

# Reinforcement Learning for Sequential Decision Making

**Peter Stone\***

Director, Learning Agents Research Group (LARG)  
Department of Computer Science  
The University of Texas at Austin

Joint work with members of LARG  
past and present

Also, Cogitai Inc.

# Who Am I?

# Who Am I?

To what degree can autonomous intelligent **agents learn** in the presence of **teammates** and/or **adversaries** in **real-time, dynamic domains**?

# Who Am I?

To what degree can autonomous intelligent **agents learn** in the presence of **teammates** and/or **adversaries** in **real-time, dynamic domains**?

## Research Areas

- Autonomous agents
- Robotics
- Multiagent systems

# Who Am I?

To what degree can autonomous intelligent **agents learn** in the presence of **teammates** and/or **adversaries** in **real-time, dynamic domains**?

## Research Areas

- Autonomous agents
- Robotics
- Multiagent systems
- Machine learning

# Who Am I?

To what degree can autonomous intelligent **agents learn** in the presence of **teammates** and/or **adversaries** in **real-time, dynamic domains**?

## Research Areas

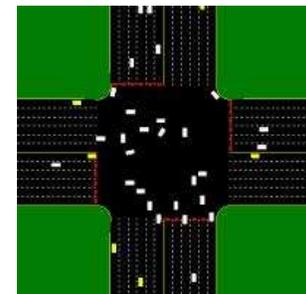
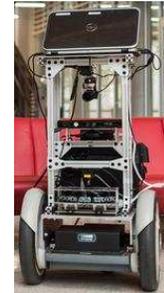
- Autonomous agents
- Robotics
- Multiagent systems
- Machine learning
  - **Reinforcement Learning**

# Who Am I?

To what degree can autonomous intelligent **agents learn** in the presence of **teammates** and/or **adversaries** in **real-time, dynamic domains**?

## Research Areas

- Autonomous agents
- Robotics
- Multiagent systems
- Machine learning
  - **Reinforcement Learning**



# AlphaGo [Silver et al., '15]: An AI milestone



# AlphaGo [Silver et al., '15]: An AI milestone



- Integrated two **existing** AI technologies

# AlphaGo [Silver et al., '15]: An AI milestone



- Integrated two **existing** AI technologies
  - Supervised learning (Deep learning)

# AlphaGo [Silver et al., '15]: An AI milestone



- Integrated two **existing** AI technologies
  - Supervised learning (Deep learning)
  - **Reinforcement learning**

# AlphaGo [Silver et al., '15]: An AI milestone



- Integrated two **existing** AI technologies
  - Supervised learning (Deep learning)
  - **Reinforcement learning**
- Does **not** solve AI

# AlphaGo [Silver et al., '15]: An AI milestone



- Integrated two **existing** AI technologies
  - Supervised learning (Deep learning)
  - **Reinforcement learning**
- Does **not** solve AI
  - The end of the line for:

# AlphaGo [Silver et al., '15]: An AI milestone

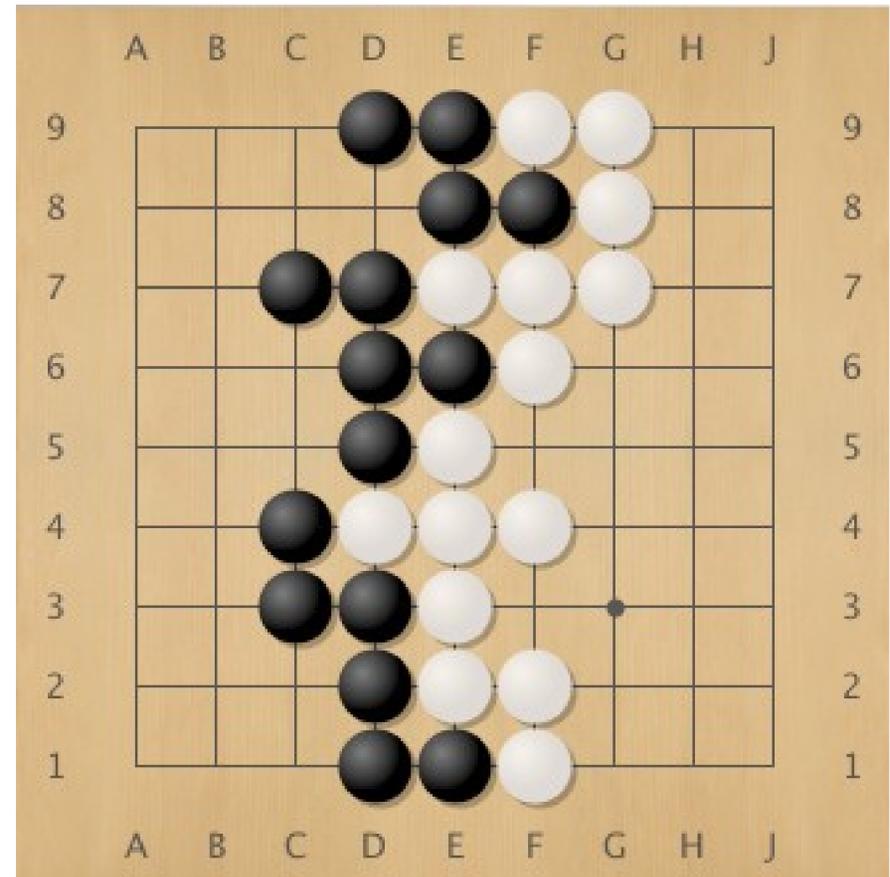


- Integrated two **existing** AI technologies
  - Supervised learning (Deep learning)
  - **Reinforcement learning**
- Does **not** solve AI
  - The end of the line for:

**2-player zero-sum discrete finite deterministic games of perfect information**

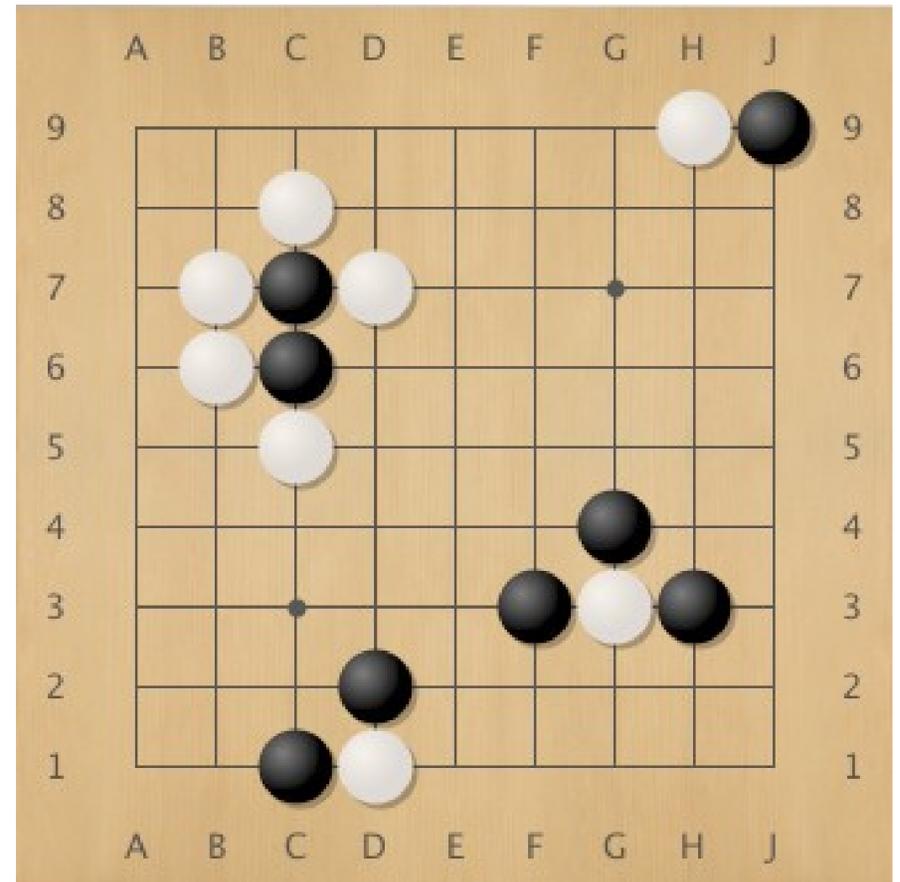
# A brief introduction to Go

- Black and white take turns to place down stones
- Once played, a stone cannot move
- The aim is to surround the most territory
- Usually played on 19x19 board



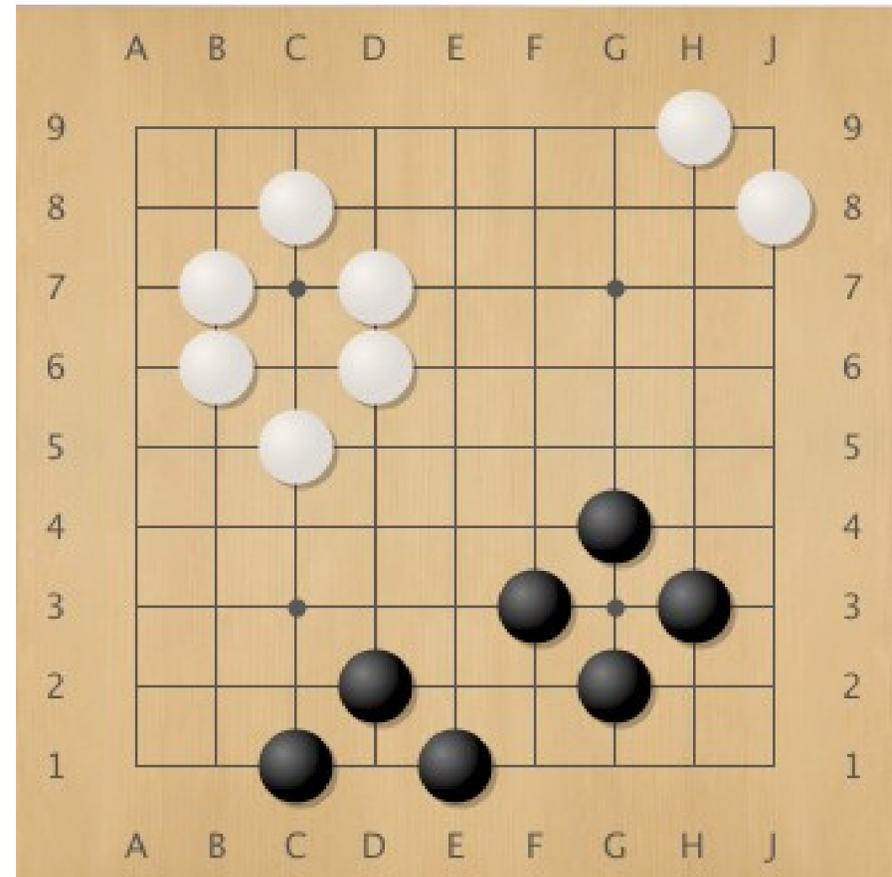
# Capturing

- The lines radiating from a stone are called *liberties*
- If a connected group of stones has all of its liberties removed then it is captured
- Captured stones are removed from the board



# Capturing

- The lines radiating from a stone are called *liberties*
- If a connected group of stones has all of its liberties removed then it is captured
- Captured stones are removed from the board

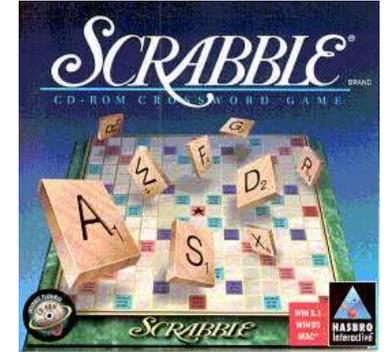
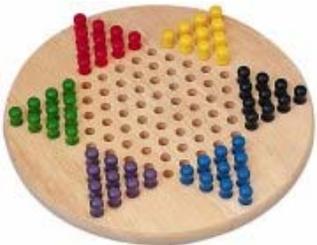
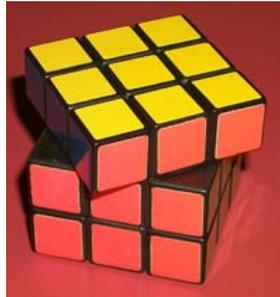


# 2-player zero-sum discrete finite deterministic games of perfect information

What do these terms mean?

- **Two player:** Duh!
- **Zero-sum:** In any outcome of any game, Player A's gains equal player B's losses. (Doesn't mean fairness: "On average, two equal players will win or lose equal amounts" not necessary for zero-sum.)
- **Discrete:** All game states and decisions are discrete values.
- **Finite:** Only a finite number of states and decisions.
- **Deterministic:** No chance (no die rolls).
- **Games:** See next page
- **Perfect information:** Both players can see the state, and each decision is made sequentially (no simultaneous moves).

# Which of these are: 2-player zero-sum discrete finite deterministic games of perfect information



- **Two player:** Duh!
- **Zero-sum:** In any outcome of any game, Player A's gains equal player B's losses.
- **Discrete:** All game states and decisions are discrete values.
- **Finite:** Only a finite number of states and decisions.
- **Deterministic:** No chance (no die rolls).
- **Games:** See next page
- **Perfect information:** Both players can see the state, and each decision is made sequentially (no simultaneous moves).

# Which of these are: 2-player zero-sum discrete finite deterministic games of perfect information



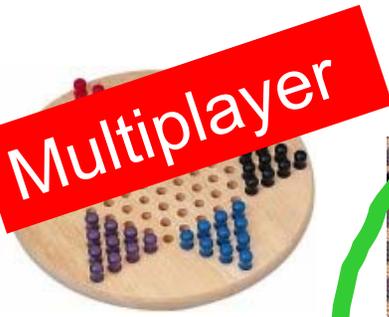
Not finite



Stochastic



One player



Multiplayer



Hidden Information

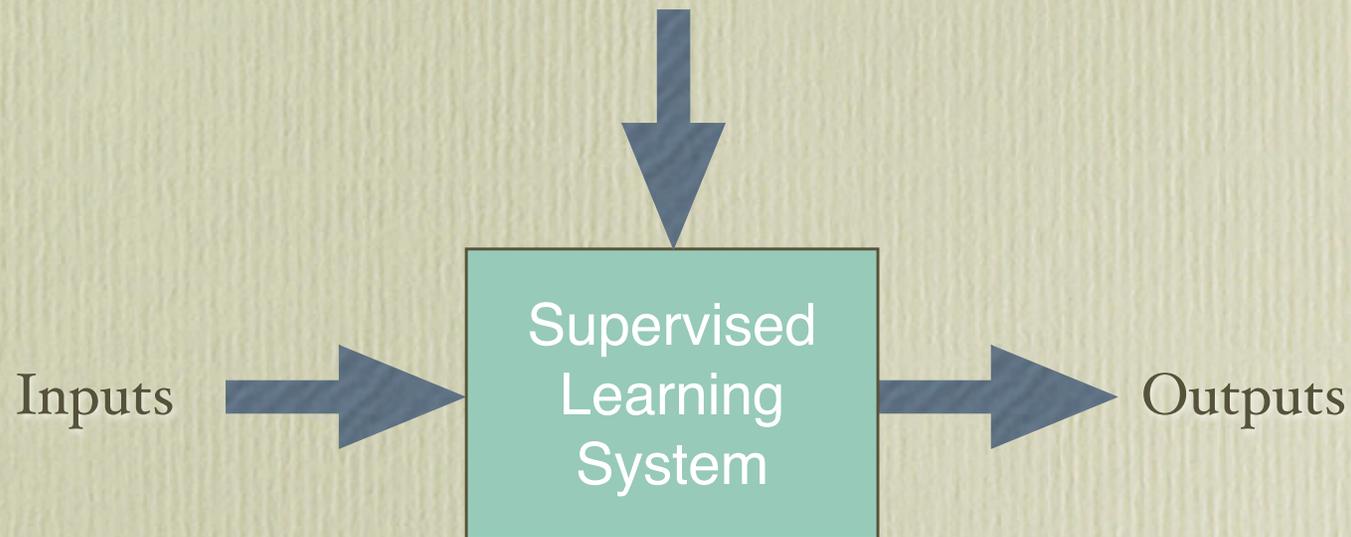


Involves Improbable Animal Behavior

- **Two player:** Duh!
- **Zero-sum:** In any outcome of any game, Player A's gains equal player B's losses.
- **Discrete:** All game states and decisions are discrete values.
- **Finite:** Only a finite number of states and decisions.
- **Deterministic:** No chance (no die rolls).
- **Games:** See next page
- **Perfect information:** Both players can see the state, and each decision is made sequentially (no simultaneous moves).

# Supervised Learning

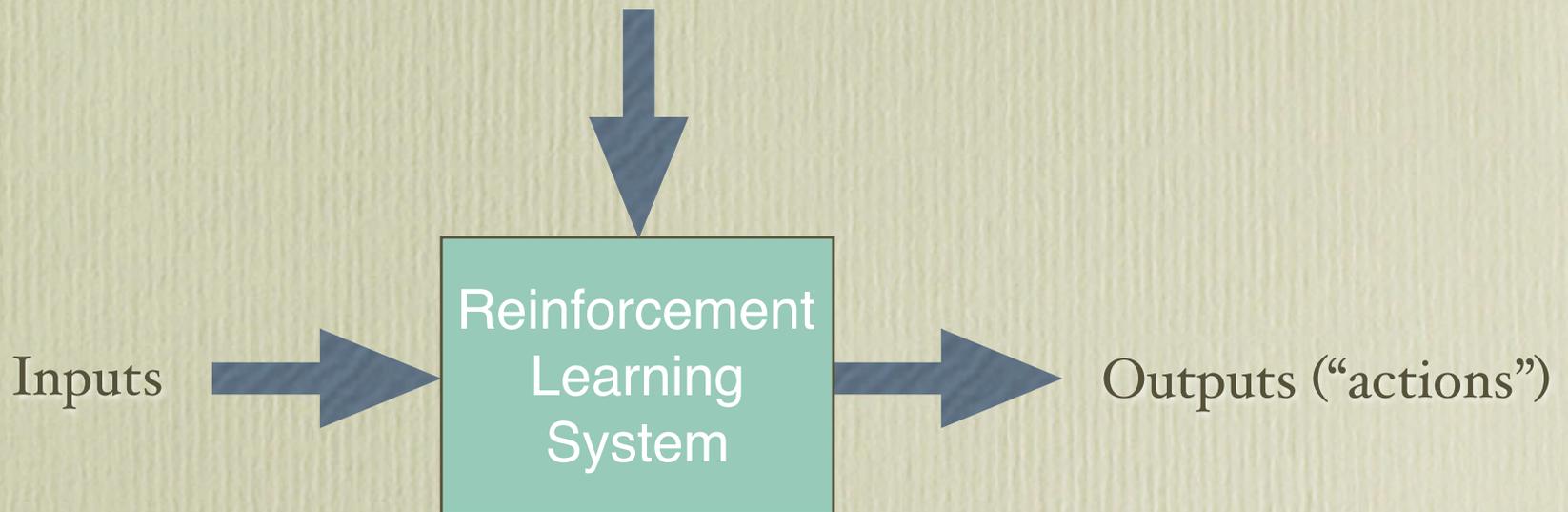
Training Info = desired (target) outputs



Error = (target output - actual output)

# Reinforcement Learning

Training Info = evaluations (“rewards” / “penalties”)



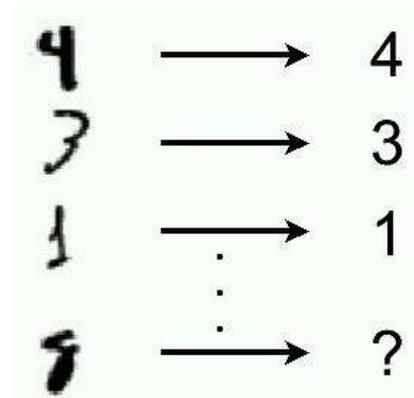
Objective: get as much reward as possible

# Key Features of RL

- Learner is not told which actions to take
- Trial-and-Error search
- Possibility of delayed reward
  - Sacrifice short-term gains for greater long-term gains
- The need to *explore* and *exploit*
- Considers the whole problem of a goal-directed agent interacting with an uncertain environment

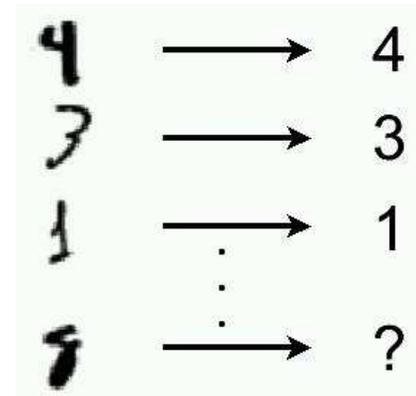
# Reinforcement Learning

Supervised learning **mature** [WEKA]



# Reinforcement Learning

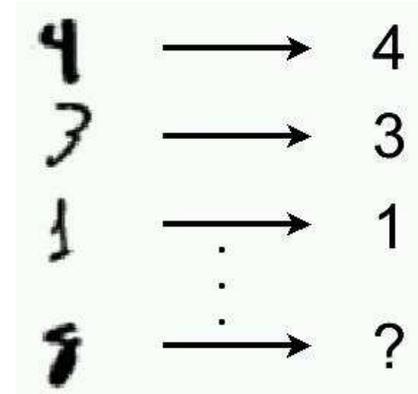
Supervised learning **mature** [WEKA]



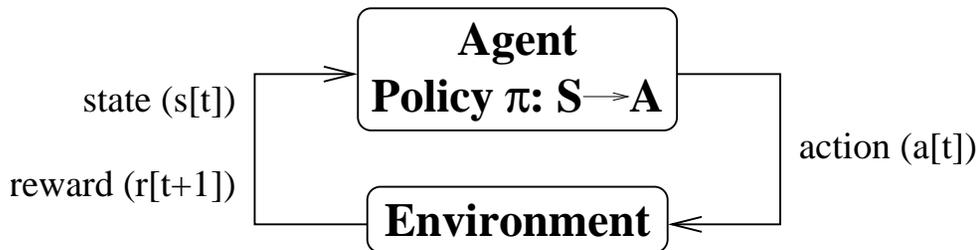
For agents, **reinforcement learning** most appropriate

# Reinforcement Learning

Supervised learning **mature** [WEKA]

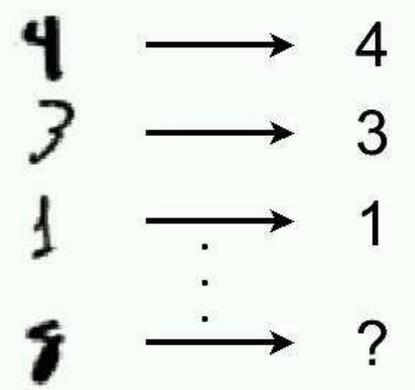


For agents, **reinforcement learning** most appropriate

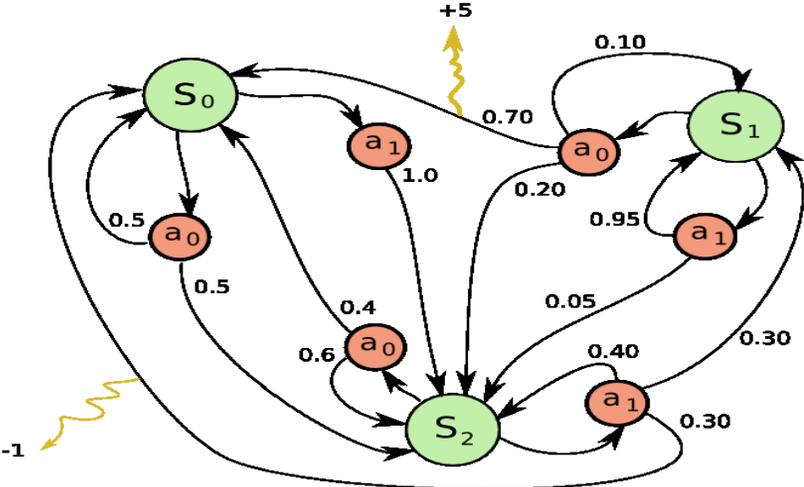
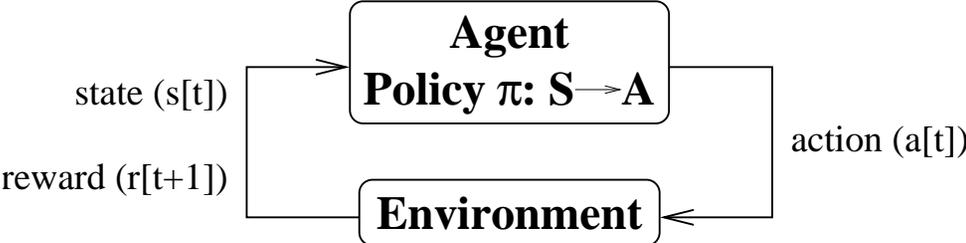


# Reinforcement Learning

Supervised learning **mature** [WEKA]

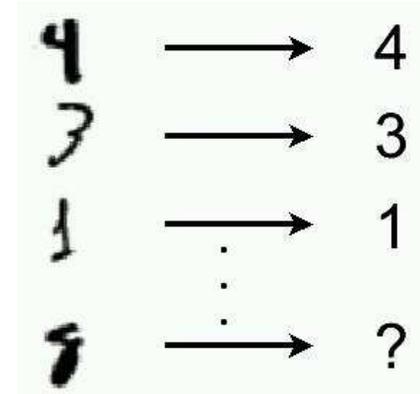


For agents, **reinforcement learning** most appropriate

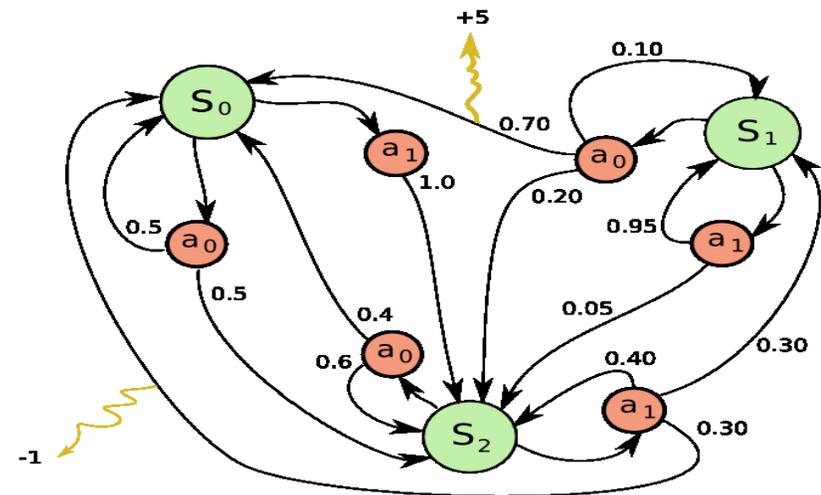
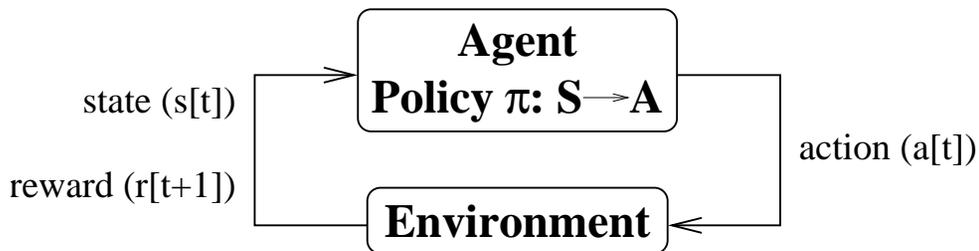


# Reinforcement Learning

Supervised learning **mature** [WEKA]



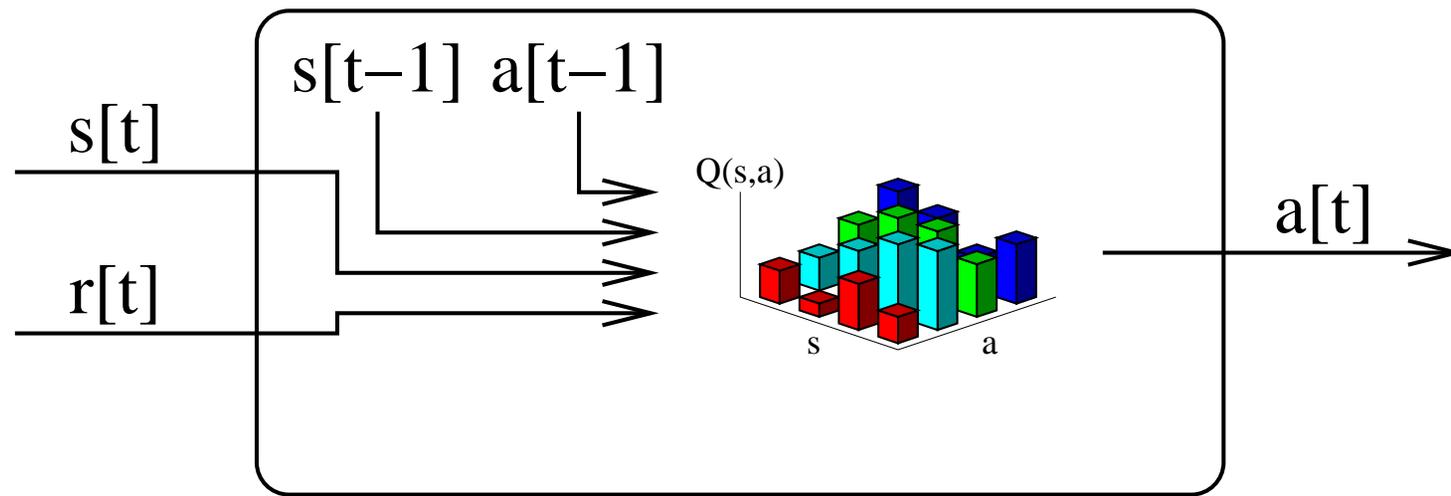
For agents, **reinforcement learning** most appropriate



- Foundational **theoretical** results
- Applications require **innovations** to scale up

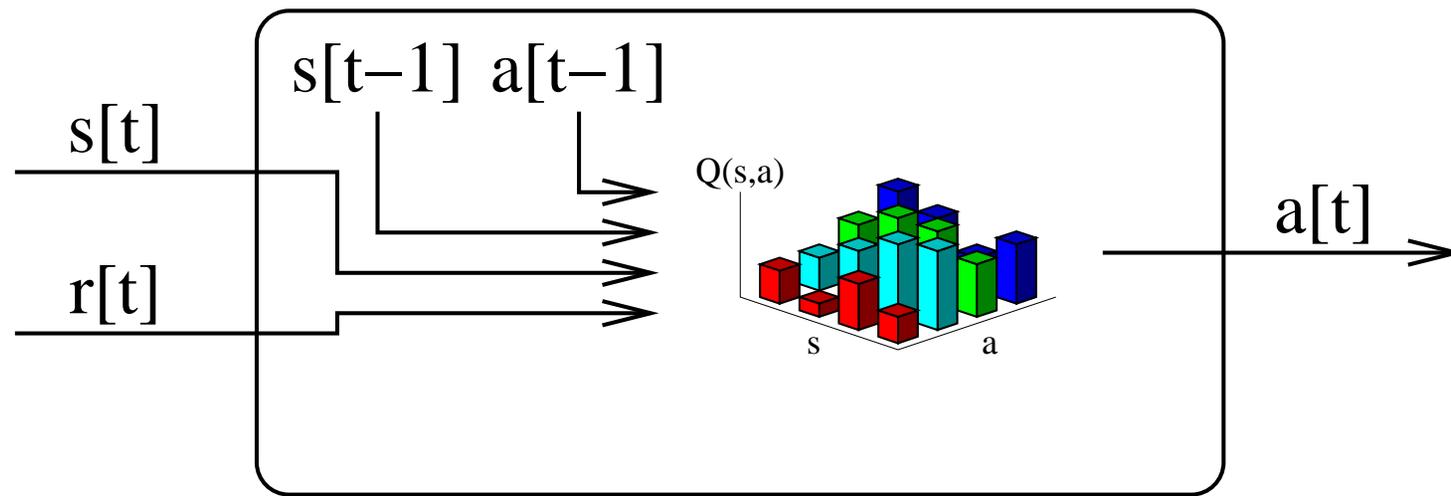
# RL Theory

Success story: Q-learning converges to  $\pi^*$  [Watkins, 89]



# RL Theory

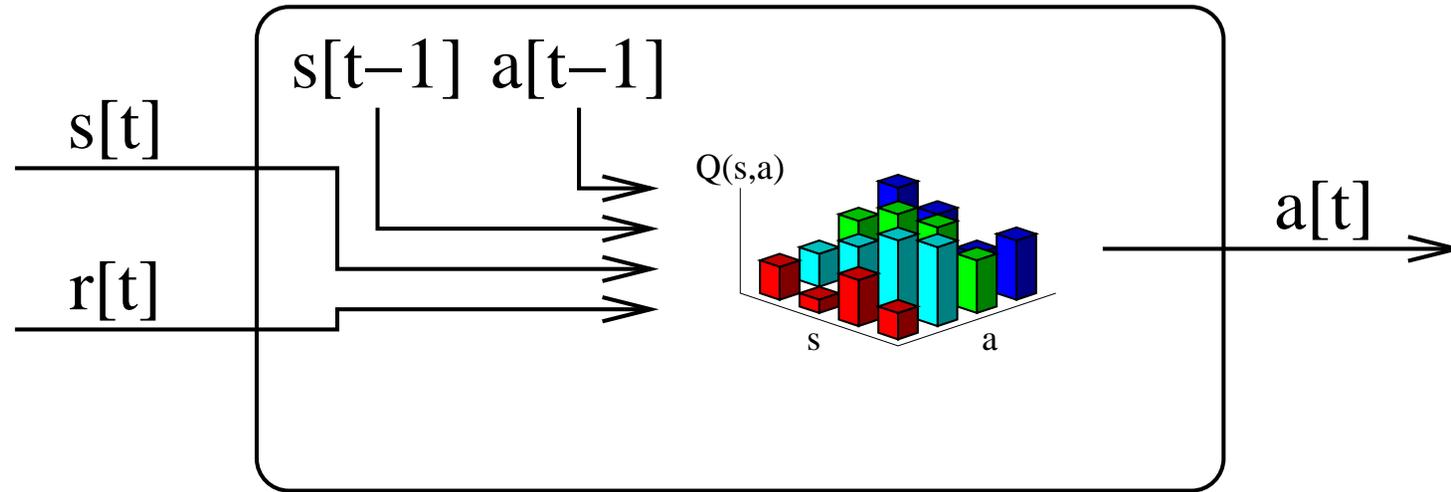
Success story: **Q-learning** converges to  $\pi^*$  [Watkins, 89]



- Table-based representation
- Visit every state infinitely often

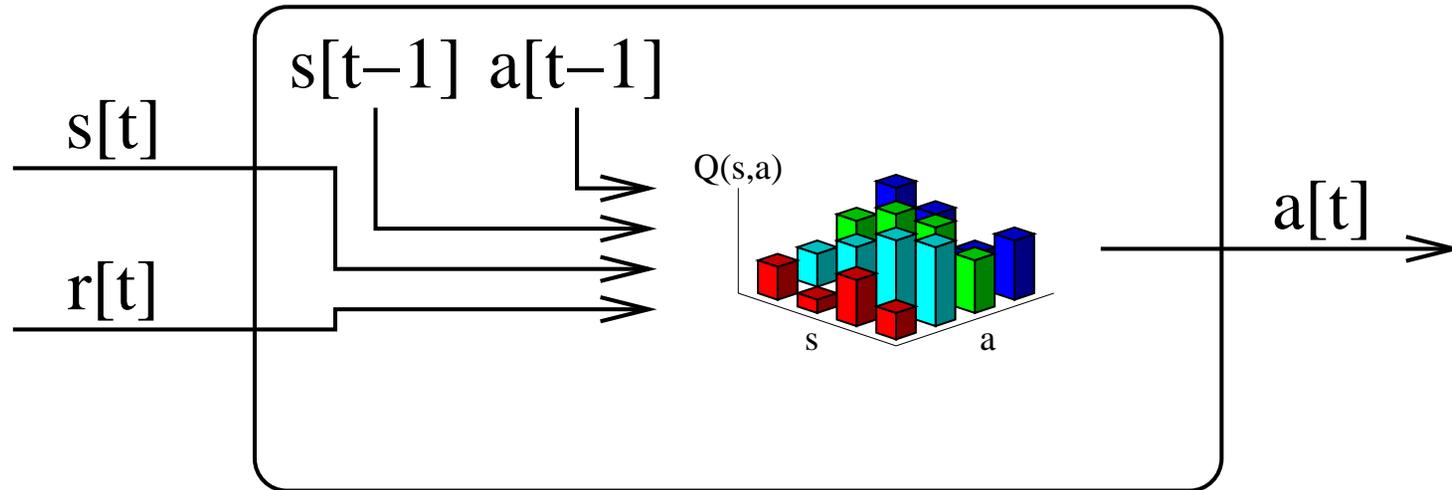
# Function Approximation

In practice, visiting **every state** impossible

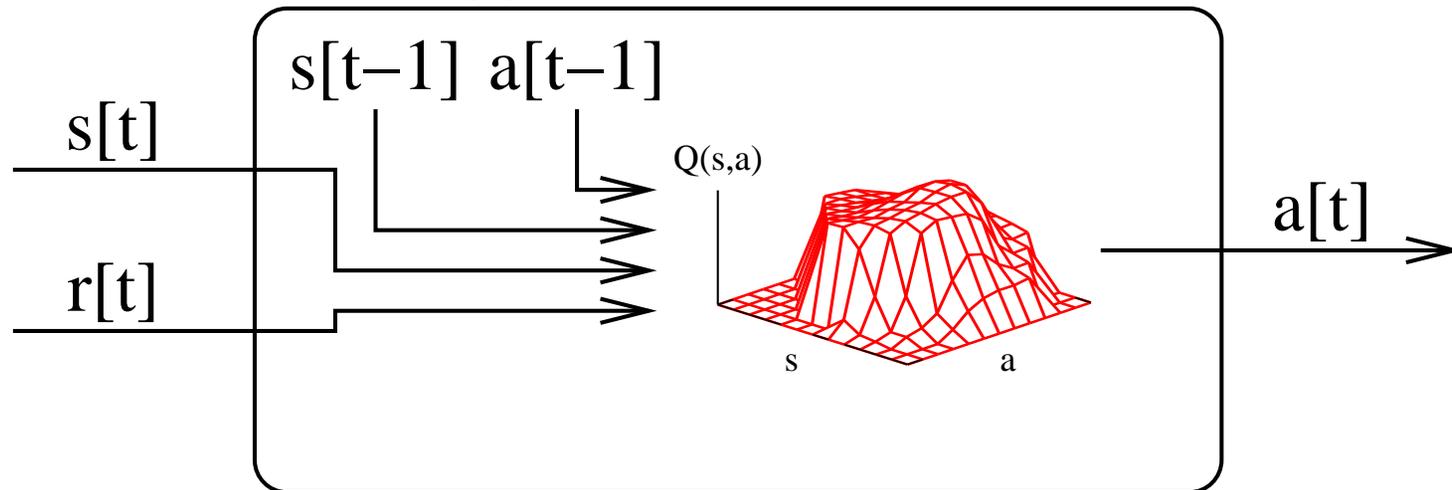


# Function Approximation

In practice, visiting **every state** impossible



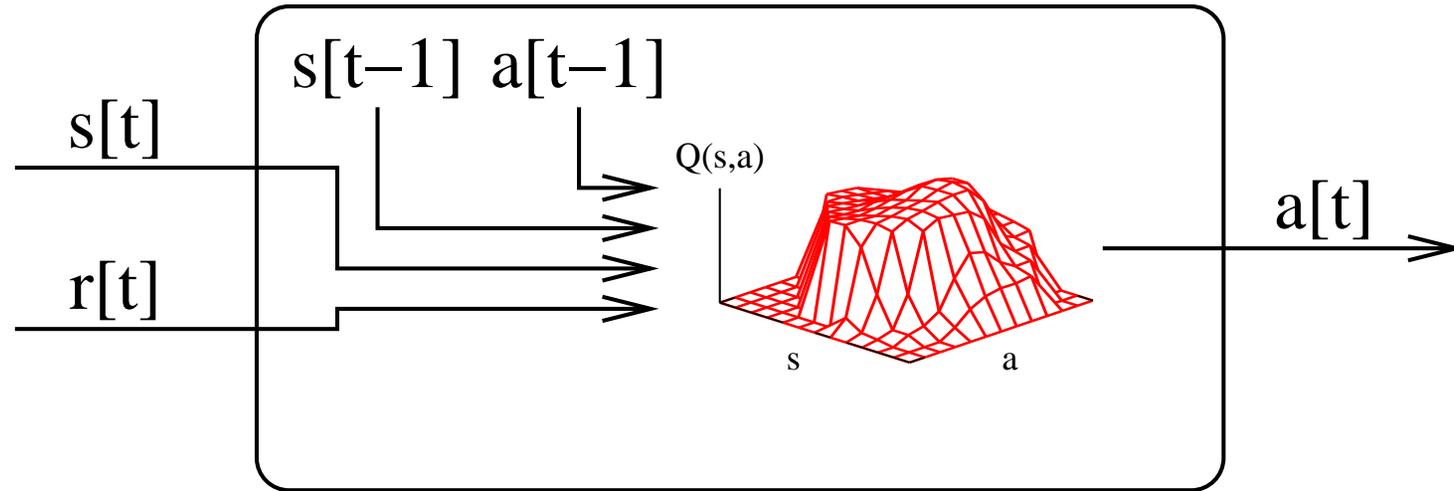
Function approximation of value function



Theoretical **guarantees** harder to come by

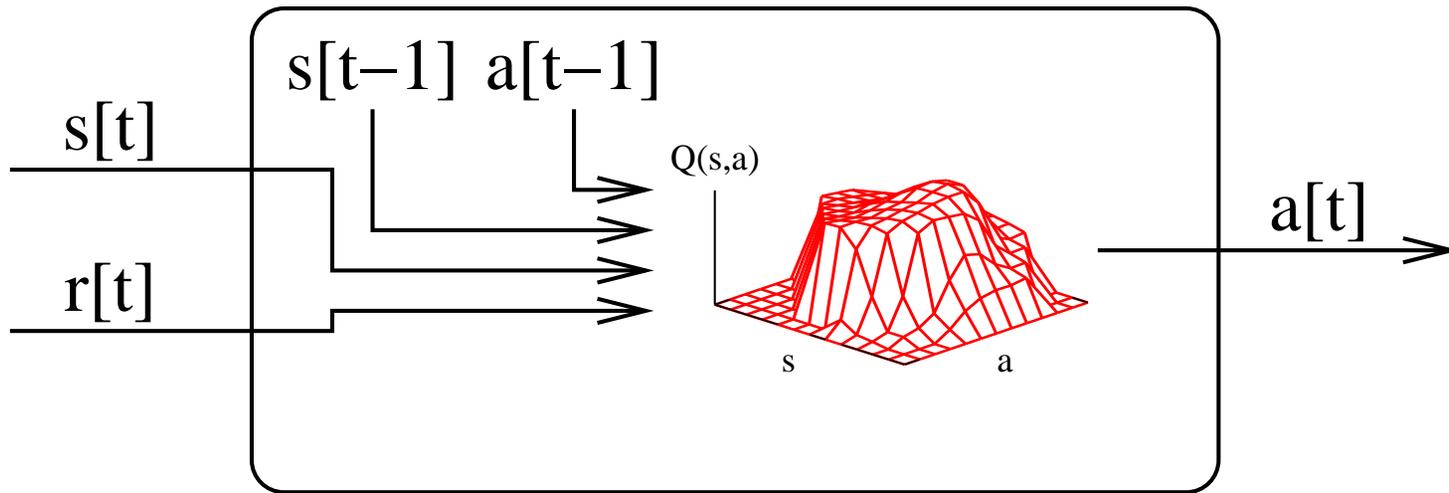
# Batch Methods

In practice, often **experience** is scarce

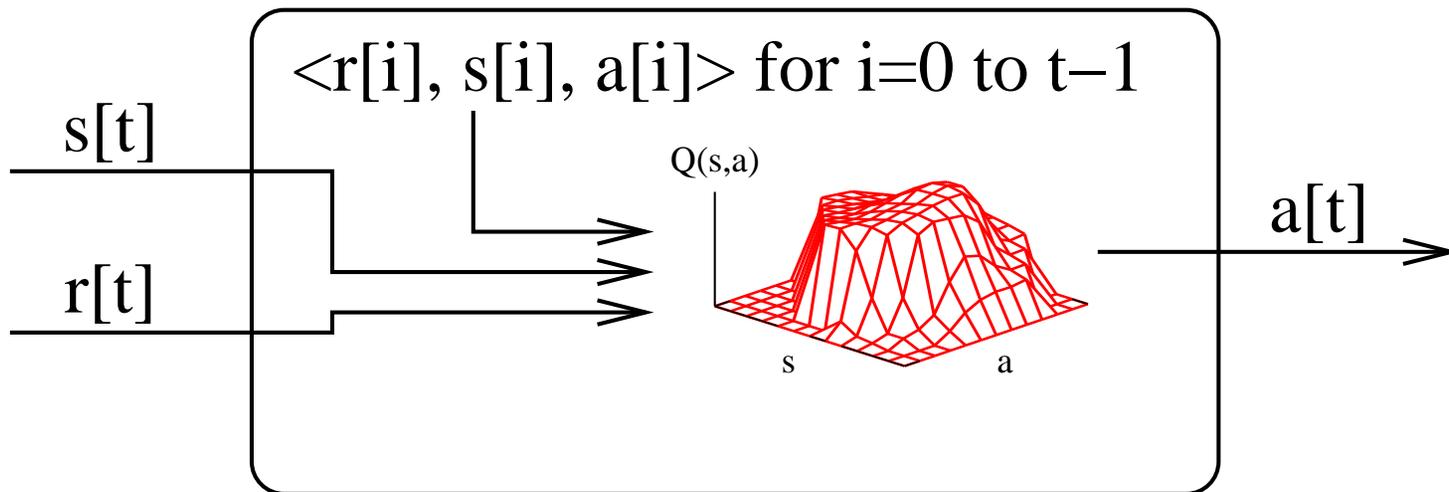


# Batch Methods

In practice, often **experience** is scarce



**Save** transitions:



# Applications: Towards a Useful Tool

- Backgammon [Tesauro, '94]
- Helicopter control [Ng et al., '03]

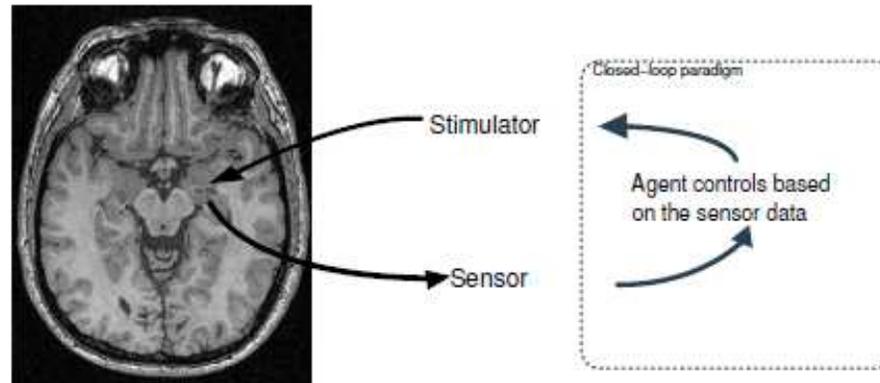


# Applications: Towards a Useful Tool

- Backgammon [Tesauro, '94]
- Helicopter control [Ng et al., '03]



- Adaptive treatment of epilepsy [Pineau et al., '08]



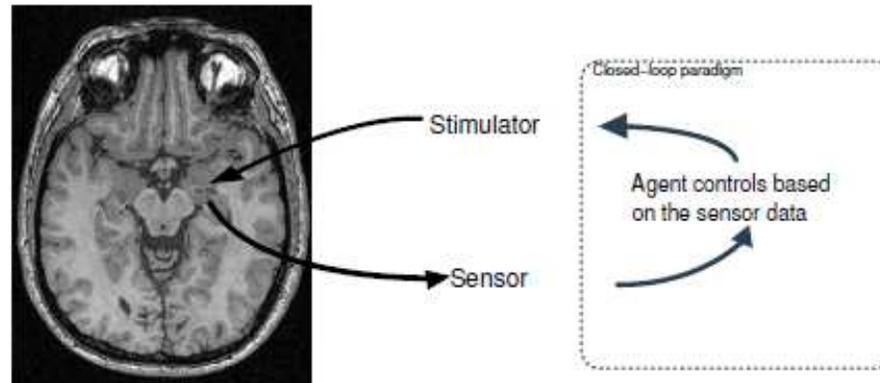
- Invasive species management, wildfire suppression [Dietterich et al., '13]

# Applications: Towards a Useful Tool

- Backgammon [Tesauro, '94]
- Helicopter control [Ng et al., '03]



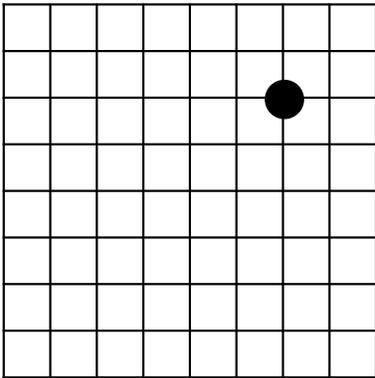
- Adaptive treatment of epilepsy [Pineau et al., '08]



- Invasive species management, wildfire suppression [Dietterich et al., '13]
- Google DeepMind beats human go champion, [Silver et al., '16]

# Computer Go AI – Definition

$d = 1$



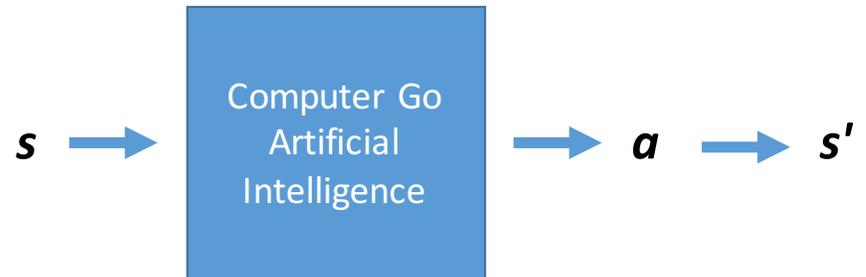
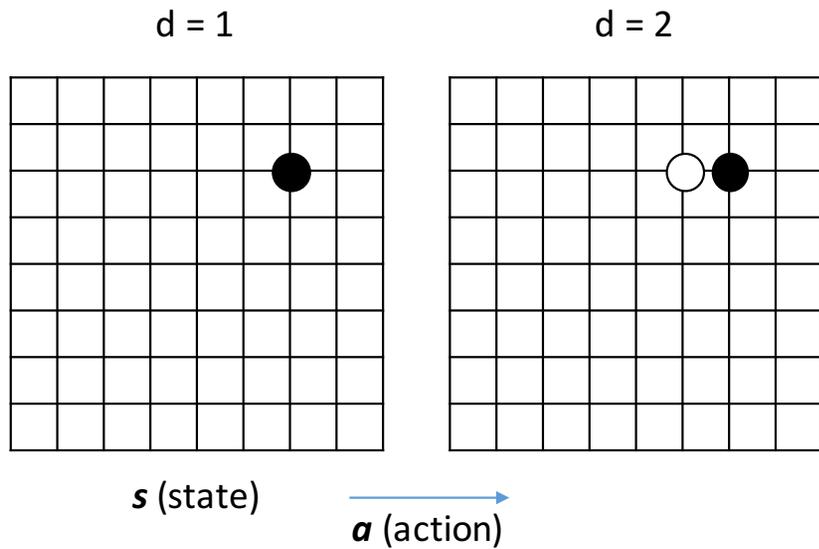
$s$  (state)

=

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

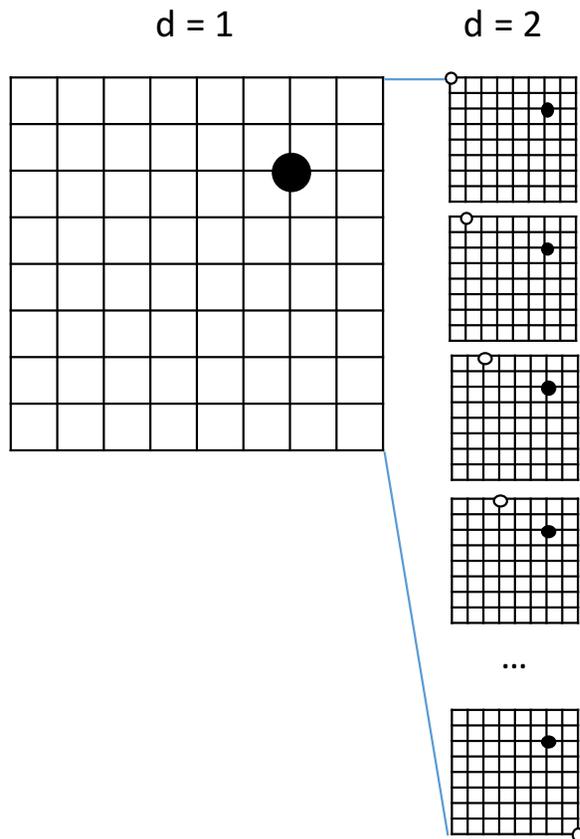
(e.g. we can represent the board into a matrix-like form)

# Computer Go AI – Definition



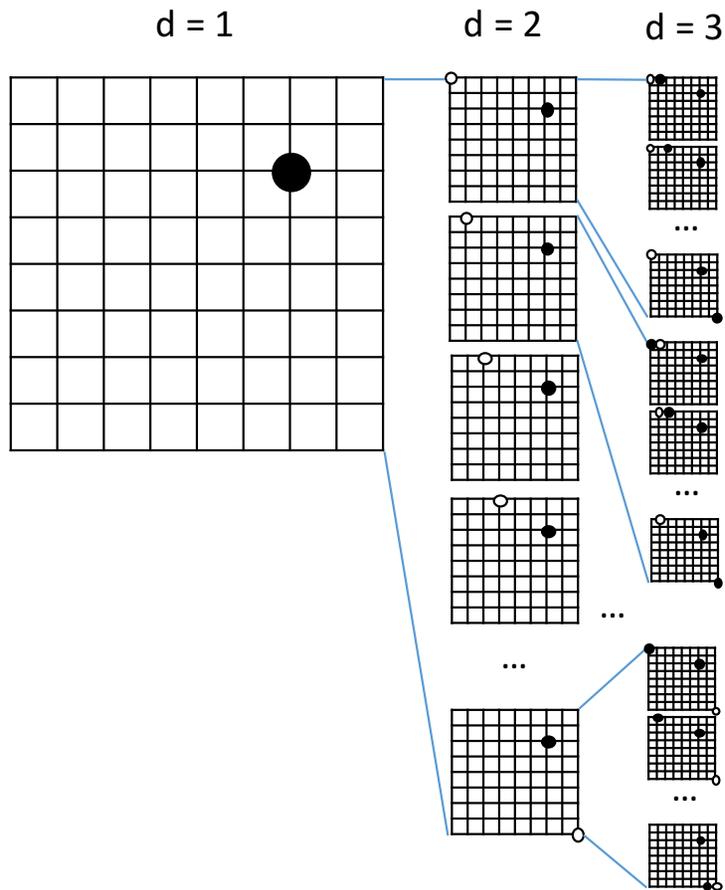
Given  $s$ , pick the best  $a$

# Computer Go AI – An Implementation Idea?

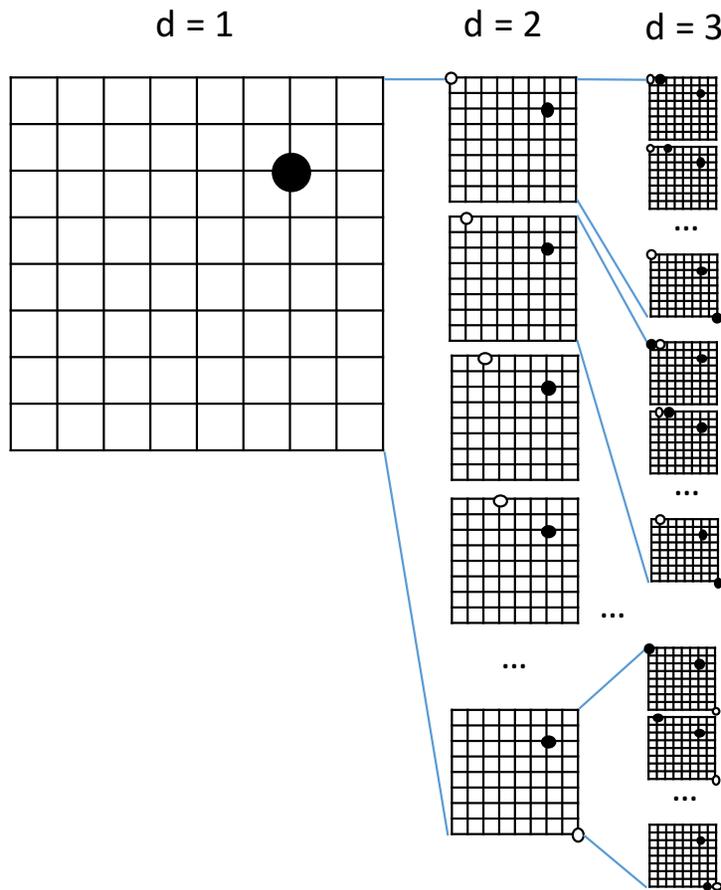


How about simulating all possible board positions?

# Computer Go AI – An Implementation Idea?

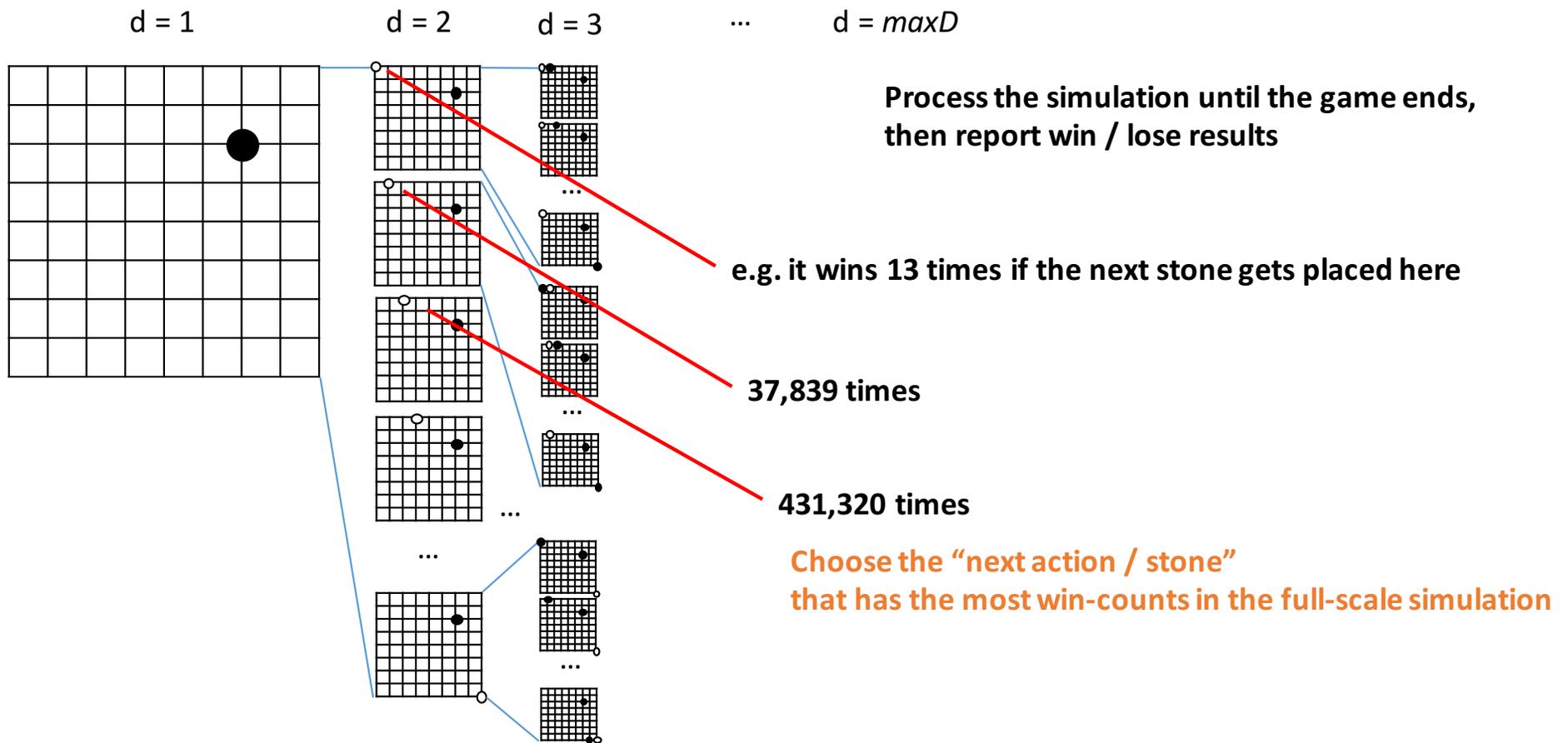


# Computer Go AI – An Implementation Idea?

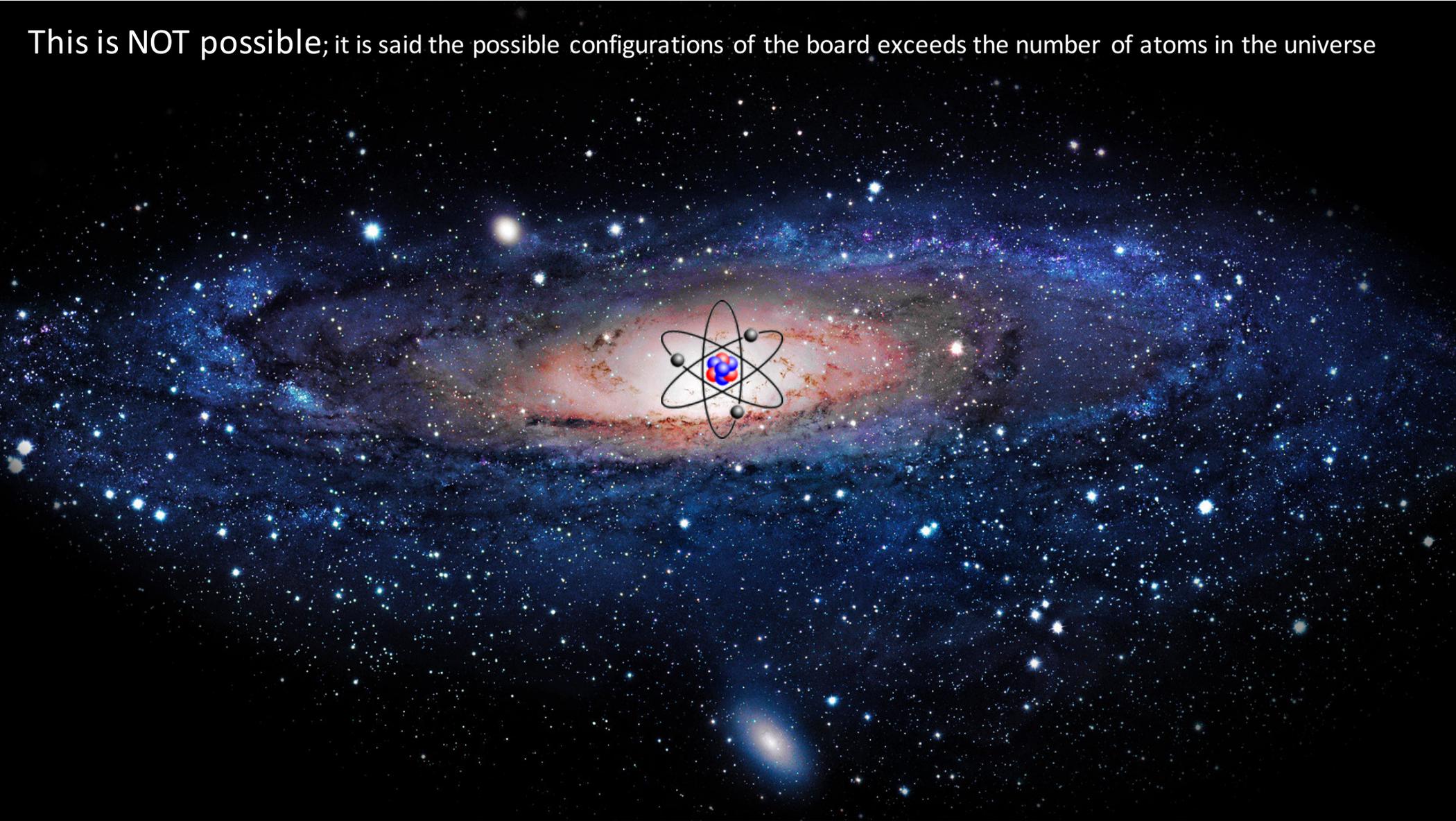


Process the simulation until the game ends,  
then report win / lose results

# Computer Go AI – An Implementation Idea?



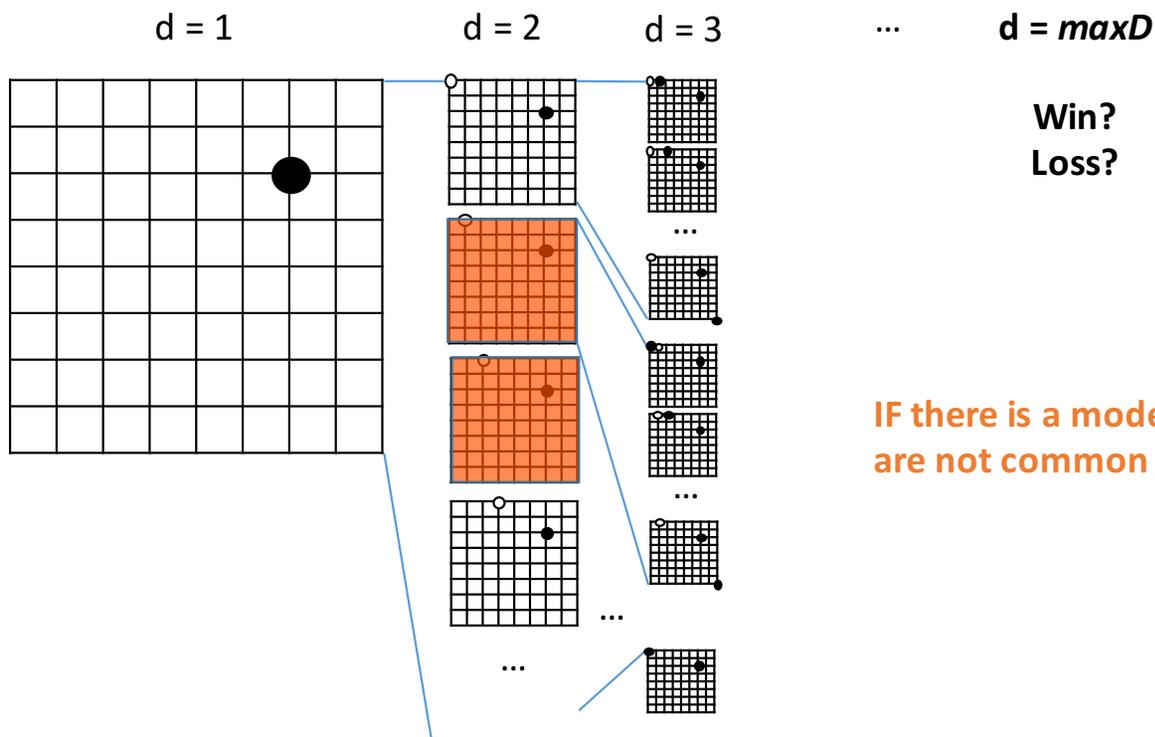
This is NOT possible; it is said the possible configurations of the board exceeds the number of atoms in the universe



**Key: To Reduce Search Space**

# Reducing Search Space

## 1. Reducing “action candidates” (Breadth Reduction)

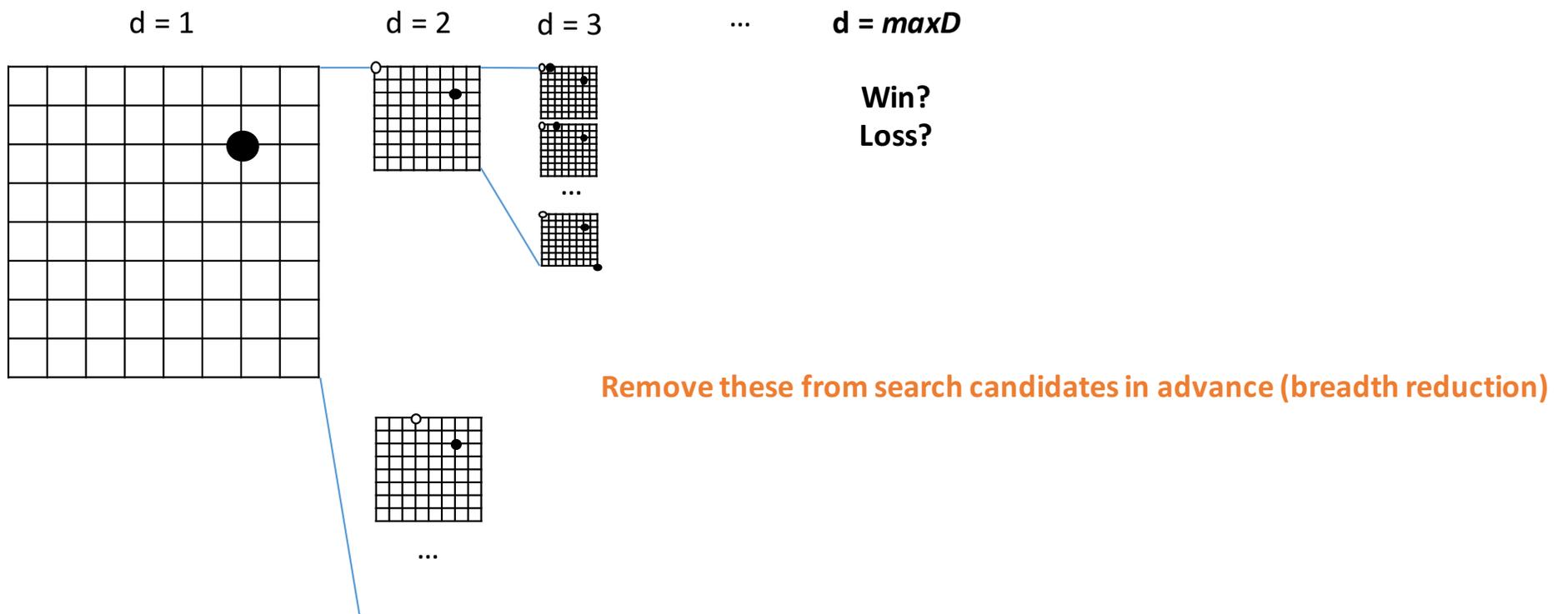


Win?  
Loss?

IF there is a model that can tell you that these moves are not common / probable (e.g. by experts, etc.) ...

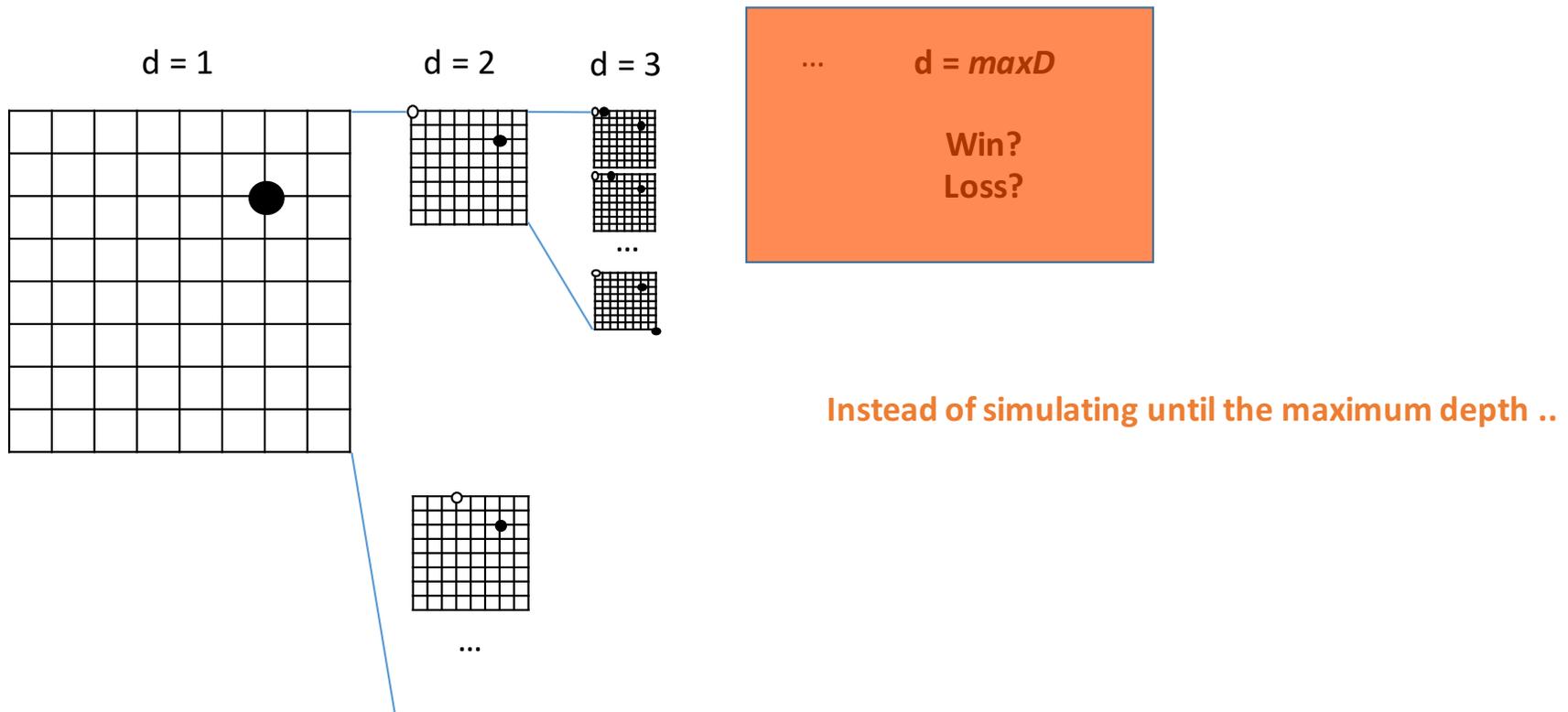
# Reducing Search Space

## 1. Reducing “action candidates” (Breadth Reduction)



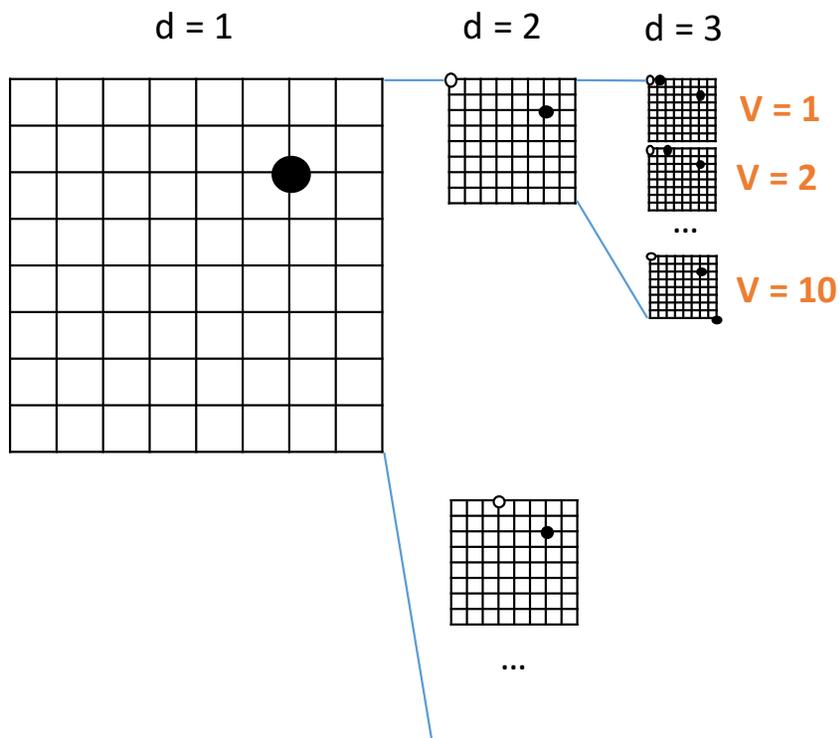
# Reducing Search Space

## 2. Position evaluation ahead of time (Depth Reduction)



# Reducing Search Space

## 2. Position evaluation ahead of time (Depth Reduction)



**IF there is a function that can measure:**  
 $V(s)$ : “board evaluation of state  $s$ ”



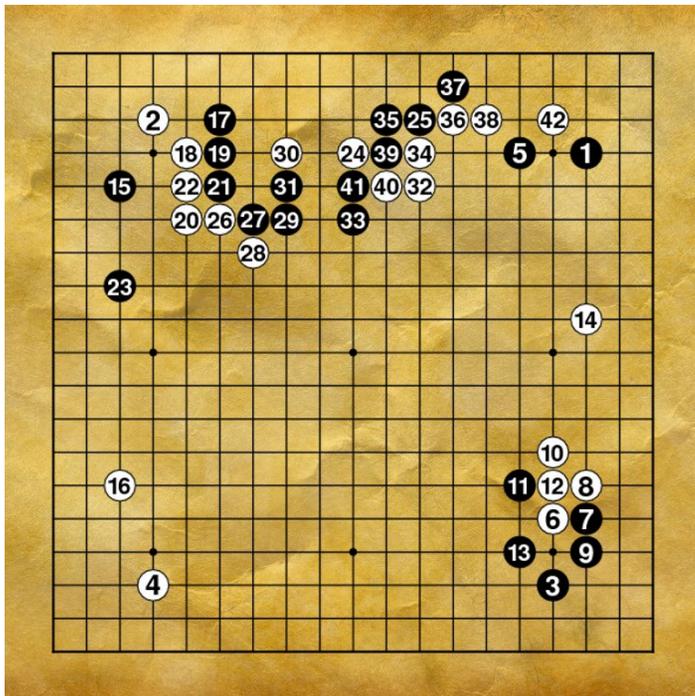
# 1. Reducing “action candidates”

Learning: **P ( next action | current state )**

$$= P ( a | s )$$

# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)



Current State

Next State

s1

s2

s2

s3

s3

s4

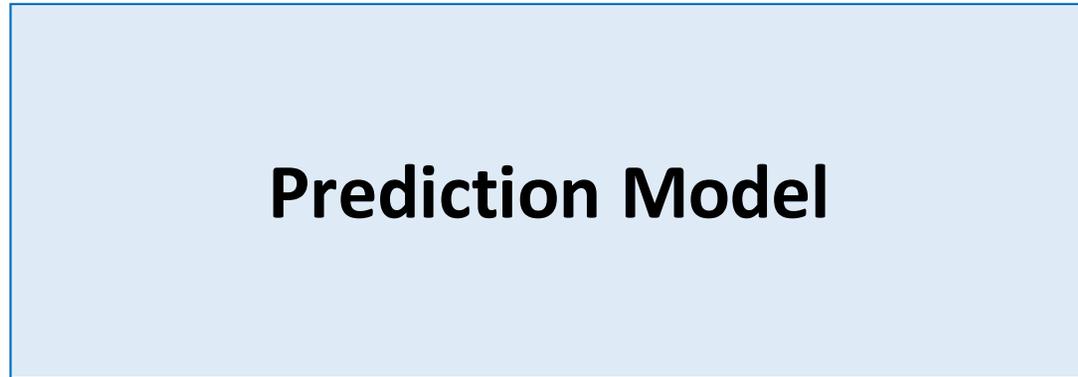
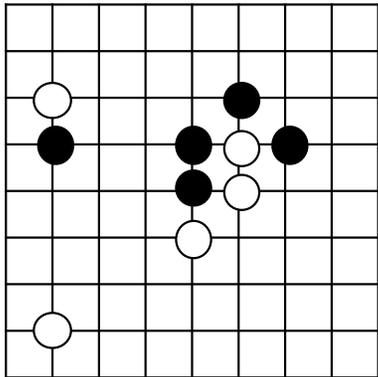
**Prediction  
Model**

**Data:** Online Go experts (5~9 dan)  
160K games, 30M board positions

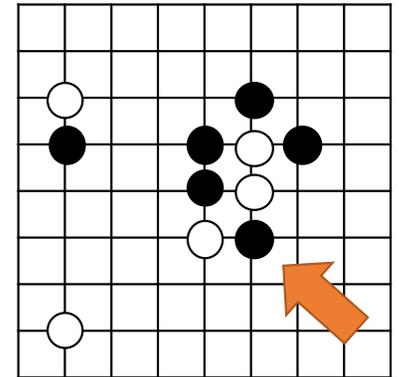
# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board



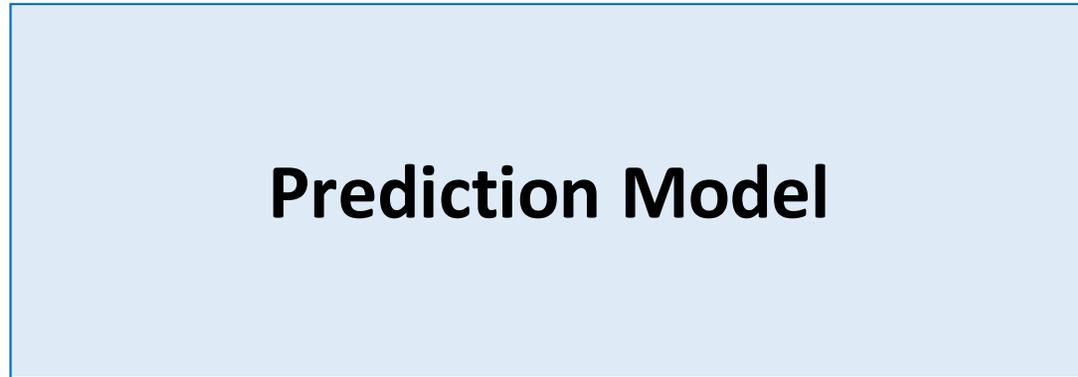
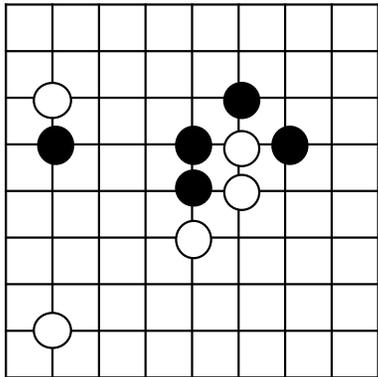
Next Board



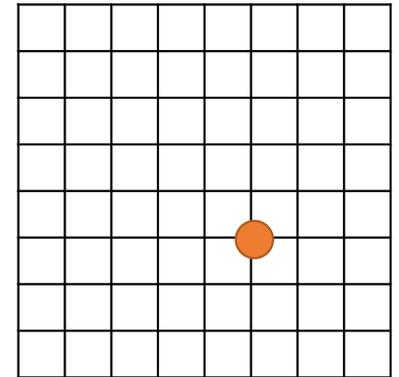
# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board



Next Action



There are  $19 \times 19 = 361$   
possible actions  
(with different probabilities)

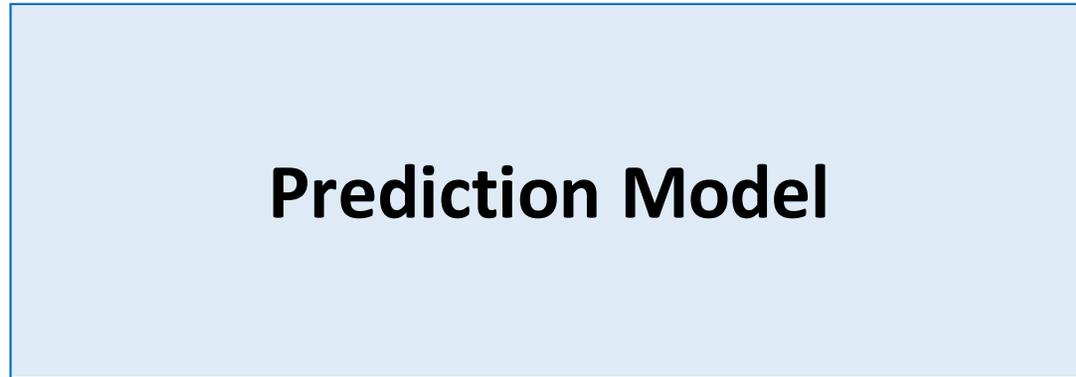
# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board

```
00 000 0000  
00 000 1000  
0-100 1-1100  
0 1 00 1-1000  
00 00-10000  
00 000 0000  
0-10000 0000  
00 000 0000
```

$s$



Next Action

```
0000000000  
0000000000  
0000000000  
0000000000  
0000000000  
0000010000  
0000000000  
0000000000  
0000000000
```

$a$

$f: s \rightarrow a$

# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board

```
00 000 0000  
00 000 1000  
0-100 1-1100  
0 1 00 1-1000  
00 00-10000  
00 000 0000  
0-10000 0000  
00 000 0000
```

$s$

**Prediction  
Model**

$g: s \rightarrow p(a|s)$

```
000000 000  
000000 000  
000000 000  
000000.2 0.100  
00000 0.4 0.200  
000000.1 000  
000000 000  
000000 000
```

$p(a|s)$

**argmax**

Next Action

```
000000000  
000000000  
000000000  
000000000  
000000000  
00000 1 000  
000000000  
000000000  
000000000
```

$a$

# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

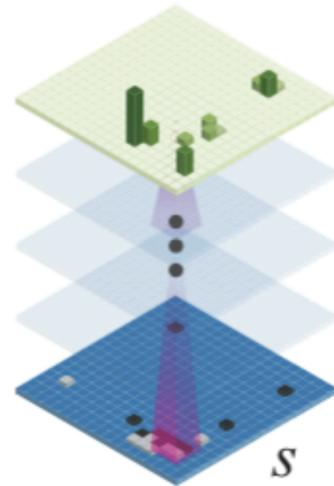
Current Board

```
00 000 0000  
00 000 1000  
0 -100 1-1100  
0 1 00 1-1000  
00 00 -10000  
00 000 0000  
0 -1000 0000  
00 000 0000
```

$s$

**Prediction  
Model**

$g: s \rightarrow p(a|s)$



$p(a|s)$

argmax

$a$

Next Action

```
000000000  
000000000  
000000000  
000000000  
000000000  
000001000  
000000000  
000000000  
000000000
```

# 1. Reducing “action candidates”

(1) Imitating expert moves (supervised learning)

Current Board

```
00 000 0000  
00 000 1000  
0-100 1-1100  
0 1 00 1-1000  
00 00-10000  
00 000 0000  
0-10000 0000  
00 000 0000
```

$s$

**Deep Learning  
(13 Layer CNN)**

$g: s \rightarrow p(a|s)$

```
000000 000  
000000 000  
000000 000  
000000.2 0.100  
00000 0.4 0.200  
000000.1 000  
000000 000  
000000 000
```

$p(a|s)$

**argmax**

Next Action

```
000000000  
000000000  
000000000  
000000000  
000000000  
00000 1 000  
000000000  
000000000  
000000000
```

$a$

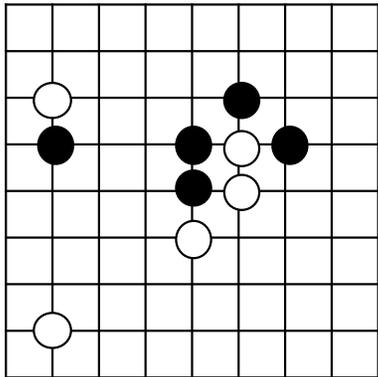
**Go:** abstraction is the key to win

**CNN:** abstraction is its *forte*

# 1. Reducing “action candidates”

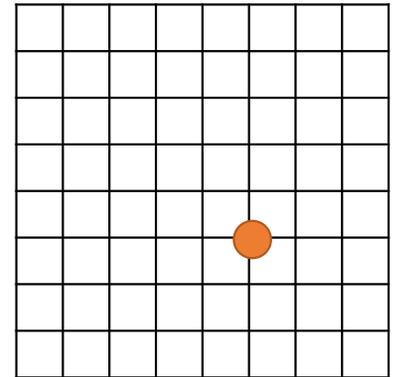
(1) Imitating expert moves (supervised learning)

Current Board



**Expert Moves Imitator Model  
(w/ CNN)**

Next Action



**Training:**  $\Delta\sigma \propto \frac{\partial \log p_\sigma(a|s)}{\partial \sigma}$

# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Improving by playing against itself

**Expert Moves  
Imitator Model  
(w/ CNN)**

vs

**Expert Moves  
Imitator Model  
(w/ CNN)**



# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

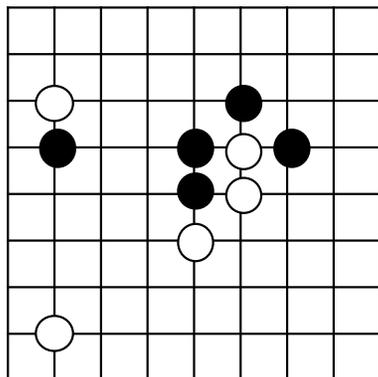


**Return:** board positions, win/lose info

# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Board position



**Expert Moves Imitator Model  
(w/ CNN)**

win/loss

**Loss**

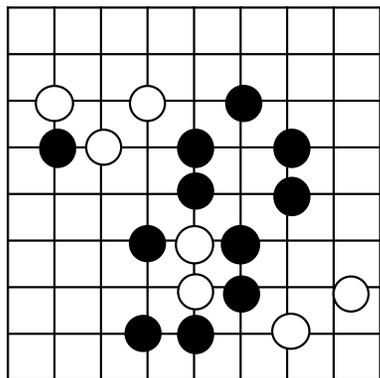
$z = -1$

**Training:**  $\Delta\rho \propto \frac{\partial \log p_\rho(a_t|s_t)}{\partial \rho} z_t$ .

# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Board position



win/loss

**Expert Moves Imitator Model  
(w/ CNN)**

**Win**

$z = +1$

**Training:**  $\Delta\rho \propto \frac{\partial \log p_\rho(a_t|s_t)}{\partial \rho} z_t$ .

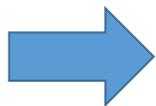
# 1. Reducing “action candidates”

(2) Improving through self-plays (reinforcement learning)

Older models vs. newer models



It uses the same topology as the expert moves imitator model, and just uses the updated parameters



**Return:** board positions, win/lose info

# 1. Reducing “action candidates”

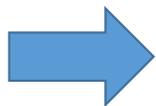
(2) Improving through self-plays (reinforcement learning)



**Return:** board positions, win/lose info

# 1. Reducing “action candidates”

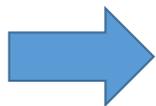
(2) Improving through self-plays (reinforcement learning)



**Return:** board positions, win/lose info

# 1. Reducing “action candidates”

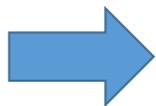
(2) Improving through self-plays (reinforcement learning)



**Return:** board positions, win/lose info

# 1. Reducing “action candidates”

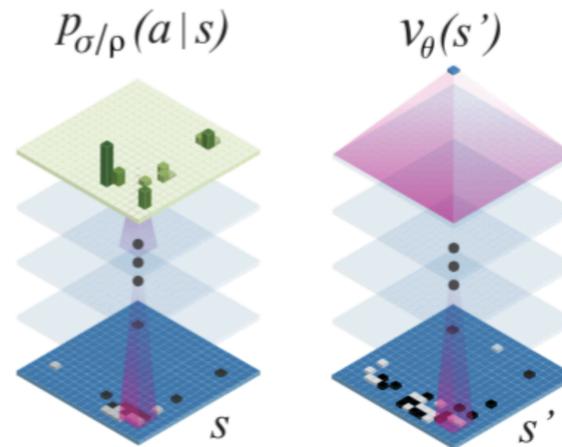
(2) Improving through self-plays (reinforcement learning)



**The final model wins 80% of the time  
when playing against the first model**

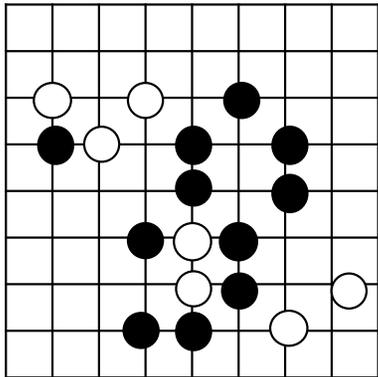
## 2. Board Evaluation

## 2. Board Evaluation



Adds a regression layer to the model  
Predicts values between 0~1  
Close to 1: a good board position  
Close to 0: a bad board position

Board Position



**Updated Model  
ver 1,000,000**

**Value  
Prediction  
Model  
(Regression)**

Win / Loss

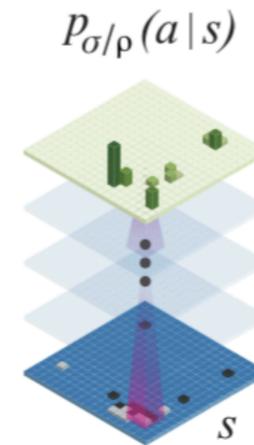
**Win  
(0~1)**

**Training:** 
$$\Delta\theta \propto \frac{\partial v_{\theta}(s)}{\partial\theta} (z - v_{\theta}(s))$$

# Reducing Search Space

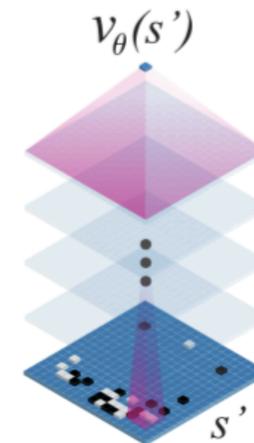
1. Reducing “action candidates”  
(Breadth Reduction)

Policy Network

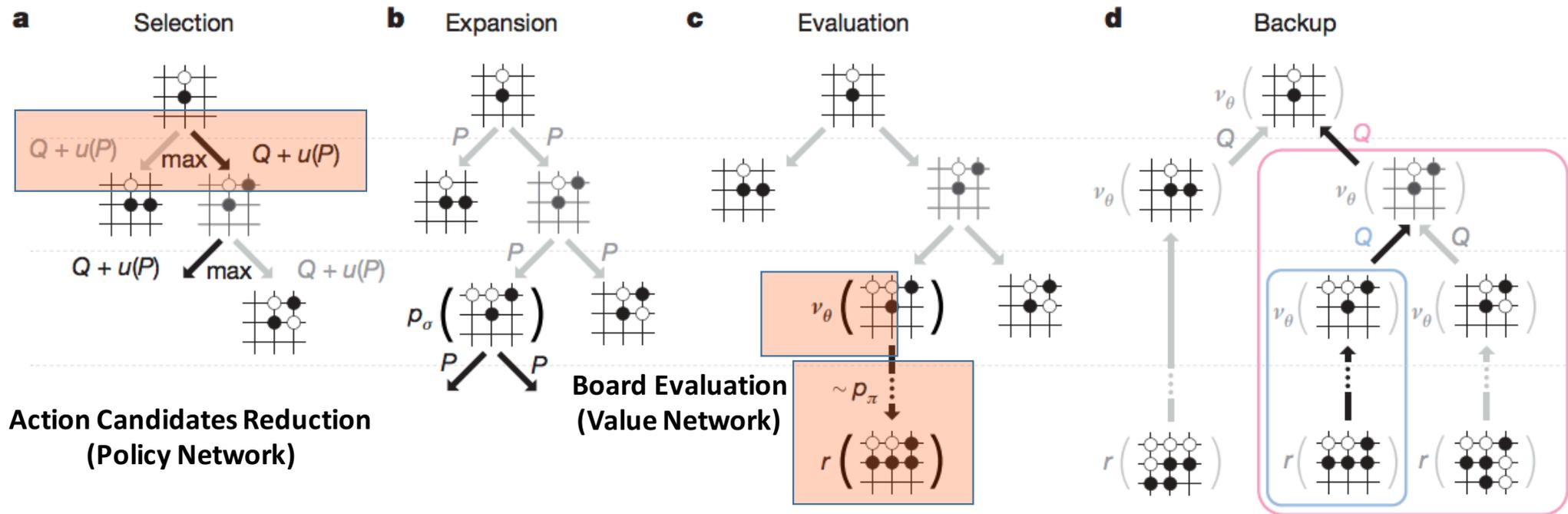


2. Board Evaluation (Depth Reduction)

Value Network

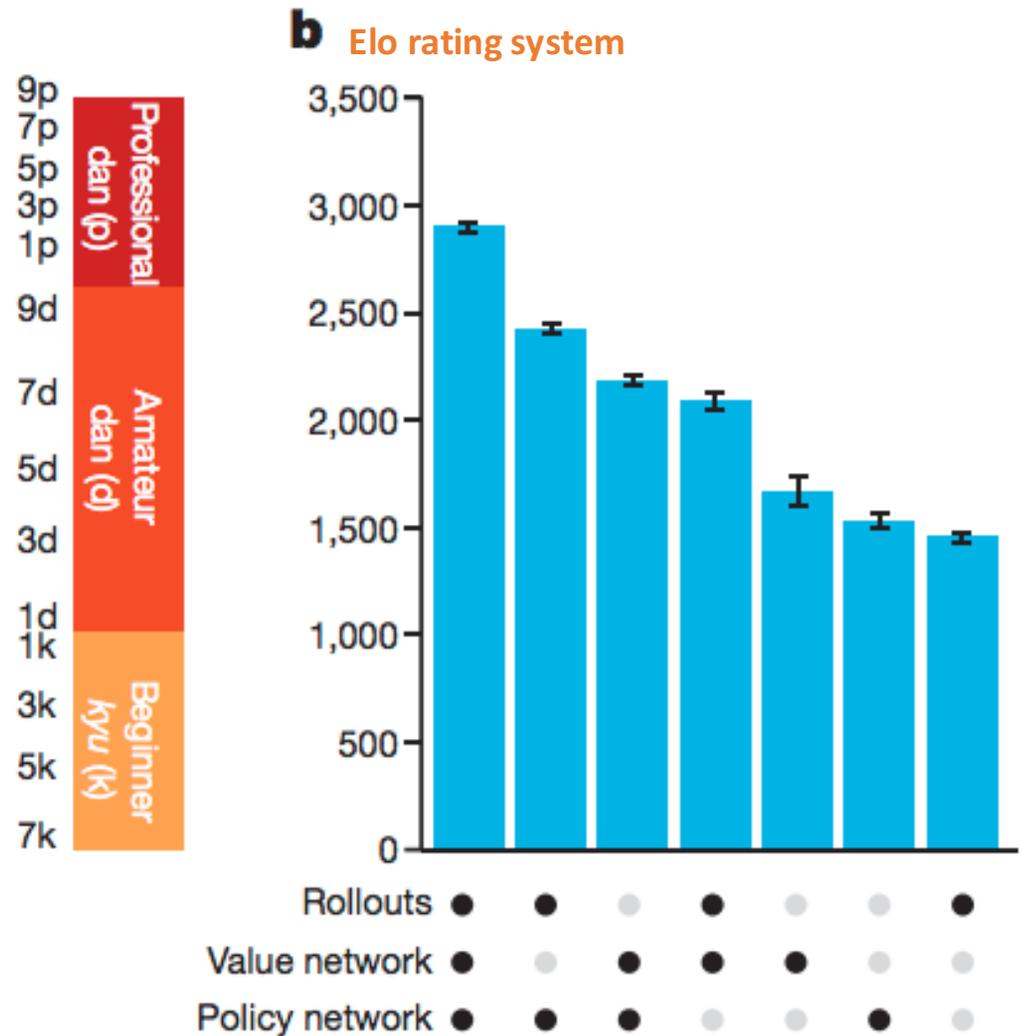


# Looking ahead (w/ Monte Carlo Search Tree)



(Rollout): Faster version of estimating  $p(a | s)$   
 $\rightarrow$  uses shallow networks ( $3 \text{ ms} \rightarrow 2 \mu\text{s}$ )

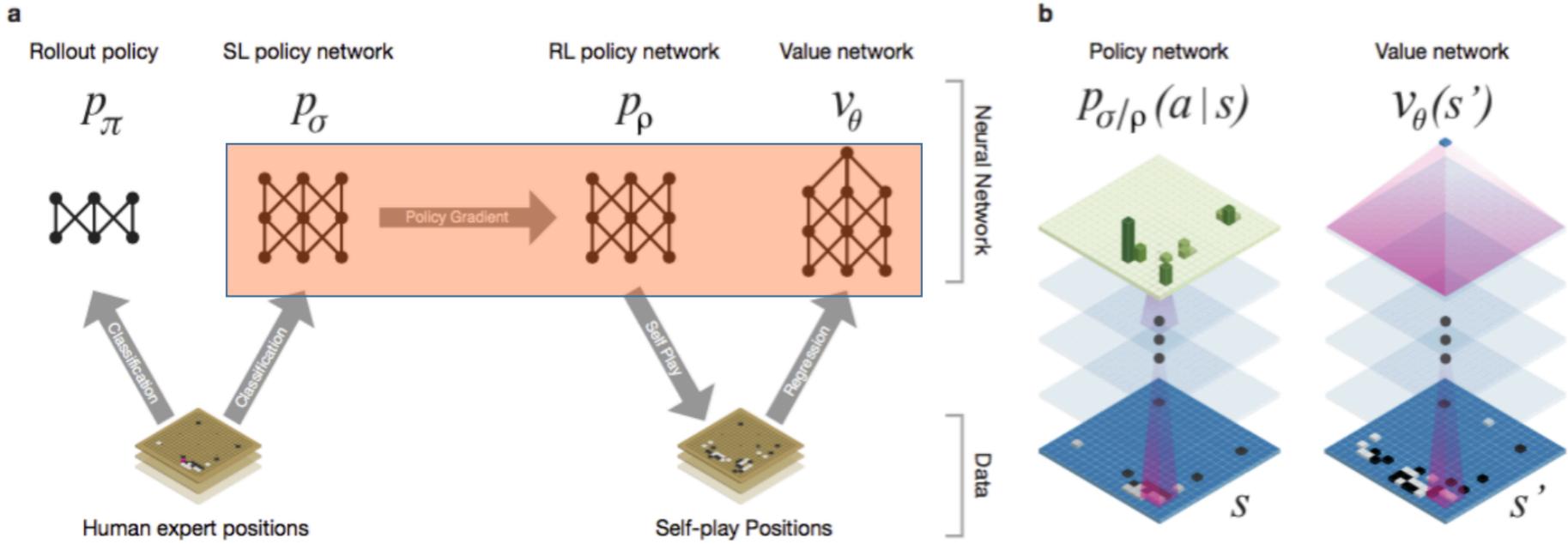
# Results



Performance with different combinations of AlphaGo components

# Takeaways

Use the networks trained for a certain task (with different loss objectives) for several other tasks



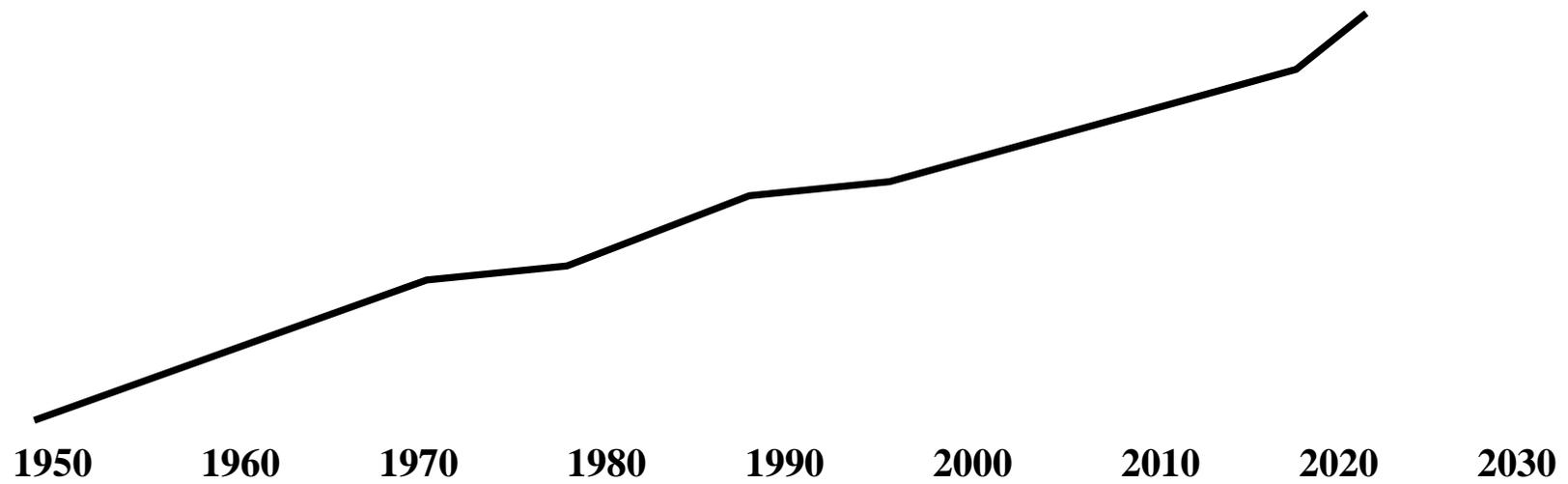
# AlphaGo [Silver et al., '15]: An AI milestone



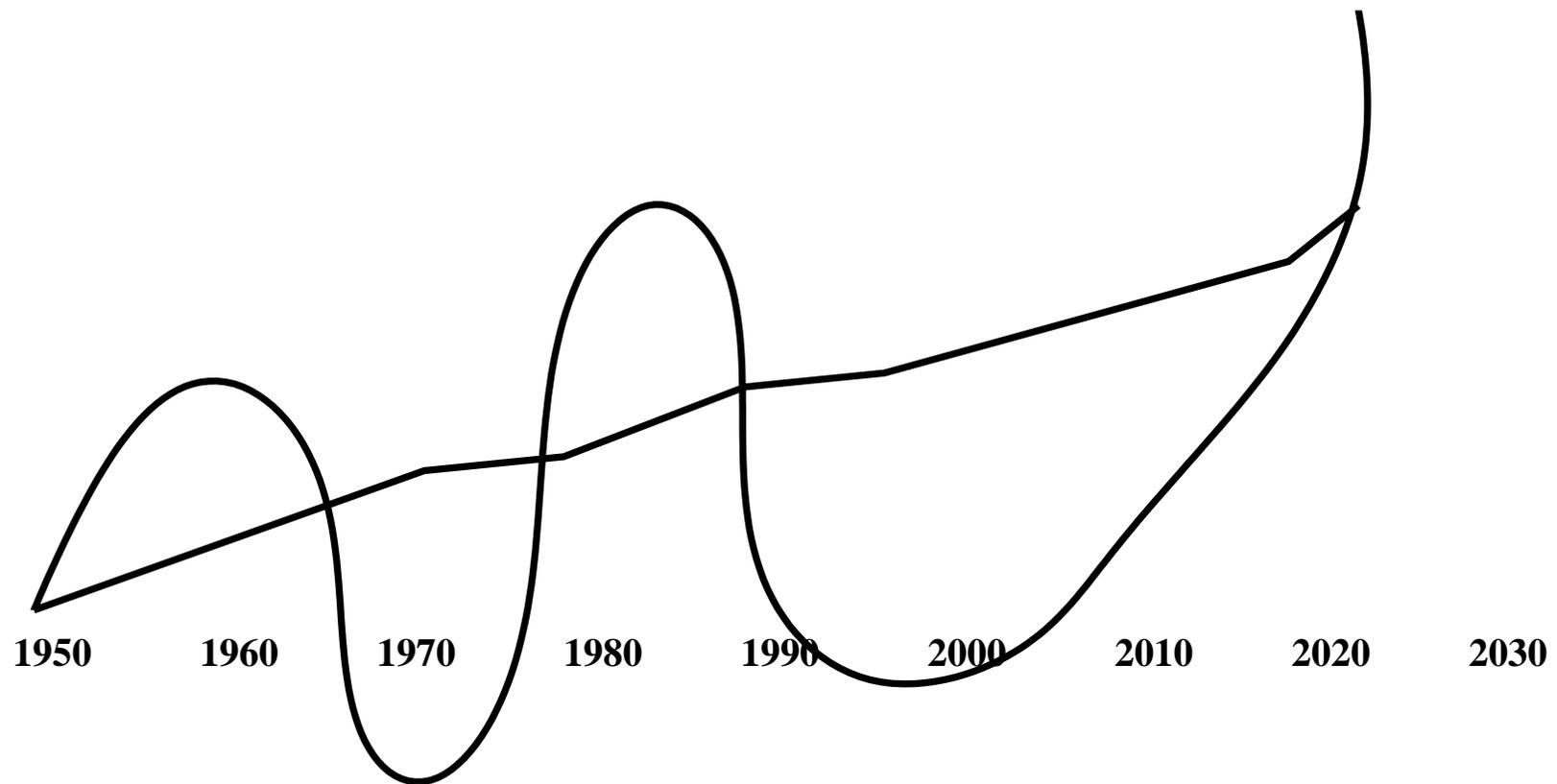
- Integrated two **existing** AI technologies
  - Supervised learning (Deep learning)
  - **Reinforcement learning**
- Does **not** solve AI
  - The end of the line for:

**2-player zero-sum discrete finite deterministic games of perfect information**

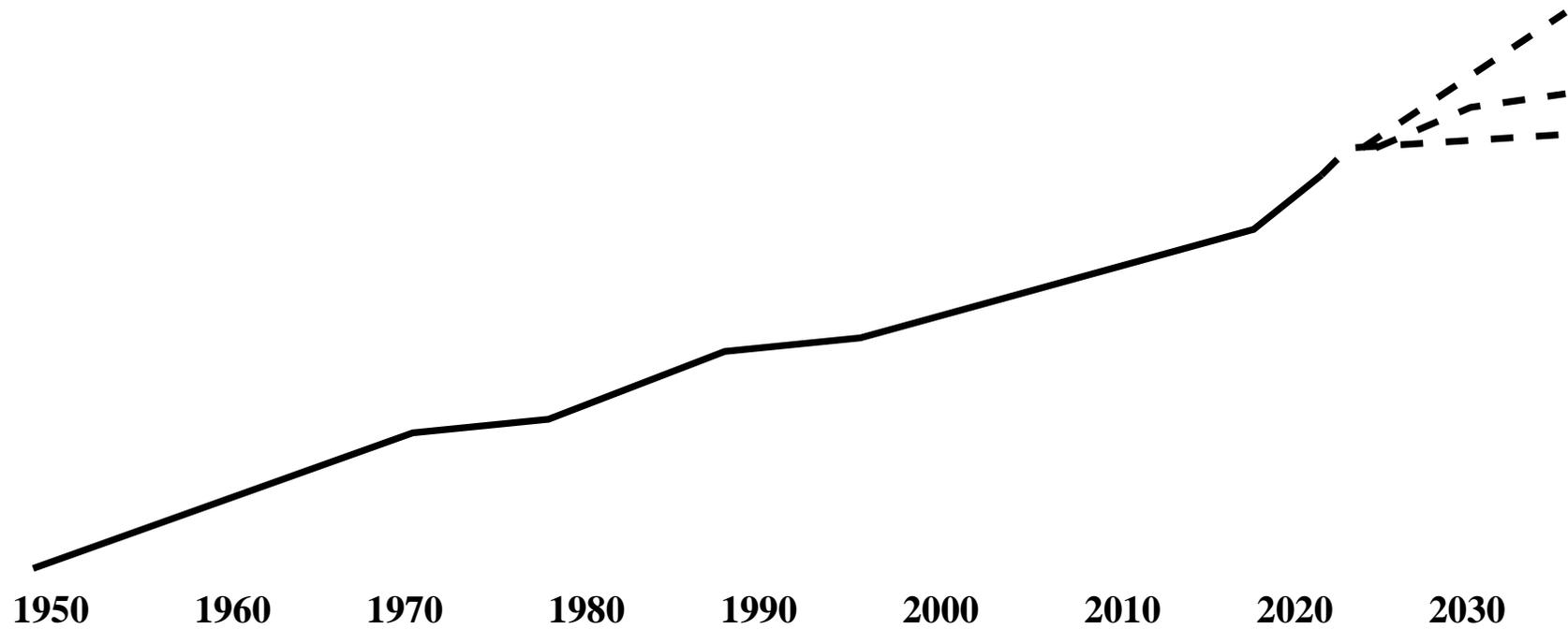
# Reality



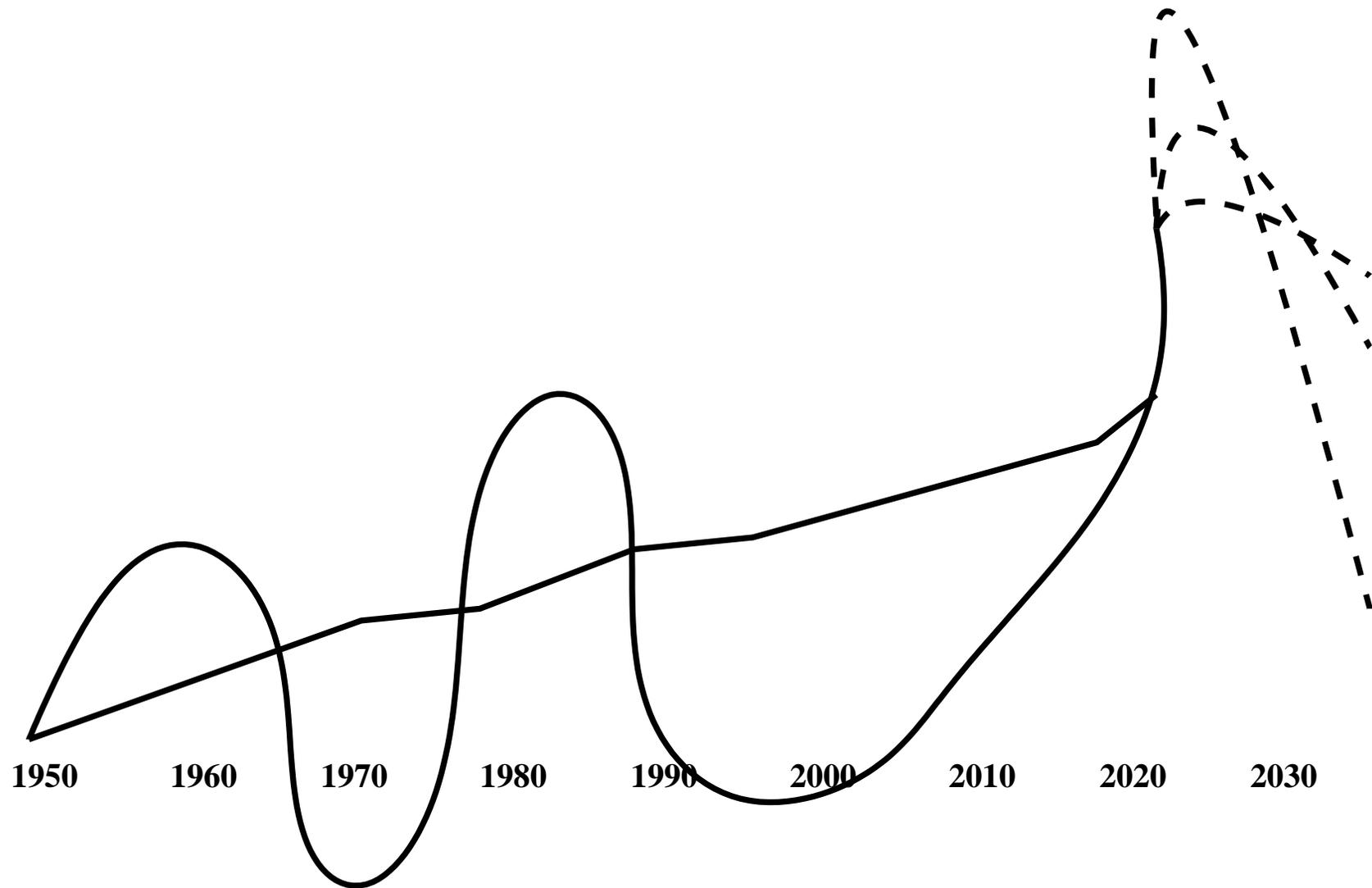
# Perceptions



# Uncertainty



# Perception Uncertainty



# Selected RL Contributions

- Human interaction

# Selected RL Contributions

- Human interaction
  - Advice, **Demonstration**



# Selected RL Contributions

- Human interaction
  - Advice, **Demonstration**
  - Positive/Negative **Feedback**



[Knox & Stone, '09]

# Selected RL Contributions

- Human interaction
  - Advice, **Demonstration**
  - Positive/Negative **Feedback**
- Transfer learning for RL



[Knox & Stone, '09]  
[Taylor & Stone, '07]

# Selected RL Contributions

- Human interaction
  - Advice, **Demonstration**
  - Positive/Negative **Feedback**



[Knox & Stone, '09]

[Taylor & Stone, '07]

- Transfer learning for RL
- Adaptive/hierarchical representations

[Whiteson & Stone, '05], [Jong & Stone, '08]

# Selected RL Contributions

- Human interaction
  - Advice, **Demonstration**
  - Positive/Negative **Feedback**



- Transfer learning for RL

[Knox & Stone, '09]  
[Taylor & Stone, '07]

- Adaptive/hierarchical representations

[Whiteson & Stone, '05], [Jong & Stone, '08]

- **TEXPLORE** for Robot RL

[Hester & Stone, '13]

# Selected RL Contributions

- Human interaction

- Advice, **Demonstration**
- Positive/Negative **Feedback**



[Knox & Stone, '09]

[Taylor & Stone, '07]

- Transfer learning for RL

- Adaptive/hierarchical representations

[Whiteson & Stone, '05], [Jong & Stone, '08]

- **TEXPLORE** for Robot RL

[Hester & Stone, '13]

- Sample efficient; real-time

# Selected RL Contributions

- Human interaction

- Advice, **Demonstration**
- Positive/Negative **Feedback**



[Knox & Stone, '09]

[Taylor & Stone, '07]

- Transfer learning for RL

- Adaptive/hierarchical representations

[Whiteson & Stone, '05], [Jong & Stone, '08]

- **TEXPLORE** for Robot RL

[Hester & Stone, '13]

- Sample efficient; real-time
- Continuous state; delayed effects

# UT Austin Villa 2014 RoboCup 3D Simulation League Champion via Overlapping Layered Learning

Patrick MacAlpine, Mike Depinet, and Peter Stone



# Layered Learning

- For domains too **complex** for tractably mapping state features  $S$   
 $\mapsto$  outputs  $O$
- Hierarchical subtask decomposition **given**:  $\{L_1, L_2, \dots, L_n\}$

# Layered Learning

- For domains too **complex** for tractably mapping state features  $S$   $\mapsto$  outputs  $O$
- Hierarchical subtask decomposition **given**:  $\{L_1, L_2, \dots, L_n\}$
- Machine learning: **exploit data** to train, adapt

# Layered Learning

- For domains too **complex** for tractably mapping state features  $S$   $\mapsto$  outputs  $O$
- Hierarchical subtask decomposition **given**:  $\{L_1, L_2, \dots, L_n\}$
- Machine learning: **exploit data** to train, adapt
- **Synthesis**: Learning in one layer feeds into next layer

# Layered Learning

- For domains too **complex** for tractably mapping state features  $S$   $\mapsto$  outputs  $O$
- Hierarchical subtask decomposition **given**:  $\{L_1, L_2, \dots, L_n\}$
- Machine learning: **exploit data** to train, adapt
- **Synthesis**: Learning in one layer feeds into next layer



# Layered Learning in Practice

First applied in **simulated** robot soccer [Stone & Veloso, '97]

	Strategic Level	Example
$L_1$	individual	ball interception
$L_2$	multiagent	pass evaluation
$L_3$	<b>team</b>	pass selection

# Layered Learning in Practice

First applied in **simulated** robot soccer [Stone & Veloso, '97]

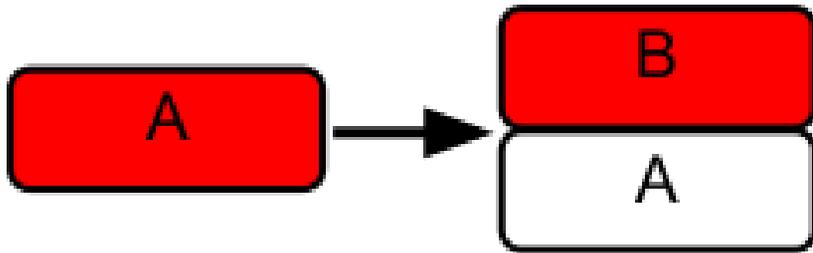
	Strategic Level	Example
$L_1$	individual	ball interception
$L_2$	multiagent	pass evaluation
$L_3$	<b>team</b>	pass selection

Later applied on **real robots** [Stone, Kohl, & Fiedelman, '06]

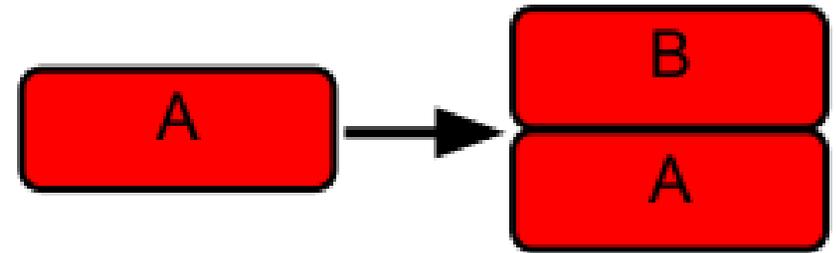
	Strategic Level	Example
$L_1$	individual	fast walking
$L_2$	individual	ball control



# Layered Learning Paradigms



Sequential Layered Learning (SLL)



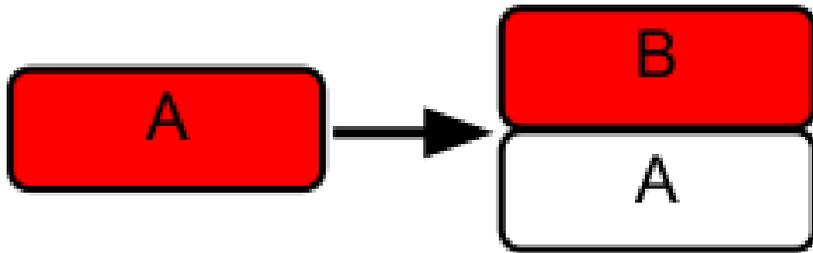
Concurrent Layered Learning (CLL)

## DESCRIPTIONS:

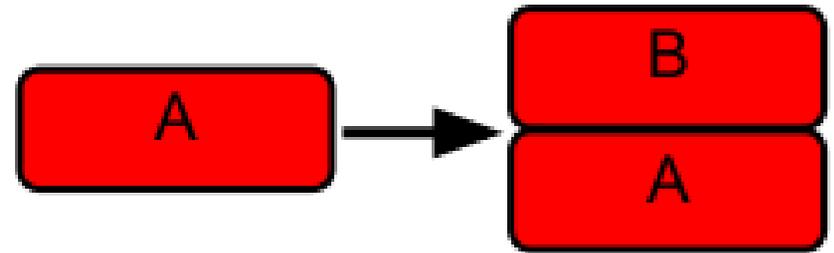
**Sequential Layered Learning:** Freeze parameters of layer after learning before learning of the next layer

**Concurrent Layered Learning:** Keep parameters of layer open during learning of the next layer

# Layered Learning Paradigms



Sequential Layered Learning (SLL)



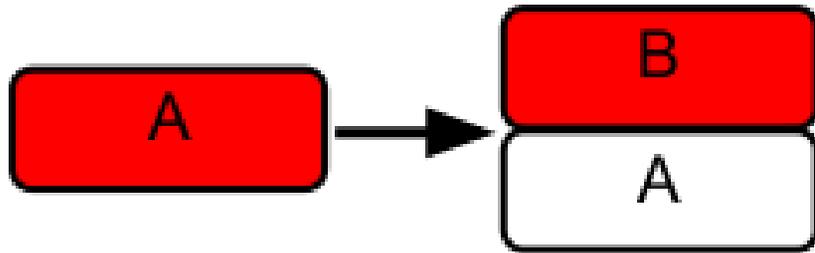
Concurrent Layered Learning (CLL)

## PROBLEMS:

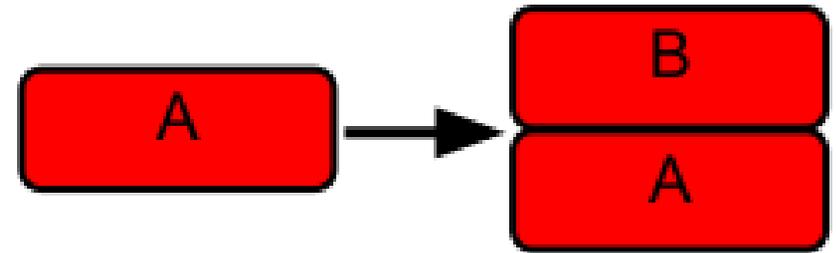
**Sequential Layered Learning:** Can be too **limiting** in the joint layer policy search space

**Concurrent Layered Learning:** The **increased dimensionality** can make learning harder or intractable

# Layered Learning Paradigms



Sequential Layered Learning (SLL)



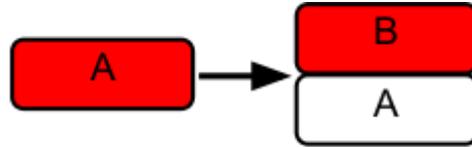
Concurrent Layered Learning (CLL)

## SOLUTION:

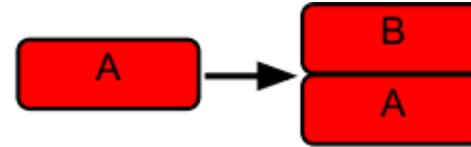
**Overlapping Layered Learning:** Tradeoff between freezing or keeping open previous learned layers

Optimizes “seam” or **overlap** between behaviors: keeps **some** parts of previously learned layers open during subsequent learning

# Overlapping Layered Learning



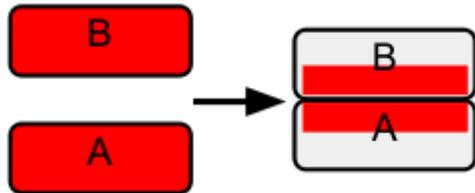
Sequential Layered Learning (SLL)



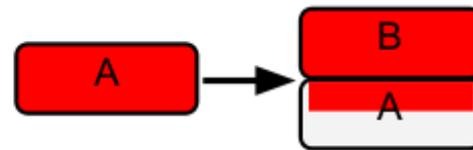
Concurrent Layered Learning (CLL)

---

## Overlapping Layered Learning



Combining Independently Learned Behaviors (CILB)

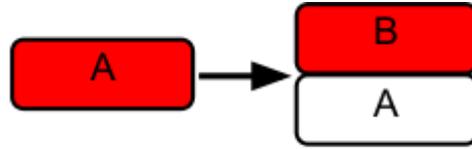


Partial Concurrent Layered Learning (PCLL)

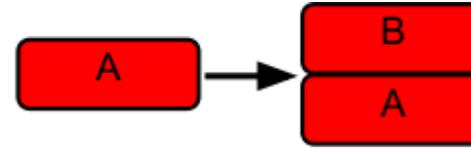


Previous Learned Layer Refinement (PLLR)

# Overlapping Layered Learning



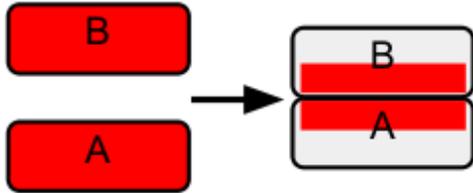
Sequential Layered Learning (SLL)



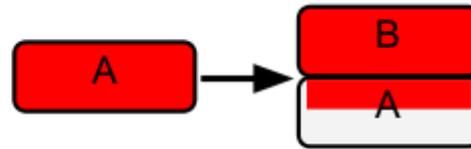
Concurrent Layered Learning (CLL)

---

## Overlapping Layered Learning



Combining Independently Learned Behaviors (CILB)



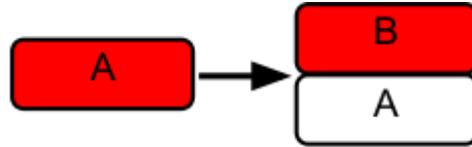
Partial Concurrent Layered Learning (PCLL)



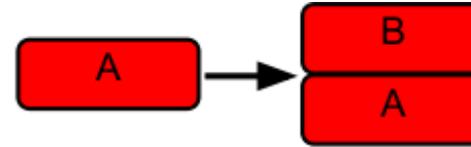
Previous Learned Layer Refinement (PLLR)

**Combining Independently Learned Behaviors:** Behaviors learned independently and then combined by relearning subset of behaviors' parameters

# Overlapping Layered Learning



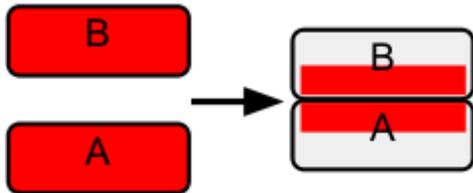
Sequential Layered Learning (SLL)



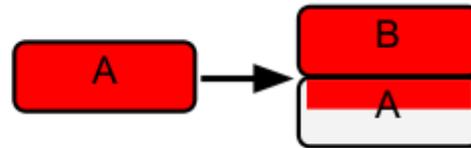
Concurrent Layered Learning (CLL)

---

## Overlapping Layered Learning



Combining Independently Learned Behaviors (CILB)



Partial Concurrent Layered Learning (PCLL)

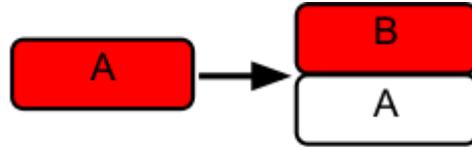


Previous Learned Layer Refinement (PLLR)

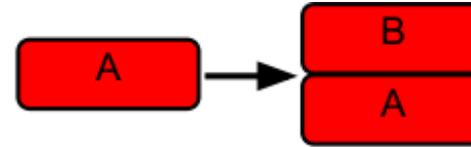
**Combining Independently Learned Behaviors:** Behaviors learned independently and then combined by relearning subset of behaviors' parameters

**Partial Concurrent Layered Learning:** Part, but not all, of a previously learned layer's behaviors are left open

# Overlapping Layered Learning



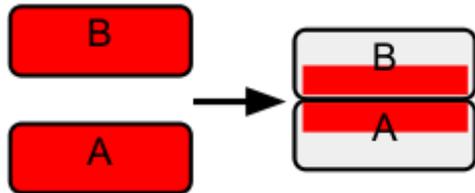
Sequential Layered Learning (SLL)



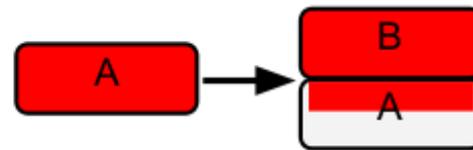
Concurrent Layered Learning (CLL)

---

## Overlapping Layered Learning



Combining Independently Learned Behaviors (CILB)



Partial Concurrent Layered Learning (PCLL)



Previous Learned Layer Refinement (PLLR)

**Combining Independently Learned Behaviors:** Behaviors learned independently and then combined by relearning subset of behaviors' parameters

**Partial Concurrent Layered Learning:** Part, but not all, of a previously learned layer's behaviors are left open

**Previous Learned Layer Refinement:** After a pair of layers is learned, part or all of the initial layer is unfrozen

# RoboCup 3D Simulation Domain

- Teams of 11 vs 11 autonomous robots play soccer
- Realistic physics using Open Dynamics Engine (ODE)
- Simulated robots modeled after Aldebaron Nao robot
- Robot receives noisy visual information about environment
- Robots can communicate over limited bandwidth channel



# RoboCup Champions 2011, 2012

# RoboCup Champions 2011, 2012

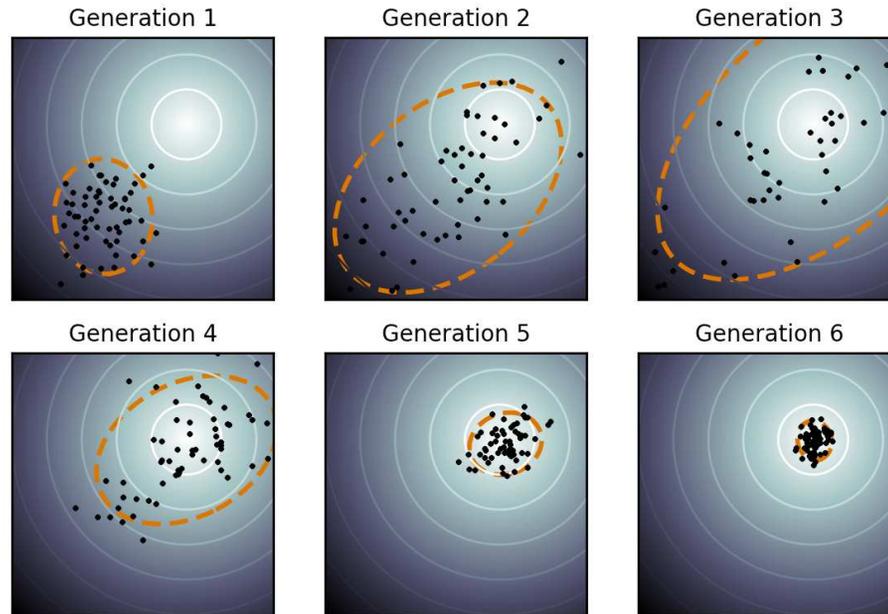
Humanoid Walk Learning via Layered Learning and CMA-ES

- Parameterized **double linear inverted pendulum model**

# RoboCup Champions 2011, 2012

## Humanoid Walk Learning via Layered Learning and CMA-ES

- Parameterized **double linear inverted pendulum model**



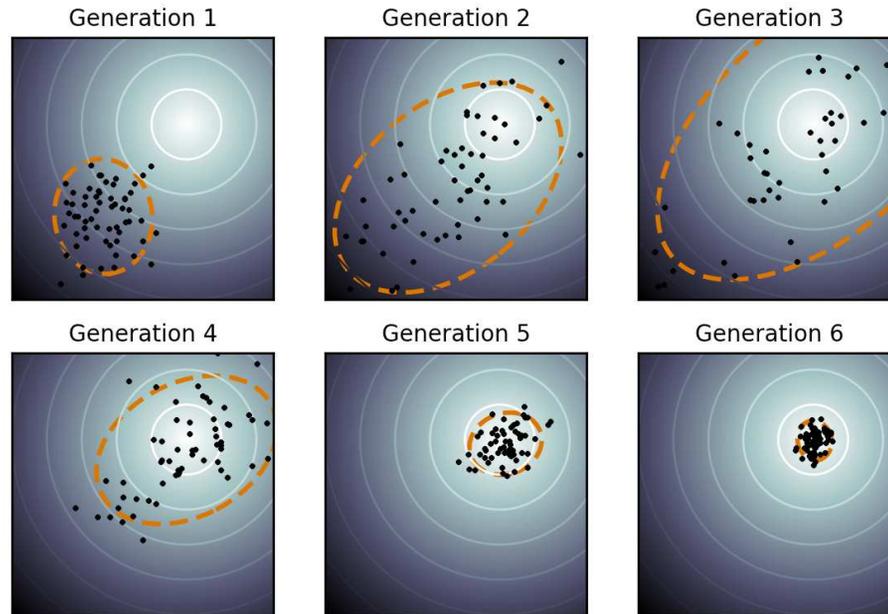
CMA-ES  
[Hansen, '09]

- Stochastic, derivative-free, numerical optimization method
- Candidates sampled from **multidimensional Gaussian**

# RoboCup Champions 2011, 2012

## Humanoid Walk Learning via Layered Learning and CMA-ES

- Parameterized **double linear inverted pendulum model**



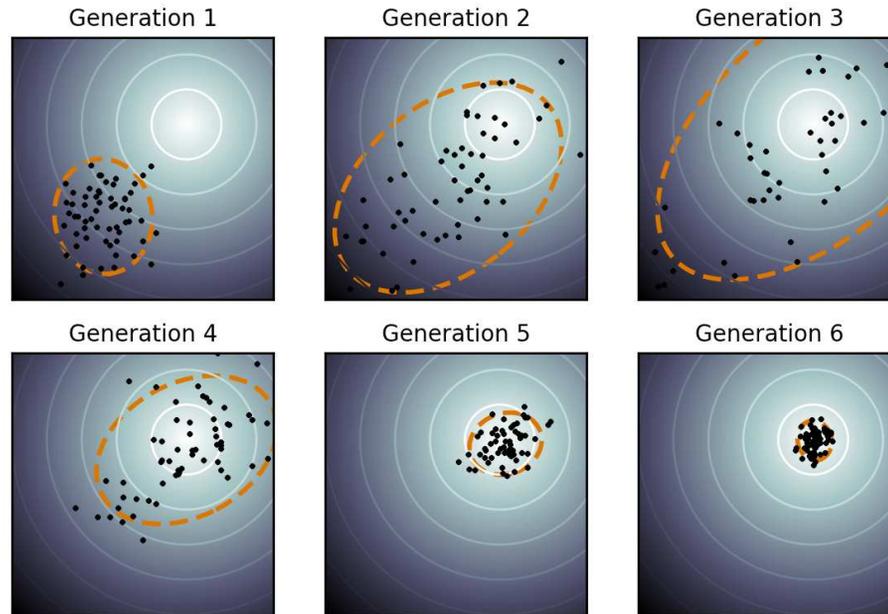
CMA-ES  
[Hansen, '09]

- Stochastic, derivative-free, numerical optimization method
- Candidates sampled from **multidimensional Gaussian**
  - **Mean** maximizes likelihood of previous successes
  - **Covariance** update controls search step sizes

# RoboCup Champions 2011, 2012

## Humanoid Walk Learning via Layered Learning and CMA-ES

- Parameterized **double linear inverted pendulum model**



CMA-ES  
[Hansen, '09]

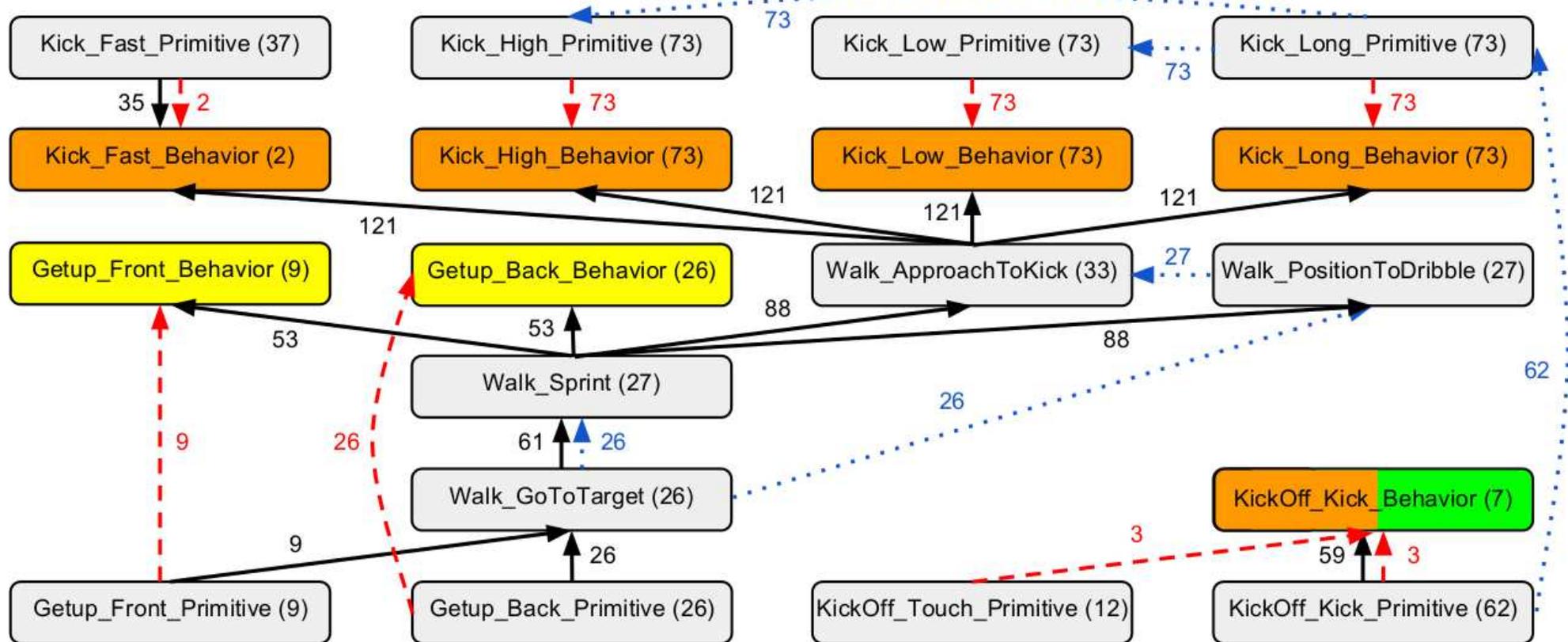
- Stochastic, derivative-free, numerical optimization method
- Candidates sampled from **multidimensional Gaussian**
  - **Mean** maximizes likelihood of previous successes
  - **Covariance** update controls search step sizes

**Initial walk**  
**3 layers**

**No layered learning**  
**Final walk**

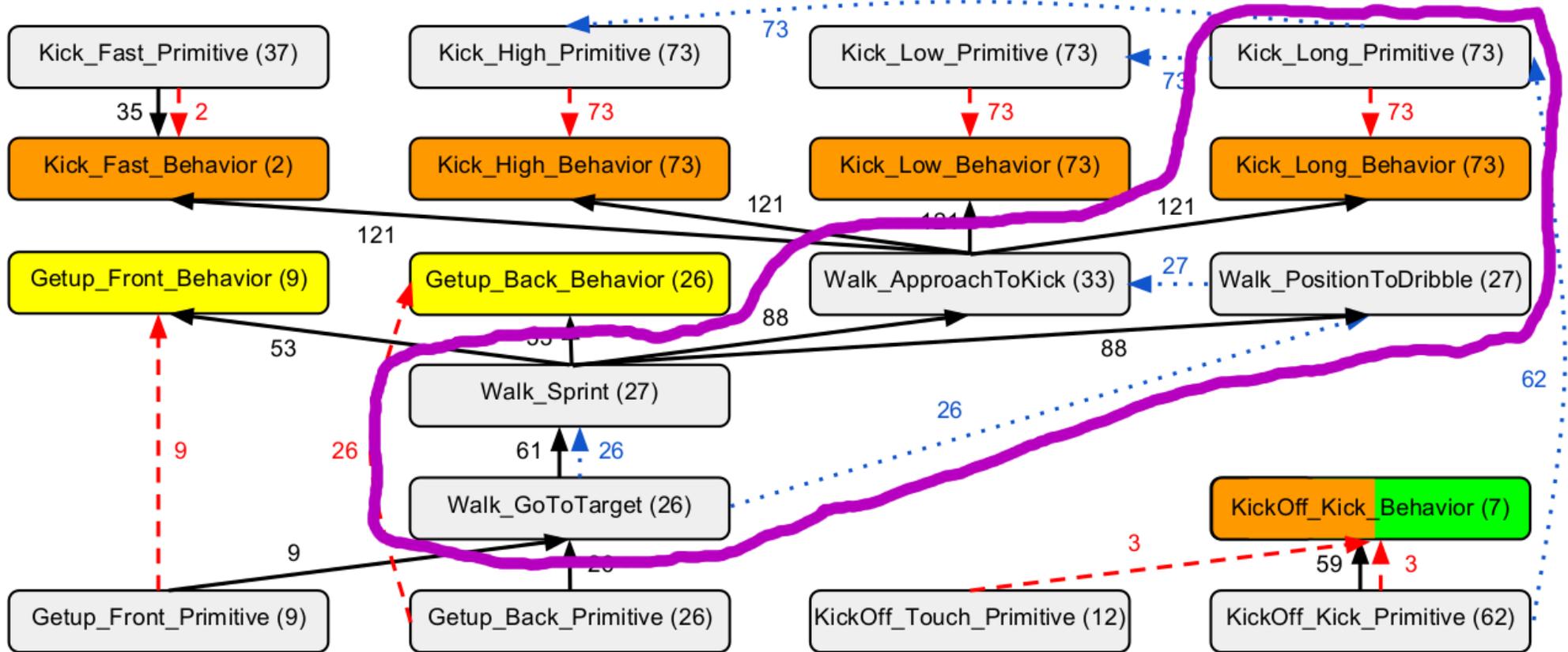
**2 layers**  
**Champs\*2**

# Learned Layers



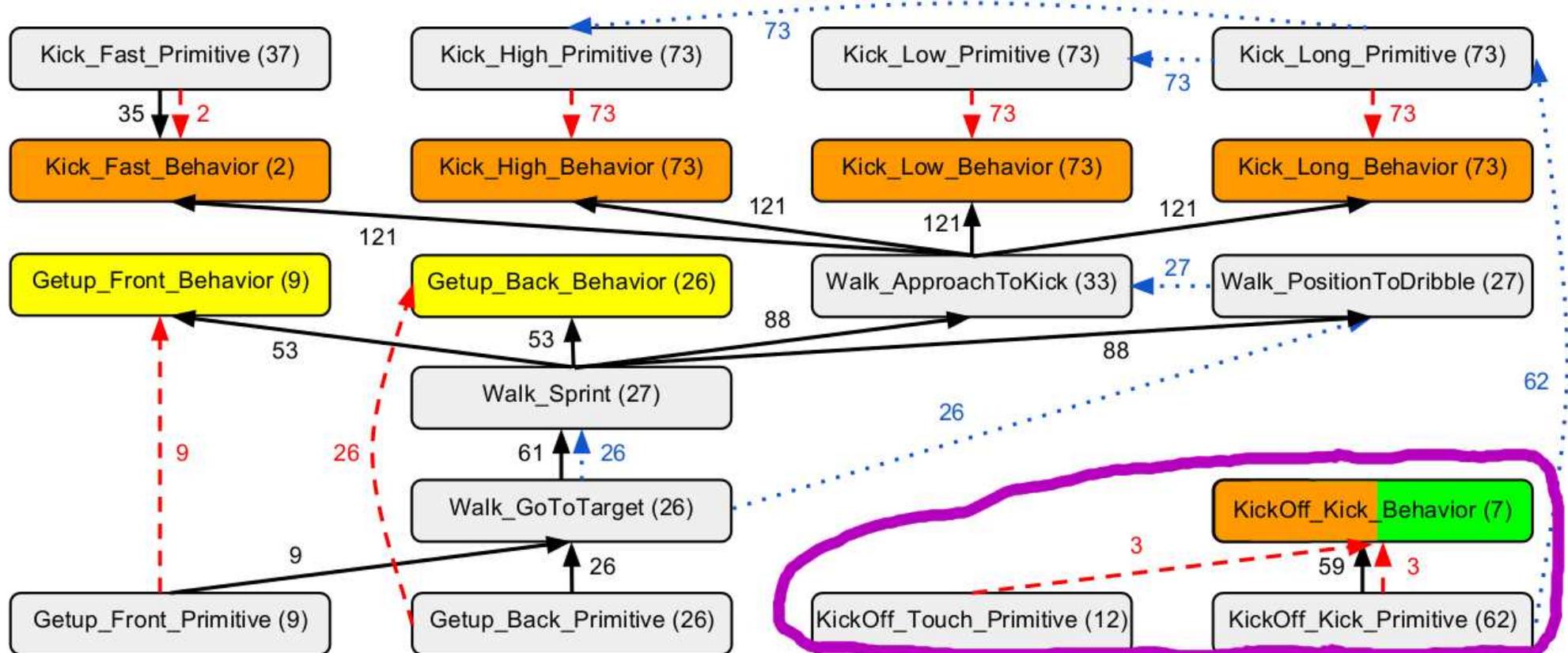
- 19 learned behaviors for standing up, walking, and kicking
  - ▶ CILB, PCLL, PLLR
- Over 500 parameters optimized during the course of learning
  - ▶ frozen, passed, seeded

# Dribbling and Kicking the Ball in the Goal



- Four different walk parameter sets
  - ▶ Target/sprint/position + **approach ball to kick**
- Learn **fixed kick**
- Combine **kick with walk**: combine independent layers (CILB)
  - ▶ **Overlap** kick parameters for positioning
- Final **walk and kick**

# Scoring on a Kickoff



- Kickoffs **indirect** (2 players must touch to score)
- Learn **fixed kick**
- Learn **touch** behavior **interferes**
- Combine **kick with touch**
  - ▶ Relearn position patterns: combine independent layers (**CILB**)
  - ▶ Learn new timing parameter: partial concurrent (**PCLL**)

# Impact of Overlapping Layered Learning

1000 games vs. top 3 teams from 2013

# Impact of Overlapping Layered Learning

1000 games vs. top 3 teams from 2013

Opponent	Average Goal Difference		
	Full Team	No Kickoff	Dribble Only
apollo3d	2.703 (0.041)	2.062 (0.038)	1.861 (0.034)
UTAustinVilla2013	1.589 (0.036)	1.225 (0.033)	0.849 (0.025)
fcportugal3d	3.991 (0.051)	3.189 (0.048)	1.584 (0.030)

**No Kickoff:** On kickoff, kick ball deep into opponent's end

**Dribble Only:** No kicking

# Repetition on Different Robot Types

Type 0: Standard Nao model

Type 1: Longer legs and arms

Type 2: Quicker moving legs

Type 3: Wider hips and longest legs and arms

Type 4: Added toes to foot

# Repetition on Different Robot Types

Type 0: Standard Nao model

Type 1: Longer legs and arms

Type 2: Quicker moving legs

Type 3: Wider hips and longest legs and arms

Type 4: Added toes to foot

	<b>Avg. Goal Difference per Robot Type</b>				
<b>Opponent</b>	<b>Type 0</b>	<b>Type 1</b>	<b>Type 2</b>	<b>Type 3</b>	<b>Type 4</b>
apollo3d	1.787	1.819	1.820	1.543	2.827
UTAustinVilla2013	0.992	0.892	1.276	0.573	1.141
fcportugal3d	2.423	3.025	3.275	2.678	4.033

# Repetition on Different Robot Types

Type 0: Standard Nao model

Type 1: Longer legs and arms

Type 2: Quicker moving legs

Type 3: Wider hips and longest legs and arms

Type 4: Added toes to foot

	<b>Avg. Goal Difference per Robot Type</b>				
<b>Opponent</b>	<b>Type 0</b>	<b>Type 1</b>	<b>Type 2</b>	<b>Type 3</b>	<b>Type 4</b>
apollo3d	1.787	1.819	1.820	1.543	2.827
UTAustinVilla2013	0.992	0.892	1.276	0.573	1.141
fcportugal3d	2.423	3.025	3.275	2.678	4.033

## Computation per type

≈ 700k parameter sets evaluated

≈ 1.5 years compute time (≈ 50 hours on condor cluster)

# RoboCup 2014

Won competition with **undefeated** record: outscored opps 52–0

Opponent	Avg. Goal Diff.	Record (W-L-T)	Goals (F/A)	KO Score %
BahiaRT	2.075 (0.030)	990-0-10	2092/17	96.2
FCPortugal	2.642 (0.034)	986-0-14	2748/106	83.4
magmaOffenburg	2.855 (0.035)	990-0-10	2864/9	88.3
RoboCanes	3.081 (0.046)	974-0-26	3155/74	69.4
FUT-K	3.236 (0.039)	998-0-2	3240/4	96.3
SEU_Jolly	4.031 (0.062)	995-0-5	4034/3	87.6
KarachiKoalas	5.681 (0.046)	1000-0-0	5682/1	87.5
ODENS	7.933 (0.041)	1000-0-0	7933/0	92.1
HfutEngine	8.510 (0.050)	1000-0-0	8510/0	94.7
Mithras3D	8.897 (0.041)	1000-0-0	8897/0	90.4
L3M-SIM	9.304 (0.043)	1000-0-0	9304/0	93.7

# RoboCup 2014

Won competition with **undefeated** record: outscored opps 52–0

Opponent	Avg. Goal Diff.	Record (W-L-T)	Goals (F/A)	KO Score %
BahiaRT	2.075 (0.030)	990-0-10	2092/17	96.2
FCPortugal	2.642 (0.034)	986-0-14	2748/106	83.4
magmaOffenburg	2.855 (0.035)	990-0-10	2864/9	88.3
RoboCanes	3.081 (0.046)	974-0-26	3155/74	69.4
FUT-K	3.236 (0.039)	998-0-2	3240/4	96.3
SEU_Jolly	4.031 (0.062)	995-0-5	4034/3	87.6
KarachiKoalas	5.681 (0.046)	1000-0-0	5682/1	87.5
ODENS	7.933 (0.041)	1000-0-0	7933/0	92.1
HfutEngine	8.510 (0.050)	1000-0-0	8510/0	94.7
Mithras3D	8.897 (0.041)	1000-0-0	8897/0	90.4
L3M-SIM	9.304 (0.043)	1000-0-0	9304/0	93.7

- After: **11,000 games**: won all by 67 (**no losses**)

# RoboCup 2014

Won competition with **undefeated** record: outscored opps 52–0

Opponent	Avg. Goal Diff.	Record (W-L-T)	Goals (F/A)	KO Score %
BahiaRT	2.075 (0.030)	990-0-10	2092/17	96.2
FCPortugal	2.642 (0.034)	986-0-14	2748/106	83.4
magmaOffenburg	2.855 (0.035)	990-0-10	2864/9	88.3
RoboCanes	3.081 (0.046)	974-0-26	3155/74	69.4
FUT-K	3.236 (0.039)	998-0-2	3240/4	96.3
SEU_Jolly	4.031 (0.062)	995-0-5	4034/3	87.6
KarachiKoalas	5.681 (0.046)	1000-0-0	5682/1	87.5
ODENS	7.933 (0.041)	1000-0-0	7933/0	92.1
HfutEngine	8.510 (0.050)	1000-0-0	8510/0	94.7
Mithras3D	8.897 (0.041)	1000-0-0	8897/0	90.4
L3M-SIM	9.304 (0.043)	1000-0-0	9304/0	93.7

- After: **11,000 games**: won all by 67 (**no losses**)
- **Highlights** from Final vs. RoboCanes (University of Miami)

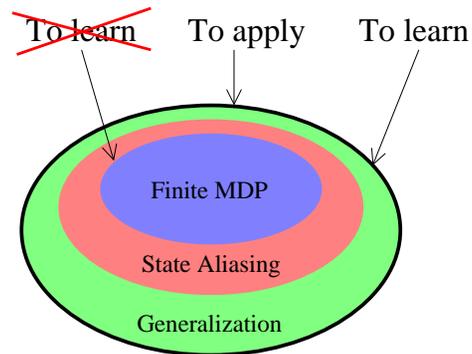
# RoboCup 2014

Won competition with **undefeated** record: outscored opps 52–0

Opponent	Avg. Goal Diff.	Record (W-L-T)	Goals (F/A)	KO Score %
BahiaRT	2.075 (0.030)	990-0-10	2092/17	96.2
FCPortugal	2.642 (0.034)	986-0-14	2748/106	83.4
magmaOffenburg	2.855 (0.035)	990-0-10	2864/9	88.3
RoboCanes	3.081 (0.046)	974-0-26	3155/74	69.4
FUT-K	3.236 (0.039)	998-0-2	3240/4	96.3
SEU_Jolly	4.031 (0.062)	995-0-5	4034/3	87.6
KarachiKoalas	5.681 (0.046)	1000-0-0	5682/1	87.5
ODENS	7.933 (0.041)	1000-0-0	7933/0	92.1
HfutEngine	8.510 (0.050)	1000-0-0	8510/0	94.7
Mithras3D	8.897 (0.041)	1000-0-0	8897/0	90.4
L3M-SIM	9.304 (0.043)	1000-0-0	9304/0	93.7

- After: **11,000 games**: won all by 67 (**no losses**)
- **Highlights** from Final vs. RoboCanes (University of Miami)
- More info: [www.cs.utexas.edu/~AustinVilla/sim/3dsimulation/](http://www.cs.utexas.edu/~AustinVilla/sim/3dsimulation/)

# Practical RL



## ● Representation

- ▶ Selecting the Algorithm: parameterized domains [K.&S., MLJ 2011]
- ▶ Adapting Representation: NEAT+Q [Whiteson & S., JMLR 2006]

## ● Interaction

- ▶ With adversaries: CMLES [Chakraborty & S., ICML 2010]
- ▶ **With ad hoc teammates: PLASTIC** [Barrett, thesis 2014]
- ▶ With people: TAMER [Knox & S., AAMAS 2010]

## ● Synthesis

- ▶ Of Algorithms: Layered Learning [S., MIT Press 2000]
- ▶ Of Concepts: Fitted R-MAXQ [Jong & S., ECML 2009]

## ● Mortality

- ▶ Leverage the Past: Transfer Learning [Taylor, S., & Liu, JMLR 2007]
- ▶ Acknowledge a Finite Future: TEXPLORE [Hester & S., MLJ 2013]

# Making Friends on the Fly: Advances in Ad Hoc Teamwork

Samuel Barrett, Katie Genter, and Peter Stone



# Ad Hoc Teamwork [Stone et al., AIJ 2013]

- Only in control of a single agent or subset of agents
- Unknown teammates
- Shared goals
- No pre-coordination

# Ad Hoc Teamwork [Stone et al., AIJ 2013]

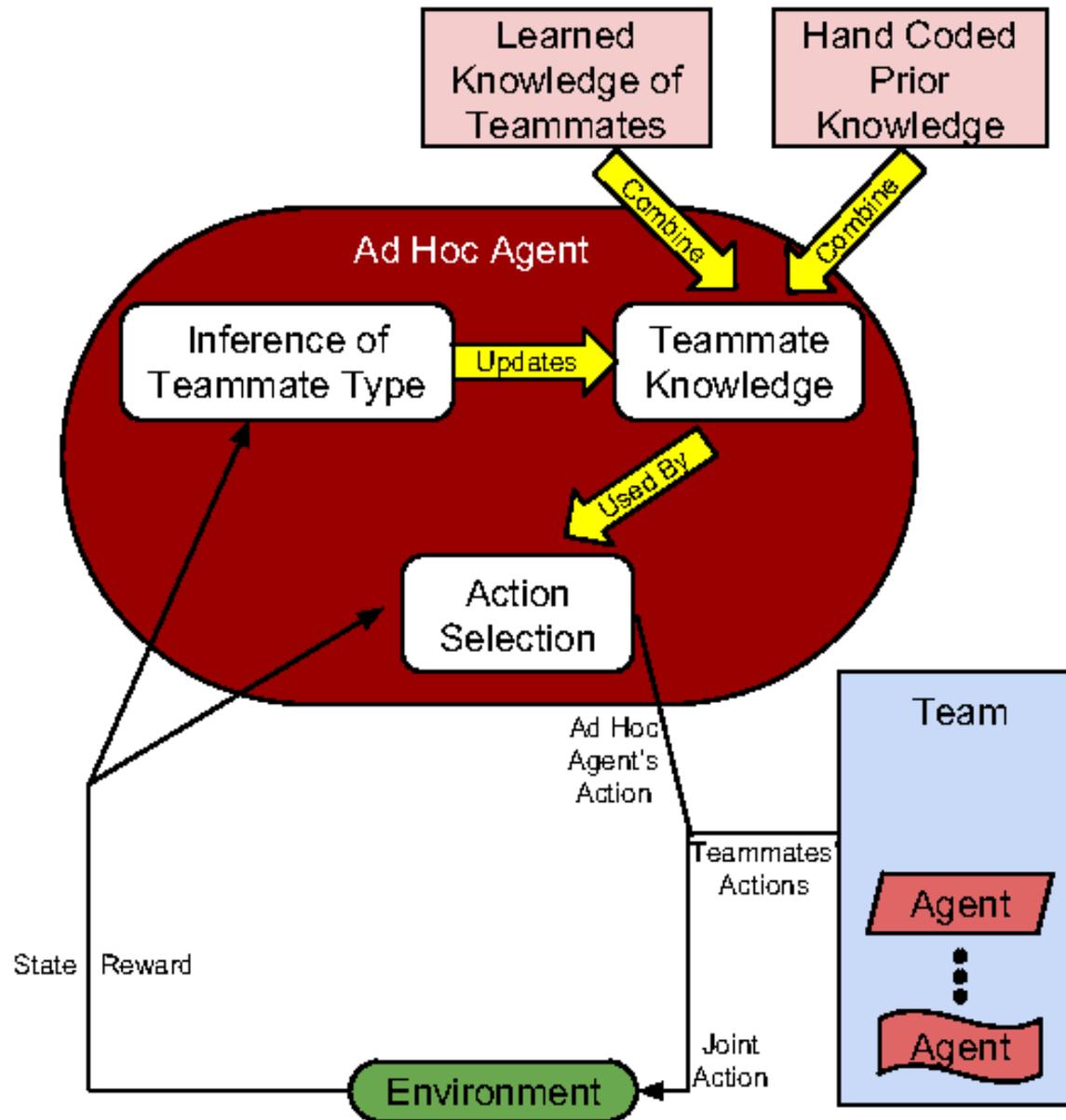
- Only in control of a single agent or subset of agents
- Unknown teammates
- Shared goals
- No pre-coordination

## Examples in humans:

- Pick up soccer
- Accident response

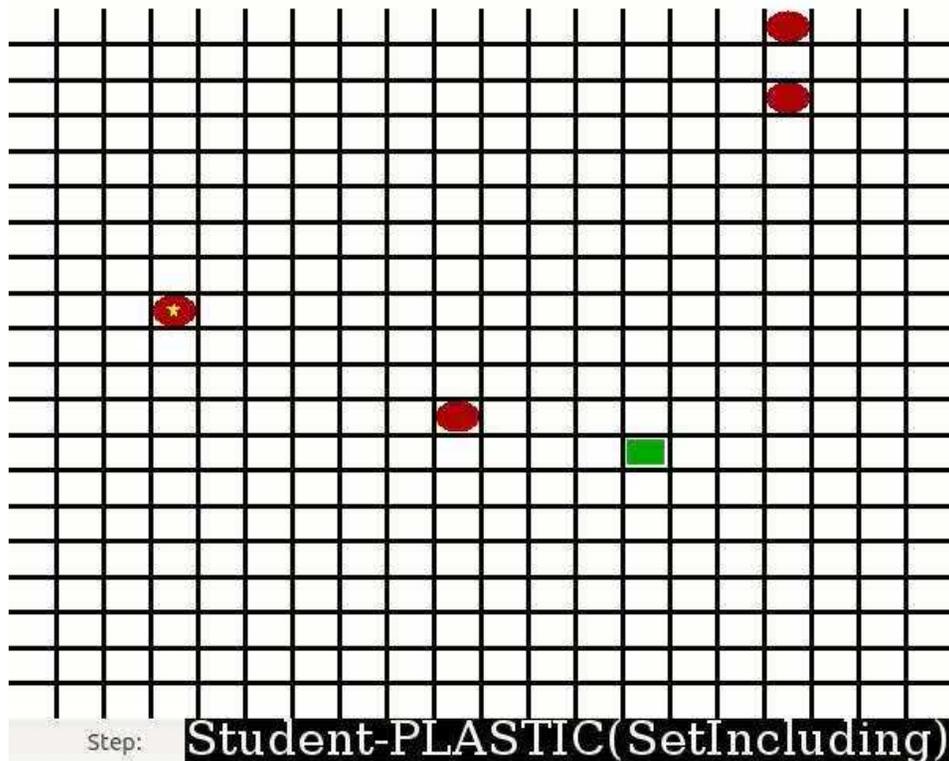


# PLASTIC: Planning and Learning to Adapt Swiftly to Teammates to Improve Cooperation



# Testbed Domains

- Agent **replaces single teammate** in otherwise coherent team
- Adapts based on knowledge **learned from previous teammates**



# Other Ad Hoc Teamwork

- Introduced as **AAAI Challenge Problem**

[Stone et al. '10]

# Other Ad Hoc Teamwork

- Introduced as **AAAI Challenge Problem**
- Theory: repeated games, bandits

[Stone et al. '10]

[Stone et al., '11]

# Other Ad Hoc Teamwork

- Introduced as **AAAI Challenge Problem** [Stone et al. '10]
- Theory: repeated games, bandits [Stone et al., '11]
- Experiments: **flocking** [Genter & Stone, '12]

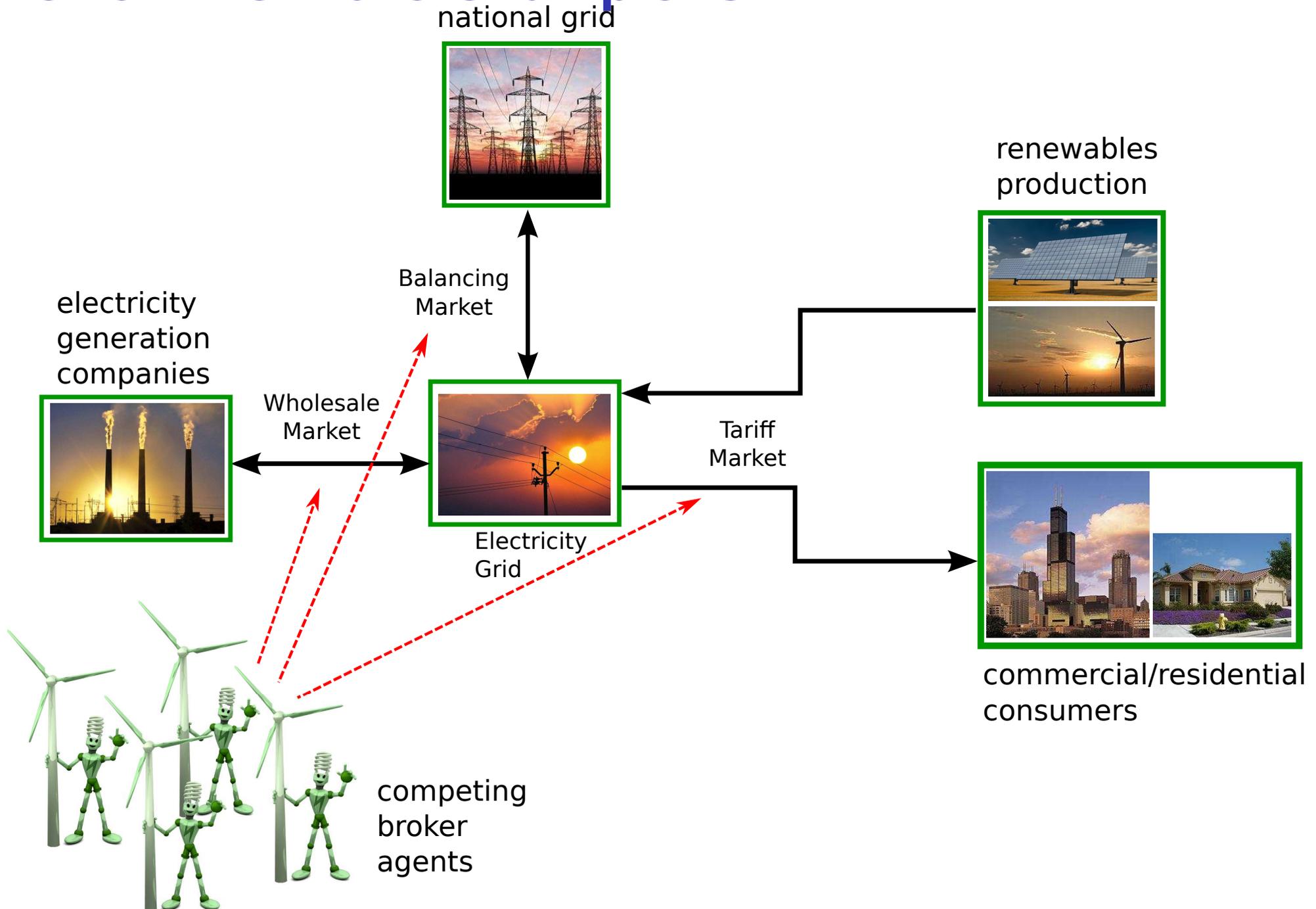
# Other Ad Hoc Teamwork

- Introduced as **AAAI Challenge Problem** [Stone et al. '10]
- Theory: repeated games, bandits [Stone et al., '11]
- Experiments: **flocking** [Genter & Stone, '12]
- **RoboCup experiments** [Genter et al., '15]

# Other Ad Hoc Teamwork

- Introduced as **AAAI Challenge Problem** [Stone et al. '10]
- Theory: repeated games, bandits [Stone et al., '11]
- Experiments: **flocking** [Genter & Stone, '12]
- **RoboCup experiments** [Genter et al., '15]
- AAI Workshops, JAAMAS special issue

# Power TAC: 2013 Champions



# Acknowledgements



Nicholas  
Jong



Matthew  
Taylor



Shimon  
Whiteson



Doran  
Chakraborty



Todd  
Hester



Shivaram  
Kalyanakrishnan



Brad  
Knox



Samuel  
Barrett



Patrick  
MacAlpine



Daniel  
Urieli



Katie  
Genter

LARG (past and present), NSF, ONR, FHWA, DARPA, IBM, Google