# Deep Learning in Tree-Based Game Solving 3

Stephen McAleer

# Outline of the next few lectures

- Deep learning in tree-based game solving 1
    - Deep learning recap
    - NFSP
    - Deep CFR
    - Policy gradient methods
- Deep learning in tree-based game solving 2
    - MCCFR
    - DREAM
    - ESCHER
    - NeuRD
- Deep learning in tree-based game solving 3
    - DeepNash for expert-level Stratego
- Deep learning in tree-based game solving 4
    - AlphaStar and OpenAI 5 for SOTA in video games
    - Double Oracle brief intro
- SOTA in double oracle algorithms
    - PSRO
    - XDO
    - SP-PSRO

# A Taxonomy of Game-Theoretic RL

- Counterfactual Regret Minimization (Zinkevich et al. 2007)
    - CFR: Zinkevich et al. 2007
    - MC-CFR: Lanctot et al. 2009
    - **Deep CFR: Brown et al. 2019**
    - DREAM: Steinberger et al. 2020                        Lecture 1
    - ESCHER: McAleer et al. 2022
- Policy Gradients
    - **Regret Policy Gradient (Srinivasan et al. 2018)**
    - OpenAI Five (OpenAI 2019)
    - Neural Replicator Dynamics (Hennes, Morrill, and Omidshafiei et al. 2020)
    - Actor Critic Hedge (Fu et al. 2022)
    - DeepNash for expert-level Stratego (Perolat, de Vylder, and Tuyls et al. 2022)
    - Magnetic Mirror Descent (Sokota et al. 2022)
- PSRO (McMahan et. al. 2003, Lanctot et al. 2017)
    - AlphaStar for expert-level Starcraft (Vinyals et al. 2019)
    - Pipeline PSRO (McAleer and Lanier et al. 2020)
    - α-PSRO (Muller et al. 2020)
    - XDO (McAleer et al. 2021)
    - Joint-PSRO (Marris et al. 2021)
    - Anytime PSRO (McAleer et al. 2022)
    - Self-Play PSRO (McAleer et al. 2022)
- **Neural Fictitious Self Play (Heinrich and Silver 2016)**

# A Taxonomy of Game-Theoretic RL

- Counterfactual Regret Minimization (Zinkevich et al. 2007)
    - CFR: Zinkevich et al. 2007
    - **MC-CFR: Lanctot et al. 2009**
    - Deep CFR: Brown et al. 2019
    - **DREAM: Steinberger et al. 2020**
    - **ESCHER: McAleer et al. 2022**
- Policy Gradients
    - Regret Policy Gradient (Srinivasan et al. 2018)
    - OpenAI Five (OpenAI 2019)
    - **Neural Replicator Dynamics (Hennes, Morrill, and Omidshafiei et al. 2020)**
    - Actor Critic Hedge (Fu et al. 2022)
    - DeepNash for expert-level Stratego (Perolat, de Vylder, and Tuyls et al. 2022)
    - Magnetic Mirror Descent (Sokota et al. 2022)
- PSRO (McMahan et. al. 2003, Lanctot et al. 2017)
    - AlphaStar for expert-level Starcraft (Vinyals et al. 2019)
    - Pipeline PSRO (McAleer and Lanier et al. 2020)
    - α-PSRO (Muller et al. 2020)
    - XDO (McAleer et al. 2021)
    - Joint-PSRO (Marris et al. 2021)
    - Anytime PSRO (McAleer et al. 2022)
    - Self-Play PSRO (McAleer et al. 2022)
- Neural Fictitious Self Play (Heinrich and Silver 2016)

Lecture 2

# A Taxonomy of Game-Theoretic RL

- Counterfactual Regret Minimization (Zinkevich et al. 2007)
    - CFR: Zinkevich et al. 2007
    - MC-CFR: Lanctot et al. 2009
    - Deep CFR: Brown et al. 2019
    - DREAM: Steinberger et al. 2020
    - ESCHER: McAleer et al. 2022

        Lecture 3 (This Lecture)

- Policy Gradients
    - Regret Policy Gradient (Srinivasan et al. 2018)
    - OpenAI Five (OpenAI 2019)
    - Neural Replicator Dynamics (Hennes, Morrill, and Omidshafiei et al. 2020)
    - Actor Critic Hedge (Fu et al. 2022)
    - **DeepNash for expert-level Stratego (Perolat, de Vylder, and Tuyls et al. 2022)**
        - **From Poincaré Recurrence to Convergence in Imperfect Information Games: Finding Equilibrium via Regularization (Perolat et al. 2021)**
    - **Magnetic Mirror Descent (Sokota et al. 2022)**
- PSRO (McMahan et. al. 2003, Lanctot et al. 2017)
    - AlphaStar for expert-level Starcraft (Vinyals et al. 2019)
    - Pipeline PSRO (McAleer and Lanier et al. 2020)
    - α-PSRO (Muller et al. 2020)
    - XDO (McAleer et al. 2021)
    - Joint-PSRO (Marris et al. 2021)
    - Anytime PSRO (McAleer et al. 2022)
    - Self-Play PSRO (McAleer et al. 2022)
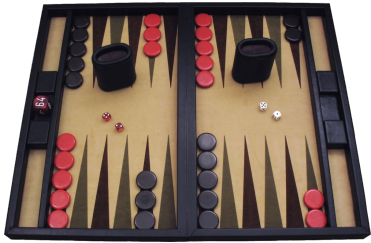- Neural Fictitious Self Play (Heinrich and Silver 2016)

# A Taxonomy of Game-Theoretic RL

- Counterfactual Regret Minimization (Zinkevich et al. 2007)
    - CFR: Zinkevich et al. 2007
    - MC-CFR: Lanctot et al. 2009
    - Deep CFR: Brown et al. 2019
    - DREAM: Steinberger et al. 2020
    - ESCHER: McAleer et al. 2022
- Policy Gradients
    - Regret Policy Gradient (Srinivasan et al. 2018)
    - **OpenAI Five (OpenAI 2019)**
    - Neural Replicator Dynamics (Hennes, Morrill, and Omidshafiei et al. 2020)
    - Actor Critic Hedge (Fu et al. 2022)
    - DeepNash for expert-level Stratego (Perolat, de Vylder, and Tuyls et al. 2022)
    - Magnetic Mirror Descent (Sokota et al. 2022)
- PSRO (McMahan et. al. 2003, Lanctot et al. 2017)
    - **AlphaStar for expert-level Starcraft (Vinyals et al. 2019)**
    - Pipeline PSRO (McAleer and Lanier et al. 2020)
    - α-PSRO (Muller et al. 2020)
    - XDO (McAleer et al. 2021)
    - Joint-PSRO (Marris et al. 2021)
    - Anytime PSRO (McAleer et al. 2022)
    - Self-Play PSRO (McAleer et al. 2022)
- Neural Fictitious Self Play (Heinrich and Silver 2016)

Lecture 4

# A Taxonomy of Game-Theoretic RL

- Counterfactual Regret Minimization (Zinkevich et al. 2007)
    - CFR: Zinkevich et al. 2007
    - MC-CFR: Lanctot et al. 2009
    - Deep CFR: Brown et al. 2019
    - DREAM: Steinberger et al. 2020
    - ESCHER: McAleer et al. 2022
- Policy Gradients
    - Regret Policy Gradient (Srinivasan et al. 2018)
    - OpenAI Five (OpenAI 2019)
    - Neural Replicator Dynamics (Hennes, Morrill, and Omidshafiei et al. 2020)
    - Actor Critic Hedge (Fu et al. 2022)
    - DeepNash for expert-level Stratego (Perolat, de Vylder, and Tuyls et al. 2022)
    - Magnetic Mirror Descent (Sokota et al. 2022)
- **PSRO (McMahan et. al. 2003, Lanctot et al. 2017)**
    - AlphaStar for expert-level Starcraft (Vinyals et al. 2019)
    - **Pipeline PSRO (McAleer and Lanier et al. 2020)**
    - **α-PSRO (Muller et al. 2020)**
    - **XDO (McAleer et al. 2021)**
    - **Joint-PSRO (Marris et al. 2021)**
    - **Anytime PSRO (McAleer et al. 2022)**
    - **Self-Play PSRO (McAleer et al. 2022)**
- Neural Fictitious Self Play (Heinrich and Silver 2016)
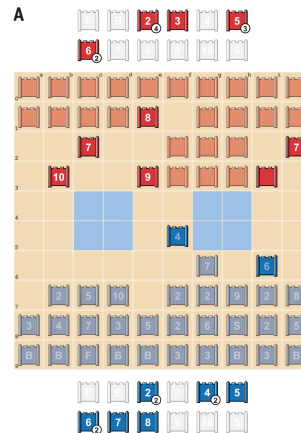
Lecture 5

# Games in AI



Backgammon
1992



Chess
1997
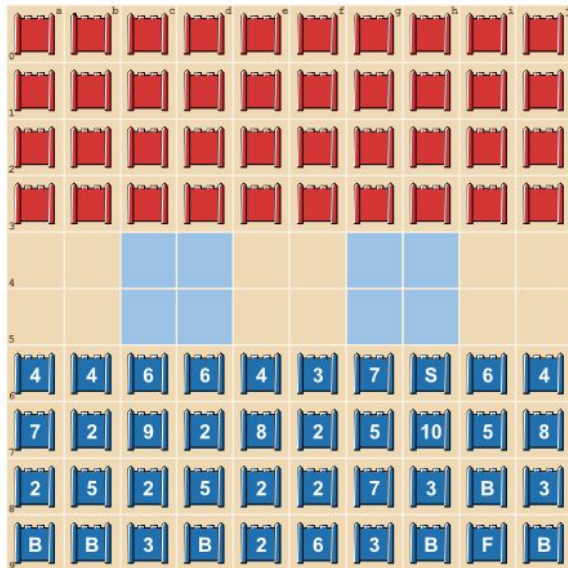


Go
2016



Poker
2017/2019
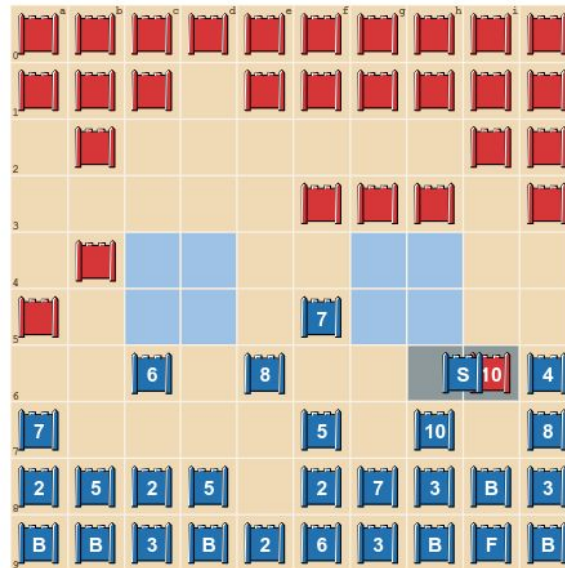


Starcraft/Dota
2019



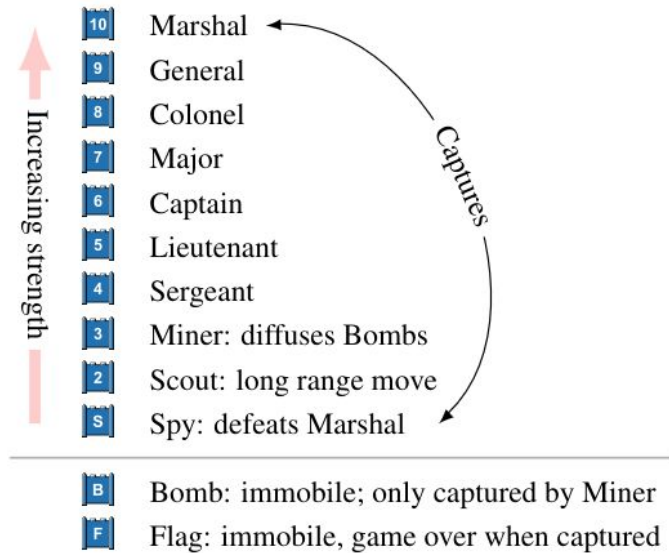Stratego
2022



Diplomacy
2022

# Stratego

- Pieces are numbered from 2 to 10 (Also a spy, bomb and **flag)**
- Higher numbers capture lower numbers (Exceptions: spys, bombs)
- First, both players place their pieces (**Can't see opponents pieces)**
- Each piece moves one square (Exception: 2)
- If your piece is captured, you see the other piece number
- Objective is to capture the opponent's flag
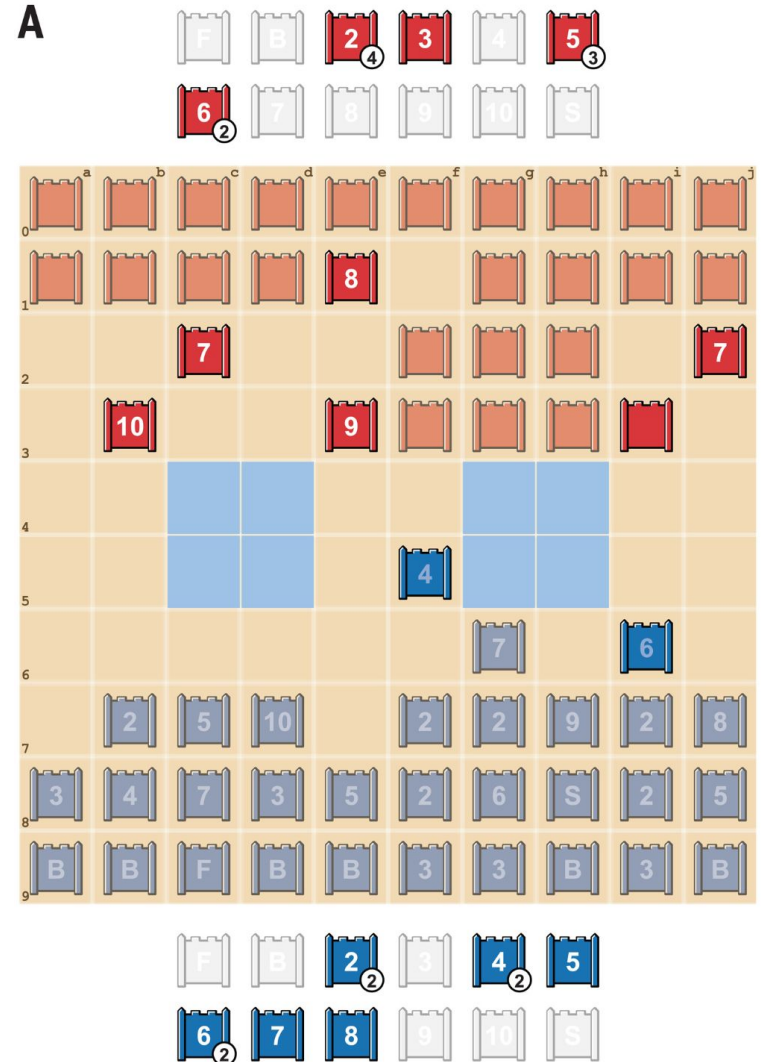


**Phase 1:** Private deployment

**Phase 2:** Game play

Piece types

| | |
|---|---|
| 10 | Marshal |
| 9 | General |
| 8 | Colonel |
| 7 | Major |
| 6 | Captain |
| 5 | Lieutenant |
| 4 | Sergeant |
| 3 | Miner: diffuses Bombs |
| 2 | Scout: long range move |
| S | Spy: defeats Marshal |
| B | Bomb: immobile; only captured by Miner |
| F | Flag: immobile, game over when captured |

Increasing strength
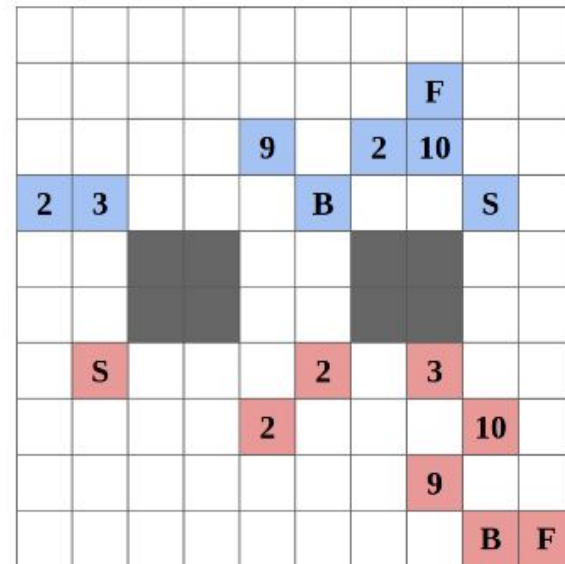
Captures

# Stratego

- Two challenges: **size** and **imperfect information**
- Size: order of $10^{535}$ nodes
    - Texas hold 'em: $10^{164}$ nodes
    - Go: $10^{360}$ nodes
- Imperfect information
    - $10^{66}$ possible deployments
    - Can't use perfect-info search
    - Bluffing, mixing are important
    - Gathering and hiding information very important
- Compared to video games, decisions are made deliberately
    - Doesn't just test reaction time and instincts

# Stratego

- Existing approaches have hand-coded rules and play at an amateur level
- PSRO-based approach got SOTA on Barrage Stratego in 2020
  - Still played at an amateur level

| Name | P2SRO Win Rate vs. Bot |
|---|---|
| Asmodeus | 81% |
| Celsius | 70% |
| Vixen | 69% |
| Celsius1.1 | 65% |
| **All Bots Average** | **71%** |

McAleer*, Lanier*, Fox, Baldi. Pipeline PSRO: A Scalable Approach for Finding Approximate Nash Equilibria in Large Games. NeurIPS 2020

# Finding Equilibrium via Regularization

- Continuous-time Follow-the-Regularized Leader (FoReL)

$$y_t^i(a^i) = \int_0^t Q_{\pi_s}^i(a^i)ds \quad \text{and} \quad \pi_t^i = \arg\max_{p \in \Delta A} \Lambda^i(p, y_t^i)$$

$$\Lambda^i(p, y) = \langle y, p \rangle - \phi_i(p)$$

$$\phi_i^*(y) = \max_p \Lambda^i(p, y)$$

- Motivation: want to get last-iterate convergence

Perolat et al. From Poincare Recurrence to Convergence in Imperfect Information Games:
Finding Equilibrium via Regularization. ICML 2021.

# Finding Equilibrium via Regularization

- In two-player zero-sum games, the Nash Gap (exploitability) is preserved, so FoReL is recurrent

$$J(y) = \sum_{i=1}^{2} \left[ \phi_i^*(y_i) - \langle y_i, \pi_i^* \rangle \right]$$



Perolat et al. From Poincare Recurrence to Convergence in Imperfect Information Games:
Finding Equilibrium via Regularization. ICML 2021.

# Finding Equilibrium via Regularization

- If we modify the game to have this new policy-dependent reward function

$$r_\pi^i(a) = r^i(a^i, a^{-i}) - \eta \log \frac{\pi^i(a^i)}{\mu^i(a^i)} + \eta \log \frac{\pi^{-i}(a^{-i})}{\mu^{-i}(a^{-i})}$$

- Then FoReL is convergent

$$\frac{d}{dt} J(y) = \sum_{i=1}^{2} \underbrace{[V_{\pi_t^i, \pi^{*-i}}^i - V_{\pi^*}^i]}_{\leq 0 \text{ because } \pi^* \text{ is a Nash}} - \eta \sum_{i=1}^{2} KL(\pi^{*i}, \pi_t^i)$$

Perolat et al. From Poincare Recurrence to Convergence in Imperfect Information Games: Finding Equilibrium via Regularization. ICML 2021.

# Finding Equilibrium via Regularization

- However, FoReL converges to a biased solution
- Plot shows eta= 0, 0.5, 1, and 10



Perolat et al. From Poincare Recurrence to Convergence in Imperfect Information Games:
Finding Equilibrium via Regularization. ICML 2021.

# Finding Equilibrium via Regularization

- Solve the original game by iteratively using last policy as the reference policy

$$r^i_{k,\pi}(h,a) = r^i(a^i, a^{-i}) - \eta \log \frac{\pi^i(a^i)}{\pi^i_{k-1}(a^i)} + \eta \log \frac{\pi^{-i}(a^{-i})}{\pi^{-i}_{k-1}(a^{-i})}$$

- This procedure monotonically gets closer to Nash

Perolat et al. From Poincare Recurrence to Convergence in Imperfect Information Games: Finding Equilibrium via Regularization. ICML 2021.

# DeepNash

- Two components
    - NeuRD
    - Regularized Nash Dynamics (R-NaD)



**Replicator dynamics:** $\frac{d}{d\tau}\pi_\tau^i(a^i) = \pi_\tau^i(a^i)\left[Q_{\pi_\tau}^i(a^i) - \sum_{b^i}\pi_\tau^i(b^i)Q_{\pi_\tau}^i(b^i)\right]$

**Reward transformation:** $r^i(\pi^i, \pi^{-i}, a^i, a^{-i}) = r^i(a^i, a^{-i}) - \eta\log\left(\frac{\pi^i(a^i)}{\pi_{\text{reg}}^i(a^i)}\right) + \eta\log\left(\frac{\pi^{-i}(a^{-i})}{\pi_{\text{reg}}^{-i}(a^{-i})}\right)$
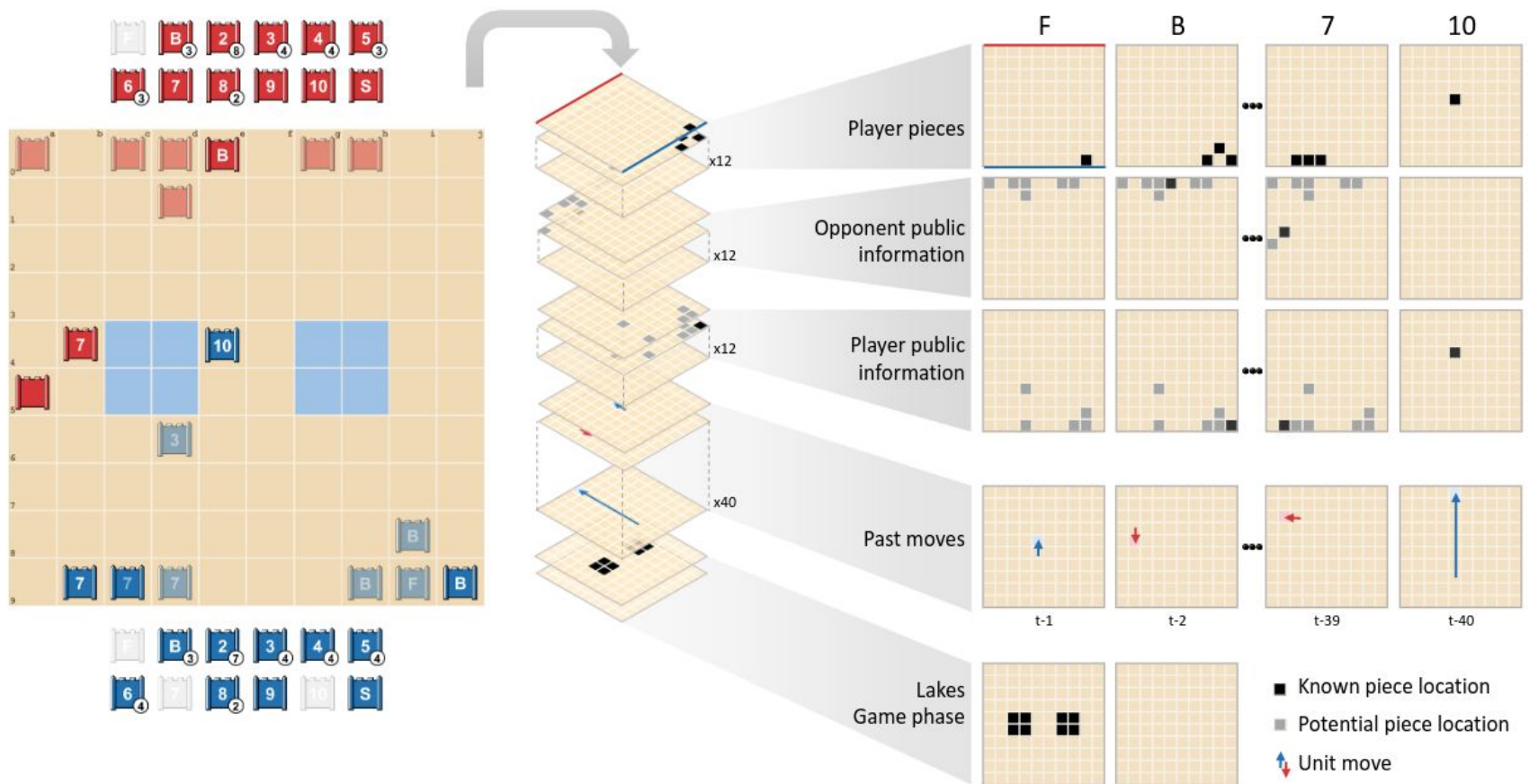
Perolat et al. Mastering the Game of Stratego with Model-Free Multiagent Reinforcement Learning. Science 2022

# DeepNash

- Regularized Nash Dynamics (R-NaD)
  - Same as in previous paper



(a) Matching pennies

(b) Algorithmic steps

(c) Dynamics and Lyapunov function

Figure 2: The R-NaD learning algorithm illustrated with the matching pennies game

# DeepNash

- Same reward transformation as before

$$r^i(\pi^i, \pi^{-i}, a^i, a^{-i}) = r^i(a^i, a^{-i}) - \eta \log\left(\frac{\pi^i(a^i)}{\pi^i_{\text{reg}}(a^i)}\right) + \eta \log\left(\frac{\pi^{-i}(a^{-i})}{\pi^{-i}_{\text{reg}}(a^{-i})}\right)$$

- Learn value function via V-Trace
- Learn policy via NeuRD

$$\Lambda_n = - \left[ \text{lr}_n \nabla l_{\text{critic}}(\theta_n) + \sum_{i=1}^{2} \frac{1}{t_{\text{effective}}} \sum_{t=0}^{t_{\text{effective}}} \sum_{a} \hat{\nabla}\theta(l_{\theta_n}(a, o_t)\text{Clip}\left(Q^{\psi_t}_{t,n}(a, o_t), c_{\text{clip NeuRD}}\right), \text{lr}_n, \beta) \right]$$

- Adapts IMPALA to parallelize

Perolat et al. Mastering the Game of Stratego with Model-Free Multiagent Reinforcement Learning.
Science 2022

# DeepNash

- Neural network input doesn't include full observation history, but a lot of it

# Results

| Opponent | Number of Games | Wins | Draws | Losses |
|---|---|---|---|---|
| Probe | 30 | 100.0% | 0.0% | 0.0% |
| Master of the Flag | 30 | 100.0% | 0.0% | 0.0% |
| Demon of Ignorance | 800 | 97.1% | 1.8% | 1.1% |
| Asmodeus | 800 | 99.7% | 0.0% | 0.3% |
| Celsius | 800 | 98.2% | 0.0% | 1.8% |
| Celsius1.1 | 800 | 97.9% | 0.0% | 2.1% |
| PeternLewis | 800 | 99.9% | 0.0% | 0.1% |
| Vixen | 800 | 100.0% | 0.0% | 0.0% |

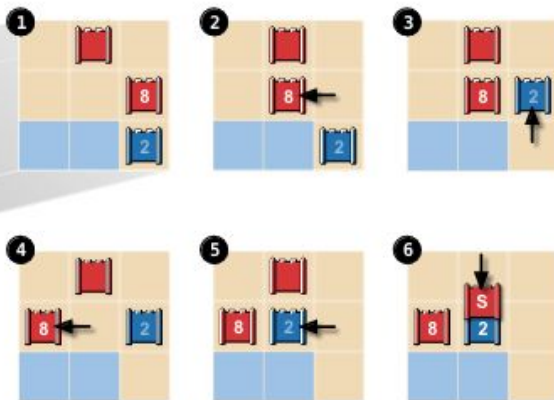**Expert-Level Performance:** Won 84% of games on online server, placing it 3rd all-time.
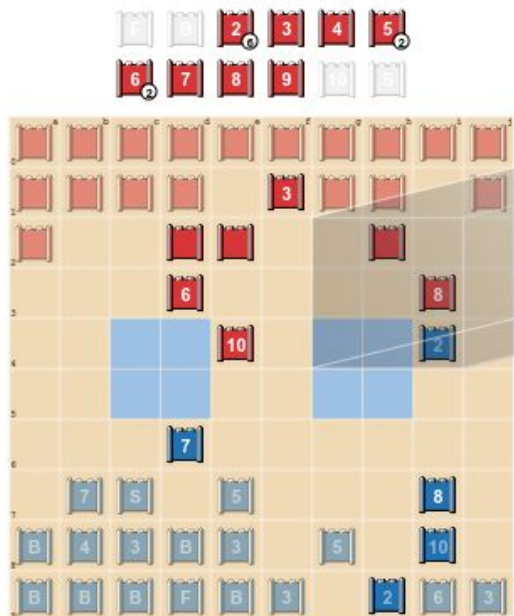
# Results



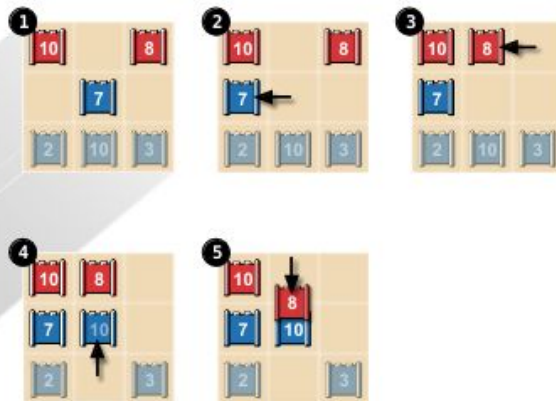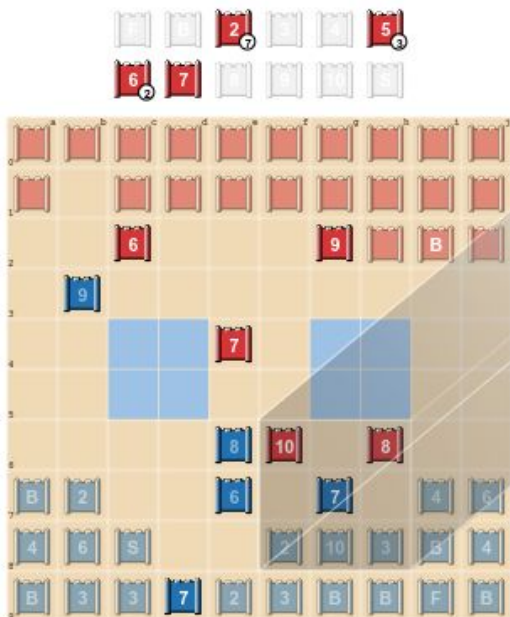(a) Four example deployments *DeepNash* played on Gravon.

(b) While Blue is behind a 7 and 8, none of its pieces are revealed and only two pieces moved. As a result *DeepNash* assesses its chance of winning to be still around 70% (Blue indeed won this match).
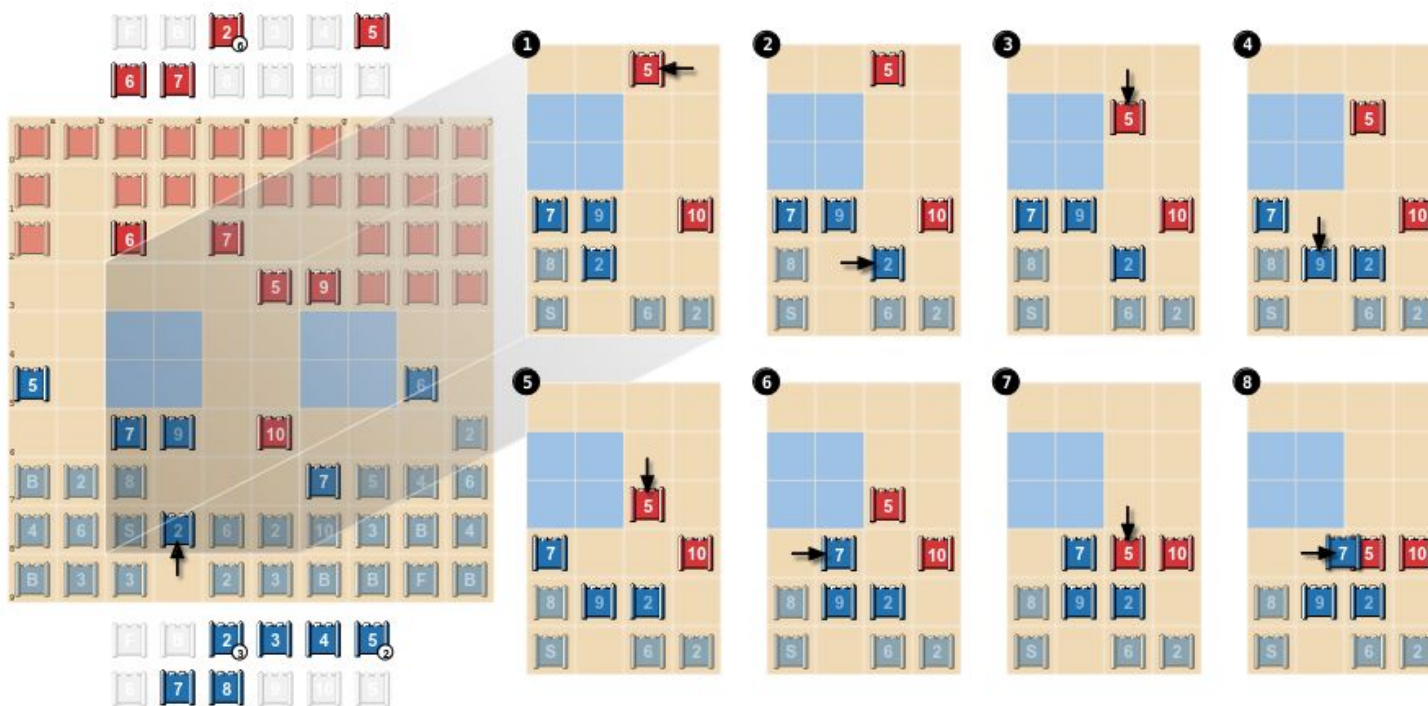
(c) Blue to move. *DeepNash*'s policy supports three moves at this state, with the indicated probabilities (the move on the right was played in the actual match). While Blue has the opportunity to capture the opponent's 6 with its 9, this move is not considered by *DeepNash*, likely because the protection of 9's identity is assessed to be more important than the material gain.

(a) Positive bluffing.

(b) Negative bluffing.

(c) *DeepNash* makes a Scout (2) behave like a Spy and gains material.

Figure 5: Illustration of *DeepNash* bluffing.

Perolat et al. Mastering the Game of Stratego with Model-Free Multiagent Reinforcement Learning. Science 2022

# What is mirror descent?

- Generalization of gradient descent to different notions of distance

$$x_{t+1} = \arg\min_x \langle g, x \rangle + \frac{1}{\eta} B(x, x_t)$$

- Negative Entropy (policy space):

$$\pi_{t+1} = \arg\max_\pi \langle q, \pi \rangle - \frac{1}{\eta} KL(\pi, \pi_t)$$

# What is **magnetic** mirror descent?

- Generalization of **regularized** gradient descent to different

  notions of distance

$$x_{t+1} = \arg \min_x \langle g, x \rangle + \frac{1}{\eta} \mathrm{B}(x, x_t) + \alpha \mathrm{B}(x, z)$$

- Negative Entropy (policy space):

$$\pi_{t+1} = \arg \max_\pi \langle q, \pi \rangle - \frac{1}{\eta} \mathrm{KL}(\pi, \pi_t) - \alpha \mathrm{KL}(\pi, \rho)$$

$$\propto \left[ \pi_t e^{\eta q} \rho^{\alpha \eta} \right]^{\frac{1}{1 + \alpha \eta}}$$
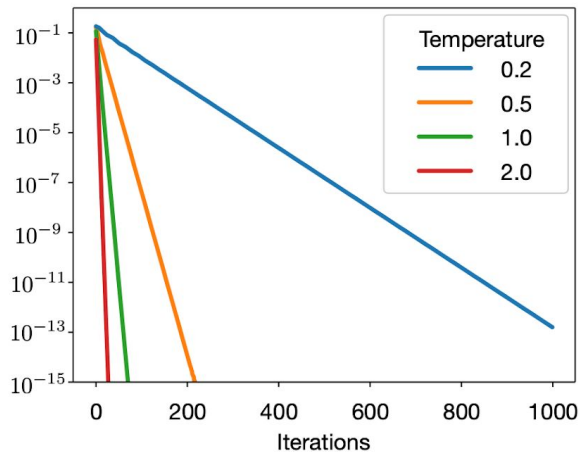
# Theoretical Grounding

In two-player zero-sum one-shot games, if
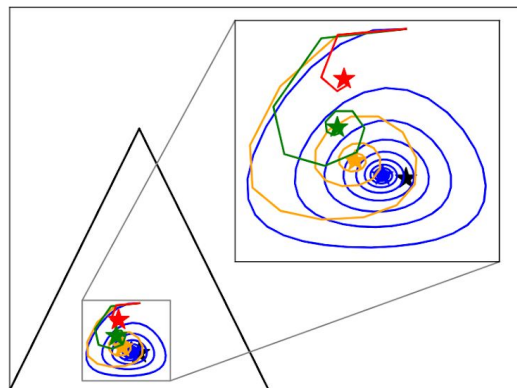$$\eta \leq \alpha/L^2$$
magnetic mirror descent converges exponentially fast to a

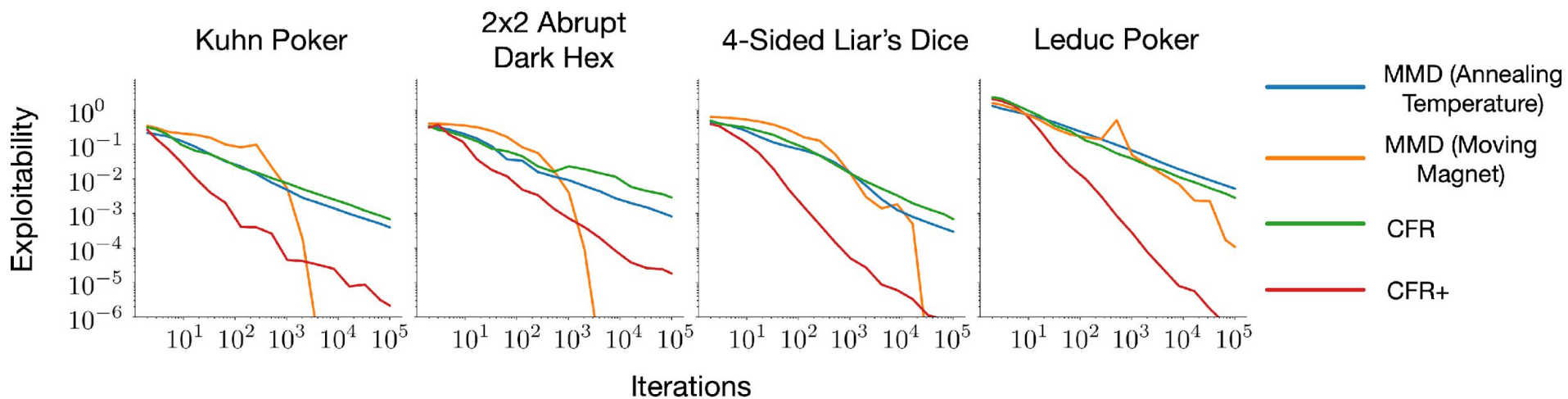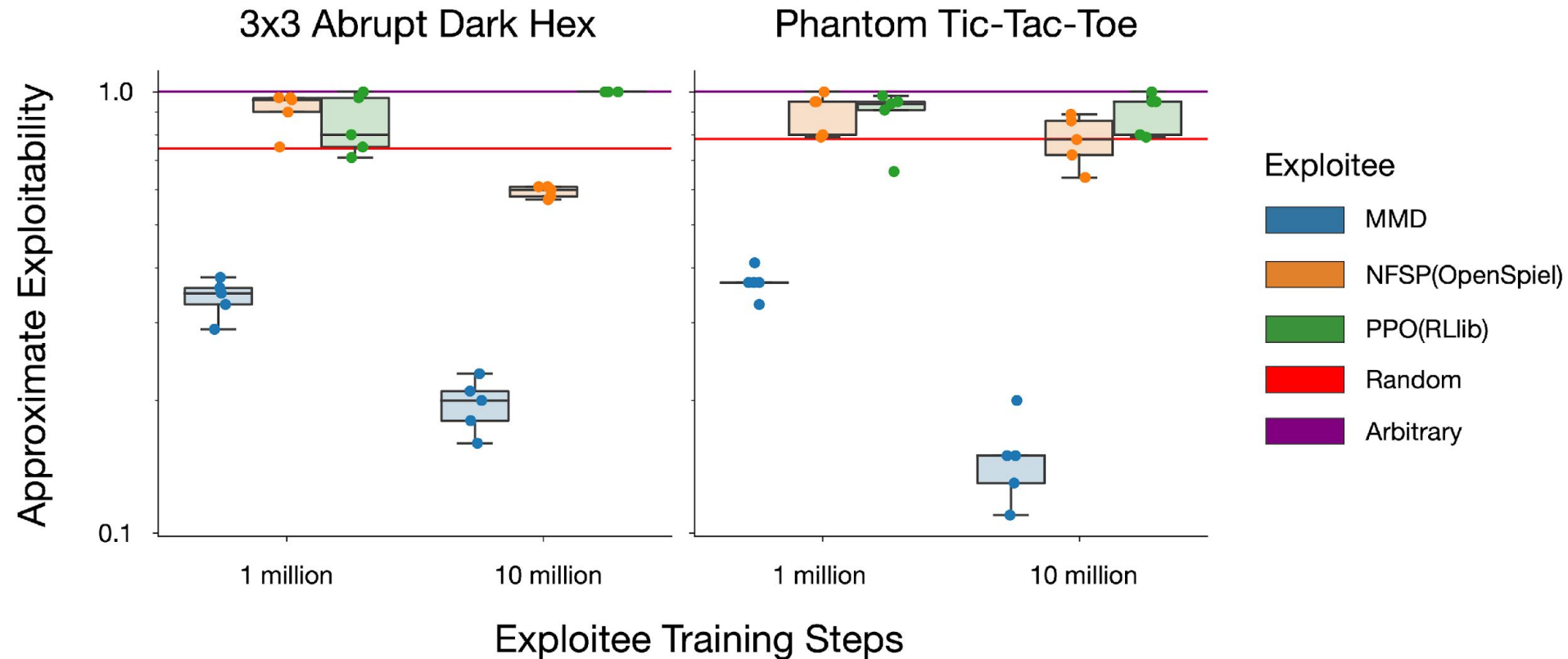regularized equilibrium in self play



KL Divergence to QRE



Simplex Trajectories



Payoff Matrix

|   | R | P | S |
|---|---|---|---|
| R | 0 | -1 | 3 |
| P | 1 | 0 | -3 |
| S | -3 | 3 | 0 |

# Comparison Against CFR

Sokota et al. A Unified Approach to Reinforcement Learning, Quantal Response Equilibria, and Two-Player Zero-Sum Games. ICLR 2023

# Deep RL Experiments: Approximate Exploitability



Sokota et al. A Unified Approach to Reinforcement Learning, Quantal Response Equilibria, and Two-Player Zero-Sum Games. ICLR 2023

# Deep RL Experiments: Head-to-Head Matchups



3x3 Abrupt Dark Hex — Phantom Tic-Tac-Toe

Opponent: MMD, NFSP(OpenSpiel), PPO(RLlib), Random, Arbitrary

Sokota et al. A Unified Approach to Reinforcement Learning, Quantal Response Equilibria, and Two-Player Zero-Sum Games. ICLR 2023