# AAMAS Tutorial

# Solving Games with Complex Strategy Spaces

# Part II: Integrating Learning with Game Theory for

# Strategic Decision Making

Fei Fang

Carnegie Mellon University

feifang@cmu.edu

# Outline

▸ Game-Theoretic Reasoning and Its Applications
  ▸ Wildlife Conservation
  ▸ Cyber Security
  ▸ Ridesharing

▸ End-to-End Learning and Decision Making in Games
  ▸ A differentiable learning framework for learning game parameters

▸ Learning-Powered Strategy Computation in Large Scale Games
  ▸ Leveraging Deep Reinforcement Learning

5/20/2019

# Security Challenges

# Sustainability Challenges



Today
≈ 3,200
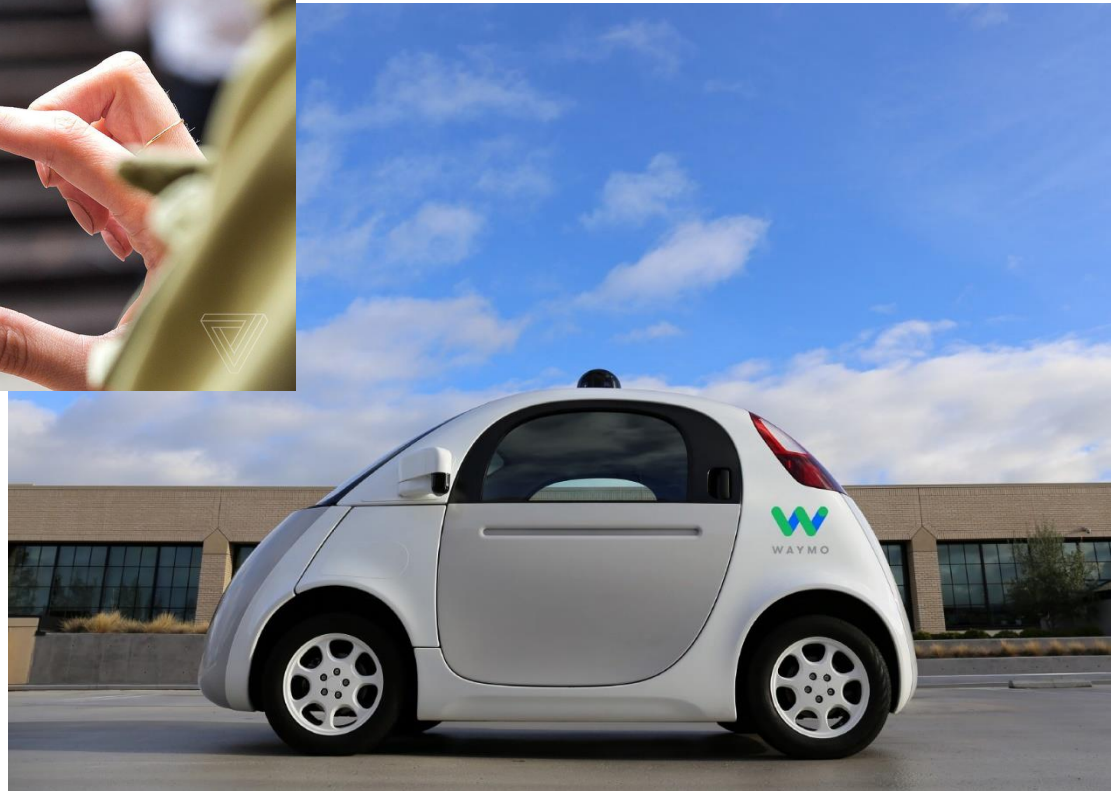
100 years ago
≈ 60,000

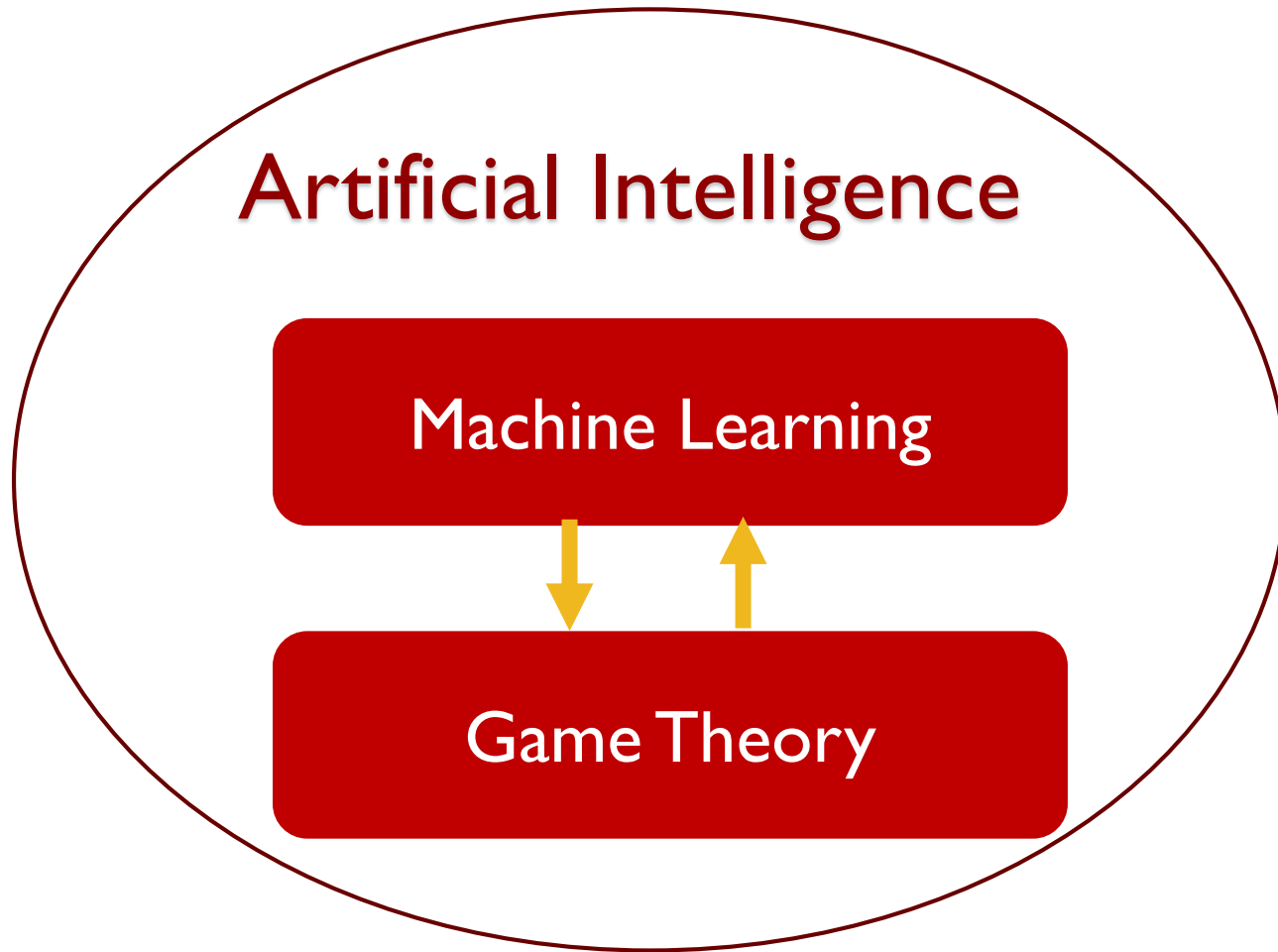# Mobility Challenges

# Societal Challenges

Security & Safety



Environmental Sustainability



Mobility

▸ Strong Stackelberg Equilibrium

  ▸ Defender: mixed strategy

  ▸ Attacker: best response, break tie in favor of defender

**Adversary**

|  | | Target #1 | Target #2 |
|---|---|---|---|
| **Target #1** | 55.6% | 5, -3 | -1, 1 |
| **Target #2** | 44.4% | -5, 4 | 2, -1 |

**Defender**

# Quiz

▸ How to get the defender's mixed strategy in SSE in this problem?

**Adversary**

|  |  | Target #1 | Target #2 |
|---|---|---|---|
| 55.6% | **Target #1** | 5, -3 | -1, 1 |
| 44.4% | **Target #2** | -5, 4 | 2, -1 |

**Defender**

# Quiz

▸ How to get the defender's mixed strategy in SSE in this problem?

  ▸ AttEU1=$p * (-3) + (1 - p) * 4 = p * 1 + (1 - p) * (-1)$=AttEU2

  ▸ Equilibrium: DefStrat=(0.556,0.444), AttStrat=(1,0)

**Adversary**

|  |  | Target #1 | Target #2 |
|---|---|---|---|
| 55.6% | Target #1 | 5, -3 | -1, 1 |
| 44.4% | Target #2 | -5, 4 | 2, -1 |

**Defender**

# Recap: SSE vs NE

▸ Zero-sum
  ▸ SSE=NE=minimax=maximin
  ▸ Approach 1: Single LP (minimax or maximin strategy)
  ▸ Approach 2: Greedy allocation for security games

▸ General-sum
  ▸ SSE≥NE
  ▸ Computing NE: PPAD Complete, LCP (linear complementarity problem) formulation, Gambit solver
  ▸ Computing SSE
    ▸ Approach 1: Multiple LPs (each solve a subproblem)
    ▸ Approach 2: A single MILP that combines all the LPs
    ▸ Approach 3: Extended greedy allocation algorithm $O(n \log n)$ for security games

# Example: Protecting Staten Island Ferry

Optimal Patrol Strategy for Protecting Moving Targets with Multiple Mobile Resources. Fei Fang, Albert Xin Jiang, Milind Tambe. In AAMAS-13

5/20/2019

# Outline

- **Game-Theoretic Reasoning and Its Applications**
  - Wildlife Conservation
  - Cyber Security
  - Ridesharing

- **End-to-End Learning and Decision Making in Games**
  - A differentiable learning framework for learning game parameters

- **Learning-Powered Strategy Computation in Large Scale Games**
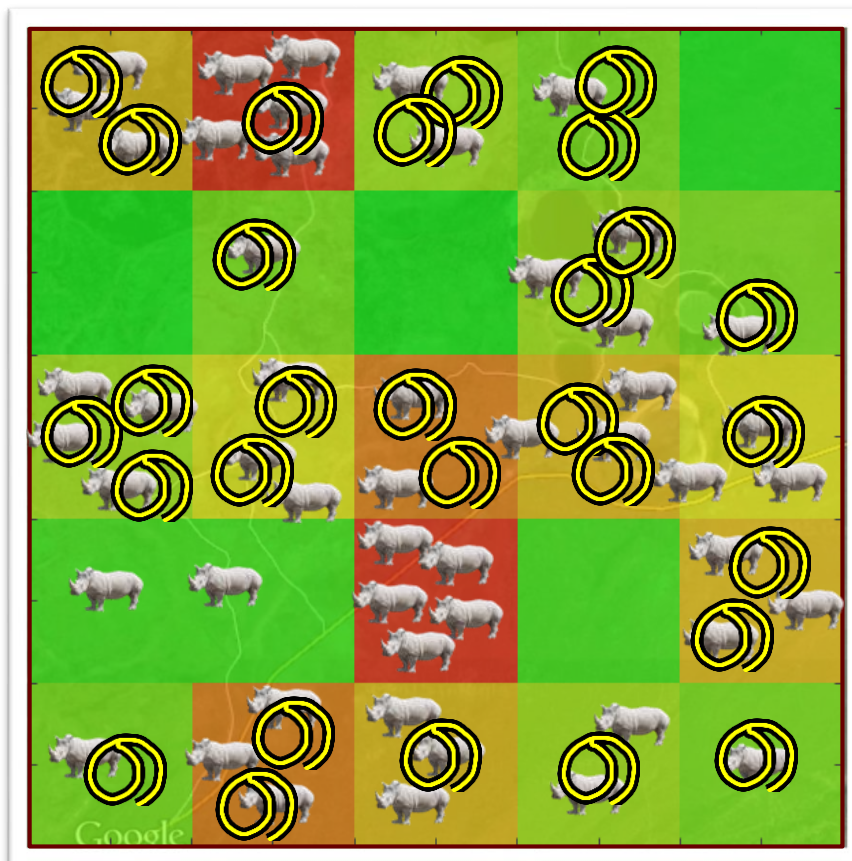  - Leveraging Deep Reinforcement Learning

# Wildlife Conservation



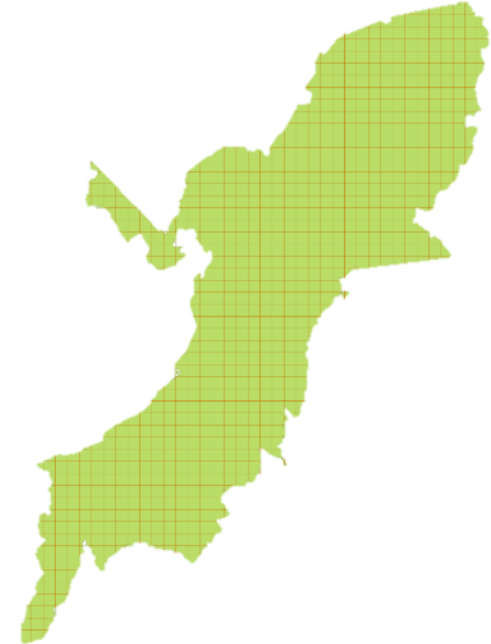Data SIO, NOAA, U.S. Navy, NGA, GEBCO
Image Landsat
Image IBCAO

Google earth

# Human Behavior in Games

- Not always perfectly rational or behave as expected!
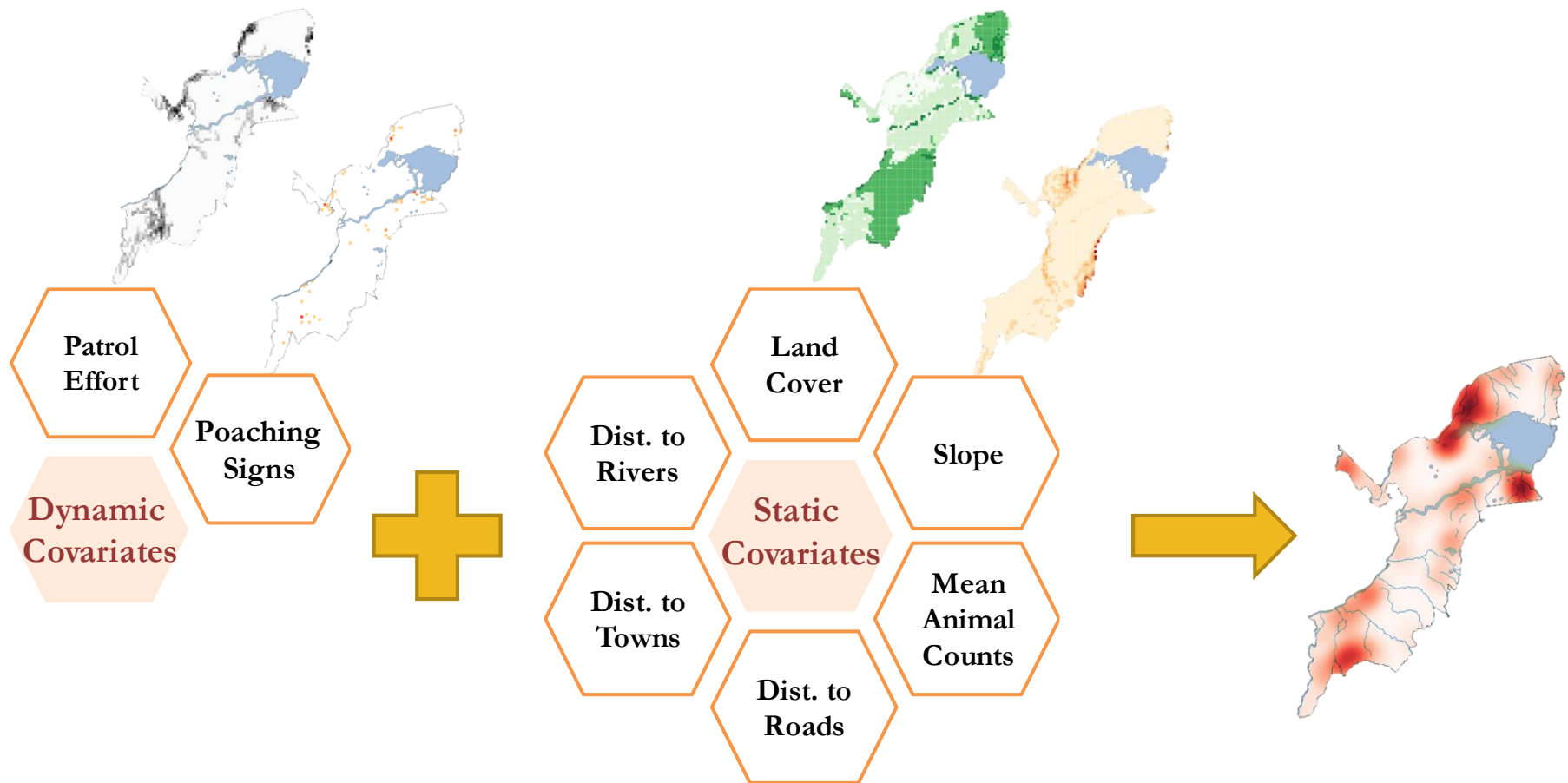- Task: Predict where the poachers place snares

# Learn from Real-World Data



▸ Raw Dataset for Queen Elizabeth National Park
  ▸ Covers 2520 sq. km
  ▸ Patrol and poaching recorded

# Learn from Real-World Data



Each data point represent a 1km×1km area in a season

# Challenge 2: Lack of Recorded Attacks



**Patrolled Cells**
(Year)

| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |
|------|------|------|------|------|------|
| 61.7 | 67.7 | 66.0 | 54.6 | 59.0 | 61.0 |

Per 100 cells

**Not Attacked Patrolled Cells**

| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |
|------|------|------|------|------|------|
| 86.0 | 91.9 | 89.5 | 92.3 | 93.0 | 88.4 |

Per 100 cells

**Attacked Patrolled Cells**

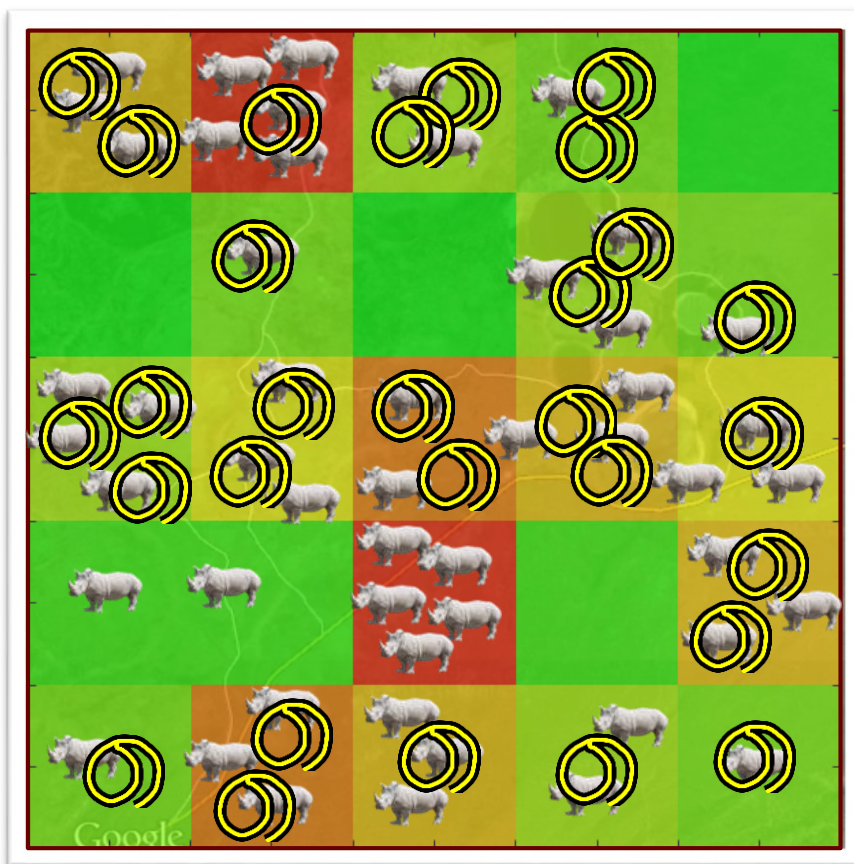| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |
|------|------|------|------|------|------|
| 14 | 8.1 | 10.5 | 7.7 | 7 | 11.6 |

Per 100 cells

# Quantal Response Model

▸ Classical model in behavioral game theory

▸ Probability of attacking target $j$

$$q_j = \frac{e^{\lambda * \text{AttEU}_j(x)}}{\sum_i e^{\lambda * \text{AttEU}_i(x)}}$$

▸ λ: represents error level (=0 means uniform random)

   ▸ Maximal likelihood estimation (λ=0.76)

   ▸ $\max_\lambda f(\lambda) = \sum_j N_j \log(q_j)$

   ▸ Solved through gradient ascent $\lambda \leftarrow \lambda + \alpha \nabla_\lambda f(\lambda)$

 McKelvey, R. D., & Palfrey, T. R. (1995). Quantal response equilibria for normal form games. Games and economic behavior, 10(1), 6-38. 5/20/2019

# Subjective Utility Quantal Response Model

$$\text{SEU}_j = \sum_k w_k f_j^k, \quad q_j = \frac{e^{\lambda * \text{SEU}_j(x)}}{\sum_i e^{\lambda * \text{SEU}_i(x)}}$$



Past Success/Failure
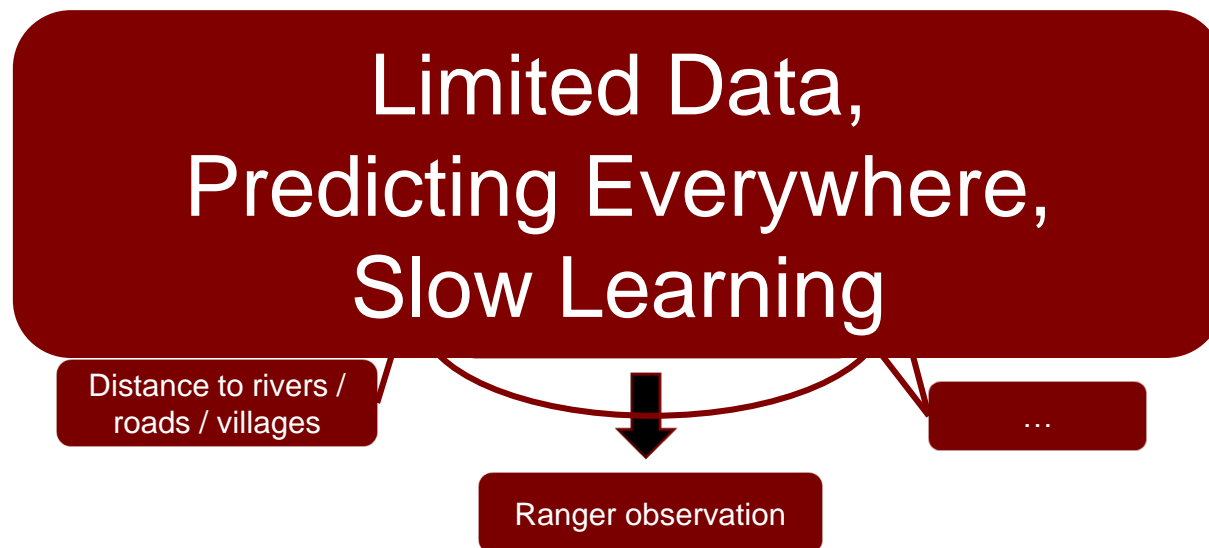Induced Features +

Coverage Probability
+ Reward/Penalty

SUQR

Attack Probability

Nguyen, T. H., Yang, R., Azaria, A., Kraus, S., & Tambe, M. Analyzing the Effectiveness of Adversary Modeling in Security Games. In AAAI, 2013.

▶ # CAPTURE

▸ Real-world Data

▸ Dynamic Bayes Net: Time Dependency & Imperfect Observation



Limited Data,
Predicting Everywhere,
Slow Learning

Distance to rivers /
roads / villages

...

Ranger observation

Thanh H. Nguyen, Arunesh Sinha, Shahrzad Gholami, Andrew Plumptre, Lucas Joppa, Milind Tambe, Margaret Driciru, Fred Wanyama, Aggrey Rwetsiba, Rob Critchlow, Colin Beale. CAPTURE: A New Predictive Anti-Poaching Tool for Wildlife Protection. In AAMAS, 2016.

# Decision Tree

▶ PROS

   ▶ High speed

   ▶ Learn global poachers behavior

   ▶ Learn nonlinearity in geo-spatial predictor

▶ CONS

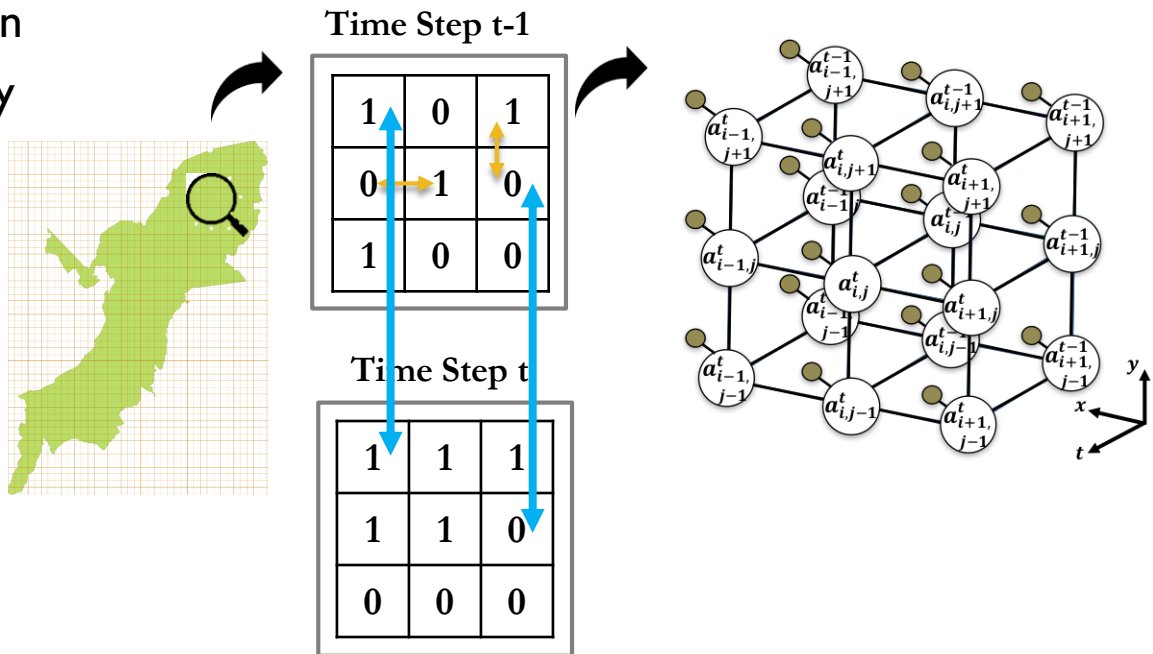   ▶ No explicit temporal dimension

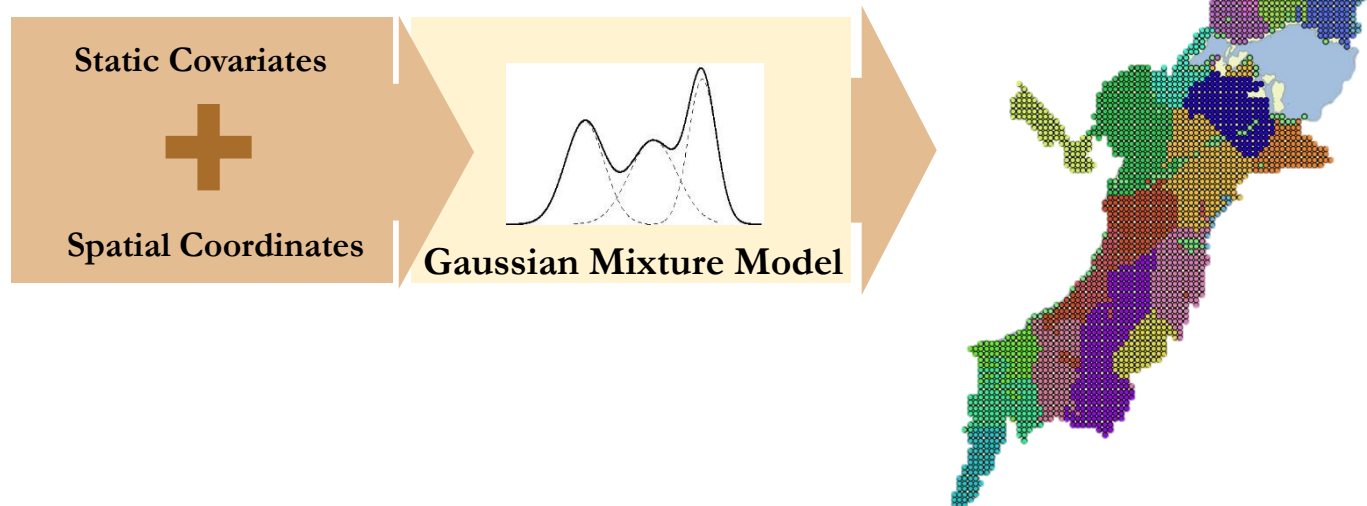   ▶ No aspect for label uncertainty

- **PROS**
  - Explicit spatial dimension
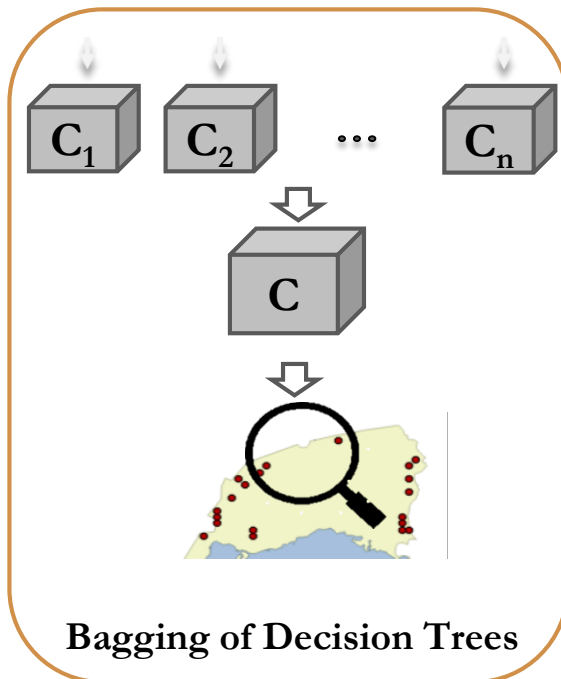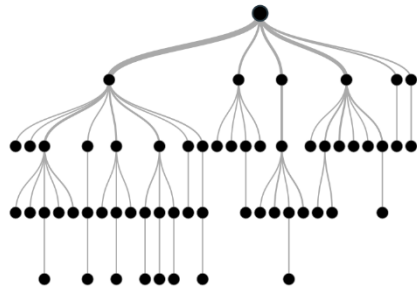  - Explicit temporal dimension
  - Addresses label uncertainty
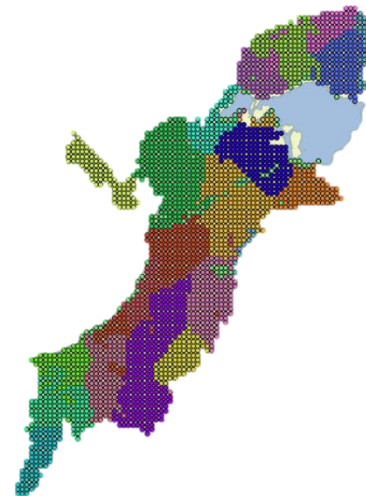
- **CONS**
  - Low speed
  - Data greedy



Time Step t-1

| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |

Time Step t

| 1 | 1 | 1 |
| 1 | 1 | 0 |
| 0 | 0 | 0 |

# Our Solution: Hybrid Model

**Static Covariates**

**+**

**Spatial Coordinates**

**Gaussian Mixture Model**

Geo-clusters

# Our Solution: Hybrid Model



**Bagging of Decision Trees**

On Intensely Monitored Regions

**Markov Random Fields**

Decision Tree
+
Markov Random Fields

Taking it for a Test Drive: A Hybrid Spatio-temporal Model for Wildlife Poaching Prediction Evaluated through a Controlled Field Test. Shahrzad Gholami, Benjamin Ford, Fei Fang, Andrew Plumptre, Milind Tambe, Margaret Driciru, Fred Wanyama, Aggrey Rwetsiba, Mustapha Nsubaga, Joshua Mabonga. In ECML-PKDD 2017

# Augment Dataset With Expert Knowledge



| Cluster Index | Cluster Value in "predict_heatmap_i10.tif" | Estimated Snaring Threat (0~10) |
|---|---|---|
| 1 | 1 | 7 |
| 2 | 2 | 3 |
| 3 | 3 | 8 |
| 4 | 4 | 7 |
| 5 | 5 | 3 |
| 6 | 6 | 2 |
| 7 | 7 | 8 |
| 8 | 8 | 3 |
| 9 | 9 | 0 |
| 10 | 10 | 0 |

▶ Negative sampling: sample from unpatrolled regions

▶ Positive sampling: Estimate from rangers' estimated scores

  ▸ Collect answers for several sets of clusters $C^1, C^2$

  ▸ Compute aggregated score a $s = \min\{s_1(C_i^1), s2(C_j^1), \dots\}$, add unlabeled points as positive points if $s \geq 6$

Exploiting Data and Human Knowledge for Predicting Wildlife Poaching. Swaminathan Gurumurthy, Lantao Yu, Chenyan Zhang, Yongchao Jin, Weiping Li, Xiaodong Zhang, Fei Fang. In COMPASS-18

5/20/2019

- ## Trespassing
  - ### 19 signs of litter, ashes, etc.
- ## Poached animals
  - ### 1 poached elephant
- ## Snaring
  - ### 1 active snare
  - ### 1 cache of 10 antelope snares
  - ### 1 roll of elephant snares
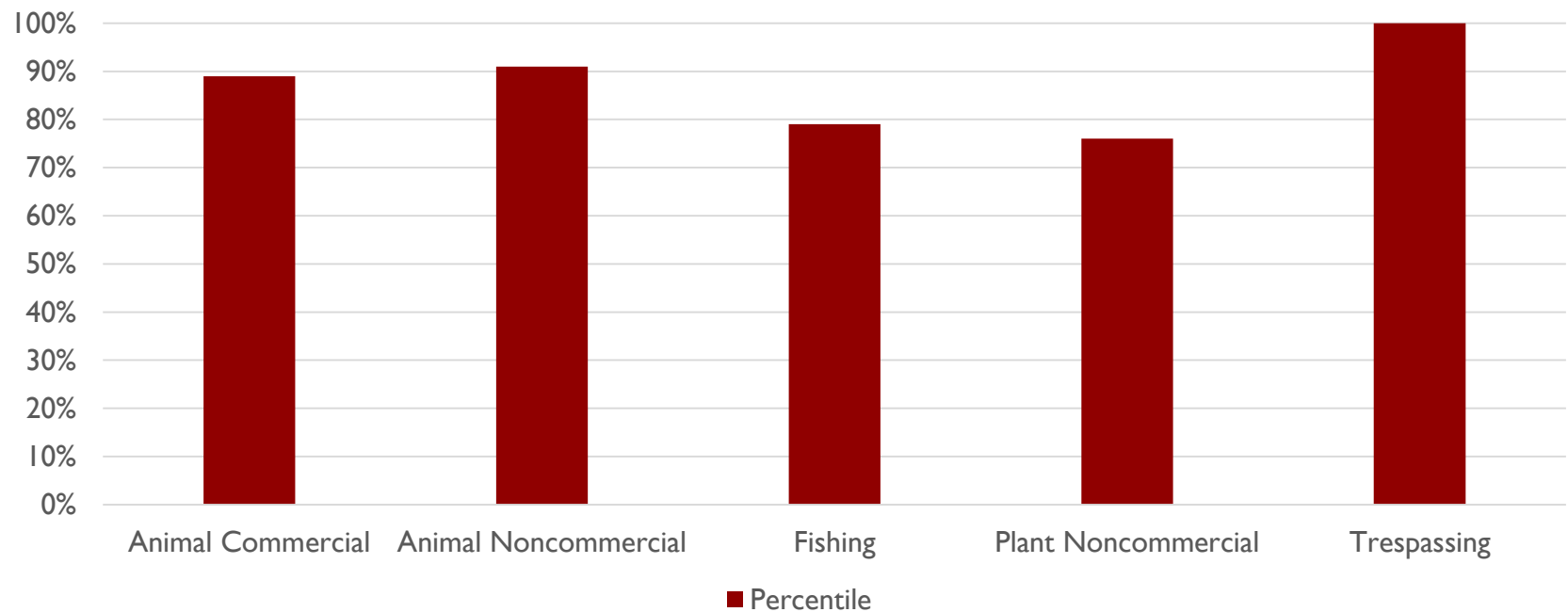- ## Snaring hit rates
  - ### Outperform 91% of months

| Historical Base Hit Rate | Our Hit Rate |
|---|---|
| Average: 0.73 | 3 |

Cloudy with a Chance of Poaching: Adversary Behavior Modeling and Forecasting with Real-World Poaching Data. Debarun Kar, Benjamin Ford, Shahrzad Gholami, Fei Fang, Andrew Plumptre, Milind Tambe, Margaret Driciru, Fred Wanyama, Aggrey Rwetsiba. In AAMAS-17

# Field Test 1 in Uganda: % Months Exceeded Historical

# Field Test 2 in Uganda (8 months)

▸ 27 areas (9-sq km each)

▸ 454 km patrolled in total

▸ No point > 5 km from patrol post

▸ No area patrolled too much/rarely

▸ No overlapping areas

▸ <= 2 areas per patrol post

0  7.5  15 km

▶ 2 experiment groups

  ▶ 1: >= 50% attack
     prediction rate

     ▸ 5 areas

  ▶ 2: < 50% attack
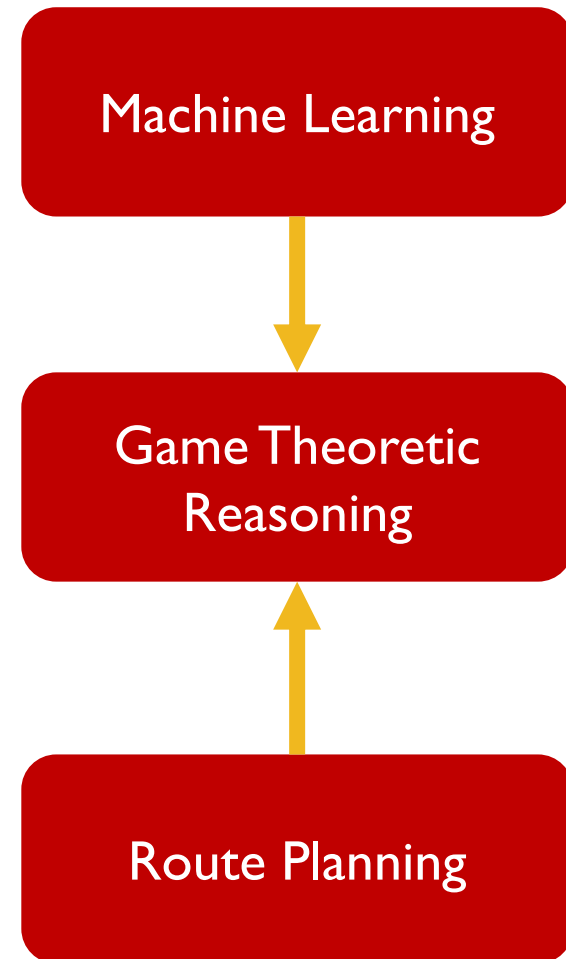     prediction rate

     ▸ 22 areas

▶ Catch Per Unit Effort
  (CPUE)

  ▶ Unit Effort = km walked

# Field Test in China

▸ Two-day field test in October 2017: 22 snares

▸ 34 patrols from November 2017 to February 2018
  ▸ 7 snares

# From Prediction to Prescription

# Game Theoretic Reasoning Based on Learned Model

▸ Find optimal patrol strategy given poachers respond to the patrol strategy according to learned model

▸ Challenges

- ▸ Learned model is hard to represent using closed form function (e.g., decision tree)

- ▸ Hard to scale up when considering scheduling constraints

# Game Theoretic Reasoning Based on Learned Model

▸ Input: A machine learning model that predicts snares

▸ Output: an optimal patrolling strategy

▸ Goal: maximize catches of snares

Optimal Patrol Planning for Green Security Games with Black-Box Attackers. Haifeng Xu, Benjamin Ford, Fei Fang, Bistra Dilkina, Andrew Plumptre, Milind Tambe, Margaret Driciru, Fred Wanyama, Aggrey Rwetsiba, Mustapha Nsubaga, Joshua Mabonga. In GameSec-17: The 8th Conference on Decision and Game Theory for Security

# Game Theoretic Reasoning Based on Learned Model

**For each cell $i$:**

$x_i$: Current patrol effort at $i$

$g_i$

$y_i$: Prob. of detecting a snare at $i$ in current period

▸ Optimization problem: $\max\limits_{x_i} \sum_i g_i(x_i)$

▸ However…



**Patrol post**
(one patroller)

# Game Theoretic Reasoning

- Observe: a pure strategy = a path from $v_{11}$ to $v_{1T}$

- Claim: a mixed strategy $\iff$ one-unit fractional flow from $v_{11}$ to $v_{1T}$

- Patrol effort at cell $i$ = the aggregated flow through cell $i$

- Build a mixed integer linear program

cells

$(N =)8$
$7$
$6$
$5$
$4$
$3$
$2$
$1$

$1$  $2$  $3$  $4$  $5$  $6(= T)$   time

$v_{11}$              $v_{1T}$

Time-unrolled Graph

▸ A MILP formulation

$$\approx \max_{x_i} \sum_i g_i(x_i)$$

$$\text{maximize } \sum_{i=1}^N \left( g_i(0) + \sum_{j=1}^m z_i^j \cdot [g_i(j) - g_i(j-1)] \right)$$

$$\text{subject to } x_i \geq \sum_{j=1}^m z_i^j \cdot [\alpha_j - \alpha_{j-1}],$$

$$x_i \leq \alpha_1 + \sum_{j=1}^m z_i^j \cdot [\alpha_{j+1} - \alpha_j],$$

$$z_i^1 \geq z_i^2 ... \geq z_i^m,$$

$$z_i^j \in \{0, 1\},$$

$$x_i = z_i^1 + z_i^2 + \cdots$$

$$x_i = \sum_{t=1}^T \left[ \sum_{e \in \sigma^+(v_{t,i})} f(e) \right],$$

Patrol effort at cell $i$ = the aggregated flow through cell $i$

$$\sum_{e \in \sigma^+(v_{t,i})} f(e) = \sum_{e \in \sigma^-(v_{t,i})} f(e),$$

$$\sum_{e \in \sigma^+(v_{T,1})} f(e) = \sum_{e \in \sigma^-(v_{1,1})} f(e) = 1$$

$f$ is a unit flow

$$0 \leq x_i \leq 1, \qquad 0 \leq f(e) \leq 1,$$

# Complex Terrain



Patrol Route (2D)

Patrol Route (3D)

▸ 8-hour patrol in April 2015: patrolling is not easy!

# Spatial Constraint

Fei Fang

# Spatial Constraint

▸ Grid based → Route based

▸ Hierarchical modeling: Focus on terrain features
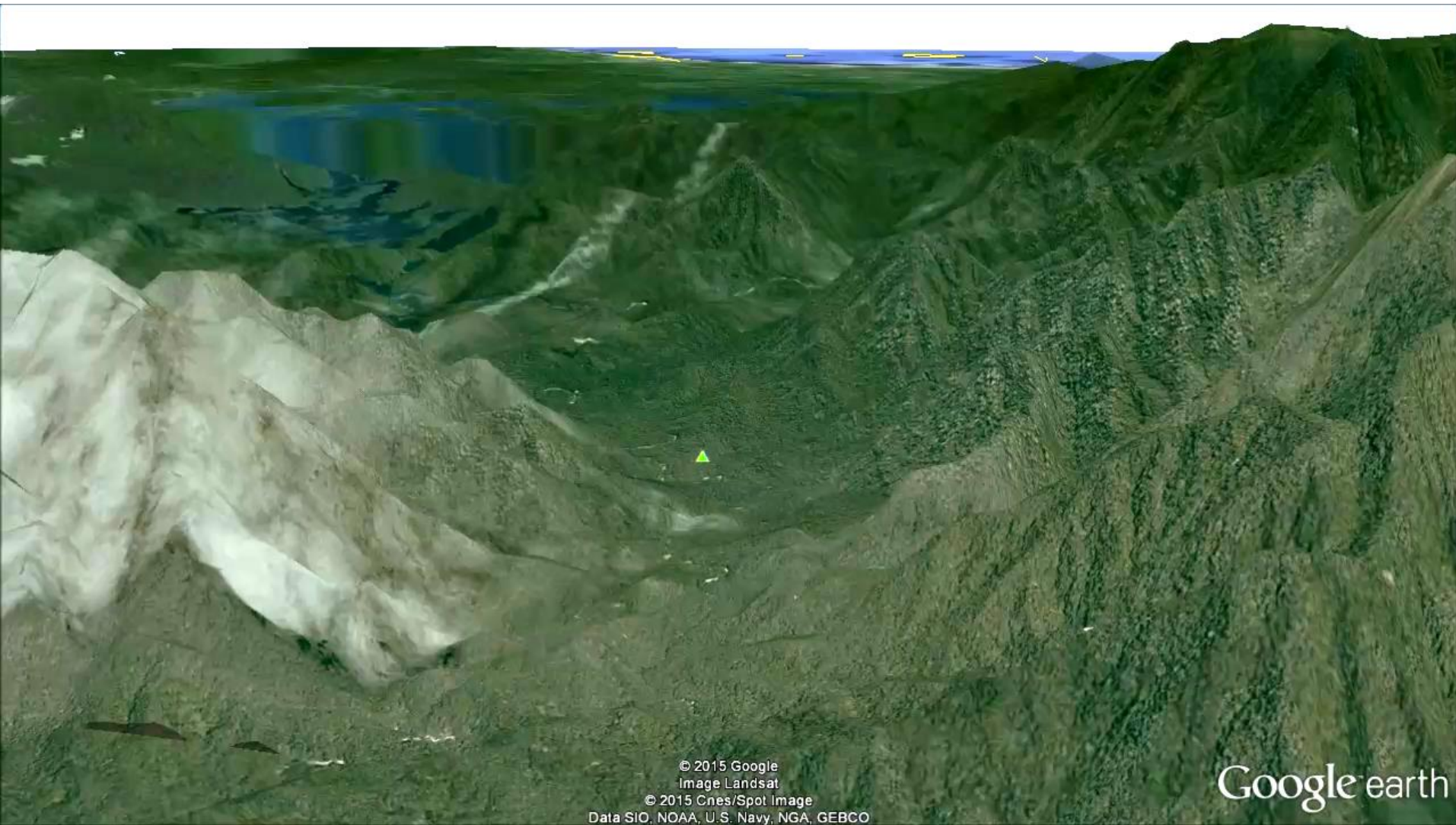
▸ Build virtual street map
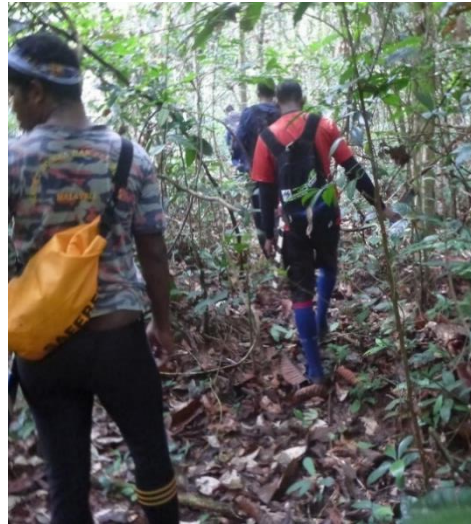
▶ Hierarchical model: Focus on terrain feature



— Ridgeline

— Stream

— Street Map

— Patrol Route

# Patrol Route Design

Deploying PAWS: Field Optimization of the Protection Assistant for Wildlife Security. Fei Fang, Thanh H. Nguyen, Rob Pickles, Wai Y. Lam, Gopalasamy R. Clements, Bo An, Amandeep Singh, Milind Tambe, Andrew Lemieux. In IAAI-16

# Field Test in Malaysia

- In collaboration with Panthera, Rimba
- Regular deployment since July 2015 (Malaysia)



Fei Fang

# Real-World Deployment

Grid Based

Route Based

Fei Fang

# Real-World Deployment

Animal Footprint

Tree Mark

Camping Sign

Tiger Sign

Lighter

# Real-World Deployment

▸ ## PAWS is deployed in the field

  ▸ ### Saved animals!

# Outline

▸ **Game-Theoretic Reasoning and Its Applications**

  ▸ Wildlife Conservation

  ▸ Cyber Security

  ▸ Ridesharing

▸ **End-to-End Learning and Decision Making in Games**

  ▸ A differentiable learning framework for learning game parameters

▸ **Learning-Powered Strategy Computation in Large Scale Games**

  ▸ Leveraging Deep Reinforcement Learning

# How Valuable is This Car?

# Deception

# Deception

# Cyber Deception

▸ What can the defender do without "patrol boats"?

▸ Use deception to confuse the attackers!



Send probes to systems to gather information

Attacker

Give information about systems on network

Enterprise Network

# Cyber Deception

▸ How should the defender disguise the systems to induce the adversary to attack the least valuable systems?

▸ Cyber Domain Challenges:

  ▸ Intelligent adversary; could perceive deception occurring

  ▸ Large number of system configurations and ways to disguise

  ▸ Arbitrary deception may not be feasible or may affect performance

▶ $K$ systems, each has true configuration (TC) $f \in F$

▶ Successful attack on system with TC $f$ yields utility $U_f$ to attacker; defender loses $U_f$ (gains $-U_f$)



$f_2 : 0$
Linux,
NGINX 1.10

$f_1 : 10$
Windows,
Apache 2.2

$f_3 : 5$
Linux,
NGINX 1.15

$f_2 : 0$
Linux,
NGINX 1.10

# Cyber Deception Game: Setting

▸ Defender disguise the systems through deceptive responses

▸ Each system gets observed configuration (OC) $\tilde{f} \in \tilde{F}$

$\boxed{\widetilde{f_2}}$
Linux, Apache 2.4

$f_2$
Linux, NGINX 1.10

$\boxed{\widetilde{f_1}}$
Windows, Apache 2.2

$f_1 : 10$
Windows, Apache 2.2

$f_3 : 5$
Linux, NGINX 1.15

$\boxed{\widetilde{f_2}}$
Linux, Apache 2.4

$\boxed{\widetilde{f_2}}$
Linux, Apache 2.4

$f_2 : 0$
Linux, NGINX 1.10

# Cyber Deception Game: Defender

▸ Know true configuration (TC) $f$

▸ Need to decides observed configuration (OC) $\tilde{f}$

▸ Systems with same TC are indifferent to the defender

▸ $N_f$ = Number of systems having TC $f \in F$



$f_1: U_{f_1} = 10$

$f_2: U_{f_2} = 0$       $\tilde{f}_1$

$f_2: U_{f_2} = 0$       $\tilde{f}_2$

$f_3: U_{f_3} = 5$

- Deception strategy encoded via integer matrix $\phi$
  - $\phi_{f,\tilde{f}}$ = number of systems with TC $f$ and OC $\tilde{f}$



| $\phi$ | $\tilde{f}_1$ | $\tilde{f}_2$ |
|--------|---------------|---------------|
| $f_1$  | 1             | 0             |
| $f_2$  | 2             | 0             |
| $f_3$  | 0             | 1             |

# Cyber Deception Game: Defender

▶ Deception strategy encoded via integer matrix $\phi$

  ▶ $\phi_{f,\tilde{f}}$ = number of systems with TC $f$ and OC $\tilde{f}$

  ▶ TC $f$ may not be masked with OC $\tilde{f}$ ($\pi_{f,\tilde{f}} = 0$)

  ▶ Showing deceptive responses incur costs $c(f, \tilde{f})$; budget $B$



$$f_1: U_{f_1} = 10$$

$$\pi_{f_1,\tilde{f}_1} = 0$$

$$f_2: U_{f_2} = 0 \quad 1$$

$$\tilde{f}_1$$

$$B = 5$$

$$f_2: U_{f_2} = 0 \quad 1$$

$$\tilde{f}_2$$

$$f_3: U_{f_3} = 5 \quad c_{f_3,\widetilde{f_2}} = 3$$

# Cyber Deception Game: Attacker

▸ Can observe OC of each system

▸ Cannot differentiate systems with same OC

▸ Uniformly randomly attacks systems with **<u>most attractive</u>** OC

⬇

**How much does the attacker know about the deception?**

▶ Powerful attacker: Knows deception strategy $\phi$

  ▸ Computes expected payoff for all OCs and best-responds

  ▸ Robust assumption to minimize worst-case loss



$$\widetilde{U}_{\tilde{f}} = \frac{\sum_{f \in F} \phi_{f,\tilde{f}} U_f}{\sum_{f \in F} \phi_{f,\tilde{f}}}$$

Expected Payoff

| $\phi$ | $\widetilde{f}_1$ | $\widetilde{f}_2$ |
|---|---|---|
| $f_1$ | 1 | 0 |
| $f_2$ | 2 | 0 |
| $f_3$ | 0 | 1 |

$$\widetilde{U}_{\tilde{f}_1} = \frac{10 + 2 * 0}{3} = 3.33$$

$$\widetilde{U}_{\tilde{f}_2} = 5/1 = 5$$

$f_1: U_{f_1} = 10$

$f_2: U_{f_2} = 0$

$f_2: U_{f_2} = 0$

$f_3: U_{f_3} = 5$

# Cyber Deception Game: Attacker

▸ Powerful attacker: Knows deception strategy $\phi$

  ▸ Computes expected payoff for all OCs and best-responds

  ▸ Robust assumption to minimize worst-case loss

▸ Naive attacker: Not aware of deception

  ▸ Believe what they observe

  ▸ Preset preferences (utilities) for attacking OCs

▸ With powerful attacker, when there are no budget constraint and feasibility constraint, what is the optimal defender strategy?

# Quiz

▸ With powerful attacker, when there are no budget constraint and feasibility constraint, what is the optimal defender strategy?

▸ Trivial case (no constraints): assign to same OC

# Against Powerful Attacker

▶ Powerful attacker: Knows deception strategy $\phi$

- ▸ Computes expected payoff for all OCs and best-responds
- ▸ Robust assumption to minimize worst-case loss

▶ When some masking infeasible or budget limited

> <u>Theorem</u>: NP-hard to compute optimal strategy for defender against powerful adversary.

- ▸ Proven via reduction to Partition problem
- ▸ NP-hard even with just feasibility or just budget constraint

▸ Solve through mathematical programming

$$\min_{u,\phi} \quad u$$

Non-linear

$$s.t. \quad u \geq \boxed{\frac{\sum_{f \in F} \phi_{f,\tilde{f}} U_f}{\sum_{f \in F} \phi_{f,\tilde{f}}}} \quad \forall \tilde{f} \in \tilde{F} \quad \boxed{\begin{array}{c} \text{Expected Utility} \\ \text{for attacking } \tilde{f} \end{array}}$$

$$\sum_{\tilde{f}} \phi_{f,\tilde{f}} = N_f$$

$$\sum_{f} \phi_{f,\tilde{f}} = N_{\tilde{f}}$$

**Feasibility Constraints**

$$\phi_{f,\tilde{f}} \leq \pi_{f,\tilde{f}}$$

$$\phi_{f,\tilde{f}} \in \mathbb{Z}_{\geq 0}$$

$$\sum_{f} \sum_{\tilde{f}} \phi_{f,\tilde{f}} c_{f,\tilde{f}} \leq B$$

**Budget Constraint**

# Against Powerful Attacker

▸ Solve through mathematical programming

▸ Reformulate to MILP: Guaranteed to find optimal solution

> ▸ Remove the non-linear constraint
>
> ▸ Adds $|K||\widetilde{F}|$ auxiliary variables
>
> ▸ Adds $4|K||\widetilde{F}|$ additional constraints

▸ Approximation algorithm: Solve sequential MILPs

▸ Heuristic algorithm: Greedy MiniMax (GMM)

> ▸ A fast heuristic which greedily minimizes attacker utility

# Against Naïve Attacker

▶ Naive attacker: Not aware of deception

- Simply believes OCs (or just not reasoning about the actual TC→OC mapping strategy used by the defender)
- Preset preferences (utilities) for attacking OCs

▶ When no budget constraints; but just the feasibility constraints

> <u>Theorem</u>: can be solved in $O(|F||\tilde{F}|)$ time

▶ When both budget and feasibility constraints present

> <u>Theorem</u>: NP-hard to compute optimal strategy for defender against naïve adversary.

▸ 20 TCs, 20 Systems

▸ Attacker Utility = 10 without deception

▸ Attacker model and belief of attacker model matters

# Outline

▸ **Game-Theoretic Reasoning and Its Applications**
  - ▸ Wildlife Conservation
  - ▸ Cyber Security
  - ▸ Ridesharing

▸ **End-to-End Learning and Decision Making in Games**
  - ▸ A differentiable learning framework for learning game parameters

▸ **Learning-Powered Strategy Computation in Large Scale Games**
  - ▸ Leveraging Deep Reinforcement Learning

▶ Surge price interface

# Evolution of Surge Pricing

▸ Coarse → Fine grained in space

# Quiz

▸ What are the potential strategic behavior of a driver (with old or new interface)?

Fei Fang

# Market Failure - 1

Fei Fang

Bad draw dispatches: "after accepting, drivers are able to contact the rider. Some may [] learn [the] destination [] and canceling if [] the trip will not be worth the time."

# Competitive Equilibrium

- Competitive Equilibrium (CE)
  - Also called Walrasian equilibrium
  - Traditional concept in economics
  - Commodity markets with flexible prices and many traders

# Competitive Equilibrium

▸ A very simple setting

  ▸ A set of items $[n] = \{1, 2, \ldots n\}$

  ▸ A set of buyers $[m] = \{1, 2, \ldots, m\}$

  ▸ Each buyer $i$ has a valuation for each item $j$: $v_{ij}$

  ▸ Given a price vector $p \in \mathbb{R}^n$, agent $i$'s utility is: $u_i(x; p) = v_i \cdot x - p \cdot x$ where $x \in \{0,1\}^n$ indicates which items the agent gets

  ▸ Each agent can get at most one item

# Competitive Equilibrium

▶ A CE consists of:
  ▸ A price vector $p \in \mathbb{R}_+^n$
  ▸ A valid allocation matrix $x$
    ▹ $x_{ij} \in \{0,1\}$ indicates whether or not item $j$ is allocated to agent $i$
    ▹ Each item is allocated at most once $\sum_i x_{ij} \leq 1, \forall j$
    ▹ Each buyer can get at most one item $\sum_j x_{ij} \leq 1, \forall i$
    ▹ Use $x_i$ to denote the binary vector for agent $i$
  ▸ $p$ and $x$ satisfy the following constraints
    ▹ Best response
      □ $x_i \in \underset{x:x\in\{0,1\}^n, \sum_j x_j \leq 1}{\mathrm{argmax}} \; u_i(x; p), \forall i$
    ▹ Market clearance
      □ $\forall j, \sum_i x_{ij} = 1$ or $p_j = 0$

# Super Bowl Example



Fei Fang

# Myopic Pricing

▸ At current time $t$, each location has a sub-market

▸ Allocate cars to the riders with highest valuations

▸ Driver-pessimal price shown in black



Fei Fang

# Quiz

▸ With Myopic Pricing, at most, how much more can the purple driver earn if he deviates from the system's assignment and all other drivers always follow the system's assignment? (Options: $100, $90, $80, $0)



Fei Fang

# Useful Deviation

▸ Purple driver rejects the assigned ride at 9:50am to earn more money



Fei Fang

▸ Model: Discrete time/location, Impatient riders, Anonymous origin-destination trip price

▸ One-shot assignment

  ▸ Assignment plan: Decompose a min-cost flow

  ▸ Pricing: Dual of flow LP

  ▸ Form competitive equilibrium (CE)

    ▸ Welfare optimal

    ▸ Maximize total payment for each driver

    ▸ Maximize utility for each rider

    ▸ Envy free

    ▸ All feasible driver payments in CE form a lattice

# ILP for Computing Optimal Assignment Plan

$$\max_{x,y} \sum_{j \in \mathcal{R}} x_j v_j - \sum_{i \in \mathcal{D}} \sum_{k=0}^{|\mathcal{Z}_i|} y_{i,k} \lambda_{i,k}$$

Dual Variables

$$\text{s.t.} \sum_{j \in \mathcal{R}} x_j \mathbb{1}\{(o_j, d_j, \tau_j) = (a,b,t)\} \leq \sum_{i \in \mathcal{D}} \sum_{k=0}^{|\mathcal{Z}_i|} y_{i,k} \mathbb{1}\{(a,b,t) \in Z_{i,k}\}, \quad p_{a,b,t} \qquad \forall (a,b,t) \in \mathcal{T}$$

$$\sum_{k=0}^{|\mathcal{Z}_i|} y_{i,k} = 1, \qquad \pi_i \qquad \text{LP Relaxation} \qquad \forall i \in \mathcal{D}$$

~~$x_j \in \{0,1\}$,~~     $x_j \leq 1 \qquad u_j$     $\forall j \in \mathcal{R}$

~~$y_{i,k} \in \{0,1\}$,~~     $x_j \geq 0$     $\forall i \in \mathcal{D}, \ k = 1, \ldots, |\mathcal{Z}_i|$

$$y_{i,k} \geq 0$$

$$\min \quad \sum_{i \in \mathcal{D}} \pi_i + \sum_{j \in \mathcal{R}} u_j$$

$$\text{s.t.} \quad \pi_i \geq \sum_{(a,b,t) \in Z_{i,k}} p_{a,b,t} - \lambda_{i,k} \qquad \forall k = 0, 1, \ldots, |\mathcal{Z}_i|, \ \forall i \in \mathcal{D}$$

$$u_j \geq v_j - p_{o_j, d_j, \tau_j}, \qquad \forall j \in \mathcal{R}$$

$$p_{a,b,t} \geq 0, \qquad \forall (a, b, t) \in \mathcal{T}$$

$$u_j \geq 0, \qquad \forall j \in \mathcal{R}$$

# Spatial-Temporal Pricing

▸ However…Drivers can deviate and trigger recomputation!

▸ Solution: Driver-Pessimal CE

  ▸ Trip price = welfare gain difference
$$p_{a,b,t} = \Phi_{a,t} - \Phi_{b,t+dist(a,b)}$$
$$\Phi_{a,t} \triangleq W(D \cup \{(t, T, a)\}, R) - W(D, R)$$

  ▸ Incentive compatible subgame perfect equilibrium

  ▸ No driver want to deviate from assigned action!

▸ SPT vs Naïve surge

# Outline

▶ **Game-Theoretic Reasoning and Its Applications**

   ▶ Wildlife Conservation

   ▶ Cyber Security

   ▶ Ridesharing

▶ **End-to-End Learning and Decision Making in Games**

   ▶ A differentiable learning framework for learning game parameters

▶ **Learning-Powered Strategy Computation in Large Scale Games**

   ▶ Leveraging Deep Reinforcement Learning

# What game are we/they playing?

▶ Common criticism: game parameters are fully known

  ▶ E.g. target importance

▶ How to learn parameters of **2-player zero sum games** from opponents' or players' actions?

Solve

Equilibrium strategies
$$u^* = v^* = [\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$$

# Inverse Problem: Game Learning



i.i.d samples from equilibrium strategies

$$a^{(1)} = ( \text{✊}, \text{🖐} )$$
$$a^{(2)} = ( \text{✊}, \text{✌} )$$
$$a^{(3)} = ( \text{🖐}, \text{✌} )$$
...

Learn

# Differentiable Learning

|  | | | |
|---|---|---|---|
| | $0$ | $-b_1$ | $-b_2$ |
| | $b_1$ | $0$ | $-b_3$ |
| | $-b_1$ | $b_3$ | $0$ |

Learn ⟵

i.i.d samples from
equilibrium strategies
$$a^{(1)} = (\;\; , \;\; )$$
$$a^{(2)} = (\;\; , \;\; )$$
$$a^{(3)} = (\;\; , \;\; )$$

▸ Guess the value of $b_i$
▸ Compute equilibrium of guessed game
▸ Check if the computed equilibrium consistent with data
▸ Adjust the value of $b_i$ to increase consistency
▸ Repeat until satisfied

$\rightarrow$ Update $b_i := b_i - \frac{\partial L}{\partial b_i}$

# NE and QRE in Zero-Sum Games

## Recall LP for computing NE

$$\min_{u,x} x$$

s.t. $x \geq \sum_i u_i P_{ij} , \forall j$

$\sum_i u_i = 1, u_i \geq 0, \forall i$

## Nash Equilibrium

‣ Assumes perfect rationality

‣ May have multiple equilibria

‣ Discontinuous w.r.t. $P$

$$\min_u \max_v u^T P v$$

s.t.

$$1^T u = 1, u \geq 0$$
$$1^T v = 1, v \geq 0$$

## Recall Quantal Response

$$q_j = \frac{e^{\lambda * \text{AttEU}_j(x)}}{\sum_i e^{\lambda * \text{AttEU}_i(x)}}$$

## Quantal Response Equilibrium

‣ Captures bounded rationality

‣ Unique

‣ Continuous w.r.t. $P$

$$\min_u \max_v u^T P v - \sum_i v_i \log v_i + \sum_i u_i \log u_i$$

s.t.

$$1^T u = 1, u \geq 0$$
$$1^T v = 1, v \geq 0$$

$$u_i^* = \frac{\exp(Pv)_i}{\sum_q \exp(Pv)_q} , v_j^* = \frac{\exp(P^T u)_j}{\sum_q \exp(P^T u)_q}$$

# Learning of normal form games

- QRE = solution of min-max convex-concave problem

$$\min_u \max_v u^T P v - \sum_i v_i \log v_i + \sum_i u_i \log u_i$$

$$1^T u = 1, \ 1^T v = 1$$

- KKT conditions:

$$Pv + \log(u) + 1 + \mu 1 = 0$$
$$P^T u - \log(v) - 1 + \nu 1 = 0$$
$$1^T u = 1, \ 1^T v = 1$$

Recall: Newton's Method for 1-D:
$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$
Generally, for nonlinear system
$$J_F(x_n)(x_{n+1} - x_n) = -F(x_n)$$

- Forward pass: Apply Newton's Method

$$\begin{bmatrix} diag(\frac{1}{u}) & P & 1 & 0 \\ P^T & -diag\left(\frac{1}{v}\right) & 0 & 1 \\ & & 0 & 0 \\ 1^T & 0 & 0 & 0 \\ 0 & 1^T & & \end{bmatrix} \begin{bmatrix} \Delta u \\ \Delta v \\ \Delta \mu \\ \Delta \nu \end{bmatrix} = - \begin{bmatrix} Pv + \log(u) + 1 + \mu 1 \\ P^T u - \log(v) - 1 + \nu 1 \\ 1^T u - 1 \\ 1^T v - 1 \end{bmatrix}$$

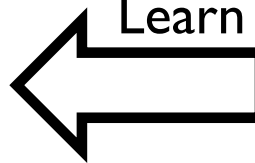▸ Backward pass: Gradients of $P$ may be obtained via the implicit function theorem

$$\nabla_P L = y_u v^T + u y_v^T,$$

where

$$
\begin{bmatrix} y_u \\ y_v \\ y_\mu \\ y_\nu \end{bmatrix} =
\begin{bmatrix}
\mathrm{diag}(\frac{1}{u}) & P & 1 & 0 \\
P^T & -\mathrm{diag}(\frac{1}{v}) & 0 & 1 \\
1^T & 0 & 0 & 0 \\
0 & 1^T & 0 & 0
\end{bmatrix}^{-1}
\begin{bmatrix} -\nabla_u L \\ -\nabla_v L \\ 0 \\ 0 \end{bmatrix}
$$

# Learning in the presence of features

|  |  |  |  |
|---|---|---|---|
|  | ✊ | ✋ | ✌ |
| ✊ | $0$ | $-b_1(x)$ | $b_2(x)$ |
| ✋ | $b_1(x)$ | $0$ | $-b_3(x)$ |
| ✌ | $-b_2(x)$ | $b_3(x)$ | $0$ |

← **Learn**

i.i.d samples from
equilibrium strategies
$$a^{(1)} = ( \ ✊ \ , \ ✋ \ )$$
$$a^{(2)} = ( \ ✊ \ , \ ✌ \ )$$
$$a^{(3)} = ( \ ✋ \ , \ ✌ \ )$$
...

Context
$$x^{(1)} = [0.1, 0.5]$$
$$x^{(2)} = [0.3, 0.7]$$
...

# Learning in the presence of features

▸ Figure out which features attract/discourage attackers

- ▸ Better understand attacker's interests
- ▸ Design better configurations which favor defenders

▸ Predict each player's mixed strategy given an *new* environment

- ▸ In practice, environment is changing over time

# Learning in the presence of features

▶ *Context* (feature) $x^{(i)}$ and payoff matrix $P_\Phi(x^{(i)})$, parameterized by $\Phi$

▶ Each player acts according to a mixed strategy $(u, v)$ given by the QRE of $P_\Phi(x^{(i)})$, giving realizations $a^{(i)}$

▶ Objective: Learn $\Phi$ from $\{x^{(i)}, a^{(i)}\}$

# End-to-end learning



**Algorithm 1:** Learning parameters $\Phi$ using SGD

**Input:** training data $\{(x^{(i)}, a^{(i)})\}$, learning rate $\eta$, $\Phi_{\text{init}}$

**for** $ep$ $in$ $\{0, \ldots, ep_{max}\}$ **do**

    Sample $(x^{(i)}, a^{(i)})$ from training data;
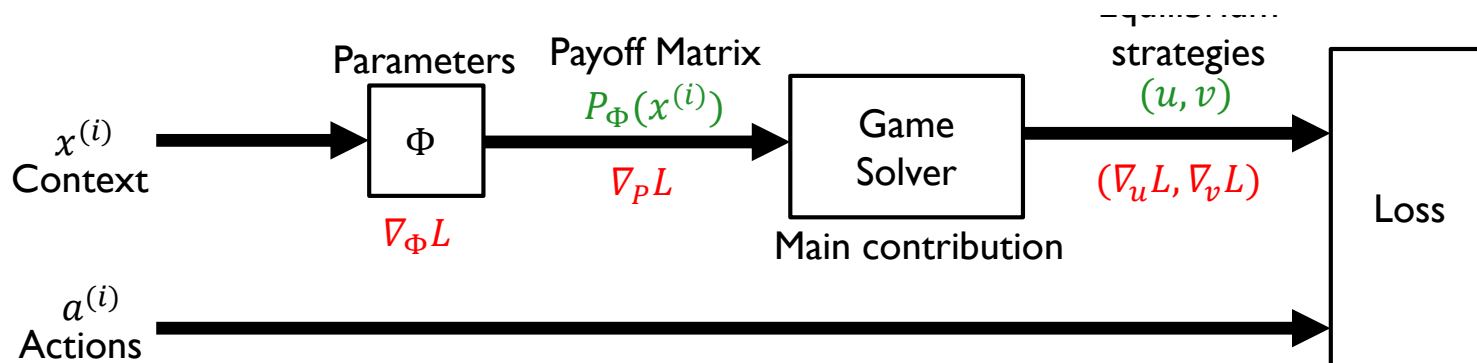
    Forward pass: Compute $P_\Phi(x^{(i)})$, QRE $(u, v)$ and loss $L(a^{(i)}, u, v)$;

    Backward pass: Compute gradients $\nabla_u L, \nabla_v L, \nabla_P L, \nabla_\Phi L$;

    Update parameters: $\Phi \leftarrow \Phi - \eta \nabla_\Phi L$;

**end**

# Extensive form Games

- Let $(u, v)$ be strategies in sequence form

- Equilibrium is expressed as solution using *dilated entropy regularization* (Equivalent to solving QRE for the reduced normal form)

$$\min_{u} \max_{v} u^T P v - \sum_{i} \sum_{a} v_a \log(\frac{v_a}{v_{p_i}}) + \sum_{i} \sum_{a} u_a \log(\frac{u_a}{u_{p_i}})$$
$$Eu = e, Fv = f$$
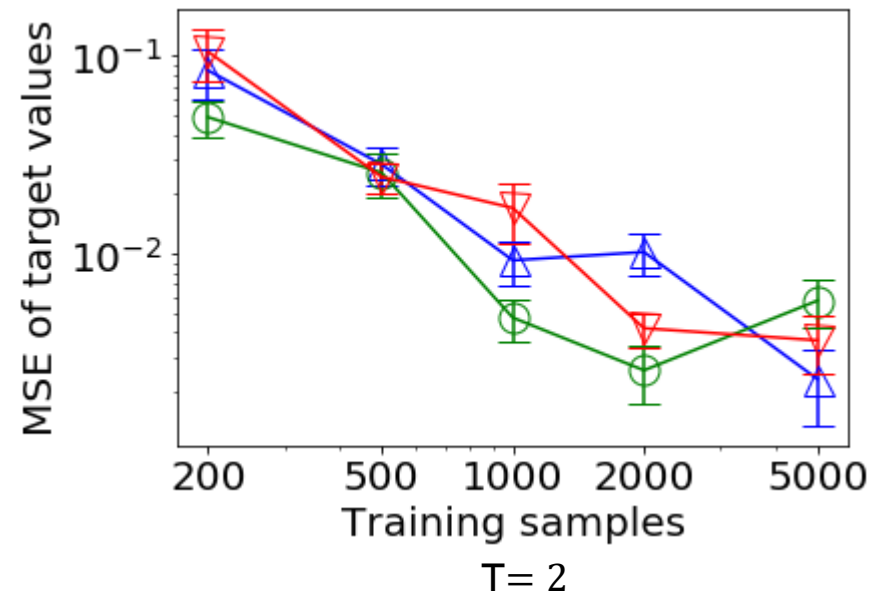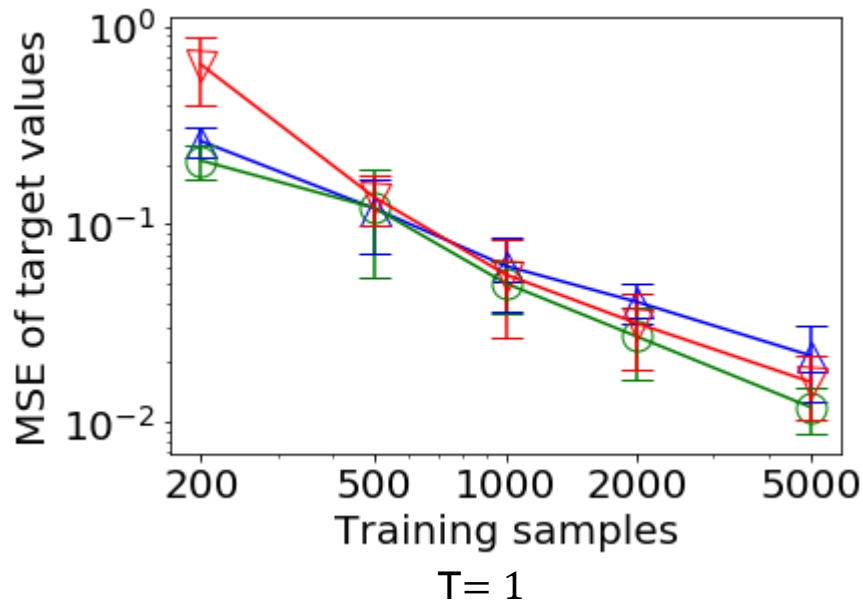
$$\nabla_P L = y_u v^T + u y_v^T,$$

where

$$\begin{bmatrix} y_u \\ y_v \\ y_\mu \\ y_\nu \end{bmatrix} = \begin{bmatrix} -\Xi(u) & P & E^T & 0 \\ P^T & \Xi(v) & 0 & F^T \\ E & 0 & 0 & 0 \\ 0 & F & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} -\nabla_u L \\ -\nabla_v L \\ 0 \\ 0 \end{bmatrix}$$

# Resource Allocation Security Game

▸ Defender: $r$ resources, $n$ targets
  ▸ Can allocate multiple resources to one target
▸ Attacker choose a target to attack
▸ Each target has value $R_i$
▸ If target $i$ is protected by $x$ resources and is attacked:
$$U_a = \frac{R_i}{2^x} = -U_d$$
▸ Attacker may learn $R_i$ from observed defender actions

▸ Extend to $T$-stage game

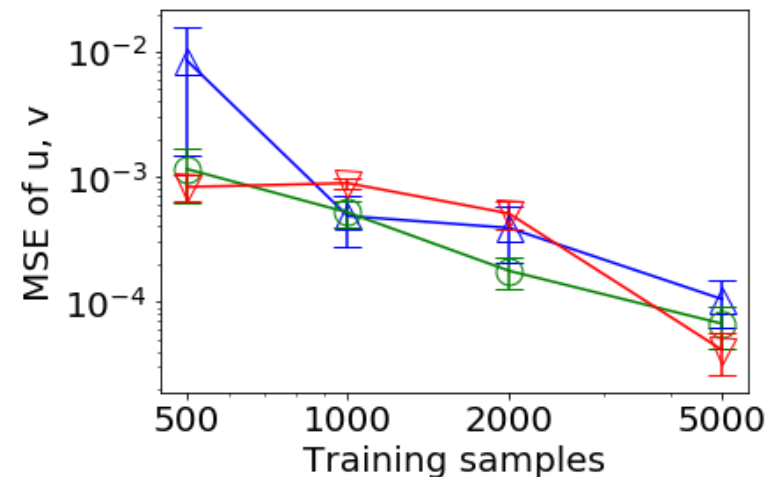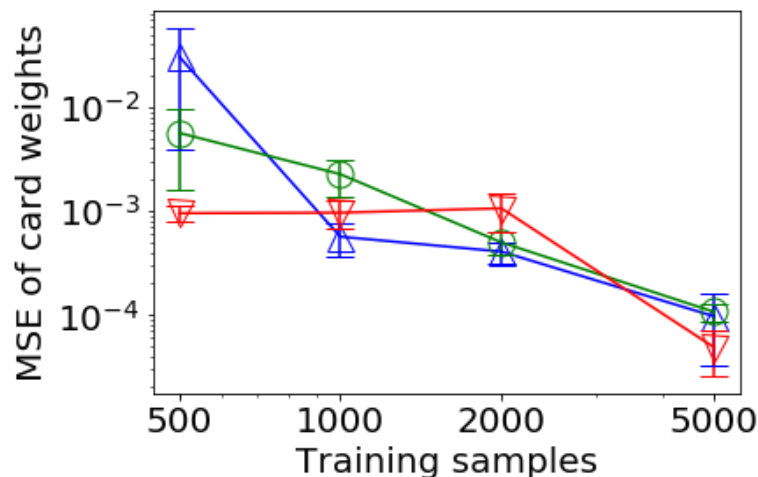# Resource Allocation Security Game

$n = 2, r = 5$



T= 1

T= 2

# One-Card Poker

▸ Learn players' belief of card distribution

▸ Variant of Kuhn Poker with 4 cards, with *non-uniform* card distributions

▸ Observe actions of each player (e.g. raise, fold)
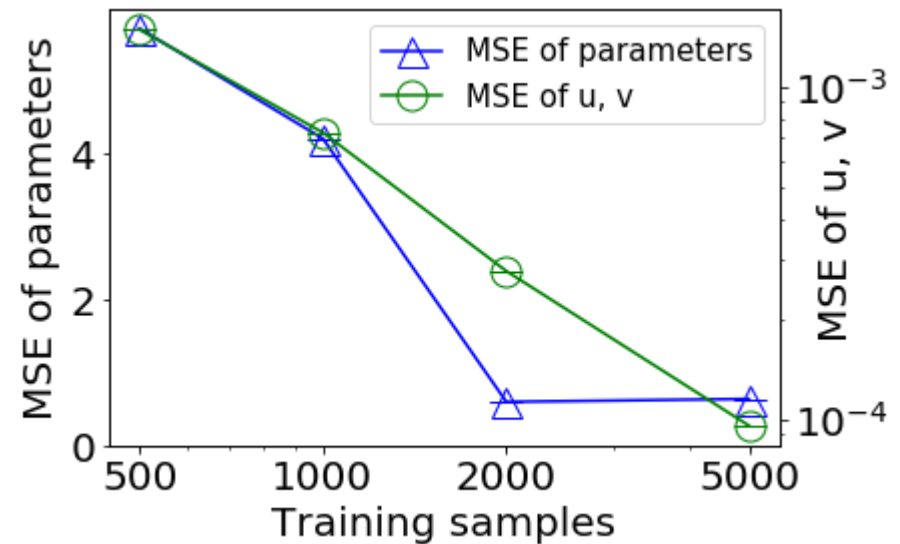
▸ Probabilities for chance nodes are embedded in $P_\Phi$

# Featurized Rock Paper Scissors

$$P = $$

|   | R | P | S |
|---|---|---|---|
| R | 0 | $-b_1$ | $b_2$ |
| P | $b_1$ | 0 | $-b_3$ |
| S | $-b_2$ | $b_3$ | 0 |

$$b = \Phi x,$$
$$x \in [0, 1]^2$$
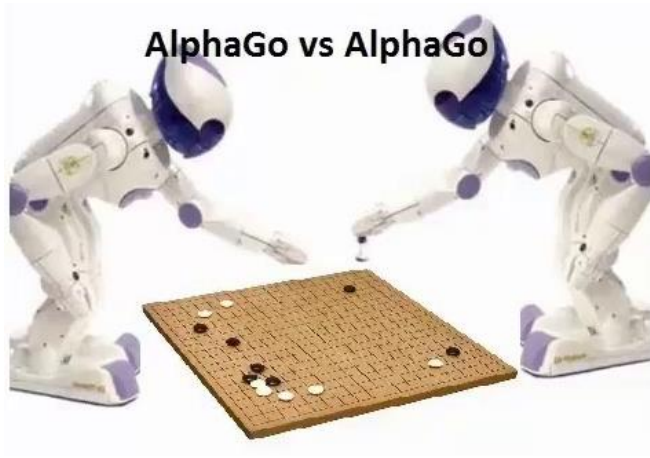$$\Phi \in [0, 10]^{3 \times 2}$$

Objective is to learn $\Phi$

# Outline

▸ Game-Theoretic Reasoning and Its Applications
  ▸ Wildlife Conservation
  ▸ Cyber Security
  ▸ Ridesharing

▸ End-to-End Learning and Decision Making in Games
  ▸ A differentiable learning framework for learning game parameters

▸ Learning-Powered Strategy Computation in Large Scale Games
  ▸ Leveraging Deep Reinforcement Learning
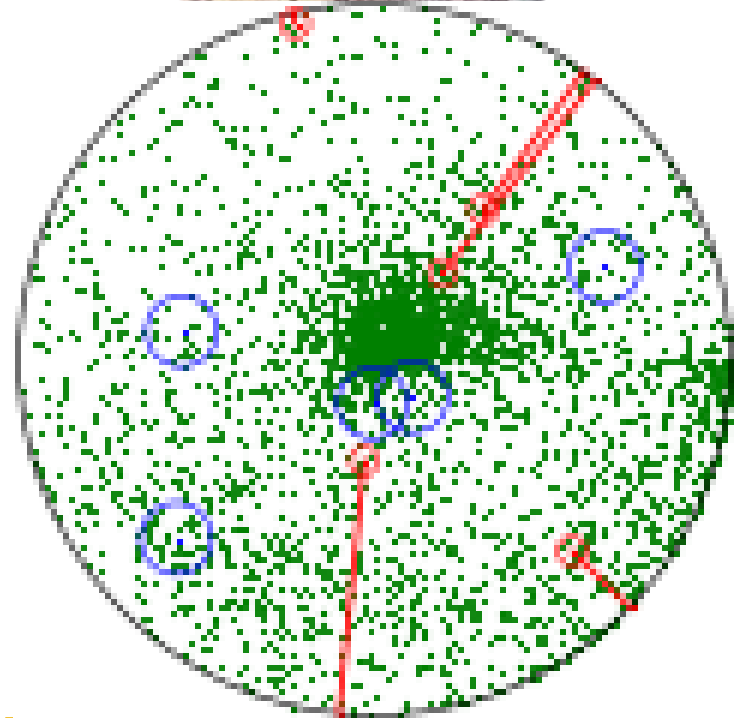
# Solving Game through Learning from Self Play



AlphaGo vs AlphaGo

https://www.youtube.com/watch?v=Ue4A2Y_i3ZQ

Self Play

Update Strategy

Compute Optimal Defender Strategy

# Solving Game through Learning from Self Play
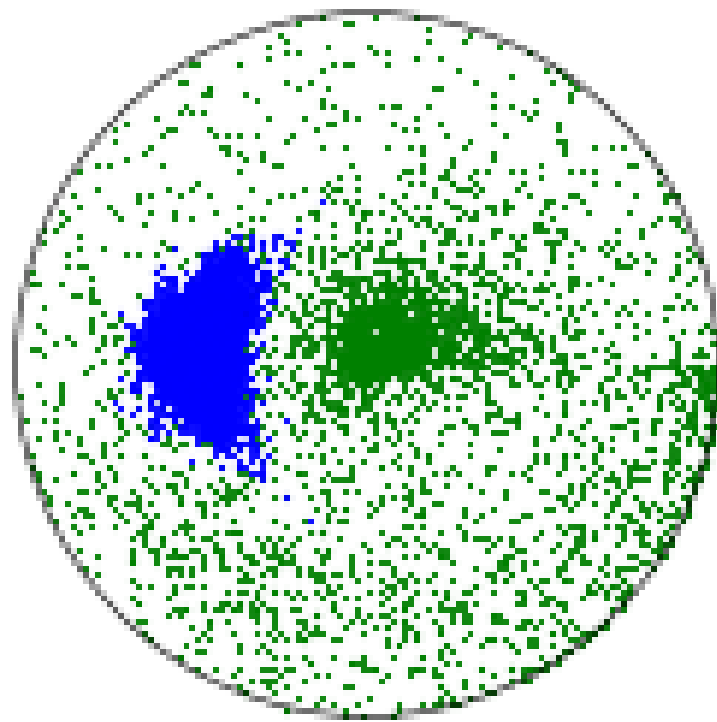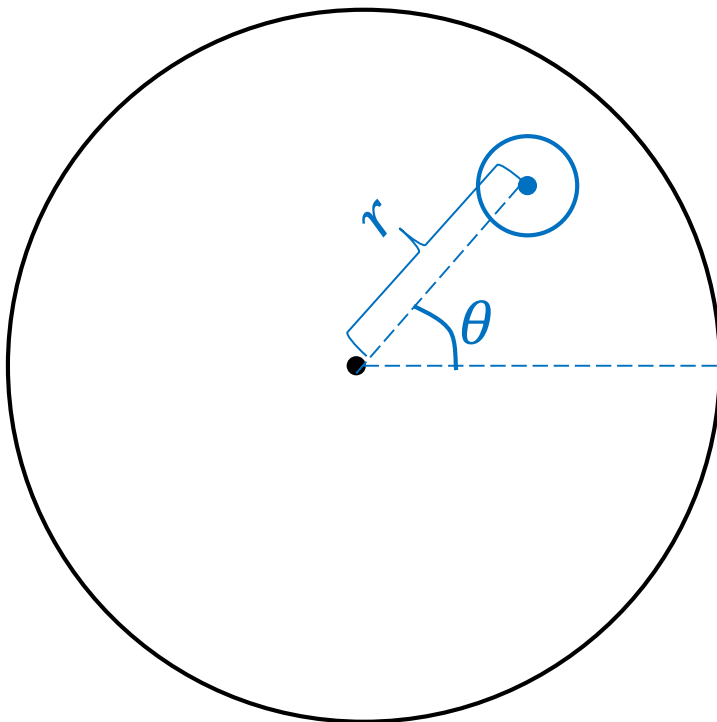


▶ Green dots: Valuable trees

▶ Blue dots: Defender location

▶ Red dots: Logging locations

▶ Zero-sum game

▶ Goal: Find defender strategy or defender policy

Policy Learning for Continuous Space Security Games using Neural Networks. Nitin Kamra, Umang Gupta, Fei Fang, Yan Liu, Milind Tambe. In AAAI-18
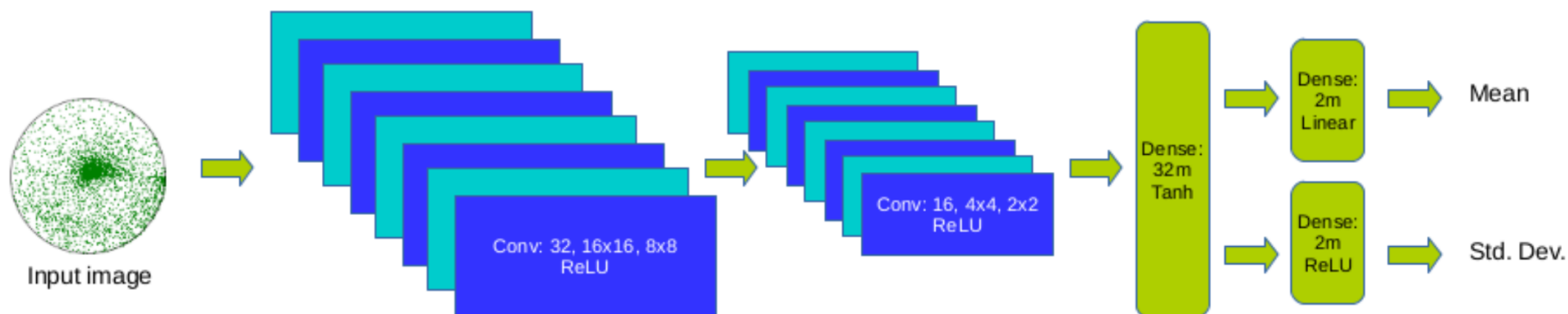
5/20/2019

▸ Key idea 1: Represent mixed strategy using logit normal distribution in polar coordinate system

$$r \sim P\left(\mathcal{N}\left(\mu_r, \sigma_r^2\right)\right)$$
$$\theta \sim P\left(\mathcal{N}\left(\mu_\theta, \sigma_\theta^2\right)\right)$$

▸ Key idea 2: Represent a "policy" with Convolutional Neural Network

  ▸ Policy: mapping from game setting to strategy
  ▸ CNN: Tree Distribution →Mean/Std of $r$ and $\theta$

▶ Key idea 3: Approximate Fictitious Play

▶ Fictitious Play: Best responds to opponent's average strategy

▶ Average strategy → Random samples from history

▶ Best response → Update neural network

# Solving Game through Learning from Self Play

▸ Put them together

---

**Algorithm 1:** OptGradFP

---

Initialization. Initialize policy parameters $w_D$ and $w_O$, replay memory *mem*;
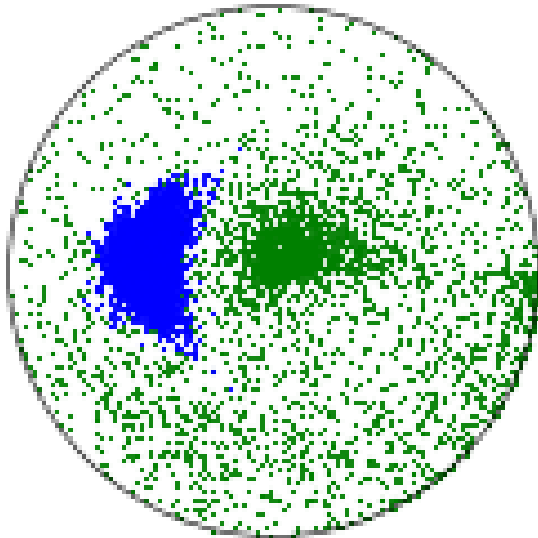
**for** *ep* in $\{0, \ldots, ep_{max}\}$ **do**

    Simulate $n_s$ game play. Sample game setting and actions from current policy $\pi_D$ and $\pi_O$ $n_s$ times, save in *mem*;

    Replay for defender. Draw $n_b$ samples from *mem*, resample defender action from current policy $\pi_D$;

    Update parameter for defender. Update defender policy parameter

    $w_D := w_D + \frac{\alpha_D}{1 + ep\, \beta_D} * \nabla_{w_D} J_D$;

    Replay for attacker. Draw $n_b$ samples from *mem*, resample attacker action from current policy $\pi_O$;

    Update parameter for attacker. Update attacker policy parameter

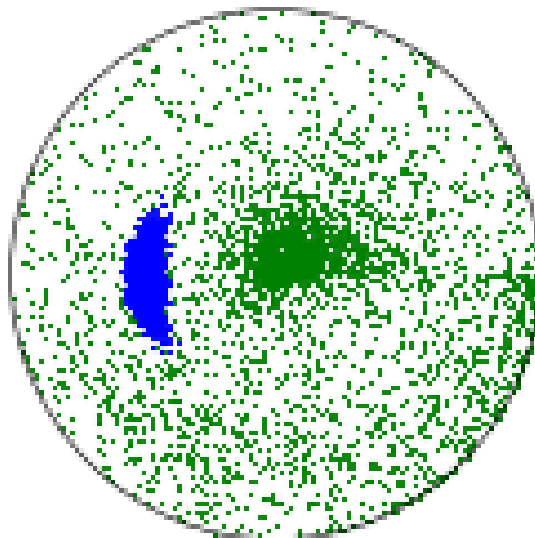    $w_O := w_O + \frac{\alpha_O}{1 + ep\, \beta_O} * \nabla_{w_O} J_O$
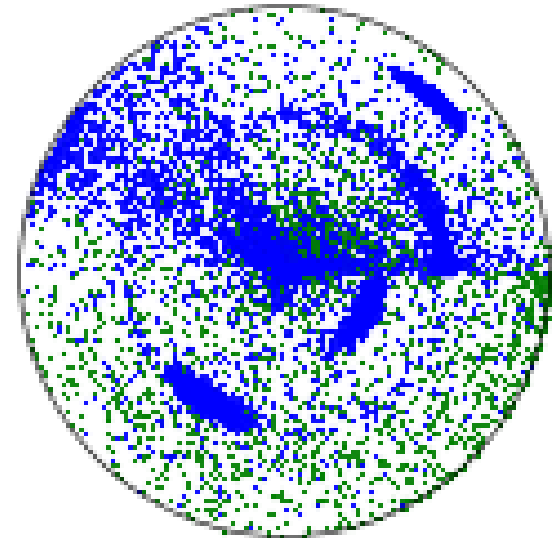
---

- ▶ **Single game setting**



Cournot Adjustment        StackGrad        OptGradFP

- ▶ **Multiple game setting**
  - ▶ Train on 1000 forest states, predict on unseen forest state
  - ▶ 7 days for training, Prediction time 90 ms
  - ▶ Shift computation from online to offline

▶ Sequential interaction

  ▸ Players make flexible decisions instead of sticking to a plan

  ▸ Players may leave traces as they take actions

▶ Example domain: Wildlife protection


Footprints


Lighters
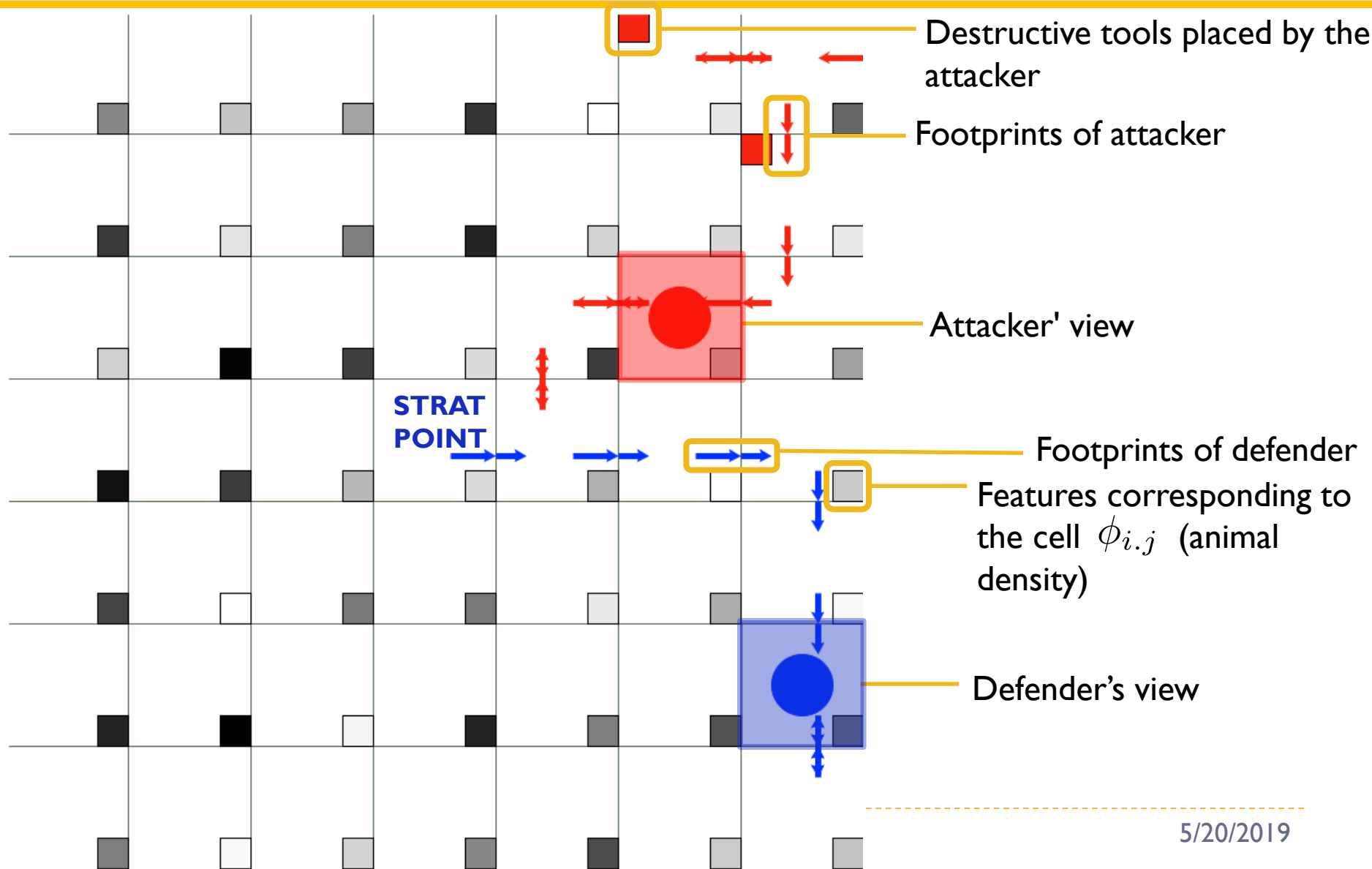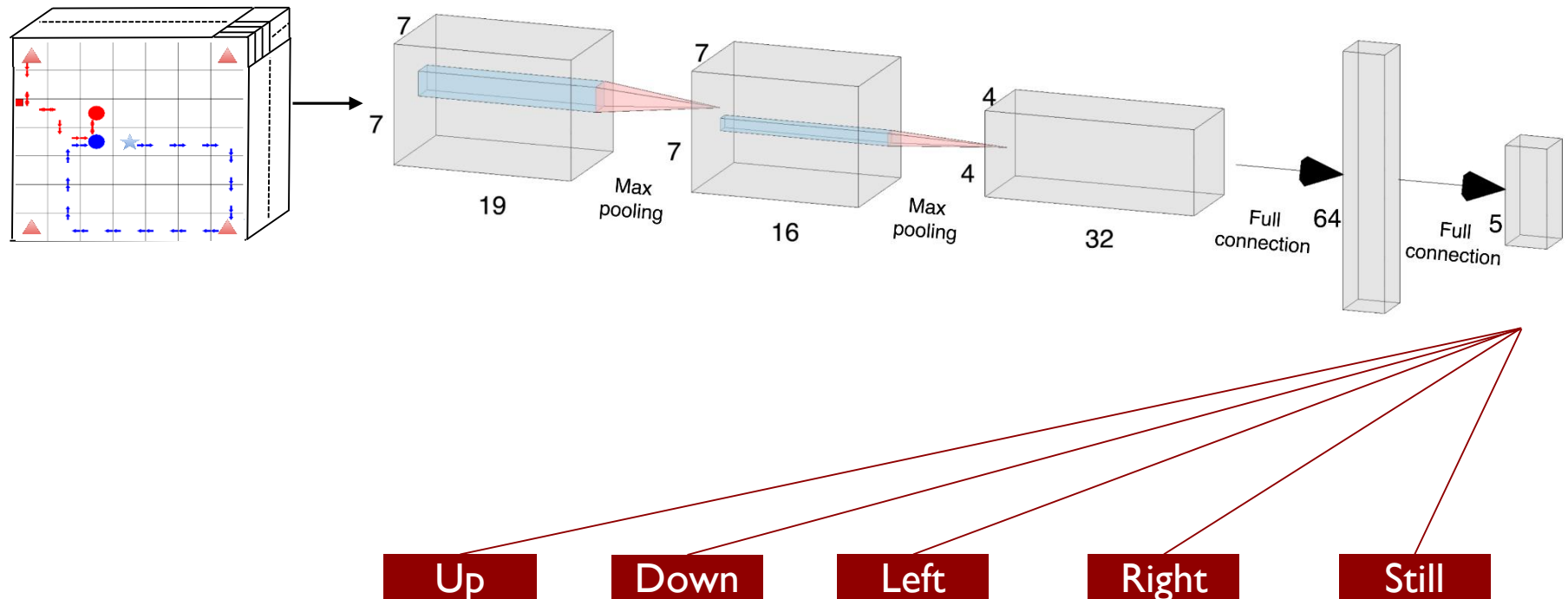

Old poacher camp
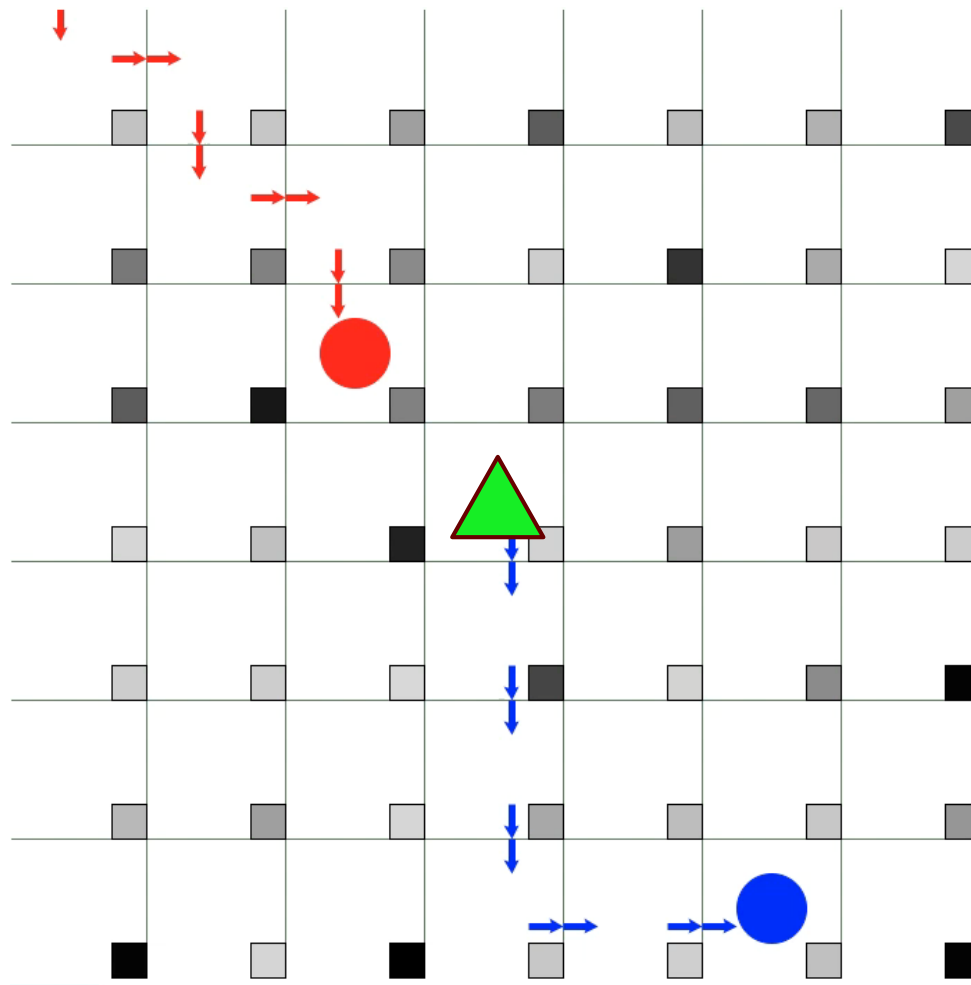

Tree marking

# Multi-Agent Reinforcement Learning



Destructive tools placed by the attacker

Footprints of attacker

Attacker' view

STRAT POINT

Footprints of defender

Features corresponding to the cell $\phi_{i.j}$ (animal density)

Defender's view

# Compute Best Response by Training a Deep Q-Network



| Up | Down | Left | Right | Still |

▸ Q Network: Game state → Q-value

▸ Use Deep Reinforcement learning to train the network and find optimal patrol policy (assuming fixed attacker)

# Compute Best Response by Training a Deep Q-Network
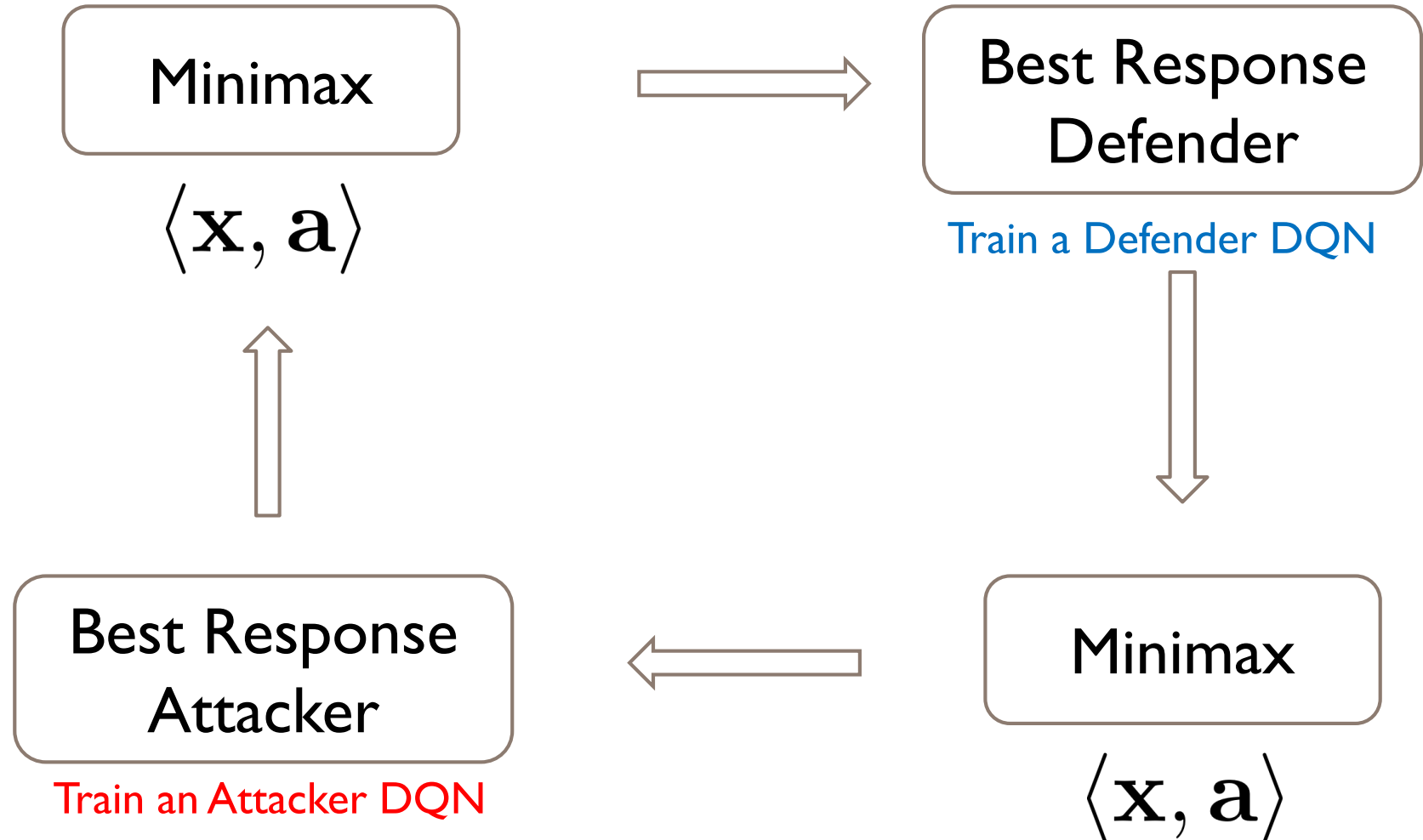


DQN Defender
vs
Non-Adaptive Attacker

Attacker

Snares

Start from one of the corners
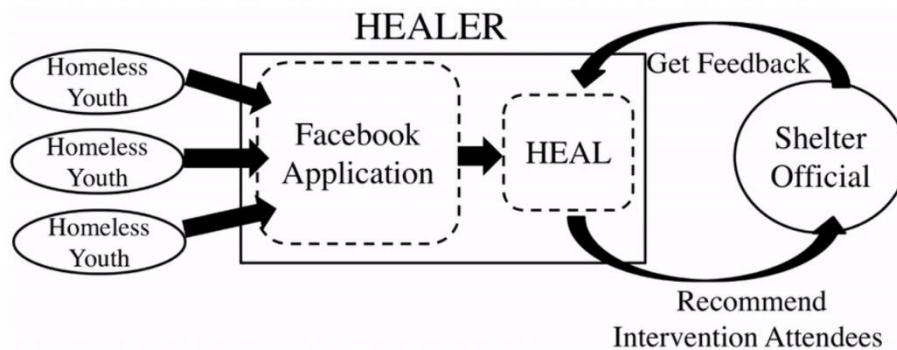
Defender

Start from
Patrol Base

# AI and Social Good

▶ AI research that can deliver societal benefits now and in the near future



http://mashable.com/2015/02/06/hiv-homeless-teens-algorithm/#..k9dRKhxaqm



https://www.pastemagazine.com/articles/2017/04/a-new-smart-technology-will-help-cities-drasticall.html

# Summary

- Game-Theoretic Reasoning and Its Applications
  - Wildlife Conservation, Cyber Security, Ridesharing
- End-to-End Learning and Decision Making in Games
- Learning-Powered Strategy Computation in Large Scale Games

# Thank you!

Fei Fang
Carnegie Mellon University
feifang@cmu.edu