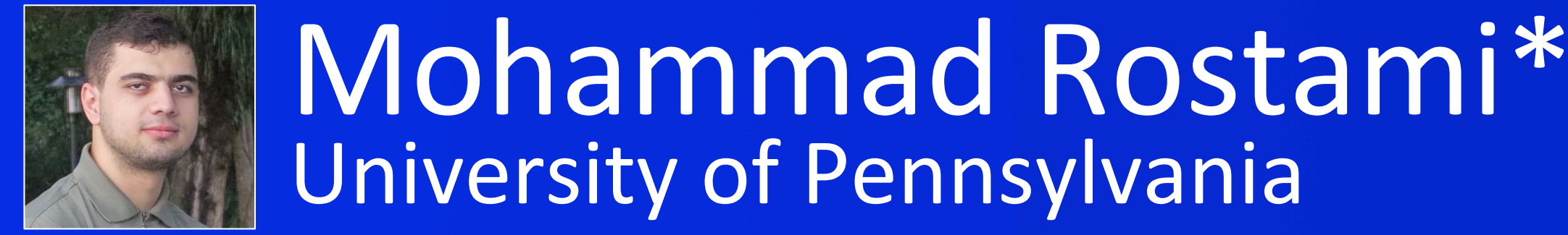
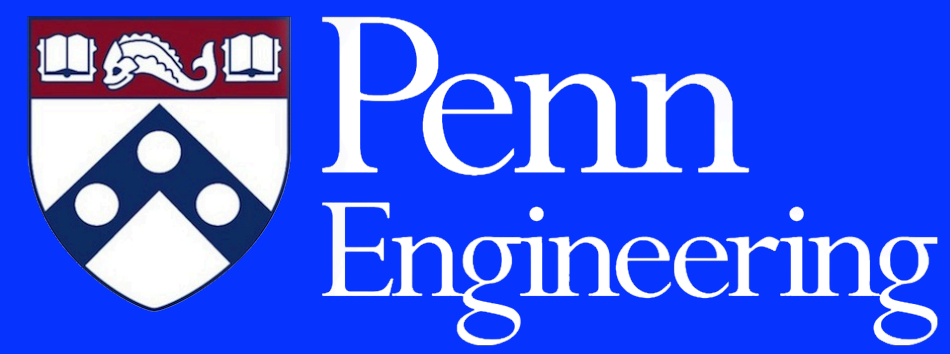


Using Task Features for Zero-Shot Knowledge Transfer in Lifelong Learning



Summary

Knowledge transfer between tasks requires an accurate estimate of the inter-task relationships, which is inefficient in lifelong learning settings. We develop a **lifelong reinforcement learning** method that incorporates **high-level task descriptors** to model the inter-task relationships.

- Improves the performance of the learned task policy
- Accurately predicts the policy for a new task via zero-shot learning, given only the task description

Motivation



Lifelong learning accelerates training of each consecutive new task by building upon previously acquired knowledge via transfer

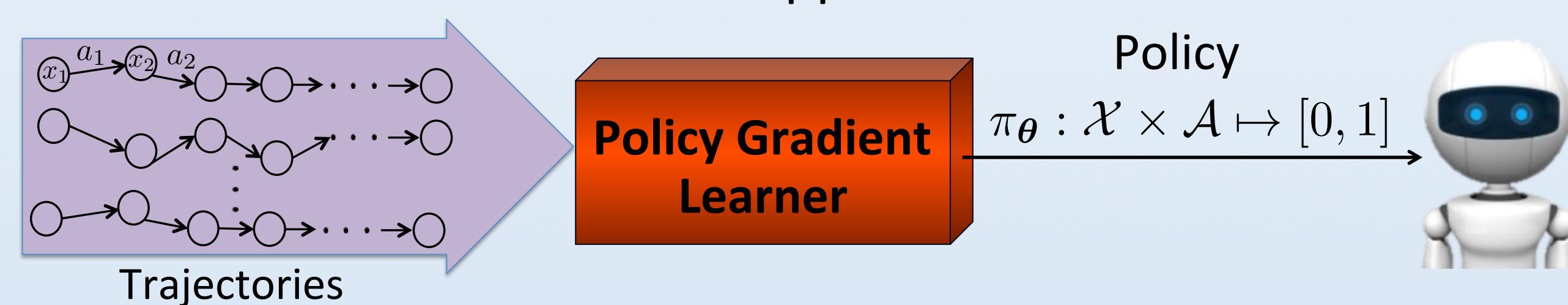
- Relevant knowledge/tasks must be identified before transfer can occur
- Requires interacting with the new task (i.e., sampling trajectories, learning, etc.) to characterize it

Alternative Idea: Can we use a high-level description of the task to identify relevant knowledge for transfer in lifelong learning?

Example task descriptor: physical specification of a quadrotor

Background: Policy Gradient (PG) Methods

- Agent interacts with environment, taking consecutive actions
- PG methods support continuous state and action spaces
- Have shown recent success in applications to robotic control



Goal: find policy π_θ that minimizes $\mathcal{J}(\theta) = \int_{\mathcal{T}} p_\theta(\tau) \mathcal{R}(\tau) d\tau$

$$p_\theta(\tau) = p_0(\mathbf{x}_0) \prod_{h=1}^H p(\mathbf{x}_{h+1} | \mathbf{x}_h, \mathbf{a}_h) \pi_\theta(\mathbf{a}_h | \mathbf{x}_h)$$

$$\mathcal{R}(\tau) = \frac{1}{H} \sum_{h=0}^H r_{h+1}$$

Sharing Knowledge Between Multiple Tasks

- Policy for task t : $\pi_{\theta^{(t)}} : \mathcal{X} \times \mathcal{A} \mapsto [0, 1]$
- Factor the policy as $\theta^{(t)} = \mathbf{L} \mathbf{s}^{(t)}$

$$\text{Obj. Fn. } e_T(\mathbf{L}) = \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{s}^{(t)}} \left[\underbrace{-\mathcal{J}(\theta^{(t)})}_{\text{fit to each task}} + \underbrace{\mu \|\mathbf{s}^{(t)}\|_1}_{\text{sparsity of coefficients}} + \underbrace{\lambda \|\mathbf{L}\|_F^2}_{\text{regularize basis complexity}} \right]$$

Batch Optimization (all tasks are given)

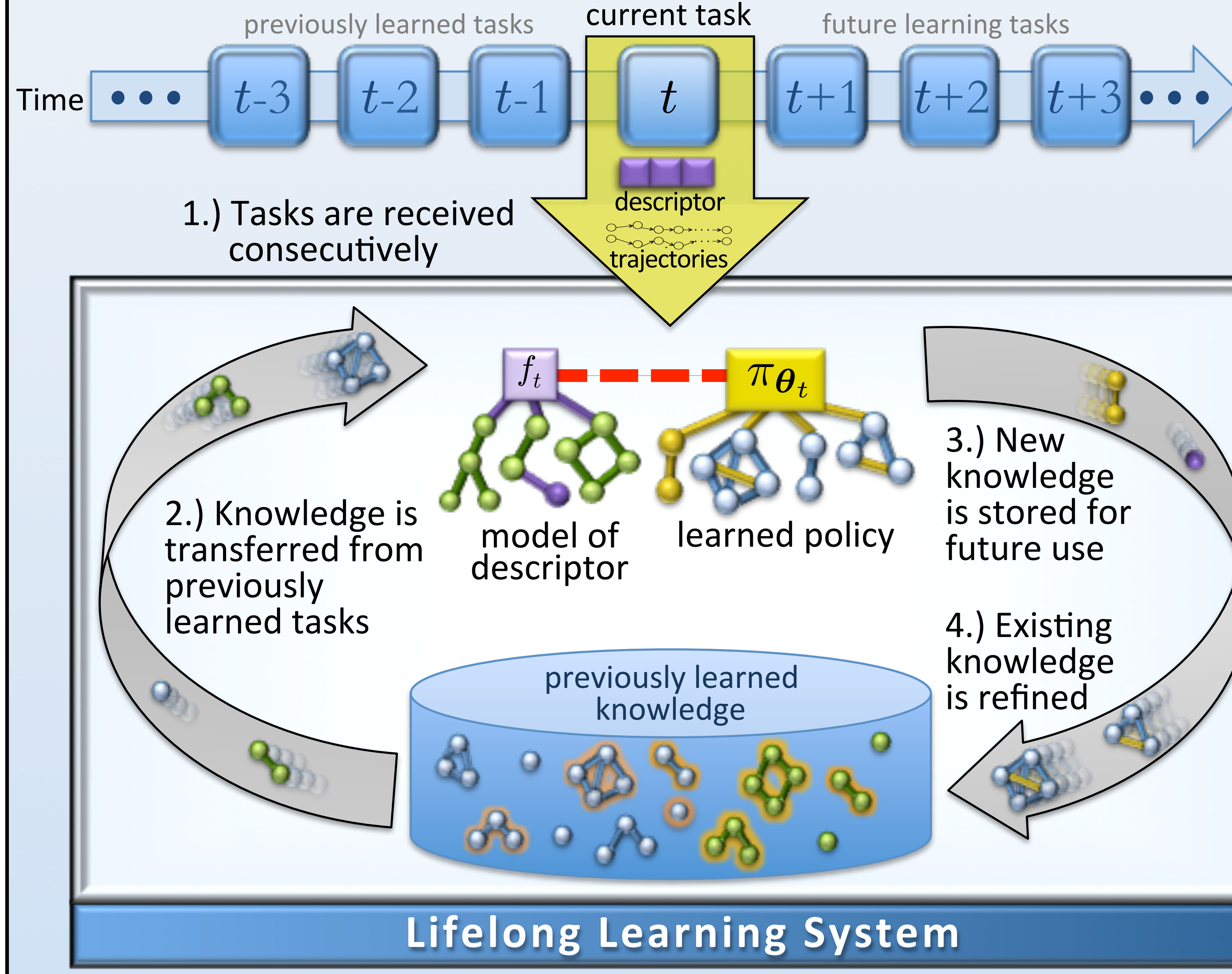
Multi-Task Learning

Online Optimization (tasks arrive consecutively)

Lifelong Learning (PG-ELLA)

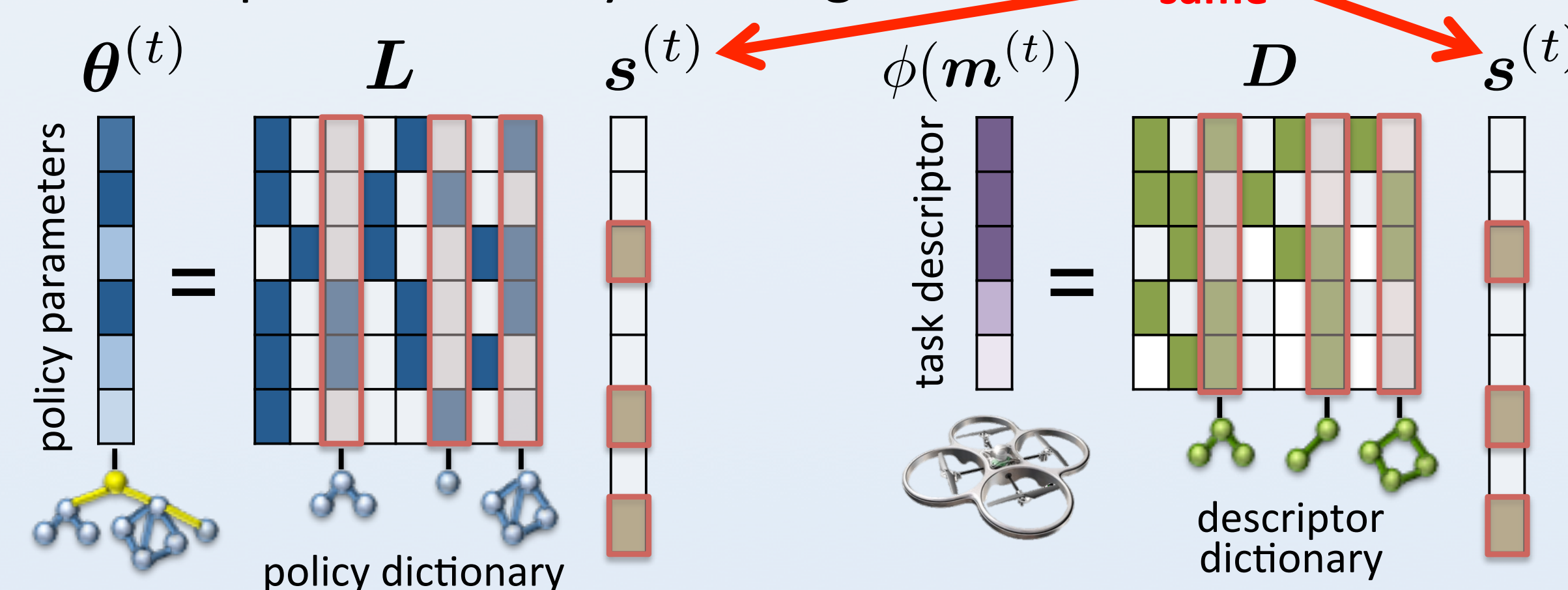
[Bou Ammar, Eaton, et al., ICML '14]

Lifelong Machine Learning with Task Descriptors



Incorporating Task Descriptors into Lifelong Learning

Key Idea: Relate policy parameters and task descriptors via coupled dictionary learning



Update the objective function to learn both dictionaries:

$$\text{Obj. Fn. } \min_{L, D, S} \frac{1}{T} \sum_t \left[\underbrace{\|\alpha^{(t)} - \mathbf{L} \mathbf{s}^{(t)}\|_{\Gamma^{(t)}}^2}_{\text{policy fit}} + \underbrace{\rho \|\phi(m^{(t)}) - \mathbf{D} \mathbf{s}^{(t)}\|_2^2}_{\text{descriptor fit}} + \underbrace{\mu \|\mathbf{s}^{(t)}\|_1}_{\text{sparsity}} + \underbrace{\lambda (\|\mathbf{L}\|_F^2 + \|\mathbf{D}\|_F^2)}_{\text{complexity}} \right]$$

Multi-Task Learning (TaDeMTL)

- Fit via alternating optimization

Lifelong Learning (TaDeLL)

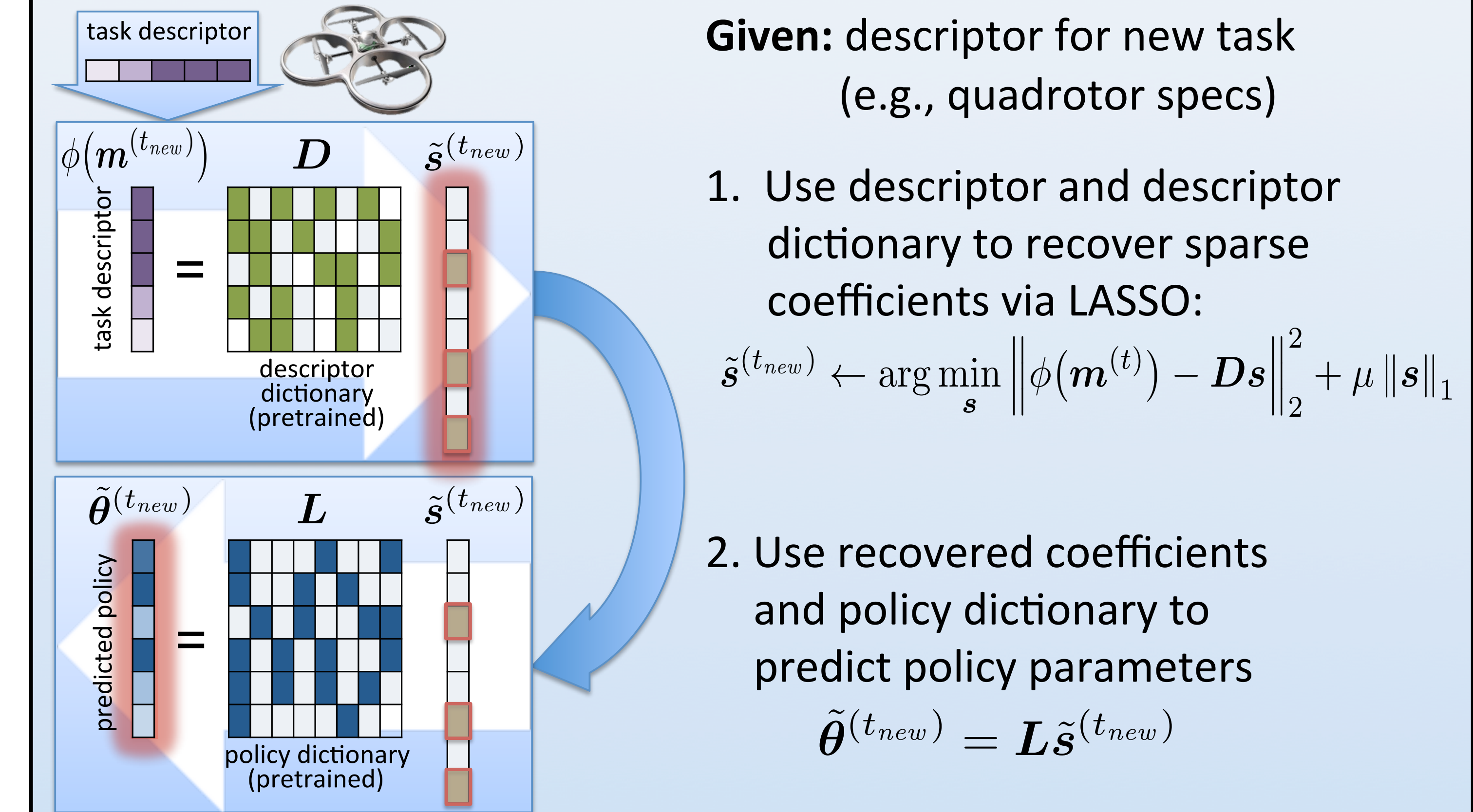
- Merge \mathbf{L} and \mathbf{D} into single dictionary \mathbf{K}
- Estimate policy $\alpha^{(t)}$ via single-task learning
- Sparse code estimated policy and descriptor in \mathbf{K}
- Update \mathbf{L} and \mathbf{D}

Algorithm 1 TaDeLL (k, λ, μ)

```

1:  $T \leftarrow 0$ 
2:  $\mathbf{L} \leftarrow \text{RandomMatrix}_{d,k}, \mathbf{D} \leftarrow \text{RandomMatrix}_{m,k}$ 
3: while some task  $(\mathcal{Z}^{(t)}, \phi(m^{(t)}))$  is available do
4:   if isNewTask( $\mathcal{Z}^{(t)}$ ) then
5:      $T \leftarrow T + 1$ 
6:      $\mathbb{T}^{(t)} \leftarrow \text{sampleRandomTrajectories}(\mathcal{Z}^{(t)})$ 
7:   else
8:      $\mathbb{T}^{(t)} \leftarrow \text{sampleTrajectories}(\mathcal{Z}^{(t)}, \pi_{\alpha^{(t)}})$ 
9:   end if
10:  Compute  $\alpha^{(t)}$  and  $\Gamma^{(t)}$  from  $\mathbb{T}^{(t)}$ 
11:   $\mathbf{s}^{(t)} \leftarrow \arg \min_{\mathbf{s}} \|\beta^{(t)} - \mathbf{K} \mathbf{s}\|_{A^{(t)}}^2 + \mu \|\mathbf{s}\|_1$ 
12:   $\mathbf{L} \leftarrow \text{updateL}(\mathbf{L}, \mathbf{s}^{(t)}, \alpha^{(t)}, \Gamma^{(t)}, \lambda)$ 
13:   $\mathbf{D} \leftarrow \text{updateD}(\mathbf{D}, \mathbf{s}^{(t)}, \phi(m^{(t)}), \rho, \mathbf{I}_{d_m}, \lambda)$ 
14:  for  $t \in \{1, \dots, T\}$  do:  $\theta^{(t)} \leftarrow \mathbf{L} \mathbf{s}^{(t)}$ 
15: end while
    
```

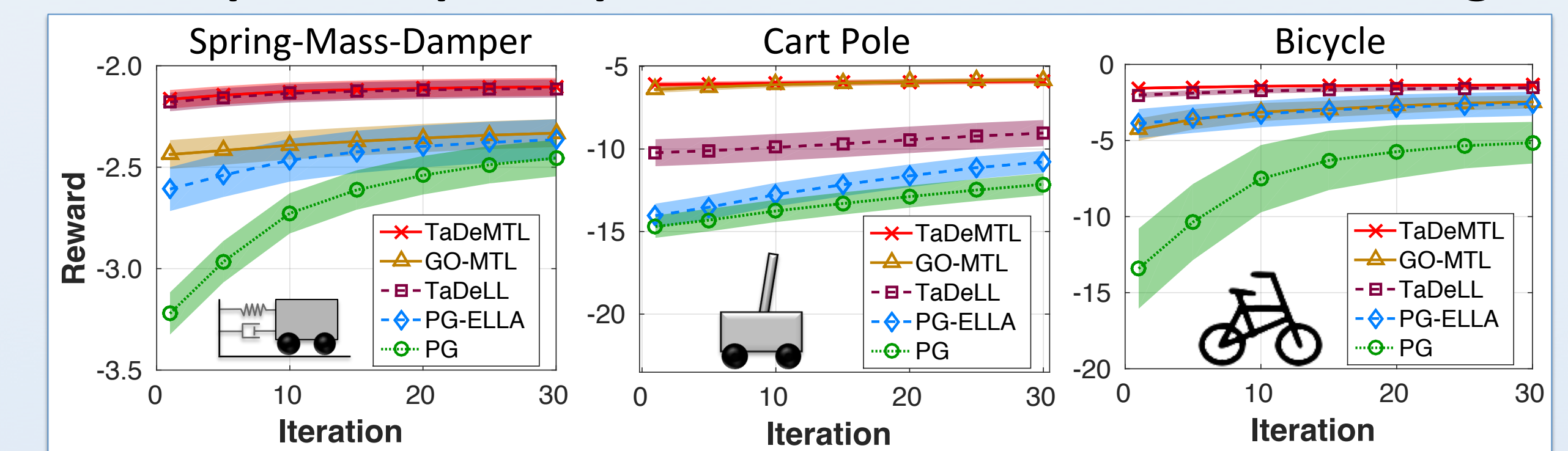
Zero-Shot Transfer via Task Descriptors



Experimental Results on Dynamical Systems

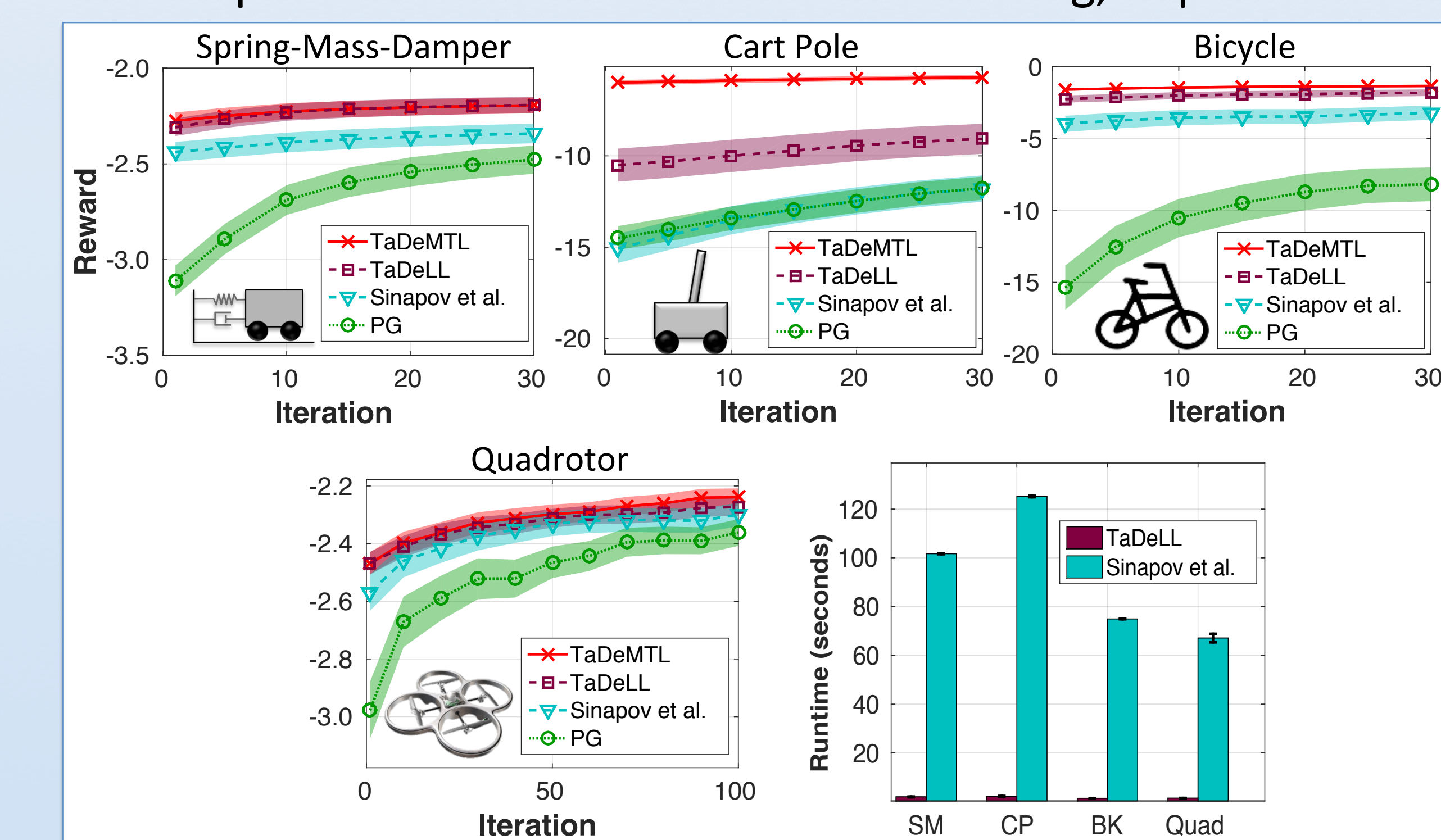
- Train on 40 different consecutive control tasks, transfer to new tasks

Task descriptors improve policies from multi-task and lifelong learning



Effective zero-shot transfer to new tasks

- Zero-shot policies used as warm-start for learning, improved via PG



Acknowledgements and Notes

This research was supported by ONR grant #N00014-11-1-0139 and AFRL grant #FA8750-14-1-0069.

* Authors contributed equally