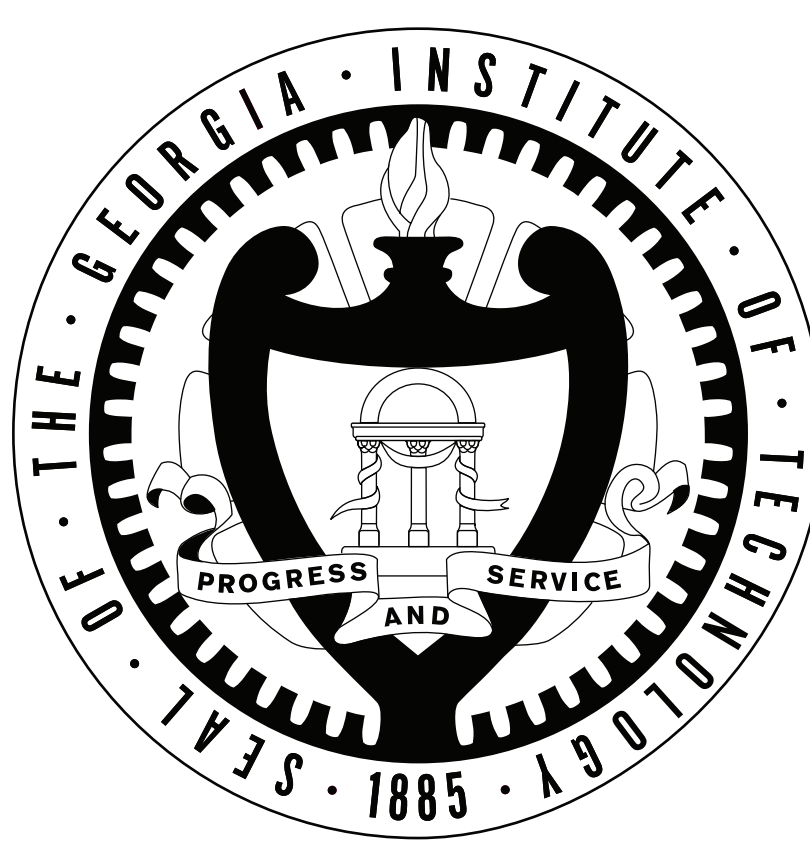


# QPACE: Quick Polytope Approximation of All Correlated Equilibria in Stochastic Games

Liam MacDermid  
Karthik Narayan  
Charles Isbell  
Lora Weiss



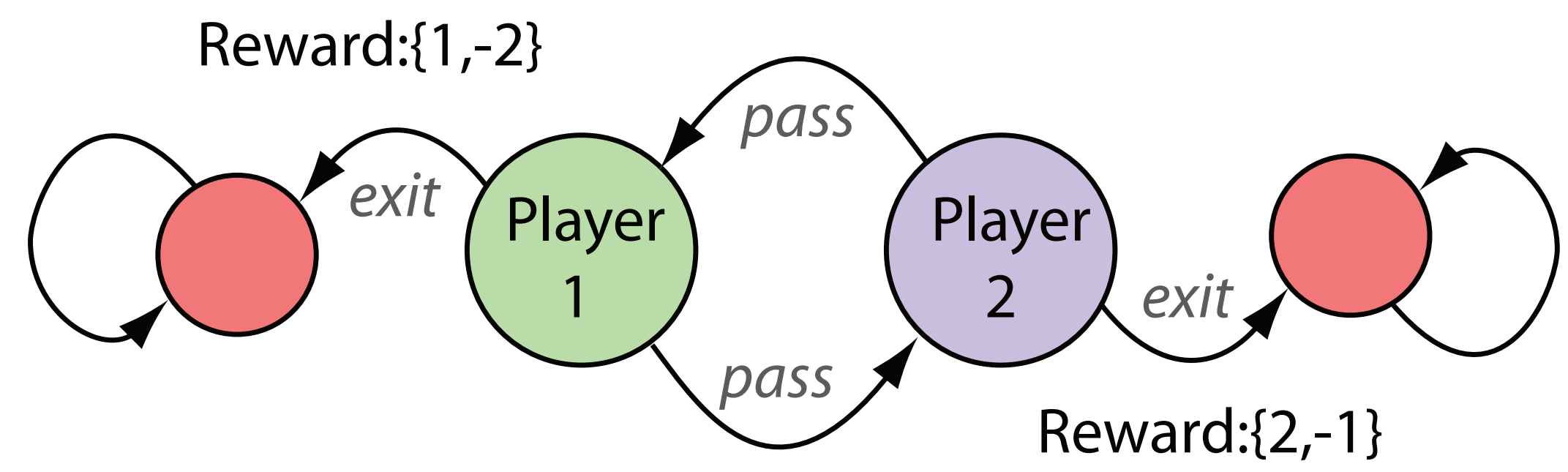
Key Ideas:

Solve multi-agent reinforcement learning by using achievable-sets instead of values (Murray & Gordon 2007).

Approximate achievable-sets using a fixed collection of half-spaces (MacDermid & Isbell 2009).

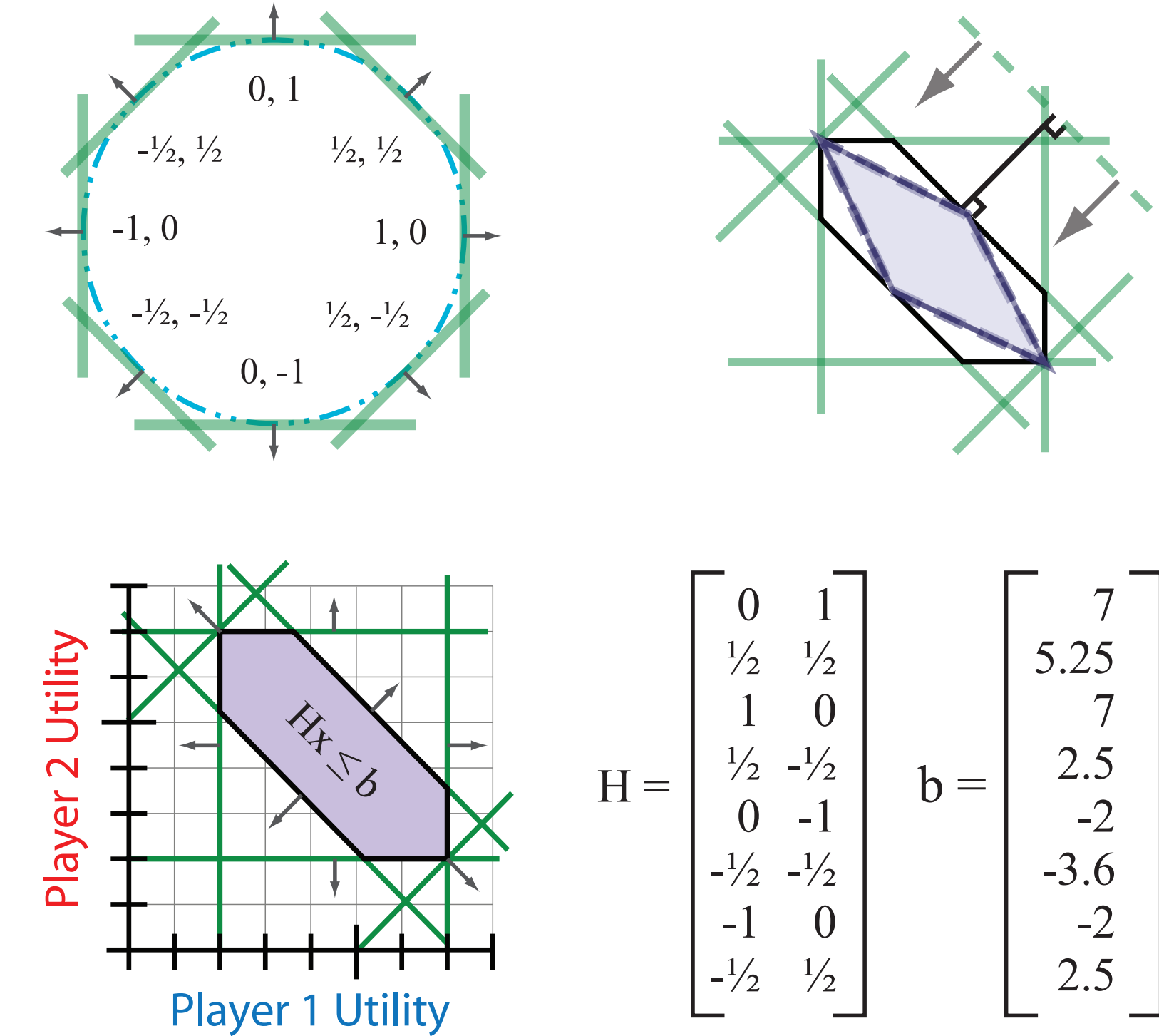
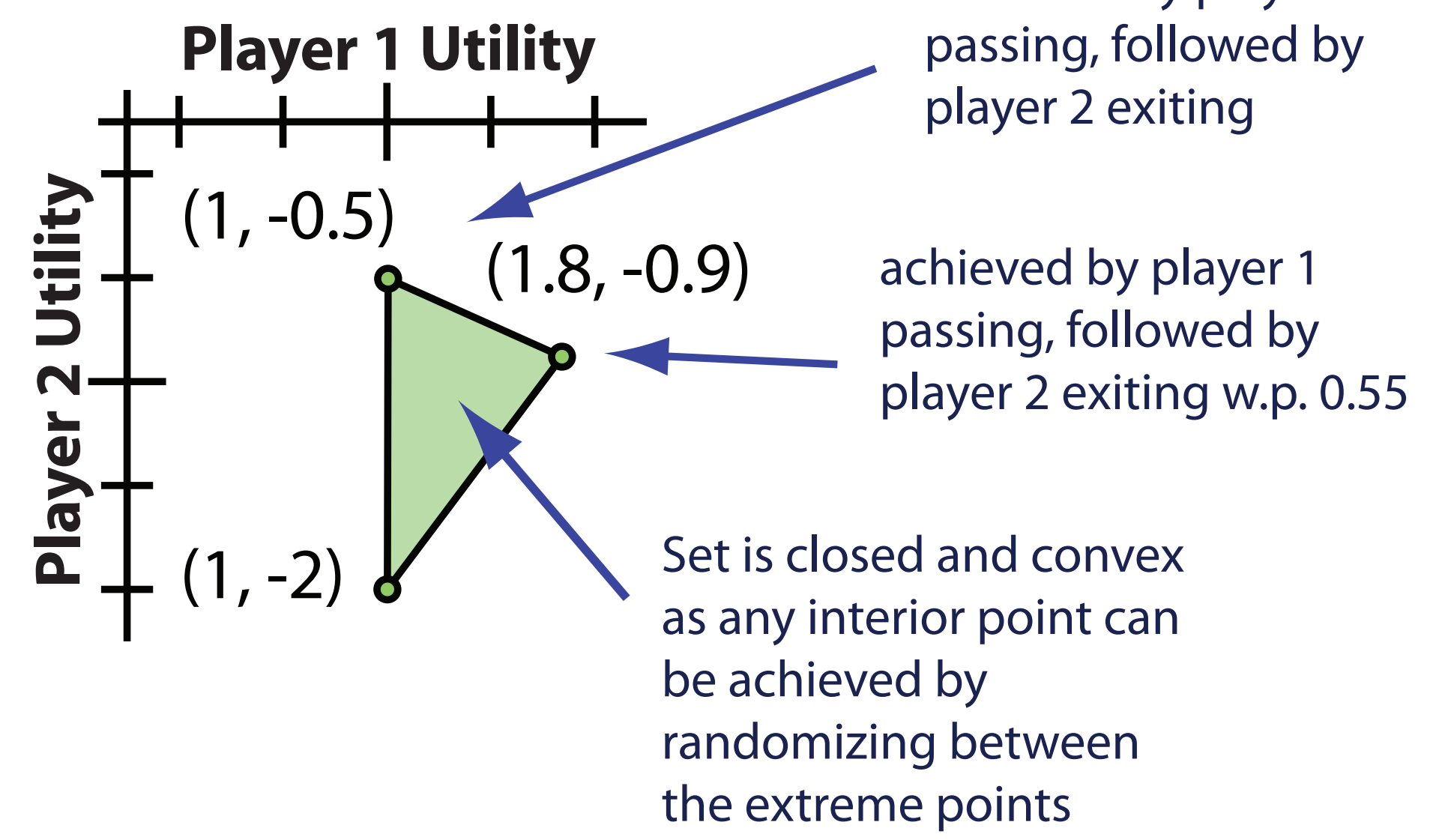
This polytope representation can be used to dramatically improve performance by (this paper):

An Example Game (The Breakup Game):



Circles represent states, outgoing arrows represent deterministic actions. Unspecified rewards are zero. Previous algorithms could not solve this game.

The final achievable-set for player 1's state



Achievable sets are approximated using regular half-spaces with a fixed set of normals.

Offsets are contracted to construct an overestimate. Normals remain constant.

- Eliminating vertex computation
- Permitting very fast Minkowski sums
- Preventing sets from changing dramatically, allowing linear program solution caching
- Transforming multi-objective LPs into a series of single-objective LPs

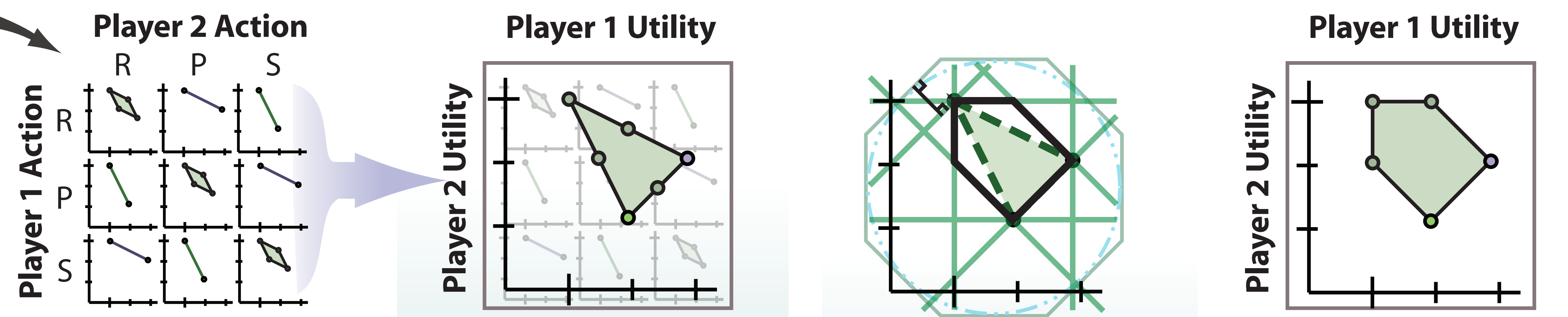
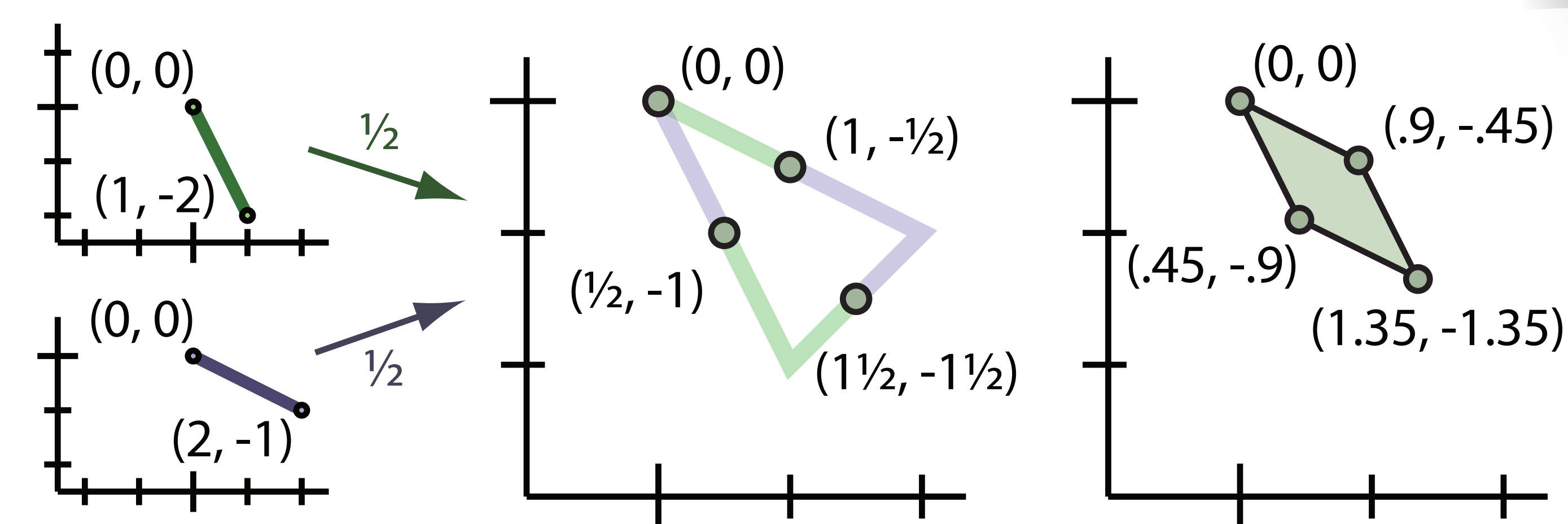
The tractable backup of achievable sets: (an iteration consists of a single backup for each state)

replaces 'max' in single agent RL

$$Q(s, \vec{a})_j = R(s, \vec{a}) \cdot H_j + \gamma \sum_{s'} P(s'|s, \vec{a}) V(s')_j$$

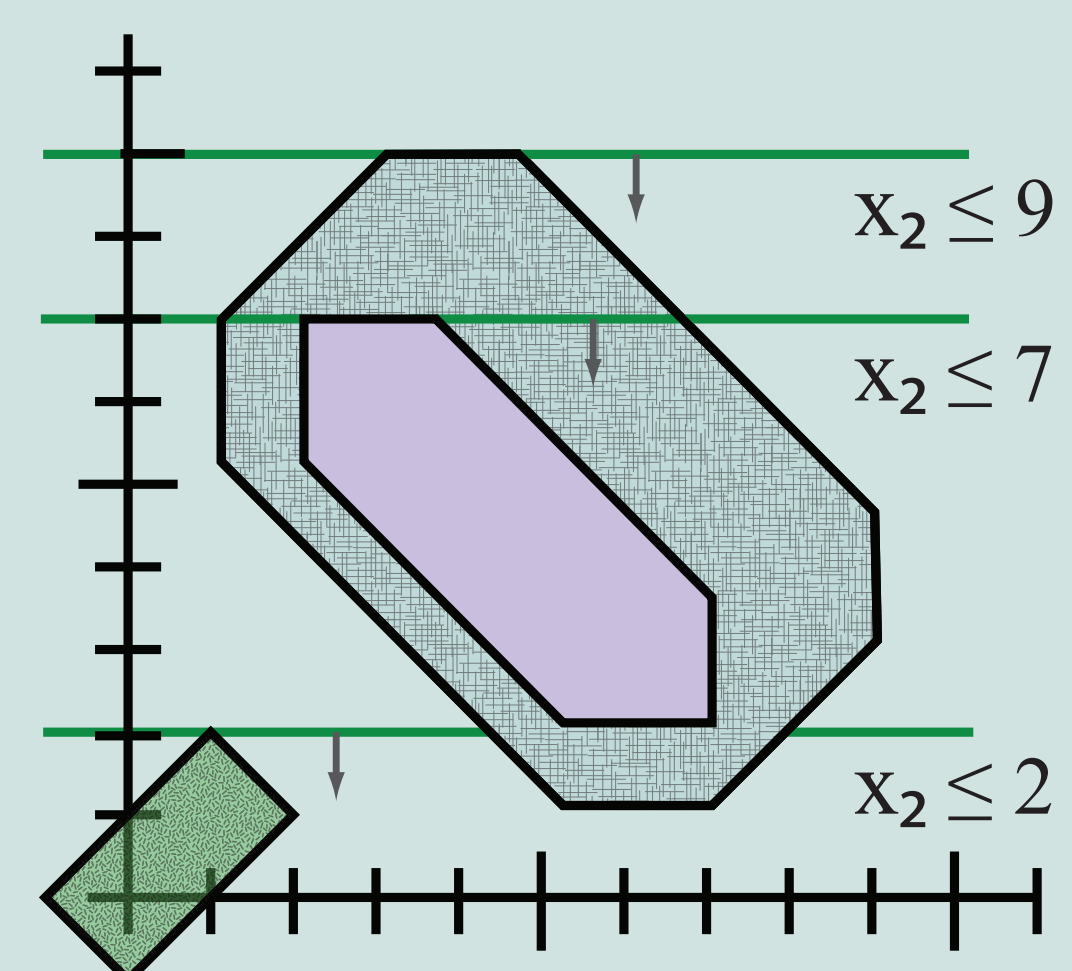
$$V^*(s) = \text{equilibrium}_{\vec{a}}(Q^*(s, \vec{a}))$$

The state shown being calculated is an initial rock-paper-scissors game played to decide who goes first in the breakup game



Value Computed:

Minkowski sums can be very efficiently computed using our regular polytope scheme



$\begin{bmatrix} 7 \\ 5.25 \\ 7 \\ 2.5 \\ -2 \\ -3.6 \\ -2 \\ 2.5 \end{bmatrix}$	$\begin{bmatrix} 2 \\ 1.5 \\ 2 \\ 0.5 \\ 1 \\ 0.5 \\ 0.5 \end{bmatrix}$	$\begin{bmatrix} 9 \\ 9 \\ 3 \\ -1 \\ -3.1 \\ -1 \\ 3 \end{bmatrix}$

For each player  $i$ , distinct actions  $\alpha, \beta \in A_i$ ,

$$\sum_{\vec{a} \in A^n} \vec{c} u_{\vec{a}(\alpha)} \geq \sum_{\vec{a} \in A^n} x_{\vec{a}(\alpha)} [\vec{g} t_{\vec{a}(\beta)} + R(s, \vec{a}(\beta))_i]$$

$$\sum_{\vec{a} \in A^n} x_{\vec{a}} = 1 \text{ and } \forall \vec{a} \in A^n, x_{\vec{a}} \geq 0$$

For each joint-action  $\vec{a} \in A^n$ ,

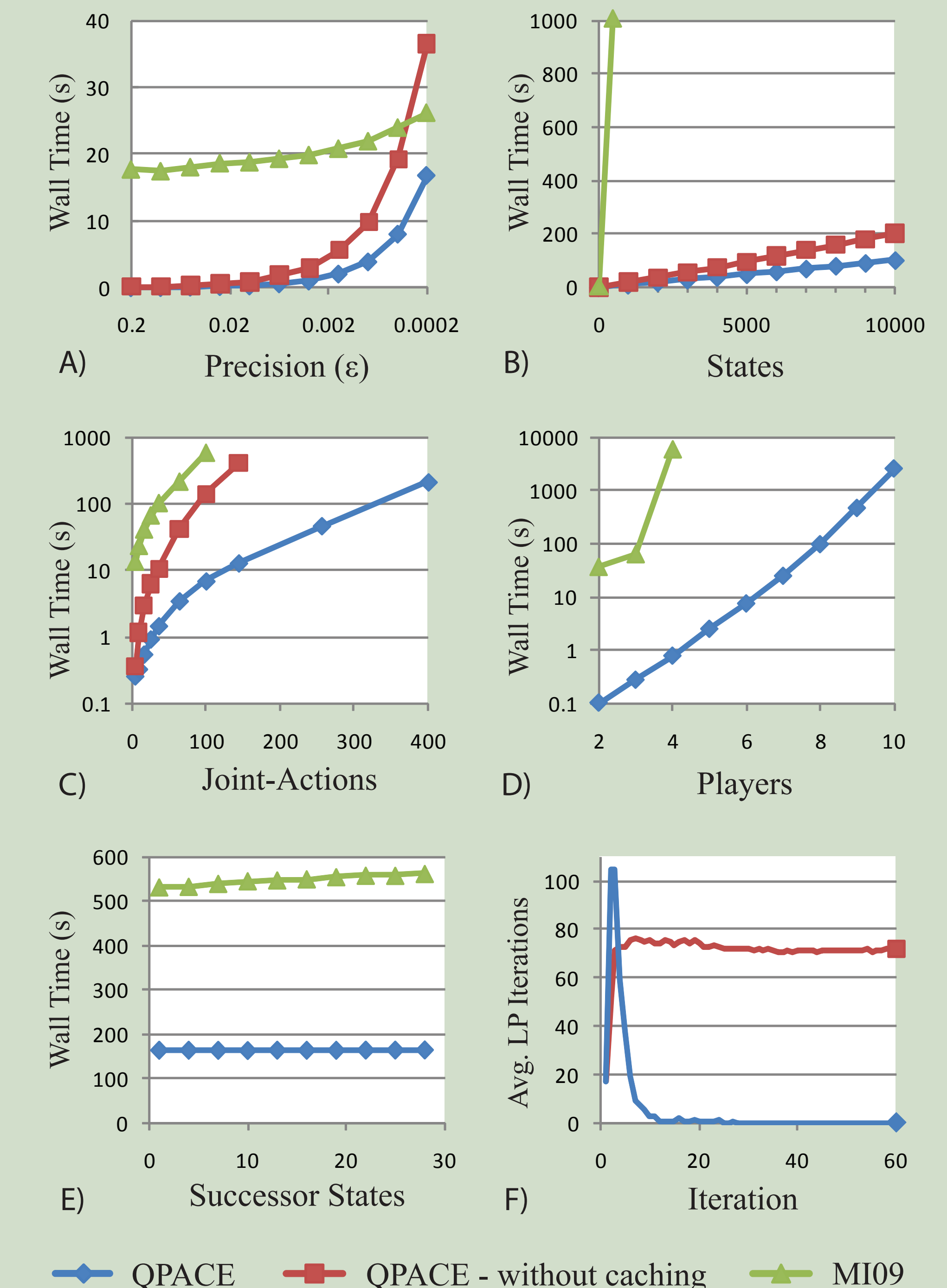
$$\vec{c} u_{\vec{a}} \in x_{\vec{a}} Q(s, \vec{a}) \text{ (i.e. } H_j \vec{c} u_{\vec{a}} \leq x_{\vec{a}} Q(s, \vec{a})_j \text{)}$$

$$\max \sum_{\vec{a}} \sum_{i \in I} H_{j,i} \cdot \vec{c} u_{\vec{a}i}$$

subject to: these inequalities

Each LP changes slowly from iteration to iteration, granting very fast results when we start the LP from the previous solution

Results: We provide the first approximation algorithm which solves stochastic games to within  $\epsilon$  absolute error of the optimal game-theoretic solution, with running time polynomial in the error bound and the number of states and joint-action.



Legend: QPACE (blue), QPACE - without caching (red), MI09 (green)