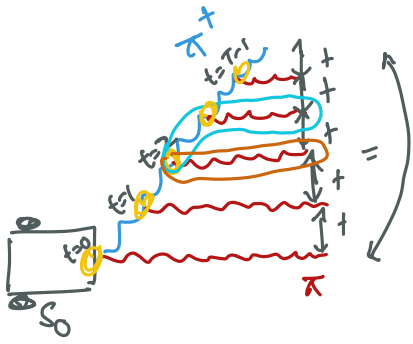


# PERFORMANCE DIFFERENCE LEMMA (PDL)



$$\begin{aligned}
 & V(s_0) - V(s_0) \\
 &= \sum_{t=0}^{T-1} E_{s_t \sim d_t^{\pi^+}} \left[ \underbrace{Q(s_t, \pi^+(s_t)) - V(s_t)}_{\text{ADVANTAGE OF } \pi^+ \text{ OVER } \pi} \right]
 \end{aligned}$$

PERFORMANCE DIFFERENCE IS

$$A(s_t, \pi^+(s_t) | \pi(s_t))$$

THE SUM OVER ON-POLICY ADVANTAGE"  
 (states  $\pi^+$  visits)  
 (OF  $\pi^+$  OVER  $\pi$ )

$A \leq 0 \Rightarrow \pi^+$  is better than  $\pi$  at that state

## POLICY ITERATION

$$\pi^+(s) = \underset{a}{\operatorname{argmin}} Q(s, a) \quad \forall s$$

$$Q(s, \pi^+(s)) \leq \underline{Q(s, \pi(s))} = V(s) \quad \forall s$$

$$Q(s, \pi^+(s)) - V(s) \leq 0 \quad \forall s$$

$$\underline{A(s, \pi^+(s))} \leq 0 \quad \forall s$$

$$V(s) - V(s)$$

$$= \sum_{t=0}^{T-1} E_{\substack{S_t \sim d_t \\ \gamma^+}} [A^{\gamma^+}(S_t, \gamma^+)] \leq 0$$

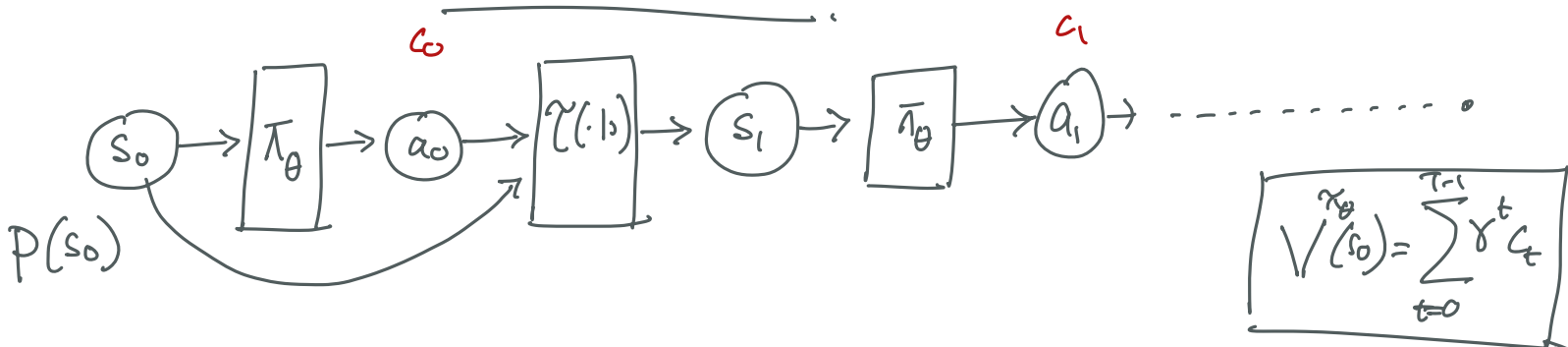
$$\Rightarrow V^{\gamma^+}(s_0) \leq V^{\gamma}(s_0)$$

$$S_t \sim d_t^{\gamma^+}$$

$$P(s_0) \pi(a_0 | s_0) \underbrace{P(s_1 | s_0, a_0)}_{\text{TRANSITION}}$$

$$\pi(a_2 | s_2) \dots P(s_T | s_{T-1}, a_{T-1})$$

### POLICY GRADIENTS



$$\nabla_{\theta} V^{\pi_\theta}(s_0) = \nabla_{\theta} J(\theta)$$