

---

# ***Game Theoretic Learning for Distributed Autonomous Systems***

**Jeff S Shamma**  
Georgia Institute of Technology

LCCC Workshop  
May 28–29, 2009

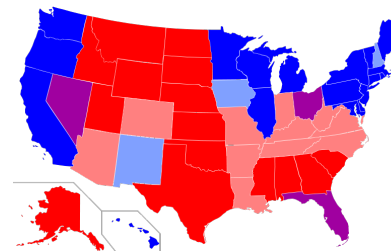
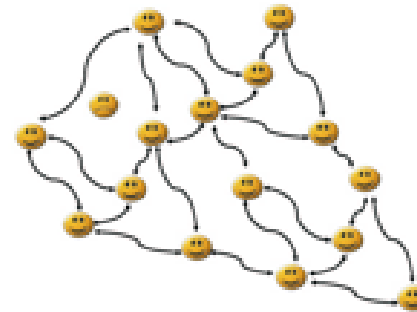


# ***Multiagent, game theoretic, cooperative, distributed, networked...***

---

- Special issues:
  - SIAM Journal on Control and Optimization, Special Issue on “Control and Optimization in Cooperative Networks”
  - ASME Journal of Dynamics Systems, Measurement, and Control, Special Issue on “Analysis and Control of Multi-Agent Dynamic Systems”
  - Journal of Intelligent and Robotic Systems, Special Issue on “Cooperative Robots, Machines, and Systems”
  - International Journal on Systems, Control and Communications, Special Issue on “Networked Control Systems”
  - Robotica, Special Issue on “Robotic Self-X Systems”
- Workshops:
  - MIT Workshop on Frontiers in Game Theory and Networked Control Systems
  - NecSys: IFAC Workshop on Estimation and Control of Networked Systems
  - GameNets: International Conference on Game Theory for Networks
  - GameComm: International Workshop on Game Theory in Communication Networks
  - 8th International Conference on Cooperative Control and Optimization (2008)

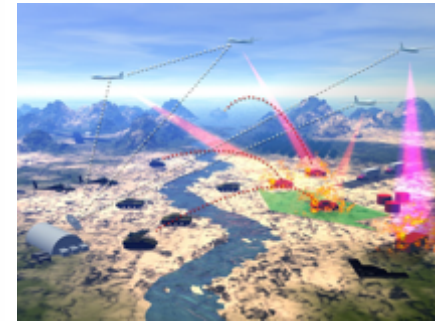
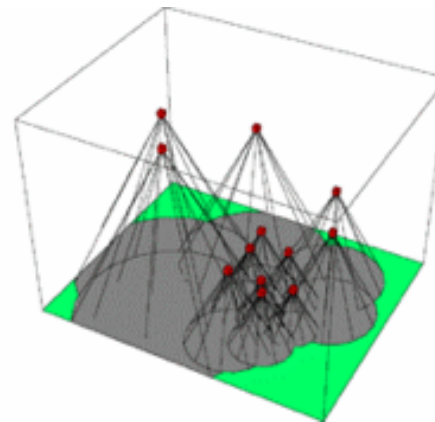
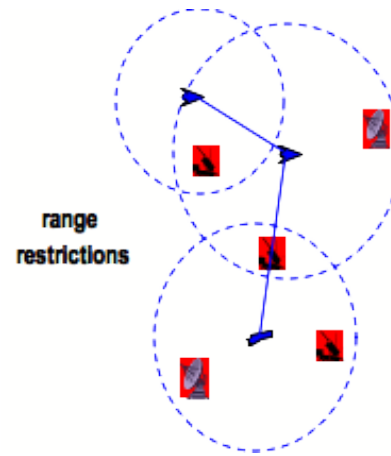
- Traffic
- Evolution of convention
- Social network formation
- Auctions & markets
- Voting
- etc
- Game elements (inherited):
  - Actors/players
  - Choices
  - Preferences



*Descriptive Agenda*

## More multiagent scenarios

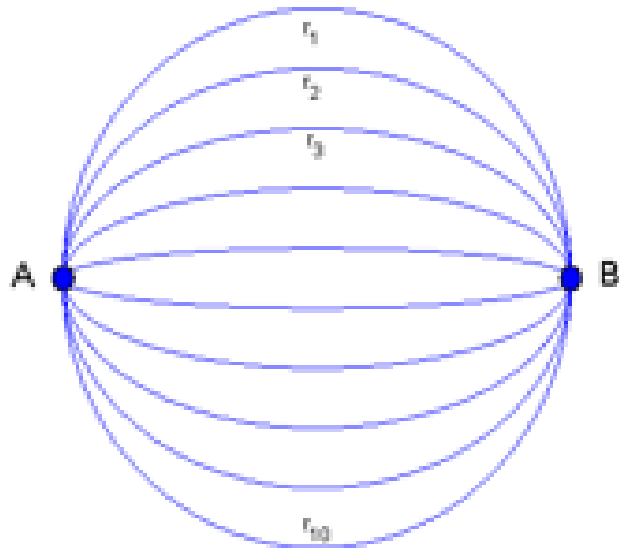
- Weapon-target assignment
- Data network routing
- Mobile sensor coverage
- Autonomous vehicle teams
- etc
- Game elements (designed):
  - Actors/players
  - Choices
  - Preferences




*Prescriptive Agenda*

- Prescriptive agenda = distributed robust optimization
- *Choose* to address cooperation as noncooperative game
- Players are *programmable components* (vs humans)
- Must *specify*
  - Elements of game (players, actions, payoffs)
  - Learning algorithm
- Metrics:
  - Information available to agent?
  - Communications/stage?
  - Processing/stage?
  - Asymptotic behavior?
  - Global objective performance?
  - Convergence rates?

- Game theoretic learning
- Special class: Potential games
- Survey of algorithms
- Illustrations



Distributed routing

							5	
4			1	8				
		7	6		3	9		
		6	9		8	3	2	
	5						7	
	8	3	4		7	5		
		5	3		6	1		
				1	2			6
	3							

Multi-agent sudoku

# Game setup & Nash equilibrium

---

- Setup:

- Players:  $\{1, \dots, p\}$

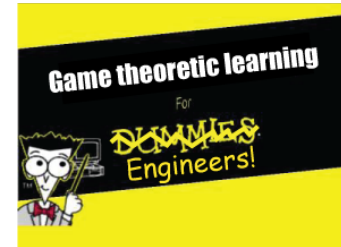
- Actions:  $a_i \in \mathcal{A}_i$

- Action profiles:

$$(a_1, a_2, \dots, a_p) \in \mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_p$$

- Payoffs:  $u_i : (a_1, a_2, \dots, a_p) = (a_i, a_{-i}) \mapsto \mathbf{R}$

- Global objective:  $G : \mathcal{A} \rightarrow \mathbf{R}$



- Action profile  $a^* \in \mathcal{A}$  is a **Nash equilibrium** (NE) if for all players:

$$u_i(a_1^*, a_2^*, \dots, a_p^*) = u_i(a_i^*, a_{-i}^*) \geq u_i(a_i', a_{-i}^*)$$

i.e., no *unilateral* incentive to change actions.



- Iterations:
  - $t = 0, 1, 2, \dots$
  - $a_i(t) = \text{rand}(s_i(t)), \quad s_i(t) \in \Delta(\mathcal{A}_i)$
  - $s_i(t) = \mathcal{F}_i(\text{available info at time } t)$
- Key questions: If NE is a descriptive outcome...
  - How could agents converge to NE?
  - Which NE?
  - Are NE efficient?

- Iterations:
  - $t = 0, 1, 2, \dots$
  - $a_i(t) = \text{rand}(s_i(t)), \quad s_i(t) \in \Delta(\mathcal{A}_i)$
  - $s_i(t) = \mathcal{F}_i(\text{available info at time } t)$
- Key questions: If NE is a descriptive outcome...
  - How could agents converge to NE?
  - Which NE?
  - Are NE efficient?
- Focus shifted away from NE towards adaptation/learning

*“The attainment of equilibrium requires a disequilibrium process”*

K. Arrow

*“Game theory lacks a general and convincing argument that a Nash outcome will occur.”*

Fudenberg & Tirole

*“...human subjects are no great shakes at thinking either [vs insects]. When they find their way to an equilibrium of a game, they typically do so using trial-and-error methods.”*

K. Binmore

Survey: Hart, “Adaptive heuristics”, 2005.

## *Game theoretic learning for prescriptive agenda?*

---

- Approach: Use game theoretic learning to steer collection towards desirable configuration
- Informational hierarchy:
  - Action based: Players can observe the actions of others.
  - Oracle based: Players receive an aggregate report of the actions of others.
  - Payoff based: Players only measure online payoffs.
- Focus:
  - Asymptotic behavior
  - Processing per stage
  - Communications per stage

- For some  $\phi : \mathcal{A} \rightarrow \mathbb{R}$

$$\begin{aligned}\phi(a_i, a_{-i}) - \phi(a'_i, a_{-i}) &> 0 \\ \Leftrightarrow \\ u_i(a_i, a_{-i}) - u_i(a'_i, a_{-i}) &> 0\end{aligned}$$

i.e., potential function increases iff unilateral improvement.

- Features:
  - Typical of “coordination games”
  - Desirable convergence properties under various algorithms
  - Need not imply “cooperation” or  $\phi = G$

- Distributed routing

- Payoff = negative congestion.  $c_r(\sigma_r)$
- Potential function:

$$\phi = \sum_r \sum_{n=1}^{\sigma_r} c_r(n)$$

- Overall congestion:

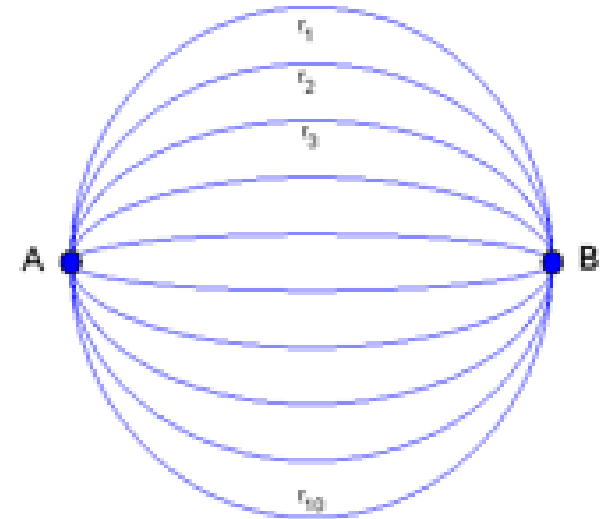
$$G = \sum_r \sigma_r c_r(\sigma_r)$$

- **Note:**  $\phi \neq G$

- Multiagent sudoku:

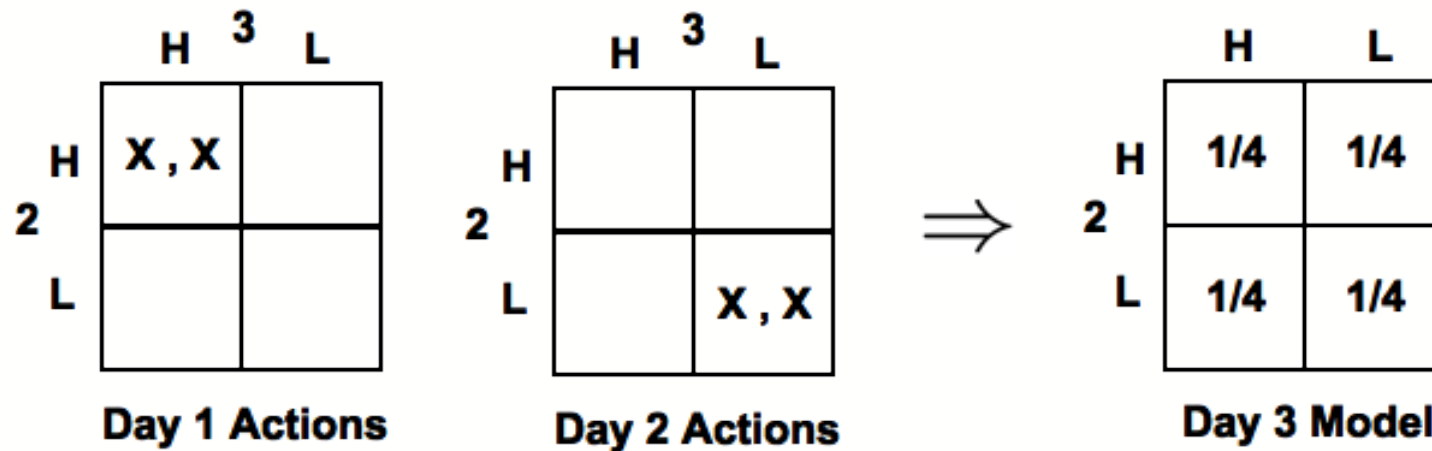
$u_i(a) = \# \text{reps in row} + \# \text{reps in column} + \# \text{reps in sector}$

$$\phi(a) = \sum_i u_i(a)$$



😡							5	
4			1	8		😡		
		7	6		3	9		😞
	😎	6	9		8	3	2	
	5			👤			7	
	8	3	4		7	5		🎉
		5	3		6	1		
		😏		1	2			6
	3							👉 😊

- Each player:
  - Maintains empirical frequencies (histograms) of other player actions
  - Forecasts (incorrectly) that others are playing randomly and independently according to empirical frequencies
  - Selects an action that maximizes expected payoff
- Bookkeeping is *action based*
- **Monderer & Shapley (1996)**: FP converges to NE in potential games.



- Viewpoint of driver 1 (3 drivers & 2 roads)
- Prohibitive-per-stage for large numbers of players with large action sets
  - Monitor all other players with IDs (cf., distributed routing)
  - Take expectation over large joint action space (cf., sudoku)

## Joint strategy fictitious play (JSFP)

---

- Virtual payoff vector

$$U_i(t) = \begin{pmatrix} u_i(1, a_{-i}(t)) \\ u_i(2, a_{-i}(t)) \\ \vdots \\ u_i(m, a_{-i}(t)) \end{pmatrix}$$

i.e., the payoffs that *could have* been obtained at time  $t$

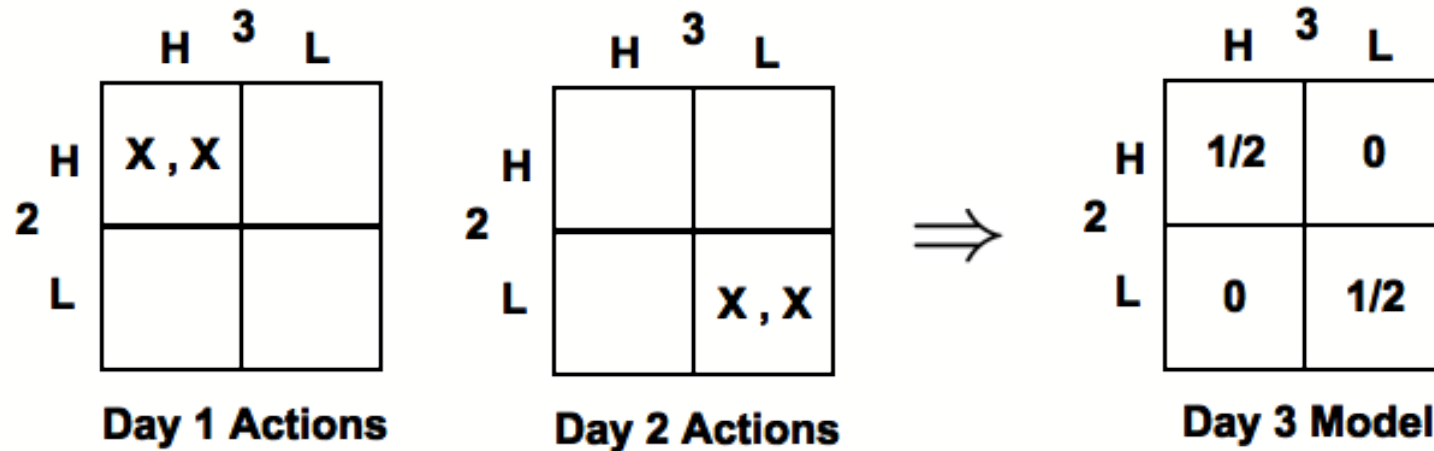
- Time averaged virtual payoff:

$$V_i(t+1) = (1 - \rho)V_i(t) + \rho U_i(t)$$

- Stepsize  $\rho$  is either
  - Constant (fading memory)
  - Diminishing (true average), e.g.,  $\rho = \frac{1}{t+1}$
- Bookkeeping is *oracle based* (cf., traffic reports)



- JSFP algorithm: Each player
  - Maintains time averaged virtual payoff
  - Selects an action with maximal virtual payoff
  - OR repeats previous stage action with some probability (i.e., inertia)
- **Marden, Arslan, & JSS (2005)**: JSFP with inertia converges to a NE in potential games.



- Equivalent to best response to *joint actions* of other players
- Related to “no regret” algorithms
- Survey: Foster & Vohra, Regret in the online decision problem, 1999.

## Equilibrium selection & Gibbs distribution

---

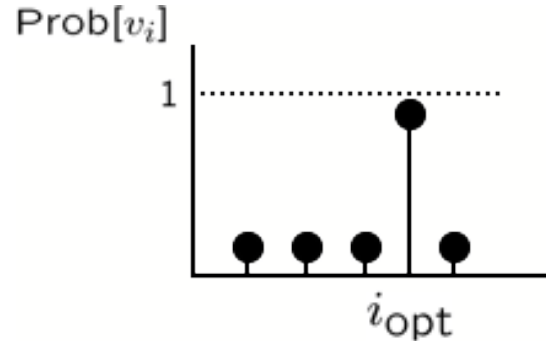
- Alternative algorithms offer more quantitative characterization of asymptotic behaviors.
- Preliminary: Gibbs distribution (cf., softmax, logit response)

$$\sigma(v; T) = \frac{1}{\mathbf{1}^T e^{v/T}} e^{v/T} \in \Delta$$

e.g.,

$$\sigma(v_1, v_2; T) = \begin{pmatrix} \frac{e^{v_1/T}}{e^{v_1/T} + e^{v_2/T}} \\ \frac{e^{v_2/T}}{e^{v_1/T} + e^{v_2/T}} \end{pmatrix}$$

- As  $T \downarrow 0$  assigns all probability to  $\arg \max \{v_1, v_2, \dots, v_n\}$



- At stage  $t$ 
  - Player  $i$  is selected at random
  - Chosen player sets

$$a_i(t) = \text{rand} \left[ \sigma \left( u_i(1, a_{-i}(t-1)), \dots, u_i(m, a_{-i}(t-1)); T \right) \right]$$

- Interpretation: Noisy best reply to previous joint actions
- Fact: SAP results in a Markov chain over joint action space  $\mathcal{A}$  with a unique stationary distribution,  $\mu$ .
- **Blume (1993)**: In (cardinal) potential games,

$$\mu(a) = \sigma(\phi(a); T) = \frac{e^{\phi(a)/T}}{\sum_{a' \in \mathcal{A}} e^{\phi(a')/T}}$$

- Implication: As  $T \downarrow 0$ , all probability assigned to potential maximizer.

- Motivation:
  - Reduced processing per stage
  - First step towards constrained actions

- At stage  $t$ :

- Player  $i$  is selected at random
- Chosen player compares  $a_i(t - 1)$  with randomly selected  $a'_i$

$$a_i(t) = \text{rand} [\sigma(u_i(a_i(t - 1), a_{-i}(t - 1)), u_i(a'_i, a_{-i}(t - 1); T))]$$

- **Arslan, Marden, & JSS (2007)**: Binary SAP results in same stationary distribution as SAP.
- Consequence: Arbitrarily high steady state probability on potential function maximizer.

- Action evolution must satisfy:  $a_i(t) \in \mathcal{C}(a_i(t-1))$

- Limited mobility
- Obstacles

- Algorithm: Same as before *except*

$$a'_i \in \mathcal{C}(a_i(t-1))$$

- **Marden & JSS (2008)**: Constrained SAP results in potential function maximizer being *stochastically stable*.
  - Arbitrarily high steady state probability on potential function maximizer
  - Does *not* characterize steady state distribution

- Action & oracle based algorithms require:
  - Explicit communications
  - Explicit representations of payoff functions
- Payoff based algorithms:
  - No (explicit) communication among agents
  - Only requires ability to *measure* payoff upon deployment

- Initialization of *baseline action* and *baseline utility*:

$$a_i^b(1) = a_i(0)$$

$$u_i^b(1) = u_i(a(0))$$

- Action selection:

$$a_i(t) = a_i^b(t) \text{ with probability } (1 - \epsilon)$$

$a_i(t)$  is chosen randomly over  $\mathcal{A}_i$  with probability  $\epsilon$

- Baseline action & utility update:

*Successful  
Experimentation*

$$a_i(t) \neq a_i^b(t)$$

$$u_i(a(t)) > u_i^b(t)$$

⇓

$$a_i^b(t+1) = a_i(t)$$

$$u_i^b(t+1) = u_i(a(t))$$

*Unsuccessful  
Experimentation*

$$a_i(t) \neq a_i^b(t)$$

$$u_i(a(t)) \leq u_i^b(t)$$

⇓

$$a_i^b(t+1) = a_i^b(t)$$

$$u_i^b(t+1) = u_i^b(t)$$

*No  
Experimentation*

$$a_i(t) = a_i^b(t)$$

⇓

$$a_i^b(t+1) = a_i^b(t)$$

$$u_i^b(t+1) = u_i(a(t))$$



- **Marden, Young, Arslan, & JSS (2007)**: For potential games,

$$\lim_{t \rightarrow \infty} \Pr [a(t) \text{ is a NE}] > p^*$$

for any  $p^* < 1$  with sufficiently small exploration rate  $\epsilon$ .

- Suitably modified algorithm admits noisy utility measurements.

- How to assign individual payoff functions?
  - Induce “localization”
  - Have desirable NE
  - Produce potential game
- Proof methods:
  - “Sticky” NE
  - Characterization of steady state distribution
  - Stochastic stability

## Illustration: Rendezvous with obstacles

- Assume undirected connected constant graph (can be generalized)
- Global objective:

$$G(a_i, a_{-i}) = -\frac{1}{2} \sum_k \sum_{j \in \mathcal{N}_k} |a_k - a_j|$$

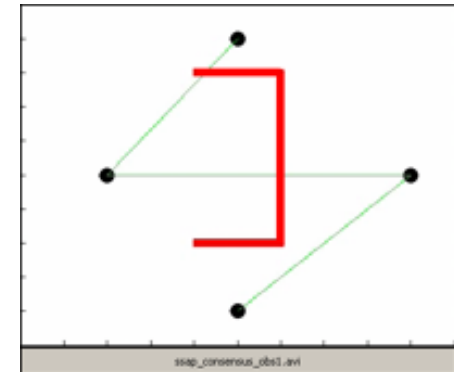
- Global objective without agent  $i$

$$G(\emptyset, a_{-i}) = -\frac{1}{2} \sum_{k \neq i} \sum_{j \in \mathcal{N}_k \setminus i} |a_k - a_j|$$

- Marginal contribution utility:

$$u_i(a_i, a_{-i}) = G(a_i, a_{-i}) - G(\emptyset, a_{-i}) = - \sum_{j \in \mathcal{N}_i} |a_i - a_j|$$

- Apply constrained SAP...

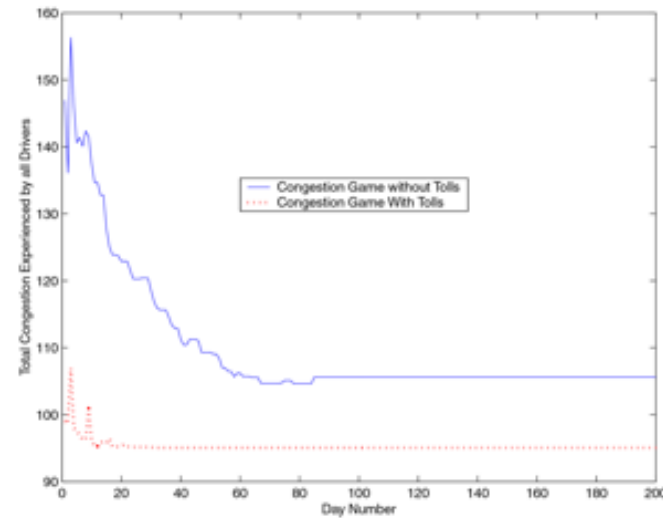
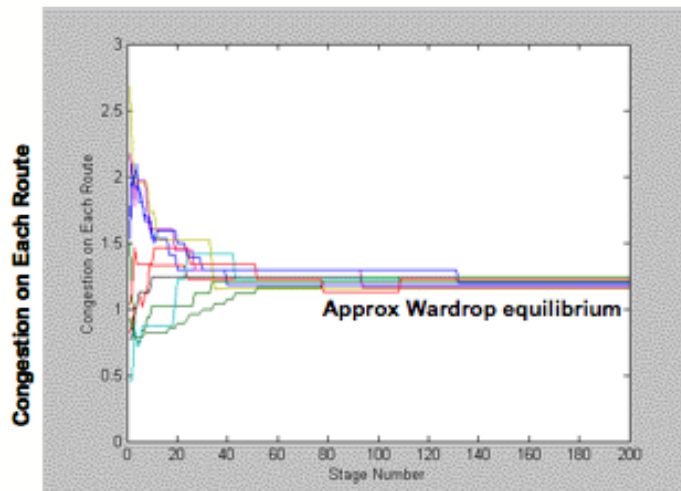
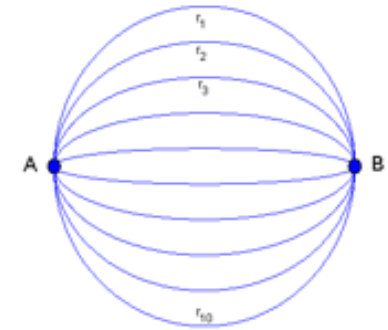


## Illustration: Distributed routing

- Setup: 10 parallel roads. 100 vehicles.
- Marginal contribution utility using overall congestion induces “tolls”

$$\tau_r(k) = (k - 1) \cdot (c_r(k) - c_r(k - 1))$$

- Apply max regret with inertia...



- Recap:

- Descriptive vs prescriptive
- Action/Oracle/Payoff based algorithms
- NE or potential function maximization
- Potential games & payoff design

- Future work:

- Convergence rates
- Exploiting prescriptive setting
- Agent dynamics
- Control theory and *descriptive* agenda

