



Update **state value function** of each agent, according to

$$L_{Encoder}^m(\varphi) = \frac{1}{T_m} \sum_{t=0}^{T_m-1} \left[ R(o_t, a_t) + \gamma V_{\bar{\varphi}}(o_{t+1}^{i_m}) - V_{\varphi}(o_t^{i_m}) \right]^2 \quad \text{on genuine step trajectories.}$$

Update **policy function**, get the recomputed advantage using GAE trick on the trajectories of blue arrow.