

Reinforcement Learning China Summer School



RLChina 2020

Multi-agent Systems

Bo AN (安波)

Nanyang Technological University

August 5, 2020

Outline

- Part 1-History and current status (10min)
- Part 2-Key research areas in MAS (30min)
- Part 3-Recent advances (40min)

- *Computer poker*
- *Game theory for security*
- *Multi-agent RL*

专题 | 中国计算机学会通讯 第10卷 第9期 2014年9月

多智能体系统研究的历史、现状及挑战

关键词：智能体 多智能体系统 人工智能

安 波¹ 史忠植²
¹新加坡南洋理工大学
²中国科学院计算机技术研究所

多智能体系统 (multi-agent systems) 由分布式人工智能演变而来, 其研究目的是解决大规模、复杂、实时和有不确定信息的现实问题, 而这类问题是单个智能体所不能解决的。多智能体系统通常具有自主性、分布性、协调性等特征, 并具有自组织能力、学习能力和推理能力, 对多智能体系统的研究既包括构建单个智能体的技术, 如建模、推理、学习及规划等, 也包括使多个智能体协调运行的技术, 例如交互通信、协调、合作、协商、调度、冲突消解等。

随着互联网技术的发展, 多智能体系统往往会出现一些“自私”的智能体(如电子商务市场的交易者), 因此需要引入博弈论分析智能体的交互策略, 其研究内容包括电子商务、拍卖、机构设计、社会选择理论等。

经过近 30 年的发展, 多智能体系统已经成为国际人工智能领域的前沿和研究热点。在近年

8



图1 多智能体系统创始人雅克·莫瑟

来的 AAAI 人工智能会议 (AAAI Conference on Artificial Intelligence) 和国际人工智能联合会与思维奖 (IJCAI Computers and Thought Award) 的获奖者有一半以上来自多智能体系统领域。

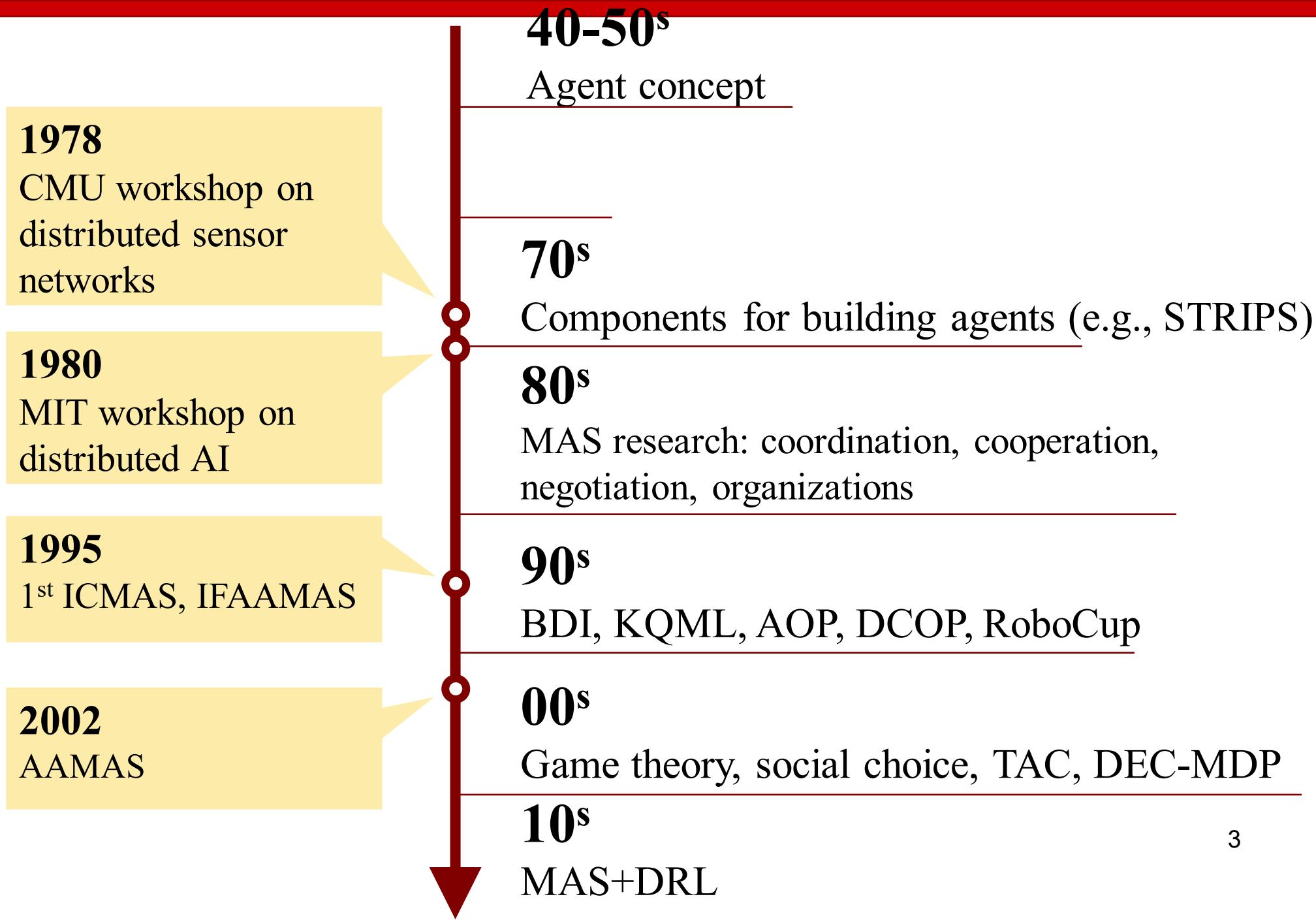
多智能体系统研究的历史

自 1956 年的翰·麦卡锡 (John McCarthy) 在著名的达特茅斯研讨会上提出“人工智能”这一概念后, “智能体”的概念便开始兴起。例如, 阿兰·图灵 (Alan Turing) 提出了用来判断一台机器是否具备人类智能的“图灵测试”。在此测试中, 测试者通过监视设备向被测试的实体 (也就是我们目前正在所说的智能体) 提问, 根据图灵的观点, 如果测试者无法区分被测试对象是计算机还是人, 那么被测试对象就是智能的。人工智能学者马文·明斯基 (Marvin Minsky) 教授于 2009 年获得了 IJCAI 杰出研究奖 (IJCAI Award for Research Excellence)。

Points I want to make:

- ✧ Mainstream AI research
- ✧ Still at the very early age: Lots of work to be done and we can make a difference in both theory and applications

MAS Research: History



MAS: One of The Most Important Field in AI

- Over 30 years' history: Initially called **distributed AI**
- Now one of the most important/active fields of AI
 - *IJCAI Award for Research Excellence: Victor Lesser (09), Barbara Grosz (15)*
 - *IJCAI Computers and Thought Award:*
 - ❖ Sarit Kraus (95), Nicholas Jennings (99), Tuomas Sandholm (03), Peter Stone (07), Vincent Conitzer (11), Ariel Procaccia (15)



- *Large number of accepted papers at AAAI/IJCAI*
- *Researchers from top universities*
- *Quality of AAMAS, JAAMAS*

MAS Research Areas

➤ Theoretic foundations

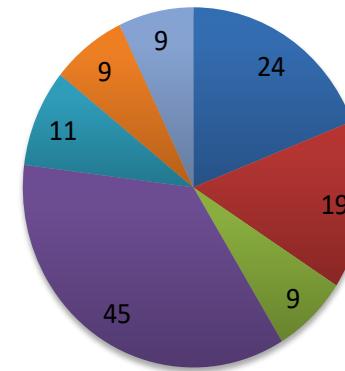
- *Coordination, DEC-(PO)MDP*
- *DCOP*
- *Game theory*

❖ mechanism design, coalition formation, social choice, negotiation

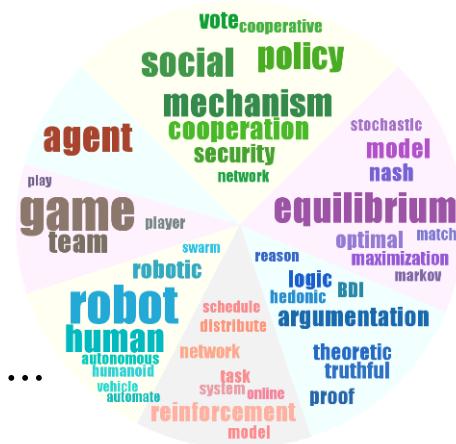
- *Multi-agent Learning*
- *Organizational design*

➤ Applications

- *Robotics, security, sustainability, rescue, sensor networks, games...*

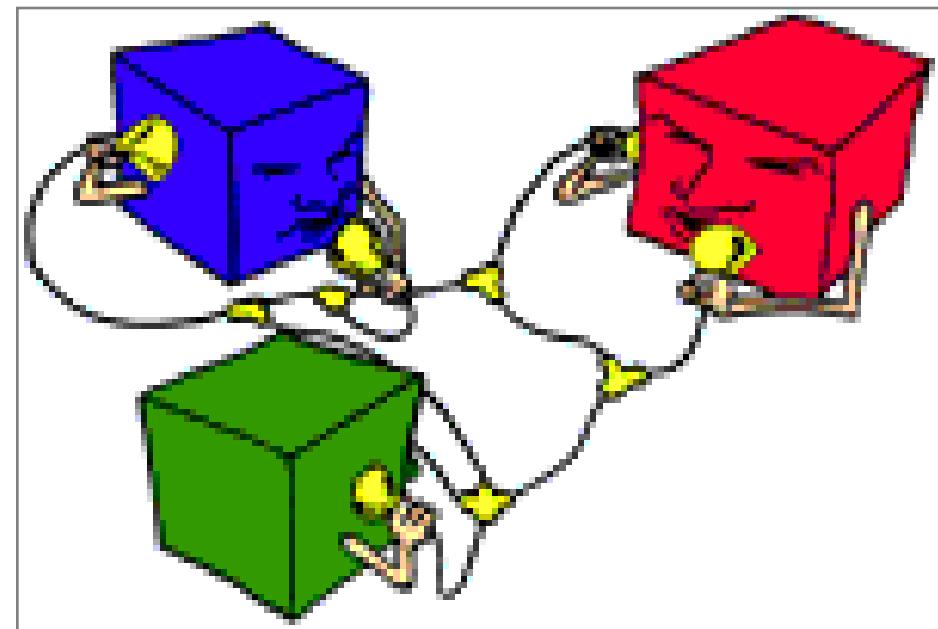


- 智能体合作 (Cooperation)
- 智能体推理 (Reasoning)
- 社群 (Societies)
- 经济学范式 (Economic paradigms)
- 人与智能体 (Humans and agents)
- 学习与自适应 (Learning and adaptation)
- 机器人学 (Robotics)

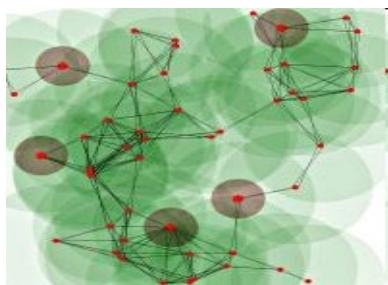


Cooperative Multi-Agent Systems:

Groups of Sophisticated AI systems that Work Together

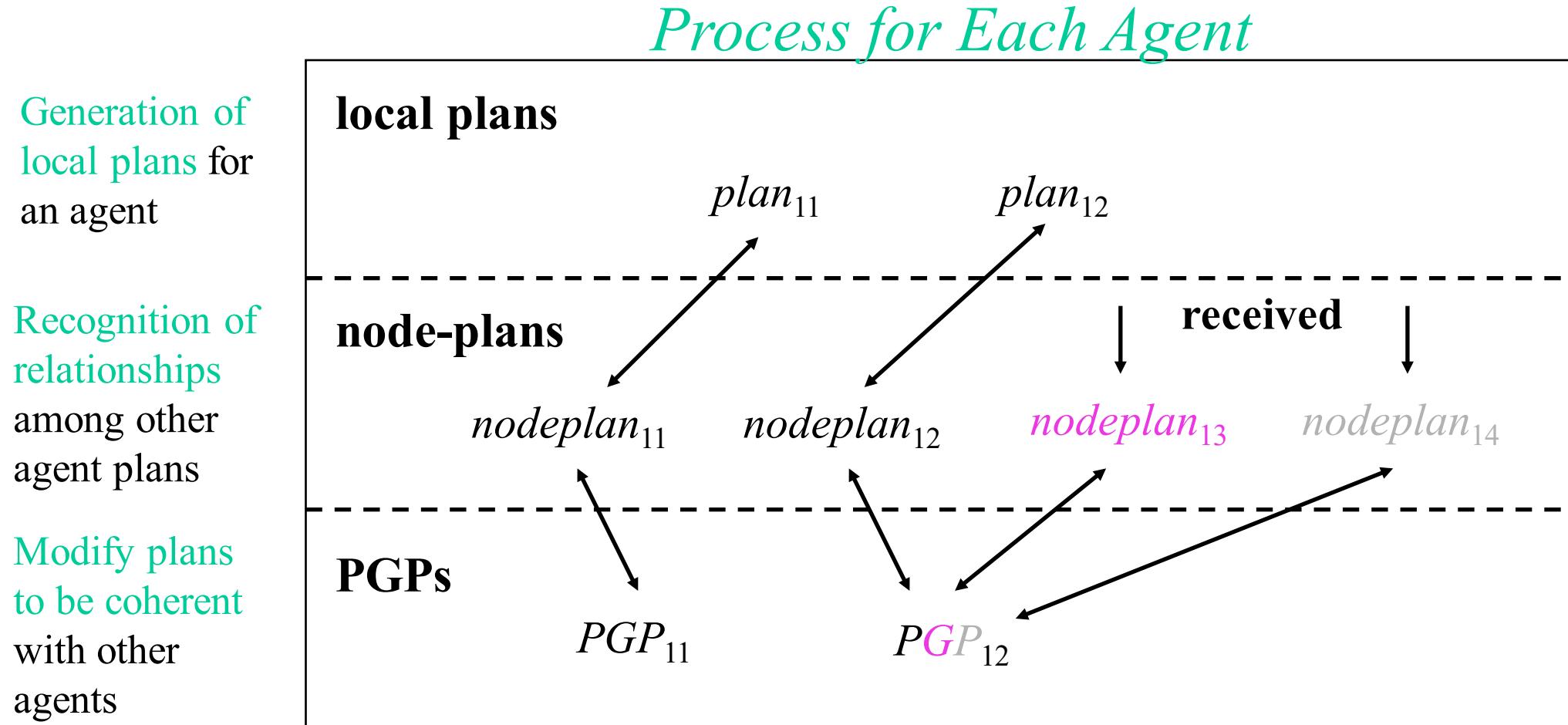


- Open, dynamic, persistent systems
- Decentralized control
- Asynchronous
- Large scale (10s to 1000s)
- Partial observability
- No real-time global reward signal
- Communication delay



Partial Global Planning (PGP) [Lesser et al]

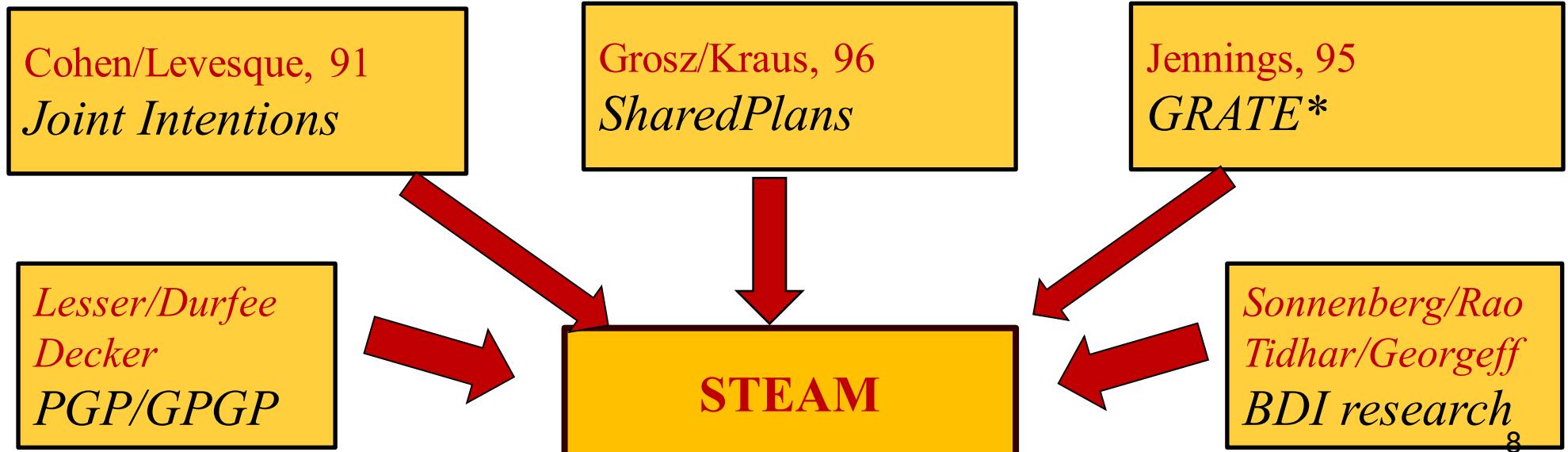
- Form a partial-global-plan (PGP) to achieve partial-global-goal
- PGP points to local models of participating plans



Towards Flexible Teamwork [Tambe]



- *Scout crashed, company waited forever*
- *Commander returned to home base alone...*
- *Ad-hoc: No Framework to anticipate failures*



Dec-(PO)MDP Models

GO-DEC-MDP Goldman et. al, 2004

OC-DEC-MDP Beynier et. al, 2005

ND-POMDP

Nair et. al, 2005

Varakantham
et. al, 2009

DPCL
IDMG

Interaction in a few
global states

Spaan et. al, 2008

Restricted
Problem

Locality of
Interaction

Factored State & Local
Observability

DEC-POMDP
Bernstein et. al, 2000

Independence
& Structured
Dependence

LO-DEC-MDP

Becker et. al, 2004

TI-DEC-MDP

Becker et. al, 2004

EDI-DEC-MDP

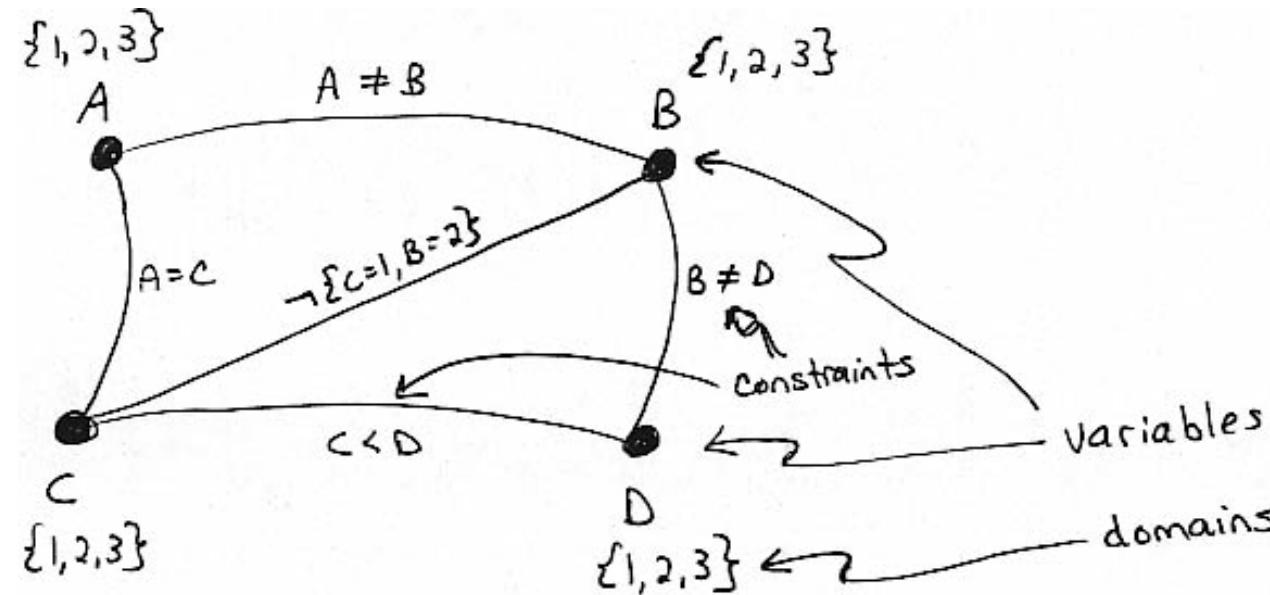
TD-POMDP

Witwicki et al., 2010

Courtesy of V. Lesser

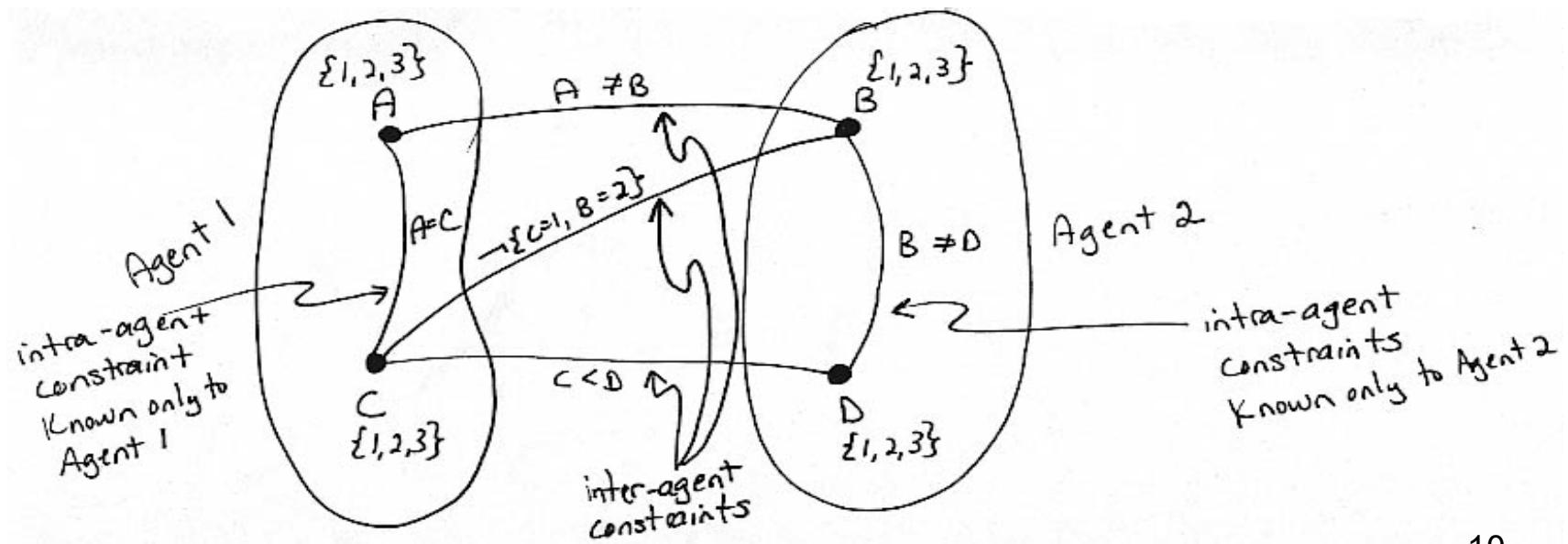
Further reading: Bernstein, Zilberstein, Immerman: The Complexity of Decentralized Control of Markov Decision Processes. UAI'00.

Distributed Constraint Satisfaction & Optimization



Applications:

- Distributed sensor networks
- Resource allocation
- Meeting scheduling
- Engineering design



Organizational Design

➤ Human organizations

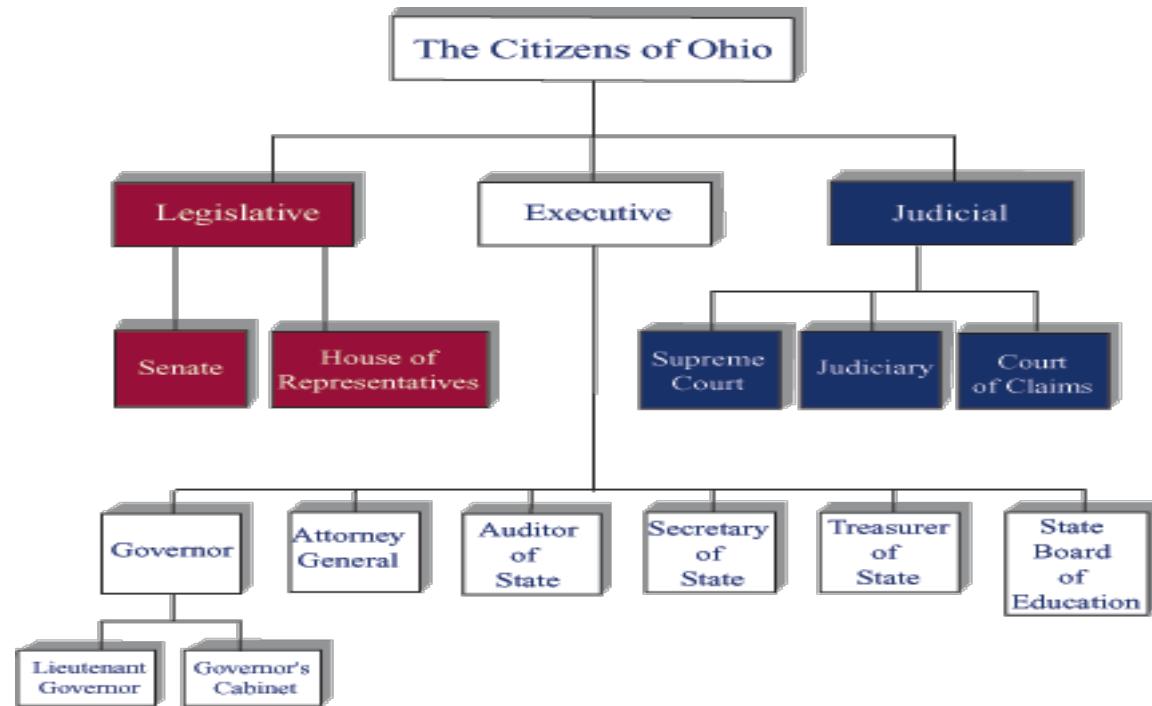
- ❑ *Businesses*
- ❑ *Governments*
- ❑ *Universities*

➤ Long-term, repetitive, multi-participant

➤ Organizations give participants a set of responsibilities, goals, incentives and guidelines

➤ MAS organizations are the same

- ❑ *Assignment of goals, incentives, roles, rights, authority, expectations, partners, rules, conventions...*



Game Theory for AI

- GT analyzes *multi-agent* interaction (1940s-)



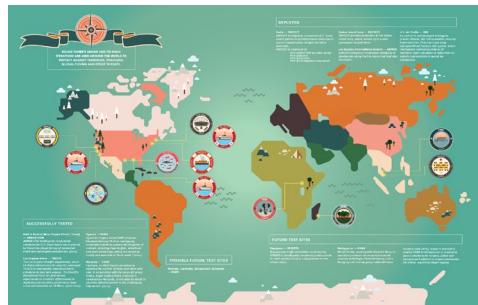
- AI: study and construction of *rational* agents [Russell & Norvig, 2003]

building a single agent
(1950s-70s)

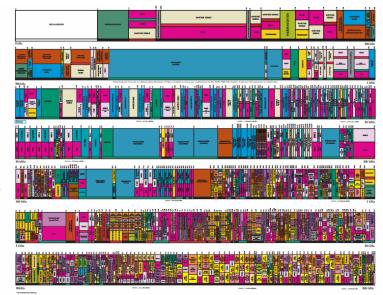
Multi-agent systems (cooperative)
(1980s-)

Multi-agent systems (competitive)
(1995-)

- GT for AI: success in computer poker, security, auction, GAN, ...

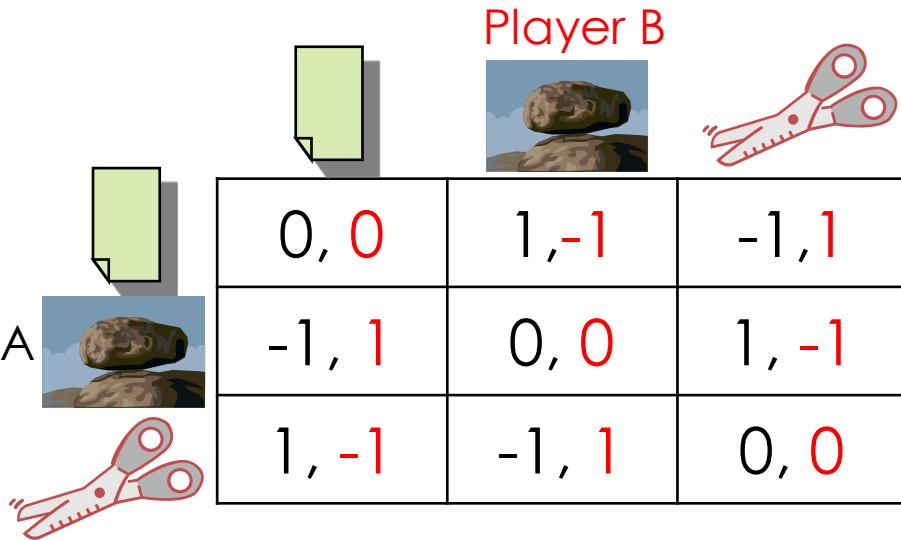


UNITED
STATES
FREQUENCY
ALLOCATIONS
THE RADIO SPECTRUM



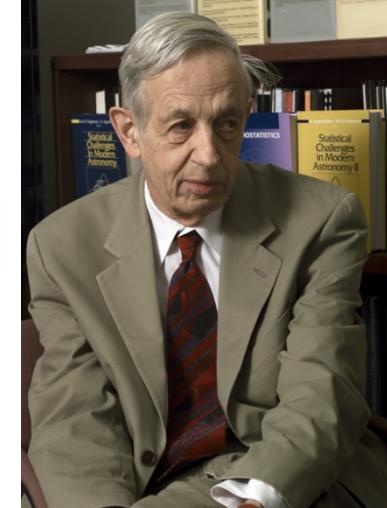
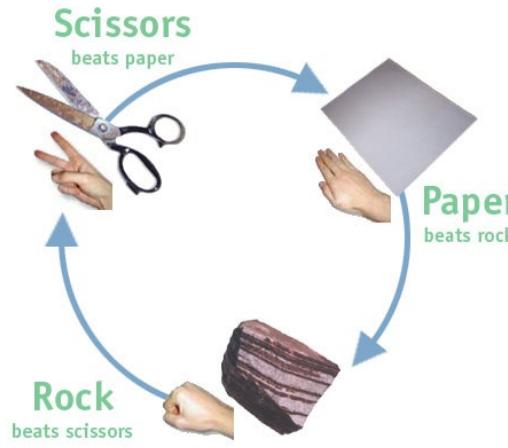
Normal-form Game and Nash Equilibrium

- Players, strategies, payoffs



A 3x3 matrix representing the Rock-Paper-Scissors game. The columns are labeled "Player B" and the rows are labeled "Player A". The payoffs are as follows:

		Player B		
		Scissors	Paper	Rock
Player A	Scissors	0, 0	1, -1	-1, 1
	Paper	-1, 1	0, 0	1, -1
	Rock	1, -1	-1, 1	0, 0



- Nash Equilibrium: no agent has incentive to unilaterally deviate

- *In any (finite) game, at least one Nash equilibrium (possibly mixed) exists* [Nash, 50]
- *In 2-player zero-sum games, a profile is an NE iff both players play minimax strategies*
- *Computing one (any) Nash equilibrium is PPAD-complete (even in 2-player games)*
[Daskalakis, Goldberg, Papadimitriou 2006; Chen, Deng 2006]
- *All known algorithms require exponential time (in the worst case)*
 - ❖ Lemke-Howson, support enumeration

- Mechanism design

Recent advances: Libratus and Pluribus



Abstraction
(offline)

- action abstraction
- card abstraction
- took the game size from 10^{161} to 10^{12}

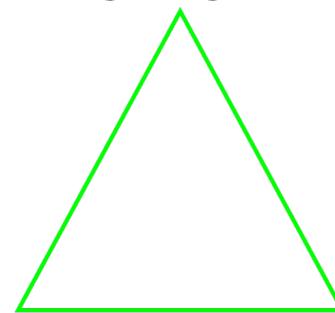
Equilibrium
Finding
(offline)

- CFR
- CFR⁺
- Monte Carlo CFR

Decomposition
and Subgame
Refinement
(online)

- endgame solving
- subgame re-solving
- max-margin subgame refinement

Original game



Automated abstraction

Abstracted game

Compute Nash

Nash equilibrium

Reverse model

Nash equilibrium

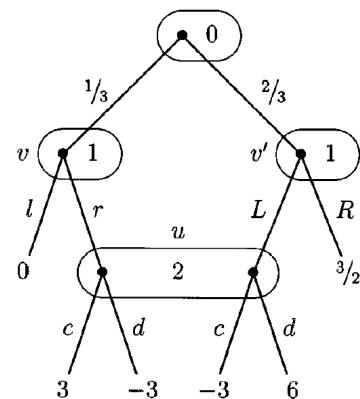
➤ Deep learning: Alberta's DeepStack, DeepMind

Regret and Regret Matching

- For each player i , action a_i and time t , define the regret $r_i(a_i, t)$ as

$$\sum_{1 \leq t' \leq t-1} (u_i(a_i, a_{-i,t'}) - u_i(a_{i,t'}, a_{-i,t'})) / (t-1)$$

- Regret matching: at time t , play an action that has positive regret $r_i(a_i, t)$ with probability proportional to $r_i(a_i, t)$
 - If none of the actions have positive regret, play uniformly at random
- Theorem [Hart & Mas-Colell 2000]:** Informally, if we pick actions according to RM in a zero-sum game, the regret for any a_i is approaching 0 as t goes to infinity
- Folk Theorem:** In a zero-sum game at time T , if both players' average overall regret is less than ε , then the average strategy is a 2ε equilibrium



	c	d
(l, L)	-2	4
(l, R)	1	1
(r, L)	-1	3
(r, R)	2	0

CFR: Regret Computing

- Key idea: minimize regret independently in each information set of an extensive-form game by maintaining the average immediate regret:

$$R_i^T(I, a) = \frac{1}{T} \sum_{t=1}^T \pi_{-i}^{\sigma^t}(I) \cdot (u_i(\sigma^t|_{I \rightarrow a}, I) - u_i(\sigma^t, I))$$

- π is the probability of reaching information set I given the strategies at time t and the player i is trying to reach that information;
 - μ is the utility of a strategy, including that we switch our action only at that information set
- Define the strategy at information set I and time $t+1$ as

$$\sigma_i^{t+1}(I, a) = \begin{cases} \frac{R_i^{t,+}(I, a)}{\sum_{a \in A(I)} R_i^{t,+}(I, a)} & \text{if } \sum_{a \in A(I)} R_i^{t,+}(I, a) > 0 \\ \frac{1}{|A(I)|} & \text{otherwise.} \end{cases}$$

Game Theory for Security

[50+@ AAAI, AAMAS, IJCAI, ICML, NeurIPS]

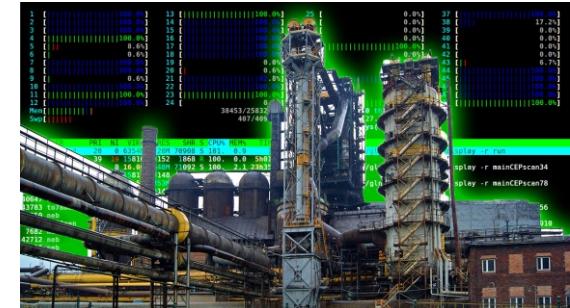
➤ Global challenges for security



Boston Marathon bombings



French oil tanker hit by a boat



Cyber physical attacks

➤ Security resource allocation

- *Limited security resources*
- *Adversary monitors defenses, exploits patterns*

➤ We pioneered the first set of applications of game theory for security resource scheduling(2007-)



- 40+ papers at premier conferences/journals, 2 best paper awards
- INFORMS Daniel H. Wagner Prize for Excellence in Operations Research Practice (2012), etc
- Operational Excellence Award from US Coast Guard (2012), etc
- Media reports: FOX News, CNN News, Federal News Radio, Defense News, The Economics Times News, Los Angeles Times, etc
- United States congressional hearing (4 times)



Security Games

- Security allocation: (i) Target weights; (ii) Opponent reaction
- Stackelberg game: Security forces commit first
- Strong Stackelberg Equilibrium
- Our contributions:
 - *Algorithms for solving large scale games*
 - *Learning adversary behavior*
 - *Applications in the real world*



Attacker

Defender

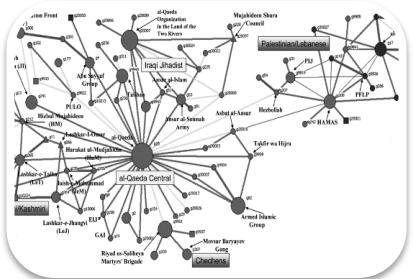


	Attacker Action #1	Attacker Action #2
Defender Action #1	7, -3	-6, 2
Defender Action #2	-3, 4	4, -6
.....
.....

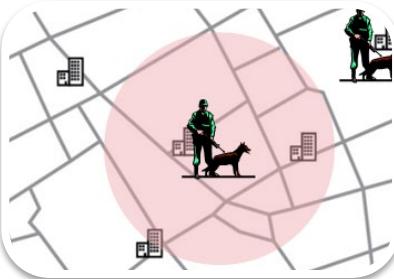
Analyzing Complex Security Games [2016-]



Dynamic Payoff



Network Games



Protection Externality



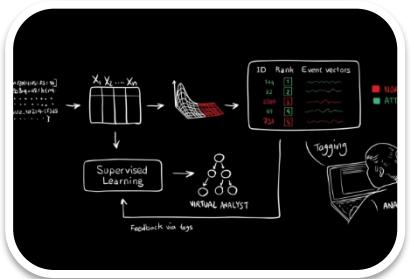
Uncertainty



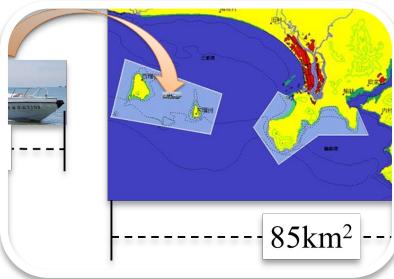
Strategic Deception



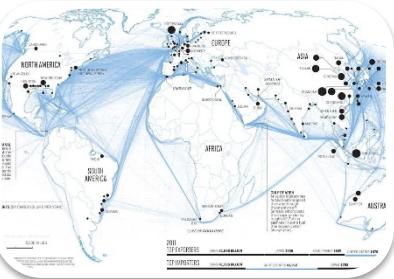
Cyber Security



Adversarial Machine Learning



Coral Reef



Nuclear Smuggling



Elections

- Combining techniques from AI, Game Theory, Operations Research ...
- Marry theory with practice
- Approaches can be applied to other domains *With proper tuning & extension*
 - ❑ Incremental strategy generation
 - ❑ Construct (multiple) equivalent games
 - ❑ Exploit compact representation
 - ❑ Abstraction
 - ❑ Tradeoff between optimality and efficiency
 - ❑ Approximation

Converging to Team-Maxmin Equilibria in Multiplayer Norm-Form Games [ICML'20]

➤ Equilibria in Multiplayer Games

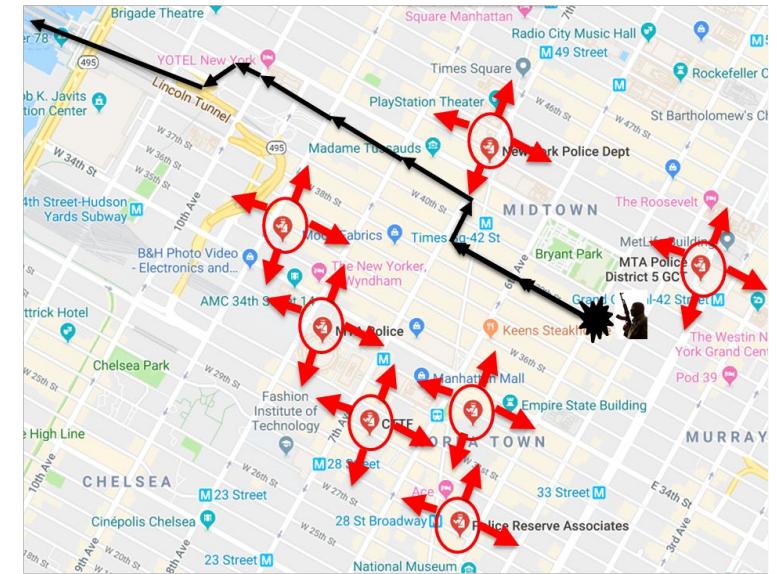
- ❑ Hard to compute: PPAD-Complete
- ❑ Hard to select: NEs are not unique
- ❑ Few results:
 - ❖ Special structure: congestion games
 - ❖ No theoretical guarantee: Pluribus (Brown and Sandholm 2019)

➤ Team-Maxmin Equilibria (von Stengel and Koller 1997)

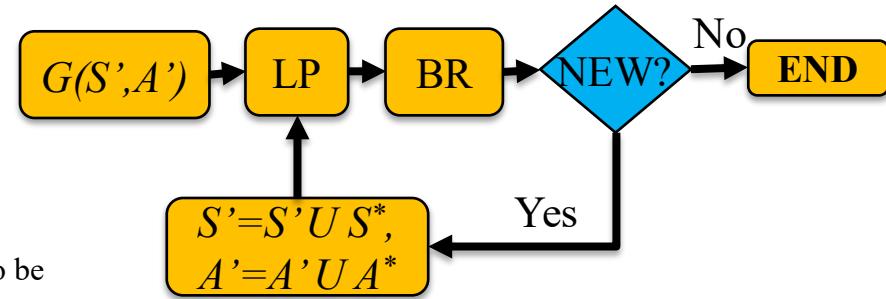
- ❑ A team of players independently plays against an adversary
- ❑ Unique in general
- ❑ FNP-hard to compute a team-maxmin equilibrium
 - ❖ Formulated as a non-convex program
 - ❖ Solved by a global optimization solver

➤ Converging to Team-Maxmin Equilibria

- ❑ Existing ISG for multiplayer games
 - ❖ Converge to an NE but many not to a TME
 - ❖ Difficult to extend the current ISG to converge to a TME
- ❑ ISGT: the first ISG guaranteeing to converging to a TME
 - ❖ Conditions in ISGT cannot be further relaxed
- ❑ CISGT: further improve the scalability
 - ❖ Initialize the strategy space by computing an equilibrium that is easier to be computed



Incremental Strategy Generation (ISG)



L×W	5×5	5×5	5×5	5×5	5×5	4×4	6×6	8×8	10×10
(p,q)	(0.8,0.6)	(0.7,0.5)	(0.6,0.4)	(0.5,0.3)	(0.4,0.2)	(0.4,0.2)	(0.4,0.2)	(0.4,0.2)	(0.4,0.2)
FullTME		∞	448s	50.4s	17.8s	0.3s	∞		
ISGT					>1000s	4s	>1000s		
CISGT	9.8s	5.9s	4.7s	3.7s	2.3s	2.2s	8.3s	24s	57s ²⁰

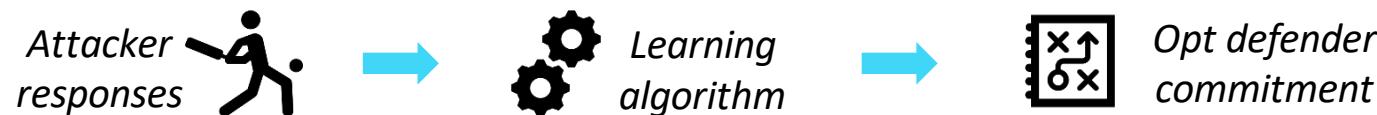
Table 1. Computing TMEs: ∞ represents out of memory.

Manipulating a Learning Defender and Ways to Counteract

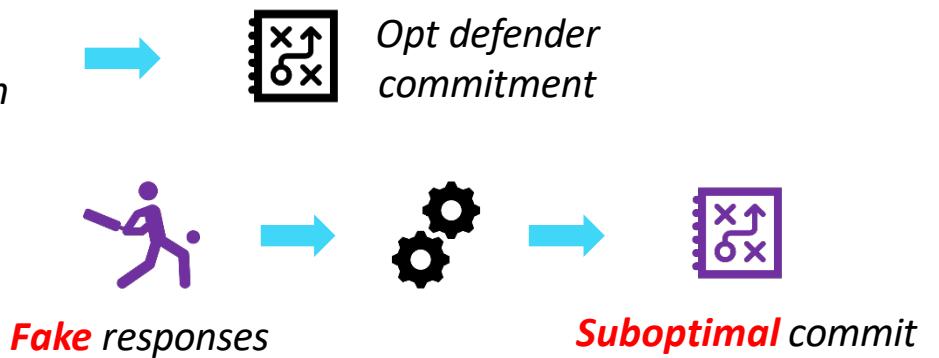
[NeurIPS'19]

- Learn to play optimally in a Stackelberg Security Game (SSG)

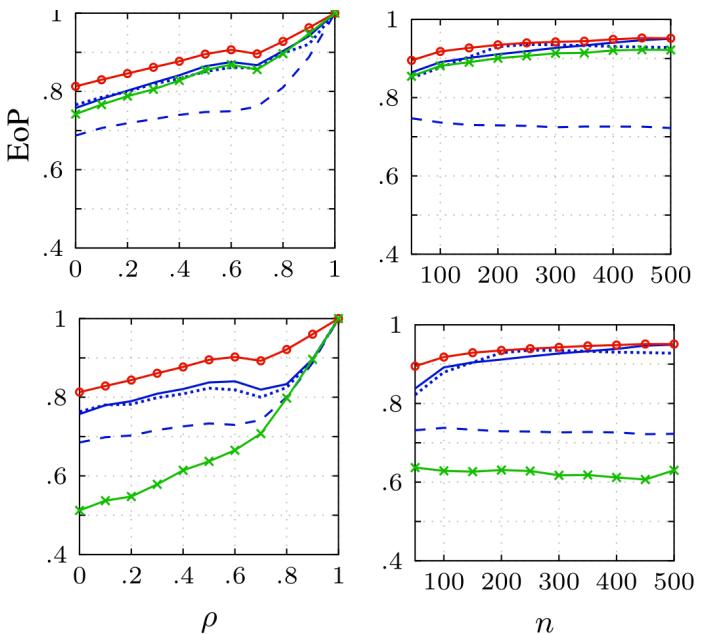
[Letchford et al., 2009; Blum et al., 2014; Haghtalab et al., 2016; Roth et al., 2016; Peng et al., 2019]



- Unrealistic assumption about *truthful* attacker responses → *What if untruthful?*



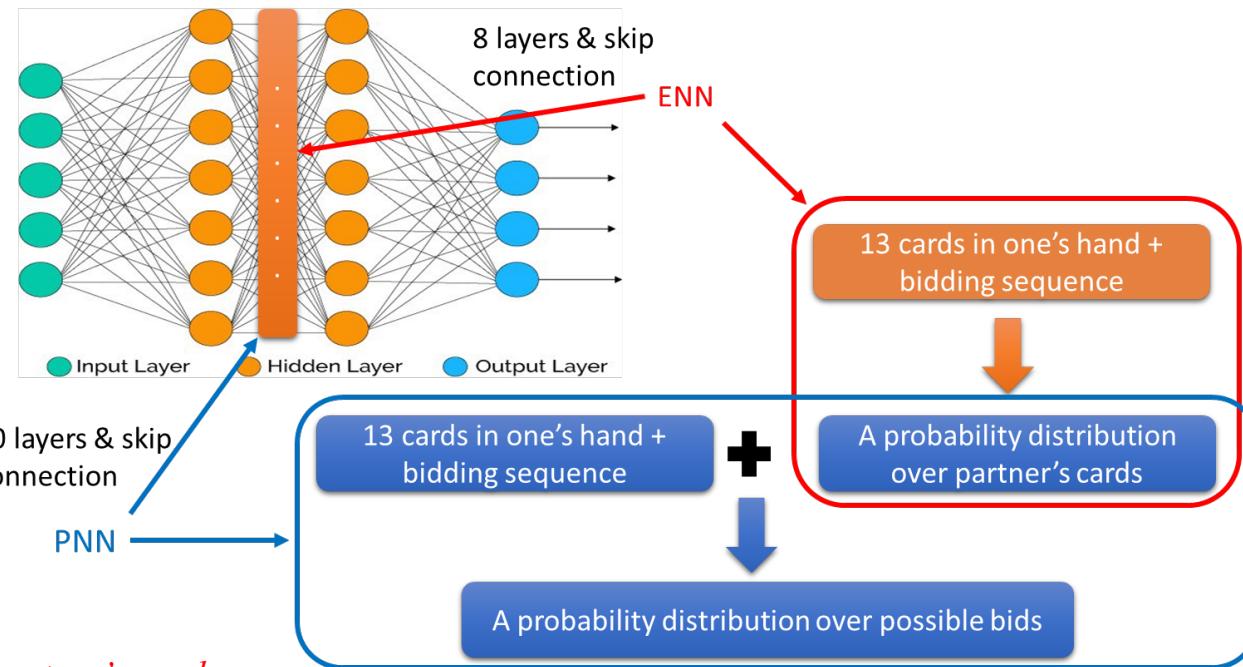
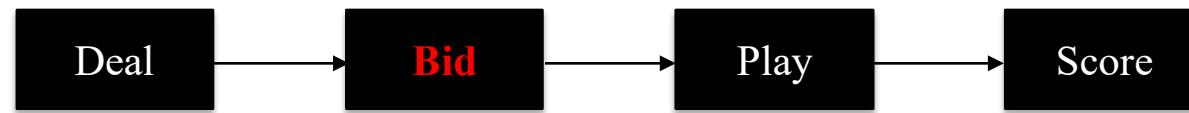
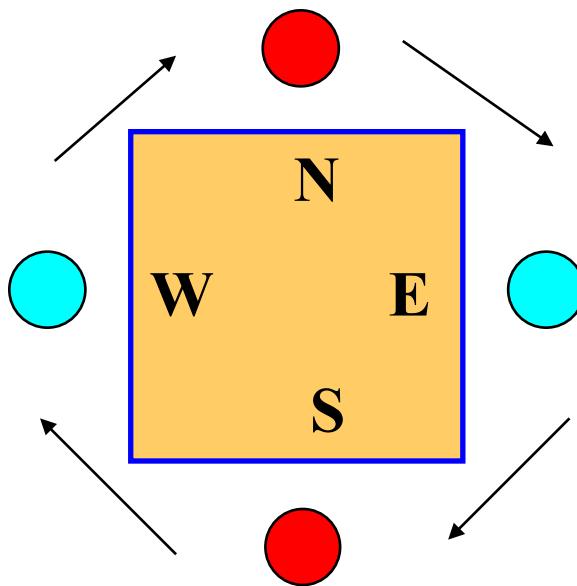
- In our work:
 - ❑ *Learning algorithms can be easily manipulated by untruthful attacker*
 - ❑ *Often optimal for attacker to deceive defender into playing a zero-sum game*
 - ❑ *A policy-based framework to play against attacker deception*
 - ❑ *A poly-time algorithm to compute optimal policy and a heuristic approach for infinite attacker types*



When do We Need RL?

- RL might be more appropriate when
 - ❑ *Problem cannot be well modelled*
 - ❑ *Large scale*
 - ❑ *Non-convex and cannot be approximated*
 - ❑ *No domain structures can be exploited*
- Multi-agent RL is receiving increasing attention
 - ❑ *Agents collaborate/compete with each other without the coordinator*
 - ❑ *Often centralized training and decentralized execution (CTDE)*
 - ❑ *MADDPG, COMA, VDN, QMIX ...*
 - ❑ *More resources:*
 - ❑ DeepMind tutorial: <https://www.karltuyls.net/wp-content/uploads/2020/06/MA-DM-ICML-ACAI.pdf>
 - ❑ Thore Graepel (DeepMind) AAMAS'20 Keynote: <https://underline.io/lecture/63-automatic-curricula-in-deep-multi-agent-reinforcement-learning>
- Does not mean multi-agent RL can always work!
- Rest of the lecture: quick overview of some of our recent works on RL
 - ❑ *Game, security, e-commerce, urban planning: competitiveness - data rich/poor*

Competitive Bridge Bidding with Deep Neural Networks [AAMAS'19]



➤ Imperfect information

- *Not necessary to infer opponents' cards*
- *Estimation neural network (ENN) to infer partner's cards*

➤ Cooperation & competition

- *Train the policy neural network (PNN) & ENN simultaneously*
- *Maintain an opponent pool for competitors*
- *Self-play with approximate reward from DDA*

➤ Large state space

- *Efficient bidding sequence representation*

	SL-PNN	RL-PNN	SL-PNN+ENN	RL-PNN+ENN	Wbridge5
SL-PNN	N/A	-8.7793	-5.653	-9.2957	-
RL-PNN	8.7793	N/A	2.1006	-1.0856	-
SL-PNN+ENN	5.653	-2.1006	N/A	-2.2854	
RL-PNN+ENN	9.2957	1.0856	2.2854	N/A	23.25

test on 10 thousand boards

test on 64 boards manually

Impression Allocation for Combating Fraud in E-Commerce via Deep RL [IJCAI'18]

➤ Fraud transactions in e-commerce:

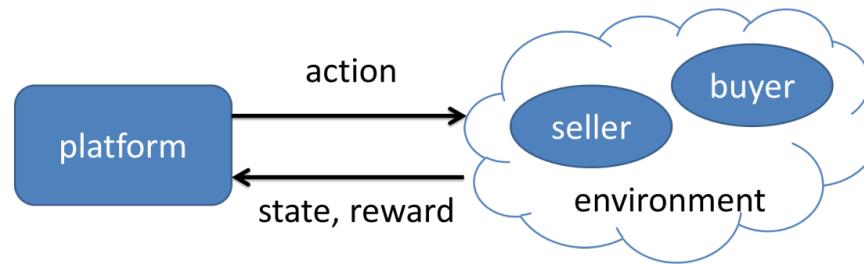
- Sellers buy their own products to fake popularity
- Deliberate, large scale fraud: *Grey industry*
- Existing approach: Machine learning for fraud detection



➤ Reducing fraud by optimal impression allocation

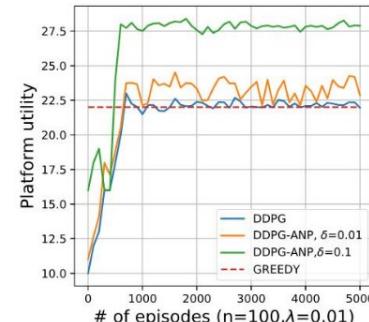
- Intuition: reward honest sellers and penalize cheating sellers

- MDP model:

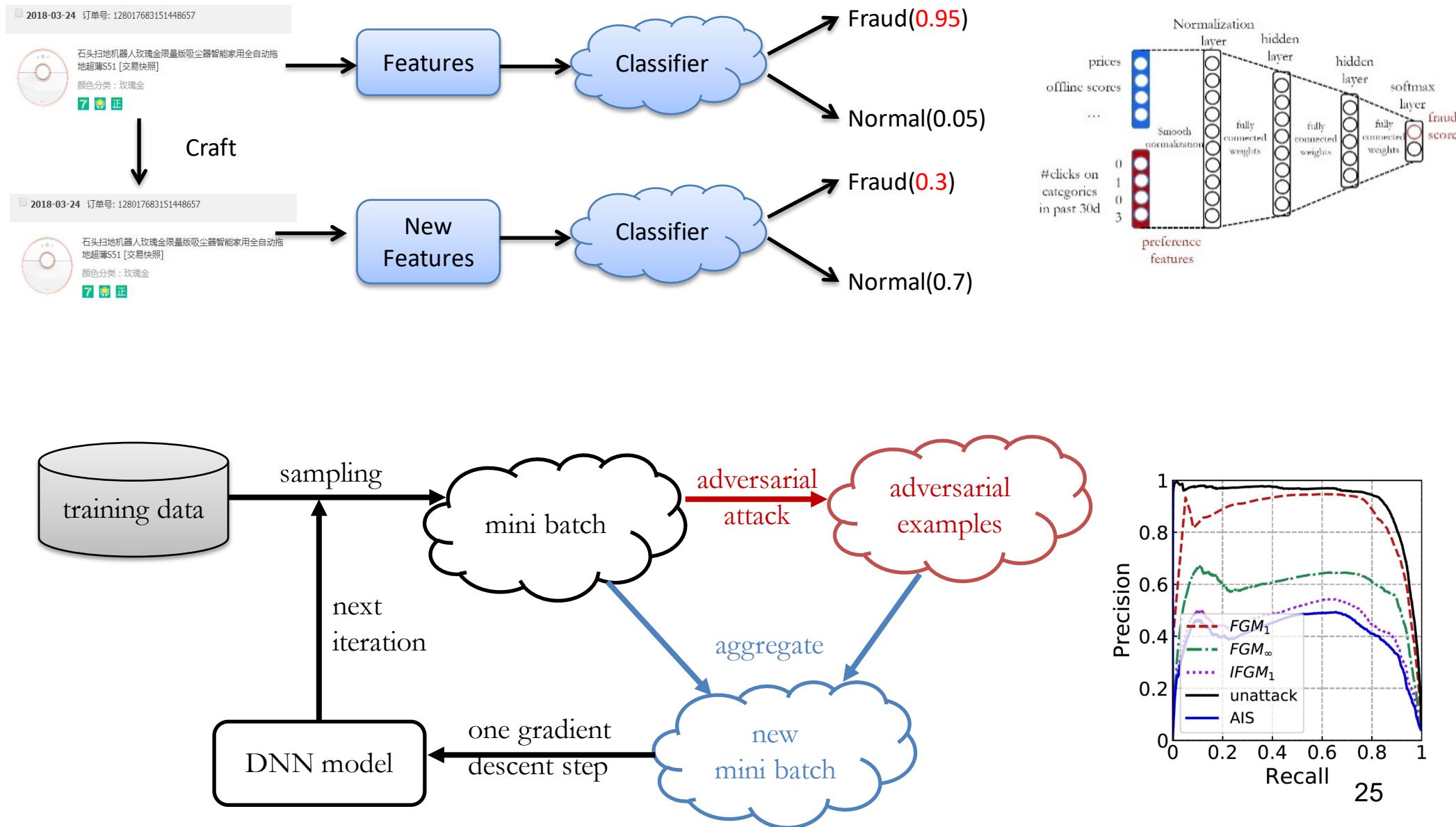


- Solving the MDP: Deep Deterministic Policy Gradient (DDPG) + reward shaping

- Shaped reward: $R(M^t, \mathbf{w}^t) = \frac{1}{n} \sum_{i=1}^n (r_i^{t+1} - \lambda f_i^{t+1}) * price_i^{t+1} - \delta \|\mathbf{w}^t\|_2$
- Avoid too large action (\mathbf{w}^t)
- Make the reward signal continuous
- Outperform all existing approaches using Alibaba's real data



Improving Robustness of Fraud Detection Using Adversarial Examples [WWW'19]



Dynamic Electronic Road Pricing (ERP) [AAAI'18]



Woodsville Tunnel (71)

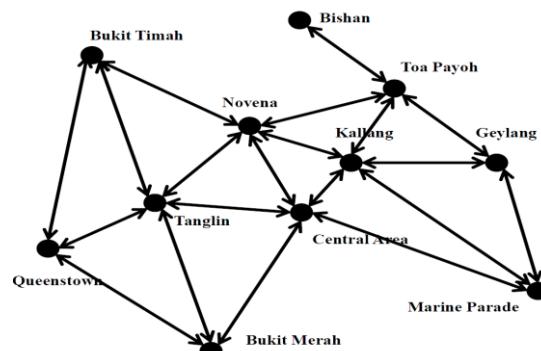
Cars/Light Goods/Taxis (Weekdays) ▾

07:00 - 07:30	\$0.00
07:30 - 08:30	\$0.50
08:30 - 08:35	\$1.00
08:35 - 08:55	\$1.50
08:55 - 09:00	\$1.00
09:00 - 09:30	\$0.50
09:30 - 22:30	\$0.00

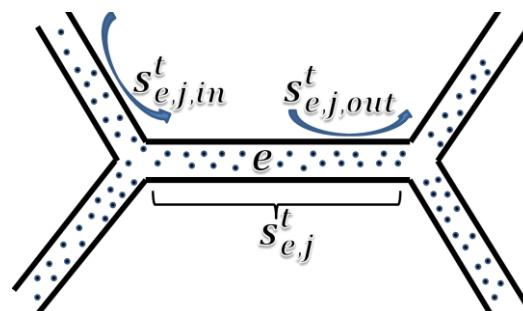
Static!!

Dynamic
ERP

Abstracted road network



An MDP formulation
(state transition)



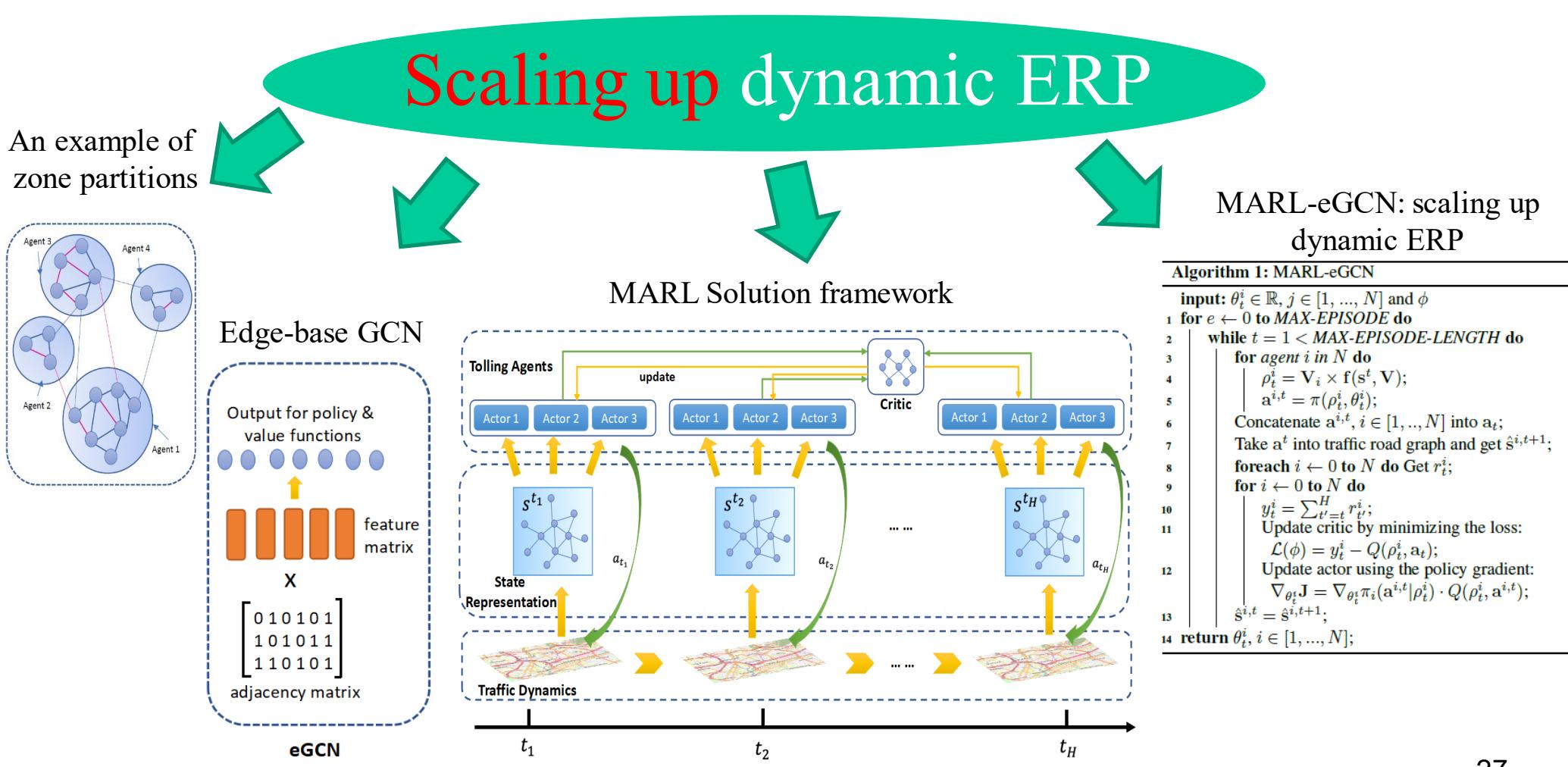
PG- β : scalable RL based on
function approximation

Algorithm 1: PG- β

```
1 Initialize  $\vartheta^t \leftarrow \vartheta_0, \theta^t \leftarrow \theta_0, \forall t = 0, 1, \dots, H-1$ ;  
2 repeat  
3   Generate an episode  $s^0, a^0, R^1, \dots, s^{H-1}, a^{H-1}, R^H$ ;  
4   for  $t = 0, \dots, H-1$  do  
5      $Q^t \leftarrow \sum_{\tau=t}^H R^\tau$ ;  
6      $\delta \leftarrow Q^t - \hat{v}(s^t, a^t, \vartheta^t)$   
7      $\vartheta^t \leftarrow \vartheta^t + \beta \delta \nabla_{\vartheta^t} \hat{v}(s^t, a^t)$   
8      $\theta^t \leftarrow \theta^t + \beta' \gamma^t \delta \nabla_{\theta^t} \log \pi^t(a^t | s^t, \theta^t)$   
9   until converge  
10  return  $\theta^t, \forall t = 0, 1, \dots, H-1$ ;
```

Scaling up Dynamic Electronic Road Pricing [IJCAI'19]

- Multi-agent reinforcement learning (MARL) for scaling up
- Edge-based graph convolutional network (GCN) for domain feature learning



Algorithm 1: MARL-eGCN

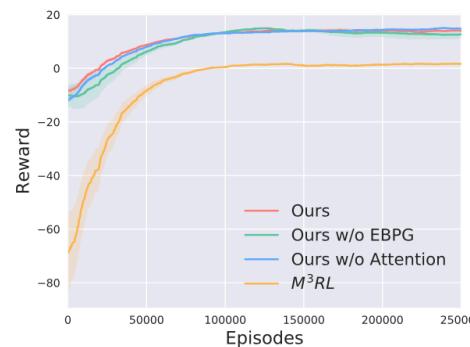
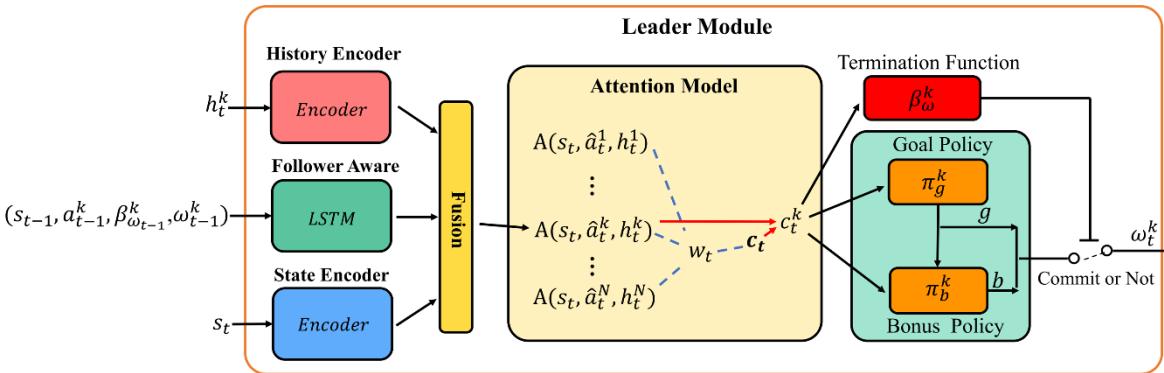
```

input:  $\theta_t^i \in \mathbb{R}$ ,  $j \in [1, \dots, N]$  and  $\phi$ 
1 for  $e \leftarrow 0$  to MAX-EPILOGUE do
2   while  $t = 1 <$  MAX-EPILOGUE-LENGTH do
3     for agent  $i$  in  $N$  do
4        $\rho_t^i = V_i \times f(s^t, V)$ ;
5        $a^{i,t} = \pi(\rho_t^i, \theta_t^i)$ ;
6       Concatenate  $a^{i,t}, i \in [1, \dots, N]$  into  $a_t$ ;
7       Take  $a_t$  into traffic road graph and get  $\hat{s}^{i,t+1}$ ;
8       foreach  $i \leftarrow 0$  to  $N$  do Get  $r_t^i$ ;
9       for  $i \leftarrow 0$  to  $N$  do
10         $y_t^i = \sum_{t'=t}^H r_{t'}^i$ ;
11        Update critic by minimizing the loss:
12         $\mathcal{L}(\phi) = y_t^i - Q(\rho_t^i, a_t)$ ;
13        Update actor using the policy gradient:
14         $\nabla_{\theta_t^i} J = \nabla_{\theta_t^i} \pi_i(a^{i,t} | \rho_t^i) \cdot Q(\rho_t^i, a^{i,t})$ ;
15    $\hat{s}^{i,t} = \hat{s}^{i,t+1}$ ;
16 return  $\theta_t^i, i \in [1, \dots, N]$ ;

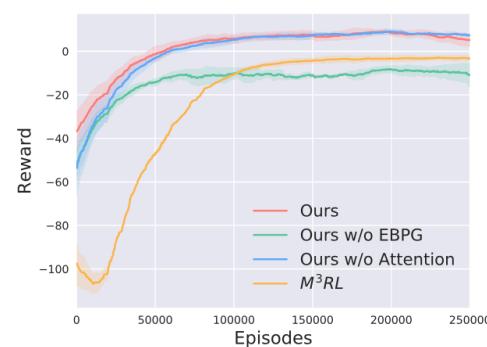
```

Learning Expensive Coordination: An Event-Based Deep RL Approach [ICLR'20]

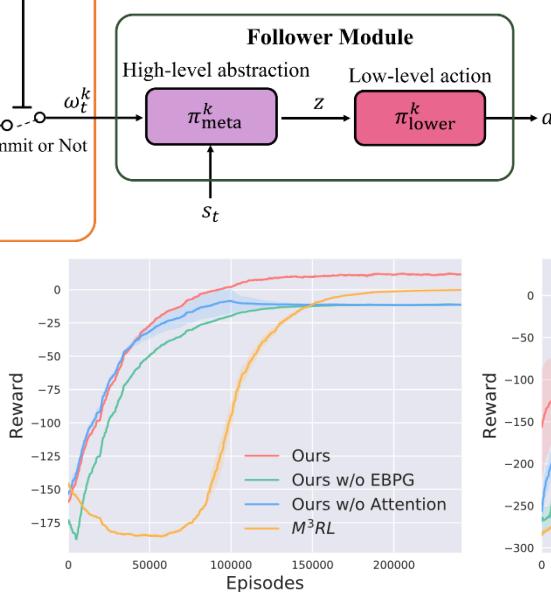
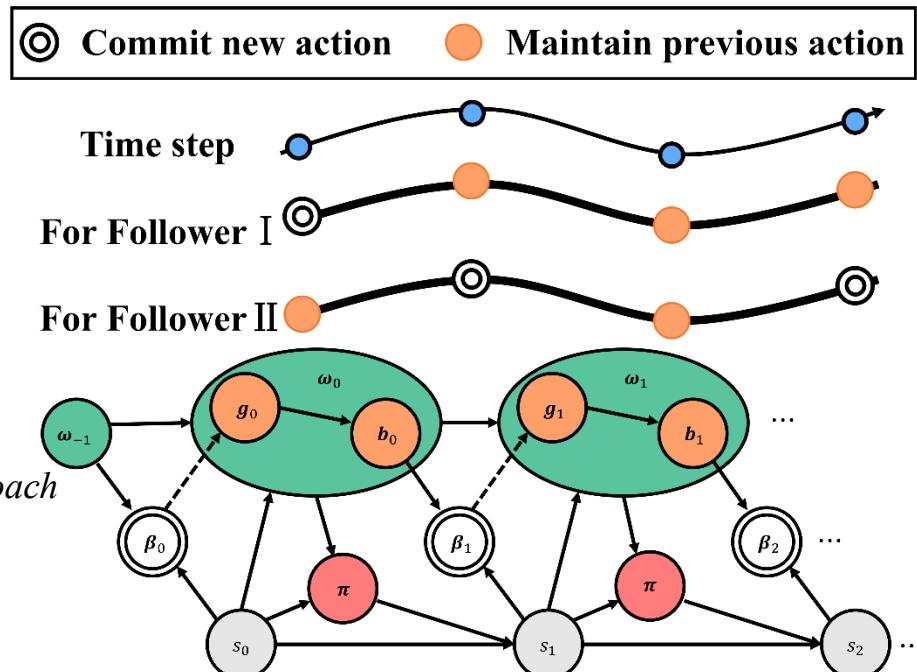
- Large state space based expensive coordination
 - Considering the followers are *self-interested*
 - The leader should coordinate them by assigning incentives
- Event-based policy gradient
 - Modeling the leader's strategy as *events*
 - A novel event-based policy is induced based on the events
- Action abstraction for followers
 - Accelerating the training process through action abstraction approach



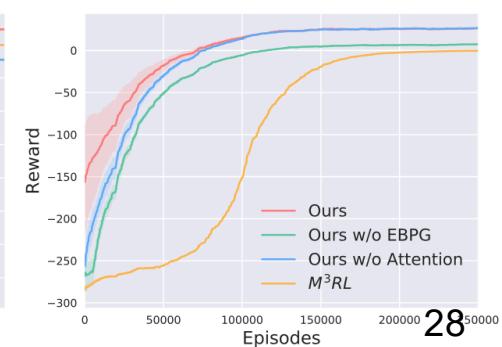
(a) Resource Collections.



(b) Multi-Bonus Resource



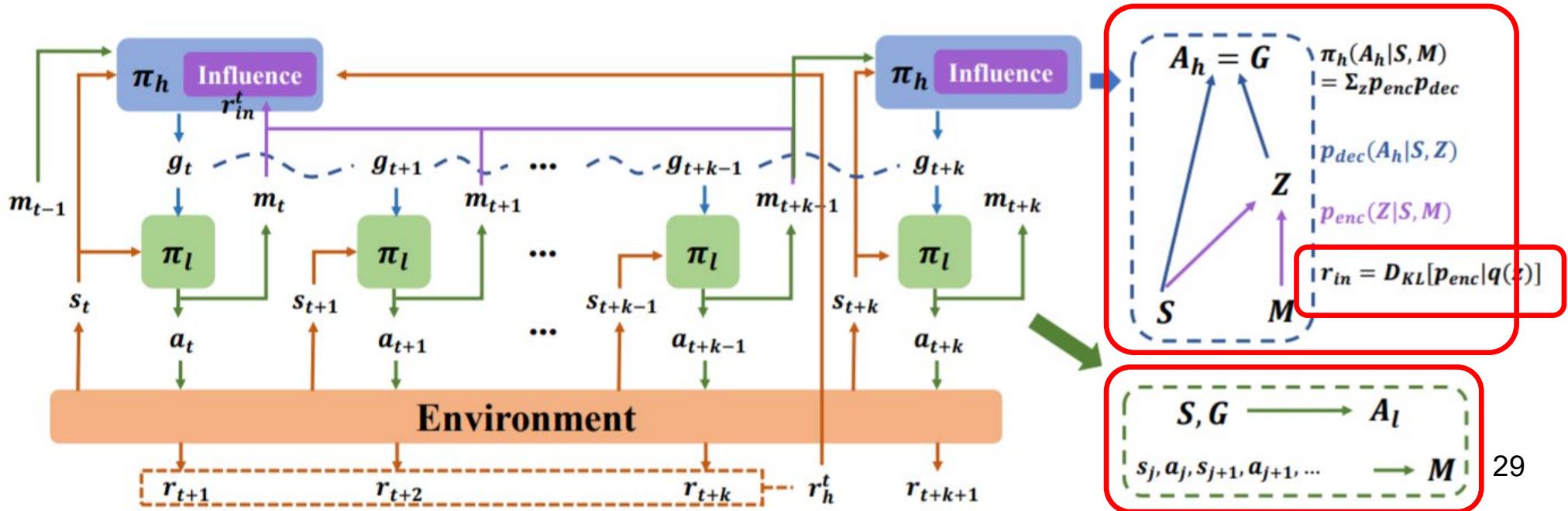
(c) Navigation.



(d) Predator-Prey.

Interactive Influence-based HRL [IJCAI'20]

- Challenge: high-level policy suffers the **non-stationarity** problem
 - *the constantly changing low-level policy leads to different state transitions and rewards of high-level policy*
- Key idea: the high-level policy makes decisions conditioned on the low-level policy representation
- Solutions:
 - *Low-level policy modeling*
 - *High-level policy stabilization via information-theoretic regularization*
 - *Influence-based adaptive exploration with auxiliary rewards*



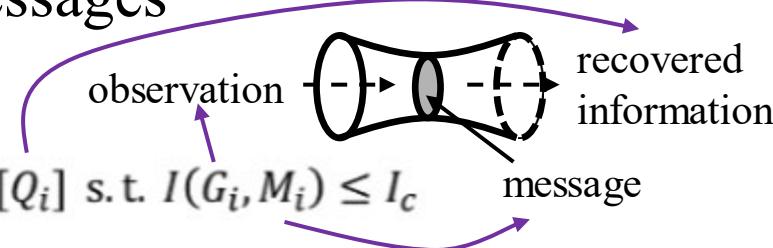
Efficient Multi-agent Communication [ICML'20]

- Limited bandwidth in multi-agent communication:



- Theorem: limited bandwidth constrains the message's entropy
- Key idea: compressing the communication messages
- Information bottleneck principle

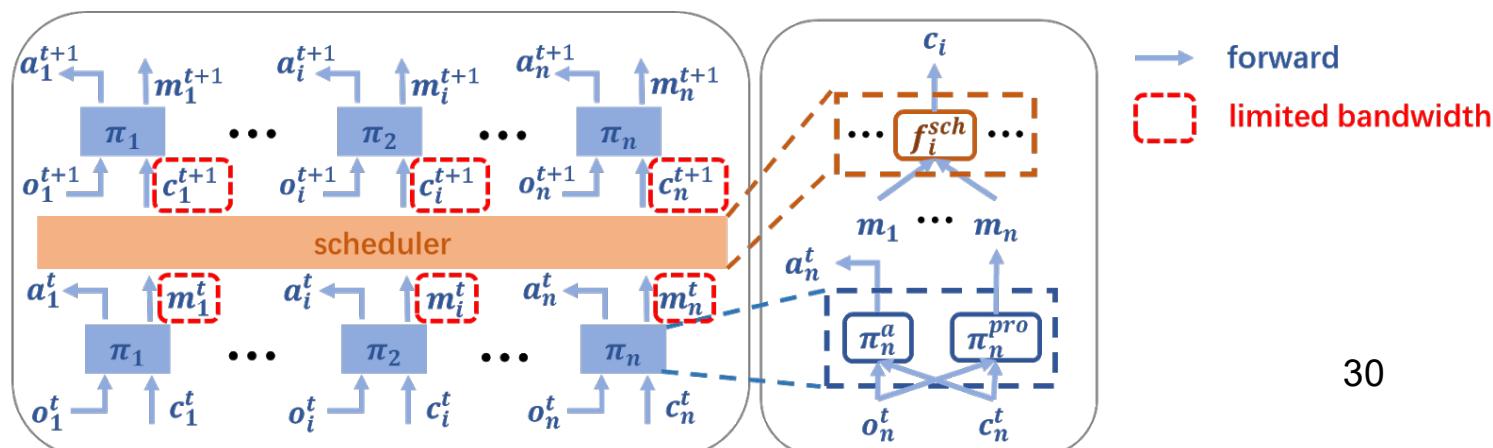
$$\max_T I(T, Y) \text{ s.t. } I(X, T) \leq I_c \rightarrow \max_{M_i} J(\theta_i) := E_{\pi_i}[Q_i] \text{ s.t. } I(G_i, M_i) \leq I_c$$



□ lower bound for control entropy

$$\tilde{J}(\theta_i) = E_{\pi_i}[Q_i] - \beta E \left[D_{KL}[\pi_i^{pro}(m_i|g_i)||z(m_i)] \right]$$

- Architecture:



Learning Behaviors with Uncertain Human Feedback

[UAI'20]

➤ Uncertain Human Feedback

- Human may not give any feedback
- Positive feedback → sub-optimal action
- Negative feedback → optimal action

$$p(a, \lambda(s); \sigma, \mu, f) = \begin{cases} p^+(a, \lambda(s); \sigma, \mu^+) & f = f^+ \\ p^-(a, \lambda(s); \sigma, \mu^-) & f = f^- \\ 1 - p^+(a, \lambda(s); \sigma, \mu^+) \\ -p^-(a, \lambda(s); \sigma, \mu^-) & f = f^0 \end{cases}$$

p^+ and p^- are functions like Gauss function where λ is the **optimal action**, μ and σ are unknown parameters controlling **mean and variance**.

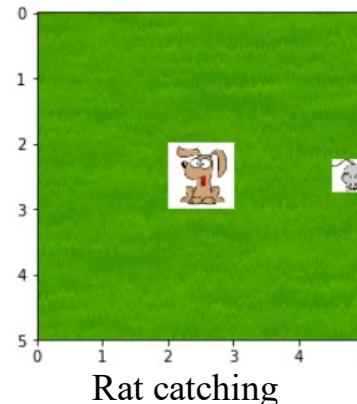
➤ Inferring Optimal Behavior

- Maximizing likelihood estimation of receiving different kinds of feedback:

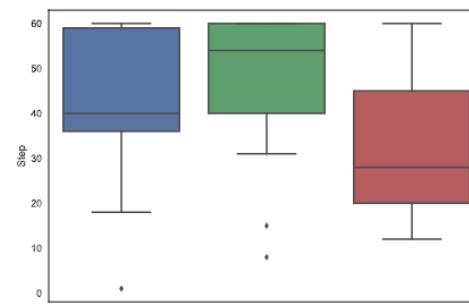
$$\arg \max_{\lambda} P(h|\lambda; \mu, \sigma)$$

- EM + GD:

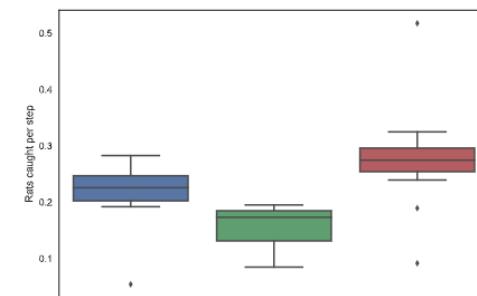
- EM updates λ with σ fixed
- μ is latent variable in integral
- GD updates σ with λ fixed
- μ is eliminated by a trick



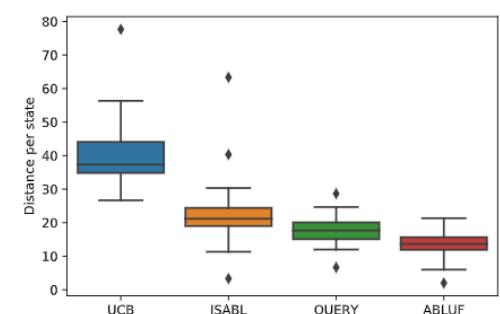
Light controlling



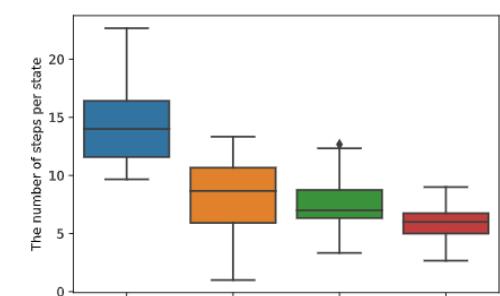
(a) Step



(b) Rats caught per step



(a) Distance/state



(b) #steps/state

We Won Microsoft Collaborative AI Challenge

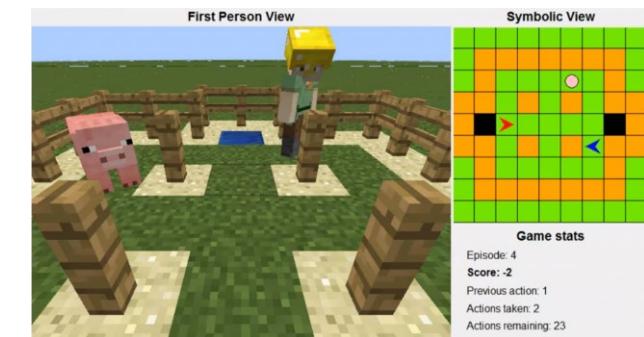
[AAAI'18]

➤ Collaborative AI

- ❑ *How can AI agents learn to recognize someone's intent (that is, what they are trying to achieve)?*
- ❑ *How can AI agents learn what behaviors are helpful when working toward a common goal?*
- ❑ *How can they coordinate or communicate with another agent to agree on a shared strategy for problem-solving?*

➤ Microsoft Malmo Collaborative AI Challenge

- ❑ *Collaborative mini-game, based on an extension "stag hunt"*
- ❑ *Uncertainty of pig movement*
- ❑ *Unknown type of the other agent*
- ❑ *Detection noise (frequency 25%)*
- ❑ *Efficient learning*



➤ Our team HogRider won the challenge (out of more than 80 teams from 26 countries)

- ❑ *learning + game theoretic reasoning + sequential decision making + optimization*

Catch the pig by HogRider

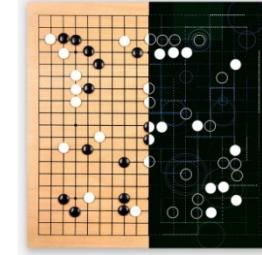
AMI Research Group
Nanyang Technological University



AI for Complex Interaction

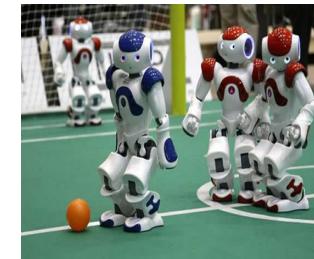
- Recent AI breakthrough

IMAGENET



- What's next: AI for *complex* interaction

- Stochastic, open environment
- Multiple players
- Sequential decision, online
- Strategic (selfish) behavior
- Distributed optimization



Google™
bing

- MARL will play a very important role...