

On wireless scheduling with partial channel-state information

Aditya Gopalan, Constantine Caramanis and Sanjay Shakkottai

Abstract— We consider a single server serving a time-slotted queued system of multiple flows, where not more than one channel can be serviced in a single time slot. Each flow has exogenous arrivals, and the service rates to the flows vary over time according to a fixed distribution. The server is allowed to observe the service rates for only a single *subset* of flows (chosen from a fixed collection of *observable* subsets) in a time slot for the purpose of making scheduling decisions. We provide a precise characterization of the stability region for such a system. We present an online scheduling algorithm that uses information about marginal distributions to pick the subset and the MaxWeight rule to pick a flow within the subset, and show that it is throughput-optimal. In the case where the observable subsets are all disjoint, we show that a simple scheduling algorithm - Max-Sum-Queue - that essentially picks subsets having the largest squared-sum of queues, followed by MaxWeight within the subset, is throughput-optimal. We show that for channels which are symmetric with respect to channel rates and distributions, and fixed-size observable subsets, Max-Sum-Queue is throughput-optimal. Finally, we demonstrate that under certain conditions, Max-Sum-Queue may not be throughput-optimal.

I. INTRODUCTION

There has been much recent interest in scheduling over wireless cellular networks where channel state information is available at the base-station [1], [2], [12]. A canonical system consists of a base-station and a collection of mobile users. Time is slotted (typically of the order of a milli-second), like in the high-speed WiMAX [9], Ultra Mobile Broadband (UMB), GSM-based HSDPA and 1xEV-DO communications technologies. In each time-slot, the channel state (the channel quality such as SINR or data rate that can be sustained over the time-slot to the mobile) is potentially available (via a feedback channel from the mobile terminals to the base-station) at the base-station. Based on the load (packets queues at the base-station) as well as the channel state, the base-station schedules users for channel access each time-slot.

However, as the capacity of the wireless system increases, it is likely that a large number of mobile users will be connected to the base-station. Thus, transmitting channel state feedback from all of the mobile to the base-station might be difficult due to feedback bandwidth constraints. A reasonable approach would be for the base-station to request channel state from a sub-collection of users (for example, users which have a large backlog of data at the base-station)

and make scheduling decisions based on this partial channel state information.

A. Main Contributions

We consider a base-station system where there are N users and channels, with each user generating data, and with channels which have an arbitrary joint distribution over a finite state-space (the channel is assumed to be independent across time but not across users), and the server *does not* have knowledge of the channel joint distribution.

In each time-slot, the base-station is allowed to get channel state¹ from one among a predefined collection of subsets of channels (for example, in a ten user system, the constraint could be that we can acquire channel state from at-most three users per time-slot). We henceforth refer to this as a system with partial channel-state information.

The scheduling task at each time-slot is to first determine the subset (of channels for which channel state will be acquired) and then determine a single user to schedule from within this subset. In this paper, we characterize the stability region for this multi-user system, and develop algorithms that achieve the stability region. The main contributions in this paper are as follows.

- (i) We derive the stability region for a system with N users and an arbitrary collection of observable subsets (i.e., a collection of subset of users for which the channel state can be simultaneously acquired), and for any joint channel distribution (across users) that are independent and identically distributed over time. The stability region corresponds to the set of arrival rates that can be sustained such that the queues at the base-station are stable (positive recurrent).

We show that the stability region is described by the convex hull of the “local” stability region for each observable subset. In other words, for each observable subset, we first consider the rates that can be sustained if none of the other users had any data. This can be characterized via a convex polytope that corresponds to the stability region (the “local” stability region for the subset) of a reduced system where users that do not belong to the subset are removed, and the base-station has *complete* channel state information for the users within the subset. The convex hull of such “local”

This work was partially supported by NSF Grants CNS-0325788, CNS-0347400, CNS-0519535, CNS-0721532, EFRI-0735905 and the Darpa ITMANET program. A. Gopalan, C. Caramanis, and S. Shakkottai are with the Department of Electrical and Computer Engineering, The University of Texas at Austin, Austin TX, and part of the Wireless Communications and Networking Group (WNCG). e-mail: {gopalan, cmcaram, shakkottai}@ece.utexas.edu.

¹At each time-slot, the complete channel state is a N dimensional vector, with the i -th component of the vector corresponding to the data rate that can be sustained to the i -th mobile user over the time-slot if this user is chosen by the scheduler. Correspondingly, the partial channel state corresponds to a sub-vector of this N dimensional vector.

stability regions over all observable subsets describe the stability region with partial channel state information. We also present a numerical example that illustrates the degradation in the stability region as the amount of channel state information decreases (i.e., fewer simultaneously observable channels). The example suggests that the stability region with partial channel information can be viewed as the intersection of the stability regions with *complete* channel information over all systems whose channels have a joint distribution that is consistent with the subset-marginal distributions (i.e., channel distributions only over the observable subsets) of the given system with partial channel-state information.

- (ii) Next, we develop a queue length based “online” policy that uses the queue-length information along with subset-marginal distributions that is throughput-optimal, i.e., the policy attains all rate points within the stability region. The policy consists of two stages: In each time-slot, (a) The base-station first determines the subset to request channel state from via using the *expected* rates over the observable subsets weighted by the *actual* queue lengths at the base-station; and (b) Within the chosen subset, the policy uses the MaxWeight rule [23], [1] which uses the product of the *actual* channel rate (received from the mobile in the chosen subset) and the *actual* queue-length to make the scheduling decision.
- (iii) We develop a simpler online policy (Max-Sum-Queue rule) where in the first stage, the subset of users chosen is determined by only the queue-lengths and does not use the expected channel rates. The Max-Sum-Queue policy chooses that subset over which the sum of the squares of the queue-lengths is largest. The second stage is the same as before, namely, the MaxWeight policy restricted to the chosen subset. We show that if the observable subsets are disjoint or the channels are symmetric, this policy is throughput-optimal. Finally, we provide a example to show that in general this policy is not throughput optimal if the channel-symmetry/disjoint-observable-subsets condition is not met.

B. Related Work

There has been much work in developing scheduling algorithms for down-link wireless systems for various performance metrics that include stability, utility maximization and probabilistic delay guarantees [23], [13], [20], [11], [5], [15], [4], [21]. However, the above studies primarily focus on the case where complete channel state information is available at the base-station.

In the context of partial channel information, related work includes that of [8] where the authors study the problem of a server (terminal) accessing N time varying channels which are independent across users and time (e.g., a multi-channel MAC). The server has a cost for (sequentially) probing channels (with a channel dependent probing cost), and gains a reward (which depends on the user and the probed state) if a packet is transmitted successfully. The authors formulate the

problem of minimizing the expected cost (reward for transmissions minus the probing cost) where the cost functions and the channel probabilities are known to the server. The authors in [8] develop constant factor (within the optimal cost) approximation algorithms that operate in polynomial time for both the saturated data case, as well as when the user (terminal) generates packets according to a Markov chain. The authors in [10], [19] have earlier considered the special cases with equal probing costs and identically distributed channels. Recent results in this context also includes [3] where the authors develop structural properties of the optimal probing strategy using a dynamic programming approach.

Further, for systems with channels that are independent across users and with infinitely backlogged data at the base-station, there have been studies on limited feedback from the mobile users to the base-station. In these studies, the mobiles use thresholds to determine if their channel quality is “good enough”, and if so, send their channel state information to the base-station [6], [17], [18], [22], [16].

II. SYSTEM MODEL AND DEFINITIONS

Consider a time-slotted model of N users serviced by a single server across N unidirectional communication channels $\{c_1, \dots, c_N\} =: C$. An integer number of data packets arrive at the input of every channel at the beginning of a time slot, to be serviced by the server. Packets get queued at the inputs of channels if they are not immediately transmitted. We assume that at most one of the channels can be activated for transmission in a single time slot.

Further, in any given time slot, t , the set of channels C (an N dimensional vector) assumes a *state* $l(t)$ from a finite set of aggregate channel states $\mathcal{L} = \{1, \dots, L\}$, with the channel state remaining constant within each time slot. In each channel state $l \in \mathcal{L}$, every channel $c_i \in C$ assumes a data service rate of μ_i^l , i.e., a maximum of μ_i^l packets can be served from queue i (corresponding to channel c_i) when the aggregate channel is in state l . The random channel state process $\mathbf{L} := (l(t) \in \mathbb{R}^N : t = 0, 1, 2, \dots)$ is assumed to be an independent and identically distributed (*iid*) discrete-time random process taking values from the finite state space \mathcal{L} . We denote the distribution $(\Pr(l(t) = i))_{i=1}^L$ by $\pi = (\pi_1, \dots, \pi_L)$. Observe that the channel state process is *iid* across time, and can have any joint distribution across users (i.e., across channels).

The packet arrival process at the input of each channel c_i is taken to be stationary and ergodic, and generated by a finite state non-negative Discrete Time Markov Chain with rate λ_i .

Our channel observations are limited to a given collection of subsets of the channel $C = \{c_1, \dots, c_N\}$ (whose union is C) called the collection of *observable subsets*. Let us denote the (finite) set of observable subsets of C by $O = \{o_1, o_2, \dots, o_K\}$. In the example of Section IV, the set O contains all subsets of size two. In a given time slot, an observable subset $o_k = \{c_{n_1}, \dots, c_{n_l}\} \subset C$ is said to be in a *sub-state* $\mu^k = (\mu_{n_1}^k, \dots, \mu_{n_l}^k)$ if $\mu_{n_j} = \mu_{n_j}^k$ for $j = 1, \dots, l$.

As in [1], we define the state of the system as the random process $S = (S(t), t = 0, 1, 2, \dots)$ where $S(t) := (Q_1(t), \dots, Q_N(t), U_{11}(t), \dots, U_{1Q_1}(t), \dots, U_{N1}(t), \dots, U_{NQ_N}(t), m(t))$. Here, $Q_j(t)$ denotes the length of the packet queue for channel $c_j \in C$ in time slot t and $U_{ik}(t)$ is the current delay of the k -th packet in queue i at time t . In this regard, a *scheduling policy* \mathcal{P} is a pair of maps $(\mathcal{G}, \mathcal{H})$, where \mathcal{G} is a map from the state of the system $S(t)$ to a fixed probability distribution on the set of observable subsets O , and \mathcal{H} is a map which takes $S(t)$ restricted to a particular observable subset, along with its sub-state, into a fixed probability distribution on the channels which comprise the subset. Such a scheduling policy \mathcal{P} is applied to select a transmitting channel using two steps. At every time slot t , in the first step, we pick an observable set randomly according to the distribution $\mathcal{G}(S(t))$ after which we are able to sample the sub-state of the chosen observable set. Then, using the distribution \mathcal{H} on the observable set and its sub-state, we pick a channel for transmission from that observable set.

A vector or point $\Lambda = (\lambda_1, \dots, \lambda_N) \in \mathbb{R}^N$ is said to be *supported* by a scheduling policy \mathcal{P} if the input packet queues at all channels in the system remain stable under scheduling using \mathcal{P} when the arrival rates at the inputs of channels c_1, \dots, c_N are $\lambda_1, \dots, \lambda_N$ respectively. Associated with each policy \mathcal{P} is its *rate region* $\mathcal{R}(\mathcal{P}) := \{\Lambda \in \mathbb{R}^N : \Lambda \text{ is supported by } \mathcal{P}\}$. The *achievable rate region* \mathcal{R} for the system described above is then defined to be the union of the rate regions for all possible scheduling policies \mathcal{P} . A rate vector Λ is said to be *achievable* if it is supported by some scheduling policy. Likewise, a set or region $A \subset \mathbb{R}^N$ is said to be achievable if all its elements are achievable. A scheduling policy is said to be *throughput-optimal* if it supports all vectors in the achievable rate region.

We wish to characterize the achievable rate region for the model we have described. Henceforth, we shall naturally confine our attention to the set of maximal observable subsets $O_M \subset O$, where the maximality is with respect to set inclusion.

III. THE ACHIEVABLE RATE REGION

In this section, we show two main results. First, we characterize the achievable rate region for any collection of maximal observable subsets O_M . Moreover, we show that this region is attained using a *Static Split Service* (SSS) [1] scheduling rule. The second part of this section characterizes all such maximal SSS scheduling rules.

Consider a maximal observable subset $U \in O_M$, $U = \{c_{k_1}, c_{k_2}, \dots, c_{k_l}\}$ where $k_1, \dots, k_l \in \{1, \dots, N\}$. Let $Q(U)$ denote the l -dimensional subspace of \mathbb{R}^N where coordinates with indices other than k_1, \dots, k_l are zero. If only users from U are served, then any stabilizable rate must lie in $Q(U)$. Denote this stabilizable rate region by $\mathcal{R}(U)$. In particular, applying Theorem 1 in [1] to the subset U , we have:

Lemma 1: There exists a scheduling rule H stabilizing a rate vector $\Lambda = (\lambda_i)_{i=1}^N \in \mathcal{R}(U) \subseteq \mathbb{R}^N$ if and only if there

exists a stochastic matrix ϕ^U such that

$$\lambda_i < v_i^U(\phi^U) := \sum_{m \in \mathcal{L}_U} \pi^{m,U} \phi_{mi}^U \mu_i^{m,U}, \quad \forall i \in U,$$

where \mathcal{L}_U is the set of sub-states of U , $\pi^{m,U}$ is the marginal probability of sub-state m of U and $\mu_i^{m,U}$ is the service rate for channel i in sub-state m .

This matrix ϕ^U defines a *Static Service Split* scheduling rule for the subset U . The rows of ϕ^U correspond to every sub-state of U and the columns of ϕ^U correspond to every channel in U . When U is in the sub-state $m = (\mu_{c_{k_1}}, \dots, \mu_{c_{k_l}})$, channel i is chosen with probability ϕ_{mi}^U .

Thus $\mathcal{R}(U)$ is a convex polytope. We can now characterize the achievable rate region for the system:

Theorem 1: The achievable region, \mathcal{C} , for the whole system is the convex hull of the stabilizable regions in each subspace $Q(\alpha)$, for $\alpha \in O_M$:

$$\mathcal{C} \triangleq \text{conv}(\{\mathcal{R}(\alpha) : \alpha \in O_M\}).$$

The proof follows from two lemmata establishing matching inner and outer bounds on the set \mathcal{C} . Achievability follows from a timesharing argument (see [7] for details):

Lemma 2: \mathcal{C} is achievable.

We next establish that this achievable region is actually tight:

Lemma 3: If $\Lambda \in \mathbb{R}^N$ is achievable, then $\Lambda \in \mathcal{C}$. In particular, Λ can be achieved by a global SSS scheduling rule given by a stochastic matrix ϕ of the form

$$\phi = \sum_{\alpha \in O_M} p_\alpha \phi^\alpha,$$

where ϕ^α are stochastic matrices as described above, and p_α is a probability distribution on the maximal observable subsets, O_M .

Similar to the notion of an SSS rule for a maximal observable subset, the matrix ϕ above defines a *global SSS rule*. A scheduling policy implementing this SSS rule for the system selects a subset α in the first step with probability p_α and subsequently uses the subset SSS rule ϕ^α to pick a queue in α . The (long-term) service rate such a rule provides to queue i is

$$v_i := \sum_{\alpha \in O_M} p_\alpha v_i^\alpha(\phi^\alpha) = \sum_{\alpha \in O_M} p_\alpha \sum_{l \in \mathcal{L}_\alpha} \pi^{l,\alpha} \mu_i^{l,\alpha} \phi_{li}^\alpha.$$

The proof is available in [7]. We note that the assumption that the channel state distribution is *iid* over time, is critically used in the proof of this lemma.

What do maximal SSS rules look like?

We conclude this subsection with a theorem [7] which provides a characterization of maximal global SSS rules. We call a global SSS rule *maximal* if no vector in \mathcal{C} dominates its vector of service rates $(v_i)_{i=1}^N$, where a vector $x \in \mathbb{R}^N$ *dominates* a vector $y \in \mathbb{R}^N$ if $x_i \leq y_i$ for all i , and $x_j < y_j$ holds for at least one j .

Theorem 2: Consider a maximal global SSS rule associated with SSS rules $\{\phi^{*\alpha} : \alpha \in O_M\}$ and a distribution $\{p_\alpha^* : \alpha \in O_M\}$ over subsets. Then, there exists a set of

strictly positive constants $\alpha_i, i = 1, \dots, N$ such that for any l, i and α ,

$$p_\alpha^* > 0, \phi_{li}^{*\alpha} > 0 \Rightarrow i \in \arg \max_{j \in \alpha} \alpha_j \mu_j^{l,\alpha}, \quad \text{and}$$

$$p_\alpha^* > 0 \Rightarrow \alpha \in \arg \max_{\beta \in O_M} \sum_{l \in \mathcal{L}_\beta} \pi^{l,\beta} (\max_{j \in \beta} \alpha_j \mu_j^{l,\beta}).$$

The result says that at time t , in the first scheduling step, a maximal global SSS rule chooses a subset α for which $\sum_{l \in \mathcal{L}_\alpha} \pi^{l,\alpha} (\max_{j \in \alpha} \alpha_j \mu_j^{l,\alpha})$ is maximized, and further picks queue i in α which maximizes $\alpha_i \mu_i^{l(t),\alpha}$, where $l(t)$ is the observed sub-state of subset α .

IV. EXAMPLE: RATE REGION FOR THREE SYMMETRIC CHANNELS

Let us determine the achievable rate region for a three-channel system $C_3 = \{c_1, c_2, c_3\}$ in which the system can take one of eight possible states $\{s_1, \dots, s_8\}$ (Table I), and where each of the channels c_i takes a rate of either a or b ($a < b$) in every state. We denote the 8 possible values of the joint distribution of all three channels by $\pi_1, \pi_2, \dots, \pi_8$ as shown in the table. Assume that we have partial information about subsets of size at most 2, i.e. we know all the joint pairwise probabilities of rates $\{\pi_{ij} : i, j \in \{1, 2, 3\}\}$. Thus the set of maximal posets is $O_M = \{\{1, 2\}, \{2, 3\}, \{3, 1\}\}$. In particular, suppose we know that $\Pr(c_i \text{ has rate } \mu_i, c_j \text{ has rate } \mu_j) = 1/4, i, j \in \{1, 2, 3\}, i \neq j, \mu_i, \mu_j \in \{a, b\}$.

Channel \ State	s_1	s_2	s_3	s_4	s_5	s_6	s_7	s_8
c_1	a	a	a	a	b	b	b	b
c_2	a	a	b	b	a	a	b	b
c_3	a	b	a	b	a	b	a	b
State probability	π_1	π_2	π_3	π_4	π_5	π_6	π_7	π_8

TABLE I
PROBABILITY ASSIGNMENTS FOR THREE-CHANNEL SYSTEM

These pairwise constraints give us a feasible set of possible channel distributions: it is the set of vectors (π_1, \dots, π_8) in the simplex that satisfy the equations $\pi_1 + \pi_2 = 1/4, \pi_1 + \pi_3 = 1/4, \pi_1 + \pi_4 = 1/4, \pi_2 + \pi_5 = 1/4$ etc. In matrix form, these constraints along with the simplex constraints become

$$\begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ & & \vdots & & & & & \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \\ \pi_5 \\ \pi_6 \\ \pi_7 \\ \pi_8 \end{pmatrix} = \begin{pmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \\ \vdots \\ 1/4 \\ 1 \end{pmatrix},$$

with $\pi_i \geq 0$ for all i . The set of possible solutions of the vector $\vec{\pi} = (\pi_1 \ \pi_2 \ \dots \ \pi_8)^T$ is the set of convex combinations of the vectors $\vec{\pi}_{(i)} = (1/4 \ 0 \ 0 \ 1/4 \ 0 \ 1/4 \ 1/4 \ 0)^T$ and $\vec{\pi}_{(ii)} = (0 \ 1/4 \ 1/4 \ 0 \ 1/4 \ 0 \ 0 \ 1/4)^T$, i.e.

$$\vec{\pi} \in \{\eta \vec{\pi}_{(i)} + (1 - \eta) \vec{\pi}_{(ii)} : 0 \leq \eta \leq 1\}.$$

The case where the three channels are independent and identically distributed (*iid*) and $\pi_i = 1/8$, is given by $\eta = 1/2$. Since $a < b$, the ‘‘worst case’’ situation for $\vec{\pi}$ is when $\eta = 0$, i.e., $\vec{\pi} = \vec{\pi}_{(ii)}$, and the best case when $\eta = 1$, and $\vec{\pi} = \vec{\pi}_{(i)}$.

Following the earlier notation, if $\phi = [\phi_{ij}]_{8 \times 3}$ denotes a stochastic matrix defining an SSS rule, where, recall, the (i, j) -th entry is the probability that channel j is chosen for service in system state s_i , then a rate vector $(\lambda_1, \lambda_2, \lambda_3)$ is stabilized by this rule iff:

$$\begin{aligned} \phi_{11}\pi_1 a + \phi_{21}\pi_2 a + \phi_{31}\pi_3 a + \dots + \phi_{81}\pi_8 b &> \lambda_1, \\ \phi_{12}\pi_1 a + \phi_{22}\pi_2 a + \phi_{32}\pi_3 b + \dots + \phi_{82}\pi_8 b &> \lambda_2, \\ \phi_{13}\pi_1 a + \phi_{23}\pi_2 b + \phi_{33}\pi_3 a + \dots + \phi_{83}\pi_8 b &> \lambda_3, \end{aligned}$$

and the stability region for the full-observation case is given by the union of all these regions over all stochastic matrices ϕ . For the case of *iid* channels, i.e., $\pi_i = 1/8$, this region is depicted in Figure 1.

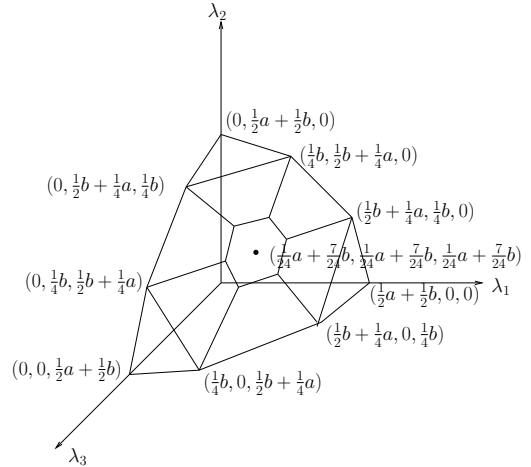


Fig. 1. Rate region for 3 channels with full knowledge of all joint probabilities.

The rate region for any pair of channels c_i and c_j for $i \in \{1, 2, 3\}$ in a 2-dimensional projection of \mathbb{R}^3 turns out to be a convex region in the first quadrant enclosed by four corner points and the origin, as shown in Figure 2.

The dotted line in Figure 2 represents the rate region if we knew only the marginal probabilities of single channels.

According to our result, the achievable rate region for the three-channel system with partial information restricted to size-2 subsets is the region enclosed within the convex hull of the corner points of every pairwise rate region (9 in all), as shown in Figure 3.

Note that this region is a strict subset of the full-information rate region depicted in Figure 1. We also observe that in this example, if we did not know the pairwise channel probabilities of say c_1 and c_2 , we would have to discard the

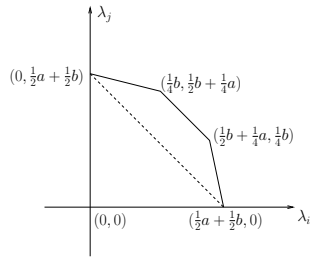


Fig. 2. Rate region for 2 channels c_i and c_j , with knowledge of joint probabilities.

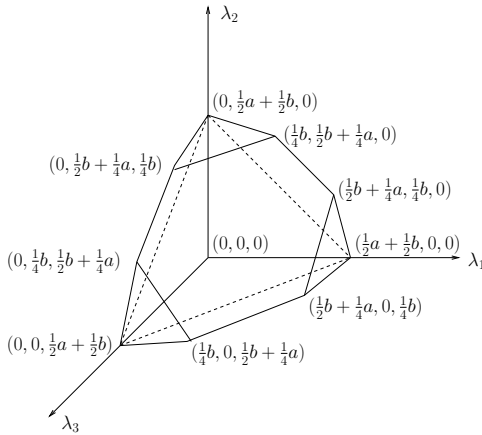


Fig. 3. Achievable rate region for 3 channels, with at most pairwise partial information

corner points $(\frac{1}{4}b, \frac{1}{2}b + \frac{1}{4}a, 0)$ and $(\frac{1}{2}b + \frac{1}{4}a, \frac{1}{4}b, 0)$ when finding the convex hull, hence the achievable region would be smaller (the dotted simplex) than in the case when we know all pairwise probabilities.

V. A THROUGHPUT-OPTIMAL SCHEDULING ALGORITHM

Motivated by the form of the result in Theorem 2, we present a scheduling algorithm which, for a system having arrival rates in the described achievable region, takes as input only the state of the system at each time slot and decides which maximal subset to observe and ultimately, which channel in that subset to schedule. Knowledge of the arrival rates is not assumed in such a case. However, it is presumed that the marginal probabilities $\pi^{l,\alpha}$ of the subset α being in the sub-state l (as in the proof of Lemma 3) are known.

Algorithm 1:

- 1) Select a poset $\delta \in O_M$, given by

$$\delta \in \arg \max_{\alpha \in O_M} \sum_{l \in \mathcal{L}_\alpha} \pi^{l,\alpha} \left(\max_{i \in \alpha} Q_i(t) \mu_i^{l,\alpha} \right),$$

where the symbols O_M , \mathcal{L}_α and $\mu_i^{l,\alpha}$ have the same meaning as in the proof of Lemma 3 and $Q_i(t)$ represents the length of the i th queue at the beginning of time slot t .

- 2) After observing the state $s \in \mathcal{L}_\delta$ of δ , schedule channel $j \in \delta$ using the *max-weighted-queue* rule, i.e.

$$j \in \arg \max_{i \in \delta} Q_i(t) \mu_i^{s,\delta}.$$

Remark: By arguments used to prove Theorem 2 [7], we have the following lemma which provides an important equivalent characterization of the above algorithm in terms of knowing the extreme points of the achievable rate region \mathcal{C} :

Lemma 4: Let E be the (finite) set of extreme points for the achievable rate region \mathcal{C} . Then, Step 1 above can be expressed as choosing $\delta \in O_M$ which satisfies

$$\delta \in \{\alpha \in O_M : u \in \arg \max_{v \in E} \langle v, Q(t) \rangle \Rightarrow u \in \mathcal{R}(\alpha)\}.$$

That is, the algorithm selects any subset whose rate region contains an extreme point maximizing the inner product $\langle x, Q(t) \rangle$ over all $x \in E$ and hence a point maximizing $\langle y, Q(t) \rangle$ over all $y \in \mathcal{C}$.

The chief result in this section is the following theorem, which says that the scheduling policy defined above is throughput-optimal for scheduling with partial channel-state information.

Theorem 3: Algorithm 1 makes the system stable if the vector of arrival rates lies in the achievable region.

The proof of stability uses fluid limit machinery. Roughly, the fluid limit of the (N dimensional) queue length corresponds to a limiting trajectory when the queue length process is “observed” over a long interval of time (by scaling and “compressing” time) and concurrently scaling down the magnitude of the queue length process. Under such a scaling, the discrete and random queue length process “looks like” a deterministic fluid process henceforth denoted by $q(t)$, which is driven by a (vector) constant rate fluid arrival process (the components corresponding to the mean arrival rates to each of the users), and whose service rate corresponds to the “average” service rate under the scheduling algorithm. For the system we are considering, showing that such a limiting fluid queue length trajectory has negative drift (as we will do so in Lemma 5) is sufficient to prove that the discrete-time stochastic queue length process is stable (positive recurrent) [14], [1], [7].

The proof requires numerous definitions and lemmata on fluid limits, their existence and their properties. We defer these to the full version ([7]) and provide here only the main intuitive Lyapunov idea.

The key step in the proof is to show that a suitably defined Lyapunov function has negative drift. So far this parallels the proof used for Theorem 3 of [1], however here we face the additional difficulty of assuring that we pick the correct observation subset $\alpha \in O_M$, in addition to picking the correct queue to serve inside that subset α . Using the Lyapunov function introduced below, we show that maximizing the negative drift of this Lyapunov function is exactly the problem of maximizing the inner product $\langle y, Q(t) \rangle$ over all $y \in \mathcal{C}$. If we pick the “wrong” subset $\alpha \in O_M$, then maximizing the linear function above becomes impossible. To side-step this problem, we rely on Lemma 4,

which guarantees that the chosen subset will indeed be one with an extreme point maximizing the linear function.

Formally, let us introduce the quadratic Lyapunov function

$$L(y) = \frac{1}{2} \sum_{i=1}^N y_i^2 \quad (1)$$

for a vector $y = (y_1, \dots, y_N)$. Let $q(t)$ denote a fluid limit of the queue-length process (this exists almost surely; see [7] for precise definitions and details). The following property establishes negative drift, and thus enables us to show stability:

Lemma 5: Consider a feasible system operating under the described scheduling discipline. For any $\delta_1 > 0$, there exists $\delta_2 > 0$ such that the following holds. With probability 1, a limiting set of functions defining the fluid limit, satisfies the following property at any regular point t :

$$L(q(t)) \geq \delta_1 \Rightarrow \frac{d}{dt} L(q(t)) \leq -\delta_2 < 0.$$

The proof relies on Lemma 4. See [7] for the full details.

As in [1], the previous lemma along with a result from [14] together imply Theorem 3.

VI. THE MAX-SUM-QUEUE ALGORITHM

In this section, we present a ‘simpler’ scheduling policy which only uses queue-length information to pick the subset to observe, and analyze its stability properties under specific assumptions:

Algorithm 2:

In each time slot t ,

- 1) Select a poset $\delta \in O_M$, given by

$$\delta = \arg \max_{\alpha \in O_M} \sum_{i \in \alpha} Q_i^2(t),$$

where $Q_i(t)$ denotes the length of the i th queue at the beginning of time slot t .

- 2) After observing the state $s \in \mathcal{L}_\delta$ of δ , schedule channel $j \in \delta$ using the *MaxWeight* rule (also known as the *Modified Largest-Weighted-Work-First (M-LWWF)* rule) [23], [1], i.e.

$$j = \arg \max_{i \in \delta} Q_i(t) \mu_i^{s, \delta}.$$

Note: A suitable rule to break ties in each case is assumed.

We shall call this algorithm the *Max-Sum-Queue* scheduling algorithm. In this section we show that it is throughput-optimal in two cases of interest: (i) when the maximal subsets in O_M are disjoint; and (ii) when the channel is symmetric in the users. In the next section, we prove by example that throughput-optimality does not hold in general.

A. Max-Sum-Queue for disjoint subsets

The following result tells us that when the collection of observable subsets is mutually disjoint, Max-Sum-Queue is throughput-optimal.

Theorem 4: Under the assumption that every pair of maximal observable subsets is disjoint, the Max-Sum-Queue

scheduling algorithm makes the system stable if the vector of arrival rates lies in the achievable region.

To prove Theorem 4, we follow a similar route as in the previous section, defining fluid limits and proving that a suitably defined Lyapunov function has negative drift. The Lyapunov function we use here is

$$L(y) = h_\beta(y) = \frac{1}{2} \sum_{i \in \beta} y_i^2,$$

where

$$\beta \in \arg \max_{\alpha \in O_M} h_\alpha(y)$$

and β is chosen according to some fixed precedence rule in $\arg \max_{\alpha \in O_M} h_\alpha(y)$, for a vector $y = (y_1, \dots, y_N)$. As with Lemma 5, the following lemma is used to establish stability.

Lemma 6: Consider a feasible system operating under the Max-Sum-Queue scheduling discipline. For any $\delta_1 > 0$, there exists $\delta_2 > 0$ such that the following holds. With probability 1, a limiting set of functions defining the fluid limit, satisfies the following additional property at any regular point t :

$$L(q(t)) \geq \delta_1 \Rightarrow \frac{d}{dt} L(q(t)) \leq -\delta_2.$$

We refer the reader to [7] for the details. There is an intuitive geometric explanation for this result. It is based on two observations: first, due to the disjoint subset assumption and the Max-Sum-Queue algorithm, if any queue is unstable, all queues are unstable; next, given an extreme point x_α in each set $\mathcal{R}(\alpha)$, the convex hull of those extreme points will always lie on an exposed face of \mathcal{C} . Note that this is not true in the general case.

B. Max-Sum-Queue for symmetric channels

It is instructive to note that the reason that the presented scheduling policies work in their respective cases is because at any point $t \in [0, \infty)$, they maximize the linear objective function $\langle q(t), u \rangle$ over all u in the convex polytope \mathcal{C} which represents the achievable rate region. The drift of the sum-of-squares Lyapunov function defined by (1) happens to be precisely the difference between $\langle q(t), \lambda \rangle$ and $\max_{u \in \mathcal{C}} \langle q(t), u \rangle$. This geometric interpretation allows us to prove the useful result that Max-Sum-Queue is actually throughput-optimal for systems of symmetric channels and subsets.

Theorem 5: Consider a symmetric system, i.e. where all the N channels have an identical distribution of service rates. Further, let the observable subsets be all subsets of a fixed cardinality K . For such a system, Max-Sum-Queue is throughput-optimal.

See [7] for the proof.

A geometric view of Max-Sum-Queue: In the following, we see why Max-Sum-Queue is throughput-optimal by examining the geometric aspect of its working. Let λ be the vector of arrival rates to the system of N channels represented by $S = \{1, \dots, N\}$, such that $\lambda \in \text{int}C$. As before, we consider the drift of the sum-of-squares Lyapunov function defined by

(1):

$$\begin{aligned} \frac{d}{dt} L(q(t)) &= \sum_{i=1}^N q_i(t) (\lambda_i - \hat{f}_i(t)) \\ &= \langle q(t), \lambda \rangle - \langle q(t), \hat{f}(t) \rangle, \end{aligned}$$

where $\hat{f}(t) \equiv (\hat{f}_i(t))_{i=1}^N$ is the instantaneous vector of service rates chosen by Max-Sum-Queue at time t in the fluid time scale. We will show that $\hat{f}(t) \in \mathcal{C}$ maximizes the inner product $\langle q(t), x \rangle$ over all $x \in \mathcal{C}$ or equivalently over all the extreme points of \mathcal{C} ; this establishes that the drift of $L(q(t))$ is strictly negative and bounded away from zero and hence Max-Sum-Queue is throughput-optimal.

We observe that the subsets which Max-Sum-Queue picks for scheduling at t are the ones of (fixed) size $K < N$, say, that contain the top K queues in the system. Without loss of generality, let $q_1(t) \geq q_2(t) \geq \dots \geq q_N(t)$, and let

$$A = \arg \max_{\beta \subset S, |\beta|=K} \sum_{i \in \beta} q_i^2(t).$$

Every set $\alpha \in A$ is picked by Max-Sum-Queue in the fluid timescale, and has the same queue values ordered in descending order. Further, since the channels are symmetric, every subset rate region $\mathcal{R}(\beta)$ for $\beta \subset S$, $|\beta| = K$, is identical up to a permutation of indices. It follows that the extreme points of \mathcal{C} maximizing $\langle \cdot, q(t) \rangle$ must lie in the rate regions $\mathcal{R}(\alpha)$ where $\alpha \in A$, since only the K heaviest queues can maximize this inner product over all permutations of extreme points.

Since these extreme points are precisely the ones picked by Max-Sum-Queue in each subset, and that $\hat{f}(t)$ lies in the convex hull of these extreme points, $\hat{f}(t)$ maximizes the inner product $\langle q(t), x \rangle$ over all $x \in \mathcal{C}$, and we are done.

VII. MAX-SUM-QUEUE APPLIED TO ARBITRARY SUBSETS

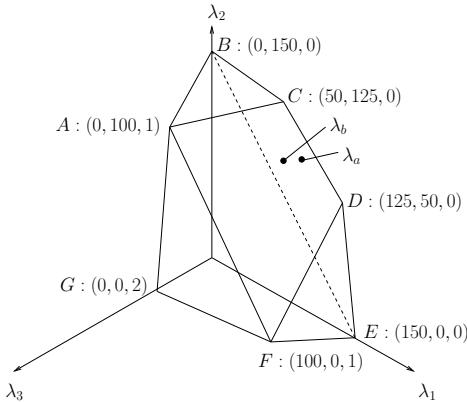


Fig. 4. Rate region for described 3-channel system

It is interesting to ask the question: Is the simple Max-Sum-Queue scheduling algorithm throughput-optimal for an arbitrary (and, in particular, non-disjoint) system of maximal observable subsets? In this section, we present an example

based on a system of three channels where under certain arrival rates in the stability region, all the queue fluid limits are seen to increase. Thus we can show [7] that the system exhibits instability in certain portions of the achievable rate region, under this policy.

Consider a system of three channels c_1 , c_2 and c_3 . The system assumes four possible states S_1 , S_2 , S_3 and S_4 with the corresponding channel rates, expressed by (rate of c_1 , rate of c_2 , rate of c_3), being (100,100,2), (100,200,2), (200,100,2) and (200,200,2) respectively. Further, each state occurs with probability $\frac{1}{4}$. The maximal observable subsets are $\alpha = \{c_1, c_2\}$, $\beta = \{c_2, c_3\}$ and $\gamma = \{c_3, c_1\}$, i.e., all pairs of channels. The achievable rate region for the system is shown in Figure 4.

Set the vector of arrival rates (shown in the figure) to be $\lambda_b \equiv (\lambda_{1b}, \lambda_{2b}, \lambda_{3b}) = (\frac{175}{2}, \frac{175}{2}, 0) - \epsilon(1, 1, 0) + \delta(0, 0, 1)$, with $\epsilon = \frac{1}{2}$ and $0 < \delta = \frac{1}{100} < \frac{1}{75}$. It is easily verifiable that λ lies in the interior of the rate region. We will show that a regular point $t \in [0, \infty)$ can exist with the fluid-limit queue-length process satisfying $q_1(t) = q_2(t) = q_3(t) > 0$, and with $\dot{q}_1(t) = \dot{q}_2(t) = \dot{q}_3(t) > 0$ (the full details are in [7]). In such a case, the fluid levels of the queues increase (linearly) at a constant rate.

Let us hypothesize that t is a regular point in $[0, \infty)$ satisfying the condition $q_1(t) = q_2(t) = q_3(t) > 0$, and attempt to find a valid set of the subset timesharing probabilities p_α , p_β and p_γ . Note that all the $q_i(t)$ being equal forces the system to be ‘serving’ all three subsets with the aforementioned timesharing probabilities, which must be positive. We can show ([7]) that the regularity hypothesis now implies $\dot{q}_1(t) = \dot{q}_2(t) = \dot{q}_3(t)$, and hence

$$\begin{aligned} \Rightarrow \lambda_{1b} - 150p_\gamma - \frac{175}{2}p_\alpha & \\ &= \lambda_{2b} - 150p_\beta - \frac{175}{2}p_\alpha \\ &= \lambda_{3b} - 0 = \delta \\ \Rightarrow p_\gamma = p_\beta, \quad \text{and} & \\ 150p_\beta + \frac{175}{2}p_\alpha = \lambda_{2b} - \delta = 86.99. & \end{aligned}$$

Together with $p_\alpha + p_\beta + p_\gamma = 1$, we get $p_\beta = p_\gamma \approx 0.02$ and $p_\alpha \approx 0.96$ which is the unique timesharing solution between the subsets α , β and γ . Hence t is valid as a regular point where all the queues are equal and increase linearly at the same rate $\delta > 0$.

Remarks:

- 1) We observe that the (mutually exclusive) conditions $q_1(t) = q_2(t) > q_3(t)$ and $q_1(t) = q_2(t) < q_3(t)$ lead to all the q_i becoming equal within finite time. Hence the state $q_1(t) = q_2(t) = q_3(t)$ is an ‘unstable attractor’ for the fluid limits in this sense.
- 2) For the arrival rate vector $\lambda_a = (87, 87, 0)$ (as in Figure 4), we can similarly show that starting from $q_1(0) = q_2(0) = q_3(0) = c > 0$ implies that $q_1(t) = q_2(t) = q_3(t) = c$ at all times $t \in [0, \infty)$.

The following proposition formalizes the case of linearly exploding fluid limits for the system defined in Section VII

when the vector of arrival rates is $\lambda = \lambda_b$ (as in Figure 4). Specifically, we show that given any large time interval, we can find a large enough initial condition for the queue lengths such that the queue lengths grow linearly within that time interval. We refer the reader to [7] for the proof.

Proposition 6: Fix any $T > 0$. Then, there exists $n_0(T)$ such that for all $n > n_0$, there exists $\epsilon > 0$ such that

$$Q_1(0) = Q_2(0) = Q_3(0) = n \\ \Rightarrow Q_i(nT) \geq (1 + \epsilon)n, \quad i = 1, 2, 3.$$

REFERENCES

- [1] M. Andrews, K. Kumaran, K. Ramanan, A. L. Stolyar, R. Vijayakumar, and P. Whiting. Scheduling in a queueing system with asynchronously varying service rates. *Probability in Engineering and Informational Sciences*, 14:191–217, 2004.
- [2] R. Berry and R. Gallager. Communication over fading channels with delay constraints. *IEEE Trans. Info. Theory*, 48:1135–1149, May 2002.
- [3] N. Chang and M. Liu. Optimal channel probing and transmission scheduling for opportunistic spectrum access. In *Proc. ACM International Conference on Mobile Computing and Networking (MobiCom)*, Montreal, Canada, September 2007.
- [4] Mung Chiang. Balancing transport and physical layers in wireless multihop networks: jointly optimal congestion control and power control. *IEEE Journal on Sel. Areas in Commun.*, 23:104–116, January 2005.
- [5] A. Eryilmaz, R. Srikant, and J. R. Perkins. Throughput optimal scheduling for broadcast channels. In *Modelling and Design of Wireless Networks, Proceedings of SPIE, E. K. P. Chong (editor)*, volume 4531, pages 70–78, Denver, CO, 2001.
- [6] D. Gesbert and M. Slim-Alouini. How much feedback is multi-user diversity really worth? In *Proc. Int. Conf. on Commun.*, pages 234–238, 2004.
- [7] Aditya Gopalan, Constantine Caramanis, and Sanjay Shakkottai. On wireless scheduling with partial channel-state information. Technical report, The University of Texas at Austin, 2007.
- [8] Sudipto Guha, Kamesh Munagala, and Saswati Sarkar. Performance guarantees through partial information based control in multichannel wireless networks. Technical report, University of Pennsylvania, 2006. <http://www.seas.upenn.edu/~swati/report.pdf>.
- [9] Arunabha Ghosh Jeffrey G. Andrews and Rias Muhamed. *Fundamentals of WiMAX: Understanding Broadband Wireless Networking*. Prentice Hall, 2007.
- [10] Z. Ji, Y. Yang, J. Zhou, M. Takai, and R. Bagrodia. Exploiting medium access diversity in rate adaptive wireless LANs. In *Proc. ACM MOBICOM*, 2004.
- [11] X. Lin, N. Shroff, and R. Srikant. A tutorial on cross-layer optimization in wireless networks. *IEEE Journal on Selected Areas in Communications*, pages 1452–1463, August 2006.
- [12] X. Liu, E. K. P. Chong, and N. B. Shroff. A framework for opportunistic scheduling in wireless networks. *Computer Networks Journal*, 41(4), 2003.
- [13] X. Liu, E. P. K. Chong, and N. Shroff. Opportunistic transmission scheduling with resource-sharing constraints in wireless networks. *IEEE Journal on Sel. Areas of Commun.*, 19(10):2053–2064, October 2001.
- [14] V. A. Malyshev and M. V. Menshikov. Ergodicity, continuity and analyticity of countable Markov chains. *Transactions of the Moscow Mathematical Society*, 39:3–48, 1979.
- [15] M. Neely, E. Modiano, and C. Rohrs. Power allocation and routing in multi-beam satellites with time varying channels. *IEEE/ACM Trans. Networking*, 11(1):138–152, February 2003.
- [16] S. Patil and G. de Veciana. Reducing feedback for opportunistic scheduling in wireless systems. *IEEE Trans. on Wireless Comm.*, January 2006. Submitted.
- [17] X. Qin and R. Berry. Opportunistic splitting algorithms for wireless networks. In *Proc. IEEE INFOCOM*, March 2004.
- [18] X. Qin and R. Berry. Opportunistic splitting algorithms for wireless networks with heterogenous users. In *Proc. Conf. on Information Sciences and Systems (CISS)*, March 2004.
- [19] A. Sabharwal, A. Khoshnevis, and E. Knightly. Opportunistic spectral usage: Bounds and a multi-band CSMA/CA protocol. *IEEE/ACM Transactions on Networking*, 2006.
- [20] S. Shakkottai, R. Srikant, and A. L. Stolyar. Pathwise optimality of the exponential scheduling rule for wireless channels. *Advances in Applied Probability*, 36(4):1021–1045, 2004.
- [21] A.L. Stolyar. Maxweight scheduling in a generalized switch: State space collapse and workload minimization in heavy traffic. *Annals of Applied Probability*, 14, No.1:1–53, 2004.
- [22] T. Tang and R. W. Heath. Opportunistic feedback for downlink multiuser diversity. *IEEE Communication Letters*, 9:948–950, 2005.
- [23] L. Tassiulas and A. Ephremides. Dynamic server allocation to parallel queues with randomly varying connectivity. *IEEE Transactions on Information Theory*, 39:466–478, March 1993.