



NeuPL: Neural Population Learning

[read paper](#)

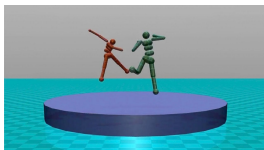
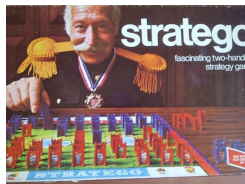
Siqi Liu^{1,2}, Luke Marris^{1,2}, Daniel Hennes², Josh Merel², Nicolas Heess², Thore Graepel^{1,2}

University College London, UK¹
DeepMind²

ICLR 2022

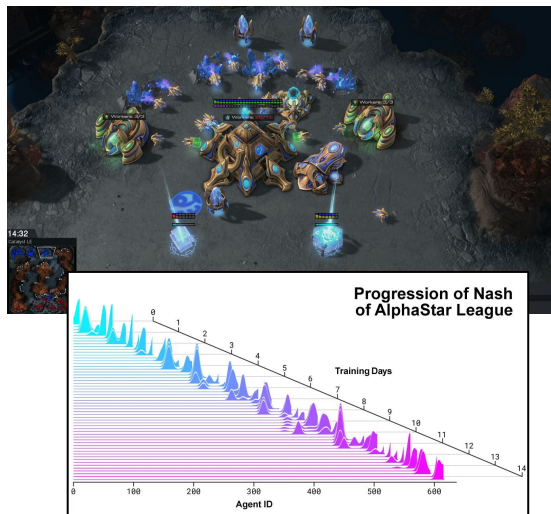
Learning in Games

Transitive Games: self-play is enough



Cyclic Games: game-theoretic reasoning

Learning in “Real-World” Games



PSRO^[1]: iterative learning of best-responses to mixture opponent strategies given by the meta-strategy solver (MSS) given expected returns between strategies.

- MSS = NE/Unif: convergence to NE.

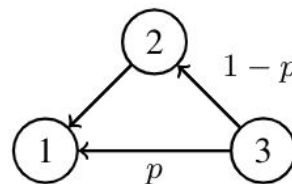
$$\sigma_2 = \Sigma_{[2, :]} = \text{SOLVE-NE}(\text{Payoffs}[:, 2, :2])$$

$$J_i = \mathbb{E}_{j \sim \Pr(\sigma_j)} [\mathbb{E}_{a \sim \pi_i, a' \sim \pi_j} [\sum_t r_t \gamma^t]]$$

Limitations:

- Independent iterative learning of policies;
- “Good-”responses due to early stopping.

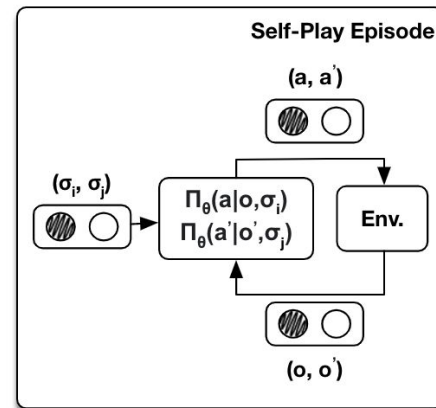
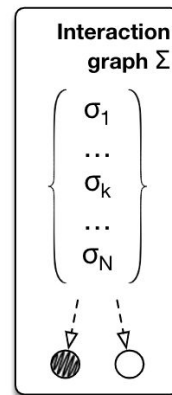
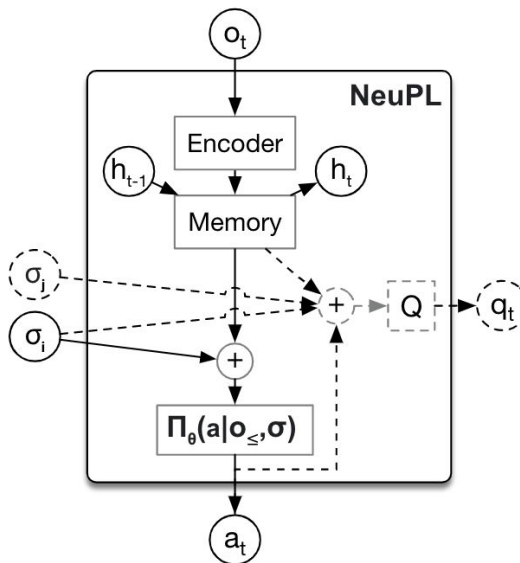
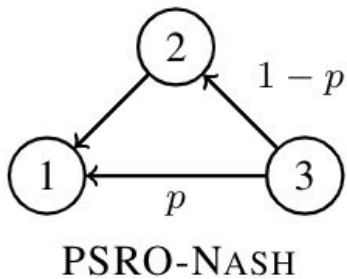
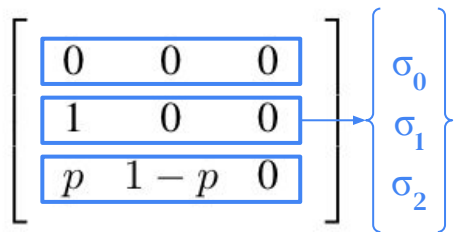
$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ p & 1-p & 0 \end{bmatrix}$$



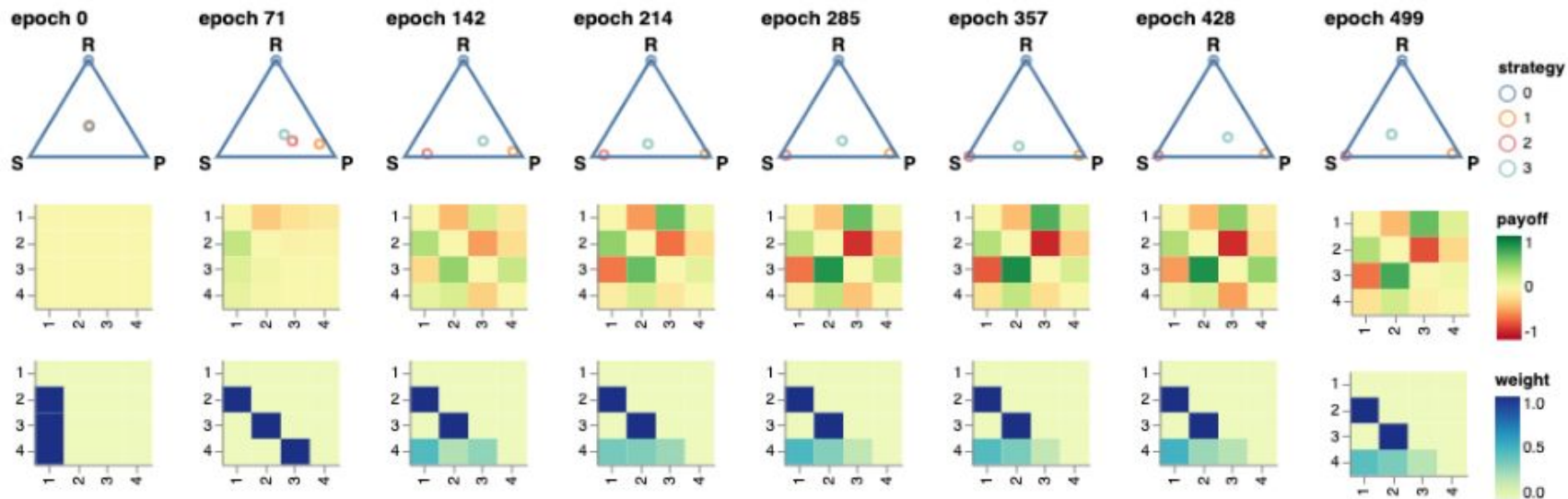
PSRO-NASH

Neural Population Learning

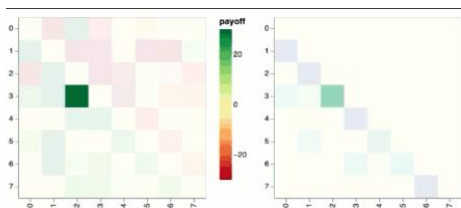
$$J_{\sigma_i} = \mathbb{E}_{\sigma_j \sim P(\sigma_i)} \left[a \sim \Pi_{\theta}(\cdot | o_{\leq t}, \sigma_i), a' \sim \Pi_{\theta}(\cdot | o'_{\leq t}, \sigma_j) \left[\sum_t r_t \gamma^t \right] \right]$$



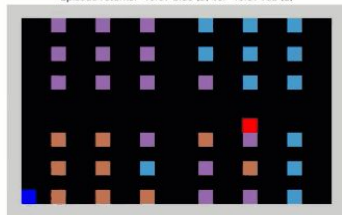
Rock-Paper-Scissors



Running-with-Scissors

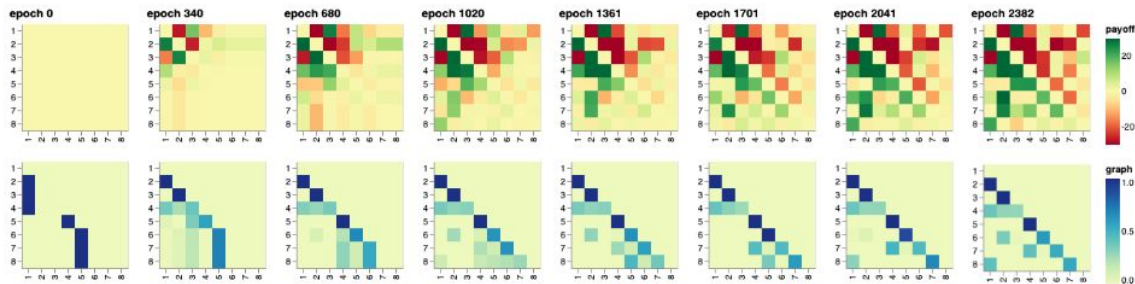
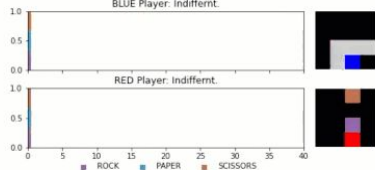


Episode returns: 46.67 blue (3) vs. -46.67 red (2)



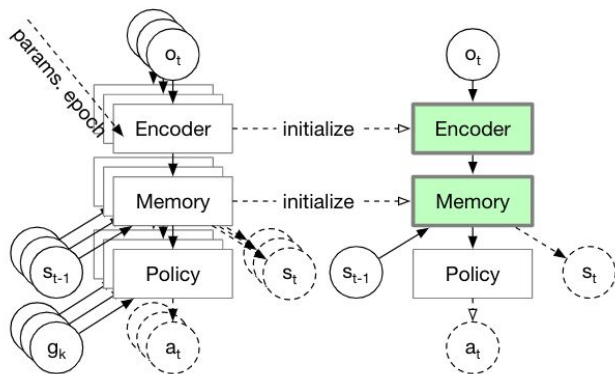
BLUE Player: Indifferent.

RED Player: Indifferent.

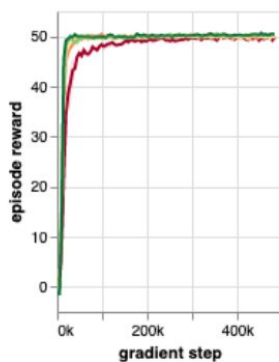


Running-with-Scissors

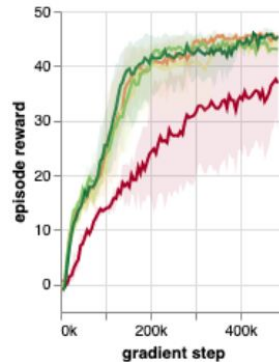
Transfer from NeuPL agent at different epochs to MPO agents



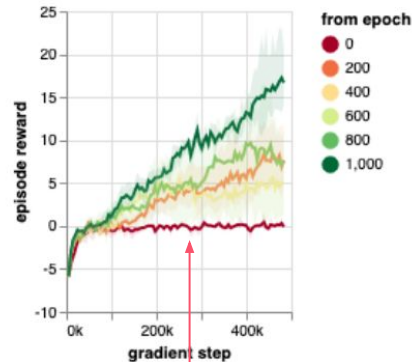
Against Nash with $n=2$



Against Nash with $n=4$



Against Nash with $n=7$



Without transfer learning, agent failed to exploit strong opponents!

Conclusion & Future Works

- **Game-Theoretic:** preserves convergence guarantees to NE (see Appendix C for proofs);
- **Transfer of skills** across strategies: learning the $N+1^{\text{th}}$ strategy becomes easier;
- **N-step Best-Response:** NeuPL yields a sequence of N -step best-responses instead of “good-”responses;
- **Efficient & Scalable:** represents a population of strategies within a single conditional network, using the compute resources of *self-play*.
- **Future Works:**
 - Beyond symmetric zero-sum games.

Algorithm 1 Neural Population Learning (Ours)

```
1:  $\Pi_\theta(\cdot|s, \sigma)$   $\triangleright$  Conditional neural population net.
2:  $\Sigma := \{\sigma_i\}_{i=1}^N$   $\triangleright$  Initial interaction graph.
3:  $\mathcal{F} : \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{N \times N}$   $\triangleright$  Meta-graph solver.
4: while true do
5:    $\Pi_\theta^\Sigma \leftarrow \{\Pi_\theta(\cdot|s, \sigma_i)\}_{i=1}^N$   $\triangleright$  Neural population.
6:   for  $\sigma_i \in \text{UNIQUE}(\Sigma)$  do
7:      $\Pi_\theta^{\sigma_i} \leftarrow \Pi_\theta(\cdot|s, \sigma_i)$ 
8:      $\Pi_\theta^{\sigma_i} \leftarrow \text{ABR}(\Pi_\theta^{\sigma_i}, \sigma_i, \Pi_\theta^\Sigma)$   $\triangleright$  Self-play.
9:    $\mathcal{U} \leftarrow \text{EVAL}(\Pi_\theta^\Sigma)$   $\triangleright$  (Optional) if  $\mathcal{F}$  adaptive.
10:   $\Sigma \leftarrow \mathcal{F}(\mathcal{U})$   $\triangleright$  (Optional) if  $\mathcal{F}$  adaptive.
11: return  $\Pi_\theta, \Sigma$ 
```
