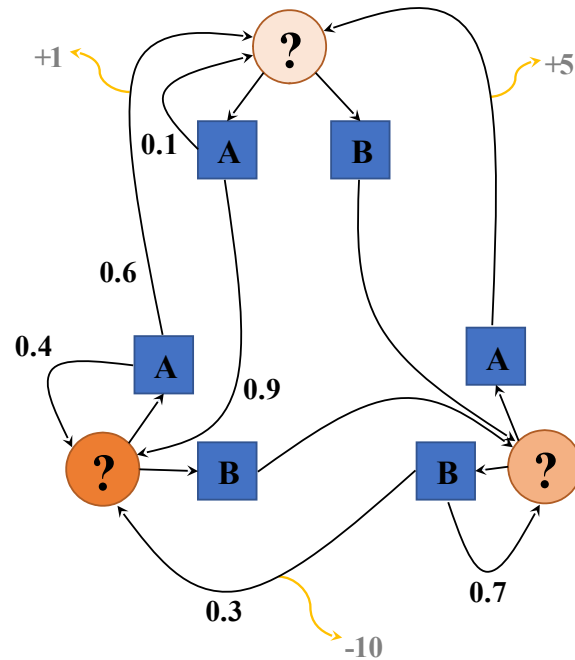


# Incomplete Information Dynamic Games

# Incomplete Information



# Partially Observable Markov Decision Process (POMDP)



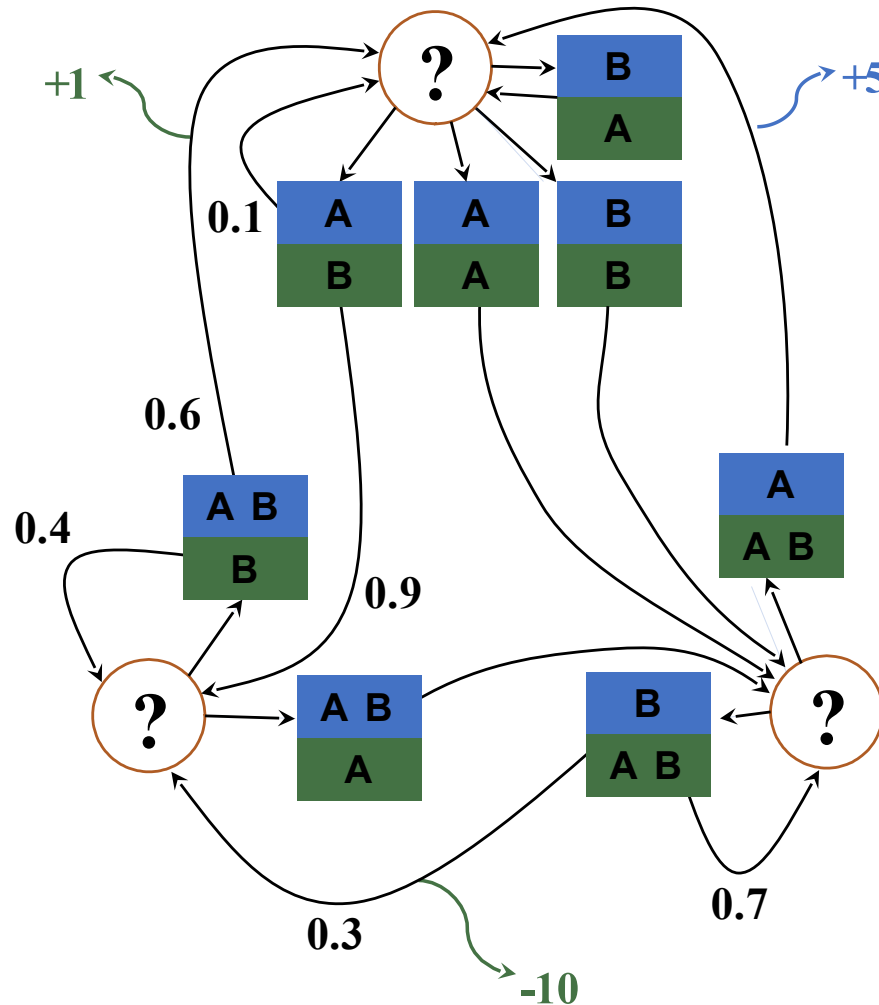
- $\mathcal{S}$  - State space
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  - Transition probability distribution
- $\mathcal{A}$  - Action space
- $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  - Reward
- $\mathcal{O}$  - Observation space
- $Z : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{O} \rightarrow \mathbb{R}$  - Observation probability distribution

Alleatory

Epistemic (Static)

Epistemic (Dynamic)

# Partially Observable Markov Game



- $\mathcal{S}$  - State space
- $T(s' | s, \mathbf{a})$  - Transition probability distribution
- $\mathcal{A}^i, i \in 1..k$  - Action spaces
- $R^i(s, \mathbf{a})$  - Reward function
- $\mathcal{O}^i, i \in 1..k$  - Observation space
- $Z(o^i | \mathbf{a}, s')$  - Observation probability distribution

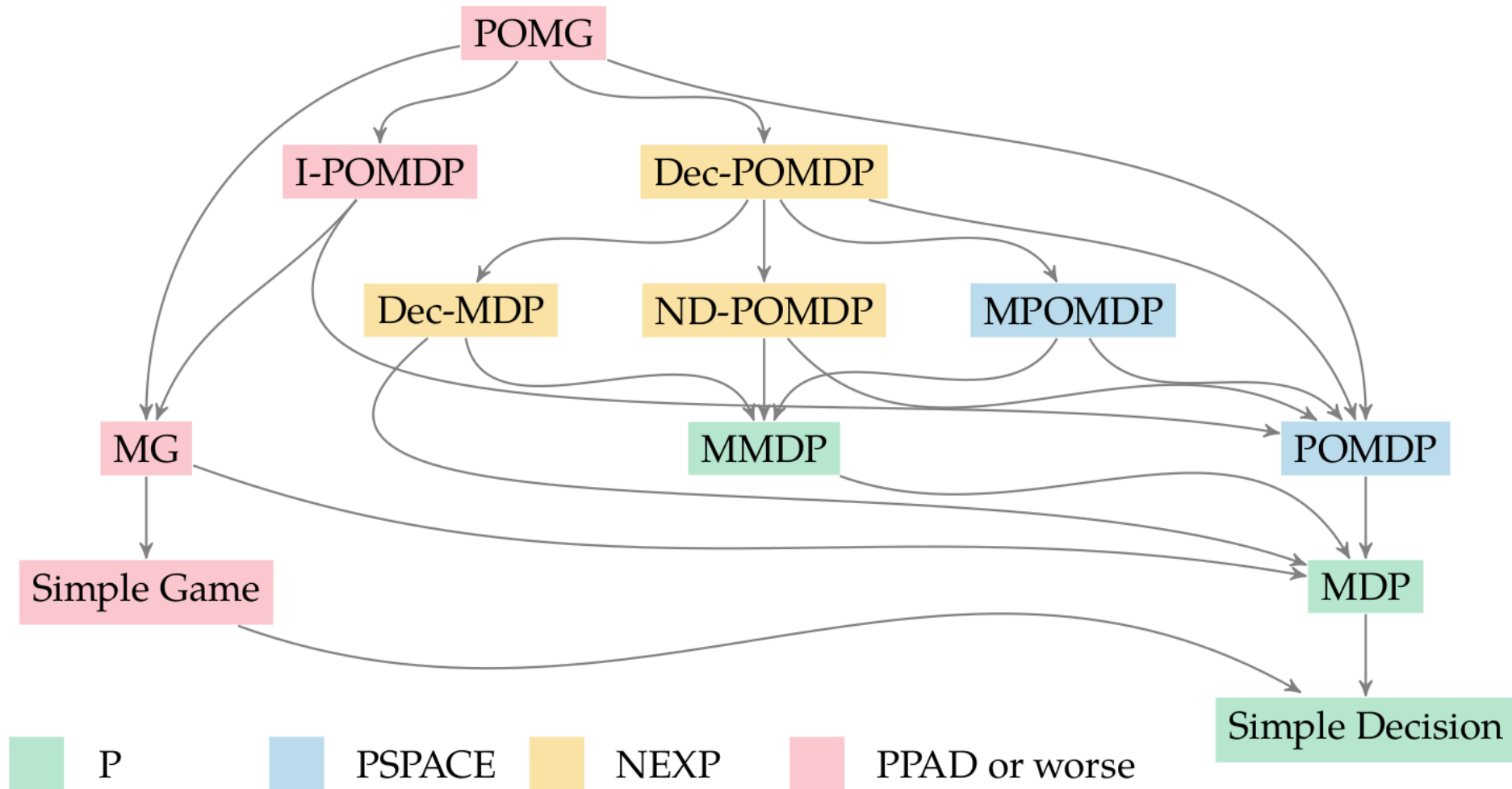
Alleatory

Epistemic (Static)

Epistemic (Dynamic)

Interaction

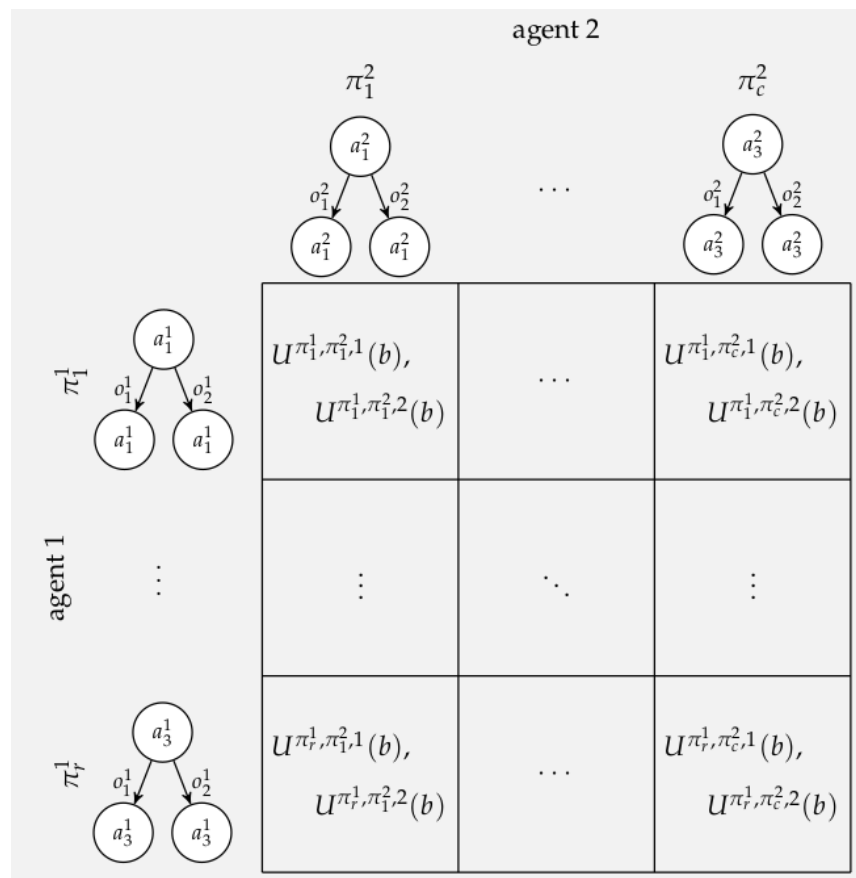
# Hierarchy of Problems



# Belief updates?

# Reduction to Simple Game

# Reduction to Simple Game





# Pruning in Dynamic Programming

$$\sum_{\pi^{-i}} \sum_s b(\pi^{-i}, s) U^{\pi^{i'}, \pi^{-i}, i}(s) \geq \sum_{\pi^{-i}} \sum_s b(\pi^{-i}, s) U^{\pi^i, \pi^{-i}, i}(s)$$

maximize  $\delta$   
 $\delta, b$

subject to  $b(\pi^{-i}, s) \geq 0$  for all  $\pi^{-i}, s$

$$\sum_{\pi^{-i}} \sum_s b(\pi^{-i}, s) = 1$$

$$\sum_{\pi^{-i}} \sum_s b(\pi^{-i}, s) \left( U^{\pi^{i'}, \pi^{-i}, i}(s) - U^{\pi^i, \pi^{-i}, i}(s) \right) \geq \delta \text{ for all } \pi^{i'}$$

# Extensive Form Game

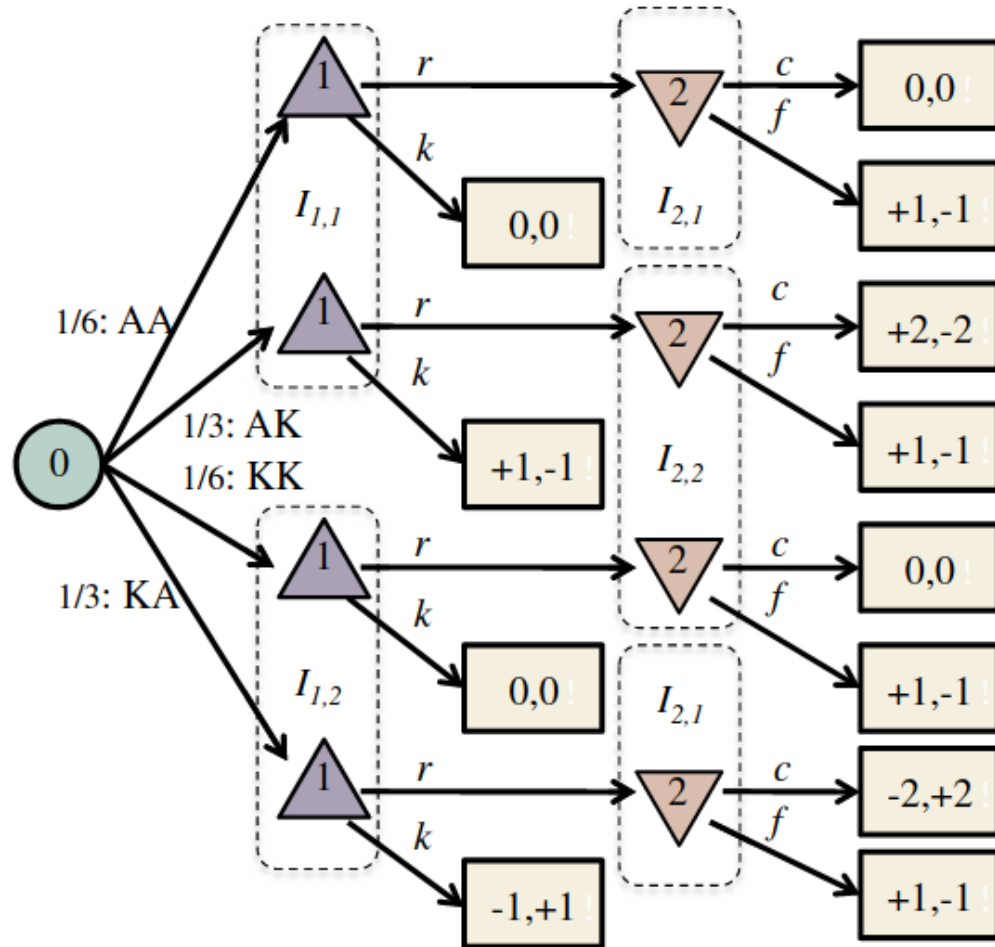
(Alternative to POMGs that is more common in the literature)

- Similar to a minimax tree for a turn-taking game
- Chance nodes
- Information sets

# King-Ace Poker Example

- 4 Cards: 2 Aces, 2 Kings
- Each player is dealt a card
- P1 can either *raise* ( $r$ ) the payoff to 2 points or *check* ( $k$ ) the payoff at 1 point
- If P1 raises, P2 can either *call* ( $c$ ) Player 1's bet, or *fold* ( $f$ ) the payoff back to 1 point
- The highest card wins

# Extensive to Matrix Form

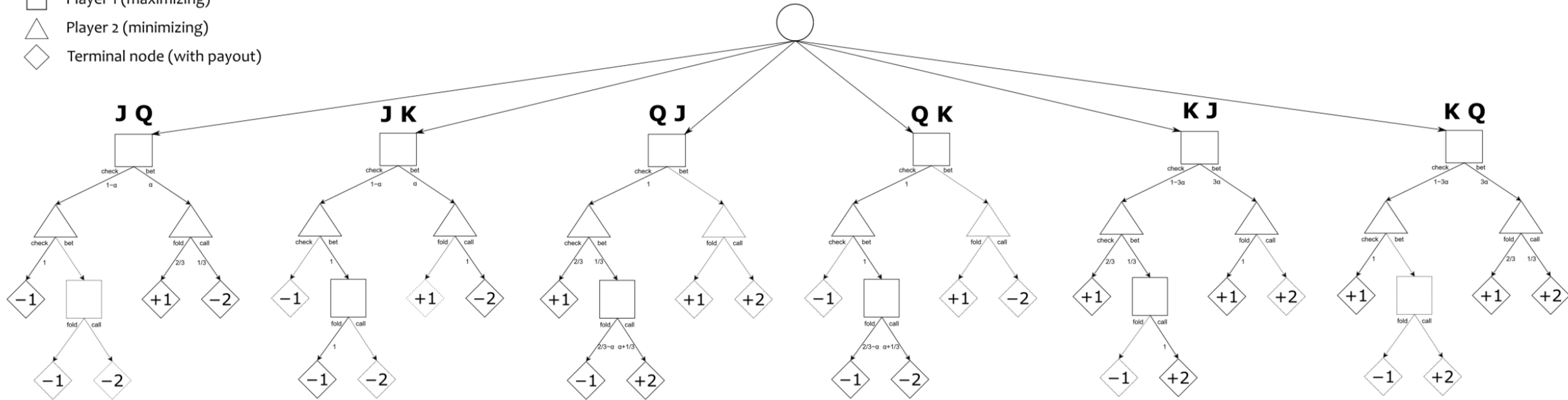


	2:cc	2:cf	2:ff	2:fc
1:rr	0	-1/6	1	7/6
1:kr	-1/3	-1/6	5/6	2/3
1:rk	1/3	<b>0</b>	1/6	1/2
1:kk	0	<b>0</b>	0	0

Exponential in number of info states!

# A more interesting example: Kuhn Poker

- Chance node (initial deal)
- Player 1 (maximizing)
- △ Player 2 (minimizing)
- ◇ Terminal node (with payout)



# Fictitious Play in Extensive Form Games

**Algorithm 2** General Fictitious Self-Play

**function** FICTITIOUSSELFPLAY( $\Gamma, n, m$ )

Initialize completely mixed  $\pi_1$

$\beta_2 \leftarrow \pi_1$

$j \leftarrow 2$

**while** within computational budget **do**

$\eta_j \leftarrow \text{MIXINGPARAMETER}(j)$

$\mathcal{D} \leftarrow \text{GENERATEDATA}(\pi_{j-1}, \beta_j, n, m, \eta_j)$

**for** each player  $i \in \mathcal{N}$  **do**

$\mathcal{M}_{RL}^i \leftarrow \text{UPDATERLMEMORY}(\mathcal{M}_{RL}^i, \mathcal{D}^i)$

$\mathcal{M}_{SL}^i \leftarrow \text{UPDATESLMEMORY}(\mathcal{M}_{SL}^i, \mathcal{D}^i)$

$\beta_{j+1}^i \leftarrow \text{REINFORCEMENTLEARNING}(\mathcal{M}_{RL}^i)$

$\pi_j^i \leftarrow \text{SUPERVISEDLEARNING}(\mathcal{M}_{SL}^i)$

**end for**

$j \leftarrow j + 1$

**end while**

**return**  $\pi_{j-1}$

**end function**

**function** GENERATEDATA( $\pi, \beta, n, m, \eta$ )

$\sigma \leftarrow (1 - \eta)\pi + \eta\beta$

$\mathcal{D} \leftarrow n$  episodes  $\{t_k\}_{1 \leq k \leq n}$ , sampled from strategy profile  $\sigma$

**for** each player  $i \in \mathcal{N}$  **do**

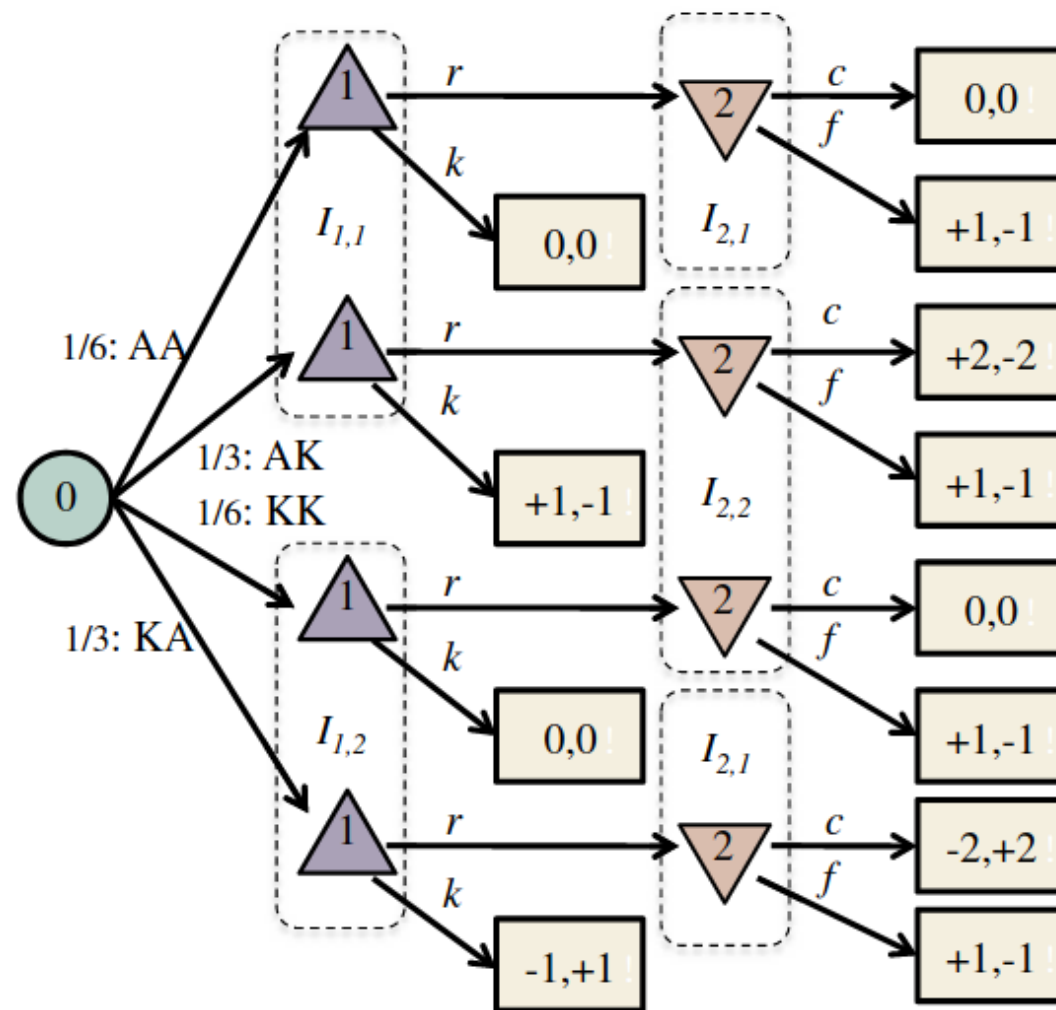
$\mathcal{D}^i \leftarrow m$  episodes  $\{t_k^i\}_{1 \leq k \leq m}$ , sampled from strategy profile  $(\beta^i, \sigma^{-i})$

$\mathcal{D}^i \leftarrow \mathcal{D}^i \cup \mathcal{D}$

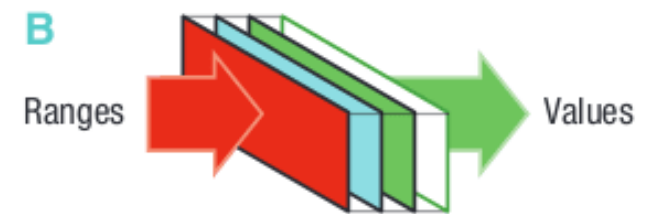
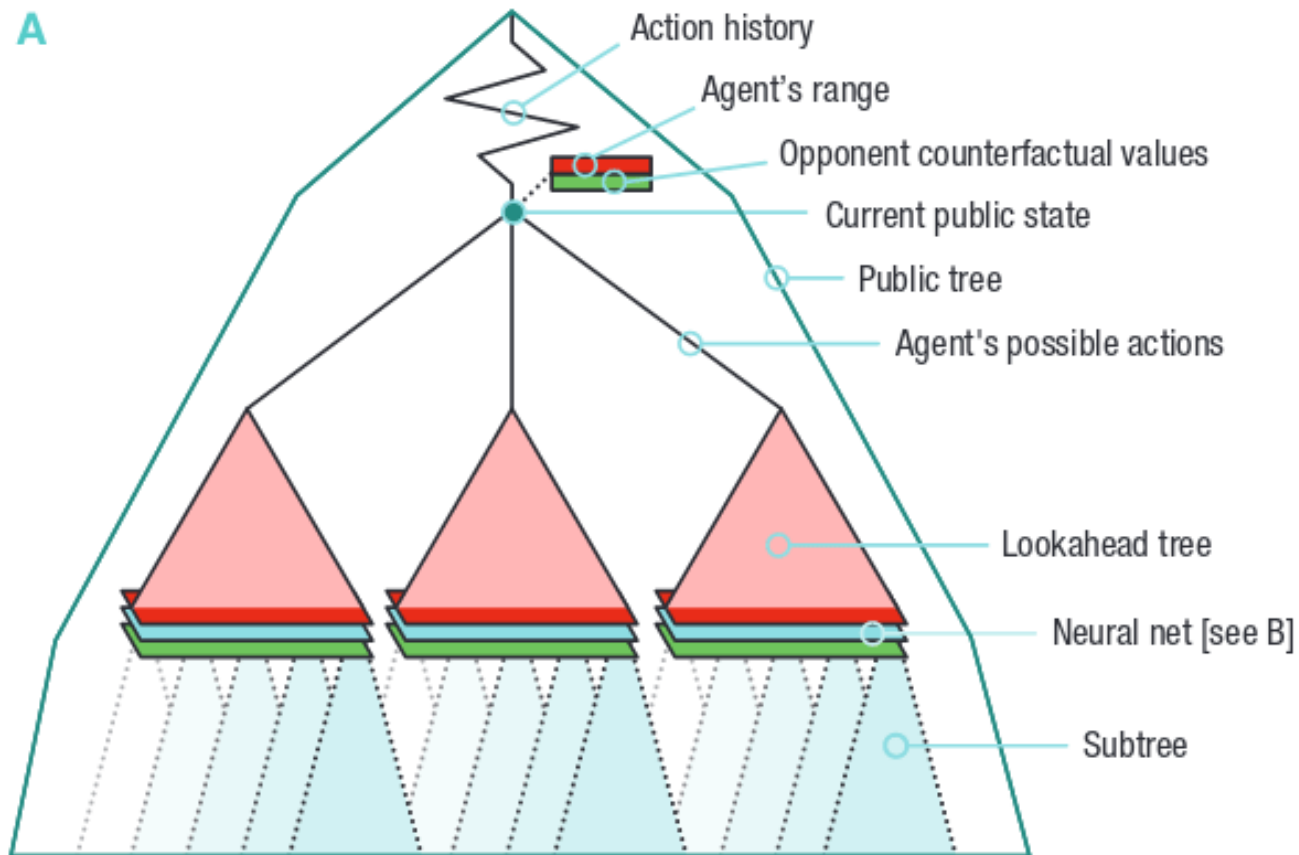
**end for**

**return**  $\{\mathcal{D}^k\}_{1 \leq k \leq N}$

**end function**

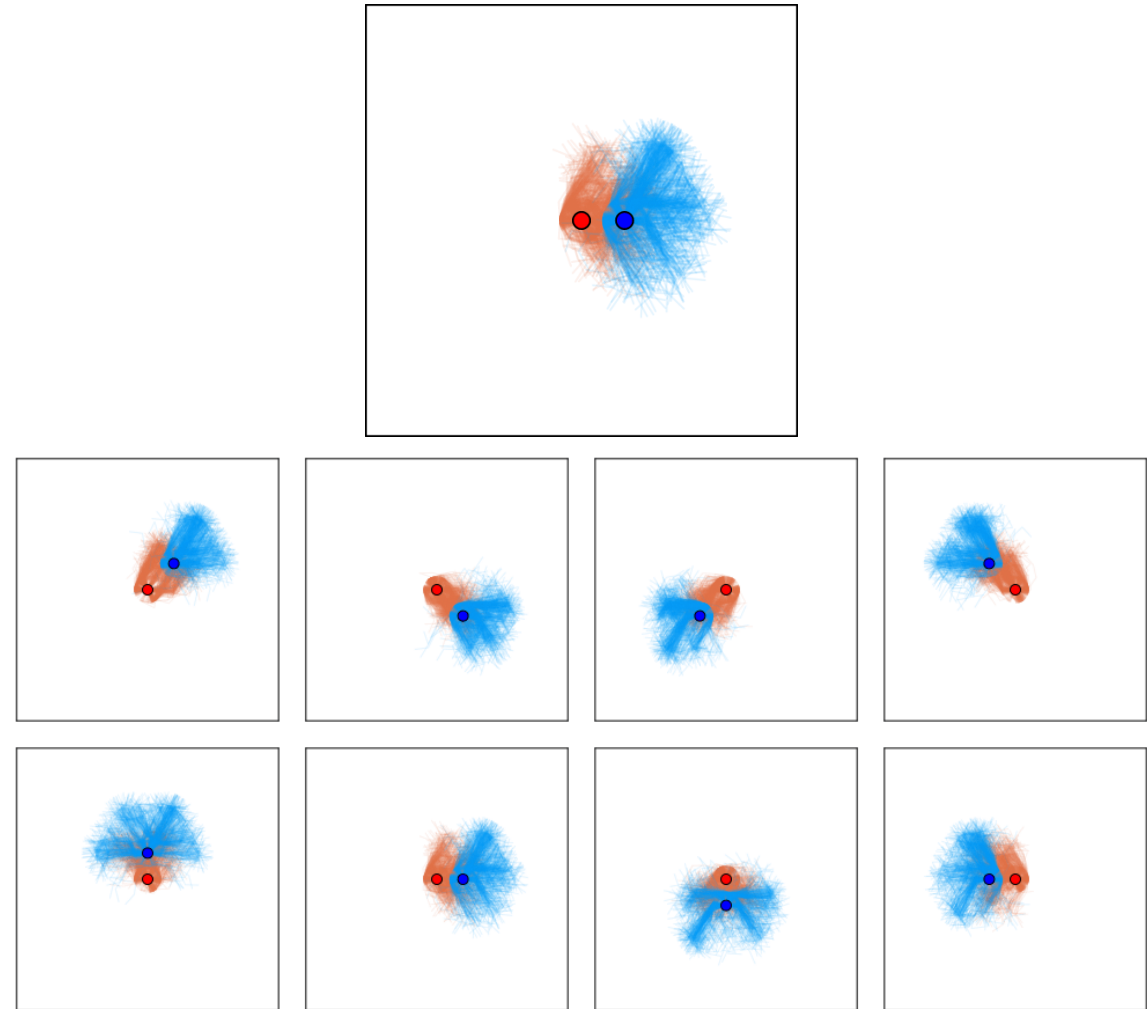
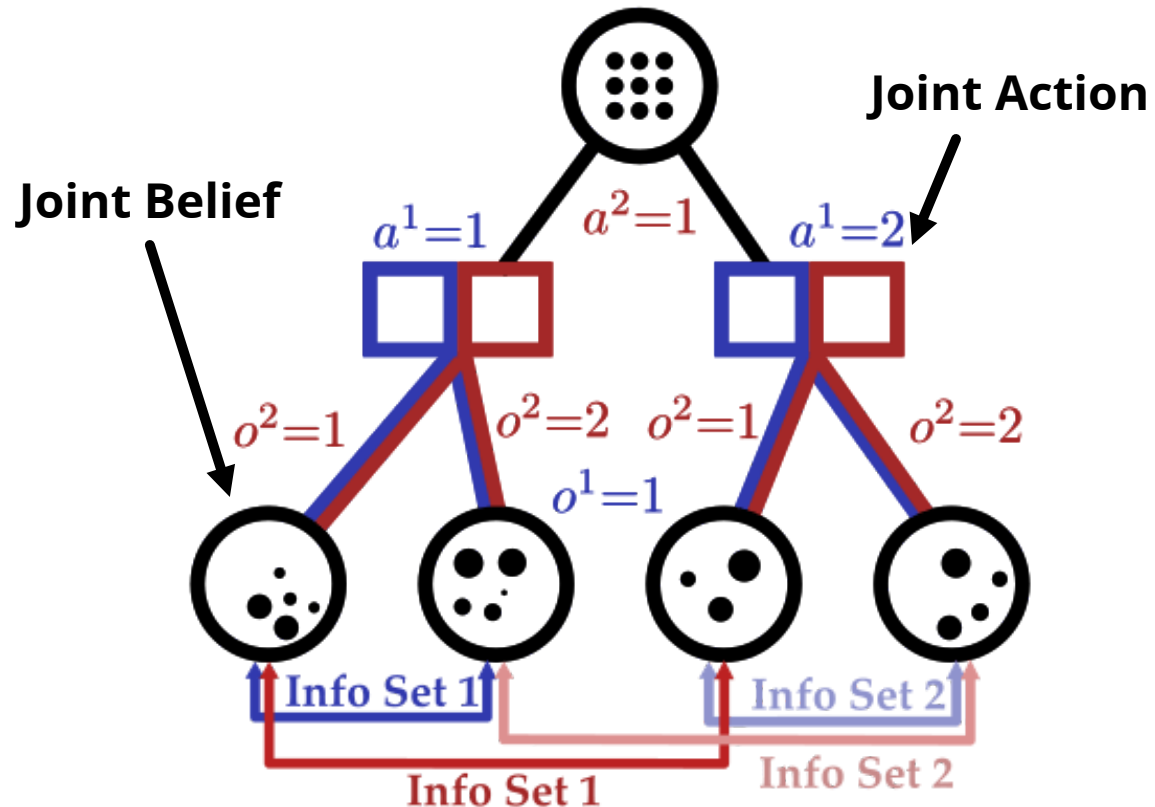


# Deep Stack: Scaling to Heads Up No Limit Texas Hold 'Em



# Tree-Based Planning in POSGs

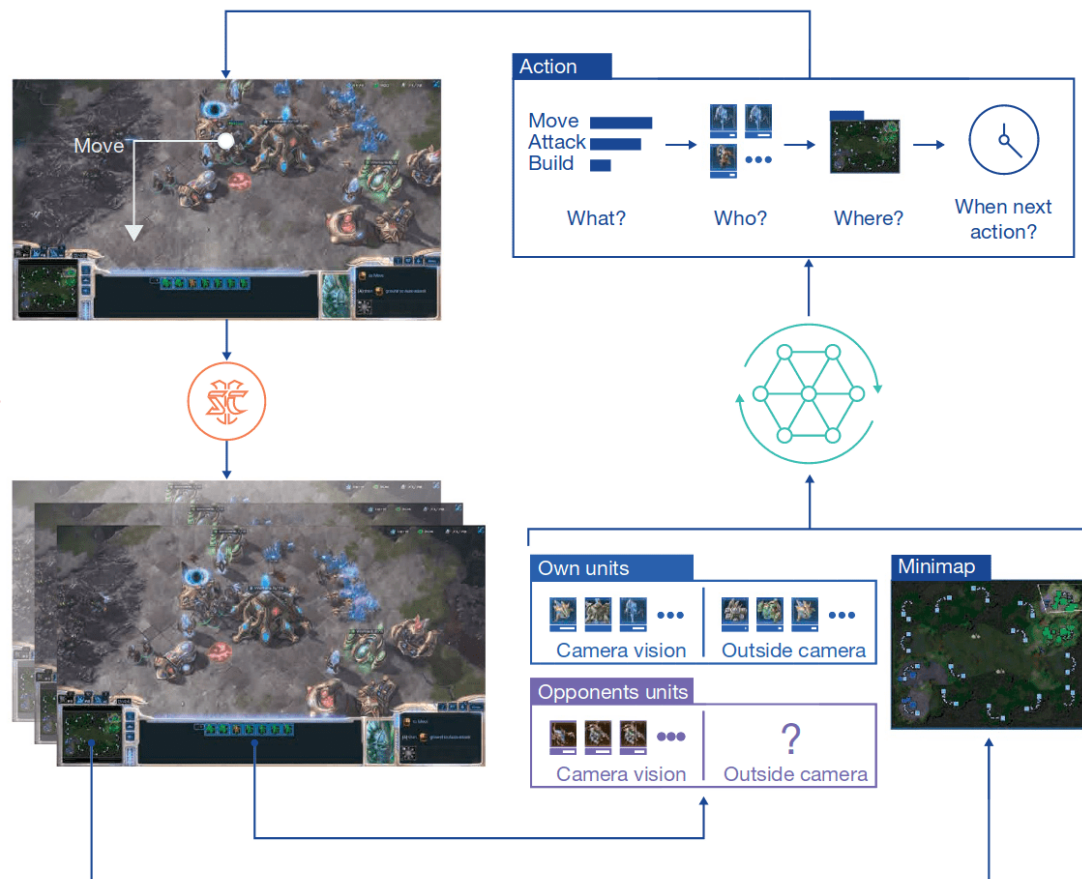
Our approach: combine particle filtering and information sets



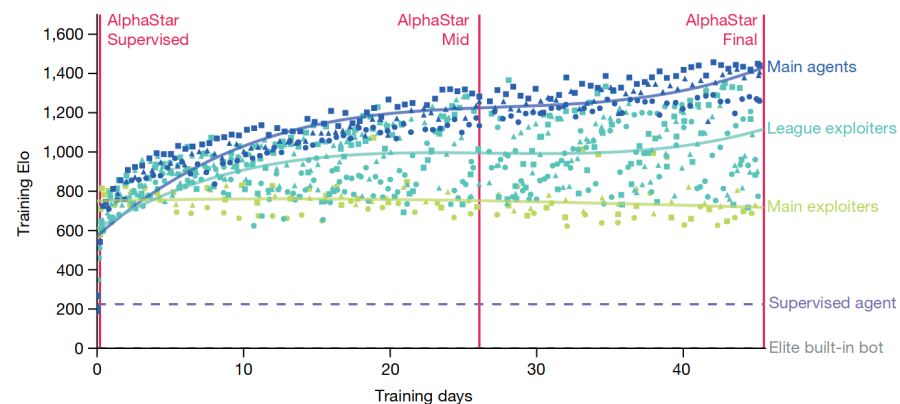
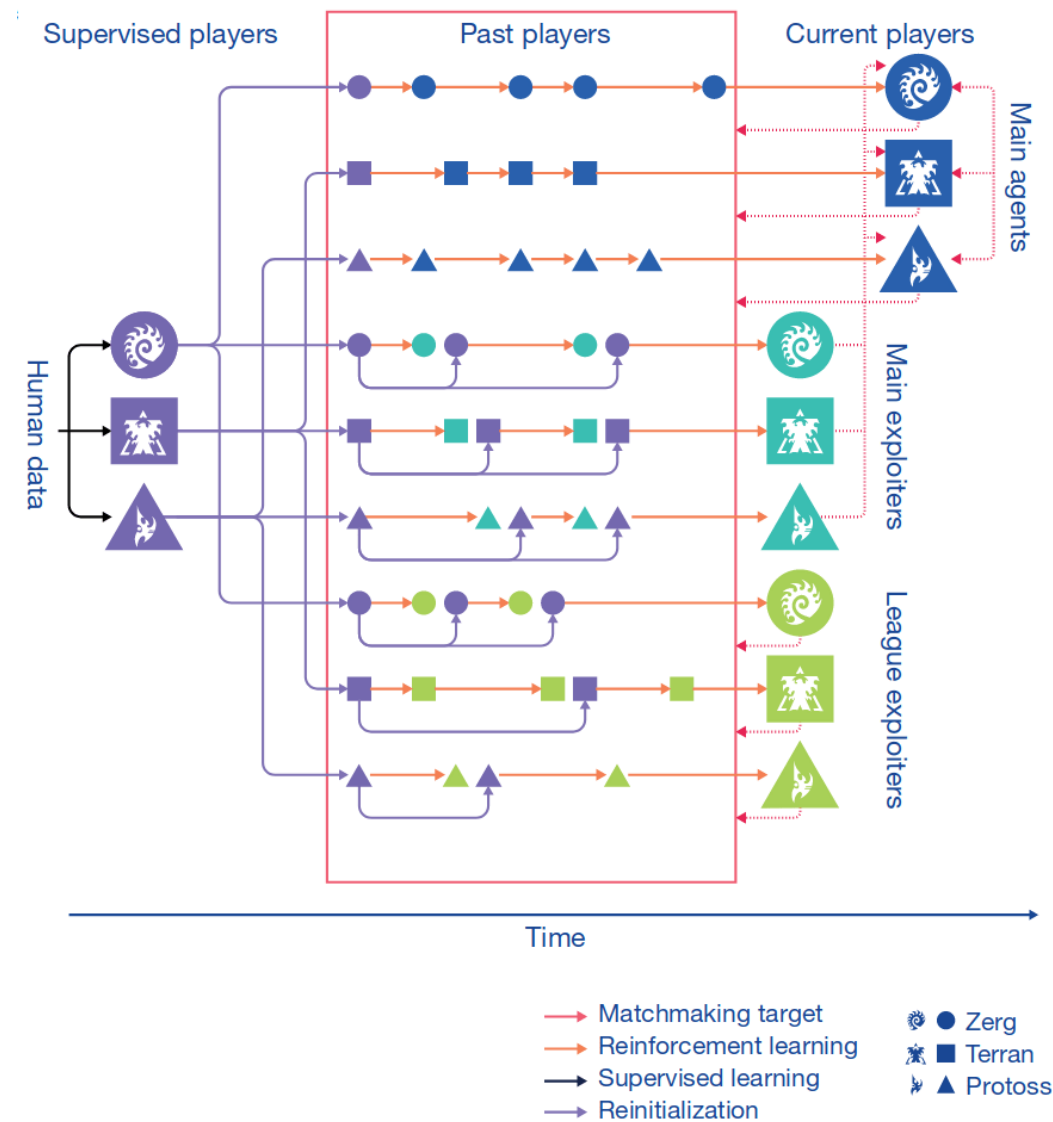


## Monitoring layer

Actions limit ~22 per 5 s  
Requested delay ~200 ms



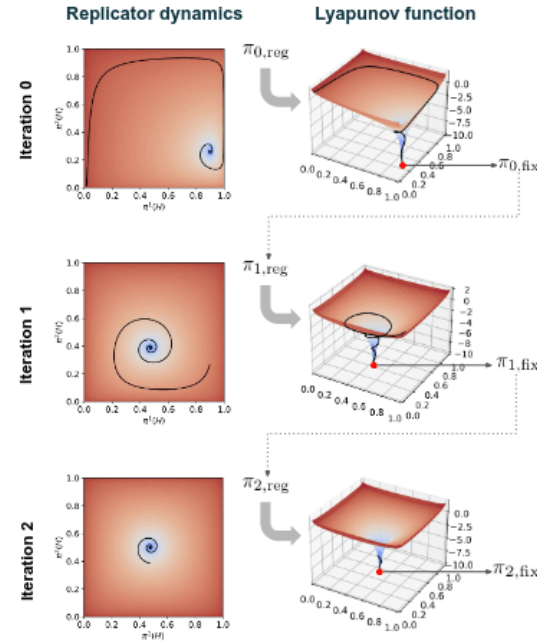
# Alpha Star



# Deep Nash

		Player 2	
		Head: $H$	Tail: $T$
Player 1	Head: $H$	1	-1
	Tail: $T$	-1	1

(a) Matching pennies



**R-NaD Iteration**

Start with an arbitrary regularization policy:  $\pi_{0,reg}$

- Reward transformation: Construct the transformed game with:  $\pi_{n,reg}$
- Dynamics: Run the replicator dynamics until convergence to:  $\pi_{n,fix}$
- Update: Set the regularization policy:  $\pi_{n+1,reg} = \pi_{n,fix}$

Repeat steps until convergence

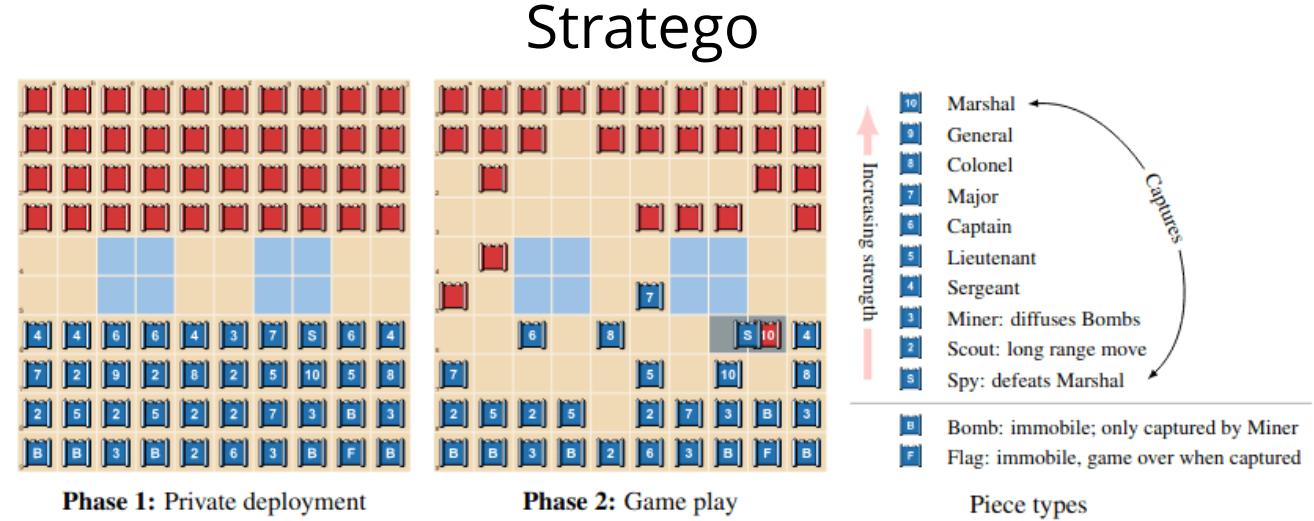
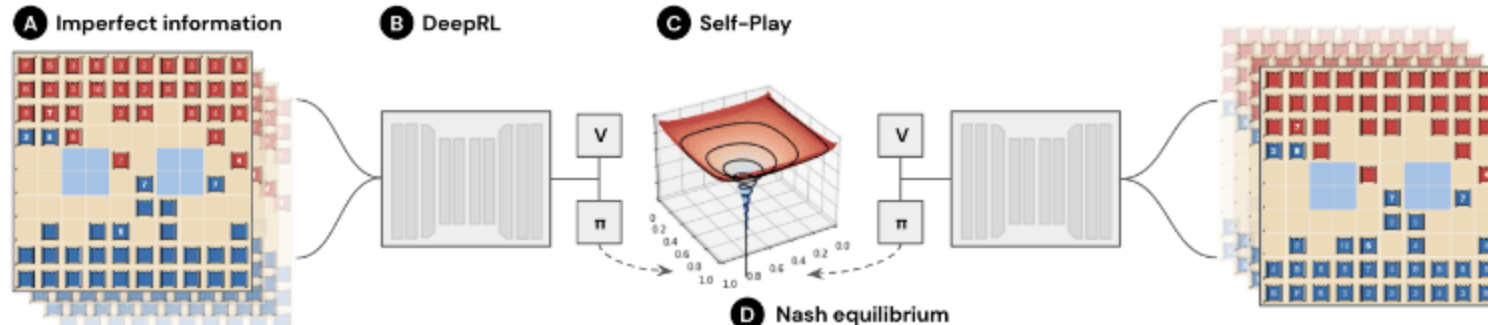


Figure 2: The R-NaD learning algorithm illustrated with the matching pennies game



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

**Reward transformation:**  $r^i(\pi^i, \pi^{-i}, a^i, a^{-i}) = r^i(a^i, a^{-i}) - \eta \log \left( \frac{\pi^i(a^i)}{\pi_{reg}^i(a^i)} \right) + \eta \log \left( \frac{\pi^{-i}(a^{-i})}{\pi_{reg}^{-i}(a^{-i})} \right)$