

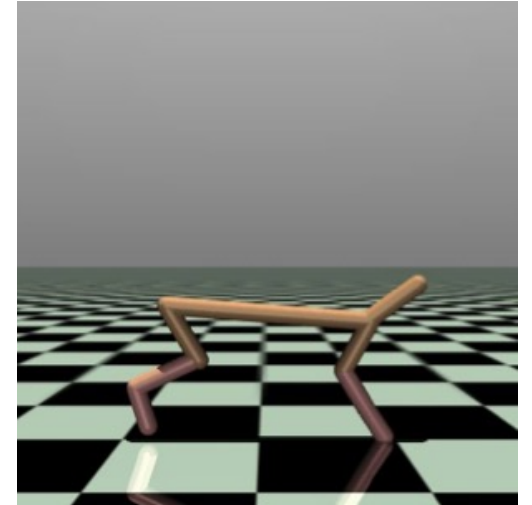
# Learning Self-Imitating Diverse Policies [Gangwani et al., 2019]

## Challenge:

- Exploration is crucial in sparse rewards settings
- Current algorithms don't decompose trajectories to steps

## Main ideas:

- $\pi(a|s)$  better reflected by  $\rho_{\pi}(s, a)$
- Optimize  $\pi$  by minimizing divergence between  $\rho_{\pi}$  and  $\rho_{\pi^*}$
- *Self-imitation from experience:*
  - Approximate  $\rho_{\pi^*}$  by high-rewarded trajectories



# Learning Self-Imitating Diverse Policies [Gangwani et al., 2019]

## Policy update:

- Approximate gradient of Jensen-Shannon divergence

