

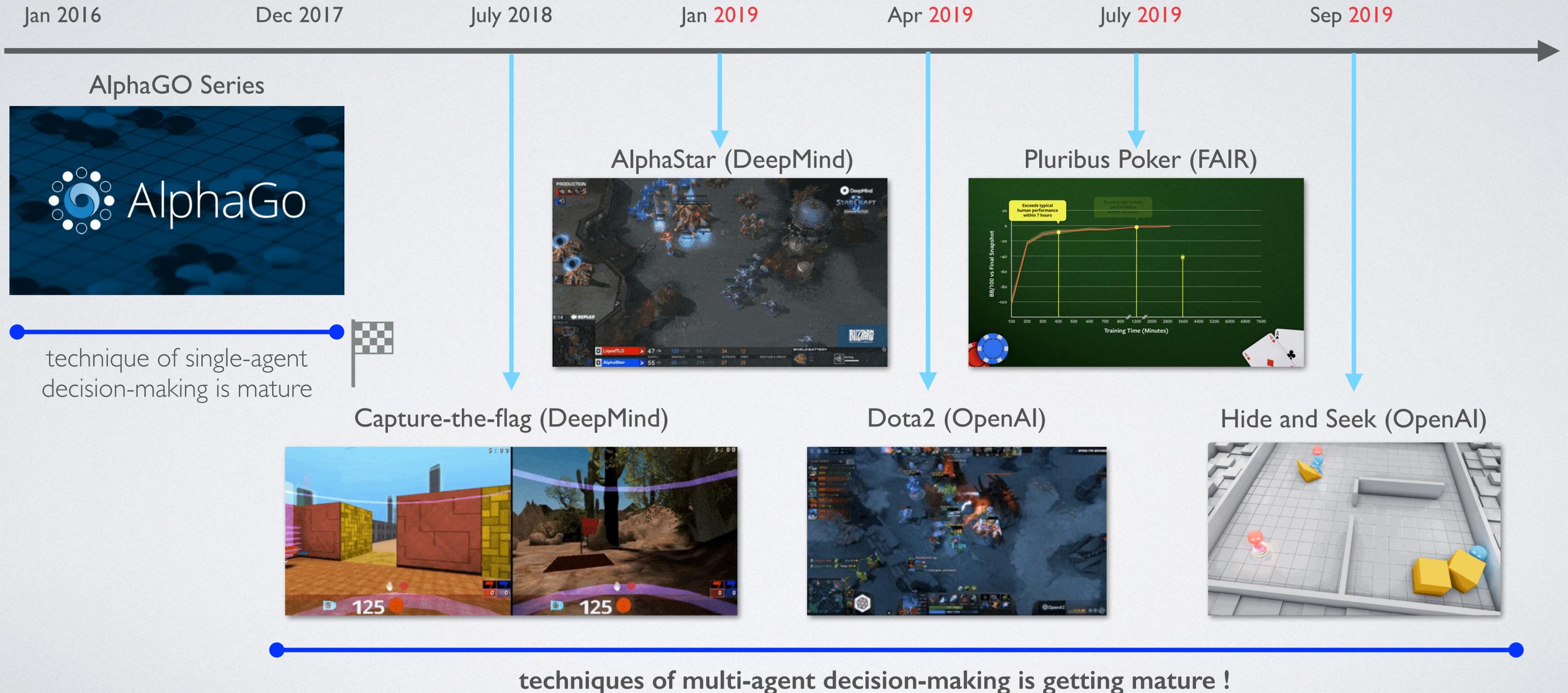


# DEALING WITH NON-TRANSITIVITY IN TWO-PLAYER ZERO-SUM GAMES

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

# Multi-Agent Reinforcement Learning in Games

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



# A General Solver to Two-Player Zero-Sum Games

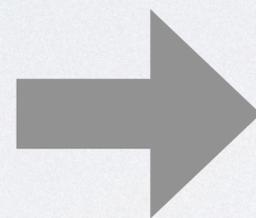
Output: the reward  $(R^1, \dots, R^N)$



Black-box multi-agent  
game engine



input



Our algorithm:



output

Low-exploitability  
strategy  
 $(\pi^{1,*}, \dots, \pi^{N,*})$

$$\mathbf{Br}^i(\pi^{-i}) = \arg \max_{\pi^i} \mathbf{E}_{a^i \sim \pi^i, a^{-i} \sim \pi^{-i}} [R^i(a^i, a^{-i})]$$

$$\text{Exploitability}(\pi) = \sum_{i=1}^2 R^i(\mathbf{Br}^i(\pi^{-i}), \pi^{-i}) - R^i(\pi)$$

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



# Not All Games Can be Solved by Tree Search

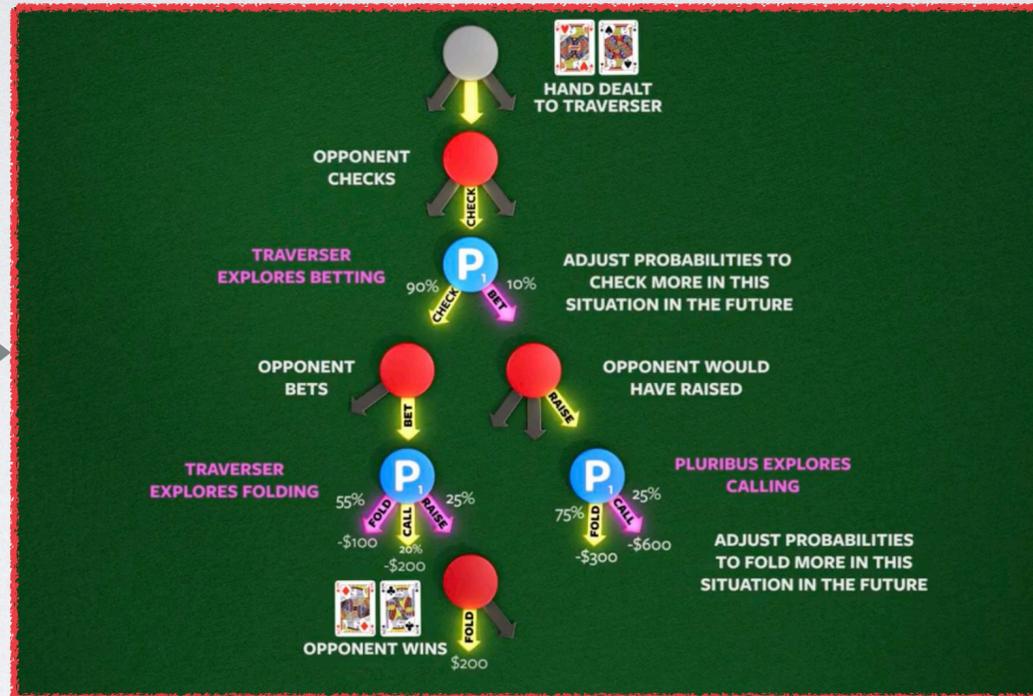
Output: the reward  $(R^1, \dots, R^N)$



**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

Planning is not always feasible. StarCraft has  $10^{26}$  choices per step (vs. the game tree size of chess  $10^{50}$ , Texas holdem  $10^{80}$ , GO  $10^{170}$ )

Enumerating all policies' actions at each state and then playing a randomise best response is infeasible (i.e. RPS can not apply)

Solution: design a game of game — meta-game, the problem problem, auto-curricula.

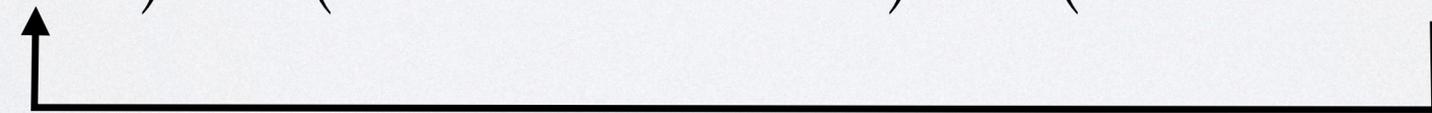
# Naive Self-play Will Not Work

Question: Can we use it as a general framework to solve any games?

**PPO + PBT + Self-play = Panacea ?**

## Algorithm 2 Self-play

```
input: agent  $v_1$   
for  $t = 1, \dots, T$  do  
     $v_{t+1} \leftarrow \text{oracle}(v_t, \phi_{v_t}(\bullet))$   
end for  
output:  $v_{T+1}$ 
```

$$(\pi^1, \pi^2) \rightarrow (\pi^1, \pi^{2,*} = \mathbf{Br}(\pi^1)) \rightarrow (\pi^{1,*} = \mathbf{Br}(\pi^{2,*}), \pi^{2,*})$$


**It depends. In most of the games, it does not work.**



self-plays

# Naive Self-play Will Not Work

- It is because of **Non-Transitivity**

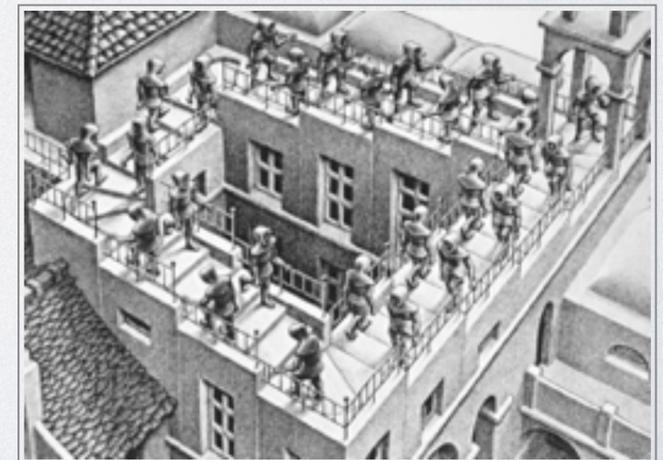
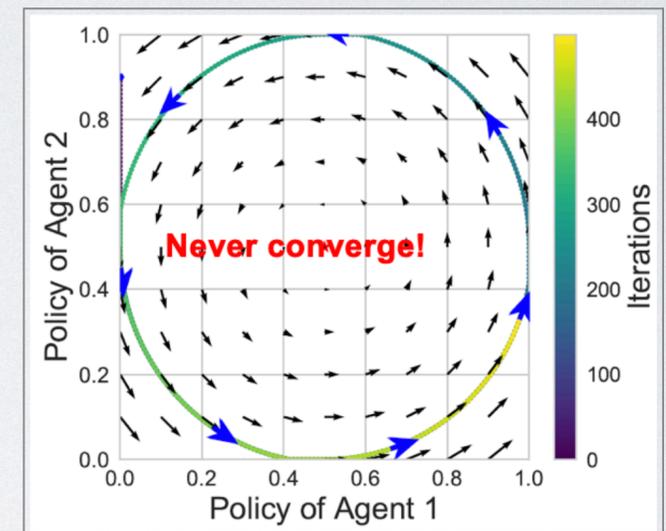
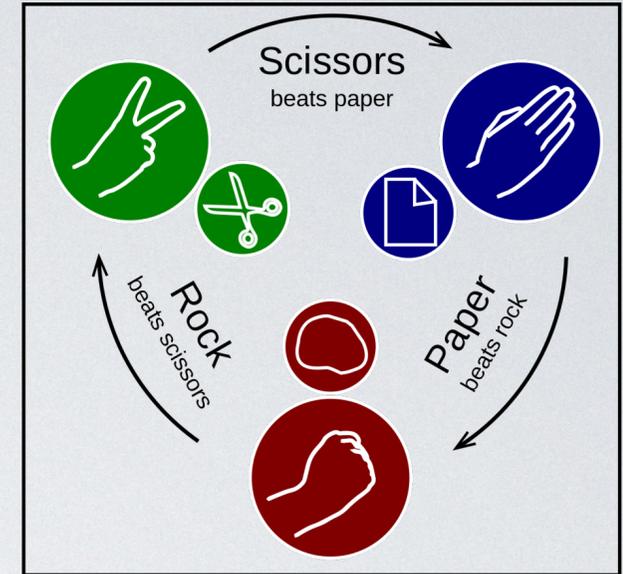
$$\int_W \phi(\mathbf{v}, \mathbf{w}) \cdot d\mathbf{w} = 0, \quad \forall \mathbf{v} \in W$$

- Rock-Paper-Scissor game:

$$\begin{bmatrix} 0 & 1 & -1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix}$$

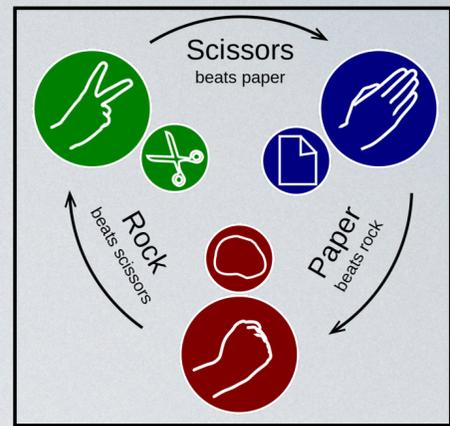
- Disc game:

$$\phi(\mathbf{v}, \mathbf{w}) = \mathbf{v}^\top \cdot \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \cdot \mathbf{w} = v_1 w_2 - v_2 w_1$$



# What is Non-Transitivity ?

- Every FFG can be decomposed into two parts



$$\text{FFG} = \text{Transitive game} \oplus \text{Non-transitive game}$$

- **Non-transitive Game:** the rules of winning are not-transitive across players.

$$v_t \text{ beats } v_{t-1}, \quad v_{t+1} \text{ beats } v_t \not\Rightarrow v_{t+1} \text{ beats } v_{t-1}$$

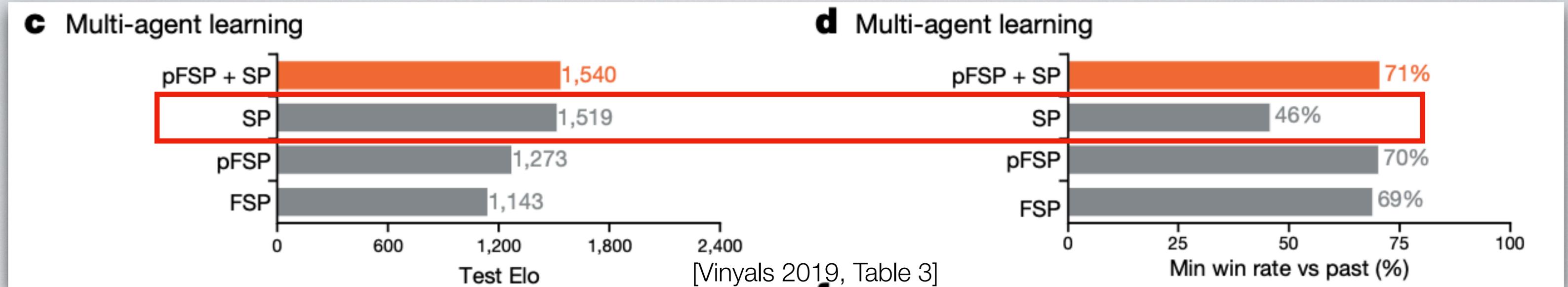
- Mutual dominance across different types of modules in a game. This is commonly observed in modern MOBA games.



- For this types of game, self-play is not helpful at all because transitivity assumption does not hold. **Self-play will lead to cyclic loops forever.**

# Non-Transitivity Harms Training !

## Example on training AlphaStar:



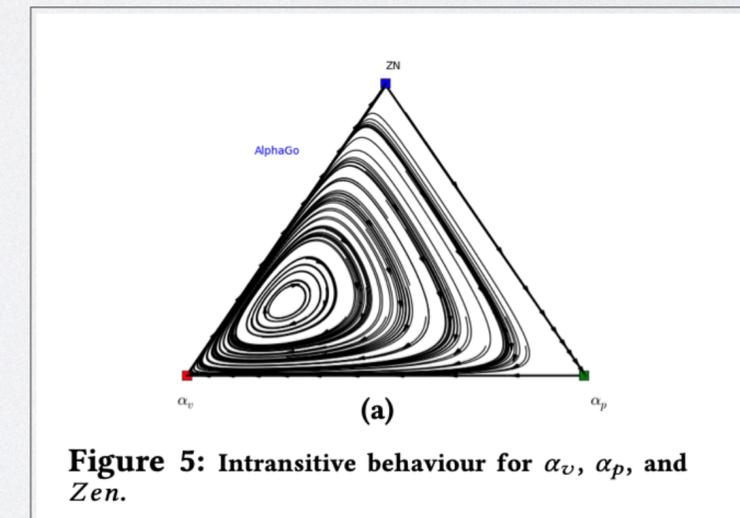
## Example on training Soccer AI:

Table 2: Average goal difference  $\pm$  one standard deviation across 5 repetitions of the experiment.

$A$ vs built-in AI	$4.25 \pm 1.72$
$B$ vs $A$	$11.93 \pm 2.19$
$B$ vs built-in AI	$-0.27 \pm 0.33$

[Karol 2020, table 2]

## Example on training AlphaGO:



[Silver 2016, table 9]

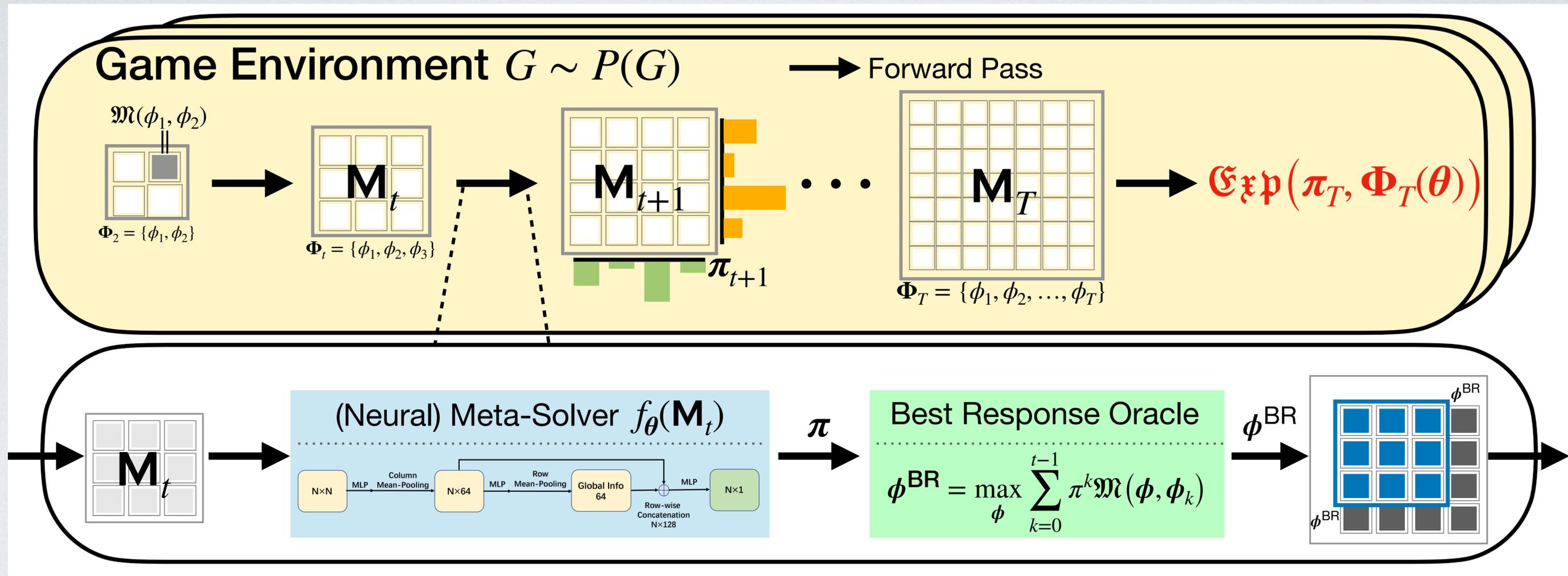
# Policy Space Response Oracle: To Build Populations of RL Agents

- A generalisation of double oracle methods on **meta-games**, with the best responder implemented through **deep RL algorithms**.
- 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-solver
- 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**:  $\sigma_{-i} = \mathbf{Nash}(\Pi^{T-1}, U)$  or  $\mathbf{RD}(\Pi^{T-1}, U)$

## 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 Incorporate Many Variants

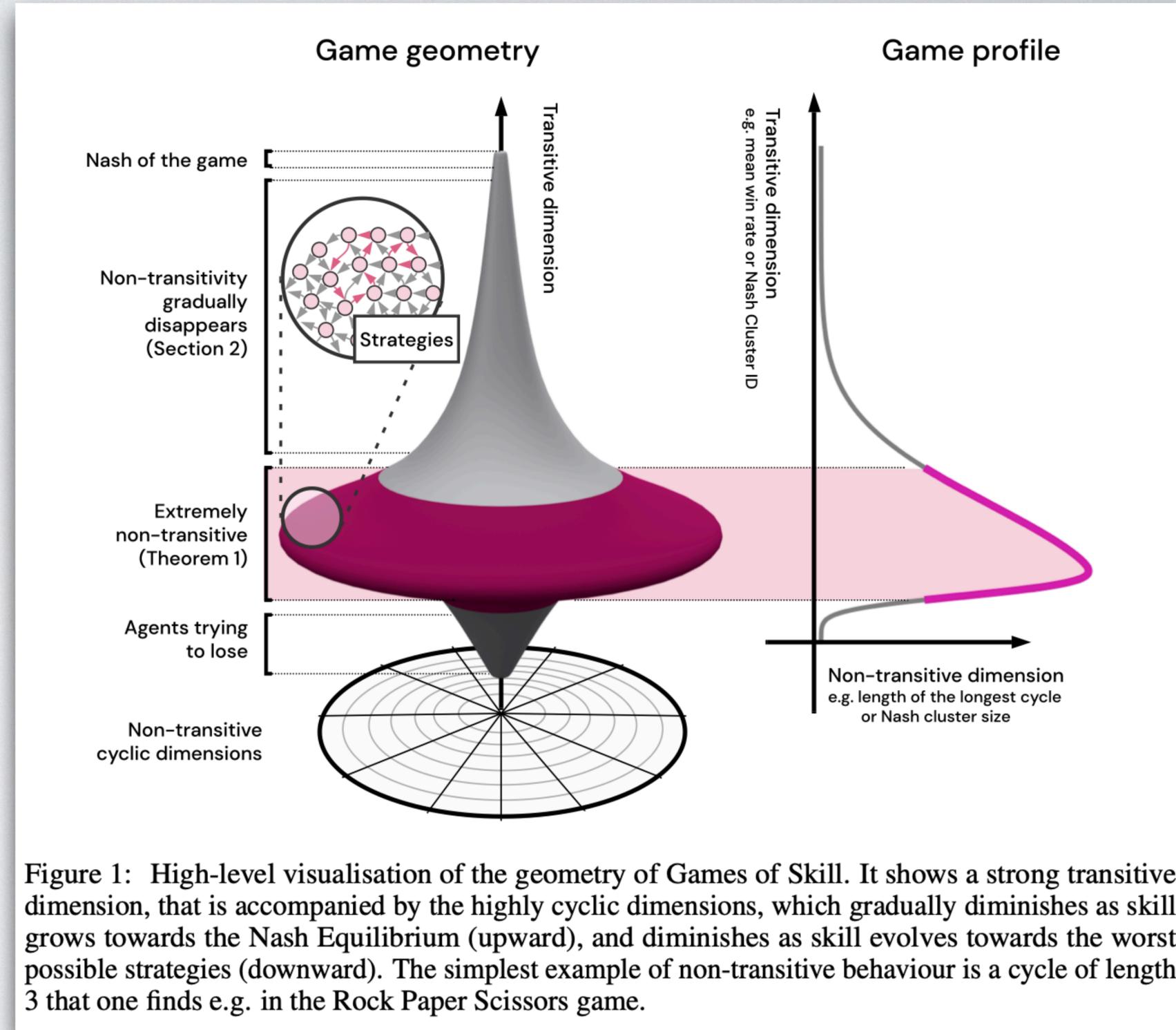


Elo rating  
 Correlated equilibrium  
 Nash equilibrium  
 Replicator dynamics  
 $\alpha$ -Rank/ $\alpha^\alpha$ -Rank

iterated best response  
 fictitious play  
 double oracle/PSRO  
 PSRO-Nash  
 $\alpha$ -PSRO  
 JPSRO

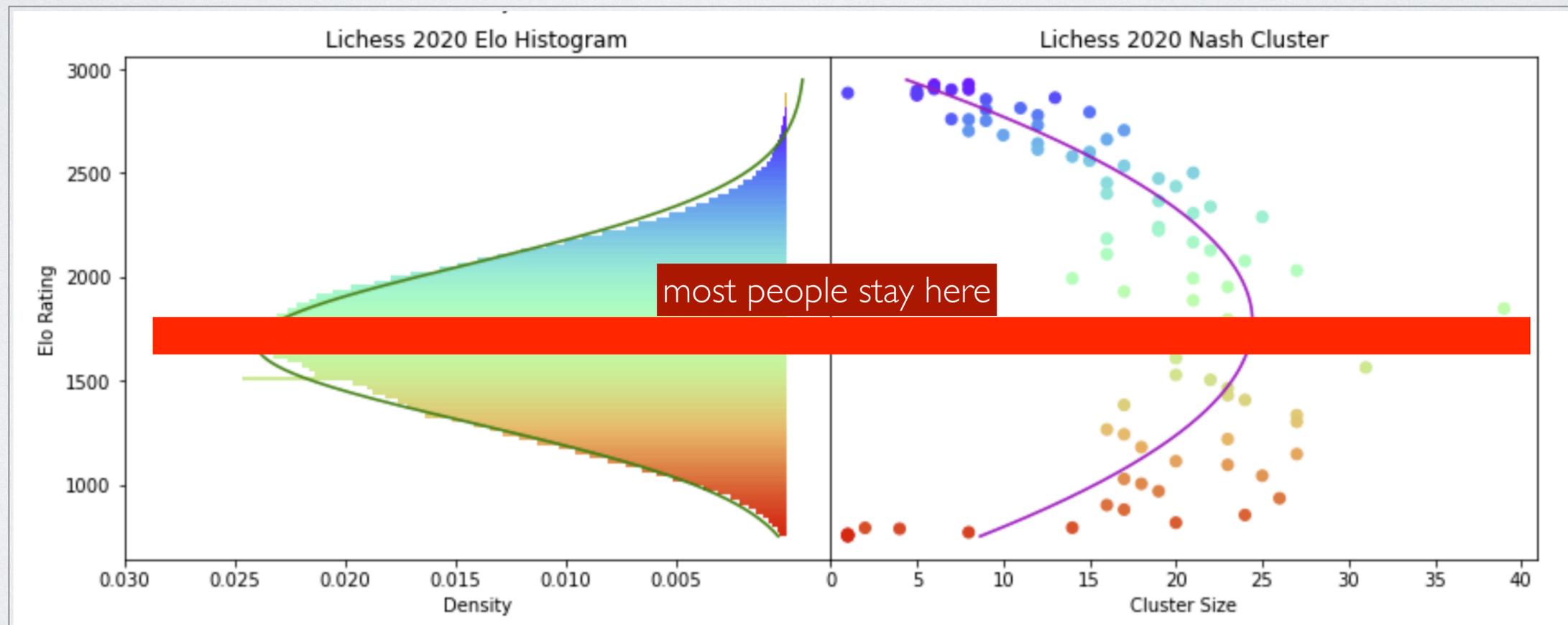
# The Spinning Top Hypothesis

- Real-world games are mixtures of both transitive and in-transitive components, e.g., Go, DOTA, StarCraft II.
- Though winning is often harder than losing a game, finding a strategy that always loses is also challenging.
- Players who regularly practice start to beat less skilled players, this corresponds to the transitive dynamics.
- At certain level (the red part), players will start to find many different strategy styles. Despite not providing a universal advantage against all opponents, players will counter each other within the same transitive group. This provide direct information of improvement.
- As players get stronger to the highest level, seeing many strategy styles, the outcome relies mostly on skill and less on one particular game styles (以不变应万变).



# Measuring the Non-Transitivity

- Real-world data set from human players on Chess
  - ◆ previous results are based on AI, now we study 1000 human players from Lichess
  - ◆ Chess presents the same spinning top pattern, which verifies the hypothesis

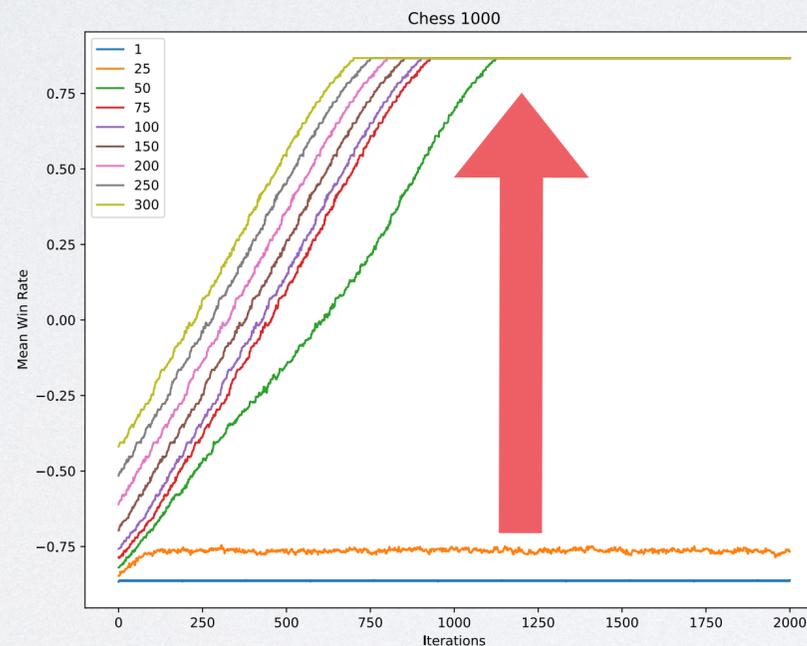
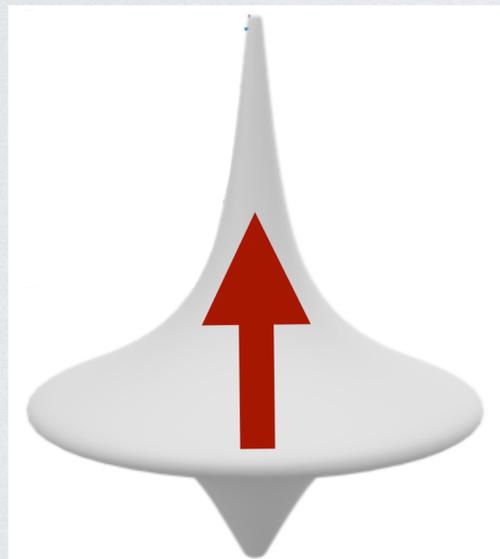


[Ricky Sanjaya]

# Why Modelling Diversity is Critical ?

- Diversity matters because **the more diverse** the strategy pool, **the less un-exploitable**. Promoting diversity can help you walk out of the in-transitive region faster.

**Theorem 3.** *If at any point in time, the training population  $\mathcal{P}^t$  includes any full Nash cluster  $C_i \subset \mathcal{P}^t$ , then training against  $\mathcal{P}^t$  by finding  $\pi$  such that  $\forall \pi_j \in \mathcal{P}^t \mathbf{f}(\pi, \pi_j) > 0$  guarantees transitive improvement in terms of the Nash clustering  $\exists_{k < i} \pi \in C_k$ .*



## Diverse Auto-Curriculum is Critical for Successful Real-World Multiagent Learning Systems\*

Blue Sky Ideas Track

Yaodong Yang<sup>†</sup>  
University College London  
Huawei R&D U.K.

Jun Luo  
Huawei Canada

Ying Wen  
Shanghai Jiao Tong University

Oliver Slumbers  
University College London

Daniel Graves  
Huawei Canada

Haitham Bou Ammar  
Huawei R&D U.K.

Jun Wang  
University College London  
Huawei R&D U.K.

Matthew E. Taylor  
University of Alberta  
Alberta Machine Intelligence Institute

- In real-world applications, you want policies to be diverse enough, covering different skill levels. This is a **realistic need** from **autonomous driving** and **gaming AI** applications.

# Recent Advance: Diverse PSRO

1. Diversity should include both response diversity (in terms of reward), and behavioural diversity (in terms of policy occupancy measure)
2. We want both the **outcomes** and the **policies that lead to those outcomes** to be diverse.

## Unifying Behavioral and Response Diversity for Open-ended Learning in Zero-sum Games

Xiangyu Liu<sup>1</sup>, Hangtian Jia<sup>2</sup>, Ying Wen<sup>1\*</sup>, Yaodong Yang<sup>3</sup>, Yujing Hu<sup>2</sup>,  
Yingfeng Chen<sup>2</sup>, Changjie Fan<sup>2</sup> and Zhipeng Hu<sup>2</sup>

<sup>1</sup>Shanghai Jiao Tong University, <sup>2</sup>Netease Fuxi AI Lab, <sup>3</sup>University College London

Method	Tool for Diversity	BD	RD	Game Type
DvD	Determinant	✓	×	Single-agent
PSRO <sub>N</sub>	None	×	×	n-player general-sum game
PSRO <sub>rN</sub>	$L_{1,1}$ norm	×	✓	2-player zero-sum game
DPP-PSRO	Determinantal point process	×	✓	2-player general-sum game
Our Methods	Occupancy measure & convex hull	✓	✓	n-player general-sum game

# Recent Advance: Diverse PSRO

1. **behavioural diversity**: assuming existing population of policy mixed by Nash distribution is  $\pi_E = (\pi_i, \pi_{E-i})$ , we want a new policy  $\pi^{M+1}$  that has a different occupancy measure

$$\rho_{\pi}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s \mid \pi) \text{ from } \pi_E:$$

$$\text{Div}_{\text{occ}}(\pi_i^{M+1}) = D_f(\rho_{\pi_i^{M+1}, \pi_{E-i}} \parallel \rho_{\pi_i, \pi_{E-i}})$$

2. in practice, one can train a neural network  $f_{\hat{\theta}}$  to fit  $(s, \mathbf{a}) \sim \rho_{\pi_E}$ , and then assign an intrinsic reward by encouraging the new policy to visit state-action pairs with large prediction error (not covered by the existing occupancy measure).

$$\max R^{\text{int}}(s, a) = \left\| f_{\hat{\theta}}(s, \mathbf{a}) - f_{\theta}(s, \mathbf{a}) \right\|^2$$

# Recent Advance: Diverse PSRO

1. **response diversity**: we want the new policy  $\pi^{M+1}$  to expand the convex hull of the existing meta-game  $A_M$  by having the new payoff vector  $\mathbf{a}_{M+1} := \left[ \phi_i(\pi_i^{M+1}, \pi_{-i}^j) \right]_{j=1}^N$  that

$$\text{Div}_{\text{rew}}(\pi_i^{M+1}) = \min_{\substack{\mathbf{1}^\top \boldsymbol{\beta} = 1 \\ \boldsymbol{\beta} \geq 0}} \left\| \mathbf{A}_M^\top \boldsymbol{\beta} - \mathbf{a}_{M+1} \right\|_2^2$$

2. the above equation has no close form, but we can optimise a lower bound

$$\text{Div}_{\text{rew}}(\pi_i^{M+1}) \geq \mathbf{F}(\pi_i^{M+1}) = \frac{\sigma_{\min}^2(\mathbf{A}) \left( 1 - \mathbf{1}^\top (\mathbf{A}^\top)^\dagger \mathbf{a}_{n+1} \right)^2}{M} + \left\| \left( \mathbf{I} - \mathbf{A}^\top (\mathbf{A}^\top)^\dagger \right) \mathbf{a}_{n+1} \right\|_2^2$$

3. **chicken-egg problem**: how can we know the payoff  $\mathbf{a}_{M+1}$  before we train the policy ?

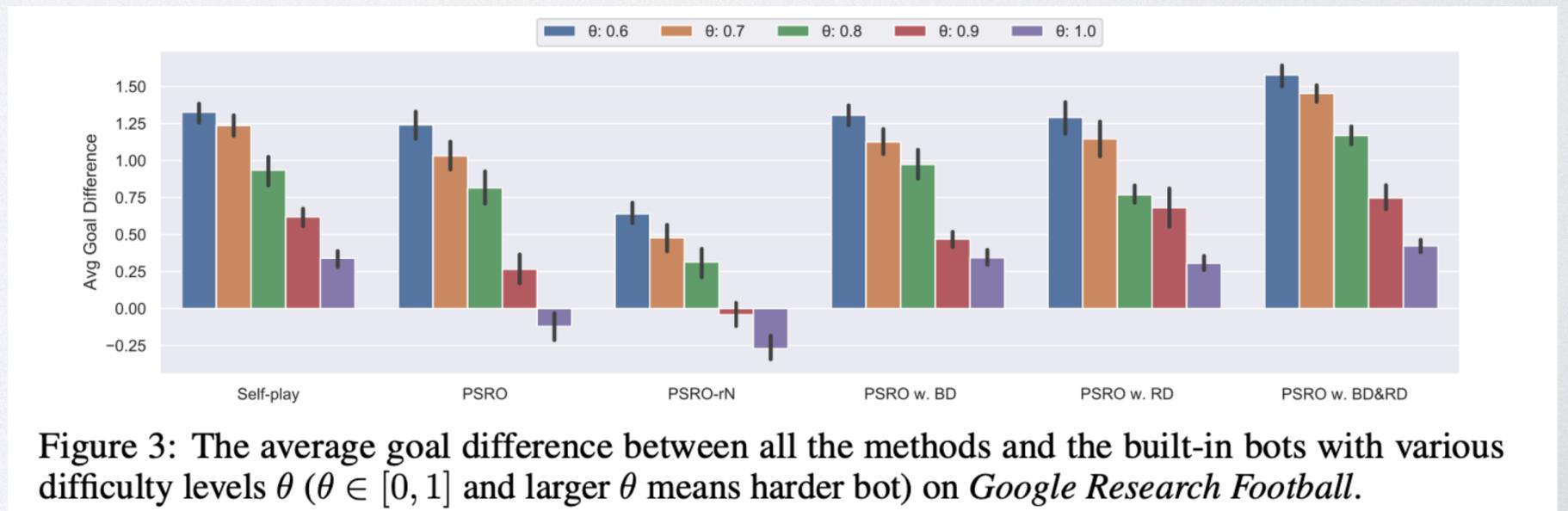
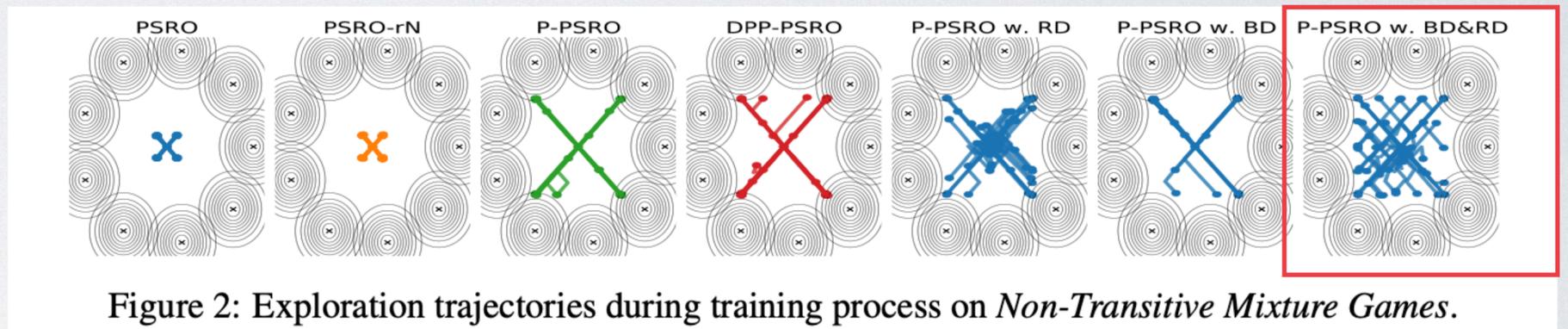
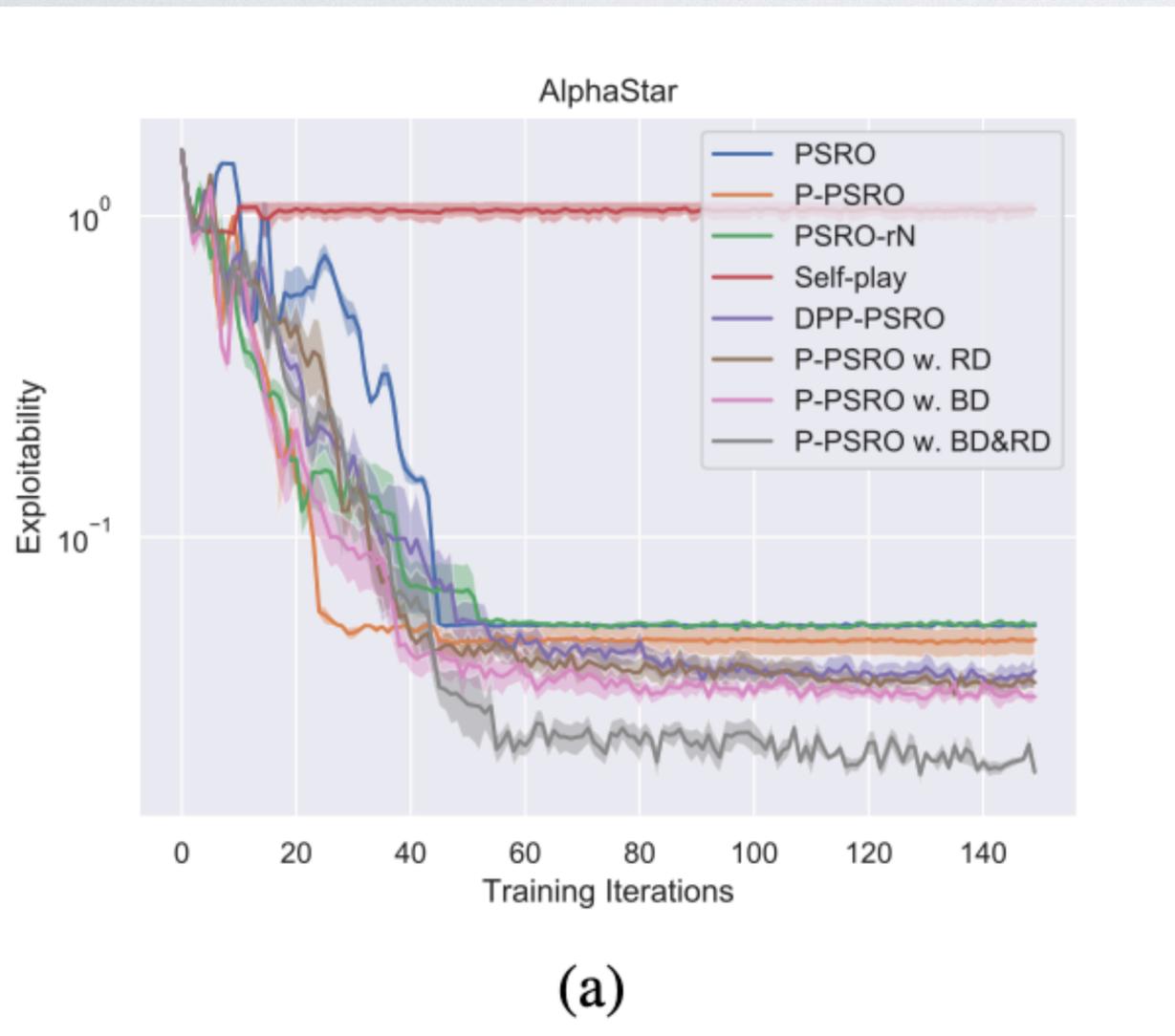
$$\frac{\partial F(\pi_i'(\theta))}{\partial \theta} = \left( \frac{\partial \phi_i(\pi_i'(\theta), \pi_{-i}^1)}{\partial \theta}, \dots, \frac{\partial \phi_i(\pi_i'(\theta), \pi_{-i}^M)}{\partial \theta} \right) \frac{\partial F}{\partial \mathbf{a}_{M+1}}$$

the answer: we can train against  $\pi_{-i}^M$  based on the weights suggested by  $\partial F / \partial \mathbf{a}_{M+1}$  !

# Recent Advance: Diverse PSRO

I. considering both diversity terms in the PSRO process

$$\arg \max_{\pi'_i} \mathbb{E}_{s, \mathbf{a} \sim \rho_{\pi'_i, \pi_{E-i}}} [r(s, \mathbf{a})] + \lambda_1 \text{Div}_{\text{occ}}(\pi'_i) + \lambda_2 \text{Div}_{\text{rew}}(\pi'_i)$$



# Diverse Behaviours Learned on Google Football

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



make offside



push and run

# Recent Advance: Auto-PSRO

## Discovering Multi-Agent Auto-Curricula in Two-Player Zero-Sum Games

Xidong Feng<sup>\*1</sup>, Oliver Slumbers<sup>\*1</sup>, Yaodong Yang<sup>†1</sup>,

Ziyu Wan<sup>2</sup>, Bo Liu<sup>3</sup>, Stephen McAleer<sup>4</sup>, Ying Wen<sup>2</sup>, Jun Wang<sup>1</sup>

1. Learning to learn: to discover multi-agent algorithms (“who to beat” and “how to beat them”) from data.
2. Maybe **game theoretical knowledge** (transitivity/non-transitivity/Nash) are not necessarily needed, the solution algorithm can be learned purely from data.
3. The idea is to **learn how to build an auto-curricula** based on the type of game provided to the meta-learning algorithm, rather than what the auto-curricula should be (e.g. PSRO/DO).
4. Why it will work better than DO/PSRO: because RL oracle can only approximate best response, and using Nash, though theoretically guaranteed, may not be the best option for a solver.
5. On single-agent RL, the discovered RL methods are proved to outperform TD learning designed by humans.

### Discovering Reinforcement Learning Algorithms

Junhyuk Oh   Matteo Hessel   Wojciech M. Czarnecki   Zhongwen Xu  
 Hado van Hasselt   Satinder Singh   David Silver  
 DeepMind

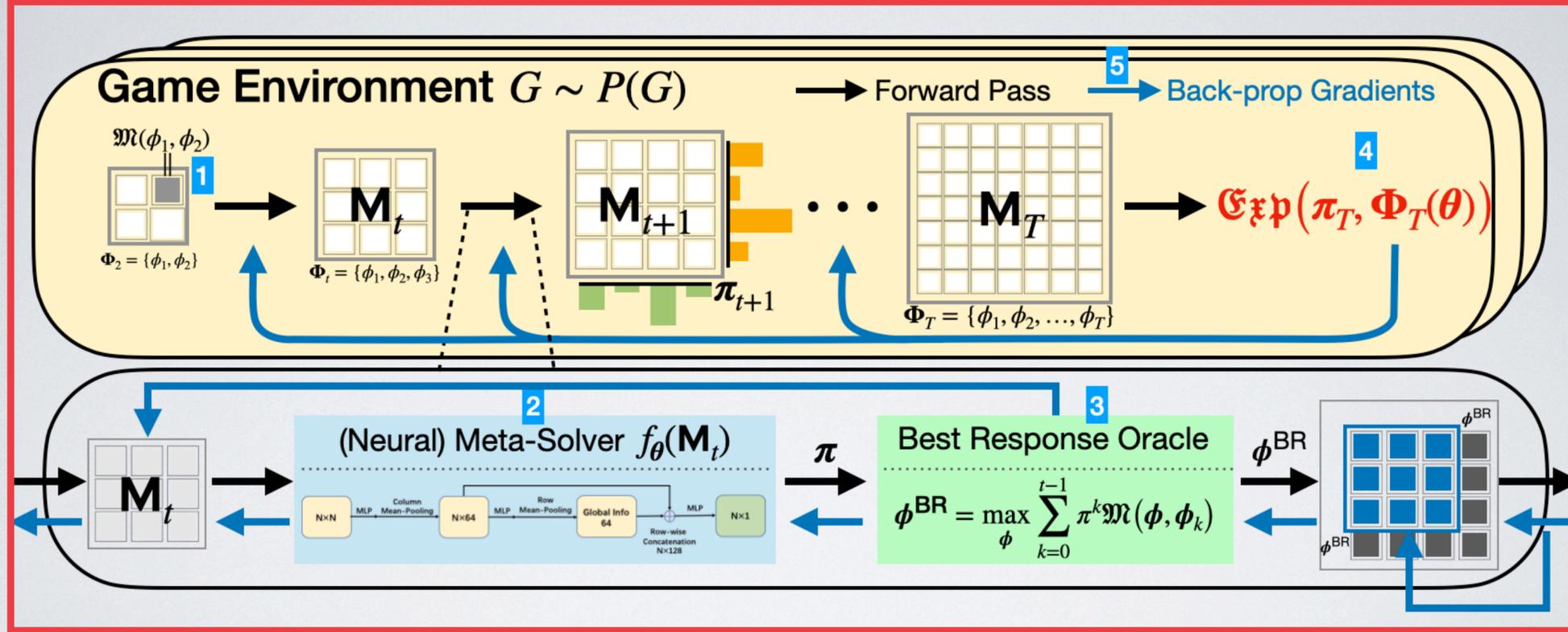
### Meta-Gradient Reinforcement Learning with an Objective Discovered Online

Zhongwen Xu, Hado van Hasselt, Matteo Hessel  
 Junhyuk Oh, Satinder Singh, David Silver  
 DeepMind  
 {zhongwen,hado,mthss,junhyuk,baveja,davidsilver}@google.com

Algorithm	Algorithm properties	What is meta-learned?
IDBD, SMD [30, 27]	† □ →	learning rate
SGD <sup>2</sup> [1]	††† ■ ←	optimiser
RL <sup>2</sup> , Meta-RL [9, 39]	††† ■ X	recurrent network
MAML, REPTILE [11, 23]	††† □ ←	initial params
Meta-Gradient [43, 46]	† □ →	$\gamma$ , $\lambda$ , reward
Meta-Gradient [38, 44, 40]	† □ ←	auxiliary tasks, hyperparams, reward weights
ML <sup>3</sup> , MetaGenRL [2, 19]	††† ■ ←	loss function
Evolved PG [16]	††† ■ X	loss function
Oh et al. 2020 [24]	††† ■ ←	target vector
This paper	† ■ ←	target

□ white box, ■ black box, † single lifetime, ††† multi-lifetime  
 ← backward mode, → forward mode, X no meta-gradient

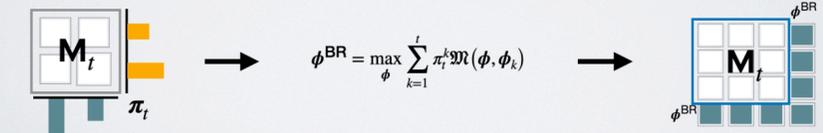
# Recent Advance: Auto-PSRO Framework



Let's differentiate the whole process using back propagation!

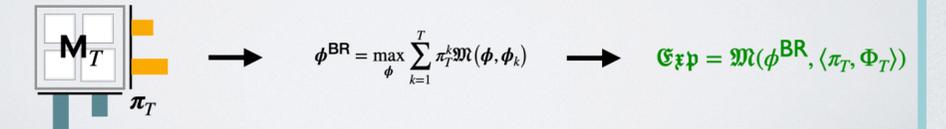
## 3 The Best-Response Oracle

- Algorithm component that controls the iterative expansion of the population
- Given a curriculum  $\pi_t \in \Delta_{|\Phi_t|}$  the goal becomes to solve a best-response to this distribution
- Goal is the following:  $\phi_t^{BR} = \operatorname{argmax}_{\phi} \sum_{k=1}^t \pi_t^k \mathfrak{M}(\phi, \phi_k)$
- Perform the optimisation in anyway desired, but this will impact the meta-gradient calculation



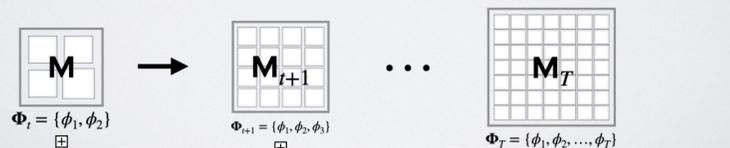
## 4 The Learning Objective

- What is the goal of the iterative update procedure?
- Given a curriculum  $\pi_T = f_\theta(M_T)$  and a population  $\Phi_T$  we want to be as close to a Nash equilibrium as possible.
- Distance to Nash measured as the exploitability:  $\mathfrak{Exp} := \max_{\phi} \mathfrak{M}(\phi, \langle \pi_T, \Phi_T \rangle)$
- i.e. How good is the best-response to the curriculum? If 0, it is a Nash equilibrium



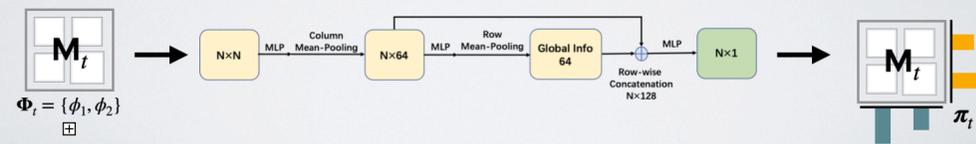
## 1 The Meta-Game

- Main component of population-based methods - **The meta-game**
- An agent is a mapping  $\phi : S \times A \rightarrow [0,1]$
- The payoff for agent  $i$  vs. agent  $j$  is defined as  $\mathfrak{M}(\phi_i, \phi_j)$
- Payoff matrix between agents in a population amenable to GT analysis
- The goal of these algorithms is to expand the populations  $\Phi$  iteratively



## 2 The Meta-Solver

- Algorithm component that controls the auto-curricula of *who to compete with*
- General examples: Nash equilibrium, Uniform distribution, Last agent
- Need to parameterise the process so that we can learn it
- A network with parameters  $\theta$  maps  $f_\theta : M_t \rightarrow [0,1]^t$  so that  $\pi_t = f_\theta(M_t)$



## 5 Optimisation through meta-gradients

- Recall the learning objective of the player:  $\mathfrak{Exp} := \max_{\phi} \mathfrak{M}(\phi, \langle \pi_T, \Phi_T \rangle)$
- Also recall that  $\pi_T = f_\theta(M_T)$ , which allows us to define the meta-solver optimisation as:

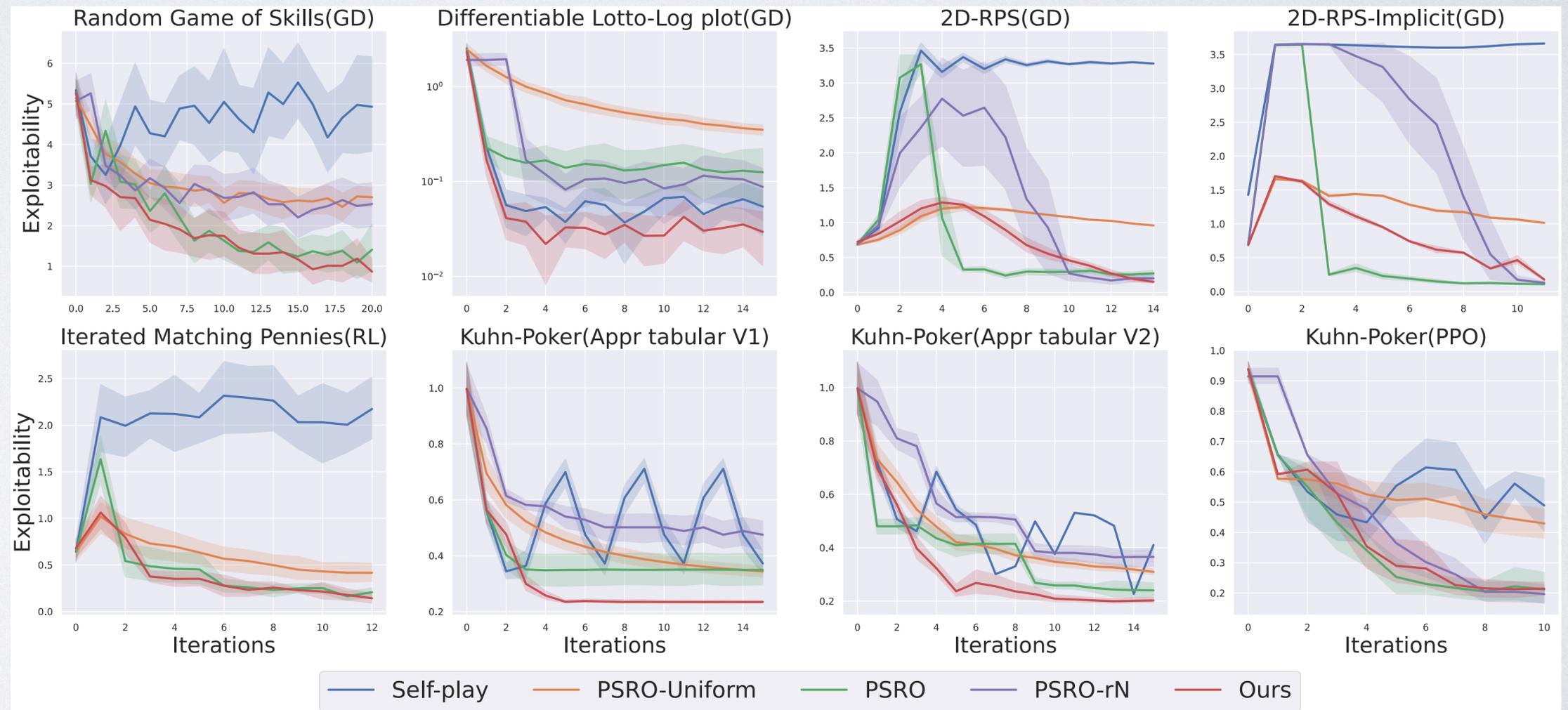
$$\nabla_{\theta} J(\theta) = \mathbb{E}_G \left[ \frac{\partial \mathfrak{M}_{T+1}}{\partial \phi_{T+1}^{BR}} \frac{\partial \phi_{T+1}^{BR}}{\partial \theta} + \frac{\partial \mathfrak{M}_{T+1}}{\partial \pi_T} \frac{\partial \pi_T}{\partial \theta} + \frac{\partial \mathfrak{M}_{T+1}}{\partial \Phi_T} \frac{\partial \Phi_T}{\partial \theta} \right]$$

Gradient of most interest decomposes to  $\frac{\partial \phi_{T+1}^{BR}}{\partial \theta} = \frac{\partial \phi_{T+1}^{BR}}{\partial \pi_T} \frac{\partial \pi_T}{\partial \theta} + \frac{\partial \phi_{T+1}^{BR}}{\partial \Phi_T} \frac{\partial \Phi_T}{\partial \theta}$

# Recent Advance: Auto-PSRO Result

- **1st question:** is our method any good on the environments where it is trained?
  - Due to long-trajectory issues, we also focus on the *approximate* best-response setting

- Performance *at least* as good as baseline measures
- Outperforms PSRO in multiple settings

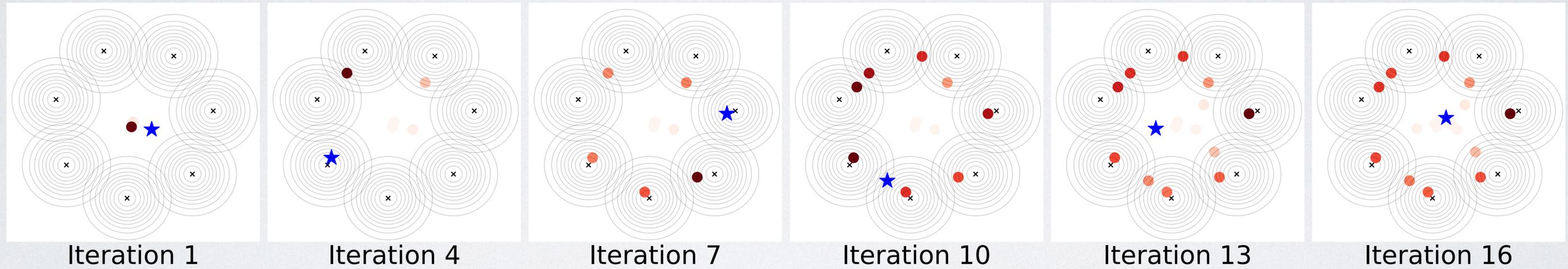


# Recent Advance: Auto-PSRO Result

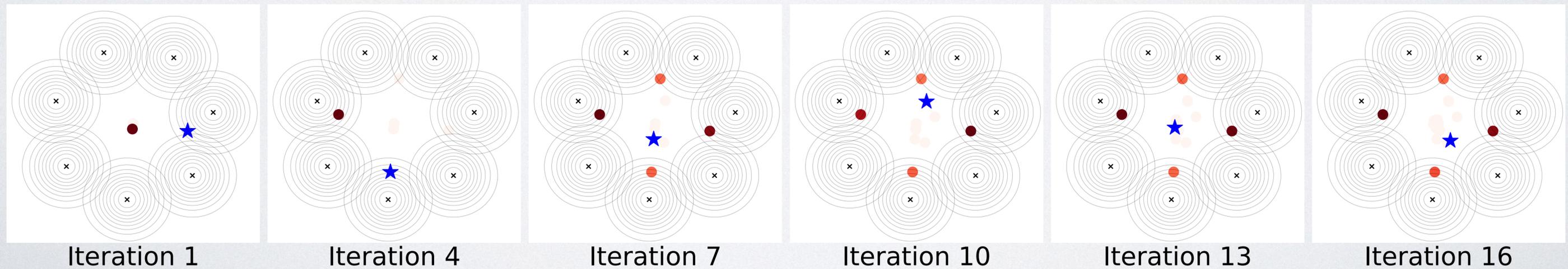
- **2nd question:** What is the learned multi-agent algorithm ?

- Compare agents found and their respective densities in the meta-distribution

Ours



PSRO



# Recent Advance: Auto-PSRO Result

- **3rd question:** Can the learned solver generalise over different games?
  - the most promising and striking aspect of LMAC - Train on small games and generalise to large game, e.g., train on Kuhn Poker and test on Leduc Poker

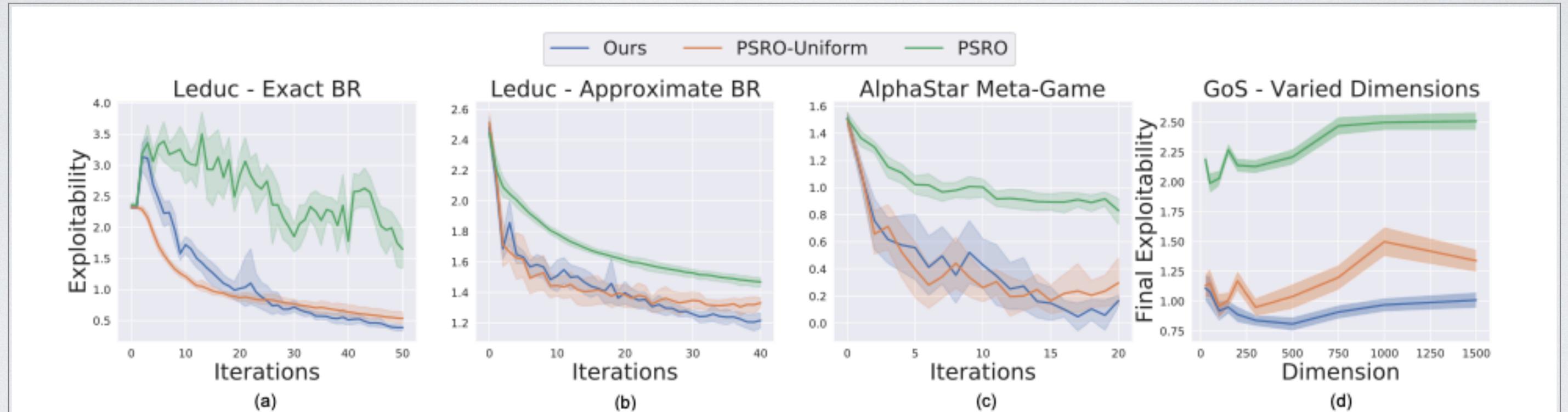


Figure 5: (a) Exploitability when trained on Kuhn Poker with an exact tabular BR oracle using ES-LMAC and tested on Leduc Poker. (b) Same as (a) with approximate tabular BR V2 (c) Exploitability when trained on GoS with a GD oracle and tested on the AlphaStar meta-game from [8] (d) Final exploitability when trained on 200 Dimension GoS and tested on a variety of dimension size GoS.

# Additional Resources:

- If you want to know more details about PSRO and its variations, please refer to
  - Talk: <https://www.bilibili.com/video/av969218959/>
  - Slides: <https://rlchina.org/lectures/lecture11.pdf>
- A self-contained MARL survey from game theoretical perspective:
  - <https://arxiv.org/abs/2011.00583>
- If you want to get hands on to solving some two-player zero-sum games, e.g., Poker/Chess
  - <https://arxiv.org/pdf/2103.00187.pdf>
  - [https://github.com/aicenter/openspiel\\_reproductions](https://github.com/aicenter/openspiel_reproductions)

Thanks!

Any Questions?