

# Learning, Equilibria, Limitations, and Robots\*

Michael Bowling  
Computer Science Department  
Carnegie Mellon University

\*Joint work with Manuela Veloso

# Talk Outline

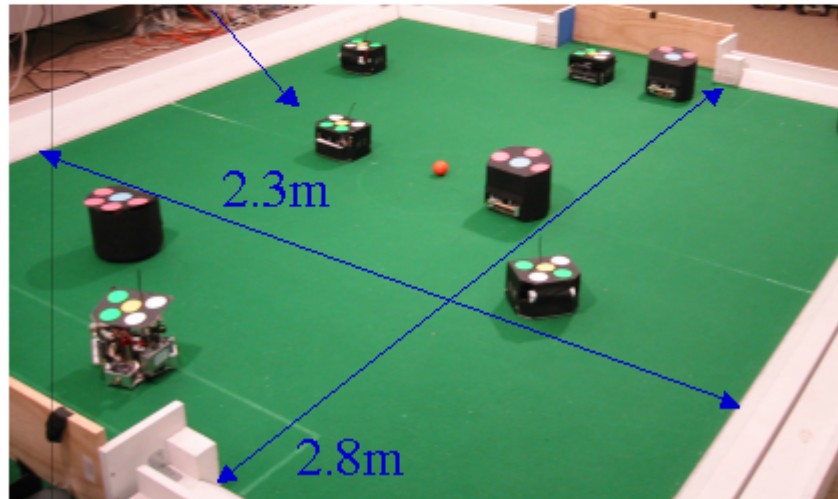
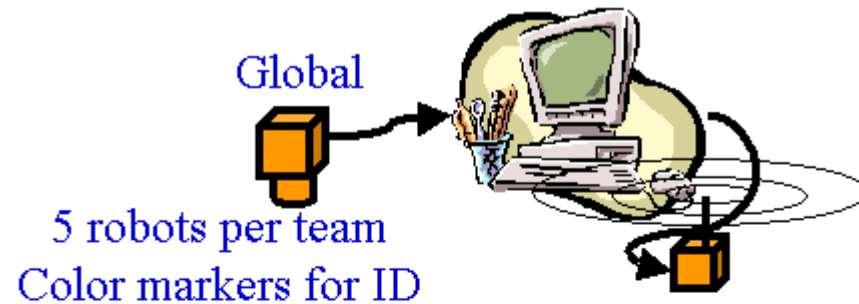
---

- Robots
  - A two robot, adversarial, concurrent learning problem.
  - The challenges for multiagent learning.
- Limitations and Equilibria
- Limitations and Learning

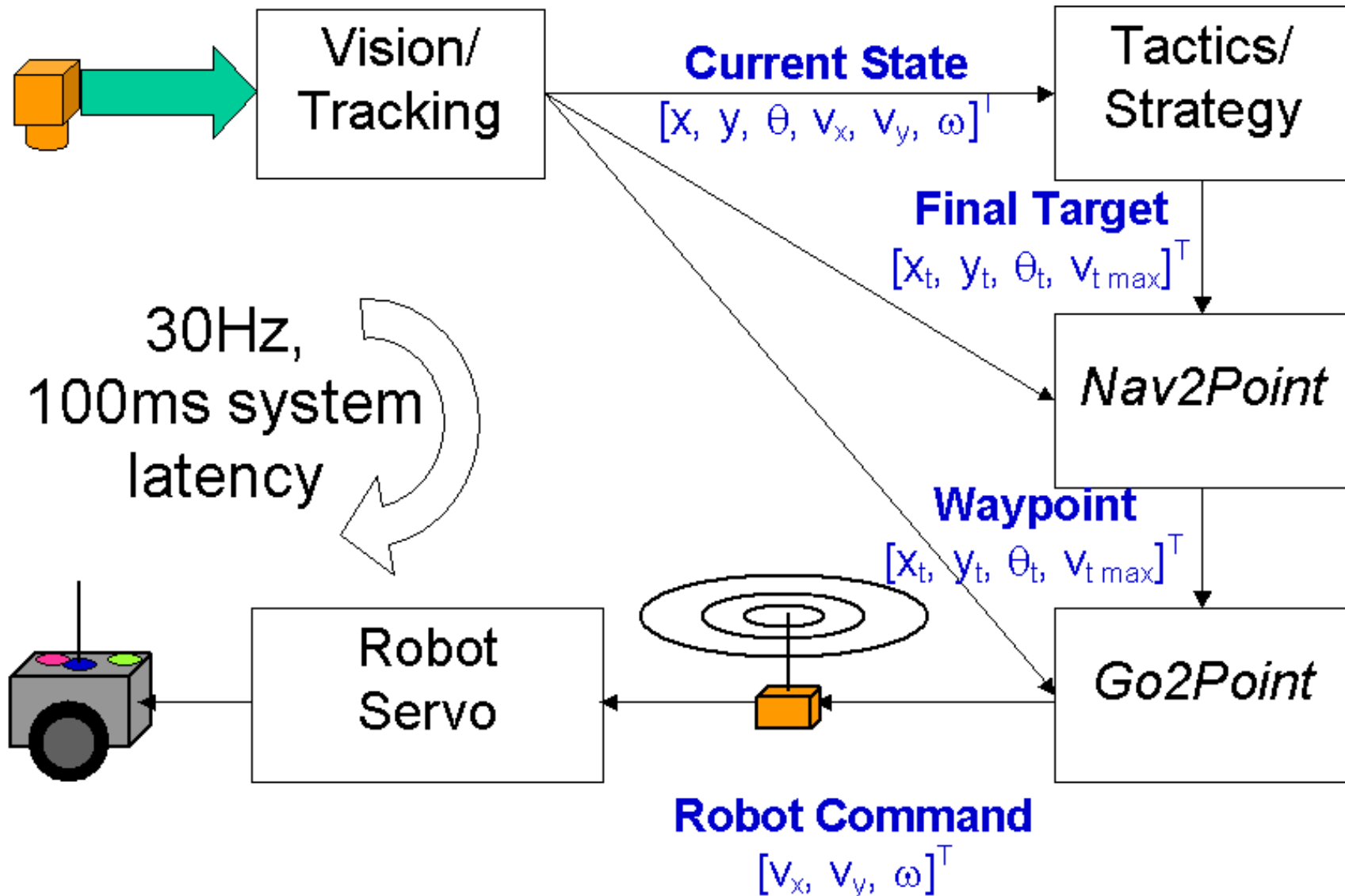


# The Domain = CMDragons = 1

---



# The Domain = CMDragons = 2



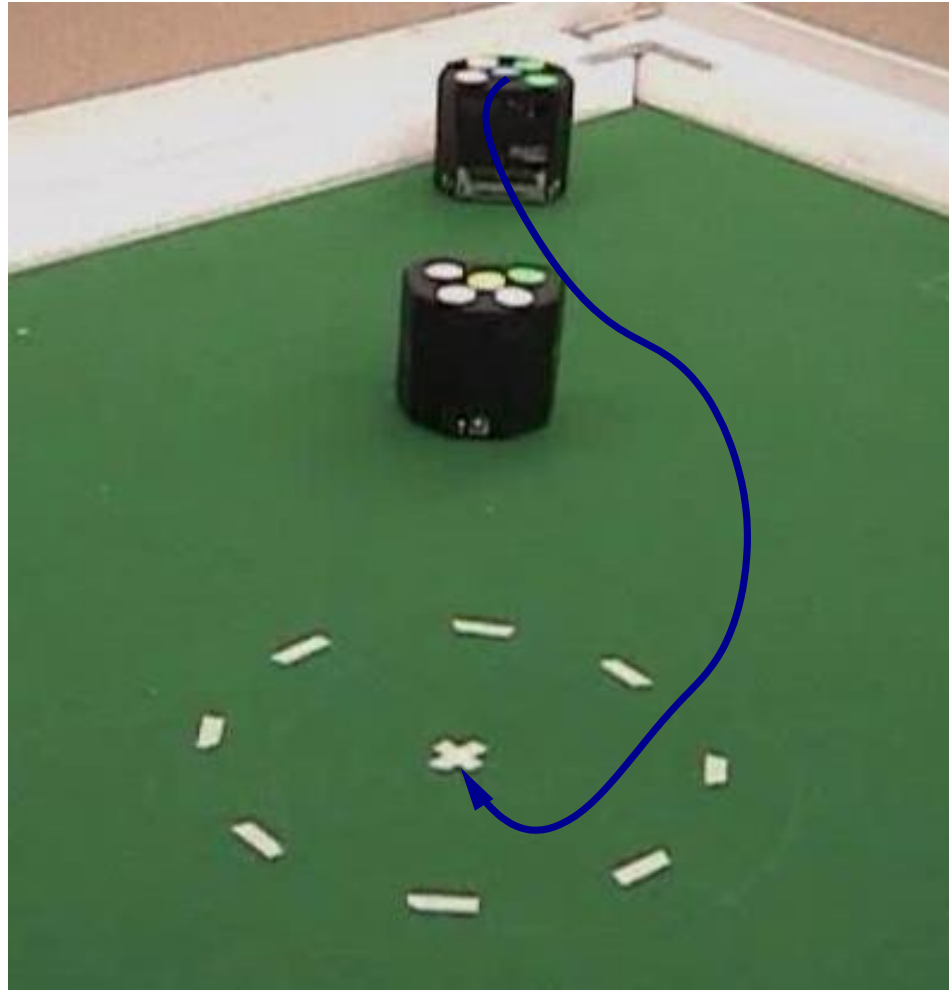
# The Task = Breakthrough

---



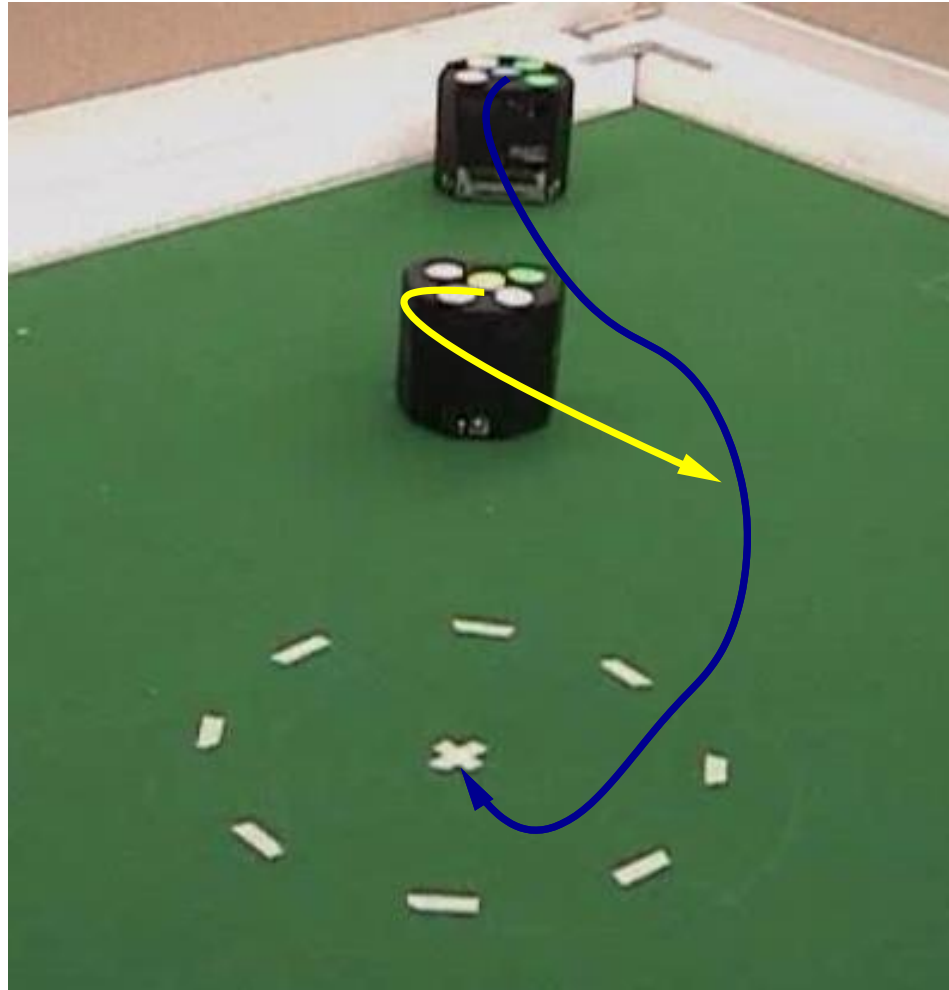
# The Task = Breakthrough

---



# The Task = Breakthrough

---



# The Challenges

---



# The Challenges

---

- Challenge #1: Continuous State and Action Spaces
  - Value function approximation, parameterized policies, state and temporal abstractions.
  - Limits agent behavior, sacrificing optimality.

# The Challenges

---

- Challenge #1: Continuous State and Action Spaces
  - Value function approximation, parameterized policies, state and temporal abstractions.
  - Limits agent behavior, sacrificing optimality.
- Challenge #2: Fixed Behavioral Components
  - Don't learn motion control or obstacle avoidance.
  - Limits agent behavior, sacrificing optimality.

# The Challenges

---

- Challenge #1: Continuous State and Action Spaces
  - Value function approximation, parameterized policies, state and temporal abstractions.
  - Limits agent behavior, sacrificing optimality.
- Challenge #2: Fixed Behavioral Components
  - Don't learn motion control or obstacle avoidance.
  - Limits agent behavior, sacrificing optimality.
- Challenge #3: Latency
  - Can predict our own state through latency, not others.
  - Asymmetric partial observability.
  - Limits agent behavior, sacrificing optimality.

# The Challenges = 1

---

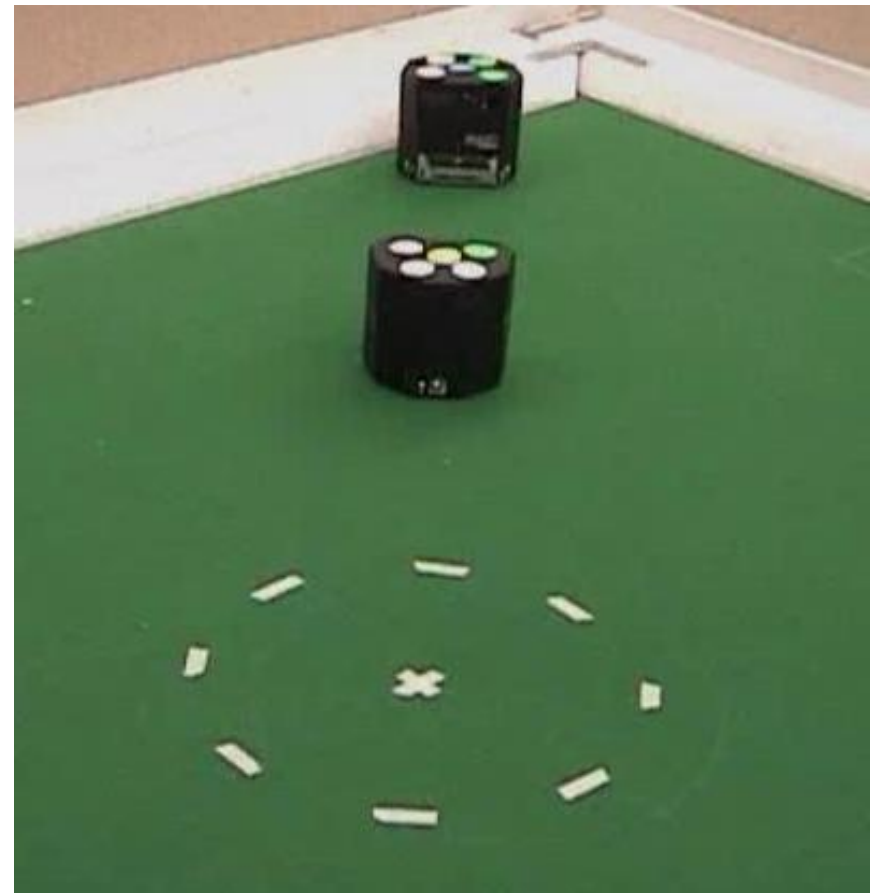
- Challenge #1: Continuous State and Action Spaces
- Challenge #2: Fixed Behavioral Components
- Challenge #3: Latency

All of these challenges involve agent **limitations**.  
... their *own* and *other's*.

# Talk Outline

---

- Robots
  - A two robot, adversarial, concurrent learning problem.
  - The challenges for multiagent learning.
- Limitations and Equilibria
- Limitations and Learning



# Limitations Restrict Behavior

---

- **Restricted Policy Space** —  $\bar{\Pi}_i \subseteq \Pi_i$

Any subset of stochastic policies,  $\pi : \mathcal{S} \rightarrow PD(\mathcal{A}_i)$ .

- **Restricted Best-Response** —  $\overline{BR}_i(\pi_{-i})$

The set of all policies from  $\bar{\Pi}_i$  that are optimal given the policies of the other players.

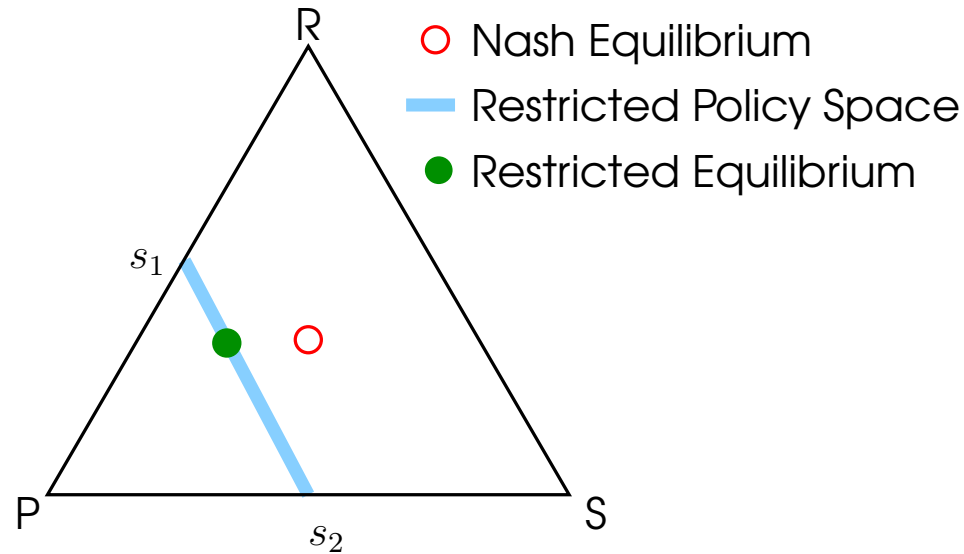
- **Restricted Equilibrium** —  $\pi_{i=1\dots n}$

$$\pi_i \in \overline{BR}_i(\pi_{-i})$$

A strategy for each player, where no player *can* and *wants* to deviate given the other players continue to play the equilibrium.

Do Restricted Equilibria Exist?

# Do Restricted Equilibria Exist? = 1

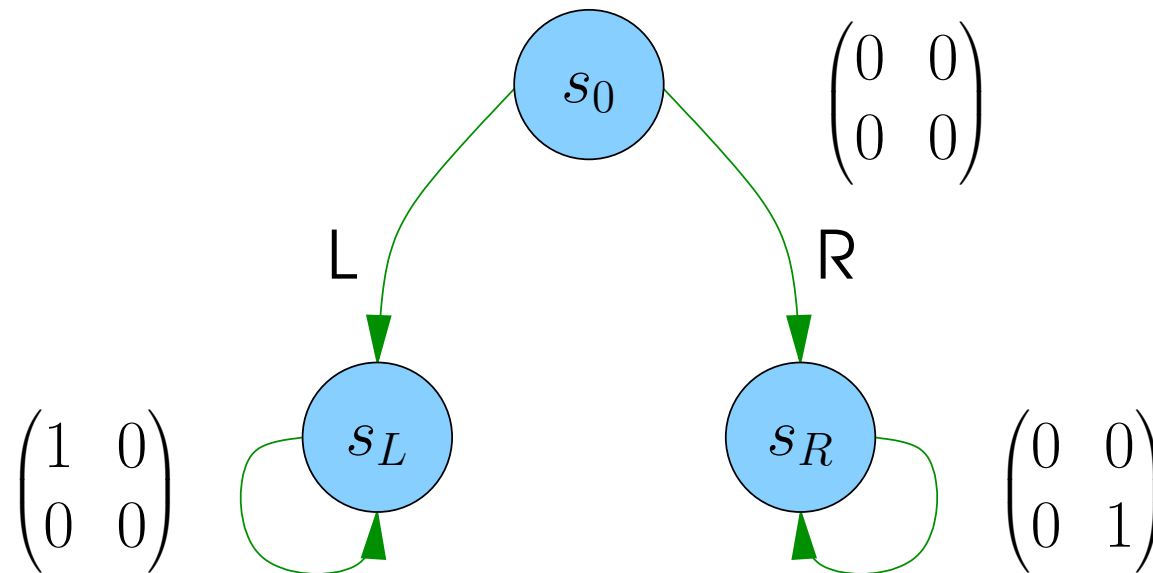


	Explicit Game	Implicit Game
Payoffs	$\begin{pmatrix} 0 & -1 & 1 \\ 1 & 0 & -1 \\ -1 & 1 & 0 \end{pmatrix}$	$\begin{pmatrix} -\frac{1}{2} & 0 \\ \frac{1}{2} & -\frac{1}{2} \\ 0 & \frac{1}{2} \end{pmatrix}$
Equilibrium	$\langle \frac{1}{3}, \frac{1}{3}, \frac{1}{3} \rangle, \langle \frac{1}{3}, \frac{1}{3}, \frac{1}{3} \rangle$	$\langle 0, \frac{1}{3}, \frac{2}{3} \rangle, \langle \frac{2}{3}, \frac{1}{3} \rangle$
Restricted Equilibrium	$\langle 0, \frac{1}{3}, \frac{2}{3} \rangle, \langle \frac{1}{3}, \frac{1}{2}, \frac{1}{6} \rangle$	

# Do Restricted Equilibria Exist? = 2

---

- Two-player, zero-sum stochastic game (Marty's Game 2).<sup>1</sup>



- Players restricted to policies that play the same distribution over actions in all states.

This game has no restricted equilibria!

---

<sup>1</sup>This counterexample is brought to you by Martin Zinkevich.



# Do Restricted Equilibria Exist? = 3

---

- In matrix games, if  $\bar{\Pi}_i$  is convex, then ...
- If  $\bar{\Pi}_i$  is statewise convex, then ...
- In no-control stochastic games, if convex  $\bar{\Pi}_i$ , then ...
- In single-controller stochastic games, if  $\bar{\Pi}_1$  is statewise convex, and  $\bar{\Pi}_{i \neq 1}$  is convex, then ...
- In team games ...

# Do Restricted Equilibria Exist? = 3

---

- In matrix games, if  $\bar{\Pi}_i$  is convex, then ...
- If  $\bar{\Pi}_i$  is statewise convex, then ...
- In no-control stochastic games, if convex  $\bar{\Pi}_i$ , then ...
- In single-controller stochastic games, if  $\bar{\Pi}_1$  is statewise convex, and  $\bar{\Pi}_{i \neq 1}$  is convex, then ...
- In team games ...

... there exists a restricted equilibrium.

**Proofs.** Uses Kakutani's fixed point theorem after showing

$$\forall \pi_{-i} \quad \overline{\text{BR}}_i(\pi_{-i}) \text{ is convex.}$$

# The Challenges = 2

---

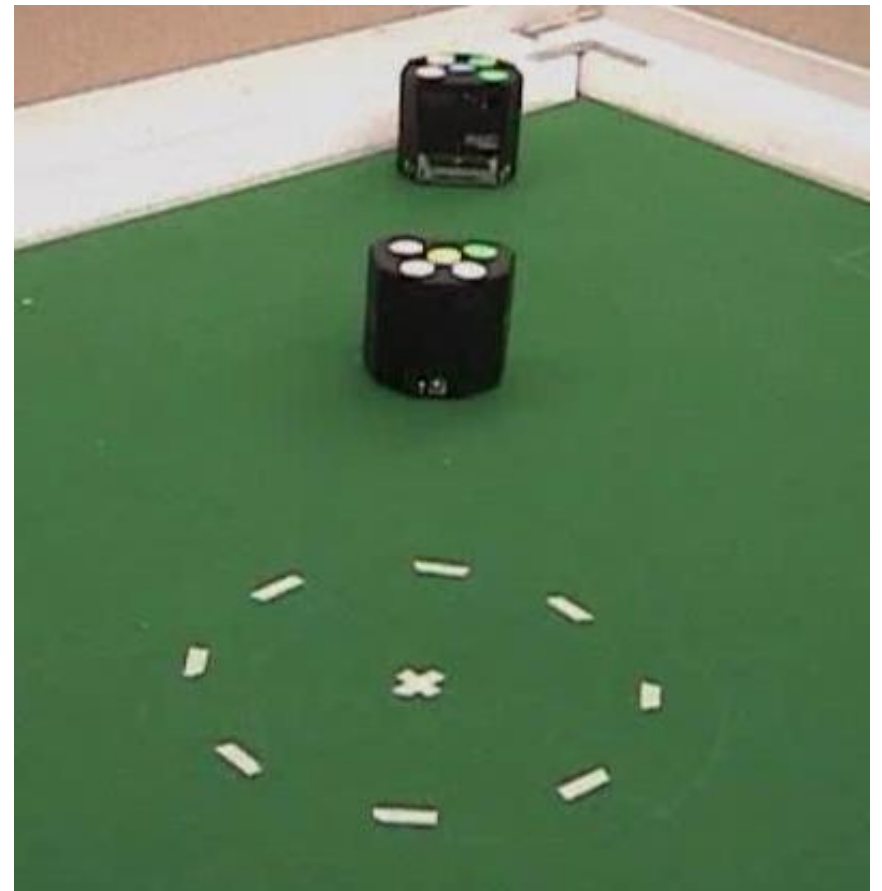
- Challenge #1: Continuous State and Action Spaces
- Challenge #2: Fixed Behavioral Components
- Challenge #3: Latency

None of these are nice enough to guarantee the existence of equilibria.

# Talk Outline

---

- Robots
  - A two robot, adversarial, concurrent learning problem.
  - The challenges for multiagent learning.
- Limitations and Equilibria
- Limitations and Learning



# Three Ideas = One Algorithm

---

- Idea #1: Policy Gradient Ascent
- Idea #2: WoLF Variable Learning Rate

GrāWoLF— Gradient-based WoLF

- Idea #3: Tile Coding

# Idea #1

---

- Policy Gradient Ascent (Sutton et al., 2000)
  - Policy improvement with parameterized policies.
  - Takes steps in direction of the gradient of the value.

$$\pi(s, a) = \frac{e^{\phi_{sa} \cdot \theta_k}}{\sum_{b \in \mathcal{A}_i} e^{\phi_{sb} \cdot \theta_k}}$$

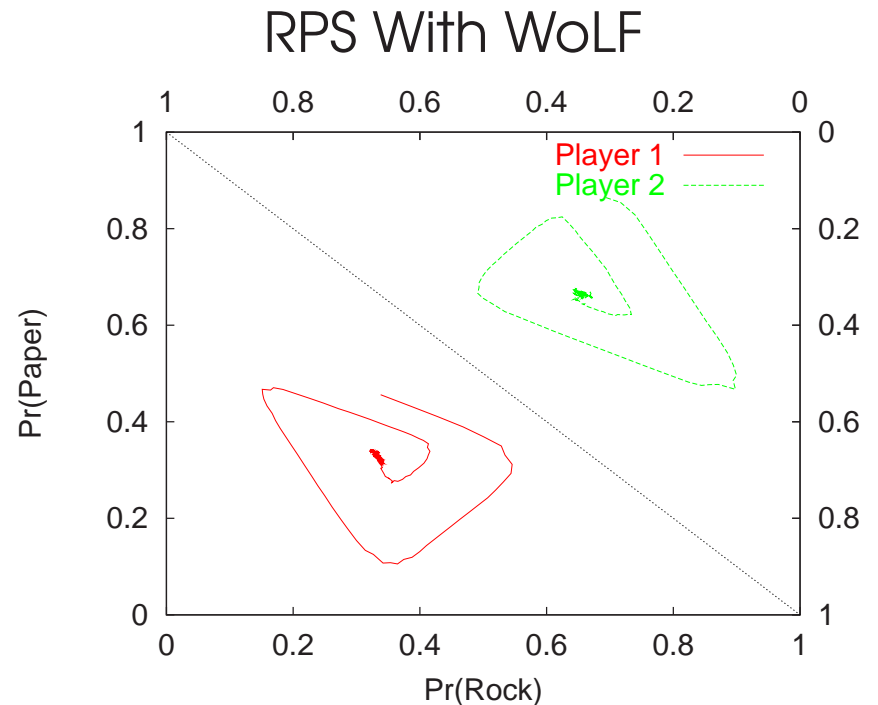
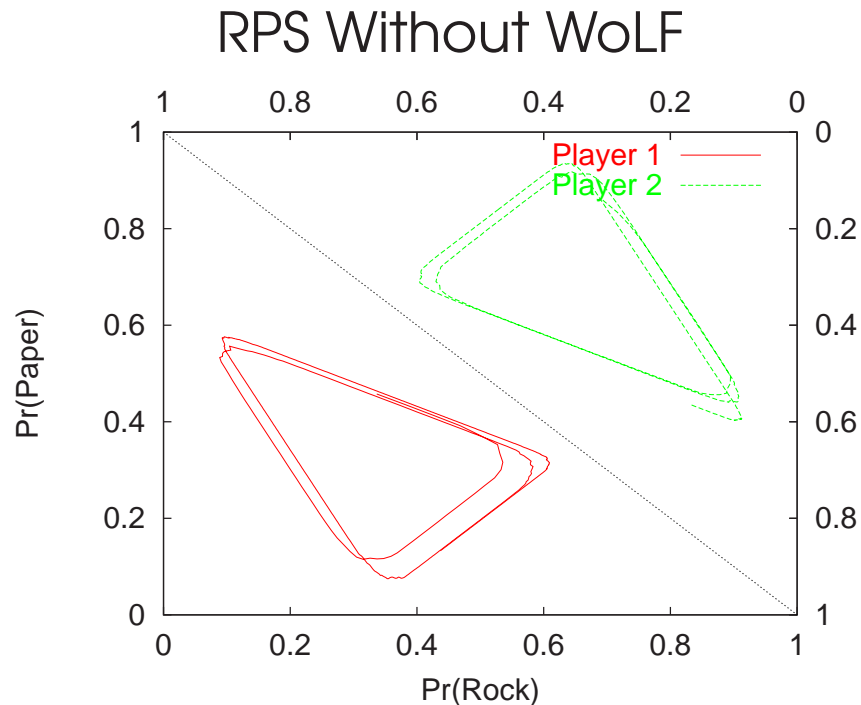
$$\theta_{k+1} = \theta_k + \alpha_k \sum_a \phi_{sa} \pi(s, a) f_k(s, a)$$

- $f_k$  is an approximation of the advantage function.

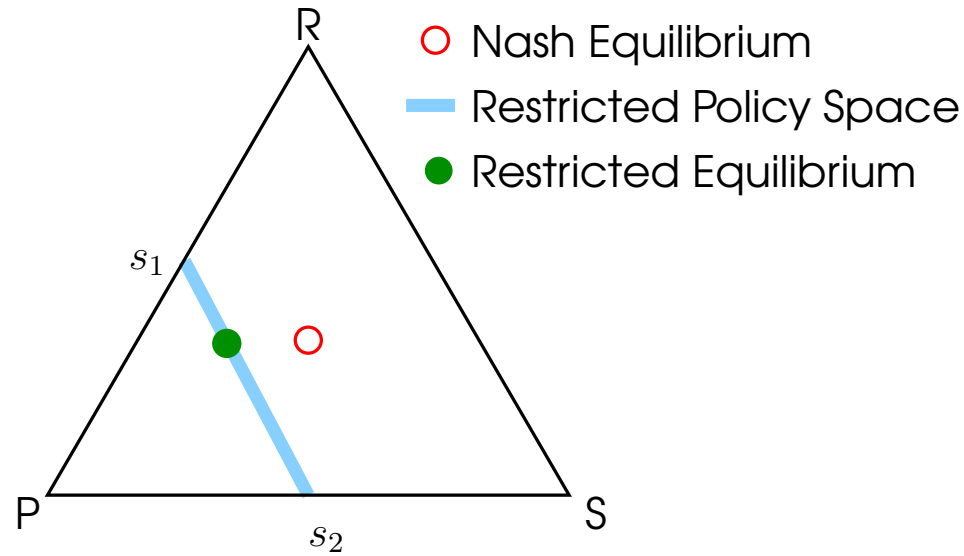
$$\begin{aligned} f_k(s, a) &\approx Q(s, a) - V^\pi(s) \\ &\approx Q(s, a) - \sum_b \pi(s, b) Q(s, b) \end{aligned}$$

# Idea #2

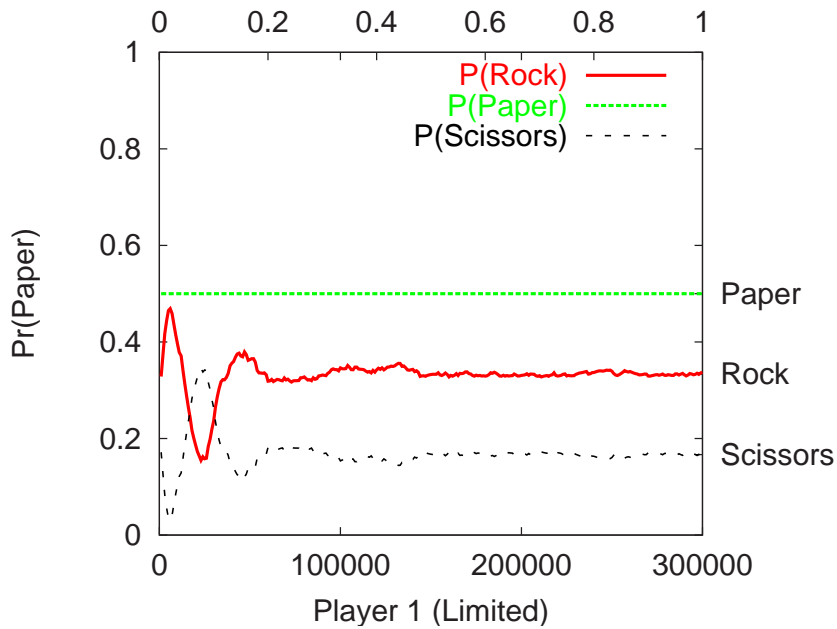
- Win or Learn Fast (WoLF) (Bowling & Veloso, 2002)
  - Variable learning rate accounts for other agents.
    - \* Learn fast when losing.
    - \* Cautious when winning, since agents may adapt.
  - Theoretical and empirical evidence of convergence.



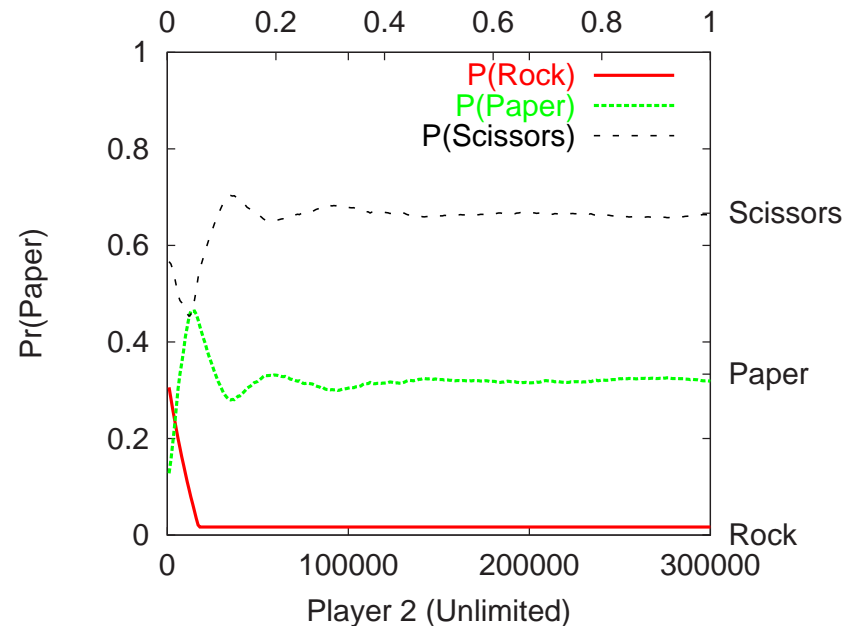
# Idea #2 = 2



### RPS Without WoLF



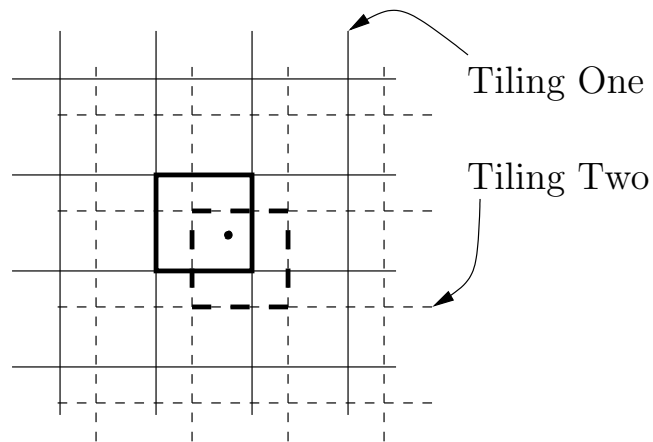
### RPS With WoLF





# Idea #3

- Tile Coding (a.k.a. CMACS) (Sutton & Barto 1998)
  - Space covered by overlapping and offset tilings.
  - Maps continuous (or discrete) spaces to a vector of boolean values.
  - Provides discretization and generalization.



# The Task

---

# *The Task = Goofspiel*

---

- A.k.a. “The Game of Pure Strategy”

# *The Task = Goofspiel*

---

- A.k.a. “The Game of Pure Strategy”
- Each player plays a full suit of cards.
- Each player uses their cards (without replacement) to bid on cards from another suit.

# The Task = Goofspiel

---

- A.k.a. “The Game of Pure Strategy”
- Each player plays a full suit of cards.
- Each player uses their cards (without replacement) to bid on cards from another suit.

$n$	$ S $	$ S \times A $	SIZEOF( $\pi$ or $Q$ )	VALUE(det)	VALUE(random)
4	692	15150	$\sim 59\text{KB}$	-2	-2.5
8	$3 \times 10^6$	$1 \times 10^7$	$\sim 47\text{MB}$	-20	-10.5
13	$1 \times 10^{11}$	$7 \times 10^{11}$	$\sim 2.5\text{TB}$	-65	-28

- The game is very large.
- Deterministic policies are very bad.
- The random policy isn't too bad.

# The Task = Goofspiel = 2

My Hand	1	3	4	5	6	8	11	13
Quartiles	*		*		*	*		*
Opp Hand	4	5	8	9	10	11	12	13
Quartiles	*		*		*	*		*
Deck	1	2	3	5	9	10	11	12
Quartiles	*		*		*	*		*
Card	11							
Action	3							

$\langle 1, 4, 6, 8, 13 \rangle,$   
 $\langle 4, 8, 10, 11, 13 \rangle,$   
 $\langle 1, 3, 9, 10, 12 \rangle,$   
 11, 3

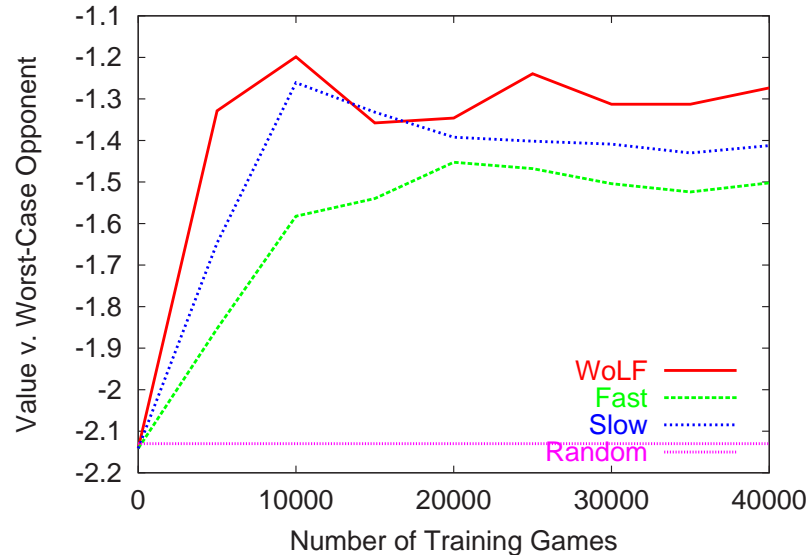
$\Downarrow$  (Tile Coding)

TILES  $\in \{0, 1\}^{10^6}$

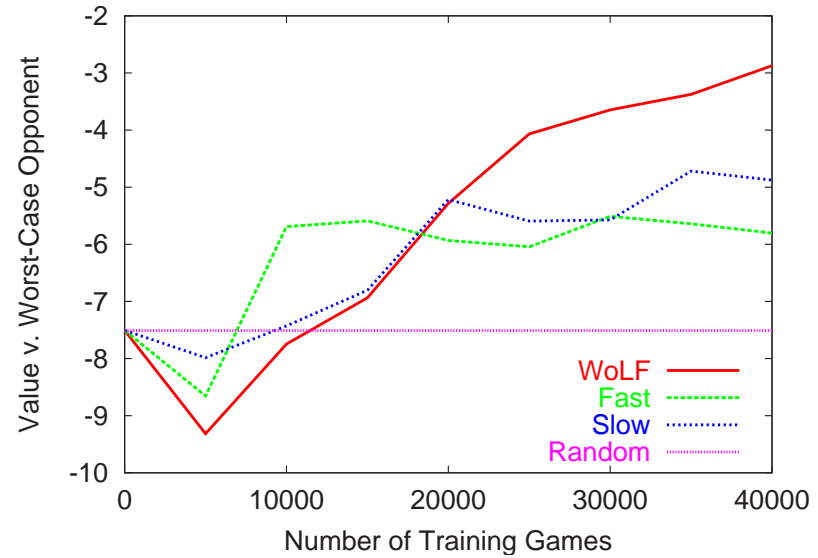
- Gradient ascent on this parameterization.
- WoLF variable learning rate on the gradient step size.

# Results = Worst-Case

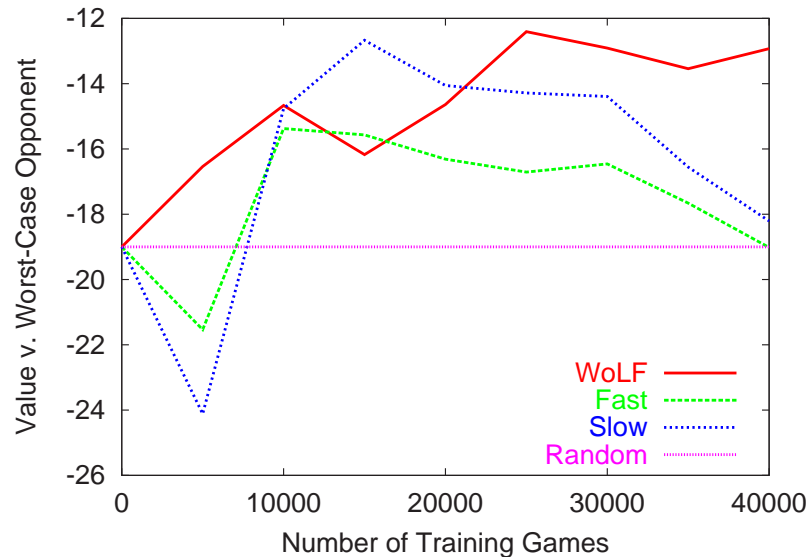
## 4 Cards



## 8 Cards



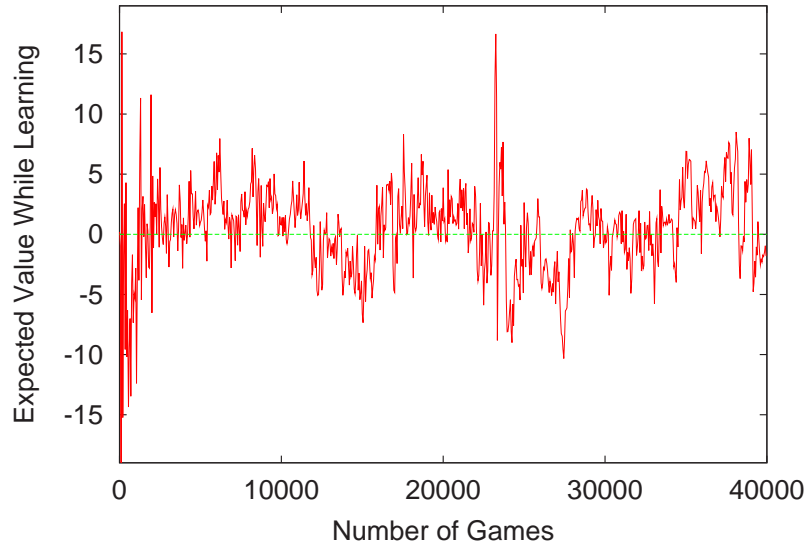
## 13 Cards



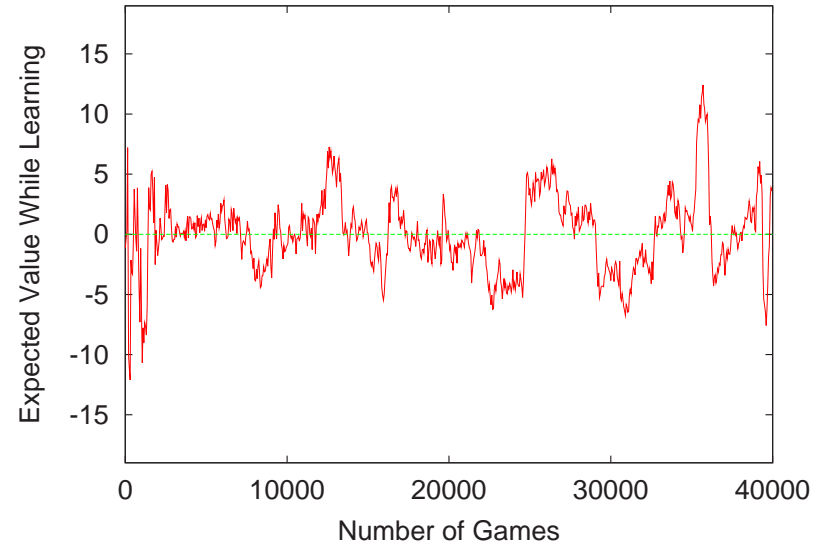
# Results = While Learning

---

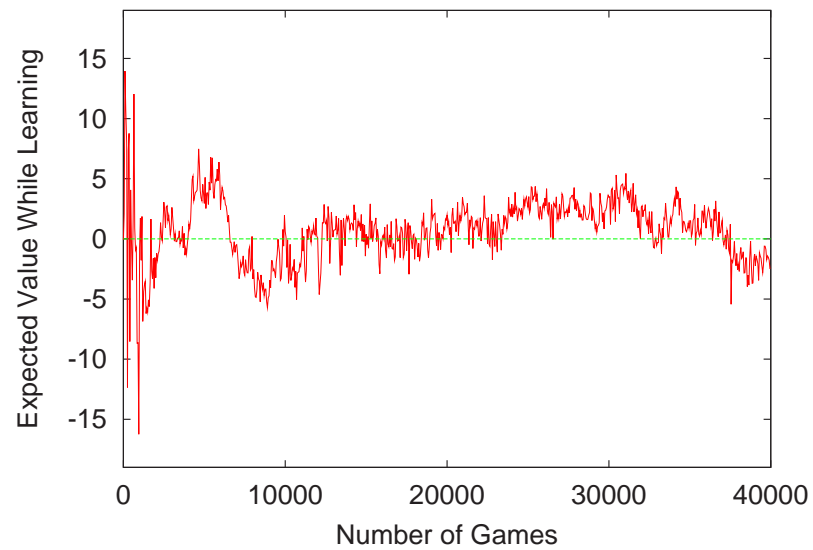
## Fast



## Slow



## WoLF





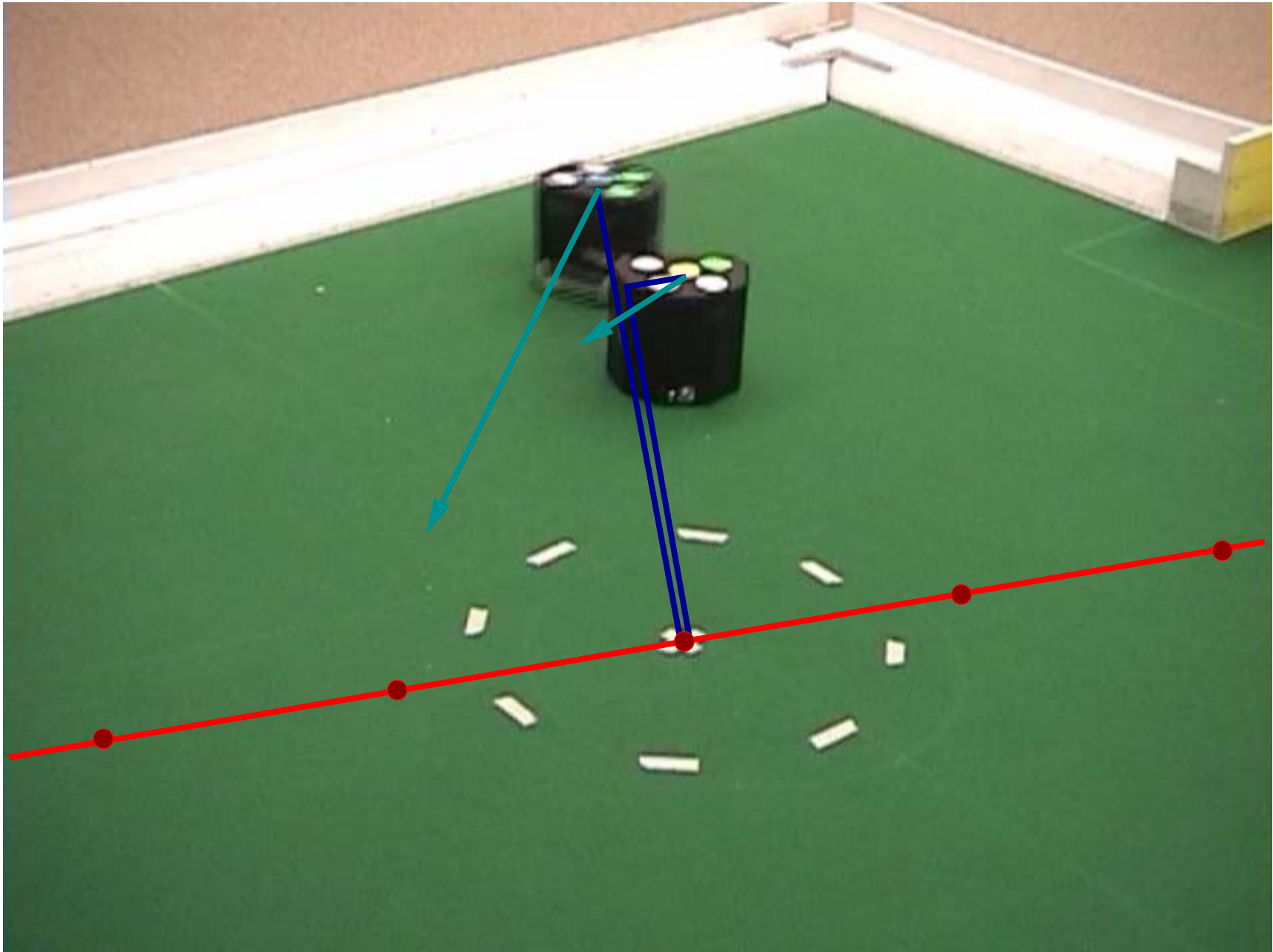
# The Task = Breakthrough

---



# The Task = Breakthrough = 2

---



# Results = Breakthrough

---

WARNING!

# Results = Breakthrough

---

## WARNING!

- These results are preliminary.... some are only hours old.
- They involve a single run of learning in a highly stochastic learning environment.
- More experiments in progress.

# Results = "To the videotape..."

---

Playback of learned policies in simulation and on the robots.

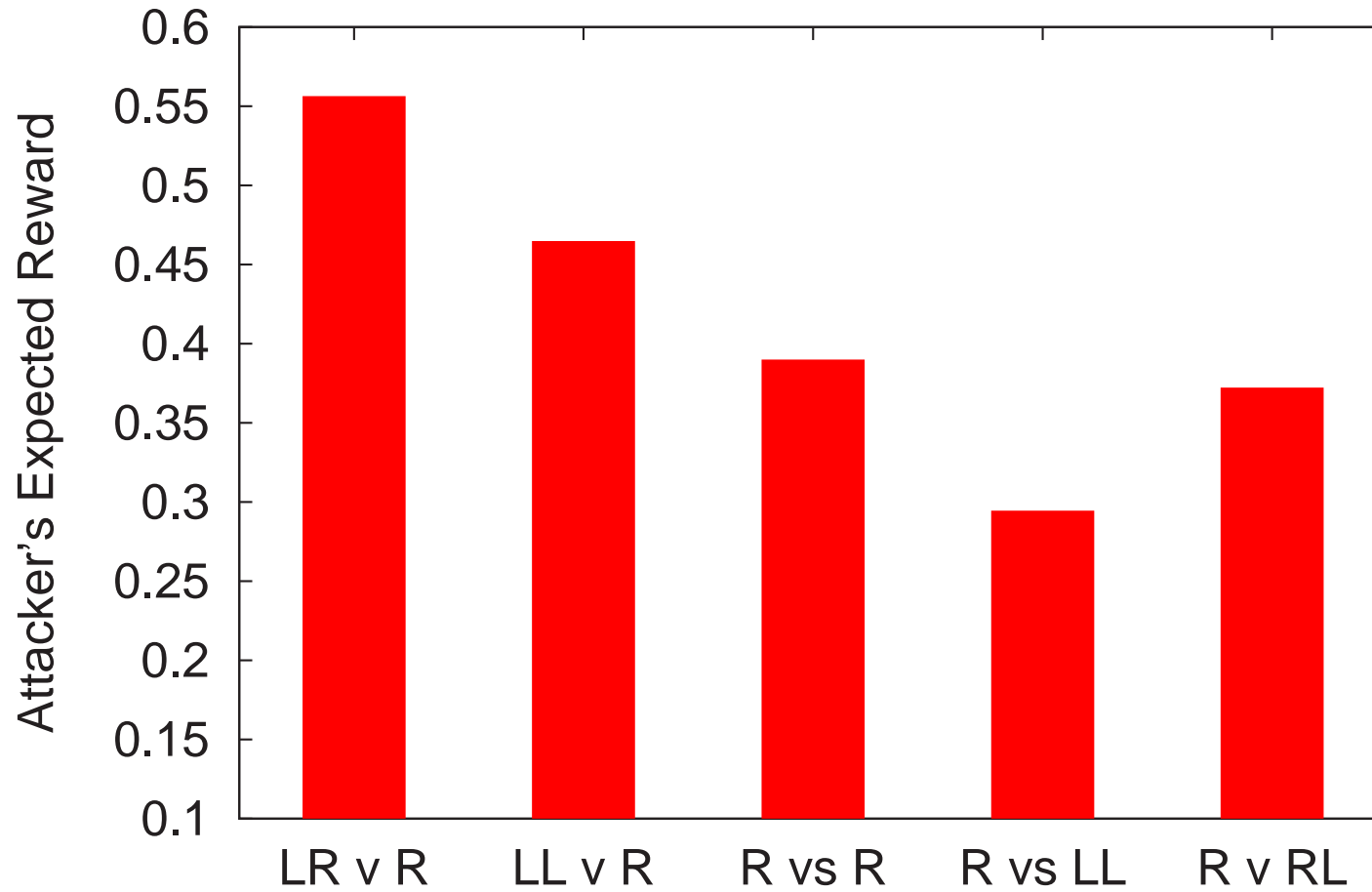
The robot video can be downloaded from...

<http://www.cs.cmu.edu/~mhb/research/>

# Results = 3

---

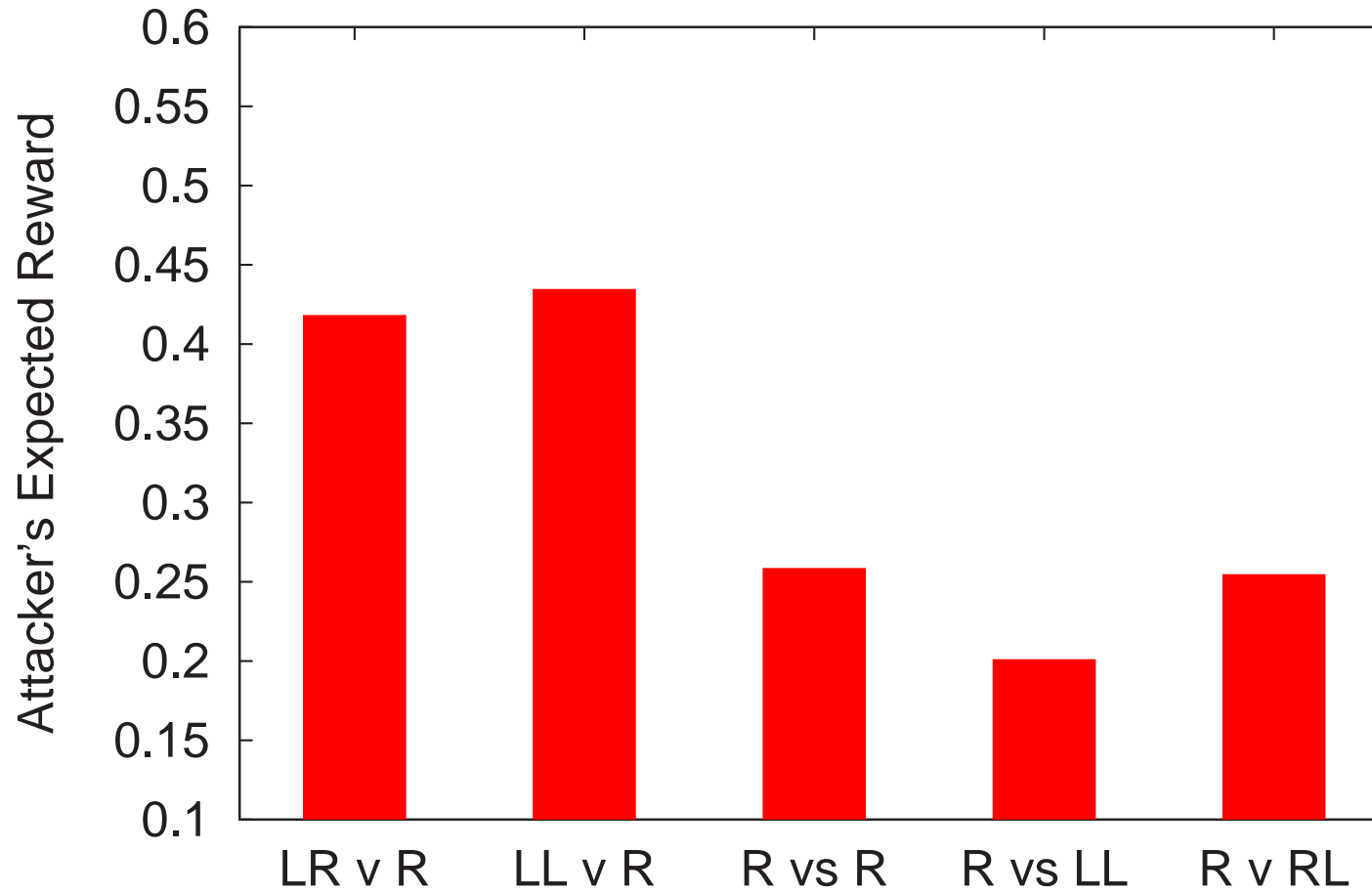
Omni vs Omni: Learned Policies



# Results = 4

---

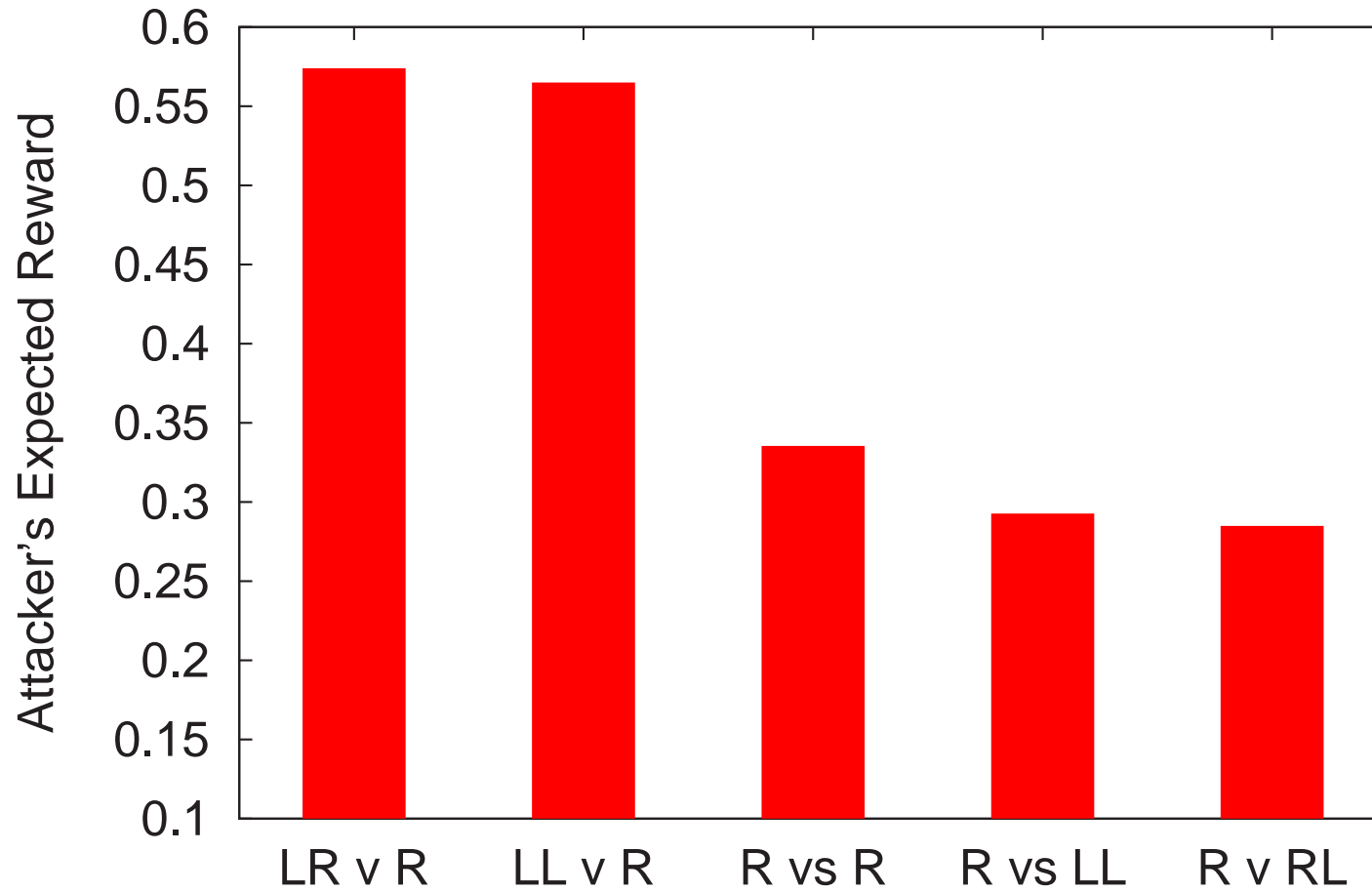
Diff vs Omni: Learned Policies



# Results = 5

---

Diff vs Diff: Learned Policies

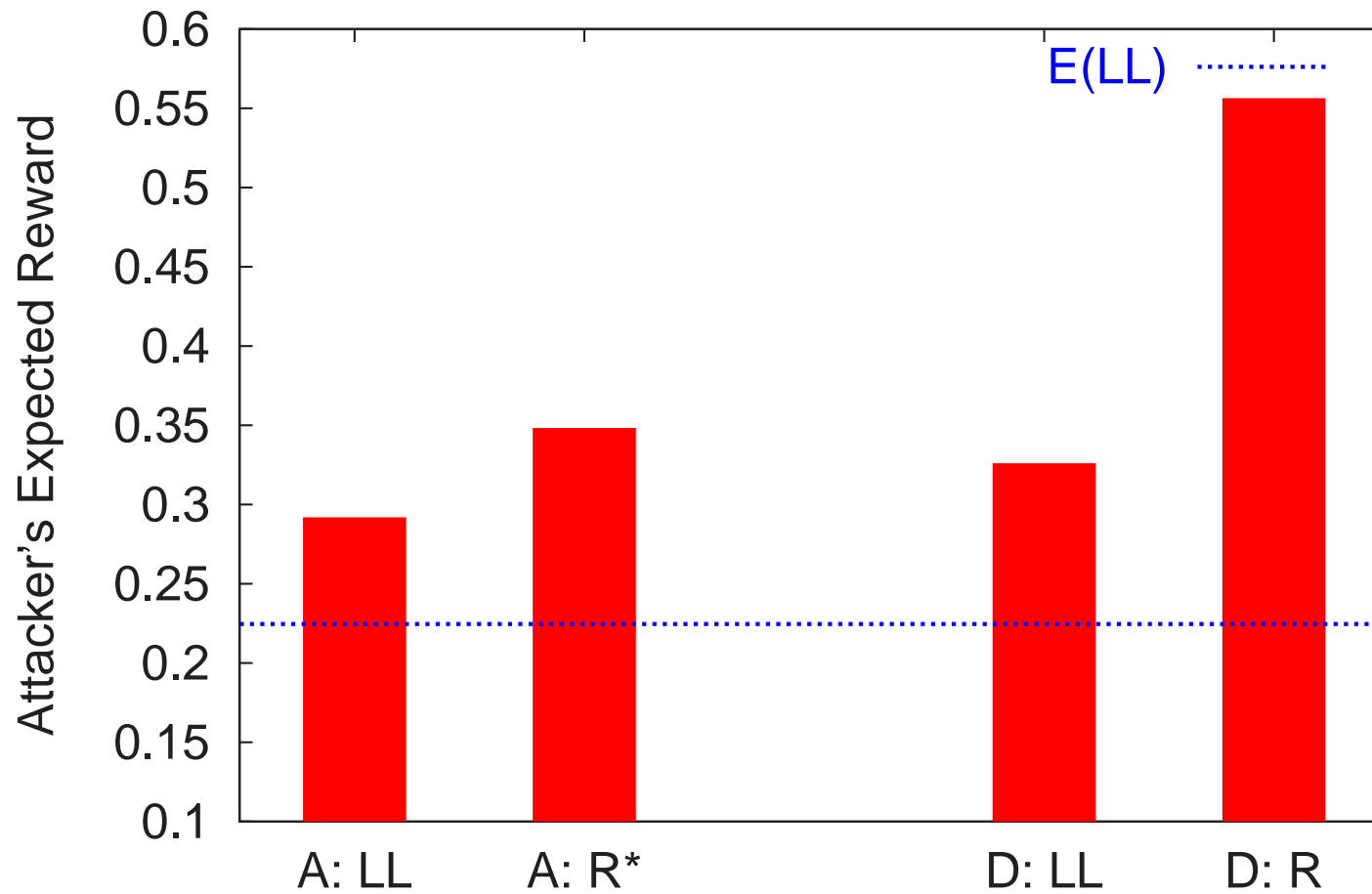




# Results = 6

---

Omni vs Omni: Worst-Case



# Results = Breakthrough

---

## WARNING!

- These results are preliminary.... some are only hours old.
- They involve a single run of learning in a highly stochastic learning environment.
- More experiments in progress.

# Big Picture

---

- How do we scale our (collective) algorithms to large problems with limited agents?

# Big Picture

---

- How do we scale our (collective) algorithms to large problems with limited agents?
  - Equilibrium learning may be in trouble.
  - Lagoudakis and Parr's approximation and minimax. (NIPS '02)
  - Correlated equilibria?

# Big Picture

---

- How do we scale our (collective) algorithms to large problems with limited agents?
  - Equilibrium learning may be in trouble.
  - Lagoudakis and Parr's approximation and minimax. (NIPS '02)
  - Correlated equilibria?
- What is the objective?

# Big Picture

---

- How do we scale our (collective) algorithms to large problems with limited agents?
  - Equilibrium learning may be in trouble.
  - Lagoudakis and Parr's approximation and minimax. (NIPS '02)
  - Correlated equilibria?
- What is the objective?
  - Performance during learning.
  - Generality of learned policies.
    - \* How can I be exploited?
    - \* What if everyone played this policy?