$$\Delta_{\mathrm{PCL}}(\tau;\theta,\phi) = -V_{\mathrm{soft}}^{\phi}(s_m) + V_{\mathrm{soft}}^{\phi}(s_n) + \sum_{t=m}^{n-1} \left( r(s_t, s_{t+1}) - \alpha \log \pi_\theta(s_{t+1} \mid s_t) \right)$$

$$\pi_\theta(s' \mid s) = P_F^\theta(s' \mid s) \qquad V_{\mathrm{soft}}^{\phi}(s) = \alpha \log F_\phi(s)$$

$$\Delta_{\mathrm{SubTB}}(\tau;\theta,\phi) = \log \frac{F_\phi(s_n) \prod_{t=m}^{n-1} P_B(s_t \mid s_{t+1})}{F_\phi(s_m) \prod_{t=m}^{n-1} P_F^\theta(s_{t+1} \mid s_t)}$$

$$\Delta_{\mathrm{SQL}}(s,s';\theta) = Q_{\mathrm{soft}}^{\theta}(s,s') - \left[ r(s,s') + \alpha \log \sum_{s'' \in \mathrm{Ch}(s')} \exp\left( \frac{1}{\alpha} Q_{\mathrm{soft}}^{\theta}(s',s'') \right) \right]$$

$$F_\theta(s) = \sum_{s'' \in \mathrm{Ch}(s)} \exp\left( \frac{1}{\alpha} Q_{\mathrm{soft}}^{\theta}(s,s'') \right) \qquad P_F^\theta(s' \mid s) \propto \exp\left( \frac{1}{\alpha} Q_{\mathrm{soft}}^{\theta}(s,s') \right)$$

$$\Delta_{\mathrm{DB}}(s,s';\theta) = \log \frac{F_\theta(s) P_F^\theta(s' \mid s)}{F_\theta(s') P_B(s \mid s')}$$

$$\Delta_{\mathrm{SQL}}(s,s';\theta) = Q_{\mathrm{soft}}^{\theta}(s,s') - \left[ r(s,s') + \alpha \log \sum_{s'' \in \mathrm{Ch}(s')} \exp\left( \frac{1}{\alpha} Q_{\mathrm{soft}}^{\theta}(s',s'') \right) \right]$$

$$\tilde{F}_\theta(s) = \sum_{s'' \in \mathrm{Ch}(s)} \exp\left( \frac{1}{\alpha} Q_{\mathrm{soft}}^{\theta}(s,s'') \right) \qquad P_F^\theta(s' \mid s) \propto \exp\left( \frac{1}{\alpha} Q_{\mathrm{soft}}^{\theta}(s,s') \right)$$

$$\Delta_{\mathrm{FL\text{-}DB}}(s,s';\theta) = \log \frac{\tilde{F}_\theta(s') P_B(s \mid s')}{\tilde{F}_\theta(s) P_F^\theta(s' \mid s)} - \frac{\mathcal{E}(s \to s')}{\alpha}$$

$$\Delta_{\pi\text{-}\mathrm{SQL}}(s,s';\theta) = \alpha \left[ \log \pi_\theta(s' \mid s) - \log \pi_\theta(s_f \mid s) + \log \pi_\theta(s_f \mid s') \right] - r(s,s')$$

$$\pi_\theta(s' \mid s) = P_F^\theta(s' \mid s)$$

$$\Delta_{\mathrm{M\text{-}DB}}(s,s';\theta) = \log \frac{\exp(-\mathcal{E}(s')/\alpha) P_B(s \mid s') P_F^\theta(s_f \mid s)}{\exp(-\mathcal{E}(s)/\alpha) P_F^\theta(s' \mid s) P_F^\theta(s_f \mid s')}$$

Terminal reward

Intermediate rewards