}G_{t+1}^\lambda$$

$$(1 - \lambda)\mathbb{E}_\pi[Q_{t+1}] + \lambda(\pi_{t+1}G_{t+1}^\lambda + \mathbb{E}_\pi[Q_{t+1}] - \pi_{t+1}Q_{t+1})$$

$$\hat{G}_t^\lambda = \sum_{k=t}^{\infty} (R_{k+1} + \gamma(\sigma_{k+1}\rho_{k+1}Q_{k+1} + (1 - \sigma_{k+1})\mathbb{E}_\pi[Q_{k+1}]) - Q_k) \prod_{i=t+1}^k \gamma\lambda(\sigma_i\rho_i + (1 - \sigma_i)\pi_i)$$

**Adjustable** action-  
dependent trace decay rate  
through  $\sigma$

# CV Q( $\sigma, \lambda$ )

$$\hat{G}_t^\lambda = R_{t+1} + \gamma(\sigma_{t+1}(\text{ACV Sarsa}) + (1 - \sigma_{t+1})(\text{TB}))$$

$$(1 - \lambda)\mathbb{E}_\pi[Q_{t+1}] + \lambda(\pi_{t+1}G_{t+1}^\lambda + \mathbb{E}_\pi[Q_{t+1}] - \pi_{t+1}Q_{t+1})$$

$$(1 - \lambda)\mathbb{E}_\pi[Q_{t+1}] + \lambda(\rho_{t+1}G_{t+1}^\lambda + \mathbb{E}_\pi[Q_{t+1}] - \rho_{t+1}Q_{t+1})$$

$$\hat{G}_t^\lambda = \sum_{k=t}^{\infty} \underbrace{(R_{k+1} + \gamma\mathbb{E}_\pi[Q_{k+1}] - Q_k)}_{\text{One-step TD error}} \prod_{i=t+1}^k \gamma\lambda(\sigma_i\rho_i + (1 - \sigma_i)\pi_i)$$

One-step TD error

Trace-decay rate

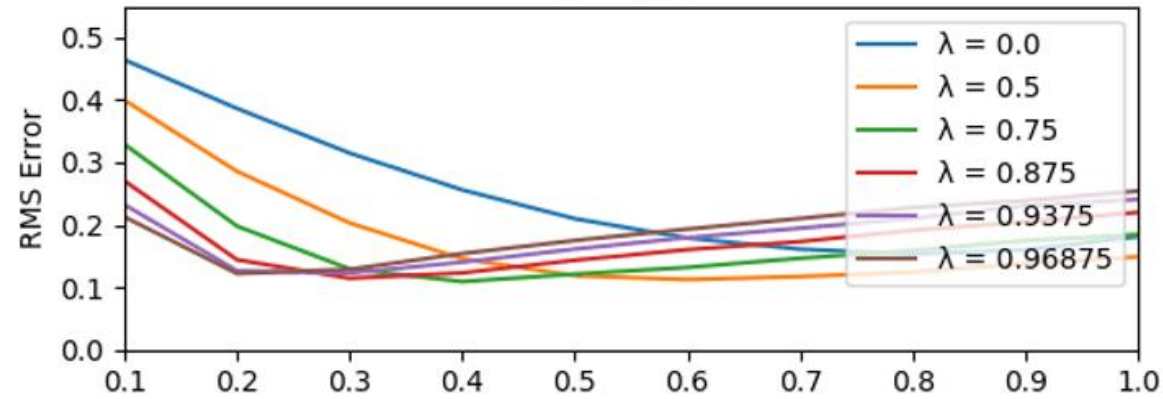


# 19-State Random Walk

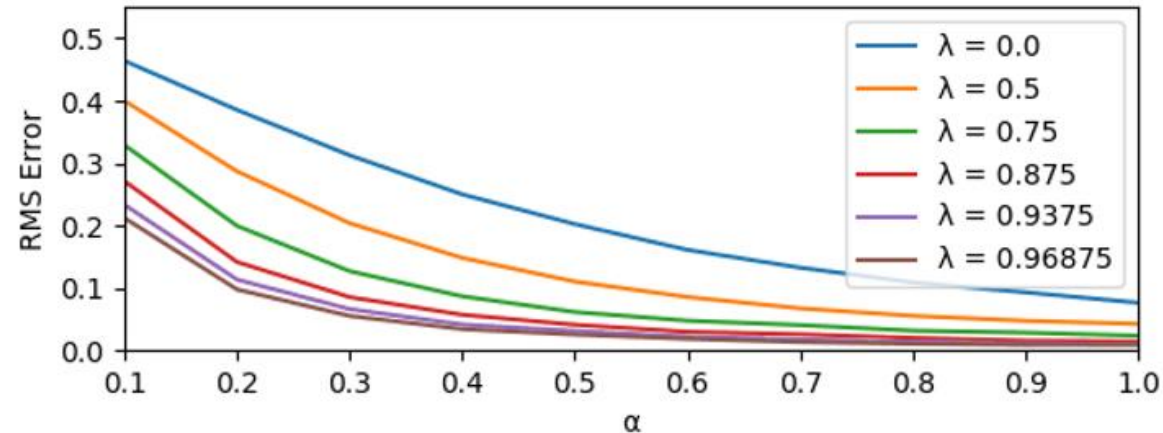
Lots of sampling

$\sigma = 0.75$

$Q(\sigma, \lambda):$



$CV Q(\sigma, \lambda):$



# Questions?

- De Asis, K.\*, Hernandez-Garcia, J. F.\*, Holland G. Z.\*, Sutton, R. S. (2018). [Multi-step Reinforcement Learning: A Unifying Algorithm.](#) *AAAI 2018*.



KRIS DE ASIS