Sarsa (on-policy TD control) for estimating $Q \approx q_*$ Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$ Initialize Q(s, a), for all $s \in S^+$, $a \in A(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode: Initialize S

Choose A from S using policy derived from Q (e.g., ε -greedy) Loop for each step of episode: Take action A, observe R, S'

Choose A' from S' using policy derived from Q (e.g., ε -greedy)

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

 $S \leftarrow S' : A \leftarrow A' :$

until S is terminal