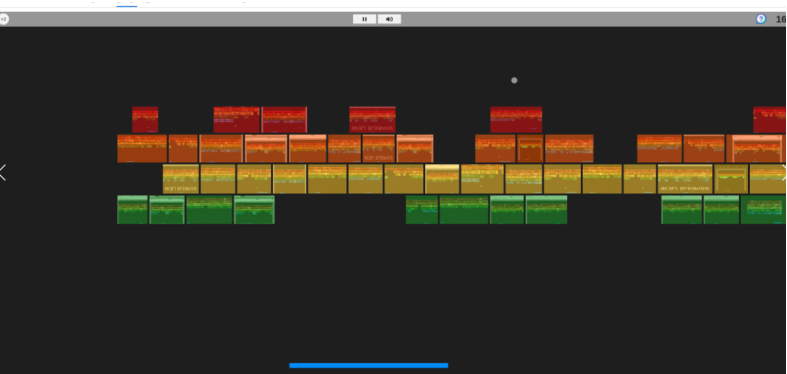
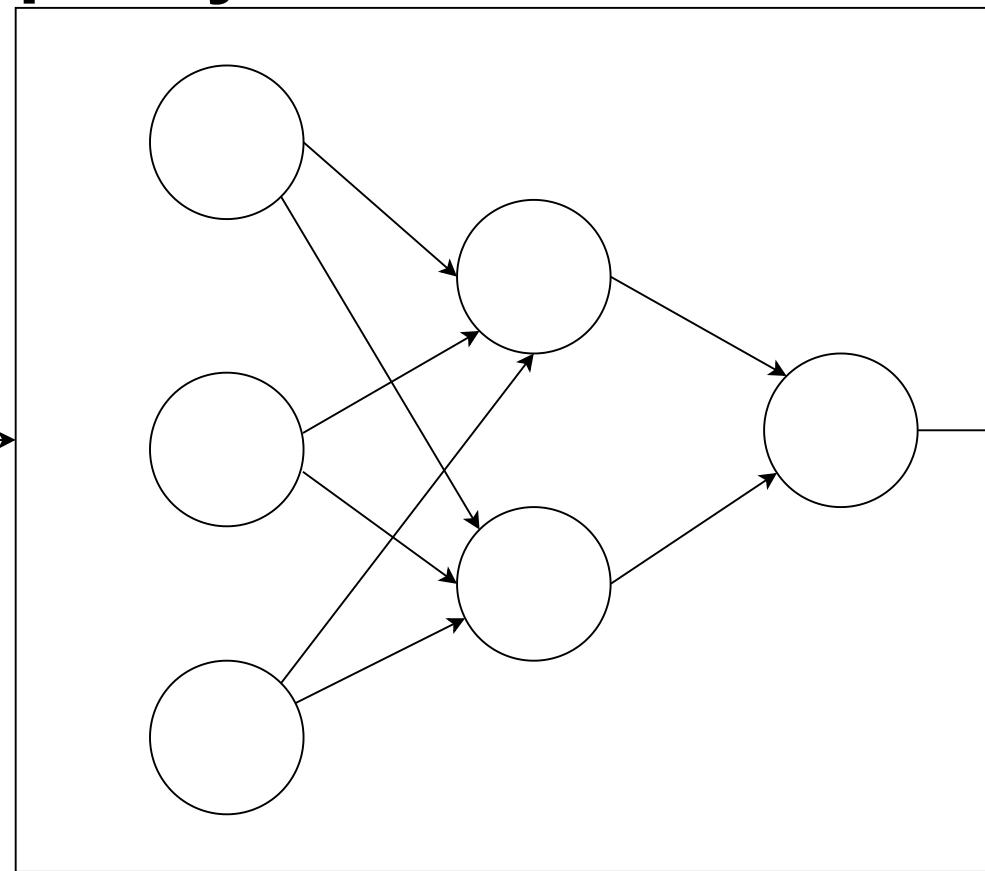


Policy Gradient: Directly Optimize the policy



State, s

Input



Agent

$\pi(s)$

$$P(a_1|s) = 0.9$$

$$P(a_2|s) = 0.1$$

$$P(a_3|s) = 0.0$$

Output

$$\sum_{a_i \in \mathcal{A}} P(a_i|s) = 1$$

$$\pi(s) \sim P(a|s)$$