# Incomplete Information Dynamic Games

# Incomplete Information



But you must have known I was not a great fool

# Partially Observable Markov Decision Process (POMDP)
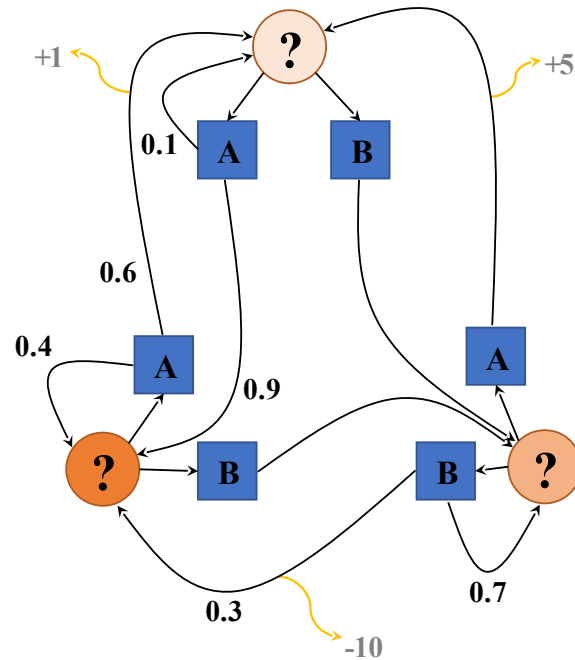


- $\mathcal{S}$ - State space
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ - Transition probability distribution
- $\mathcal{A}$ - Action space
- $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ - Reward
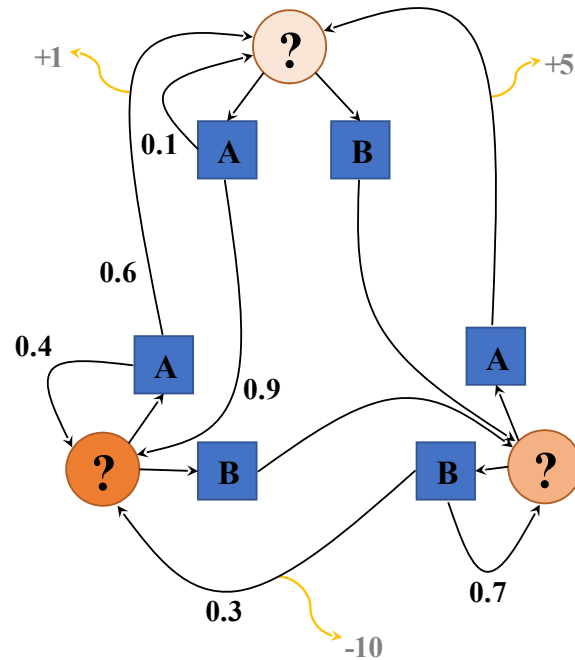
**Alleatory**

# 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
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ - Transition probability distribution
- $\mathcal{A}$ - Action space
- $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ - Reward
- $\mathcal{O}$ - Observation space
- $Z : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{O} \to \mathbb{R}$ - Observation probability distribution
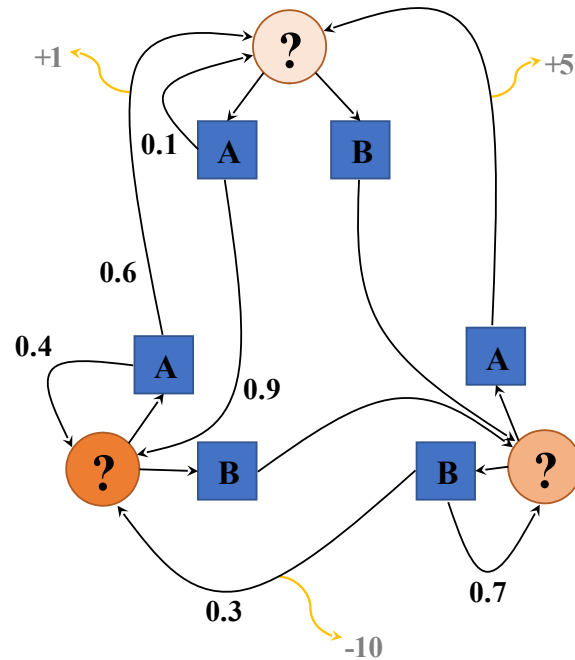
**Alleatory**

# Partially Observable Markov Decision Process (POMDP)



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

**Alleatory**
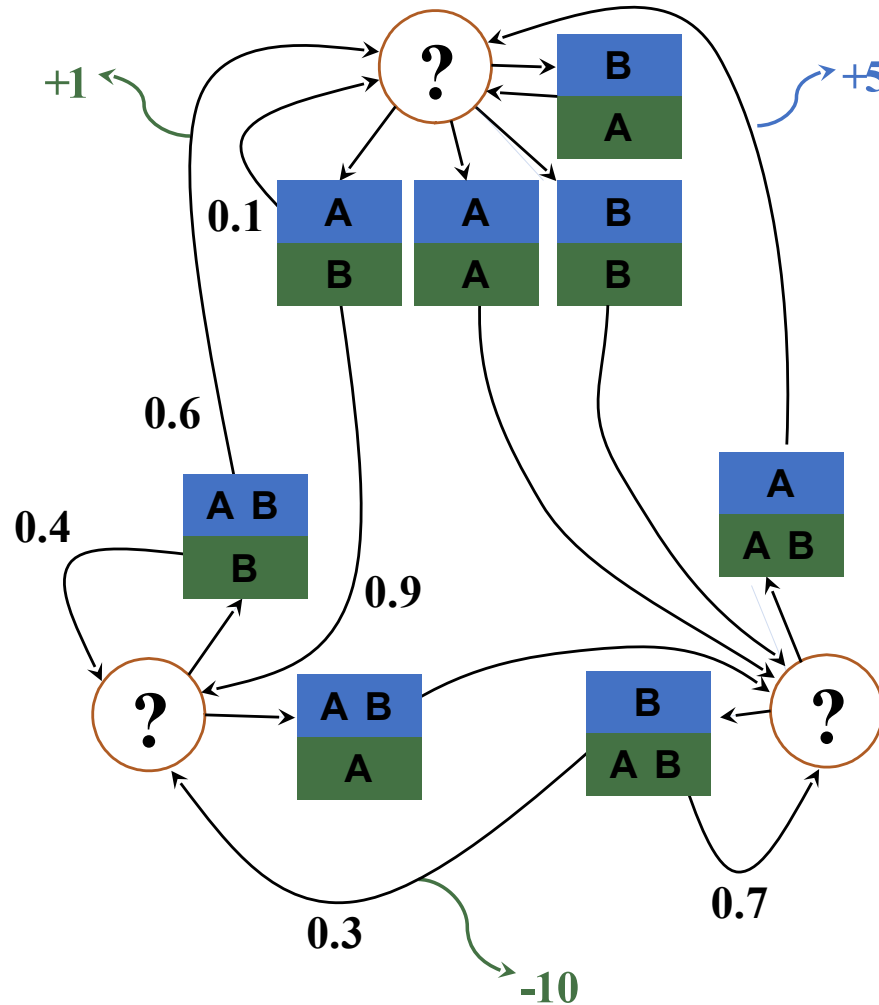
**Epistemic (Static)**

**Epistemic (Dynamic)**

# Partially Observable Markov Game

$$\left( \mathcal{S}, T, \{\mathcal{A}^i\}, \{R^i\}, \{\mathcal{O}^i\}, \{Z^i\}, \right)$$



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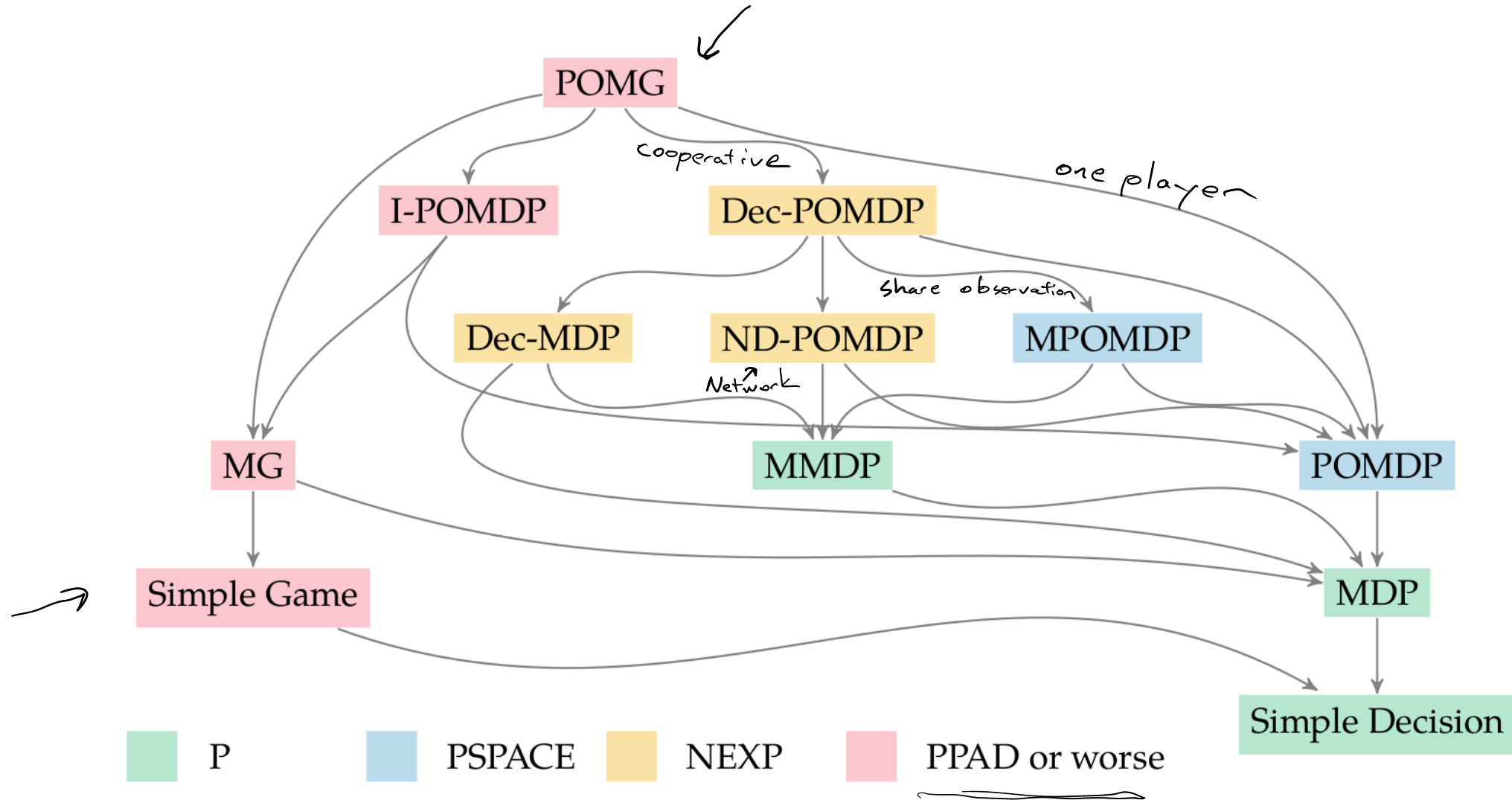- $Z(o^i \mid \boldsymbol{a}, s')$ - Observation probability distribution

| Alleatory | Epistemic (Static) |
| --- | --- |
| Epistemic (Dynamic) | Interaction |

# Hierarchy of Problems

# Belief updates?

POMDP: $\qquad b'(s') \propto Z(o|a,s') \sum_s T(s'|s,a) b(s)$
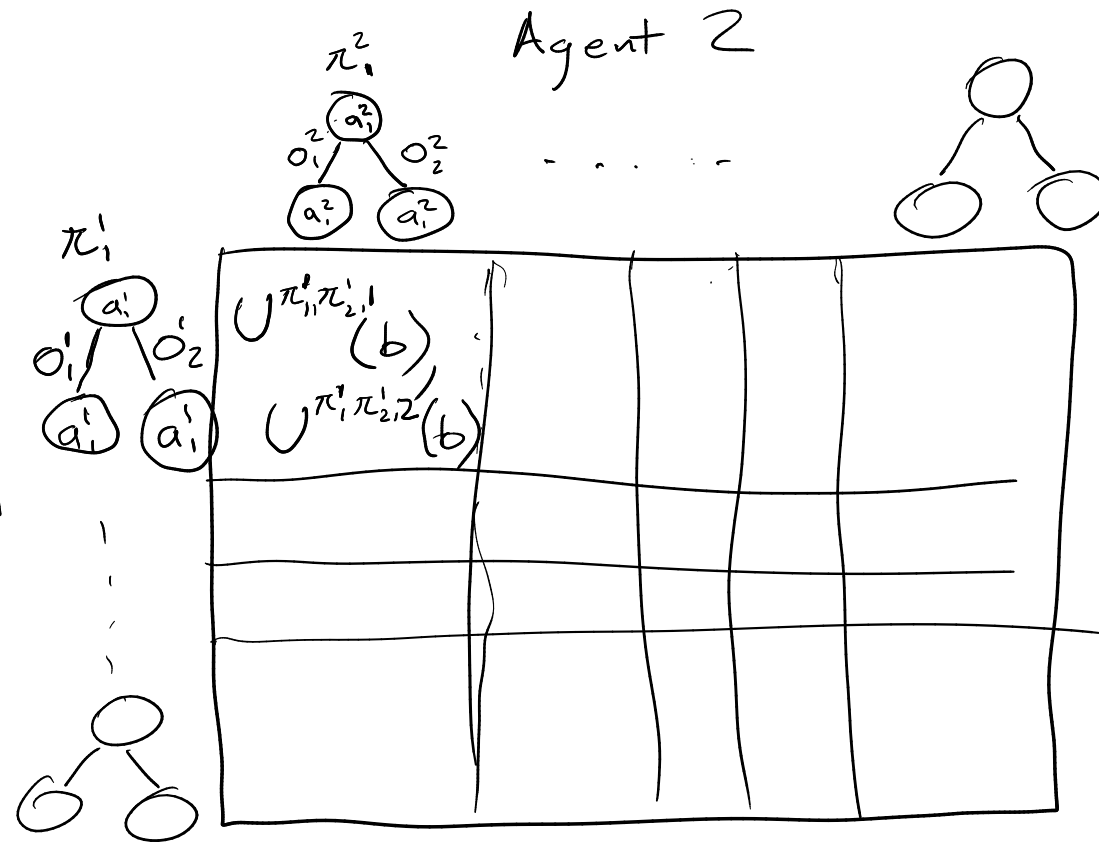
POMG $\qquad Z(o|\vec{a},s') \sum_s T(s'|s,\vec{a}) b(s)$ $\qquad$ X

$\underset{\text{joint action}}{\curvearrowleft}$ $\qquad$ $\underset{\text{joint actions}}{\curvearrowleft}$

$\pi^{-i}$

Problem: Usually trying to solve for $\pi^{-i}$ at the same time as choosing our actions

# Reduction to Simple Game

Agent 2

$\pi_1^2$

$a_1^2$

$o_1^2$   $o_2^2$

$a_1^2$   $a_1^2$

$\pi_1^1$

$a_1^1$

$o_1^1$   $o_2^1$

$a_1^1$   $a_1^1$

Agent 1

$U^{\pi_1^1, \pi_{2,1}^1}(b)$

$U^{\pi_1^1, \pi_{2,2}^1}(b)$

2 ways

1. Dynamic Programming
   with Pruning

2. Best responses
   "Double Oracle"
   - start with strategy
   - compute best response
   - add best response
     to matrix game
   - solve matrix game

# Pruning in Dynamic Programming

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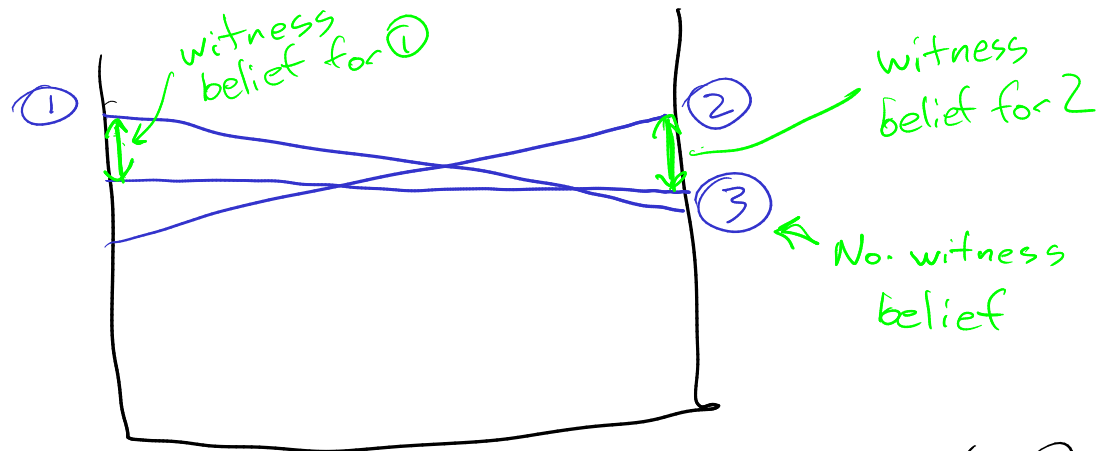
Loop

    Evaluate $N+1$ step policies

    Prune dominated policies

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 $\pi^i$.

POMDPs



witness belief for ①

witness belief for 2

No. witness belief

maximize $\delta$
$\delta, b$

subject to $\quad b(\pi^{-i}, s) \geq 0 \qquad \forall \pi^{-i}, s$

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

$$\sum_{\pi^{-i}} \sum_{s} b(\pi^{-i}, s) \left( U^{\pi^{i'}, \pi^{-i}_i}(s) - U^{\pi^{i}, \pi^{-i}_i}(s) \right) \geq \delta \quad \forall \pi^{i'}$$

If $\delta > 0$
keep $\pi^i$

# 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

# 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
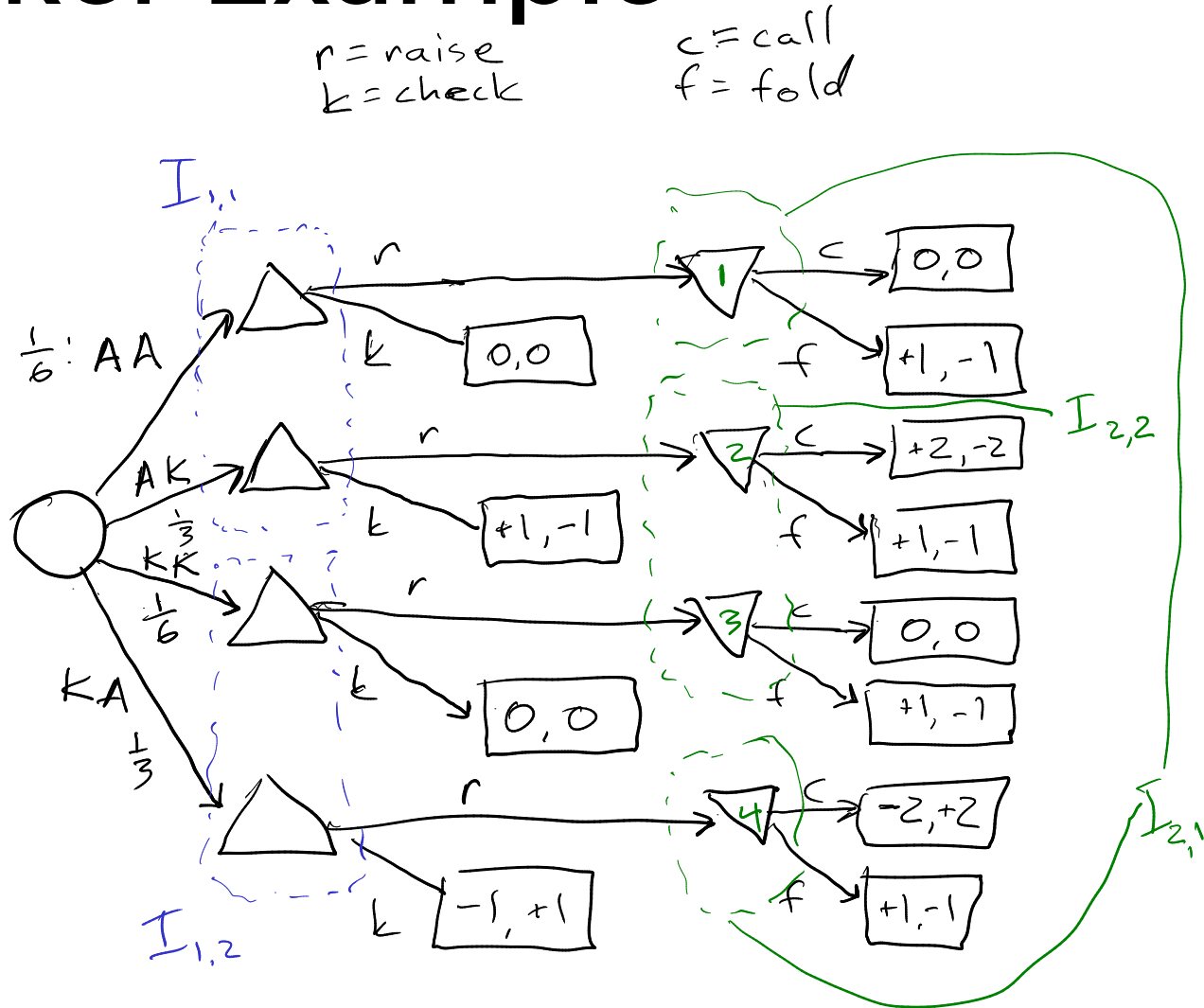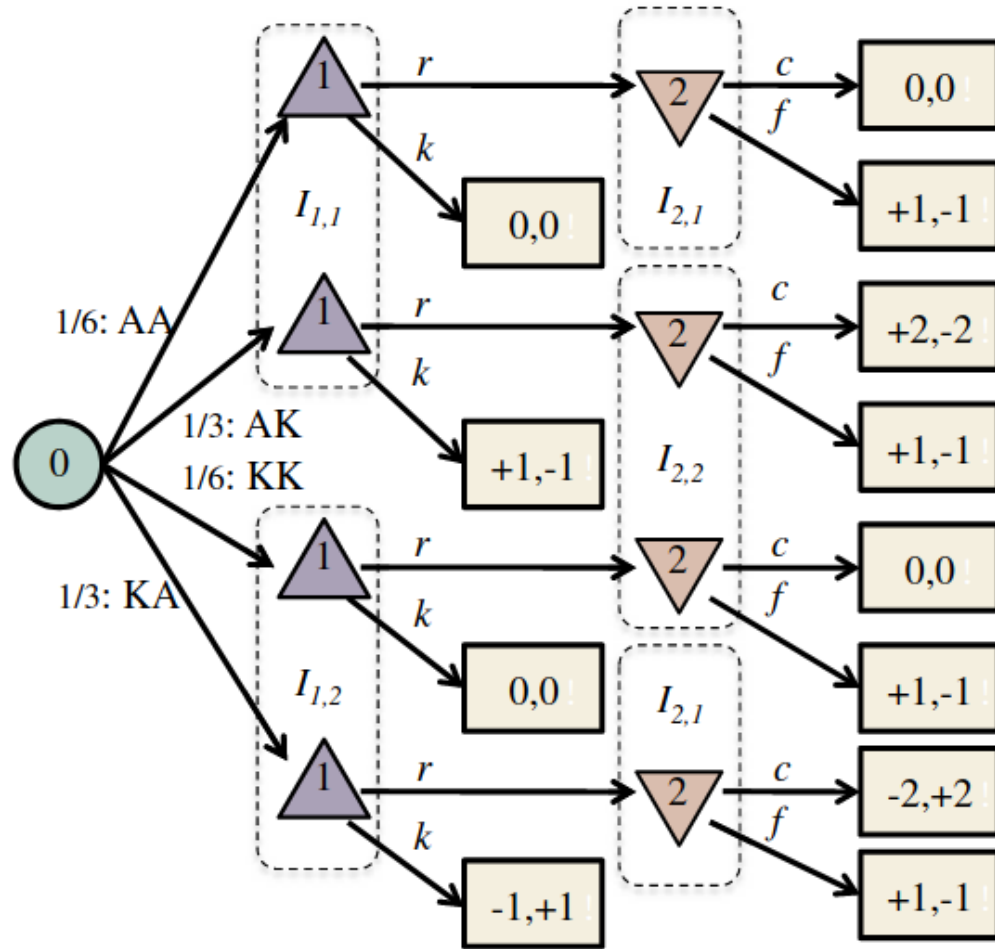
# King-Ace Poker Example

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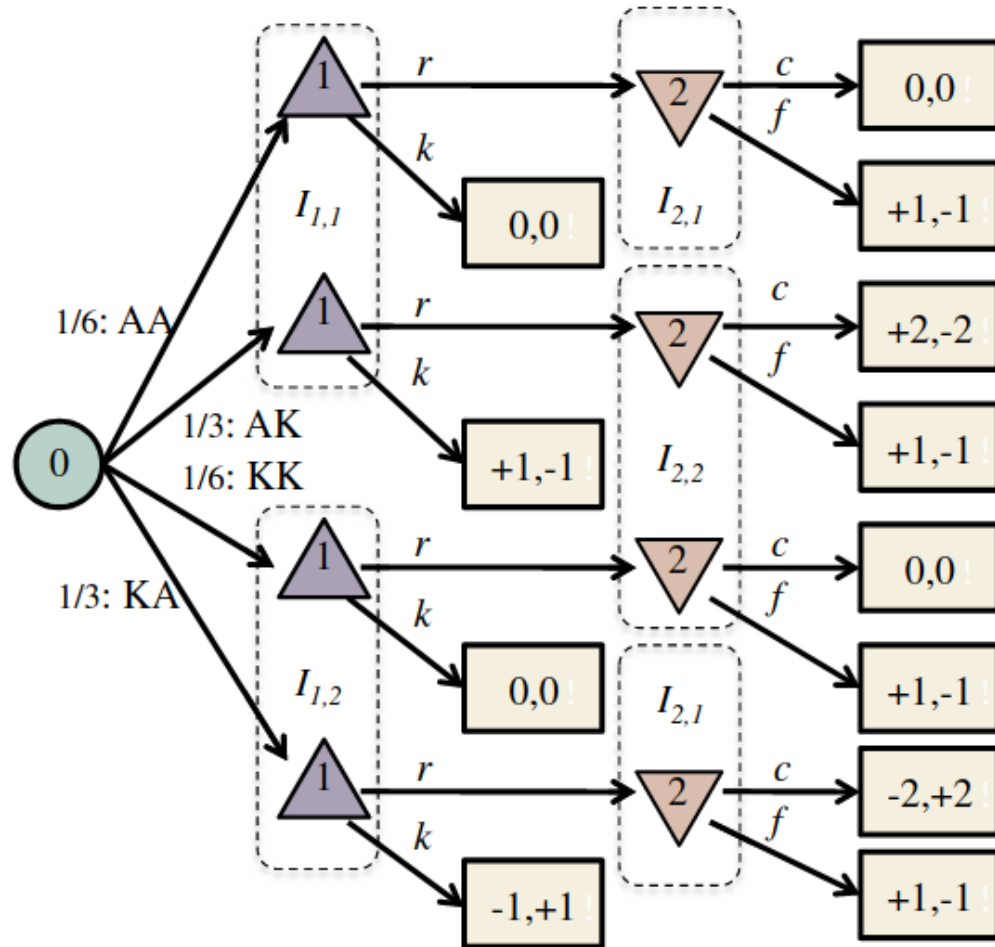
- 4 Cards: 2 Aces, 2 Kings
- Each player is dealt a card
- P1 can either *raise* ($r$) the payoff to 2 points or *check* ($k$) the payoff at 1 point
- If P1 raises, P2 can either *call* ($c$) Player 1's bet, or *fold* ($f$) the payoff back to 1 point

# King-Ace Poker Example

- 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
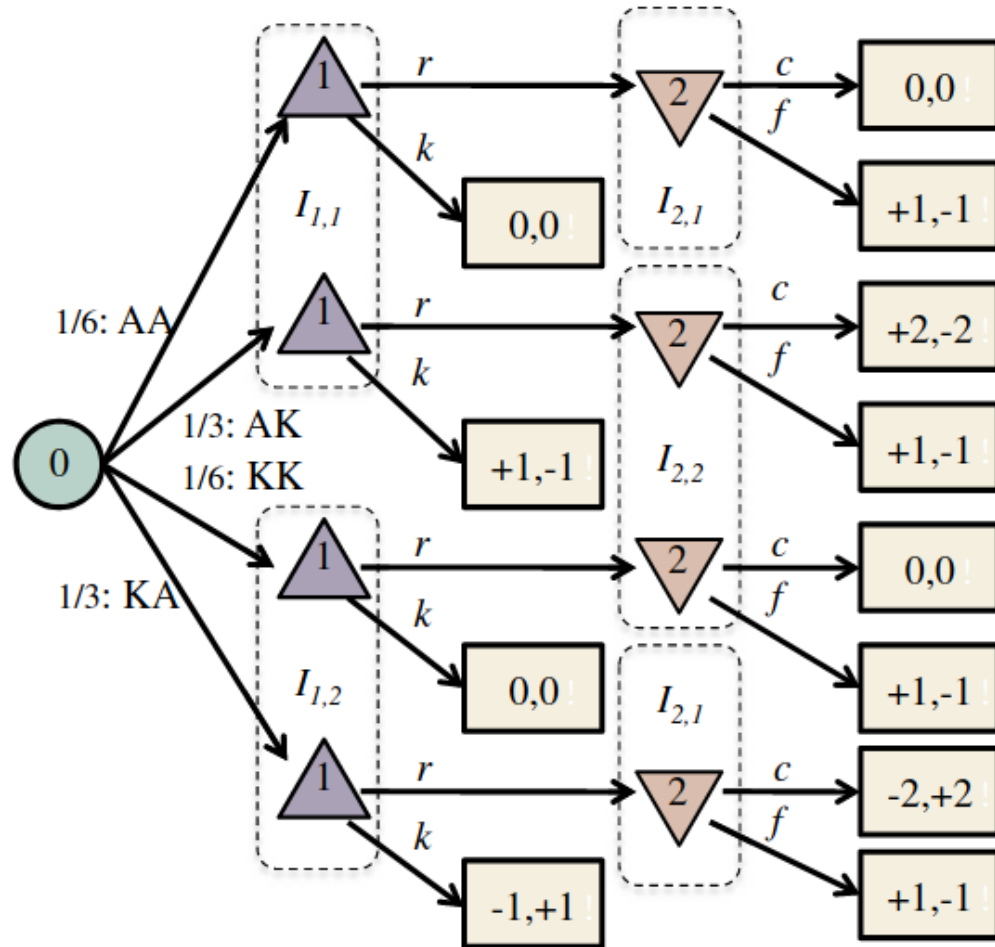
# Extensive to Matrix Form

# 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 |

# 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

# Fictitious Play in Extensive Form Games



Algorithm 2 General Fictitious Self-Play
```
function FICTITIOUSSELFPLAY(Γ, n, m)
    Initialize completely mixed π₁
    β₂ ← π₁
    j ← 2
    while within computational budget do
        ηⱼ ← MIXINGPARAMETER(j)
        𝒟 ← GENERATEDATA(πⱼ₋₁, βⱼ, n, m, ηⱼ)
        for each player i ∈ 𝒩 do
            ℳⁱ_RL ← UPDATERLMEMORY(ℳⁱ_RL, 𝒟ⁱ)
            ℳⁱ_SL ← UPDATESLMEMORY(ℳⁱ_SL, 𝒟ⁱ)
            βⁱⱼ₊₁ ← REINFORCEMENTLEARNING(ℳⁱ_RL)
            πⁱⱼ ← SUPERVISEDLEARNING(ℳⁱ_SL)
        end for
        j ← j + 1
    end while
    return πⱼ₋₁
end function

function GENERATEDATA(π, β, n, m, η)
    σ ← (1 − η)π + ηβ
    𝒟 ← n episodes {tₖ}₁≤ₖ≤ₙ, sampled from strategy
    profile σ
    for each player i ∈ 𝒩 do
        𝒟ⁱ ← m episodes {tⁱₖ}₁≤ₖ≤ₘ, sampled from strat-
        egy profile (βⁱ, σ⁻ⁱ)
        𝒟ⁱ ← 𝒟ⁱ ∪ 𝒟
    end for
    return {𝒟ᵏ}₁≤ₖ≤ₙ
end function
```

Heinrich et al. 2015 "Fictitious Self Play in Extensive-Form Games"

This slide not covered on exam

13

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

Counterfactual Regret Minimization + Deep Learning

# Can game learning methods like CFR be used in Large POMGs?