



Every time-step:

Policy improvement:
 ϵ -greedy policy improvement

Policy evaluation:
 SARSA, $Q \approx q_\pi$