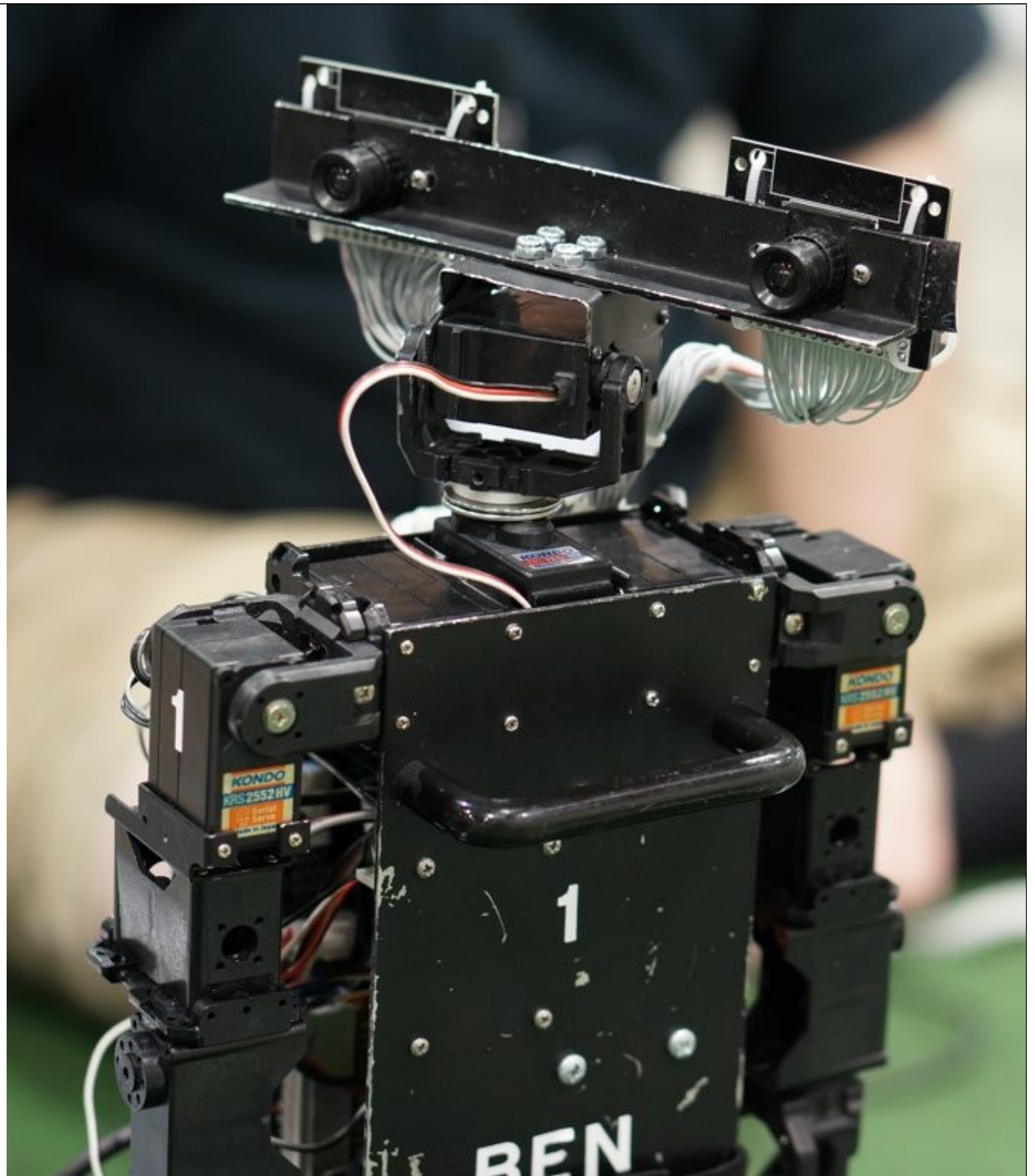


Online Training of Robots and Multirobot Teams

Sean Luke

陸子荀

Department of
Computer Science
George Mason University



About Me

- Associate Professor
Department of Computer Science
George Mason University
- **Interests**
Multiagent Systems
Machine Learning
Multirobotics
Stochastic Optimization and Evolutionary Computation
Simulation
- **Software (and Hardware)**
ECJ *Evolutionary Computation Toolkit*
MASON *Multiagent Simulation Toolkit*
RoboPatriots and FlockBots *Robot Architectures*



My Current Multiagent Systems Problem

Topics in This Talk

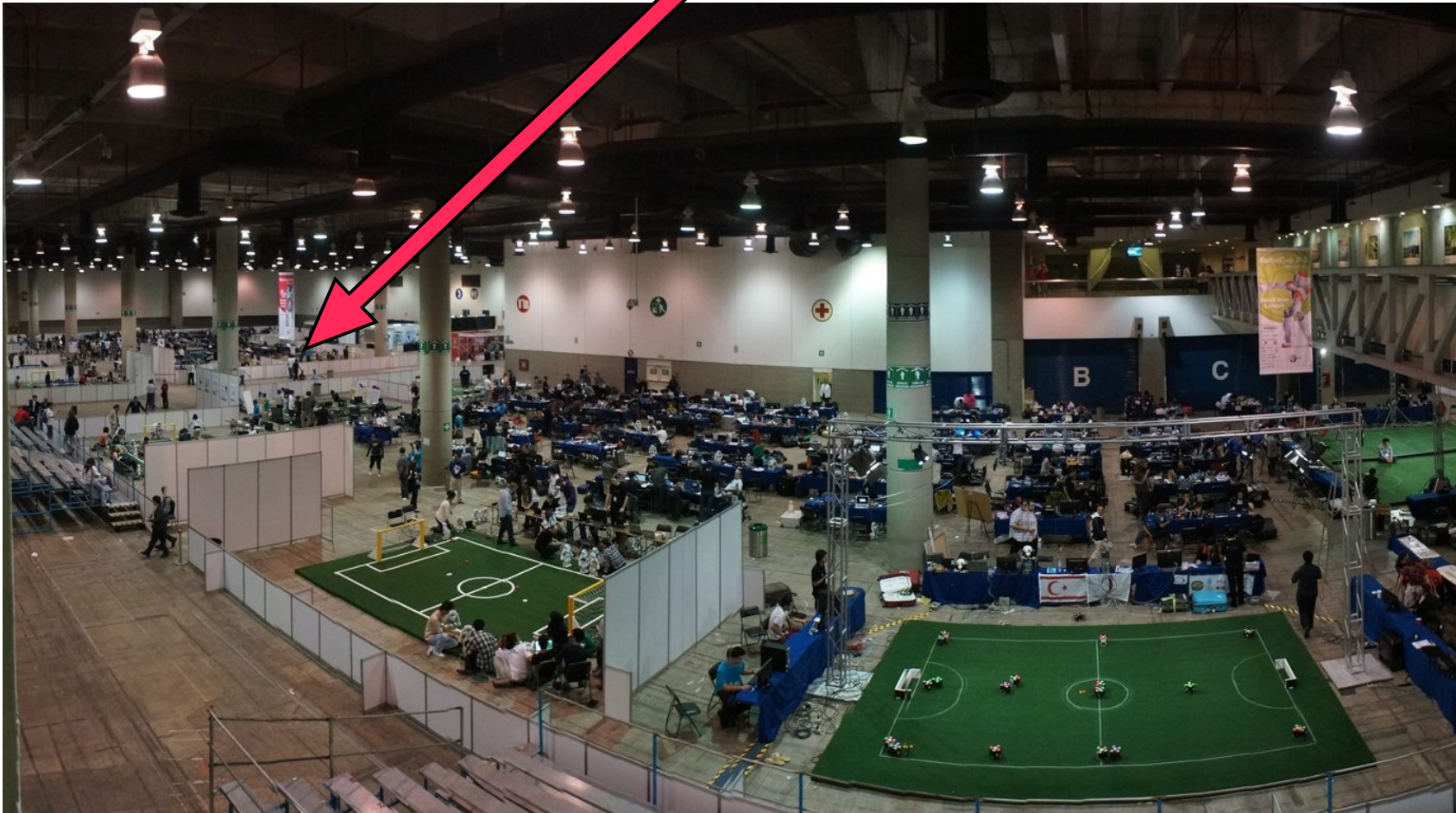
- **RoboCup**
- **Multiagent and Multi-robot Systems**
- **Pheromone-based Robotics:** An Example of Emergent Behavior
- **HiTAB:** Single-Agent and Single-Robot Training
- **Unlearning:** Dealing with noise in single-agent training
- **Behavioral Bootstrapping:** training a flat (leaderless) swarm
- **M-HiTAB:** Hierarchical Multiagent and Multi-Robot Training

RoboCup 2012 Mexico City



RoboCup 2012

George Mason University

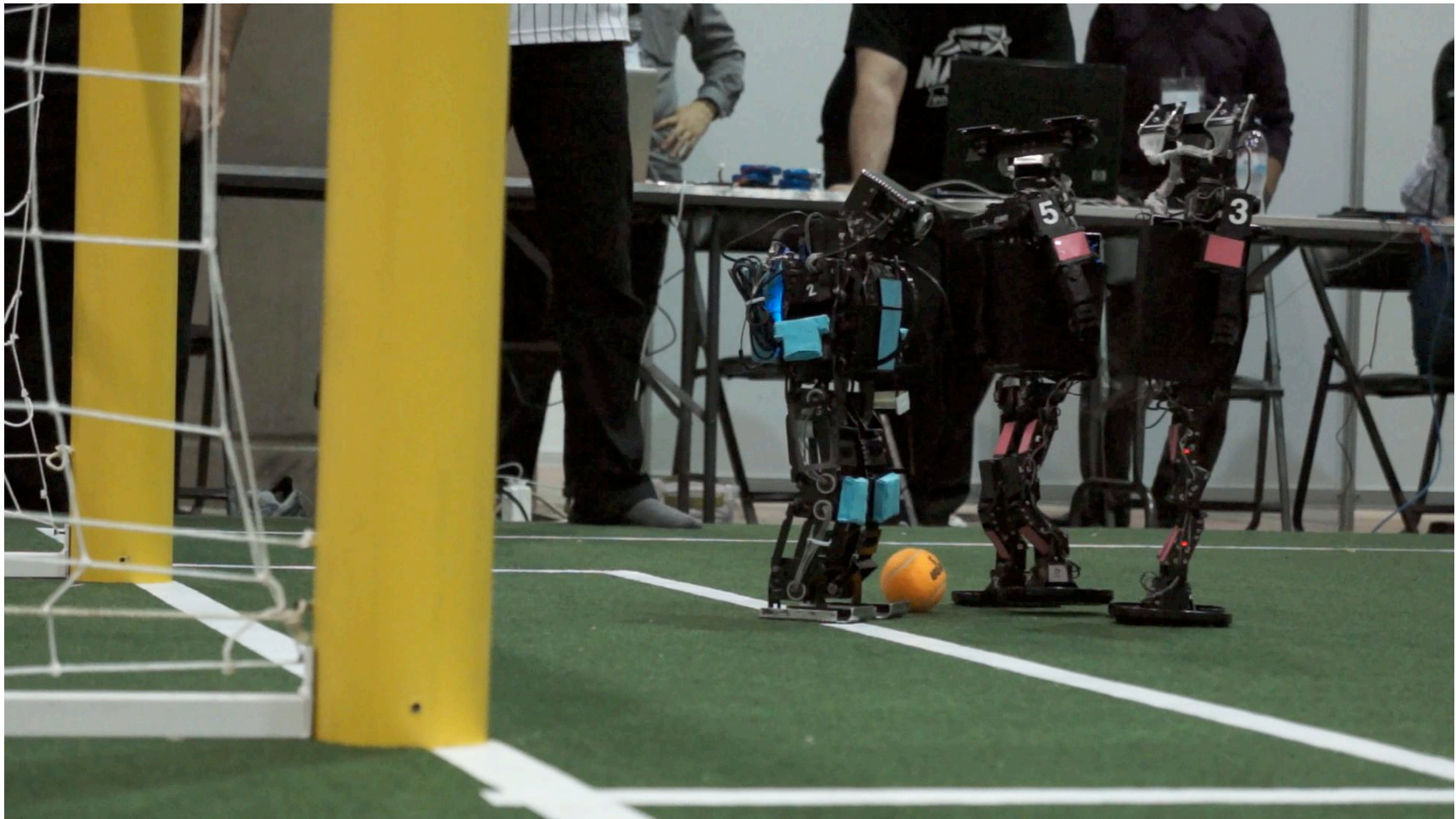




RoboCup 2012

GMU: Pink

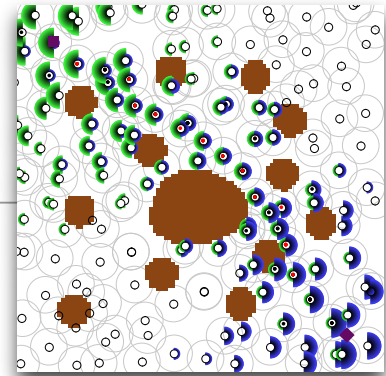
Osaka: Blue



A Multiagent System (or MAS)

- **Agent:** an autonomous entity which iteratively manipulates its environment in response to feedback received from the environment.
- **Multiagent System:** a system of ... you know ... multiple agents.
 - Agent interaction
 - Emergence
- **Distributed Systems Problem:** given multiple processors and resources under your control, solve a given task.
- **Multiagent Systems Problem:** given multiple agents with **major constraints on communication or mutual knowledge**, solve a given task.

Why Develop / Simulate MAS?



- **Science:** MAS **models** can help us make predictions and test hypotheses when it would be impossible, immoral, or unrealistic to perform real-world tests.
 - Biology, Physics, **Social Sciences**
 - *Goal:* accurate replication of existing phenomena
- **Engineering:** MAS **methods** help us test new techniques or inventions.
 - Games, Animation, Networked Agents, **Multirobotics**
 - *Goal:* optimization or demonstration of new methods

Multiagent Systems (for Engineering)

- **Agent or Robot Teams**

Small Numbers,
Often Heterogeneous
Lots of Communication/Interaction
Global Communication



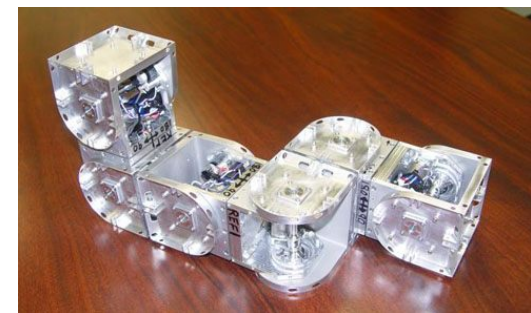
- **Agent or Robot Swarms**

Large Numbers

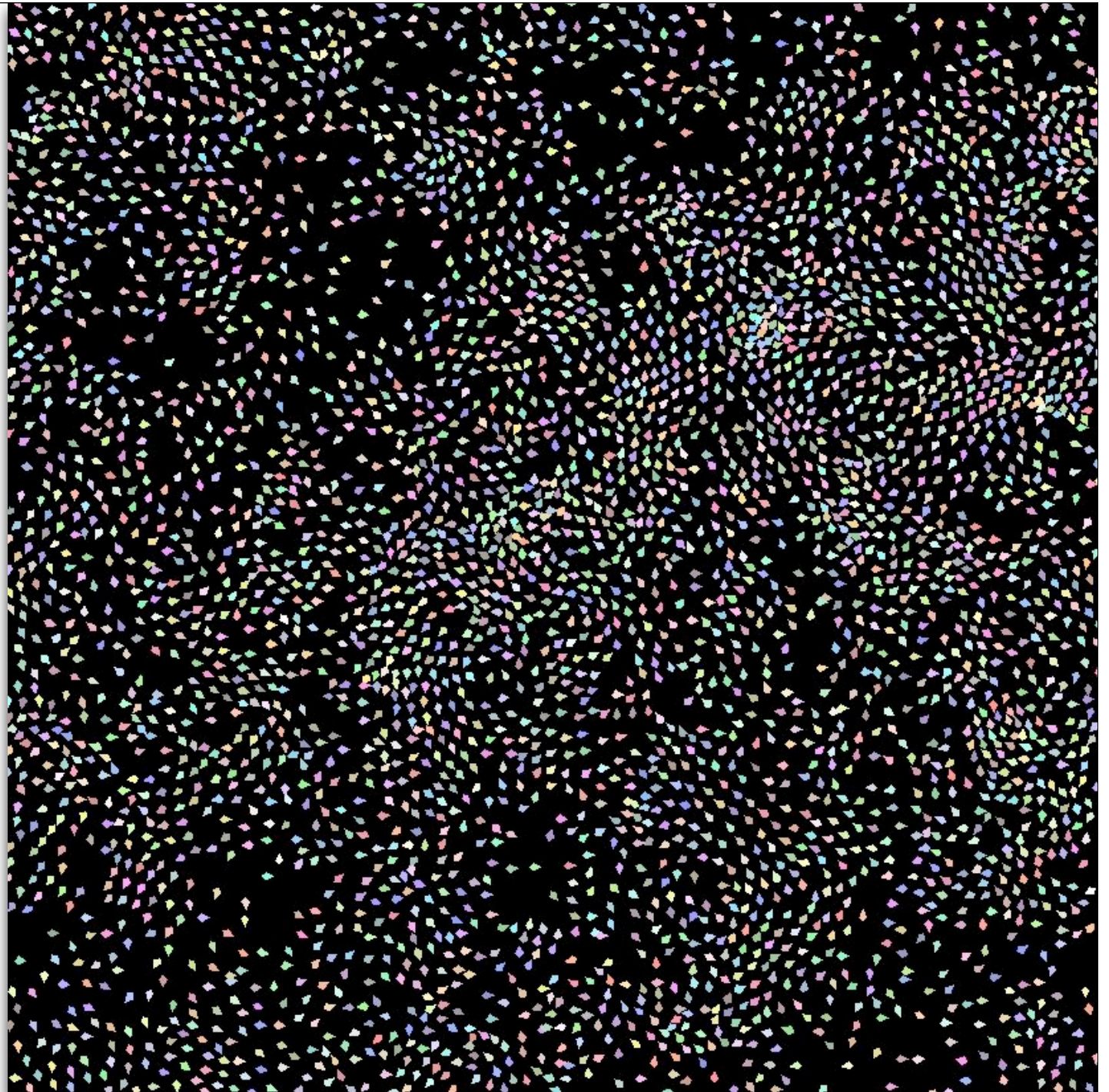


- **Modular Robots**

A Robot Consists of *Modules* (the “Agents”)
Moderate Numbers, Usually Homogeneous
Communication via Internal Network
Is this really a multiagent system?



Multiagent
Systems
Are Very
Complex



The Multiagent Systems Design Space is **Big**

- **Factors in the complexity of a Multiagent Systems Design:**
 - Number of Agents
 - Complexity of Agent Behavior and Capability
 - Heterogeneity of Agents
 - Degree of Agent Interaction
 - Communication Complexity
 - Designing Robust and Cost-Effective Designs
- **This becomes very complicated very quickly**

Tradeoffs (in Multirobotics)

- **Agent or Robot Teams**
Small Numbers (often 2 or 3!)

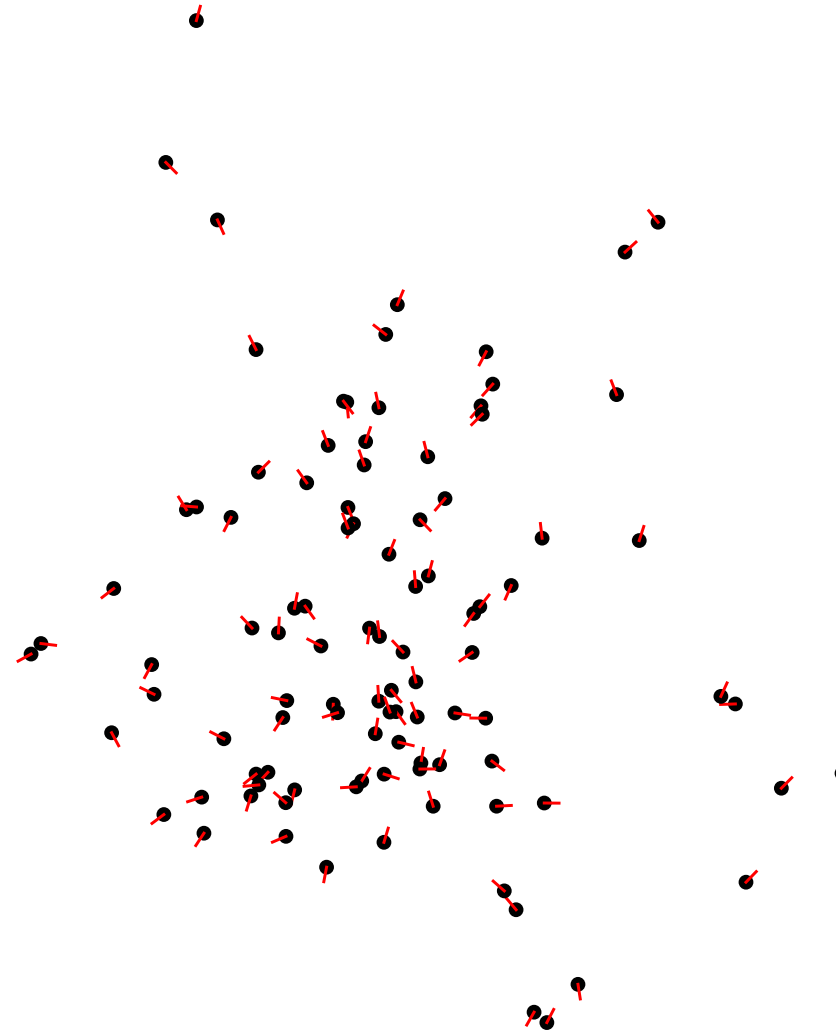


- **Agent or Robot Swarms**
Homogeneous
Little Communication/Interaction
Local Communication
Very Simple Behaviors
- **The more agents,
the simpler they get!**



Emergent Behavior

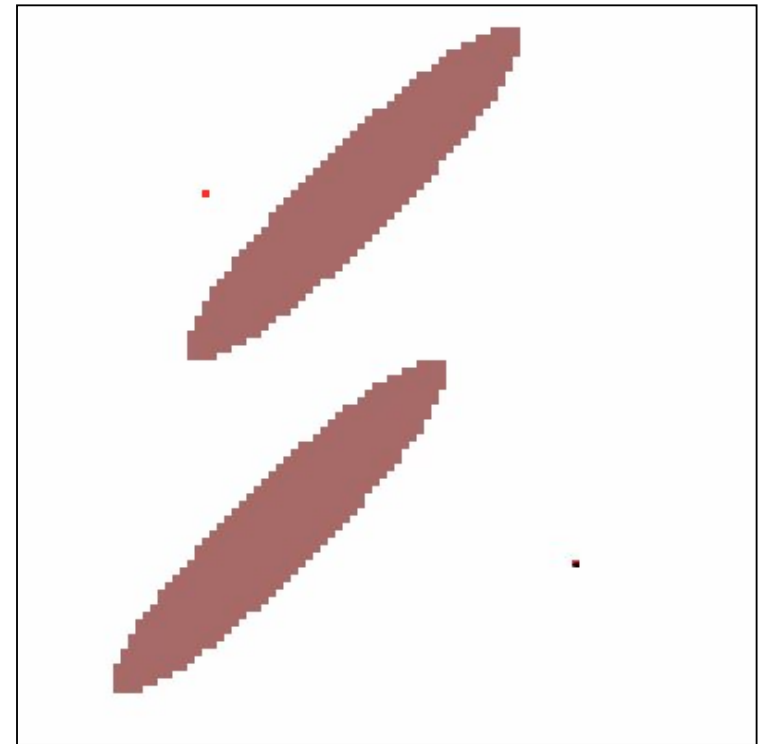
- **Simple Micro-Level Behaviors**
- **Complex Emergent Macrophenomena**
- **Can you Predict the Macrophenomena given the Micro-level Behaviors?**
- **Complexity Theorists Love Emergence**
- **Multiagent / Multirobot Designers Hate Emergence**



Can you predict this?

Example: Ant Pheromone Foraging

- **Most ant pheromone literature uses a single pheromone**
(Biologically plausible, but bad algorithms)
- **We use multiple pheromones**
2 in this example: *Food* and *Nest*
- **Each ant *follows* one pheromone but *updates* another.**
- **Each ant is in a *state*, which determines which pheromones it follows / updates.**



Example: Ant Pheromone Foraging

- | | | |
|-------------------------|--------------------------|--------------------------|
| • States: | Follow Pheromone: | Update Pheromone: |
| <i>Looking for Food</i> | Food | Nest |
| <i>Looking for Nest</i> | Nest | Food |

• **Following:**

An ant is in state s'

Go to square s'' with highest pheromone $U_p(s'')$

$$s' = \operatorname{argmax}_{s'' \in S''} U_p(s'')$$

• **Updating:**

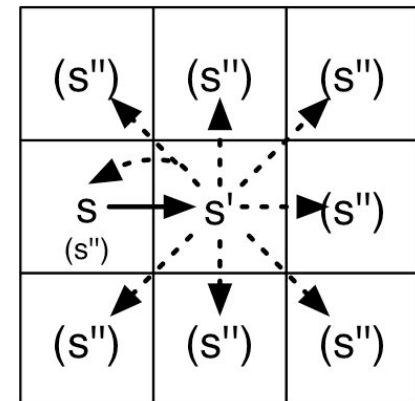
An ant is in state s'

Update $U_p(s')$

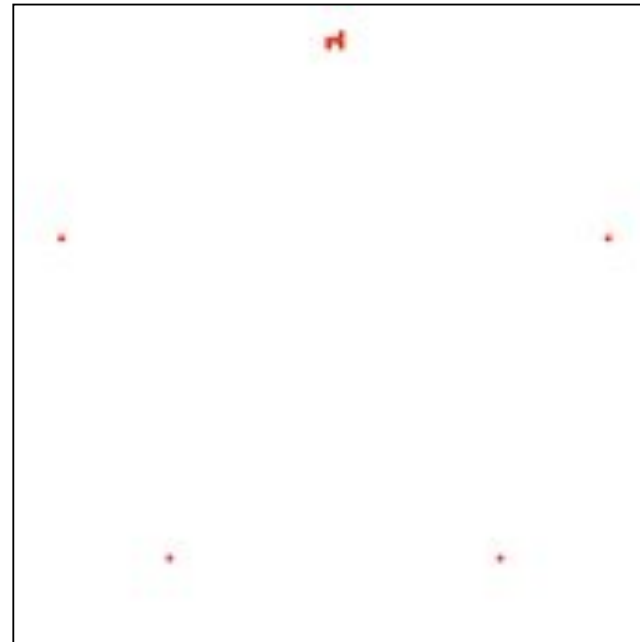
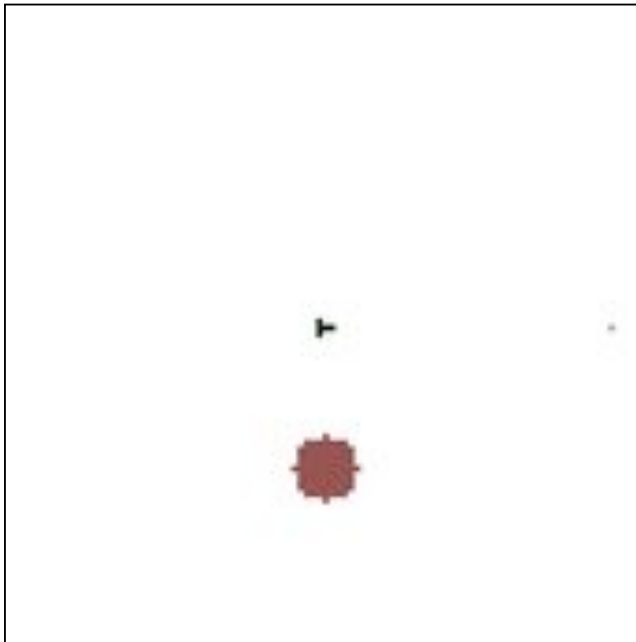
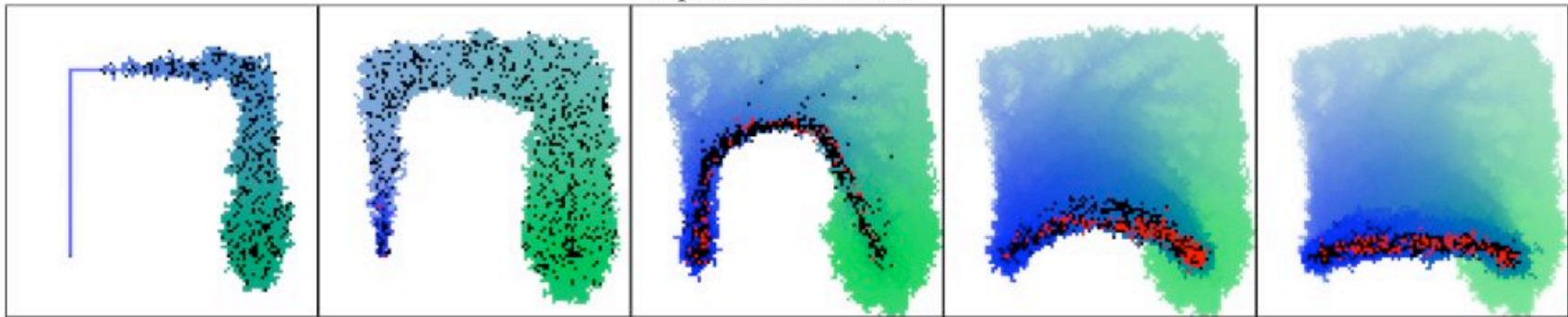
Reward $R(s')$ is received only if at nest / food

$$U_p(s') = R(s') + \gamma \max_{s'' \in S''} U_p(s'')$$

• **Form of multi-utility value iteration**

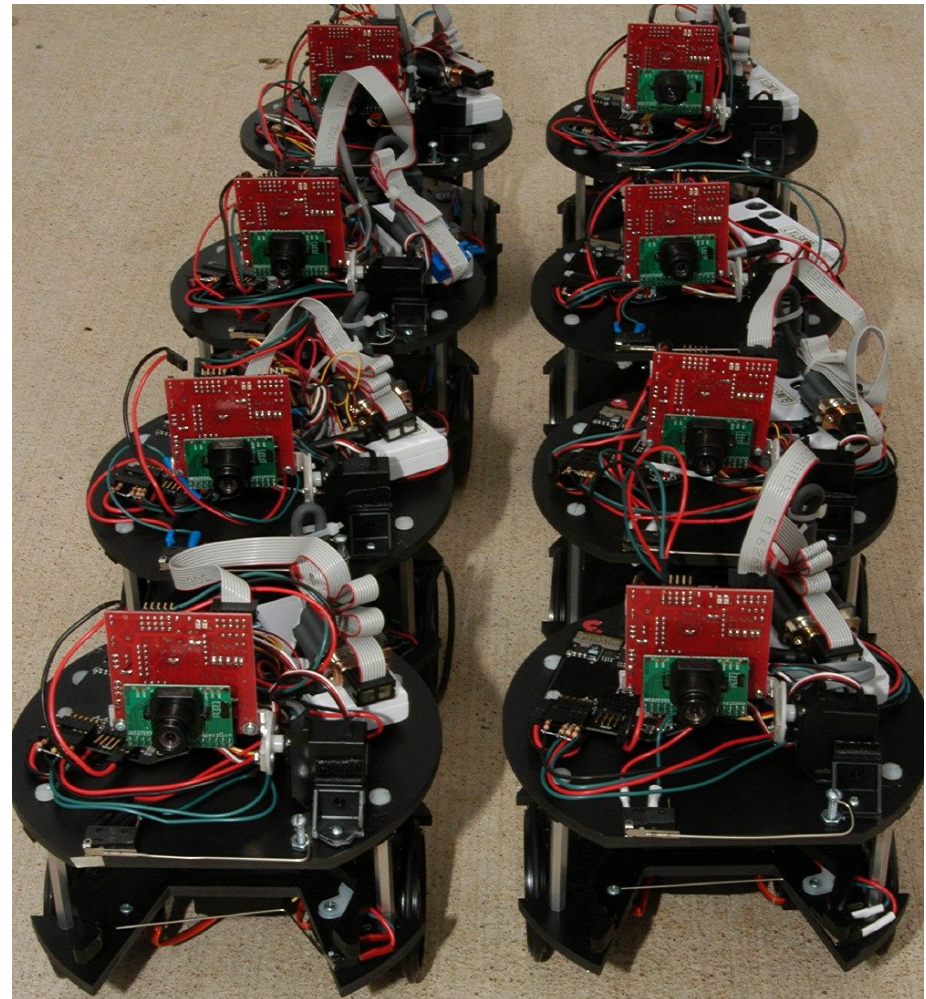


Example: Ant Pheromone Foraging

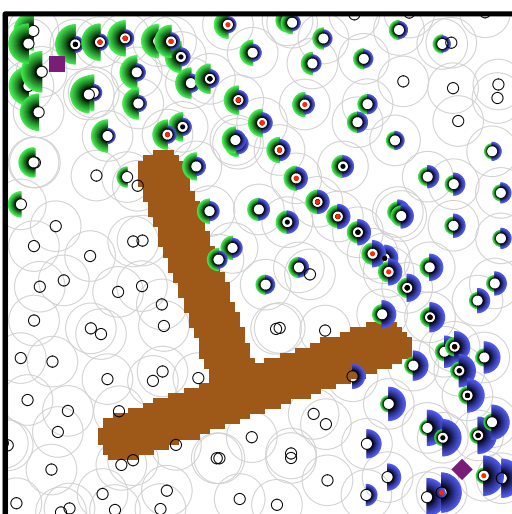
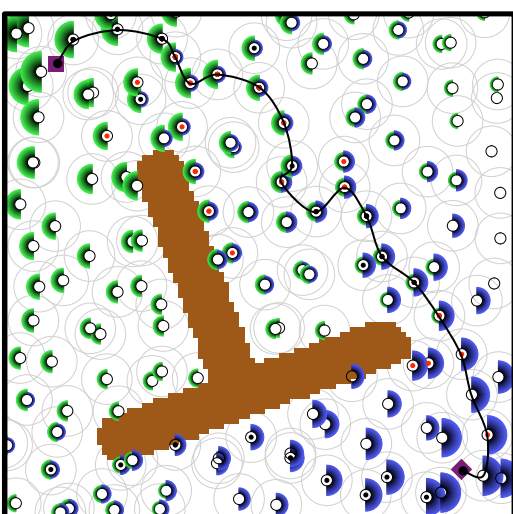
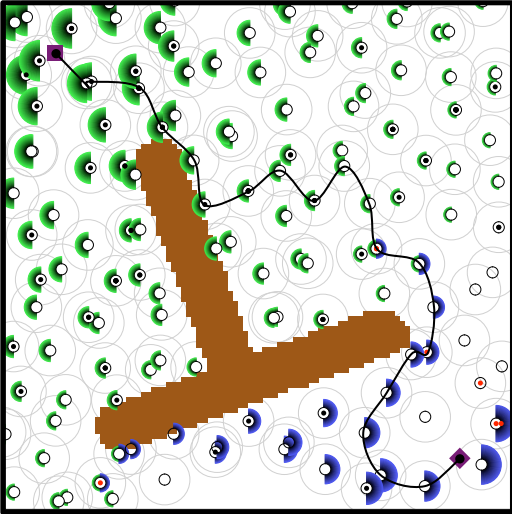
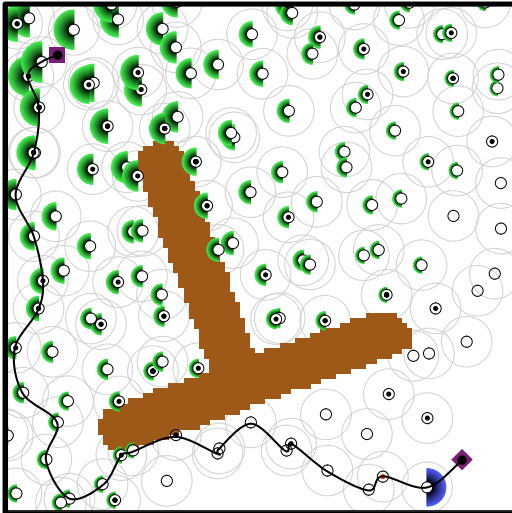
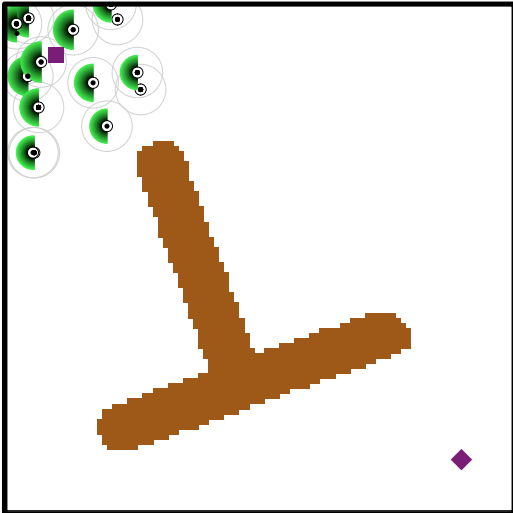
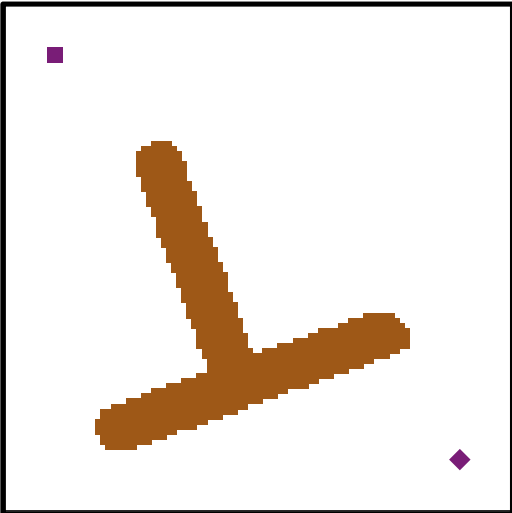


Example: Ant Pheromone Foraging **With Beacons**

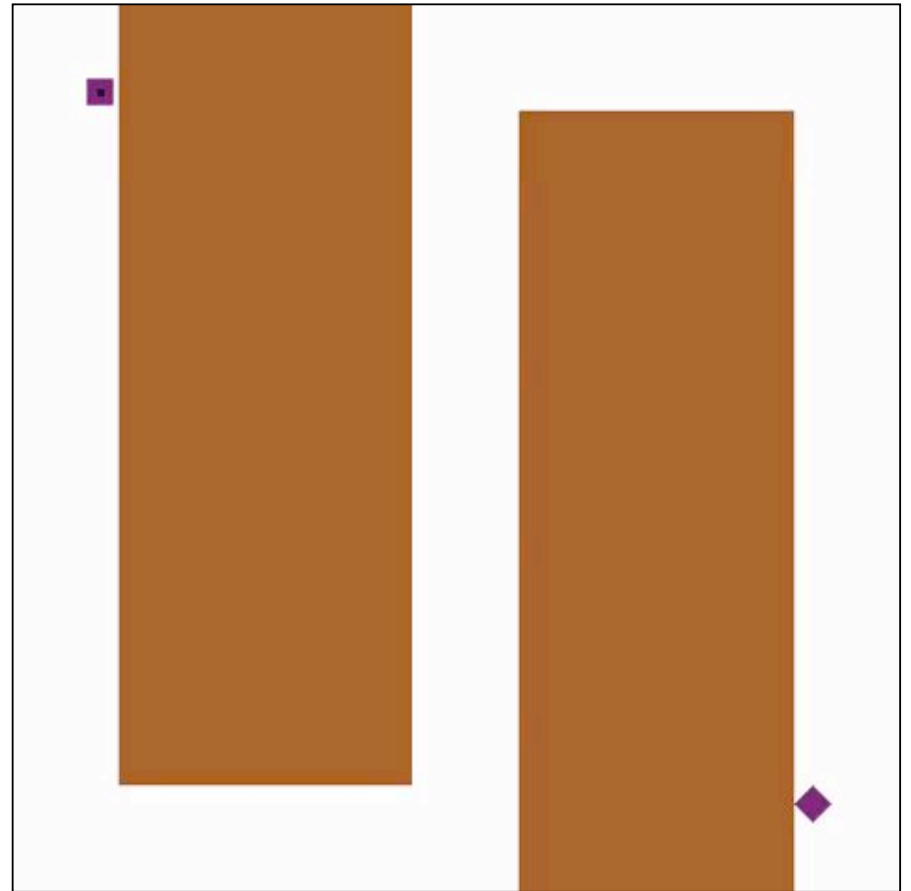
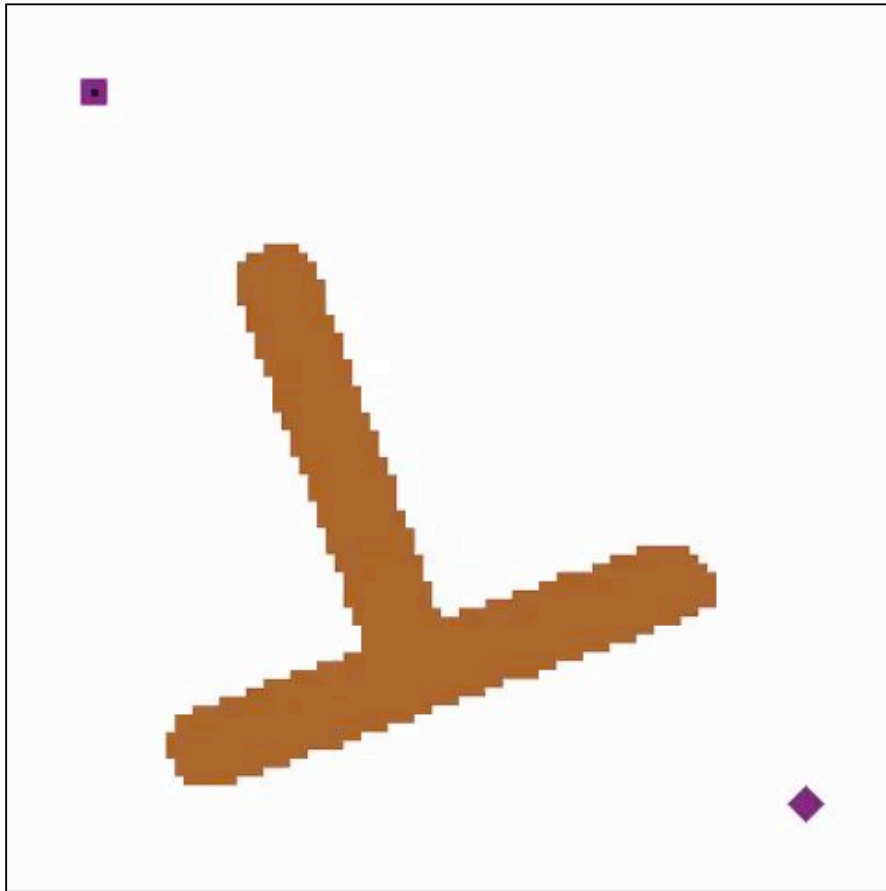
- **The Flockbots**
- **Small (15cm diameter) differential drive robots capable of deploying, moving, and removing cans**
- **Cans contain Sensor Motes which act as movable pheromon**



Example: Ant Pheromone Foraging **With Beacons**



Example: Ant Pheromone Foraging **With Beacons**



Agent Learning and Training

- **Machine Learning**

Given a sample of data drawn from an environment, construct a model which explains the environment.

- **Agent Training**

An agent is using machine learning, but there is a **trainer** present who observes the agent build and use its model, and suggests corrections.

- **Learning from Demonstration**

A robot learns to do a task after being given sample data by a human. This is **training** only if the human iteratively updates the sample data to provide corrections or suggestions. It is also **very expensive**.

- **Our Research**

1. Develop methods to do training of nontrivial **single agent** behaviors.
2. Develop methods to do training of nontrivial **multiagent** behaviors.

Single and Multi-Agent Training with Few Samples

- **Single-Agent Training Challenge**

The **Curse of Dimensionality**. The size of the training / learning space can be very large for complex behaviors, but the number of samples is very small.

- **Multi-Agent Training Challenge**

The **Multiagent Inverse Problem**. Training multiple agents presents a difficult inverse problem which gets worse and worse with more agents, more interactions, and more complex behaviors.

Current Learning from Demonstration Systems

- **Learning Paths or Trajectories**

Large numbers of samples

Machine learning is easy

- **Learning Behaviors or Plans**

Small numbers of samples

Machine learning is very difficult

- We want to learn **sophisticated** behaviors based on a very **small** number of samples.

HiTAB

(Single-Agent Training)

- **Goal**

Train **complex, stateful** behaviors from a very **small** number of samples in **real time** on simulated agents or robots.

- **Difficulty**

Curse of dimensionality. Robot behaviors can be complex, but we only have to train on a **small** number of samples.

- **Solution: Behavioral Decomposition**

Manually break complex behaviors into simpler behaviors. Learn the simpler behaviors. Then learn their composition into the complex behaviors.

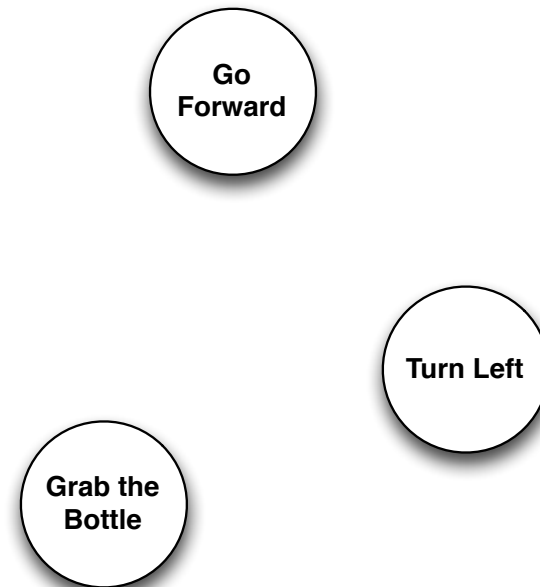
This projects the complex behaviors' joint space into smaller, simpler spaces that are much easier to learn with few samples.

HiTAB Single-Agent Model

- **Hierarchical Finite-State Automata (HFA) as Moore Machines**
 - Each **Behavior** is a **State**
 - **Recursive** Behaviors may themselves be other automata
 - **Transitions** from State to State based on environment **Features**
 - **Parameterizable** “Go to X” rather than “Go to the Ball”
- **Each timestep**
 - Transition function is queried based on current environment features, possibly resulting in a new current state
 - Current state’s behavior is pulsed one iteration

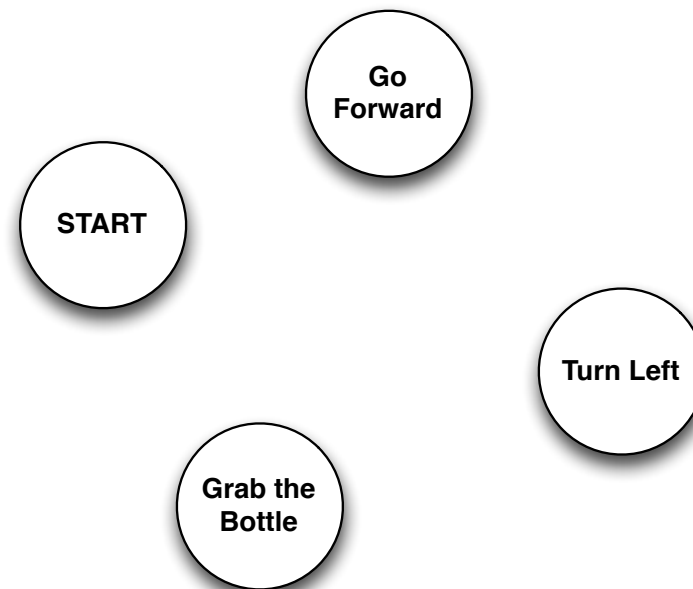
Moore Machines

- A Moore Machine is a **Finite-State Automaton** with:
 - A set of **states** corresponding to **behaviors**



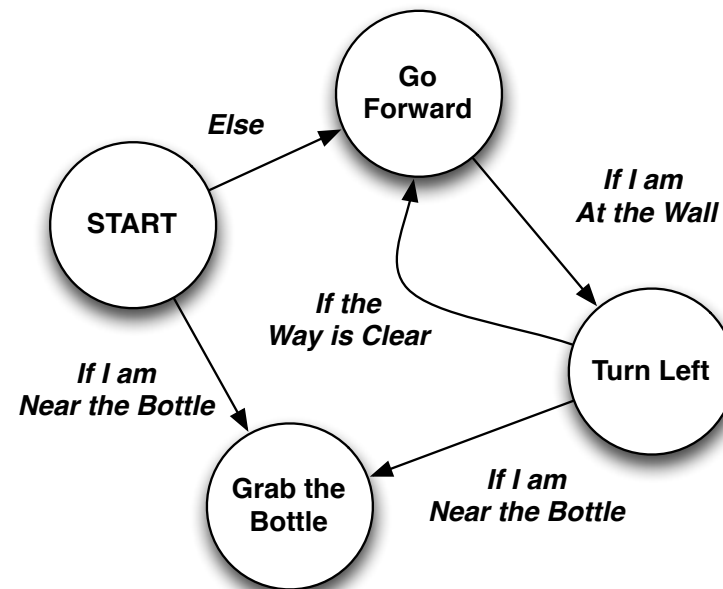
Moore Machines

- A Moore Machine is a **Finite-State Automaton** with:
 - A set of **states** corresponding to **behaviors**
 - A special **START state** (there are no end states)



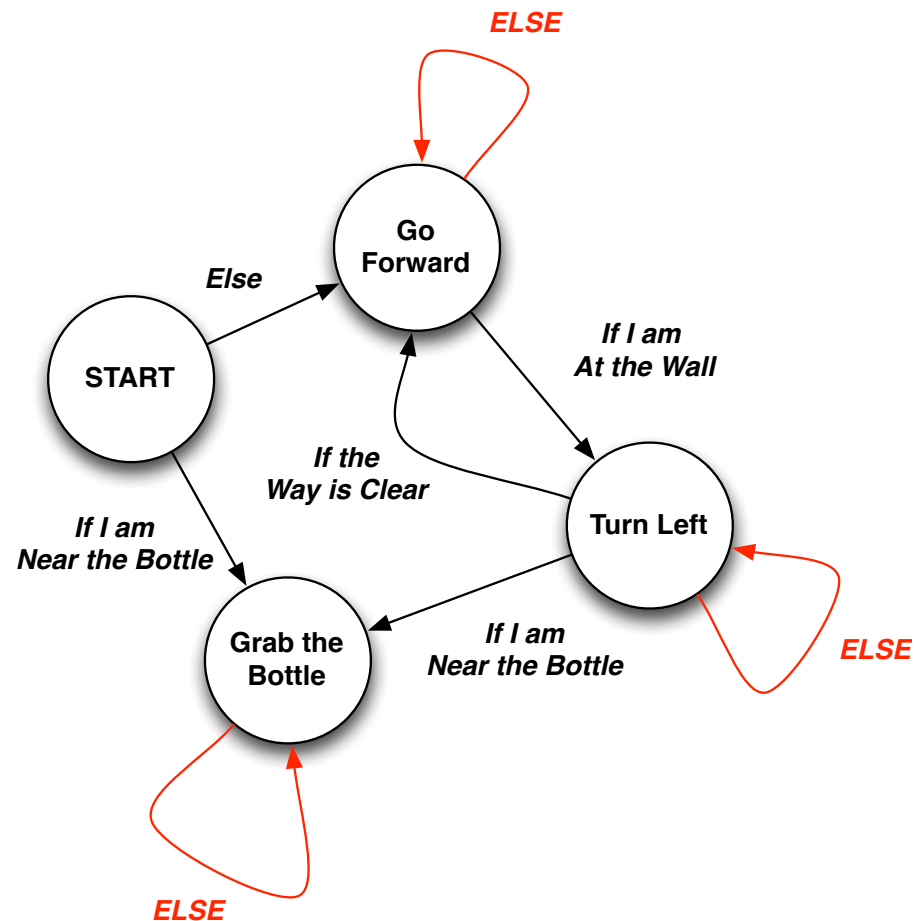
Moore Machines

- A Moore Machine is a **Finite-State Automaton** with:
 - A set of **states** corresponding to **behaviors**
 - A special **START** state (there are no end states)
 - A set of **directed edges**
 - All edges leaving a state are called its **transition function**



Moore Machines

- A Moore Machine is a **Finite-State Automaton** with:
 - A set of **states** corresponding to **behaviors**
 - A special **START** state (there are no end states)
 - A set of **directed edges**
 - All edges leaving a state are called its **transition function**
 - **No self-edges** (they are implied and mean “else”)

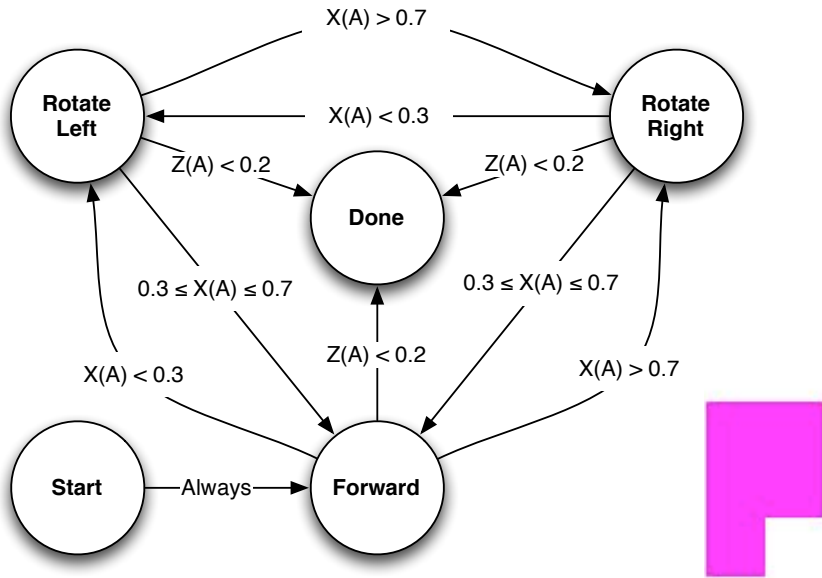




Home Base



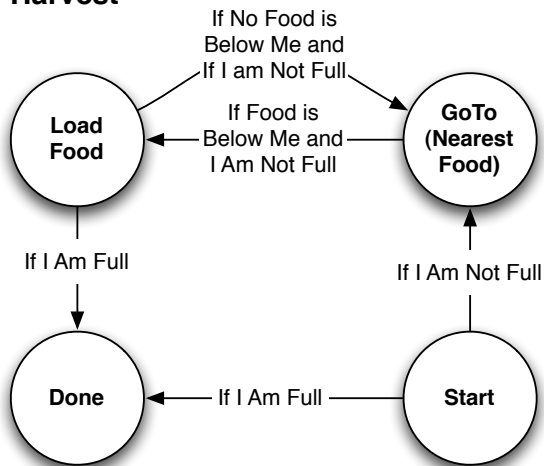
GoTo (A)



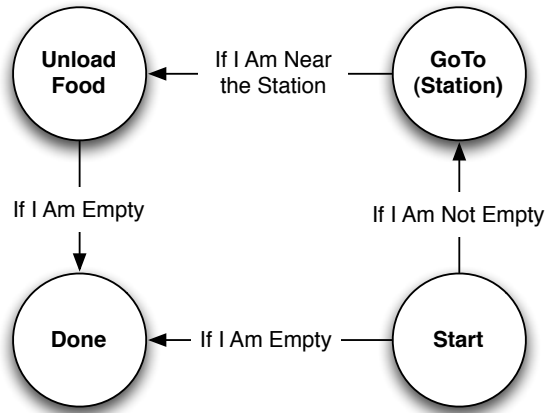
Home Base



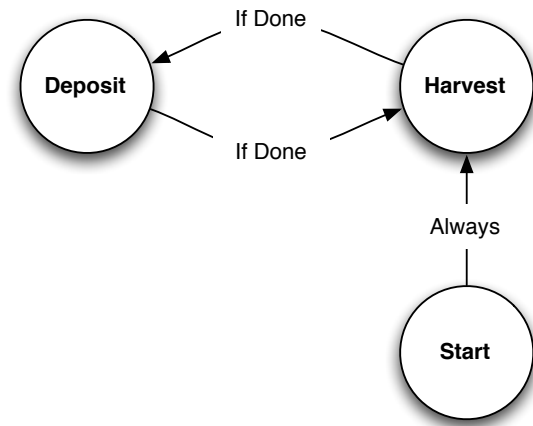
Harvest



Deposit



Forage



Training a HiTAB Automaton

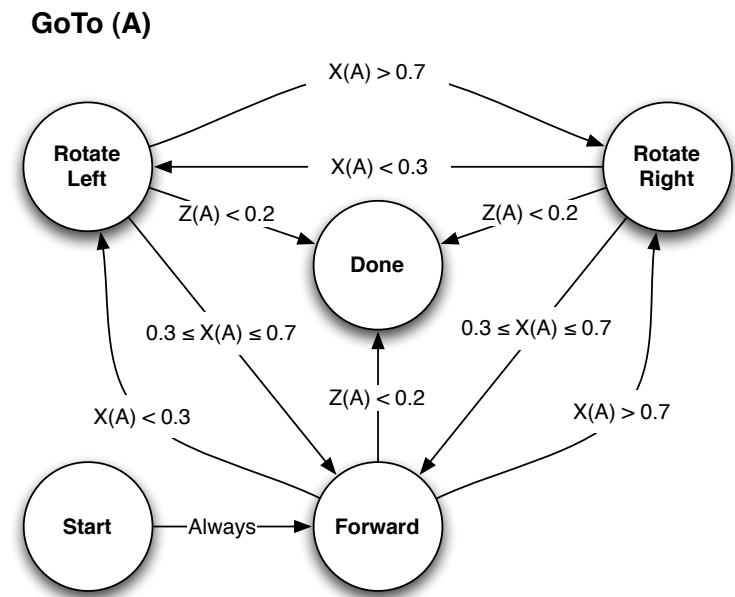
- For each state s , we learn the **transition function** $T(s, f)$ for edges leaving s .
- **Gather Data.** When the user transitions to a new state/behavior, log: [*old behavior, current feature vector, new behavior*]

- **Build $T(s, f) \Rightarrow s'$ for each state s**

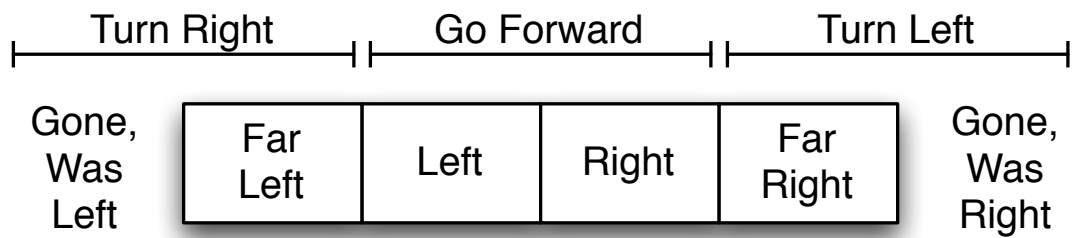
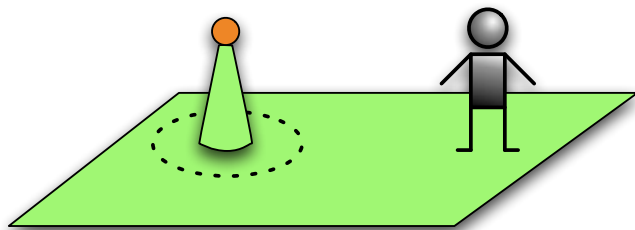
Gather all samples $[s, f, s']$ starting with s
Reduce to just $f \Rightarrow s'$

This is just a classification task

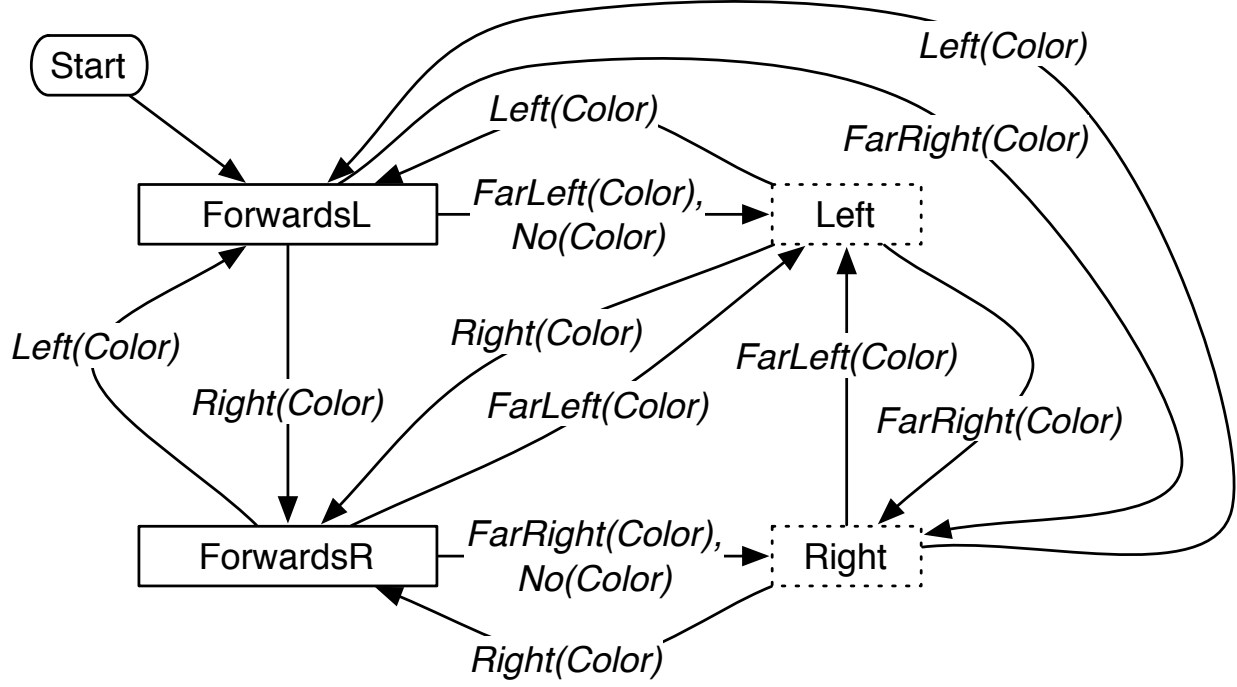
- **Delete all unused states, add to library**



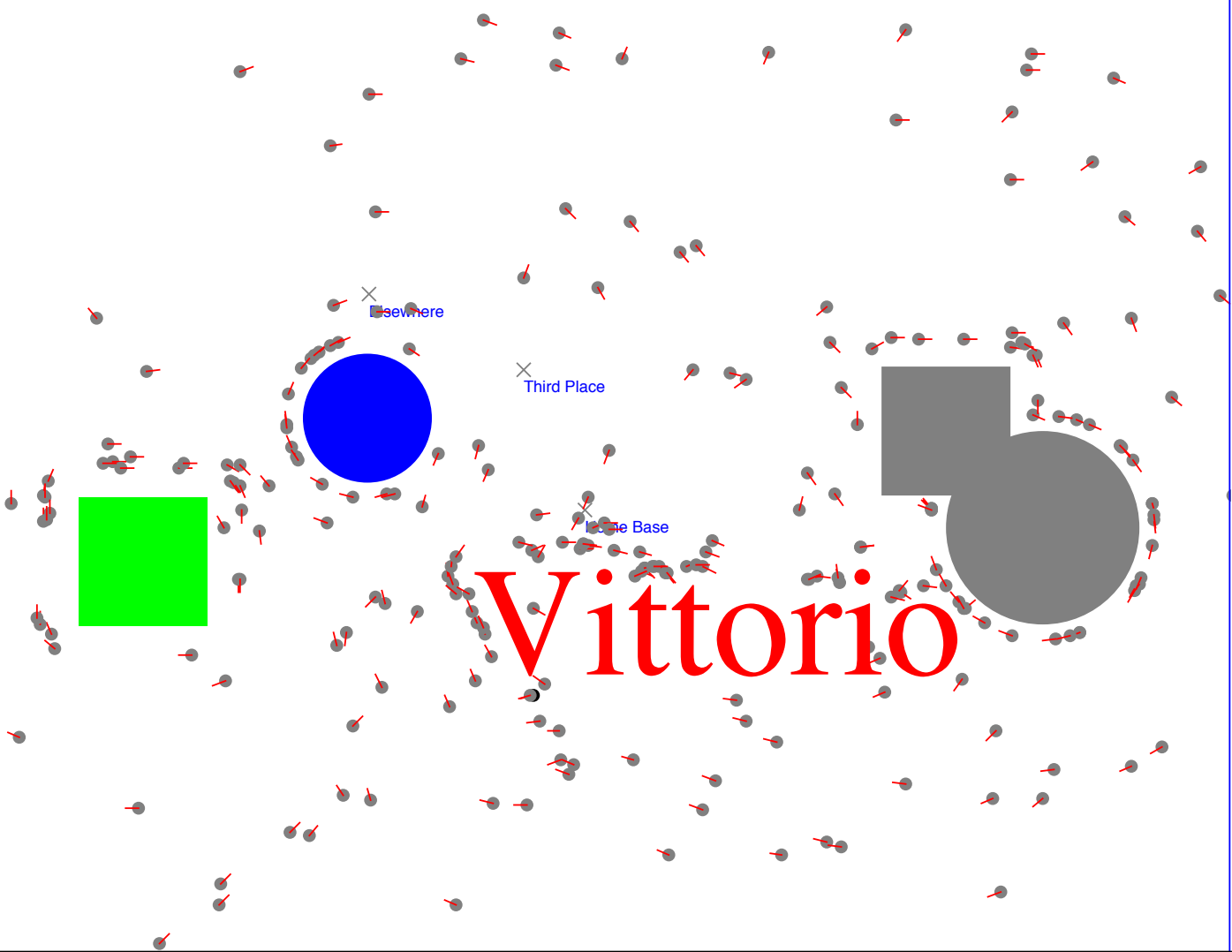
Statefulness Is Important



- A **Policy** $\pi(f) \rightarrow a$ is not a sufficiently rich representation to learn many robot behaviors.
- We learn a finite state machine transition function $T(s, f) \rightarrow s$



Demonstration...



Unlearning: Training Despite Noise (IJCAI 2013)

- **Situation: Training**

When the agent performs its learned behavior incorrectly, the trainer **corrects the behavior.**

- **Problem**

How do we use the corrective information to update the model?

- **Complication**

We have a **very small number of samples.** (Samples are precious).

- In typical machine learning (with many samples), we'd just add the corrective samples to our sample set and re-learn the model.

- In **unlearning**, we use the corrective samples to **detect and remove noisy sample data.**

Unlearning

- **We have:**

- S** Original sample set (with some noisy samples)

- M** Original learned model from S

- C** Set of corrective samples

- **We produce:**

- S'** Revised sample set (identifying/removing some noisy samples)

- M'** Revised learned model from S'

- **Approach**

- Identify the samples $B \subseteq S$ which caused M to misclassify C

- Determine which samples in $N \subseteq B$ are *likely* to be noise

- Remove N from S, producing S'

Identifying Noise in Samples

- **Identifying B requires algorithms customized for your particular model algorithm**
C4.5, K-NN, SVMs
- **A sample in $b \in B$ caused M to misclassify $c \in C$ for two reasons:**
 1. b is **noisy** *or*
 2. The sample space in S is too **sparse**, so b was inappropriately made responsible too large a region.
- Based on the model M and the algorithm which produced b , we determine if it's **probably** #1 or #2
 - How many other samples are misclassifying c ? [if many, it's likely #2]
 - How far is b from c ? [if far, it's likely #2]

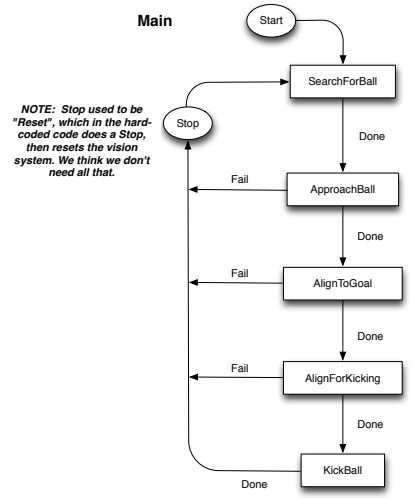
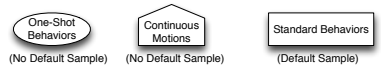
Typical Results

Dataset	<i>Noise = 1/5</i>				<i>Noise = 1/20</i>				<i>Noise = 1/100</i>			
	<i>U+C</i>	<i>U+C+E</i>	<i>Metric</i>	<i>Non-Metric</i>	<i>U+C</i>	<i>U+C+E</i>	<i>Metric</i>	<i>Non-Metric</i>	<i>U+C</i>	<i>U+C+E</i>	<i>Metric</i>	<i>Non-Metric</i>
<i>1-NN</i>												
Iris	0.9553	0.9131	0.9307	0.9255	0.9553	0.8002	0.8901	0.8601	0.9553	0.7519	0.9461	0.8490
Glass	0.6921	0.6707	0.6810	0.6822	0.6921	0.6441	0.6816	0.6705	0.6921	0.5653	0.6887	0.6421
Wine	0.9533	0.9370	0.9464	0.9442	0.9533	0.7998	0.9506	0.8722	0.9533	0.7566	0.9520	0.8488
<i>3-NN</i>												
Iris	0.9537	0.9409	0.9468	0.9492	0.9537	0.8887	0.9361	0.9295	0.9537	0.8539	0.9370	0.9331
Glass	0.7008	0.6734	0.6895	0.6980	0.7008	0.6615	0.6927	0.6971	0.7008	0.6193	0.6866	0.6828
Wine	0.9615	0.9524	0.9607	0.9594	0.9615	0.8895	0.9511	0.9472	0.9615	0.8548	0.9462	0.9408
<i>Decision Tree (Unpruned)</i>												
Iris	0.9459	0.8705	0.8915	0.8877	0.9459	0.8029	0.8497	0.8535	0.9459	0.8014	0.8765	0.8616
Glass	0.6701	0.6379	0.6577	0.6572	0.6701	0.6355	0.6544	0.6514	0.6701	0.6306	0.6591	0.6492
Wine	0.9332	0.8321	0.8638	0.8636	0.9332	0.7375	0.8103	0.7956	0.9332	0.7206	0.8365	0.8079
<i>Decision Tree (Pruned)</i>												
Iris	0.9427	0.9135	0.9213	0.9226	0.9427	0.8761	0.9081	0.9094	0.9427	0.8799	0.9250	0.9213
Glass	0.6711	0.6330	0.6520	0.6529	0.6711	0.6274	0.6460	0.6426	0.6711	0.6301	0.6501	0.6496
Wine	0.9340	0.8591	0.8811	0.8846	0.9340	0.8185	0.8749	0.8715	0.9340	0.8093	0.8892	0.8844
<i>Support Vector Machine</i>												
Iris	0.9102	0.3886	0.4280	0.9070	0.9102	0.7389	0.8649	0.8705	0.9102	0.7374	0.8695	0.8668
Glass	0.3346	0.3311	0.3163	0.3393	0.3346	0.3329	0.3313	0.3284	0.3346	0.3249	0.3259	0.3350
Wine	0.9329	0.3906	0.3991	0.9350	0.9329	0.6400	0.8828	0.8861	0.9329	0.6544	0.8834	0.8867

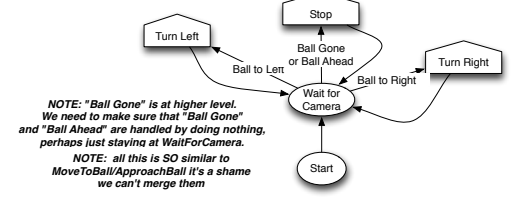
RoboCup 2012

- Use HiTAB to train a humanoid robot team **at the competition**
- **Learn 17 Finite-State Automata**

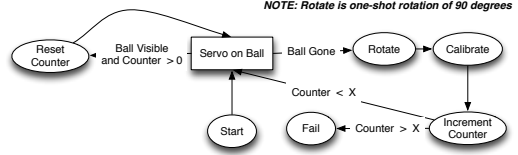




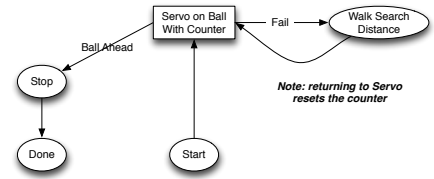
Servo on Ball



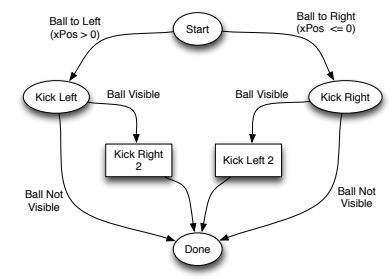
Servo on Ball With Counter



Search for Ball

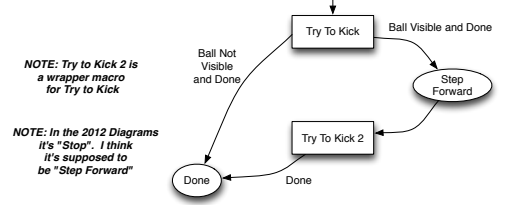


Try to Kick

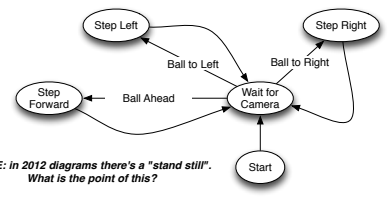


NOTE: Kick Right 2 and Kick Left 2 are wrapper macros for Kick Right and Kick Left, or alternatively are just separately saved-out kick-right and kick-left states

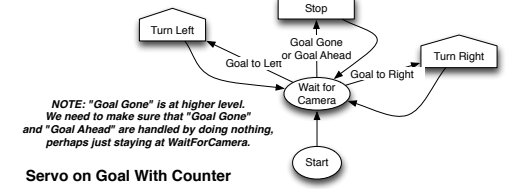
Kick Ball



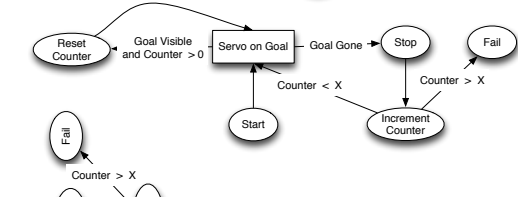
Aim for Kick



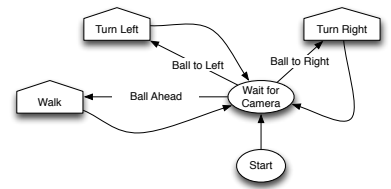
Servo on Goal



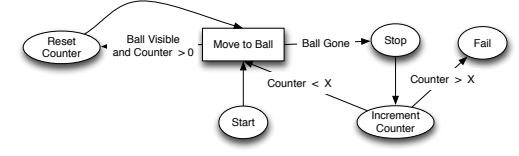
Servo on Goal With Counter



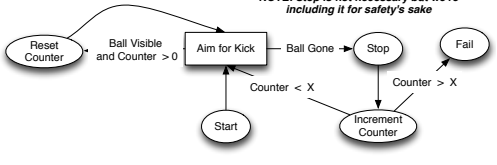
Move to Ball



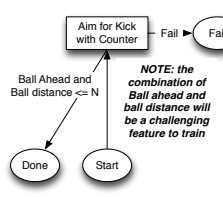
Move to Ball With Counter



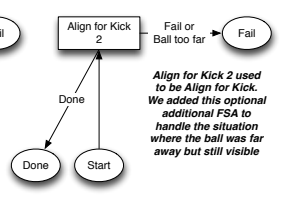
Aim for Kick with Counter



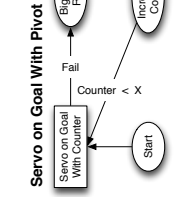
Align for Kick 2



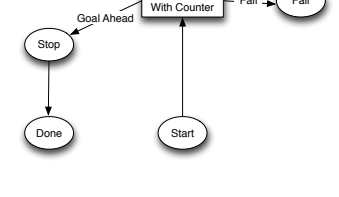
Align for Kick



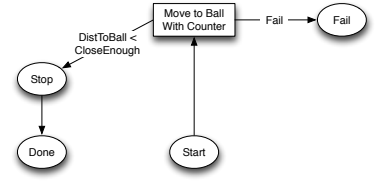
Servo on Goal With Pivot



Align to Goal



Approach Ball



Simple “Flat” Swarms with HiTAB

- **Homogenous Case:** Every agent uses the same behavior.
*This is **not just parallel**: the agents interact.*
- **Heterogeneous Case:** Agents belong to disjoint classes. Only agents in the same class use the same behavior.
- If the **interesting** behaviors require interaction, how do you train agents simultaneously?
- Example: to passing behaviors, you must teach **two** robots at the same time how to coordinate passing and receiving.

Behavioral Bootstrapping

- If you have **multiple agents that must be trained simultaneously**
- ... and you only have **one trainer ... ?**

- **Homogeneous Case**
 1. Set all agents to empty behaviors (doing nothing)
 2. Select an Agent and train a **slightly better behavior** in the context of the agents' existing behaviors
 3. Distribute this behavior to all the agents
 4. Go to 1

Behavioral Bootstrapping

- **Heterogeneous Case** *(2-agent example)*

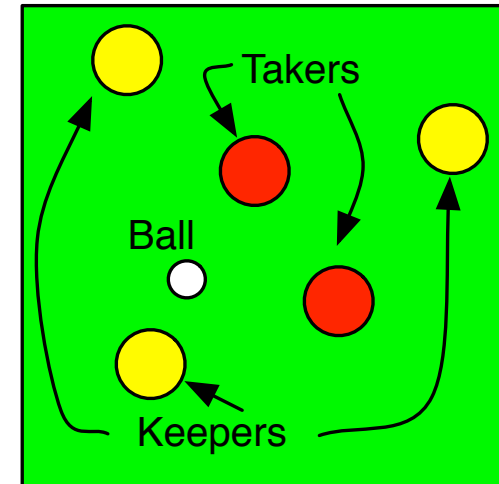
1. Set both agents to empty behaviors (doing nothing)
2. Select Agent A and train a **slightly better behavior** in the context of Agent B's existing behavior
3. Select Agent B and train a **slightly better behavior** in the context of Agent A's existing behavior
4. Go to 1

Behavioral Bootstrapping: **Keepaway Soccer**

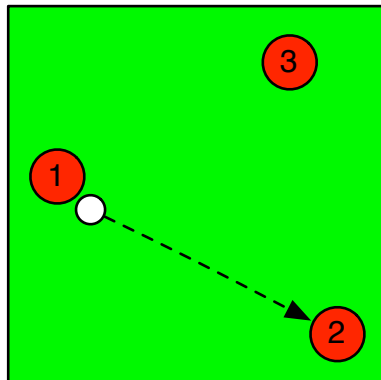
- **Three Keepers, Two Takers**

The Keepers have control of the ball
The Takers are trying to take the ball

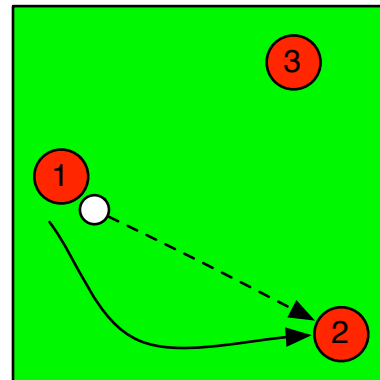
The Takers are hard-coded
We are training the Keepers (Homogeneous)



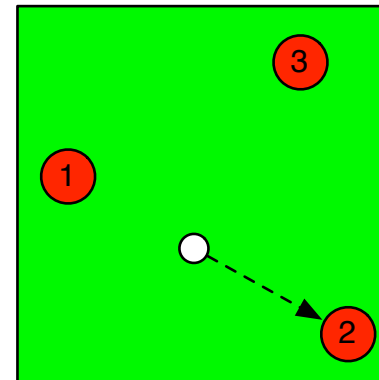
- **Passing Requires coordination between a passer and a receiver**



Player 1 decides to pass to Player 2



As Player 1 passes, it also yells to Player 2



Player 2 stops trying to Get Open and prepares to Receive

Behavioral Bootstrapping: **Keepaway Soccer**



Behavioral Bootstrapping: **Keepaway Soccer**

- **Results**

- **University of Texas, Austin Hard-Coded Team**

 - 5.6 Seconds On Average (before takers take the ball)

- **George Mason University Bootstrapped Team**

 - 7 Seconds on Average

 - 9 Seconds on Average if using “yelling”

Multiagent Training

- **Techniques for Multiagent Training** are nearly always optimizers.
 - Multiagent Reinforcement Learning, Stochastic Optimization
 - Supervised Learning is **extremely rare** for multiagent training. Yet training is a supervised task!
- **User Modeling** The team learns about one another
- **Training (or Demonstration)** The team learns to do a task set by you

The MAS Inverse Problem

- **Emergence** Given the micro-behaviors, we can't guess the emergent macro-phenomenon without simulation.
- **The MAS Inverse Problem** Given a desired emergent macro-phenomenon, we can't guess the micro-behaviors **at all**.
- **How this Affects Training:**
 - The trainer can tell the agents “in situation X, the macro-phenomenon should be Y” (when it's dark, storm the castle)
 - To learn, an agent needs to know “in situation X, my micro-behavior should be Z” (when it's dark, stay to the left of Bob)
 - **We can't easily compute the micro-behaviors to achieve the desired macro-phenomena**

Optimization Solves Inverse Problems

- **Training With an Optimizer:**

- Create a new candidate solution consisting of micro-behaviors.
- Test in the simulator to observe the resulting macro-phenomenon.
- Assess the error in the macro-phenomenon.
- Repeat.

Optimization Solves Inverse Problems

- **Supervised Learning Doesn't Work**

Multiagent Systems Inverse Problem. The separation between the micro-behaviors and macro-level phenomenon is too large

- **Stochastic Optimization**

- Simulated Annealing, Hill-Climbing, *etc.*: test one solution at a time
- Evolutionary Computation: test many solutions at a time
(very good for multiagent systems)

- **Reinforcement Learning**

- Q-Learning, Policy Search

- **BUT:** optimization requires many trials to gather samples. In robotics, a trial is *very expensive*.

Multi-Agent HiTAB: Training Hierarchies of Swarms

- **Goal**

Train **complex**, stateful behaviors from a very small number of samples in real time in **arbitrarily large swarms of agents**.

- **Difficulties**

1. Curse of dimensionality. [like single-agent]
2. The **Multiagent Inverse Problem**.

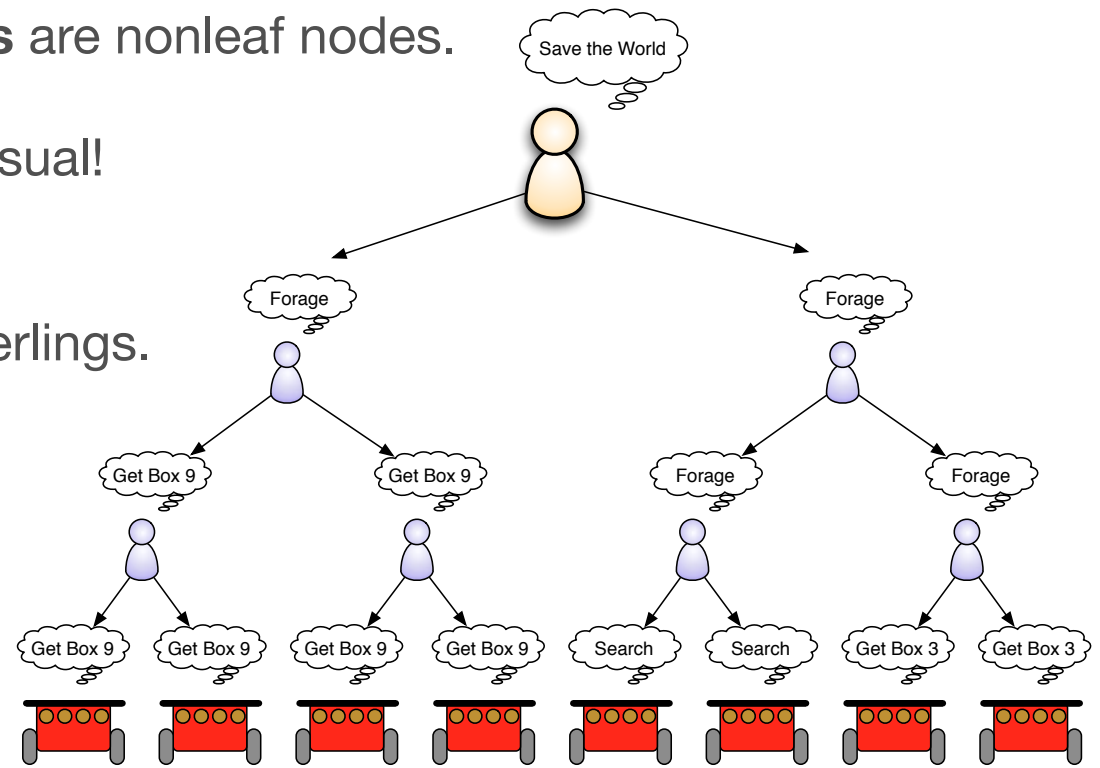
- **Solution: Swarm Decomposition**

Manually break the joint multiagent behaviors into simpler behaviors for smaller sub-swarms.

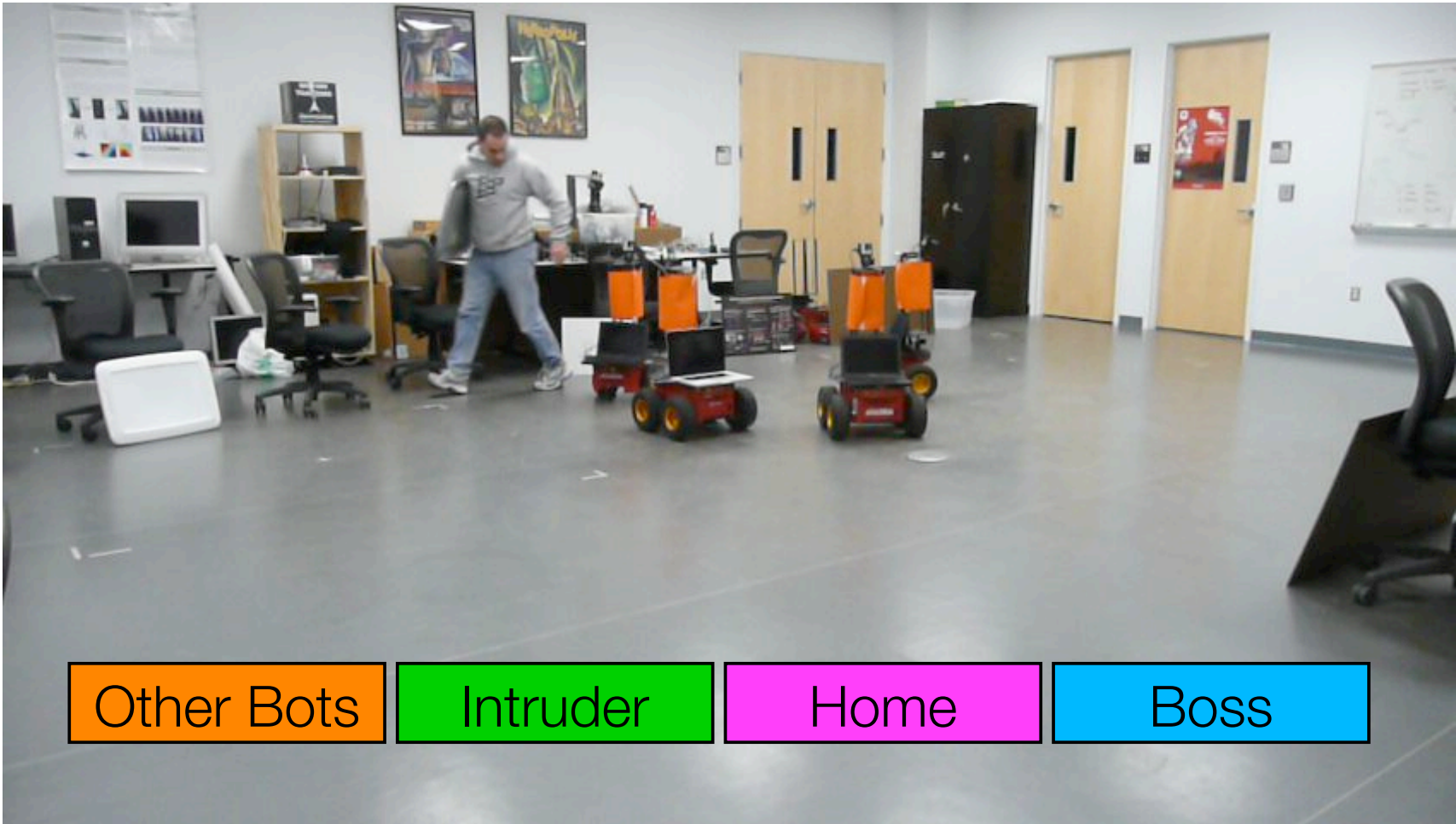
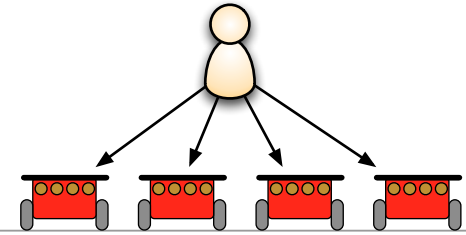
Train the simpler behaviors on small swarms, then train composed behaviors on larger swarms.

HiTAB Multi-Agent Model

- Decompose the swarm into a **hierarchy of subswarms**.
- “Regular” (real) agents are leaf nodes.
- **Controller (“boss”) agents** are nonleaf nodes.
- Train controller agents as usual!
- **Basic Behaviors**
Top-level behaviors of underlings.
- **Features**
Statistics about underlings.

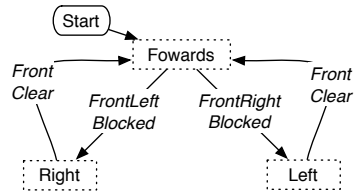


Simple Multiagent Example

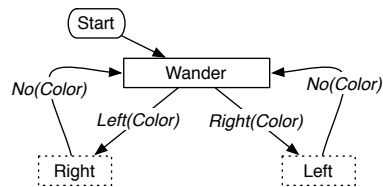


Simple Multi-Agent Example

1. Wander

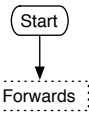


2. Disperse(Color)

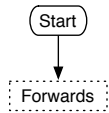


3. Various Cover FSAs

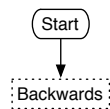
3A. ForwardsL



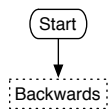
3B. ForwardsR



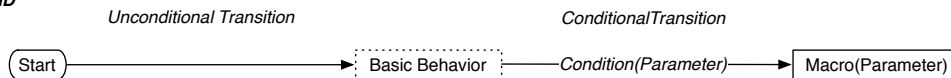
3C. BackwardsL



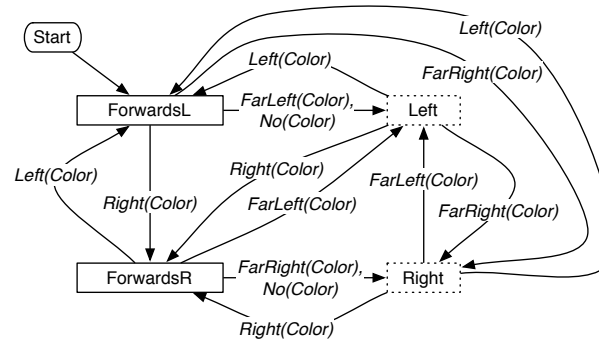
3D. BackwardsR



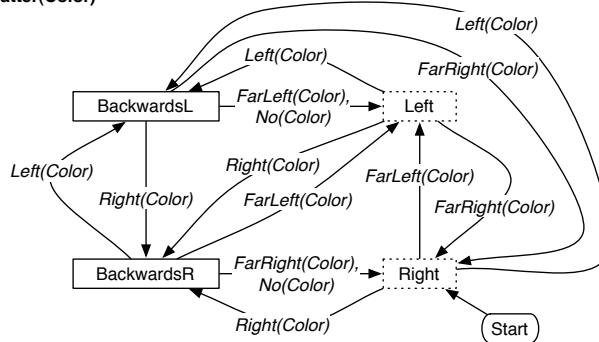
LEGEND



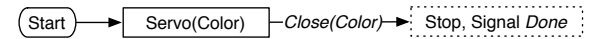
4. Servo(Color)



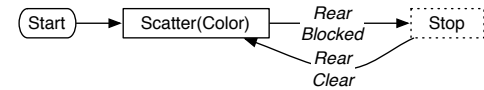
5. Scatter(Color)



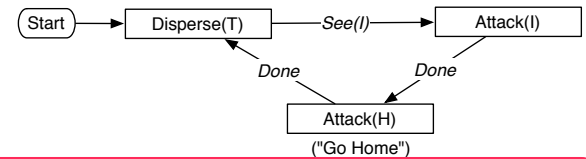
6. Attack(Color)



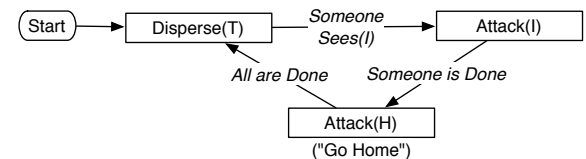
7. RunAway(Color)



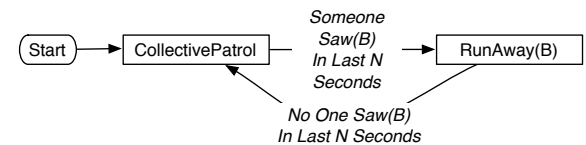
8. Patrol



9. CollectivePatrol



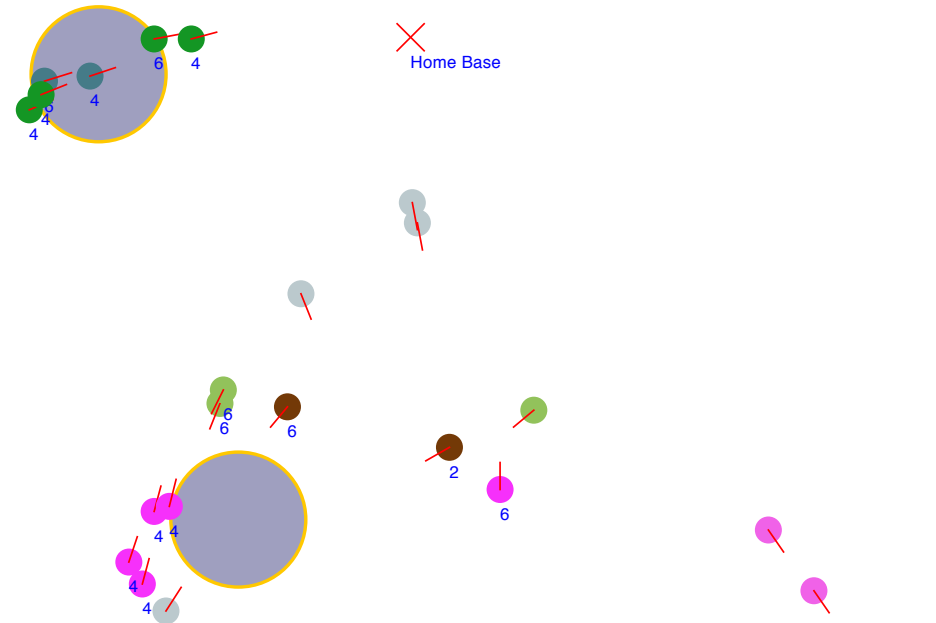
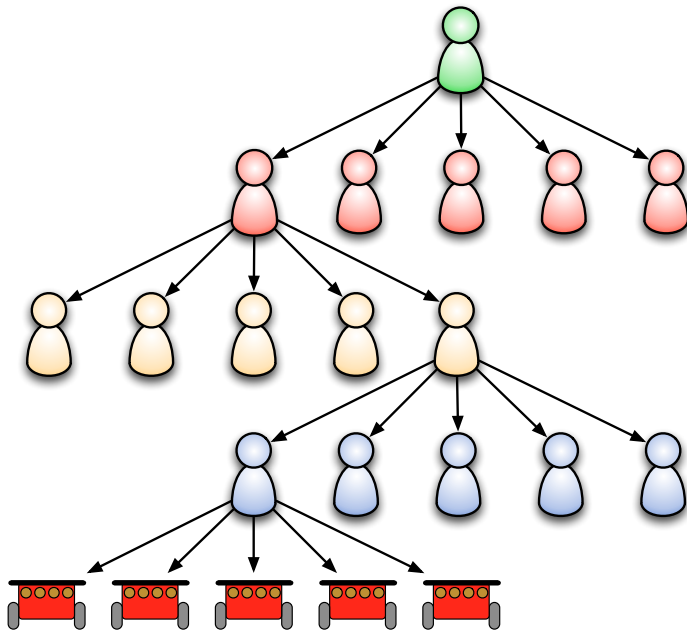
10. CollectivePatrolAndDefer



COLORS
 T Team Color
 I Intruder Color
 H Home Base Color
 B Boss Color

Larger Multi-Agent Model

- **Box Collecting**
Boxes require 5, 25, or 125 agents to retrieve
- **We've trained up to 625 agents**



Collaborators

HiTab

Daniele Nardi

Vittorio Ziparo

University of Rome, La Sapienza

Students

Ant Pheromones

Brian Hrolenok

Liviu Panait

Gabriel Balan

Katherine Russell

Single-Agent HiTab

Katherine Russell

Khaled Talukder

Ahmed ElMolla

Kevin Andrea

Multi-Agent HiTab, Unlearning, Behavioral Bootstrapping

Keith Sullivan

Bill Squires

