# OMRL - Identifiability Challenges and Effective Data Collection Strategies [Dorfman et al. 2021]

- Extension to VariBAD Algorithm

  - Off-Policy Version

- Introducing Data-Collection/Augmenting Strategy for offline data

- Identifiability Problem: MDP ambiguity

  - Special to OMRL

  - Problem of identifying MDP

  - Tackle it with **policy replaying**: online

    - Use data-policy of another MDP in current to augment data

  - Tackle it with **reward relabeling**: offline

    - Use reward function of different MDP to augment data of current one
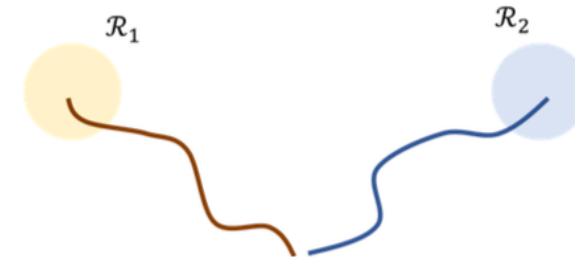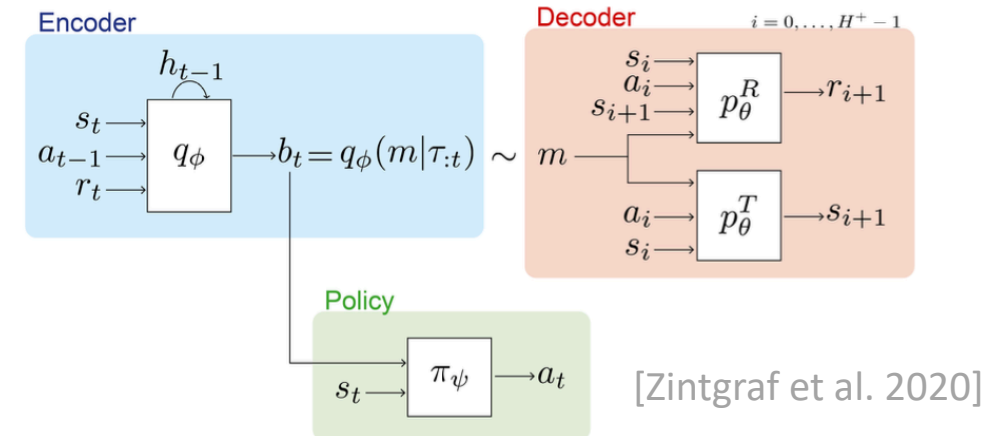
$\mathcal{R}_1$ $\mathcal{R}_2$

Figure 2: Reward ambiguity: from the two trajectories, it is impossible to know if there are two MDPs with different rewards (blue and yellow circles), or one MDP with rewards at both locations.

Encoder

$h_{t-1}$

$s_t \rightarrow$
$a_{t-1} \rightarrow \boxed{q_\phi} \rightarrow b_t = q_\phi(m|\tau_{:t}) \sim m$
$r_t \rightarrow$

Decoder $\quad i = 0, \ldots, H^+ - 1$

$s_i \rightarrow$
$a_i \rightarrow \boxed{p_\theta^R} \rightarrow r_{i+1}$
$s_{i+1} \rightarrow$

$a_i \rightarrow \boxed{p_\theta^T} \rightarrow s_{i+1}$
$s_i \rightarrow$

Policy

$s_t \rightarrow \boxed{\pi_\psi} \rightarrow a_t$

[Zintgraf et al. 2020]

# OMRL - Identifiability Challenges and Effective Data Collection Strategies [Dorfman et al. 2021]
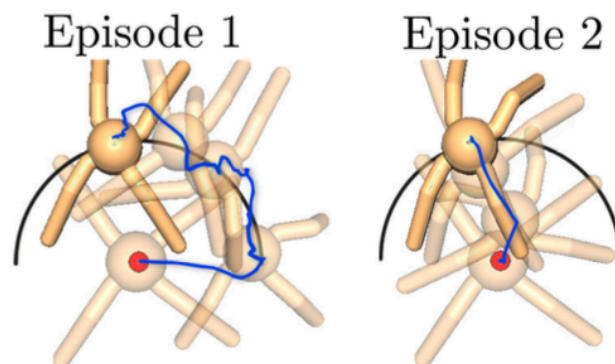


Figure 5: Ant-Semi-circle: trajectories from trained policy on a new goal. In the first episode the ant searches for the goal, and in the second one it directly moves toward the goal it has previously found. This search behavior is different from the goal-reaching behaviors that dominate the data.
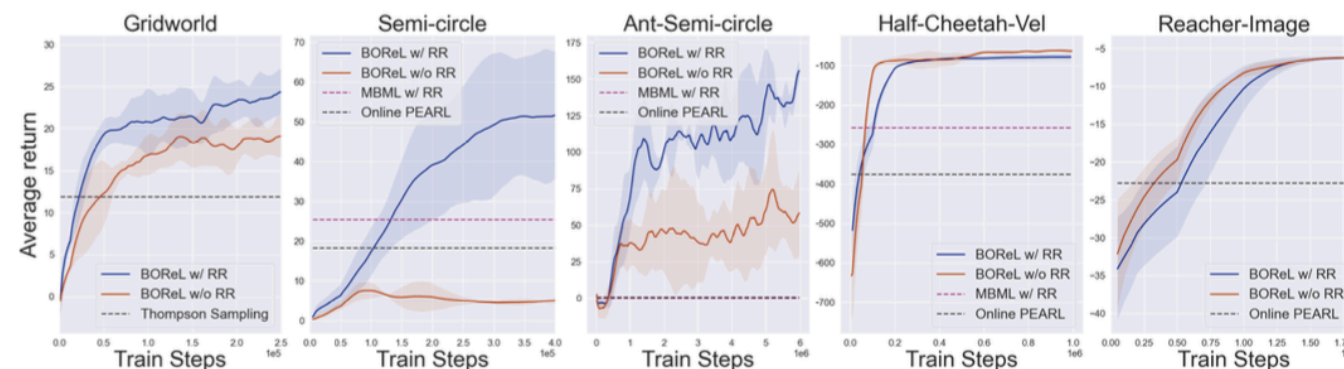


Figure 3: Offline performance on domains with varying rewards. We compare BOReL with and without reward relabeling (blue and red, respectively) with Thompson sampling baselines – calculated exactly in Gridworld, and using online PEARL and offline MBML for the other domains. Full training curves for baselines appear in the supplementary; here we plot only the best performance.