



Learning for Control of Synthetic and Cyborg

PI' s: Pieter Abbeel (PI), Ron Fearing (co-PI),
Students: Nimbus Goehausen, Woody Hoburg,

Summary

Objective

Development of learning and adaptation capabilities that will enable operation of synthetic and cyborg insects in complicated environments, such as collapsed buildings.

Platforms



Cyborg Beetle



Synthetic Crawler

Year 1 Results

Cyborg beetle:

Setup, flight initiation, cessation

Synthetic crawler:



On-board electronics: video, gyro, accel.

Learning and adaptation:

Faster-learning policy gradient method

Year 2 Results

Cyborg beetle:

Wing-folding muscle, Conditioning

Synthetic crawler:



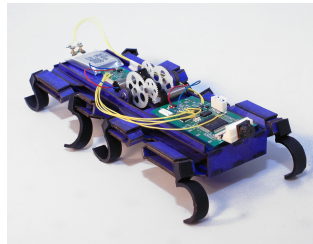
Complete crawler, Preliminary control

Learning and adaptation:

Risk-sensitive reinforcement learning

Synthetic Crawler

OctoRoACH



Specs

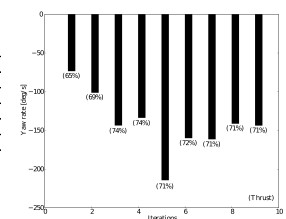
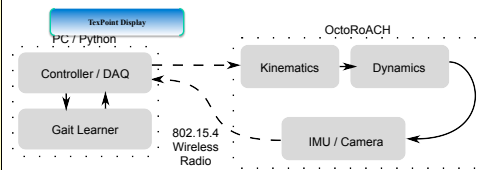
Sensing:

- Inertial Measurements Unit
- Back-EMF
- Cell-phone camera

Control:

- 40 MHz dsPIC microcontroller
- Differential drive-train
- Manually-tuned PID controllers
- 802.15.4 Wireless Radio

Learning to turn – Policy Search



Preliminary results were obtained using the policy search algorithm PEGASUS [Ng & Jordan 2000] for getting the optimal policy for turning on carpet.

From running gait at 10% thrust, robot turned fastest when other leg was run at 71% thrust.

Risk-Sensitive RL (1)

Motivation

- Expected return is often optimized
- Reasons: Law of large numbers + Mathematical convenience
- Not necessarily good criterion when only getting to execute once
- Motivating example: consider buying tickets to fly to a very important meeting, for which it would be disastrous if you arrived late. Some flights arrive on time more often than others, and these delays might be amplified if you miss connecting flights. With these risks in mind, would you rather take a route with an expected travel time of 13:48:09 and no further guarantees, would you choose a route that takes 13:53:28 on average with a standard deviation of 0:58:28, or would you prefer a route that takes less than 17:32:04 with 99% probability?

Preliminaries

- Markov Decision Process

- State space S
- Action space A
- Stochastic dynamics
- Stochastic reward function R

- Return from time t onwards

$$V^t(\pi) = \sum_{h=1}^{\infty} \gamma^{h-1} R^h(S_t, \pi^h(S_t), S_{t+1})$$

- "Optimize" $V(\pi) := V^0(\pi)$ (random var.)
- Traditionally: optimize expected return:

$$\max_{\pi} E_s[V(\pi)] \quad \text{where} \quad E_s[\diamond] := E[\diamond | S_0 = s]$$

Insects in Uncertain Dynamic Environments

Michel Maharbiz (co-PI)

Amol Jadhav, Svet Kolev, Hirotaka Sato, Jie Tang



Cyborg Beetle

All beetles were handled according to long-established institutional and scientific guidelines for experimenting with animals.
All insects in our lab are treated with respect and care.
Animals are NOT devices.

Goal

Create a reliable micro air vehicle (MAV) from a live insect

Minimum control functions:

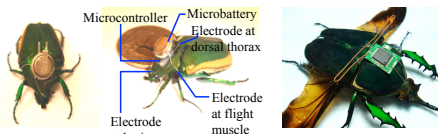
- On/Off
- Throttle
- Turning

Why Beetles?

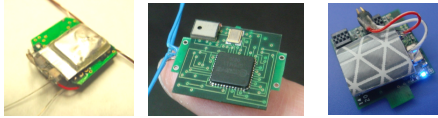
	Beetle	Moth	Locust	Fly
Payload capacity / gram	3.0	1.0	0.5	0.001
Adult life span / days	180	10	65	30
Survival under starving / days	30	-	-	3

Variety, size range: 1mm to 10cm
Strong flyers
maximum velocities: 7 ~ 14 km/hr
flight durations: 10 min ~ 3 hours
Easy to rear and breed
Generally harmless to humans

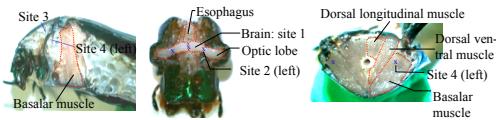
Hardware



V3.0 specs: 1.6 x 1.6 cm, 1.22 grams (0.67 g board + 0.55 g battery)
- controller, radio, oscillator, antenna, wingbeat mic, battery (4 v, 8.5 mAh)
- 1 day lifetime (sleep), 30 min lifetime flight operation (sleep/wake)



Implant Sites

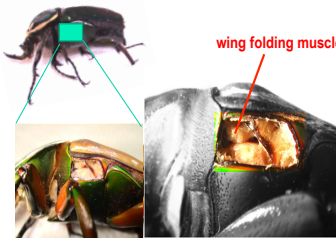


Flight Initiation, Cessation

- Initiation success rate: **97% (N= 30)**
- Cessation success rate: **100% (N=100)**
- Mean response time, cessation: **73 ms**
- Mean response time, initiation: **540 ms**

Wing-Folding Muscle

Anatomy

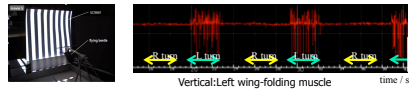


RE Snodgrass, Principles of Insect Morphology (1935), FW Darwin & JWS Pringle, Proc RSL (1959), K Ikeda et al., J Ins Physiol (1965)

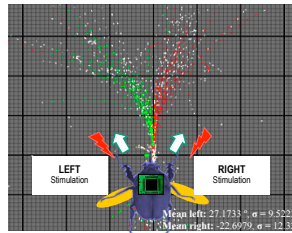
Findings

The beetles missing wing folding muscles still folded and unfolded the wings but were losing steerage.

→ Hypothesis: "wing-folding muscle" could be key to steering rather than wing-folding



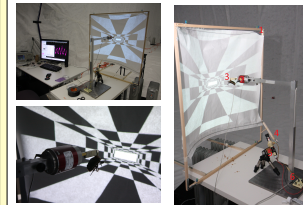
Train of the EMG spikes appeared only during the turn in the ipsilateral direction.



— left stimulation — right stimulation

Conditioning

Setup



1. Oval screen
2. Torque meter
3. Beetle attachment point
4. Servo shutter
5. 1000 mW infrared laser
6. Control signals + power

Findings

- Beetles seem to try and center a bright pattern when provided with the visual-torque feedback loop
- Attempts at overriding natural preferences (e.g., reaction to optical flow from striped patterns) through conditioning have not been a success

Plans

- Set up virtual environment with tunnel-like surroundings for beetle.
- Hypothesis: beetle will center bright end of tunnel
- Set up splits in the tunnel, with the two sides of the split naturally being equally good for the beetle. Condition the beetle to prefer one type.
- Move to stimuli that can be "mounted" onto the beetle itself

Where We're Headed

Low-level feedback

Graded turning control based on either or both of wing-folding muscle stimulation and conditioning

Miniaturization to fly



Risk-Sensitive RL (2)

Proposed Objective

- The Chernoff functional:

$$C_\delta^2[V(\pi)] = \sup_{\theta > 0} \theta^{-1} (\log \delta - \log E_\pi [e^{-\theta V(\pi)}])$$

Theorem 1. Let X be a random variable which has a moment generating function, that is $E[\exp(\theta X)]$ is finite in a neighborhood of $\theta = 0$, and $\delta \in [0, 1]$. Then, the Chernoff functional of this random variable, $C_\delta^2[X]$, is well defined and has the following properties:

- $P[X \leq C_\delta^2[X]] \leq \delta$
- $C_\delta^2[X] = E[X]$
- $\lim_{\delta \rightarrow 0} C_\delta^2[X] = \inf\{x : P\{X \leq x\} > 0\}$ which could be $-\infty$.
- As $\delta \rightarrow 1$, $C_\delta^2[X] \approx E[X] - \sqrt{2 \log(1/\delta) \text{Var}[X]}$
- $C_\delta^2[X] = E[X] - \sqrt{2 \log(1/\delta) \text{Var}[X]}$ if X is Gaussian.
- $C_\delta^2[X]$ is a smooth, increasing function of δ .

Cumulant Generating Function Iteration

Theorem 2. Consider an MDP with a finite number of states and actions, and finite horizon, such that rewards are independent conditional on states, $R^i(S_t, \pi^i(S_t), S_{t+1}) \perp V^{i+1}(\pi^i) | S_t, S_{t+1}$, and they have finite conditional moment generating functions, $E[\exp(\theta R_{t+1}^i) | S_t, S_{t+1}]$, in a neighborhood of $\theta = 0$. Then we can find $\pi^* = \arg \min_\pi \log E_\pi[\exp(-\theta V(\pi))]$ by the iteration below, and the moment generating function $E_\pi[\exp(-\theta V(\pi))]$ is finite in a neighborhood of $\theta = 0$.

$$q^i(s, a) := \log \sum_{\pi} \exp(\log p_{\pi, a, s} + \log E_\pi [e^{-\theta V(\pi)}] + u^{i+1}(s)) \quad (5)$$

$$u^i(s) := \min_a q^i(s, a) = \min_a \log E [e^{-\theta V(\pi)} | S_t = s], \quad \pi^i(s) := \arg \min_\pi q^i(s, a)$$

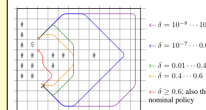
Finding the Global Optimum

$$\max_{\pi} C_\delta^2[V(\pi)] = \max_{\theta > 0} \sup_{\pi} \frac{\log \delta - \log E_\pi [e^{-\theta V(\pi)}]}{\theta} = \sup_{\theta > 0} \frac{\log \delta - \min_a \log E_\pi [e^{-\theta V(\pi)}]}{\theta}$$

The problem is not always convex in θ , but we will see that it becomes convex after a change of variables to $z = \theta^{-1}$. Recall the notation $f(\theta) := \min_a \log E_\pi [e^{-\theta V(\pi)}]$.

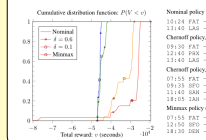
$$\max_{\pi} C_\delta^2[V(\pi)] = \sup_{\theta > 0} \theta^{-1} (\log \delta - f(\theta)) = \inf_{z > 0} (-z \log \delta + z f(z^{-1}))$$

Experiment 1: Gridworld



Each state corresponds to a square in the grid, and the actions, {N, NE, E, SE, S, SW, W, NW}, cause a move in the respective direction. In unmarked squares, the actor's intention is executed with probability 93. Each of the seven remaining actions might be executed instead, each with probability 0.01. Squares marked with # and \$ are absorbing states. The latter gives a reward of 35 when entered, and the former gives a penalty of 30. Any other state transitions cost 1. The horizon is 35.

Experiment 2: Air Travel Planning



- We assume that, in the case of a missed connecting flight due to delays, the airline will re-issue a ticket for the route of your choice leading to the original destination.
- We use historical data for February 2011, from the Office of Airline Information, Bureau of Transportation Statistics (BTS), available at www.bts.gov
- In general, we observed that lowering δ will produce policies with higher expected travel time, but lower standard deviation.