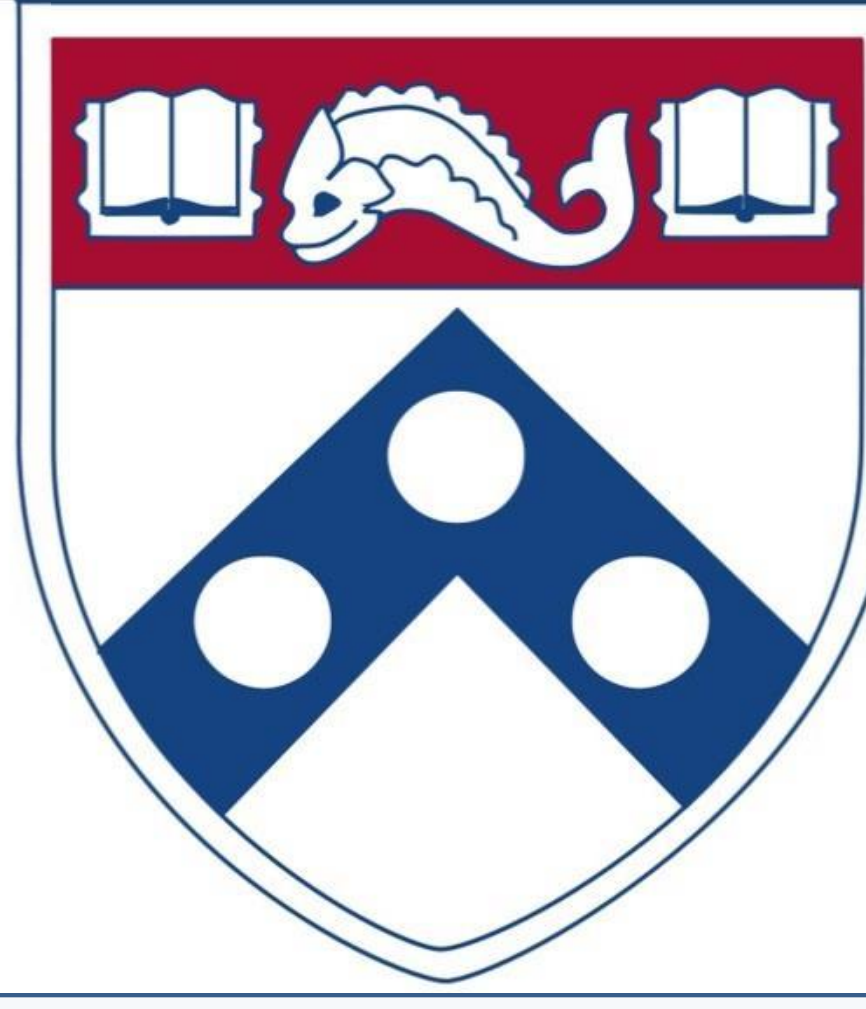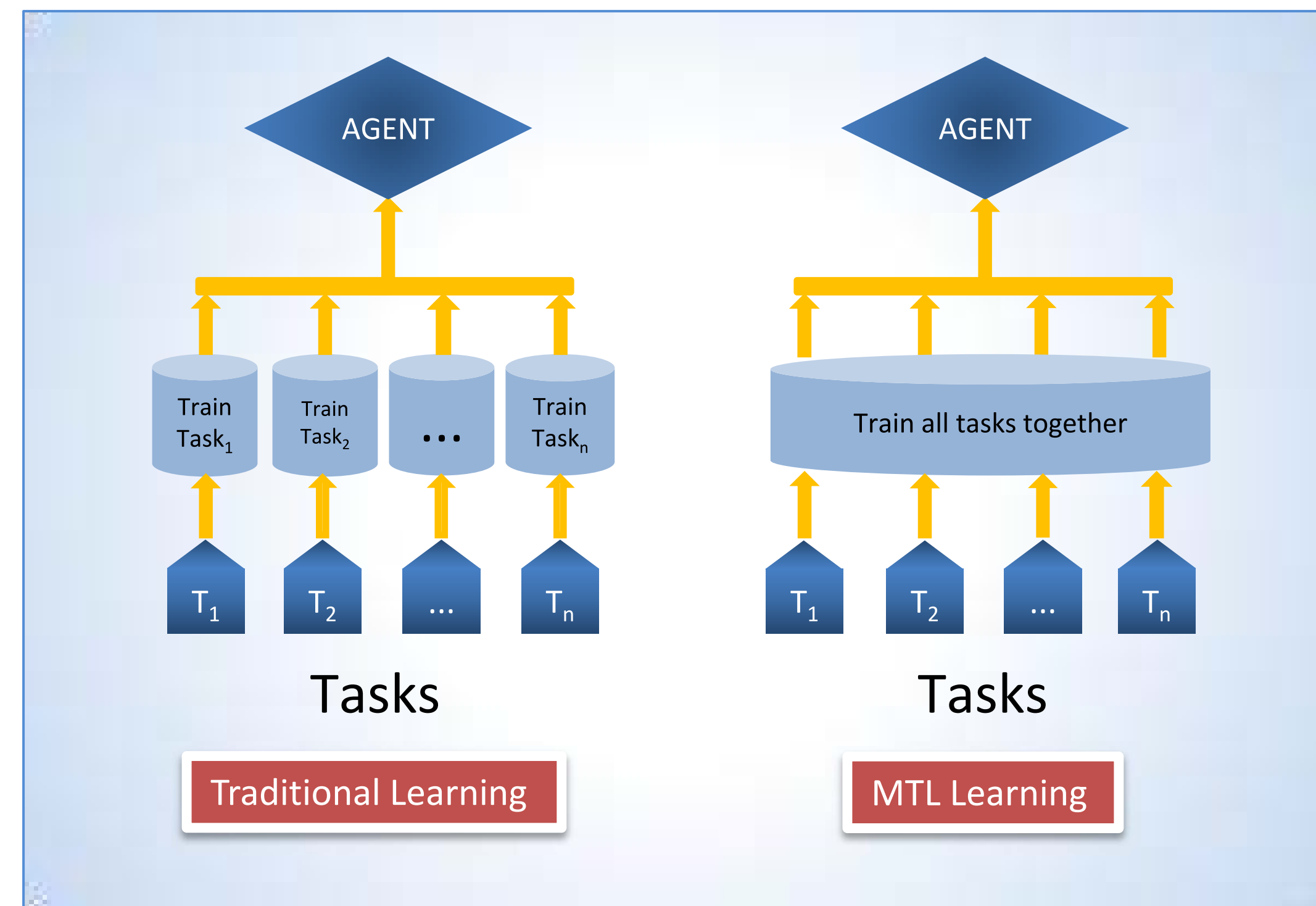# Online Multi-Task Gradient Temporal-Difference Learning

## Vishnu Purushothaman Sreenivasan, Haitham Bou Ammar, and Eric Eaton
## University of Pennsylvania, Computer and Information Science Department
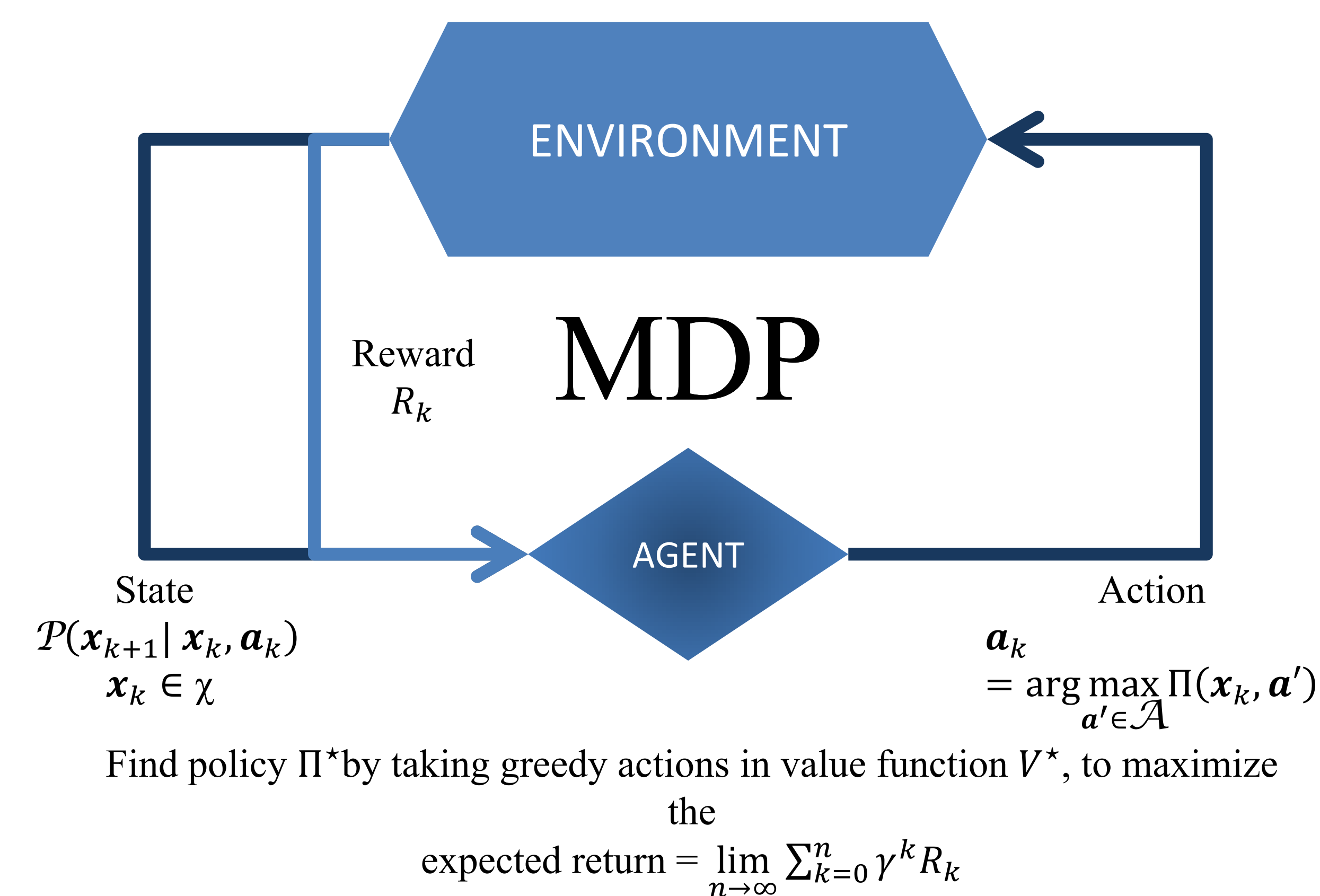
## Motivation

- Reinforcement learning is widely used for design of autonomous systems, but RL agents often require extensive experience to achieve optimal behavior.
- Policies for multiple tasks are often required to be learnt by the agent in order to achieve the overall objective.
- In such scenarios, learning multiple tasks models jointly (called multi-task learning or MTL) produces improved performance but at a large computational cost.



Tasks — Traditional Learning

Tasks — MTL Learning

## Goal

To design a MTL formulation for RL that
1. reduces the required overall interaction time of the agent with the environment,
2. allows the agent to rapidly learn new tasks by building on prior knowledge.

## Background: Reinforcement Learning



ENVIRONMENT

Reward $R_k$

MDP

AGENT

State $\mathcal{P}(x_{k+1} \mid x_k, a_k)$
$x_k \in \chi$

Action $a_k = \arg\max_{a' \in \mathcal{A}} \Pi(x_k, a')$

Find policy $\Pi^*$ by taking greedy actions in value function $V^*$, to maximize the expected return $= \lim_{n \to \infty} \sum_{k=0}^{n} \gamma^k R_k$

## Background: Gradient Temporal Difference Learning

- Value function is approximated by a linear combination of a set of basis functions $\Phi(x)$ representing the state space.

$$V = \theta^T \Phi(x)$$

V – Value function
$\Theta$ – Parameter vector $\Theta \in R^n$
$\Phi$ – State basis function $\Phi : \chi \to R^n$

- Value function estimated from the set $\{(\Phi(x_k), \Phi(x_k'), R_k)\}_{k=1,2,\ldots}$ where,

$$x_k - \text{Current state}, \quad x_k' - \text{Successor state}$$
$$\Phi = \Phi(x_k), \Phi' = \Phi(x_k')$$

- GTD minimizes the L2 norm of the temporal difference error:

$$J(\theta) = E[\delta\phi]^T E[\delta\phi]$$

by following the gradient of the objective function:

$$\nabla_\theta J(\theta) = E[\phi(\phi - \gamma\phi')^T]^T E[\delta\phi]$$
$$\delta = R_k + \gamma\theta^T\phi' - \theta^T\phi$$

## Problem Definition

- Agent learns a series of RL tasks $\mathcal{Z}^{(1)},\ldots,\mathcal{Z}^{(Tmax)}$, each of which is an MDP.

$$\mathcal{Z}^{(t)} = \langle \chi^{(t)}, \mathcal{A}^{(t)}, \mathcal{P}^{(t)}, R^{(t)}, \gamma^{(t)} \rangle$$

  - Tasks may be revisited any number of times and in any order.
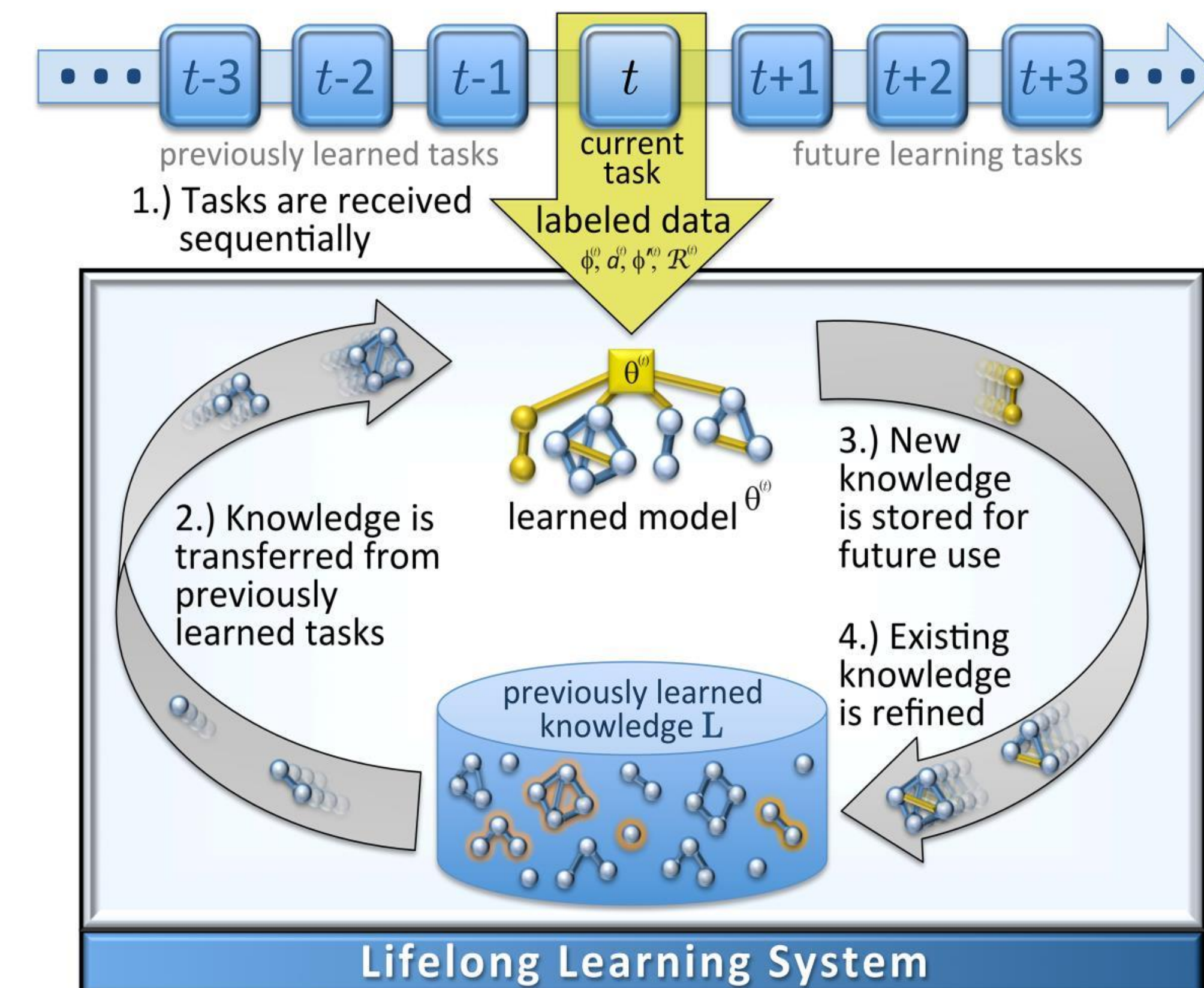  - Agent does not know the total number of tasks a priori.
- The goal is to learn an optimal set of value functions

$$V^* = \left\{ V^*_{\{\theta^{((1))}\}}, \ldots, V^*_{\{\theta^{((Tmax))}\}} \right\}$$

with corresponding parameter vectors $\theta^{(1)}, \ldots, \theta^{(Tmax)}$.

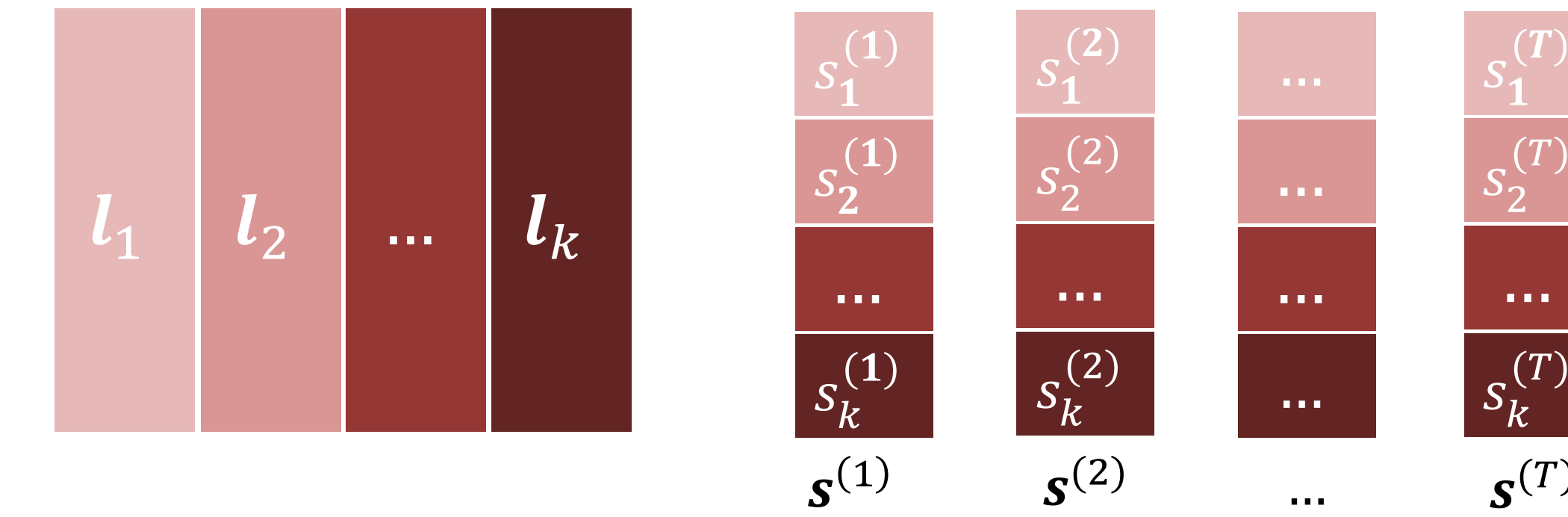- We consider the model-based RL setting, which can be readily extended to a model-free scenario.

## Online Multi-Task Learning Process



$\cdots$ $t-3$ $t-2$ $t-1$ $t$ $t+1$ $t+2$ $t+3$ $\cdots$

previously learned tasks — current task labeled data — future learning tasks

1.) Tasks are received sequentially

learned model $\theta^t$

2.) Knowledge is transferred from previously learned tasks

3.) New knowledge is stored for future use

4.) Existing knowledge is refined

previously learned knowledge L

**Lifelong Learning System**

## Approach

Maintain a library of k latent components $L \in R^{\{d \times k\}}$ that is shared among all the tasks and forms a basis for representing the parameter vector of the task models.

$$\theta^{(t)} = L\, s^{(t)}$$



$l_1$ $l_2$ $\cdots$ $l_k$

$s^{(1)}$ $s^{(2)}$ $\cdots$ $s^{(T)}$

Given T tasks, the MTL objective function is

$$e_T(L) = \frac{1}{T}\sum_{t=1}^{T} \min_{s^{(t)}} \left[ J(\theta^{(t)}) + \mu\left\| s^{(t)} \right\|_1 \right] + \lambda \|L\|_F^2$$

### Eliminating Dependence on All Trajectories

- The above equation is not jointly convex in $L$ and $s^{(t)}$'s.
  - Approximating the loss function $J(\theta^{(t)})$ with the second order Taylor expansion around the optimal single-task solution $\alpha^{(t)}$.
  - Computation of $\alpha^{(t)}$ is performed using GTD.

$$e_T(L) = \frac{1}{T}\sum_{t=1}^{T} \min_{s^{(t)}} \left[ \left\| \alpha^{(t)} - L s^{(t)} \right\|_{\Gamma^{(t)}}^2 + \mu\left\| s^{(t)} \right\|_1 \right] + \lambda \|L\|_F^2$$

$$\alpha^{(t)} = \arg\min_\theta J(\theta^{(t)}) \qquad \Gamma^{(t)} = \nabla_{\theta^{(t)},\theta^{(t)}} J(\theta^{(t)})$$
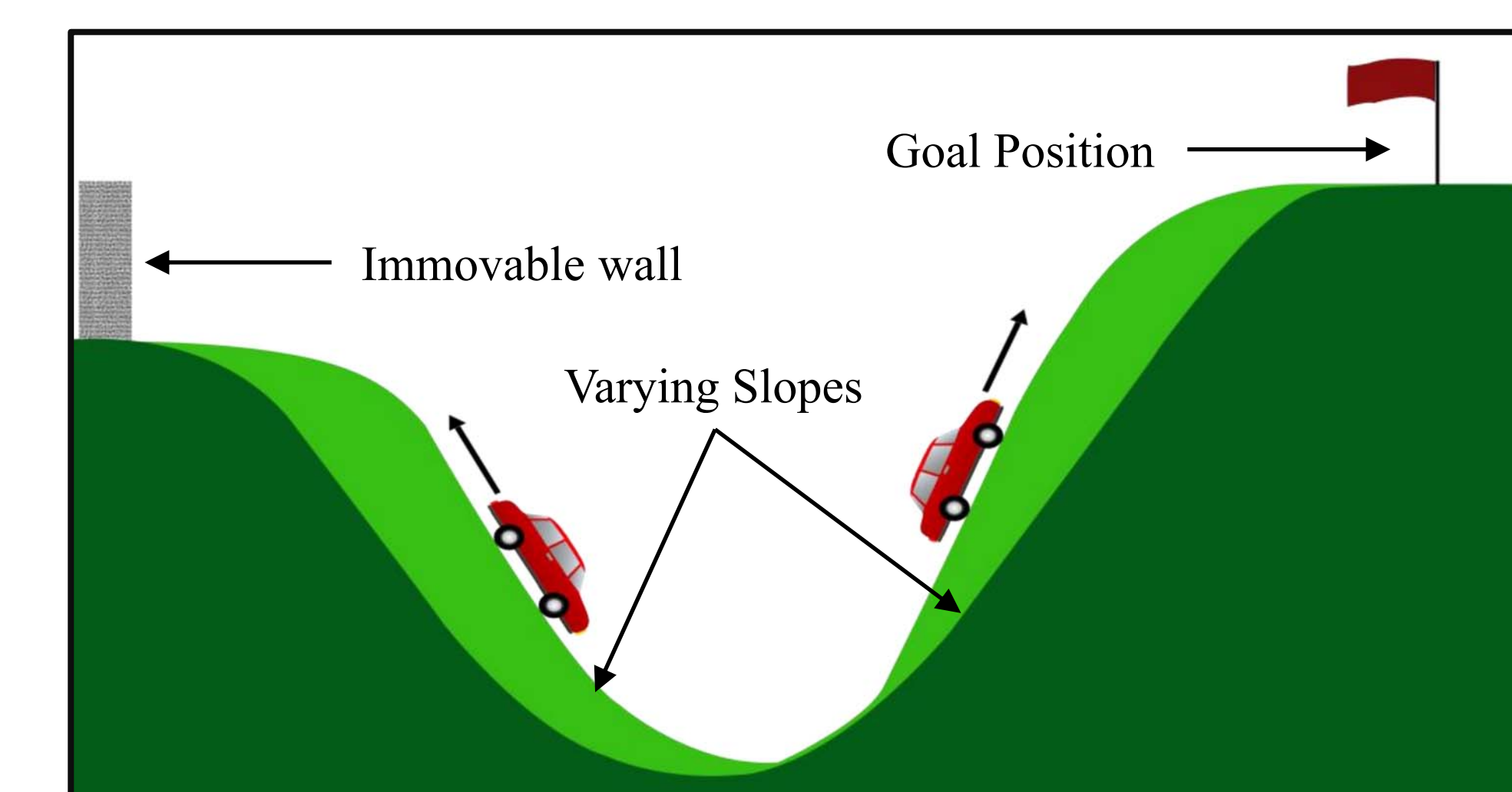
### Eliminating the Reoptimization of Other Tasks

- Modify the MTL objective function by eliminating minimization over all $s^{(t)}$'s.
- Updating $s^{(t)}$'s only when training on task t.

$$s^{(t)} \leftarrow \arg\min_{s^{(t)}} \left( L_m, s^{(t)}, \alpha^{(t)}, \Gamma^{(t)} \right)$$

$$L_{m+1} \leftarrow \arg\min_L \frac{1}{T}\sum_{t=1}^{T} l\left( L, s^{(t)}, \alpha^{(t)}, \Gamma^{(t)} \right) + \lambda \|L\|_F^2$$
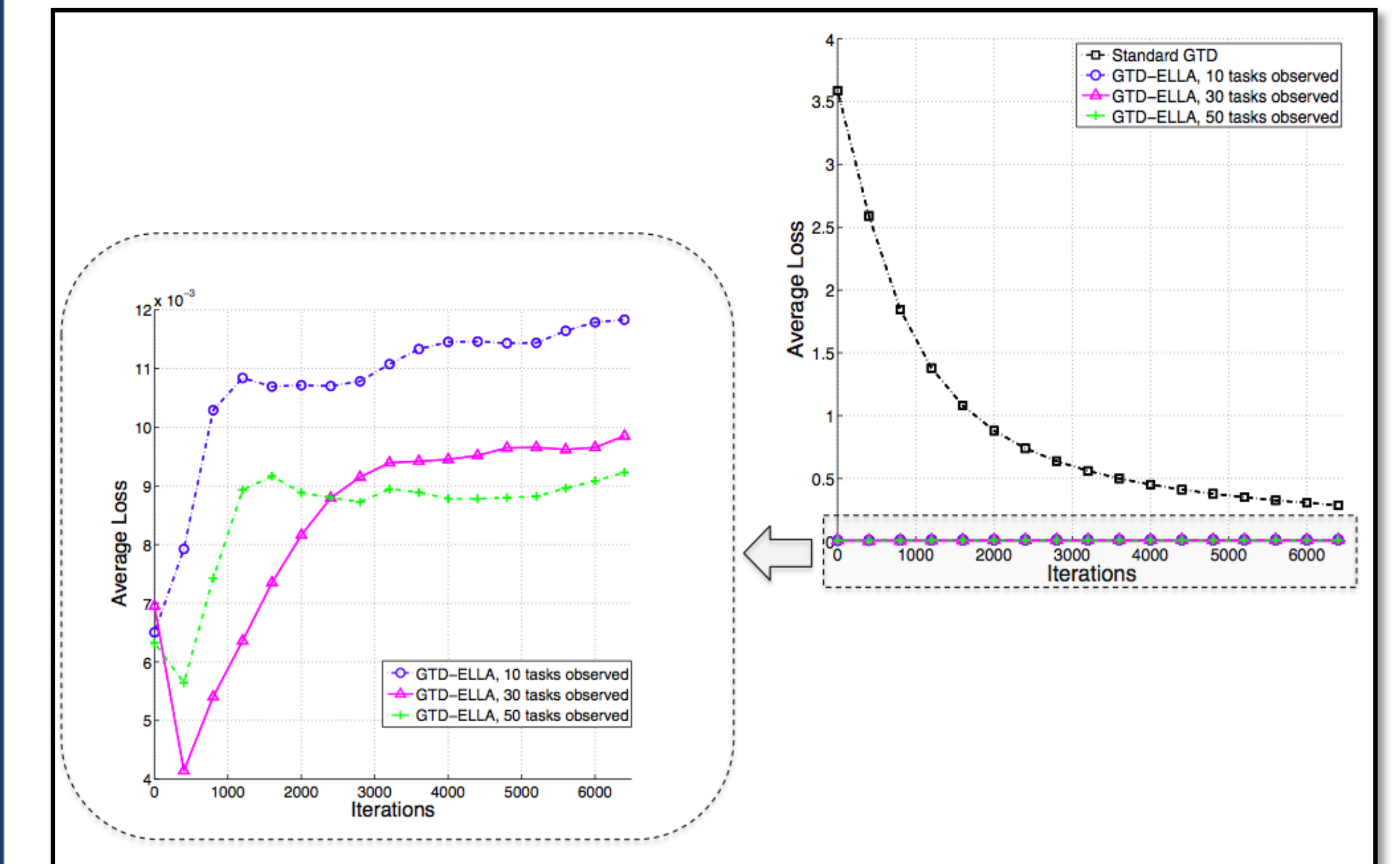
where $l(L, s^{(t)}, \alpha^{(t)}, \Gamma^{(t)}) = \mu\|s\|_1 + \|\alpha - Ls\|^2$ and $L_m$ corresponds to the value of the latent basis at the $m^{th}$ iteration.

## Mountain Car Tasks



Goal Position

Immovable wall

Varying Slopes

## Preliminary Results

- We evaluated GTD-ELLA on multiple tasks in the mountain car (MC) domain.
- State is given by position and velocity, represented by 6 radial basis functions linearly spaced across both the dimensions.
- Parameters:
  - Position is bounded between 1.2 and 0.6.
  - Velocity is bounded between -0.07 to 0.07.
  - Rewards of -1 in all states except goal state at which reward is 0.
- Generated 75 tasks by randomizing the valley slope which also changes the valley position.
- We trained GTD-ELLA on different number of task to learn L and evaluation was conducted on 25 unobserved MC tasks using either GTD-ELLA or standard GTD(0).
- The results indicate GTD-ELLA significantly improves RL performance when training on new tasks. Further, as the agent learns more tasks, its overall performance improves.



## Future Work

- Extend the GTD-ELLA algorithm to a model-free RL setting.
- Support transfer between tasks with different feature spaces

## Acknowledgements

## Contact

Vishnu Purushothaman Sreenivasan
2nd Year, Robotics MSE, University of Pennsylavania
email : visp@seas.upenn.edu