

A UNIFIED GAME-THEORETIC APPROACH TO MULTIAGENT REINFORCEMENT LEARNING

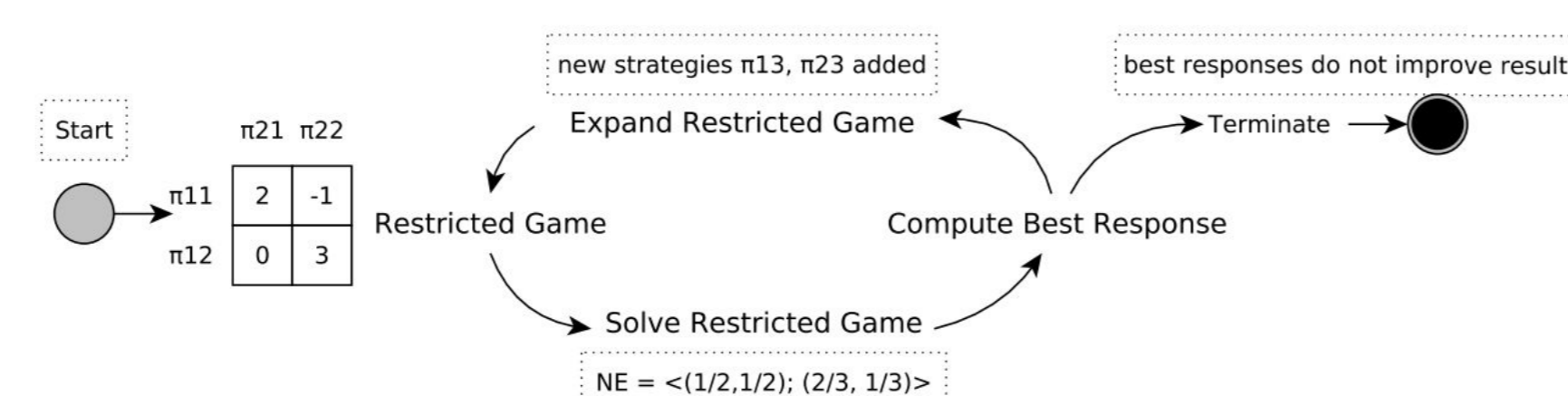


Marc Lanctot, Vinicius Zambaldi, Audrūnas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, Thore Graepel

INTRODUCTION & MOTIVATION

- Renewed interest in multiagent reinforcement learning (MARL)
 - Naive solution: Independent RL (InRL):
 - Every agent treats other agents as part of their local environment
 - Good:** simple
 - Bad:** no convergence guarantees
 - Ugly:** overfitting, lack of generalization
- New metric to quantify the **joint policy correlation (JPC)**
- We present a new algorithm:
 - Based on **empirical game-theoretic analysis** (Wellman '06, Walsh et al. '02)
 - Uses SoTA **deep reinforcement learning**: The Reactor (Gruslys et al. '17)
 - Generalizes fictitious play, iterated best resp., InRL, double oracle
 - Reduces JPC & finds trade-off between best resp./Nash

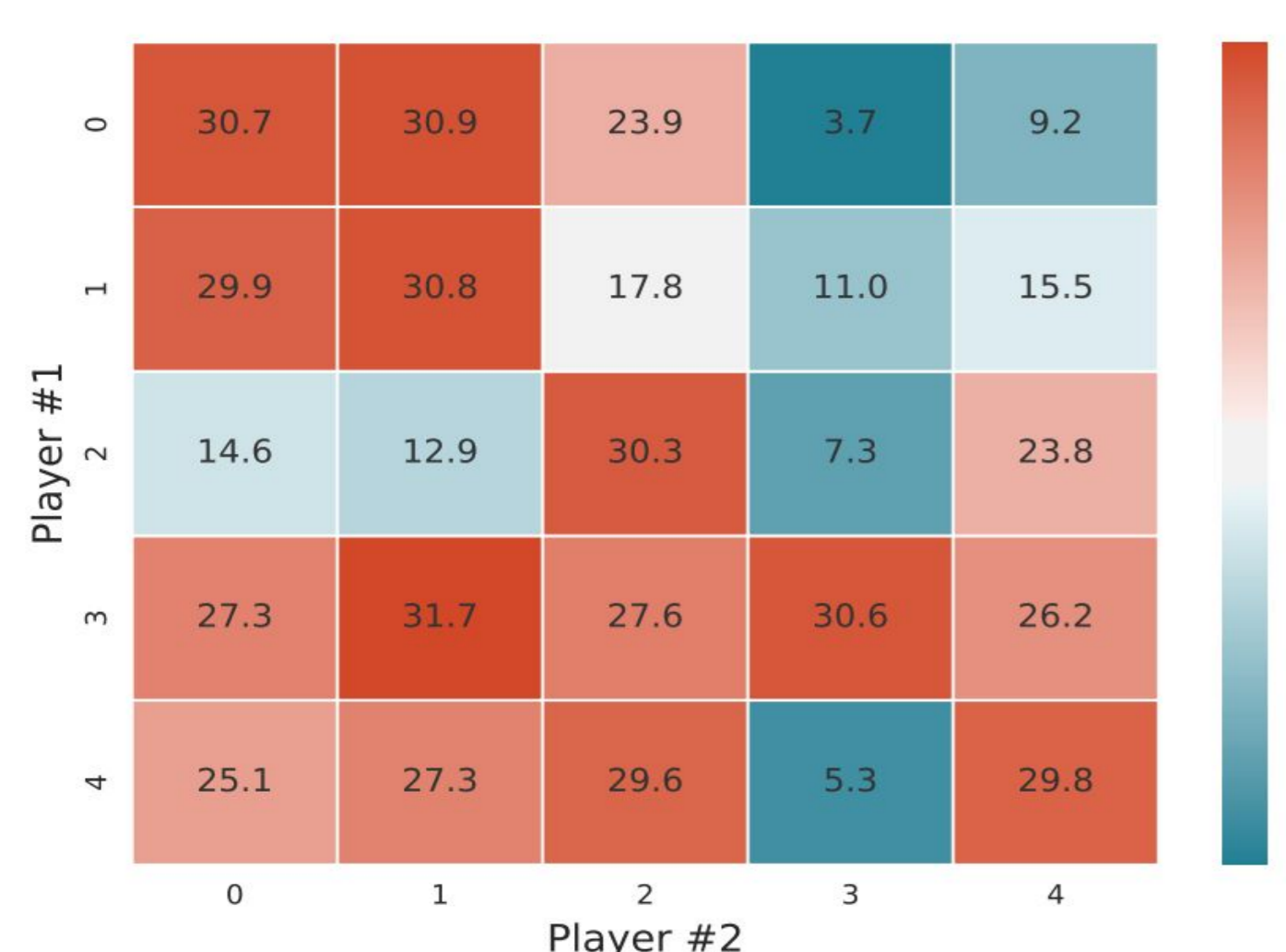
Our algorithm is largely inspired by double oracle (McMahan, Gordon, and Blum '03)



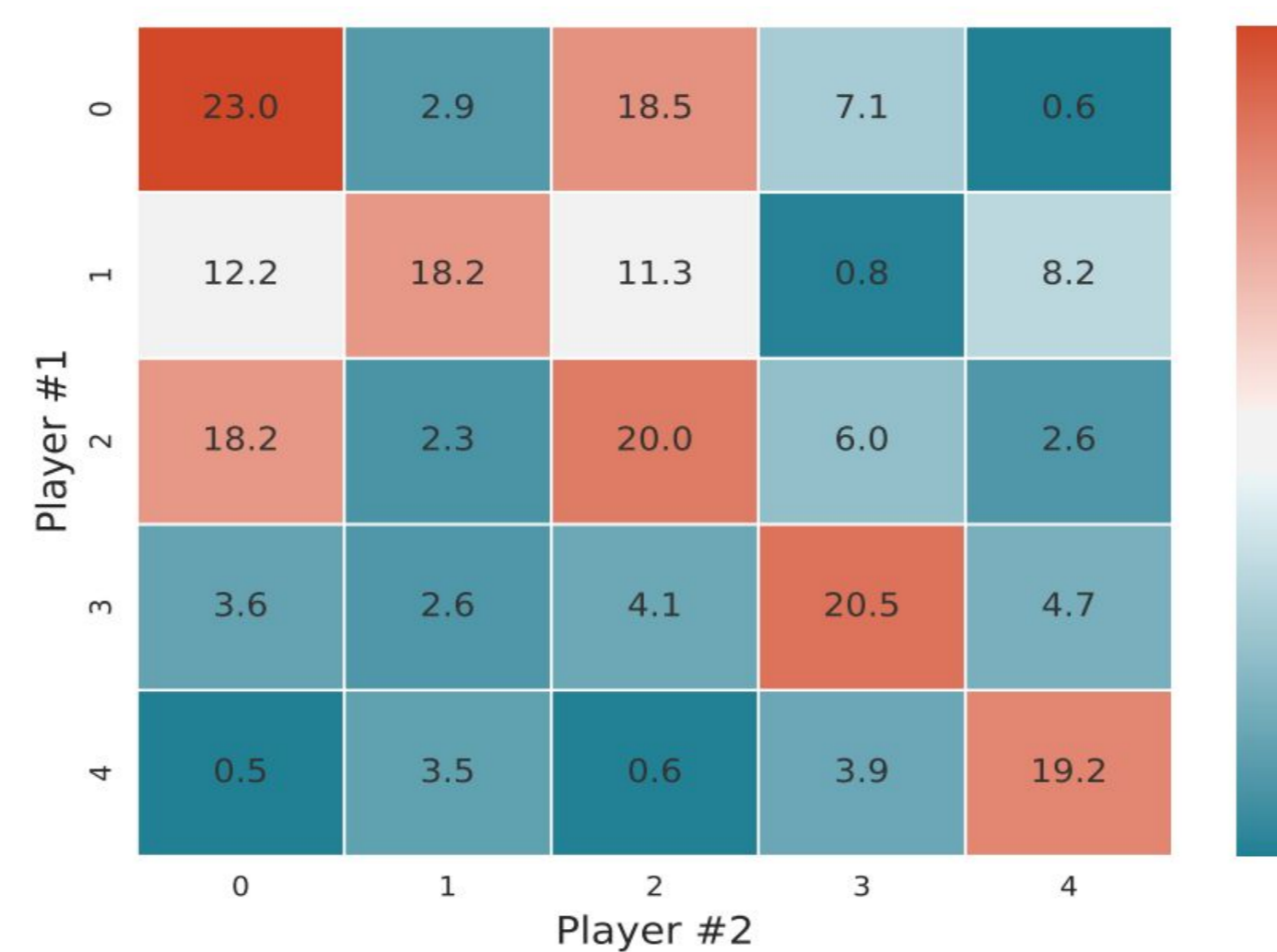
JOINT POLICY CORRELATION (JPC) IN INDEPENDENT RL

- Start D independent training instances $\rightarrow (\pi_1^d, \pi_2^d)$, for $d \in \{1, \dots, D\}$
- Produce **JPC Matrix**: $a_{i,j} = \mathbb{E}_{\rho \sim \pi_1^i, \pi_2^j} [\text{RETURN}(\rho)]$
- Define:
 - Avg. prop. loss** (from miscoordination): $R_- = \frac{\bar{D} - \bar{O}}{D}$

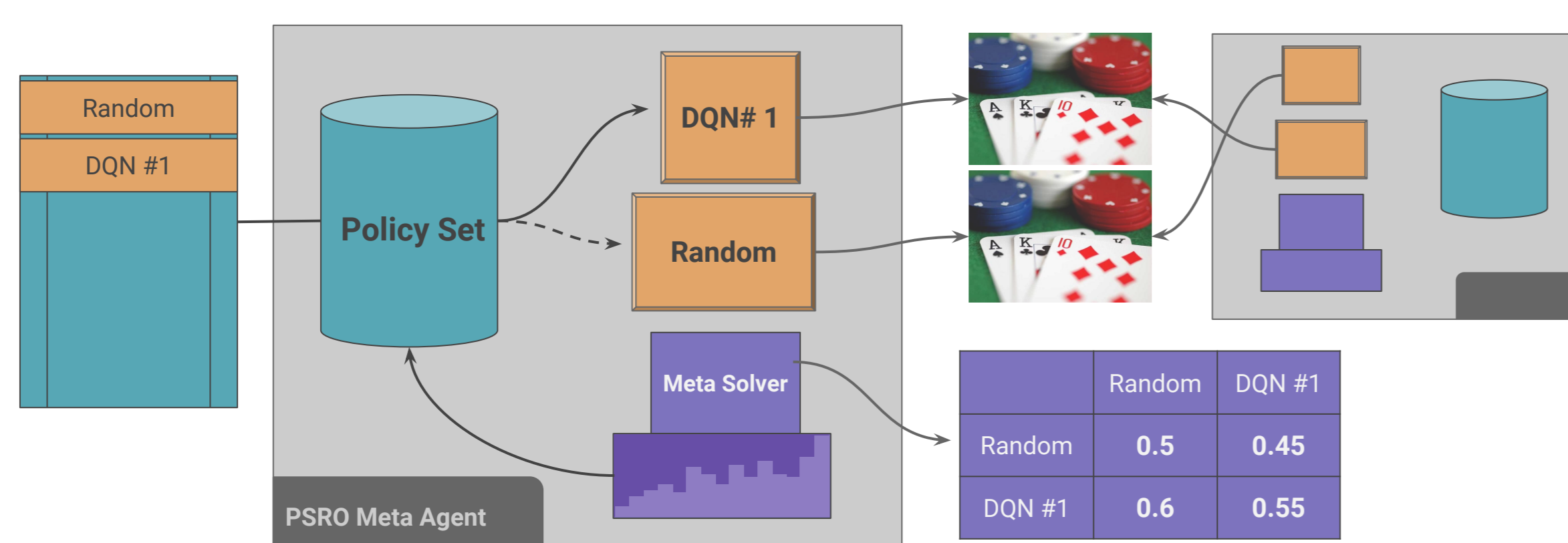
InRL - Laser Tag (small2)



InRL - Laser Tag (small4)



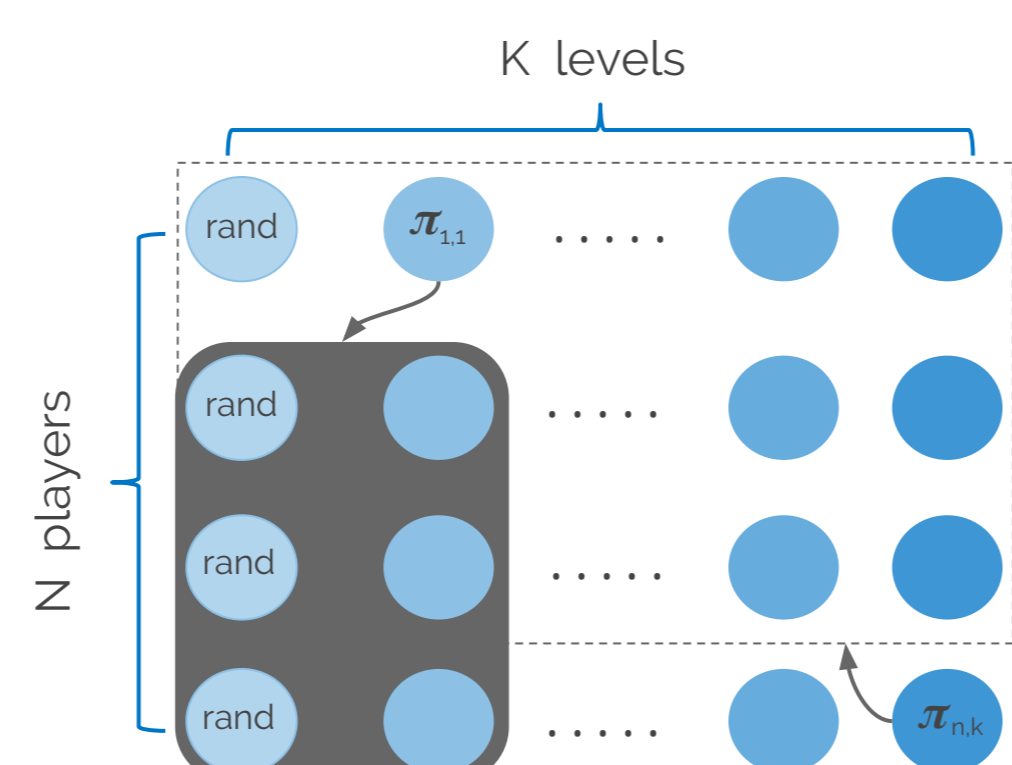
POLICY-SPACE RESPONSE ORACLES (PSRO)



Player i produces approx. best resp. **oracles**: $\Pi_i = (\pi_{i,0}, \pi_{i,1}, \dots, \pi_{i,K})$ and maintains a **meta-strategy** $\sigma_i = (\sigma_{i,0}, \sigma_{i,1}, \dots, \sigma_{i,K}) \in \Delta(\Pi)$
All sets of oracles together induce an **empirical game** (Π^K, U^{Π^K})

DEEP COGNITIVE HIERARCHIES (DCH)

- Parallel PSRO
- Choose K levels
- Start $N \times K$ processes (dist. oracles)
- Async. save + refresh $\pi_{i,k}, \sigma_{i,k}$
- Similar to Camerer et al.
- Space $O(N^2 K^2) < O(K^N)$



META-SOLVERS

Main idea: "solve" the empirical game using learning from expert advice (see Cesa-Bianchi & Lugosi '06)

Limited feedback? No problem: use bandit versions!

- Hedge / EXP3 (Auer et al. 1995)
 - Regret Matching (Hart & Mas-Colell '00, '01)
 - Projected Replicator Dynamics (New!)
- Force **exploratory** meta-strategies

$$U^{\Pi} = (\mathbf{A}, \mathbf{B}) \quad \frac{dx_k}{dt} = x_k[(\mathbf{A}\mathbf{y})_k - \mathbf{x}^T \mathbf{A}\mathbf{y}], \quad \frac{dy_k}{dt} = y_k[(\mathbf{x}^T \mathbf{B})_k - \mathbf{x}^T \mathbf{B}\mathbf{y}],$$

$$\sigma = (\mathbf{x}, \mathbf{y}) \quad \mathbf{x} \leftarrow P(\mathbf{x} + \delta \frac{dx}{dt}) \quad \mathbf{y} \leftarrow P(\mathbf{y} + \delta \frac{dy}{dt})$$

$$\sigma_{i,k} \geq \frac{\gamma}{K+1}$$

Algorithm

- Fictitious Play
- Iterated Best Response
- Independent RL
- Double Oracle

Meta-Strategy

$$\sigma_i = \left(\frac{1}{K}, \frac{1}{K}, \dots, \frac{1}{K}, 0 \right)$$

$$\sigma_i = (0, 0, \dots, 1, 0)$$

$$\sigma_i = (0, 0, \dots, 0, 1)$$

$$\sigma_i \in \text{NashEq}(U^{\Pi})$$

RESULTS

First-Person Gridworld Games

Game	\bar{D}	\bar{O}	R- (InRL)	R- (DCH)	Delta
LT small2	30.4	20.1	34.2 %	5.5 %	-28.7 %
LT small3	23.1	9.1	62.5 %	8.2 %	-54.3 %
LT small4	20.2	5.7	71.7 %	15.0 %	-56.7 %
Gathering	147.3	146.8	-	-	-
Pathfind	108.7	106.3	-	-	-

- Laser Tag:
- Moves: Turn left, right, fwd, back, strafe left, strafe right, timeout beam
 - +1 for a tag
 - Agents have local views
 - Agents re-enter at spawn points

JPC reduction using DCH @ 10 levels ($\sigma_{i,10}$)

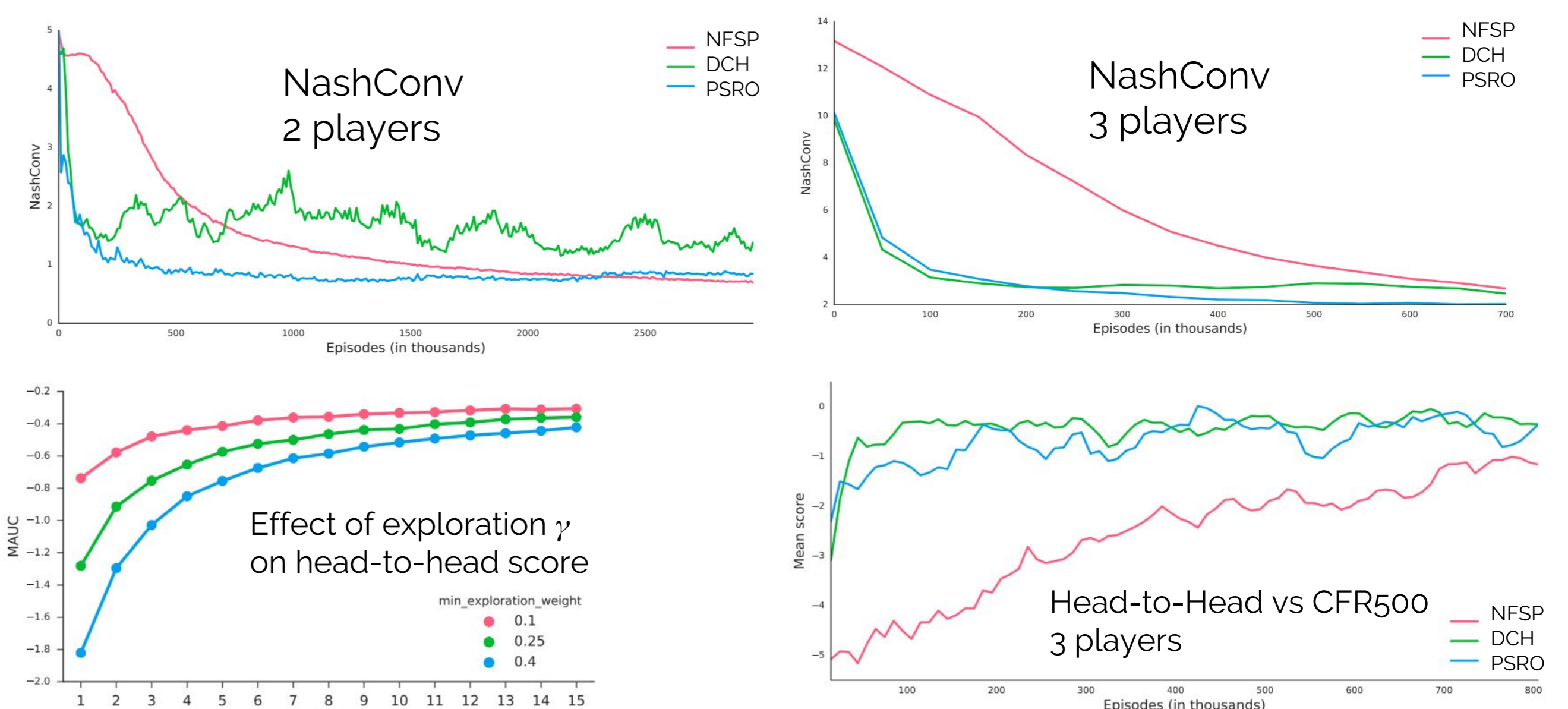
- No meta-strategy ($\pi_{i,10}$)? **14.7%, 27%, 11.8%**
- Five/three levels ($\sigma_{i,5}, \sigma_{i,3}$) small4: **15.6, 24.6**

Leduc Poker

- 1 chip ante: 2 rounds
- 6-card deck, 2 suits, 1 priv card
- Limit raises: 2 per rnd (2,4)
- Public card after rnd 2

- Compare to:
- NFSP (Heinrich & Silver '16)
 - CFR (Zinkevich et al. '07)

$$\text{NASHCONV}(\sigma) = \sum_i^n \max_{\sigma'_i} u_i(\sigma'_i, \sigma_{-i}) - u_i(\sigma)$$



Effect of levels/hierarchy on NashConv

