

Clemens Lode, Urban Richter, Hartmut Schmeck
Karlsruhe Institute of Technology (Germany)
Institute AIFB
July, GECCO 2010

Adaption of XCS to Multi-Learner Predator/Prey Scenarios

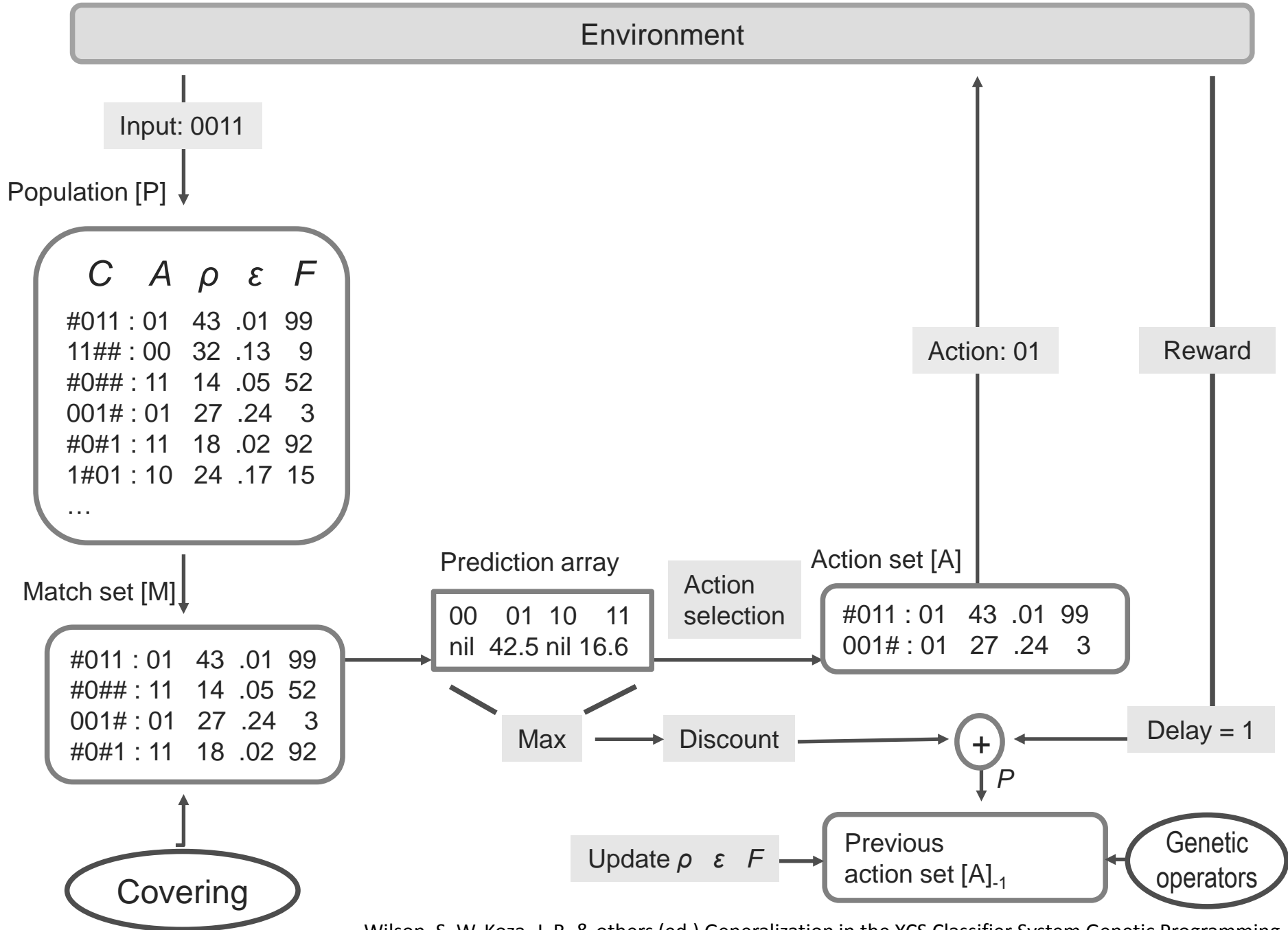


Outline

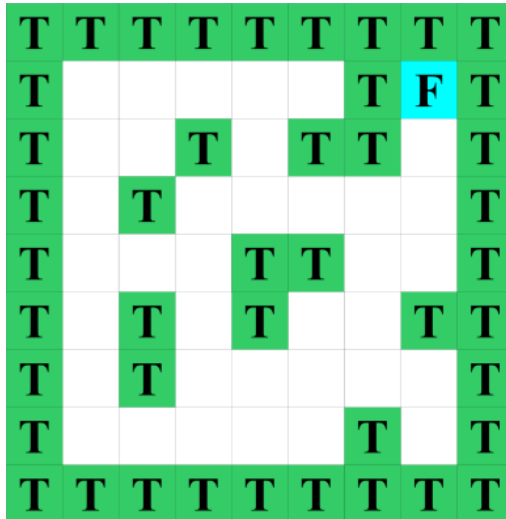


<http://www.flickr.com/photos/yathin>

- Learning Classifier Systems
- XCS in Predator/Prey Scenarios
- Adapting the Reward Function
- Experimental Results



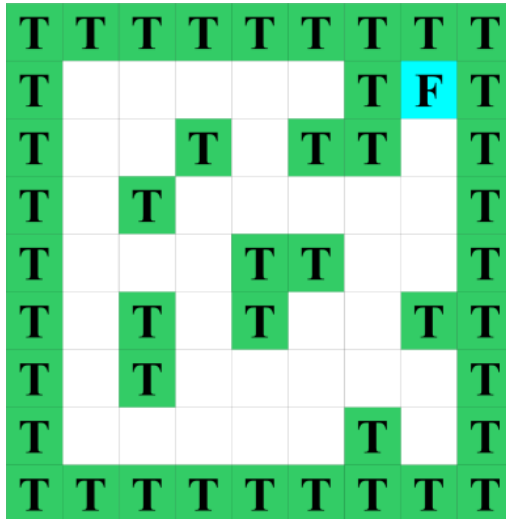
Learning Classifier Systems



T: Tree
F: Food

- Standard (Multi-Step) Problem:
 - Maze6
- Goal:
 - Find the shortest path to from a random position to food

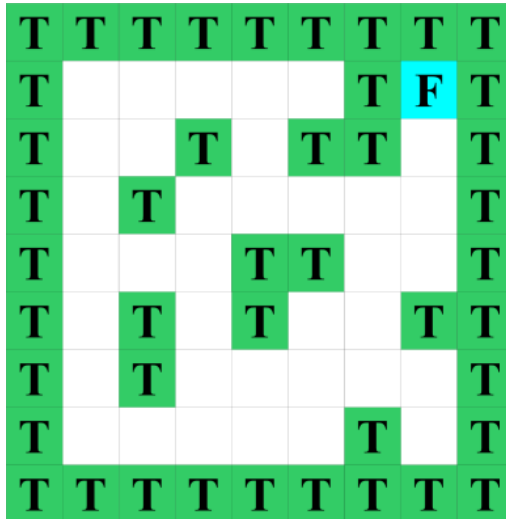
Learning Classifier Systems



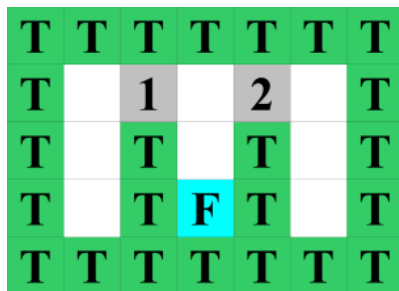
T: Tree
F: Food

- Problem:
 - Limited sensors, no global knowledge
 - Partially observable Markov decision process
- Solution:
 - Iterations, back-propagation of reward

Learning Classifier Systems



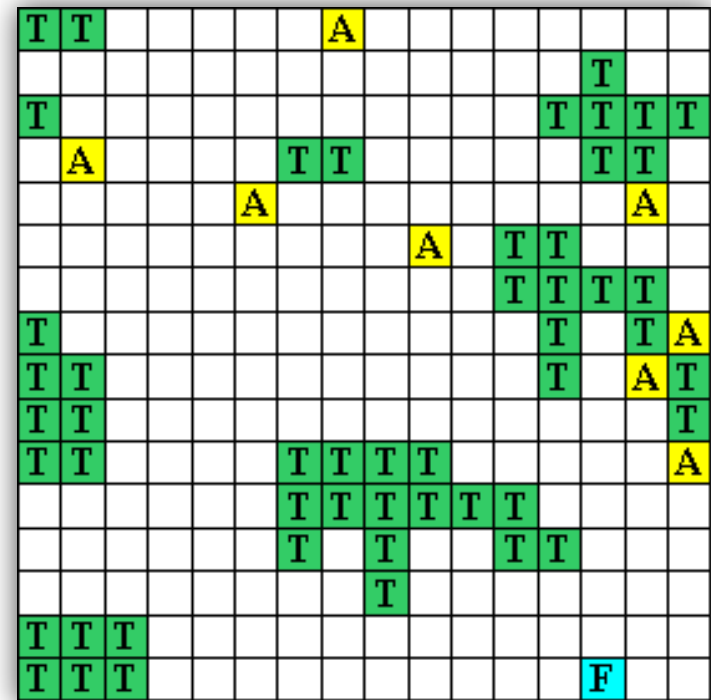
T: Tree
F: Food



- Problem:
 - Limited sensors, no global knowledge
 - Partially observable Markov decision process
- Solution:
 - Iterations, back-propagation of reward
- Aliasing positions:
 - Handle by using memory

Predator/Prey Scenarios

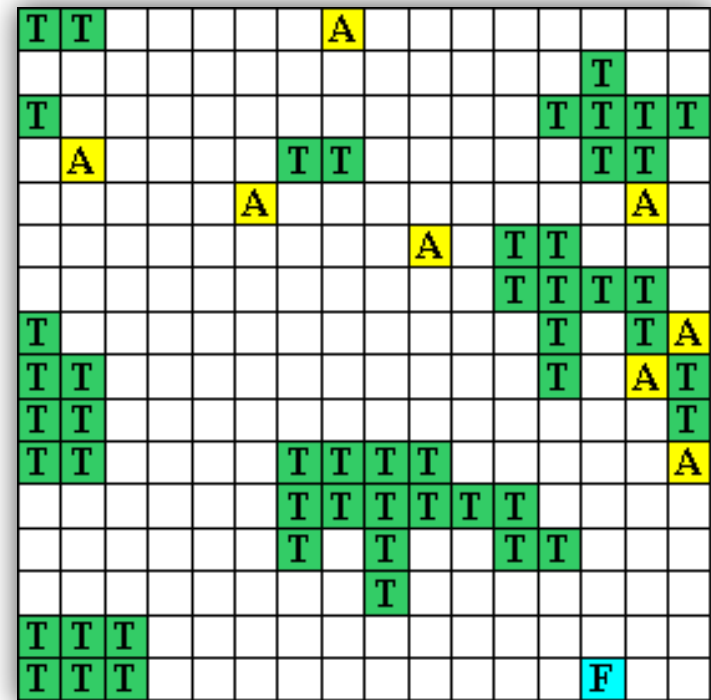
- Many aliasing positions
- Other agents present
 - Dynamic world
 - food and other agents move
- Limited sensors
 - ca. 5 cells range



T: Trees
F: Food
A: Agent

Predator/Prey Scenarios

- Terminology:
 - Obstacles, prey, predator
- Goal: Try to stay near the prey
 - Global observation task
 - Runs continuously
 - Maximize average quality



T: Trees/Obstacles
F: Food/Prey
A: Agent/Predator

Classification of Predator/Prey Scenarios

- (1) Access to local information only
 - (2) Open areas with some obstacles
 - (3) Internal state unknown to others
- No standard MDP
 - Limited sensors (1, 3)
 - Aliasing positions (2)

Classification of Predator/Prey Scenarios

- (1) Access to local information only
 - (2) Open areas with some obstacles
 - (3) Internal state unknown to others
 - (4) Dynamic scenario
- No standard MDP
 - Limited sensors (1, 3)
 - Aliasing positions (2)
 - No POMDP
 - Non-static scenario (4)

Classification of Predator/Prey Scenarios

- (1) Access to local information only
 - (2) Open areas with some obstacles
 - (3) Internal state unknown to others
 - (4) Dynamic scenario
 - (5) Predators share global observation task
 - (6) Runs continuously
- No standard MDP
 - Limited sensors (1, 3)
 - Aliasing positions (2)
 - No POMDP
 - Non-static scenario (4)
 - XCS has to be adapted
 - No “final” reward (5), no iterations (6)

Adaption of the Standard XCS Reward Function

- Standard implementation:
 - Reward:
 - Prey is in a neighboring cell
- Adapted implementation
 - Reward:
 - Prey is in observation range
 - "XCS obs"
 - Prey is in sight range
 - "XCS sight"

Adaption of the Standard XCS Reward Function

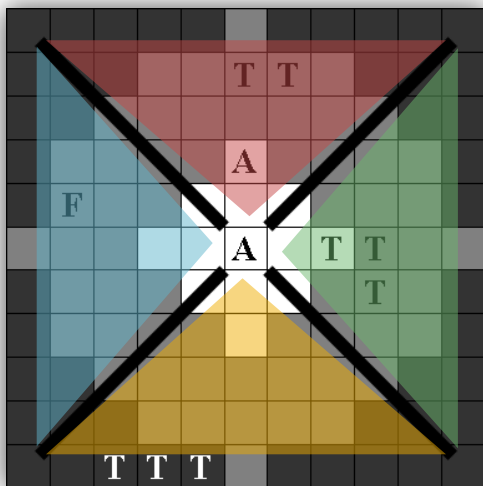
- Standard implementation:
 - Reward:
 - Prey is in a neighboring cell
 - Reward action:
 - Assign reward
 - Restart scenario
 - Explore/exploit phase

- Adapted implementation
 - Reward:
 - Prey is in observation range
 - "XCS obs"
 - Prey is in sight range
 - "XCS sight"
 - Reward action:
 - Assign reward
 - Continue scenario
 - Always use exploit phase

Sensors

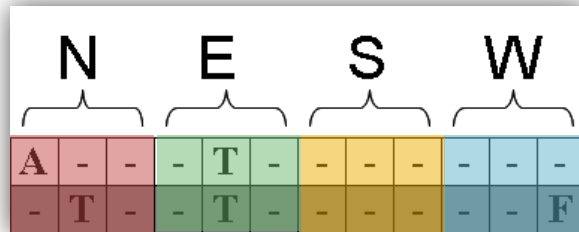
One sensor array for each direction

N			E			S			W		
A	-	-	-	T	-	-	-	-	-	-	-
-	T	-	-	T	-	-	-	-	-	-	F

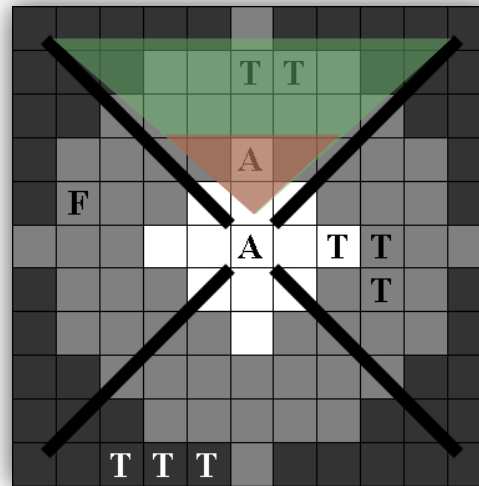
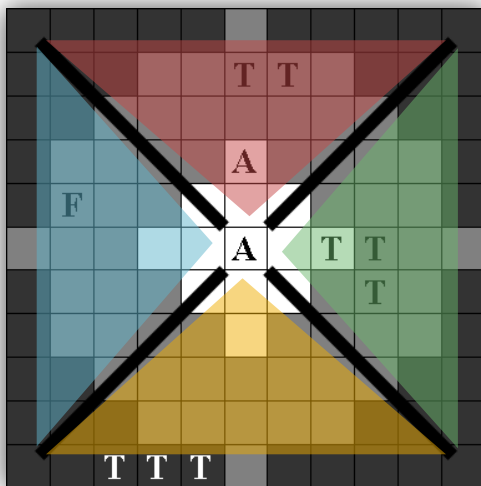
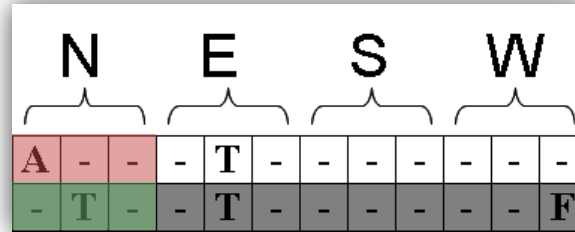


Sensors

One sensor array for each direction

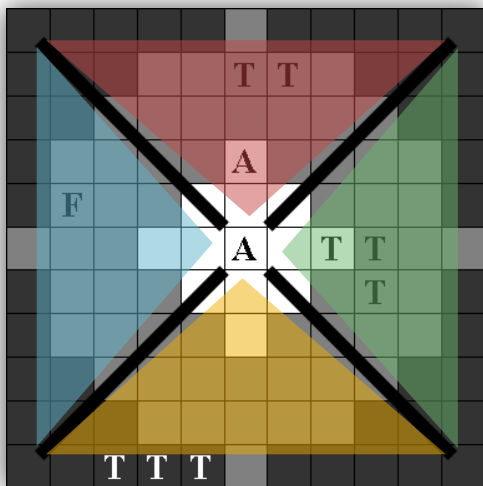
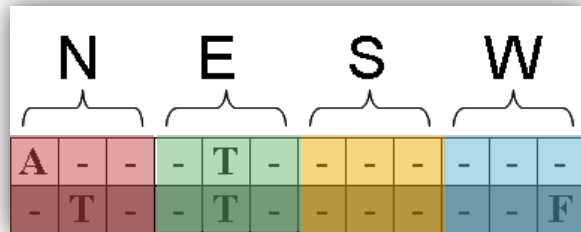


Sensors can sense either far or near (observation range / sight range)

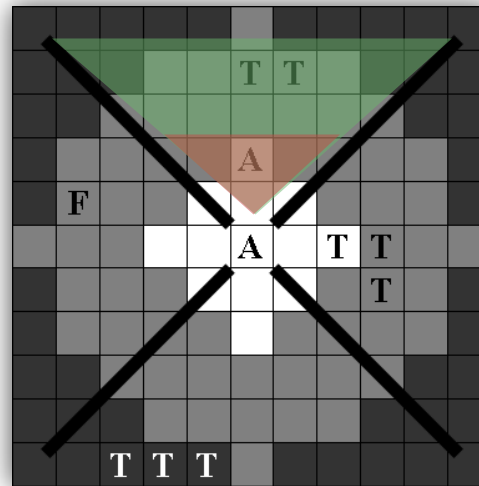
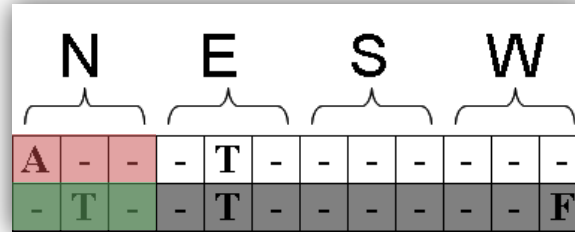


Sensors

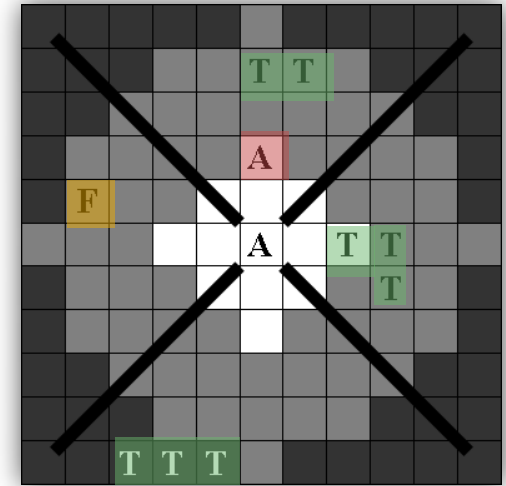
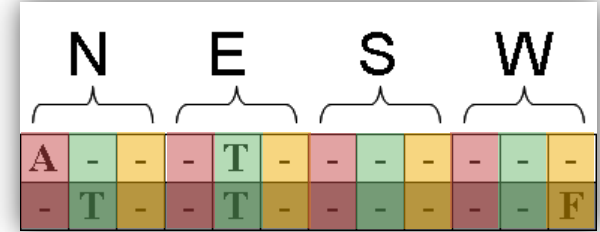
One sensor array for each direction



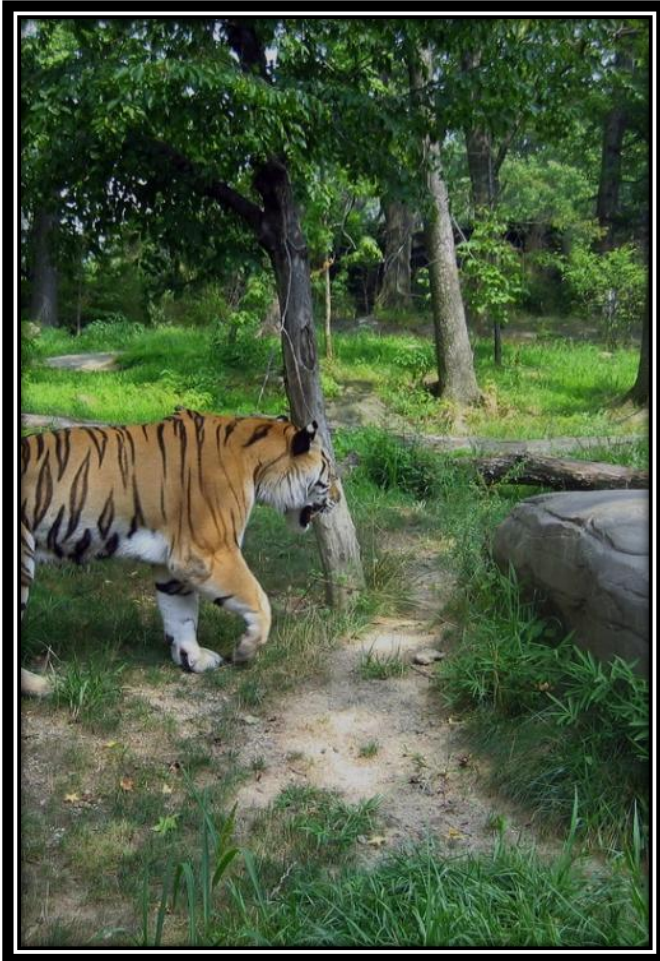
Sensors can sense either far or near (observation range / sight range)



Sensors can distinguish between predators, prey, and obstacles



Testing Methodology

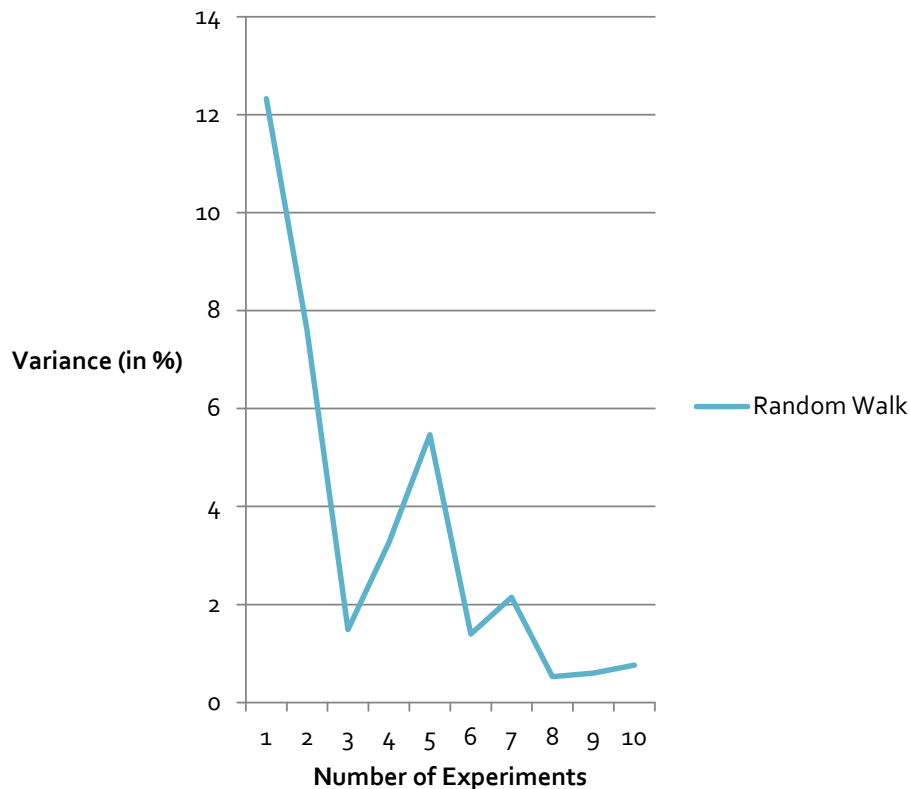


http://www.flickr.com/photos/james_crowley

- Behavior
 - 1 Prey
 - "Obstacle-evading prey"
 - "Predator-evading prey"
 - "Blinded Prey"
 - Speed 2 cells / step
 - 8 Predators
 - Speed 1 cell / step
- Standard XCS parameter settings

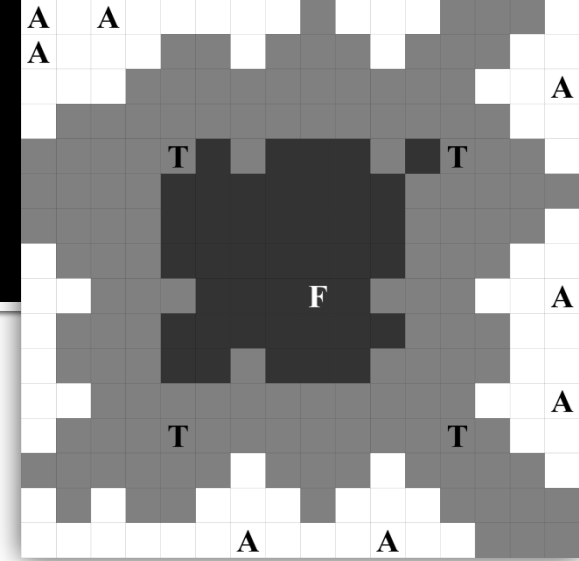
Testing Methodology

Variance in the
16x16 Predator/Prey Scenario



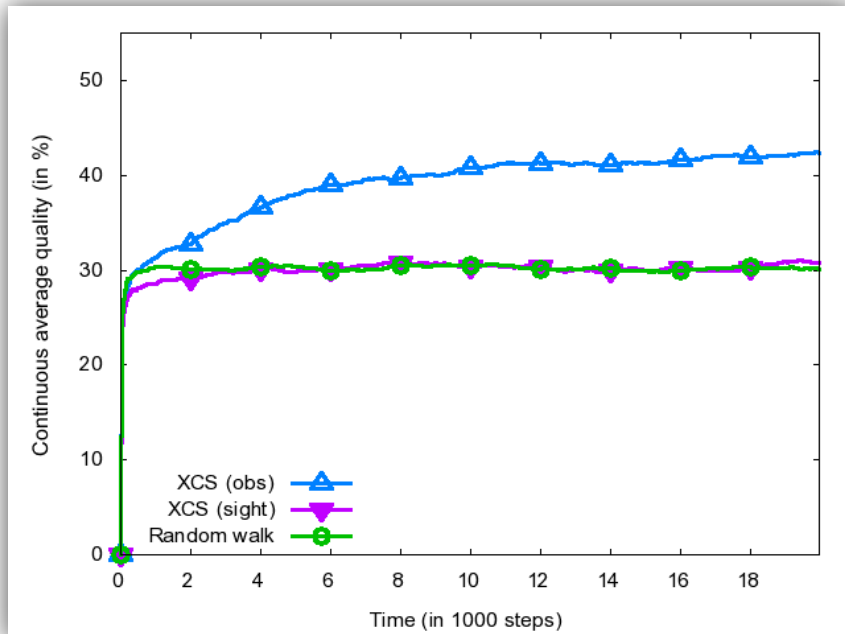
- 2,000,000 steps
- Reset of XCS every 20,000 steps (=“experiment”)
- Reset of scenario (new random positions) every 2,000 steps

XCS Experimental Results “Pillar Scenario”

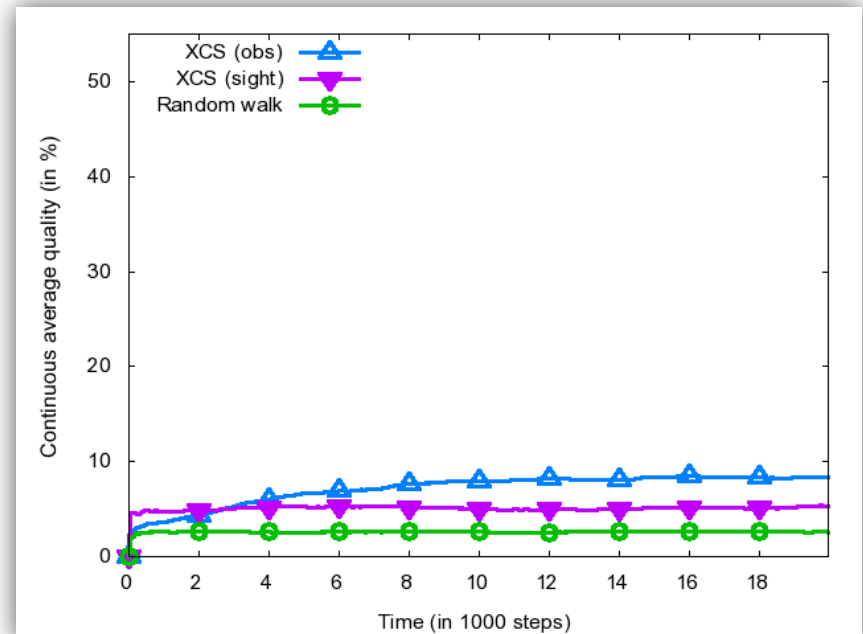


- XCS (obs) shows some learning

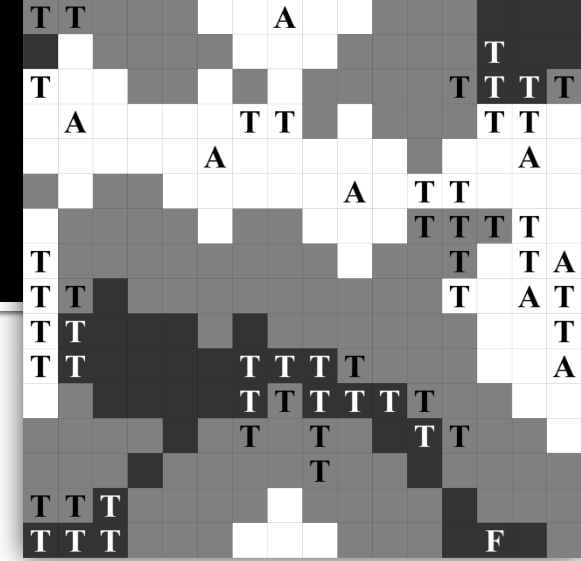
Obstacle-evading prey



Predator-evading prey

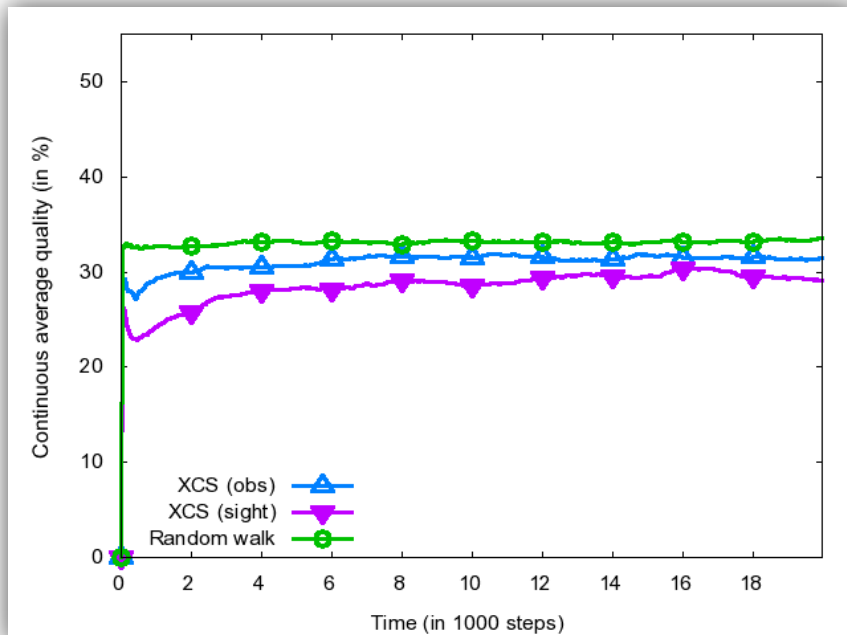


XCS Experimental Results “Random Scenario”

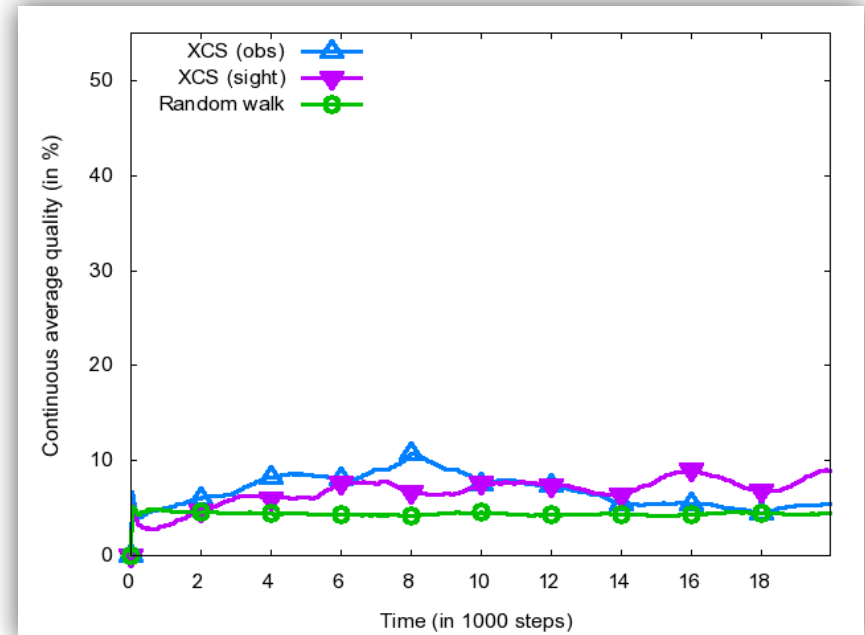


- XCS shows very little learning

Obstacle-evading prey



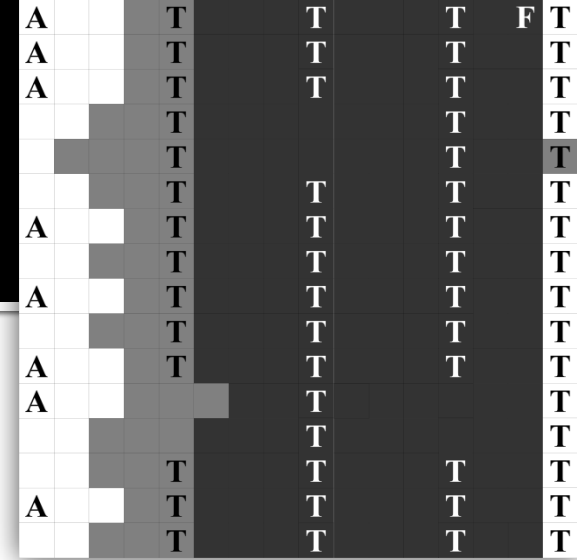
Predator-evading prey



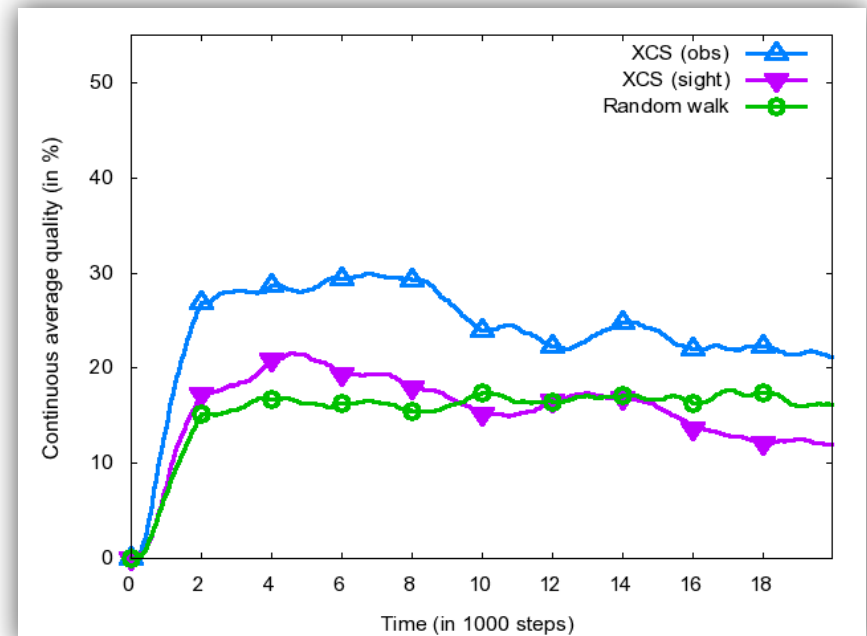
XCS Experimental Results

“Difficult Scenario”

- XCS shows significant learning
 - But also unlearning after 8,000 steps
- “Difficult Scenario” is a maze-like scenario, this result was expected



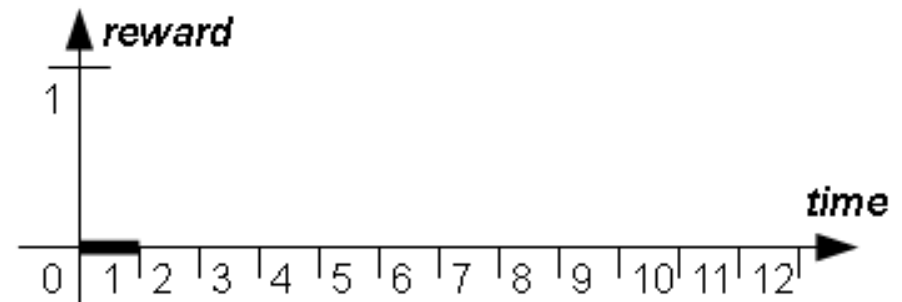
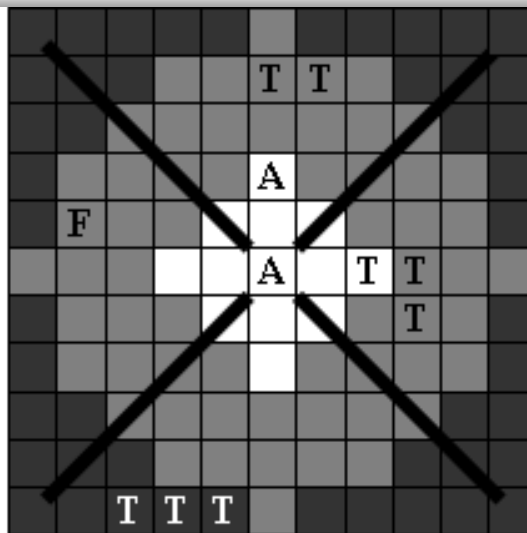
Blinded prey



Reward Events "eventXCS"

- Move West

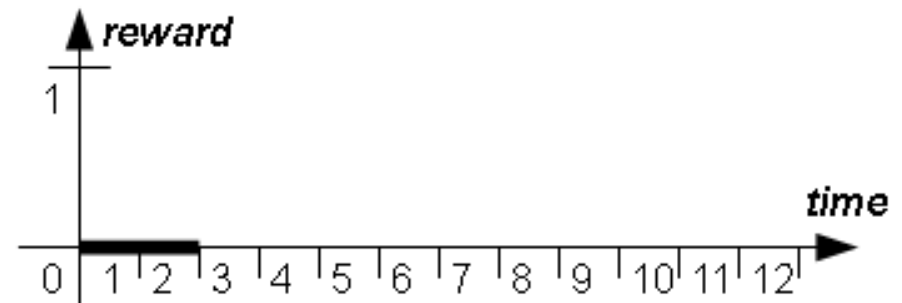
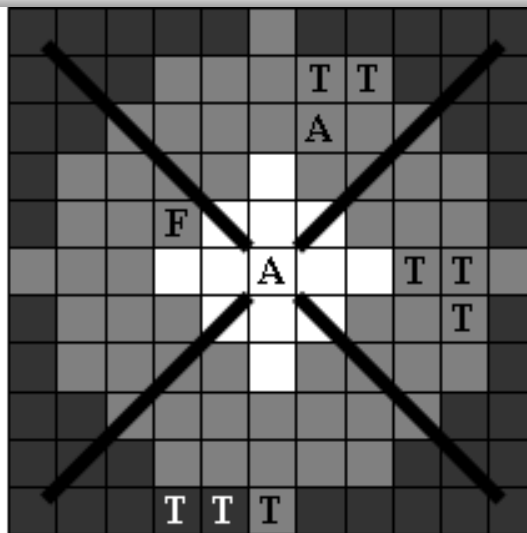
	N	E	S	W
A	1	0	0	0
T	0	1	0	0
F	0	0	0	0
A	0	0	0	0
T	1	1	0	0
F	0	0	0	1



Reward Events "eventXCS"

- Move West

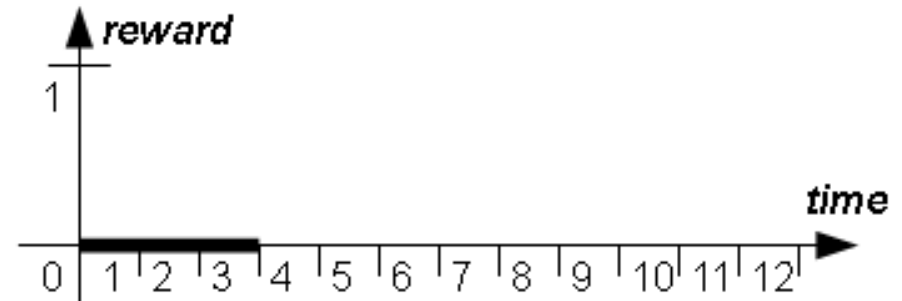
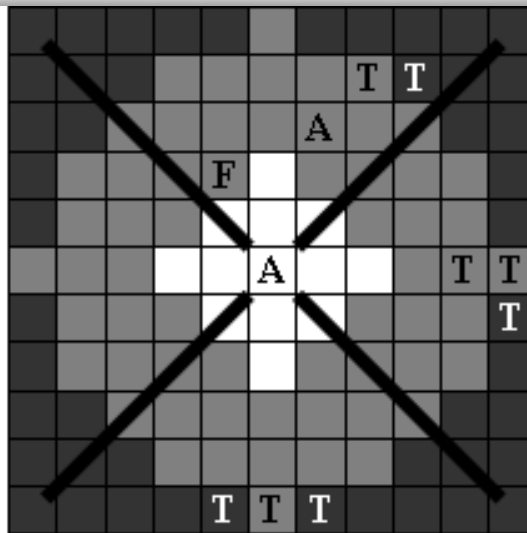
	N	E	S	W
A	0	0	0	0
T	0	0	0	0
F	0	0	0	0
A	1	0	0	0
T	1	1	1	0
F	0	0	0	1



Reward Events "eventXCS"

- Move North

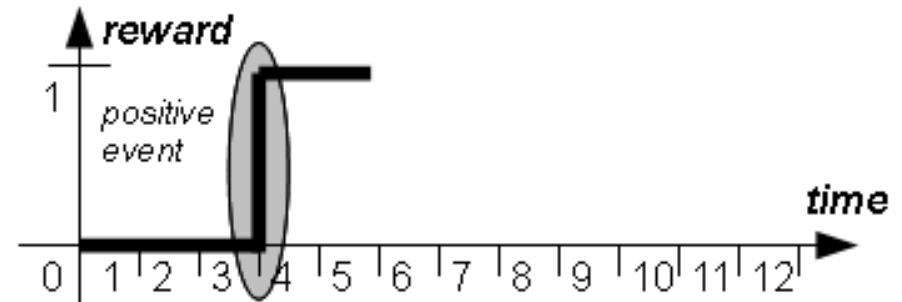
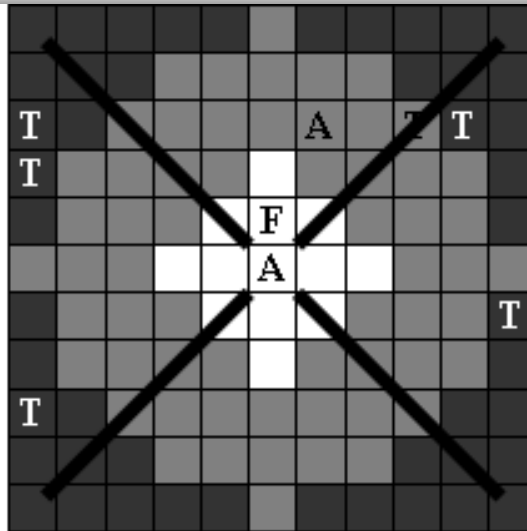
	N	E	S	W
A	0	0	0	0
T	0	0	0	0
F	0	0	0	0
A	1	0	0	0
T	1	1	1	0
F	1	0	0	0



Reward Events "eventXCS"

- Move North

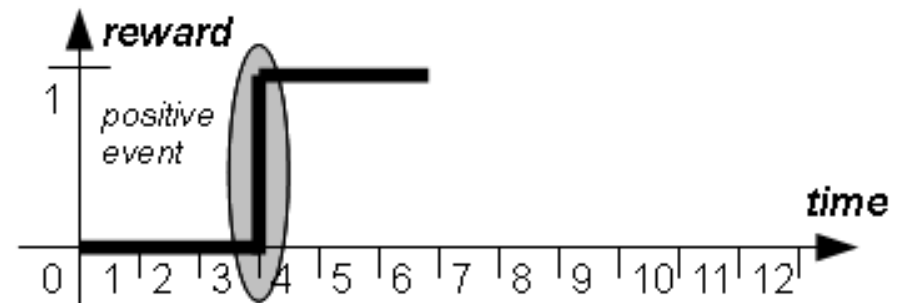
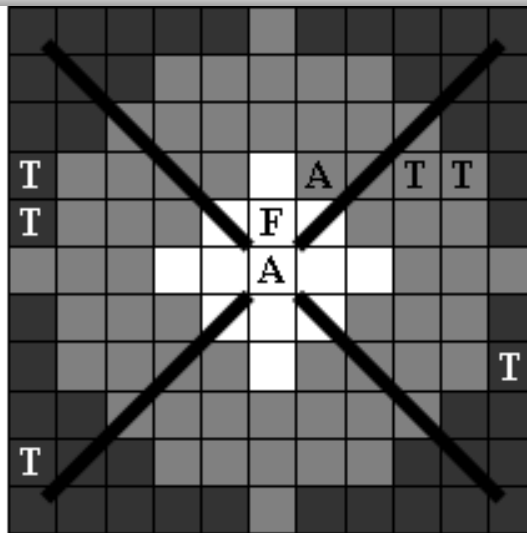
	N	E	S	W
A	0	0	0	0
T	0	0	0	0
F	1	0	0	0
A	1	0	0	0
T	1	0	0	0
F	0	0	0	0



Reward Events "eventXCS"

- Move North

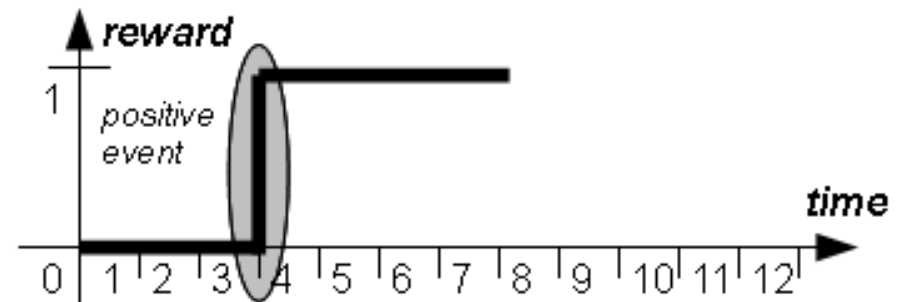
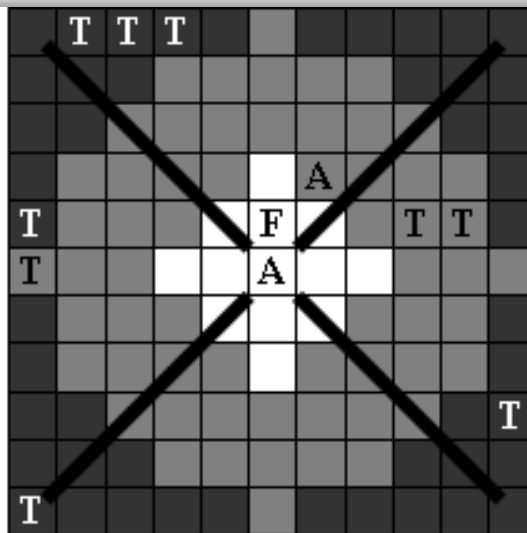
	N	E	S	W
A	0	0	0	0
T	0	0	0	0
F	1	0	0	0
A	1	0	0	0
T	0	1	0	0
F	0	0	0	0



Reward Events "eventXCS"

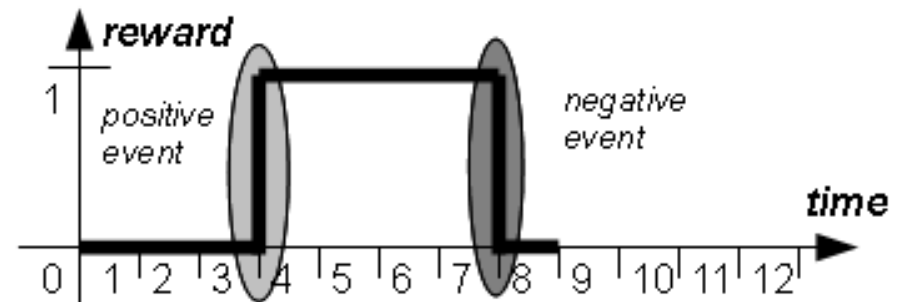
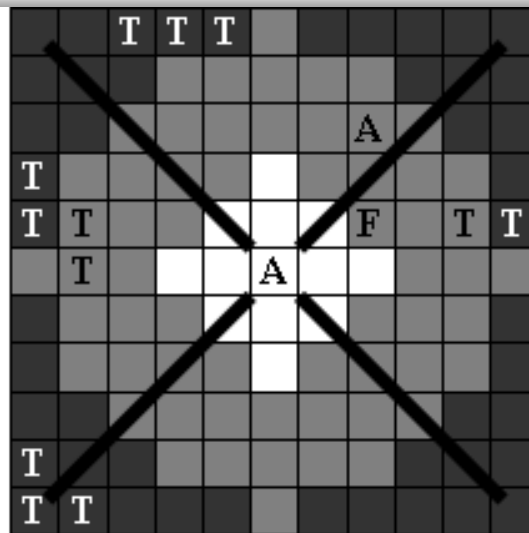
- Move West

	N	E	S	W
A	0	0	0	0
T	0	0	0	0
F	1	0	0	0
A	1	0	0	0
T	0	1	0	1
F	0	0	0	0



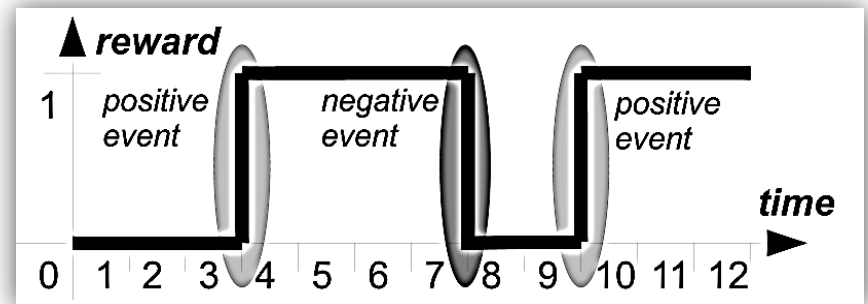
Reward Events "eventXCS"

	N	E	S	W
A	0	0	0	0
T	0	0	0	0
F	0	0	0	0
A	1	0	0	0
T	0	1	0	1
F	0	1	0	0



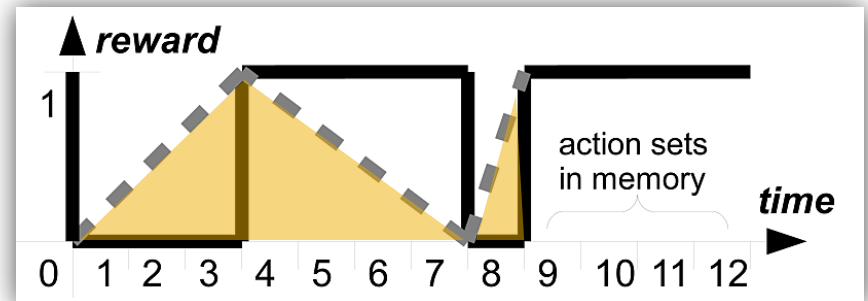
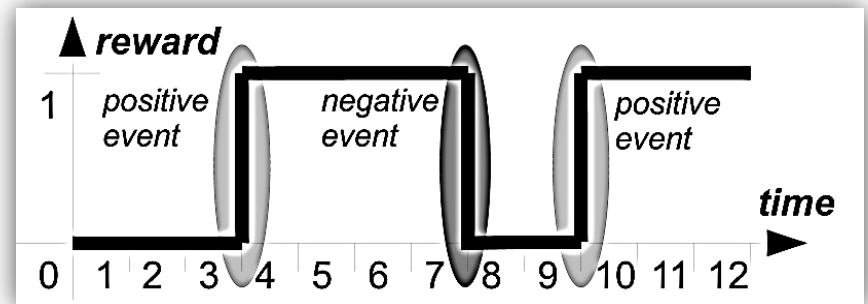
Reward Distribution “eventXCS”

- Analyze succession of positive and negative events

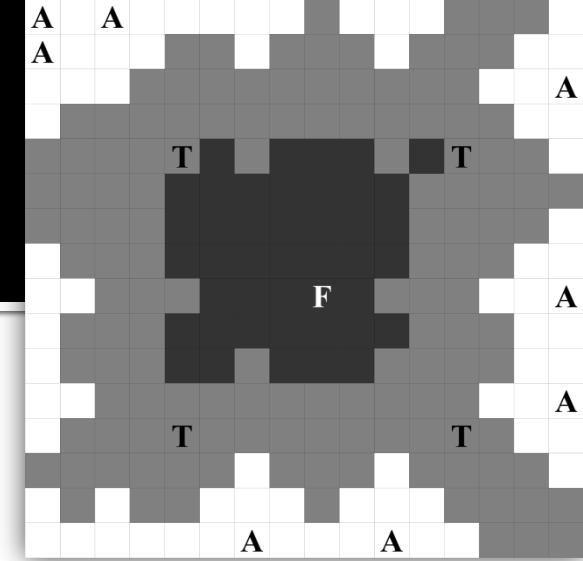


Reward Distribution “eventXCS”

- Analyze succession of positive and negative events
- Distribute the reward as soon as possible (i.e. at each event)
- Idea:
 - Action sets close to an event probably contributed more

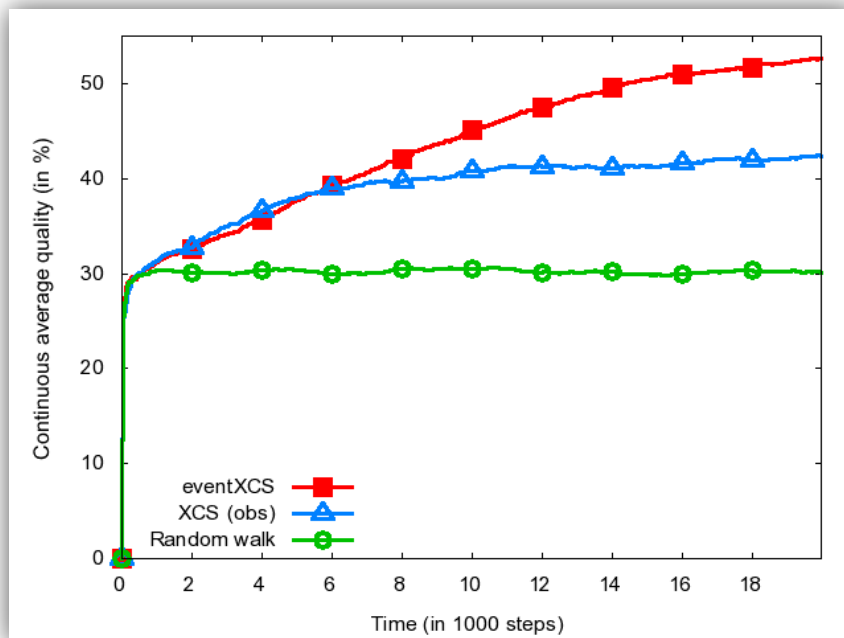


Experimental Results “Pillar Scenario”

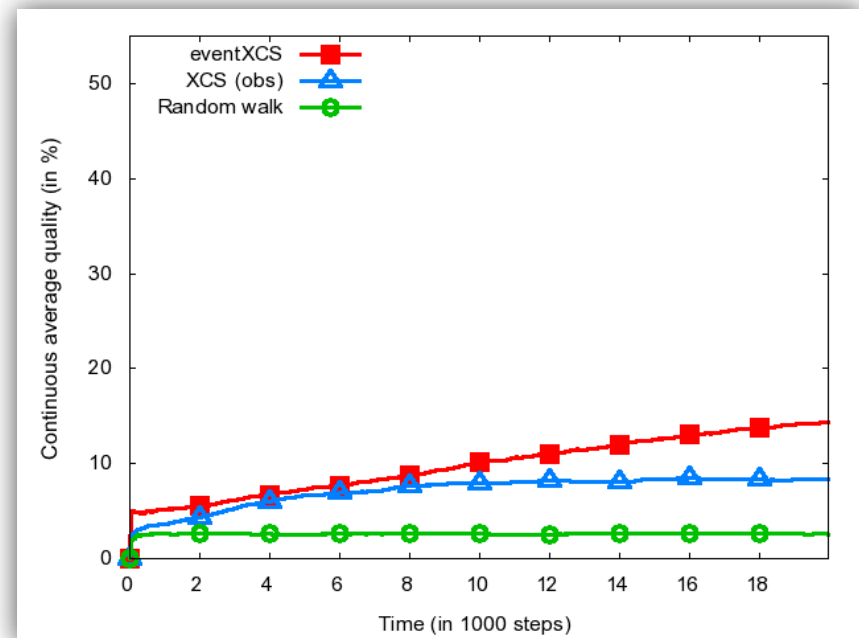


- eventXCS clearly outperforms XCS

Obstacle-evading prey



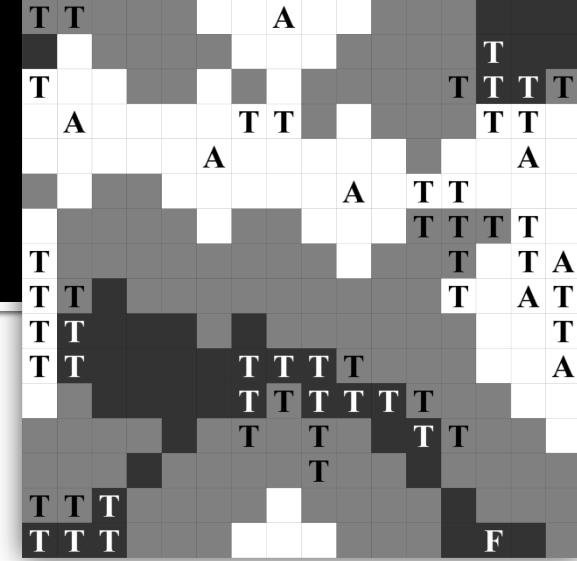
Predator-evading prey



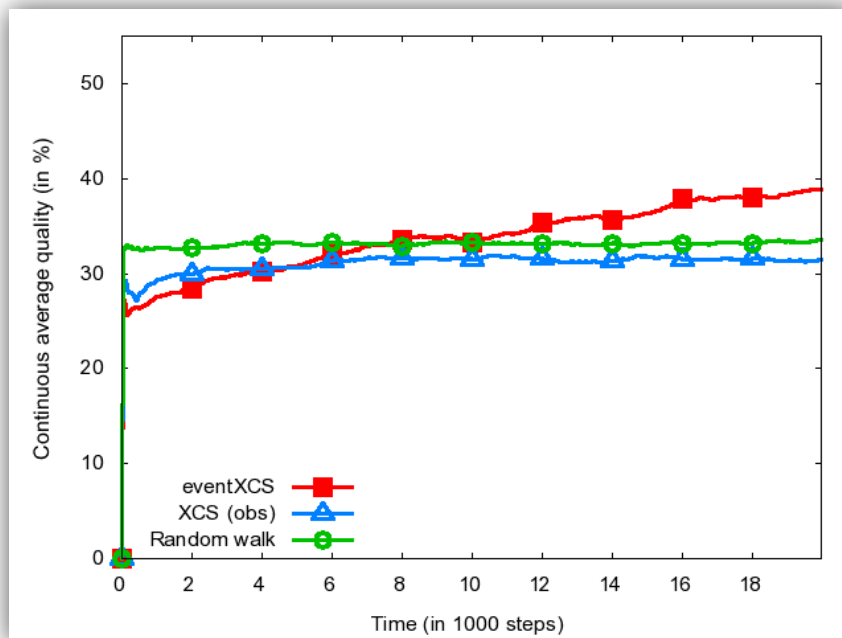
Experimental Results

“Random Scenario”

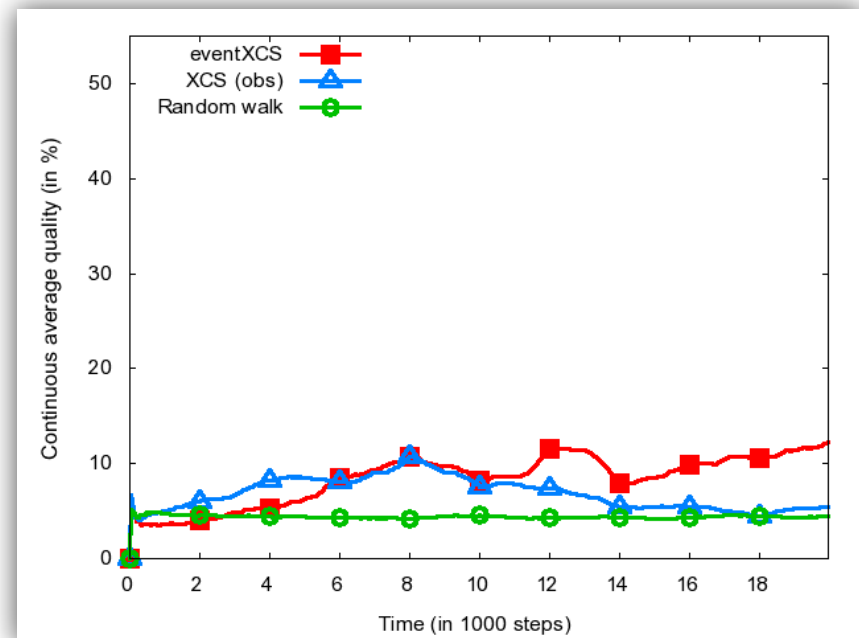
- eventXCS shows slow but steady learning with an obstacle-evading prey



Obstacle-evading prey



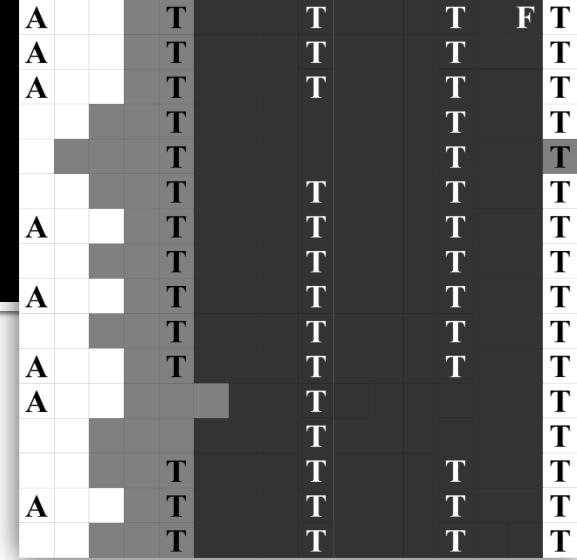
Predator-evading prey



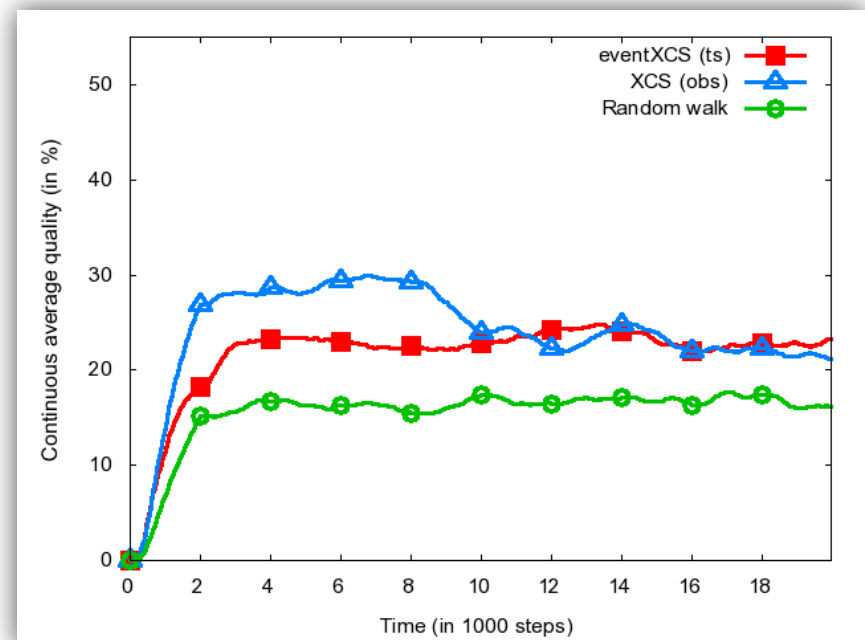
Experimental Results

“Difficult Scenario”

- eventXCS fails in this scenario (not displayed, fitness = ~0)
- Using “tournament selection” shows acceptable results with no sign of unlearning



Blinded prey



Conclusion

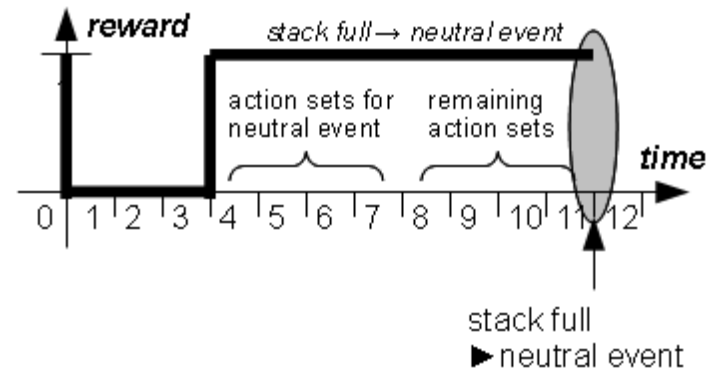


- Predator/Prey scenarios are NOMDP
- XCS can learn (with minimal adaptations) in some P/P scenarios
- Using event handling and reward distribution (eventXCS) much better learning can be observed
- But: Might need some improvement in difficult scenarios

Backup slides

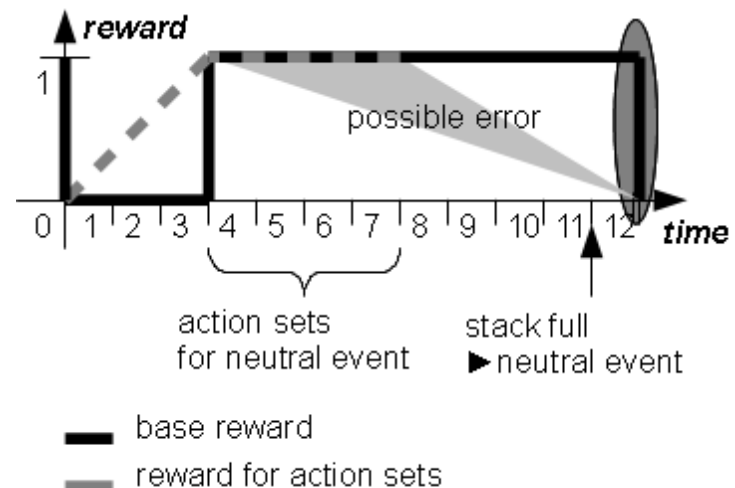
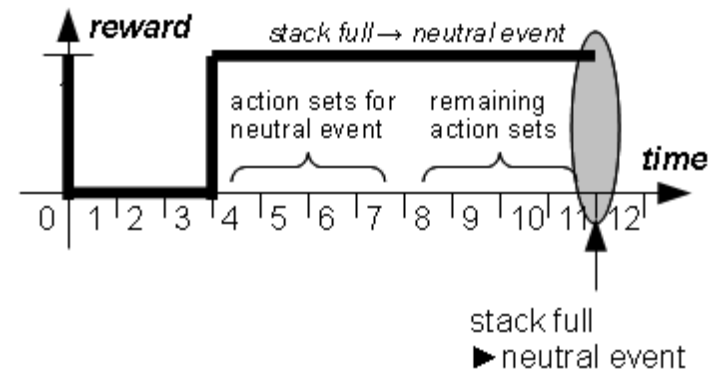
Neutral Events

- Neutral Event
 - No positive or negative event for a number of steps
 - Half of the action sets is discarded and receives reward
 - Idea:
 - Good actions are rewarded earlier
 - Preventing of dead ends



Neutral Events

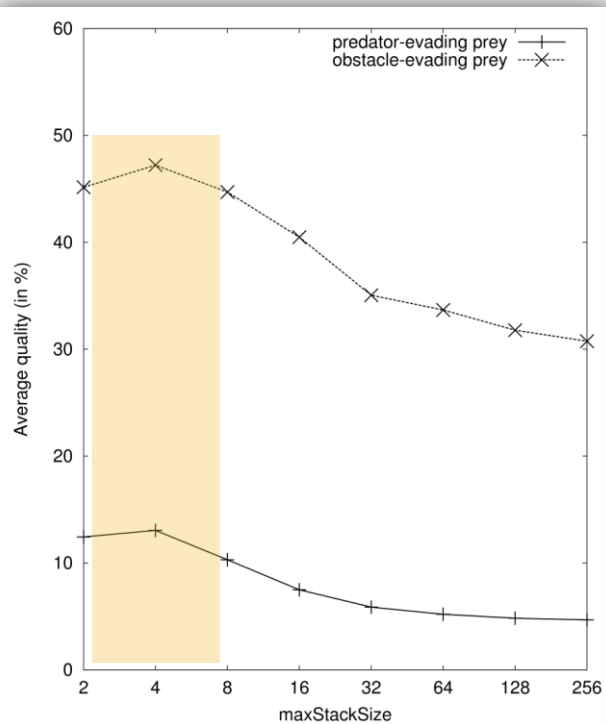
- Neutral Event
 - No positive or negative event for a number of steps
 - Half of the action sets is discarded and receives reward
 - Idea:
 - Good actions are rewarded earlier
 - Preventing of dead ends
- Problem:
 - Error possibility high if directly followed by an event.



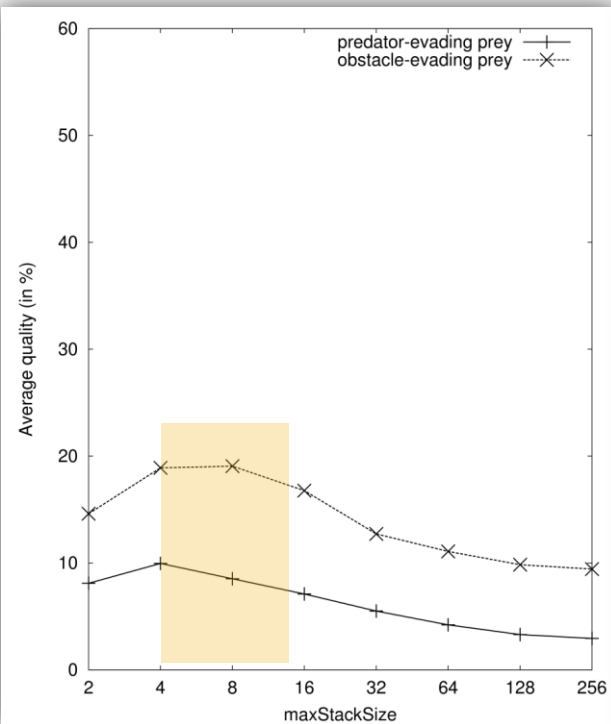
Neutral Events

- Tests have shown that a stack size of 8 is generally good for all three scenarios

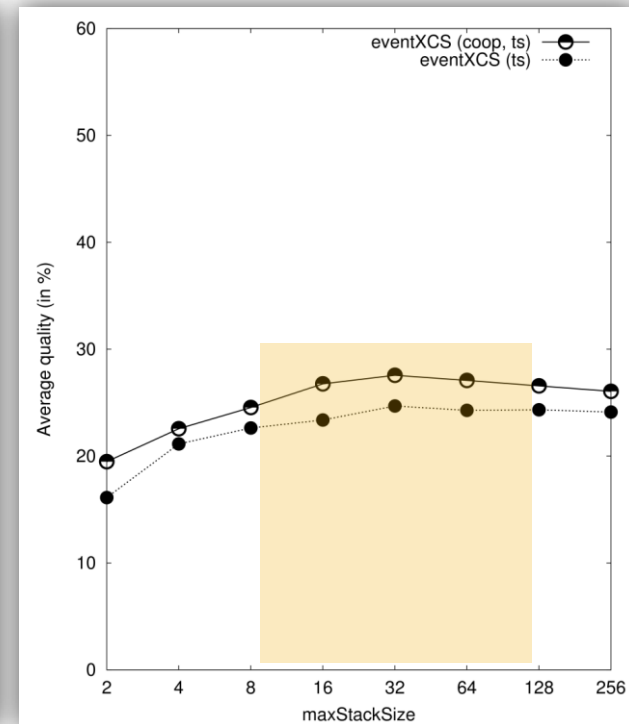
"Pillar Scenario"



"Random Scenario"

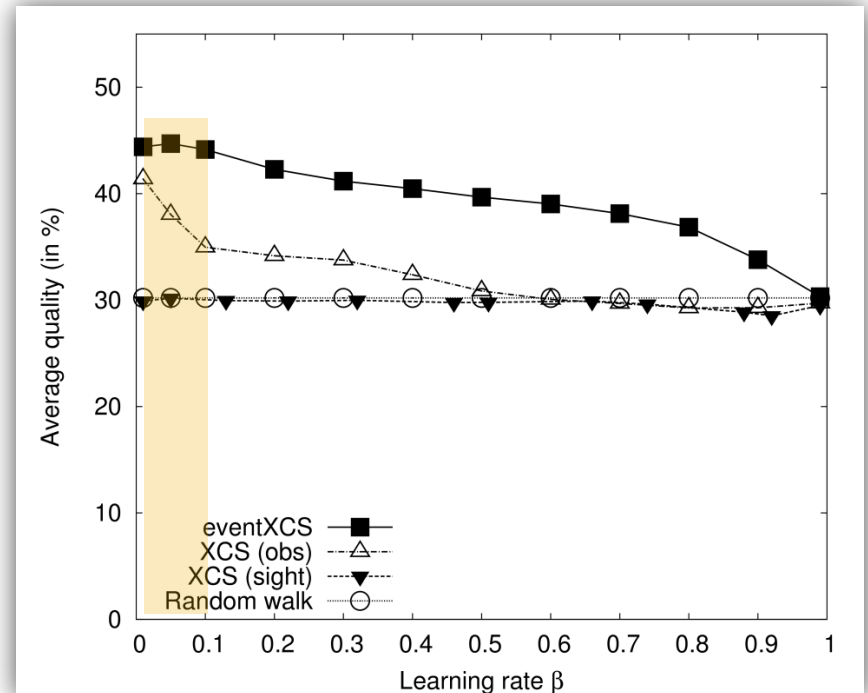


"Difficult Scenario"



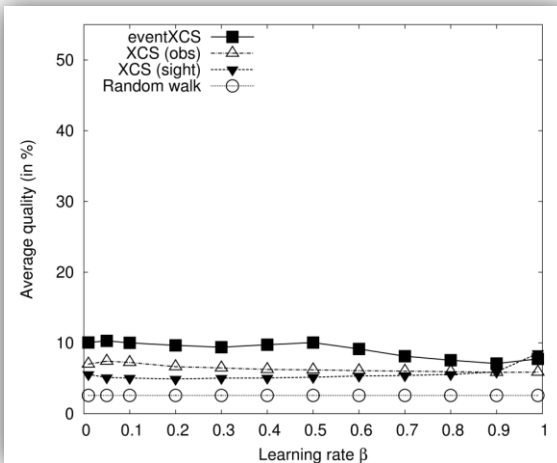
Learning Rate β

- Pillar Scenario
 - Obstacle-evading prey
- Low learning rate (0.05) good, eventXCS very stable

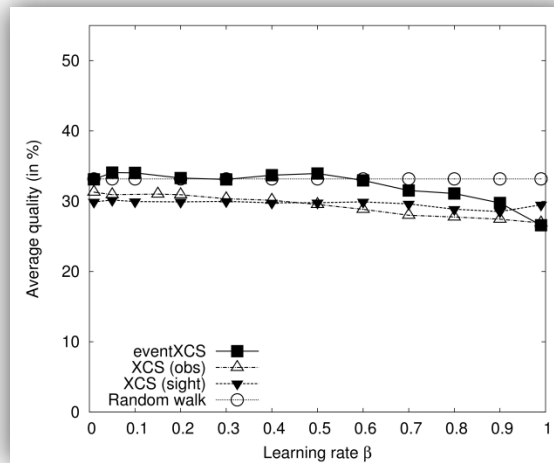


Learning Rate β

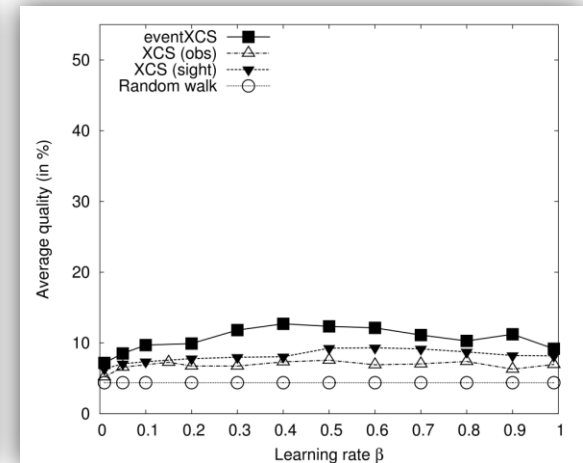
Pillar Scenario
Predator-evading
prey



Random Scenario,
Obstacle evading prey



Random Scenario,
Predator evading



Learning Rate β

- Difficult Scenario
 - Blind prey
- High learning rates show an advantage because of long distance to the prey

