

Policy Gradient RL to learn fast walk

Goal: Enable an Aibo to walk as fast as possible

Policy Gradient RL to learn fast walk

Goal: Enable an Aibo to walk as fast as possible

- Start with a **parameterized walk**
- **Learn** fastest possible parameters

Policy Gradient RL to learn fast walk

Goal: Enable an Aibo to walk as fast as possible

- Start with a **parameterized walk**
- **Learn** fastest possible parameters
- **No simulator** available:
 - Learn entirely on robots
 - Minimal human intervention

Walking Aibos

- Walks that “come with” Aibo are **slow**
- **RoboCup** soccer: **25+ Aibo teams** internationally
 - Motivates faster walks

Walking Aibos

- Walks that “come with” Aibo are **slow**
- **RoboCup** soccer: **25+ Aibo teams** internationally
 - Motivates faster walks

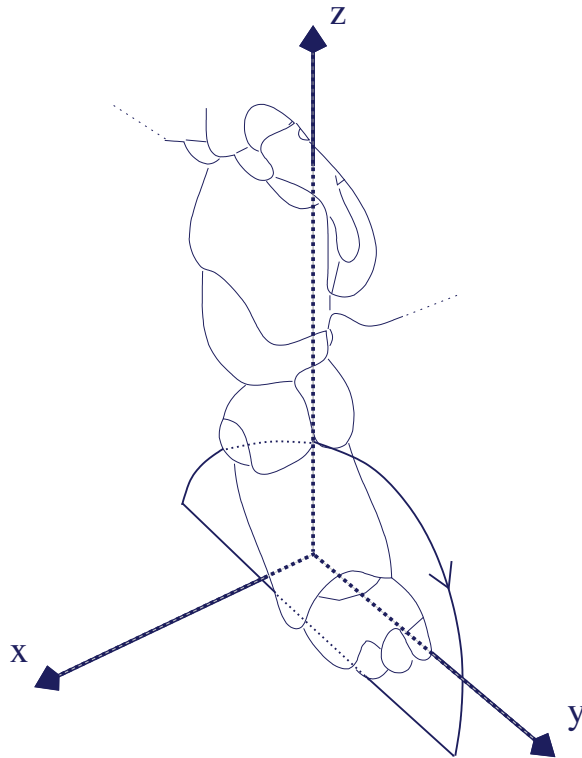
Hand-tuned gaits (2003)			Learned gaits	
German Team	UT Austin Villa	UNSW	Hornby et al. (1999)	Kim & Uther (2003)
230 mm/s	245	254	170	270 (± 5)

A Parameterized Walk

- Developed from scratch as part of **UT Austin Villa** 2003
- **Trot gait** with elliptical locus on each leg



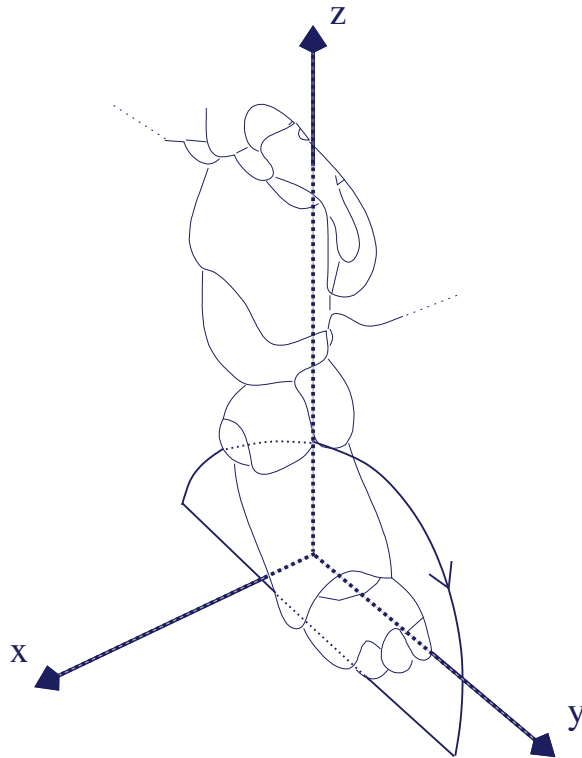
Locus Parameters



- Ellipse length
- Ellipse height
- Position on x axis
- Position on y axis
- Body height
- Timing values

12 continuous parameters

Locus Parameters



- Ellipse length
- Ellipse height
- Position on x axis
- Position on y axis
- Body height
- Timing values

12 continuous parameters

- Hand tuning by April, '03: **140 mm/s**
- Hand tuning by July, '03: **245 mm/s**

Experimental Setup

- Policy $\pi = \{\theta_1, \dots, \theta_{12}\}$, $V(\pi) = \text{walk speed}$ when using π

Experimental Setup

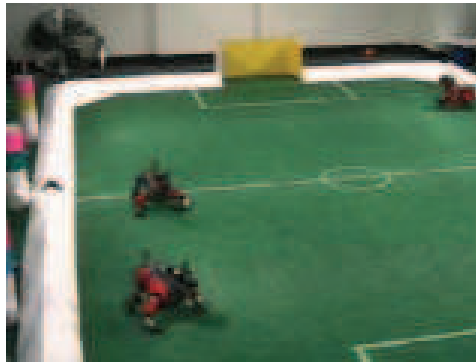
- Policy $\pi = \{\theta_1, \dots, \theta_{12}\}$, $V(\pi)$ = walk **speed** when using π
- Training Scenario
 - Robots **time themselves** traversing fixed distance
 - Multiple traversals (3) per policy to account for **noise**

Experimental Setup

- Policy $\pi = \{\theta_1, \dots, \theta_{12}\}$, $V(\pi)$ = walk **speed** when using π
- Training Scenario
 - Robots **time themselves** traversing fixed distance
 - Multiple traversals (3) per policy to account for **noise**
 - **Multiple robots** evaluate policies simultaneously
 - Off-board computer **collects results, assigns policies**

Experimental Setup

- Policy $\pi = \{\theta_1, \dots, \theta_{12}\}$, $V(\pi) = \text{walk speed}$ when using π
- Training Scenario
 - Robots **time themselves** traversing fixed distance
 - Multiple traversals (3) per policy to account for **noise**
 - **Multiple robots** evaluate policies simultaneously
 - Off-board computer **collects results, assigns policies**



No human intervention except battery changes

Policy Gradient RL

- From π want to move in direction of **gradient** of $V(\pi)$

Policy Gradient RL

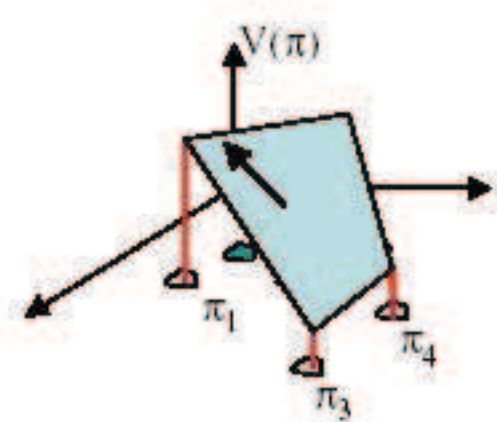
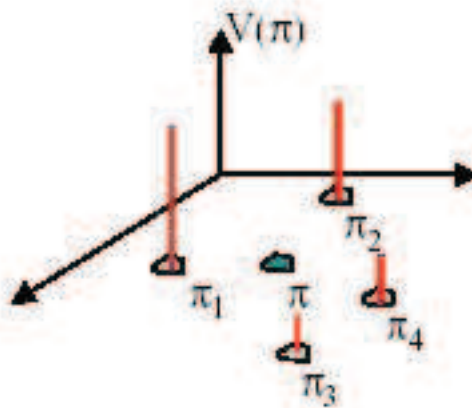
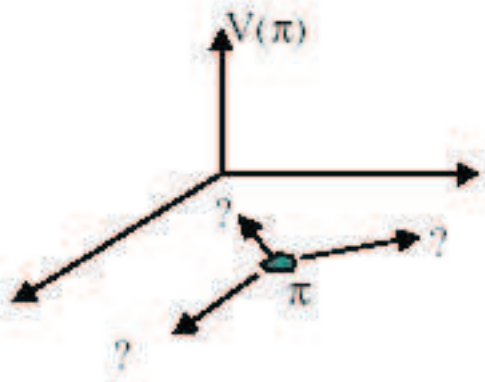
- From π want to move in direction of **gradient** of $V(\pi)$
 - Can't compute $\frac{\partial V(\pi)}{\partial \theta_i}$ directly: **estimate** empirically

Policy Gradient RL

- From π want to move in direction of **gradient** of $V(\pi)$
 - Can't compute $\frac{\partial V(\pi)}{\partial \theta_i}$ directly: **estimate** empirically
- Evaluate **neighboring policies** to estimate gradient
- Each trial randomly varies **every parameter**

Policy Gradient RL

- From π want to move in direction of **gradient** of $V(\pi)$
 - Can't compute $\frac{\partial V(\pi)}{\partial \theta_i}$ directly: **estimate** empirically
- Evaluate **neighboring policies** to estimate gradient
- Each trial randomly varies **every parameter**



Experiments

- Started from **stable**, but fairly slow gait
- Used **3 robots** simultaneously
- Each iteration takes 45 traversals, $7\frac{1}{2}$ minutes

Experiments

- Started from **stable**, but fairly slow gait
- Used **3 robots** simultaneously
- Each iteration takes 45 traversals, $7\frac{1}{2}$ minutes

Before learning

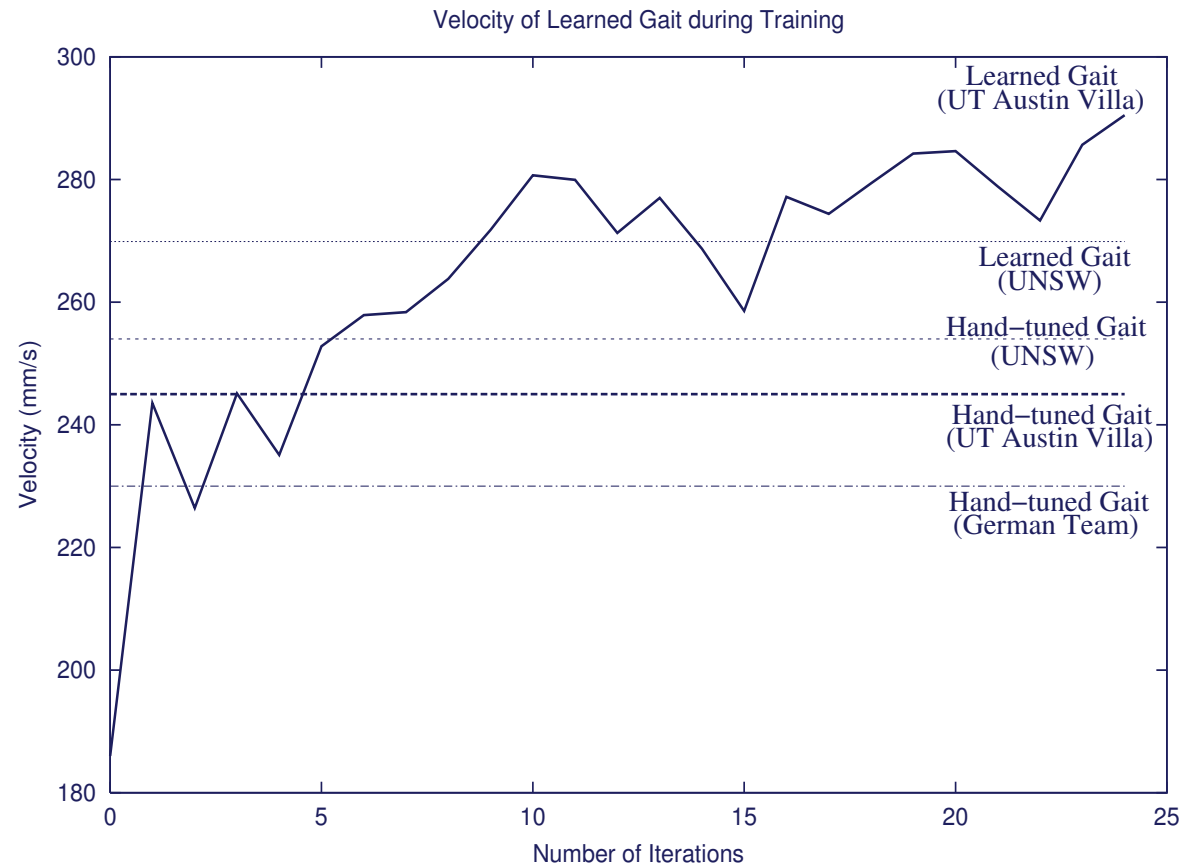


After learning

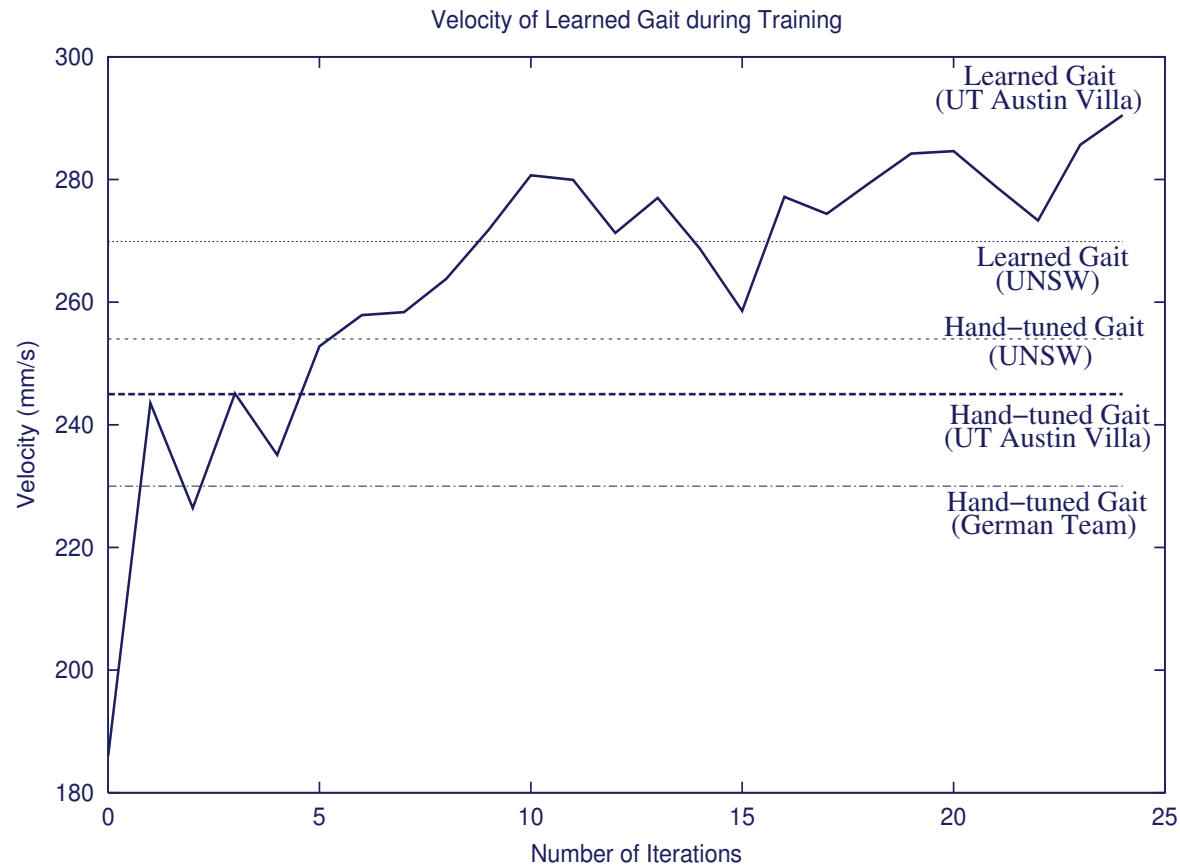


- 24 iterations = **1080 field traversals**, \approx **3 hours**

Results



Results

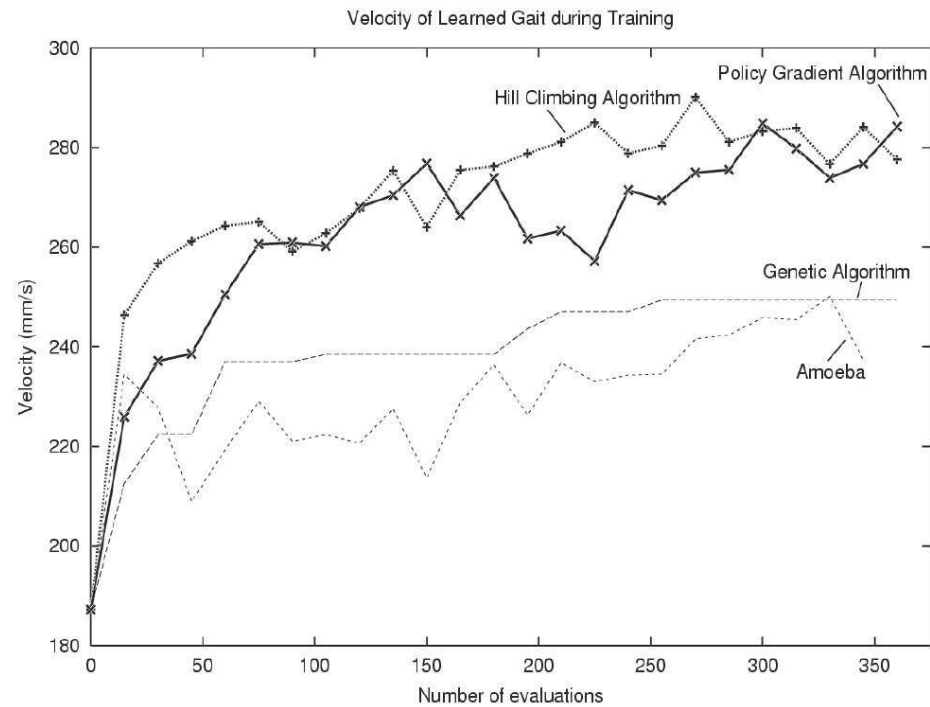


- Additional iterations didn't help
- Spikes: evaluation **noise**? large **step size**?

Learned Parameters

Parameter	Initial Value	ϵ	Best Value
Front ellipse:			
(height)	4.2	0.35	4.081
(x offset)	2.8	0.35	0.574
(y offset)	4.9	0.35	5.152
Rear ellipse:			
(height)	5.6	0.35	6.02
(x offset)	0.0	0.35	0.217
(y offset)	-2.8	0.35	-2.982
Ellipse length	4.893	0.35	5.285
Ellipse skew multiplier	0.035	0.175	0.049
Front height	7.7	0.35	7.483
Rear height	11.2	0.35	10.843
Time to move through locus	0.704	0.016	0.679
Time on ground	0.5	0.05	0.430

Algorithmic Comparison, Robot Port



Before learning



After learning



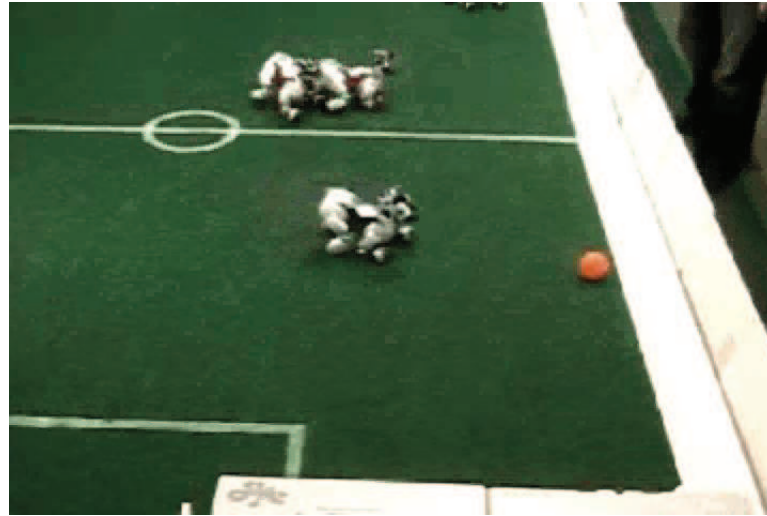
Summary

- Used policy gradient RL to **learn fastest Aibo walk**
- All learning done **on real robots**
- **No human intervention** (except battery changes)

Outline

- Learning sensor and action models (Stronger, S, '06)
- Machine learning for fast walking (Kohl, S, '04)
- **Learning to acquire the ball** (Fidelman, S, '06)
- Color constancy on mobile robots (Sridharan, S, '05)
- Autonomous Color Learning (Sridharan, S, '06)

Grasping the Ball



- **Three stages:** walk to ball; slow down; lower chin
- Head proprioception, IR chest sensor \mapsto ball distance
- Movement specified by **4 parameters**

Grasping the Ball

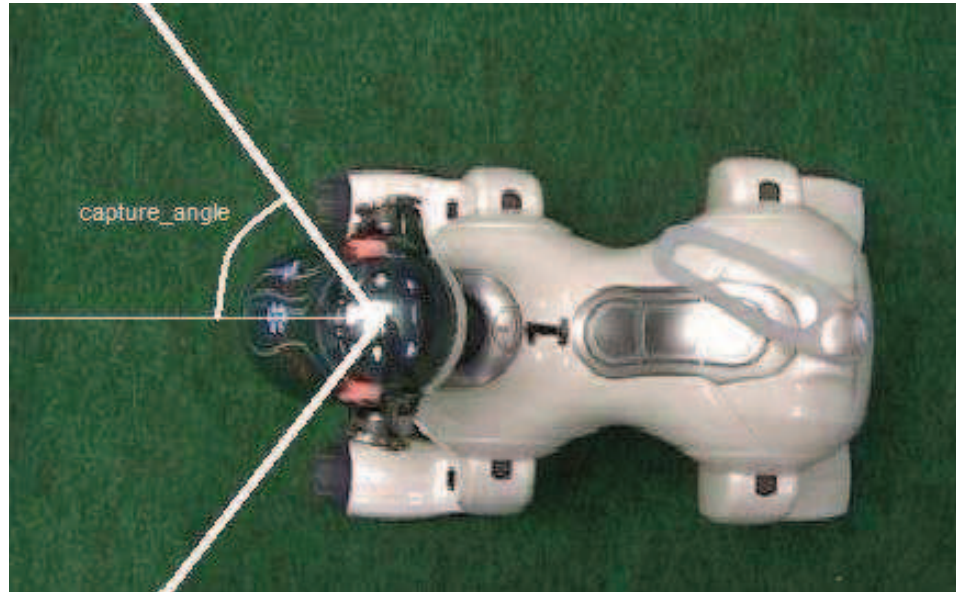


- **Three stages:** walk to ball; slow down; lower chin
- Head proprioception, IR chest sensor \mapsto ball distance
- Movement specified by **4 parameters**

Brittle!

Parameterization

- **slowdown_dist:** when to slow down
- **slowdown_factor:** how much to slow down
- **capture_angle:** when to stop turning



- **capture_dist:** when to put down head

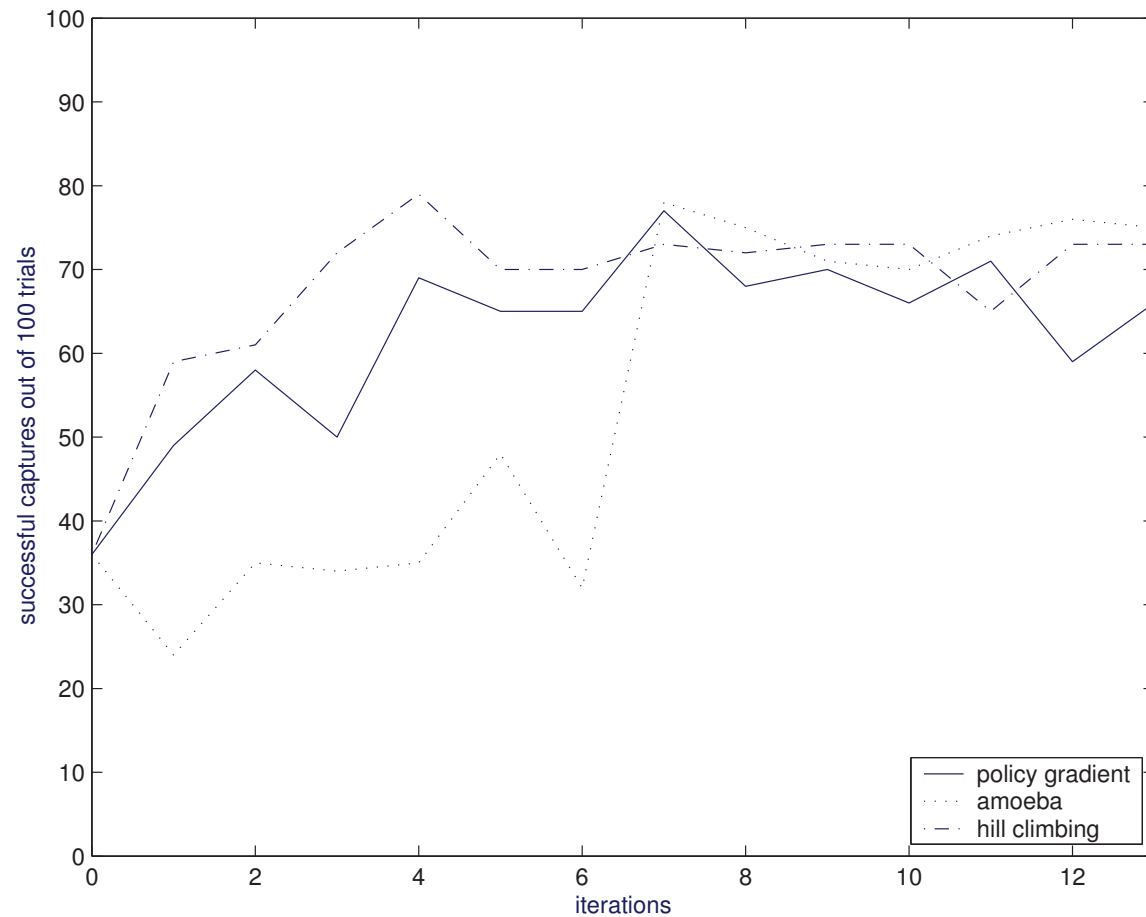
Learning the Chin Pinch

- **Binary, noisy** reinforcement signal: multiple trials
- Robot evaluates self: **no human intervention**



Results

- Evaluation of **policy gradient**, **hill climbing**, **amoeba**



What it learned



Policy	slowdown dist	slowdown factor	capture angle	capture dist	Success rate
Initial	200mm	0.7	15.0°	110mm	36%
Policy gradient	125mm	1	17.4°	152mm	64%
Amoeba	208mm	1	33.4°	162mm	69%
Hill climbing	240mm	1	35.0°	170mm	66%