

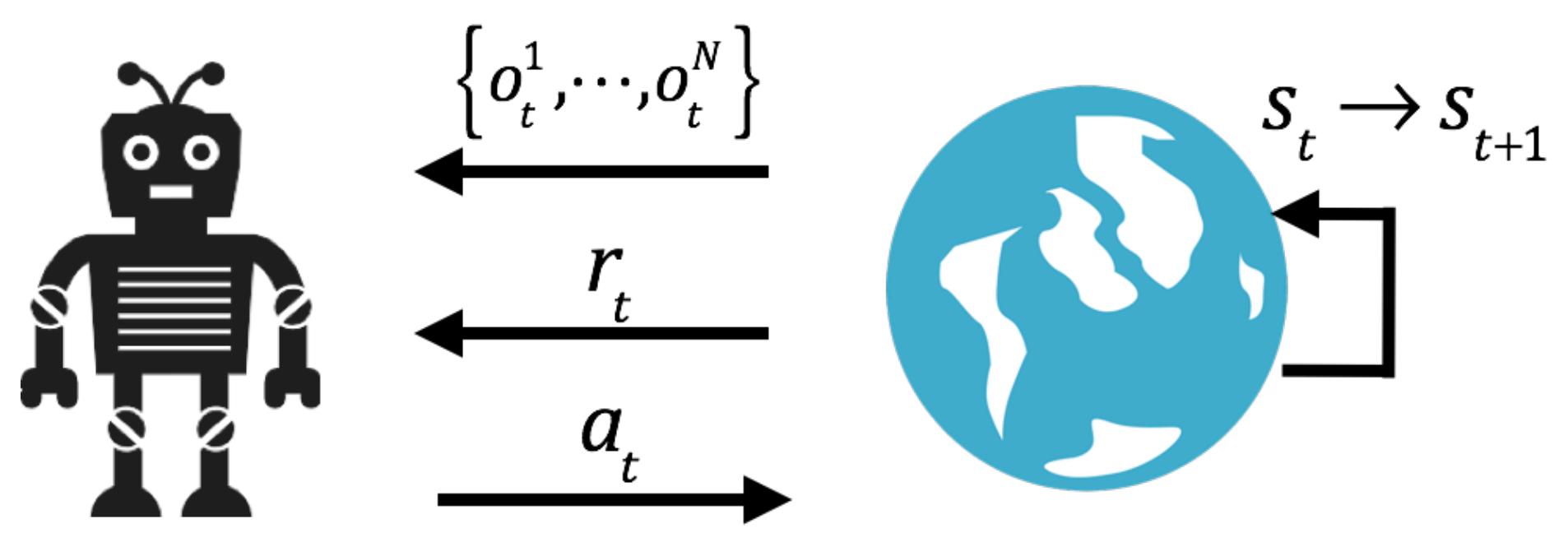
Multi-View Reinforcement Learning

Minne Li*, Lisheng Wu*, Haitham Bou Ammar[^], Jun Wang
 University College London
 *Equal contributions. [^]Honorary Lecturer at University College London

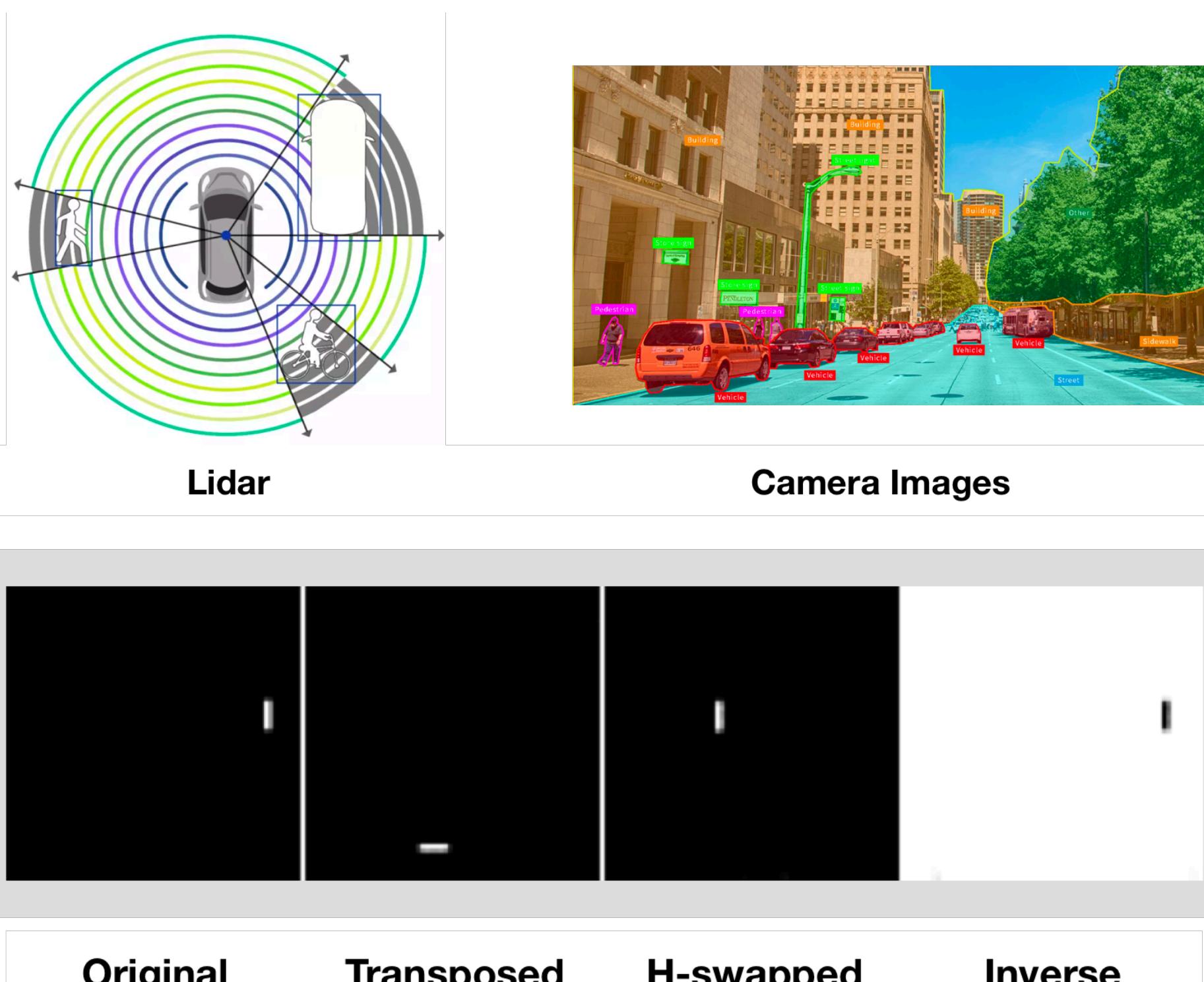


Motivation

This work is concerned with multi-view reinforcement learning (MVRL), which allows for decision making when agents share common dynamics but adhere to different observation models.



Examples



In Atari Pong, balls and paddles always follow the same physical principles of persistence, continuity, cohesion and solidity, although they are observed through totally different frame of reference.

Multi-view MDP

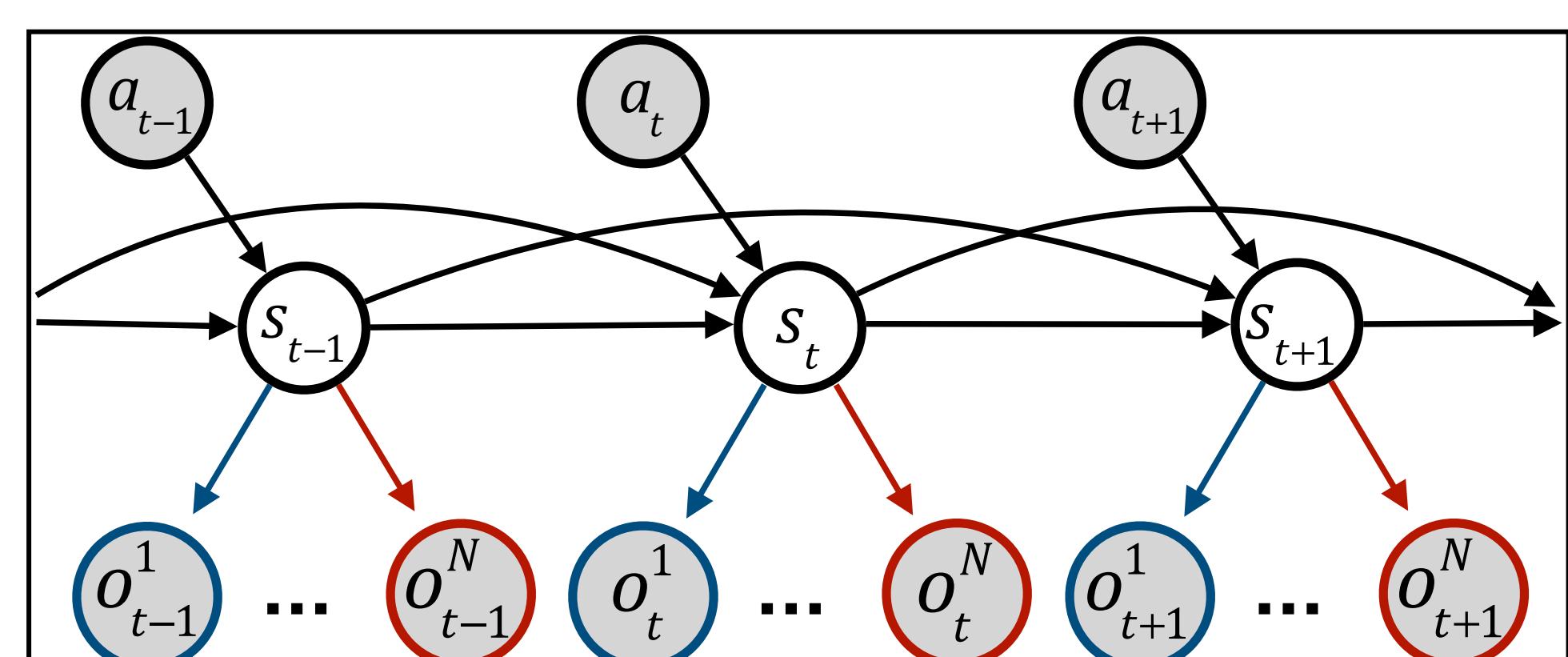
Given $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{O}_1, \mathcal{P}_{obs}^1, \dots, \mathcal{O}_N, \mathcal{P}_{obs}^N \rangle$, a history of heterogeneous observations (and potentially actions) $\mathcal{H}_t = \{\mathbf{o}_1^{i_1}, \dots, \mathbf{o}_1^{i_1}\}$, where $\mathbf{o}_k^{i_k} \sim \mathcal{P}_{obs}^{i_k}(\cdot | \mathbf{s}_k)$ for $k \in \{1, \dots, T\}$ and $i_k \in \{1, \dots, N\}$, the goal is to find optimal action policy π^* to maximize the return:

$$G_T(\tau^M) = \sum_t \gamma^t \mathcal{R}(s_t, a_t),$$

where $\tau^M = [s_1, \mathbf{o}_1^{i_1}, \mathbf{a}_1, \dots, s_T, \mathbf{o}_T^{i_T}]$,

by $\max_{\pi^M} \mathbb{E}_{\tau^M} [G_T(\tau^M)]$, where $\tau^M \sim p_{\pi^M}(\tau^M)$ and

$$p_{\pi^M}(\tau^M) = \mathcal{P}_0(s_1) \mathcal{P}_{obs}^{i_1}(o_1^{i_1} | s_1) \pi^M(a_1 | \mathcal{H}_1) \prod_{t=2}^T \mathcal{P}_{obs}^{i_t}(o_t^{i_t} | s_t) \mathcal{P}(s_t | s_{t-1}, a_{t-1}) \pi^M(a_{t-1} | \mathcal{H}_{t-1})$$



Graphical model of multi-view learning.

Model-free Solutions (MV-MF)

Updating the policy parameters ω as

$$\omega^{k+1} = \omega^k + \eta^k \frac{1}{M} \sum_{j=1}^M \sum_{t=1}^T \nabla_{\omega} \log \pi^M(a_t^j | \mathcal{H}_t^j) (\mathcal{R}(s_t^j, a_t^j) - \mathcal{B}_{\phi}(\mathcal{H}_t^j))$$

where $\mathcal{B}_{\phi}(\mathcal{H}_t)$ is the observation-based baseline.

Problems of MV-MF: sample-inefficient.

Model-based Solutions

By maximizing the model evidence of the observation outcome

$$p(\mathbf{o}_{1:T}^{MC}) = \prod_{i_t=1}^N p(\mathbf{o}_{i_t}^{i_t}) = \prod_{i_t=1}^N \int_{s_1} \dots \int_{s_T} \underbrace{\dots}_{\text{latent variables}} \underbrace{\dots}_{\text{marginalization}} d\mathbf{s}_1 \dots d\mathbf{s}_T,$$

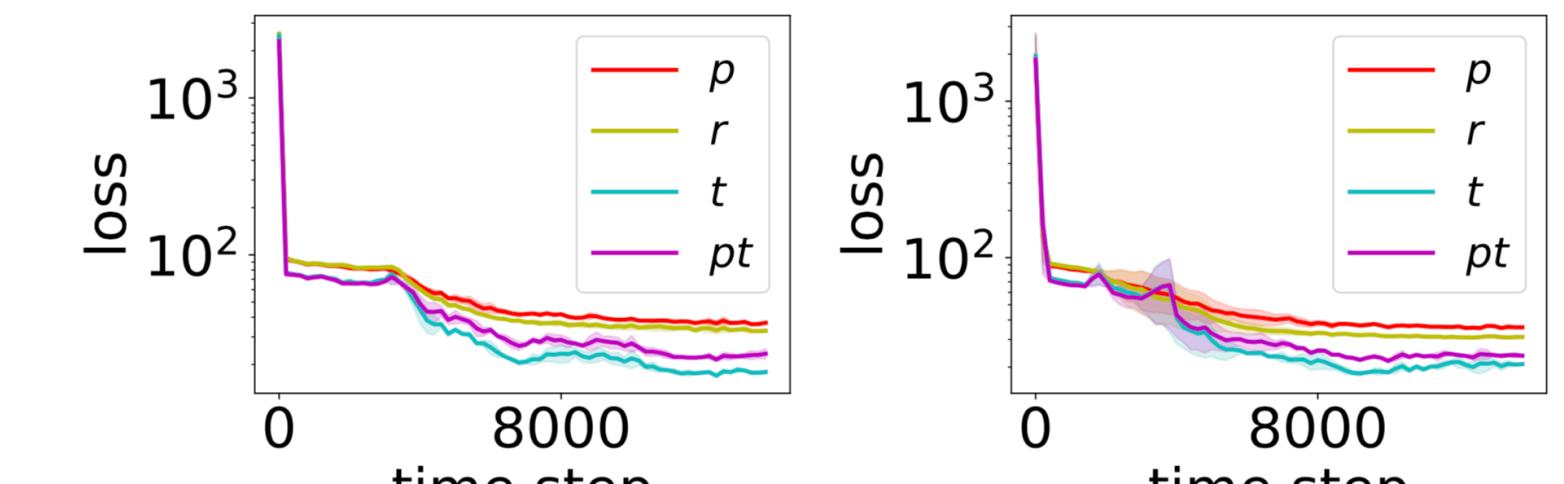
where $\mathbf{o}_{1:T}^{MC}$ collects all multi-view observations and actions across time-steps.

The above can be derived as maximizing the multi-view ELBO (see paper for details)

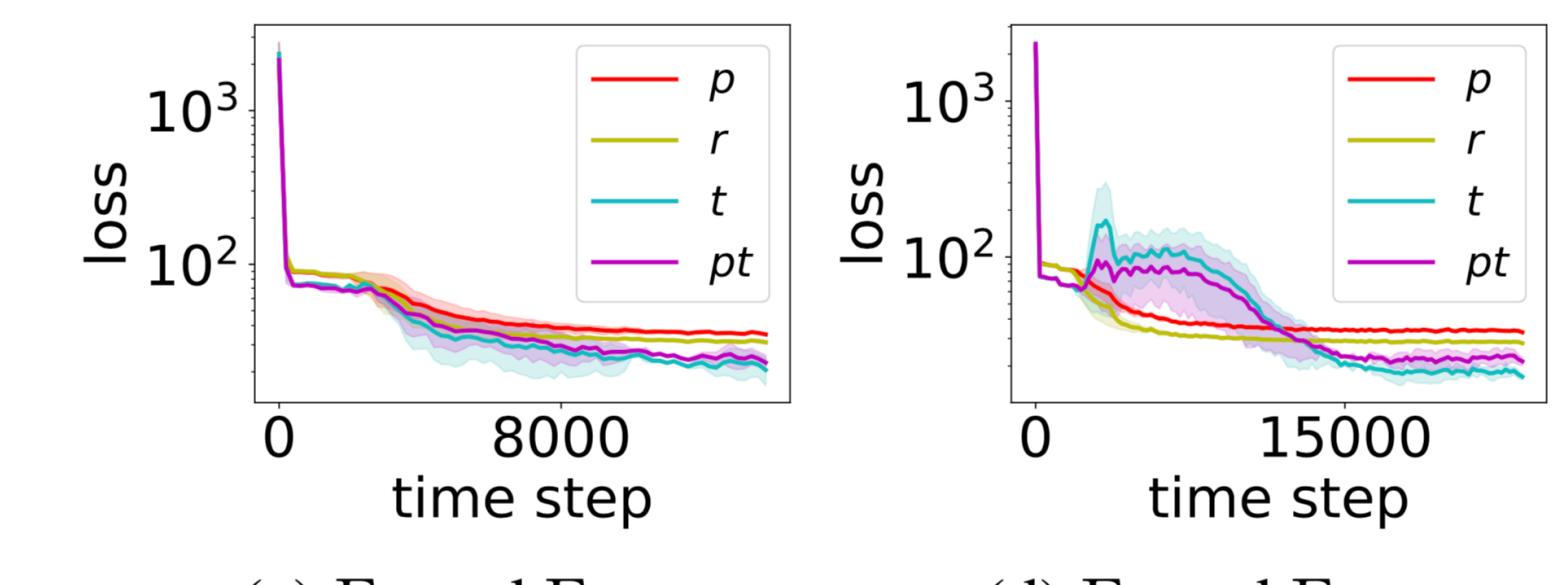
$$\max_{\theta_m, \phi} \sum_{i_t=1}^N \sum_{t=1}^T \left[\mathbb{E}_{q_{\phi^{i_t}}(s_t | \mathcal{H}_t)} [\log p_{\theta^{i_t}}(o_{i_t}^{i_t} | s_t)] - \text{KL} (q_{\phi^{i_t}}(s_t | \mathcal{H}_t) || p_{\theta^{i_t}}(s_t | \mathcal{H}_t)) \right]$$

We can then have model-based MVRL through cross-view policy transfer (MV-PT), which only require very few data from a specific view to train the multi-view model. This can then be used for action selection by: 1) inferring the corresponding latent state, and 2) feeding the latent state into the policy learned from another view with greater accessibility.

Modeling Results

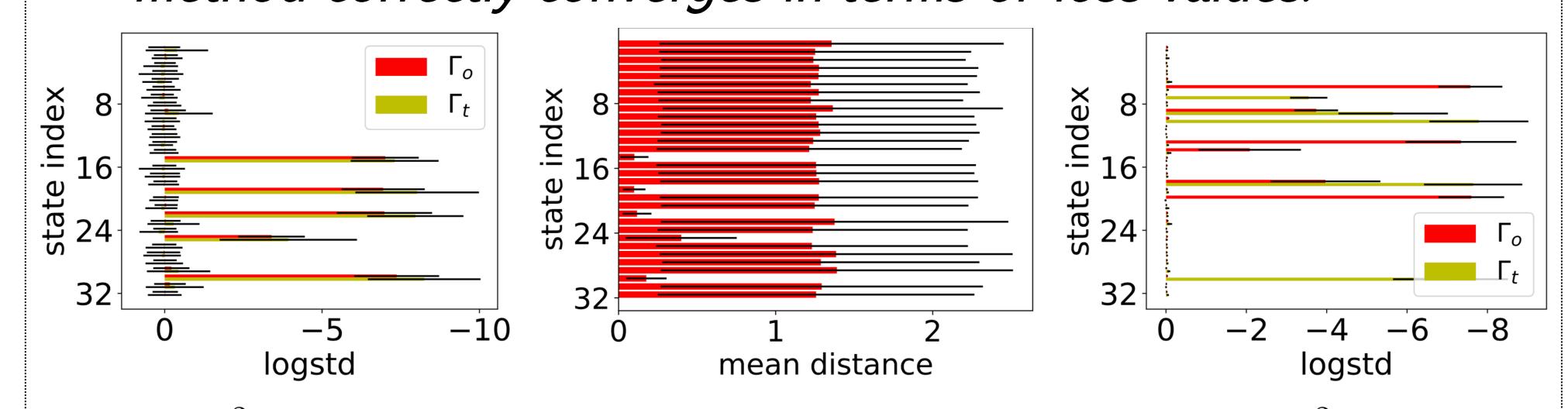


(a) Γ_o and Γ_t (b) Γ_o and Γ_h



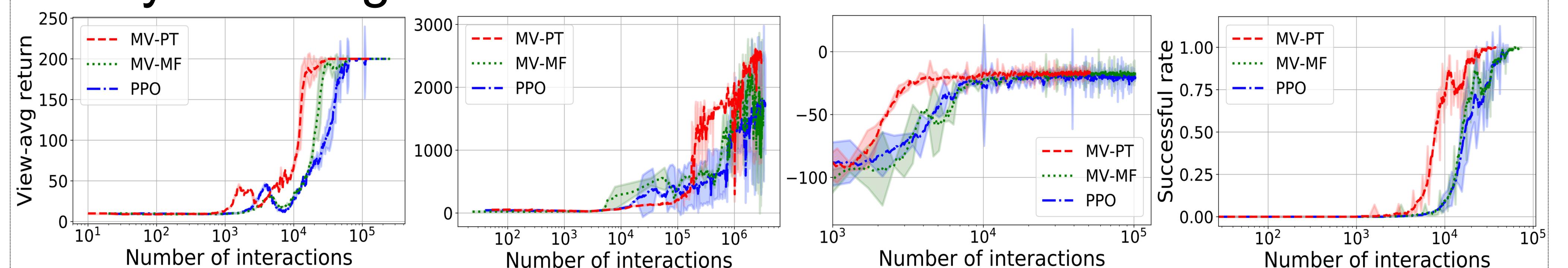
(c) Γ_o and Γ_c (d) Γ_o and Γ_m

Training multi-view models on Atari Pong with four view variations (t, h, c, m). Legend: p : prediction loss; r : reconstruction loss; t : transformation loss; pt : predicted transformation loss. These results demonstrate that our method correctly converges in terms of loss values.



Difference between inferred latent states from the original view and the transposed view. Results demonstrating that our method is capable of learning key elements -- a property essential for multi-view dynamics learning. These results also demonstrate that extracting such key-elements is challenging for world-models.

Policy Learning Results



(a) Cartpole Policy Learning (b) Hopper Policy Learning (c) RACECAR Policy Learning (d) Parking Policy Learning
 Policy learning results demonstrating that our method outperforms others in terms of sample complexities. We demonstrate these results on 4 different dynamical systems. In all the above, our method can outperform state-of-the-art significantly.