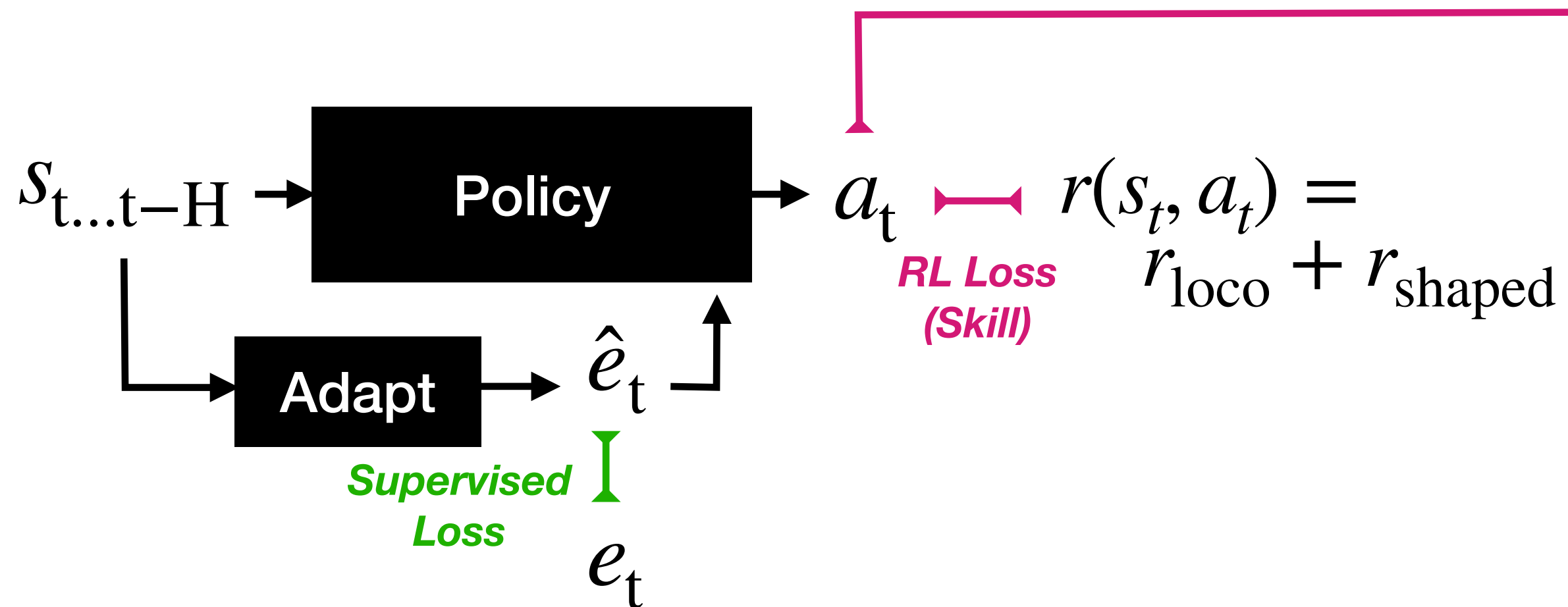


## Stage 1: Motor Skill Learning

Policy Optimization with Concurrent State Estimation

$$L(\theta) = \hat{\mathbb{E}}_t \left[ \log \pi_{\theta}(a_t | s_t, \hat{e}_t) \hat{A}_t \right] + \|e_t - \hat{e}_t\|^2$$



## Stage 2: Adapt for Active Estimation

Policy Optimization with State Estimation Reward + Imitation Loss

$$L(\theta) = \hat{\mathbb{E}}_t \left[ \log \pi_{\theta}^{\text{est}}(a_t | s_t, \hat{e}_t) \hat{A}_t^{\text{est}} \right] + \|e_t - \hat{e}_t\|^2 + \beta D_{\text{KL}} [\pi_{\theta}^{\text{est}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]$$

