

# Game Theory, Evolutionary Dynamics, and Multi-Agent Learning

Prof. Nicola Gatti  
[\(nicola.gatti@polimi.it\)](mailto:nicola.gatti@polimi.it)

# Game theory

# Game theory: basics

## Normal form

- Players
- Actions
- Outcomes
- Utilities
- Strategies
- Solutions


# Game theory: basics

- Normal form
- Players
  - Actions
  - Outcomes
  - Utilities
  - Strategies
  - Solutions

		Player 2		
		Player I		

# Game theory: basics

- Normal form
- Players
- Actions
- Outcomes
- Utilities
- Strategies
- Solutions

		Player 2		
		Rock	Paper	Scissors
		Rock		
Player 1		Paper		
		Scissors		

# Game theory: basics

- Normal form
- Players
  - Actions
  - Outcomes
  - Utilities
  - Strategies
  - Solutions

		Player 2		
		Rock	Paper	Scissors
		Rock	Tie	Player 2 wins
Player 1		Paper	Player 1 wins	Tie
		Scissors	Player 2 wins	Player 1 wins

# Game theory: basics

- Normal form
- Players
  - Actions
  - Outcomes
  - Utilities
  - Strategies
  - Solutions

		Player 2			
		Rock	Paper	Scissors	
		Rock	0,0	-1,1	1,-1
		Paper	1,-1	0,0	-1,1
		Scissors	-1,1	1,-1	0,0

# Game theory: basics

- Normal form
- Players
  - Actions
  - Outcomes
  - Utilities
  - Strategies
  - Solutions

		Player 2			
		Rock	Paper	Scissors	
		Rock	0,0	-1,1	1,-1
		Paper	-1,1	0,0	-1,1
		Scissors	1,-1	-1,1	0,0
Player 1			$\sigma_2(\text{Rock})$	$\sigma_2(\text{Paper})$	
				$\sigma_2(\text{Scissors})$	
				$\sigma_1(\text{Rock})$	
				$\sigma_1(\text{Paper})$	
				$\sigma_1(\text{Scissors})$	

# Game theory: basics

- Normal form
- Players
  - Actions
  - Outcomes
  - Utilities
  - Strategies
  - Solutions

		Player 2				
		Rock	Paper	Scissors		
		Rock	0,0	-1,1	1,-1	1/3
		Paper	-1,1	0,0	-1,1	1/3
		Scissors	1,-1	1,-1	0,0	1/3
Player 1			1/3	1/3	1/3	

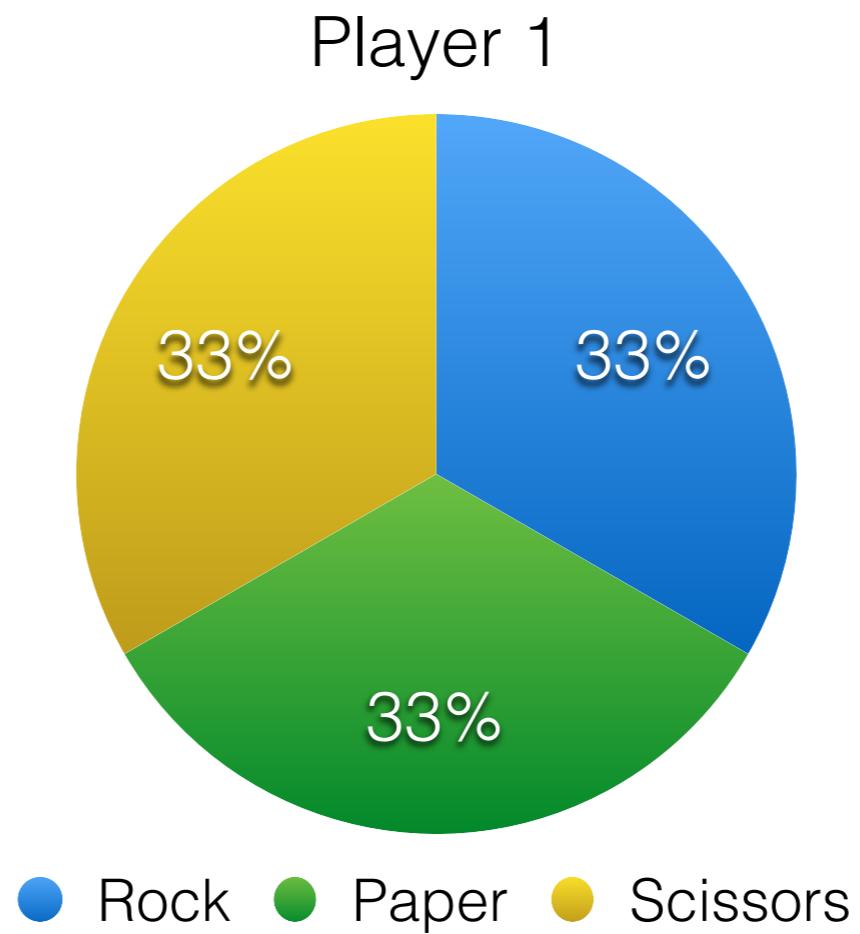
# Nash Equilibrium

A strategy profile  $(\sigma_1^*, \sigma_2^*)$  is a Nash equilibrium if and if:

- $\sigma_1^* \in \arg \max_{\sigma_1} \left\{ \sigma_1 U_1 \sigma_2^* \right\}$
- $\sigma_2^* \in \arg \max_{\sigma_2} \left\{ \sigma_1^* U_2 \sigma_2 \right\}$

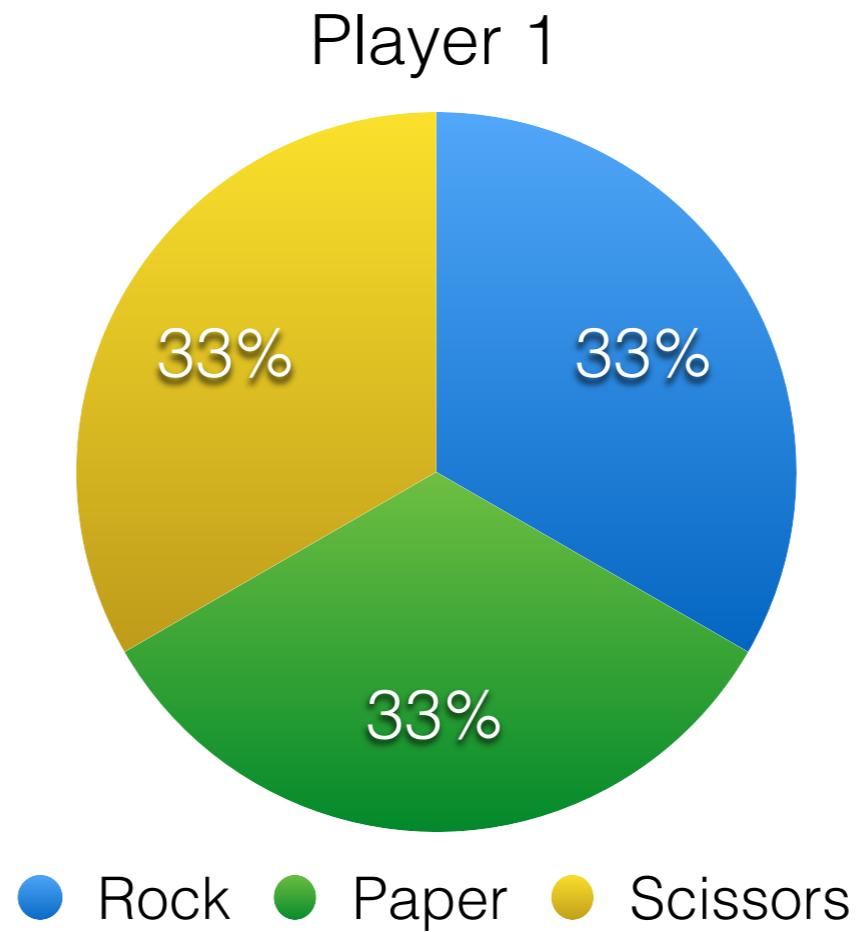
# Evolutionary dynamics

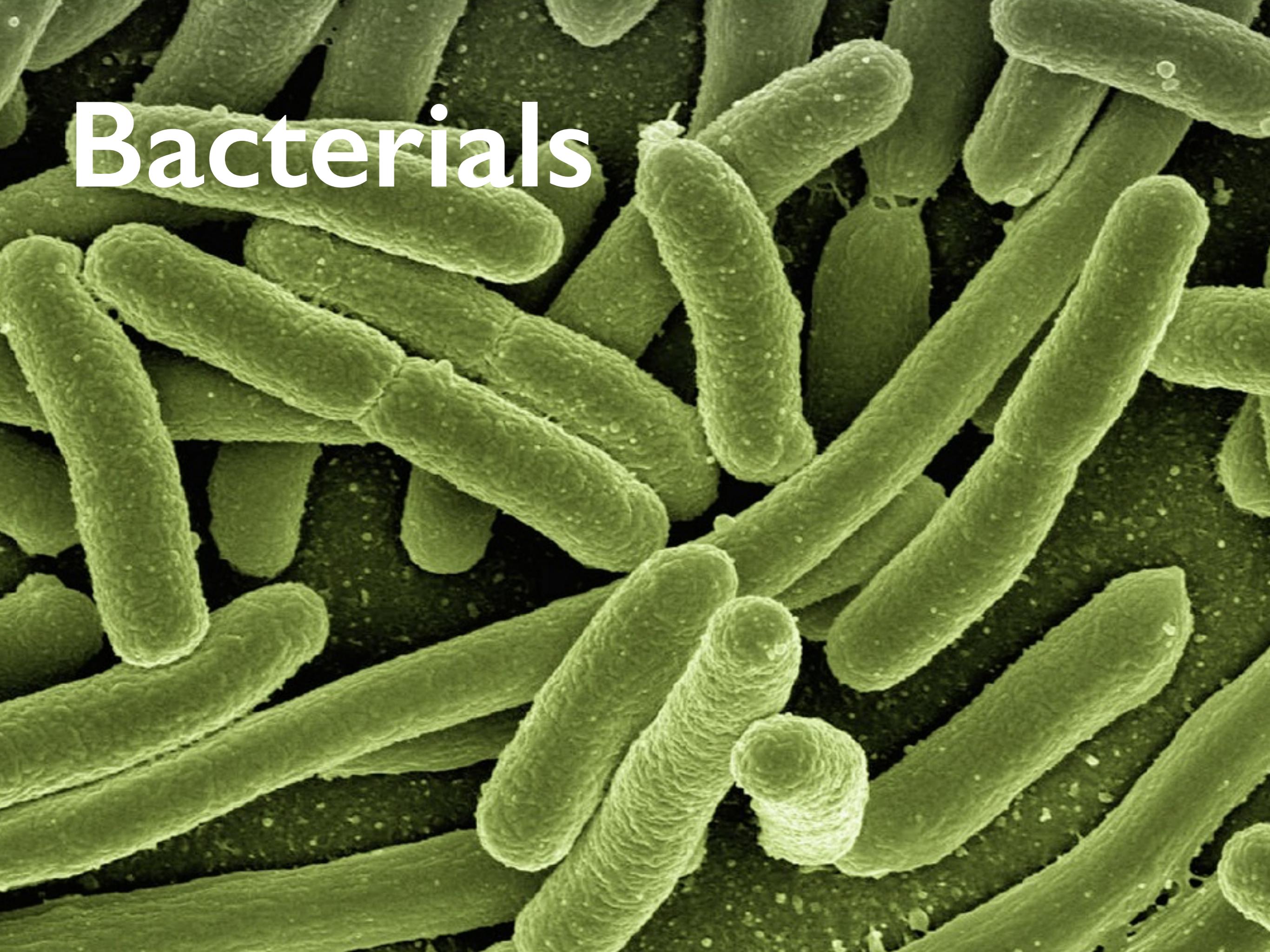
# An evolutionary interpretation



# An evolutionary interpretation

Infinite population  
of individuals

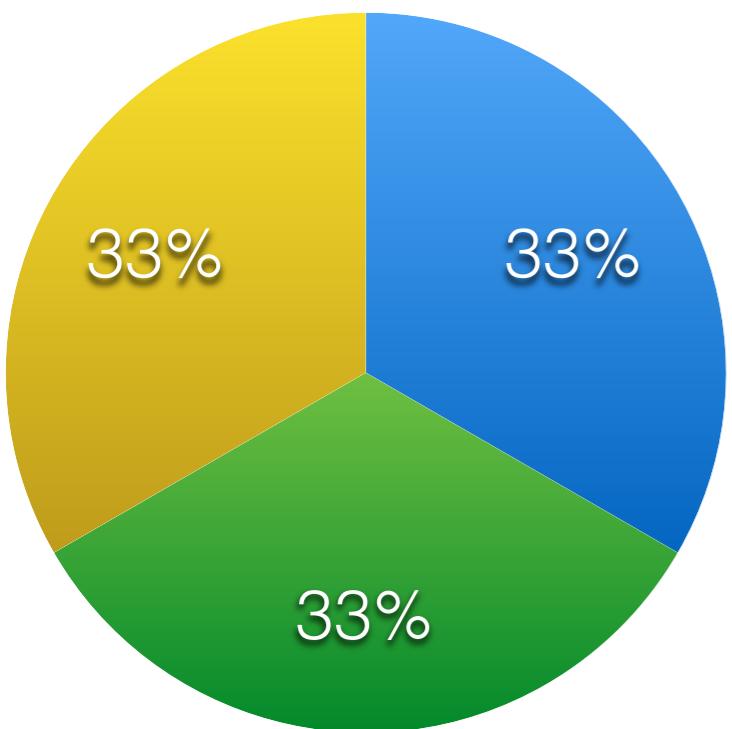


A scanning electron micrograph (SEM) showing a dense cluster of bacterial cells. The cells exhibit various morphologies: some are long and thin rods, while others are shorter and rounded cocci. The surface texture of the bacteria is visible, appearing somewhat granular or finely wrinkled. They are set against a dark, textured background that looks like a porous or fibrous material.

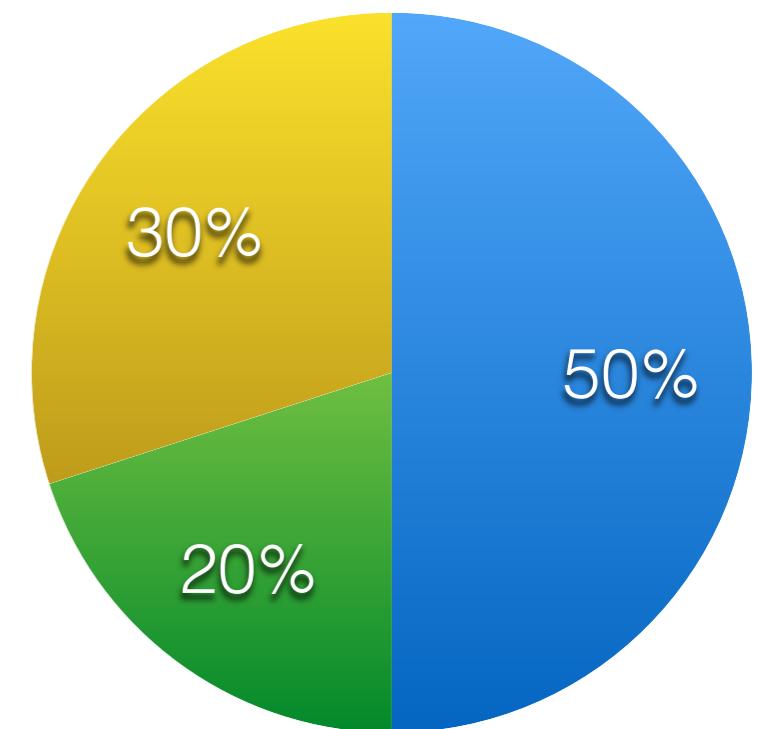
Bacterials

# An evolutionary interpretation

Player 1



Player 2

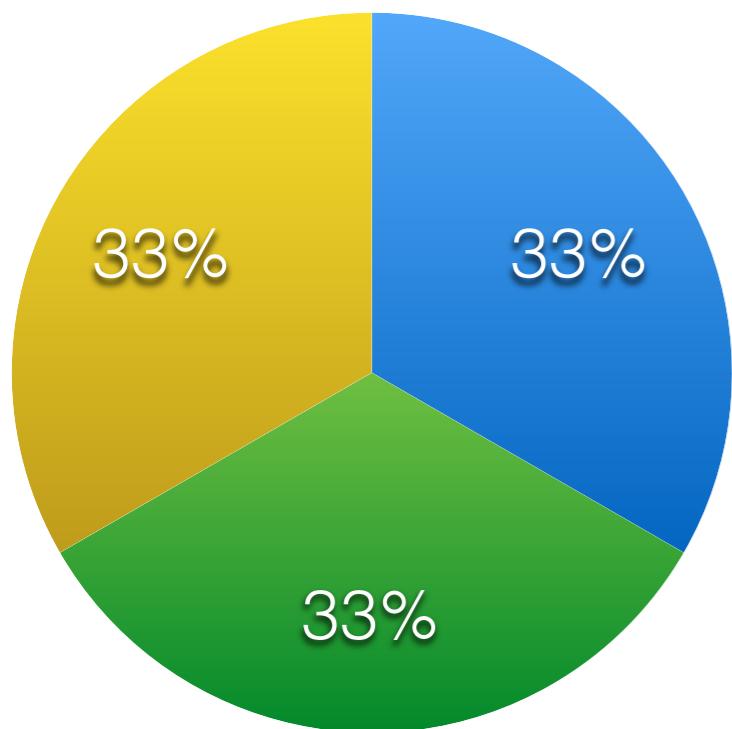


● Rock ● Paper ● Scissors

● Rock ● Paper ● Scissors

# An evolutionary interpretation

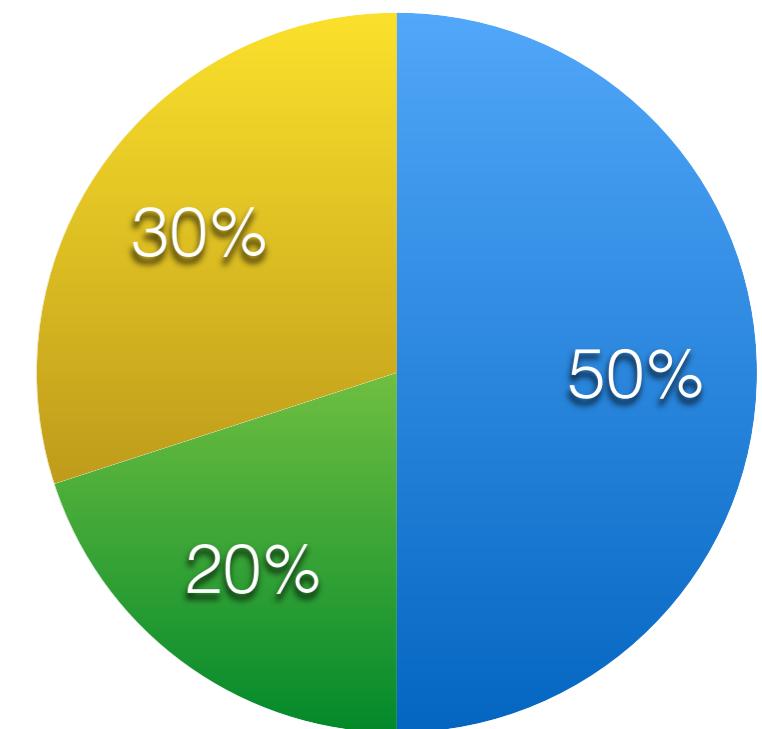
Player 1



● Rock ● Paper ● Scissors

	Rock	Paper	Scissors
Rock	0,0	-1,1	1,-1
Paper	1,-1	0,0	-1,1
Scissors	-1,1	1,-1	0,0

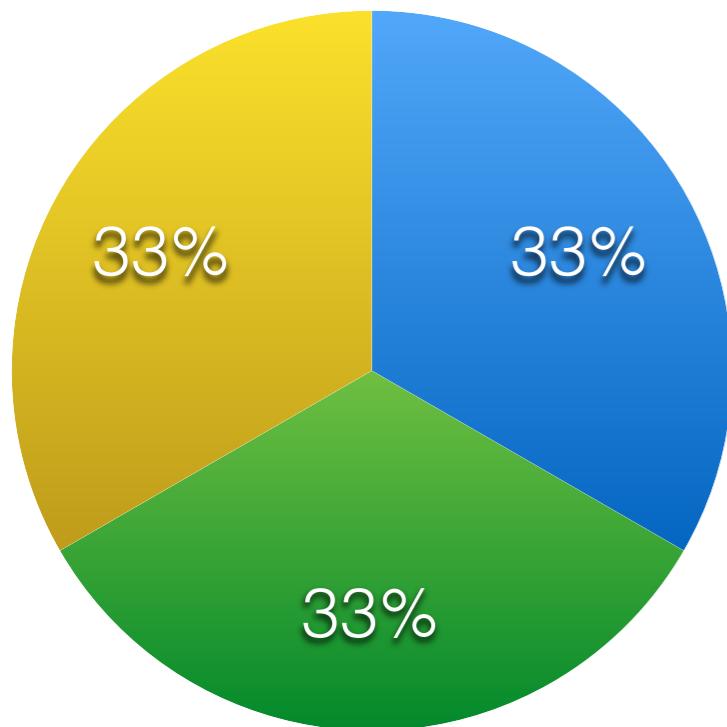
Player 2



● Rock ● Paper ● Scissors

# An evolutionary interpretation

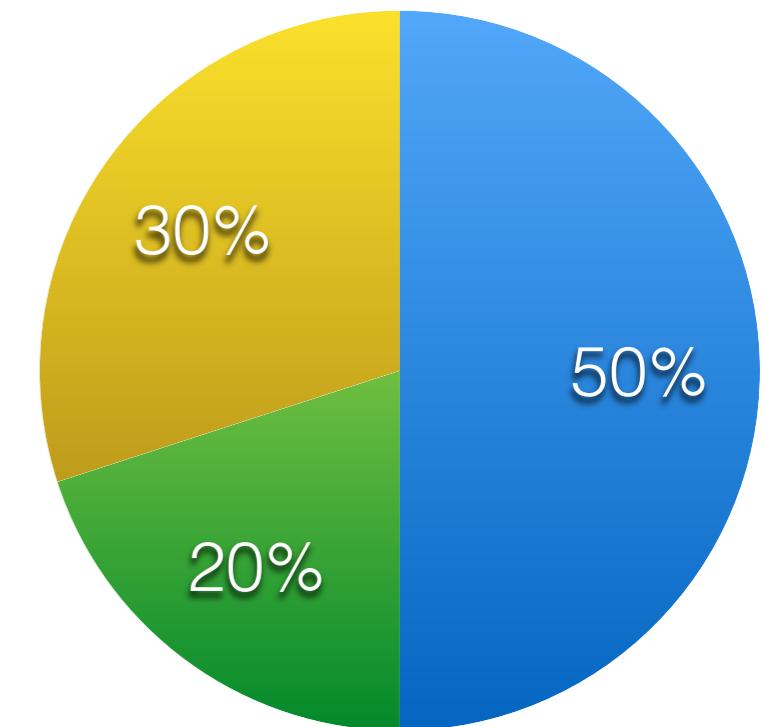
Player 1



● Rock ● Paper ● Scissors

	Rock	Paper	Scissors
Rock	0,0	-1,1	1,-1
Paper	1,-1	0,0	-1,1
Scissors	-1,1	1,-1	0,0

Player 2

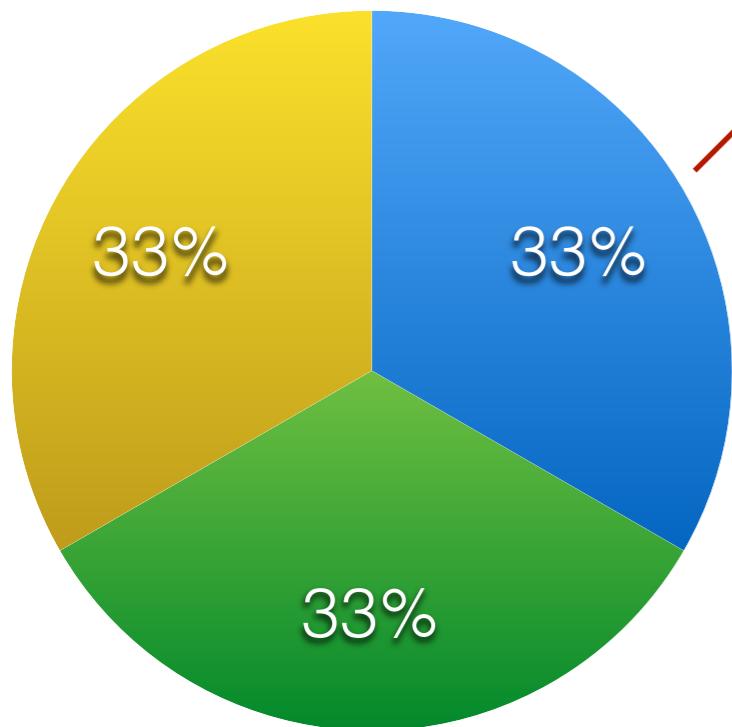


● Rock ● Paper ● Scissors

At each round, each individual of a population plays against each individual of the opponent's population weighted by the corresponding probability

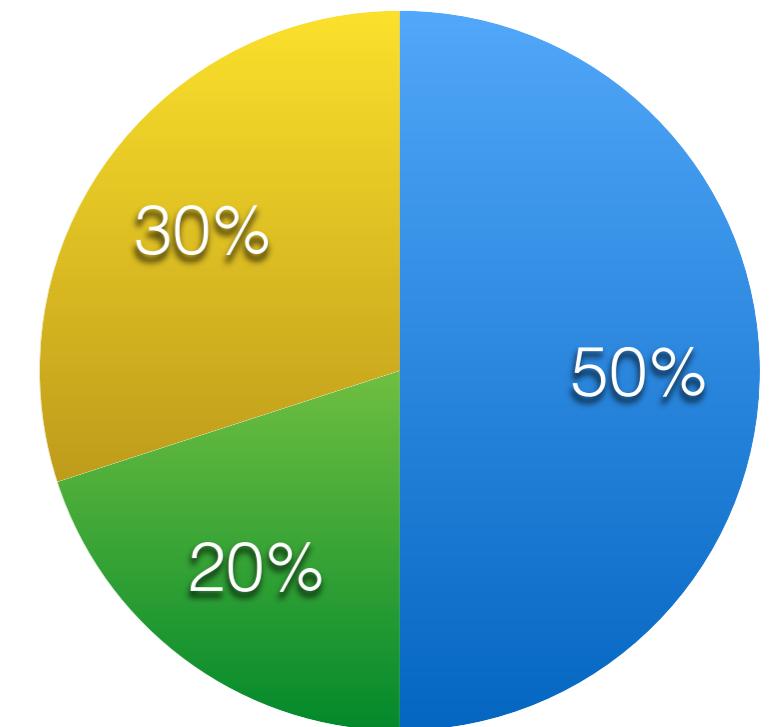
# An evolutionary interpretation

Player 1



$$\text{Utility} = 0 * \frac{1}{2} - 1 * \frac{1}{5} + 1 * \frac{3}{10} = 0.1$$

Player 2



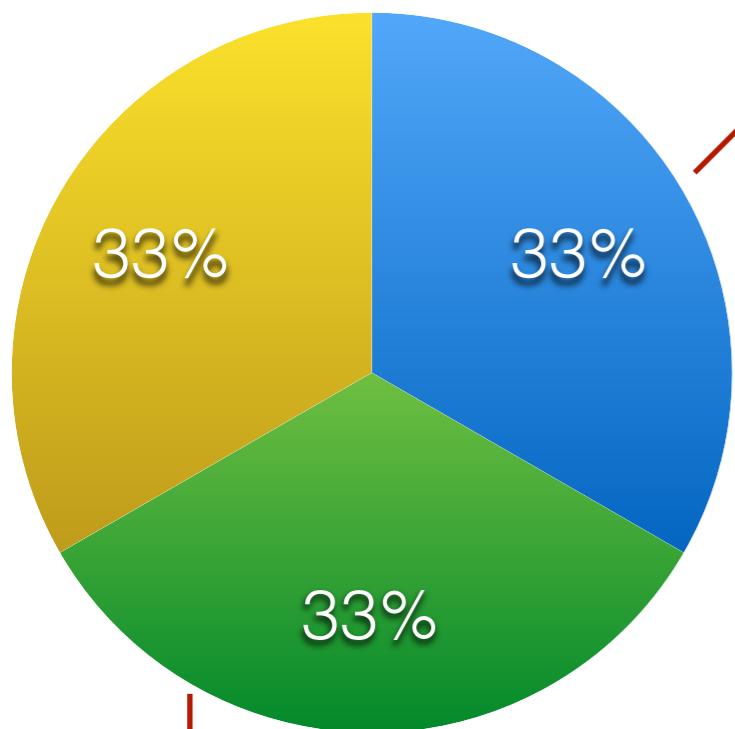
		Rock	Paper	Scissors
Rock	0,0	-1,1	1,-1	
Paper	1,-1	0,0	-1,1	
Scissors	-1,1	1,-1	0,0	

$\frac{1}{2}$     $\frac{1}{5}$     $\frac{3}{10}$

- Rock
- Paper
- Scissors

# An evolutionary interpretation

Player 1



Utility = 0.1

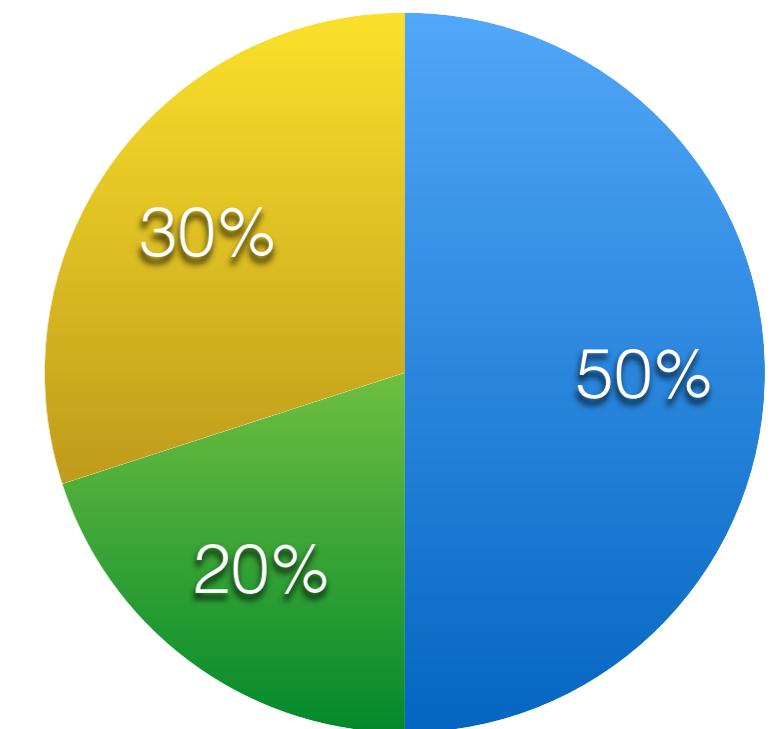
	Rock	Paper	Scissors
Rock	0,0	-1,1	1,-1
Paper	1,-1	0,0	-1,1
Scissors	-1,1	1,-1	0,0

Below the table are the values 1/2, 1/5, and 3/10.

- Rock
- Paper
- Scissors

$$\text{Utility} = 1 * 1/2 - 0 * 1/5 - 1 * 3/10 = 0.2$$

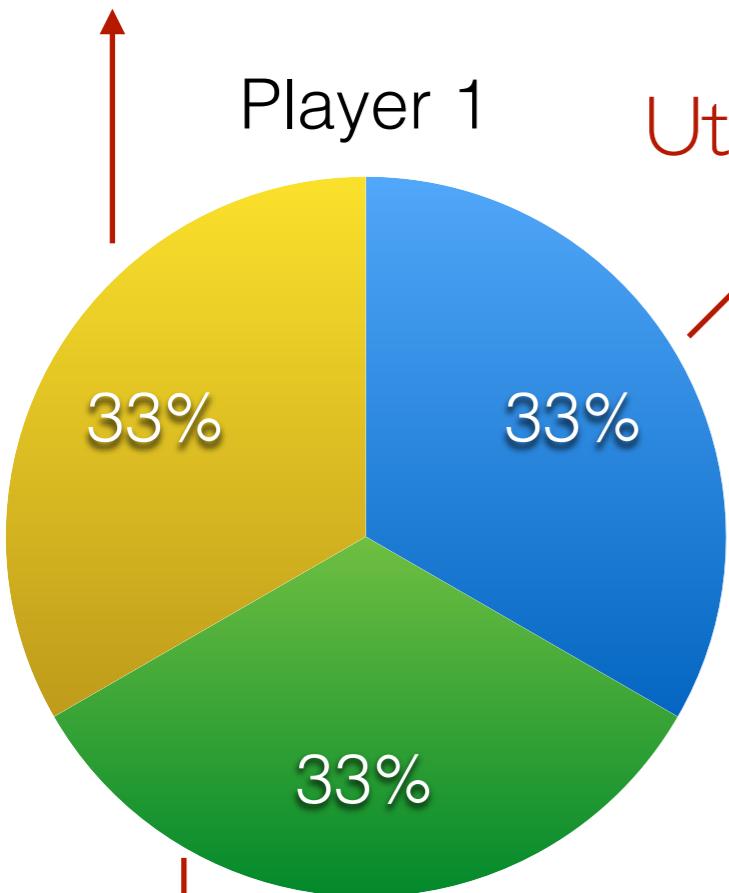
Player 2



- Rock
- Paper
- Scissors

# An evolutionary interpretation

$$\text{Utility} = -| * | / 2 - | * | / 5 - 0 * 3 / 10 = -0.3$$



Player 1

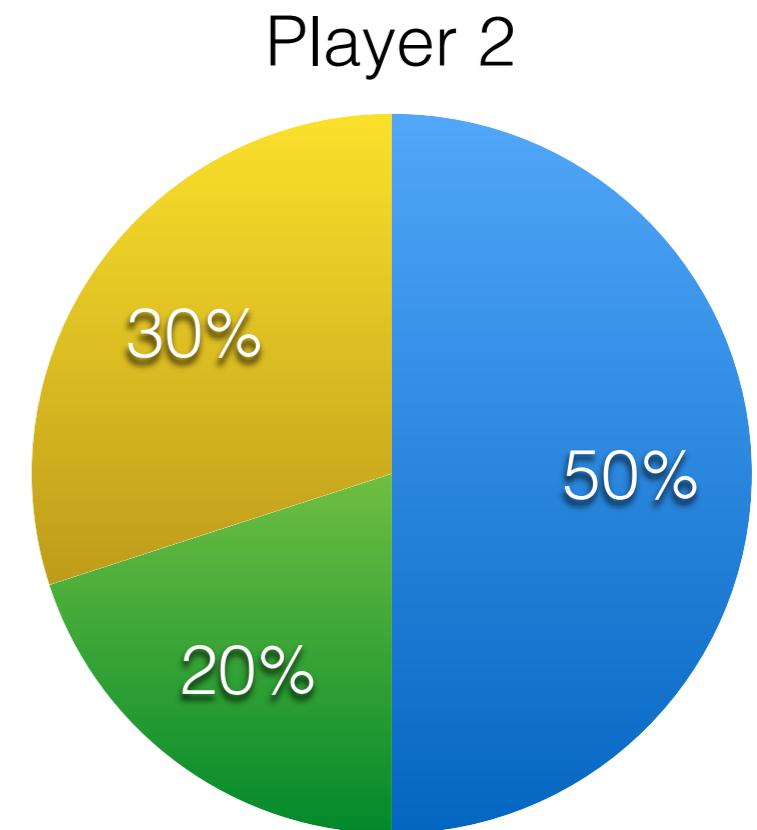
$$\text{Utility} = 0.1$$

		Rock	Paper	Scissors
Rock	0,0	-1,1	1,-1	
Paper	1,-1	0,0	-1,1	
Scissors	-1,1	1,-1	0,0	

$1/2$     $1/5$     $3/10$

- Rock
- Paper
- Scissors

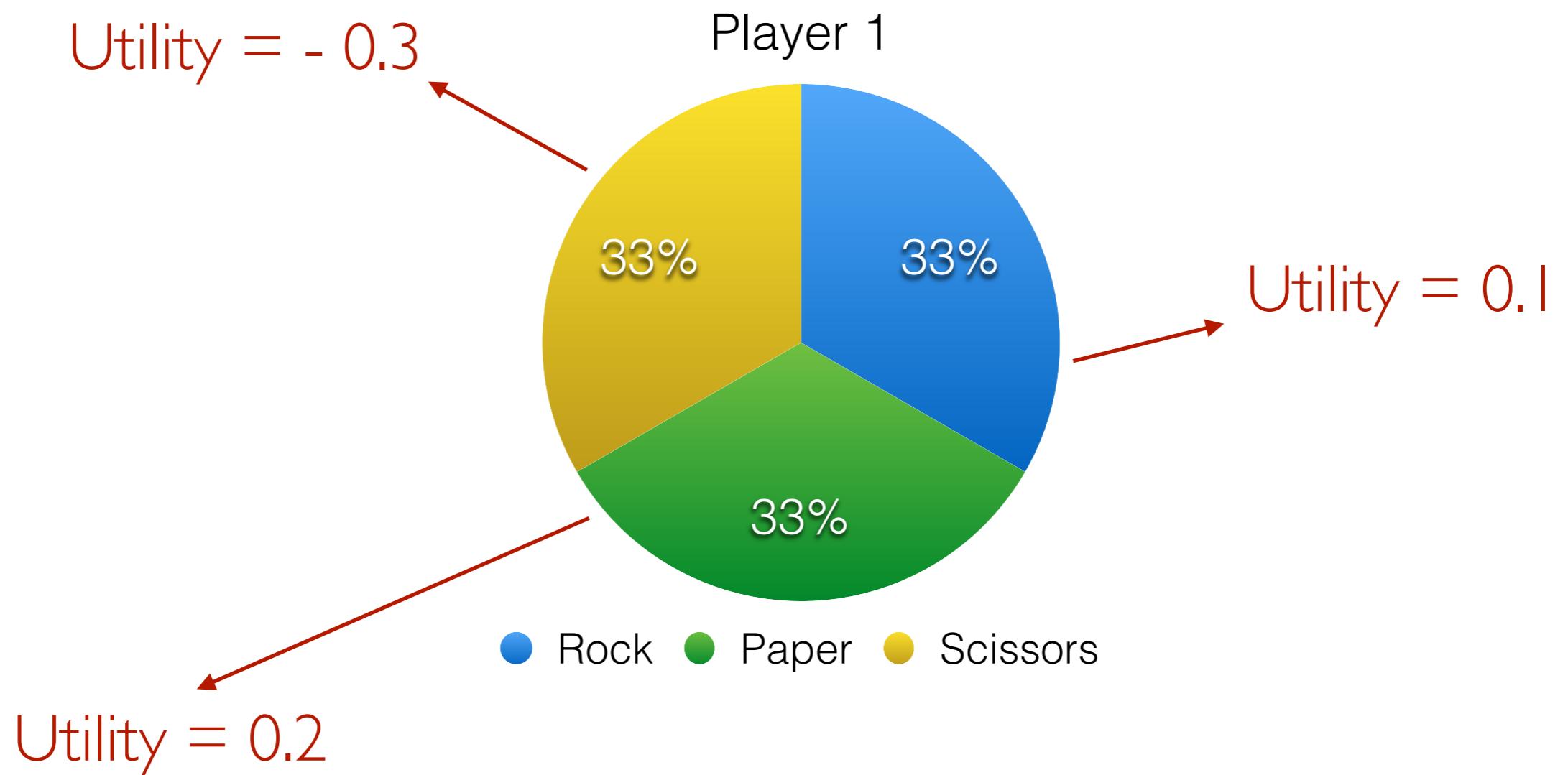
$$\text{Utility} = 0.2$$



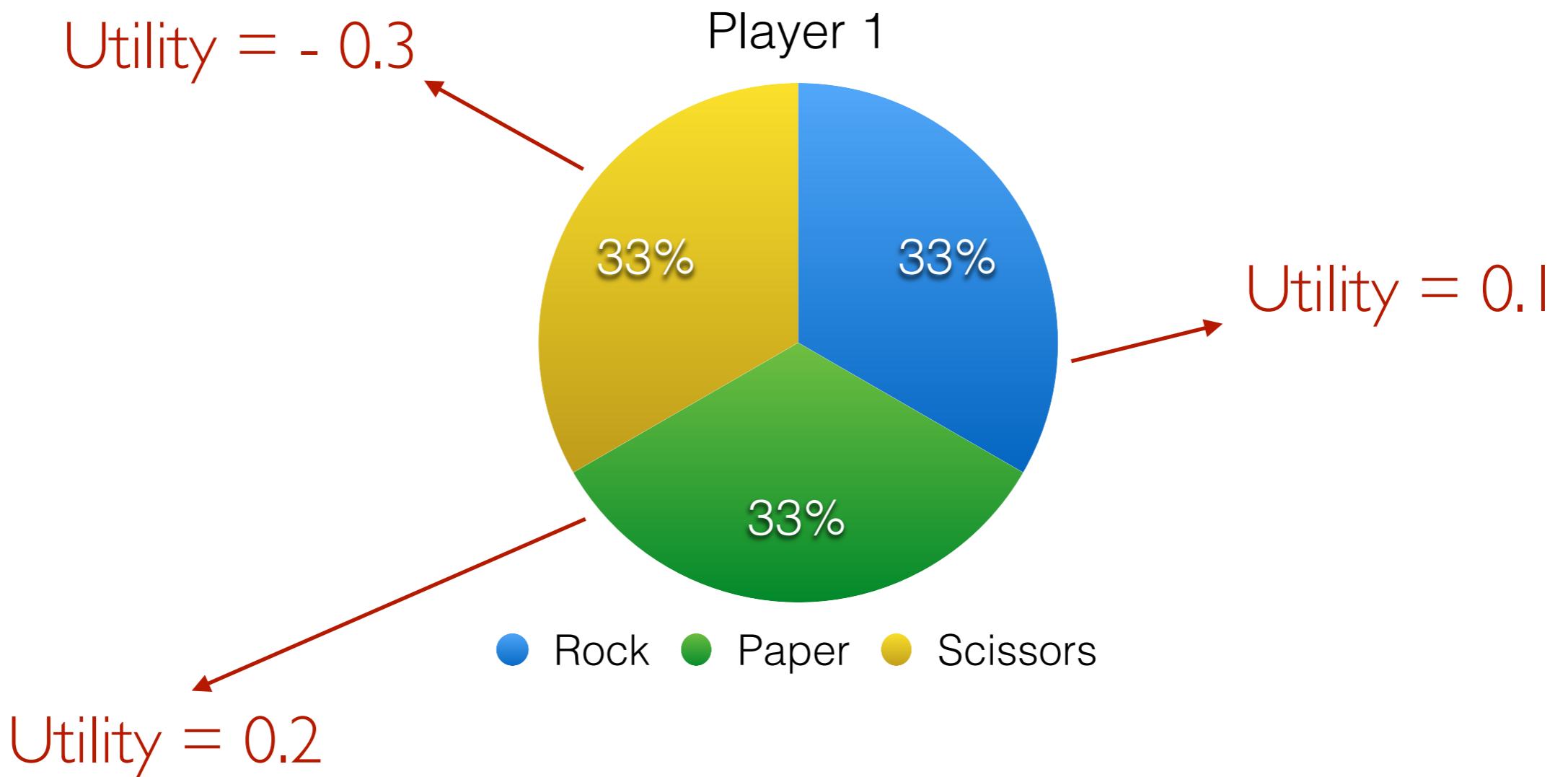
Player 2

- Rock
- Paper
- Scissors

# An evolutionary interpretation

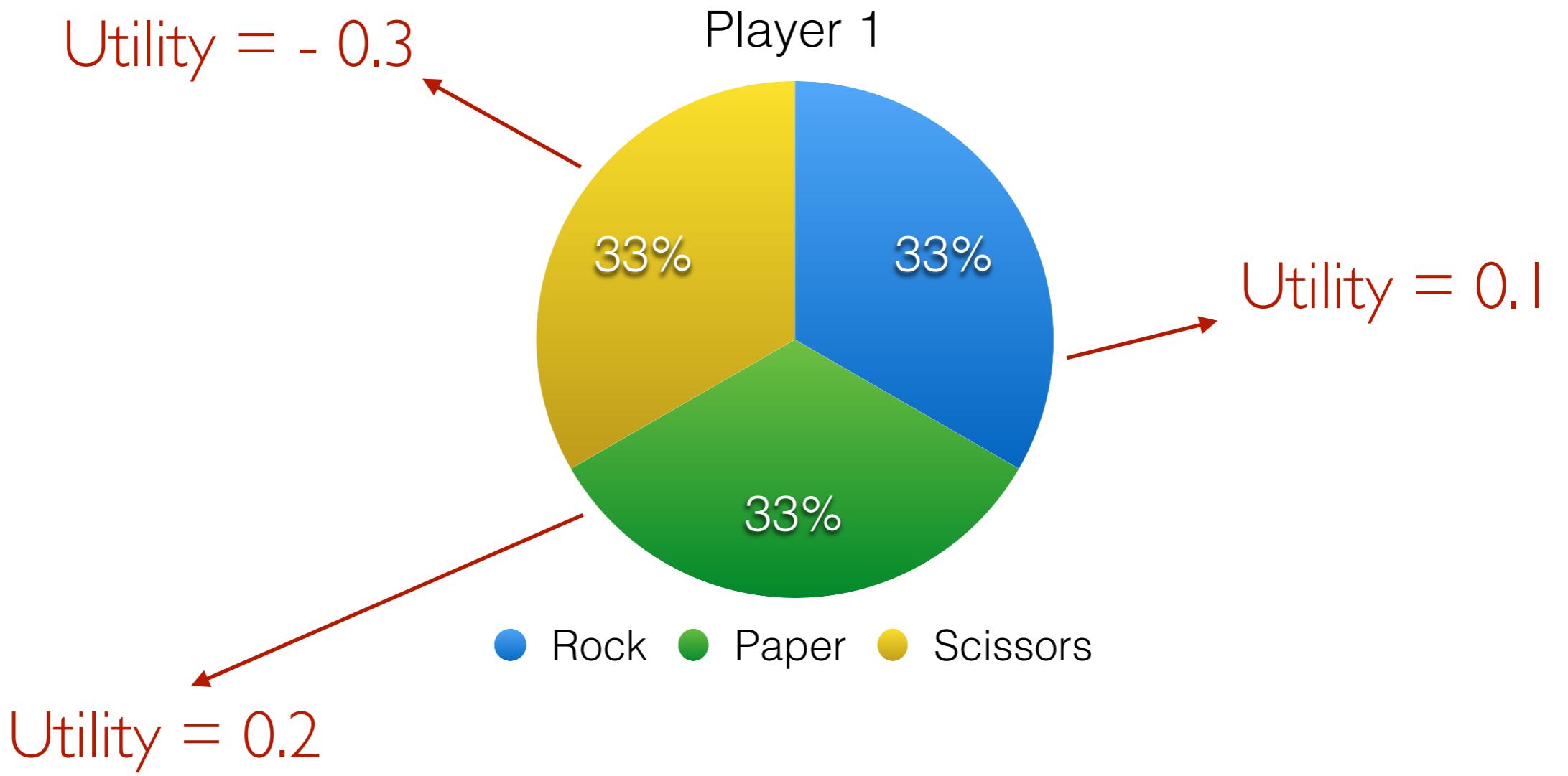


# An evolutionary interpretation



The average utility is  $1/3 * 0.1 + 1/3 * 0.2 - 0.3 * 1/3 = 0$

# An evolutionary interpretation

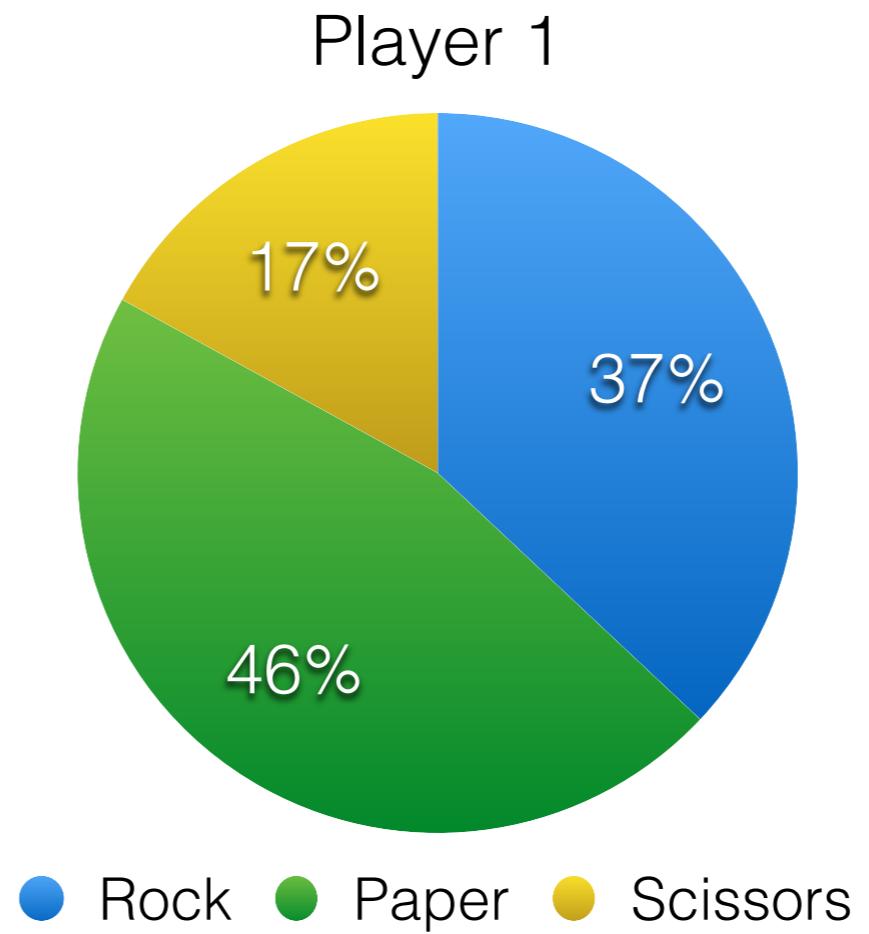


The average utility is 0

A **positive** (utility - average utility) leads to an **increase** of the population

A **negative** (utility - average utility) leads to a **decrease** of the population

# An evolutionary interpretation



New population after the replication

# Revision protocol

**Question:** how the populations change?

Replicator dynamics

$$\dot{\sigma}_1(a, t) = \sigma_1(a, t) \left( e_a U_1 \sigma_2(t) - \sigma_1(t) U_1 \sigma_2(t) \right)$$

Utility given by playing  $a$   
with a probability of  $\sigma_1$



Average population utility



# Example (I)

$$\sigma_1(t) = \begin{bmatrix} \sigma_1(\text{R}, t) & \sigma_1(\text{P}, t) & \sigma_1(\text{S}, t) \end{bmatrix}$$

$$\sigma_2(t) = \begin{bmatrix} \sigma_2(\text{R}, t) & \sigma_2(\text{P}, t) & \sigma_2(\text{S}, t) \end{bmatrix}$$

$$U_1 = \begin{bmatrix} 0 & -1 & 1 \\ 1 & 0 & -1 \\ -1 & 1 & 0 \end{bmatrix}$$

$$U_2 = \begin{bmatrix} 0 & 1 & -1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix}$$

$$U_1(\text{R}, \cdot) = \begin{bmatrix} 0 & 1 & -1 \end{bmatrix}$$

$$U_1(\text{P}, \cdot) = \begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$$

$$U_1(\text{S}, \cdot) = \begin{bmatrix} 1 & -1 & 0 \end{bmatrix}$$

# Example (2)

$$\dot{\sigma}_1(\mathbf{R}, t) = \sigma_1(\mathbf{R}, t) \begin{pmatrix} U_1(\mathbf{R}, \cdot) \sigma_2(t) & -\sigma_1(t) U_1 \sigma_2(t) \end{pmatrix}$$

$$\dot{\sigma}_1(\mathbf{P}, t) = \sigma_1(\mathbf{P}, t) \begin{pmatrix} U_1(\mathbf{P}, \cdot) \sigma_2(t) & -\sigma_1(t) U_1 \sigma_2(t) \end{pmatrix}$$

$$\dot{\sigma}_1(\mathbf{S}, t) = \sigma_1(\mathbf{S}, t) \begin{pmatrix} U_1(\mathbf{S}, \cdot) \sigma_2(t) & -\sigma_1(t) U_1 \sigma_2(t) \end{pmatrix}$$

# Evolutionary Stable Strategies

A strategy is an ESS if it is immune to invasion by mutant strategies, given that the mutants initially occupy a small fraction of population

Every ESS is an asymptotically stable fixed point of the replicator dynamics

# Evolutionary Stable Strategies (I)

A strategy is an *ESS* if it is immune to invasion by mutant strategies, given that the mutants initially occupy a small fraction of population

Every ESS is an asymptotically stable fixed point of the replicator dynamics



While a NE always exists,  
an ESS may not exist

# Evolutionary Stable Strategies (2)

Given a Nash equilibrium, a small perturbation could make the equilibrium unstable

		Player 2	
		Action 1	Action 2
		Action 1	1, 1
Player 1		Action 2	0, 0
		Action 2	0, 0

Nash equilibrium

# Evolutionary Stable Strategies (2)

Given a Nash equilibrium, a small perturbation could make the equilibrium unstable

		Player 2	
		Action 1	Action 2
		Action 1	1, 1
Player 1		0, 0	$\epsilon$
		0, 0	$1 - \epsilon$

# Evolutionary Stable Strategies (2)

Given a Nash equilibrium, a small perturbation could make the equilibrium unstable

The only best  
response is  
Action 1

		Player 2	
		Action 1	Action 2
		Action 1	1, 1
		Action 2	0, 0
Player 1		Action 1	$\epsilon$
		Action 2	$1 - \epsilon$

# Evolutionary Stable Strategies (2)

Given a Nash equilibrium, a small perturbation could make the equilibrium unstable

		Player 2		
		Action 1	Action 2	
		ESS		
Player 1	Action 1	1, 1	0, 0	$1 - \epsilon$
	Action 2	0, 0	0, 0	$\epsilon$

$1 - \epsilon$        $\epsilon$

# Prisoner's dilemma

		Player 2	
		Cooperate	Defeat
		Cooperate	3,3
		Defeat	5,0
		Defeat	1,1

# Prisoner's dilemma

Player 2

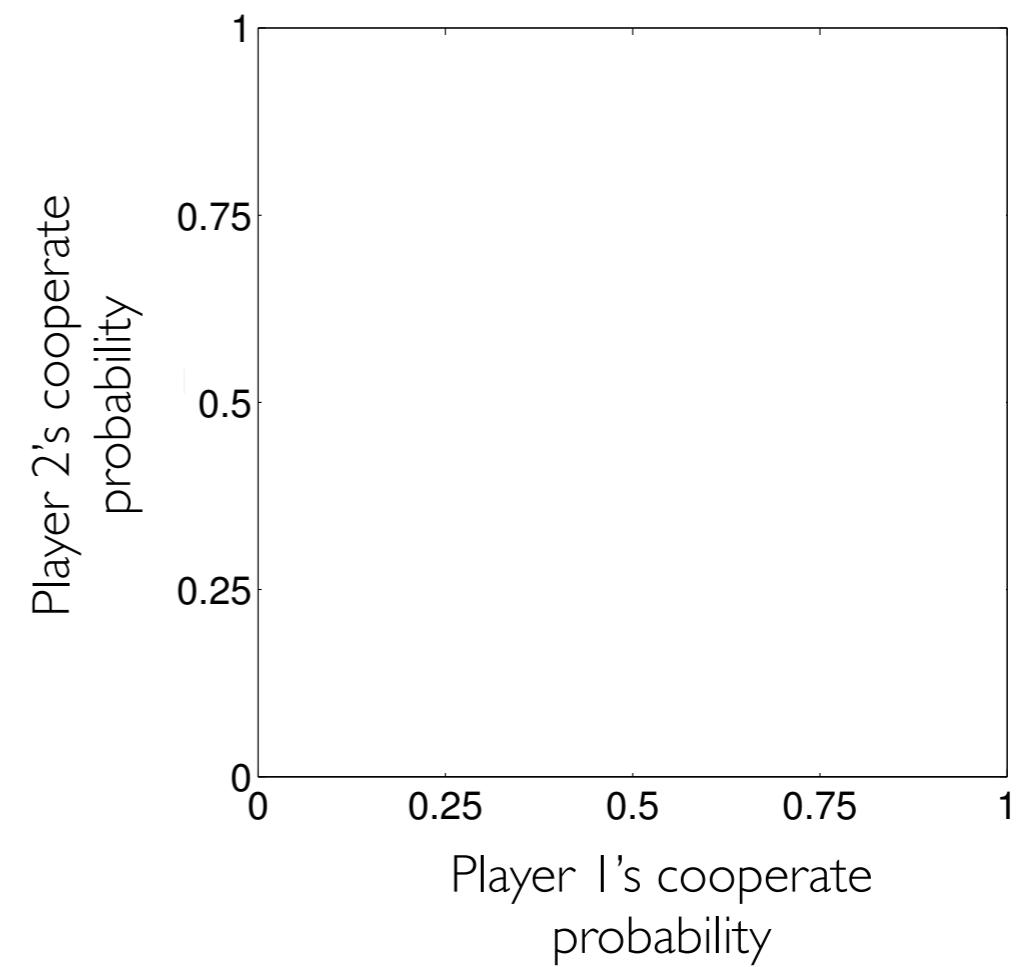
		Cooperate	Defeat
		Cooperate	0,5
Player 1	Cooperate	3,3	
	Defeat	5,0	1,1



# Prisoner's dilemma

Player 2

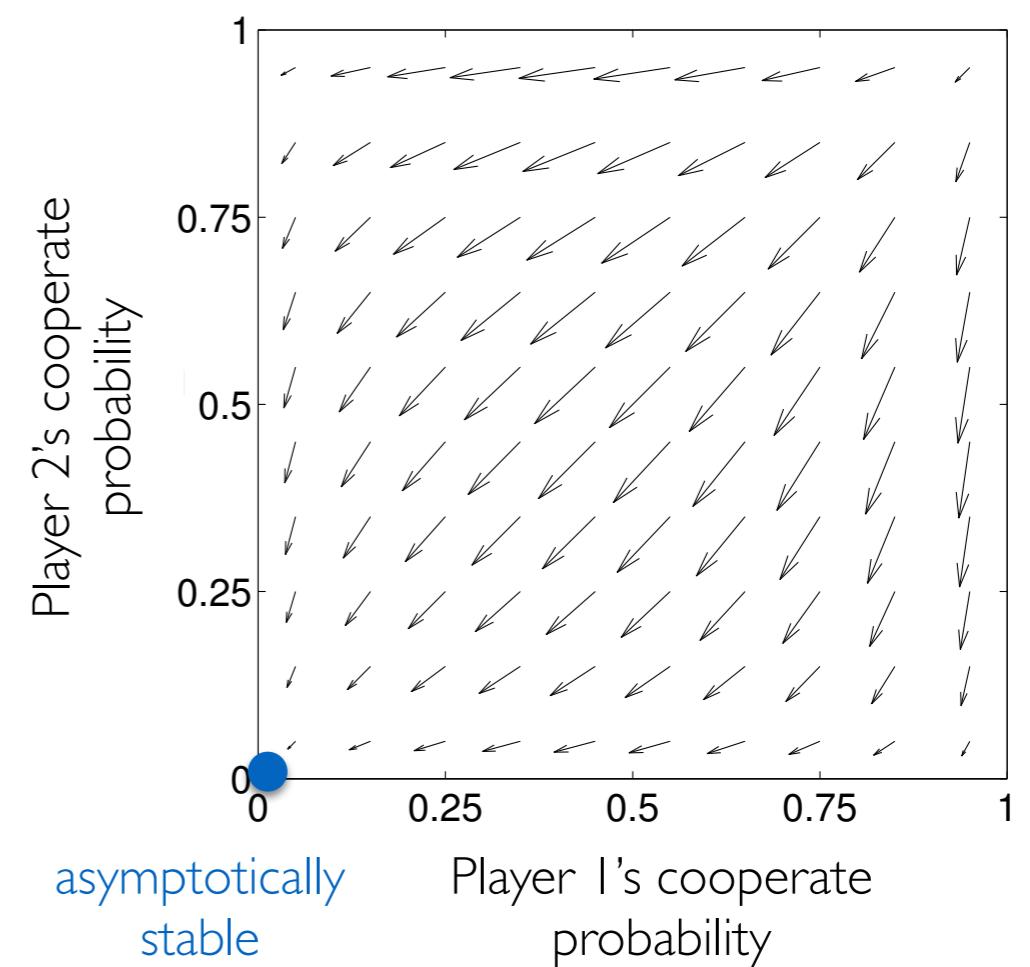
		Cooperate	Defeat
		Cooperate	3,3
Player 1	Cooperate	5,0	0,5
	Defeat	1,1	



# Prisoner's dilemma

Player 2

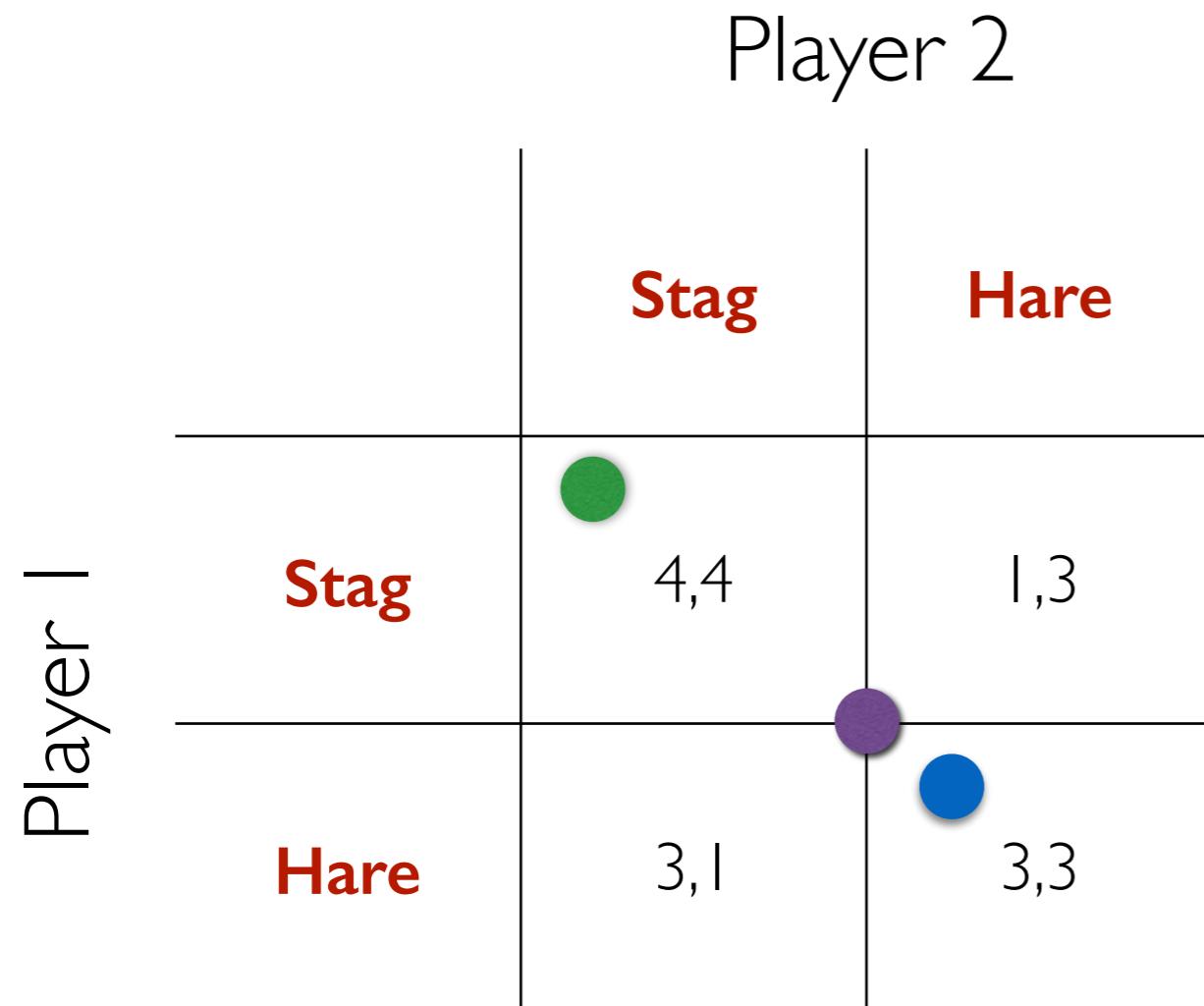
		Cooperate	Defeat
		Cooperate	3,3
Player 1	Cooperate	0,5	
	Defeat	5,0	1,1



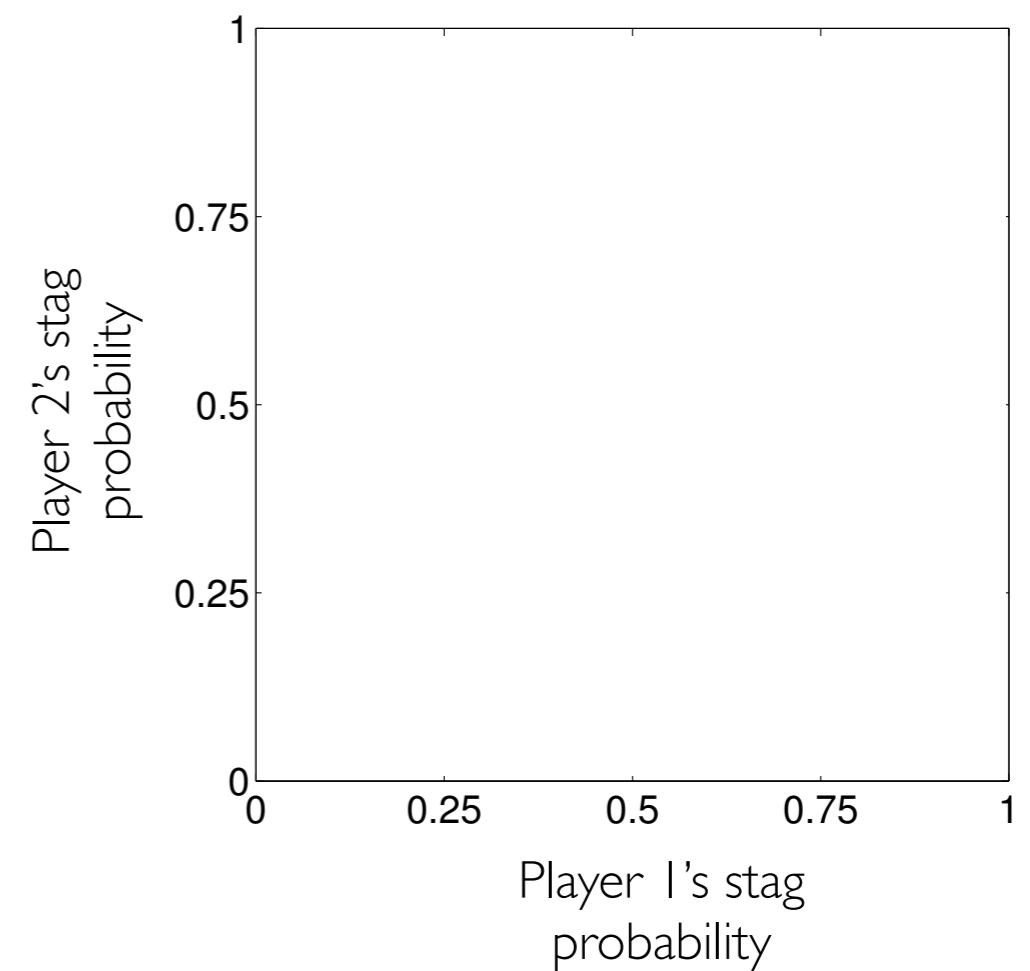
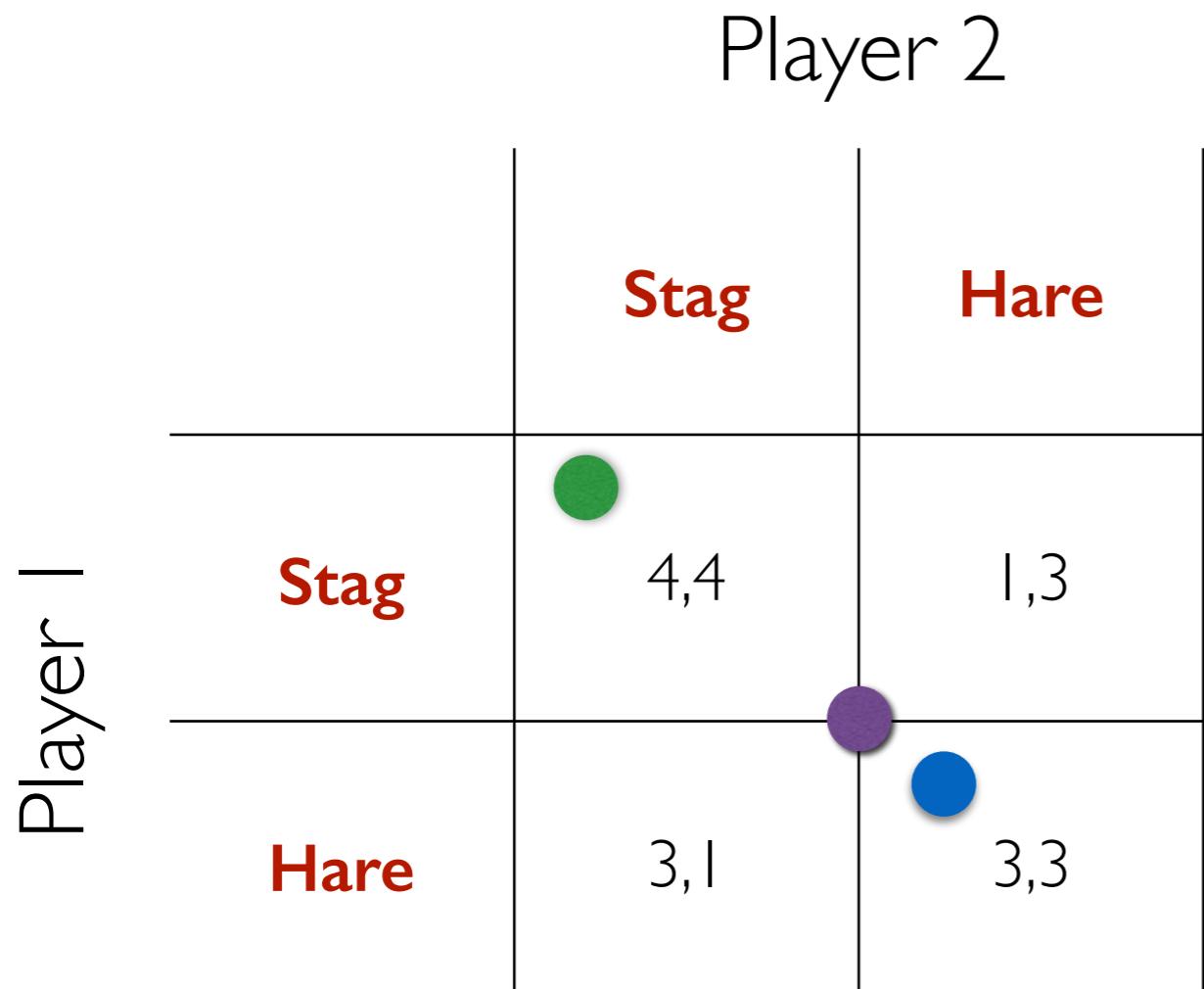
# Stag hunt

		Player 2	
		Stag	Hare
		Stag	4,4
		Hare	1,3
		Hare	3,1
		Hare	3,3

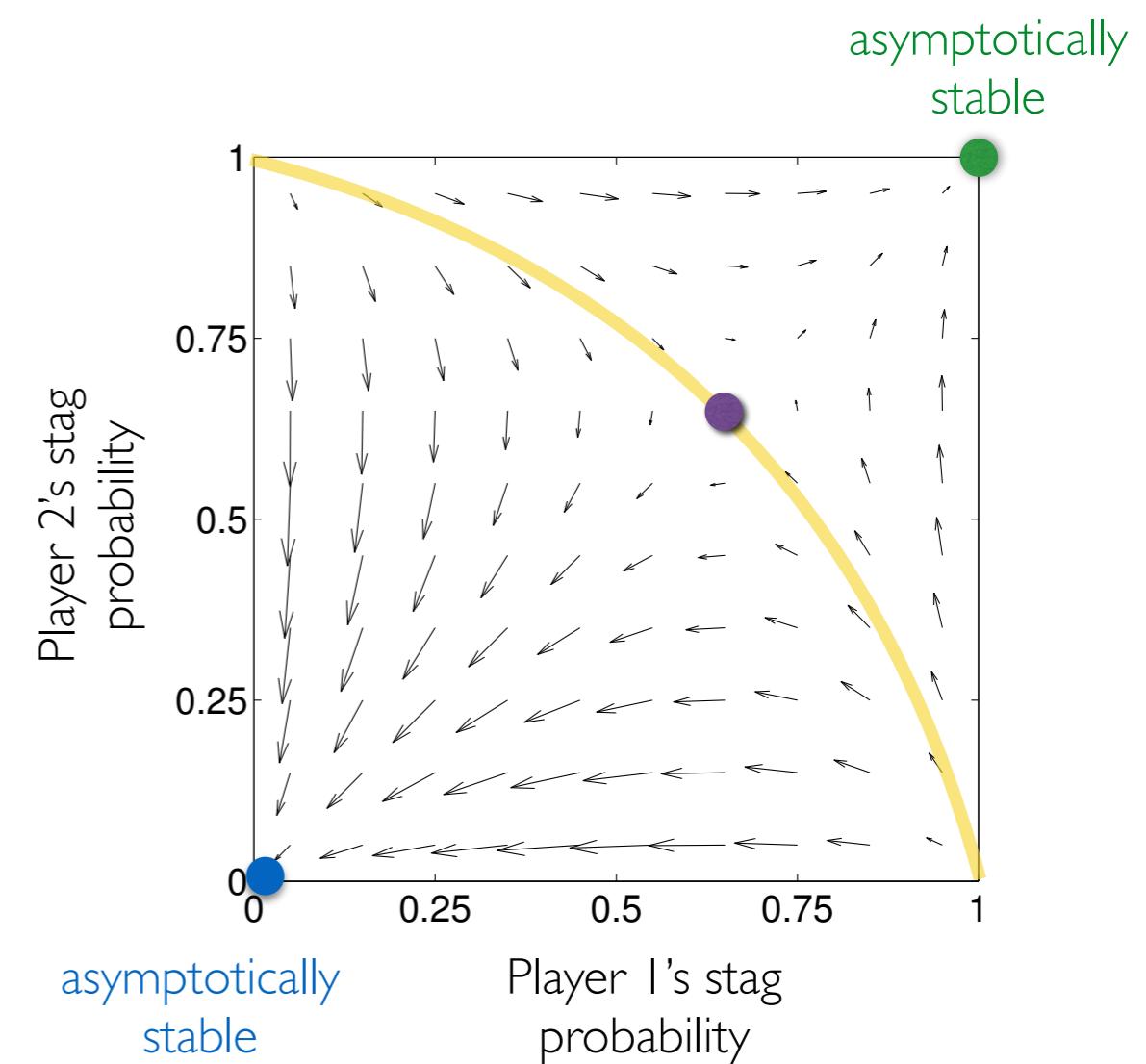
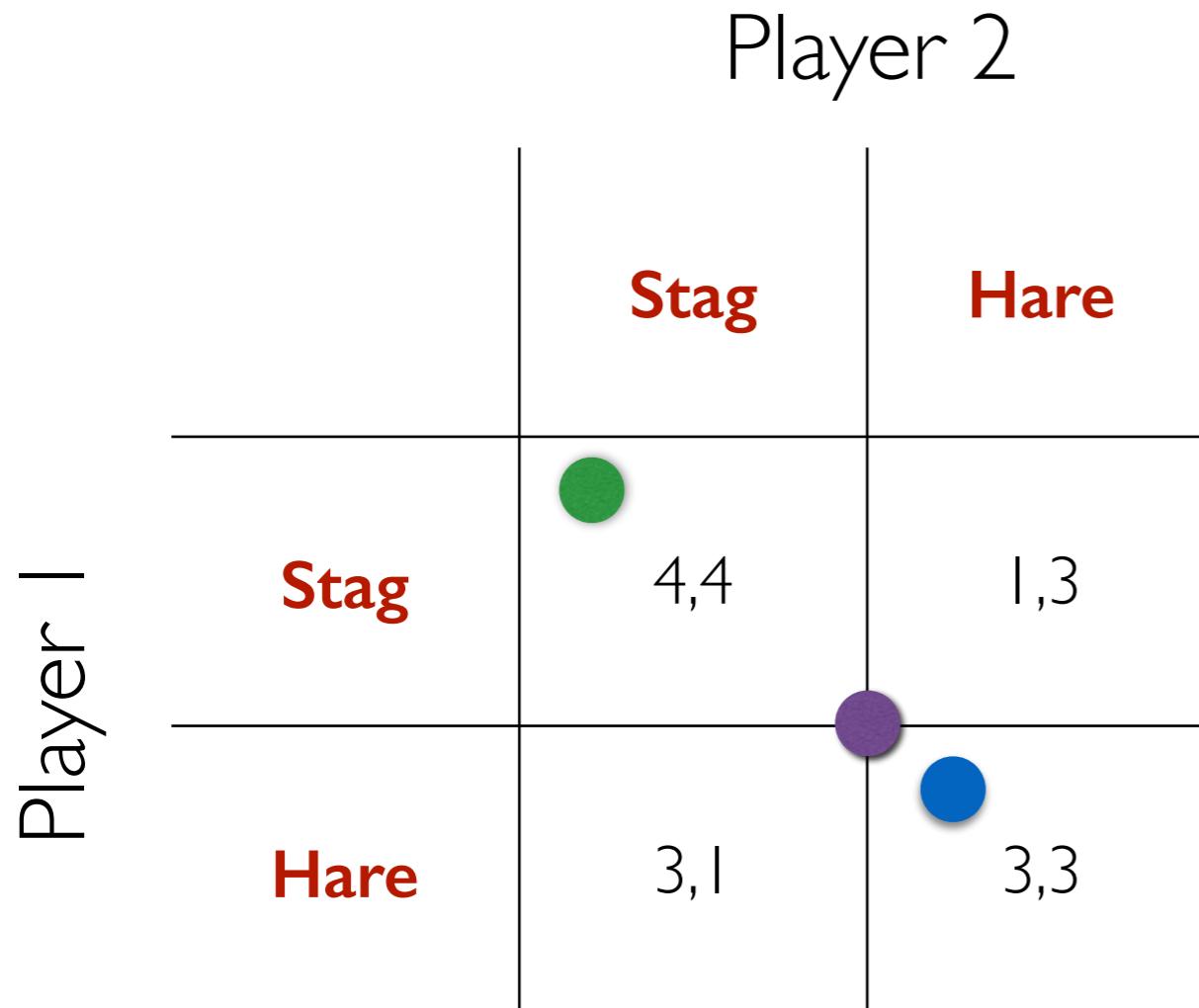
# Stag hunt



# Stag hunt



# Stag hunt



# Matching pennies

		Player 2	
		Head	Tail
		Head	0,1
		Tail	1,0
		Head	1,0
		Tail	0,1

# Matching pennies

Player 2

	Head	Tail
Head	0,1	1,0
Tail	1,0	0,1

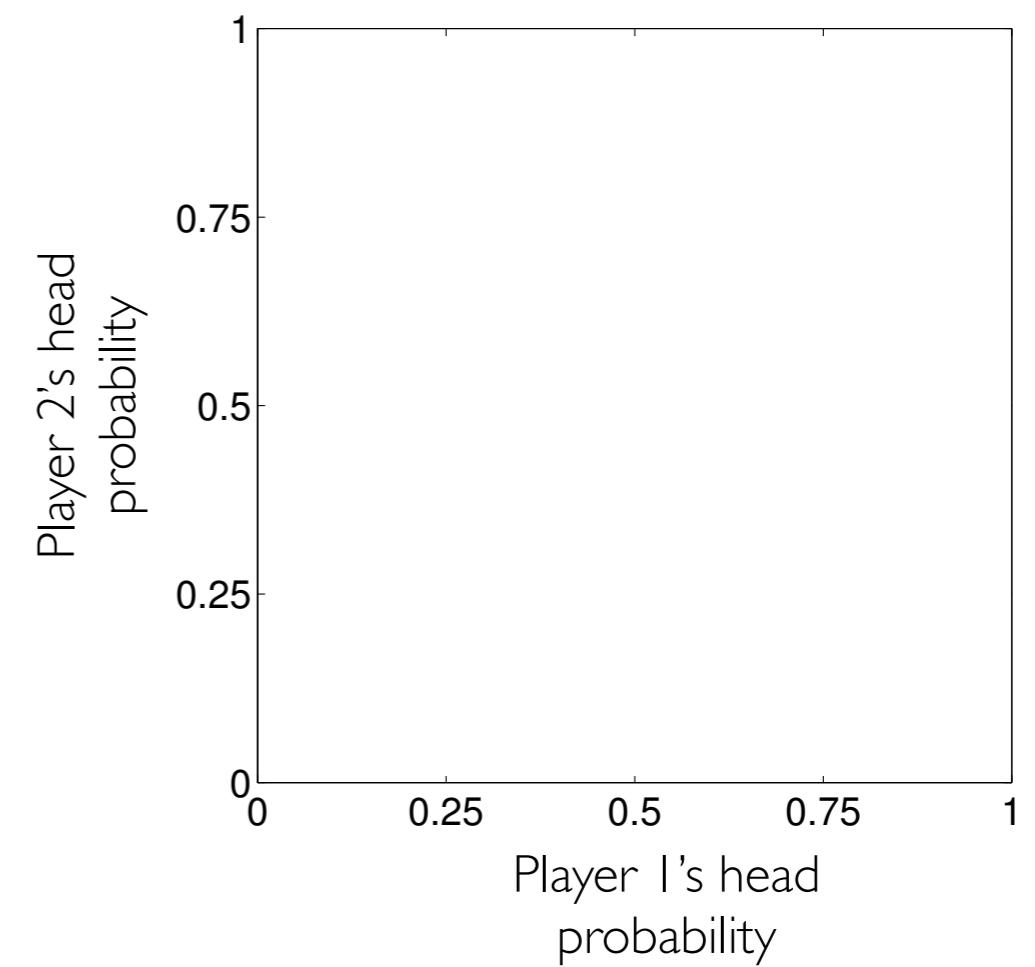
Player 1



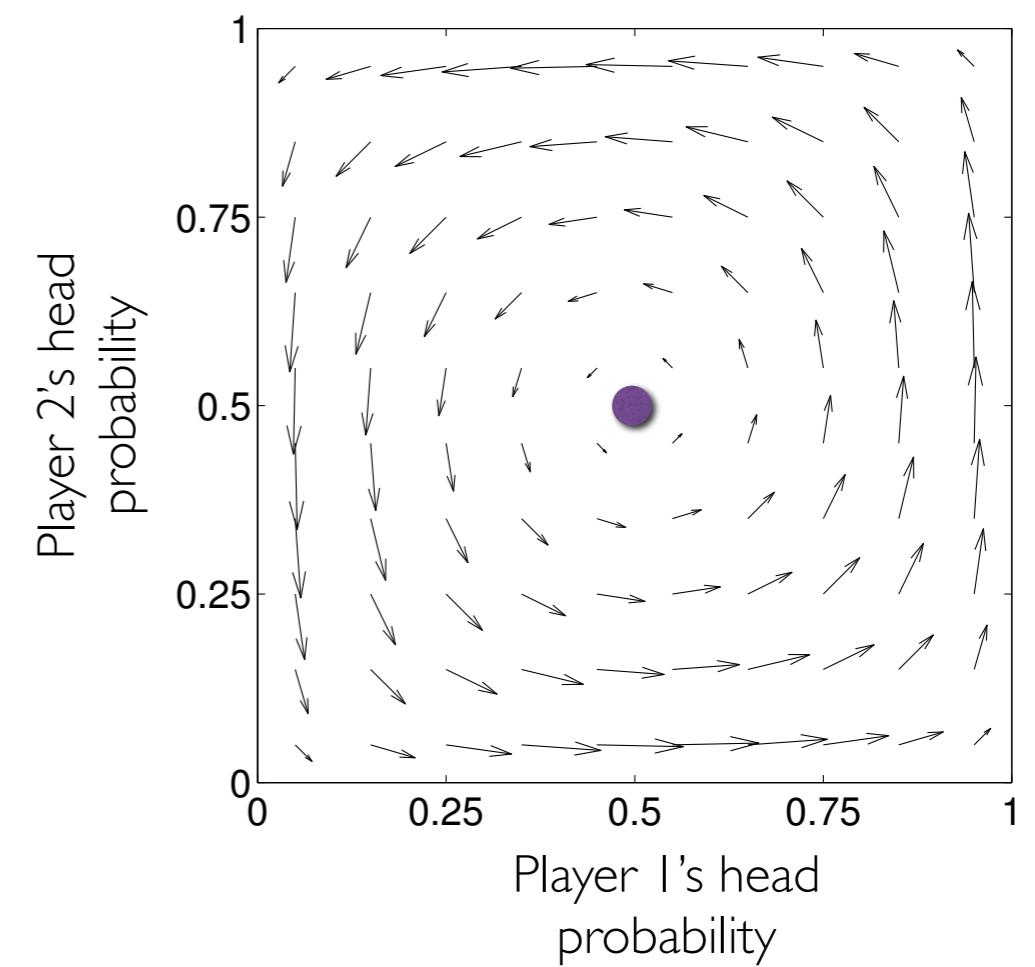
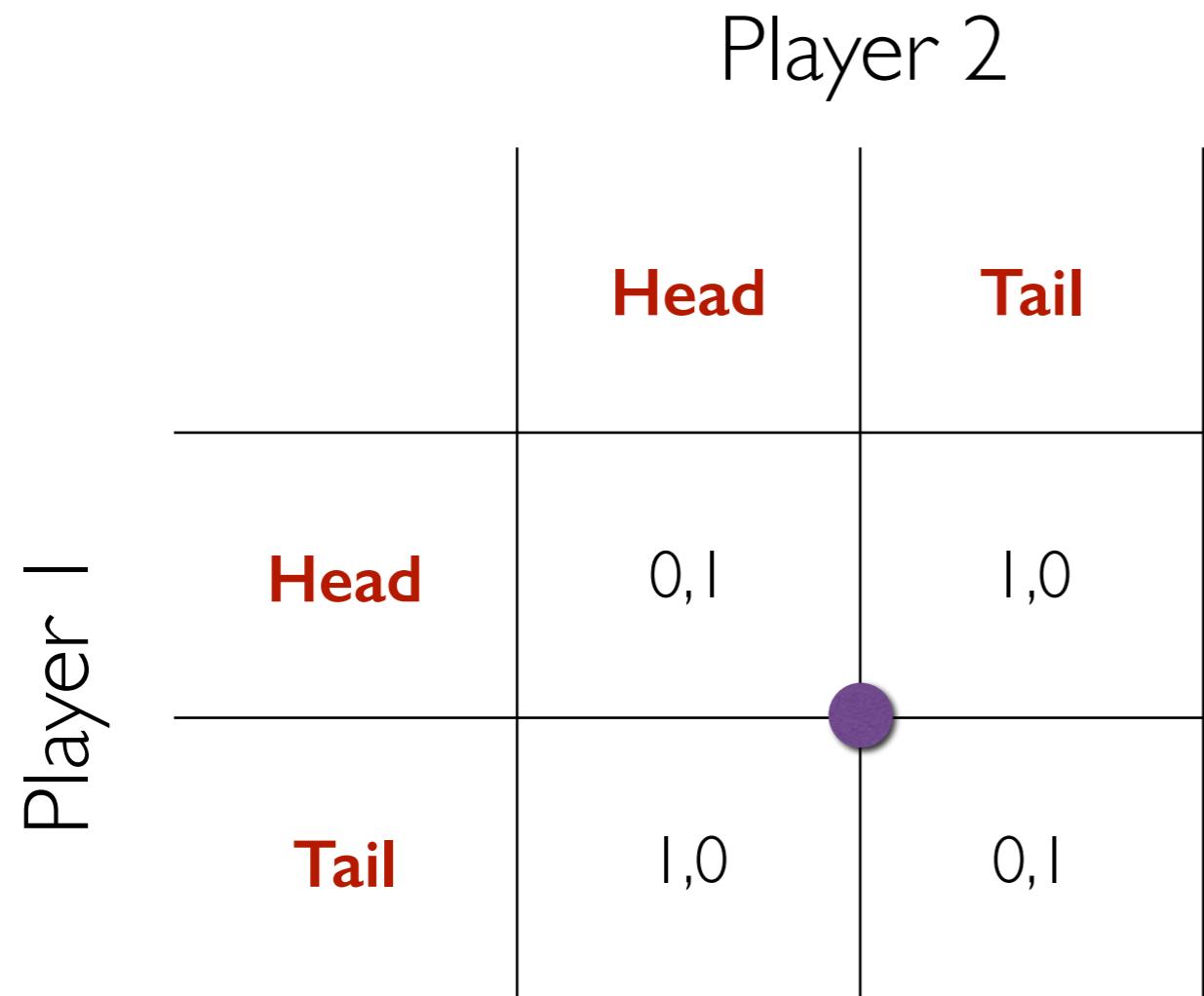
# Matching pennies

Player 2

	Head	Tail
Head	0,1	1,0
Tail	1,0	0,1



# Matching pennies



# Starting point

$$\dot{\sigma}_1(a, t) = \boxed{\sigma_1(a, t)} \left( e_a U_1 \sigma_2(t) - \sigma_1(t) U_1 \sigma_2(t) \right)$$

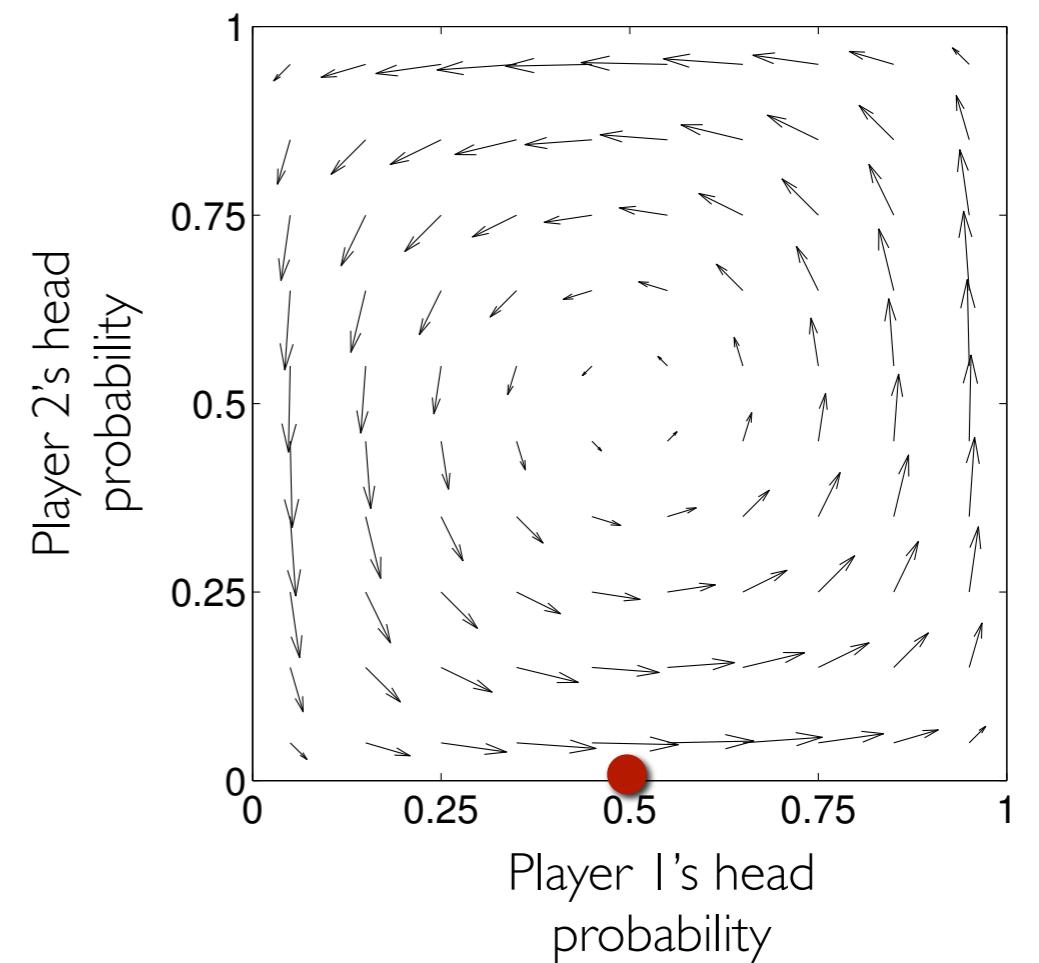
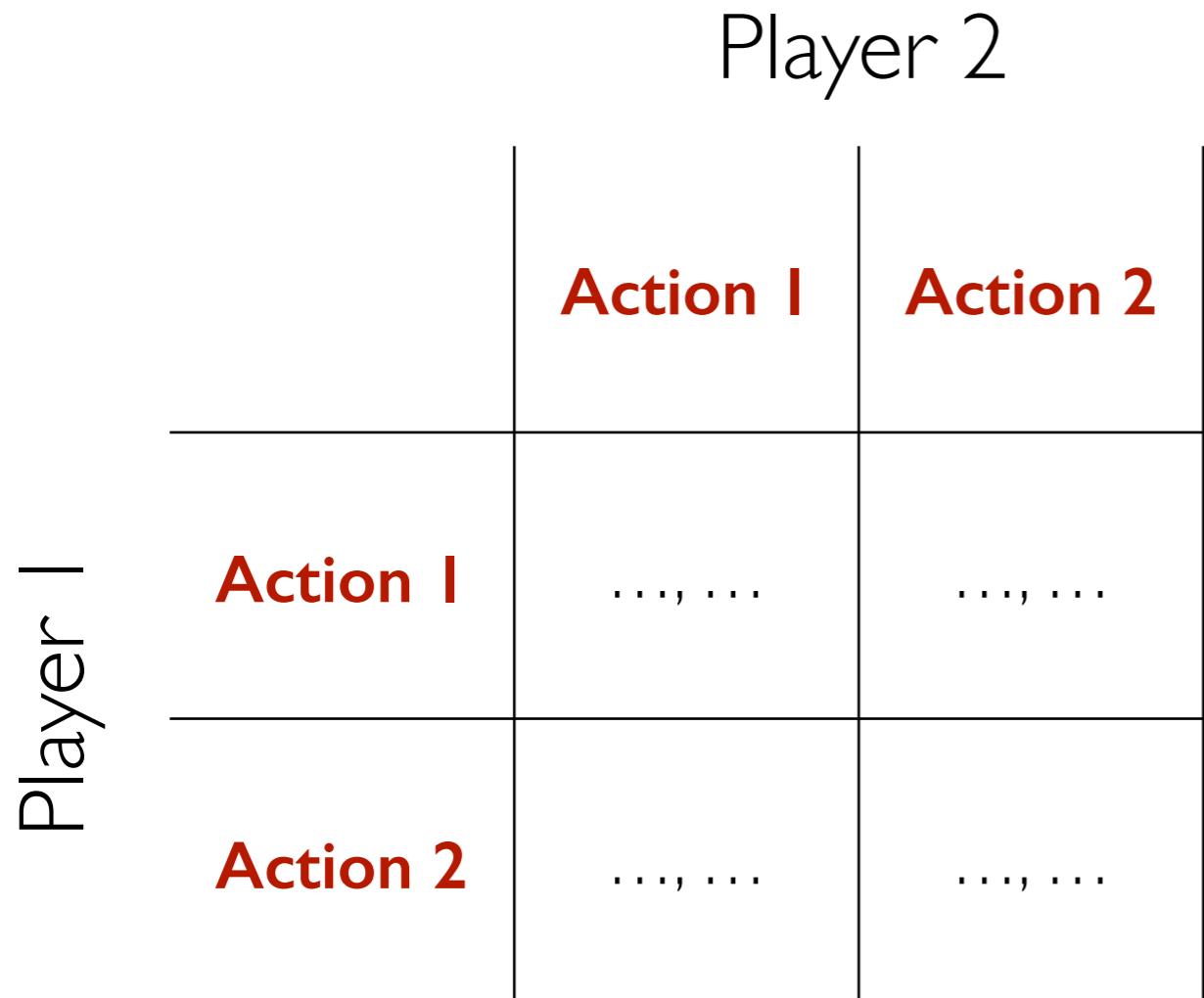


If the starting point is 0, then it will be always zero



Having a starting point with zero strategies does not make sense

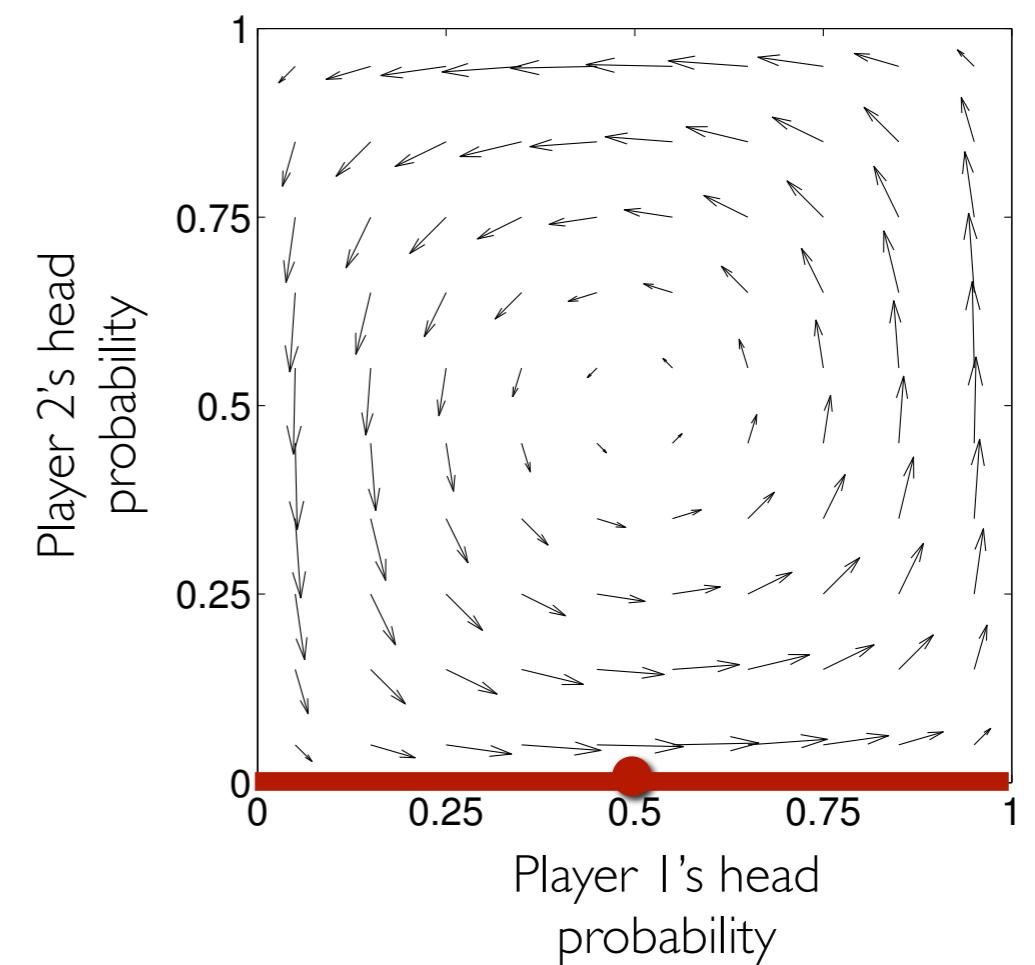
# Starting point



# Starting point

Player 2

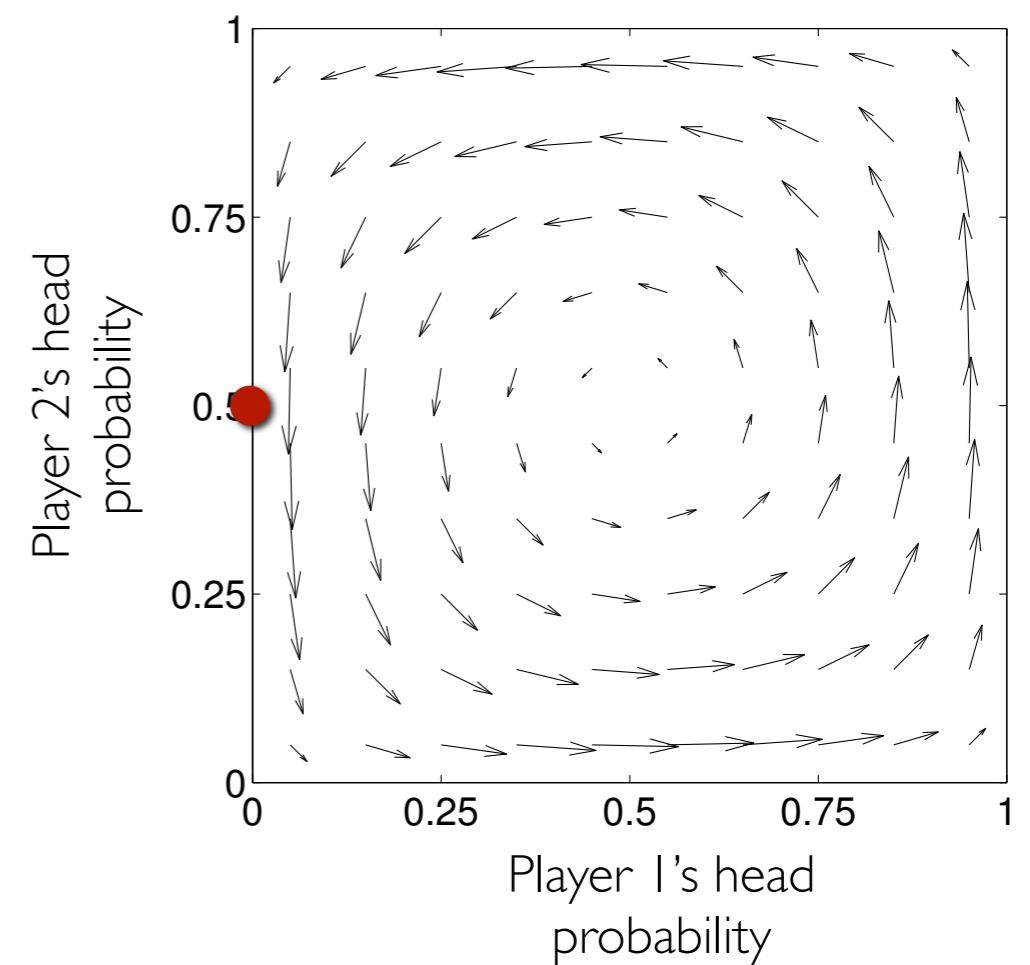
	Action 1	Action 2
Action 1	..., ...	..., ...
Action 2	..., ...	..., ...



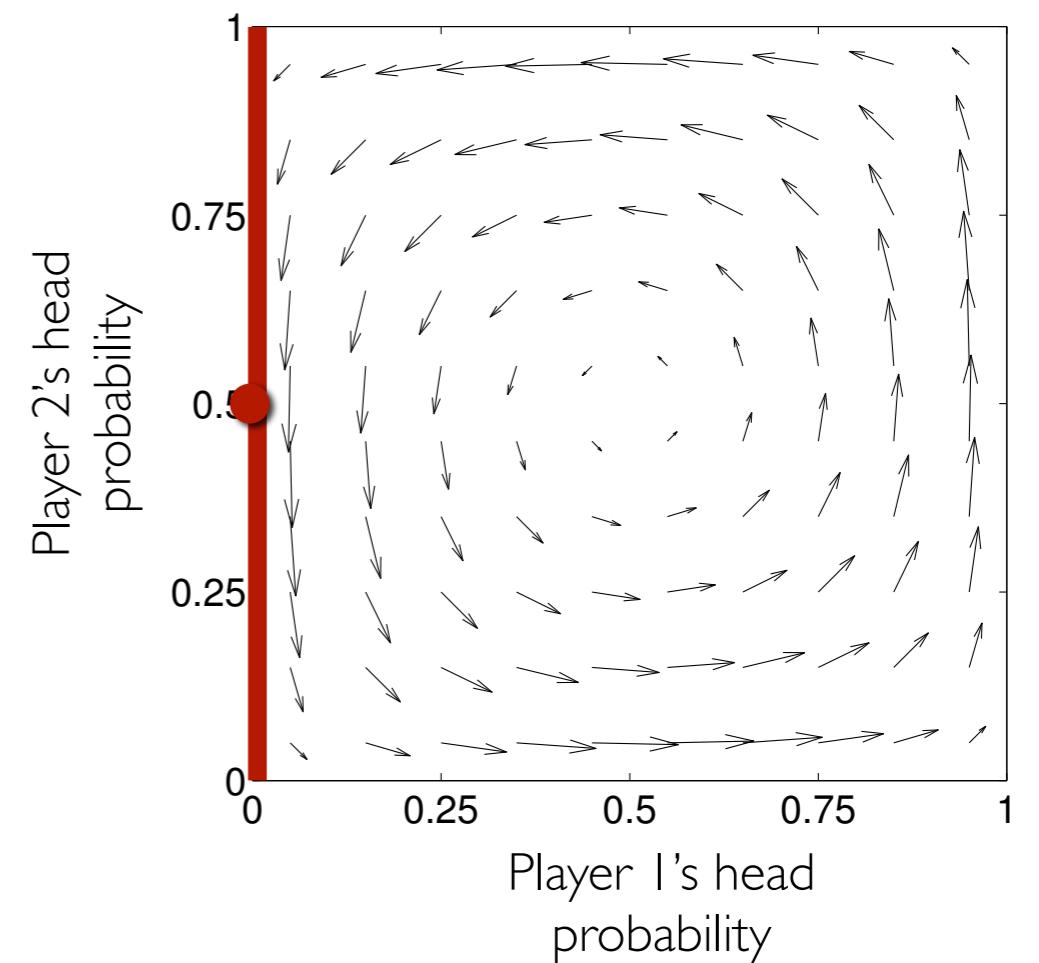
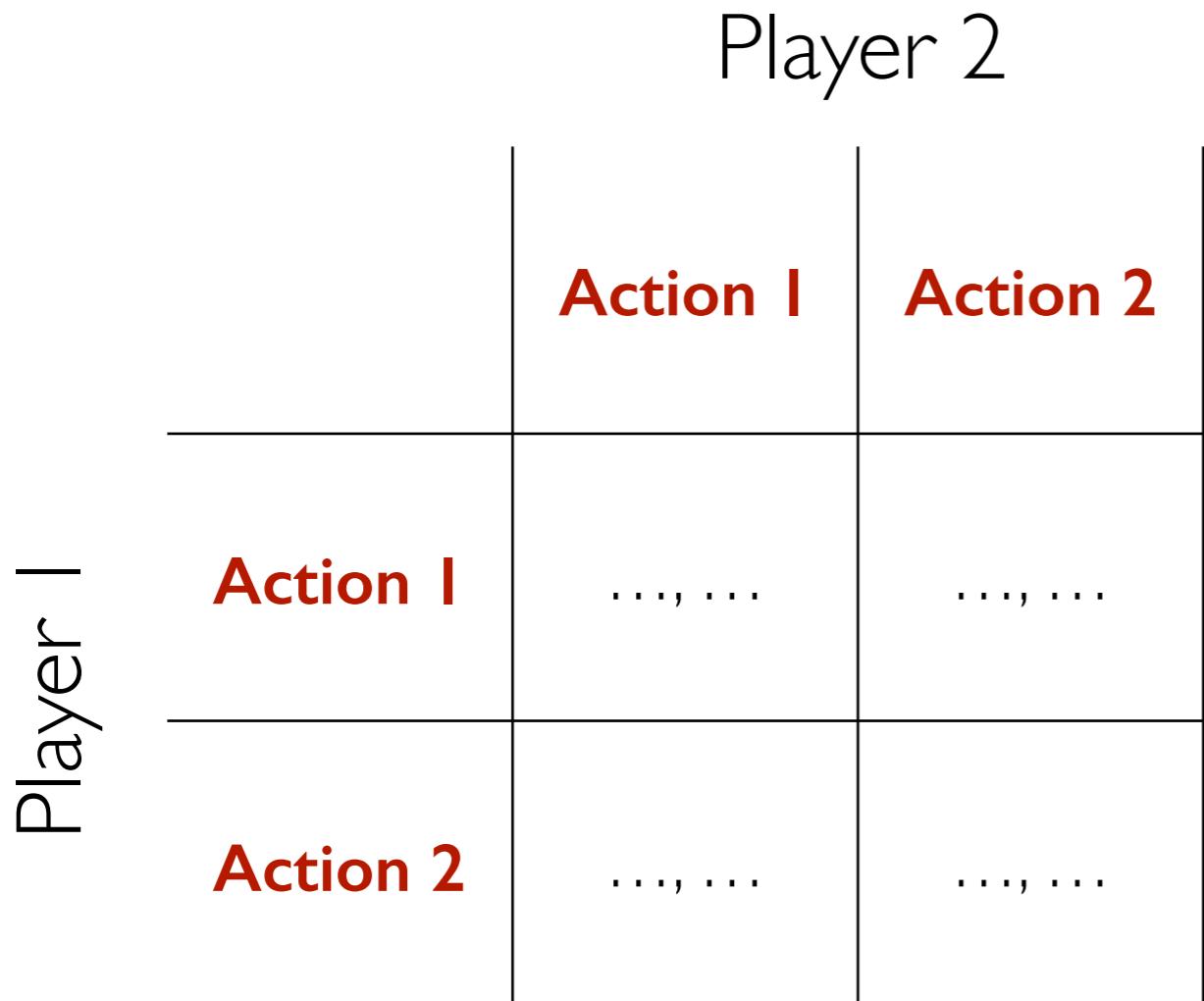
# Starting point

Player 2

	Action 1	Action 2
Action 1	..., ...	..., ...
Action 2	..., ...	..., ...



# Starting point



# Single-population games

- There is a single population of individuals that is playing against itself
- An example: bacterials



# Example

## Single population

Single population

	$\sigma(\text{TYPE 1})$	$\sigma(\text{TYPE 2})$	$\sigma(\text{TYPE 3})$
$\sigma(\text{TYPE 1})$	Type 1	Type 2	Type 3
$\sigma(\text{TYPE 2})$	1	2	0
$\sigma(\text{TYPE 2})$	0	1	2
$\sigma(\text{TYPE 3})$	2	0	1

# Convergence

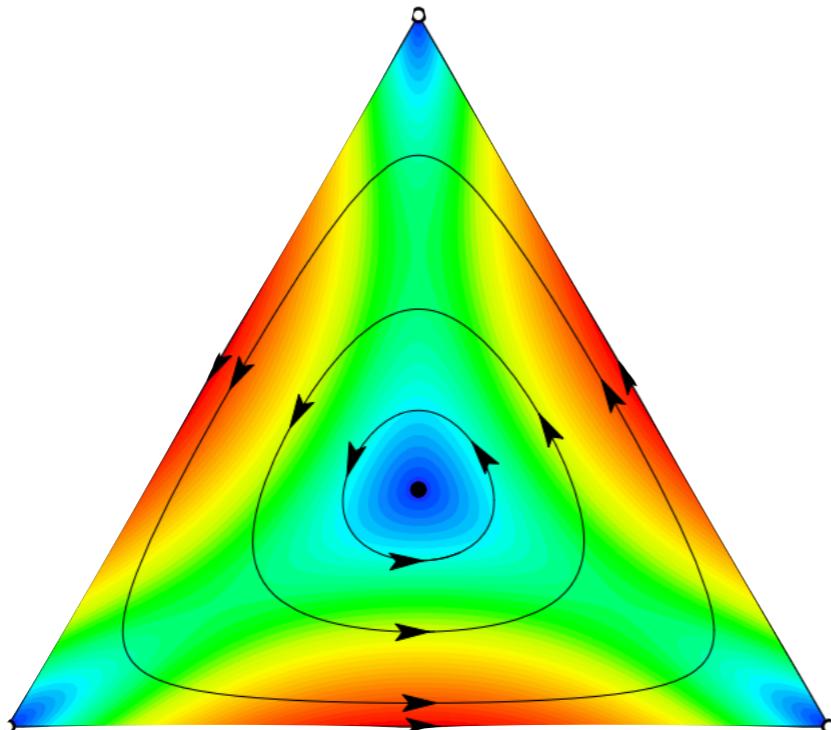
**Theorem.** The replicator equation for a matrix game (single-population) satisfies:

1. A stable rest point is a NE
2. A convergent trajectory in the interior of the strategy space evolves to a NE
3. A strict NE is locally asymptotically stable

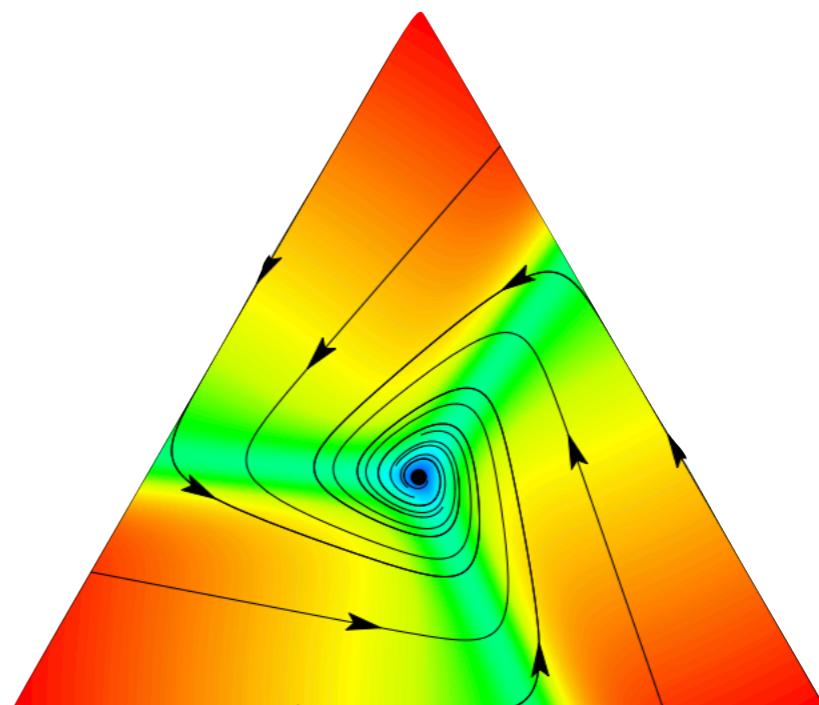
# Multi-population games

- Every player is associated with a different population of individuals
- **Theorem.**  $(p^*, q^*)$  is a two-species ESS if and only if  $(p^*, q^*)$  is a strict Nash in a neighbourhood and
  - is a locally asymptotically stable rest point of the replicator equation
  - if it is in the interior, then it is a globally asymptotically stable rest point of the replicator equation

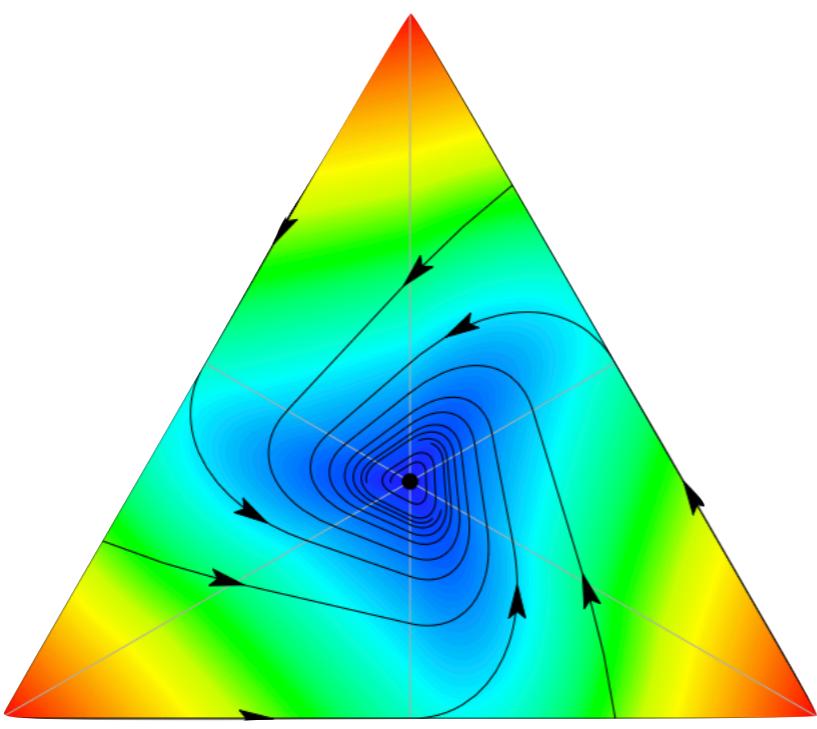
# Known dynamics



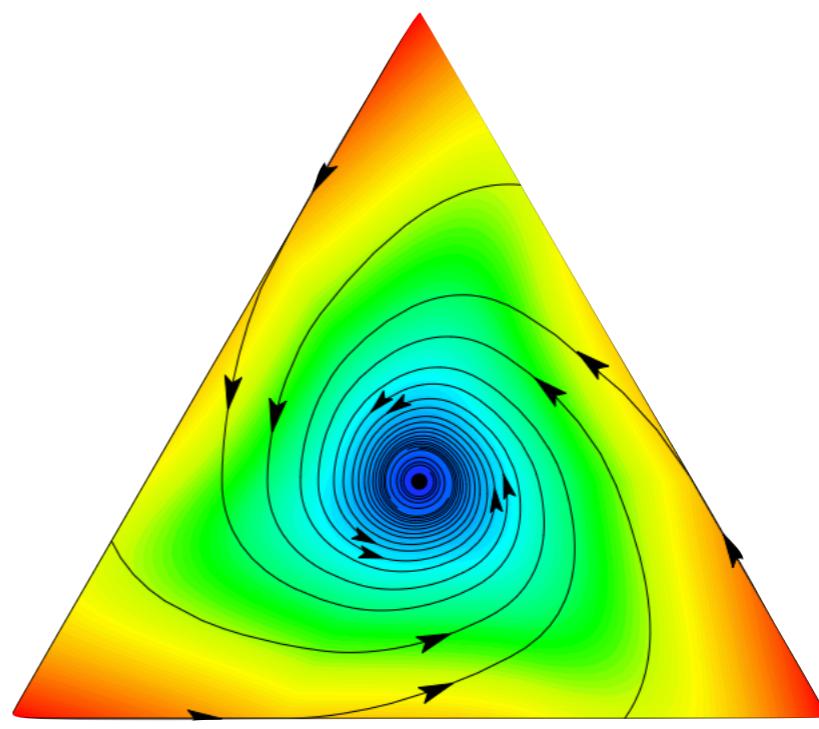
Replicator dynamics



Logit



BNN



Smith

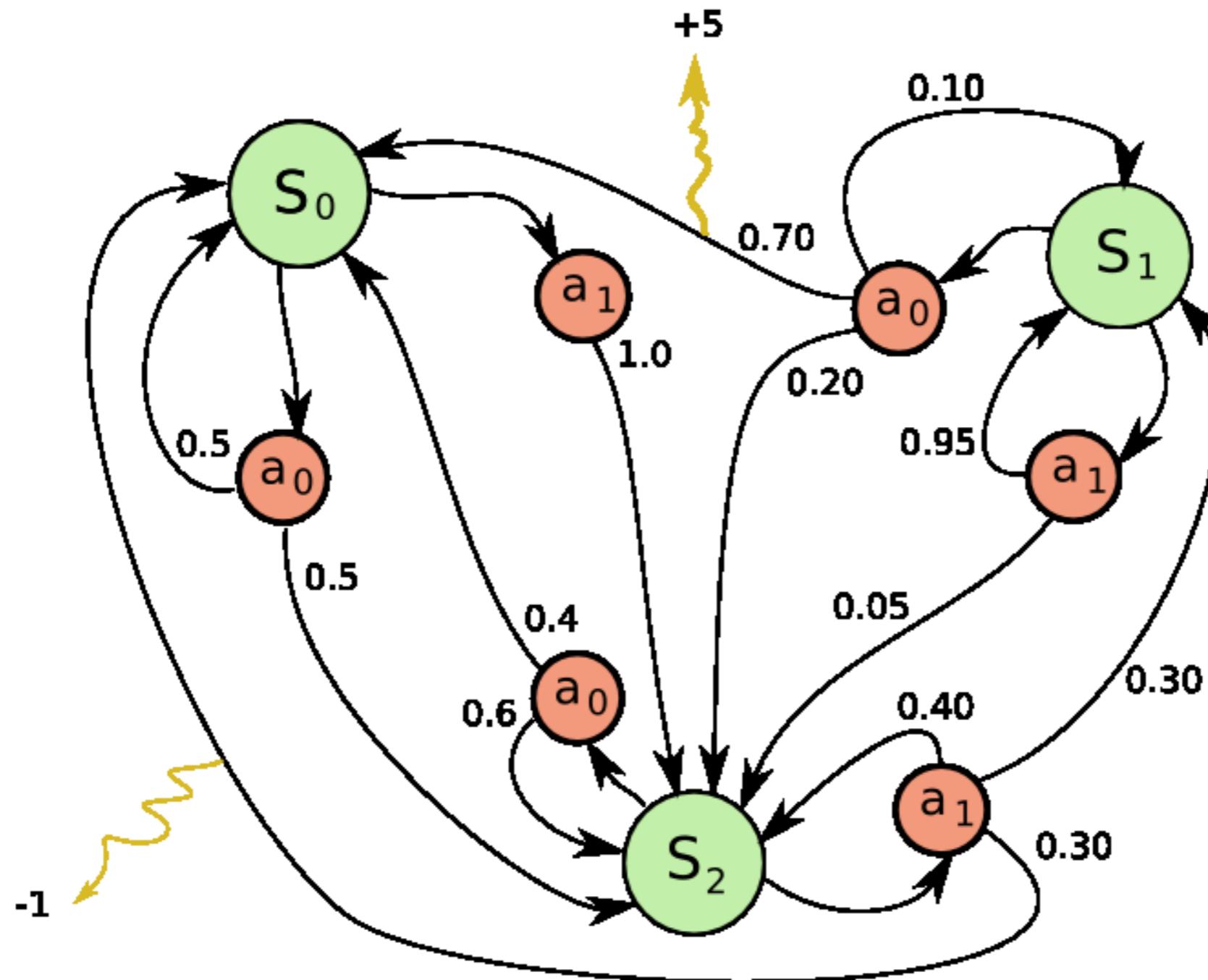
# Properties

Dynamic	Family	(C)	(SD)	(NS)	(PC)
replicator	imitation	yes	yes	no	yes
logit	perturbed best response	yes	yes*	no	no
BNN	excess payoff	yes	no	yes	yes
Smith	pairwise comparison	yes	yes	yes	yes

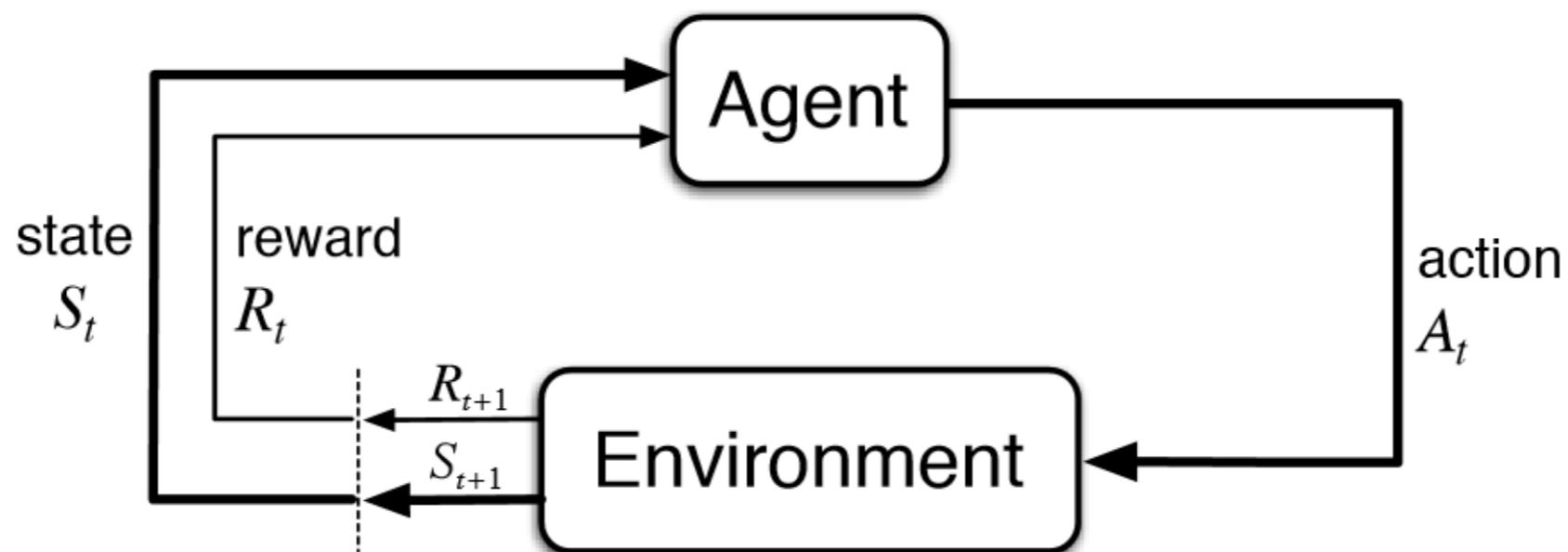
- (C) *Continuity:*  $\rho$  is Lipschitz continuous
- (SD) *Scarcity of data:*  $\rho_{ij}$  only depends on  $\pi_i, \pi_j$ , and  $x_j$
- (NS) *Nash stationarity:*  $V^F(x) = \mathbf{0}$  if and only if  $x \in NE(F)$
- (PC) *Positive correlation:*  $V^F(x) \neq \mathbf{0}$  implies that  $V^F(x)'F(x) > 0$

# Multi-agent learning

# Markov decision problem



# Reinforcement learning



# Q-learning (I)

For every pair state/action:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s, a') - Q(s, a) \right]$$



# Example: normal-form games

Player 1

- 1 state
- 2 actions (Cooperate, Defeat)

		Player 2	
		Cooperate	Defeat
		Cooperate	3,3
		Defeat	5,0
Player 1			1,1

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

		Player 2	
		Cooperate	Defeat
		Cooperate	3,3
		Defeat	5,0
Player 1			,

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

		Player 2	
		Cooperate	Defeat
		Cooperate	Defeat
Player I	Cooperate	3,3	0,5
	Defeat	5,0	1,1

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

		Player 2	
		Cooperate	Defeat
		Cooperate	Defeat
Player I	Cooperate	3,3	0,5
	Defeat	5,0	1,1

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

		Player 2	
		Cooperate	Defeat
		3,3	0,5
Player I	Cooperate	5,0	1,1
	Defeat		

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

		Player 2	
		Cooperate	Defeat
		3,3	0,5
Player I	Cooperate	3,3	0,5
	Defeat	5,0	1,1

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

		Player 2	
		Cooperate	Defeat
		3,3	0,5
Player I	Cooperate	3,3	0,5
	Defeat	5,0	1,1

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

		Player 2	
		Cooperate	Defeat
		Cooperate	Defeat
Player I	Cooperate	3,3	0,5
	Defeat	5,0	1,1

# Example: normal-form games

$$Q(a) \leftarrow Q(a) + \alpha (r - Q(a)) \quad \alpha = 0.2$$

$$\sigma_1(a) = \begin{cases} 1.0 & a = \text{Cooperate} \\ 0.0 & a = \text{Defeat} \end{cases}$$

$$\sigma_1(a) = \begin{cases} 0.2 & a = \text{Cooperate} \\ 0.8 & a = \text{Defeat} \end{cases}$$

round	Player 2's action	Player 1's $Q$ function
$t = 0$	—	$Q(\text{Cooperate}) = 0$
$t = 1$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.6$
$t = 2$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.48$
$t = 3$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.384$
$t = 4$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.3072$
$t = 5$	$a = \text{Defeat}$	$Q(\text{Cooperate}) = 0.24576$
$t = 6$	$a = \text{Cooperate}$	$Q(\text{Cooperate}) = 0.496608$

		Player 2	
		Cooperate	Defeat
		Cooperate	Defeat
Player I	Cooperate	3,3	0,5
	Defeat	5,0	1,1

# Q-learning (2)

Softmax (a.k.a. Boltzam exploration)

$$\sigma_i(a) = \frac{\exp(Q(s, a)/\tau)}{\sum_{a'} \exp(Q(s, a')/\tau)}$$

# Q-learning (2)

Softmax (a.k.a. Boltzam exploration)

$$\sigma_i(a) = \frac{\exp(Q(s, a)/\tau)}{\sum_{a'} \exp(Q(s, a')/\tau)}$$

temperature

# Q-learning (2)

Softmax (a.k.a. Boltzam exploration)

$$\sigma_i(a) = \frac{\exp(Q(s, a)/\tau)}{\sum_{a'} \exp(Q(s, a')/\tau)}$$

temperature

Every action is played with strictly positive probability

The larger the temperature, the smoother the function

If the temperature is 0, we would have a best response

# Example: normal-form games

$Q(\text{Cooperate})$	$Q(\text{Defeat})$	$\sigma_1(\text{Cooperate})$	$\sigma_1(\text{Cooperate})$
0	0	0.5	0.5
1	0	0.731	0.269
5	0	0.99331	0.00669
10	0	0.999955	0.000045

# Example: normal-form games

$Q(\text{Cooperate})$	$Q(\text{Defeat})$	$\sigma_1(\text{Cooperate})$	$\sigma_1(\text{Cooperate})$
0	0	0.5	0.5
1	0	0.731	0.269
5	0	0.99331	0.00669
10	0	0.999955	0.000045

# Example: normal-form games

$Q(\text{Cooperate})$	$Q(\text{Defeat})$	$\sigma_1(\text{Cooperate})$	$\sigma_1(\text{Cooperate})$
0	0	0.5	0.5
1	0	0.731	0.269
5	0	0.99331	0.00669
10	0	0.999955	0.000045

# Example: normal-form games

$Q(\text{Cooperate})$	$Q(\text{Defeat})$	$\sigma_1(\text{Cooperate})$	$\sigma_1(\text{Cooperate})$
0	0	0.5	0.5
1	0	0.731	0.269
5	0	0.99331	0.00669
10	0	0.999955	0.000045

# Self-play Q-learning dynamics

# Self-play learning

Q-learning algorithm

Player 2

		Cooperate	Defeat
		Cooperate	3,3      0,5
Player 1	Cooperate	3,3	0,5
	Defeat	5,0	1,1

Q-learning algorithm

# Learning dynamics

Assumptions:

- Time is continuous
- All the actions can be selected simultaneously

# Learning dynamics

Assumptions:

- Time is continuous
- All the actions can be selected simultaneously

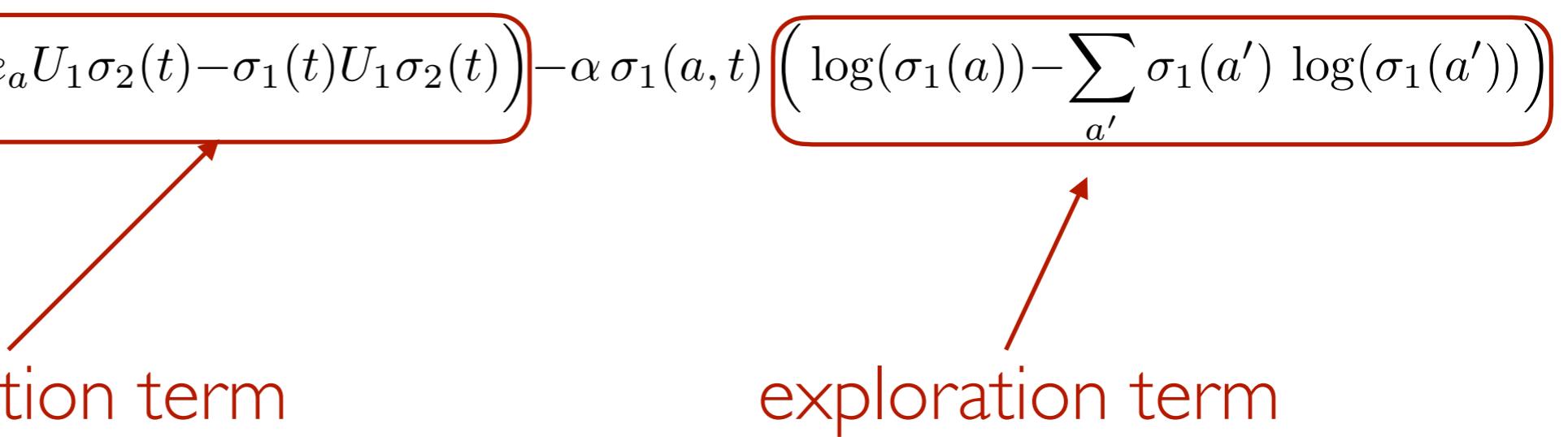
$$\dot{\sigma}_1(a, t) = \frac{\alpha \sigma_1(a, t)}{\tau} \left( e_a U_1 \sigma_2(t) - \sigma_1(t) U_1 \sigma_2(t) \right) - \alpha \sigma_1(a, t) \left( \log(\sigma_1(a)) - \sum_{a'} \sigma_1(a') \log(\sigma_1(a')) \right)$$

# Learning dynamics

Assumptions:

- Time is continuous
- All the actions can be selected simultaneously

$$\dot{\sigma}_1(a, t) = \frac{\alpha \sigma_1(a, t)}{\tau} \left( e_a U_1 \sigma_2(t) - \sigma_1(t) U_1 \sigma_2(t) \right) - \alpha \sigma_1(a, t) \left( \log(\sigma_1(a)) - \sum_{a'} \sigma_1(a') \log(\sigma_1(a')) \right)$$



exploitation term

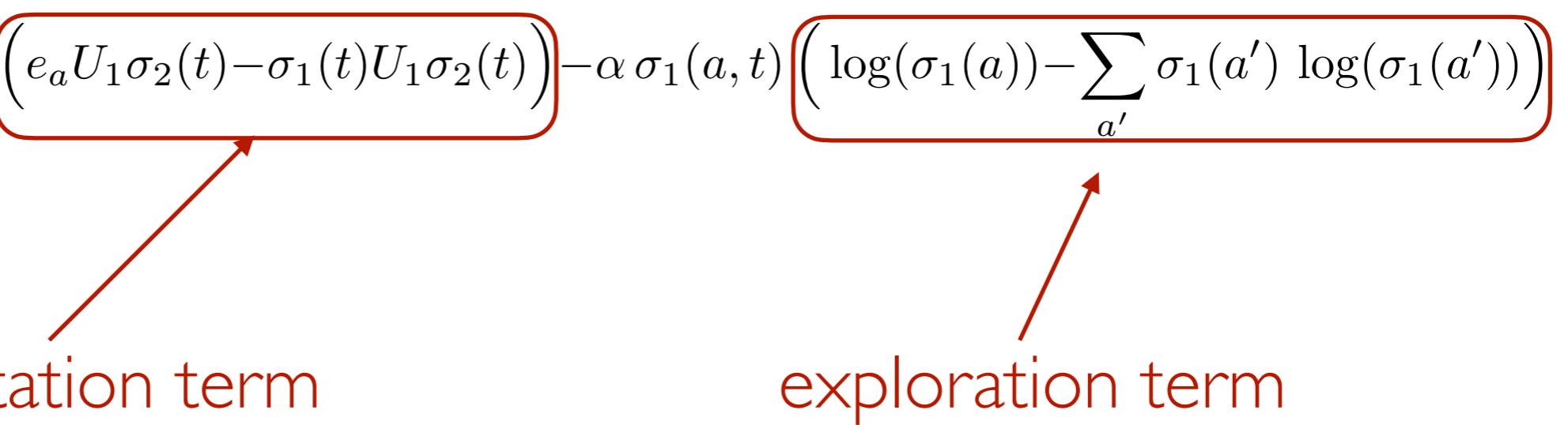
exploration term

# Learning dynamics

Assumptions:

- Time is continuous
- All the actions can be selected simultaneously

$$\dot{\sigma}_1(a, t) = \frac{\alpha \sigma_1(a, t)}{\tau} \left( e_a U_1 \sigma_2(t) - \sigma_1(t) U_1 \sigma_2(t) \right) - \alpha \sigma_1(a, t) \left( \log(\sigma_1(a)) - \sum_{a'} \sigma_1(a') \log(\sigma_1(a')) \right)$$

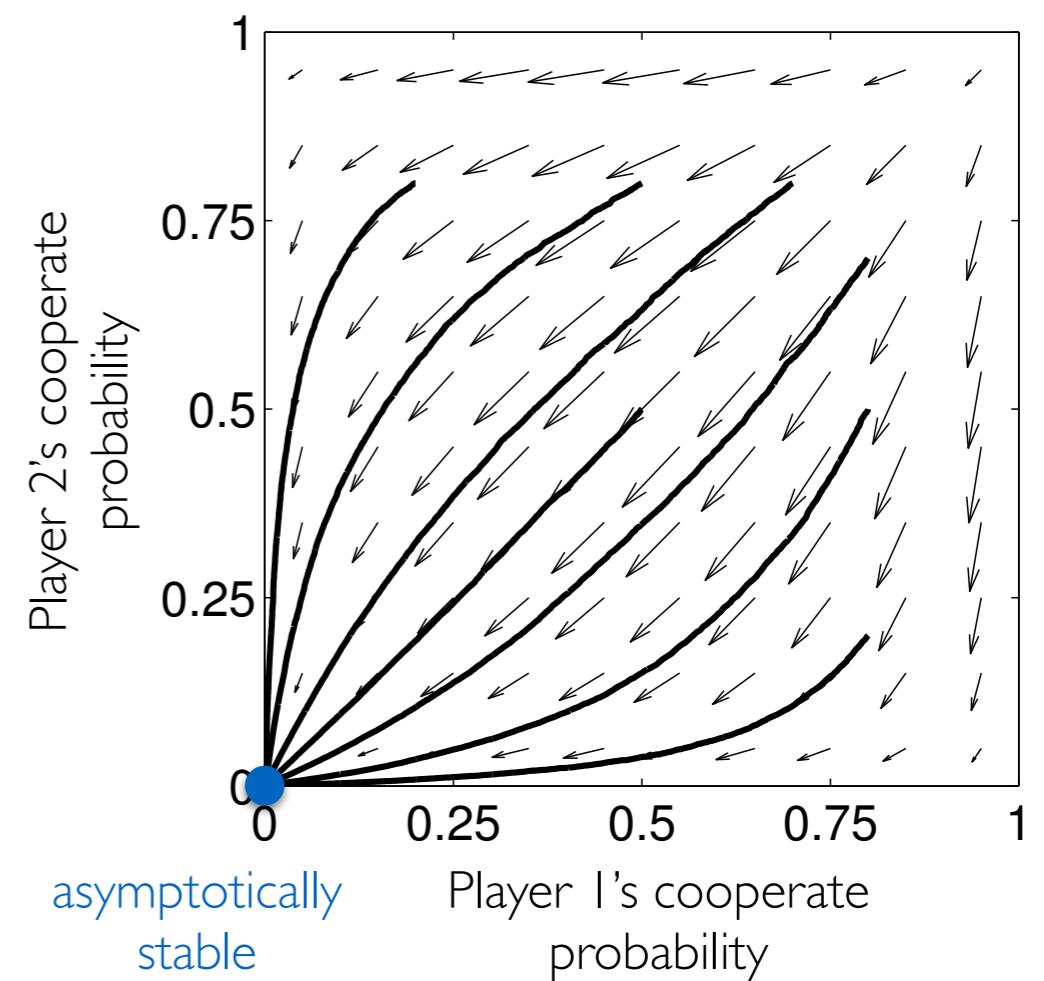
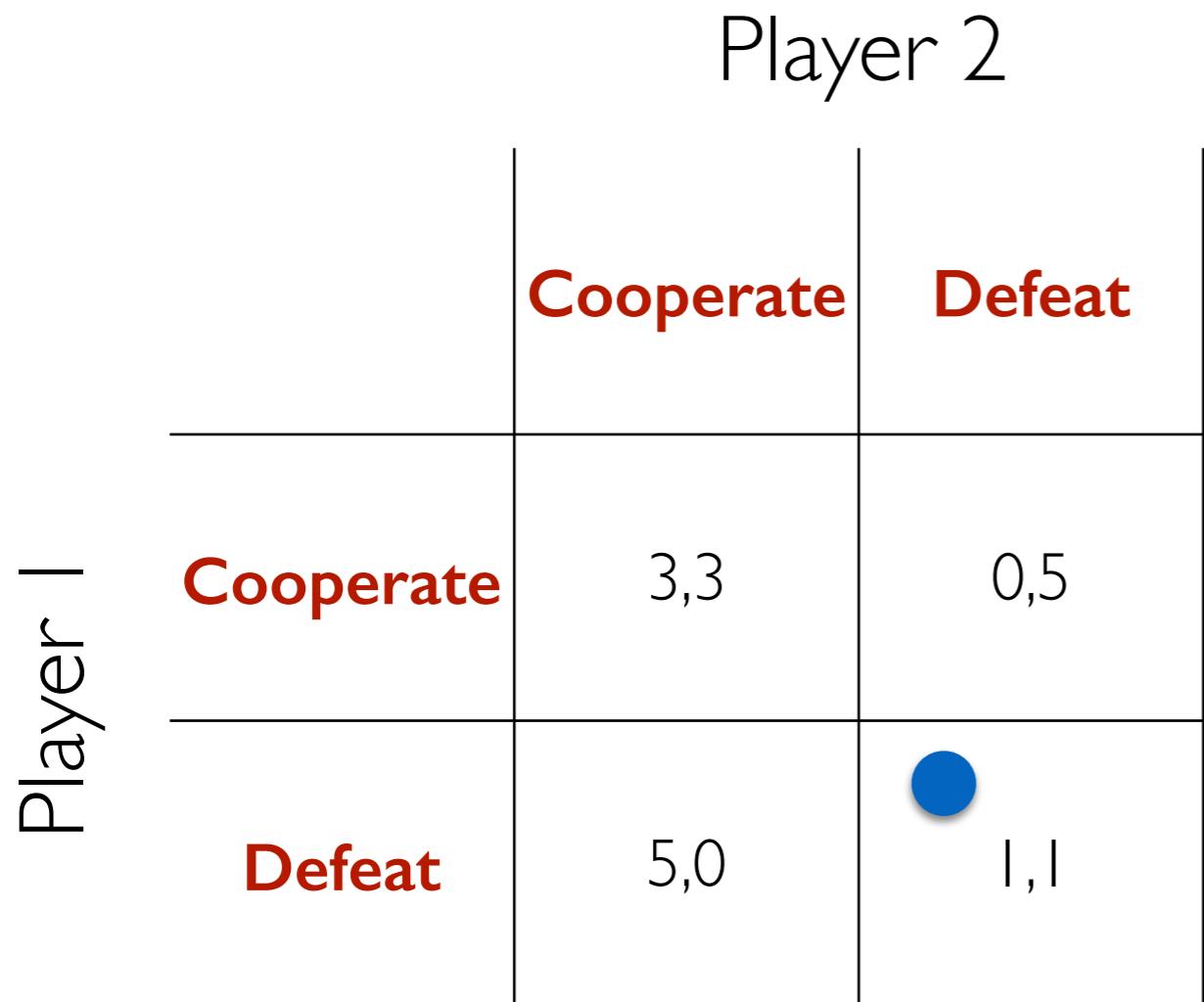


exploitation term

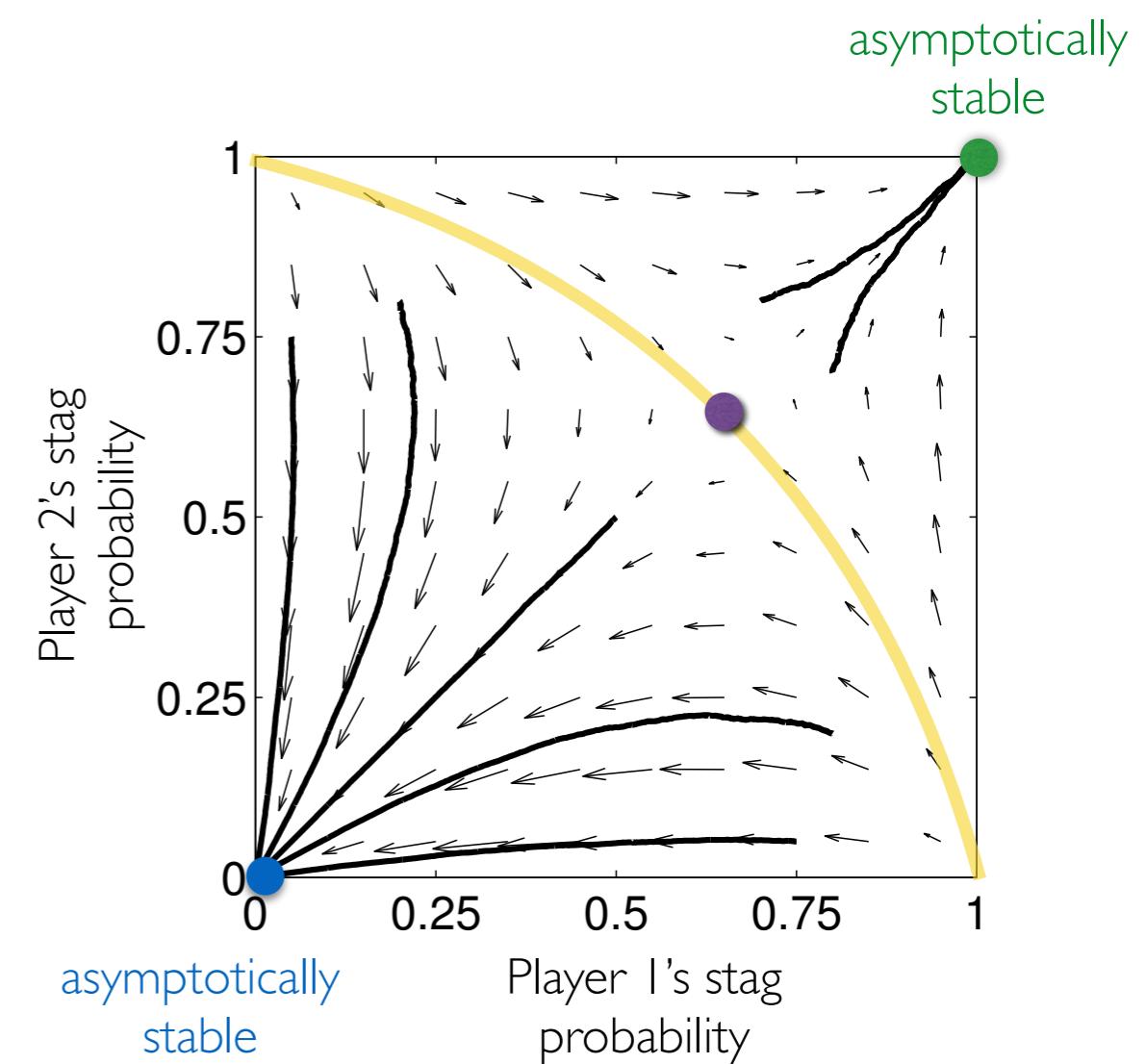
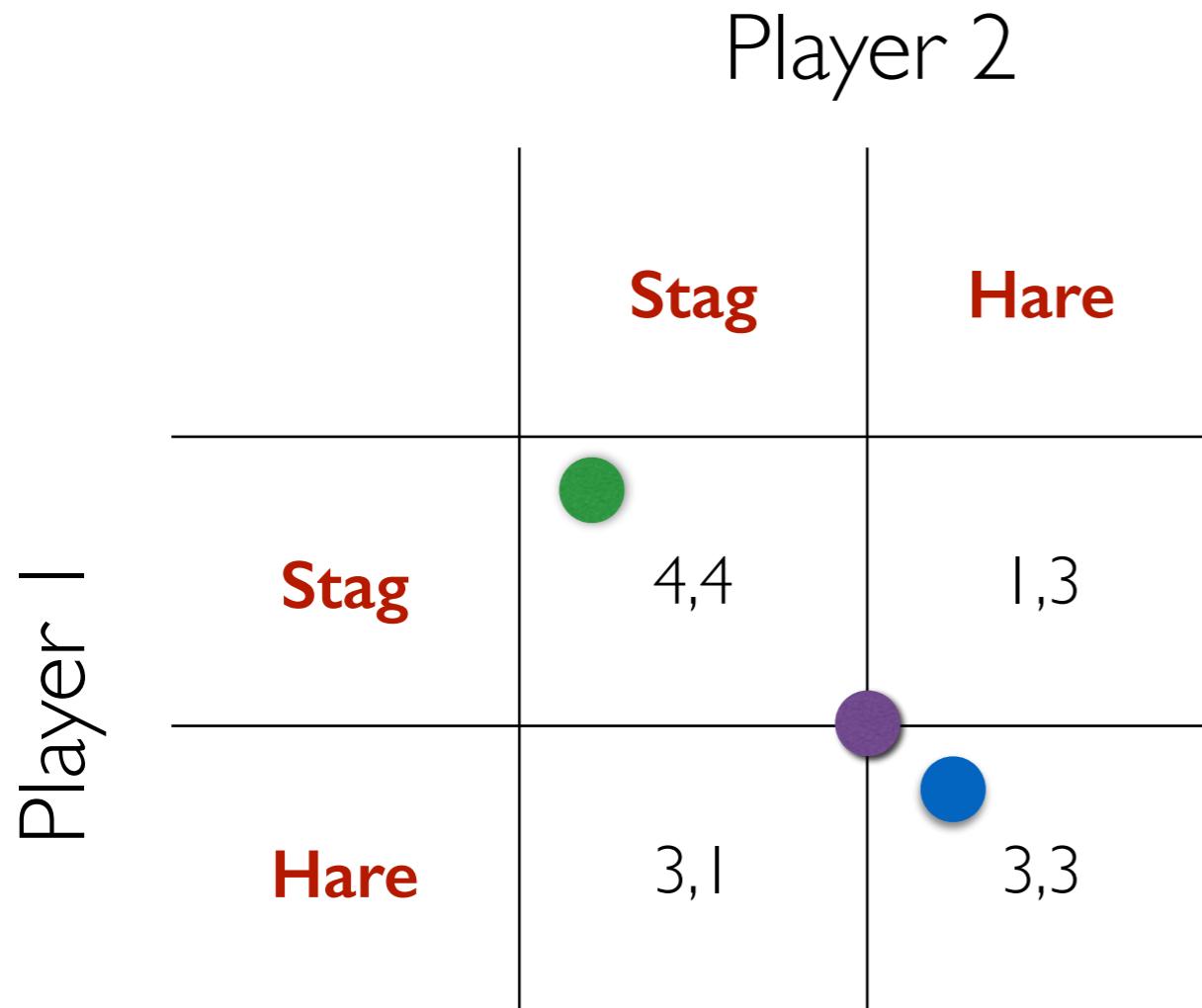
exploration term

When the temperature is 0, the Q-learning behaves as the replicator dynamics

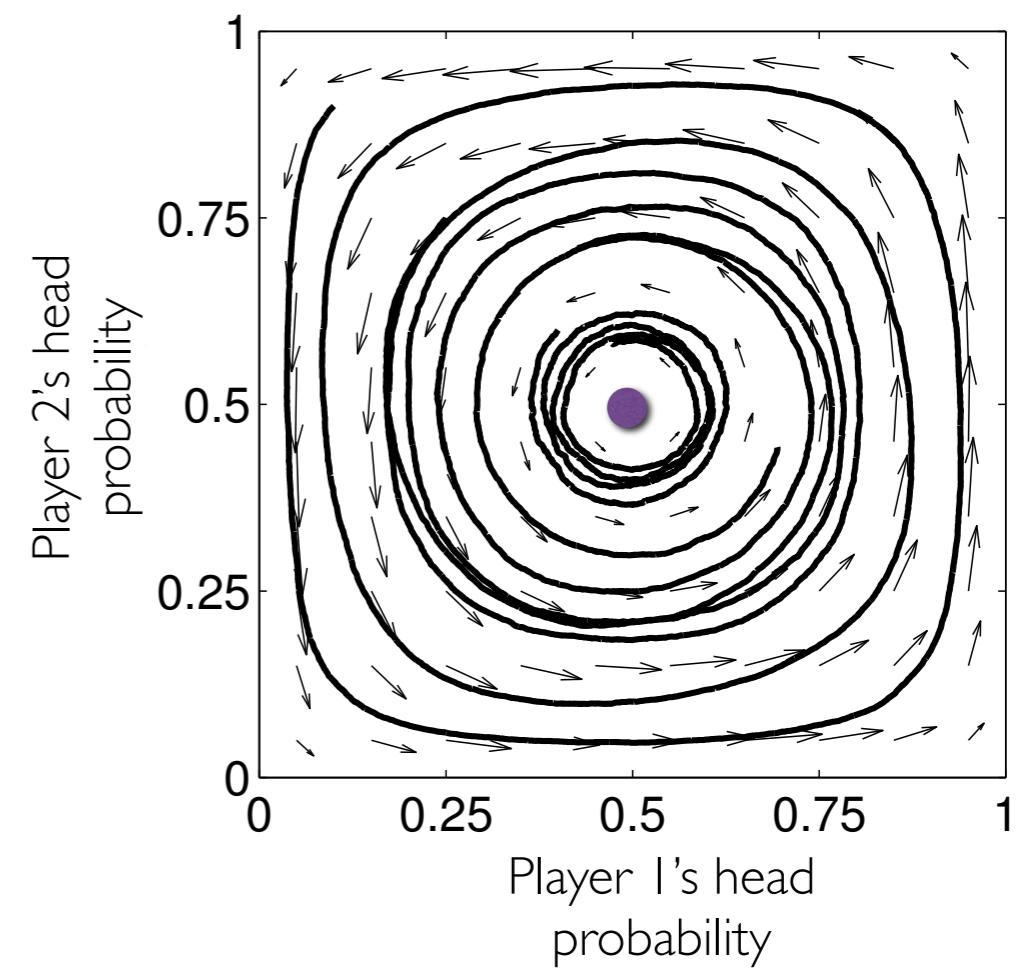
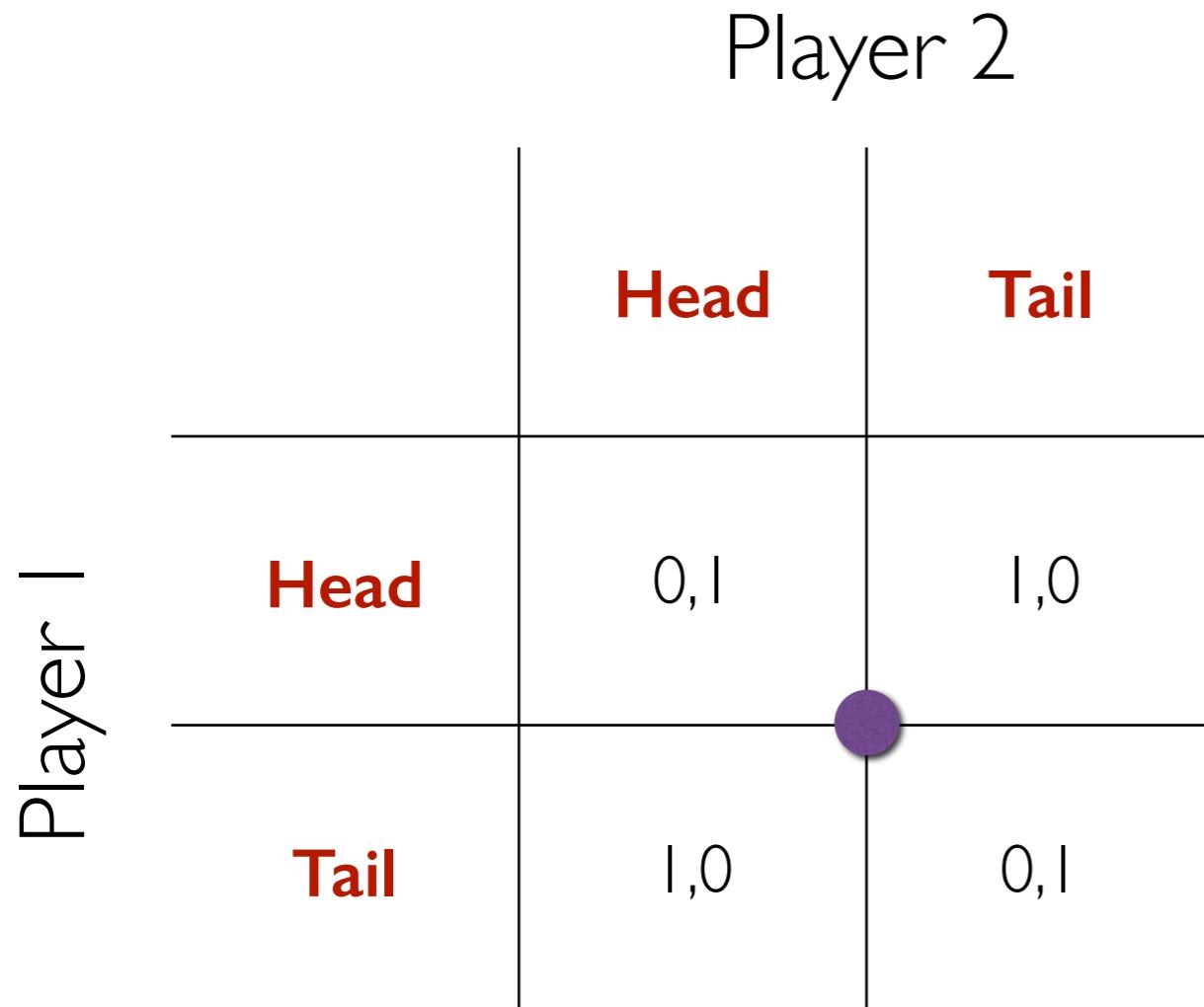
# Prisoner's dilemma



# Stag hunt



# Matching pennies



# Other MAL algorithms dynamics based on replicator equation

	Algorithm	$\dot{x}$
Infinitesimal Gradient Ascent	IGA	$\alpha \vec{\delta}$
IGA-Win or Learn Fast	IGA-WoLF	$\vec{\delta} \cdot \begin{cases} \alpha_{min} & \text{if } V(\mathbf{x}) > V(\mathbf{x}^*) \\ \alpha_{max} & \text{otherwise} \end{cases}$
Weighted Policy Learning	WPL	$\alpha \vec{\delta} \cdot \begin{cases} x & \text{if } \vec{\delta} < 0 \\ (1 - x) & \text{otherwise} \end{cases}$
Cross Learning	CL	$x(1 - x) \vec{\delta}$
Frequency-Adjusted Q-learning	FAQ	$\alpha x(1 - x) [\vec{\delta} \cdot \tau^{-1} - \log \frac{x}{1-x}]$
Regret Minimization (Polynomial Weights)	RM	$\alpha x(1 - x) \vec{\delta} \cdot \begin{cases} (1 + \alpha x \vec{\delta})^{-1} & \text{if } \vec{\delta} < 0 \\ (1 - \alpha(1 - x) \vec{\delta})^{-1} & \text{otherwise} \end{cases}$