

User Scheduling in Pilot-Assisted TDD Multiuser MIMO Systems

Zhiyuan Jiang, *Student member, IEEE*, Sheng Zhou, *Member, IEEE*, and Zhisheng Niu, *Fellow, IEEE*
 Tsinghua National Laboratory for Information Science and Technology,
 Department of Electronic Engineering, Tsinghua University, Beijing, 100084, China
 jiang-zy10@mails.tsinghua.edu.cn, {sheng.zhou, niuzhs}@tsinghua.edu.cn

Abstract—User scheduling in multiuser multiple-input-multiple-output (MU-MIMO) systems is fundamentally different with single-user systems¹, in the sense that without spatial multiplexing, users in single-user systems are sharing the time-frequency degree-of-freedom (DoFs), whereas in MU-MIMO systems, due to the fact that the number of spatial DoFs scales with the number of users (assuming sufficient base station (BS) antennas), users are not sharing the DoFs, but rather creating additional DoFs for their own use. However, instead of limited by the available DoFs, the number of simultaneous users are limited by the channel state information (CSI) acquisition overhead in pilot-assisted MU-MIMO systems. In this paper, we investigate the user scheduling scheme in pilot-assisted time-division-duplex (TDD) MU-MIMO systems. Leveraging the Lyapunov optimization techniques, we derive the throughput-optimal scheduling policy which serves as a performance bound due to its non-causality and high complexity. We then propose a heuristic scheme, which is causal and substantially decreases the complexity. Moreover, it performs fairly close to the optimum.

I. INTRODUCTION

Multiuser multiple-input-multiple-output (MU-MIMO) technology enables simultaneous (on the same time-frequency resource) data transmissions to a multiplicity of user terminals via distinguishable spatial modes. With perfect channel state information at transmitter (CSIT) and at receiver (CSIR), the capacity of the MU-MIMO system is significantly larger than that of single-user systems [1].

Despite the vital importance of CSIT, existing literature usually assumes that the CSIT overhead is negligible, and thus the MU-MIMO system can accommodate a large number of simultaneous users, especially in the massive MIMO system, where a vast excess number of base station (BS) antennas are deployed and the time-division-duplex (TDD) mode is leveraged to explore the channel reciprocity [2]. However, even with the TDD mode, the CSIT acquisition process can entail significant overhead when a large number of users are scheduled concurrently in a limited-coherence channel. The problem has been scantily treated in the literature, which motivates our work.

In this paper, we investigate the design of user scheduling scheme, also referred to as dynamic channel acquisition (DCA) scheme, to maximize the throughput of the TDD MU-MIMO

system downlink, on account of the user queue information (UQI) and the channel acquisition overhead. A unique issue that we address is that users have distinct channel coherence times². The major contribution is that based on the Lyapunov-drift optimization, we formulate the generic user-scheduling problem as the genie-aided optimization problem (GAP). The corresponding user scheduling scheme, referred to as the GAP-rule, is proved to be throughput-optimal, i.e., it stabilizes the system as long as the arrival rates are inside the capacity region. In view of the fact that the GAP-rule is practically infeasible due to its complexity and non-causality, we propose a heuristic scheme, namely the queue-based quantized-block-length user scheduling scheme (QQS), which substantially reduces the complexity and also is practically feasible. It is shown by simulations that the QQS is asymptotically throughput-optimal under the conditions that the system dimension is large and the user coherence times can be approximately grouped such that in each group the user coherence times are identical.

The paper is organized as follows. Section II describes the system model and gives some preliminary knowledge. In Section III, the GAP is formulated to maximize the Lyapunov-drift. In Section IV, we propose the QQS. Section V gives the numerical results. Finally, in Section VI, we draw the conclusions. Throughout the paper, we use boldface uppercase letters, boldface lowercase letters and lowercase letters to designate matrices, column vectors and scalars respectively.

II. SYSTEM MODEL AND PRELIMINARIES

We consider the downlink (forward-link) of a single cell where an M -antenna BS serves N single-antenna users, and one channel use³ is characterized as

$$\mathbf{y}(t) = \mathbf{H}(t)\mathbf{x}(t) + \mathbf{z}(t), \quad (1)$$

where $\mathbf{x}(t) \in \mathbb{C}^M$ denotes the complex transmit signal vector of M antennas at the BS, t in the bracket denotes the index of channel use, $\mathbf{y}(t) \in \mathbb{C}^N$ denotes the receive signal vector of N single-antenna users, $\mathbf{z}(t)$ denotes the cyclic symmetric zero mean complex Gaussian additive noise, i.e., $\mathbf{z}(t) \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$, and $\mathbf{H}(t) \in \mathbb{C}^{N \times M}$ denotes identically independently distributed (i.i.d.) Rayleigh fading

¹Single-user systems in this paper refer to systems wherein each time-frequency resource is dedicated to one user, e.g., in 3G systems, users share system resources by using different code-sequences, and single-user MIMO systems.

²This is due to different user mobilities and scattering environments.

³A channel use, or a time slot, corresponds to an independent complex signal-space dimension in the time-frequency domain.

coefficients with unit-norm entries. In particular, we consider linear precoding, where⁴

$$\mathbf{x}(t) = \zeta(t)\mathbf{W}(t)\mathbf{s}(t), \quad (2)$$

where $\zeta(t)$ is the power normalization factor with $\zeta(t)^2 = \frac{P}{\text{tr}(\mathbf{W}(t)\mathbf{W}(t)^\dagger)}$, P is the total transmit power, $\mathbf{W}(t)$ is the precoding matrix, and $\mathbf{s}(t)$ denotes the i.i.d. user data streams. The signal-to-interference-noise ratio (SINR) of user- n is written as,

$$\gamma_n(t) = \frac{\zeta(t)^2 |\mathbf{h}_n^\dagger(t)\mathbf{w}_n(t)|^2}{\sum_{j \neq n} \zeta(t)^2 |\mathbf{h}_n^\dagger(t)\mathbf{w}_j(t)|^2 + z_n(t)^2}, \quad (3)$$

where we write $\mathbf{H}(t) = [\mathbf{h}_1(t), \mathbf{h}_2(t), \dots, \mathbf{h}_N(t)]^\dagger$ and $\mathbf{W}(t) = [\mathbf{w}_1(t), \mathbf{w}_2(t), \dots, \mathbf{w}_N(t)]$. Furthermore, let $Q_n(t)$ denote the queue length in bits of user n at the beginning of t -th channel use, let $a_n(t)$ denote the number of arrival bits from upper layer to the physical layer between $(t-1)$ -th and t -th channel uses, and let $\mu_n(t)$ denote the allocated number of service bits to Queue- n , which equals the allocated service rate (bits/channel use) in this case. Then the queuing dynamics are written as

$$Q_n(t+1) = Q_n(t) - \tilde{\mu}_n(t) + a_n(t), \quad (4)$$

where $\tilde{\mu}_n(t) = \min\{Q_n(t), \mu_n(t)\}$ denotes the actual service number of bits, considering the circumstances that sometimes the queue is emptied given the amount of allocated service bits.

Definition 1: Queue- n is said to be strongly stable if [3]

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[Q_n(t)] < \infty, \quad (5)$$

when there is no bound on the buffer size for any n .

When all queues are strongly stable in the system, the time-average actual service rate equals the arrival rate, i.e.,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \tilde{\mu}_n(t) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T a_n(t), \quad \forall n. \quad (6)$$

Notice that the left-right-side is the time average of the realizations of the actual service rate, thus we do not need the expectation to hold (6).

The achievable ergodic rate region \mathcal{R} is defined as the convex hull of all achievable rate points of n users. Denote all the feasible transmission schemes as \mathcal{X} . Under a fixed transmission scheme⁵ $\pi_s \in \mathcal{X}$, the user achievable rate point is defined as the time-average of user rates

$$\bar{R}_n = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T R_n(\mathbf{H}(t), \pi_s(t)), \quad \forall n, \quad (7)$$

⁴Notice that $\mathbf{W}(t)$ can be any general linear precoding matrix, whereas we adopt the zero-forcing precoding matrix, i.e., $\mathbf{W}(t) = \mathbf{H}(t)^\dagger (\mathbf{H}(t)\mathbf{H}(t)^\dagger)^{-1}$ in the simulations.

⁵A transmission scheme here means a realization of resource allocation, or user scheduling scheme, with a fixed precoding scheme which is zero-forcing in the simulations in Section V.

where $R_n(\mathbf{H}(t), \pi_s(t))$ is defined as the transmission rate of user- n with channel realization $\mathbf{H}(t)$ and user scheduling decision $\pi_s(t)$. Based on the ergodicity, (7) equals

$$\bar{R}_n = \mathbb{E}\{R_n(\mathbf{H}, \pi_s)\}, \quad \forall n, \quad (8)$$

where the expectation is taken over all possible channel gain $\mathbf{H}(t)$ and possibly $\pi_s(t)$ ⁶. The achievable ergodic rate region can be characterized as

$$\mathcal{R} = \text{coh} \bigcup_{\pi_s \in \mathcal{X}} \{\bar{\mathbf{R}} : 0 \leq \bar{R}_n \leq \mathbb{E}[R_n(\mathbf{H}, \pi_s)]\}, \quad (9)$$

where $\bar{\mathbf{R}}$ is a N -dimensional region, \bar{R}_n is its n -th component, and ‘‘coh’’ denotes ‘‘the closure of the convex hull’’.

Definition 2: (Throughput-Optimality) A scheduling scheme is throughput-optimal if for any arrival rate point inside the achievable ergodic rate region, the system can be stabilized by the scheduling scheme.

Note that the throughput in this paper refers to the downlink throughput, not concerning uplink throughput. We consider a generic scenario where each user has its distinct block length, which denotes the number of consecutive channel uses that the user-channel stays static, or referred to as channel coherence time for brevity. The block fading channel model is adopted in this paper, where every user’s channel stays constant for T_n consecutive channel uses, and changes to another constant according to an i.i.d. (over time and users) random process. Denote by T_n the channel coherence time of user- n , and let $\mathcal{T} = \{T_1, T_2, \dots, T_N\}$ ⁷.

In this work, we assume that for every T_n time slots, the BS uses the uplink channel training symbol (TDD mode is considered here) to estimate the channel of user- n . By doing this, we assume the BS obtains the *perfect CSIT* of user- n . Note that this is actually an optimistic assumption on the MU-MIMO system, since normally we can only get a noisy version of the CSIT, and the system has to do the channel training more than once to obtain a more accurate estimation. Nonetheless our results can be extended to this scenario directly by multiplying the training length with a predefined factor, taking into account the imperfection of channel estimation [1].

III. GENIE-AIDED DYNAMIC CHANNEL ACQUISITION

In this section, we first formulate the generic DCA optimization problem, the GAP, which maximizes the Lyapunov-drift every scheduling step with the aid of a genie who provides the BS the instantaneous channel coefficients before channel estimation. The resulting scheduling scheme, albeit practically infeasible, is termed as the *GAP-rule*. Due to the throughput-optimality, the GAP-rule serves as a performance bound in this paper. In the next section, we will propose a heuristic

⁶When a randomized control policy is considered.

⁷The distinction of user block lengths is due to the fact that there are several factors that can affect the block length of each user, such as distinct user-mobility, scattering environment nearby, frequency offset and etc. We assume the BS knows the channel coherence time *a priori*, since channel coherence time is a second-order channel statistics, which can be regarded to be static for a relatively long period.

algorithm, which is the practically feasible version of the GAP-rule and also shows near-optimal performance. Note that we assume the coherence times of all users are known to the BS, since the channel coherence time is usually changing slowly, about seconds to tens of seconds, and thus it can be estimated efficiently.

A. GAP

The generic user scheduling optimization problem, namely GAP, is formulated based on the framework of [3]. To stabilize the system whenever the arrival rate is inside the achievable rate region, the optimization boils down to select the users which optimize the Lyapunov drift in each scheduling step, i.e.,

$$\underset{S \subseteq \mathcal{N}}{\text{maximize}} \quad \left[\sum_{n \in S} \frac{Q_n(t_k) \beta_n(t_k)}{T_k} \right], \quad (10)$$

where S is the optimization variable which denotes the set of scheduled users, \mathcal{N} is the overall user set.

$$\beta_n(t_k) = \begin{cases} (T_k - |S|) r_n^{\text{SM}}(t_k) & \text{if } |S| > 1, \\ T_k r_n^{\text{STC}}(t_k) & \text{if } |S| = 1, \end{cases} \quad (11)$$

where $\beta_n(t_k)$ denotes the allocated service bits of user n at time t_k , and

$$r_n^{\text{SM}}(t_k) = \log \left(1 + \gamma_n^{(t_k)} \right), \quad (12)$$

$$r_n^{\text{STC}}(t_k) = \log \left(1 + \frac{\|\mathbf{h}_n(t_k)\|^2 P}{M \sigma^2} \right), \quad (13)$$

and

$$T_k = \begin{cases} \min_{n \in S} [T_n], & \text{if } |S| > 1, \\ T_{\text{STC}}, & \text{if } |S| = 1. \end{cases} \quad (14)$$

The objective in (10) can be seen as the queue-size-weighted sum of the user service rates. Since we assume the users each occupies one uplink training channel use to obtain a perfect CSIT, the time/frequency resource dedicated to channel estimation is the number of simultaneous transmission users and the remaining resource is $T_k - |S|$ in (11). T_{STC} is a predefined constant. When the number of selected users is larger than one, spatial multiplexing is enabled with user rate $r_n^{\text{SM}}(t_k)$ and channel estimation overhead $|S|$. Otherwise, STC is leveraged to serve one user at a time without channel estimation overhead with rate $r_n^{\text{STC}}(t_k)$.

It is clear that the rates in (12) and (13) *cannot* be evaluated to proceed the optimization in practice unless we have a genie who provides the BS all the channel coefficients without having to do the channel estimation. Most existing literature assumes the CSI is known *a priori* [4] without considering the acquisition overhead, or coarse knowledge of CSI is available [5], neither of which is practical when the channel coherence time is short. Even supposing the genie is available, the algorithm is still NP-hard under generic linear precoding, meaning that we have to exhaustively search all the user sets to obtain the optimum. Therefore, in the next section, we will propose heuristic algorithms, which are much more practical and meanwhile with little to none performance degradation.

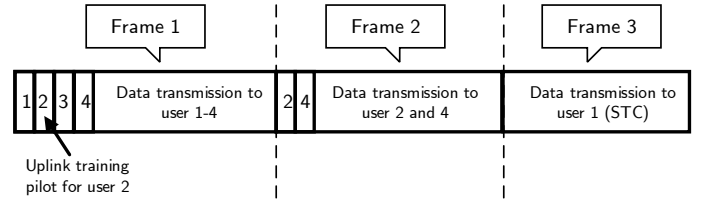


Fig. 1. A sample path of the control scheme. Only three frame transmission is shown for simplicity.

Notice that the scheduling scheme adopts a variable frame-length design. A sample path of the control scheme is shown in Fig. 1. The scheduling scheme decides which set of users to serve. If multiple users are chosen, the chosen users have to send training pilots first to let the BS have the CSIT. The frame length when multiple users are chosen is set to be the *minimum channel coherence time* of the selected users. In this way, during one frame, the channel estimations of all scheduled users are meaningful. On the other hand, if only one user is chosen, the BS will use the STC scheme, with *no channel estimation* needed. Note that the assumption of the frame length when multiple users are selected renders a lower bound of the achievable rates since some users may have remaining channel coherence times. However, this assumption makes the Lyapunov drift analysis tractable. Otherwise, the decisions of different scheduling stages would be dependent due to the possible remaining channel coherence times of some users, which makes the analysis much more difficult.

Specifically, the frame-by-frame queueing dynamics are written as

$$Q_n(t_{k+1}) = \max[Q_n(t_k) - \beta_n(t_k), 0] + \alpha_n(t_k), \quad \forall n \quad (15)$$

where t_k denotes the beginning of frame- k , $\beta_n(t_k)$ and $\alpha_n(t_k)$ denote the allocated service bits and arrival bits during the time interval $[t_k, t_{k+1})$, respectively.

Theorem 1: (Throughput-Optimality of the GAP-Rule) Suppose $a_n(t_k)$ is i.i.d. over time and satisfies

$$0 \leq a_n(t_k) \leq A_{\max}, \quad \forall n, k \quad (16)$$

under the frame design described in Fig. 1, the GAP-rule is throughput-optimal.

Proof: The proof is based on the framework of [3], with the difference that our scheme adopts a variable frame-length structure. We need to specify that the proof is still effective in this circumstance. The detailed proof can be found in our journal version [6]. ■

Corollary 1: (The Ergodic Sum Capacity of the GAP-Rule) The ergodic sum capacity can be computed by running the following admission control and scheduling schemes:

Admission control: Before each frame, for every queue with queue size $Q_n(t_k) < V$, the number of arrival bits is set to be W_{\max} , where V and W_{\max} are constants⁸. Otherwise, there are no arrival bits during this frame.

Scheduling: Schedule the users with the GAP-rule.

⁸For Typical values, V and W_{\max} can be approximately 100-fold of the arrival rate.

Then calculate the time-average sum arrival rate A_{avg} . We have

$$A_{\text{avg}} \geq R^* - \frac{B}{V}, \quad (17)$$

and the system is stable, where B is a finite constant, and R^* is the maximum ergodic sum rate, i.e., the ergodic sum capacity.

Remark 1: Leveraging Corollary 1, by letting V be sufficiently large, we can find the maximum ergodic sum rate of the GAP-rule. Theorem 1 establishes a throughput-optimal scheduling scheme, whereas it is still unknown how to characterize the achievable throughput, as well as the achievable rate region explicitly. To this end, we provide Corollary 1 to calculate the maximum ergodic sum rate, which will help us demonstrate the performance gain of the DCA scheme in Section V.

IV. QUEUE-BASED QUANTIZED-BLOCK-LENGTH SCHEDULING SCHEME (QQS)

Due to the fact that the GAP-rule described above requires genie-aided CSIT, and that it is NP-hard, the scheme is practically infeasible. To this end, we propose the QQS, corresponding to the practical version of GAP-rule. In this section, we will first specify the QQS, which bases its scheduling decision solely upon the UQI, neglecting the channel fading fluctuations. Additionally, to reduce the complexity, we divide the users into groups according to their respective channel coherence times, and schedule among different groups, based on the intuition that serving users with significant channel coherence time difference is *not* a good choice since the users with longer coherence time will be encumbered. The QQS is specified as

- Step 1) Initialization:

Denote the overall user set by \mathcal{N} . Divide the users into K groups, each of which denoted by \mathcal{N}_k , $k = 1, 2, \dots, K$, based on a uniform channel coherence time quantization

$$\mathcal{N}_k = \left\{ n \in \mathcal{N} \left| \frac{k-1}{K} T_{\max} \leq T_n \leq \frac{k}{K} T_{\max} \right. \right\}, \quad (18)$$

where $T_{\max} = \max[T_n]$, $\forall n$, and the users are indexed by

$$\mathcal{N}_k = \{k_1, k_2, \dots, k_{|\mathcal{N}_k|}\}, \quad (19)$$

such that $Q_{k_1} \geq Q_{k_2} \geq \dots \geq Q_{k_{|\mathcal{N}_k|}}$. And

$$\bar{T}_k = \mathcal{M}[T_n, n \in \mathcal{N}_k], \quad (20)$$

$$\mathcal{F}_k = \{k_1\}, \quad (21)$$

$$i=1, \quad (22)$$

where $\mathcal{M}(\cdot)$ denotes the empirical mean.

- Step 2) Group Selection:

For $k = 1 : K$,

For $i = 1 : \mathcal{N}_k$,

If

$$\left(1 - \frac{i+1}{\bar{T}_k}\right) Q_{k_{i+1}} - \frac{1}{\bar{T}_k} \sum_{n=1}^i Q_{k_n} > 0, \quad (23)$$

let

$$\mathcal{F}_k = \mathcal{F}_k \cup \{k_{i+1}\} \quad (24)$$

$$i = i + 1, \quad (25)$$

Else, break for.

End for.

End for.

- Step 3): For each group k , compute

$$\mathcal{P}_k = \max \left[\left\{ \left(1 - \frac{|\mathcal{F}_k|}{\bar{T}_k}\right) \sum_{n \in \mathcal{F}_k} Q_n \right\} \cup \{Q_{k_i} : \forall i \in [1, |\mathcal{N}_k|]\} \right], \quad (26)$$

and set

$$\mathcal{F}_k = k_j, \quad (27)$$

only if the maximization in (26) finds its maximum at a single queue length, Q_{k_j} .

Let

$$k^* = \text{argmax}[\mathcal{P}_k]. \quad (28)$$

- Step 4): Output the scheduled user set \mathcal{F}_{k^*} . ■

Several technical details of the QQS should be mentioned. By treating the block lengths of users in each group as a constant, as in (20), and neglecting the channel fading impact, Eq. (23) stems from the fact that it is sufficient to check whether it is worth adding the user with the largest queue size to the scheduled set. For (26), \mathcal{P}_k denotes the pre-log factor of the queue-weighted sum rate for the k -th group after we select the scheduled set \mathcal{F}_k , considering the possibility that scheduling one user with STC mode is the better choice, i.e., the union with Q_{k_i} . By selecting the maximal $k^* = \text{argmax}[\mathcal{P}_k]$, we find the optimal scheduled group of users, within the heuristic of the algorithm.

Remark 2: It is clear that the computational complexity of the QQS algorithm is $\mathcal{O}(N)$ because it only involves running a sequential test of all the users, the GAP-based algorithms are $\mathcal{O}(2^N)$ because they do an exhaustive search over the user set. The reason for the dramatic complexity decrease compared with the GAP-rule is two-fold. First we group the users based on their channel coherence times, and treating the users in each group with identical channel coherence time. Note that in practice, such a quantization is reasonable since the users are usually categorized into several *mobility states*, see e.g. [7, Section 5.2.4.3] for standardizations in the Long-Term Evolution (LTE) system.

Secondly, we neglect the impact of rate fluctuations due to channel fading. Nevertheless, it can be anticipated that when the number of BS antennas becomes large, i.e., in massive MIMO systems, the user rates are no longer affected by the small-scale channel fading, which is called the *channel hardening effect* [2]. Therefore, the QQS is expected to be *asymptotically throughput-optimal* in the large system regime. This effect will be shown in the numerical results in Section V.

TABLE I
SYSTEM PARAMETERS

Parameters	Value
Carrier frequency f_c	2.6 GHz
Cell radius	1000 m
Bandwidth	15 KHz
Downlink SNR	15 dB
Total time slots	20000
Precoder	Zero-forcing
Channel model	i.i.d. Rayleigh fading model

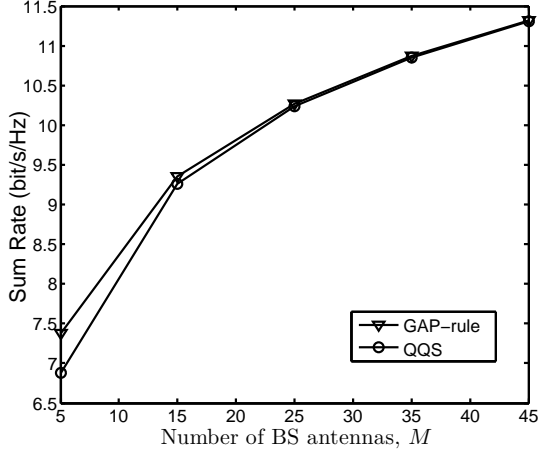


Fig. 2. Cell sum rate of GAP-rule and QQS with user channel coherence time given in Table II. The number of users $N = 5$. The number of user-groups is $K = 2$.

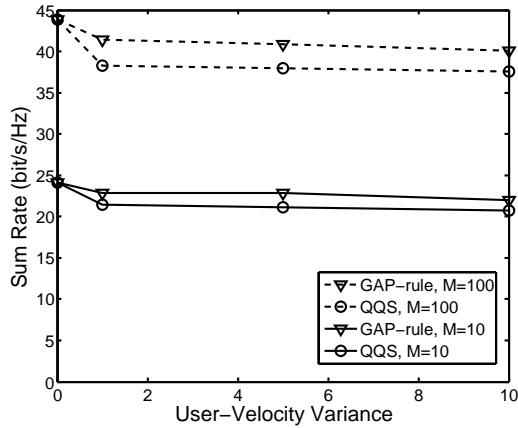


Fig. 3. Cell sum rate of GAP-rule and QQS with Gaussian random user-velocity. The number of users $N = 10$. The number of user-groups, i.e., K , is optimized by exhaustive search.

V. NUMERICAL RESULTS

In this section, we show the performance of our proposed schemes through computer simulations. The parameters used in the simulations are shown in Table I. First, to illustrate the performance of the QQS, we compare the QQS with the throughput-optimal GAP-rule. The sub-optimality of the QQS is shown, which stems from the fact that the QQS heuristically makes two simplifications of the GAP-rule, namely ignoring the channel fluctuations and quantizing the user channel co-

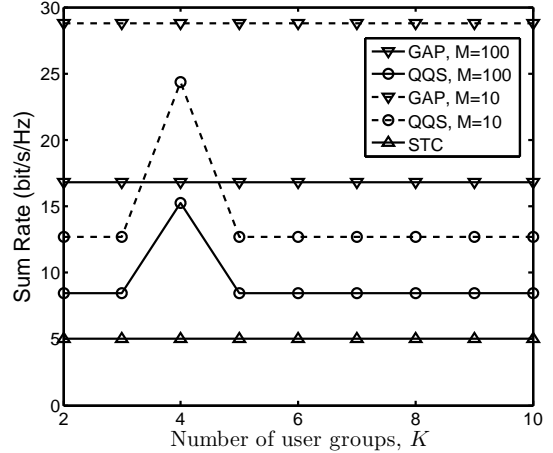


Fig. 4. Cell sum rate of GAP-rule and QQS with uniformly-distributed random user-velocity. The number of users $N = 10$.

TABLE II
USER COHERENCE TIMES

User 1	User 2	User 3	User 4	User 5
100	100	100	5	5

herence time. In Fig. 2, we demonstrate the sub-optimality due to neglecting the channel fluctuations, by letting the user channel coherence times be naturally quantized, thus eliminating the sub-optimality due to coherence-time quantization. We consider a scenario when two types of users coexist: 2 high-mobility (60 km/h) vs. 3 low-mobility (3 km/h) users. The user channel coherence times⁹ in terms of the number of channel uses are shown in Table II, according to

$$T = B_c \times T_c, \quad (29)$$

where $B_c = \frac{c}{4\Delta d}$ and $T_c = \frac{1}{8f_c v c}$, c denotes the light speed, Δd is related to the cell radius and v denotes the user velocity [8]. We run the simulation of the algorithms for 20000 time slots, which is 1.3 seconds under these parameters, and compute the sum rate by averaging the service rate and Corollary 1. It is observed that the QQS is asymptotically throughput-optimal, in the large-system regime. Despite of the sub-optimality when the number of BS antennas is limited, the rate-loss is marginal. Notice that even when $M = 5$, the sum rate loss is approximately 0.5 bit/s/Hz, since it is well-known that in MU-MIMO systems, the user-rates with linear beamforming converge to the so-called deterministic equivalents quite fast, when the system dimension goes up [9], even with moderate M . Therefore it is reasonable to put aside the channel fluctuations and focus on the queue information as in the QQS.

Fig. 3 shows the QQS performance when the user coherence times are random. We let the user-velocities be truncated Gaussian distributed with mean 3 and 60 km/h¹⁰, and the

⁹Note that we refer to the channel coherence time as the block length in the block fading model in this paper. Because users in one cell usually have identical channel coherent bandwidths, different block lengths of users are mainly due to different coherence times. Therefore we use channel coherence time instead of block length in the paper for better illustration.

¹⁰The negative velocities are eliminated and re-generated.

variance σ_v^2 is given as the x-axis of the figure. It is observed that when $\sigma_v^2 = 0$, i.e., the channel coherence times are naturally quantized, the QQS performs as good as the GAP, confirming the intuition of avoiding scheduling users with dramatically distinct channel coherence times together is reasonable. Furthermore, when the user-velocities vary, the rate loss of the QQS is fairly acceptable, given the fact that it not only dramatically decreases the complexity, but also makes the algorithm practical comparing with the GAP-rule. It is also observed that the rate gap of $M = 100$ is larger than that of $M = 10$, whereas relatively, the rate gaps are similar given the relative difference, which indicates that the analytic expression of the rate gap may involve a term that scales with M , possibly as $\log(M)$ since this is the form of the power gain.

The impact of the number of user-groups in the QQS is shown in Fig. 4. It is important to set the number of user groups K in the QQS, since it allows transmission to users in the same group exclusively. Specifically, the channel coherence time approximation in each group will be inaccurate when K is too small. On the other hand, over-grouping the users, i.e., large K simply leads to time-sharing among different users. It is observed that there exists an optimum K , resulting the performance of the QQS is fairly close to the GAP. The analytic analysis of the optimal number of K is not given due to the heuristics of the QQS algorithm. However, given the limited searching space of K , which at most scales linearly with the number of users, and the fact that the search is only required as often as the user channel coherence time changes, which is shown to be on the order of seconds to tens of seconds, an exhaustive search is acceptable. Nonetheless, the exact analysis is left to be an interesting future work. Note that the design of better user-grouping schemes, i.e., user channel coherence time quantization schemes, rather than uniform and fixed quantization, is also worth studying in the future. For comparison, the sum rate of simply time-division-multiple-access (TDMA) among users is also plotted, which is evidently inferior compared to QQS or GAP-rule. It is noteworthy to mention that multiplexing all the users renders a *zero* throughput, given the user coherence times in this simulation, due to the CSIT acquisition overhead occupies all the available time-frequency resources.

VI. CONCLUSIONS

In this paper, we have investigated the user-scheduling problem in MU-MIMO downlinks considering CSIT acqui-

sition overhead. It is shown that the performance of a system regardless of the CSIT acquisition overhead is very poor when the channel coherence time is comparable with the CSIT acquisition overhead. Therefore, a CSIT-overhead aware scheduling scheme is in great need. To this end, we formulate the generic Lyapunov-drift optimization scheduling problem, namely GAP, according to which the GAP-rule is established and proved to be throughput-optimal. In view of the NP-hardness and non-causality of the GAP-rule-based schemes, the QQS is proposed, which only requires $\mathcal{O}(N)$ complexity, and possesses practical feasibility. It is shown that the QQS performs fairly close to the GAP-rule, when the system dimension is large and the user coherence times can be grouped. The QQS suffers reasonable rate loss when either of the conditions is not met exactly. Nevertheless, the performance improvement comparing with non-DCA schemes or simple TDMA is substantial.

ACKNOWLEDGMENT

This work is sponsored in part by the National Basic Research Program of China (973 Program: 2012CB316001), the National Science Foundation of China (NSFC) under grant No. 61201191 and No. 61321061, and Hitachi R&D Headquarter.

REFERENCES

- [1] G. Caire, N. Jindal, M. Kobayashi, and N. Ravindran, "Multiuser mimo achievable rates with downlink training and channel state feedback," *IEEE Trans. Inform. Theory*, vol. 56, pp. 2845–2866, Jun 2010.
- [2] T. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, pp. 3590–3600, Nov 2010.
- [3] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," *Synthesis Lectures on Communication Networks*, vol. 3, no. 1, pp. 1–211, 2010.
- [4] T. Yoo and A. Goldsmith, "On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming," *IEEE Trans. Select. Areas Commun.*, vol. 24, pp. 528–541, Mar 2006.
- [5] R. Zakhour and D. Gesbert, "A two-stage approach to feedback design in multi-user mimo channels with limited channel state information," in *IEEE PIMRC 2007.*, pp. 1–5.
- [6] Z. Jiang, S. Zhou, and Z. Niu, "Dynamic channel acquisition in mmimo," *submitted to IEEE Trans. Commun.*
- [7] E. U. T. R. Access, "User equipment (ue) procedures in idle mode (release 9)," *3GPP TS*, vol. 9, p. V9, 2009.
- [8] D. Tse and P. Viswanath, *Fundamentals of wireless communication*. Cambridge university press, 2005.
- [9] S. Wagner, R. Couillet, M. Debbah, and D. T. M. Slock, "Large system analysis of linear precoding in correlated miso broadcast channels under limited feedback," *IEEE Trans. Inform. Theory*, vol. 58, pp. 4509–4537, Jul 2012.