# Strategic Learning and Dynamics in Networking and Computing Games

**Mihaela van der Schaar**

**Electrical Engineering, UCLA**

**Multimedia Communications and Systems Lab**

**http://medianetlab.ee.ucla.edu/**

**UCLA**

# Challenges for next-generation networks

- **Current status**
  - PHY layer innovations – significant capacity improvements
  - Source coding innovations – proliferation of a variety of applications
  - MAC, network and transport layers – often based on simplistic assumptions about users, ad-hoc rules, available information etc.

- **Key observations**
  - Collaborative communication/networking - OK for sensor nets, but most applications lack incentives for collaboration
  - Network and computing resources - shared among heterogeneous, intelligent users
  - Strategic behaviors of users - try to maximize their own utilities (even if this impacts the performance of other users)
  - Dynamic environment - not only channels/paths, but also source characteristics, application requirements, and ….
  - Informational decentralization: information required for resource management is decentralized (info is private to the users)

UCLA

# Illustrative example – Resource management in Current WLANs

- MAC protocol in IEEE 802.11a/b/g and e
  - Distributed Coordination Function (DCF)
  - Point Coordination Function (PCF)
- Underlying assumption for protocol design
  - **Users are not strategic**
    - Users have to follow protocol rules (e.g. CSMA/CA)
    - Users have to declare their resource requirements truthfully (e.g. in polling-based channel access – 802.11e HCF or 802.11a PCF)
    - Users have to collaborate with each other

**Problem 1:** Violates individual rationality of users and there are no incentives for users to adhere to these rules

**Problem 2:** Rules can be easily violated by simply adjusting communication algorithms' parameters, while still being protocol compliant; Often impossible to differentiate between users experiencing high traffic load/bad channel conditions, dumb users, and malicious users -> private information

# Existing resource management solutions assume non-strategic wireless users

Emphasis on fairness, not on incentives for truthful declaration

- Generalized Processor Sharing (GPS) [Gallager, 1993]
- Air-fair polling
- Cross-layer resource allocation schemes
  - Longest Queue receives Highest Possible Rate (LQHPR) [Yeh, 2003]
  - Cross-layer resource allocation by exploiting the packet priority and channel diversity [Zakhor, 2002][Scaglione, vdSchaar, 2005]
  - Utility-based resource allocation for multimedia applications [Girod 2006][Park, vdSchaar, 2006][Su, vdSchaar, 2007]

## Consequence: Tragedy of commons
- 802.11e Resource Allocation [vdSchaar, 2004, 2006]
- CSMA/CA in 802.11 WLAN [Cagalj, 2005][Konorski, 2006]
- R. W. Lucky, IEEE Spectrum, 2006

# What happens if users are strategic?

- CSMA/CA is vulnerable to selfish users using back-off attacks

- Selfish users gain significantly higher utility in IEEE 802.11e HCF protocol
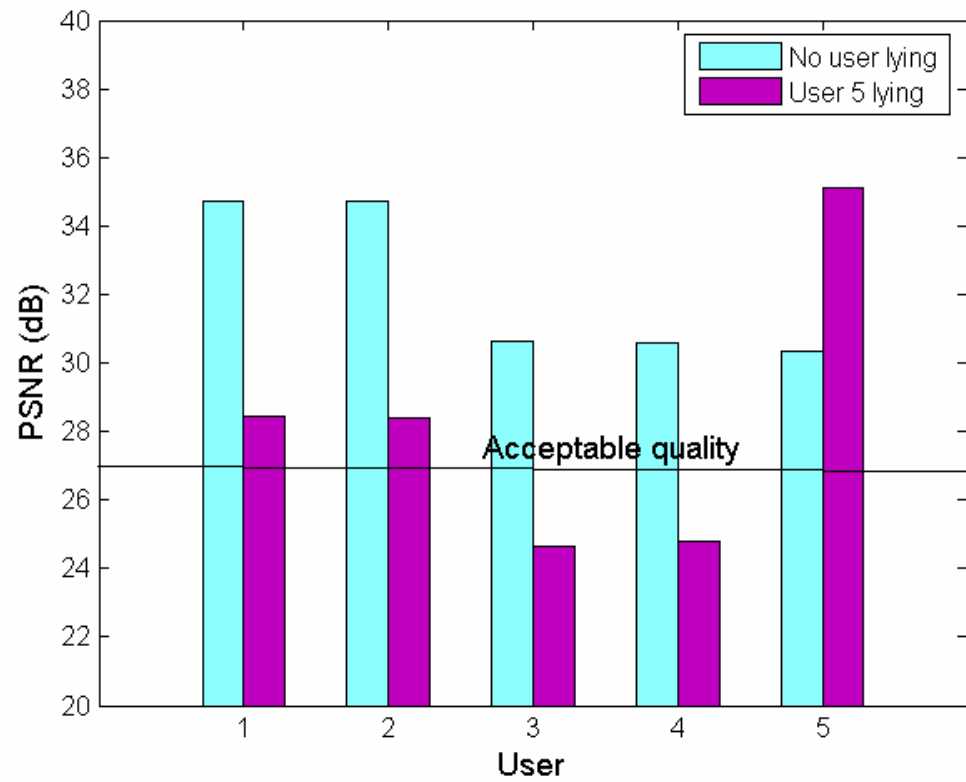
User 1: Foreman
User 2: Foreman
User 3: Coastguard
User 4: Coastguard
User 5: Mobile

Channel:
average SNR=23dB with variation 5dB



**UCLA** **Solutions?**

# Related research work (brief summary, not complete)

- Distributed power control [Cioffi][Poor][Pottie][Bambos]

- Dynamic spectrum access [Honig][Berry][Jordan][Liu]

- Routing/networking games [Lazar][Low]

- Mechanism design for wired networks [Lazar][Johari][Parkes]

- Bargaining games [Liu][MacKenzie]

- Network utility maximization for collaborative and/or homogeneous users [Chiang][Srikant]

- Existence of equilibriums in communication games, i.e. descriptive rather than constructive [Goodman][Poor][Johari, Goldsmith][Liu]

- Equilibrium selection/design and methods for getting to that equilibrium are key [Lazar][Altman]

UCLA

# Limitations of existing works/ Issues considered in our research

- **Information decentralization**
  - private information
  - information history – depends on the user's observations/protocols
  - strategic message exchanges
  - common knowledge – may differ
- **Different types of non-collaborative behavior**
  - self-interested users
  - malicious users
  - dumb users (do not optimize their cross-layer strategies efficiently)
- **Different strategies for playing the game**
  - foresighted vs. myopic users
  - risk neutral, adverse
- **Dynamics**
  - environment, but also other users (coupling between users)
- **Heterogeneity**
  - utility, experienced dynamics (traffic/loading, channels), complexity, (bounded) rationality
- **Users can learn -> not single-agent, but multi-agent learning**

# Design space for next-generation networks

## Rules

**Currently** ☹
Fairness – no consideration of resulting utility
Homogeneous users considered
No incentives to truthfully declare resource requirements/ rewards
No jamming prevention

**Desired** ☺
Resource management policies
- should be adapted based on the available resources, participating users, social decisions
- should consider the environment dynamics and users' heterogeneity– actions, strategies, utilities

**UCLA**

# Design space for next-generation networks

## Actions & Strategies

**Actions**
- protocol compliant
- unique algorithms in various layers allow users' differentiation

**Strategies – for selecting actions**
- depend on the available info
- foresighted/myopic
- malicious/altruistic
- risk-loving/risk-adverse

*Strategies* – probability of selecting various actions
In non-collaborative network environments, users do not want to use pure strategies, but rather use *mixed strategies*!

# Design space for next-generation networks
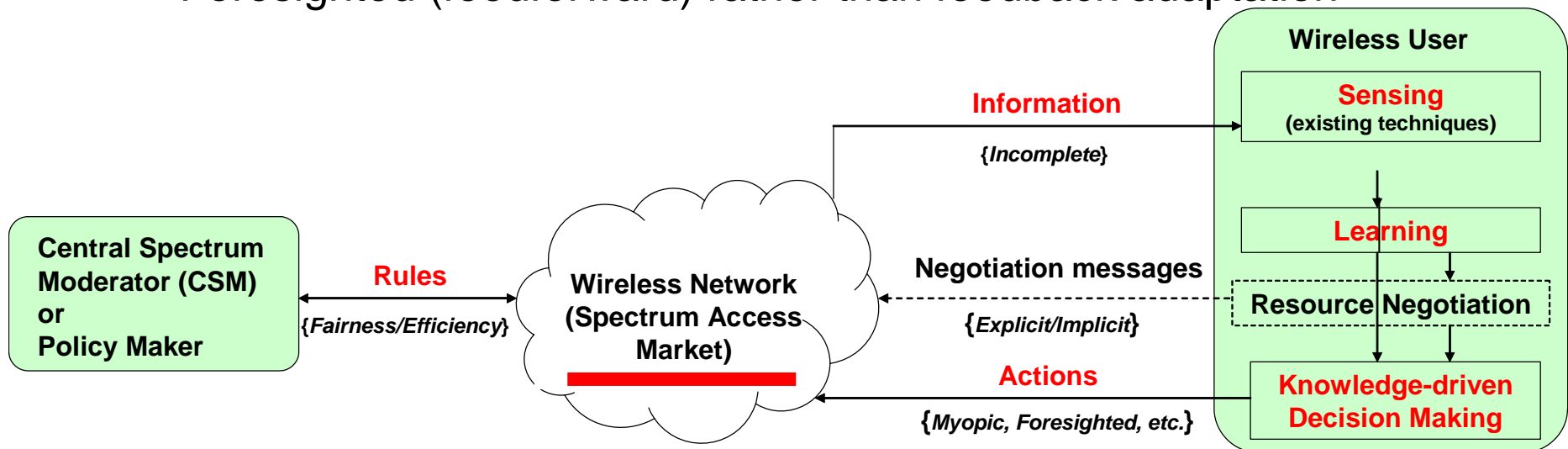
## Available information

**Information (heterogeneous)**
-information is private
-incomplete information about other network entities
(their actions, strategies, utilities, beliefs, etc.)
- common knowledge
-dynamic environment => time-varying information

UCLA

# Next-generation network design
## (NSF Career 2004)

- Model users as strategic agents playing a dynamic, stochastic game aimed at dividing network and/or computing resources

- The game is played with incomplete information

- Users can learn their environment (source and channel characteristics, but also competing users!) based on the available information, their utilities and limited computational abilities ->
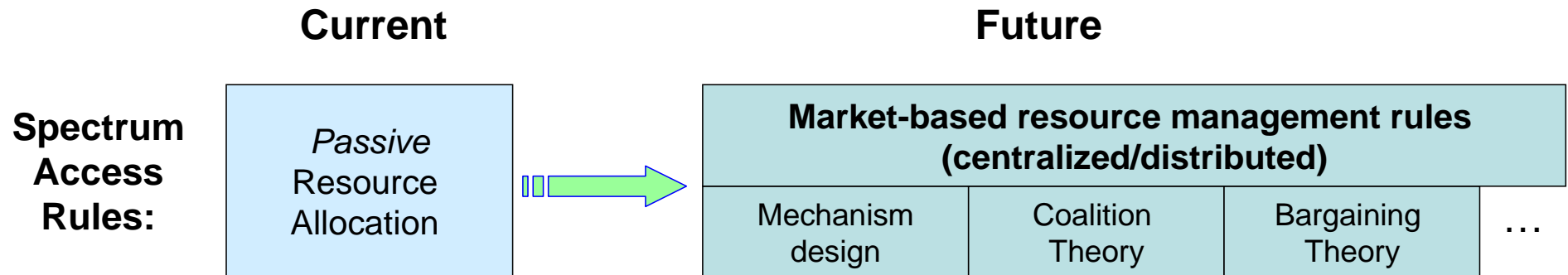
  Foresighted (feedforward) rather than feedback adaptation



## Why knowledge-driven?

Knowledge acquisition involves complex cognitive processes: sensing, learning, communication, association and reasoning. [Wikipedia]

# Creating dynamic resource markets/games

**Current**

**Future**

**Spectrum Access Rules:**

| *Passive* Resource Allocation |
| :---: |

| Market-based resource management rules (centralized/distributed) | | |
| :---: | :---: | :---: |
| Mechanism design | Coalition Theory | Bargaining Theory |

…

**Mechanism design** [Fu, vdSchaar, 2006, 2007]
**Coalition theory** [Park, vdSchaar, 2007]
**Bargaining** [Park, vdSchaar, 2006]
**Utility-driven resource allocation** [Scaglione, vdSchaar, 2005][Chen, vdSchaar, 2006][Su, vdSchaar, 2007]

**UCLA**

# Criteria for design & construction of dynamic resource markets/games

- Resource types
- One-shot versus multi-stage games
- Stochastic vs. repeated games
- Centralized vs. decentralized (who enforces the rules?)
- Social decisions (fairness rules)
- Budget-balanced vs. money-making resource allocation
- Consider what information the users' possess
- Selection/design of suitable equilibrium concepts
- Implementation

**UCLA**

# Proposed generalized stochastic game
## [Fu, vd Schaar, 2006, 2007]

Formally, the generalized stochastic game is defined as a tuple $(\mathcal{I}, \mathcal{S}, \mathcal{W}, \mathcal{A}, \mathcal{B}, \boldsymbol{P_s}, P_w, \mathcal{R})$, where

$\mathcal{I}$ is the set of agents (SUs), i.e. $\mathcal{I}=\{1,...,M\}$,

$\mathcal{S}$ is the set of state profiles of all SUs, i.e. $\mathcal{S}=\mathcal{S}_1 \times \cdots \times \mathcal{S}_M$ with $\mathcal{S}_i$ being the state set of SU $i$,

$\mathcal{W}$ is the set of network resource states,

$\mathcal{A}$ is the joint external action space $\mathcal{A}=\mathcal{A}_1 \times \cdots \times \mathcal{A}_M$, with $\mathcal{A}_i$ being the external action set of SU $i$,

$\mathcal{B}$ is the joint internal action space $\mathcal{B}=\mathcal{B}_1 \times \cdots \times \mathcal{B}_M$, with $\mathcal{B}_i$ being the internal action set of SU $i$ to transmit delay-sensitive data,

$\boldsymbol{P_s}$ is a transition probability function defined as a mapping from the current state profile $s \in \mathcal{S}$, corresponding joint external actions $a \in \mathcal{A}$ and internal actions $b \in \mathcal{B}$ and the next state profile $s' \in \mathcal{S}$ to a real number between 0 and 1, i.e. $\boldsymbol{P}: \mathcal{S} \times \mathcal{A} \times \mathcal{B} \times \mathcal{S} \mapsto [0,1]$,

$P_w$ is a transition probability function defined as a mapping from the current resource state $w \in \mathcal{W}$ and the next state $w' \in \mathcal{W}$ to a real number between 0 and 1, i.e. $P: \mathcal{W} \times \mathcal{W} \mapsto [0,1]$.

$\mathcal{R}$ is a reward vector function defined as a mapping from the current state profile $s \in \mathcal{S}$ and corresponding joint external and internal actions $a \in \mathcal{A}$ and $b \in \mathcal{B}$ to an $M$-dimensional real vector with each element being the reward to a particular agent, i.e. $\mathcal{R}: \mathcal{S} \times \mathcal{A} \times \mathcal{B} \mapsto \mathbb{R}^M$.
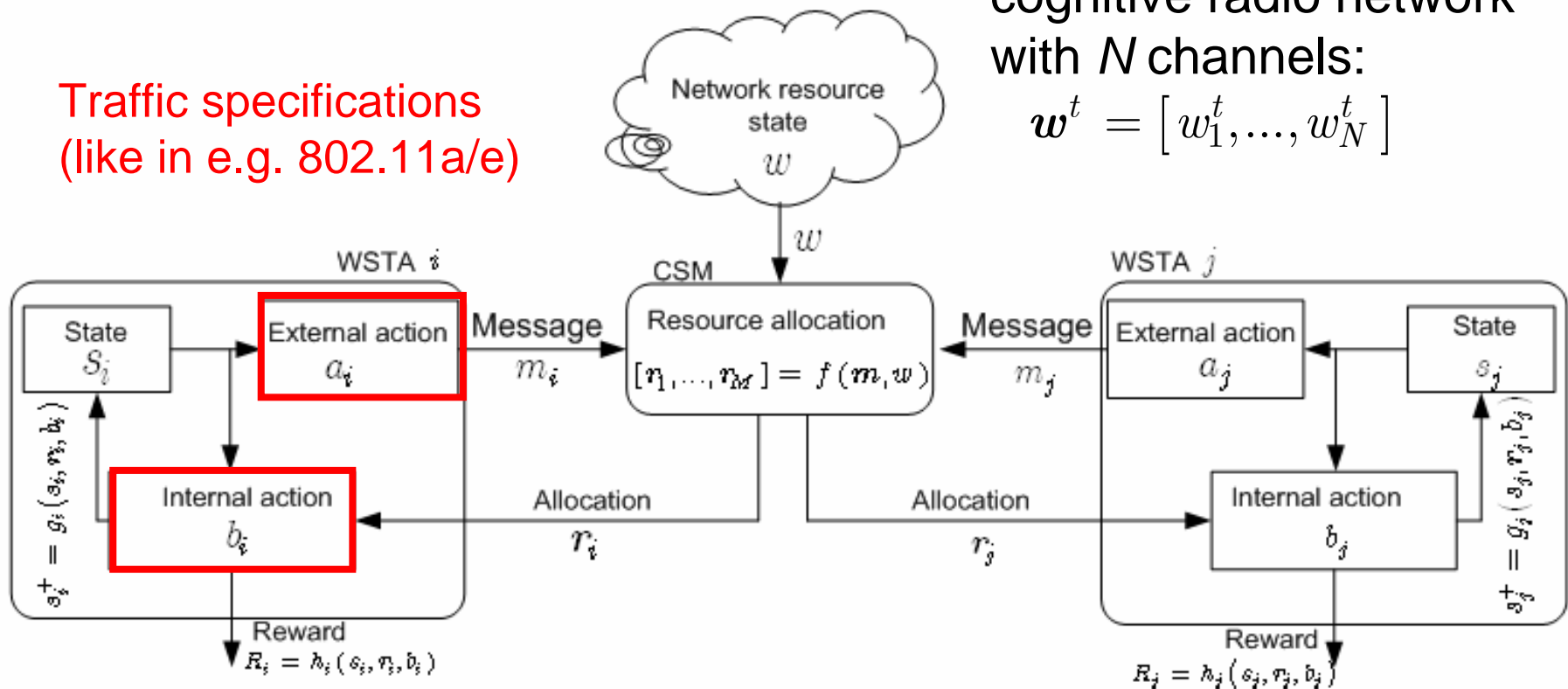
UCLA

# Centralized general stochastic game model

**Numerous networking/computing games:**
- Networks: 802.11 nets – polling based, Cellular nets, Cognitive radio nets
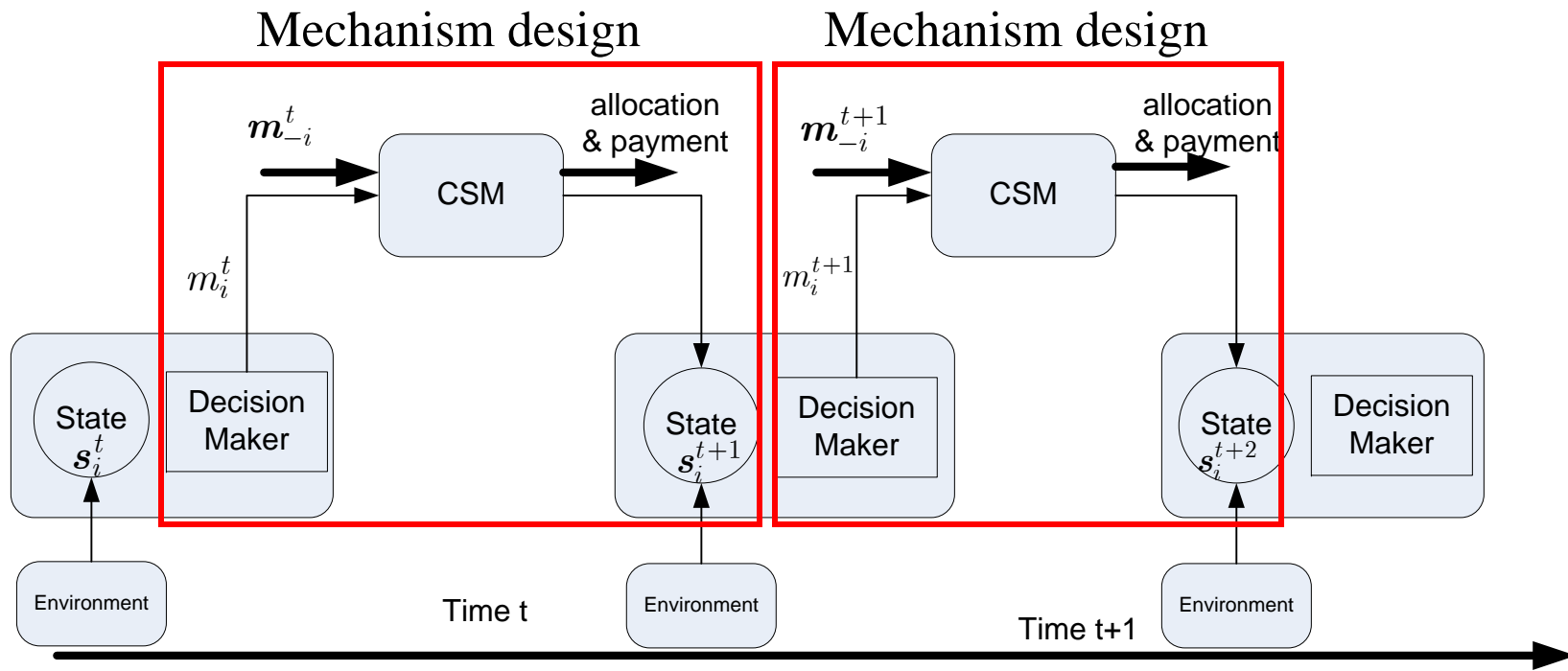- Computing systems: multi-tasks systems etc.

cognitive radio network with *N* channels:
$$\boldsymbol{w}^t = \left[ w_1^t, ..., w_N^t \right]$$

Traffic specifications
(like in e.g. 802.11a/e)



Retransmission limits, scheduling strategies, FEC etc.

UCLA

# Evolution of multi-user interaction



**Mechanism design**
-Solutions: VCG, pricing mechanism, generalized auctions etc.
-Informational and complexity requirements
-Equilibrium selection: Nash, dominant etc.
-Incentives for truthful revelation

# Centralized general stochastic game – moderator side (example)

- After each wireless user submits a bid vector $m_i^t = a_i^t$, and CSM performs two computations:

  (i) channel allocation and (ii) payment computation
  $$\boldsymbol{r}^t = \left( \boldsymbol{z}^t, \boldsymbol{\tau}^t \right) = \Omega \left( \boldsymbol{a}^t, w^t \right)$$

- Social welfare (fairness):
  $$\boldsymbol{z}^{t,opt} = \arg \max_{\boldsymbol{z}^t} \sum_{i=1}^{M} \tilde{h}_i \left( a_i^t, \boldsymbol{z}_i^t, w \right)$$

  utility function of user i as known by the CSM

- Taxation – assume second price auction

  [Klemperer,1999][Sun, Modiano, Zheng, 2006]
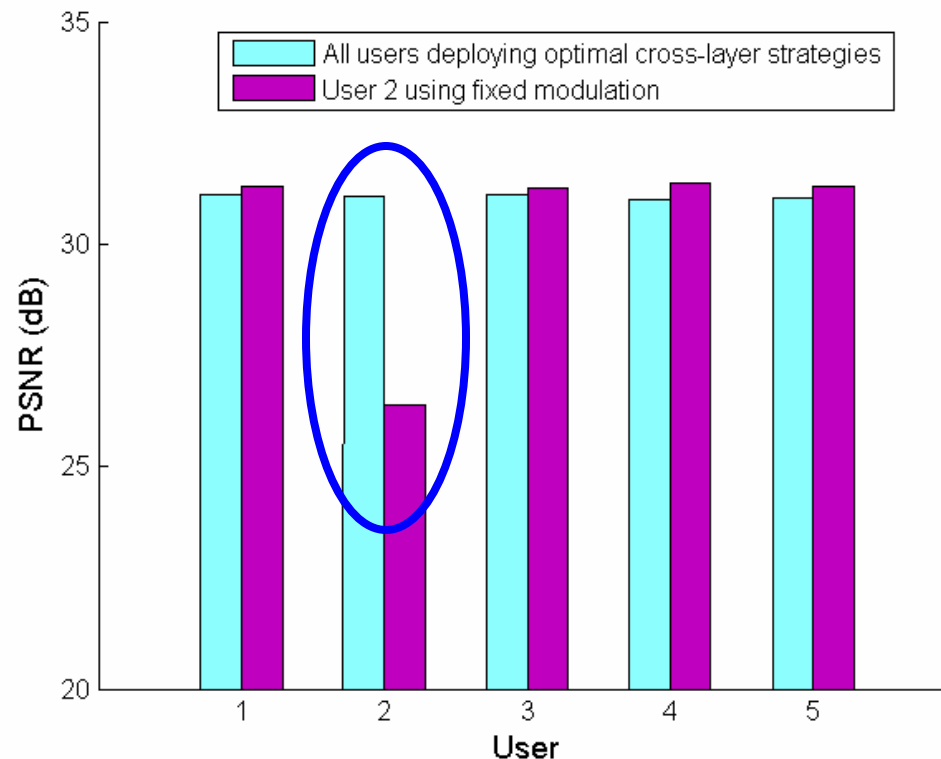  $$\tau_i^t = \sum_{\substack{j=1, \\ j \neq i}}^{M} \tilde{h}_j \left( a_j^t, \boldsymbol{z}_j^{t,opt}, w \right) - \max_{\boldsymbol{z}_{-i}^t} \sum_{\substack{j=1, \\ j \neq i}}^{M} \tilde{h}_j \left( a_i^t, \boldsymbol{z}_i^t, w \right)$$

UCLA

# Truthful revelation?

For one-shot games in wireless communication games (e.g. one-time resource allocation, like in 802.11e HCF), we proved that [F.Fu, vdSchaar, 2006]

- – Optimal strategy is to adopt the best anticipated cross-layer strategy and reveal the "true" type (utility function)
- – Optimal strategy is dominant ☺ , and thus, it can be chosen without knowing other users' strategies

- Why is dominant strategy equilibrium desirable?
  - – No need to know other users' actions/strategies –> can use single agent learning
- For multi-stage games – everything gets more interesting ☺

*UCLA*

# Illustrative Results –
## Impact of wireless users "smartness" (selected algorithms and cross-layer optimization)



All users are transmitting Foreman video sequences.

Channel: average SNR=23dB with variation 5dB

**UCLA**

# How to play the stochastic game?

- History & observation
  - History: $h^t = \{s^0, w^0, a^0, b^0, z^0, \tau^0, ..., s^{t-1}, w^{t-1}, a^{t-1}, b^{t-1}, z^{t-1}, \tau^{t-1}, s^t\} \in \mathcal{H}^t$
  - Observation :

$$o_i^t = \{s_i^0, w^0, a_i^0, b_i^0, z_i^0, \tau_i^0, ..., s_i^{t-1}, w^{t-1}, a_i^{t-1}, b_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, s_i^t\} \subset h^t , \quad o_i^t \in \mathcal{O}_i^t$$

- Policy $\quad \pi_i^t : \mathcal{O}_i^t \mapsto \mathcal{A}_i \times \mathcal{B}_i \quad \left[a_i^t, b_i^t\right] = \pi_i^t(o_i^t)$

$$\boldsymbol{\pi}_i = \left(\pi_i^0, ..., \pi_i^t, ...\right) , \quad \boldsymbol{\pi}^t = \left(\pi_1^t, ..., \pi_M^t\right) = \left(\pi_i^t, \boldsymbol{\pi}_{-i}^t\right)$$

- Reward: $R_i^t(s_i^t, \boldsymbol{r}_i^t, b_i^t) = u\left(s_i^t, z_i^t, b_i^t\right) + \tau_i^t \quad \longrightarrow \quad R_i^t(s_i^t, o_i^t, b_i^t)$

- Discounted reward: $\quad Q_i^t((\pi_i^t, \boldsymbol{\pi}_{-i}^t) \mid s^t, w^t) = \sum_{k=t}^{\infty} (\alpha_i)^{k-t} R_i^k(s_i^k, o_i^k, b_i^k) ,$

- Best response: $\quad \beta_i(\boldsymbol{\pi}_{-i}^t) = \arg\max_{\pi_i} Q_i^t((\pi_i^t, \boldsymbol{\pi}_{-i}^t) \mid s^t, w^t)$

# Key challenge

- An SU *may not exactly know* the other SUs' actions and models, and it cannot know their private information
- Thus, an SU *can only predict the dynamics (uncertainties)* caused by the competing SUs based on its observations from past interactions

For instance, in wireless networks:

*Private information* (e.g. characteristics of the application traffic, channel gain or channel conditions - SINR, etc.)

*Network information* (e.g. network resource states, primary users etc.)

*Opponents information* (e.g. states and possible actions of the opponents)

## How to solve this problem? Multi-agent learning!

UCLA

# What information should be learnt?

$$\pi_i^* = \arg\max_{\pi_i} Q_i\left(\pi_i, \boldsymbol{\pi}_{-i} \mid s_i, \boldsymbol{s}_{-i}, \boldsymbol{w}\right)$$

To solve this optimization, the following information is required by SU $i$:

1. the state transition model of SU $i$, $p\left(s_i^{t+1} \mid s_i^t, a_i^t, \boldsymbol{a}_{-i}^t, b_i\right)$;

2. the state transition model of other SUs, $p\left(s_j^{t+1} \mid s_j^t, a_j^t, \boldsymbol{a}_{-j}^t, b_j\right), \forall j \neq i$;

3. the state of other SUs, $\boldsymbol{s}_{-i}$;

4. the policy of other SUs, $\boldsymbol{\pi}_{-i}$;

5. the network resource state $w$.

**UCLA**

# Multi-agent learning - definition

We define a **learning algorithm** $\mathcal{L}_i$ as:

$$\left[a_i^t, b_i^t\right] = \pi_i^t\left(s_i^t, B_{\boldsymbol{s}_{-i}}^t, B_{\boldsymbol{\pi}_{-i}}^t, B_w^t\right)$$

**Output of the multi-user interaction game**:

$$\Omega^t = Game\left(\boldsymbol{s}^t, \boldsymbol{a}^t, w^t\right)$$

**Observation** of SU $i$

$$o_i^t = O\left(s_i^t, \Omega_i^t, b_i^t\right),$$

where $O$ is the observation function which depends on the current state, the current game output and the current internal action taken.

**Policy update:**

$$\pi_i^{t+1} = \mathcal{F}_i\left(\pi_i^t, o_i^t, I_{-i}^t\right)$$

$\mathcal{F}$ is the update function about the belief and policies
$I_{-i}^t$ is the exchanged information with the other SUs

**Beliefs** about the other SUs' states $s_{-i}$, policies $\pi_{-i}$ and the network resource state $w$:

$$B_{\boldsymbol{\pi}_{-i}}^{t+1} = \mathcal{F}_{\boldsymbol{\pi}_{-i}}\left(B_{\boldsymbol{\pi}_{-i}}^t, o_i^t, I_{-i}^t\right) \quad, \ B_w^{t+1} = \mathcal{F}_w\left(B_w^t, o_i^t, I_{-i}^t\right), \ B_{\boldsymbol{s}_{-i}}^{t+1} = \mathcal{F}_{\boldsymbol{s}_{-i}}\left(B_{\boldsymbol{s}_{-i}}^t, o_i^t, I_{-i}^t\right)$$

UCLA

# Value of Learning [F.Fu,vdSchaar, 2007]

$$\mathcal{V}^{\pi_i^{\mathcal{L}_i(o_i,I_{-i})}}(T) = \frac{1}{T}\sum_{t=1}^{T} R_i^t\left(\pi_i^{\mathcal{L}_i(o_i,I_{-i})}\right)$$

where the reward $R_i^t$ depends on both the learning approach $\mathcal{L}_i$ and on the observation $o_i^t$ and information exchanged $I_{-i}^t$
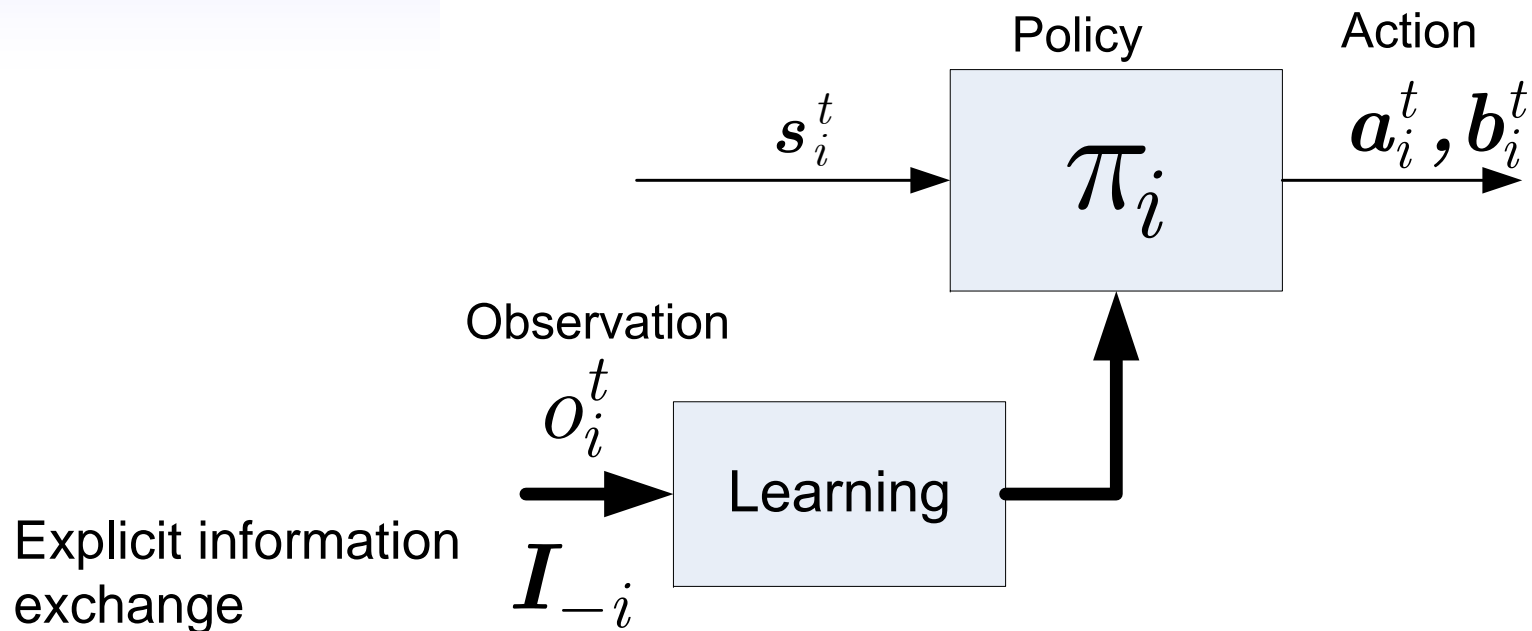
For instance, given the same observation $o_i^t$ and exchanged information $I_{-i}^t$, if the time average rewards of two algorithms $\mathcal{L}_i'$ and $\mathcal{L}_i''$ satisfy $\mathcal{V}^{\pi_i^{\mathcal{L}_i'(o_i,I_{-i})}}(T) > \mathcal{V}^{\pi_i^{\mathcal{L}_i''(o_i,I_{-i})}}(T)$, then we say that learning algorithm $\mathcal{L}_i'$ is better than $\mathcal{L}_i''$

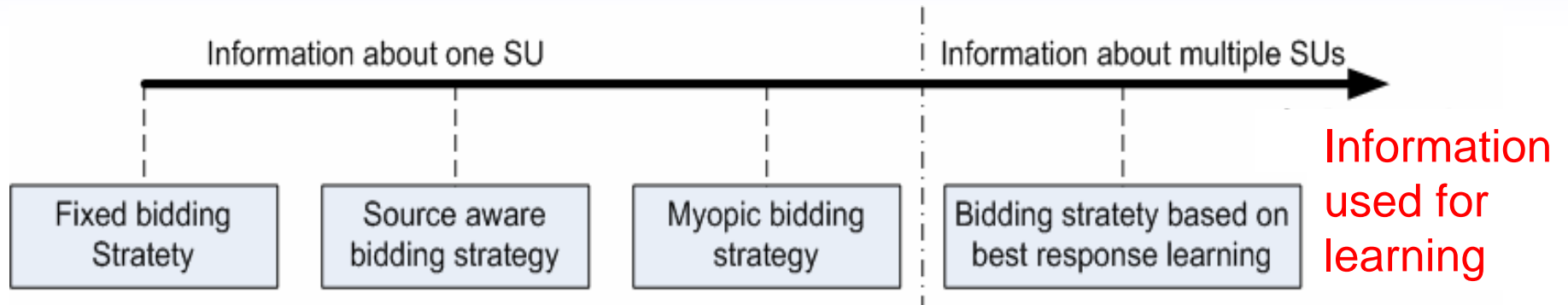**How much to learn for a desired performance (utility)?**
[Y. Su, vdSchaar, 2008]

**UCLA**

# Multi-agent learning - illustration

Policy       Action

$s_i^t \longrightarrow$   $\pi_i$   $\longrightarrow a_i^t, b_i^t$

Observation

$o_i^t \longrightarrow$   Learning

Explicit information
exchange   $I_{-i}$

Solutions depend on the information availability:
- Reinforcement learning (no explicit modeling of other users)
[Fu, vdSchaar, 2007]
- Fictitious Play (explicit modeling of other users – needs to
know what actions opponents took, but not their strategies)
[Shiang, vdSchaar, 2007]

UCLA

# Illustrative results for bidding and learning strategies

Information about one SU | Information about multiple SUs



Information used for learning

- Fixed bidding strategy $\pi_i^{fixed}$: this strategy generates a constant bid vector during each stage of the auction game, irrespective of the state that SU $i$ is currently in and of the states other SUs are in.

- Source-aware bidding strategy $\pi_i^{source}$: this strategy generates various bid vectors by considering the dynamics in source characteristics (based on the current buffer state), but not the channel dynamics.

- Myopic bidding strategy $\pi_i^{myopic}$: this strategy takes into account both the environmental disturbances and the impact caused by other SUs. However, it does not consider the impact on its future rewards.

- Bidding strategy based on best response learning $\pi_i^{\mathcal{L}_i}$: This strategy is produced using the presented learning, which considers both the environmental dynamics and the impact on the future reward.

UCLA

# Illustrative results
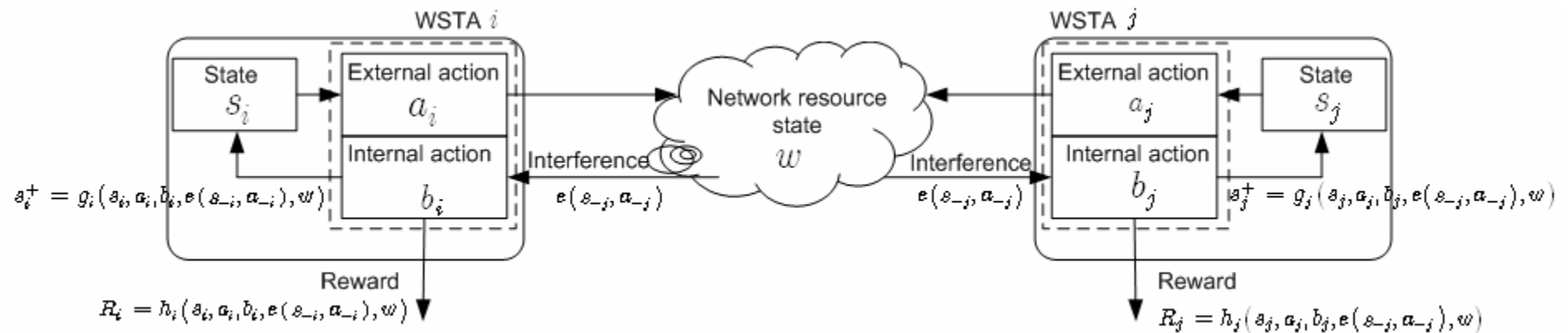
Coastguard video sequence, 500 ms delay

**Performance of competing SUs with various bidding strategies**

| | Bidding Strategies | SU 1 | | | SU 2 | | |
|---|---|---|---|---|---|---|---|
| | | Video Quality (PSNR) | Average tax | Average reward | Video Quality (PSNR) | Average tax | Average reward |
| Scenario 1 | $\pi_1^{fixed}, \pi_2^{myopic}$ | 25 dB | 0.1222 | 2.6337 | 36 dB | 0.5495 | 1.5105 |
| Scenario 2 | $\pi_1^{source}, \pi_2^{myopic}$ | 26 dB | 0.3147 | 2.4915 | 33 dB | 0.6048 | 1.6116 |
| Scenario 3 | $\pi_1^{myopic}, \pi_2^{myopic}$ | 29 dB | 0.4669 | 1.9767 | 30 dB | 0.3763 | 1.7837 |
| Scenario 4 | $\pi_1^{\mathcal{L}}, \pi_2^{myopic}$ | 35 dB | 0.6923 | 1.7428 | 27 dB | 0.4197 | 2.2967 |

UCLA

# Distributed stochastic games

## Numerous networking/computing games:
- Networks: power control games, contention games etc.
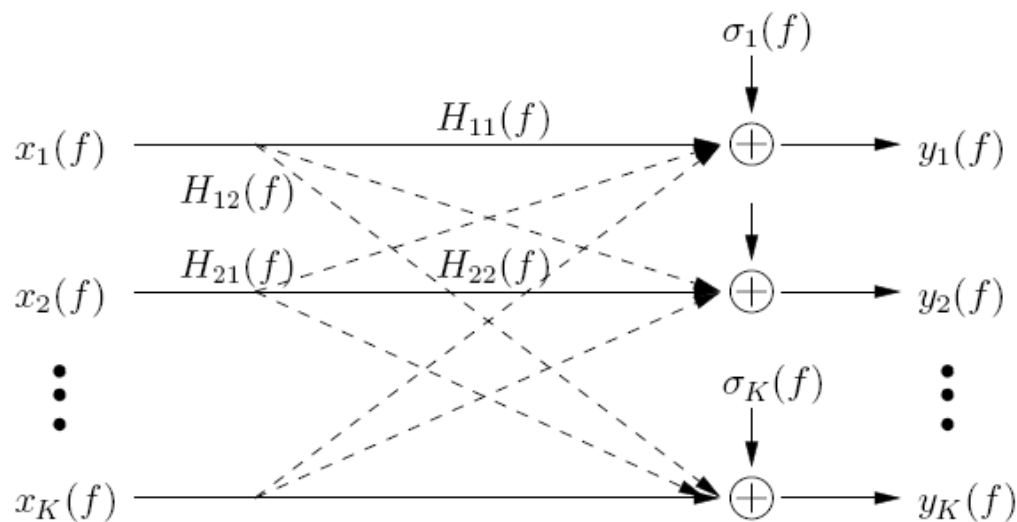- Computing systems: peer-to-peer, multi-tasks systems etc.



E.g. in power control games:
external action can be the selected power allocation,
internal action can be the selected modulation and channel coding scheme
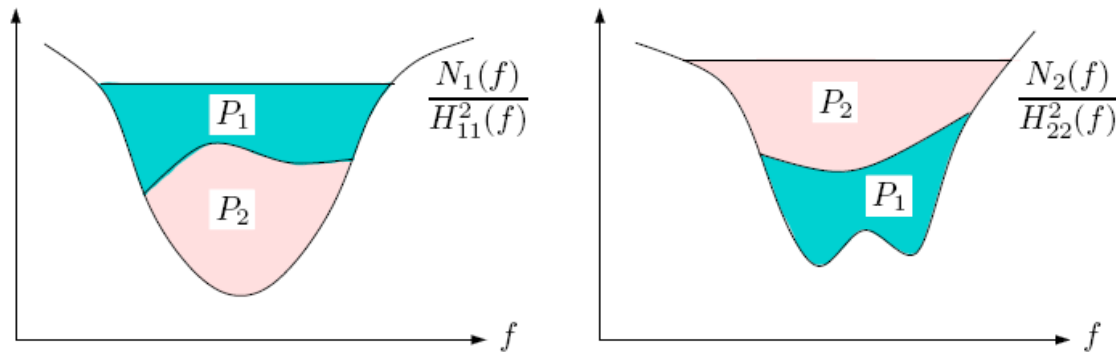
**UCLA**

# Distributed games - Illustrative results

- Multi-user power control problem
  - Interference-limited multi-user communication systems
  - Frequency-selective channels
  - Transmit PSD design
- Goal
  - Maximize selfish users' rates



UCLA

# Existing solutions

**Solution – Iterative waterfilling (W. Yu, J. Cioffi, 2002)**



$$P_1^{(0)}(f) \rightarrow P_2^{(0)}(f) \rightarrow P_1^{(1)}(f) \rightarrow P_2^{(1)}(f) \rightarrow \cdots$$

- Nash equilibrium: competitive optimal  ☹
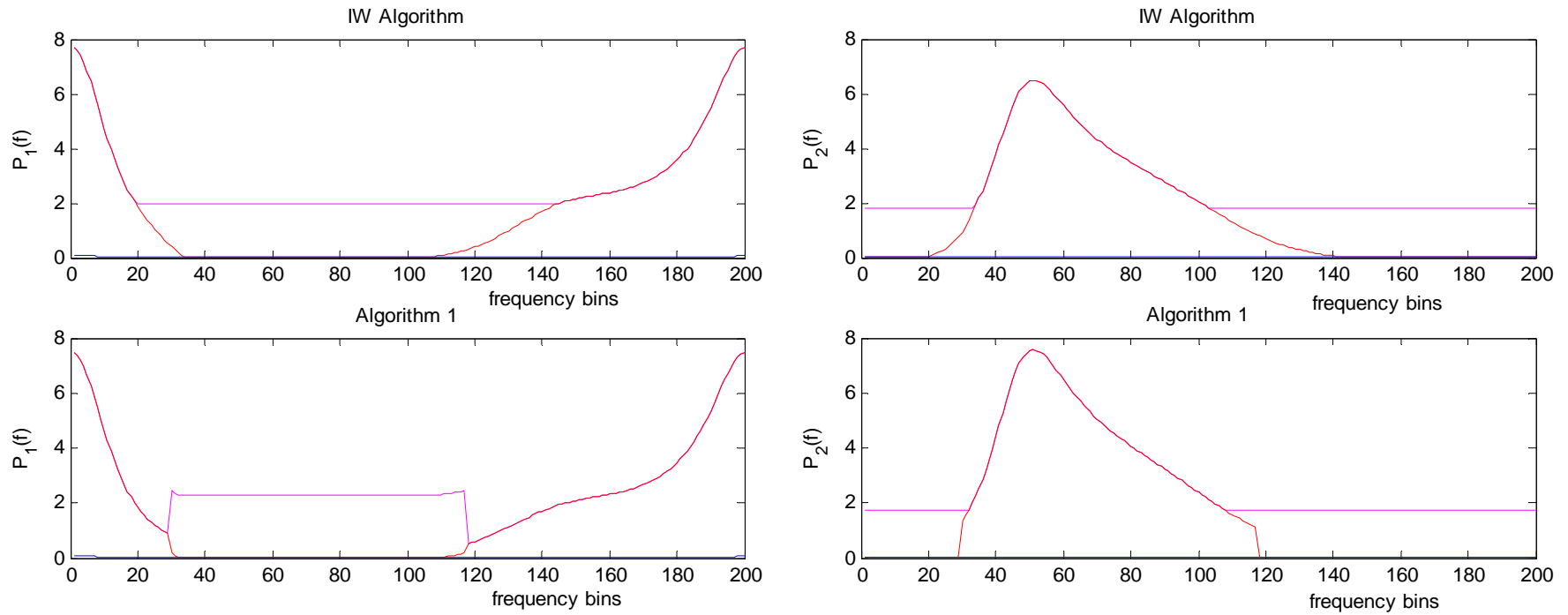- Convergence is achieved by iterative water-filling

Can we do better? How?

UCLA

# A New Perspective on Multi-user Power Control Games in Interference Channels [Y. Su, vdSchaar,2007]

- Iterative Waterfilling =>Myopic users -> Nash equilibrium
- Foresighted strategy in determining the transmit PSD -> Stackelberg equilibrium
  - Bi-level programming formulation
  - A low-complexity sub-optimal approach based on the necessary KKT conditions

# Illustrative results



- **Substantial performance improvements for both foresighted and myopic users ! ☺**

- **How to achieve this result using learning?**

**UCLA**

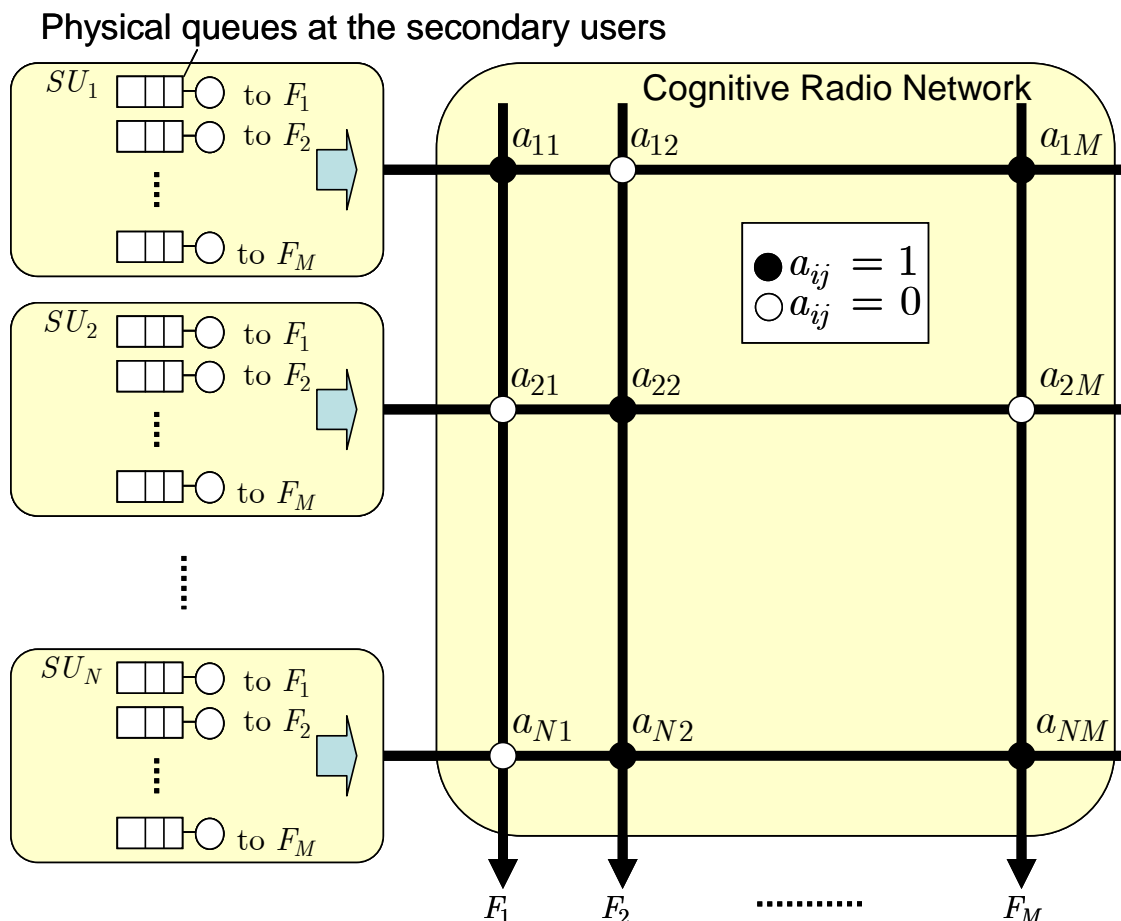# Preliminary results for different learning schemes in repeated power control games

Simulation results using different learning techniques

| Adopted schemes | SU | Reward (Kbit/joule) | Average reward |
|---|---|---|---|
| Myopic scheme | 1 | 519.0 | |
| | 2 | 195.2 | |
| | 3 | 530.6 | 890.15 |
| | 4 | 2073.0 | |
| | 5 | 1132.9 | |
| AR learning scheme | 1 | 555.2 | |
| | 2 | 113.5 | |
| | 3 | 345.6 | 1005.6 |
| | 4 | 2830.2 | |
| | 5 | 1183.7 | |
| AA learning scheme | 1 | 529.3 | |
| | 2 | 475.6 | |
| | 3 | 476.8 | 1069.3 |
| | 4 | 2831.2 | |
| | 5 | 1033.3 | |

**Adaptive Reinforcement (AR)**

**Adaptive Action Learning (AA)**

UCLA    **Stackelberg (perfect info.) Average Reward: 1250**

# Distributed and dynamic resource management
## *with information exchanges* [H. Shiang, vdSchaar, 2007]

Physical queues at the secondary users
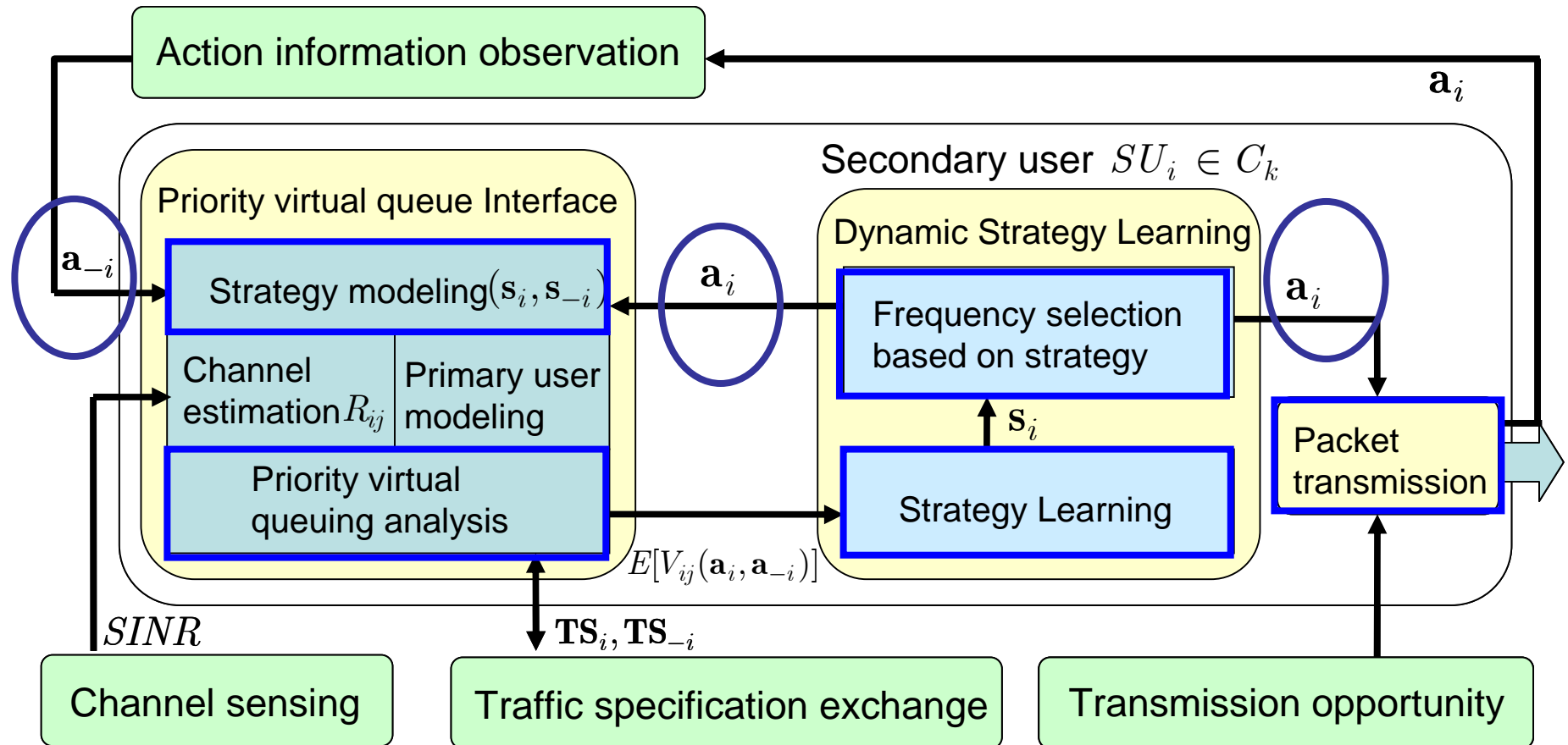


Delay-sensitive users

$$\mathbf{s}_i^* = \arg \max_{\mathbf{s}_i \in \mathcal{S}^M} E_{(\mathbf{s}_i, \mathbf{s}_{-i})}[u_i(\mathbf{a}_i, \mathbf{a}_{-i})]$$

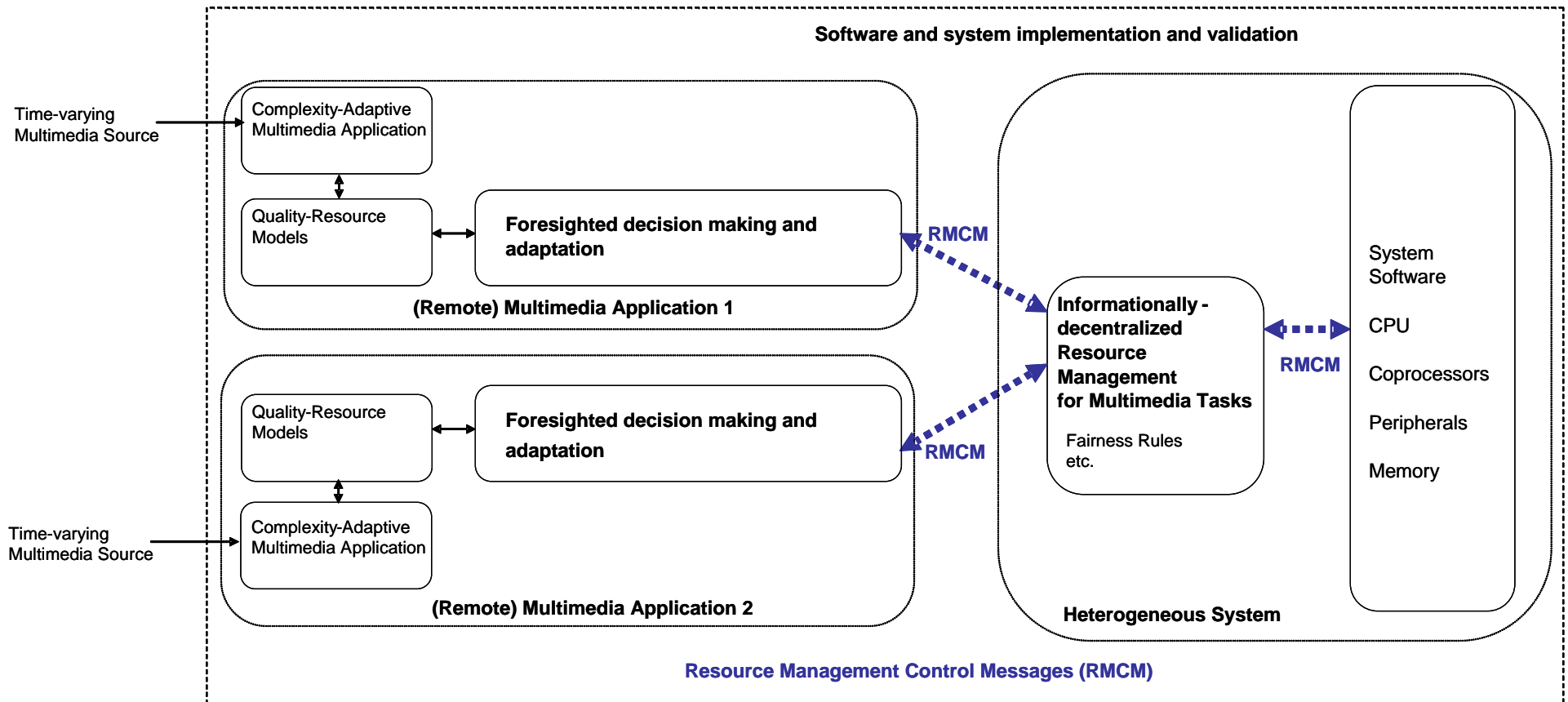Virtual queue interface for modeling inter-user communication

$$\mathbf{s}_i^* = \arg \max_{\mathbf{s}_i \in \mathcal{S}^M} \sum_{j=1}^{M} s_{ij} \cdot E[V_{ij}(\mathbf{a}_i, \mathbf{a}_{-i})]$$

Model-based learning
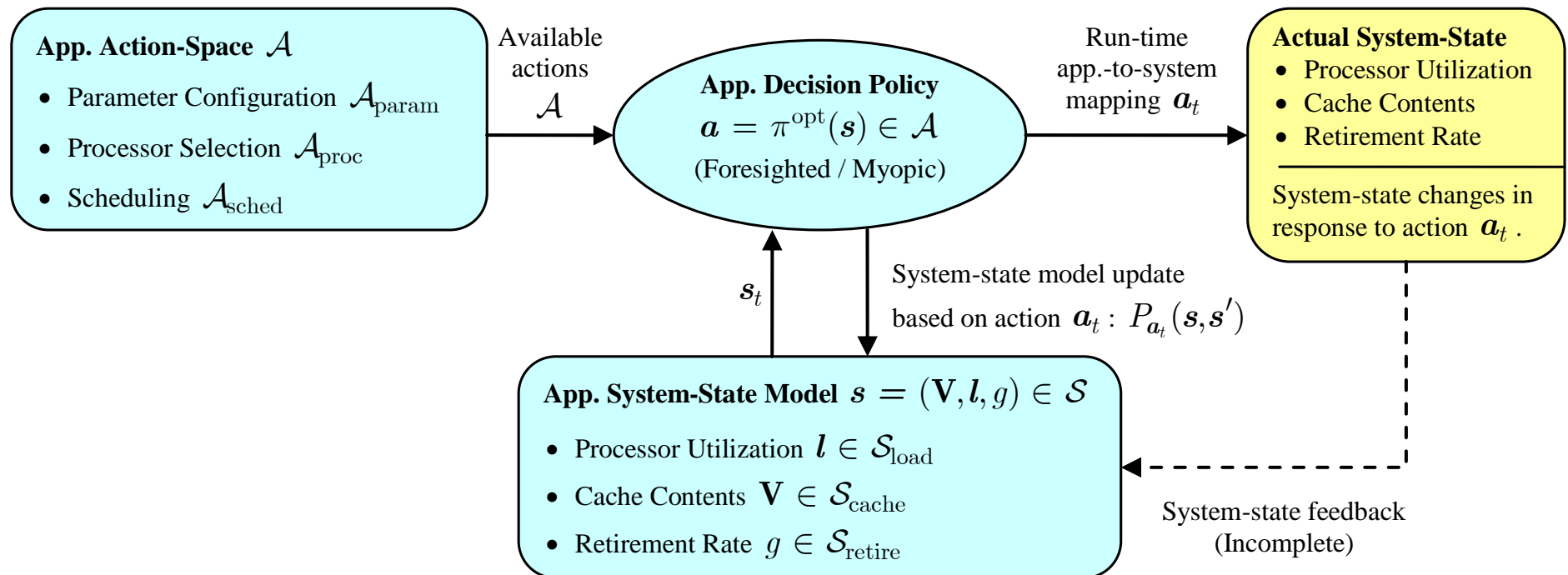
**UCLA**

# Dynamic Strategy Learning



**UCLA**

# Foresighted adaptation and learning in computing games



[Foo, vdSchaar, 2006,2007,2008][Akyol, vdSchaar, 2006]
[vdSchaar, Andreopoulos, 2005]

UCLA

# Illustration of how the application decision policy takes actions based on the system-state model,

# and how these actions impact the actual system state



**App. Action-Space** $\mathcal{A}$

- Parameter Configuration $\mathcal{A}_{\text{param}}$
- Processor Selection $\mathcal{A}_{\text{proc}}$
- Scheduling $\mathcal{A}_{\text{sched}}$

Available actions $\mathcal{A}$

**App. Decision Policy**

$$a = \pi^{\text{opt}}(s) \in \mathcal{A}$$

(Foresighted / Myopic)

Run-time app.-to-system mapping $a_t$

**Actual System-State**

- Processor Utilization
- Cache Contents
- Retirement Rate

System-state changes in response to action $a_t$.

$s_t$

System-state model update based on action $a_t$: $P_{a_t}(s, s')$

**App. System-State Model** $s = (\mathbf{V}, l, g) \in \mathcal{S}$

- Processor Utilization $l \in \mathcal{S}_{\text{load}}$
- Cache Contents $\mathbf{V} \in \mathcal{S}_{\text{cache}}$
- Retirement Rate $g \in \mathcal{S}_{\text{retire}}$

System-state feedback (Incomplete)

**UCLA**

# Our Goal

**Add a new dimension to multi-user networks/systems by explicitly considering strategic users, dynamics, heterogeneity and information availability**

- Opens opportunities for new theoretical foundations and algorithm designs, new metrics needed

- Significant performance improvements

- Backwards compatible with existing protocols

- Simple system designs for building next-generation dynamic, robust and trustable networks

**UCLA**

# Multimedia Communications and Systems Laboratory

## See our research at:
### http://medianetlab.ee.ucla.edu



**Current Ph.D. Students**
Fangwen Fu
Hyunggon Park
Hsien-Po Shiang
Brian Foo
Nicholas Mastronarde
Zhichu Lin
Yi Su

**Current M.Sc. Students**
Wenchi Tu

UCLA