

Deep and Cooperative Learning for Approximate Dynamic Programming

MAX GREENE AND WARREN DIXON

Dept. of Mechanical and Aerospace Engineering, Univ. of Florida





- Deep Neural Network Identifier for Approximate Optimal Control
 - Online System ID using DNN
 - Update the DNN in Real-Time
 - Approximate System Dynamics
 - Use System Model to Approximate Optimal Controller
- Cooperative Approximate Optimal Control
 - Leverage Multiple Agents Learning the Same Function

Deep Neural Network-based Approximate Optimal Tracking for Unknown Nonlinear Systems

MAX GREENE¹, ZACHARY BELL², SCOTT NIVISON², WARREN DIXON¹

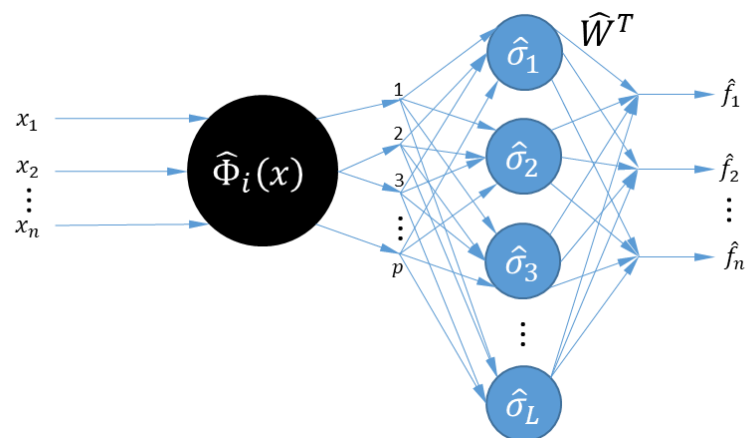
¹Dept. of Mechanical and Aerospace Engineering, Univ. of Florida

²Munitions Directorate, Eglin AFB, Air Force Research Laboratory



Real-Time DNN Approximation

- Multi-Timescale DNN
 - Online system identifier
 - Better than single hidden-layer
 - Inner features updates iteratively (switched)
 - Output-layer updates in real-time
 - Can we use a DNN to approximate an optimal control policy online?



Inner Features (Trained Iteratively) Fully Connected Layer (Updated in real-time)



Real-Time DNN Approximation

Dynamical System

Given a control affine nonlinear dynamical system:

$$\dot{x} = f(x) + g(x)u$$

DNN Approximation

Approximate the drift dynamics with a DNN:

$$\text{Dynamics:} \quad f(x) = \theta^T \phi^*(\Phi^*(x)) + \epsilon(x)$$

$$i^{\text{th}} \text{ DNN Approximation:} \quad \hat{f}_i(x, \hat{\theta}) = \hat{\theta}^T \hat{\phi}_i(\hat{\Phi}_i(x))$$

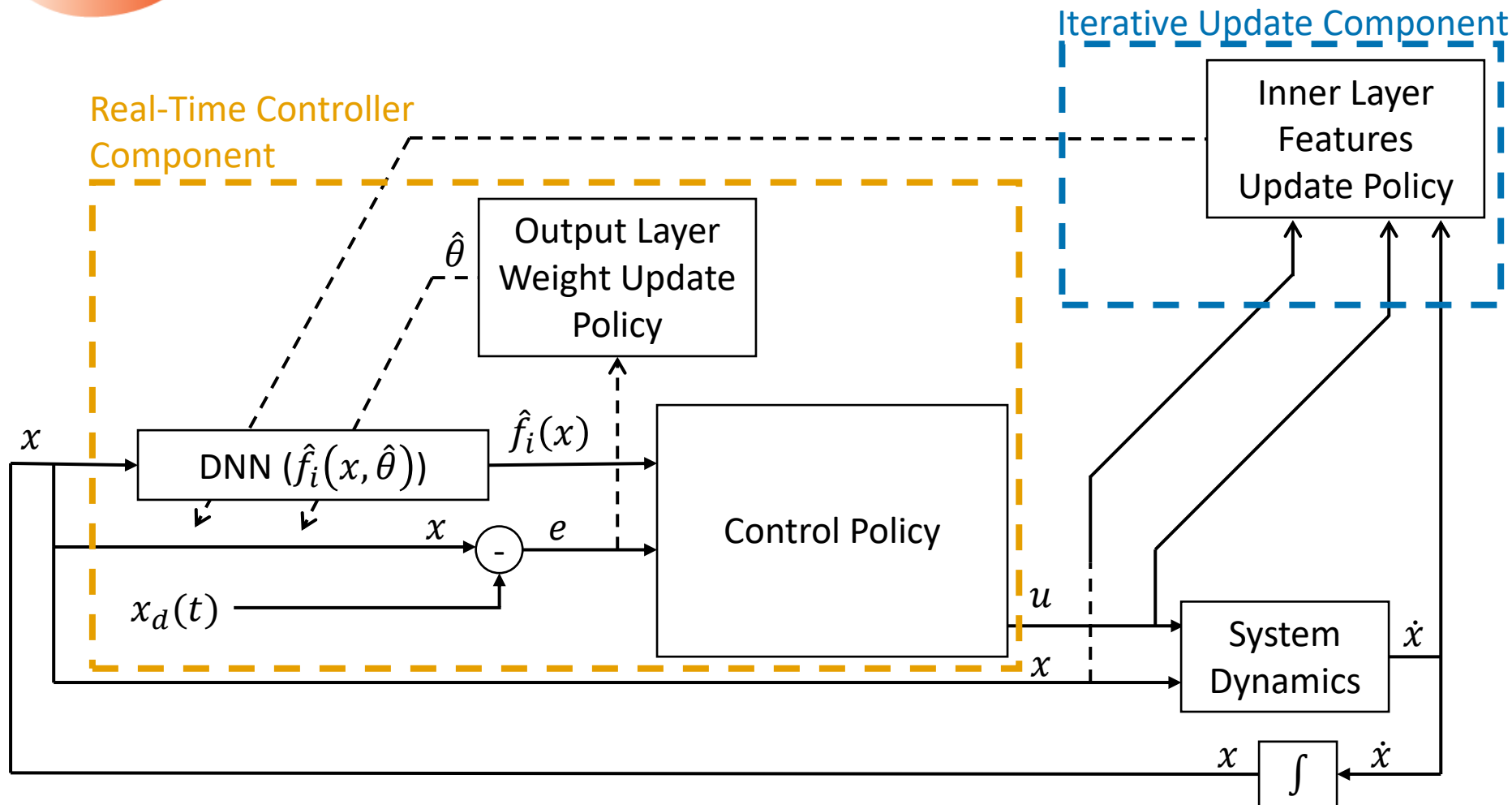
Looks Like:

Reminds us of linearly parameterizable dynamics:

$$f(x) = \theta^T \phi^*(\Phi^*(x)) + \epsilon(x) \approx \theta^T Y(x) + \epsilon(x)$$



Real-Time DNN Approximation





Dynamical System

Given a control affine nonlinear dynamical system:

$$\dot{\zeta} = F(\zeta) + G(\zeta)\mu$$

Control Objective

Design a controller, u , which minimizes a cost function:

$$J(\zeta, \mu) = \min_{u(\tau) \in U} \int_0^{\infty} Q(\zeta(\tau)) + \mu(\tau)^T R \mu(\tau) d\tau$$

Cost-to-Go

Optimal value function:

$$V^*(x) = \min_{u(\tau) \in U} \int_0^{\infty} Q(\zeta(\tau)) + \mu(\tau)^T R \mu(\tau) d\tau$$



Hamilton Jacobi Bellman Equation

Hamilton Jacobi Bellman (HJB) equation:

$$0 = Q(\zeta) + \mu^*(\zeta)^T R \mu^*(\zeta) + \nabla_{\zeta} V^*(\zeta) (F(\zeta) + G(\zeta) \mu^*(\zeta))$$

Optimal Controller

From the HJB equation:

$$\mu^*(\zeta) = -\frac{1}{2} R^{-1} G(\zeta)^T \left(\nabla_{\zeta} V^*(\zeta) \right)^T$$

- Cannot solve HJB analytically.
- Approximate the Value Function (V^*)
 - Universal Function Approximation Property
 - Neural Networks



Optimal Value Function and Optimal Control Policy:

$$V^*(\zeta) = W^T \sigma(\zeta) + \varepsilon(\zeta) \quad \mu^*(\zeta) = -\frac{1}{2} R^{-1} G(\zeta)^T (\nabla_{\zeta} \sigma(\zeta)^T W + \nabla_{\zeta} \varepsilon(\zeta)^T)$$

Unknown: Neural weights

$$\widehat{W}_c, \widehat{W}_a \rightarrow W$$

\widehat{W}_c : Critic weight
 \widehat{W}_a : Actor weight

Value Function and Optimal Control Policy Approximation

$$\widehat{V}(\zeta, \widehat{W}_c) = \widehat{W}_c^T \sigma(\zeta) \quad \widehat{\mu}(\zeta, \widehat{W}_a) = -\frac{1}{2} R^{-1} G(\zeta)^T (\nabla_{\zeta} \sigma(\zeta)^T \widehat{W}_a)$$

Bellman Error (BE): Residual from HJB

$$\widehat{\delta}(\zeta, \widehat{W}_c, \widehat{W}_a) \triangleq Q(\zeta) + \widehat{\mu}(\zeta, \widehat{W}_a)^T R \widehat{\mu}(\zeta, \widehat{W}_a) + \nabla_{\zeta} \widehat{V}(\zeta, \widehat{W}_c) \left(\widehat{F}_i(\zeta) + G(\zeta) \widehat{\mu}(\zeta, \widehat{W}_a) \right)$$



Instantaneous Bellman Error: Residual from HJB

$$\hat{\delta}(t) \triangleq \hat{\delta}(\zeta(t), \hat{W}_c(t), \hat{W}_a(t))$$

Weight Update Laws using R-MBRL

$$\dot{\hat{W}}_c(t) = -\eta_{c1} \Gamma \frac{\omega(t)}{\rho(t)} \hat{\delta}(t) - \eta_{c2} \Gamma \frac{1}{N} \sum_{i=1}^N \frac{\omega_i(t)}{\rho_i(t)} \hat{\delta}_i(t)$$

Minimize the Bellman Error

$$\dot{\Gamma}(t) = \left(\lambda \Gamma(t) - \eta_{c1} \frac{\Gamma(t) \omega(t) \omega(t)^T \Gamma(t)}{\rho(t)} - \eta_{c2} \Gamma(t) \left(\frac{1}{N} \sum_{i=1}^N \frac{\omega_i(t) \omega_i^T(t)}{\rho_i(t)} \right) \Gamma(t) \right) \mathbf{1}$$

Adjust the “learning rate”

$$\begin{aligned} \dot{\hat{W}}_a(t) = & -\eta_{c1} (\hat{W}_a(t) - \hat{W}_c(t)) - \eta_{a2} \hat{W}_a(t) + \eta_{c1} \frac{G_\sigma(t)^T \hat{W}_a(t) \omega(t)^T}{4\rho(t)} \hat{W}_c(t) \\ & + \eta_{c2} \left(\frac{1}{N} \sum_{i=1}^N \frac{G_{i\sigma}^T \hat{W}_a(t) \omega_i(t)^T}{4\rho_i(t)} \right) \hat{W}_c(t) \end{aligned}$$

Follow the critic weights



Instantaneous Bellman Error: Residual from HJB

$$\hat{\delta}(t) \triangleq \hat{\delta}(\zeta(t), \widehat{W}_c(t), \widehat{W}_a(t))$$

Weight Update Laws using R-MBRL

$$\dot{\widehat{W}}_c(t) = -\eta_{c1} \Gamma \frac{\omega(t)}{\rho(t)} \hat{\delta}(t) - \eta_{c2} \Gamma \frac{1}{N} \sum_{i=1}^N \frac{\omega_i(t)}{\rho_i(t)} \hat{\delta}_i(t)$$

On-Trajectory
Point

$$\dot{\Gamma}(t) = \left(\lambda \Gamma(t) - \eta_{c1} \frac{\Gamma(t) \omega(t) \omega(t)^T \Gamma(t)}{\rho(t)} - \eta_{c2} \Gamma(t) \left(\frac{1}{N} \sum_{i=1}^N \frac{\omega_i(t) \omega_i^T(t)}{\rho_i(t)} \right) \Gamma(t) \right) \mathbf{1}_{\{\underline{\Gamma} \leq \|\Gamma\| \leq \bar{\Gamma}\}}$$

$$\begin{aligned} \dot{\widehat{W}}_a(t) = & -\eta_{c1} (\widehat{W}_a(t) - \widehat{W}_c(t)) - \eta_{a2} \widehat{W}_a(t) + \eta_{c1} \frac{G_\sigma(t)^T \widehat{W}_a(t) \omega(t)^T}{4\rho(t)} \widehat{W}_c(t) \\ & + \eta_{c2} \left(\frac{1}{N} \sum_{i=1}^N \frac{G_{i\sigma}^T \widehat{W}_a(t) \omega_i(t)^T}{4\rho_i(t)} \right) \widehat{W}_c(t) \end{aligned}$$

Off-Trajectory
Points



Theorem 1:

Provided the necessary assumptions, gain conditions, and initial conditions hold, then the tracking error $e(t)$, weight estimation errors $\tilde{W}_c(t)$ and $\tilde{W}_a(t)$, system ID DNN output weight estimation errors $\tilde{\theta}(t)$, state estimator error $\tilde{x}(t)$, and transient control policy $\mu(t)$ are uniformly ultimately bounded (UUB).

- While iteratively retraining (switching) DNN-based system ID terms:
 - $\|e\|$ converges to a neighborhood around zero
 - $\|\tilde{\theta}\|$, $\|\tilde{x}\|$ converge to a neighborhood around zero
 - $\|\tilde{W}_c\|$, $\|\tilde{W}_a\|$ converge to a neighborhood around zero
 - $\mu(t)$ is an approximation of the optimal control policy



- Contribution of Theorem 1

- Trajectory Tracking
- Approximate Optimal Control Policy
- Allows iteratively updating the DNN inner-layer features.

- Obstacles in Theorem 1

- Updating the DNN introduces discontinuities

- $\hat{\delta}(t) \triangleq \nabla_{\zeta} \hat{V}(\zeta, \hat{W}_c) \left(\hat{F}_i(\zeta) + G(\zeta) \hat{\mu}(\zeta, \hat{W}_a) \right) + \dots$

- $\hat{W}_c(t) = -\eta_{c1} \Gamma \frac{\omega(t)}{\rho(t)} \hat{\delta}(t) - \eta_{c2} \frac{1}{N} \sum_{i=1}^{N_e} \frac{\omega_i(t)}{\rho_i(t)} \hat{\delta}_i(t)$

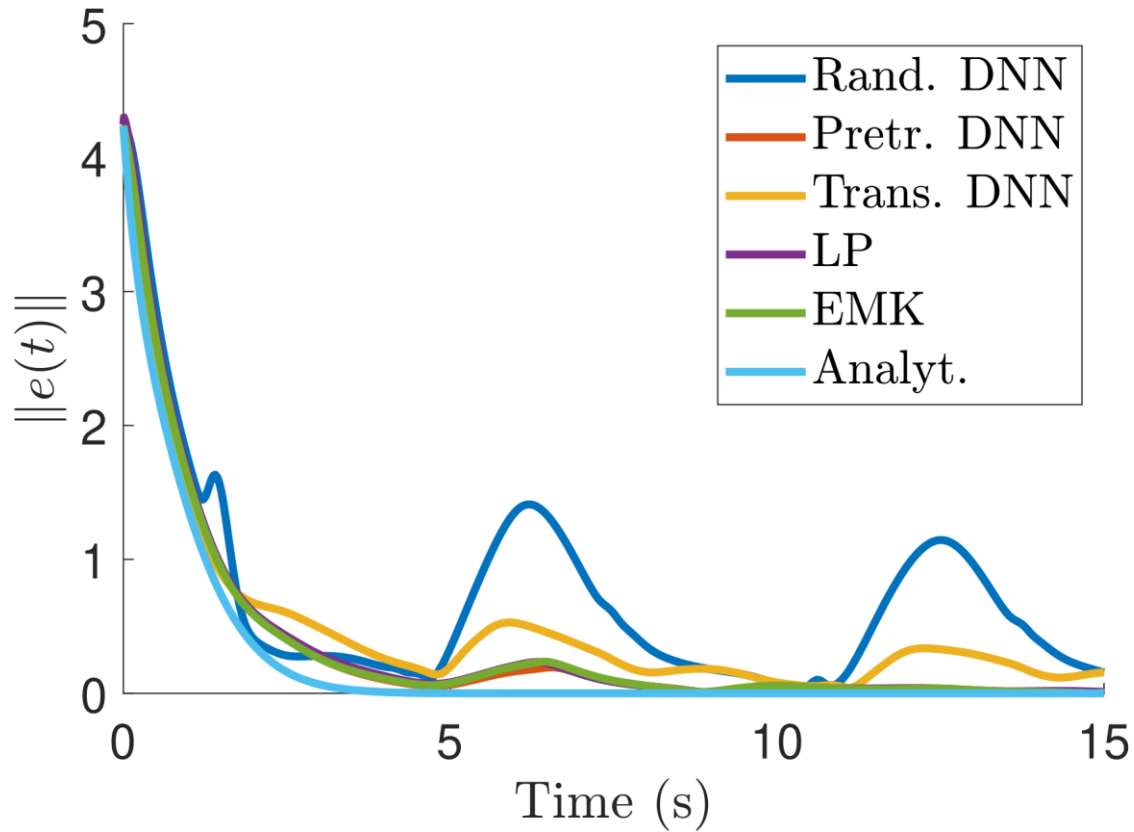
- $\hat{W}_c(t)$ is discontinuous

- $\hat{W}_c(t)$ is nonsmooth

- Requires use of generalized solutions to show stability



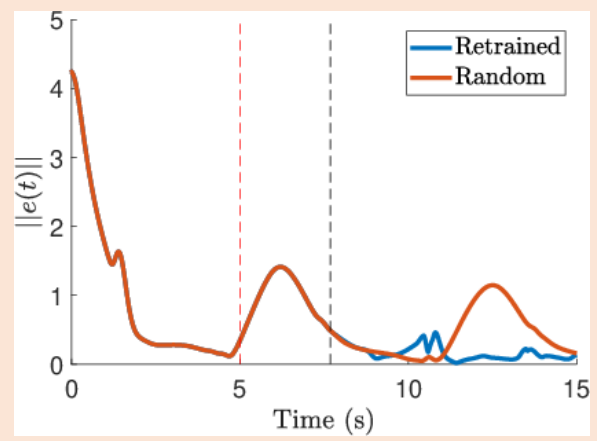
DNN ADP Simulation Results



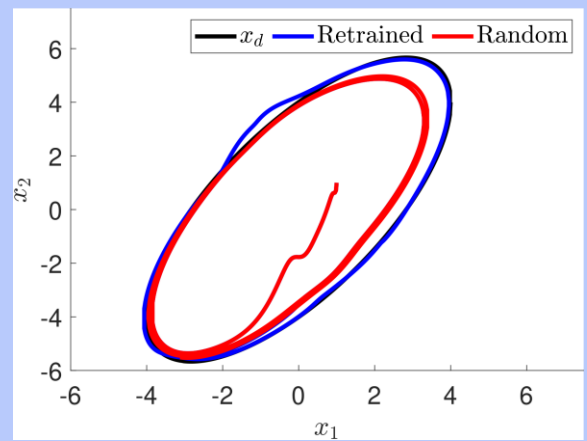


DNN ADP Simulation Results

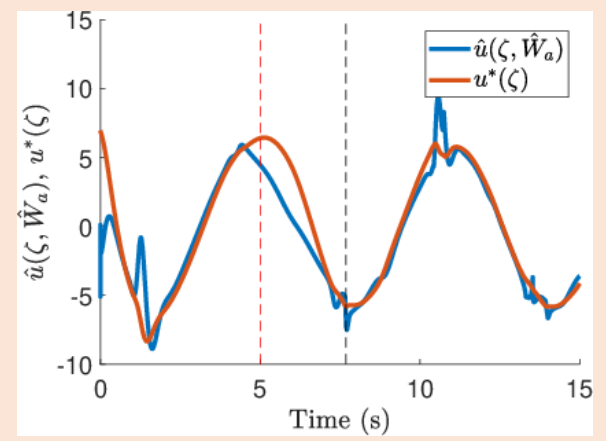
Norm of Error



Phase Plot

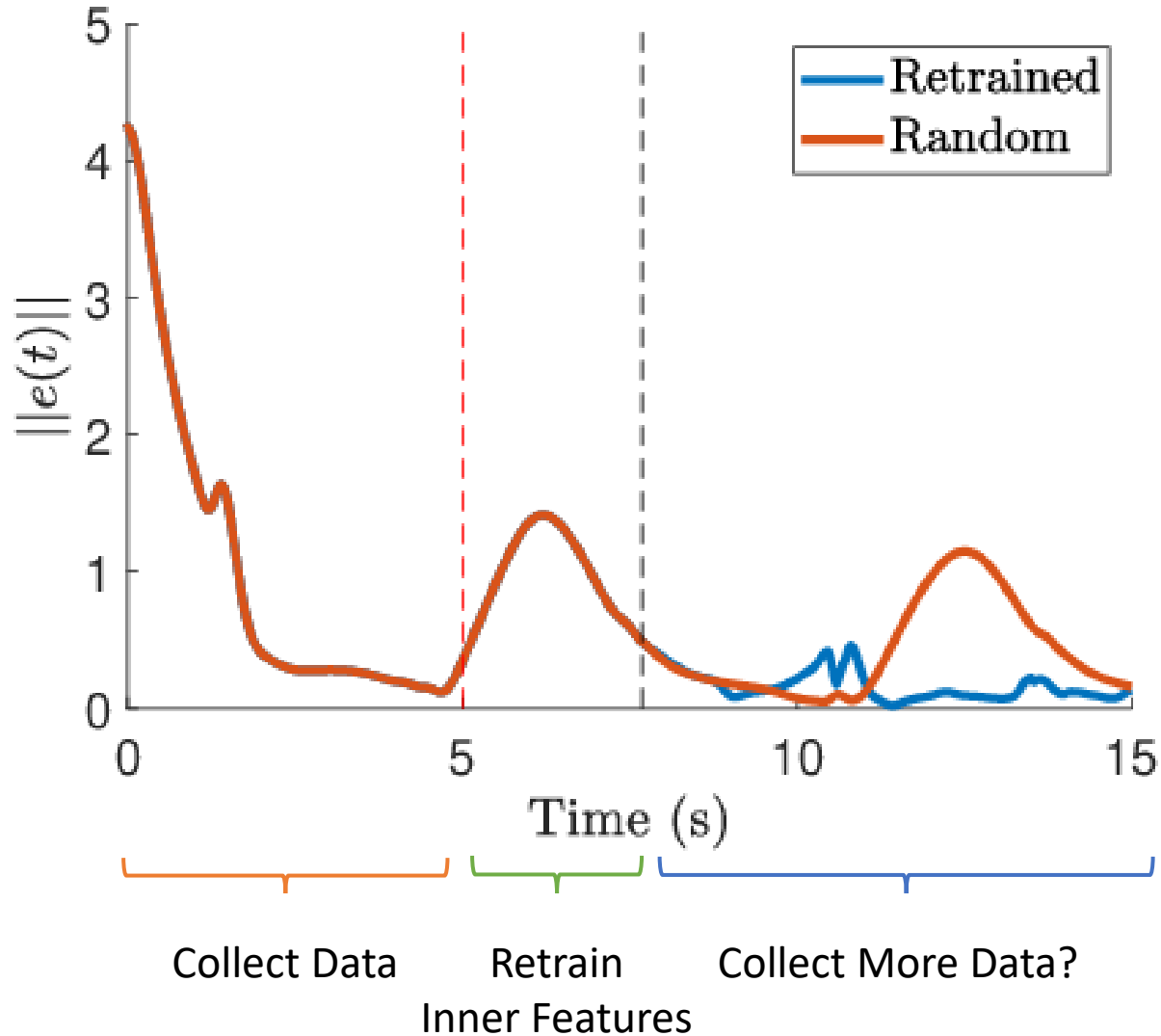


Optimal Control Comparison





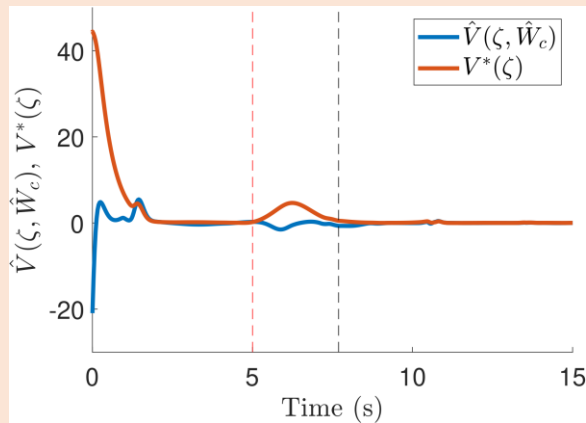
DNN ADP Simulation Results



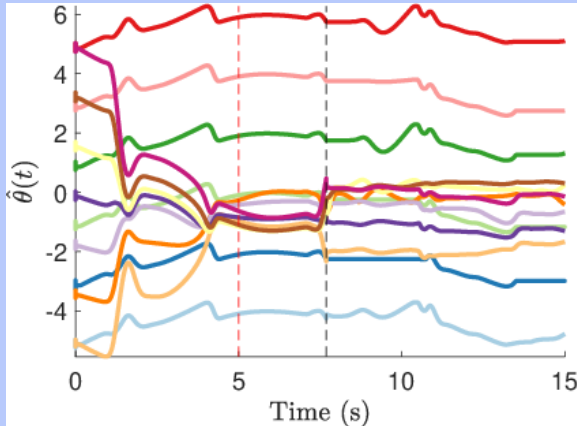


DNN ADP Simulation Results

Value Function Comparison



DNN Output Weights



Cooperative Model-Based Reinforcement Learning for Approximate Optimal Tracking

MAX GREENE¹, ZACHARY BELL², SCOTT NIVISON², JONATHAN HOW³, WARREN DIXON¹

¹Dept. of Mechanical and Aerospace Engineering, Univ. of Florida

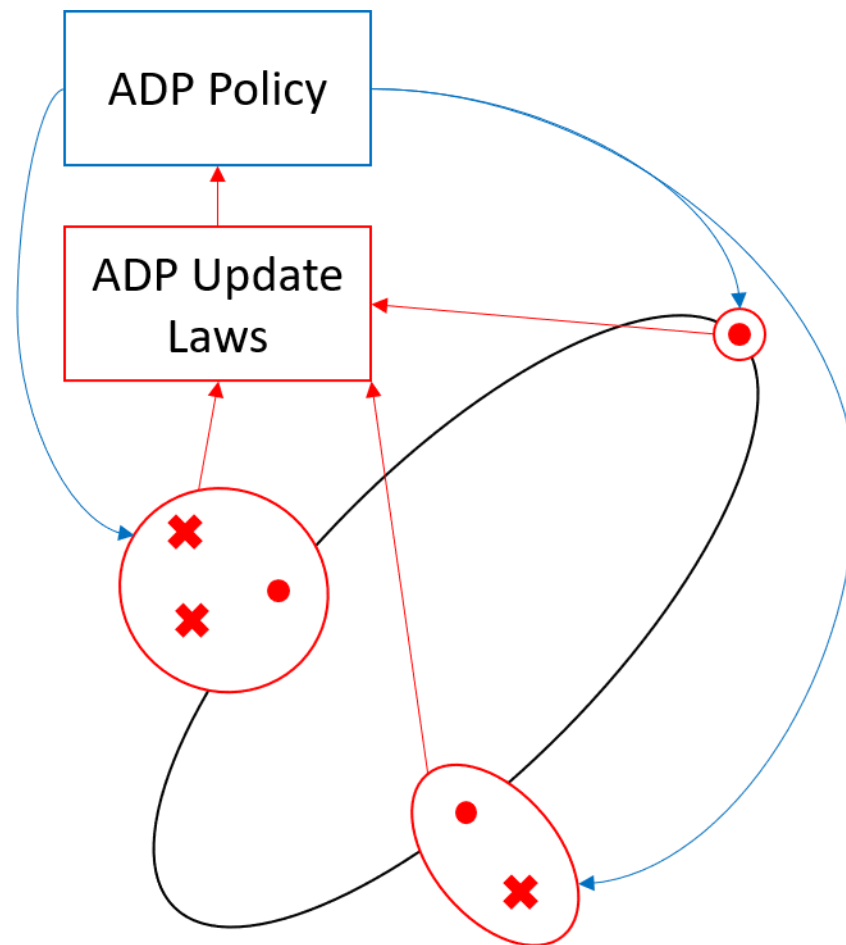
²Munitions Directorate, Eglin AFB, Air Force Research Laboratory

³Department of Aero. and Astro., Massachusetts Institute of Technology



Cooperative BE Extrapolation

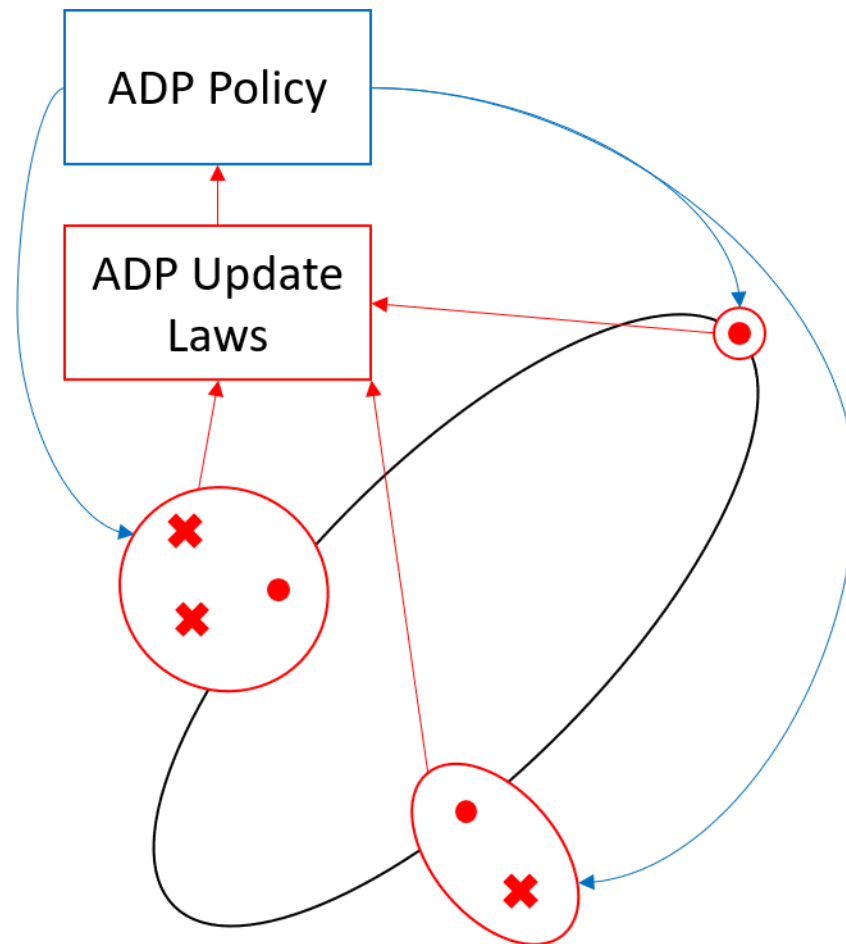
- Developed Method
 - Homogeneous agents
 - i.e., identical dynamics
 - Each agent computes one on- and some off-trajectory BE extrapolation points.
 - BE extrapolation is used in update laws.
 - Divides computational expense of BE extrapolation.





Cooperative BE Extrapolation

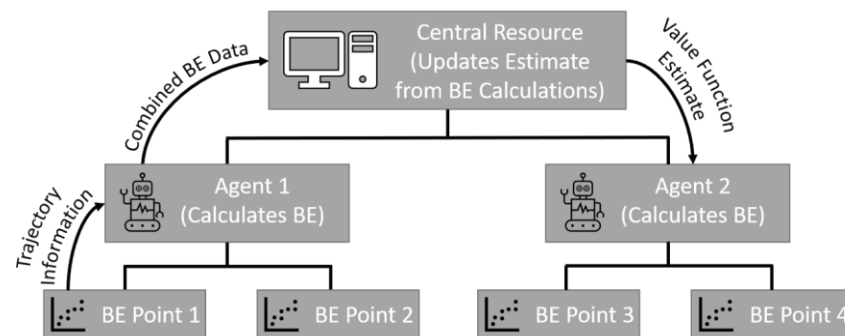
- Why Cooperative?
 - Share data between agents
 - Train DNN system ID
 - Collect data faster
 - More agents means more data
 - Value Function Appx.
 - Without excessive computational loads on one agent

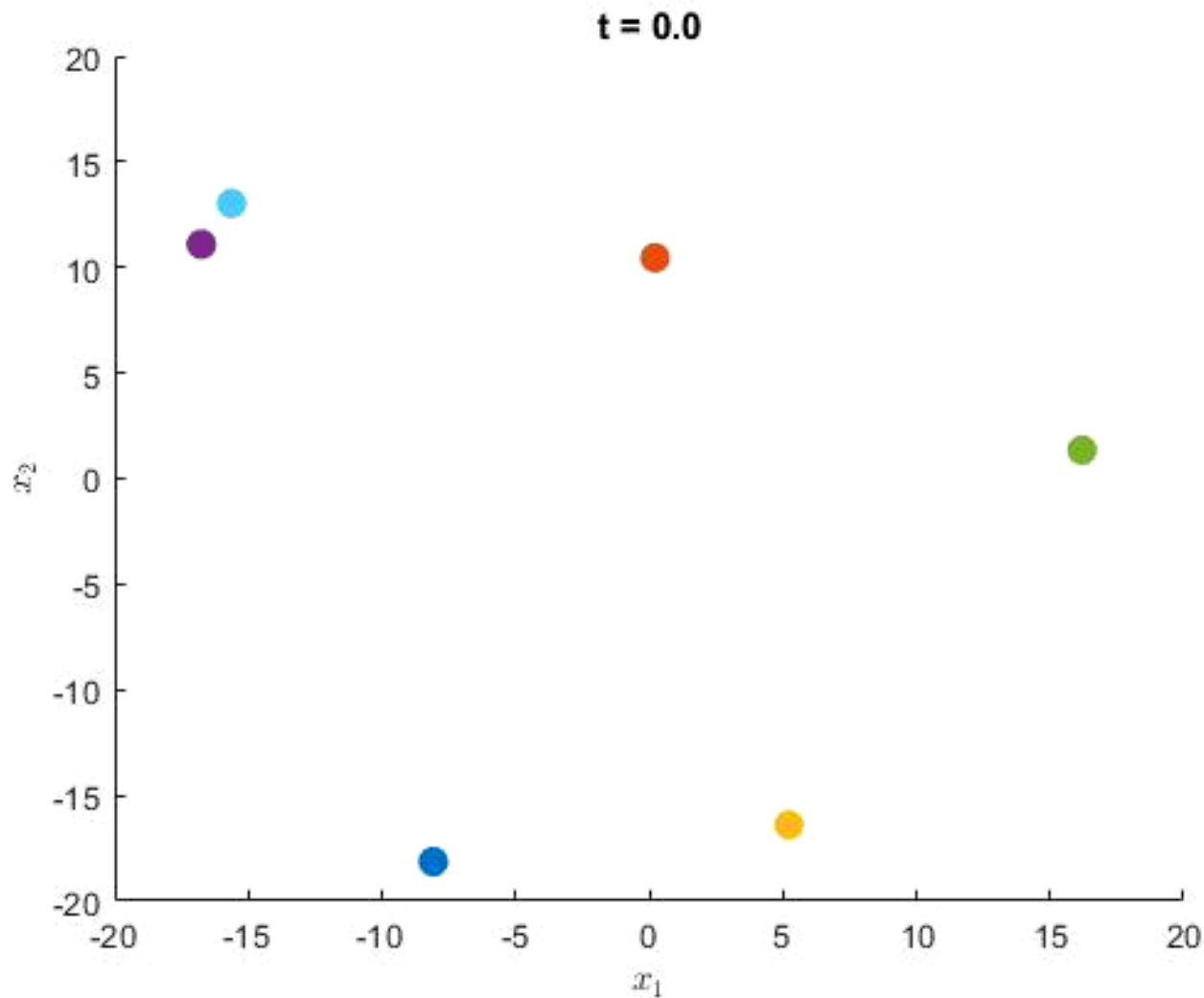




Cooperative BE Extrapolation

- Developed Method
 - Alternative view
 - Each agent calculates BE trajectory data.
 - Agent communicates BE data to central resource
 - Central resource communicates ADP policy to agents





- Future Work
 - DNN ADP
 - Perform DNN-based value function approximation
 - Currently single-layer NN
 - Cooperative ADP
 - Extend to heterogeneous agents
 - How do we use share data between agents to cooperatively accomplish a global objective?

Questions?

MAXGREENE12@UFL.EDU

