

Imperfect

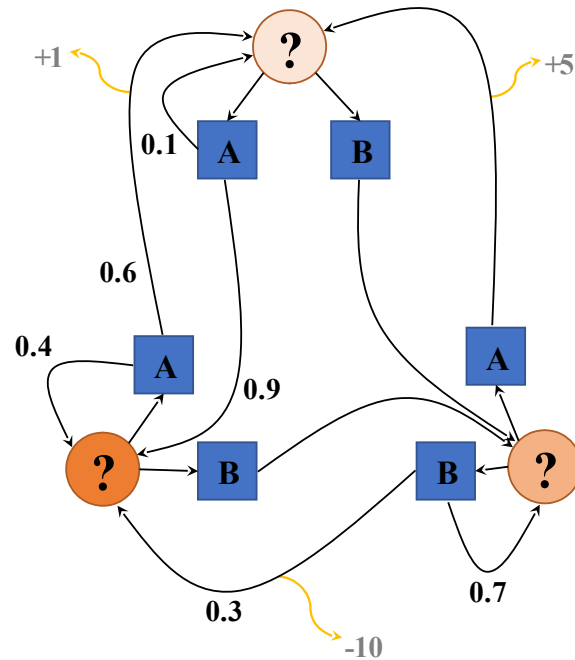
~~Incomplete~~ Information Dynamic Games

Imperfect ~~Incomplete~~ Information



Partially Observable Markov Decision Process (POMDP)

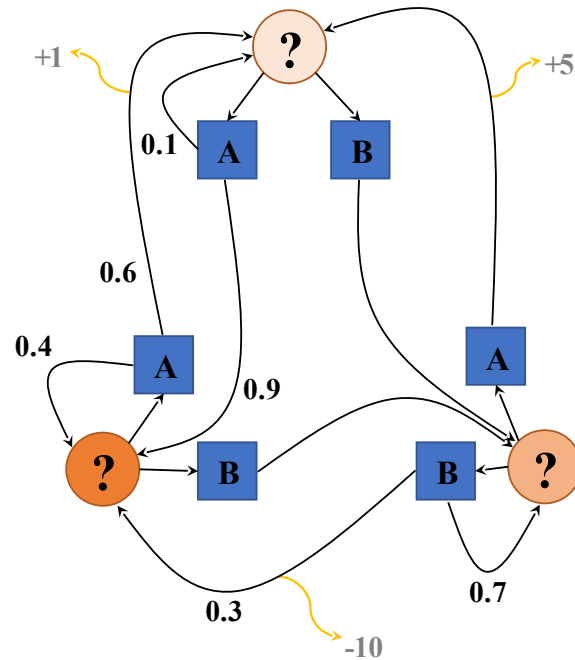
- \mathcal{S} - State space
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ - Transition probability distribution
- \mathcal{A} - Action space
- $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ - Reward



Alleatory

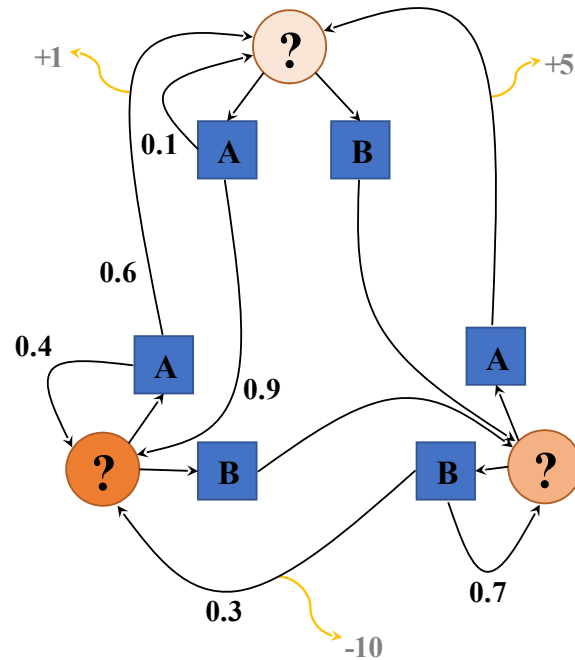
Partially Observable Markov Decision Process (POMDP)

- \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



Alleatory

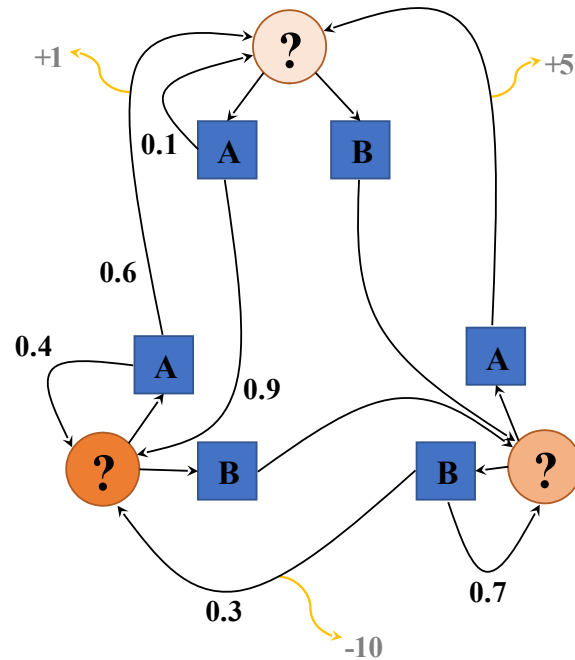
Partially Observable Markov Decision Process (POMDP)



- \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

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

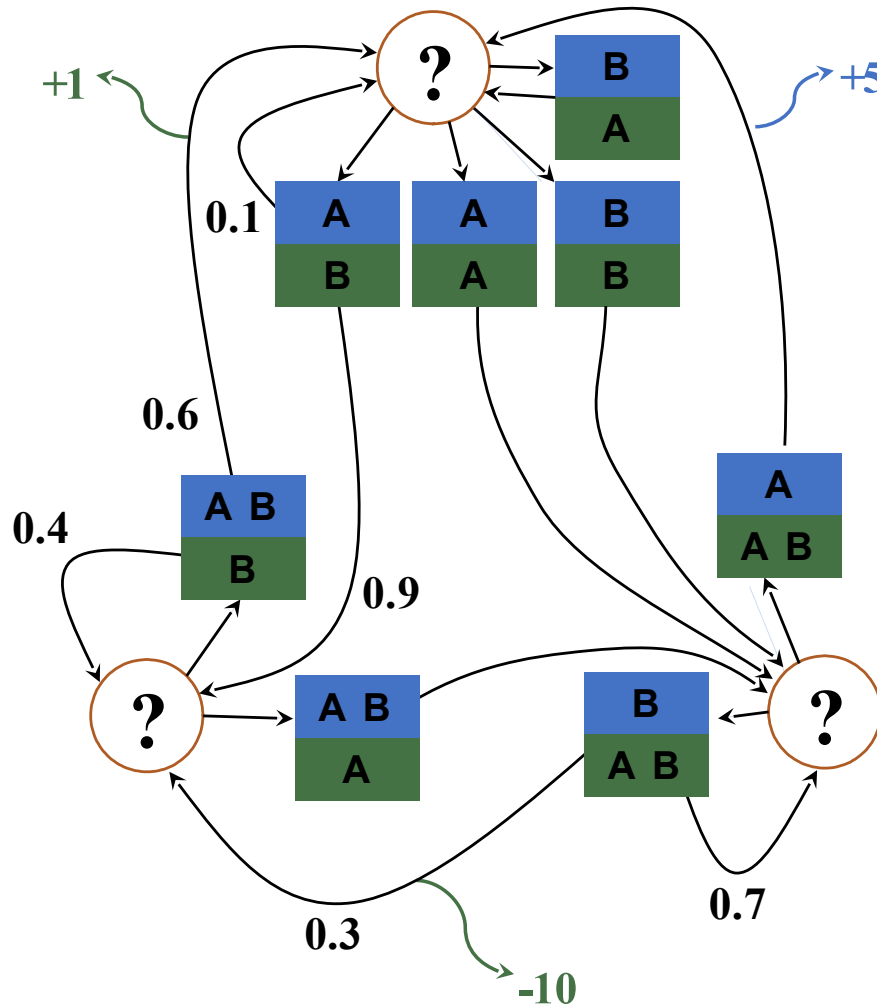
Alleatory

Epistemic (Static)

Epistemic (Dynamic)

Partially Observable Markov Game

Stochastic



- \mathcal{S} - State space
- $T(s' | s, \mathbf{a})$ - Transition probability distribution
joint action
- $\mathcal{A}^i, i \in 1..k$ - 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

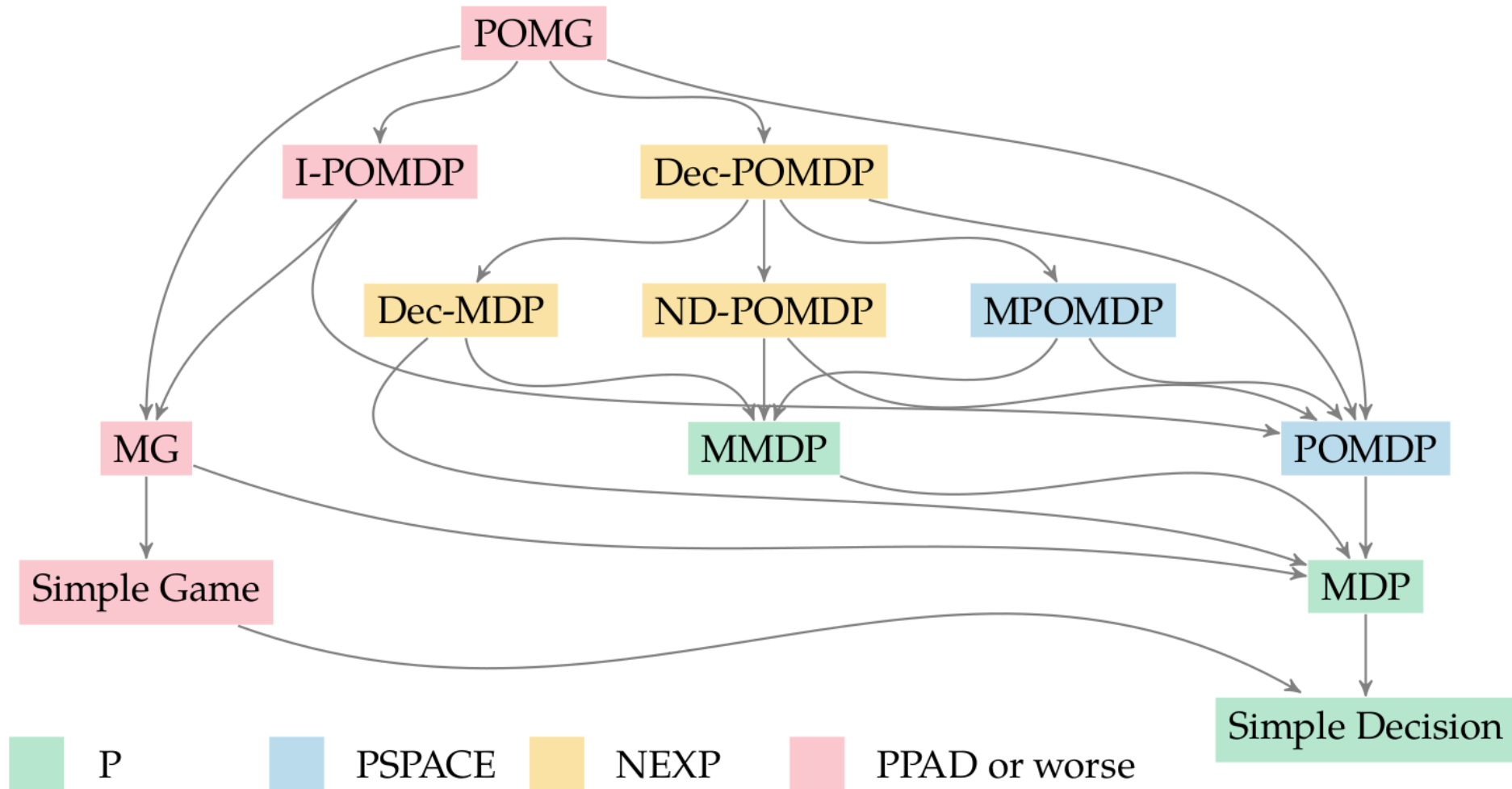
Alleatory

Epistemic (Static)

Epistemic (Dynamic)

Interaction

Hierarchy of Problems



Belief updates?

POMDP :

$$b'(s') \propto Z(o|a, s') \sum_{s \in S} T(s'|s, a) b(s)$$

POMG :

$$b'(s') \propto Z(o^i | \overset{\pi^i}{\vec{a}}, s') \sum_{s \in S} T(s' | s, \overset{\text{joint action}}{\vec{a}}) b(s)$$

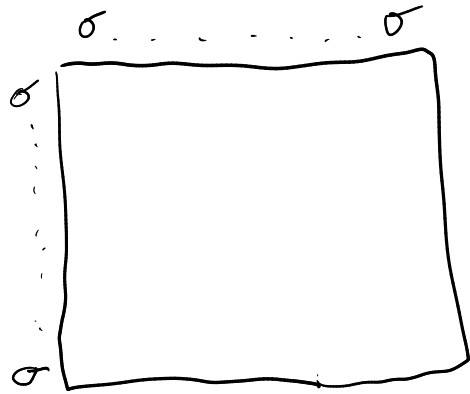
X

Problems

1. Requires reasoning about other player's actions
2. Requires reasoning about other player's observations
3. Usually trying to solve for π^i at the same time as choosing our actions

Reduction to Simple Game

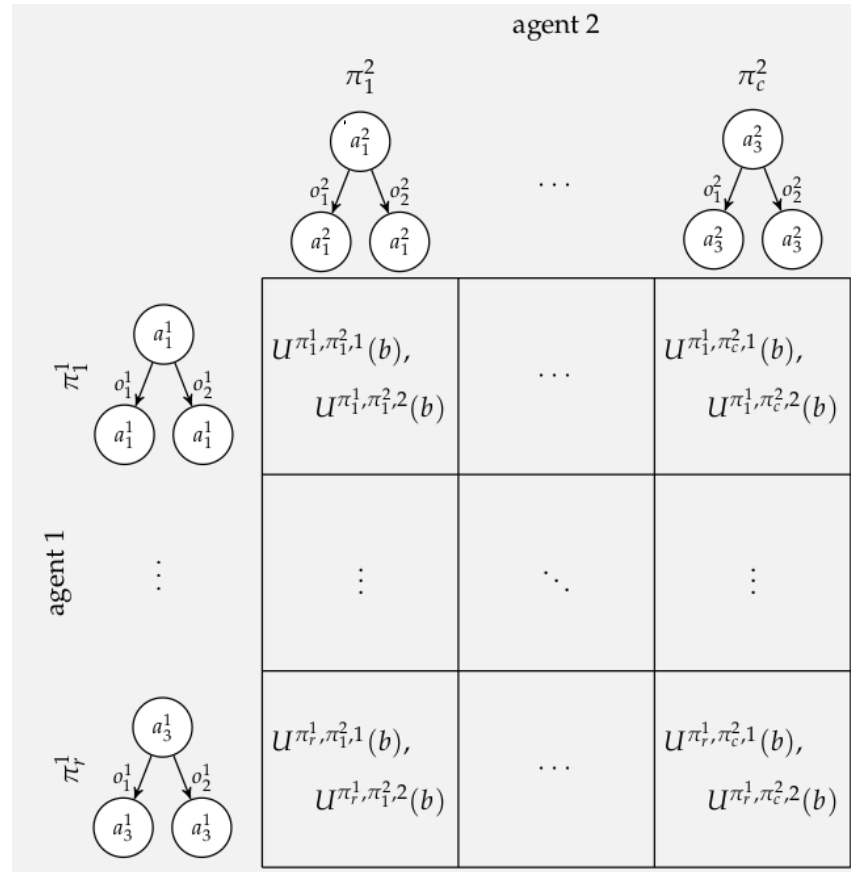
MDP



Reduction to Simple Game

Reduction to Simple Game

2 step



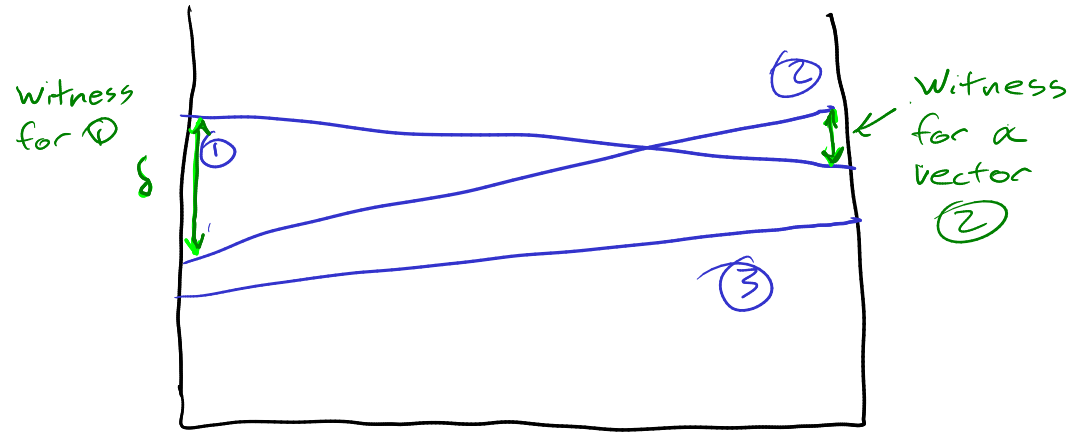
Pruning in Dynamic Programming

Pruning in Dynamic Programming

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

Pruning in Dynamic Programming

From POMDP



Dynamic Programming for POMGs

Start with all possible $N=1$ -step policies

Loop until $N=\bar{N}$

Evaluate $N+1$ step policies $\rightarrow U$

Prune Dominated Policies

Solve Matrix Game for all \bar{N} -step policies

Pruning

If there exists $\pi^{i'}$ such that

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

for all "beliefs", we can prune π^i .

$$\left[\begin{array}{l} \underset{\delta, b}{\text{maximize}} \quad \delta \\ \text{subject to} \quad \left. \begin{array}{l} b(\pi^{-i}, s) \geq 0 \text{ for all } \pi^{-i}, s \\ \sum_{\pi^{-i}} \sum_s b(\pi^{-i}, s) = 1 \end{array} \right\} \text{probability} \\ \sum_{\pi^{-i}} \sum_s b(\pi^{-i}, s) \left(U^{\pi^{i'}, \pi^{-i}, i}(s) - U^{\pi^i, \pi^{-i}, i}(s) \right) \geq \delta \text{ for all } \pi^{i'} \end{array} \right]$$

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

King-Ace Poker Example

- 4 Cards: 2 Aces, 2 Kings

King-Ace Poker Example

- 4 Cards: 2 Aces, 2 Kings
- Each player is dealt a card

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

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

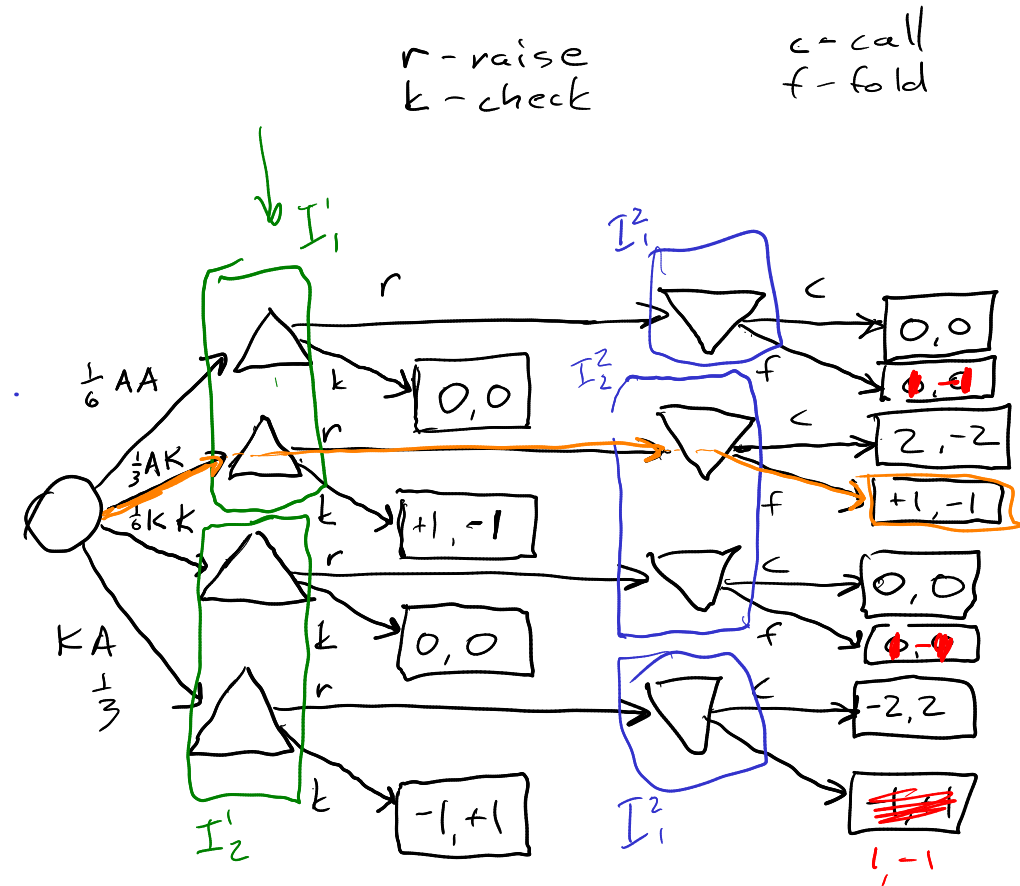
King-Ace Poker Example

Aces High

- 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

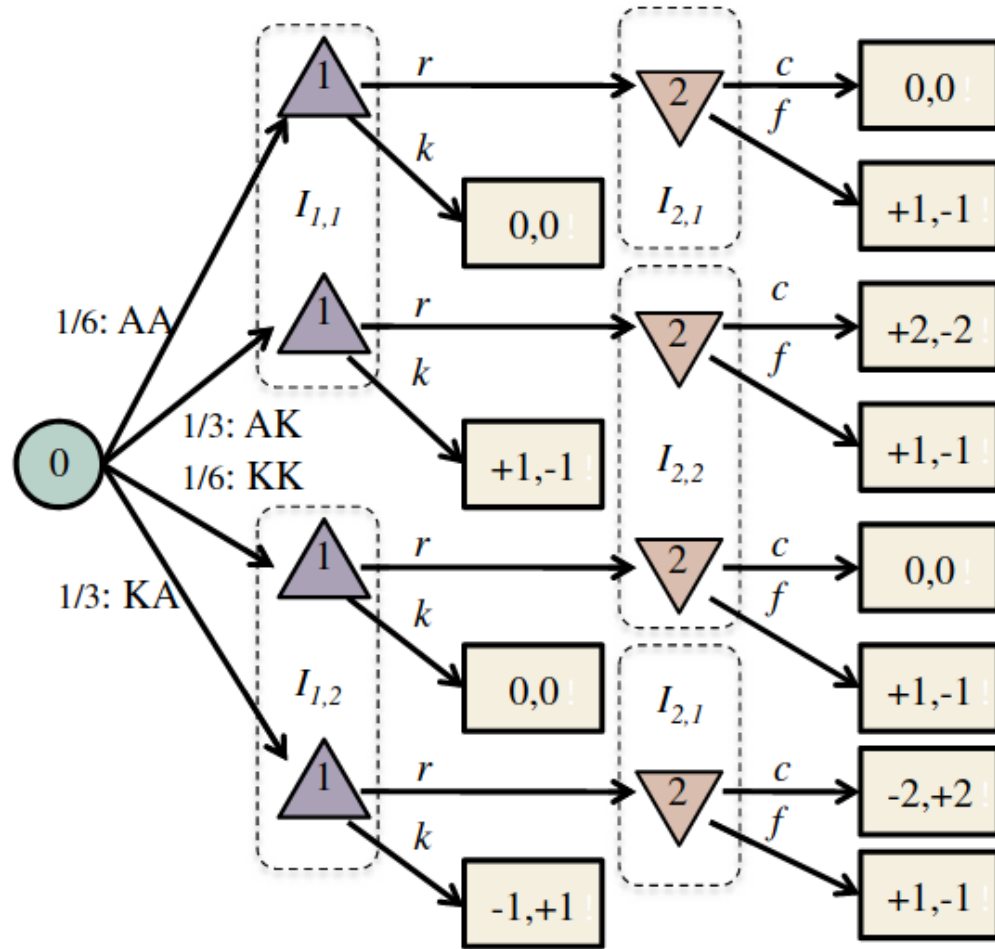
Player 1
wins 1 point

Chance
node

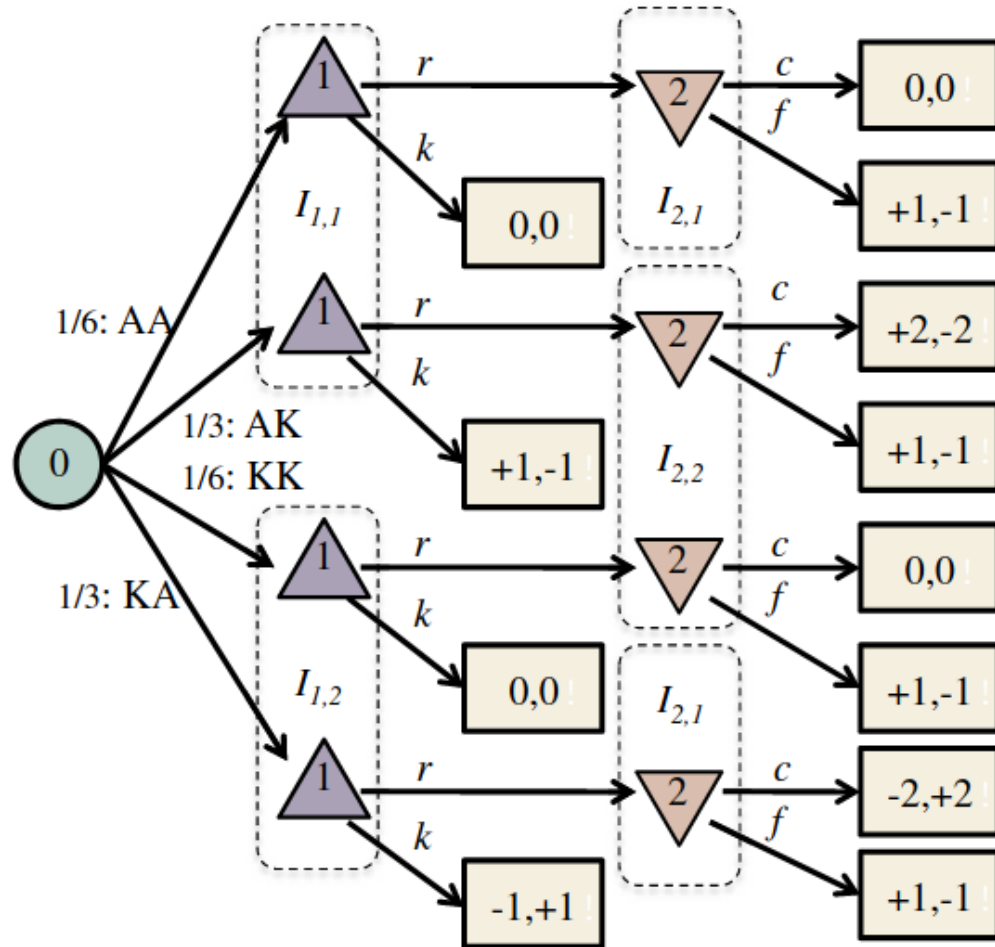


Information sets group nodes that are indistinguishable to a player

Extensive to Matrix Form



Extensive to Matrix Form



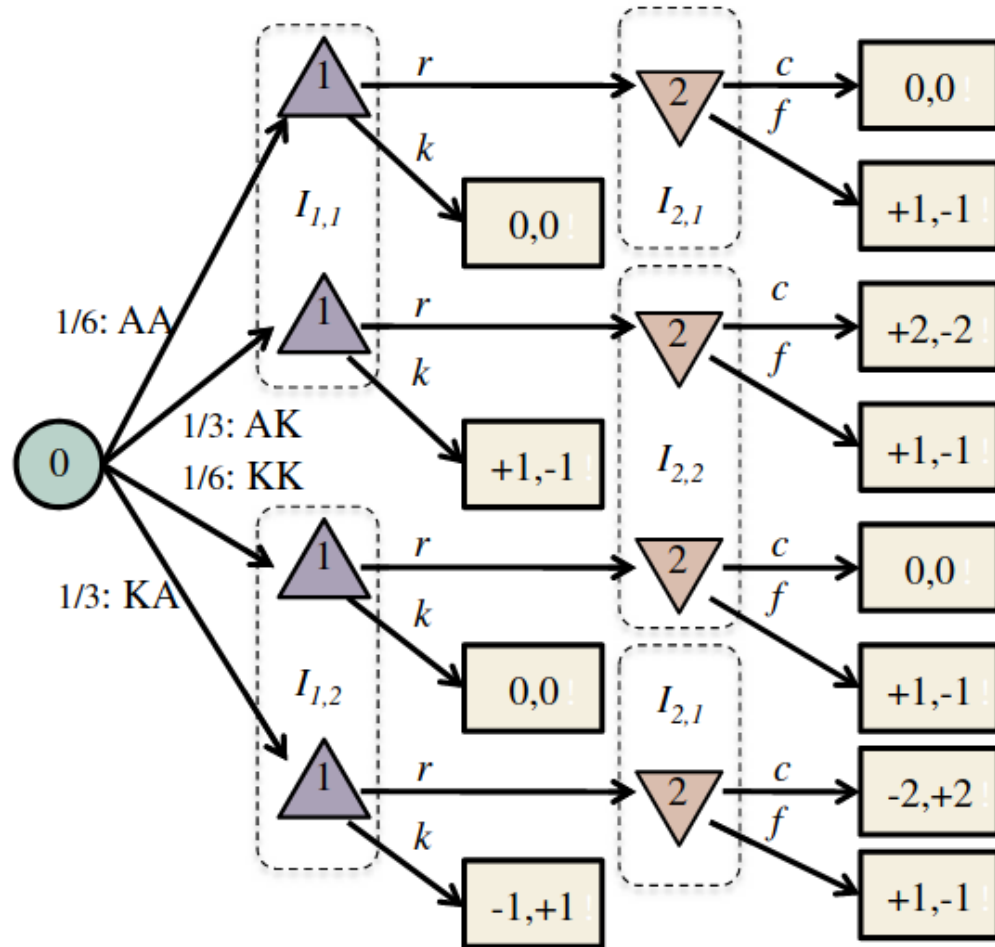
Arrows point to a best response in a row or column

Bars indicate responses are equally the best

If a cell has no arrows leaving, it is a pure NE

		\overbrace{AK}	\overbrace{AK}	\overbrace{AK}
	$2:cc$	$2:cf$	$2:ff$	$2:fc$
\overbrace{AK} $1:rr$	0	$-1/6$	1	$7/6$
\overbrace{AK} $1:kr$	$-1/3$	$-1/6$	$5/6$	$2/3$
\overbrace{AK} $1:rk$	$1/3$	0	$1/6$	$1/2$
\overbrace{AK} $1:kk$	0	0	0	0

Extensive to Matrix Form

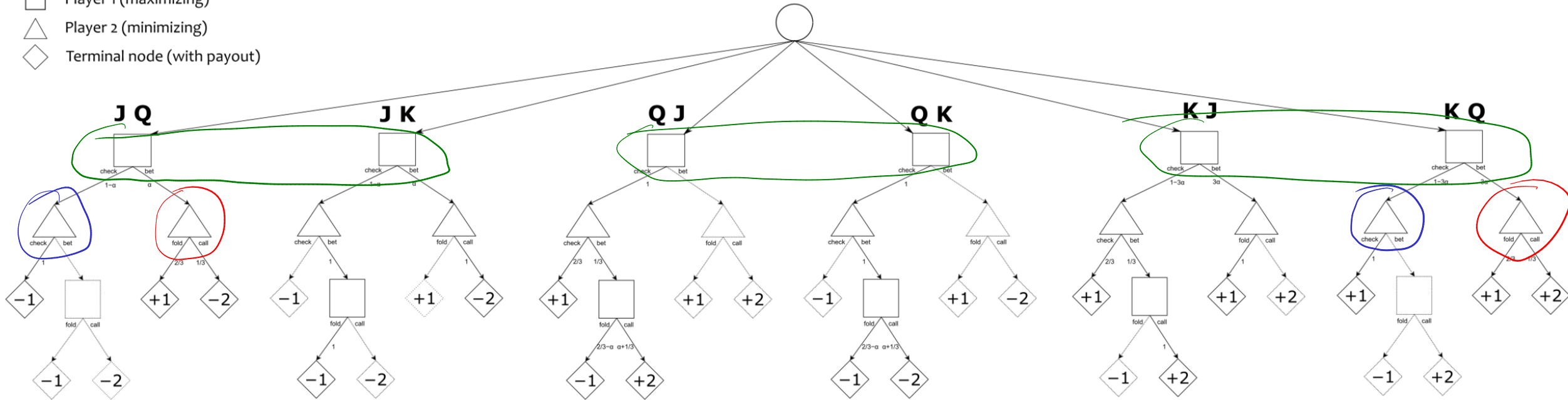


	2:cc	2:cf	2:ff	2:fc
1:rr	0	-1/6	1	7/6
1:kr	-1/3	-1/6	5/6	2/3
1:rk	1/3	0	1/6	1/2
1:kk	0	0	0	0

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(Γ, n, m)

Initialize completely mixed π_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 π_{j-1}

end function

function GENERATEDATA(π, β, n, m, η)

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

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

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

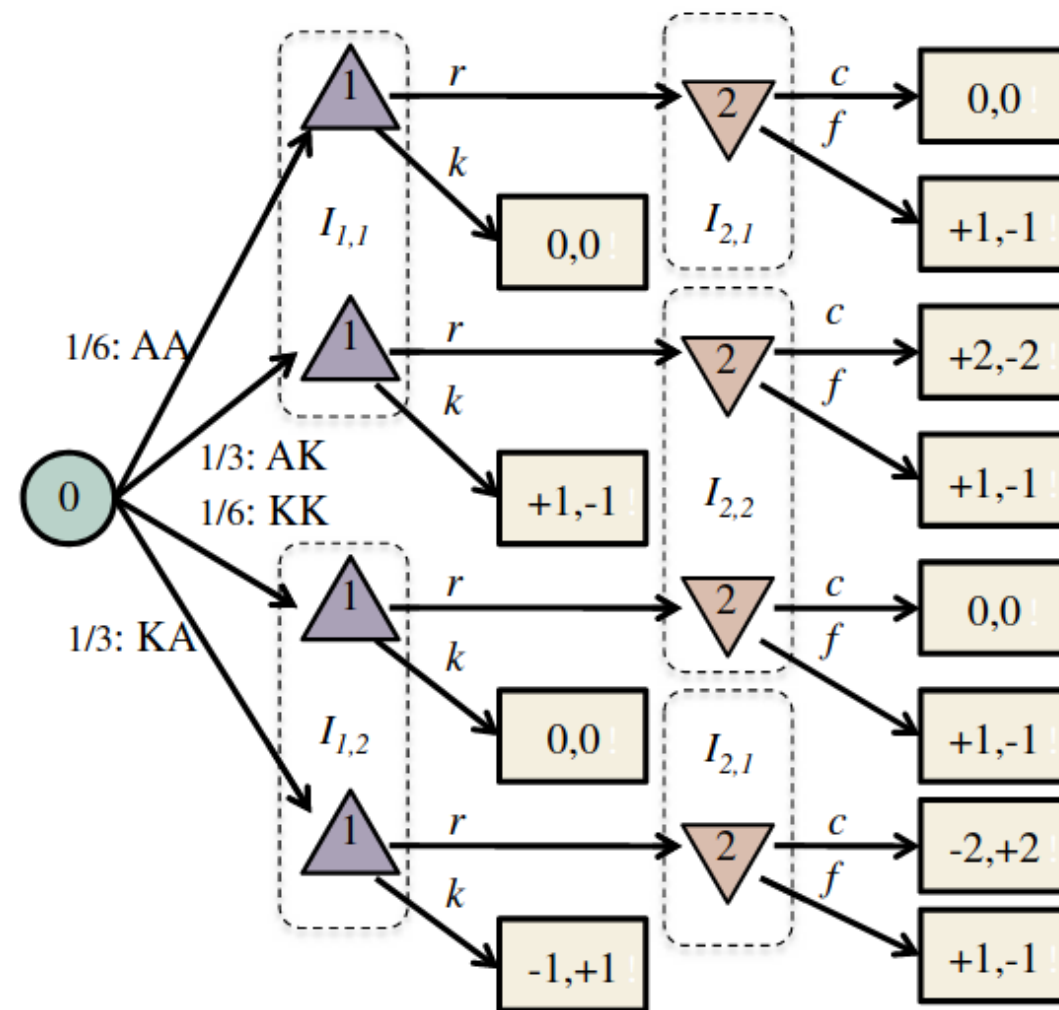
$\mathcal{D}^i \leftarrow m$ episodes $\{t_k^i\}_{1 \leq k \leq m}$, sampled from strategy profile (β^i, σ^{-i})

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

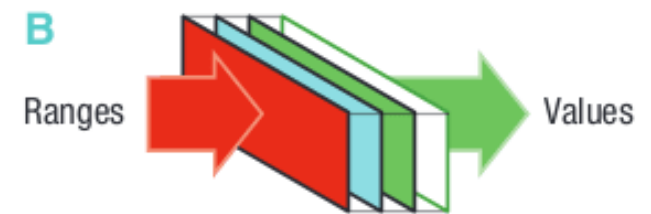
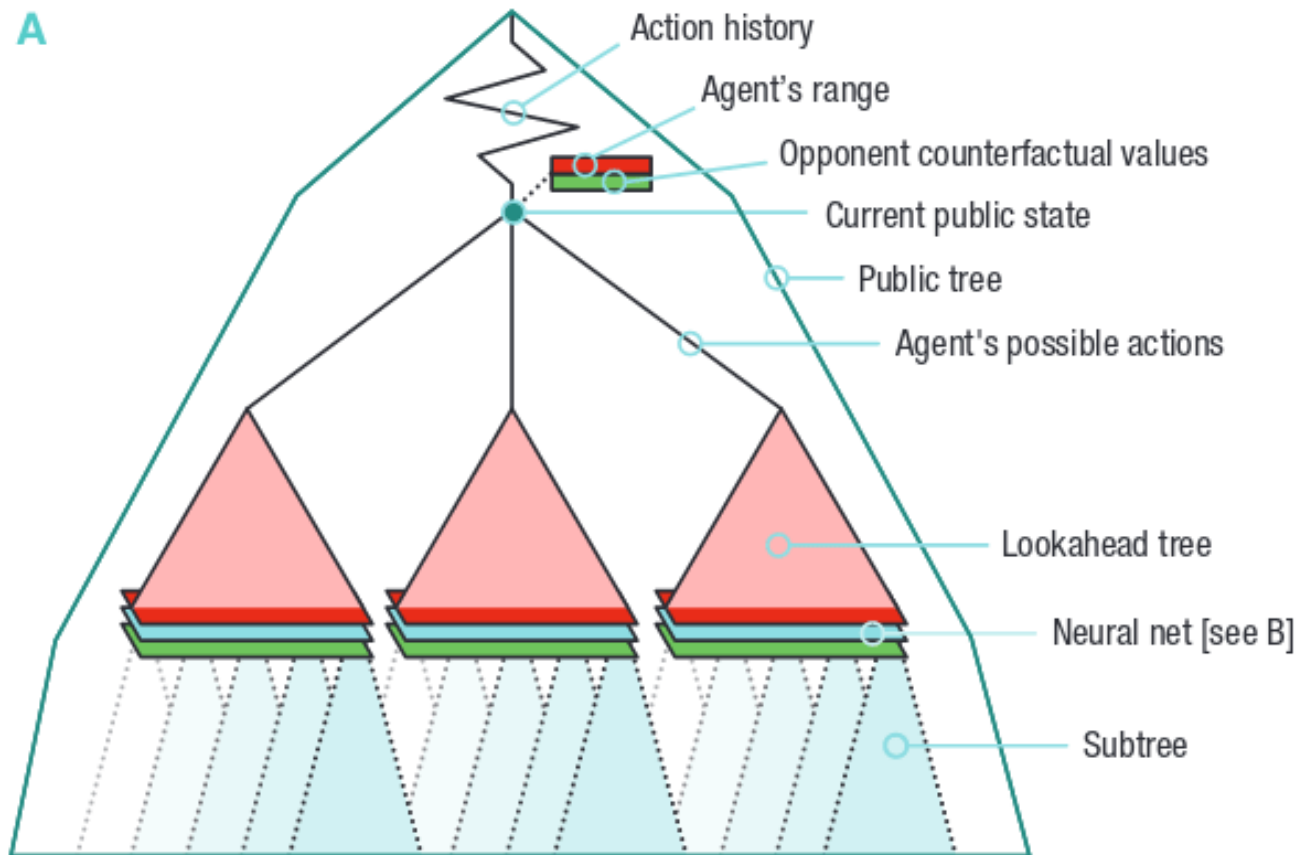
end for

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

end function



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



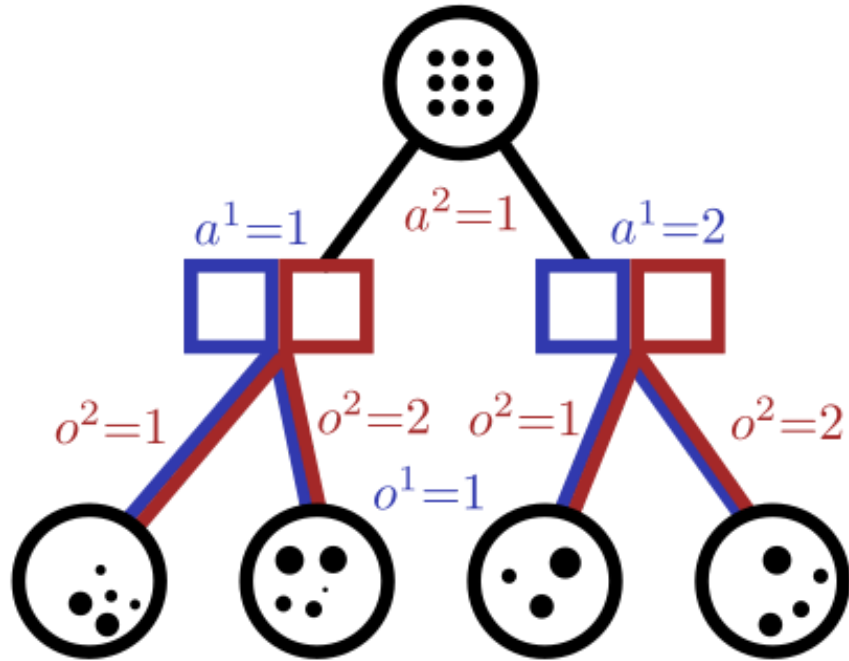
Tree-Based Planning in POSGs

Tree-Based Planning in POSGs

**Our approach: combine particle
filtering and information sets**

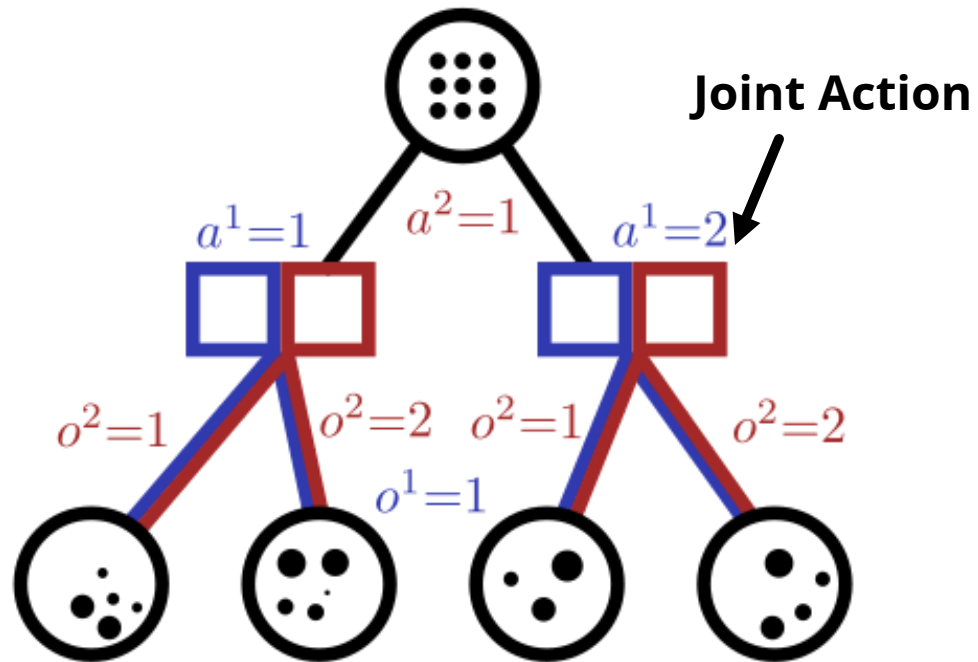
Tree-Based Planning in POSGs

Our approach: combine particle filtering and information sets



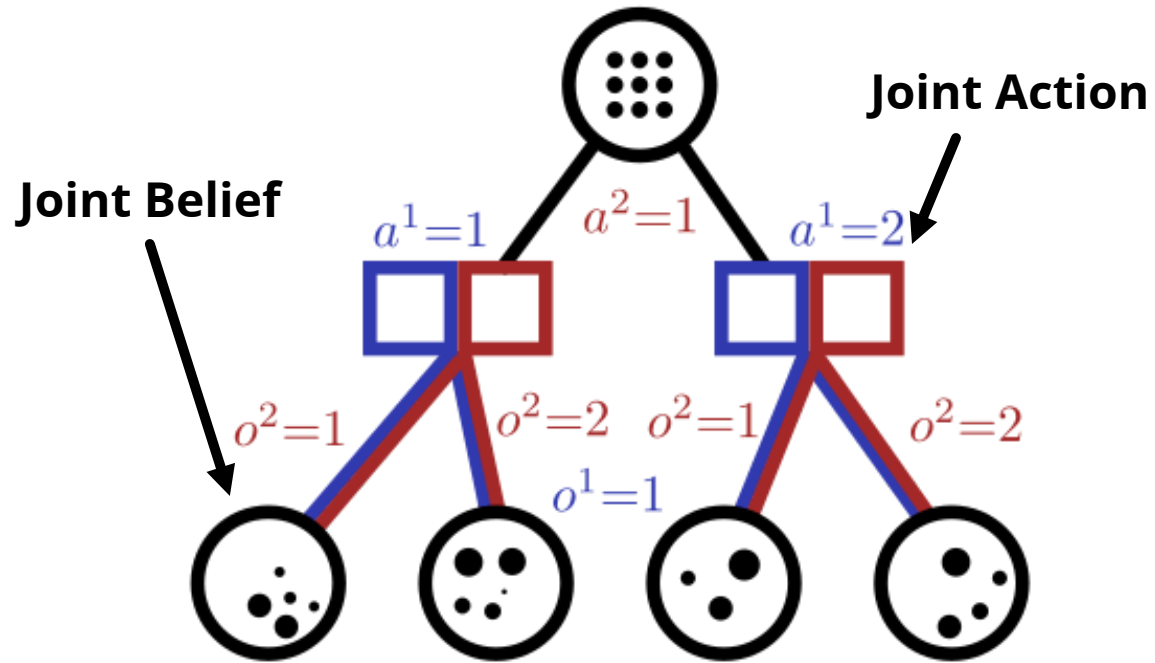
Tree-Based Planning in POSGs

Our approach: combine particle filtering and information sets



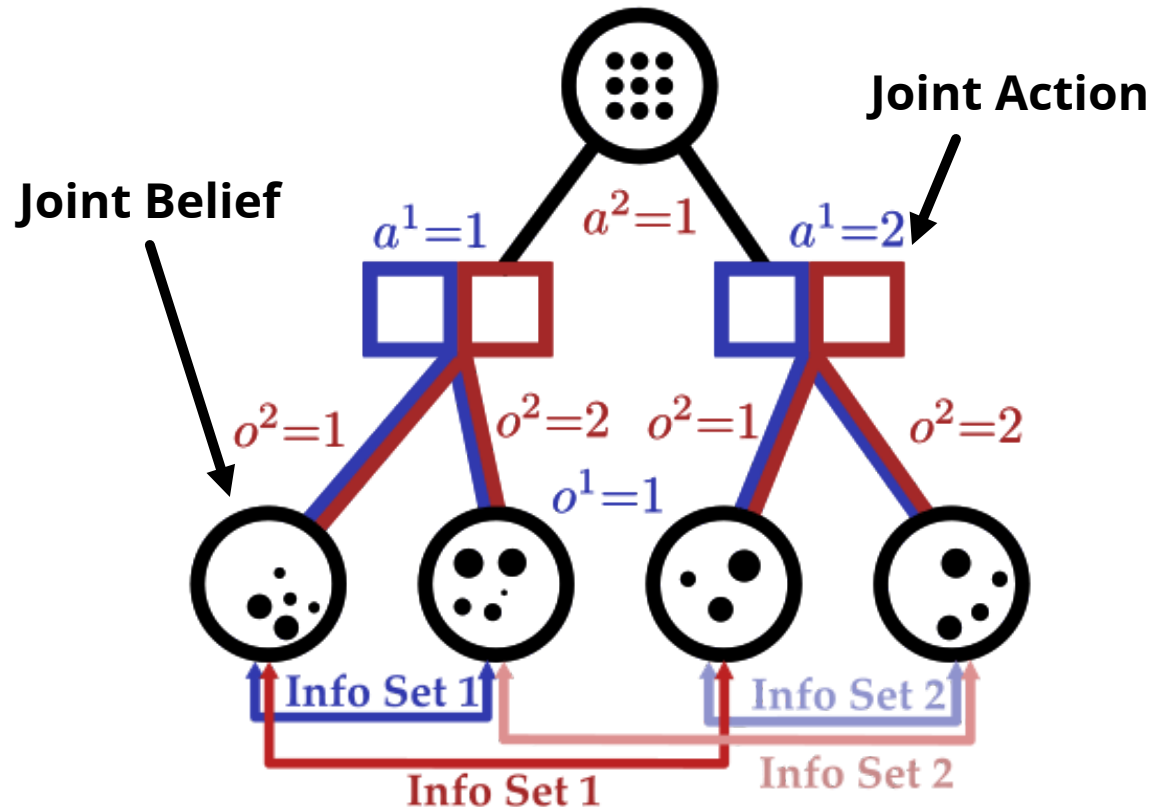
Tree-Based Planning in POSGs

Our approach: combine particle filtering and information sets



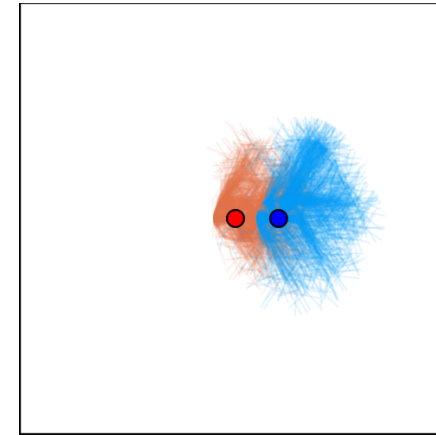
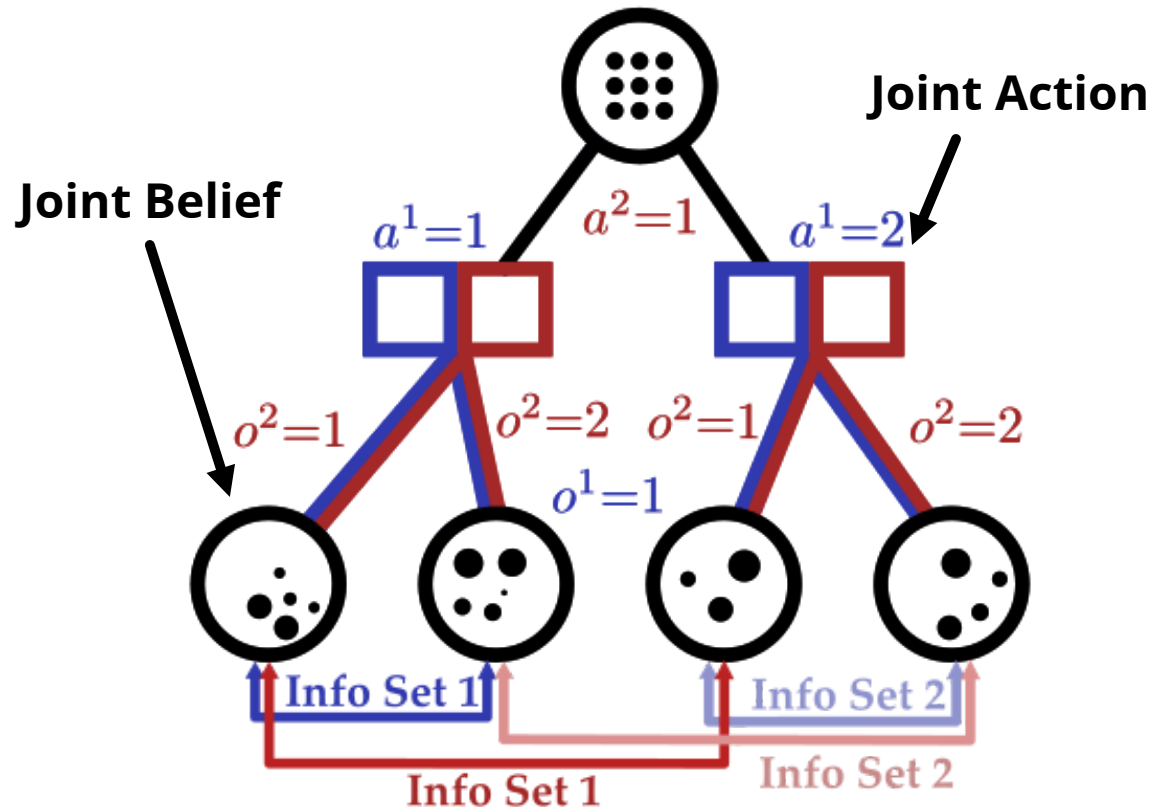
Tree-Based Planning in POSGs

Our approach: combine particle filtering and information sets



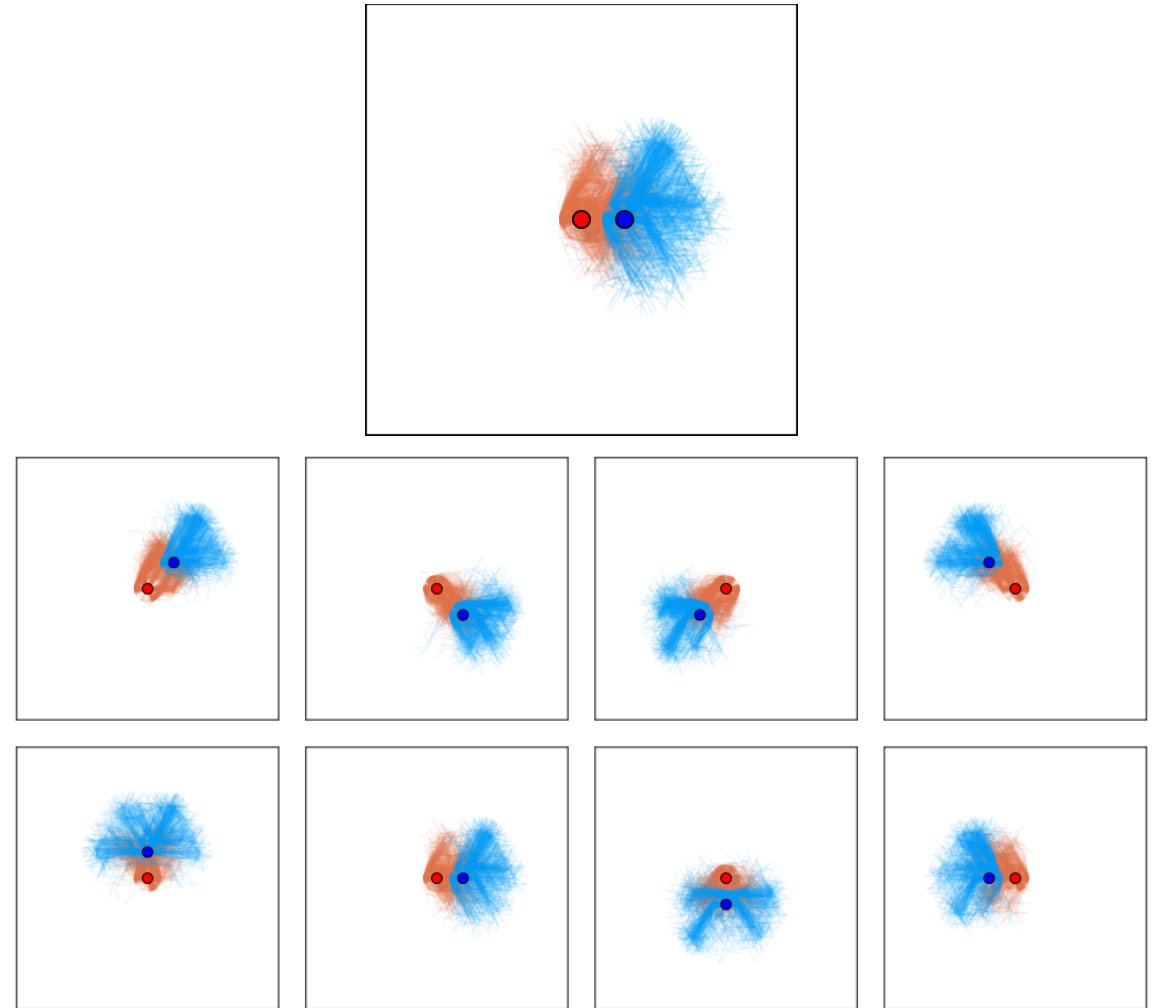
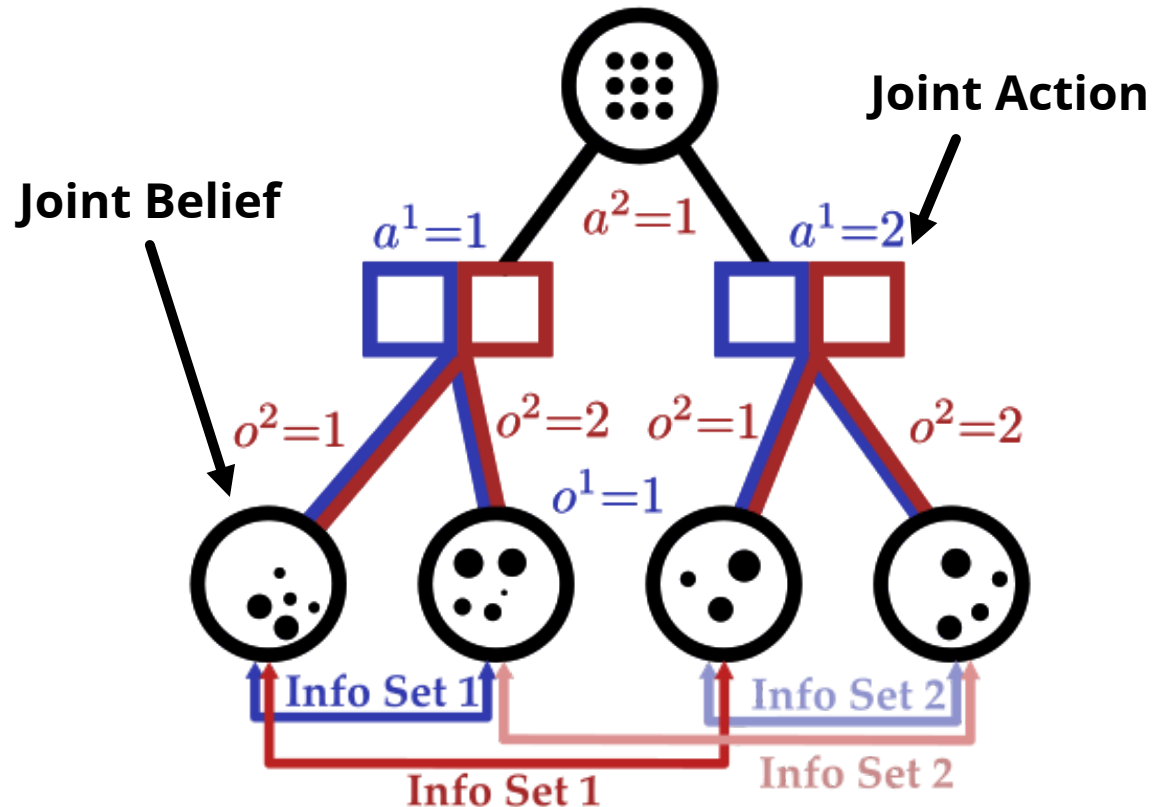
Tree-Based Planning in POSGs

Our approach: combine particle filtering and information sets



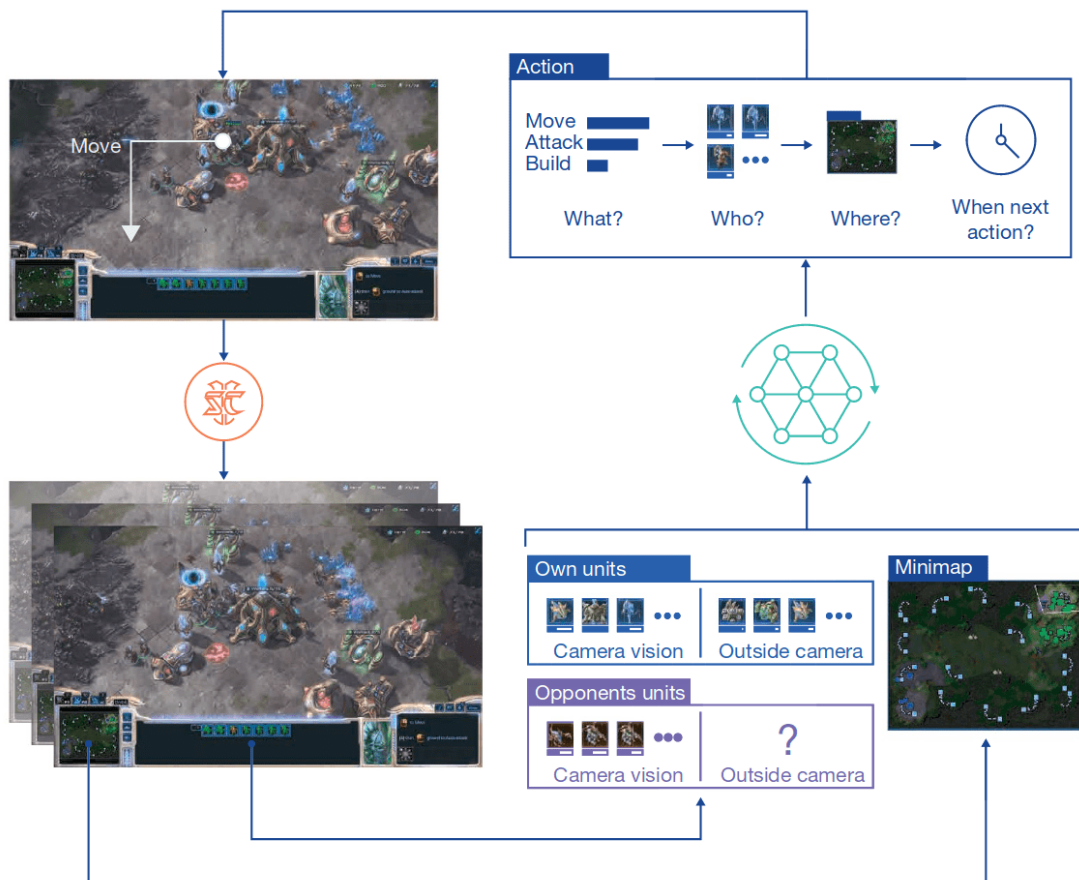
Tree-Based Planning in POSGs

Our approach: combine particle filtering and information sets



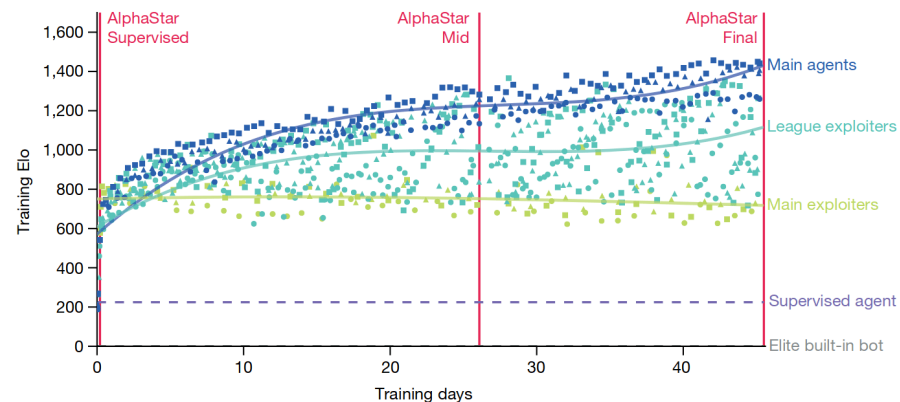
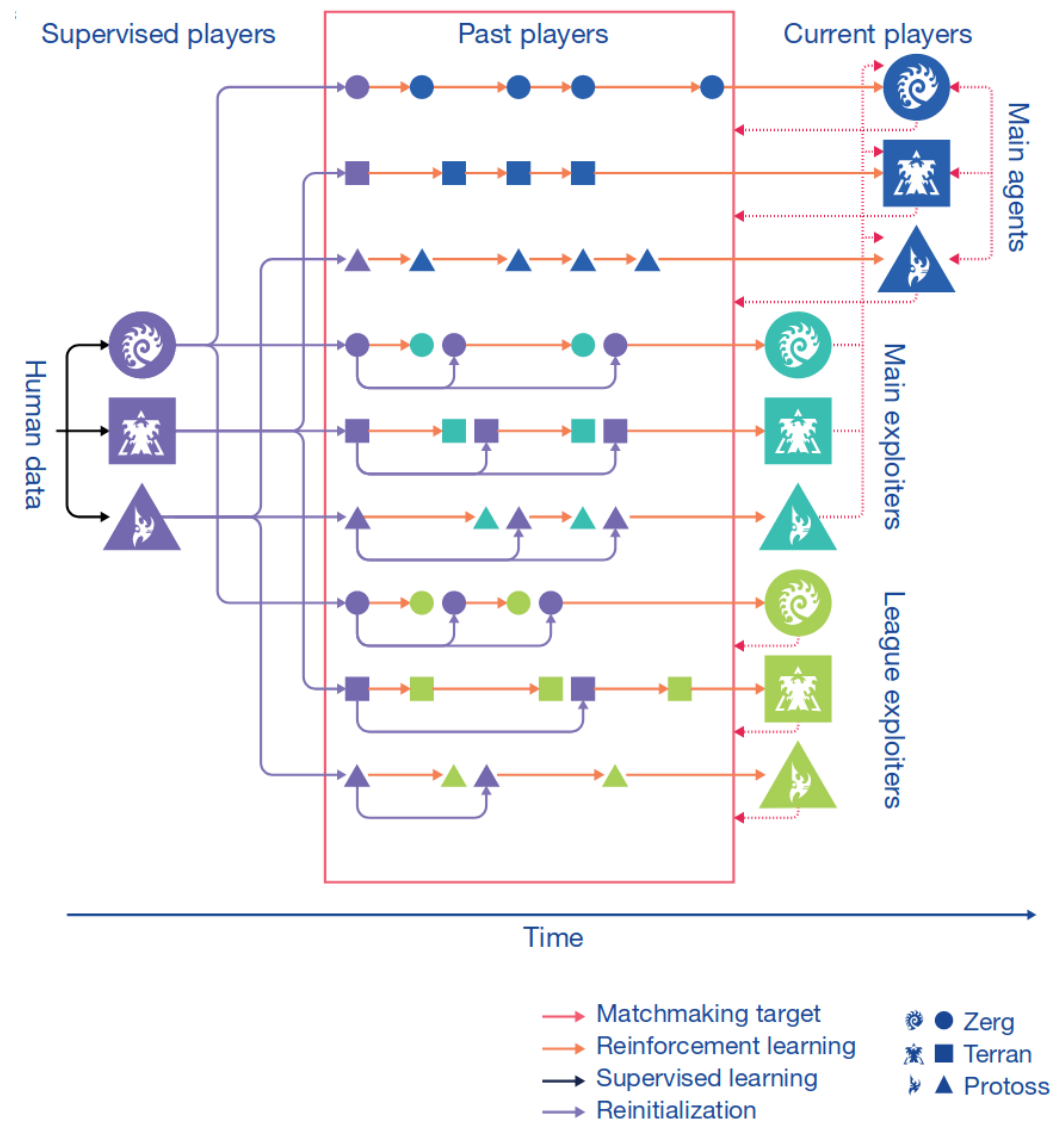
Monitoring layer

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



Real-time processing delay 30 ms

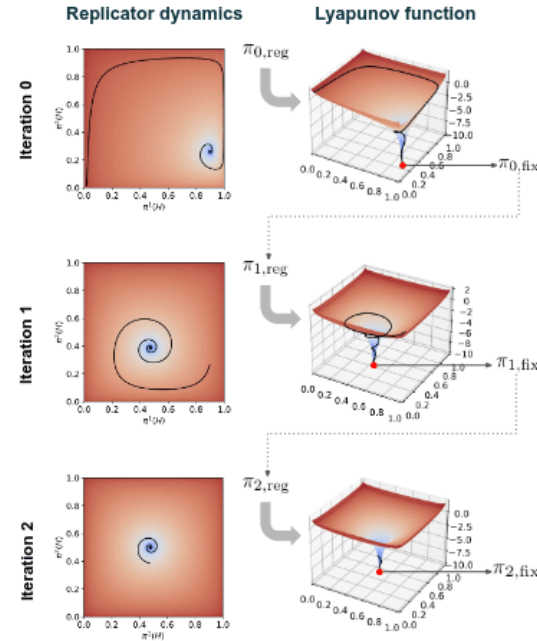
Alpha Star



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: $\pi_{0,reg}$

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

$$\pi_{n+1,reg} = \pi_{n,fix}$$

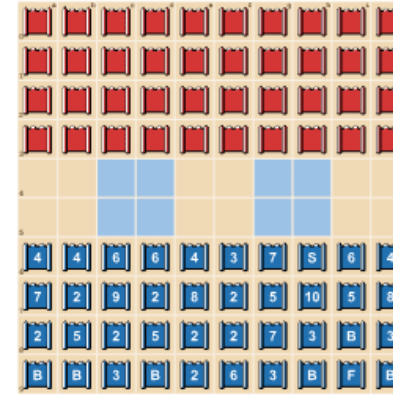
Repeat steps until convergence

(b) Algorithmic steps

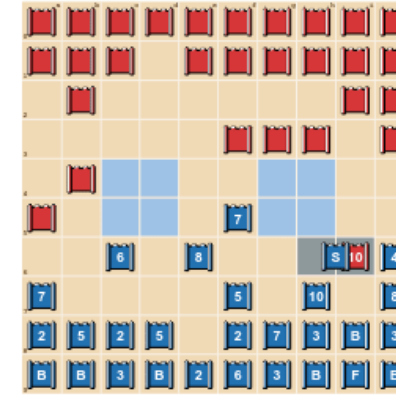
(c) Dynamics and Lyapunov function

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

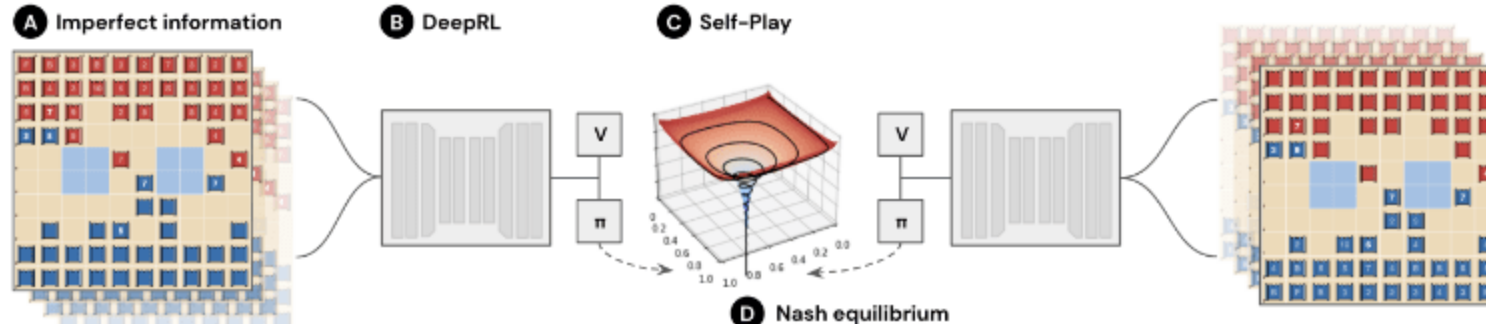
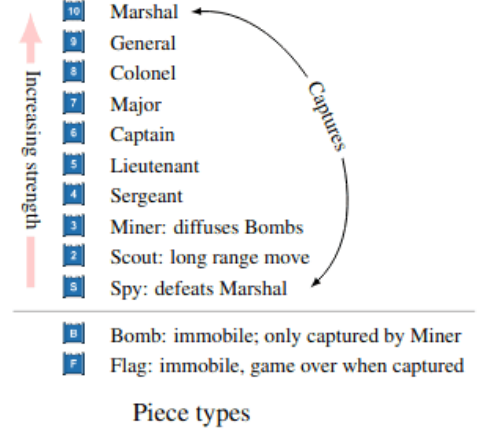
Stratego



Phase 1: Private deployment



Phase 2: Game play



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

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