

Introduction

We developed an autonomous framework that uses unsupervised manifold alignment to learn inter-task mappings and effectively transfer samples between different task domains. Our results demonstrate the success of our approach for transfer between highly dissimilar control tasks (e.g., from cart-poles to quadrotors), and show that transfer quality is positively correlated with manifold alignment quality.

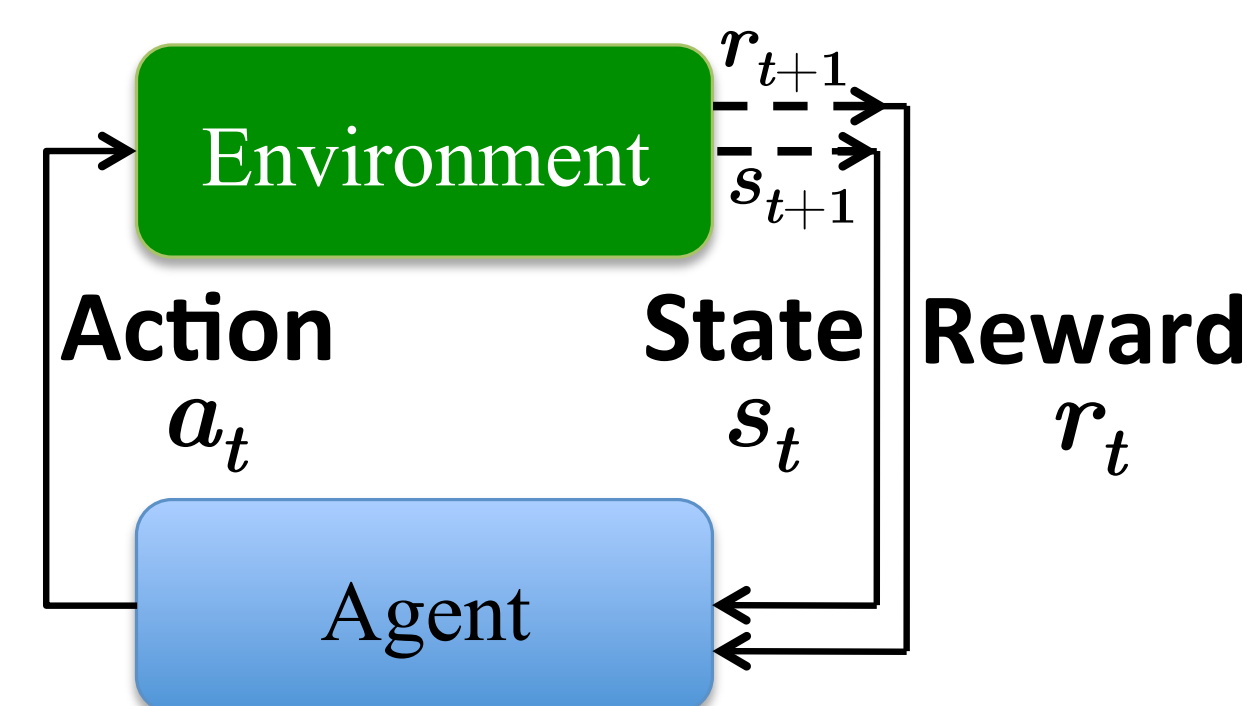
Motivation:

- Transfer learning enables rapid training of a control policy for a new target task by reusing knowledge from other source tasks.
- In the case of multiple task domain, an **inter-task mapping** χ is needed to map knowledge between tasks.
 - χ maps state-action-next-state triplets from the source task to the target task, which can be used for policy initialization.

Background: Reinforcement Learning

Reinforcement Learning (RL) problems are formalized as Markov Decision Processes (MDPs): $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}_0, \mathcal{P}, r \rangle$, where

- $\mathcal{S} \in \mathbb{R}^d$ is the state space
- $\mathcal{A} \in \mathbb{R}^m$ is the action space
- \mathcal{P}_0 is the initial state distribution
- $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the transition probability function
- $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward function.



Goal: Learn an optimal policy $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$ that maximizes the total discounted reward.

Background: Policy Gradient RL

In Policy Gradient (PG) methods, the policy is parameterized by $\theta \in \mathbb{R}^d$ and a vector of state features Φ . The goal is to maximize

$$\mathcal{J}(\theta) = \int_{\mathbb{T}} p_{\theta}(\tau) \mathcal{R}(\tau) d\tau, \text{ where}$$

$$p_{\theta}(\tau) = \mathcal{P}_0(s_1) \prod_{t=1}^H \mathcal{P}(s_{t+1}|s_t, a_t) \pi(a_t|s_t) \quad \leftarrow \text{Probability of trajectory}$$

$$\mathcal{R}(\tau) = \frac{1}{H} \sum_{t=1}^H r(s_{t+1}, a_t, s_t) \quad \leftarrow \text{Reward of trajectory}$$

Problem: PG suffers from high computational and sample complexities.

Unsupervised Manifold Alignment for Learning the Inter-Task Mapping χ_S

Phase I: Learning the inter-task mapping χ_S via unsupervised manifold alignment

1. Sample *a.) optimal* trajectories from the source task using $\pi_{(S)}^*$ and *b.) random* trajectories from the target task.
2. Flatten all trajectories and construct a k -NN graph to capture the local geometry of the states in both the source and target tasks.
3. Identify a shared representation between the source and target tasks that captures local state transition dynamics by optimizing

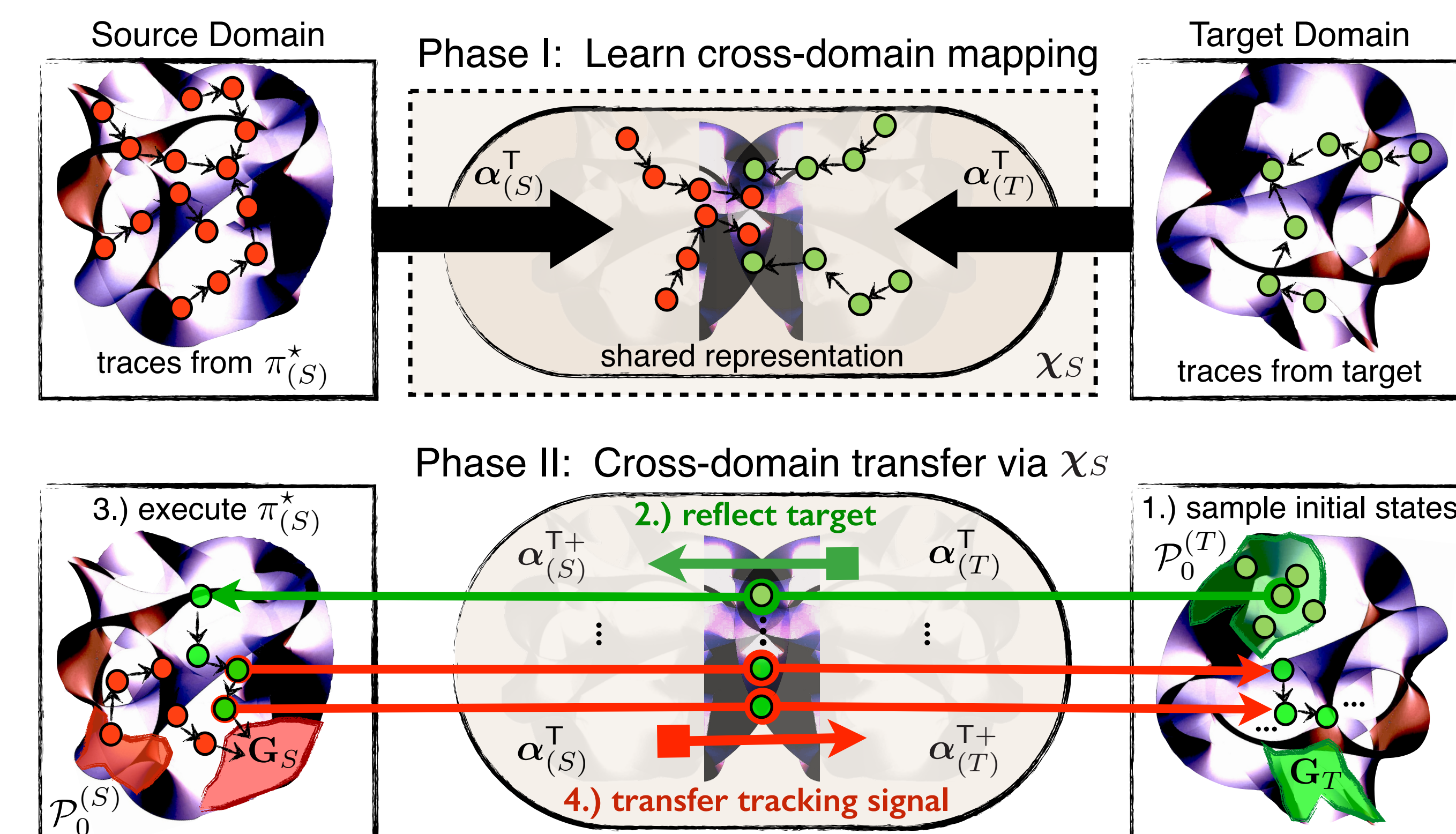
$$\mathcal{J}(\alpha_{(S)}, \alpha_{(T)}) = \underbrace{\mu \sum_{i,j} \left(\alpha_{(S)}^T s_i^{(S)*} - \alpha_{(T)}^T s_j^{(T)*} \right)^2 W_{i,j}}_{\text{cross-task connection}} + \underbrace{0.5 \sum_{i,j} \left(\alpha_{(S)}^T s_i^{(S)*} - \alpha_{(S)}^T s_j^{(S)} \right)^2 W_{S(S)}^{i,j}}_{\text{source task geometry}} + \underbrace{0.5 \sum_{i,j} \left(\alpha_{(T)}^T s_i^{(T)*} - \alpha_{(T)}^T s_j^{(T)} \right)^2 W_{S(T)}^{i,j}}_{\text{target task geometry}}$$

where the W 's are the weighted adjacency matrices, s are the states, the α 's are the projections into the shared latent space, and the superscripts or subscripts of (S) and (T) denote whether these variables correspond to the source or target task, respectively.

4. The inter-task mapping $\chi_S = \alpha_{(T)}^+ \alpha_{(S)}^T [\cdot]$.

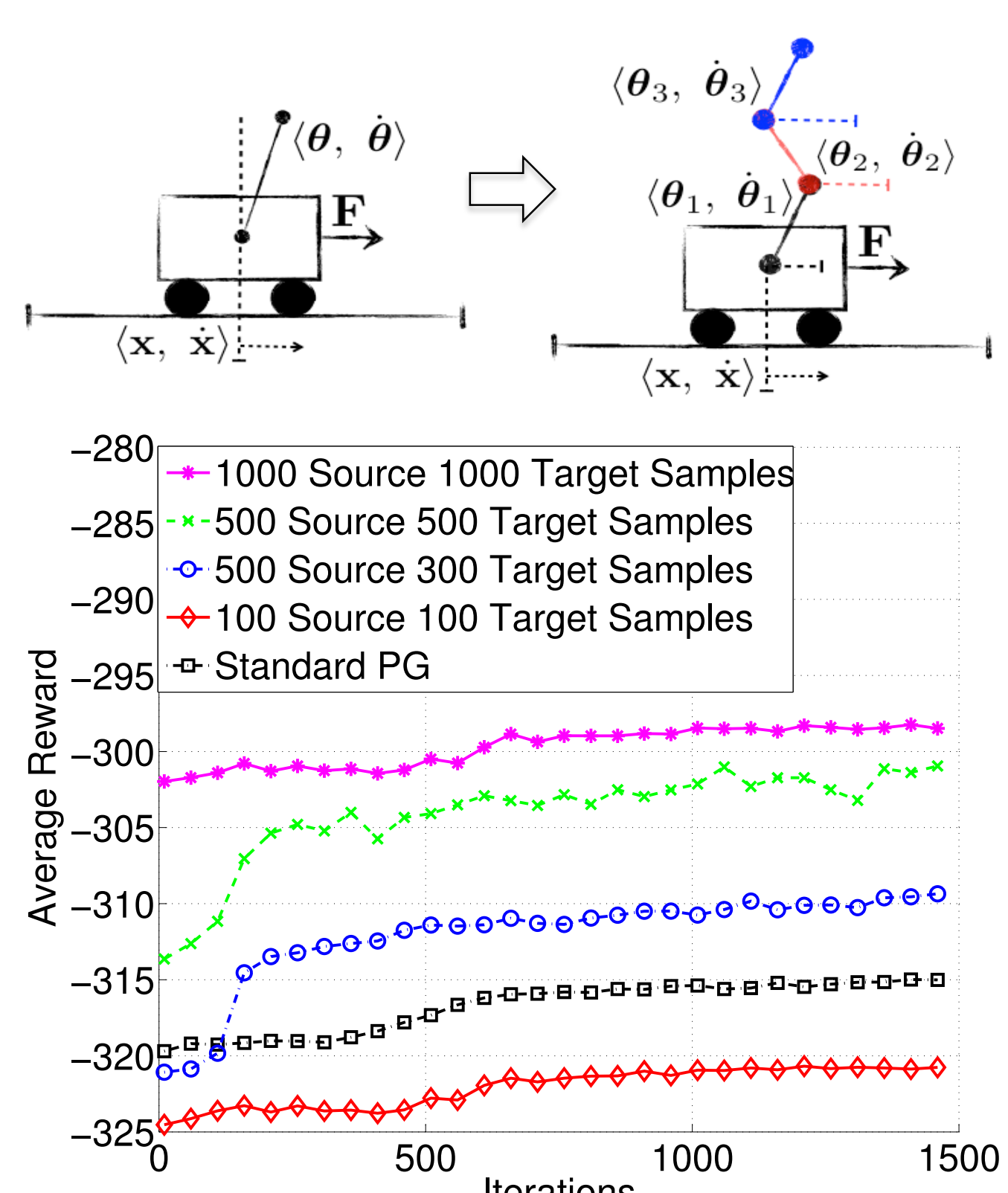
Phase II: Initialize the target task's policy via transfer

1. Sample initial target states $s_0^{(T)} \sim \mathcal{P}_0^{(T)}$.
2. Project initial target states $s_0^{(T)}$ to the source task via χ_S .
3. Execute $\pi_{(S)}^*$ from these projected states, yielding optimal trajectories $\tilde{\tau}_{(S)}$.
4. Transfer optimal source trajectories $\tilde{\tau}_{(S)}$ to the target task via χ_S^+ , yielding target trajectories $\tilde{\tau}_{(T)}$.
5. Initialize target task policy $\pi_{(T)}$ from $\tilde{\tau}_{(T)}$, yielding $\theta_{(T)}^{(0)}$. Improve $\pi_{(T)}$ using standard policy gradient methods.

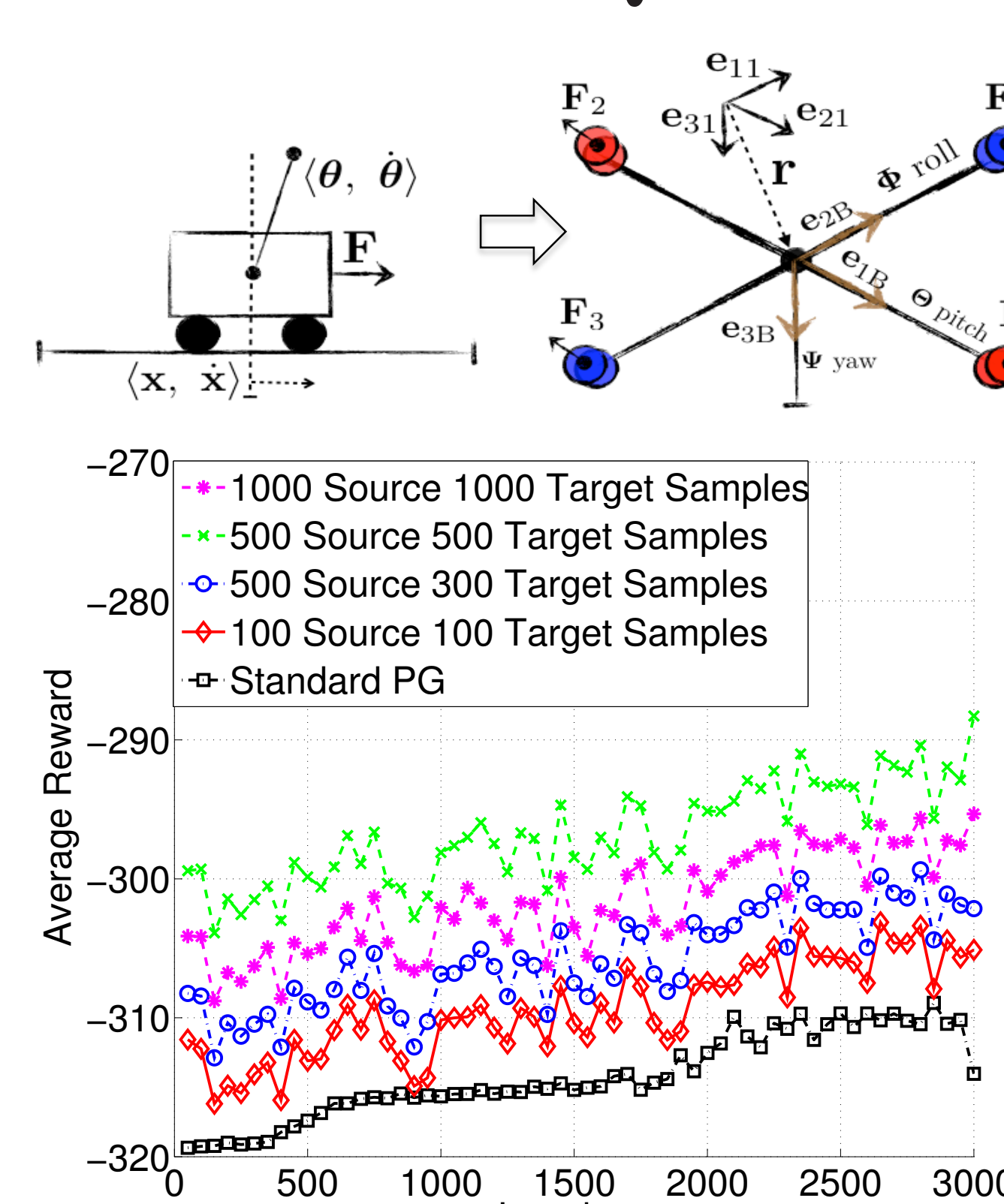


Selected Results of Transfer Between Different Dynamical Systems

Cart-Pole to 3-link Cart-Pole

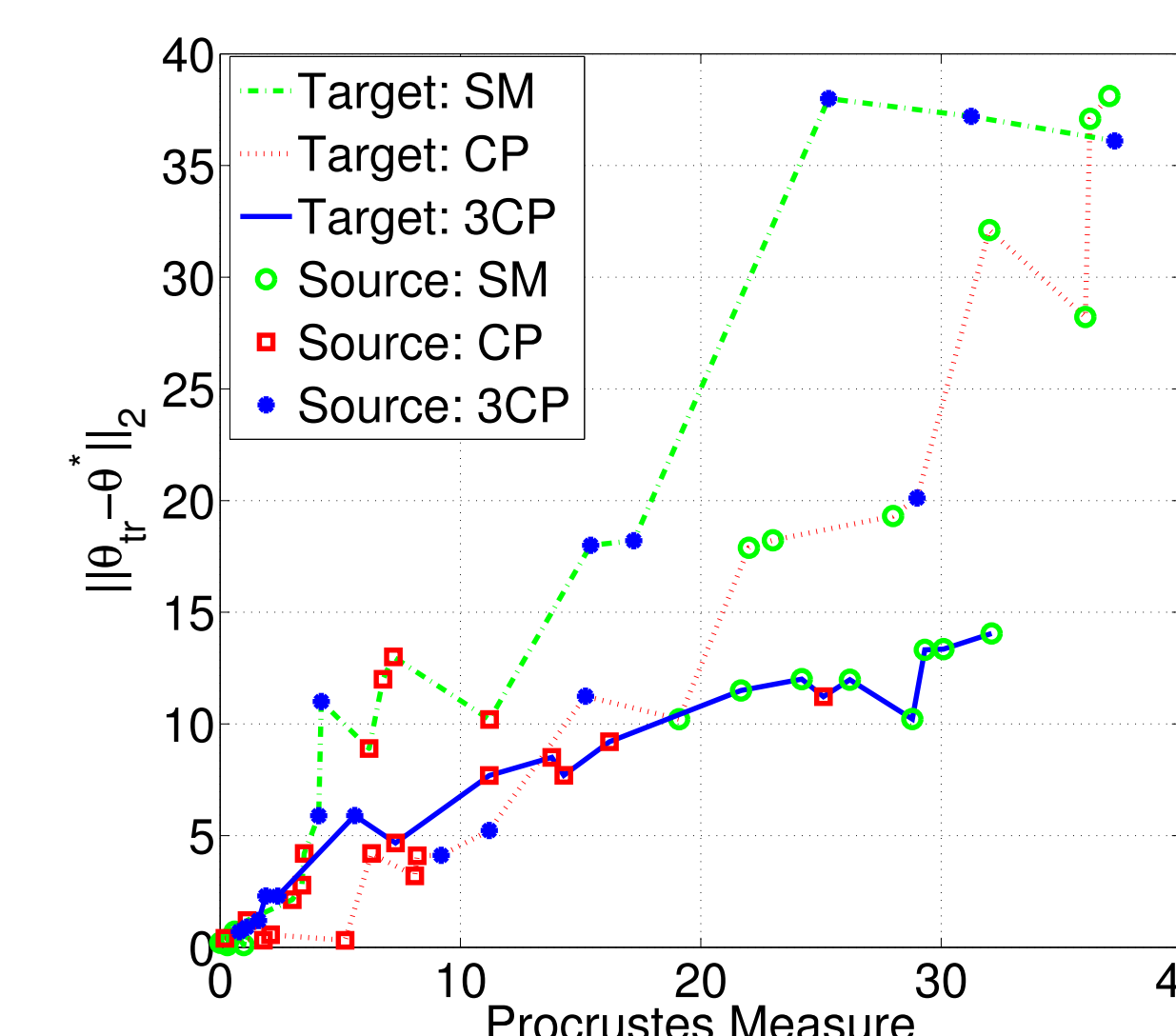


Cart-Pole to Quadrotor



Predicting Success of Cross-Domain Transfer

- Transfer quality ($\|\theta_{tr} - \theta^*\|_2$) is positive correlated with manifold alignment quality (Procrustes measure).
- Manifold alignment quality may indicate when our approach to cross-domain transfer is likely to succeed.



Unsupervised manifold alignment enables robust cross-domain transfer between highly dissimilar tasks