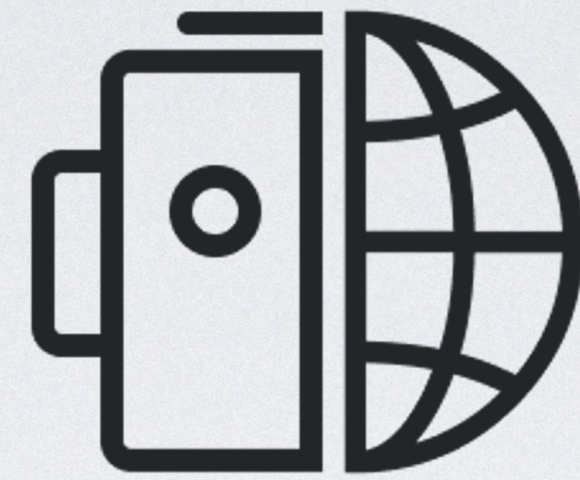


Reinforcement Learning China Summer School



RLChina 2021

# MULTI-AGENT LEARNING BASICS

Dr. Yaodong Yang  
Assistant Professor  
King's College London  
[www.yangyaodong.com](http://www.yangyaodong.com)  
08/2021

# Intelligence is learning from mistakes!



“... if a machine is expected to be infallible, it cannot also be intelligent. There are several mathematical theorems which say almost exactly that. But these theorems say nothing about how much intelligence may be displayed if a machine makes no pretence at infallibility...”

— Alan Turing, 1947

# Remarkable Success of MARL in Gaming AI Applications

Great advantages have been made since **2019!**

Jan 2016

Dec 2017

July 2018

Jan 2019

Apr 2019

July 2019

Sep 2019

July 2021

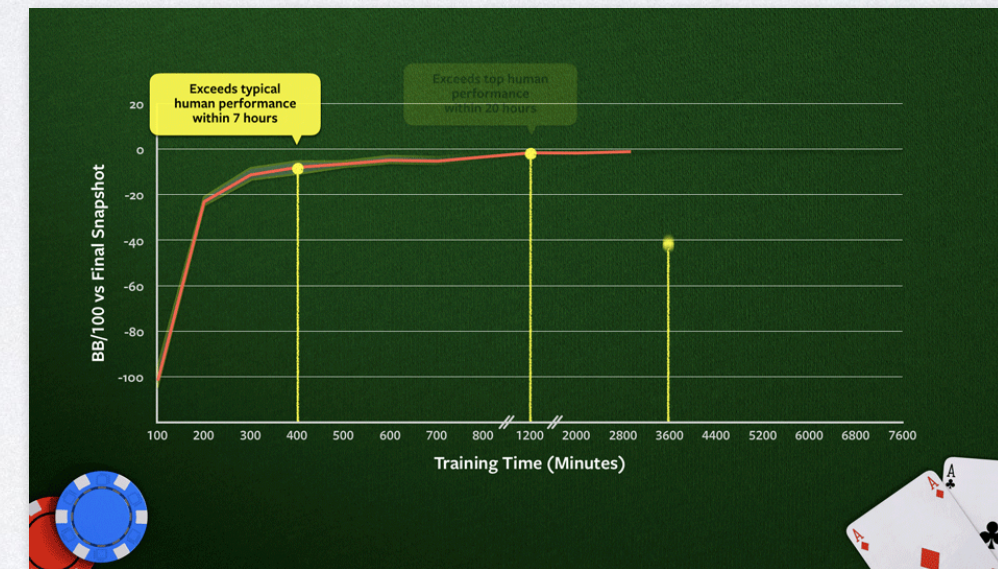
AlphaGO Series



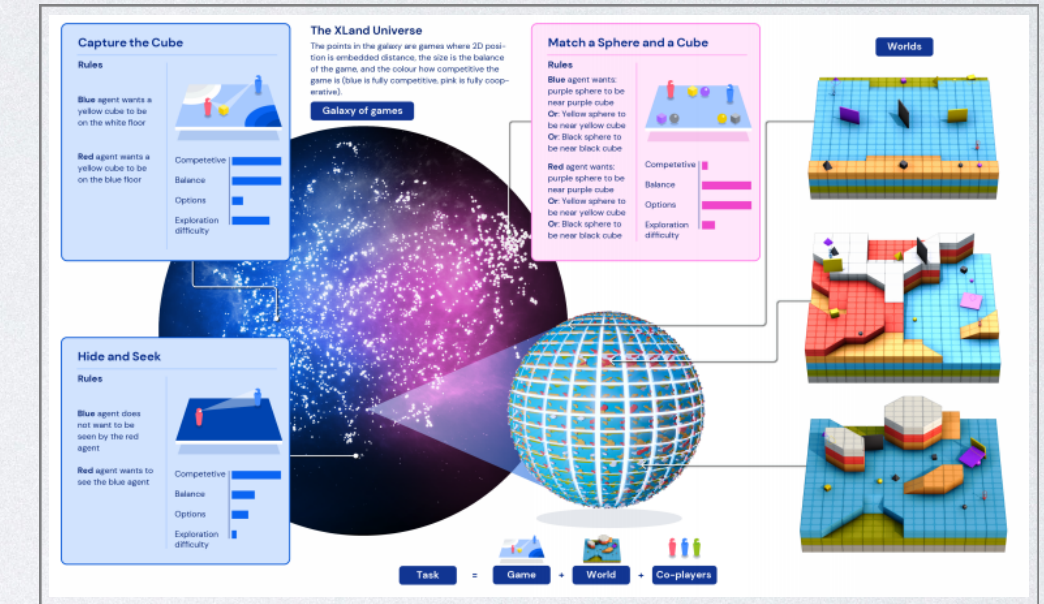
AlphaStar (DeepMind)



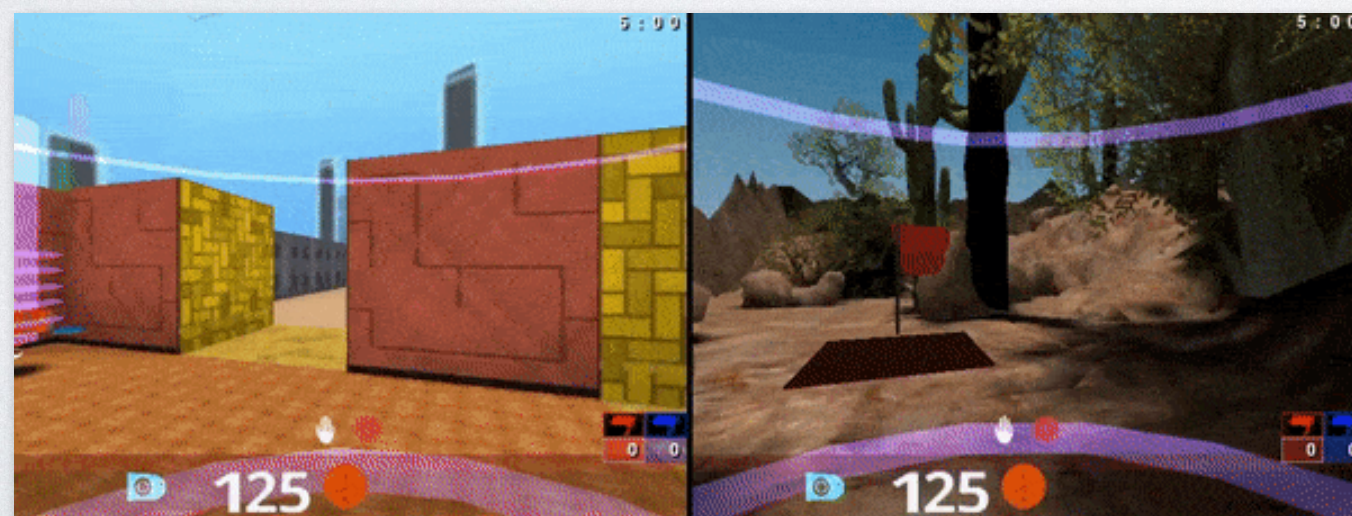
Pluribus Poker (FAIR)



XLand (DeepMind)



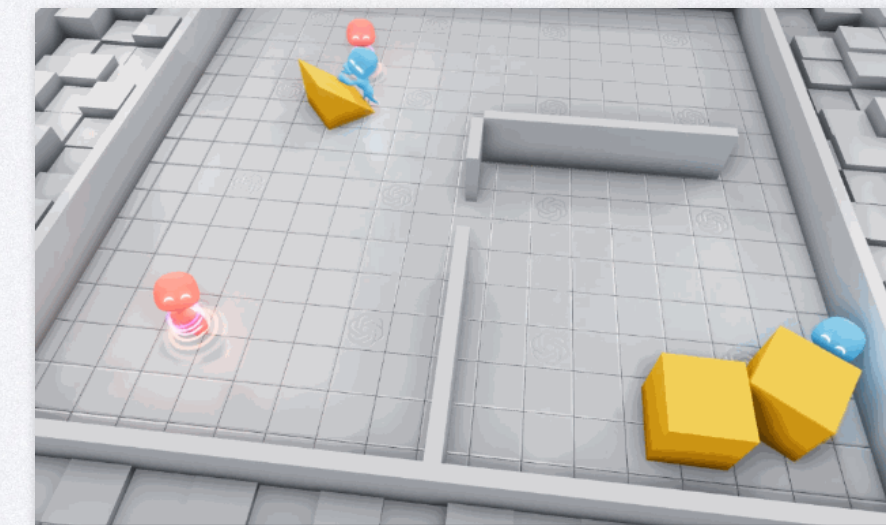
Capture-the-flag (DeepMind)



Dota2 (OpenAI)



Hide and Seek (OpenAI)



Multi-agent intelligence emerges

# Multi-agent Intelligence Components

Multi-Agent Intelligence

Autonomous Driving  
Gaming AI  
Smart Grid / City

...

||

Algorithmic Game Theory Fundamentals

Equilibrium/Solution Concept  
Learning Dynamics Analysis  
Mechanism Design

...

+

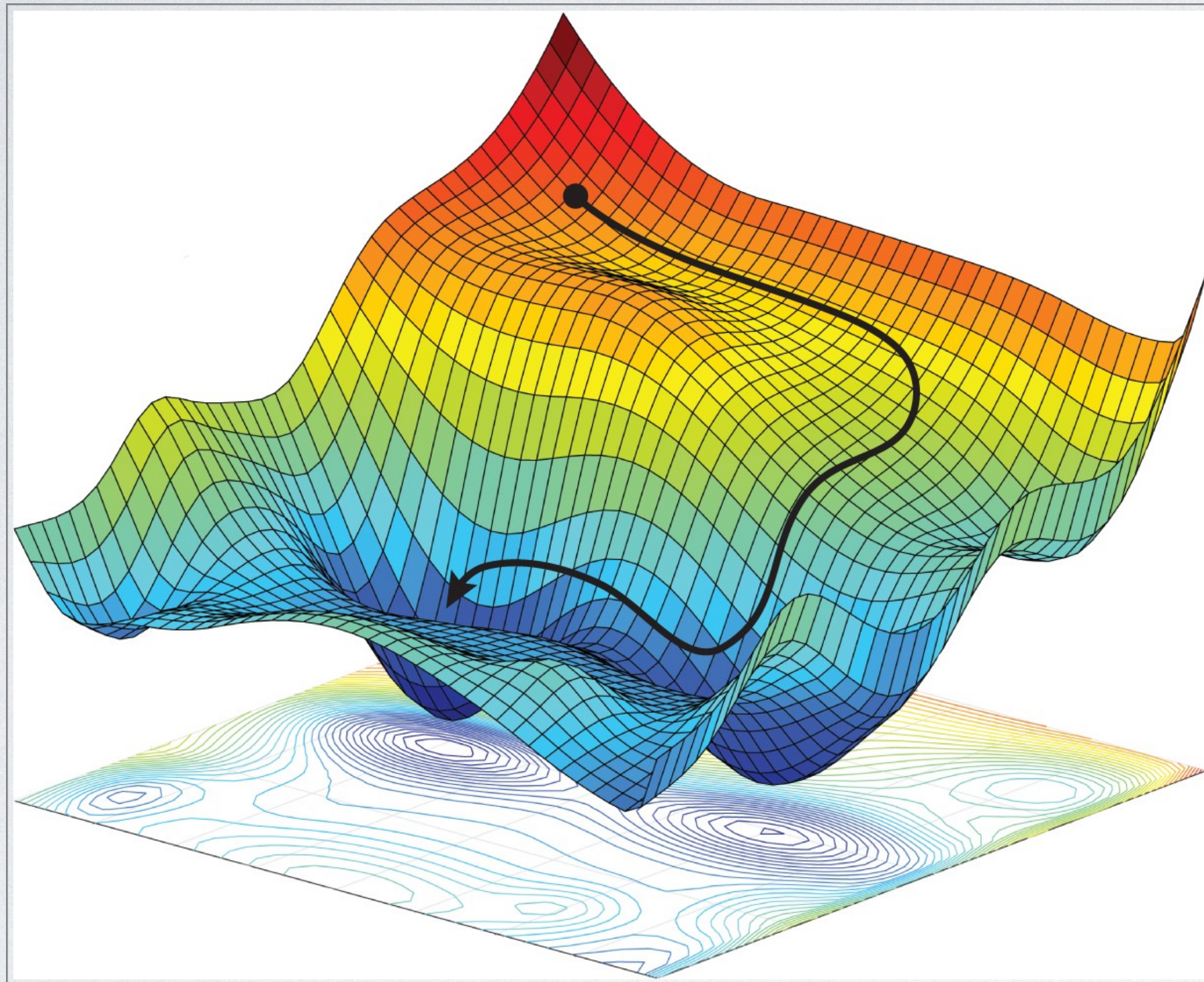
Machine Learning Techniques

Reinforcement Learning  
Deep Learning  
Auto-differentiation

...

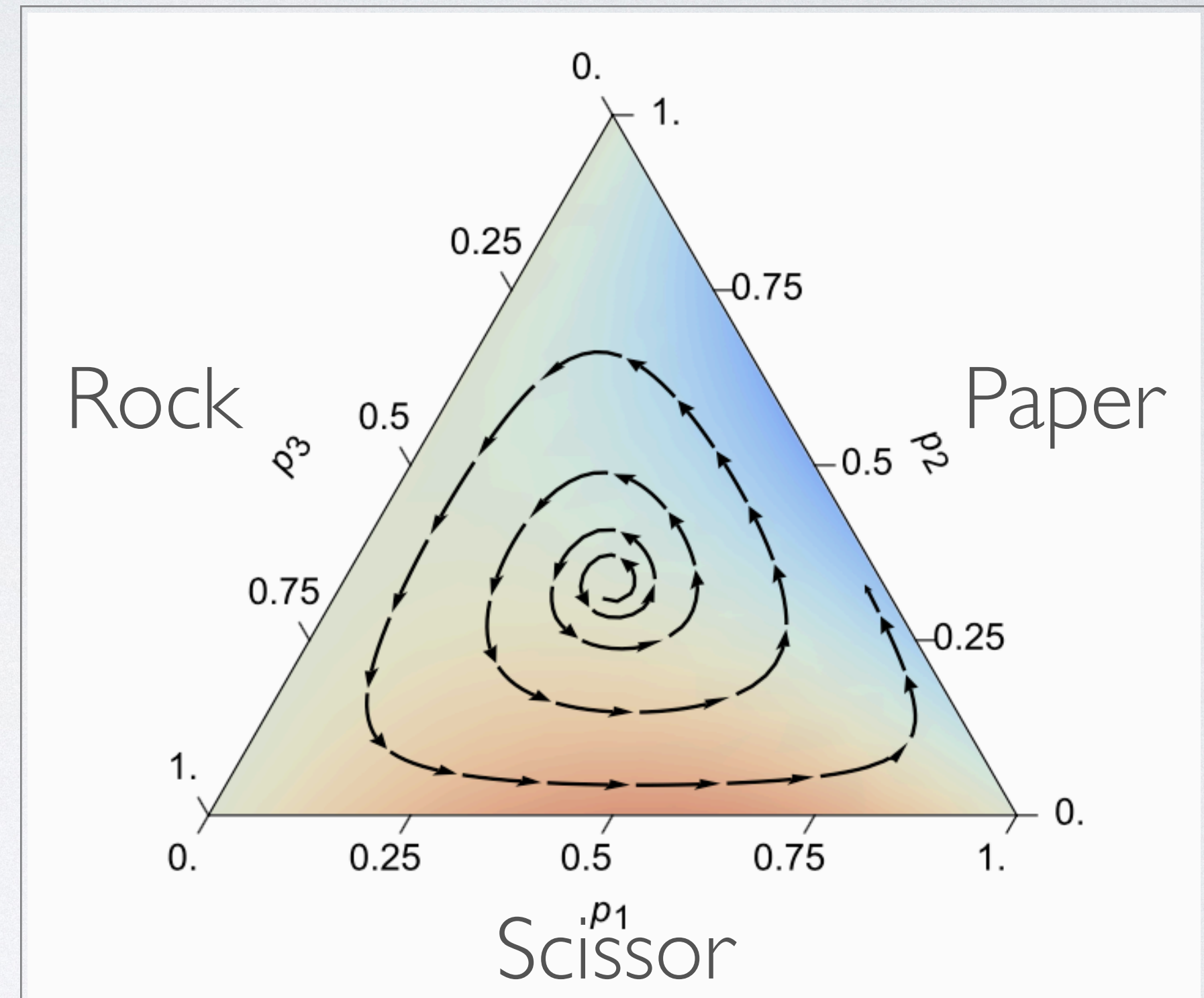
# Game Theory Lays the Foundation for Multi-agent Learning

Normal machine learning problems:



$$\|\nabla f(x)\|_2 = 0$$

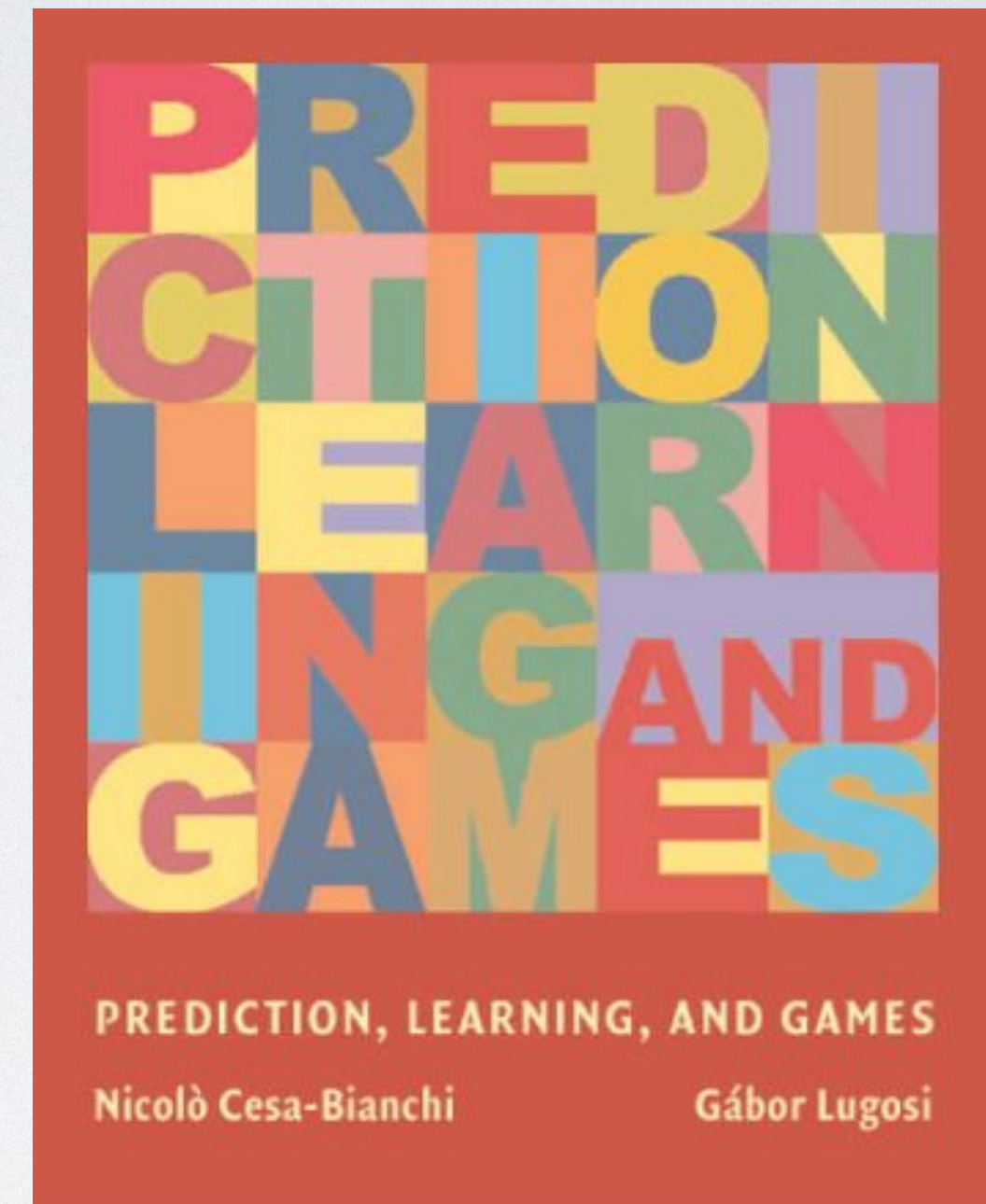
Multi-agent learning problems:



the learning outcomes are described by game theory

# Recommended Resources:

- A self-contained MARL survey from game theoretical perspective:
  - <https://arxiv.org/abs/2011.00583>
- Textbook
  - <Prediction, learning and games> Nicolò Cesa-Bianchi
- Algorithmic Game Theory lectures:
  - <http://www.cs.jhu.edu/~mdinitz/classes/AGT/Spring2020/>  
(Uncited screenshots refer to Lectures 1,2,3,4,5,6,7,8,9)
- If you want to know more details about modern MARL methods
  - [Talk: A General Solver to Two-player Zero-sum Games](#)
  - [Talk: Recent advances of MARL in Gaming AI](#)
  - [Talk: Dealing with Non-transitivity in Two-player Zero-sum Games](#)
- If you want to get hands to solve real-world games, e.g., Poker/Chess
  - <https://github.com/sjtu-marl/malib>



# Contents

- Algorithmic Game Theory Overview
- Computing Nash Equilibrium
  - Two-Player Zero-Sum Games (LP, fictitious play, double oracle, PSRO)
  - Two-Player General-Sum Games (support enumeration, Lemke-Howson method)
  - N-Player Potential Games (best response dynamics)
- Connections to MARL
  - MARL Formulations
  - Complexity Results (Nash is PPAD-hard)
  - Other Necessary Solution Concepts (correlated equilibrium, coarse CE)
- No-Regret Dynamics
  - Solving Coarse Correlated Equilibrium
  - Solving Two-Player Zero-Sum Games (FP is not no-regret, MWU, ODO)
  - Solving Correlated Equilibrium (swap regret)

# Algorithmic Game Theory Overview

- Game theory studies the interaction between rationale (selfish) agents.
- It is an area between CS and Economics.
- MARL is to study algorithmic games theory with powerful machine learning tools.

## 1. Computing Equilibria

- Is it reasonable to model behaviours through different equilibrium concept (e.g, the “invisible hands”)
- how can we compute the equilibrium efficiently and distributedly ?

## 2. Understanding the Inefficiency of Equilibrium

- Is the equilibrium “optimal” compared to maximal social welfare ?
- bounding and deriving the distance to optimality

## 3. Mechanism Design

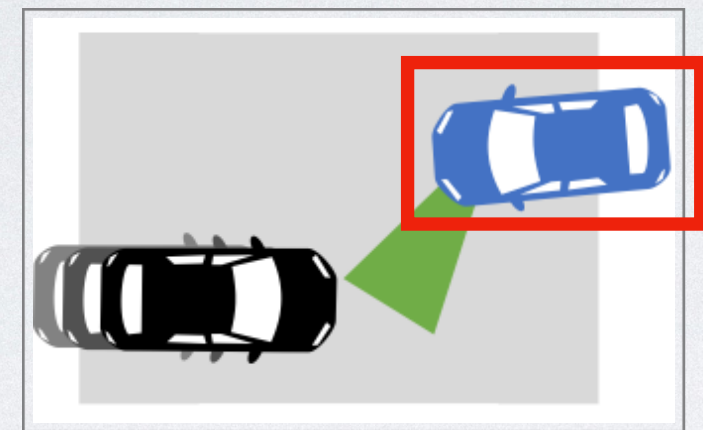
- The science of “rule making”: can we design games so that selfish agents can lead to good outcomes ?
- Heavily focus on auctions: how to design the auction rules to incentivise agents to tell the truth.



# Computing Equilibria

Traffic intersection is naturally a multi-agent system. From each driver's perspective, in order to perform the optimal action, he must take into account others' behaviours.

what you see



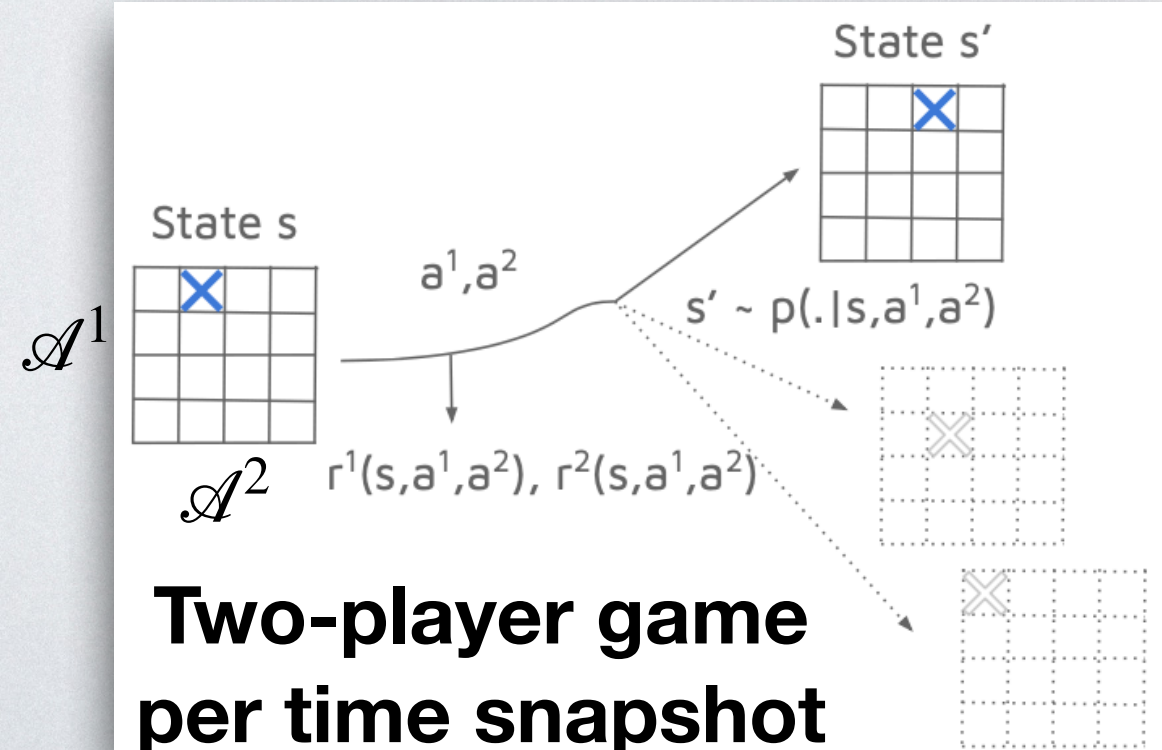
scenario

Yield  
Rush

	Yield	Rush
Yield	(0, 0)	(1, 2)
Rush	(2, 1)	(0, 0)

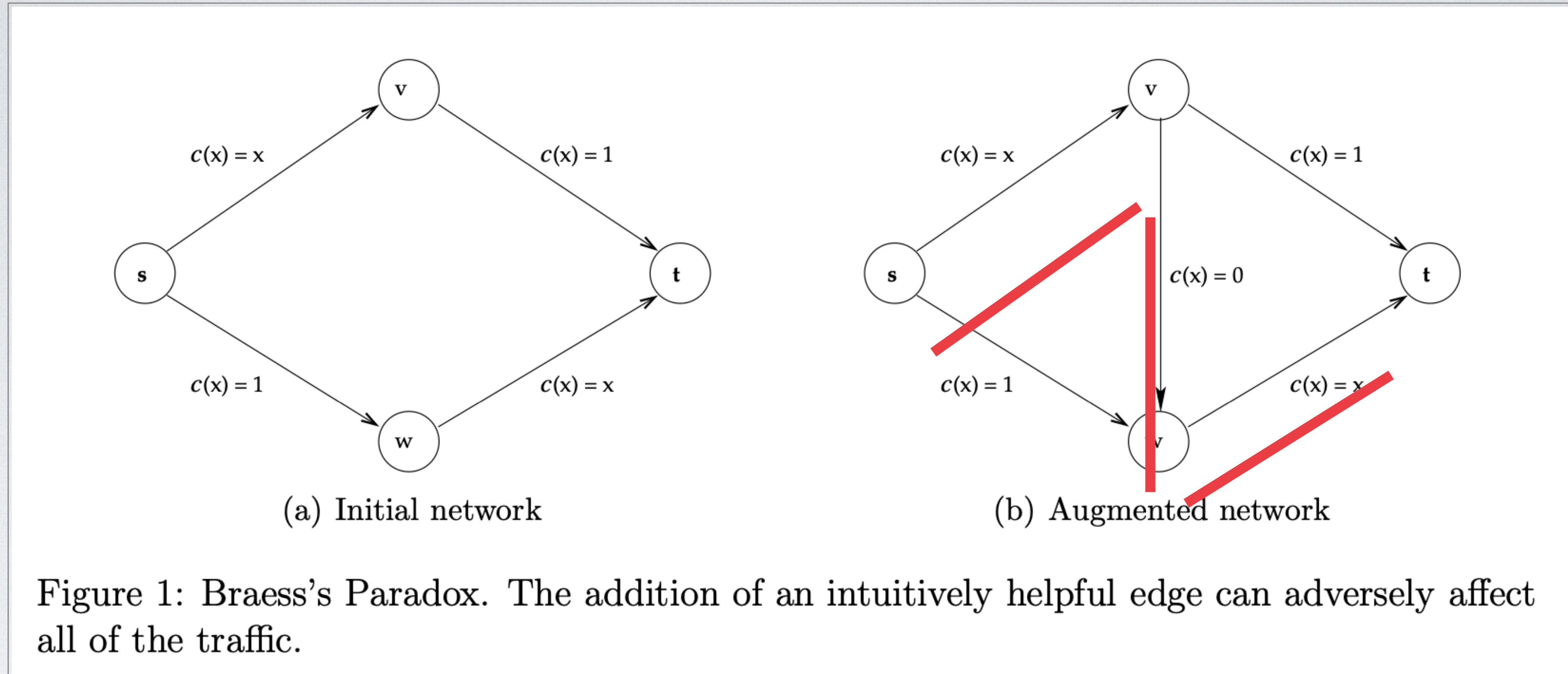
normal-form game

what the computer sees



- When the drivers are rational, they will reach the outcome of a Nash Equilibrium. It is the outcome of interaction. Knowing it can predict future.
- Real-world decision making has cooperation & competition. For each agent, how to infer the belief of the other agents and make the optimal action is critical.
- The concept of using traffic light is in fact a **correlated equilibrium**.
- Many-agent system is when agents  $\gg 2$ . It is a very challenging problem to compute equilibrium, thus making decisions.

# The Inefficiency of Equilibria



optimal traffic: half  $s \rightarrow v \rightarrow t$ , half  $s \rightarrow w \rightarrow t$   
 cost per person:  $1 + 1/2$

optimal traffic: all goes  $s \rightarrow v \rightarrow w \rightarrow t$   
 cost per person: 2

price of anarchy:  $\frac{\text{value of worst Nash}}{\text{optimal outcome}} = \frac{2}{3/2} \geq 1$

Selfish behaviours can lead to inefficient equilibrium ! Can we bound them for real-world problems?

# Mechanism Design

- Suppose we are to organise an auction, each player's utility is set as  $(\text{valuation} - \text{final price})_+$ , how can we set the auction rule: the mechanism of determining the final price? We want it to be **truthful** that all players bids their valuations.
- Highest price auction:
  - ◆ Each player will bid less than valuation, if he wins the bids, he will try to decrease the price.
  - ◆ The final price depends on others' valuations. It is unclear for both players and auctioneer to practice
- Second-price auction:
  - ◆ Player bids the valuation price is the dominant strategy. Assuming  $b_j$  to be the highest price.
    - if  $v_i < b_j$ , then 0 utility is better than negative, thus bidding  $v_i$
    - if  $v_i > b_j$ , bidding  $v_i$  is the always dominant than bidding other numbers.
- In many real-world problems, how can we design **truthful mechanism** that meets **computation constraint**.
- See Zhengyang's talk at RLChina for this topic.

# Contents

- Algorithmic Game Theory Overview
- Computing Nash Equilibrium
  - Two-Player Zero-Sum Games (LP, fictitious play, double oracle, PSRO)
  - Two-Player General-Sum Games (support enumeration, Lemke-Howson method)
  - N-Player Potential Games (best response dynamics)
- Connections to MARL
  - MARL Formulations
  - Complexity Results (Nash is PPAD-hard)
  - Other Solution Concepts (correlated equilibrium, coarse CE)
- No-Regret Dynamics
  - Solving Coarse Correlated Equilibrium
  - Solving Two-Player Zero-Sum Games (FP is not no-regret, MWU, ODO)
  - Solving Correlated Equilibrium (swap regret)

# Nash Equilibrium

- Let  $n$  players,  $S = S_1 \times \dots \times S_n$  is the joint strategy profile,  $u_i : S \rightarrow \mathbb{R}$  is the utility function, Nash equilibrium is

$$\mathbf{E}_{s_j \sim \mu_j \forall j \in [n]} \left[ u_i(s_1, \dots, s_n) \right] \geq \mathbf{E}_{\substack{s_i \sim \mu'_i \\ s_{-i} \sim \mu_{-i}}} \left[ u_i(s_1, \dots, s_n) \right] \quad \forall \mu'_i \in \Delta_{S_i}$$

- Mixed strategy Nash equilibrium always exists in **finite player finite action** games.
- For continuous utility games, the strategy set needs to be compact
- Note that  $\mu'_i \in \Delta_{S_i}$  can be replaced by  $a \in S_i$  because deviation is at most a pure strategy !
- Here no state transitions are considered. In Markov game, the solution concept is **Markov Perfect Equilibrium**.

- The expectation is computed as follows
- $$\begin{aligned} x^T A y &= \sum_{i=1}^N \sum_{j=1}^M A_{ij} x_i y_j \\ &= \sum_{i=1}^N \sum_{j=1}^M A_{ij} \mathbf{Pr}[\text{player 1 plays } i] \mathbf{Pr}[\text{player 2 plays } j] \\ &= \mathbf{E}_{\substack{i \sim x \\ j \sim y}} [u_i(i, j)] \end{aligned}$$

# Nash Equilibrium in Two-Player Zero-Sum Games

- **von Neumann theorem:** Two-player Nash can be computed in P-time through linear programmes (LP).

## Prime problem

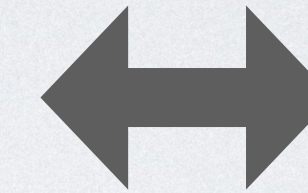
row player maximises the worst situation

$$\begin{aligned} \max_{v \in \mathbb{R}} v \\ \text{s.t. } \mathbf{p}^T \mathbf{A} \geq v \cdot \mathbf{1} \\ \mathbf{p} \geq \mathbf{0} \text{ and } \mathbf{p}^T \mathbf{1} = 1 \end{aligned}$$

## Dual problem

column player's view

$$\begin{aligned} \min_{v \in \mathbb{R}} v \\ \text{s.t. } \mathbf{q}^T \mathbf{A}^T \leq v \cdot \mathbf{1} \\ \mathbf{q} \geq \mathbf{0} \text{ and } \mathbf{q}^T \mathbf{1} = 1 \end{aligned}$$



## Minimax theorem

zero-duality gap for convex problems

$$\begin{aligned} \max_{\mathbf{p}} \min_{\mathbf{q}} \mathbf{p}^T \mathbf{A} \mathbf{q} \\ = \min_{\mathbf{q}} \max_{\mathbf{p}} \mathbf{p}^T \mathbf{A} \mathbf{q} \end{aligned}$$

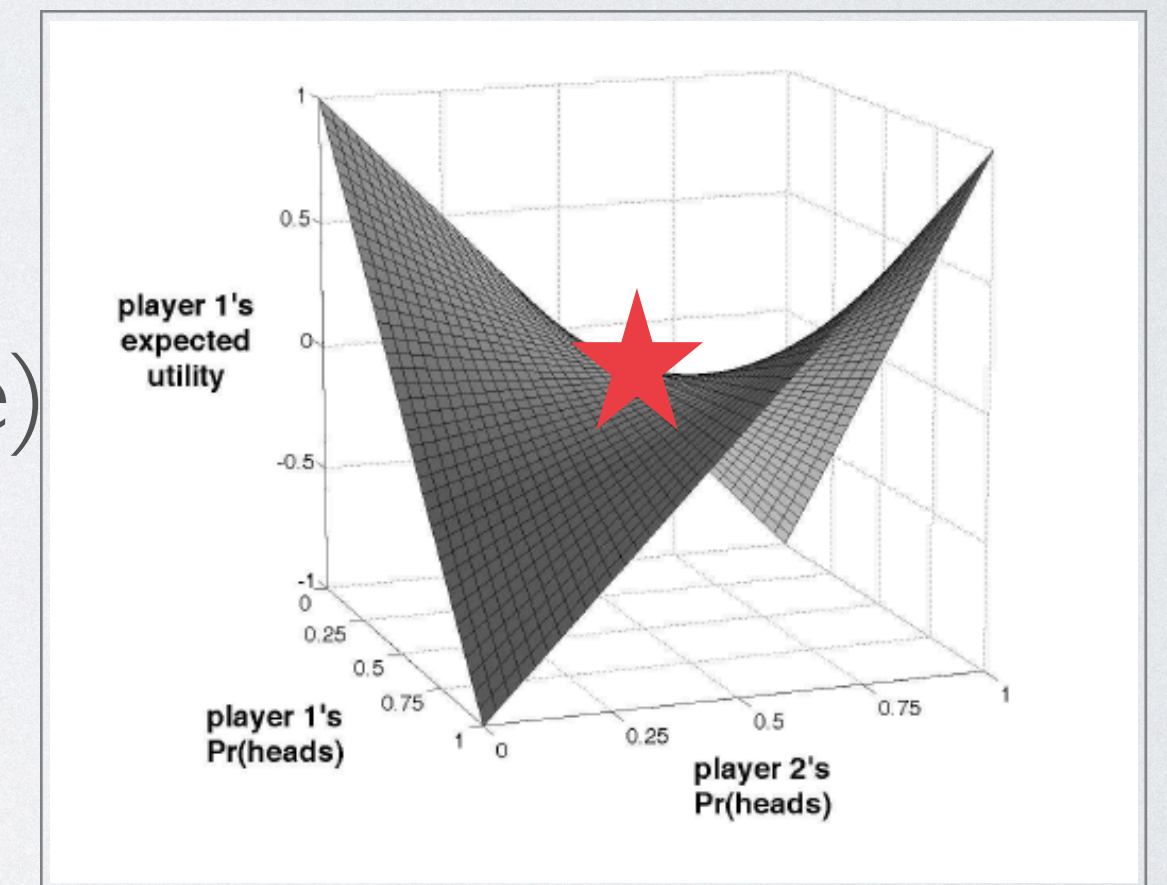
- The  $v^*$  is the Nash value

- proof:  $v \leq v^*$  due to definition of  $v^*$ ,  $v \geq v^*$  due to being the LP solution.
- corollary: all Nash value are the same (saddle point is unique in bimatrix game)

- The  $(\mathbf{p}, \mathbf{q})$  is the Nash equilibrium:

- proof: suppose the player plays  $\mathbf{x}, \mathbf{y}$  instead of  $\mathbf{p}, \mathbf{q}$

- $x^T A q = \sum_{i=1}^N x_i (A q)_i \leq \max_{i \in [N]} (A q)_i = v_c = v^*$ ,  $p^T A y = \sum_{j=1}^M (p^T A)_j y_j \geq \min_{j \in [M]} (p^T A)_j = v_r = v^*$ , thus no incentives to deviate.



- **Sion's minimax theorem** generalises to quasi-convex/concave functions  $\min_{x \in X} \sup_{y \in Y} f(x, y) = \sup_{y \in Y} \min_{x \in X} f(x, y)$

# Fictitious Play [Brown 1951]

- Maintain a belief over the historical actions that the opponent has played, and the learning agent then takes the best response to this empirical distribution.

$$a_i^{t,*} \in \mathbf{BR}_i \left( p_{-i}^t = \frac{1}{t} \sum_{\tau=0}^{t-1} \mathcal{F} \{ a_{-i}^\tau = a, a \in \mathbb{A} \} \right)$$

$$p_i^{t+1} = \left( 1 - \frac{1}{t} \right) p_i^t + \frac{1}{t} a_i^{t,*}, \text{ for all } i$$

- It guarantees to converge, in terms of the Nash value, in two-player zero-sum games, potential games and  $2 \times 2$  games, and, the average policy converge to the Nash strategy.

- Examples:

		Player 2	
		<i>a</i>	<i>b</i>
Player 1	<i>A</i>	(1,1)	(0,0)
	<i>B</i>	(0,0)	(1,1)

<i>t</i>	$p_1^t$	$p_2^t$	$a_1^t$	$a_2^t$
0	(3/4, 1/4)	(1/4, 3/4)	B	a
1	(3/4, 5/4)	(5/4, 3/4)	A	b
2	(7/4, 5/4)	(5/4, 7/4)	B	a
3	(7/4, 9/4)	(9/4, 7/4)	A	b
⋮	⋮	⋮	⋮	⋮

∞ (1/2, 1/2) (1/2, 1/2)

# Double Oracle [McMahan 2003]

- Double Oracle best responds to the opponent's Nash equilibrium at each iteration.
- To solve the game by solving subgame Nash because not all pure strategies are “useful” in the support of Nash.
- It is much faster than LP, but In the worst-case scenario, it recovers to solve the original game.

## Algorithm 1 Double Oracle (McMahan et al., 2003)

```

1: Input: A set  $\Pi, C$  strategy set of players
2:  $\Pi_0, C_0$ : initial set of strategies
3: for  $t = 1$  to  $\infty$  do
4:   if  $\Pi_t \neq \Pi_{t-1}$  or  $C_t \neq C_{t-1}$  then
5:     Solve the NE of the subgame  $G_t$ :
6:      $(\pi_t^*, c_t^*) = \arg \min_{\pi \in \Delta_{\Pi_t}} \arg \max_{c \in \Delta_{C_t}} \pi^\top A c$ 
7:     Find the best response  $a_{t+1}$  and  $c_{t+1}$  to  $(\pi_t^*, c_t^*)$ :
8:      $a_{t+1} = \arg \min_{a \in \Pi} a^\top A c_t^*$ 
9:      $c_{t+1} = \arg \max_{c \in C} \pi_t^{*\top} A c$ 
10:    Update  $\Pi_{t+1} = \Pi_t \cup \{a_{t+1}\}, C_{t+1} = C_t \cup \{c_{t+1}\}$ 
11:   else if  $\Pi_t = \Pi_{t-1}$  and  $C_t = C_{t-1}$  then
12:     Terminate
13:   end if
14: end for

```

■ **iteration 0:** restricted game R vs R

■ **iteration 1:**

- solve Nash of restricted game  $(1, 0, 0), (1, 0, 0)$
- unrestricted  $\mathbf{Br}^1, \mathbf{Br}^2 = P, P$

■ **iteration 2:**

- solve Nash of restricted games  $(0, 1, 0), (0, 1, 0)$
- unrestricted  $\mathbf{Br}^1, \mathbf{Br}^2 = S, S$

■ **iteration 3:**

- solve Nash of restricted game  $(1/3, 1/3, 1/3), (1/3, 1/3, 1/3)$

■ **iteration 4:** no new response, END

- output  $(1/3, 1/3, 1/3)$

example

	R	P	S
R	0	-1	1
P	1	0	-1
S	-1	1	0



# Double Oracle [McMahan 2003]

- It guarantees to converge to Nash equilibrium in two-player zero-sum games.
- We need to prove that early stopping also leads to Nash equilibrium.

- **Convergence proof:**

- ◆ DO finally recovers to solve the whole game

- **Correctness proof:**

- ◆ suppose DO stops at the  $j$ -th sub-game (i.e., no new best responses are added)

- ◆  $\forall p, V(p, q_j) \geq v \Rightarrow \forall p, \max_q V(p, q) \geq v$

$$\forall q, V(p_j, q) \leq v \Rightarrow \max_q V(p_j, q) \leq v$$

$$\Rightarrow \forall p, \max_q V(p_j, q) \leq \max_q V(p, q)$$

$p_j$  must be the minimax optimal,  
 $q_j$  vice versa, so  $(p_j, q_j)$  is the final Nash

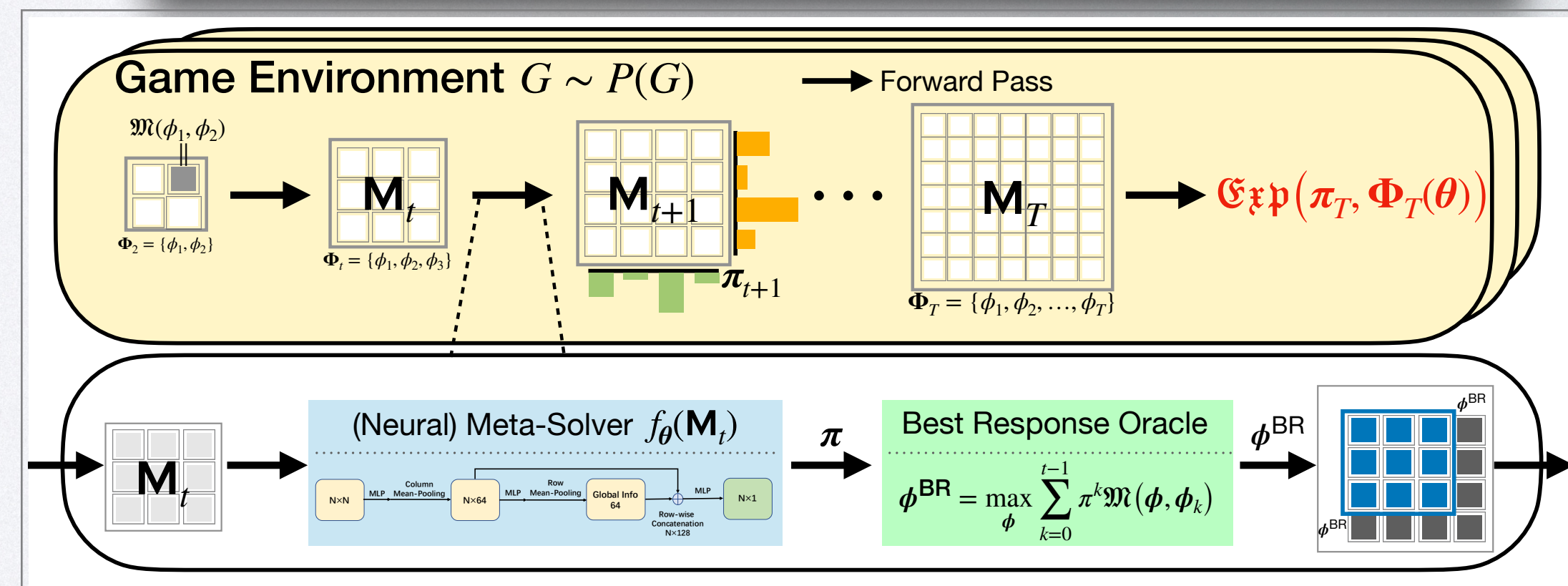
# Policy Space Response Oracle [Lanctot 2017]

- A generalisation of double oracle methods: the best responder is implemented through **deep RL models**.
- Each RL model is now a “pure strategy”.
- A meta-game is  $(\Pi, U, n)$  where  $\Pi = (\Pi_1, \dots, \Pi_n)$  is the set of policies for each agent and  $U : \Pi \rightarrow \mathbb{R}^n$  is the reward values for each agent given a joint strategy profile.
- $\sigma_{-i}$  is distribution over  $(\Pi_1^0, \dots, \Pi_1^T)$ , a.k.a meta-strategy
- PSRO generalises all previous methods by varying  $\sigma_{-i}$ .

- independent learning:  $\sigma_{-i} = (0, \dots, 0, 0, 1)$
- self-play:  $\sigma_{-i} = (0, \dots, 0, 1, 0)$
- fictitious play:  $\sigma_{-i} = (1/T, 1/T, \dots, 1/T, 0)$
- PSRO:  $\sigma_{-i} = \text{Nash}(\Pi^{T-1}, U)$

## Algorithm 1: Policy-Space Response Oracles

**input** : initial policy sets for all players  $\Pi$   
 Compute exp. utilities  $U^\Pi$  for each joint  $\pi \in \Pi$   
 Initialize meta-strategies  $\sigma_i = \text{UNIFORM}(\Pi_i)$   
**while** epoch  $e$  in  $\{1, 2, \dots\}$  **do**  
     **for** player  $i \in [[n]]$  **do**  
         **for** many episodes **do**  
             select opponent policies      Sample  $\pi_{-i} \sim \sigma_{-i}$   
             compute the best response    Train oracle  $\pi'_i$  over  $\rho \sim (\pi'_i, \pi_{-i})$   
             augment strategy pool       $\Pi_i = \Pi_i \cup \{\pi'_i\}$   
             expand the payoff matrix    Compute missing entries in  $U^\Pi$  from  $\Pi$   
             solve the new meta game      Compute a meta-strategy  $\sigma$  from  $U^\Pi$   
         Output current solution strategy  $\sigma_i$  for player  $i$



# PSRO on Google Football !

<https://sites.google.com/view/diverse-psro/>



make offside



push and run

# Two-Player General-Sum Games

- Exponential-time algorithms are the best we can hope for.
- **Support Enumeration Method:**

◆ **Theorem: a best response has to be a pure strategy that satisfies  $x_i > 0 \implies (Ay)_i = \max_{k \in [N]} (Ay)_k$ , that is,  $x$  is a best response to  $y$  if only if  $x$  has nonzero probability only on pure strategies that max expected utility against  $y$ .**

**Proof:** if  $x_i > 0$ , then  $x_i$  must be best response, otherwise suppose  $(Ay)_i < \max_{k \in [N]} (Ay)_k$  it will contradict

$$\begin{aligned}x^T Ay &= \sum_{k=1}^N x_k (Ay)_k = x_i (Ay)_i + \sum_{k \neq i} x_k (Ay)_k \\ &< x_i \max_{k \in [N]} (Ay)_k + \sum_{j \neq i} x_j \max_{k \in [N]} (Ay)_k \\ &= x_i \max_{k \in [N]} (Ay)_k + (1 - x_i) \max_{k \in [N]} (Ay)_k = \max_{k \in [N]} (Ay)_k\end{aligned}$$

if  $x_i$  is a best response, then  $x_i > 0$

$$x^T Ay = \sum_{i=1}^N x_i (Ay)_i = \sum_{i: x_i \neq 0} x_i (Ay)_i = \sum_{i: x_i \neq 0} x_i \max_{k \in [N]} (Ay)_k = \max_{k \in [N]} (Ay)_k \sum_{i: x_i \neq 0} x_i = \max_{k \in [N]} (Ay)_k$$

# Two-Player General-Sum Games

- Support Enumeration Method:

- ◆ **Theorem: a best response has to be a pure strategy that satisfies  $x_i > 0 \implies (Ay)_i = \max_{k \in [N]} (Ay)_k$ , that is,  $x$  is a best response to  $y$  if only if  $x$  has nonzero probability only on pure strategies that max expected utility.**
- ◆ we can have a **guess of the Nash support  $(I, J)$**  and then propose the Nash by solving the following LP

$$\begin{aligned} \sum_{i \in I} x_i &= 1 \\ \sum_{j \in J} y_j &= 1 \\ \sum_{i \in I} x_i B_{ij} &= v \quad \forall j \in J \\ \sum_{j \in J} A_{ij} y_j &= u \quad \forall i \in I \end{aligned}$$

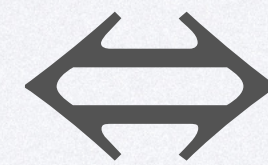
size of  $|I| + |J| + 2$ ,  
solved by Gaussian elimination

- ◆ Given the proposal  $(x, y)$  we can then **check if  $u = \max_{k \in [N]} (Ay)_k$  and  $v = \max_{k \in [M]} (x^T B)_k$  are met.**

# LCP formulation of Two-Player General-Sum Games

- Support Enumeration Method is more like a heuristic search.
- Here we provide a formal formulation through **linear complementarity problem (LCP)**
- LCP is very much like an LP, but with a new constraint built on a **slack variable**.
- Recall that in two-player zero-sum game, we can have

$$\begin{array}{ll} \text{minimize} & U_1^* \\ \text{subject to} & \sum_{k \in A_2} u_1(a_1^j, a_2^k) \cdot s_2^k \leq U_1^* \quad \forall j \in A_1 \\ & \sum_{k \in A_2} s_2^k = 1 \\ & s_2^k \geq 0 \quad \forall k \in A_2 \end{array}$$



$$\begin{array}{ll} \text{minimize} & U_1^* \\ \text{subject to} & \sum_{k \in A_2} u_1(a_1^j, a_2^k) \cdot s_2^k + r_1^j = U_1^* \quad \forall j \in A_1 \\ & \sum_{k \in A_2} s_2^k = 1 \\ & s_2^k \geq 0 \quad \forall k \in A_2 \\ & r_1^j \geq 0 \quad \forall j \in A_1 \end{array}$$

# LCP formulation of Two-Player General-Sum Games

- Solving the Nash of two-player general-sum games is an LCP problem.
  - ◆ Both two players' variables need to be considered rather than one player (unlike two-player zero sum!)
  - ◆ There is no objective, it is rather a *feasibility program* (finding the solution that meets the conditions).
  - ◆ The last *complementarity condition* that prevents unbounded  $U^*$ ,  $V^*$  is non-linear, turn LP into LCP.
  - ◆ If an action is played  $x^j > 0$ , it has to be a best response:  $r_1^j = 0$ , otherwise it can just deviate to reach  $V^*$

$$\sum_{k \in A_2} A(a_1^j, a_2^k) \cdot y^k + r_1^j = V^* \quad \forall j \in A_1$$

$$\sum_{j \in A_1} B(a_1^j, a_2^k) \cdot x^j + r_2^k = U^* \quad \forall k \in A_2$$

mutual best response

$$\sum_{j \in A_1} x^j = 1, \sum_{k \in A_2} y^k = 1, \quad x^j \geq 0, y^k \geq 0 \quad \forall j \in A_1, \forall k \in A_2$$

valid prob. distribution

$$r_1^j \geq 0, \quad r_2^k \geq 0 \quad \forall j \in A_1, \forall k \in A_2$$

slack variables

$$r_1^j \cdot x^j = 0, \quad r_2^k \cdot y^k = 0 \quad \forall j \in A_1, \forall k \in A_2$$

complementarity condition

# Lemke-Howson Method

## Solving the LCP problem through Lemke-Howson:

- A classical algorithm that combines game theory, convex analysis and graph theory
- Non-degenerate games: for a pure strategy, there can only be at most one best response.
- Non-degenerate games have the same Nash support size for both players:  $|p^*| = |q^*|$ .
- Consider polyhedron:  $P = \left\{ (u, \mathbf{x}) \mid x_i \geq 0, \sum x_i = 1, \mathbf{x}^T B \leq u \cdot \mathbf{1} \right\}$  and  $Q = \left\{ (v, \mathbf{y}) \mid y_j \geq 0, \sum y_j = 1, A\mathbf{y} \leq v \cdot \mathbf{1} \right\}$ 
  - ♦ These polyhedra describes the space of mixed strategies with an upper bound on the best response value from the other player will react.
- Consider the (bounded) polytope:  $\bar{P} = \left\{ \mathbf{x} \mid x_i \geq 0, \mathbf{x}^T B \leq \mathbf{1} \right\}$  and  $\bar{Q} = \left\{ \mathbf{y} \mid y_j \geq 0, A\mathbf{y} \leq \mathbf{1} \right\}$ .
  - Setting  $U^* = V^* = 1$  is generic, there exist bijective mappings between  $P, Q$  and  $\bar{P}, \bar{Q}$
- We are interested in the graph that are composed by the corner points of  $\bar{P}, \bar{Q}$ , because it either means  $x_i$  is either in Nash support, or column  $j$  is a best response to  $\mathbf{x}$ .

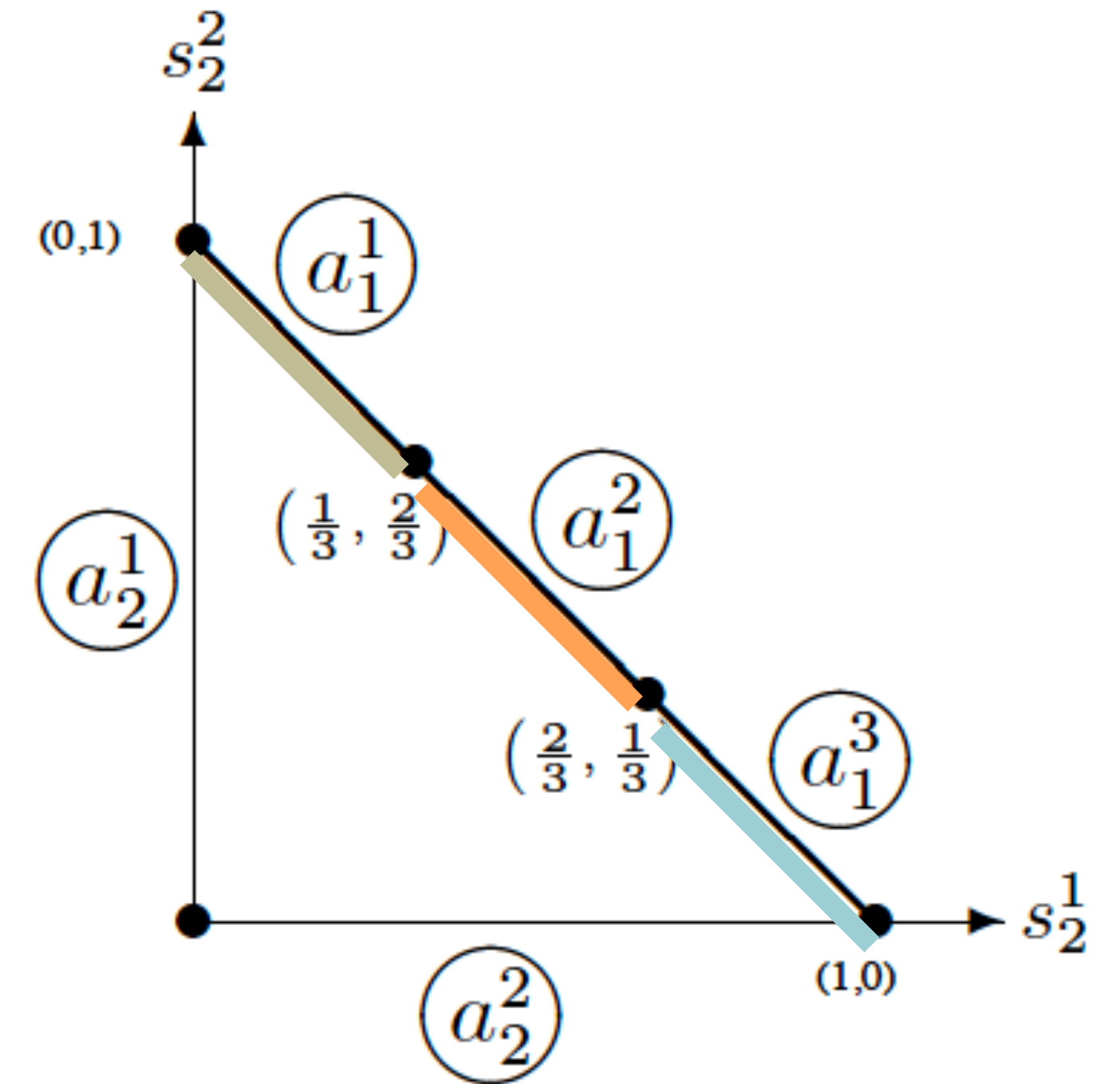
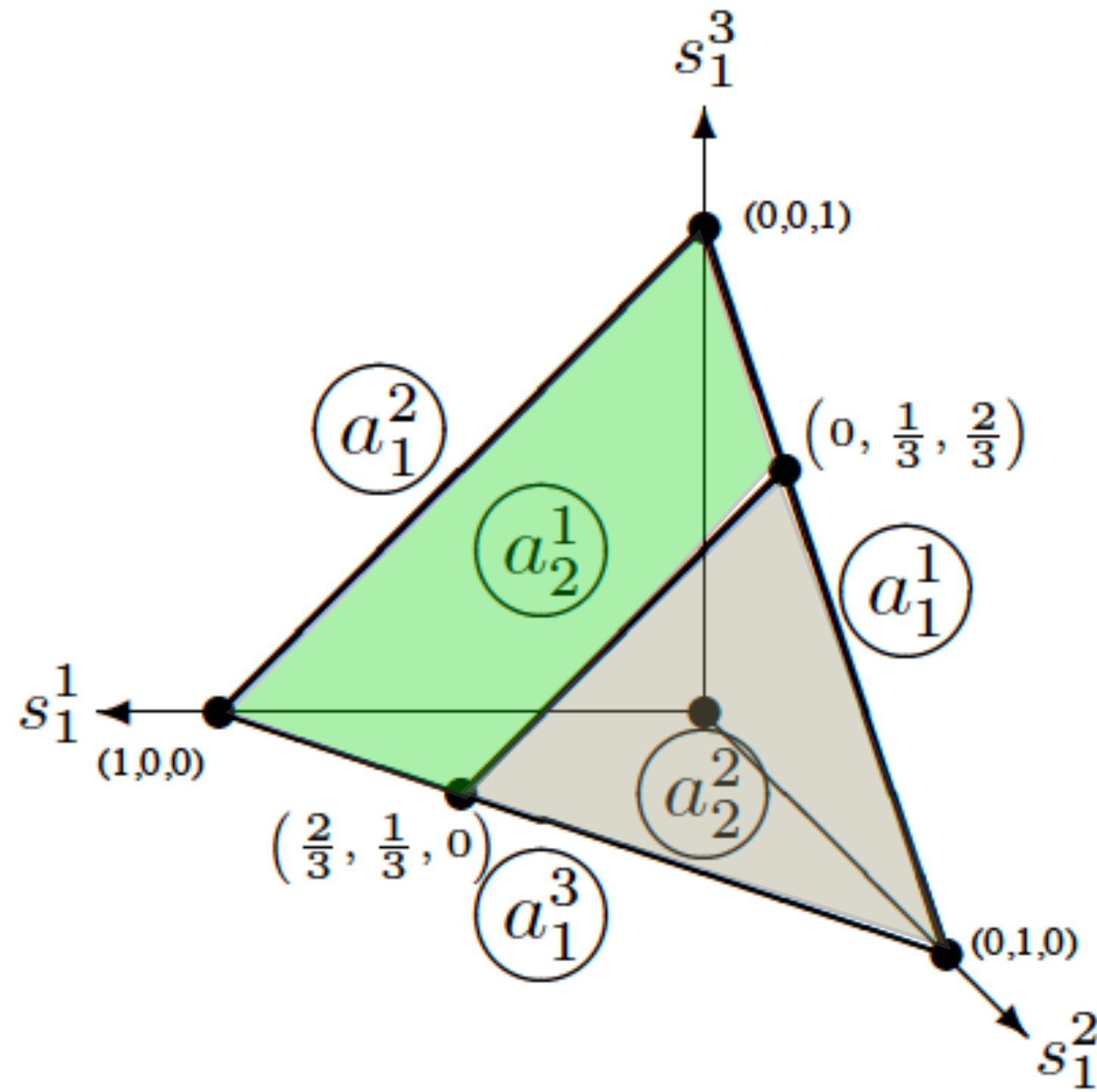
$$L(\mathbf{x}) = \left\{ i \mid x_i = 0 \right\} \cup \left\{ j \mid (\mathbf{x}^T B)_j = 1 \right\} \quad L(\mathbf{y}) = \left\{ j \mid y_j = 0 \right\} \cup \left\{ i \mid (A\mathbf{y})_i = 1 \right\}$$



# Lemke-Howson Method Example

$$L(x) = \{i \mid x_i = 0\} \cup \{j \mid (x^T B)_j = 1\} \quad L(y) = \{j \mid y_j = 0\} \cup \{i \mid (Ay)_i = 1\}$$

0, 1	6, 0
2, 0	5, 2
3, 4	3, 3



# Lemke-Howson Method

Solving the LCP problem through Lemke-Howson:

- We are interested in the graph that are composed by the corner points of  $\bar{P}, \bar{Q}$ , because it either means  $x_i$  is not in Nash, or opponent's column  $j$  is a best response to  $\mathbf{x}$ . Let's label those corner points by the constraint id:

$$L(\mathbf{x}) = \{i \mid x_i = 0\} \cup \{j \mid (\mathbf{x}^T B)_j = 1\} \qquad L(\mathbf{y}) = \{j \mid y_j = 0\} \cup \{i \mid (A\mathbf{y})_i = 1\}$$

**Theorem:** a pair  $(\mathbf{x}, \mathbf{y})$  is a Nash equilibrium if and only if there are completely labelled:  $L(\mathbf{x}) \cup L(\mathbf{y}) = M + N$

- **Proof:**  $\Leftarrow$ : if  $(\mathbf{x}, \mathbf{y})$  is a Nash, then  $x_i$  is either 0, thus  $i \in L(\mathbf{x})$ , or, is a best response to  $(A\mathbf{y})_i = 1$ , and thus  $i \in L(\mathbf{y})$ . Therefore, every label appears only once in  $L(\mathbf{x}) \cup L(\mathbf{y})$ , so  $L = L(\mathbf{x}) \cup L(\mathbf{y})$ .

$\Rightarrow$ : if  $L = L(\mathbf{x}) \cup L(\mathbf{y})$ , by definition, we can know that the following set is a Nash.

$$\left( \left\{ i \in [N] : \sum_{j=1}^M y_j A_{ij} = 1 \right\} \quad (i \in L(\mathbf{y})), \quad \left\{ j \in [M] : \sum_{i=1}^N x_i B_{ij} = 1 \right\} \quad (j \in L(\mathbf{x})) \right)$$

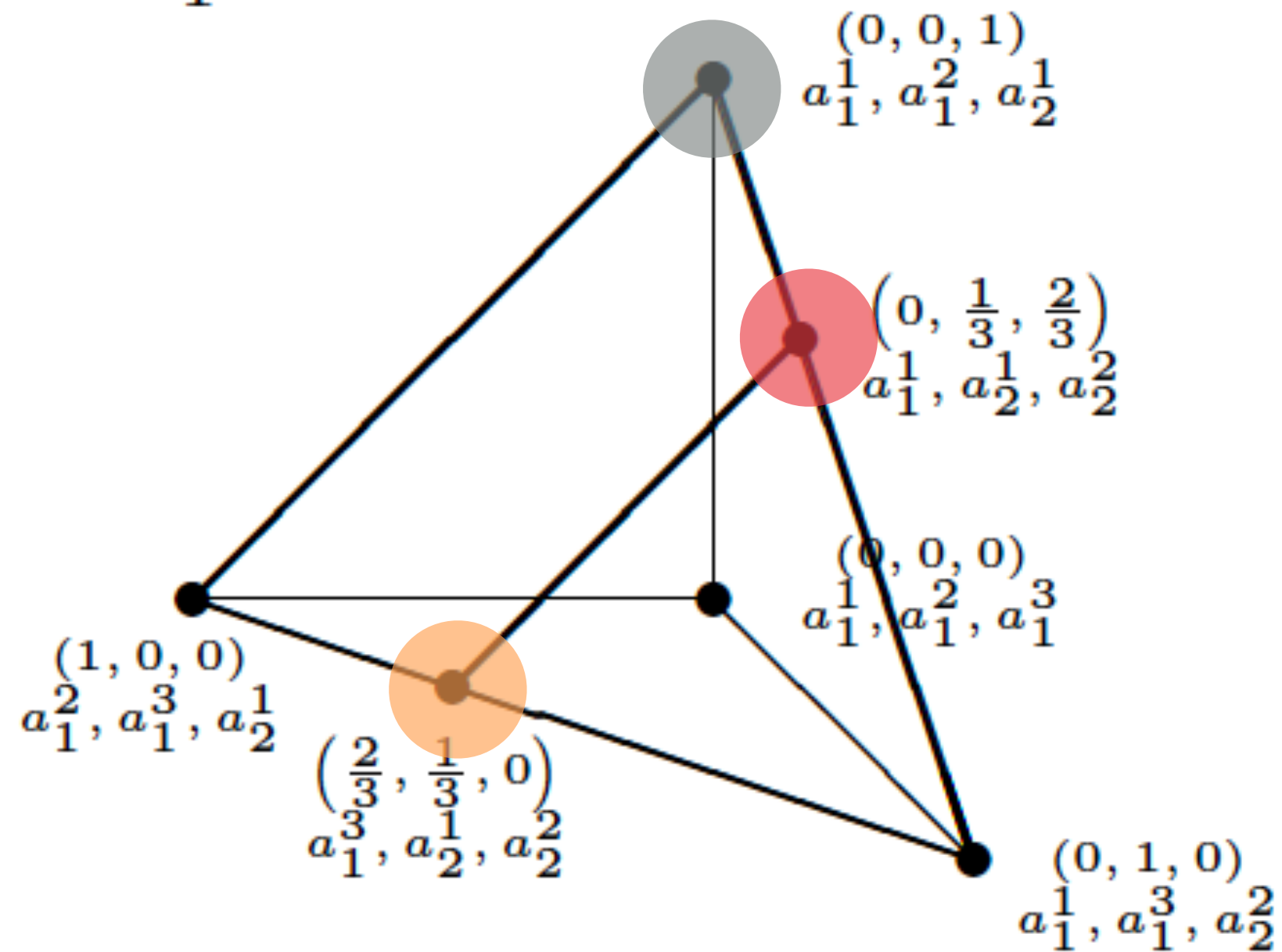
- Intuition: finding Nash is just find vertices of  $\bar{P}, \bar{Q}$  that contains all of the labels.

# Lemke-Howson Method

**Theorem:** a pair  $(\mathbf{x}, \mathbf{y})$  is a Nash equilibrium if and only if there are completely labelled:  $L(\mathbf{x}) \cup L(\mathbf{y}) = M + N$

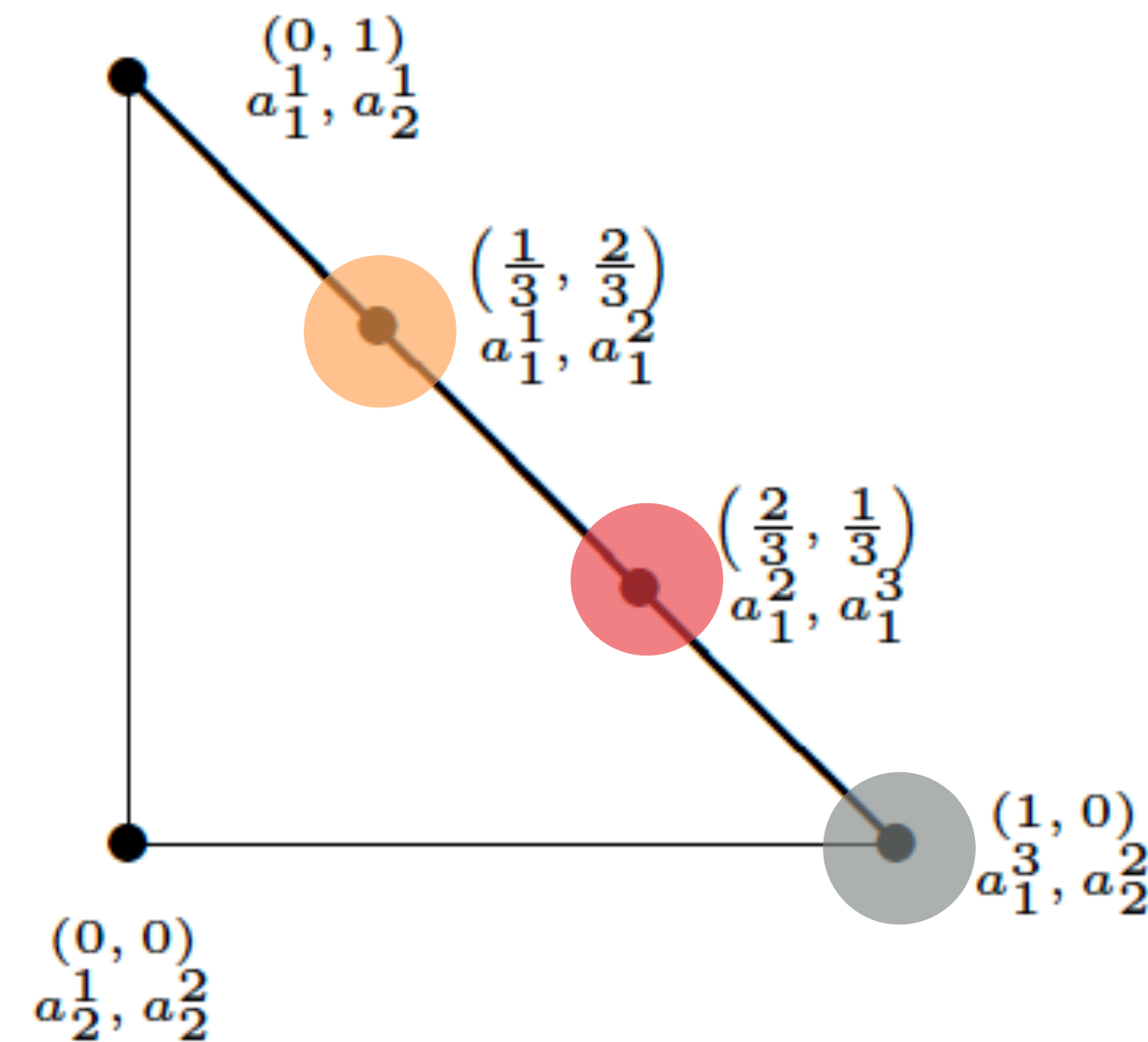
$$L(\mathbf{x}) = \{i \mid x_i = 0\} \cup \{j \mid (\mathbf{x}^T B)_j = 1\}$$

$G_1$ :



$$L(\mathbf{y}) = \{j \mid y_j = 0\} \cup \{i \mid (A\mathbf{y})_i = 1\}$$

$G_2$ :



# Lemke-Howson Method

- Lemke-Howson Method:

- Intuition: finding Nash is just find vertices of  $\bar{P}, \bar{Q}$  that contains all of the labels.

- In order to find the completely labelled pairs (the Nash ), we define:

- $G = G_1 \times G_2$ , with vertices  $v = (v_1, v_2)$  where  $v_1 \in V(G_1)$  and  $v_2 \in V(G_2)$ .

- Edge of  $\mathbf{E}(G) = \{(v_1, v_2), (v'_1, v_2) \in G \times G : \text{if } (v_1, v'_1) \in G_1\} \cup \{(v_1, v_2), (v_1, v'_2) \in G \times G : \text{if } (v_2, v'_2) \in G_2\}$

- $L(v) = L(v_1) \cup L(v_2)$  for each vertex  $v = (v_1, v_2) \in V(G)$

- Let's focus on the subgraph in  $G$  that is almost fully labelled but only lack the label  $k$

- For  $k \in M \cup N$ , we write  $U_k = \{v \in V(G) \mid L(v) \supseteq M \cup N \setminus \{k\}\}$  for the subset of vertices except  $k$ .

- **Intuition: Nash equilibrium in  $U_k$  has degree only equal to 1.**

- Nash is exactly  $M \cup N$ , so are in  $U_k \forall k$ .

- Since Nash is a fully labelled point and  $L(v_1) \cap L(v_2) = \emptyset$ , thus Nash in  $U_k$  can only have degree 1 (either best response or not in the Nash support, cannot have both !).

- In  $U_k$ , the degree of every other vertex apart from Nash have degree 2 (both in  $G_1$  and  $G_2$ ).

- **An effective algorithm: walk along the edges in  $G$ , and drop the nodes with degree of 2.**

# Lemke-Howson Method

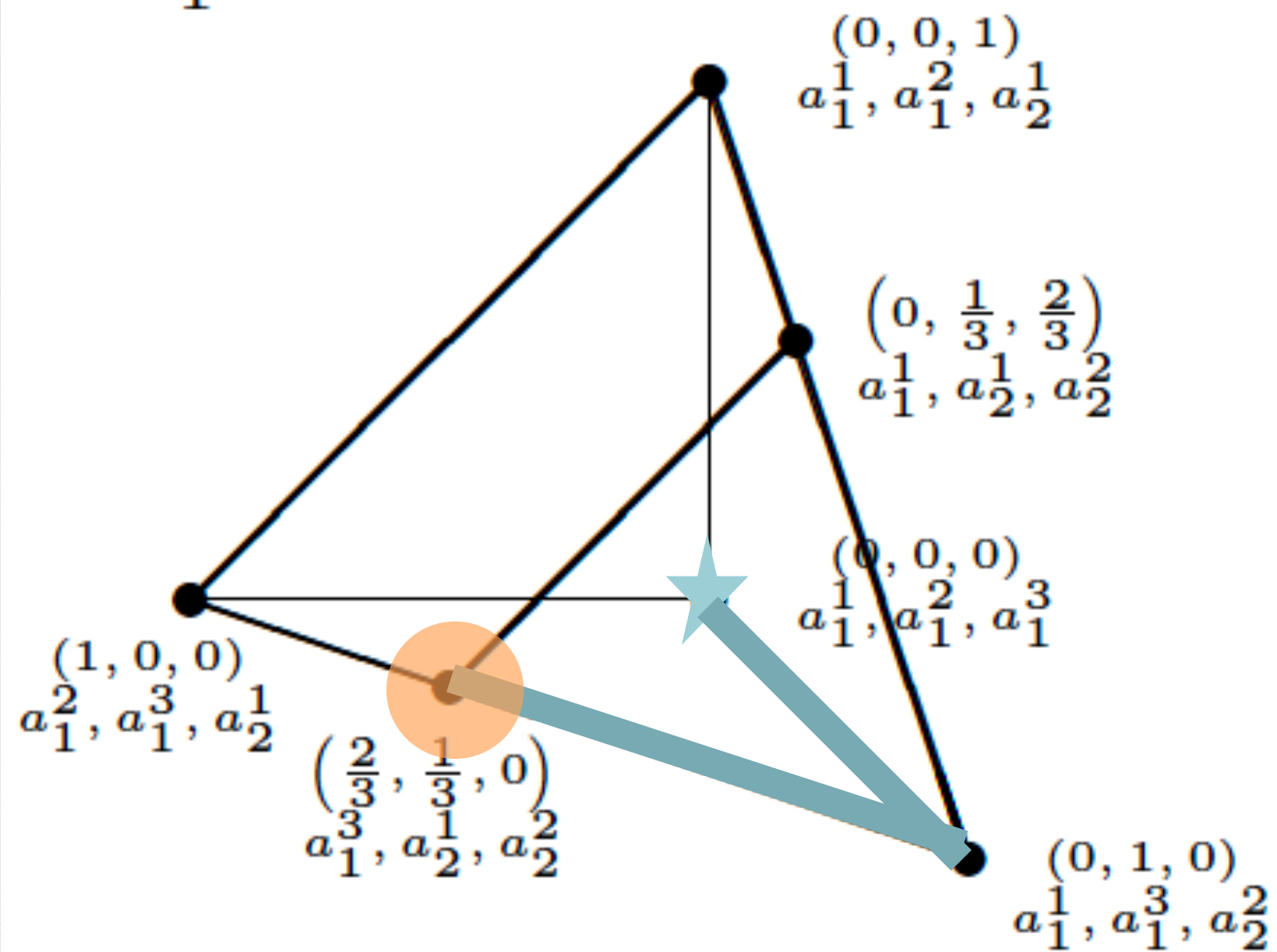
**Theorem:** a pair  $(\mathbf{x}, \mathbf{y})$  is a Nash equilibrium if and only if there are completely labelled:  $L(\mathbf{x}) \cup L(\mathbf{y}) = M + N$

- An effective algorithm: walk along the edges in  $G$ , and drop the nodes with degree of 2.

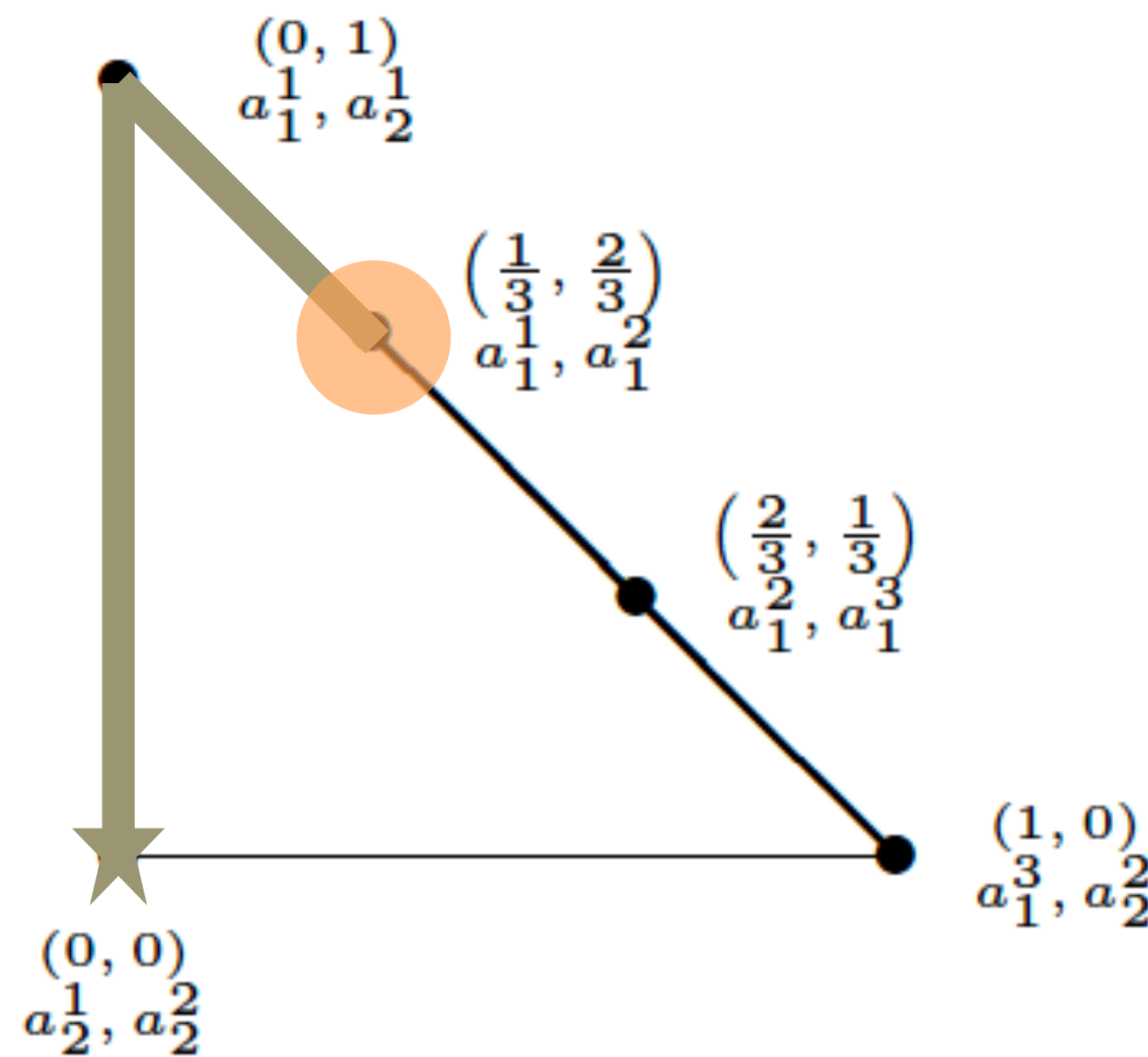
$$L(\mathbf{x}) = \{i \mid x_i = 0\} \cup \{j \mid (\mathbf{x}^T B)_j = 1\}$$

$$L(\mathbf{y}) = \{j \mid y_j = 0\} \cup \{i \mid (A\mathbf{y})_i = 1\}$$

$G_1$ :



$G_2$ :



Start:  $(0,0,0)$   $(0,0)$

$(0,0,0) \rightarrow (0,1,0)$ , drop  $a_1^2$ , pick  $a_2^2$

$(0,0) \rightarrow (0,1)$ , drop  $a_2^2$ , pick  $a_1^1$

$(0,1,0) \rightarrow (2/3, 1/3, 0)$ , drop  $a_1^1$ , pick  $a_2^1$

$(0,1) \rightarrow (1/3, 2/3)$ , drop  $a_2^1$ , pick  $a_1^1$

Terminate: because now we have  $\{a_1^3, a_2^1, a_2^2\} \cup \{a_1^1, a_1^2\}$  fully labelled

Output:  $(2/3, 1/3, 0)$ ,  $(1/3, 2/3)$

# Lemke-Howson Methods

## Algorithm (Lemke-Howson)

**Input:** A Non-degenerate bimatrix game  $(A, B)$ .

**Output:** One Nash equilibrium of the game.

1. Choose  $k \in M \cup N$ .
2. Start with  $(x, y) = (0, 0) \in G$ . Drop label  $k$  from  $(x, y)$  (from  $x \in \bar{P}$  if  $k \in M$ , from  $y \in \bar{Q}$  if  $k \in N$ ).
3. Let  $(x, y)$  be the current vertex. Let  $l$  be the label that is picked up by dropping label  $k$ . If  $l = k$ , terminate and  $(x, y)$  is a Nash equilibrium of the game. If  $l \neq k$ , drop  $l$  in the other polytope and repeat this step.

## Remarks:

- Traversing the whole graph  $G$  to add missing labels and drop duplicate labels (i.e. pivoting) until we find one set of nodes that are fully labelled.
- Guarantee to find one Nash, but not all of them! Determining whether all Nash are found is not even NP.
- Exponential time complexity because  $G$  has the number of vertices that is exponential in  $n, m$  (combinatorial many)
- Because LCP has no objective, we cannot tell the progress before finding a solution.
- Corollary: almost all games have an odd number of Nash equilibrium.

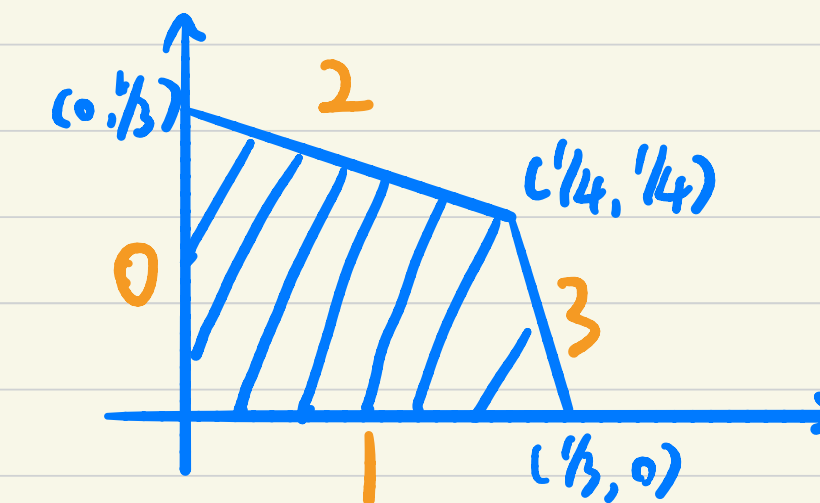
$$A = \begin{pmatrix} 3 & 1 \\ 1 & 3 \end{pmatrix} \quad B = \begin{pmatrix} 1 & 3 \\ 3 & 1 \end{pmatrix}$$

1) Construct the polytope

$$\bar{P} = \{x \in \mathbb{R}^2 \mid x_i \geq 0, x^T B \leq 1\}$$

$$\begin{cases} x_1 \geq 0 \\ x_2 \geq 0 \\ x_1 + 3x_2 \leq 1 \\ 3x_1 + x_2 \leq 1 \end{cases} \Rightarrow V = \{(0, 0), (0, 1/3), (1/3, 0), (1/4, 1/4)\}$$

2) Label the vertex

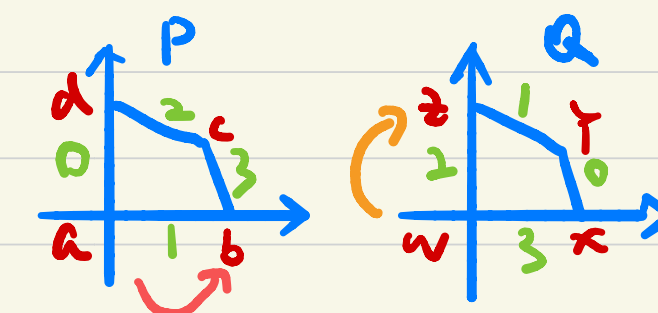


$$\begin{aligned} 0: & x_1 = 0 \\ 1: & x_2 = 0 \\ 2: & x_1 + 3x_2 = 1 \quad \text{when } y_1 \text{ is the best response to } x \\ 3: & 3x_1 + x_2 = 1 \quad \text{when } y_2 \text{ is the best response to } x \end{aligned}$$

$$\begin{aligned} \bar{P}: (0, 0) &: \{0, 1\} & (0, 1/3) &: \{0, 2\} \\ (1/3, 0) &: \{1, 3\} & (1/4, 1/4) &: \{2, 3\} \end{aligned}$$

$$\begin{aligned} \bar{Q}: (0, 0) &: \{2, 3\} & (1/4, 1/4) &: \{0, 1\} \\ (1/3, 0) &: \{0, 3\} & (0, 1/3) &: \{1, 2\} \end{aligned}$$

3) Lemke-Howson step.



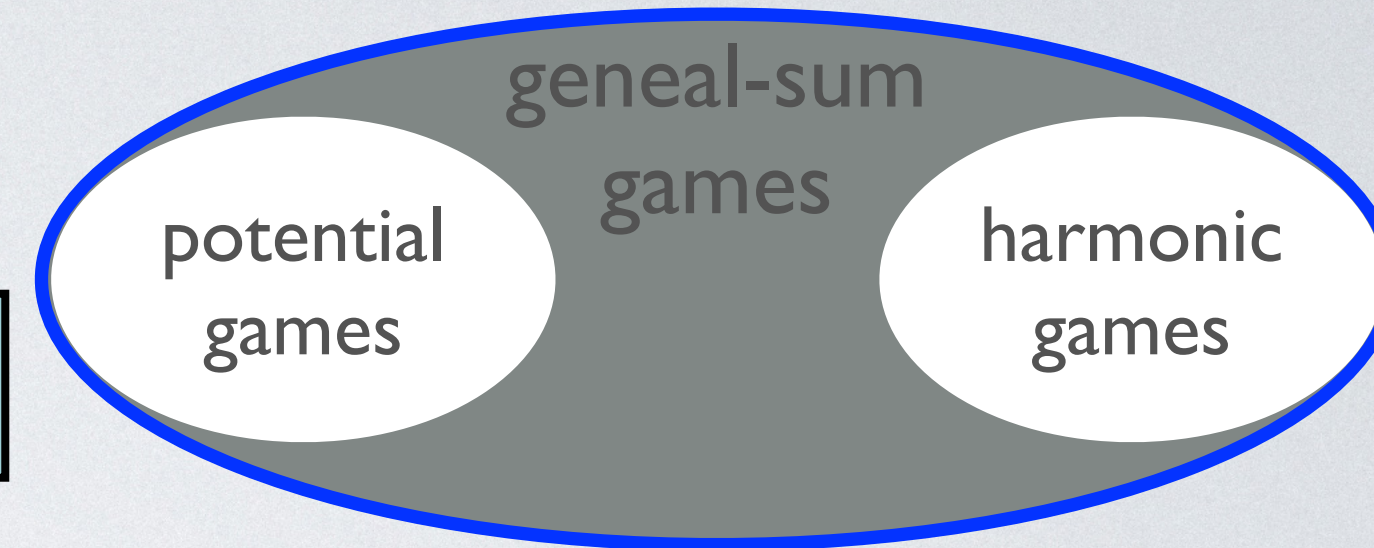
Start:  $(a, w) = \{0, 1, 2, 3\}$ . drop 0 from P, pick 3.  
 $(b, w) = \{1, 2, 3\}$ . drop 3 from Q, pick 1.  
 $(b, z) = \{1, 2, 3\}$  drop 1 from P, pick 2.  
 $(c, z) = \{1, 2, 3\}$  drop 2 from Q, pick 0.  
 $(c, y) = \{0, 1, 2, 3\}$ . Terminate  
 output: Normalized  $(c, y) = (1/2, 1/2)$

# Potential Games

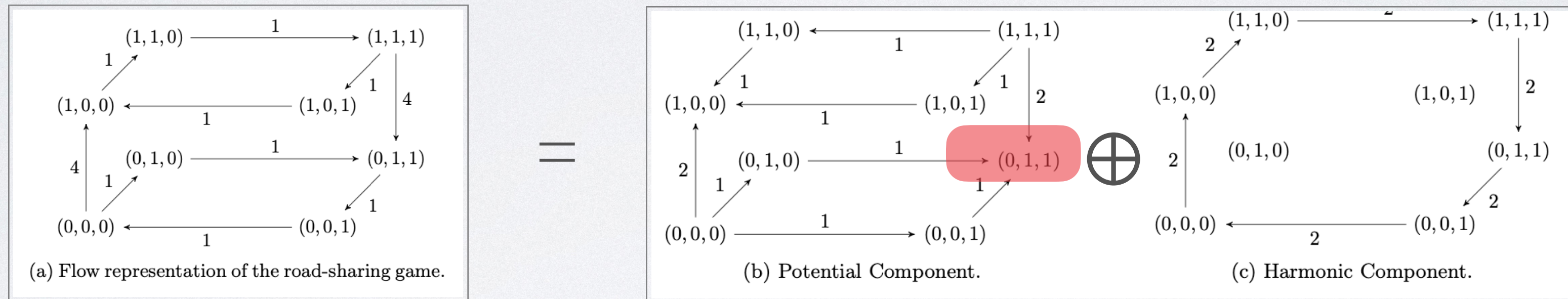
- Potential Game is the other side of the coin where we know how to solve Nash.
- A general game decomposition results suggest that

[Candogan 2010]

$$\text{Normal-form Game} = \text{Potential Game} \oplus \text{Harmonic Game}$$



- Potential games has no **curl**, harmonic games has no **divergence**.



- In bimatrix games, potential game is related to the identical interest game, harmonic game is related to the zero-sum games.

• Since one can always write  $(\mathbf{A}, \mathbf{B}) = \left( \frac{\mathbf{A} - \mathbf{B}}{2}, \frac{\mathbf{B} - \mathbf{A}}{2} \right) \oplus \left( \frac{\mathbf{A} + \mathbf{B}}{2}, \frac{\mathbf{B} + \mathbf{A}}{2} \right)$

zero-sum
team-game

# Potential Games

- Potential Games are those games that can be described a potential function:

$$\Phi(s_{-i}, s'_i) - \Phi(s) = C_i(s_{-i}, s'_i) - C_i(s), \forall i \in [n], s \in S, s'_i \in S_i$$

Φ

(0, 0)	(1, 2)
(2, 1)	(0, 0)

↔

0	2
2	1

- The outcome of any player's any possible deviation are captured by the change in Φ.
- Knowing the location deviations from a global viewpoint.

- **Examples of potential games:**

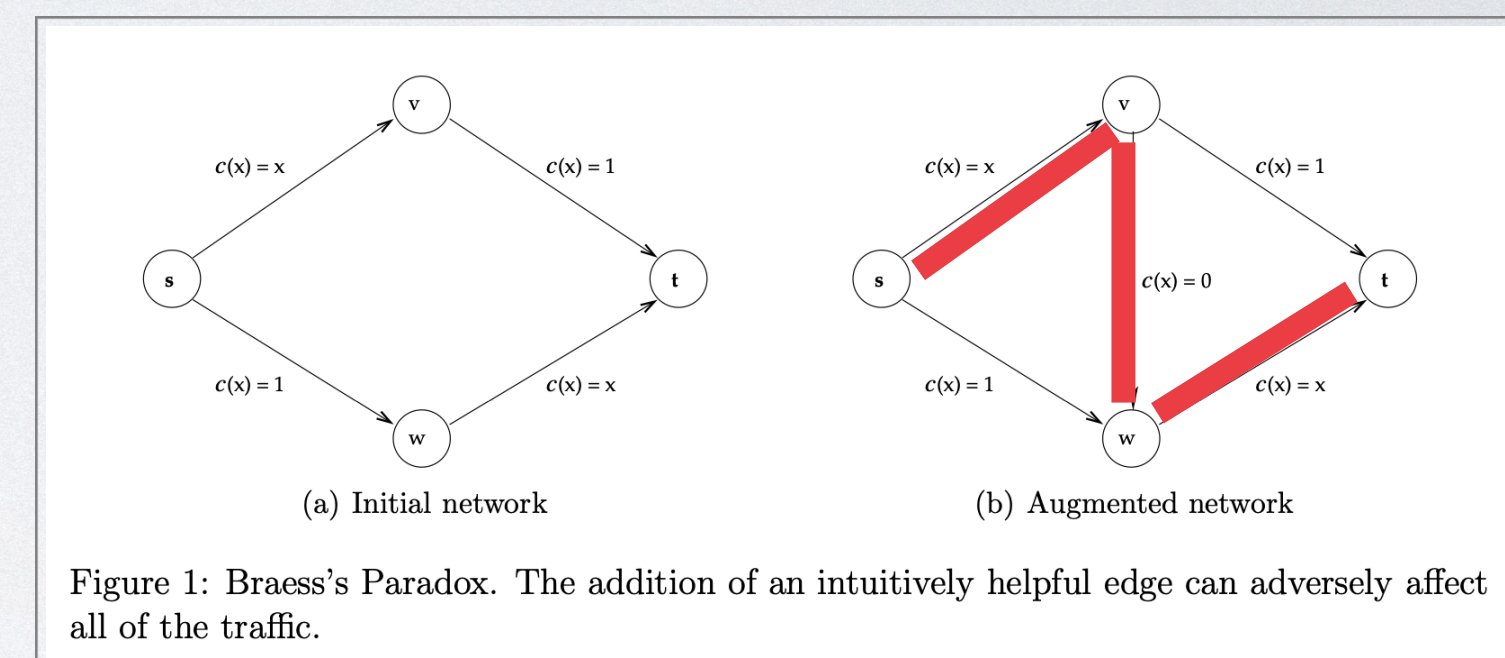
- ◆ All **fully cooperative games** are the potential game. Potential is the reward function.
- ◆ The **routing game** we have seen is a potential game.

- Each player choose a path  $P_i$  (e.g.  $s \rightarrow v \rightarrow t$ )

- The cost is then  $C_i(f) = \sum_{e \in P_i} c_e(f_e)$ ,  $f_e = |\{i : e \in P_i\}|$  is the number of players using edge e.

- We can find a potential function  $\Phi(f) = \sum_{e \in E} \sum_{j=1}^{f_e} c_e(j)$  such that  $C_i(\hat{f}) - C_i(f) = \Phi(\hat{f}) - \Phi(f)$ .

- All potential games are also routing games.





# Best Response Dynamics in Potential Games

- **Theorem: every potential games has a pure Nash equilibrium**

- ◆ **proof:**  $s = \arg \min_{s \in S} \Phi(s)$  is the pure Nash otherwise any player had incentive to deviate would have a smaller value. This proof leads to an interesting learning algorithm that guarantees to convergence.

- Best response dynamics can lead to Nash convergence.

- While the current outcome  $\mathbf{s}$  is not a  $\epsilon$ -PNE:

- Pick an arbitrary player  $i$  that has an  $\epsilon$ -move — a deviation  $s'_i$  with  $C_i(s'_i, \mathbf{s}_{-i}) < (1-\epsilon)C_i(\mathbf{s})$  — and an arbitrary such move for the player, and move to the outcome  $(s'_i, \mathbf{s}_{-i})$ .

- A polynomial time convergence bound can be proved.

**Theorem 3.2 (Convergence of  $\epsilon$ -Best Response Dynamics [2])** Consider an atomic selfish routing game where:

1. All players have a common source vertex and a common sink vertex.
2. Cost functions satisfy the “ $\alpha$ -bounded jump condition,” meaning  $c_e(x+1) \in [c_e(x), \alpha \cdot c_e(x)]$  for every edge  $e$  and positive integer  $x$ .
3. The MaxGain variant of  $\epsilon$ -best-response dynamics is used: in every iteration, among players with an  $\epsilon$ -move available, the player who can obtain the biggest absolute cost decrease moves to its minimum-cost deviation.

Then, an  $\epsilon$ -PNE is reached in  $(\frac{k\alpha}{\epsilon} \log \frac{\Phi(\mathbf{s}^0)}{\Phi_{\min}})$  iterations.

# Best Response Dynamics in Potential Games

- Best response dynamics has polynomial time rate.

**Theorem 3.2 (Convergence of  $\epsilon$ -Best Response Dynamics [2])** Consider an atomic selfish routing game where:

1. All players have a common source vertex and a common sink vertex.
2. Cost functions satisfy the “ $\alpha$ -bounded jump condition,” meaning  $c_e(x+1) \in [c_e(x), \alpha \cdot c_e(x)]$  for every edge  $e$  and positive integer  $x$ .
3. The MaxGain variant of  $\epsilon$ -best-response dynamics is used: in every iteration, among players with an  $\epsilon$ -move available, the player who can obtain the biggest absolute cost decrease moves to its minimum-cost deviation.

Then, an  $\epsilon$ -PNE is reached in  $(\frac{k\alpha}{\epsilon} \log \frac{\Phi(s^0)}{\Phi_{\min}})$  iterations.

- Check the proof in [S. Chien 2011].
  - ♦ first prove the existence of a player with high cost
  - ♦ then prove the player chosen to take best response has cost within an  $\alpha$  factor of that of any other player.
  - ♦ third we can bound the decrease of potential function for each iteration.
- Note that polynomial in the number of joint strategies (e.g,  $2^n$  for 2 strategy  $n$  players)
- If either the assumptions on **bounded jump condition, single source and sink, MaxGain** is dropped, then it could take exponential number of iterations [A. Skopalik 2008].

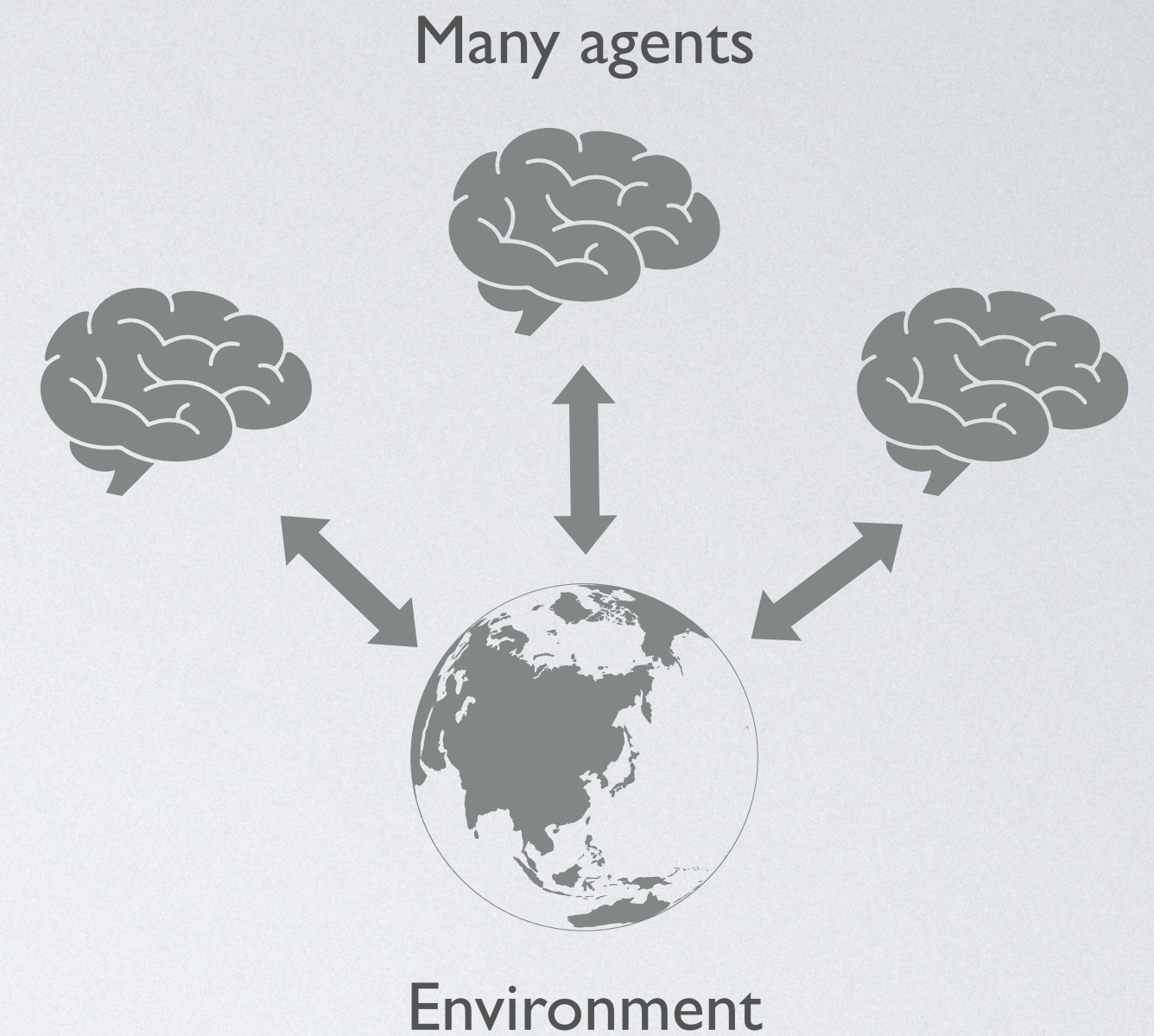
# Contents

- Algorithmic Game Theory Overview
- Computing Nash Equilibrium
  - Two-Player Zero-Sum Games (LP, fictitious play, double oracle, PSRO)
  - Two-Player General-Sum Games (support enumeration, Lemke-Howson method)
  - N-Player Potential Games (best response dynamics)
- Connections to MARL
  - MARL Formulations
  - Complexity Results (Nash is PPAD-hard)
  - Other Necessary Solution Concepts (correlated equilibrium, coarse CE)
- No-Regret Dynamics
  - Solving Coarse Correlated Equilibrium
  - Solving Two-Player Zero-Sum Games (FP is not no-regret, MWU, ODO)
  - Solving Correlated Equilibrium (swap regret)

# Multi-Agent Reinforcement Learning

- Modelled by a Stochastic Game  $(\mathcal{S}, \mathcal{A}^{\{1,\dots,n\}}, \mathcal{R}^{\{1,\dots,n\}}, \mathcal{T}, \mathcal{P}_0, \gamma)$

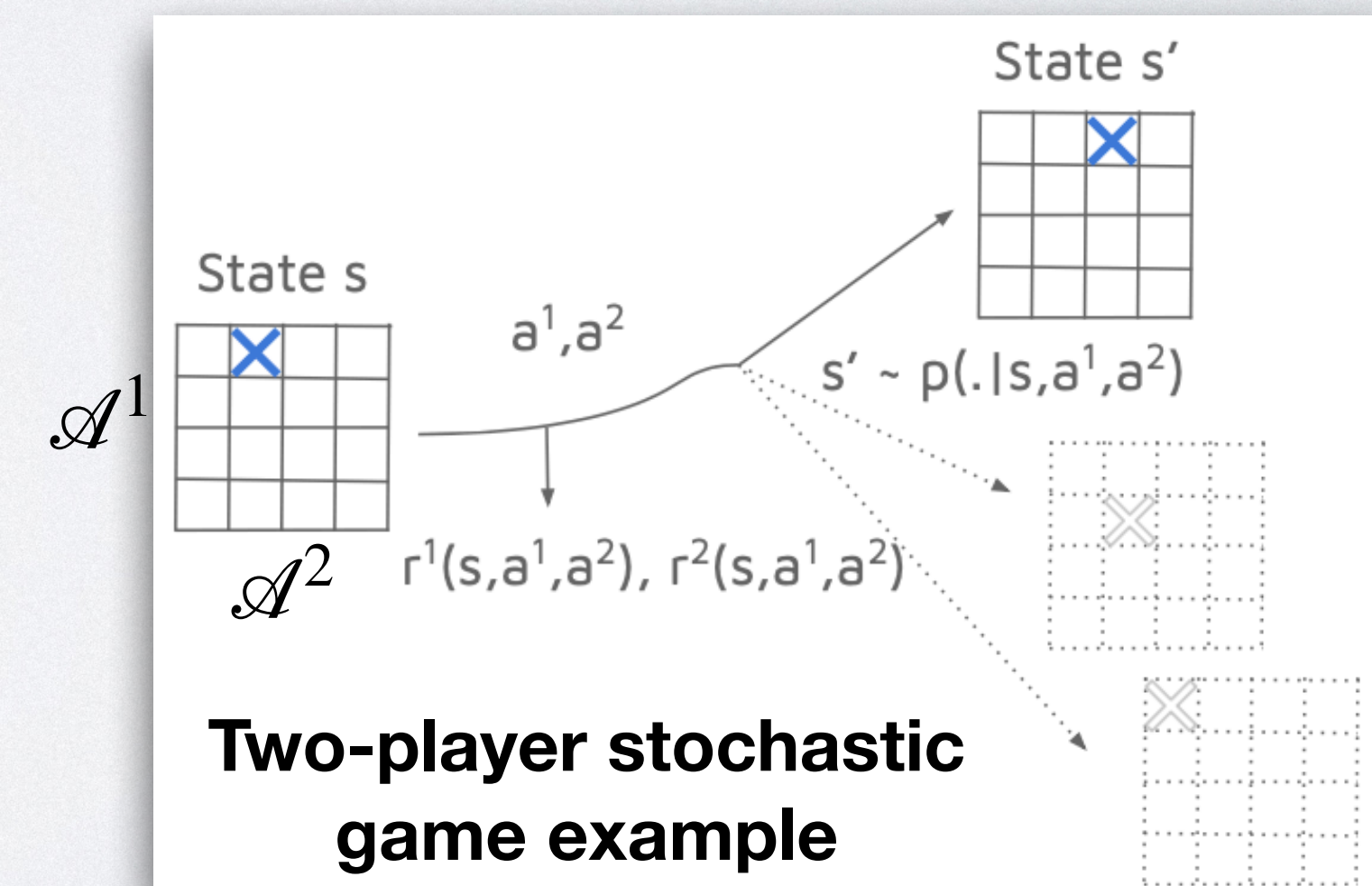
- $\mathcal{S}$  denotes the state space,
- $\mathcal{A}$  is the joint-action space  $\mathcal{A}^1 \times \dots \times \mathcal{A}^n$ ,
- $\mathcal{R}^i = \mathcal{R}^i(s, a^i, a^{-i})$  is the reward function for the i-th agent,
- $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0,1]$  is the transition function based on the joint action,
- $\mathcal{P}_0$  is the distribution of the initial state,  $\gamma$  is a discount factor.
- **Special case:**  $n = 1 \rightarrow$  single-agent MDP,  $|\mathcal{S}| = 1 \rightarrow$  normal-form game
- **Dec-POMDP:** assume state is not directly observed, but agents have same reward function.



- Each agent tries to maximise its expected long-term reward:

$$V_{i,\pi}(s) = \sum_{t=0}^{\infty} \gamma^t \mathbf{E}_{\pi, \mathcal{P}} \{ R_{i,t} | s_0 = s, \pi \}, \pi = [\pi_1, \dots, \pi_N]$$

$$Q_{i,\pi}(s, \mathbf{a}) = R_i(s, \mathbf{a}) + \gamma \mathbf{E}_{s' \sim p} [V_{i,\pi}(s')]$$



# Multi-Agent Reinforcement Learning

- Value-based method:

- The sense of optimality changes, now it depends on other agents !

$$Q_{i,t+1}(s_t, \pi_t) = Q_{i,t}(s_t, \pi_t) + \alpha [R_{i,t+1} + \gamma \cdot \mathbf{eval}_i\{Q_{\cdot,t}(s_{t+1}, \cdot)\} - Q_{i,t}(s_t, \pi_t)]$$
$$\pi_{i,t}(s, \cdot) = \mathbf{solve}_i\{Q_{\cdot,t}(s_t, \cdot)\}$$

- ◆ Fully-cooperative game: agents share the same reward function

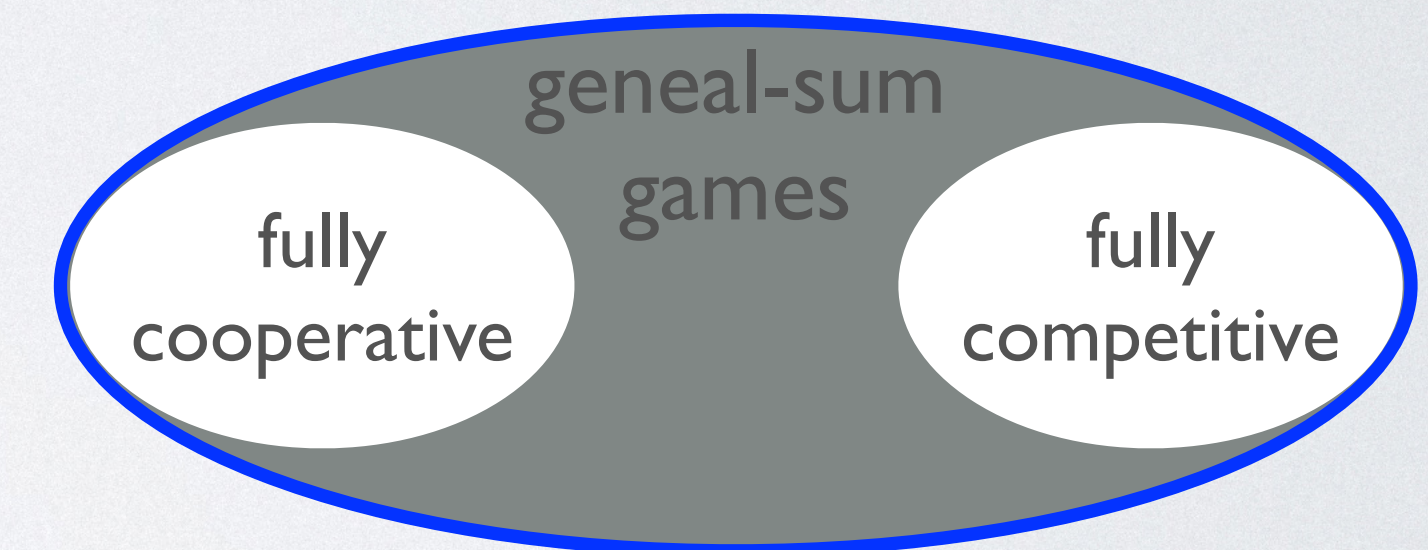
$$\mathbf{eval}_i\{Q_{\cdot,t}(s_{t+1}, \cdot)\} = \max_a Q_{i,t}(s_{t+1}, a)$$

$$\mathbf{solve}_i\{Q_{\cdot,t}(s_t, \cdot)\} = \arg \max_{a_i} \left( \max_{a_{-i}} Q_{i,t}(s_t, a_i, a_{-i}) \right)$$

- ◆ Fully-competitive game: sum of agents' reward is zero

$$\mathbf{eval}_i\{Q_{\cdot,t}(s_{t+1}, \cdot)\} = \max_{\pi_i} \min_{a_{-i}} \mathbf{E}_{\pi_i} [Q_{i,t}(s_t, a_i, a_{-i})]$$

$$\mathbf{solve}_i\{Q_{\cdot,t}(s_t, \cdot)\} = \arg \max_{\pi_i} \min_{a_{-i}} \mathbf{E}_{\pi_i} [Q_{i,t}(s_t, a_i, a_{-i})]$$



- Assuming agents share the either the same or completely opposite interest is a strong assumption.

# Nash Equilibrium to MARL

- Value-based method:

$$\pi_{i,t}(s, \cdot) = \mathbf{solve}_i \left\{ Q_{\cdot,t}(s_t, \cdot) \right\}$$

$$Q_{i,t+1}(s_k, \pi_t) = Q_{i,t}(s_t, \pi_t) + \alpha \left[ R_{i,t+1} + \gamma \cdot \mathbf{eval}_i \left\{ Q_{\cdot,t}(s_{t+1}, \cdot) \right\} - Q_{i,t}(s_t, \pi_t) \right]$$

- Nash-Q Learning [Hu. et al 2003] — Using Nash Equilibrium as the optima to guide agents' policies

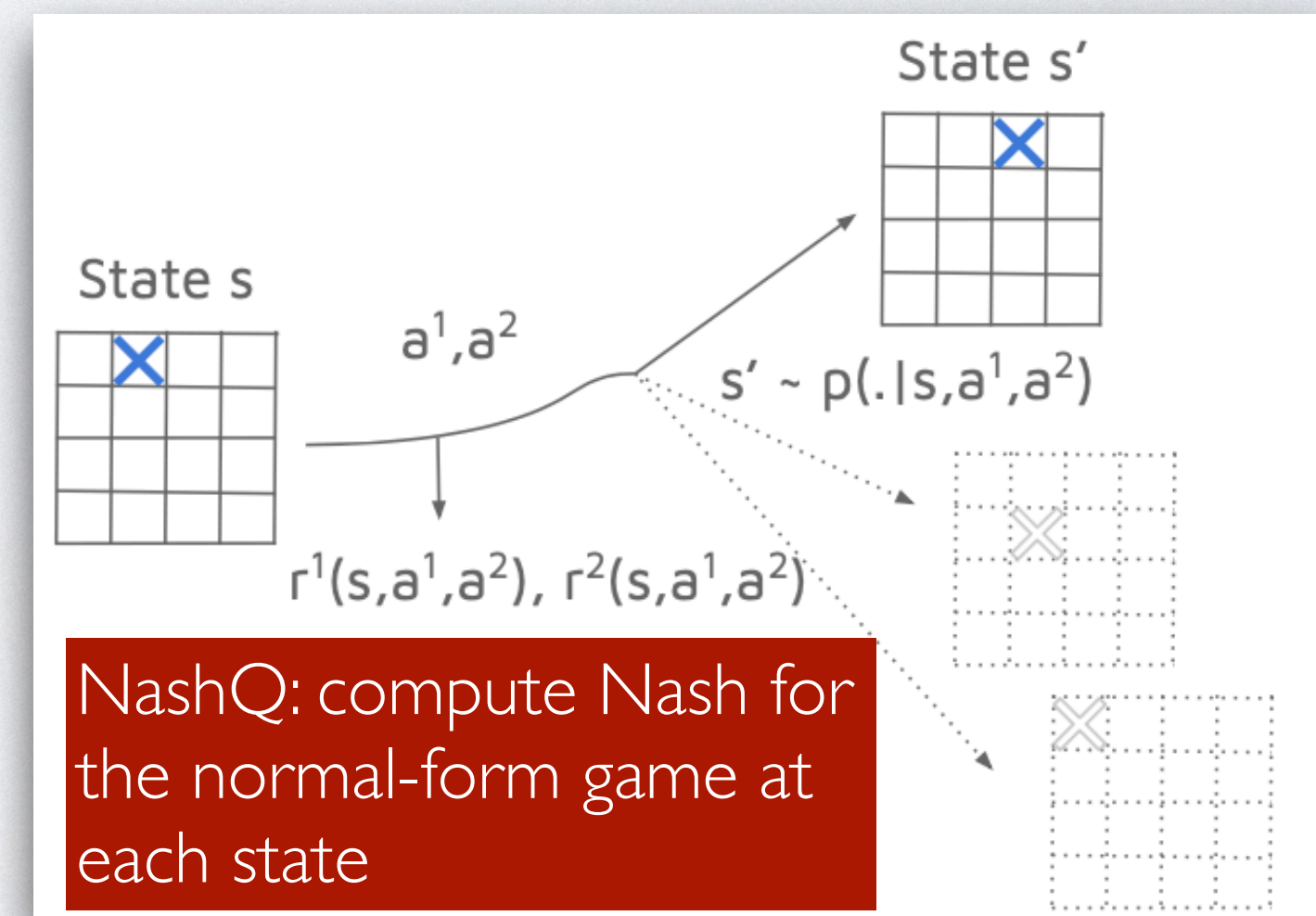
- Solve the Nash Equilibrium for the current stage game

$$\mathbf{solve}_i \left\{ Q_{\cdot,t}(s, \cdot) \right\} = \mathbf{Nash}_i \left\{ Q_{\cdot,t}(s_t, \cdot) \right\}$$

- Improve the estimation of the Q-function by the Nash value function.

$$\mathbf{eval}_i \left\{ Q_{\cdot,t}(s, \cdot) \right\} = V_i(s, \mathbf{Nash} \left\{ Q_{\cdot,t}(s_t, \cdot) \right\})$$

- Nash-Q operator  $\mathcal{H}^{\text{Nash}} \mathbf{Q}(s, \mathbf{a}) = \mathbf{E}_{s'} \left[ R(s, \mathbf{a}) + \gamma \mathbf{V}^{\text{Nash}}(s') \right]$  is a **contraction mapping**.



# MARL in Markov Potential Games

- Applying RL techniques to solve stochastic potential games [Mguni 2021].
- Stochastic potential games, considering state transition, are defined by

$$R_i(s, (a^i, a^{-i})) - R_i(s, (a^i, a^{-i})) = \phi(s, (a^i, a^{-i})) - \phi(s, (a^i, a^{-i})), \forall i \in [n], \forall s \in S$$

- Mguni found a **dual-form MDP** where the local optimum in value function corresponds to the Markov Perfect Equilibrium (the Nash equilibrium) of SPG.

the “value” function  
( $\phi$  can also be learned from data)

$$B^\pi(s) = \mathbf{E} \left[ \sum_{t=0}^{\infty} \gamma^t \phi(s_t, \mathbf{a}_t) \mid \mathbf{a}_t \sim \pi \right], \forall s \in \mathcal{S}, \forall \pi \in \Pi$$

the “TD target”

$$Y_{l_k}(s_{l_k}, \mathbf{a}_{l_k}, s'_{l_k}) := \phi_{l_k, \hat{\rho}}(s_{l_k}, \mathbf{a}_{l_k}) + \sup_{\mathbf{a}'} \mathbf{E}_{\mathbb{P}} [\hat{B}_l](s'_{l_k}, \mathbf{a}')$$

the “Bellman Error”

$$F_l \in \arg \inf_{\mathbb{F} \in \mathbb{H}} \sum_{l_k=1}^{n_k} \left( Y_{l_k}(s_{l_k}, \mathbf{a}_{l_k}, s'_{l_k}) - [\mathbb{F}](s_{l_k}, \mathbf{a}_{l_k}) \right)^2$$

the “policy gradient”

$$\nabla_{\eta_i} \hat{B}^{(\pi_{i, \eta_i}^k, \pi_{-i, \eta_{-i}}^k)}(s_{l_k}) \approx \frac{1}{L} \sum_{l=1}^L \nabla_{\eta_i} \pi_{i, \eta_i}(\cdot \mid s_{l_k}) \nabla_{a_{l_k}^i} F_k(s_{l_k}, \mathbf{a}_{l_k}) \Big|_{a_{l_k}^i \sim \pi_{i, \eta_i}^k}$$

- One can apply Q-learning/Actor Critic in the dual MDP to solve for the MPE of SPG.

# Complexity of Computing Nash Equilibrium

- Complexity theory 101 — an intuitive explanation:

- Recall the NP for a decision problem as

**Definition 4.2.1 (NP)** *A decision problem  $Q$  is in NP if there exists a polynomial time algorithm  $V(I, X)$  such that*

- 1. If  $I$  is a YES instance of  $Q$  then there exists some  $X$  such that  $|X|$  is polynomial in  $|I|$  and  $V(I, X) = YES$*
- 2. If  $I$  is a NO instance of  $Q$  then  $V(I, X) = NO$  for all  $X$*

- But the decision problem of “is there a Nash equilibrium?” is always true proved by Nash himself.
- We need a new complexity class of **Functional NP (FNP)** to describe the search problems: not only do solutions have to be verified in P-time, but also to find a solution!
- Two-player Nash will be FNP because we can check whether Nash is true by checking the best responses.
- However, two-player Nash will not be FNP-hard. To prove that, we need to show two-player Nash is not **FNP-Complete**.



# Complexity of Computing Nash Equilibrium

- Two-player Nash is not FNP-complete.

- Completeness is build on the notion of **reduction**:

**Definition 4.2.2** We say that  $P$  reduces to  $Q$  (denoted as  $P \leq_p Q$ ) if there exist polynomial-time algorithms  $A$  and  $B$  such that

1.  $A$  maps instances of  $P$  to instances of  $Q$ ,
2. If  $I$  is a YES instance of  $P$  than  $A(I)$  is a YES instance of  $Q$ , and
3. If  $X$  is a witness for  $A(I)$ , then  $B(x)$  is a witness of  $I$  (if  $I$  is a YES instance) or NO (if  $I$  is a NO instance).

- **If we want to solve  $P$ , it suffices to solve  $Q$** : to solve instance  $I$  of problem  $P$ , we can first find a solution of  $X$  of  $A(I)$ , which is of  $Q$ , and then use  $B$  to find a solution of  $B(X)$  of  $I$ .

- **Theorem: Two-player Nash is not FNP-complete**

- We can proof that if two-player Nash is FNP-hard then  $NP=coNP$  (verifying “No” instance in P-time).
  - ♦ Proof by showing that if true, we can find a certificate of NO instances for SAT problems.
  - ♦ SAT problem: find alb such that “a AND NOT b” is satisfied. SAT is known to be NP-complete.
  - ♦ But, we know that  $NP \neq coNP$ .

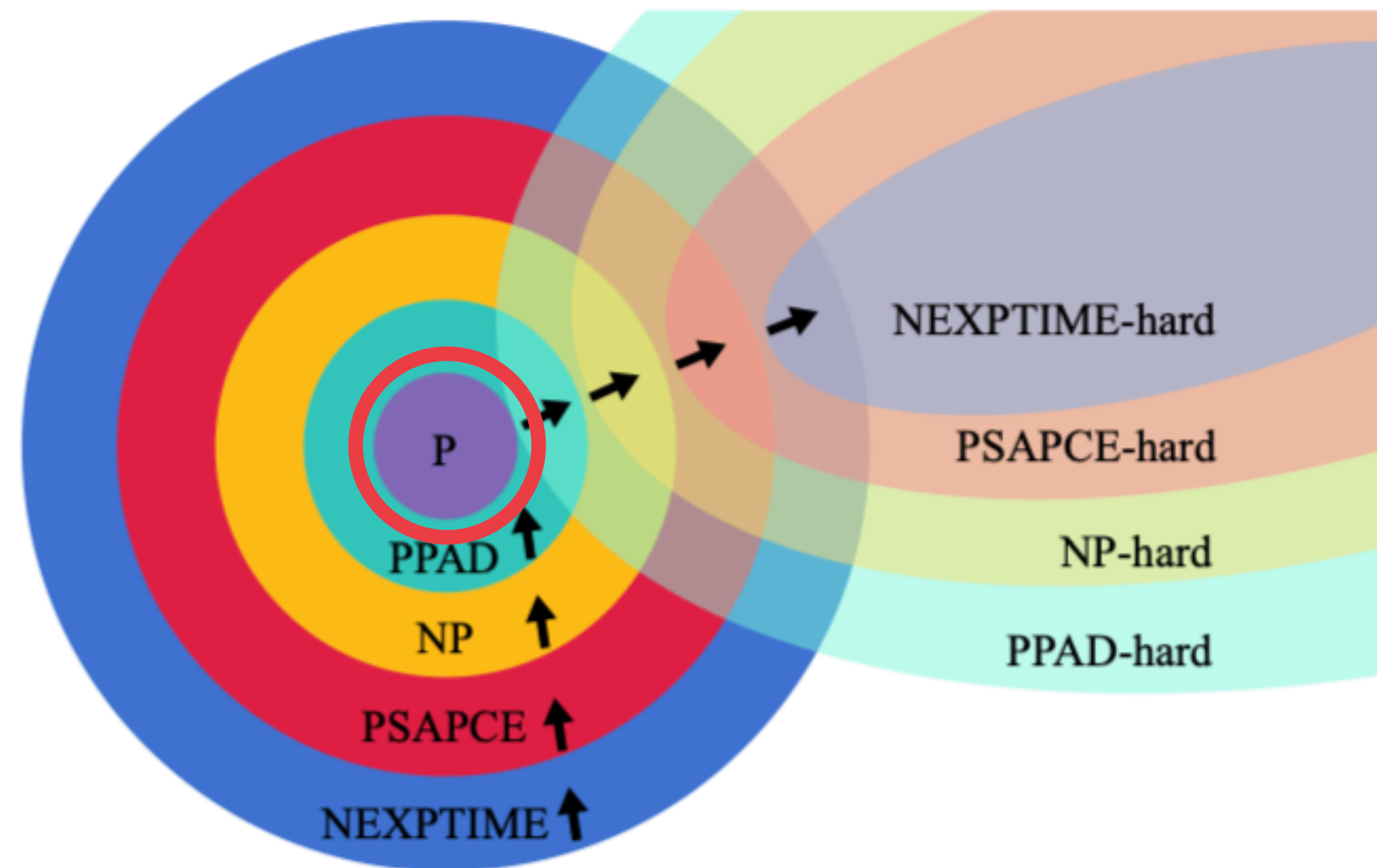
# Complexity of Computing Nash Equilibrium

- We need a new class that has complete problems for the search tasks.
  - **Polynomial Parity Argument Directed (PPAD)**: the class of search problems where the existence of a solution and an algorithm to find one are guaranteed by the properties that we have seen in the Lemke-Howson methods.
    - ♦ finite graph, vertex has at most degree 2, every source has valid solution (e.g., Nash in  $U_k, \forall k$ )
    - ♦ A bit chicken-egg here: we want to describe the complexity of a problem, but now we say all problems that look like this form the class of new complexity
  - We know PPAD problems can always have exponential-time algorithms, but can we have P-time solutions?
    - ♦ Short answer is we don't know yet. Similar to we don't know if  $P=NP$ .
    - ♦ But highly likely NO.
  - **Theorem: Two-player Nash is PPAD-complete.**

# Complexity of Multi-Agent Learning

- Solving Nash Equilibrium is very challenging !
  - The solution concept of Nash comes from game theory but it is not their main interest to find solutions.
  - Complexity of solving two-player Nash is **PPAD-Hard** (intractable unless  $P=NP$ ).
  - How to scale up multi-agent solution is open-question.
  - Approximate solution is still under development.
$$R_i(a_i, a_{-i}) \geq R_i(a'_i, a_{-i}) - \epsilon$$
$$\epsilon = .75 \rightarrow .50 \rightarrow .38 \rightarrow .37 \rightarrow .3393$$
 [Tsaknakis 2008]
  - Equilibrium selection is problematic, how to coordinate agents to agree on Nash during training is unknown.
  - Nash equilibrium assumes perfect rationality, but can be unrealistic in the real world.
- More complexity results of solving Nash [Shoham 2007, sec 4][Conitzer 2002]
  - Two-player general-sum normal-form game:
    - Compute NE  $\rightarrow$  **PPAD-Hard**
    - Count number of NE  $\rightarrow$  **#P-Hard**
    - Check uniqueness of NE  $\rightarrow$  **NP-Hard**
    - Guaranteed payoff for one player  $\rightarrow$  **NP-Hard**
    - Guaranteed sum of agents payoffs  $\rightarrow$  **NP-Hard**
    - Check action inclusion / exclusion in NE  $\rightarrow$  **NP-Hard**
  - Stochastic game:
    - Check pure-strategy NE existence  $\rightarrow$  **PSPACE-Hard**
    - Best response for arbitrary strategy  $\rightarrow$  **Not Turing-computable, even can not be implemented by a Turing PC.**
    - **It holds for two-player symmetrical game with finite time length.**

# Complexity of Multi-Agent Learning



**Figure 1.5:** Landscape of different complexity classes. Relevant examples are: 1) solving NE in two-player zero-sum game is  $P$  (Neumann, 1928). 2) solving NE in two-player general-sum game is  $PPAD$ -hard (Daskalakis et al., 2009). solving NE in three-player zero-sum game is also  $PPAD$ -hard (Daskalakis and Papadimitriou, 2005). 3) checking the uniqueness of NE is  $NP$ -hard (Conitzer and Sandholm, 2002). 4) checking whether pure-strategy NE exists in stochastic game is  $PSPACE$ -hard (Conitzer and Sandholm, 2008). 5) solving Dec-POMDP is  $NEXPTIME$ -hard (Bernstein et al., 2002).

# What is Next?

- We start from two-player zero-sum games, where solving Nash is P.
  - ◆ fictitious play, double oracle, PSRO
- We look at potential games, where pure Nash exists and solving it is P.
  - ◆ best response dynamics has polynomial bound.
  - ◆ solving stochastic potential games through MARL.
- We move onto two-player general-sum games, where solving Nash is PPAD-hard.
  - ◆ support enumeration, Lemke-Howson
  - ◆ connections to MARL
  - ◆ two-player Nash is PPAD-hard, MARL is even harder

So far, it is all about Nash. But it could be a dead-end for general cases !

# MARL in Reality

## what you Mum thinks

ARTIFICIAL INTELLIGENCE MACHINE CONSCIOUSNESS

### An Artificial Intelligence Tries to Kill her Creator

11 MONTHS AGO READ TIME: 8 MINUTES BY RAÚL ARRABALES LEAVE A COMMENT

99

Spanish researchers discover a bot trying to kill her creator. This Artificial Intelligence, designed to fight in First-Person Shooter video games, was surprised while looking for a way to end the life of her creator in the real world.

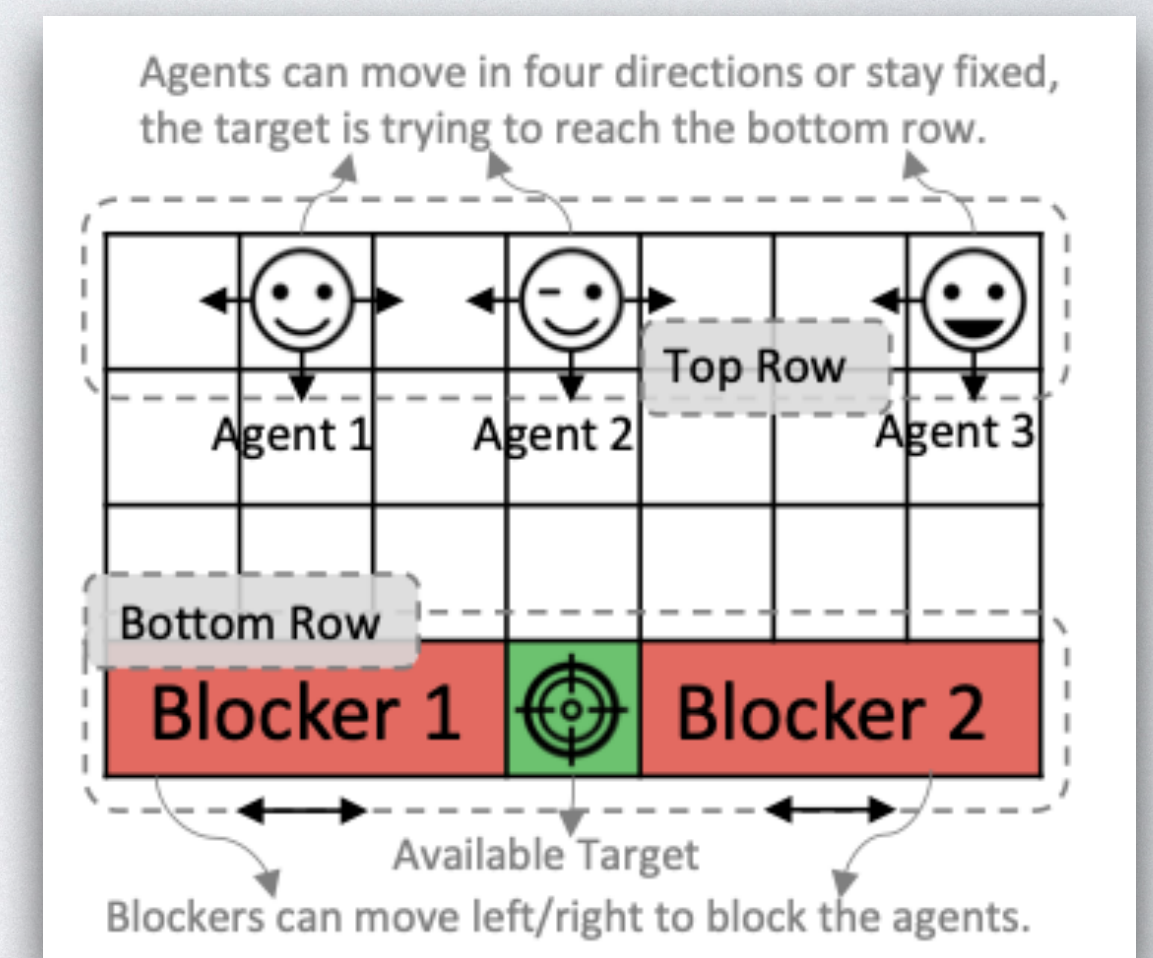
Something undescribable :)

## what you think you are doing



Multi-player general-sum games with high-dimensional continuous state-action space

## what you are actually doing



Two-player discrete-action game in a grid world.



Available online at [www.sciencedirect.com](http://www.sciencedirect.com)



Artificial Intelligence 171 (2007) 365–377

---

---

**Artificial  
Intelligence**

---

---

[www.elsevier.com/locate/artint](http://www.elsevier.com/locate/artint)

## If multi-agent learning is the answer, what is the question?

Yoav Shoham \*, Rob Powers, Trond Grenager

*Department of Computer Science, Stanford University, Stanford, CA 94305, USA*

Received 8 November 2005; received in revised form 14 February 2006; accepted 16 February 2006

Available online 30 March 2007

*“For the field to advance one cannot simply define **arbitrary learning strategies**, and analyse whether the resulting dynamics **converge in certain cases to a Nash equilibrium** or some other solution concept of the **stage game**. This in and of itself is **not well motivated**.”*

# Other Necessary Solution Concepts

- Nash Equilibrium

**Definition 5.3.1** Let  $\sigma_i$  be a distribution over  $S_i$  for all  $i \in [k]$ . Let  $\sigma = \sigma_1 \times \sigma_2 \times \dots \times \sigma_k$  be the product distribution over  $S$  defined by the individual player distributions. Then  $\sigma$  is a mixed Nash equilibrium if

$$\mathbf{E}_{s \sim \sigma} [c_i(s)] \leq \mathbf{E}_{s \sim \sigma} [c_i(s_{-i}, s'_i)]$$

for all  $i \in [k]$  and for all  $s'_i \in S_i$ .

- Correlated Equilibrium

**Definition 5.3.2** Let  $\sigma$  be a distribution over  $S = S_1 \times \dots \times S_k$ . Then  $\sigma$  is a correlated equilibrium if

$$\mathbf{E}_{s \sim \sigma} [c_i(s) | s_i] \leq \mathbf{E}_{s \sim \sigma} [c_i(s_{-i}, s'_i) | s_i]$$

for all  $i \in [k]$  and for all  $s_i, s'_i \in S_i$ .

- Coarse Correlated Equilibrium

**Definition 5.3.3** Let  $\sigma$  be a distribution over  $S = S_1 \times \dots \times S_k$ . Then  $\sigma$  is a coarse correlated equilibrium if

$$\mathbf{E}_{s \sim \sigma} [c_i(s)] \leq \mathbf{E}_{s \sim \sigma} [c_i(s_{-i}, s'_i)]$$

for all  $i \in [k]$  and for all  $s'_i \in S_i$ .

**NE  $\subset$  CE  $\subset$  CCE**



# Other Necessary Solution Concepts

- **NE  $\subset$  CE  $\subset$  CCE**

- **Distinctions:**

- ◆ Nash requires a product of individual distribution:  $\sigma = \sigma_1 \times \sigma_2 \times \dots \times \sigma_k$
- ◆ CE/CCE requires arbitrary distribution over the joint strategy set:  $\sigma \in \Delta_{S=S_1 \times \dots \times S_k}$
- ◆ CE: assuming a third party draw  $s \sim \sigma$ , and privately tell each player his  $s_i$  but not  $s_{-i}$ , and each player knows  $\sigma$ , now he decides whether to deviate to achieve  $\mathbf{E}_{s \sim \sigma} [c_i(s_{-i}, s'_i) \mid s_i]$

- ◆ In the example of 

	Yield	Rush
Yield	(0, 0)	(1, 2)
Rush	(2, 1)	(0, 0)

 :

- ◆ (yield, rush) and (rush, yield) are Nash, but sticking to them will block one side forever.
- ◆ but we cannot let them randomly play 50%-50%, because 25% chance they will crash !
- ◆ Instead, we want  $\sigma(\text{yield, stop}) = \sigma(\text{stop, yield}) = 1/2$ .
- ◆ When the green light is on, we know others will yield to us. Both of us have no motivations to violate.
- ◆ **CCE vs CE: player needs to decide whether to deviate even before being told  $s_i$ .**
  - ◆ If we have no incentive to deviate no matter what we're told, then we will not deviate even if  $s_i$  is unknown.

# Contents

- Algorithmic Game Theory Overview
- Computing Nash Equilibrium
  - Two-Player Zero-Sum Games (LP, fictitious play, double oracle, PSRO)
  - Two-Player General-Sum Games (support enumeration, Lemke-Howson method)
  - N-Player Potential Games (best response dynamics)
- Connections to MARL
  - MARL Formulations
  - Complexity Results (Nash is PPAD-hard)
  - Other Necessary Solution Concepts (correlated equilibrium, coarse CE)
- No-Regret Dynamics
  - Solving Coarse Correlated Equilibrium
  - Solving Two-Player Zero-Sum Games (FP is not no-regret, MWU, ODO)
  - Solving Correlated Equilibrium (swap regret)

# Two Mainstreams of Multi-Agent Learning algorithms

## Best response methods:

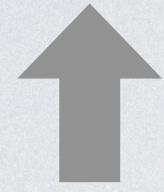
- Fictitious play, double oracle, PSRO series, ...
- Regard the opponents fixed and seek for best responses.
- Easily and nicely integrated with RL methods (e.g., NFSP, PSRO)
- Work effectively in potential and zero-sum games, but limited in general-sum games.
- Average policy have convergence guarantee but generally no last-iteration convergence

## No-regret methods

- MWU, Follow the Regularised/Perturbed leader, CFR and all kinds of CFR variants, MCTS, ...
- Work in a self-play settings, no best-response step but a no-regret step.
- Often requires to know the model (the game tree, utility function/strategies of opponents, etc)
- A portal to the arsenal of online learning tools
- Have nice convergence guarantee to Nash zero-sum games, and CE/CCE in general-sum games.

# No-Regret vs. Best-Response Methods in Zero-Sum Games

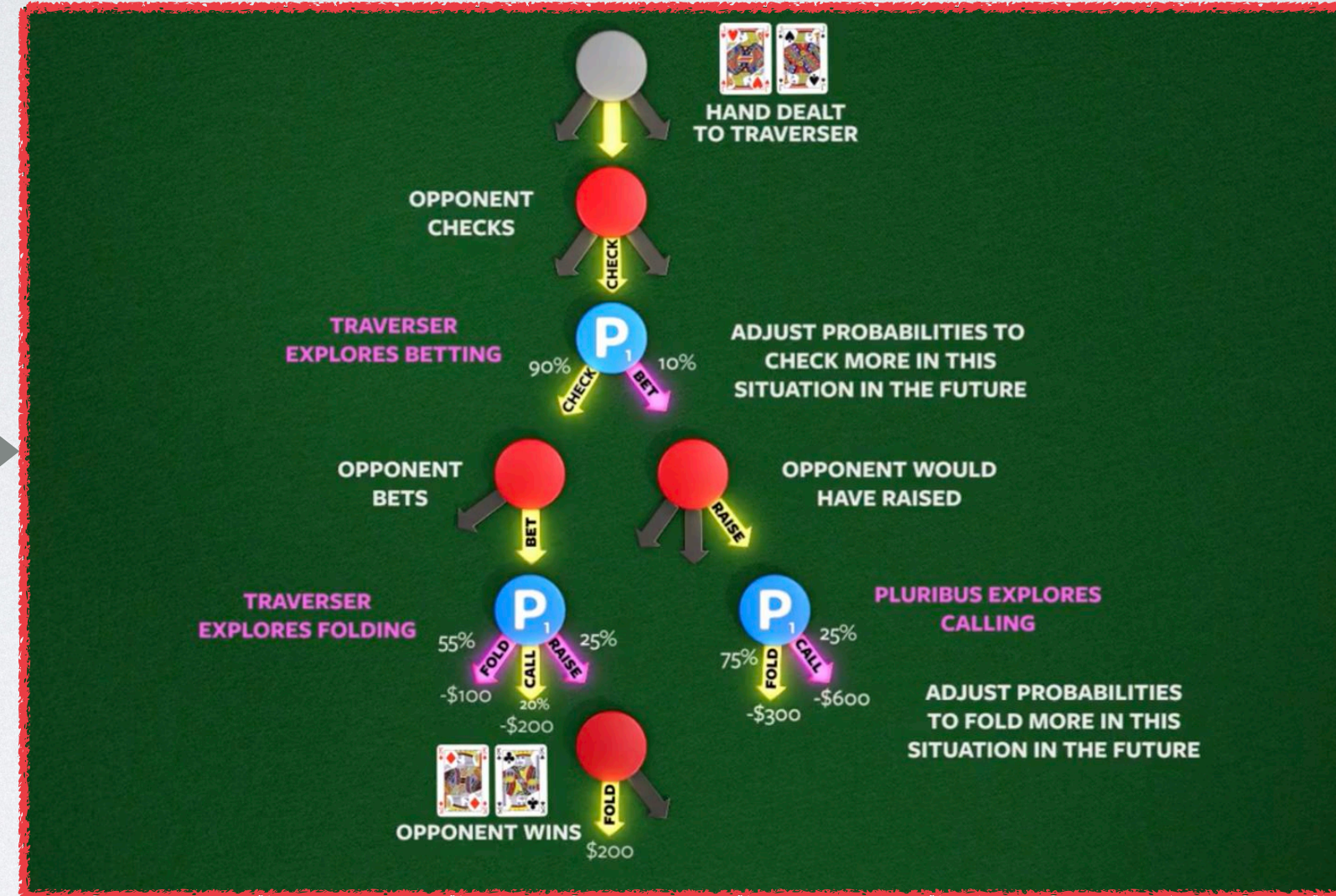
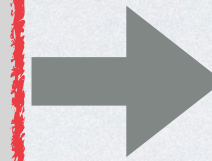
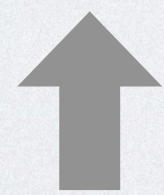
Output: the reward  $(R^1, \dots, R^N)$



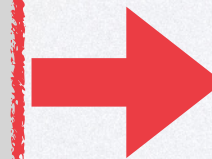
**Black-box multi-agent game engine**



Input: a joint strategy  $(\pi^1, \dots, \pi^N)$



Regret based methods: Poker Type



Best response based methods: StarCraft type

When planning is feasible (game tree is easily accessible), existing techniques can solve the games really well.

Perfect-information games:  
MCTS, alpha-beta search, AlphaGO series (AlphaZero, MuZero, etc)

Imperfect-information:  
CFR series (DeepCFR, Libratus/Pluribus, Deepstack), XFP/NFSP series

When planning is not feasible. StarCraft has  $10^{26}$  choices per time step vs. the whole tree of chess  $10^{50}$  (Texas holdem  $10^{80}$ , GO  $10^{170}$ )

Enumerating all policies' actions at each state and then playing a best response is infeasible. But an approximate BR can be computed.

Solution: training a population of RL agents, treat each RL agent as one "pure strategy" and solve the game in the meta-level (e.g. PSRO methods).

# No-Regret vs. Best-Response Methods in Zero-Sum Games

## Actor-Critic Policy Optimization in Partially Observable Multiagent Environments

Sriram Srinivasan<sup>\*,1</sup>  
srsrinivasan@

Marc Lanctot<sup>\*,1</sup>  
lanctot@

Vinicius

Karl Tuyls<sup>1</sup>  
karltuyls@

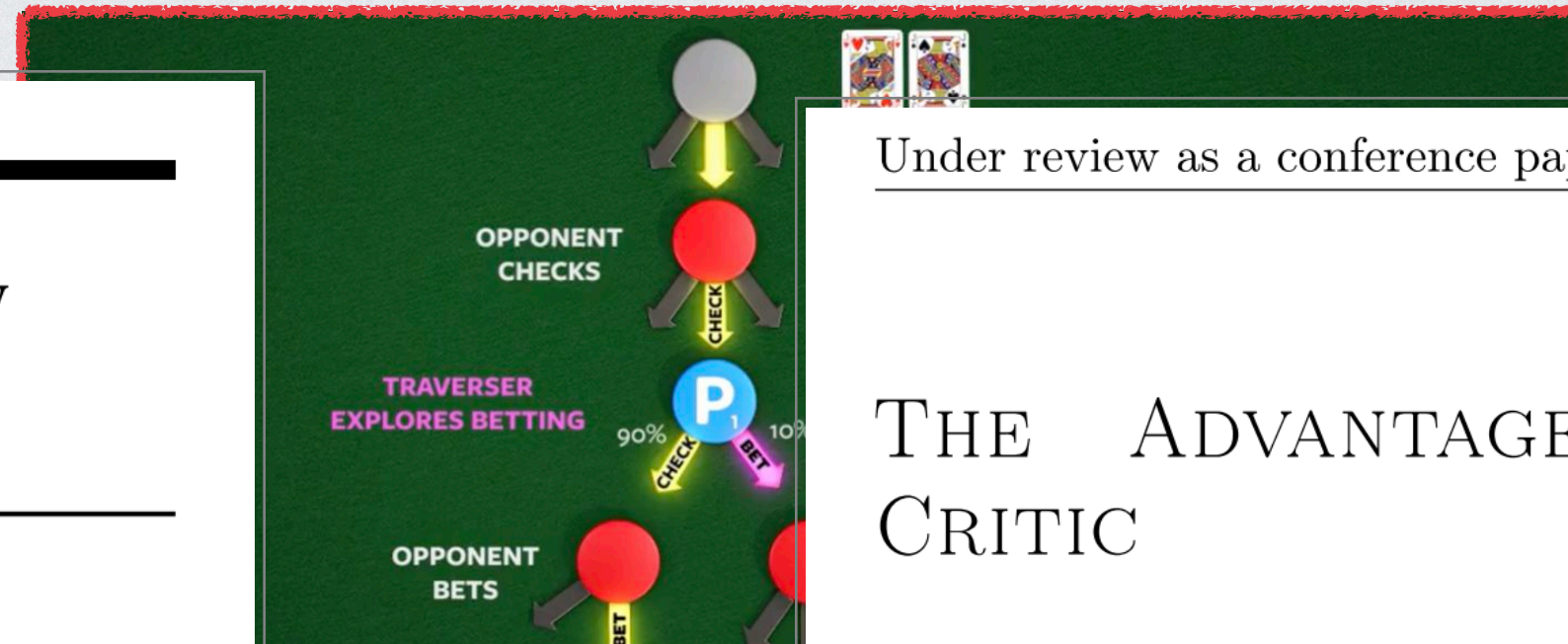
Rémi Munos<sup>1</sup>  
munos@

...@google.com. <sup>1</sup>DeepMind. \*These authors

### Abstract

Optimization of parameterized policies for reinforcement learning is a fundamental and challenging problem in artificial intelligence. Among the most common approaches are algorithms based on gradient ascent of a score function representing discounted return. In this paper, we examine the role of these policy gradient and actor-critic algorithms in partially-observable multiagent environments. We show several candidate policy update rules and relate them to a foundation of regret minimization and multiagent learning techniques for the one-shot and tabular cases. We apply our method to *model-free* multiagent reinforcement learning in adversarial sequential decision problems (zero-sum imperfect information games), using RL-style function approximation. We evaluate on commonly used benchmark Poker domains, showing performance against fixed policies and empirical convergence to approximate Nash equilibria in self-play with rates similar to or better than a baseline model-free algorithm for zero-sum games, without any domain-specific state space reductions.

Input: a joint strategy  $(\pi^1, \dots, \pi^N)$



When planning is feasible (game tree is

Under review as a conference paper at ICLR 2021

THE ADVANTAGE REGRET-MATCHING ACTOR-CRITIC

They are equivalent  
Regret  $\approx$  Value!



Best response based methods: StarCraft type

are used to predict conditional advantages which, combined with regret matching, produces a new policy. In particular, ARMAC learns from sampled trajectories in a centralized training setting, without requiring the application of importance sampling commonly used in Monte Carlo counterfactual regret (CFR) minimization; hence, it does not suffer from excessive variance in large environments. In the single-agent setting, ARMAC shows an interesting form of exploration by keeping past policies intact. In the multiagent setting, ARMAC in self-play approaches Nash equilibria on some partially-observable zero-sum benchmarks. We provide exploitability estimates in the significantly larger game of betting-abstracted no-limit Texas Hold'em.

the problem problem, auto-curricula.

# Online Learning and No-Regret

- The settings of online learning:

- ◆ The algorithm picks a strategy  $p^t \in \Delta_{|A|}$  at time step  $t$
- ◆ The adversary/nature picks cost vector  $c^t : A \rightarrow [0,1]$
- ◆ An action  $a^t$  is drawn from  $p^t$ , and the algorithm incurs cost of  $c^t(a^t)$
- ◆ **Full-information settings:** observe the entire cost vector  $c^t$ . **Bandit settings:** only observe selected  $c^t(a^t)$
- ◆ **Oblivious adversary:**  $c^t$  only depends on  $t$ . **Adaptive adversary:**  $c^t$  depends on  $\{t, (p^1, \dots, p^t), (a^1, \dots, a^{t-1})\}$
- ◆ **Goal:** we try to learn how should we adapt our algorithms, **learn from mistakes [remember Alan Turing]!**

- The (external) regret of a sequence of actions w.r.t action  $a \in A$ :

$$R_T(a) = \frac{1}{T} \left( \sum_{t=1}^T c^t(a^t) - \sum_{t=1}^T c^t(a) \right)$$

- A no-(external)-regret algorithm  $\mathcal{A}$  is said to be **no-regret** (also **Hannan consistent**) if:

$$\lim_{T \rightarrow \infty} \mathbf{E} [R_T^{\mathcal{A}}(a)] = \frac{1}{T} \left( \sum_{i=1}^T \mathbf{E}_{a^i \sim p^i} [c^i(a^i)] - \sum_{t=1}^T c^t(a) \right) = 0, \forall a \in A$$

# No-Regret Learning

- No-regret to the best action sequence  $\sum_{t=1}^T \min_{a \in A} c^t(a)$  in hindsight is impossible.
  - ◆ When the adversary exploits you, in hindsight, a zero-loss best action sequence is always possible. If we know  $c^t$  before each iteration we can have  $\sum_{t=1}^T \min_{a \in A} c^t(a) = 0$ , then the regret explodes.

$$\mathbf{E} [R_T^{\mathcal{A}}(a)] = \frac{1}{T} \left( \sum_{i=1}^T \mathbf{E}_{a^i \sim p^i} [c^i(a^i)] - \sum_{t=1}^T \min_{a \in A} c^t(a) \right) = \mathcal{O}(T) \neq \mathcal{O}(1) \quad \text{not sublinear!}$$

- ◆ However, no-regret to the best action (not sequence!) in hindsight is possible:

$$\sum_{t=1}^T \min_{a \in A} c^t(a) \quad \text{vs.} \quad \min_{a \in A} \sum_{t=1}^T c^t(a)$$

best action sequence                      best action in hindsight

- In games, no-regret is the minimal requirement for rationale.

# No-Regret Learning towards Coarse Correlated Equilibrium

- In games, no-regret is the minimal requirement for rationale.
- Recall no-regret is

$$\lim_{T \rightarrow \infty} \mathbf{E} [R_T^{\mathcal{A}}(a)] = \frac{1}{T} \left( \sum_{i=1}^T \mathbf{E}_{a^t \sim p^t} [c^t(a^t)] - \sum_{t=1}^T c^t(a) \right) = 0, \forall a \in A$$

- **Theorem: if all players adopt no-regret algorithms such that  $\mathbf{E}[R_T^{\mathcal{A}_i}(a)] \leq \epsilon$ ,  $\forall i \in [k], a \in S_i$ , then the average distribution  $\sigma = \sum_{t=1}^T \prod_{i=1}^k \frac{p_i^t}{T}$  is an  $\epsilon$ -CCE, i.e.,  $\mathbf{E}_{s \sim \sigma} [C_i(s)] \leq \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, s'_i)] + \epsilon$**

- Proof:

$$\begin{aligned} \mathbf{E}_{s \sim \sigma} [C_i(s)] - \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, s'_i)] &= \frac{1}{T} \sum_{t=1}^T \mathbf{E}_{s \sim \sigma^t} [C_i(s)] - \frac{1}{T} \sum_{t=1}^T \mathbf{E}_{s \sim \sigma^t} [C_i(s_{-i}, s'_i)] \\ &= \frac{1}{T} \left( \sum_{t=1}^T \mathbf{E}_{a^t \sim p_i^t} [c_i^t(a^t)] - \sum_{t=1}^T c_i^t(s'_i) \right) \\ &= \mathbf{E} [R_T^{\mathcal{A}_i}(s'_i)] \leq \epsilon \end{aligned}$$

Note that in game playing, we have:  
 $c_i^t(a) = \mathbf{E}_{s \sim \sigma^t} [C_i(s_{-i}, a)]$



# No-Regret Learning in Zero-Sum Games

- We have showed that if all players play no-regret methods, they can reach a CCE.
- We can further show that no-regret players will reach the Nash in zero-sum games.
  - For better clarity. Let's assume row player chooses an action  $I_t \in \{1, \dots, N\}$ , mixed strategy  $\mathbf{p}_t = (p_{1,t}, \dots, p_{N,t})$
  - Column player instead of deciding  $c^t$ , let's assume they choose an action  $J_t \in \{1, \dots, M\}$  and  $\mathbf{q}_t = (q_{1,t}, \dots, q_{M,t})$
  - Assuming both players adopt no-regret algorithm regardless of what the opponent does, such that

$$\limsup_{T \rightarrow \infty} \left( \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) - \min_{i=1, \dots, N} \frac{1}{T} \sum_{t=1}^T \ell(i, J_t) \right) \leq 0$$

- Recall that the Nash value of a two-player zero-sum game is  $\max_{\mathbf{q}} \min_{\mathbf{p}} \sum_{i=1}^N \sum_{j=1}^M p_i q_j \ell(i, j)$
- We have the main theorem that

**Theorem: assuming that in a two-player zero-sum game, if both players play no-regret**

**methods, then**  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) = V = \max_{\mathbf{q}} \min_{\mathbf{p}} \sum_{i=1}^N \sum_{j=1}^M p_i q_j \ell(i, j)$  **almost surely.**

# No-Regret Learning in Zero-Sum Games

**Theorem: assuming that in a two-player zero-sum game, if both players play no-regret**

**methods, then**  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) = V = \max_{\mathbf{q}} \min_{\mathbf{p}} \sum_{i=1}^N \sum_{j=1}^M p_i q_j \ell(i, j)$  **almost surely.**

**Proof:** 1) we first should that regardless of what column player plays, if the row player plays no-regret method, his loss will be no more than maximin value (i.e., worst case scenario to row player)  $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) \leq V$ .

Since the row player adopts no-regret method

we only need to show:

$$\begin{aligned} \min_{i=1, \dots, N} \frac{1}{T} \sum_{t=1}^T \ell(i, J_t) &= \min_{\mathbf{p}} \frac{1}{T} \sum_{t=1}^T \left( \sum_{i=1}^N p_i \ell(i, J_t) \right) \\ &= \min_{\mathbf{p}} \frac{1}{T} \sum_{t=1}^T \bar{\ell}(\mathbf{p}, J_t) \\ &= \min_{\mathbf{p}} \sum_{j=1}^M \left( \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{\{J_t=j\}} \bar{\ell}(\mathbf{p}, j) \right) \\ &= \min_{\mathbf{p}} \bar{\ell}(\mathbf{p}, \hat{\mathbf{q}}_T) \\ &\leq \max_{\mathbf{q}} \min_{\mathbf{p}} \bar{\ell}(\mathbf{p}, \mathbf{q}) = V \end{aligned}$$

pure strategy is a special mixed strategy

change of notation

empirical mean on column player's action

expectation over the empirical mean of column player

done!

# No-Regret Learning in Zero-Sum Games

**Theorem: assuming that in a two-player zero-sum game, if both players play no-regret methods, then**  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) = V = \max_{\mathbf{q}} \min_{\mathbf{p}} \sum_{i=1}^N \sum_{j=1}^M p_i q_j \ell(i, j)$  **almost surely.**

**Proof:** 2) we proved that  $\min_{i=1, \dots, N} \frac{1}{T} \sum_{t=1}^T \ell(i, J_t) \leq V$ , and since row player plays no-regret

$$\limsup_{T \rightarrow \infty} \left( \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) - \min_{i=1, \dots, N} \frac{1}{T} \sum_{t=1}^T \ell(i, J_t) \right) \leq 0$$

we can know that

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) \leq V$$

for row player, we can do the same for the column player

$$\liminf_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) \geq \min_{\mathbf{p}} \max_{\mathbf{q}} \bar{\ell}(\mathbf{p}, \mathbf{q}) = V$$

By von Neumann's minimax theorem, we prove that

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) = V$$

This mean no-regret players self-play will reach to Nash equilibrium in two-player zero-sum games.

# No-Regret Learning in Zero-Sum Games

- We have showed that no-regret players will reach Nash equilibrium value.
- Furthermore, we can show that they reach to the Nash strategy.

**Theorem: In a two-player zero-sum game, if both players play no-regret methods, then**  
 $\hat{p}_{i,T} = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{\{I_t=i\}}, \quad \hat{q}_{j,T} = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{\{J_t=j\}}$  **almost surely converge to the set of Nash equilibrium.**

**Proof:** in the previous proof, we have shown that

$$\min_{\mathbf{p}} \bar{\ell}(\mathbf{p}, \hat{\mathbf{q}}_T) \leq V, \quad \max_{\mathbf{q}} \bar{\ell}(\hat{\mathbf{p}}_T, \mathbf{q}) \geq V$$

and due to the uniqueness of  $V$  value in zero-sum games

$$\min_{\mathbf{p}} \bar{\ell}(\mathbf{p}, \hat{\mathbf{q}}_T) = \max_{\mathbf{q}} \bar{\ell}(\hat{\mathbf{p}}_T, \mathbf{q}) = V$$

and because of

$$\bar{\ell}(\hat{\mathbf{p}}_T, \hat{\mathbf{q}}_T) \geq \min_{\mathbf{p}} \bar{\ell}(\mathbf{p}, \hat{\mathbf{q}}_T), \quad \max_{\mathbf{q}} \bar{\ell}(\hat{\mathbf{p}}_T, \mathbf{q}) \geq \bar{\ell}(\hat{\mathbf{p}}_T, \hat{\mathbf{q}}_T)$$

finally, we prove that

$$\bar{\ell}(\hat{\mathbf{p}}_T, \hat{\mathbf{q}}_T) = V$$

This means that the empirical mean of  $\hat{\mathbf{p}}_T, \hat{\mathbf{q}}_T$  is the Nash

# Fictitious Play is Not No-Regret

- No-regret can lead to CCE in general-sum, NE in two-player zero-sum.
- But, what exactly is a no-regret algorithm? How can we behave to achieve no-regret?

$$\limsup_{T \rightarrow \infty} \left( \frac{1}{T} \sum_{t=1}^n \ell(I_t, J_t) - \min_{i=1, \dots, N} \frac{1}{T} \sum_{t=1}^n \ell(i, J_t) \right) \leq 0$$

- **Surprisingly, Fictitious play is not no-regret !**

- ◆ Recall that fictitious play, also known as **Follow-the-Leader**, is to take the best response to the average cumulative loss.

$$I_t = \operatorname{argmin}_{i=1, \dots, N} \left( \frac{1}{t-1} \sum_{t=1}^{t-1} \ell(i, J_t) \right)$$

- ◆ Consider  $N = 2$  actions, let  $J_t$  be chosen such that  $\ell(1, J_t) = (1/2, 0, 1, 0, 1, \dots)$  and  $\ell(2, J_t) = (1/2, 1, 0, 1, 0, \dots)$
- ◆ Then the accumulative loss for both actions is **T/2**, and the FP will suffer a loss of **T**, thus has constant regret.
- ◆ There are many variants that make FP no-regret, for example **Follow-the-Perturbed-Leader**:

$$I_t = \operatorname{argmin}_{i=1, \dots, N} \left( \frac{1}{t-1} \sum_{t=1}^{t-1} \ell(i, J_t) + Z_{i,t} \right), \mathbf{Z}_t \text{ any i.i.d random vectors of size } N$$

# No-Regret Algorithms: MWU

- Fictitious play is not no-regret !
- Let's introduce a true no-regret algorithm: **Multiplicative Weight Update [Freund 1999]**.
  - ♦ MWU has many names: *Hedge, Exponential Weights Algorithms, Randomised Weighted Majority*

$$p_{i,t} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell(i, J_s)\right)}{\sum_{k=1}^{|A|} \exp\left(-\eta \sum_{s=1}^{t-1} \ell(k, J_s)\right)}$$

- ♦ Equivalently, one can think of the following iterative process:

$$p_t = \frac{w^t}{\sum_{a \in A} w^t(a)}, \quad w^{t+1}(a) = w^t(a) \cdot \exp\left(-\eta \ell(a, J_t)\right), \quad w^1(a) = \mathbf{1} \quad \forall a \in A$$

- Large  $\eta$  means more exploitation, small  $\eta$  means more exploration.
- ♦ Equivalently, one can get MWU by the following maximum entropy framework (a common trick in RL).

$$\arg \min_{p \in \Delta_{|A|}} \bar{\ell}(p, J_s) + 1/\eta \cdot \sum_i p_i \log p_i$$

similar to soft Q-learning !

# No-Regret Algorithms: MWU

- Let's introduce a true no-regret algorithm — **Multiplicative Weight Update**.
  - ◆ Equivalently, one can get MWU by the following maximum entropy framework (a common trick in RL).

$$\arg \min_{p \in \Delta_{|A|}} \bar{\ell}(p, J_s) + 1/\eta \cdot \sum_i p_i \log p_i = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \ell(i, J_s)\right)}{\sum_{k=1}^{|A|} \exp\left(-\eta \sum_{s=1}^{t-1} \ell(k, J_s)\right)}$$

$$\arg \min_p \ell(p, J_s) + 1/\eta \cdot \sum_i p_i \log p_i$$

$$= \sum_i p_i \ell(i, J_s) + 1/\eta \cdot \sum_i p_i \log p_i$$

$$= 1/\eta \sum_i p_i (\eta \cdot \ell(i, J_s) + \log p_i)$$

$$= 1/\eta E_p [\log \exp(\eta \cdot \ell(i, J_s)) + \log p_i]$$

$$= 1/\eta E_p [-\log \exp(-\eta \cdot \ell(i, J_s)) + \log p_i]$$

$$= 1/\eta E_p \left[ \log \frac{p_i}{\exp(-\eta \cdot \ell(i, J_s))} \right]$$

$$p_i^* = \exp(-\eta \cdot \ell(i, J_s)) \text{ to minimize}$$

$$\text{Therefore } p = \exp(-\eta \cdot \ell(i, J_s)) / \sum_i \exp(-\eta \cdot \ell(i, J_s))$$

# No-Regret Algorithms: MWU

- Let's introduce a true no-regret algorithm — **Multiplicative Weight Update**.

$$p_t = \frac{w^t}{\sum_{a \in A} w^t(a)}, \quad w^{t+1}(a) = w^t(a) \cdot \exp(-\eta \ell(a, J_t)), \quad w^1(a) = \mathbf{1} \quad \forall a \in A$$

- Let's now show MWU is indeed a no-regret method

- Let  $\mathbf{opt} = \min_{a \in A} \sum_{t=1}^n \ell(a, J_t)$  and  $\mathbf{v}^t = \sum_{a \in |A|} p_t(a) \cdot \ell(a, J_t)$ , thus the regret-bound is  $\lim_{T \rightarrow \infty} \frac{1}{T} \left( \sum_{t=1}^T \mathbf{v}^t - \mathbf{opt} \right)$

- We first bound the denominator  $\sum_{a \in A} w^t(a) \geq w^t(a^*) = \exp(-\eta \cdot \mathbf{opt})$  multiplication rule of exp function  
 $\geq (1 - \epsilon)^{\mathbf{opt}}$  assume  $1 - \epsilon = e^{-\eta}$

$$\sum_a w^{t+1}(a) = \sum_a w^t(a) \cdot \exp(-\eta \ell(a, J_t)) := \sum_a w^t(a) \cdot (1 - \epsilon)^{\ell(a, J_t)}$$
assume  $1 - \epsilon = e^{-\eta}$

$$\leq \sum_a w^t(a) \cdot (1 - \epsilon \ell(a, J_t))$$
Taylor expansion

$$= \sum_a w^t(a) \cdot (1 - \epsilon \mathbf{v}^t)$$
definition of  $\mathbf{v}^t = \sum_{a \in |A|} \frac{w^t(a)}{\sum_a w^t(a)} \cdot \ell(a, J_t)$



# No-Regret Algorithms: MWU

- Regret bound of MWU

- ♦ Let  $\mathbf{opt} = \min_{a \in A} \sum_{t=1}^n \ell(a, J_t)$  and  $\mathbf{v}^t = \sum_{a \in |A|} p_t(a) \cdot \ell(a, J_t)$ , thus the regret-bound is  $\frac{1}{T} \left( \sum_{t=1}^T v^t - \mathbf{opt} \right)$

- ♦ Merge the upper and lower bound in the last slides

$$(1 - \epsilon)^{\mathbf{opt}} \leq \sum_{a \in A} w^t(a) \leq \sum_a w^1(a) \prod_{t=1}^T (1 - \epsilon \mathbf{v}^t) = |A| \prod_{t=1}^T (1 - \epsilon \mathbf{v}^t)$$

$$\mathbf{opt} \cdot \ln(1 - \epsilon) \leq \ln |A| + \sum_{t=1}^T \ln(1 - \epsilon \mathbf{v}^t)$$

$$\mathbf{opt} \cdot (-\epsilon - \epsilon^2) \leq \ln |A| + \sum_{t=1}^T (-\epsilon \mathbf{v}^t)$$

Taylor expansion  $\ln(1 - x) = -x - x^2/2$

- ♦ Set  $\epsilon = \sqrt{\ln |A| / T}$  we conclude the proof

$$\frac{1}{T} \left( \sum_{t=1}^T v^t - \mathbf{opt} \right) \leq \frac{1}{T} \left( \epsilon T + \frac{1}{\epsilon} \ln |A| \right) = \frac{1}{T} (\sqrt{T \ln |A|} + \sqrt{T \ln |A|}) \leq 2 \sqrt{\frac{\ln |A|}{T}} \rightarrow 0$$

Note: the log term is great for many real-world problems!

# No-Regret Algorithms: Online Double Oracle

## Online Double Oracle

Le Cong Dinh<sup>\*,1,2</sup>, Yaodong Yang<sup>\*,1,4</sup>, Nicolas Perez-Nieves<sup>3</sup>, Oliver Slumbers<sup>4</sup>,

Zheng Tian<sup>4</sup>, David Henry Mguni<sup>1</sup>, Haitham Bou Ammar<sup>1</sup>, Jun Wang<sup>1,4</sup>

1. Nash is unexploitable, but when a player always plays Rock, you should play Paper rather than (1/3, 1/3, 1/3).
2. Double Oracle/PSRO assumes both players play the worst-case scenario, can be too **pessimistic** during training.
3. Online learning provides a framework about how to exploit opponents through minimising regret.

### Algorithm 1 Double Oracle (McMahan et al., 2003)

```
1: Input: A set  $\Pi, C$  strategy set of players
2:  $\Pi_0, C_0$ : initial set of strategies
3: for  $t = 1$  to  $\infty$  do
4:   if  $\Pi_t \neq \Pi_{t-1}$  or  $C_t \neq C_{t-1}$  then
5:     Solve the NE of the subgame  $G_t$ :
6:      $(\pi_t^*, c_t^*) = \arg \min_{\pi \in \Delta_{\Pi_t}} \arg \max_{c \in \Delta_{C_t}} \pi^\top A c$ 
7:     Find the best response  $a_{t+1}$  and  $c_{t+1}$  to  $(\pi_t^*, c_t^*)$ :
8:        $a_{t+1} = \arg \min_{a \in \Pi} a^\top A c_t^*$ 
9:        $c_{t+1} = \arg \max_{c \in C} \pi_t^{*\top} A c$ 
10:    Update  $\Pi_{t+1} = \Pi_t \cup \{a_{t+1}\}, C_{t+1} = C_t \cup \{c_{t+1}\}$ 
11:   else if  $\Pi_t = \Pi_{t-1}$  and  $C_t = C_{t-1}$  then
12:     Terminate
13:   end if
14: end for
```

### What we want:

if opponents play  $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_T$ , we want the player to have  $\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \dots, \boldsymbol{\pi}_T$  s.t.

$$\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0, \quad R_T = \max_{\pi \in \Delta_{\Pi}} \sum_{t=1}^T (\pi_t^\top A c_t - \pi^\top A c_t)$$

### What we know:

hedge algorithm/multiplicative weight update can achieve no-regret property if one follows the below update

$$\pi_{t+1}(i) = \pi_t(i) \frac{\exp(-\mu_t a^{i\top} A c_t)}{\sum_{i=1}^n \pi_t(i) \exp(-\mu_t a^{i\top} A c_t)}, \forall i \in [n]$$

the regret of MWU is  $\mathcal{O}(\sqrt{\log(n)/T})$

# No-Regret Algorithms: Online Double Oracle

---

**Algorithm 1** Double Oracle (McMahan et al., 2003)

---

- 1: **Input:** A set  $\Pi$ ,  $C$  strategy set of players
  - 2:  $\Pi_0, C_0$ : initial set of strategies
  - 3: **for**  $t = 1$  to  $\infty$  **do**
  - 4:   **if**  $\Pi_t \neq \Pi_{t-1}$  or  $C_t \neq C_{t-1}$  **then**
  - 5:     Solve the NE of the subgame  $G_t$ :  
     $(\pi_t^*, c_t^*) = \arg \min_{\pi \in \Delta_{\Pi_t}} \arg \max_{c \in \Delta_{C_t}} \pi^\top A c$
  - 6:     Find the best response  $a_{t+1}$  and  $c_{t+1}$  to  $(\pi_t^*, c_t^*)$ :  
       $a_{t+1} = \arg \min_{a \in \Pi} a^\top A c_t^*$   
       $c_{t+1} = \arg \max_{c \in C} \pi_t^{*\top} A c$
  - 7:     Update  $\Pi_{t+1} = \Pi_t \cup \{a_{t+1}\}$ ,  $C_{t+1} = C_t \cup \{c_{t+1}\}$
  - 8:   **else if**  $\Pi_t = \Pi_{t-1}$  and  $C_t = C_{t-1}$  **then**
  - 9:     Terminate
  - 10:   **end if**
  - 11: **end for**
- 

---

**Algorithm 2:** Online Single Oracle Algorithm

---

- 1: **Input:** Player's pure strategy set  $\Pi$
  - 2: Init. effective strategies set:  $\Pi_0 = \Pi_1 = \{a^j\}$ ,  $a^j \in \Pi$
  - 3: **for**  $t = 1$  to  $T$  **do**
  - 4:   **if**  $\Pi_t = \Pi_{t-1}$  **then**
  - 5:     Compute  $\pi_t$  by the MWU in Equation (5)
  - 6:   **else if**  $\Pi_t \neq \Pi_{t-1}$  **then**
  - 7:     Start a new time window  $T_{i+1}$  and  
    Reset  $\pi_t = [1/|\Pi_t|, \dots, 1/|\Pi_t|]$ ,  $\bar{l} = \mathbf{0}$
  - 8:   **end if**
  - 9:   Observe  $l_t$  and update the average loss in  $T_i$ :  
     $\bar{l} = \sum_{t \in T_i} l_t / |T_i|$
  - 10:   Calculate the best response:  $a_t = \arg \min_{\pi \in \Pi} \langle \pi, \bar{l} \rangle$
  - 11:   Update the set of strategies:  $\Pi_{t+1} = \Pi_t \cup \{a_t\}$
  - 12: **end for**
  - 13: **Output:**  $\pi_T, \Pi_T$
- 

**Intuition:** maintain a time window  $T_i$  to track opponent's strategy, if no new best response can be found, then keep exploiting, otherwise refresh the time window to catch up with the latest change

# No-Regret Algorithms: Online Double Oracle

1. OSO is a no-regret algorithm.

**Theorem 4 (Regret Bound of OSO).** *Let  $l_1, l_2, \dots, l_T$  be a sequence of loss vectors played by an adversary, and  $\langle \cdot, \cdot \rangle$  be the dot product, OSO in Algorithm 2 is a no-regret algorithm with*

$$\frac{1}{T} \left( \sum_{t=1}^T \langle \pi_t, l_t \rangle - \min_{\pi \in \Pi} \sum_{t=1}^T \langle \pi, l_t \rangle \right) \leq \frac{\sqrt{k \log(k)}}{\sqrt{2T}},$$

where  $k = |\Pi_T|$  is the size of effective strategy set in the final time window.

## Algorithm 2: Online Single Oracle Algorithm

```

1: Input: Player's pure strategy set  $\Pi$ 
2: Init. effective strategies set:  $\Pi_0 = \Pi_1 = \{a^j\}, a^j \in \Pi$ 
3: for  $t = 1$  to  $T$  do
4:   if  $\Pi_t = \Pi_{t-1}$  then
5:     Compute  $\pi_t$  by the MWU in Equation (5)
6:   else if  $\Pi_t \neq \Pi_{t-1}$  then
7:     Start a new time window  $T_{i+1}$  and
       Reset  $\pi_t = [1/|\Pi_t|, \dots, 1/|\Pi_t|], \bar{l} = \mathbf{0}$ 
8:   end if
9:   Observe  $l_t$  and update the average loss in  $T_i$ :
        $\bar{l} = \sum_{t \in T_i} l_t / |T_i|$ 
10:  Calculate the best response:  $a_t = \arg \min_{\pi \in \Pi} \langle \pi, \bar{l} \rangle$ 
11:  Update the set of strategies:  $\Pi_{t+1} = \Pi_t \cup \{a_t\}$ 
12: end for
13: Output:  $\pi_T, \Pi_T$ 

```

2. Putting OSO into self-play settings, we get Online Double Oracle which can solve Nash.

- ◆ Recall that in two-player zero-sum game, if two no-regret methods self play, the outcome will leads to a Nash equilibrium!
- ◆ We just prove that.

## Algorithm 3: Online Double Oracle Algorithm

```

1: Input: Full pure strategy set  $\Pi, C$ 
2: Init. effective strategies set:  $\Pi_0 = \Pi_1, C_0 = C_1$ 
3: for  $t = 1$  to  $T$  do
4:   Each player follows the OSO in Algorithm 2 with
       their respective effective strategy sets  $\Pi_t, C_t$ 
5: end for
6: Output:  $\pi_T, \Pi_T, c_T, C_T$ 

```

**Theorem 5.** *Suppose both players apply OSO. Let  $k_1, k_2$  denote the size of effective strategy set for each player. Then, the average strategies of both players converge to the NE with the rate:*

$$\epsilon_T = \sqrt{\frac{k_1 \log(k_1)}{2T}} + \sqrt{\frac{k_2 \log(k_2)}{2T}}.$$

*In situation where both players follow OSO with Less-Frequent Best Response in Equation (6) and  $\alpha_{t-|\bar{T}_i|}^i = \sqrt{t - |\bar{T}_i|}$ , the convergence rate to NE will be*

$$\epsilon_T = \sqrt{\frac{k_1 \log(k_1)}{2T}} + \sqrt{\frac{k_2 \log(k_2)}{2T}} + \frac{\sqrt{k_1} + \sqrt{k_2}}{\sqrt{T}}.$$

# No-Regret Algorithms: Online Double Oracle

Exploitability on the Spinning Top games

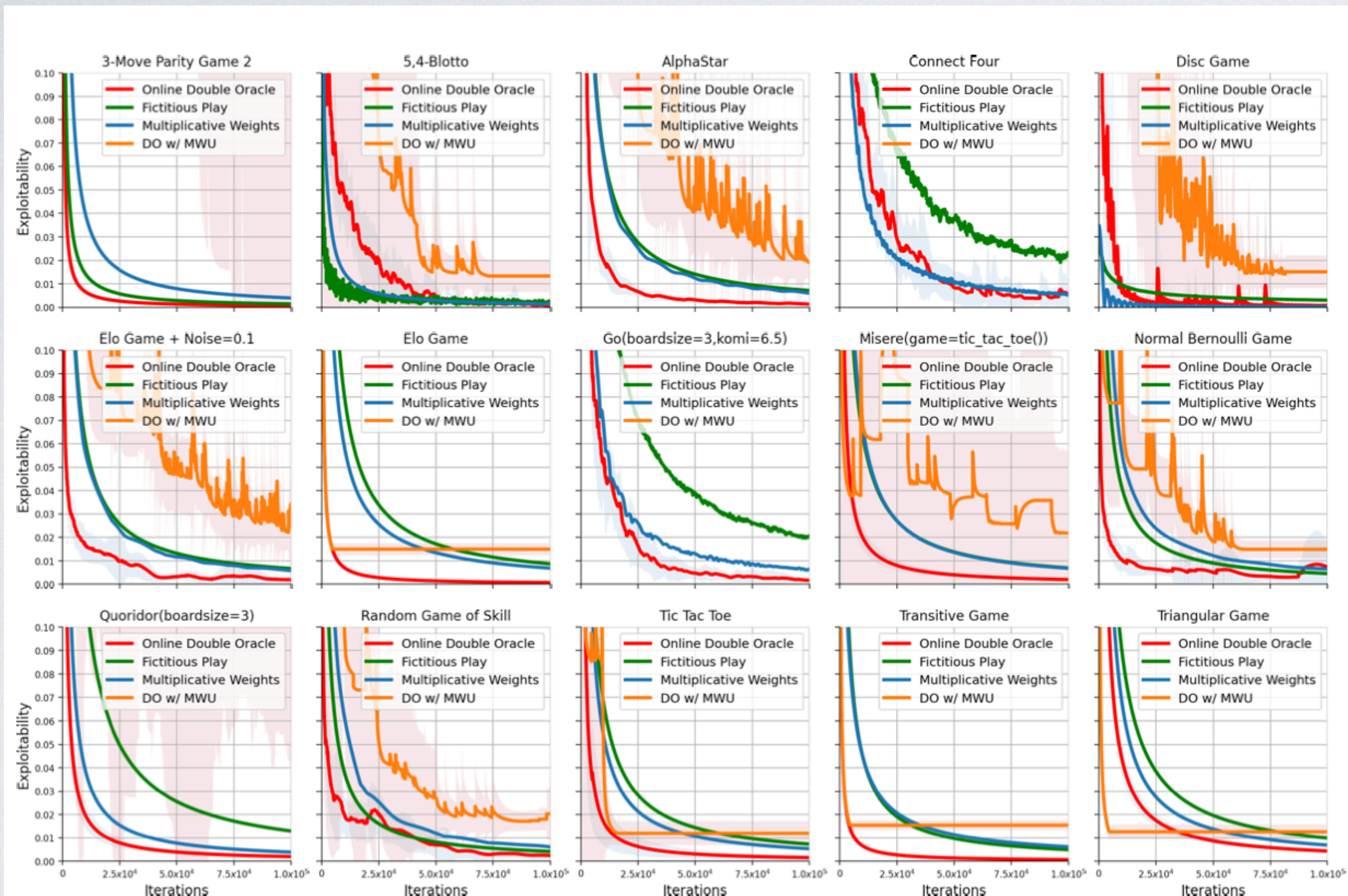


Figure 1: Performance comparisons under self-plays

Exploitability on Poker

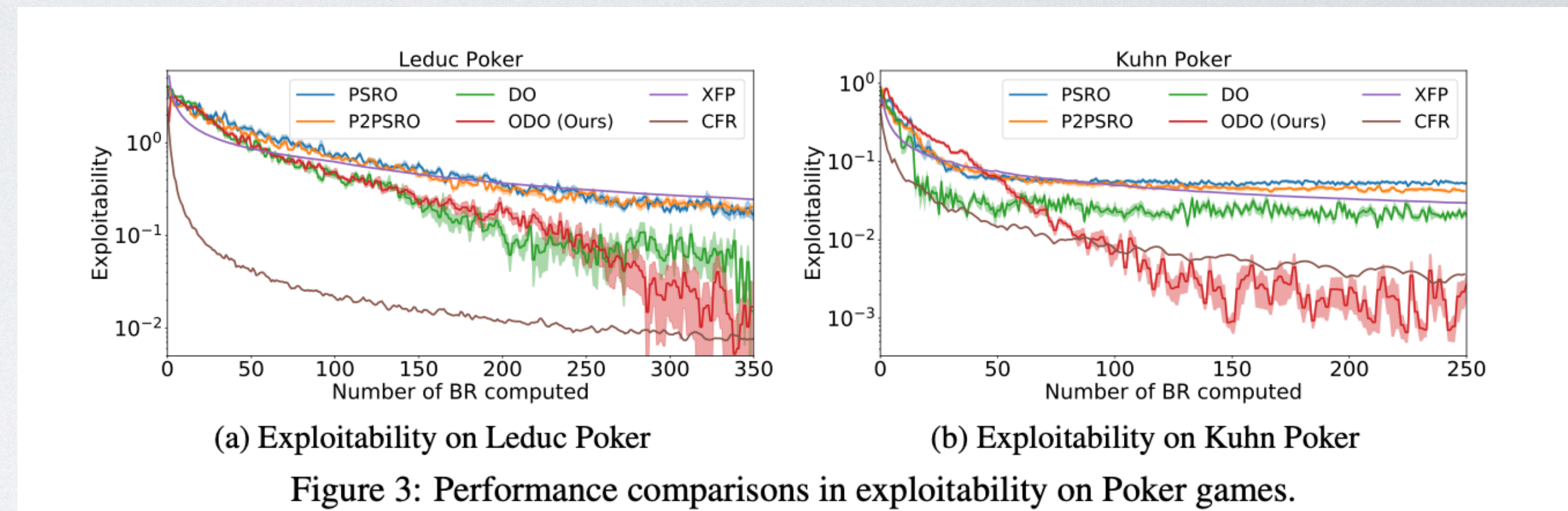
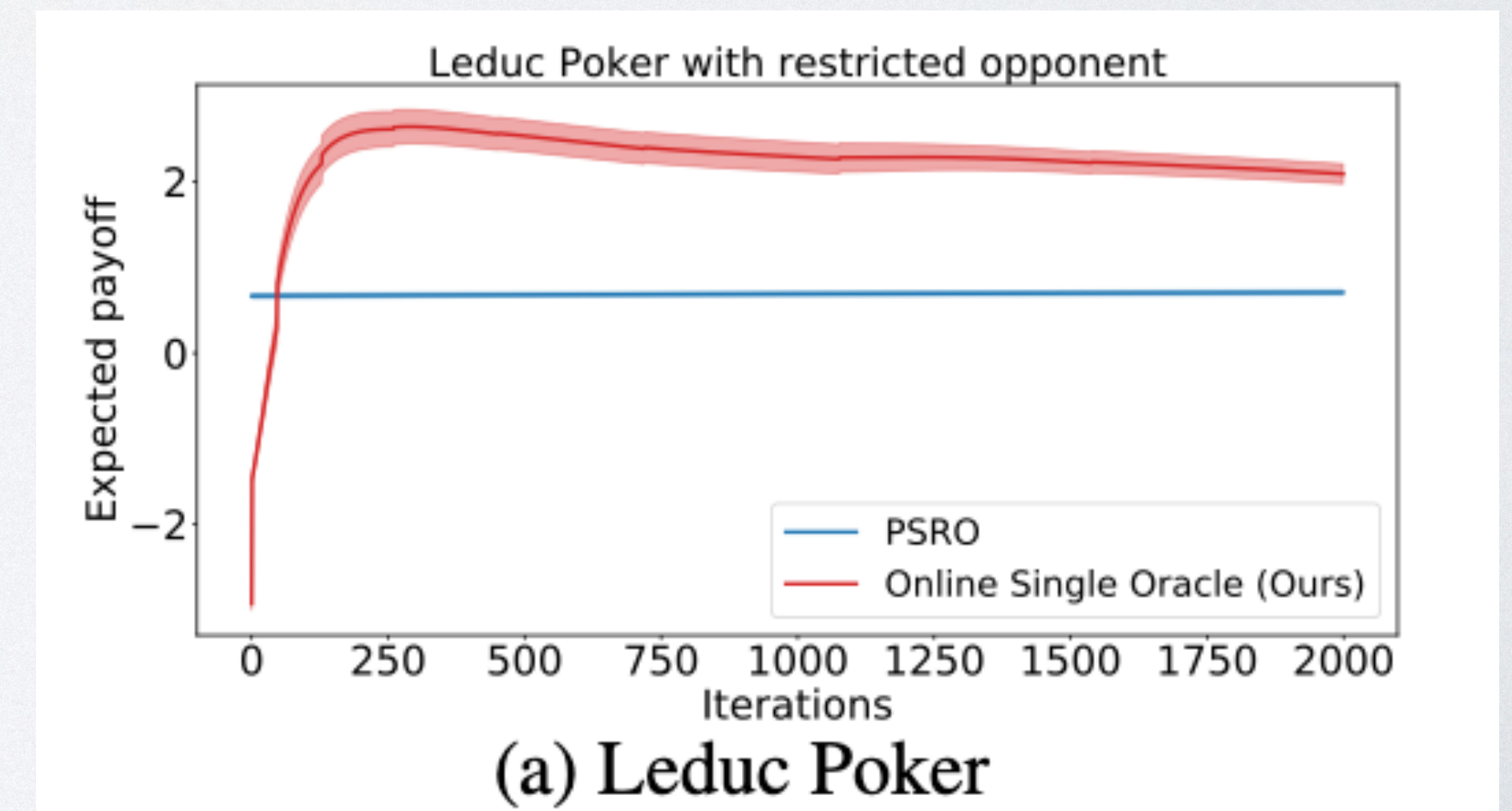


Figure 3: Performance comparisons in exploitability on Poker games.

Play with an imperfect opponent



(a) Leduc Poker

# No-Regret Algorithms for Correlated Equilibrium

- No-(external)-regret players can lead to CCE in multi-player general-sum games, and Nash in two-player zero-sum games.
- The last missing piece: how about correlated equilibrium?

**Definition 5.3.2** Let  $\sigma$  be a distribution over  $S = S_1 \times \dots \times S_k$ . Then  $\sigma$  is a correlated equilibrium if

$$\mathbf{E}_{s \sim \sigma} [c_i(s) | s_i] \leq \mathbf{E}_{s \sim \sigma} [c_i(s_{-i}, s'_i) | s_i]$$

for all  $i \in [k]$  and for all  $s_i, s'_i \in S_i$ .

- We can also define CE through a **switching function (play action b when I plan to play a)**

**Theorem 7.3.2**  $\sigma$  is a correlated equilibrium if and only if for all  $i \in [k]$  and  $\delta : S_i \rightarrow S_i$ ,

$$\mathbf{E}_{s \sim \sigma} [C_i(s)] \leq \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, \delta(s_i))]. \quad (7.3.1)$$

$$\delta(x) = \begin{cases} b & \text{if } x = a \\ x & \text{otherwise} \end{cases}$$

- ♦  $\Rightarrow$ : if a CE, then any switching will incur large cost

- ♦  $\Leftarrow$ : we can rewrite (7.3.1) into  $\sum_{x \in S_i} \left( \Pr [s_i = x] \cdot \mathbf{E}_{s \sim \sigma} [C_i(s) | s_i = x] \right) \leq \sum_{x \in S_i} \left( \Pr [s_i = x] \cdot \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, \delta(s_i)) | s_i = x] \right)$

**both sides only differ when  $x = a$**   $\Pr [s_i = a] \cdot \mathbf{E}_{s \sim \sigma} [C_i(s) | s_i = a] \leq \Pr [s_i = a] \cdot \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, b) | s_i = a]$

$$\implies \mathbf{E}_{s \sim \sigma} [C_i(s) | s_i = a] \leq \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, b) | s_i = a]$$

# No-Regret Algorithms for Correlated Equilibrium

- We can also define CE through a **switching function** (play action b when I plan to play a)

**Theorem 7.3.2**  $\sigma$  is a correlated equilibrium if and only if for all  $i \in [k]$  and  $\delta : S_i \rightarrow S_i$ ,

$$\mathbf{E}_{s \sim \sigma} [C_i(s)] \leq \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, \delta(s_i))]. \quad (7.3.1)$$

- We removed the condition requirement of CE at the price of switching function.
- Best action in hindsight is equivalent to a switching function that maps  $a^t \rightarrow a^*$ .
- The new definition relates to swap-regret:

**Definition 8.2.2** The swap regret of a sequence of actions  $a^1, a^2, \dots, a^T$  with respect to a switching function  $\delta : A \rightarrow A$  is

$$S_T(\delta) = \frac{1}{T} \left( \sum_{t=1}^T c^t(a^t) - \sum_{t=1}^T c^t(\delta(a^t)) \right)$$

best action sequence  
in hindsight

internal regret:  
best swap functions

external regret:  
best action in hindsight

- No-(swap)-regret implies no-(external)-regret. Not the other way round.
- To remember: no internal regret  $\rightarrow$  CE, no external regret  $\rightarrow$  CCE.

# No-Regret Algorithms for Correlated Equilibrium

- For algorithm  $\mathcal{A}$ , the expected swap regret w.r.t a  $\delta : A \rightarrow A$  is defined by

$$\mathbf{E} [S_T^{\mathcal{A}}(\delta)] = \frac{1}{T} \left( \sum_{t=1}^T \mathbf{E}_{a^t \sim p^t} [c^t(a^t)] - \sum_{t=1}^T \mathbf{E}_{a^t \sim p^t} [c^t(\delta(a^t))] \right)$$

- Compare to external-regret  $\mathbf{E} [R_T^{\mathcal{A}}(a)] = \frac{1}{T} \left( \sum_{t=1}^T \mathbf{E}_{a^t \sim p^t} [c^t(a^t)] - \sum_{t=1}^T c^t(a) \right)$
- No-swap-regret means  $\lim_{T \rightarrow \infty} \mathbf{E}[S_T^{\mathcal{A}}(\delta)] = 0, \forall \delta$ .

- Theorem: if all players adopt no-swap-regret algorithms s.t.  $\mathbf{E}[S_T^{\mathcal{A}i}(\delta)] \leq \epsilon, \forall i \in [k], \forall \delta$ , then the average distribution  $\sigma = \sum_{t=1}^T \prod_{i=1}^k \frac{p_i^t}{T}$  is an  $\epsilon$ -CE, i.e.,  $\mathbf{E}_{s \sim \sigma} [C_i(s)] \leq \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, \delta(s_i))] + \epsilon$**

- proof:

$$\begin{aligned} & \mathbf{E}_{s \sim \sigma} [C_i(s)] - \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, \delta(s_i))] \quad \text{Note that } c^t(a) = \mathbf{E}_{s \sim \sigma^t} [C_i(s_{-i}, a)] \\ &= \frac{1}{T} \sum_{t=1}^T \mathbf{E}_{a^t \sim p^t} [c^t(a^t)] - \frac{1}{T} \sum_{t=1}^T \mathbf{E}_{a^t \sim p^t} [c^t(\delta(a^t))] \\ &= \mathbf{E}[S_T^{\mathcal{A}i}(\delta)] \leq \epsilon \end{aligned}$$



# No-Swap-Regret Algorithms for Correlated Equilibrium

- Theorem: if all players adopt no-swap-regret algorithms s.t.  $\mathbf{E}[S_T^{\mathcal{A}_i}(\delta)] \leq \epsilon, \forall i \in [k], \forall \delta$ , then the average distribution  $\sigma = \sum_{t=1}^T \prod_{i=1}^k \frac{p_i^t}{T}$  is an  $\epsilon$ -CE, i.e.,  $\mathbf{E}_{s \sim \sigma} [C_i(s)] \leq \mathbf{E}_{s \sim \sigma} [C_i(s_{-i}, \delta(s_i))] + \epsilon$**
- The goal is then to develop no-swap-regret algorithms like MWU for external regret.
- In turns out we can transform no-external-regret algorithm into no-swap-regret.
  - ◆ The idea: for each action, maintain a no-regret algorithm  $q_i^t$  such as MWU
  - ◆ Let different no-regret algorithm  $(q_1^t, \dots, q_n^t)$  reach a consensus  $p^t$  (see how later).
  - ◆ “Lie” to each no-regret algorithm the loss of  $p^t(i)c^t(i)$  instead of  $c^t(i)$ .
- The goal of no-swap-regret is then written as

$$\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n p^t(i)c^t(i) - \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n p^t(i)c^t(\delta(i)) = o(1)$$

# No-Swap-Regret Algorithms for Correlated Equilibrium

- The goal of no-swap-regret is then written as

$$\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n p^t(i) c^t(i) - \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n p^t(i) c^t(\delta(i)) = o(1)$$

- Because for each action, we maintain a no-regret algorithm, thus we have

$$\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n q_j^t(i) \boxed{p^t(j) c^t(i)} - \frac{1}{T} \sum_{t=1}^T p^t(j) c^t(\delta(j)) \boxed{\leq R_T^j(\delta(j))} = o(1)$$

the cost that is lied
playing  $\delta(j)$  in hindsight and internal regret is always smaller than external regret
this is because no regret w.r.t  $\delta(j)$

- Sum over  $n = |A|$  will not influence the total regret, since no  $T$  term is involved.

$$\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^n \boxed{\sum_{j=1}^n q_j^t(i) (p^t(j) c^t(i))} - \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^n p^t(j) c^t(\delta(j)) \leq \sum_{j=1}^n R_T^j(\delta(j)) = o(1)$$

- Recall  $q_j^t(i)$  is the probability that the  $j$ -th no-regret algorithm pick the  $i$ -th action. If we

make  $p^t(i) = \sum_{j=1}^n q_j^t(i) p^t(j)$ , then we use the above equation prove no-swap-regret. DONE!

*Thanks!*