stay
in ⟶ (in, stay)

quit

$\frac{2}{3} : +\$4$   $\frac{1}{3} : +\$4$

(in, quit) ⟶ end

$1 : -\$10$

states $s$
chance nodes
actions $a$
transition $T(s' \mid s, a) \in \mathbb{R}$
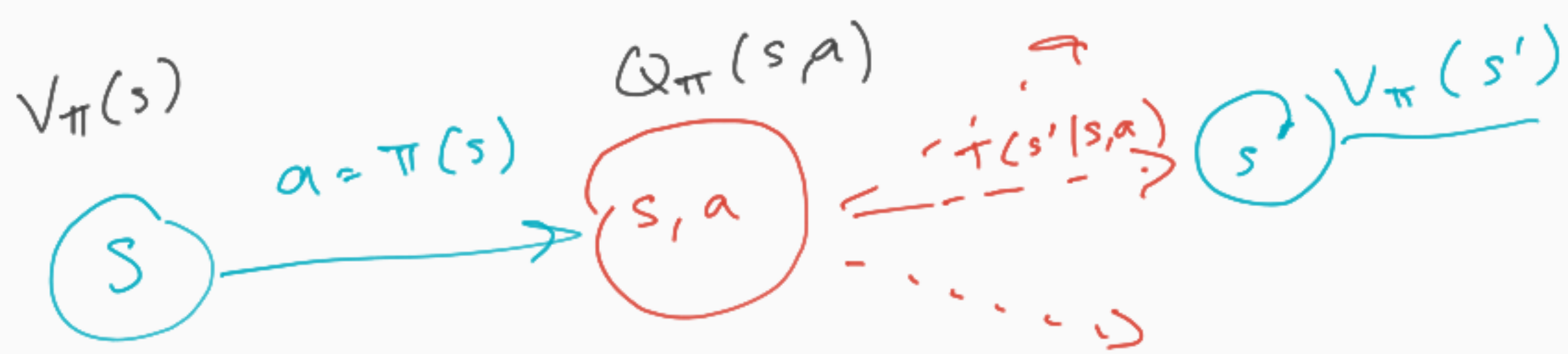reward $R(s, a, s')$
discount factor $0 \le \gamma \le 1$
policy $\pi(s) \rightarrow a$
utility $\sum_{t=0}^{T} \gamma^t r_t$
Value $E[U] = V_\pi(s)$
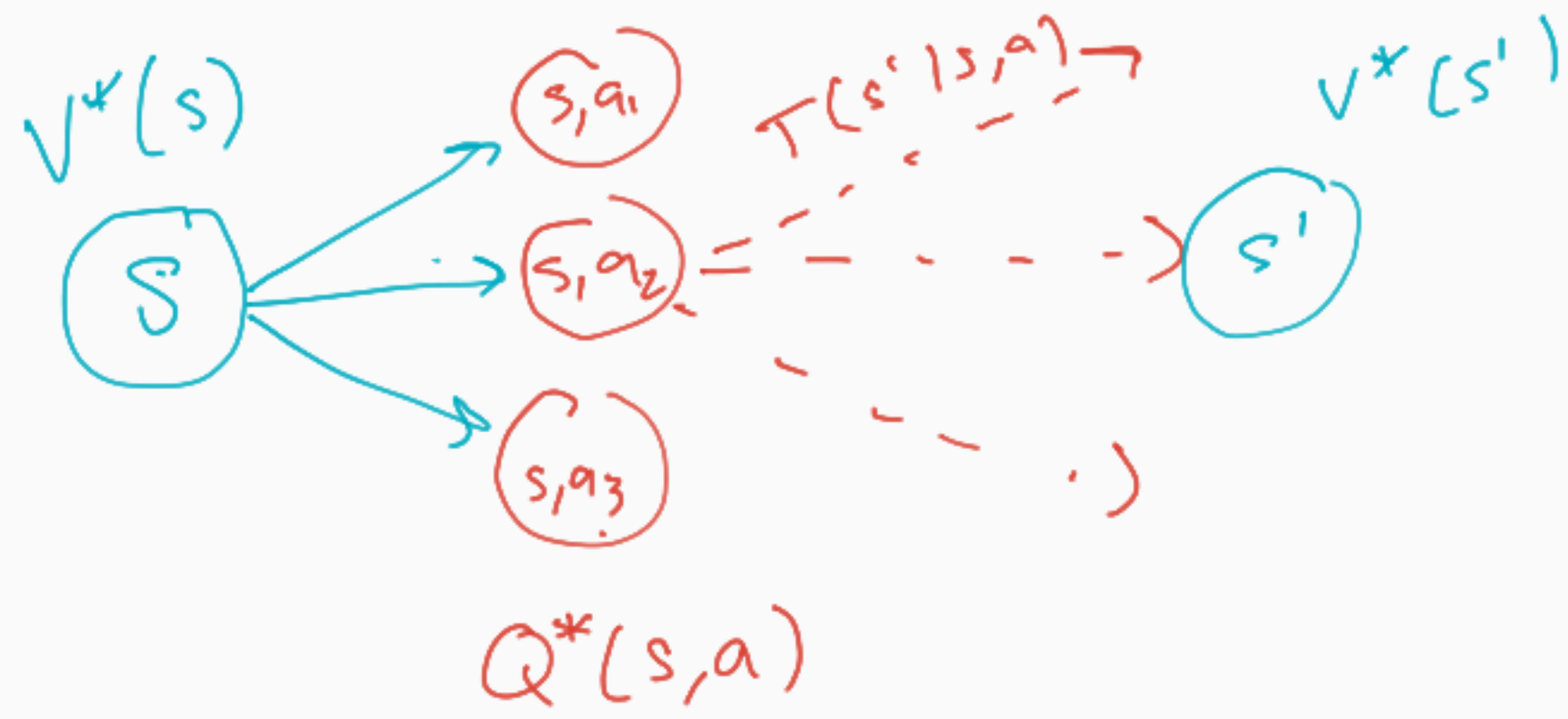
$V_\pi(s)$

$Q_\pi(s,a)$

$V_\pi(s')$

$a = \pi(s)$

$s$

$s, a$

$T(s'|s,a)$

$s'$

$$V_\pi(s) = \begin{cases} 0 & \text{IsEnd}(s) = \text{True} \\ Q_\pi(s, \pi(s)) \end{cases}$$

$$Q_\pi(s,a) = \sum_{s'} T(s'|s,a) \left[ R(s,a,s') + \underbrace{\gamma V_\pi(s')}_{\gamma Q_\pi(s', \pi(s'))} \right]$$

Bellman Equation

$V^*(s)$

$(s, a_1)$     $T(s'|s, a) \rightarrow$     $V^*(s')$

$S$     $(s, a_2) = - \cdot - - \rightarrow$     $s'$

$(s, a_3)$

$Q^*(s, a)$

$$V^*(s) = \begin{cases} 0 & \text{isEnd}(s) = \text{True} \\ \max_{a \in \text{Actions}(s)} Q^*(s, a) \end{cases}$$

$$Q^*(s, a) = \sum_{s'} T(s'|s, a) \left[ R(s, a, s') + \gamma V^*(s') \right]$$

$\left.\begin{array}{l}\text{Bellman} \\ \text{Optimality} \\ \text{Equation.}\end{array}\right\}$

$$\pi^*(s) = \underset{a \in \text{Actions}(s)}{\text{argmax}} \, Q^*(s, a)$$