# Lifelong Learning for Disturbance Rejection on Mobile Robots
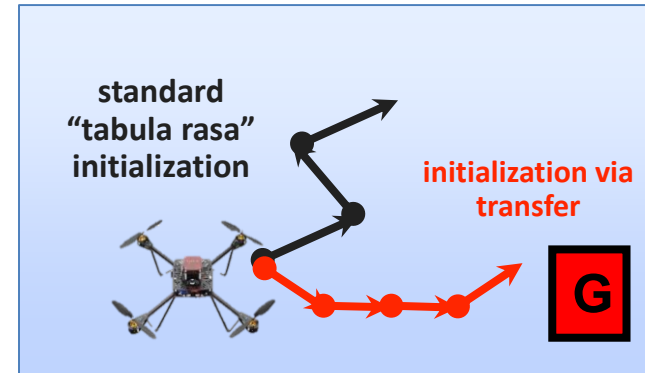
**David Isele, José Marcio Luna, Eric Eaton,**
**Gabriel V. de la Cruz, James Irwin,**
**Brandon Kallaher, Matthew E. Taylor**

# Motivation

**Problem 1:** Without prior knowledge, RL in a new task is slow
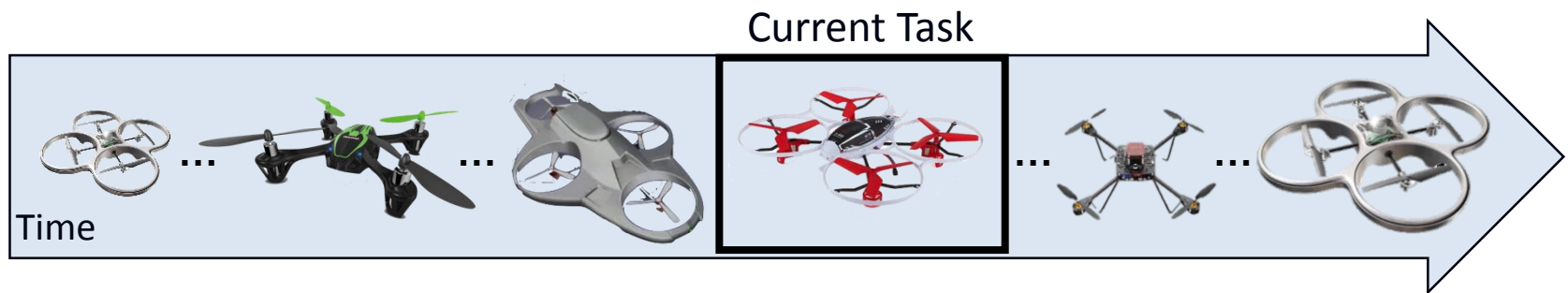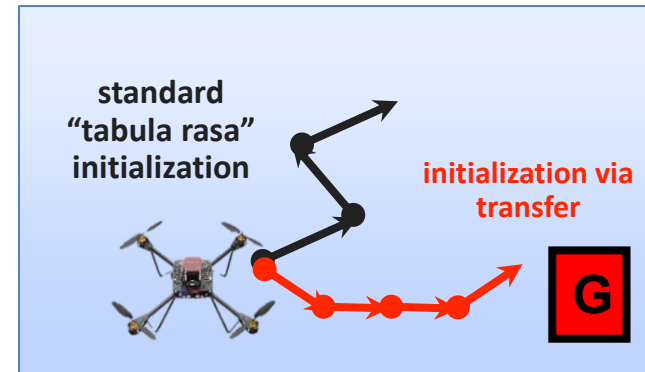
**Idea:** Reuse knowledge from previously learned tasks

# Motivation

**Problem 1:** Without prior knowledge, RL in a new task is slow

**Idea:** Reuse knowledge from previously learned tasks



Current Task



Time

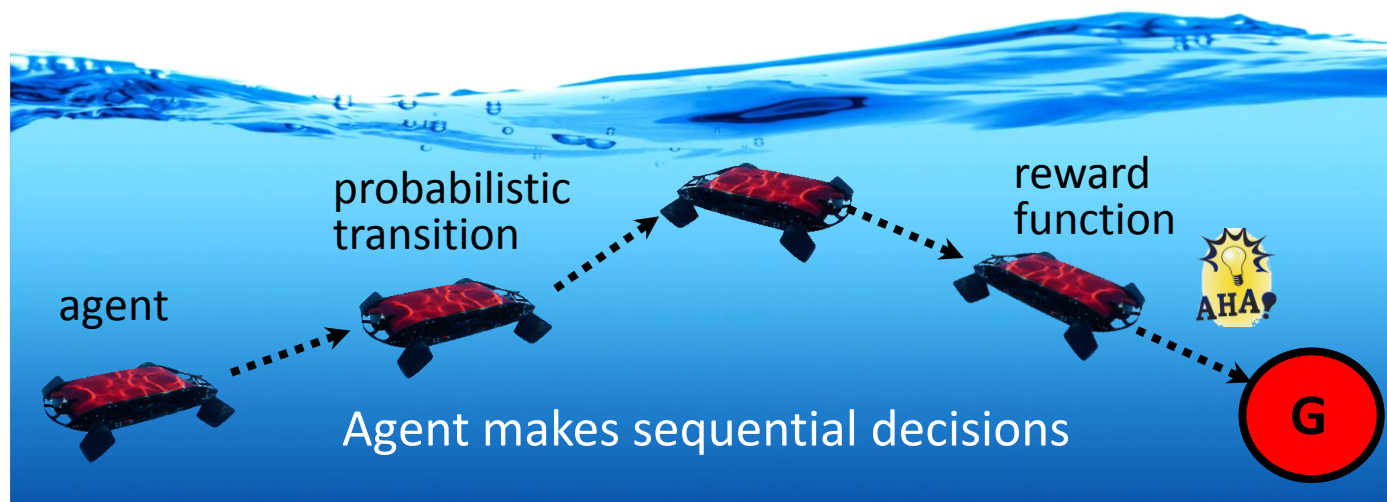We focus on the **lifelong learning** case:
Agent learns multiple tasks consecutively
Want stability guarantees as the number of tasks grows large

# Background

# Background: Policy Gradient Methods for Control

- Agent interacts with environment, taking consecutive actions
- PG methods support continuous state and action spaces
  - Have shown recent success in applications to robotic control [Kober & Peters 2011; Peters & Schaal 2008; Sutton et al. 2000]



Agent makes sequential decisions

- Formalized as a Markov Decision Process (MDP)
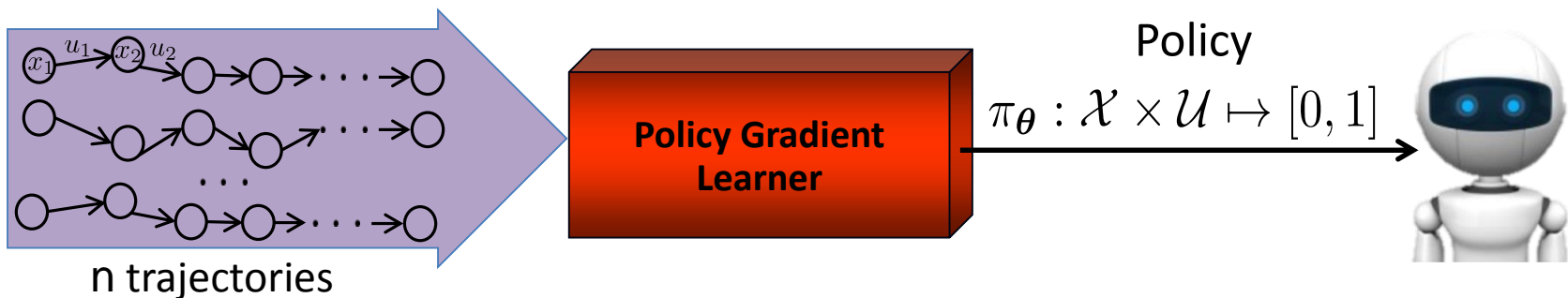
# Background: Policy Gradient Methods for Control

- Agent interacts with environment, taking consecutive actions
- PG methods support continuous state and action spaces
  - Have shown recent success in applications to robotic control
    - [Kober & Peters 2011; Peters & Schaal 2008; Sutton et al. 2000]



$$\pi_{\boldsymbol{\theta}} : \mathcal{X} \times \mathcal{U} \mapsto [0, 1]$$
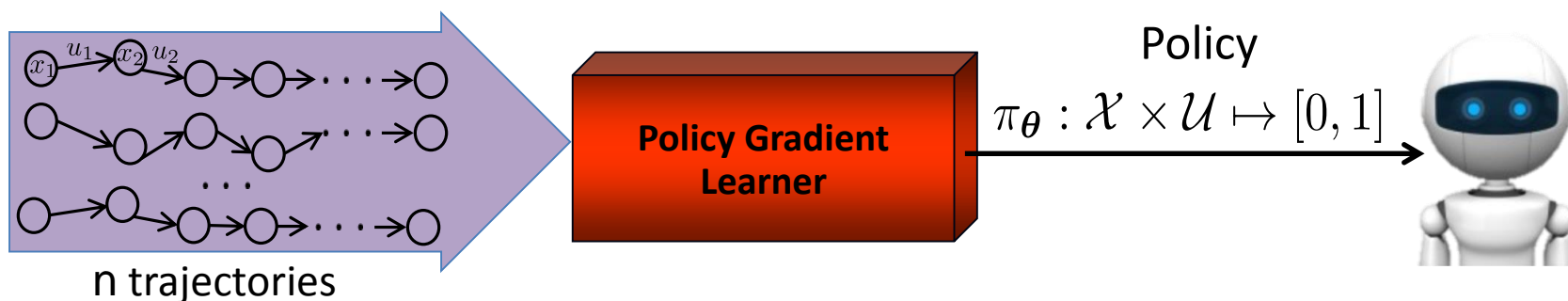
Policy

n trajectories

Policy Gradient Learner

# Background: Policy Gradient Methods for Control

- Agent interacts with environment, taking consecutive actions
- PG methods support continuous state and action spaces
  - Have shown recent success in applications to robotic control
    - [Kober & Peters 2011; Peters & Schaal 2008; Sutton et al. 2000]



n trajectories

Policy Gradient Learner

Policy

$$\pi_{\boldsymbol{\theta}} : \mathcal{X} \times \mathcal{U} \mapsto [0, 1]$$

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

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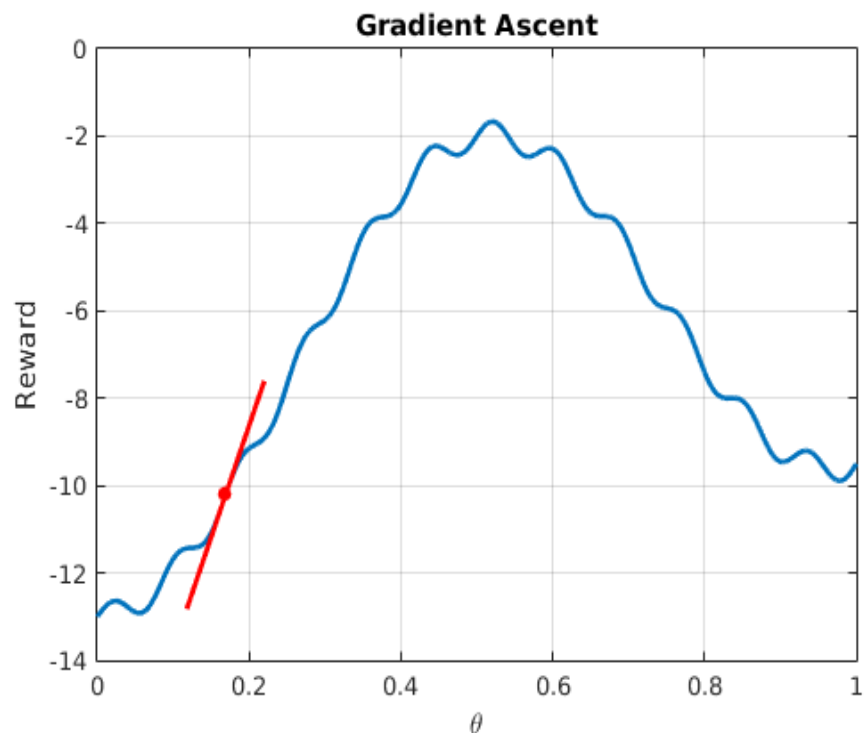
probability of trajectory

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

reward function

# Background: Finite Difference Policy Gradients

Approximate the change in reward with sampled disturbances

$$\Delta \mathcal{J} \approx \mathcal{J}(\theta + \Delta\theta) - \mathcal{J}(\theta)$$
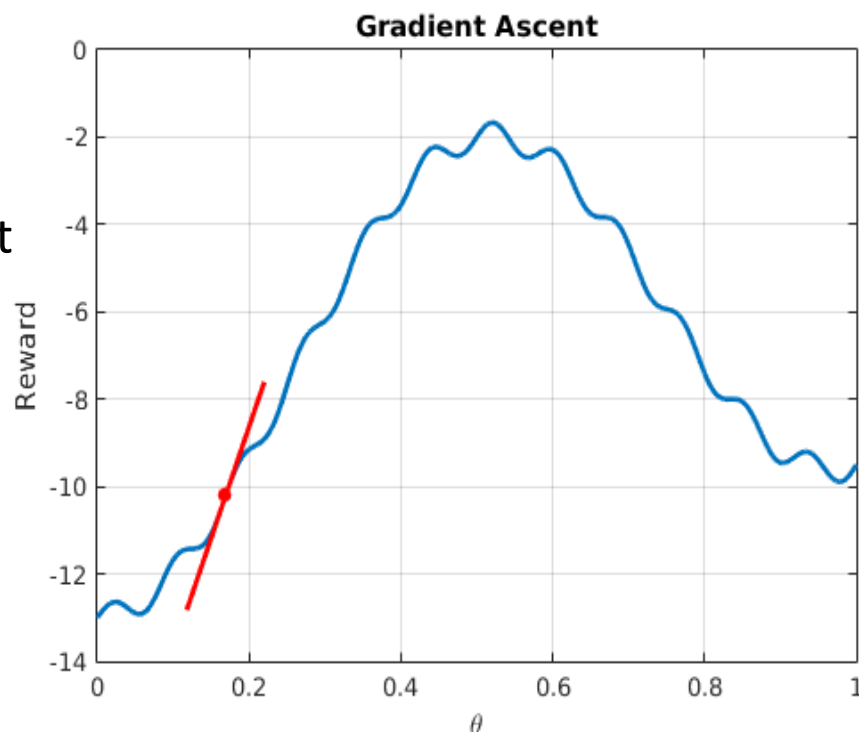


Gradient Ascent

# Background: Finite Difference Policy Gradients

Approximate the change in reward with sampled disturbances

$$\Delta \mathcal{J} \approx \mathcal{J}(\theta + \Delta\theta) - \mathcal{J}(\theta)$$

Use the pseudo-inverse to find the gradient

$$\frac{d\mathcal{J}}{d\theta} = (\Delta\theta^\top \Delta\theta)^{-1} \Delta\theta^\top \Delta\mathcal{J}$$

**Gradient Ascent**

# Background: Finite Difference Policy Gradients

Approximate the change in reward with sampled disturbances
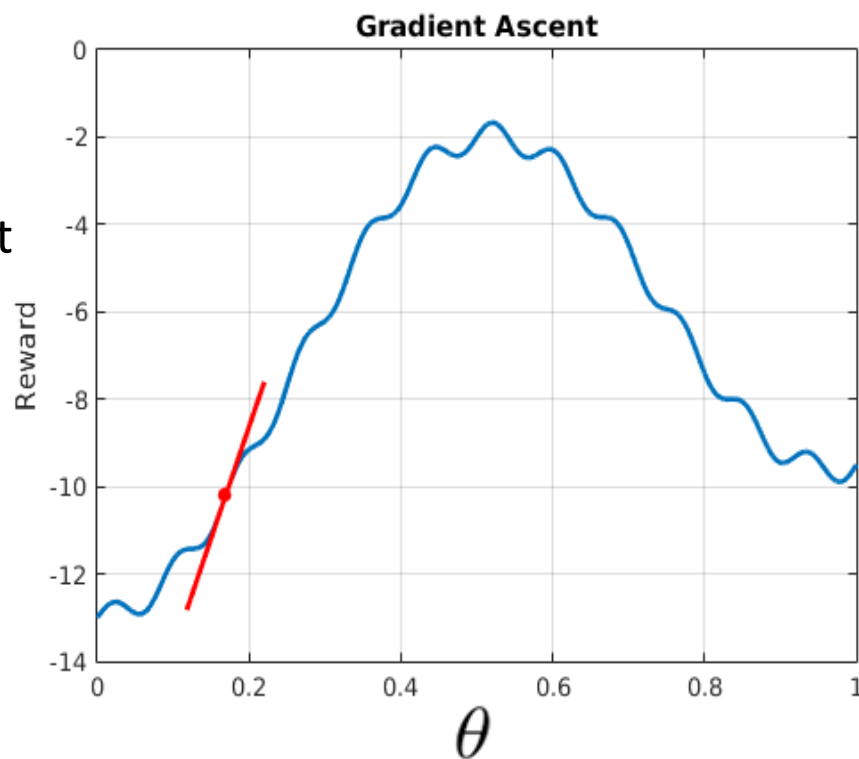
$$\Delta \mathcal{J} \approx \mathcal{J}(\theta + \Delta\theta) - \mathcal{J}(\theta)$$

Use the pseudo-inverse to find the gradient

$$\frac{d\mathcal{J}}{d\theta} = (\Delta\theta^\top \Delta\theta)^{-1} \Delta\theta^\top \Delta\mathcal{J}$$
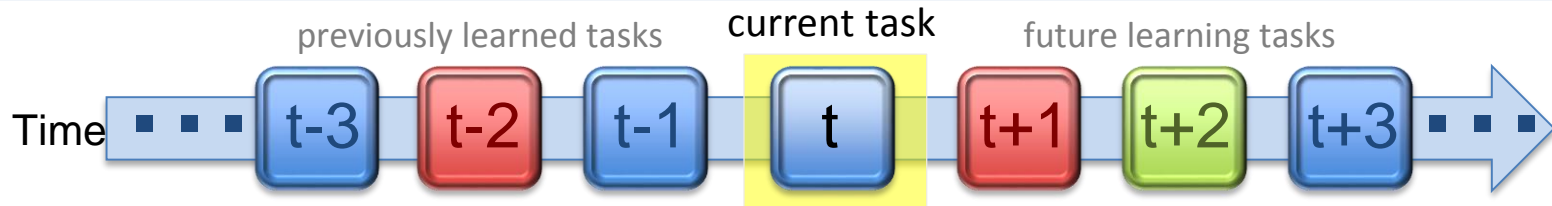
Update the current policy

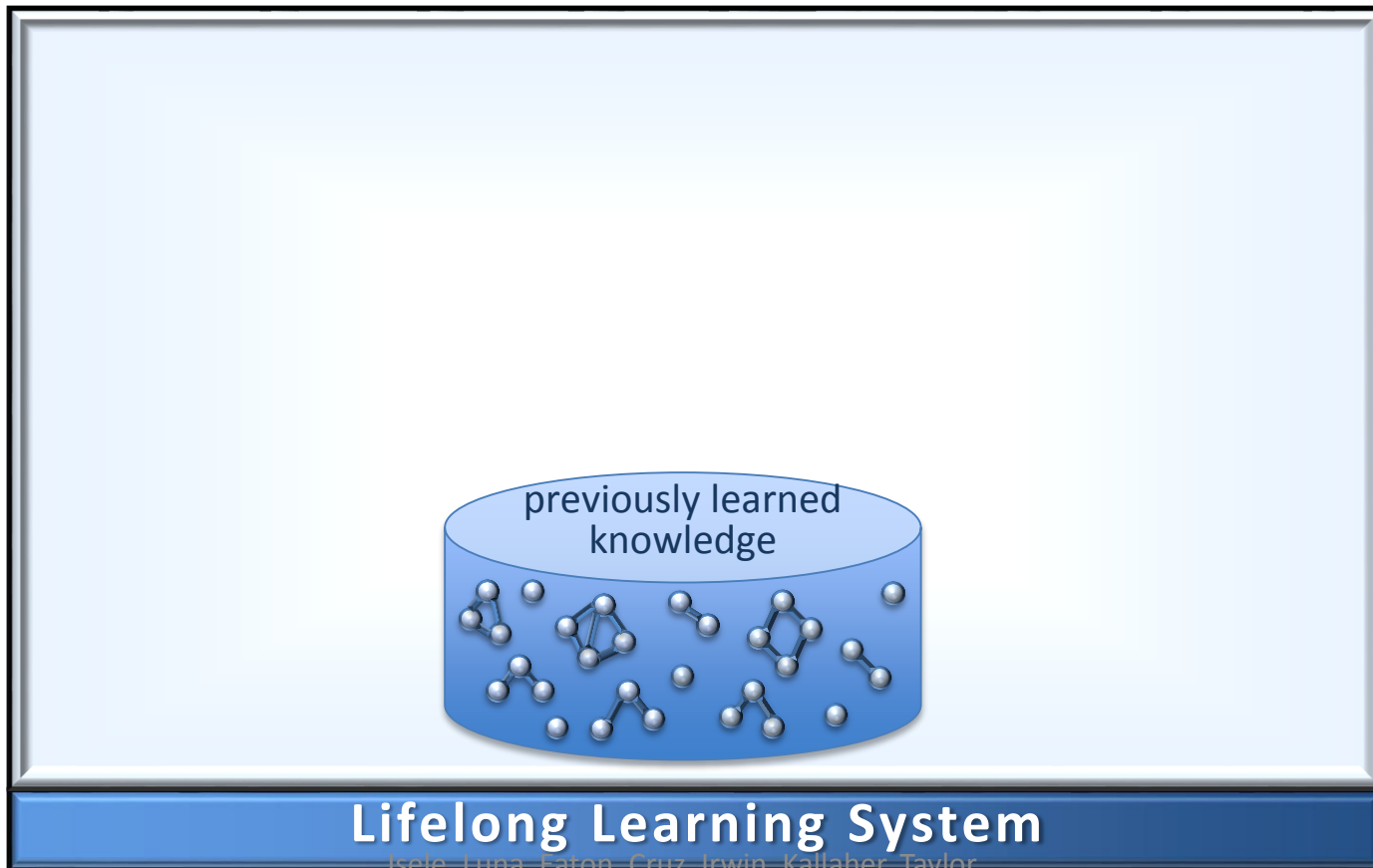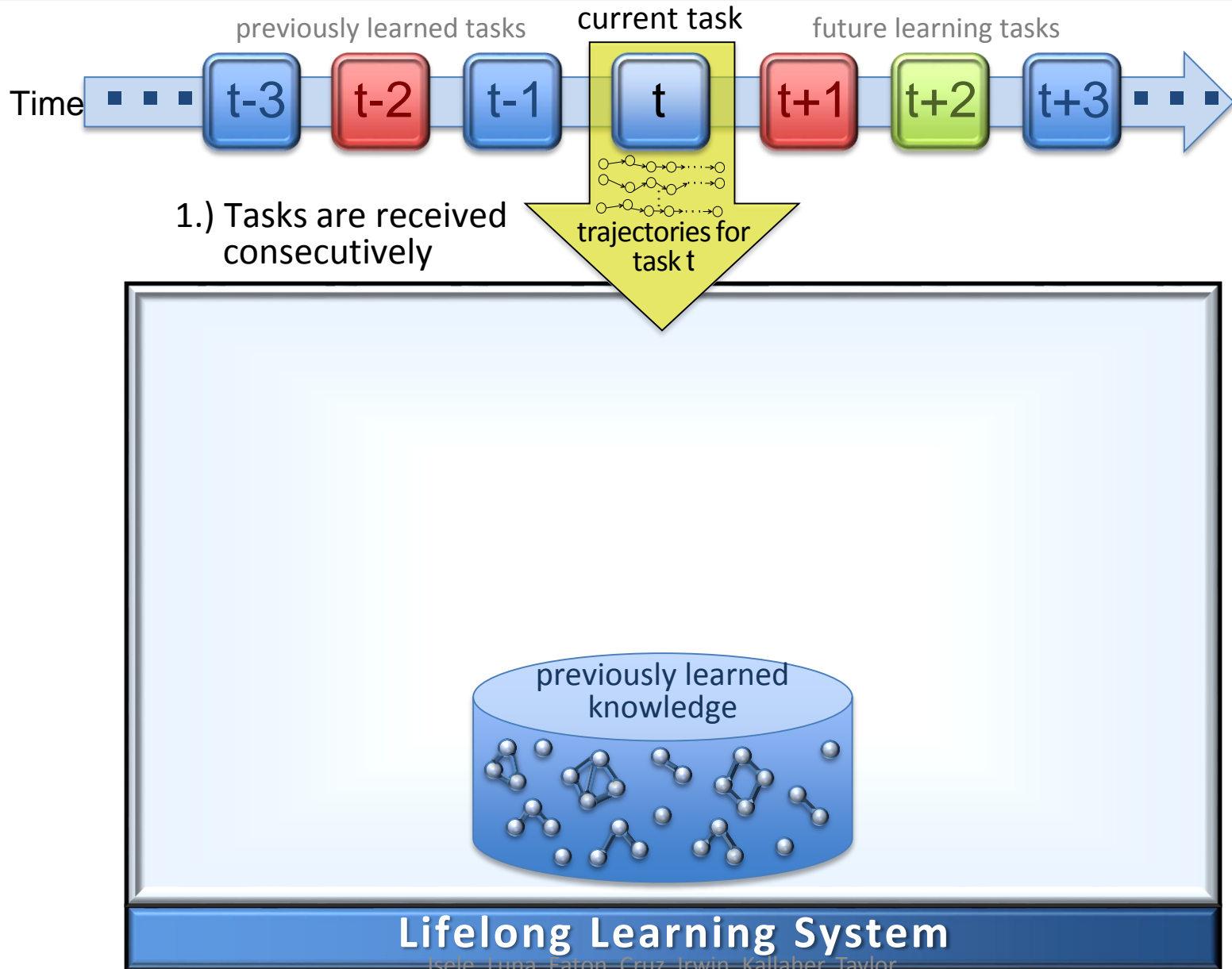$$\theta \leftarrow \theta + \alpha \frac{d\mathcal{J}}{d\theta}$$



**Gradient Ascent**

# Lifelong PG Learning

# Lifelong Machine Learning

previously learned tasks     current task     future learning tasks

Time   ▪ ▪ ▪ ▪  t-3  t-2  t-1  t  t+1  t+2  t+3  ▪ ▪ ▪ ▶

1.) Tasks are received consecutively

previously learned knowledge

**Lifelong Learning System**

Isele, Luna, Eaton, Cruz, Irwin, Kallaher, Taylor

# Lifelong Machine Learning

previously learned tasks | current task | future learning tasks

Time

t-3  t-2  t-1  t  t+1  t+2  t+3

trajectories for task t

1.) Tasks are received consecutively

previously learned knowledge

**Lifelong Learning System**

# Lifelong Machine Learning



previously learned tasks

current task

future learning tasks

Time    t-3   t-2   t-1   t   t+1   t+2   t+3

trajectories for task t

1.) Tasks are received consecutively

previously learned knowledge

**Lifelong Learning System**

# Lifelong Machine Learning

# Lifelong Machine Learning



previously learned tasks  current task  future learning tasks

Time  t-3  t-2  t-1  t  t+1  t+2  t+3

1.) Tasks are received consecutively

trajectories for task t

$\pi_{\boldsymbol{\theta}_t}$

learned policy

2.) Knowledge is transferred from previously learned tasks

3.) New knowledge is stored for future use

previously learned knowledge

**Lifelong Learning System**

# Lifelong Machine Learning

previously learned tasks    current task    future learning tasks

Time    t-3    t-2    t-1    t    t+1    t+2    t+3

trajectories for task t

1.) Tasks are received consecutively

$\pi_{\boldsymbol{\theta}_t}$

learned policy

2.) Knowledge is transferred from previously learned tasks

3.) New knowledge is stored for future use

4.) Existing knowledge is refined

previously learned knowledge

**Lifelong Learning System**

Isele, Luna, Eaton, Cruz, Irwin, Kallaher, Taylor

# PG-ELLA Objective

**Issue:** the objective is dependent on <u>all</u> trajectories

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

# PG-ELLA Objective

**Issue:** the objective is dependent on <u>all</u> trajectories

$$e_T\left(\mathbf{L}\right) = \frac{1}{T}\sum_{t=1}^{T}\min_{\mathbf{s}^{(t)}}\left[-\mathcal{J}\left(\boldsymbol{\theta}^{(t)}\right) + \mu\left\|\mathbf{s}^{(t)}\right\|_1\right] + \lambda\|\mathbf{L}\|_{\mathsf{F}}^2$$

$$\hat{e}_T\left(\mathbf{L}\right) = \frac{1}{T}\sum_{t=1}^{T}\min_{\mathbf{s}^{(t)}}\left[\left\|\boldsymbol{\alpha}^{(t)} - \mathbf{L}\mathbf{s}^{(t)}\right\|_{\boldsymbol{\Gamma}^{(t)}}^2 + \mu\left\|\mathbf{s}^{(t)}\right\|_1\right] + \lambda\|\mathbf{L}\|_{\mathsf{F}}^2$$

Hessian

# **Experiments**

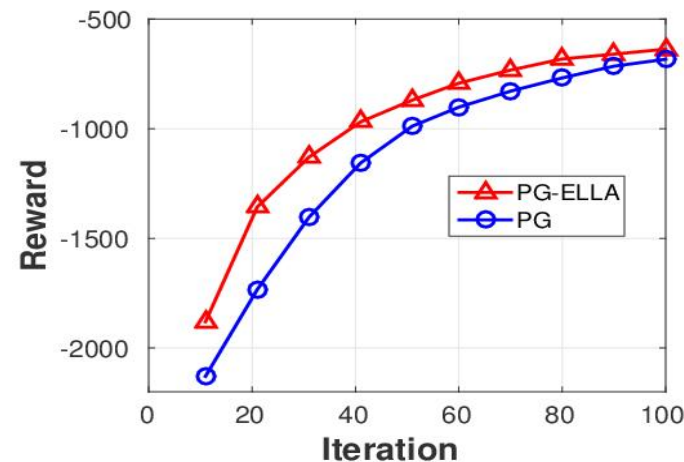## Verification on Robots

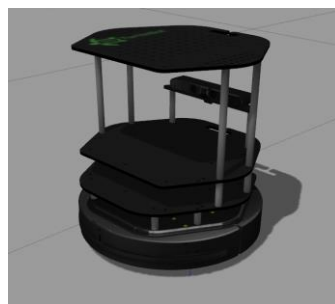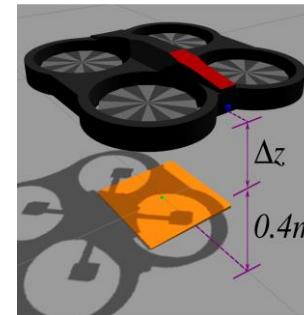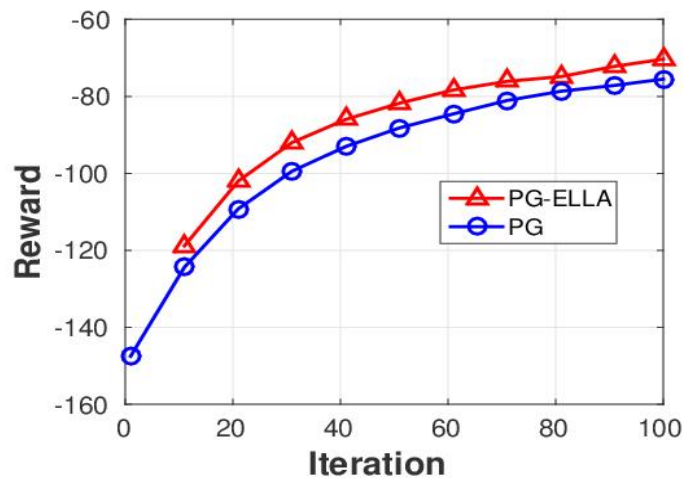# Results for Robot Go-to-Goal Task

- Run RL on a new robot (goal and disturbance) for a small number of iterations
- Use PG-ELLA to adjust policy according to known solutions
- Continue training



## PG-ELLA improves Learning

# Better Results Incorporating Prior

- Initialization with average policy of other robots improves benefit









## PG-ELLA improves Learning

# Lifelong Learning for Disturbance Rejection on Mobile Robots

**David Isele, José Marcio Luna, Eric Eaton,**
**Gabriel V. de la Cruz, James Irwin, Brandon Kallaher, Matthew E. Taylor**

# Thank you!

# Questions?