

Preference-balancing Motion Planning under Stochastic Disturbances

Aleksandra Faust¹, Nick Malone², and Lydia Tapia²

¹Sandia National Laboratories

²University of New Mexico

Preference Balancing Tasks (PBTs)

- ▶ Robotic motion
 - ▶ Control-affine system
 - ▶ Continuous states and actions
 - ▶ High-dimensional
- ▶ Complicated dynamics
 - ▶ Difficult demonstration
 - ▶ Lack of motion primitives
- ▶ Described with preferences



Image: SpaceX

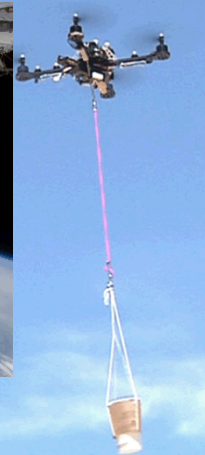
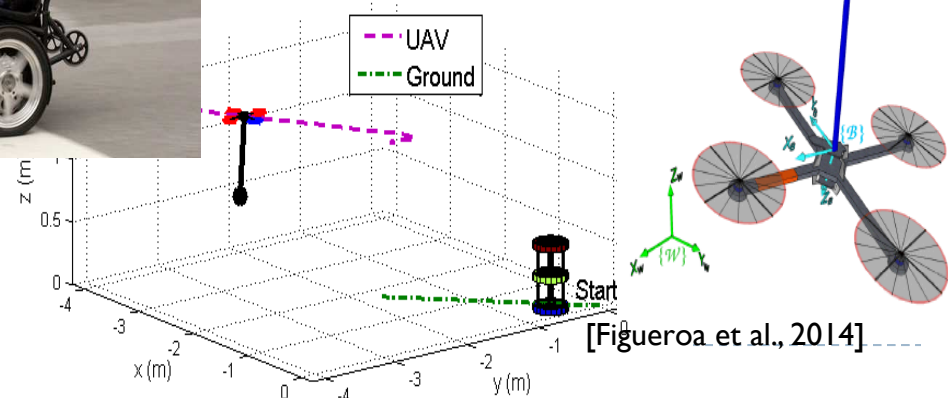


Image:
NASA/JPL/Caltech



Image: GM

[Faust et al., 2013]



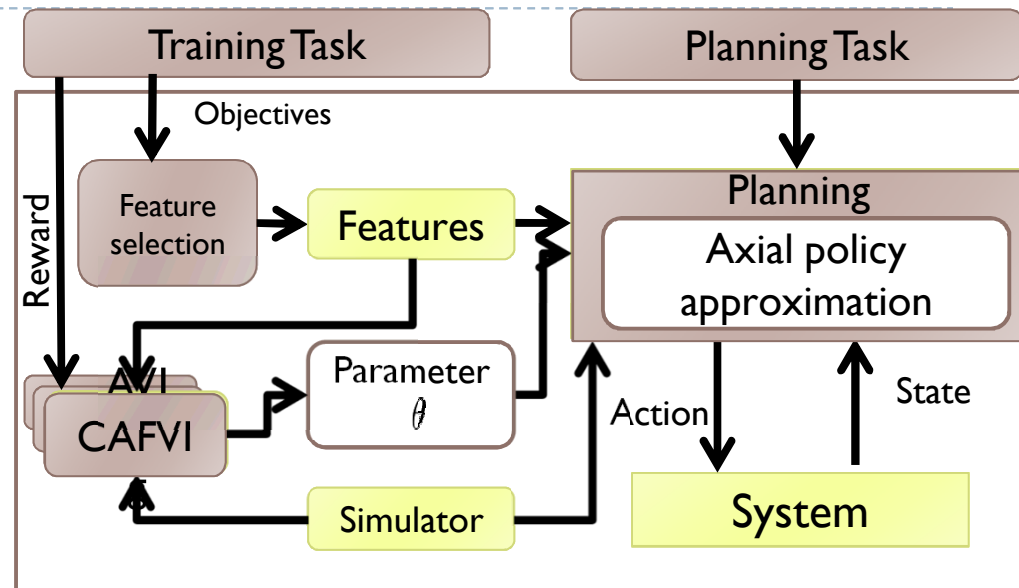
[Figueroa et al., 2014]

[Faust et al., in press]

PrEference Appraisal Reinforcement Learning (PEARL)

[Faust, 2014]

- ▶ Learns to perform PBTs*
- ▶ Batch reinforcement learning (RL)
- ▶ Learning
 - ▶ Continuous action fitted value iteration (CAFVI) [Faust, et al, in press]
 - ▶ Linear map state-value function approximation
 - ▶ Features are squared preferences $V(\mathbf{s}) = \boldsymbol{\theta}^T \mathbf{F}(\mathbf{s})$
- ▶ Planning
 - ▶ Generates trajectory
 - ▶ Real-time, one step at the time
 - ▶ Axial sum policy approximation [Faust, et al, in press]

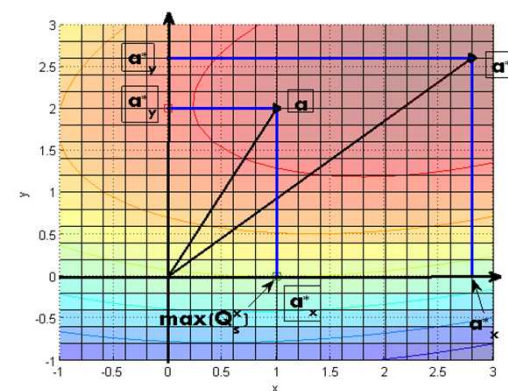
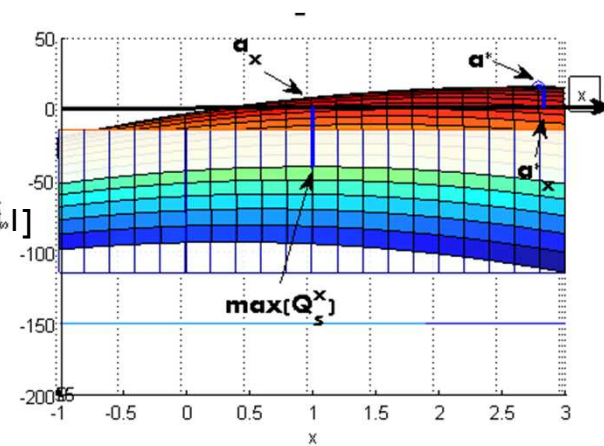


Works only with deterministic systems

*Preference Balancing Tasks

Axial sum policy approximation

- ▶ Greedy policy $h^*(x) = \operatorname{argmax}_{u \in U} V(D(x, u))$ in continuous spaces
 - ▶ Sampling-based search space narrowing [Busoniu et al. 2013] [Mansley et al. 2011] [Walsh et al. 2010]
 - ▶ Gradient descent [Hasselt et al. 2012]
- ▶ Deterministic axial sum policy (DAS) [Faust, et al, in press]
 - ▶ Interpolate action-value function along each the axis
 - ▶ Find maximum
 - ▶ Combine with a vector or convex sum
- ▶ Sufficient conditions for convergence to goal
 - ▶ Control-affine system with bounded drift [Faust, et al, in press]
 - ▶ Squared-features
 - ▶ Negative weights



Extend DAS to work under external disturbances.

Least Squares Axial Policy Approximation (LSAPA)

► Problem

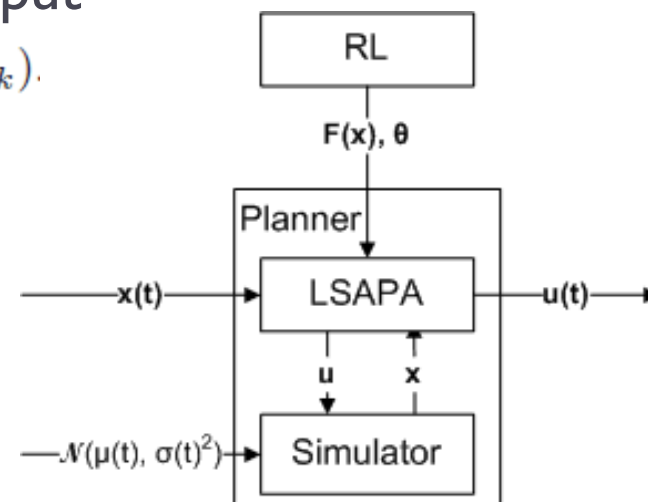
- PBTs $V(x) = \sum_{i=1}^{d_g} \theta_i F_i(x).$
- Control-affine system with external input disturbance $x_{k+1} = f(x_k) + g(x_k)(u_k + \eta_k).$

► Learning

- Deterministic CAFVI

► Planning

- Estimate disturbance in real-time
- Least Squares Axial Policy Approximation selects an action at every time step
- Adapts to observed disturbance



Least Squares Axial Policy Approximation Continued

► On each input axis

► Action-value function $Q_{\mathbf{x},i}(u) = \mathbf{p}_i^T [u^2 \ u \ 1]^T$

► Collect d_n dynamic samples $U_i = [u_{1,i} \ \dots \ u_{d_n,i}]^T$

► Calculate $X_i = [\mathbf{x}'_{1,i} \ \dots \ \mathbf{x}'_{d_n,i}]^T$

$$Q_i = [Q_{\mathbf{x},1}(u_{1,i}) \ \dots \ Q_{\mathbf{x},d_n}(u_{d_n,i})]^T$$

$$Q_{\mathbf{x},j}(u_{j,i}) = \theta^T F(\mathbf{x}'_{j,i}) \quad j = 1, \dots, d_n$$

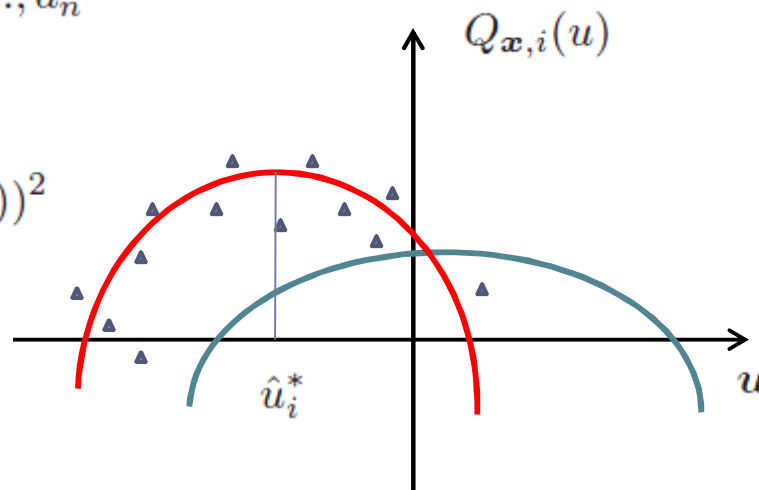
► Solve supervised ML problem

$$\hat{\mathbf{p}}_i = \underset{\mathbf{p}_i}{\operatorname{argmin}} \sum_{j=1}^{d_n} (C_{j,i} \mathbf{p}_i - Q_{\mathbf{x},j}(u_{j,i}))^2$$

$$C_i = \begin{bmatrix} (u_{1,i})^2 & u_{1,i} & 1 \\ (u_{2,i})^2 & u_{2,i} & 1 \\ \vdots & \vdots & \vdots \\ (u_{d_n,i})^2 & u_{d_n,i} & 1 \end{bmatrix}$$

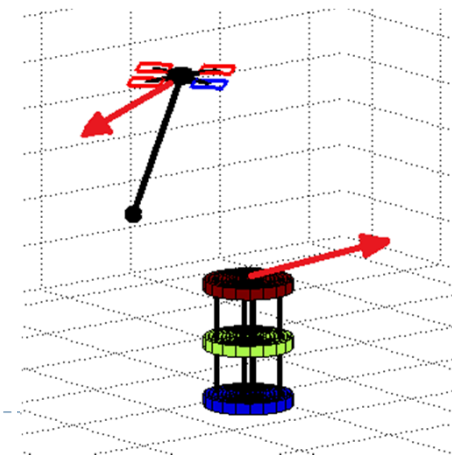
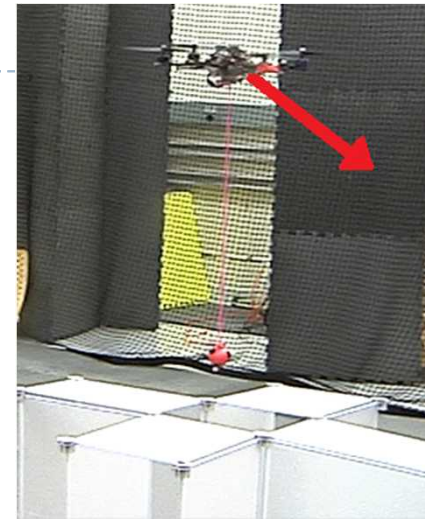
► Find maximum

$$\hat{u}_i^* = -\frac{\hat{p}_{1,i}}{2\hat{p}_{2,i}}$$



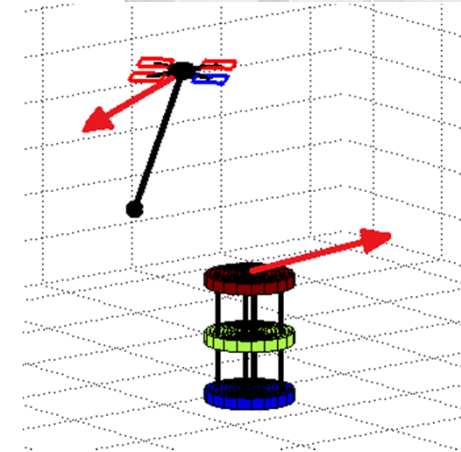
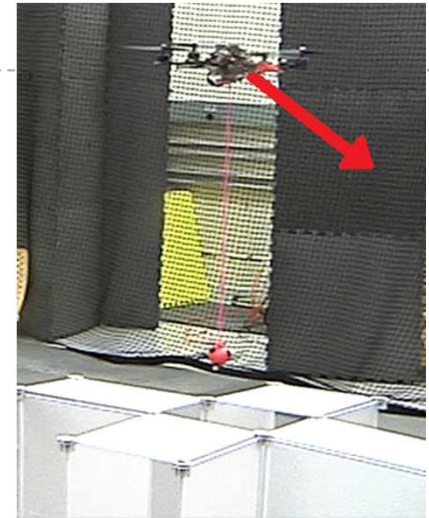
Results

- ▶ Related work on UAV robust control
 - ▶ Trajectory tracking [Alexis et al. 2010]
 - ▶ Harmonic potential fields [Masoud 2011]
 - ▶ Low-level controllers [DeCastro and Kess-Gazit 2013]
 - ▶ Trajectory libraries [Majumdar and Tedrake 2013]
 - ▶ Blimp path planning with dynamic programming [Kawano 2011]
- ▶ Coffee-delivery Tasks
 - ▶ Swing-free aerial cargo delivery
 - ▶ Rendezvous task

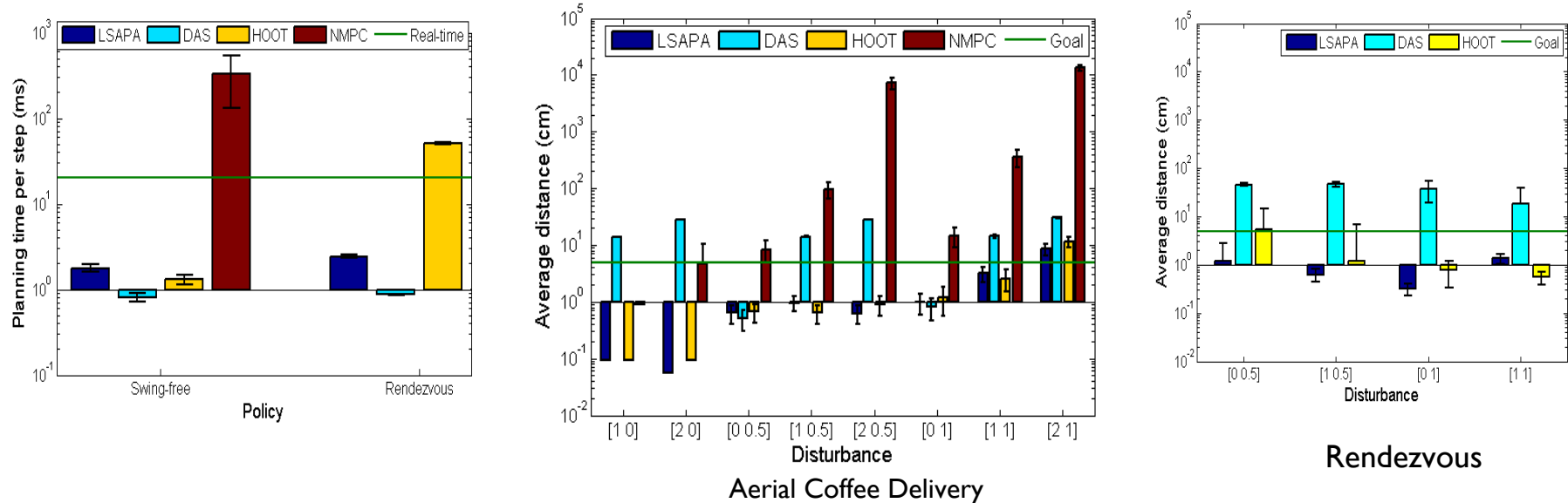


Coffee Delivery: Setup

- ▶ **Aerial Problem**
 - ▶ Holonomic cargo-bearing UAV
 - ▶ Bring the suspended load to the destination
 - ▶ Minimal residual load oscillations at arrival
- ▶ **Rendezvous Problem**
 - ▶ Holonomic cargo-bearing UAV and ground robot
 - ▶ Bring the suspended load to the ground robot
 - ▶ Minimal residual load oscillations at arrival
- ▶ **Preferences, reduce**
 - ▶ Distance from the destination
 - ▶ Vehicle's velocity
 - ▶ Load displacement
 - ▶ Load's velocity
- ▶ **MDP**
 - ▶ Aerial: 10-dimensional vector states, 3-dimensional actions
 - ▶ Rendezvous: 16-dimensional vector states, 5-dimensional actions



Trajectory Characteristics under Varying Disturbance



- Least Squares Axial Policy Approximation (LSAPA) [Faust et al., 2015]
- Deterministic Axial Sum (DAS) [Faust et al., 2014]
- HOOT [Mansley et al. '11]
- Nonlinear Model Predictive Control (NMPC) [Grune and Pannek, 2011]

LSAPA and DAS perform decision-making in real-time.

LSAPA reaches the goal for non-zero mean disturbances.

Swing-free aerial cargo delivery with disturbances

Preference-balancing Motion Planning under Stochastic Disturbances

Aleksandra Faust, Nick Malone, and Lydia Tapia
Department of Computer Science
University of New Mexico
<https://www.cs.unm.edu/amprg>

<http://www.cs.unm.edu/~afaust/movies/afaustlcra15.mp4>



Questions

- ▶ Thank you!

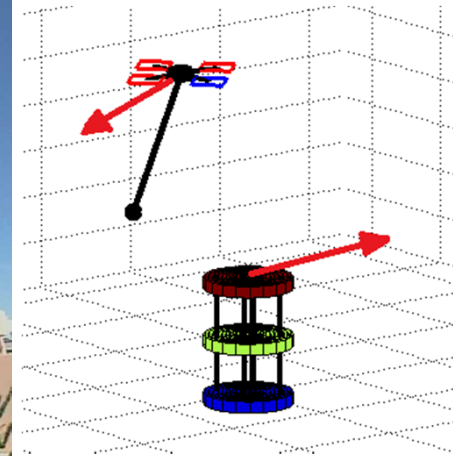
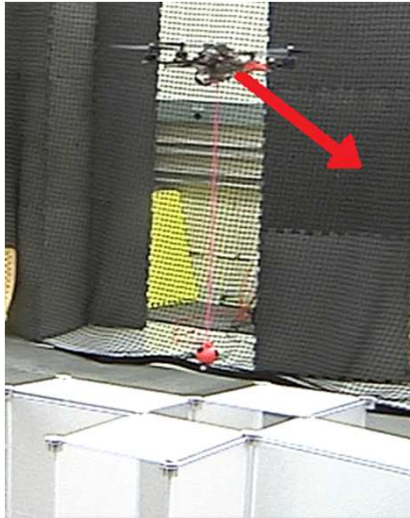
- ▶ Acknowledgements
 - ▶ Dr. Peter Ruymgaart for discussing disturbance modelling
 - ▶ Patricio Cruz for assisting with experiments
 - ▶ Reviewers for very helpful and constructive feedback

This work was in part supported by NM Space Grant and National Institutes of Health (NIH) Grant P20RR018754 to the Center for Evolutionary and Theoretical Immunology.



----- Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. -----





Preference-balancing Motion Planning under Stochastic Disturbances

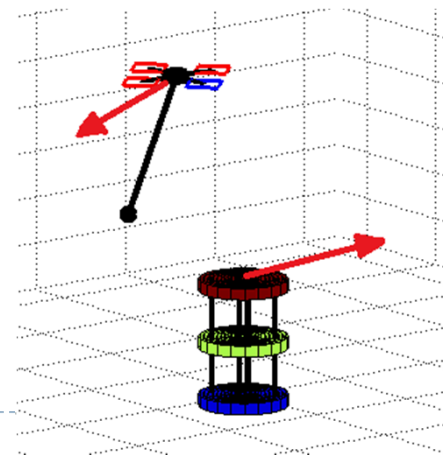
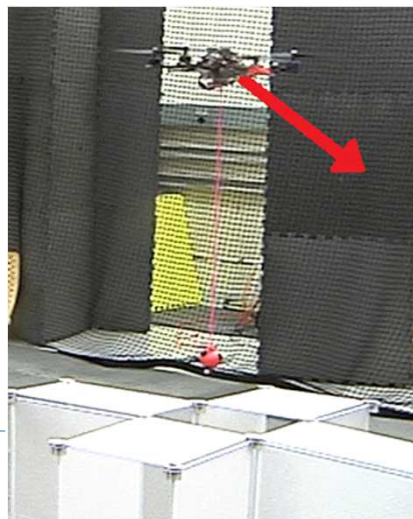
Aleksandra Faust¹, Nick Malone², and Lydia Tapia²

¹Sandia National Laboratories

²University of New Mexico

Preference-balancing Motion Planning under Stochastic Disturbances Summary

- ▶ Reinforcement learning with no disturbances
- ▶ Online planning in the presence of disturbances
- ▶ Applicable for continuous actions
- ▶ Linear in the input dimensionality
- ▶ Works through a Least Squares Axial Sum Policy Approximation



Swing-free aerial cargo delivery with disturbances

Preference-balancing Motion Planning under Stochastic Disturbances

Aleksandra Faust, Nick Malone, and Lydia Tapia
Department of Computer Science
University of New Mexico
<https://www.cs.unm.edu/amprg>

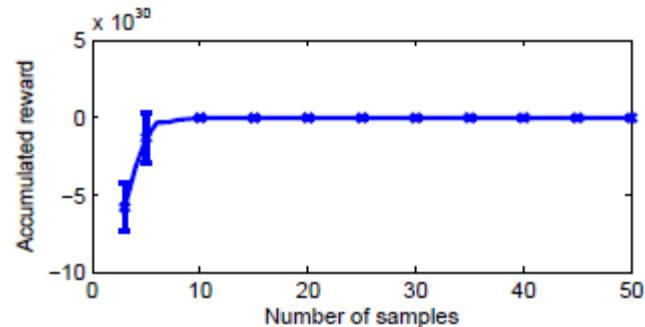
<http://www.cs.unm.edu/~afaust/movies/afaustlcra15.mp4>



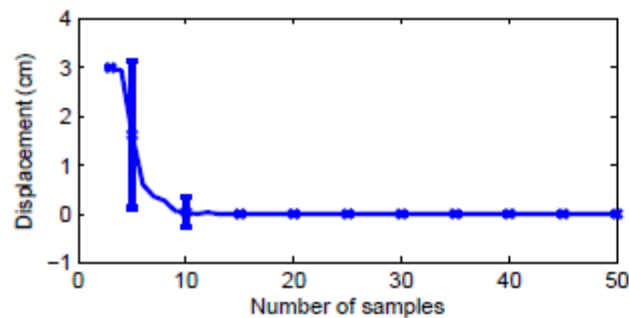


Flying Inverted Pendulum

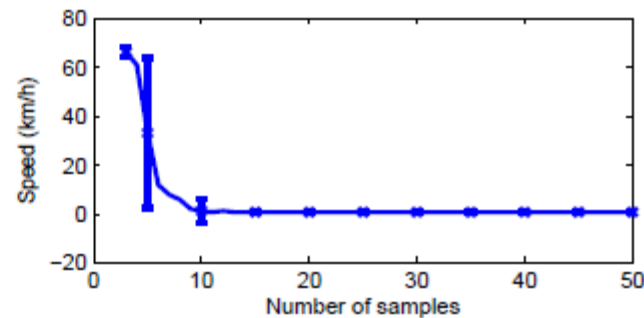
Number of samples need



(a) Accumulated reward



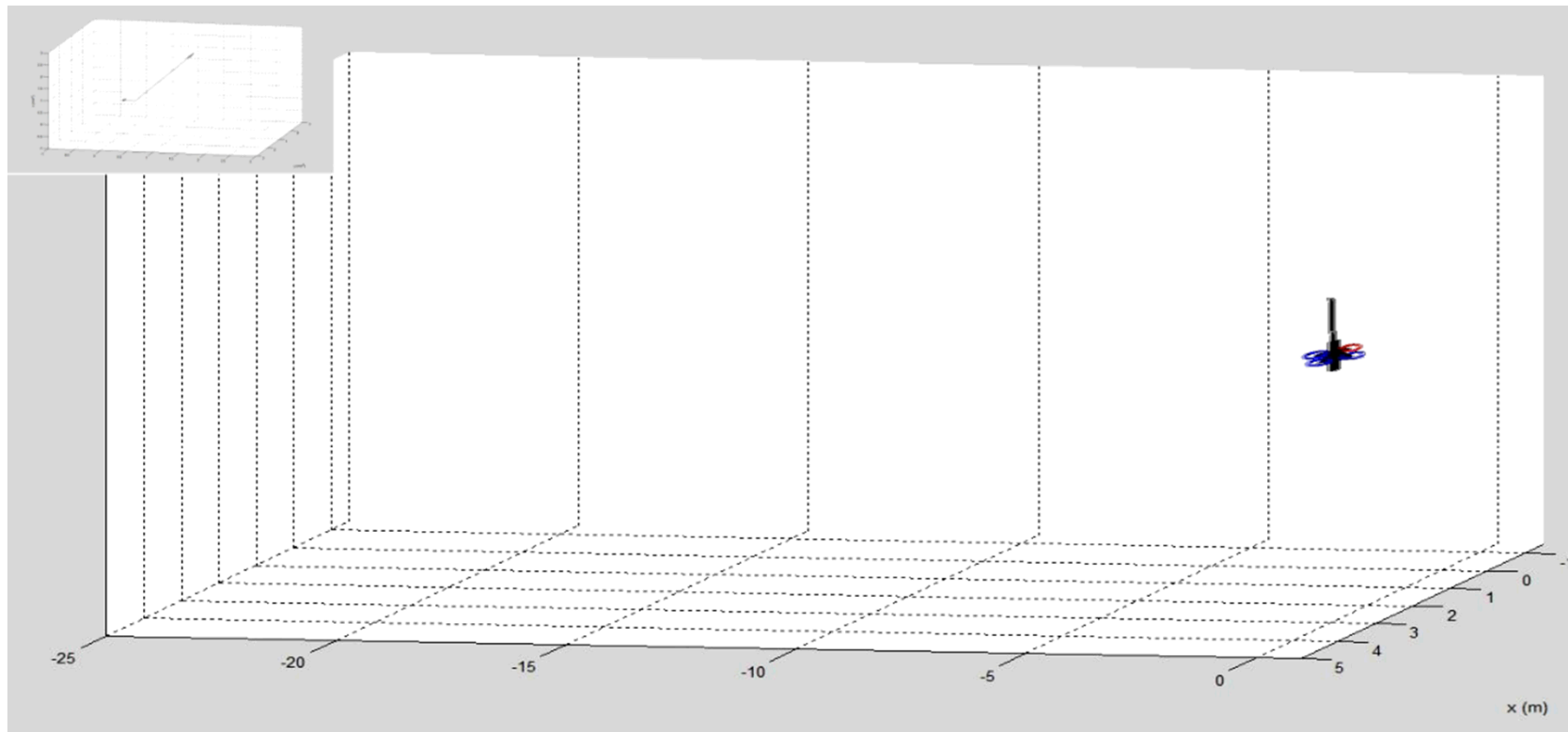
(b) Displacement



(c) Speed

Trajectory characteristics improve exponentially with number of samples.

Flying Inverted Pendulum



Task: *Flying inverted pendulum*

Planning: *Least squares axial policy approximation (LSAPA)*

Upper left: *Stochastic disturbance $\sim N(1,1)$*

Related work

▶ UAV Robust control

- ▶ Trajectory tracking [Alexis et al. 2010]
- ▶ Harmonic potential fields [Masoud 2011]
- ▶ Low-level controllers [DeCastro and Kess-Gazit 2013]
- ▶ Trajectory libraries [Majumdar and Tedrake 2013]
- ▶ Blimp path planning with dynamic programming [Kawano 2011]

▶ Greedy policy approximation

- ▶ Sampling based planning search space narrowing [Mansley et al. 2011]
[Busoniu et al. 2013] [Walsh et al. 2010] [Bubek et al. 2011]
- ▶ Gradient descent [Hasselt et al. 2012]