Auto-Tuning for Cellular Scheduling Through Bandit-Learning and Low-Dimensional Clustering

Isfar Tariq[®], Graduate Student Member, IEEE, Rajat Sen, Member, IEEE, Thomas Novlan, Member, IEEE, Salam Akoum, Member, IEEE, Milap Majmundar, Senior Member, IEEE,

Gustavo de Veciana[®], Fellow, IEEE, Sanjay Shakkottai[®], Fellow, IEEE

Abstract— We propose an online algorithm for clustering channel-states and learning the associated achievable multiuser rates. Our motivation stems from the complexity of multiuser scheduling. For instance, MU-MIMO scheduling involves the selection of a user subset and associated rate selection each time-slot for varying channel states (the vector of quantized channels matrices for each of the users) — a complex integer optimization problem that is different for each channel state. Instead, our algorithm clusters the collection of channel states to a much lower dimension, and for each cluster provides achievable multiuser capacity trade-offs, which can be used for user and rate selection. Our algorithm uses a bandit approach, where it learns both the unknown partitions of the channel-state space (channel-state clustering) as well as the rate region for each cluster along a pre-specified set of directions, by observing the success/failure of the scheduling decisions (e.g. through packet loss). We propose an epoch-greedy learning algorithm that achieves a sub-linear regret, given access to a class of classifying functions over the channel-state space. We empirically validate our approach on a high-fidelity 5G New Radio (NR) wireless simulator developed within AT&T Labs. We show that our epoch-greedy bandit algorithm learns the channel-state clusters and the associated rate regions. Further, adaptive scheduling using this learned rate-region model (map from channel-state to the set of feasible rates) outperforms the corresponding hand-tuned static maps in multiple settings. Thus, we believe that auto-tuning cellular systems through learning-assisted scheduling algorithms can significantly improve performance in real deployments.

Manuscript received March 25, 2020; revised January 2, 2021; accepted April 1, 2021; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor R. La. Date of publication May 19, 2021; date of current version October 15, 2021. This work was supported in part by the NSF under Grant CNS-1731658, Grant CNS-1718089, and Grant CNS-1910112; in part by the Army Futures Command under Grant W911NF-19-2-0333; in part by the Wireless Networking and Communications Group Industrial Affiliates Program; and in part by the U.S. Department of Transportation (DoT) supported the Data-Supported Transportation Operations and Planning (D-STOP) Tier 1 University Transportation Center. This article was presented in part at the 2019 IEEE International Conference on Computer Communications. (Corresponding author: Isfar Tariq.)

Isfar Tariq, Gustavo de Veciana, and Sanjay Shakkottai are with the Department of Electrical and Computer Engineering, The University of Texas at Austin, Austin, TX 78712 USA (e-mail: isfartariq@gmail.com; gustavo@ece.utexas.edu; shakkott@mail.utexas.edu).

Rajat Sen is with Google Research, Mountain View, CA 94043 USA (e-mail: rajat.sen@utexas.edu).

Thomas Novlan, Salam Akoum, and Milap Majmundar are with the AT&T Labs, Austin, TX 78795 USA (e-mail: thomas_novlan@labs.att.com; salam_akoum@labs.att.com; milap_majmundar@labs.att.com).

This article has supplementary downloadable material available at https://doi.org/10.1109/TNET.2021.3077455, provided by the authors.

Digital Object Identifier 10.1109/TNET.2021.3077455

Index Terms—Online learning, bandit algorithms, wireless networks, scheduling, capacity region, auto-tuning.

I. INTRODUCTION

WIRELESS cellular networks have become increasingly more complex to operate – the aggregate number of parameters available for optimization at various layers can range in the thousands (e.g. MIMO antenna weights, power levels, coding and modulation rates, and frequency/sub-frame allocation to users), and the choice of which depend on the channel-states of the users. Thus, when scheduling users (e.g. in MU-MIMO scheduling [2]), a channel-state dependent combinatorial optimization problem needs to be solved each time-slot, where a subset of users need to be selected, and transmission rates and power levels jointly determined for each of these users from among the allowable parameters. This problem however has a latent low-dimensionality that can be exploited, namely that for channel-states that are "near" each other, the optimal solution (user and rate selection) is likely to be the same. Thus, if we cluster channel-states, and determine the effective rate region trade-offs for each cluster, these cluster-dependent rate regions can be used for user and rate selection, and thus significantly reduce the complexity of user and rate scheduling.

However, these clusters are unlikely to be universal, meaning that different scenarios (e.g. indoor, outdoor urban, outdoor rural) would lead to different channel-state clusterings. Indeed, it is also likely that the clusters and associated rate-regions will also vary with the time of day depending on different loading/use-case scenarios. This then, motivates an *online* clustering and multi-user rate region learning approach.

Such learned rate regions are useful in practice. In literature, scheduling algorithms typically assume that they have access to the set of available rates that are feasible for each channel-state. For instance, MaxWeight-like rules [3] that schedule based on the product of the queue-length and channel-rate implicitly assume that the map from the channel-state to the set of feasible channel rates is known, and thus solve an optimization each time-slot (over all feasible rates) resulting in the scheduling decision. These maps from channel-state to feasible rates, in reality, are hand-tuned by operators based on experiments, and these static maps are chosen such that they are good across several deployment scenarios. Instead, our approach of learning the feasible rate-region in each scenario,

 1 In a MIMO setting, the channel-state for *each* of the users is the channel H matrix, and in practice the base-station would have access to an approximation of this (e.g. quantized version).

1558-2566 © 2021 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

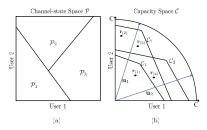


Fig. 1. An illustrative example of the channel-state space $\mathcal P$ and the corresponding capacity classes for n=2 users, K=3 capacity classes and d=1. $\{\mathbf r_{(i)}\}_{i\in[4]}$ are different rate vectors that can be scheduled. The vectors $\{\mathbf u_i\}_{i\in[2]}$ correspond to the directions along which we need to maximize user rates.

and thus having a *different* map for each scenario (effectively, an auto-tuning approach for the PHY/MAC scheduler) permits improved performance in deployments. We refer to Simulations Section VI and Section VI-D for more discussion in the setting of the AT&T Labs cellular network simulator.

Main Contributions: The contributions in this paper are two-fold. From a modeling and algorithm development perspective, we develop a clustering model for the wireless downlink, and develop an epoch-greedy bandit algorithm that learns the clusters, the associated capacity regions and schedules users using learned parameters to minimize regret. From a system simulation perspective, we study the benefits of auto-tuning cellular scheduling using such bandit learning, and demonstrate significant benefits on a high-fidelity cellular simulator developed by AT&T Labs.

We consider a system where the channel-state space \mathcal{P} clusters into K (unknown) classes, with a corresponding multiuser rate-region for each class. Our goal is to develop *online* strategies that can learn clusterings of different channel-states that have similar multiuser rate regions along with the boundaries of these regions. Simultaneously while learning, we need to schedule users based on the observed channel-states to maximize the user rates along pre-specified directions (see Figure 1(a) precise definition is given in Section III). Our contributions are:

(i) We propose an *epoch-greedy* bandit algorithm for our problem setting. The algorithm assumes access to a class of experts/classifying functions Π , where an expert in Π is a mapping from the channel-state space to $\{0,1\}$. We also assume that the class of experts is rich enough, such that there exists a set of functions, which when composed together can yield a function from the channel-state space to $\{1, 2, \dots, K\}$ which correctly identifies the class in which each channel-state belongs in. Similar assumptions have been made in the realizable setting in stochastic contextual bandits [4]. Our approach achieves a balance between three objectives: (i) Class Explorelearning the clustering of the channel-state space using the class of experts and feedback obtained by scheduling different rates in an exploratory manner (ii) Capacity Explore- learning the boundaries of the capacity regions in the specified directions for the different channel-state region clusters (iii) Exploit - finally, exploiting the knowledge learned, by scheduling the rate vector of maximum possible magnitude in the specified direction, that lies within the capacity region corresponding to the channel-state observed in a time-slot.

(ii) We consider a notion of *cumulative regret*, where the regret in our setting is the difference between the total effective rate obtained by a learning policy in T time-slots

and the total rate obtained by a *genie policy* which knows the capacity clusters and the corresponding capacity regions and given a channel-state, always schedules the rate vector of maximum possible magnitude in the specified direction, which lies within the capacity region corresponding to the channel-state. We provide a rigorous definition of regret for our problem in Equation (4). We analyze our algorithm and prove that it has a regret scaling of $\mathcal{O}(T^{2/3}\log T)$ at time T.

(iii) We perform extensive simulations on a high-fidelity simulator – Wireless Next-Generation Simulation (WiNGS) - developed withing AT&T Labs. WiNGS includes a fully dynamic, event-driven system-level simulator which models both the 5G New Radio (NR) physical layer as well as the air interface protocols including the MAC, RLC, and PDCP sublayers. First, we note that the epoch-greedy bandit algorithm is able to learn channel-state clusters and the associated capacity regions over a short time-scale (30 seconds or less in our settings). Further, the associated scheduler using this learned model is able to match and/or outperform hand-tuned policies in multi-user MIMO settings (both with and without out-of-cell interference). This is because a learning-based algorithm is able to auto-tune to each specific scenario (user locations, interference environment, etc.), whereas static hand-tuned policies are chosen for good average performance across many scenarios. Thus, we believe that such bandit algorithms can play a significant role in auto-tuning wireless cellular systems in real deployments. We refer to Section VI-D for additional details.

Finally, we circle back to one of our motivations – understanding the channel-state-dependent capacity regions. Note that since our algorithms focus on optimizing along a pre-specified set of directions, the resulting capacity region that can be constructed for each channel-state class will be an approximation (because we can potentially miss some of the faces of the capacity region). However, if the capacity regions are "nice", then the direction vectors can be designed in order to get an almost exact estimate of the capacity regions. For instance in [5], it has been shown that convex polytopes formed by the intersection R half-spaces (the hyper-planes should have rational coefficients) and for which the vertex enumeration problem is efficient [6], can be learned with O(poly(R,d')) noiseless membership queries, where d' is the dimension of the space.

II. RELATED WORK

Over the last few decades, there has been a lot of work on opportunistic scheduling for wireless networks. This has led to a powerful framework of algorithms that utilize channel feedback and the queue lengths to achieve objectives like system stability, optimization of a utility function or average delay [3]. In the setting of multi-user MIMO wireless networks (MU-MIMO), scheduling algorithms need to optimize over user selection, beamforming (antenna weight selection), power allocations, physical layer modulation and coding parameters [2], [7], [8]. Here, the user selection sub-problem (choosing a subset of users for transmission from among all the possible users) renders leads to a combinatorial explosion in complexity, and several approximations have been used as guidelines for complexity reduction [9]–[11].

We approach dimensionality reduction through online clustering, and our algorithmic approach is related to the contextual multi-armed bandit problem [12]–[14]. The stochastic

contextual bandits with experts problem [4], [13], [15], [16] is especially relevant to our problem. This problem has been studied in the literature starting with the epoch-greedy policy in [13] leading to the more powerful and essentially statistically optimal policies in [4], [15], [16]. Our problem is somewhat similar to this setting as the channel-states observed is analogous to the context and the feedback received after scheduling a rate vector is similar to the stochastic reward observed after pulling an arm. We also assume access to a class of experts that map the space of channel-states to $\{1, 2, \dots, K\}$, where K is the number of capacity classes. However, it should be noted that the feedback received in our setting is much more challenging, as it does not provide direct information about the capacity classes unlike the rewards received from the arms in contextual bandits, which directly reflects the utility of that arm under the given context. Moreover, in our problem there is an additional task of learning the boundary of the capacity regions, even after the clustering of the channel-state region into K classes has been learned.

In the context of learning the capacity regions, there is a line of related work on learning convex polytopes which are formed by the intersection of a finite number of half-spaces, from noiseless membership queries [5], [17]. In [5] binary search type strategies have been used to provide efficient algorithms for learning a class of convex polytopes that are formed by the intersection of half-spaces defined by hyper-planes with rational coefficients and for which the vertex enumeration problem can be solved efficiently. Finally, an earlier version of this paper appeared in [1].

III. SYSTEM MODEL AND DEFINITIONS

We consider a discrete time scheduling system with n users and a single scheduler. At each time t, the scheduler observes a channel-state vector $\mathbf{q}(t) = \{\mathbf{q}_1(t), \mathbf{q}_2(t), \dots, \mathbf{q}_n(t)\}$ where $\mathbf{q}_i(t) \in \mathcal{Q}^d$ is the channel-state for user $i \in [n]$, where $[n] \triangleq \{1, 2, ..., n\}$. The set \mathcal{Q} can be a bounded subset of \mathbb{R} or a discrete alphabet set. We denote the set of all channel-state vectors as $\mathcal{P}(=(\mathcal{Q}^d)^n)$. At any time t we observe the channel-state vector \mathbf{q} from a time-invariant distribution $f_{\mathcal{Q}}$ over \mathcal{P} (this distribution depends on the wireless channel between the user and the base-station).

A. Scheduling a Rate Vector

Corresponding to each channel-state, there is a unique capacity region that the system can support. The capacity region corresponding to a channel-state is defined as the set of all user rate vectors $\mathbf{r} \in \mathbb{R}^n_+$ that can be achieved with probability close to one, potentially by time-sharing. Strictly speaking, we are really considering the rate region, i.e., the set of user rate vectors that are achievable using the available physical layer strategies at the base-station (convex hull of the data rates that be generated using the available physical layer coding/modulation/antenna-beamforming choices), as opposed to an information-theoretic characterization. We however use the term capacity region instead of rate region for clarity of description.

In our subsequent discussion, when using the phrase "schedule a rate vector **r**", it means that we notify the PHY/MAC parameter selection algorithm that **r** needs to be scheduled. Then, this algorithm tries to achieve the rate **r** potentially by time-sharing among various allowable physical layer rates, and over a block of several physical layer time-slots, in which

the channel-state remains the same. Finally, at the end of this time-share block, a notification is received which tells us whether the requested rate **r** is achieved or not. Therefore, in the subsequent discussion we use 'time-slot' as an abstraction for one trial by the PHY/MAC parameter selection algorithm to achieve a rate over a block of physical layer time-slots. Note that we use a finite length block of physical layer time-slots to judge whether a rate **r** can be achieved and therefore the notification is bound to be noisy. This noise is captured in our *noise model* which is described later in this section.

B. Channel-State Partitions and Capacity Regions

We assume that the channel-state space \mathcal{P} can be partitioned into K sets denoted by $\mathcal{P}_1, \mathcal{P}_2, \ldots, \mathcal{P}_K$ with their corresponding unique capacity regions $\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_K$ respectively 2 such that for any $\mathbf{q} \in \mathcal{P}_i$, the capacity region is \mathcal{C}_i . In the case where \mathcal{Q} is discrete and finite, it is reasonable to assume that $K \ll |\mathcal{P}| = |\mathcal{Q}|^{nd}$. The capacity regions $\mathcal{C}_1, \mathcal{C}_2, \ldots, \mathcal{C}_K \subseteq \mathcal{C} \subset \mathbb{R}^n$ are convex polytopes that lie in the positive quadrant. Further, for non-negative vectors \mathbf{x}, \mathbf{y} , if $\mathbf{x} \leq \mathbf{y}$ (element-wise) and $\mathbf{y} \in \mathcal{C}_i$, then $\mathbf{x} \in \mathcal{C}_i$ for any $i \in [K]$. We also assume that all the capacity regions lie inside the positive quadrant of the ball with radius C centered at the origin, i.e $\mathcal{C}_i \subset \mathcal{B}(0, C)^+$ for all $i \in [K]$. Here, $\mathcal{B}(0, x) = \{\mathbf{u} \in \mathbb{R}^n : \|\mathbf{u}\|_2 \leq x\}$ and \mathcal{A}^+ denotes the subset of \mathcal{A} that lies in the positive quadrant.

We provide an illustrative example in Fig. 1, with n=2 users and K=3 capacity classes. In our example each user provides a one-dimensional channel-state vector, therefore the dimensions of both $\mathcal P$ and $\mathcal C$ are two. The partitions of the channel-state space $\mathcal P$ is shown in Fig. 1(a), which correspond to K different capacity regions in Fig. 1(b). We shall define an index function relating any channel-state vector $\mathbf q \in \mathcal P$ to the channel-state partition and the capacity region as follows.

Definition 1 (Index Function $\mathcal{I}(.)$): Given a channel-state vector \mathbf{q} , $\mathcal{I}(\mathbf{q})$ is the index of the element of the partition that contains \mathbf{q} , i.e. $\mathbf{q} \in \mathcal{P}_{\mathcal{I}(\mathbf{q})}$.

The following assumption states that channel-state's from each element of the partition are observed sufficiently often.

Assumption 1 (Class Probabilities): We assume that $\mathbb{P}(\mathcal{I}(\mathbf{Q}) = i) > \beta = O\left(\frac{1}{K}\right) \ \forall \ i \in [K]$, where $\mathbf{Q} \in \mathcal{P}$ is a random variable with distribution $f_{\mathcal{Q}}$ capturing variability in the system.

C. Separation of Capacity Regions

We assume that the capacity regions are sufficiently different from each other. For instance, in Fig. 1 if \mathcal{C}_1 and \mathcal{C}_2 were almost identical to each other, then it would be better to merge $\mathcal{P}_1, \mathcal{P}_2$ and treat it as a system with K=2. The following assumption says that for any two capacity regions $\mathcal{C}_i, \mathcal{C}_j$ a sufficient fraction of the volume lies outside of their intersection.

Assumption 2 (Separability): We assume the capacity regions are well separated, i.e., for all $i, j \in [K]$

$$d(\mathcal{C}_i, \mathcal{C}_j) \triangleq \frac{|(\mathcal{C}_i \setminus \mathcal{C}_j) \cup (\mathcal{C}_j \setminus \mathcal{C}_i)|}{|\mathcal{B}(0, C)^+|} \geq \lambda > 0,$$

where |A| denotes the volume of the set A.

²Our theoretical guarantees require the channel-state regions corresponding to the different capacity regions be disjoint, however our algorithm can also handle cases where the channel-state classes are not disjoint.

D. Noise Model

Let $Y(\mathbf{q}, \mathbf{r}) \in \{0, 1\}$ denote a random variable modeling the observed notification when a rate $\mathbf{r} \in \mathbb{R}^n$ is scheduled when the observed channel-state vector is \mathbf{q} . Here, $Y(\mathbf{q}, \mathbf{r}) = 1$ signifies a successful transmission and $Y(\mathbf{q}, \mathbf{r}) = 0$ signifies a failure to achieve that rate vector.³ The success or failure to transmit a rate vector \mathbf{r} under channel-state \mathbf{q} is assumed to be an i.i.d random variable $Y(\mathbf{q}, \mathbf{r})$ with distribution given by

$$\mathbb{P}(Y(\mathbf{q}, \mathbf{r}) = 1) = \begin{cases} 1 - \rho(\mathbf{q}, \mathbf{r}), & \text{if } \mathbf{r} \in \mathcal{C}_{\mathcal{I}(\mathbf{q})}, \\ \rho(\mathbf{q}, \mathbf{r}), & \text{if } \mathbf{r} \in \mathcal{B}(0, C)^+ \setminus \mathcal{C}_{\mathcal{I}(\mathbf{q})}, \\ 0, & \text{otherwise.} \end{cases}$$

where $\rho(\mathbf{q}, \mathbf{r})$ can be viewed as a noise parameter (essentially the packet error rate) which depends on the channel-state \mathbf{q} and the rate \mathbf{r} .

Assumption 3 (Noise Rate): We assume that $\rho(\mathbf{q}, \mathbf{r}) \leq \rho < 1/8$, $\forall \mathbf{q}, \mathbf{r}$. We further assume that for all \mathbf{p}, \mathbf{q} and i such that $\mathbf{p}, \mathbf{q} \in \mathcal{P}_i$, $\rho(\mathbf{p}, \mathbf{r}) = \rho(\mathbf{q}, \mathbf{r})$. For notational convenience, for all $\mathbf{q} \in \mathcal{P}_i$, let $\rho_i(\mathbf{r}) \triangleq \rho(\mathbf{q}, \mathbf{r}) = \rho_{\mathcal{I}(q)}(\mathbf{r})$.

Given a channel-state and the corresponding capacity region, when $\bf r$ approaches the boundary of the capacity region (from inside) the probability of successful transmission is close to 1 but decreases slightly near the boundary. The success probability drops significantly after $\bf r$ crosses the boundary (there is a discontinuous jump in success probability at the boundary). After crossing the boundary, $\rho(\bf q, \bf r)$ decreases till $|\bf r|=C$, beyond which $\rho(\bf q, \bf r)=0$.

E. Bandit Feedback and Objectives

Let $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_D\}$ be a set of unit vectors such that $\mathbf{u}_i \in \mathbb{R}^n_+$. This set is fixed a priori. The broad objective is to discover the maximum possible service rates (i.e. Pareto-optimal points) in these directions,⁴ given a particular channel-state.

Since we use 'time-slot' as an abstraction for a block of several physical layer time-slots where a rate vector **r** is attempted to be scheduled potentially by time-sharing. Therefore, a wide-range of direction vectors within the capacity region can be supported.

Concurrently with the channel-state, a direction vector \mathbf{u} is chosen uniformly at random from the set \mathcal{U} . The task is to schedule a rate vector within the capacity region $\mathcal{C}_{\mathcal{I}(\mathbf{q})}$, of maximum possible magnitude in the direction \mathbf{u} . In other words, we would ideally like to schedule a rate vector $c\mathbf{u}$ such that

$$c(\mathbf{q}) = \underset{d}{\operatorname{arg\,max}} \{d | d\mathbf{u} \in \mathcal{C}_{\mathcal{I}(\mathbf{q})}\}.$$

The precise order of events at a given time-step is as follows:

- A channel-state $\mathbf{q}(t)$ from the distribution $f_{\mathcal{Q}}$ is observed. A direction $\mathbf{u}(t)$ drawn uniformly at random from \mathcal{U} is also specified.
- The policy optionally selects a magnitude $c(\mathbf{q}(t), \mathbf{u}(t)) \in [0, C]$ to be scheduled in the direction $\mathbf{u}(t)$ and the rate vector $\mathbf{r}(t) = c(\mathbf{q}(t), \mathbf{u}(t))\mathbf{u}(t)$ is scheduled. On the other hand, the policy may choose any other rate vector $\mathbf{r}(t)$

that does not lie in the specified direction. In this case the reward obtained is zero in the time-step.⁵

• A notification $Y(\mathbf{q}(t), \mathbf{r}(t)) \in \{0, 1\}$ is then observed.

Remark 1: Consider a 5G-NR setting where two users (aka user-pair) can be scheduled in each time-slot using MU-MIMO (note in general that more than two users could be co-scheduled; in this example, we restrict to two users for discussion clarity). In a typical 5G-NR system, the scheduler and the resource manager are separate entities (see Section V for details). The scheduler's task is to select the user-pair, based on the set of available/active user-pairs and other scheduling constraints. The resource manager's task is to allocate rates to the selected user-pair. In current systems, the resource manager is hard-coded, i.e., a fixed mapping between channel quality to the selected rate (more details in Section V). Furthermore, this map is usually selected in a conservative manner so that the map is "universal" across many deployment scenarios. Our system model focuses on a learning-based resource manager that enables us to instead learn the map between channel quality and user rates in a scenario-dependent manner.

In such a MU-MIMO setting, we use the abstraction of direction vector to represent the scheduling decisions made by the scheduler. This vector direction corresponds to: (a) the choice of user-pair, and (b) their data rate ratio. However, the magnitude of this vector is determined by the resource manager. As we will see in later sections, our system model permits the resource manager to learn the boundary of the capacity region (meaning, the largest permissible magnitude for this vector) by probing the system and learning from success/failure notifications. As a numerical example, for the system shown in Figure 1, let us suppose that at time step t we observe the channel-state from \mathcal{P}_1 and the scheduler provides us the direction vector $(\frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}})$. Let us further suppose that we schedule the rate vector $\mathbf{r} = (\mathbf{10}, \mathbf{20})$ using MU-MIMO, with the corresponding rates for User 1 being 10 and for User 2 being 20. If $\mathbf{r} = (\mathbf{10}, \mathbf{20})$ lies inside the capacity region C_1 , then we will likely receive a success notification and perhaps rarely a failure notification.

Finally, this paper focuses on the rate allocation problem (i.e., the resource manager) where we have access to a scheduler which selects the direction vector to schedule according to allocation constraints of the system at any given time step.

F. Expected Reward Function

Recall $Y(\mathbf{q}, \mathbf{r})$ is the notification received for transmitting rate vector \mathbf{r} when the observed channel-state was \mathbf{q} . Let us define the reward $r(\mathbf{q}, \mathbf{r}, \mathbf{u})$ for a rate vector \mathbf{r} , channel-state \mathbf{q} and direction vector \mathbf{u} to be

$$r(\mathbf{q}, \mathbf{r}, \mathbf{u}) = |\mathbf{r}| \mathbb{1} \{ \mathbf{r}.\mathbf{u} = |\mathbf{r}| \} Y(\mathbf{q}, \mathbf{r}), \tag{1}$$

where $\mathbb{1}\{\}$ is the indicator function.

Note that for any $\mathbf{p}, \mathbf{q} \in \mathcal{P}_i$, we have $\mathbb{E}[r(\mathbf{q}, \mathbf{r}, \mathbf{u})] = \mathbb{E}[r(\mathbf{p}, \mathbf{r}, \mathbf{u})] \triangleq \mathbb{E}[r_i(\mathbf{r}, \mathbf{u})]$. Therefore, we define the expected reward function $f_{\mathbf{u},i}(c)$ for direction vector \mathbf{u} , capacity region \mathcal{C}_i and magnitude c, as follows:

$$f_{\mathbf{u},i}(c) = \mathbb{E}[r_i(c\mathbf{u}, \mathbf{u})].$$
 (2)

 $^{^3}$ A failed transmission can be caused by several factors, such as transmission errors due to fast fading channel or incorrect channel-state sent by users (i.e. estimation and quantization error in $\mathbf{q}(t)$).

⁴The direction vector \mathbf{u} controls the ratio of data that is transmitted to the users. Furthermore, the direction vector \mathbf{u} , also controls the set of users that are scheduled, i.e. user i is scheduled if $\mathbf{u}(i) > \mathbf{0}$.

⁵Note that this is a conservative estimate of the reward. In general, there is some non-zero value in scheduling any rate vector in the capacity region corresponding to the observed channel-state. However, our theoretical guarantees will be under this conservative reward model, and in practice the performance observed will only be better.

The function $f_{\mathbf{u},i}(c)$ is the expected rate achieved if we schedule a rate vector $c\mathbf{u}$ when the channel-state observed belongs to capacity class i. It can be evaluated as follows,

$$f_{\mathbf{u},i}(c) = \begin{cases} c(1 - \rho_i(c\mathbf{u})), & \text{if } c\mathbf{u} \in \mathcal{C}_i, \\ c\rho_i(c\mathbf{u}), & \text{if } c\mathbf{u} \in \mathcal{B}(0,C)^+ \setminus \mathcal{C}_i, \\ 0, & \text{otherwise.} \end{cases}$$
 (3)

Since we have assumed that $\rho_i(\mathbf{r}) < \frac{1}{8} \, \forall \, \mathbf{r}$ therefore $f_{\mathbf{u},i}(c)$ is a discontinuous function of c, and the discontinuity is located at the point where a ray in the direction \mathbf{u} meets the boundary of \mathcal{C}_i . We make the following assumption on the expected rate function.

Assumption 4 (Maxima of Rate Function): Let us define

$$\hat{c}_{\mathbf{u},i} = \arg\max_{c} f_{\mathbf{u},i}(c)$$

and $c_{\mathbf{u},i}^* = \max_c \{c | c\mathbf{u} \in \mathcal{C}_i\}$. We assume that the noise function $\rho_i(\mathbf{r})$ is such that $c_{\mathbf{u},i}^* = \hat{c}_{\mathbf{u},i}$.

This assumption basically implies that for all $i \in [K]$ the maximum of the rate function $f_{\mathbf{u},i}(c)$ is achieved at the point where a ray in the direction \mathbf{u} meets the boundary of \mathcal{C}_i .

G. Class of Experts

We assume access to a class of binary experts/classifiers Π , where each expert $\hat{\pi} \in \hat{\Pi}$ is a function mapping the space of channel-states to $\{0,1\}$ i.e $\hat{\pi}: \mathcal{P} \to \{0,1\}$.

Assumption 5 (Classifying Functions): Let κ be a proper subset of [K]. Let us define the following binary function $\hat{\mathcal{I}}_{\kappa}(\mathbf{q}) = \sum_{i \in \kappa} \mathbb{1}\{\mathbf{q} \in \mathcal{P}_i\}$. We assume that the set of binary experts/classifiers $\hat{\Pi}$ is such that for all $\kappa \subset [K]$, $\hat{\mathcal{I}}_{\kappa}(.) \in \hat{\Pi}$. We further assume that the VC dimension [18] of our class of experts is V.

The above assumption states that the binary functions from \mathcal{P} to $\{0,1\}$ that are induced by labeling the channel-state's belonging to a set $\kappa \subset [K]$ of capacity classes as 1 and the rest as 0, are a part of our class of experts, for all such proper subsets κ . Consider the example exhibited in Figure 1. Suppose $\kappa = \{1,2\}$. Then, $\hat{\mathcal{I}}_{\kappa}(\mathbf{q})$ divides \mathcal{P} into two regions $\mathcal{P}_1 \cup \mathcal{P}_2$ and \mathcal{P}_3 . Note that both these regions can be represented as the intersection of at most two half-spaces, as the boundaries of the partitions are linear. This is true for all such proper subsets κ . Therefore, if our class of binary classifiers contains all the separators that are intersections of at most two half-spaces (i.e. $\mathbb{1}\{\bigcap_{i=1}^2 \boldsymbol{\beta}_i^T \mathbf{q} + \beta_{i,0} > 0\}$, parameterized by $(\boldsymbol{\beta}_i, \beta_{i,0})$, then Assumption 5 is valid.

Assumption 5 basically implies that there exists a group of binary functions in Π , which when composed together can yield the true index function, i.e., $\mathcal{I}_{\kappa}(\mathbf{q})$ for different κ 's can be composed together to recover $\mathcal{I}(\mathbf{q})$. For example in Figure 1, $\hat{I}_{[2,3]}(\mathbf{q})$ differentiates Class 1 from 2,3 and $\mathcal{I}_{[3]}(\mathbf{q})$ separates Class 3 from the rest. Given a channel-state \mathbf{q} , if for instance $\mathcal{I}_{[3]}(\mathbf{q}) = 0$ and $\mathcal{I}_{[2,3]}(\mathbf{q}) = 1$ then we can infer that $\mathcal{I}(\mathbf{q}) = 2$. Note that this is similar to the *realizable* setting in the contextual bandits with experts problem [4], where it is assumed that the true behavior of the system can be represented by one of the expert function. However, finding the correct expert in an online setting is an algorithmic challenge. For instance for the example discussed above we can recover the index function $\mathcal{I}(\mathbf{q})$ by using $\mathcal{I}_{[3]}(\mathbf{q})$ and $\hat{\mathcal{I}}_{[2,3]}(\mathbf{q})$. For channel states from channel-state cluster 1, 2 and 3 the value for $(\hat{\mathcal{I}}_{[3]}(\mathbf{q}), \hat{\mathcal{I}}_{[2,3]}(\mathbf{q}))$ will be (0,0), (0,1)

and (1,1) respectively. Therefore, we can construct a function $\hat{\mathcal{I}}(q) = f(\hat{\mathcal{I}}_{[3]}(\mathbf{q}), \hat{\mathcal{I}}_{[2,3]}(\mathbf{q}))$ such that f(0,0) = 1, f(0,1) = 2 and f(1,1) = 3. We also refer to Section IV for an example on constructing $\hat{\mathcal{I}}(q)$ in an online setting.

Remark 2: In practical settings, the class of experts/ classifiers can be provided by implementing a module/function for the classifier directly on the base-station (this is similar to how the scheduler or resource manager module is implemented in the 5G-NR setting we discussed in Remark 1). For instance, if we select logistic regression, then our class of experts would consist of functions $\hat{\pi}(x) = \mathbb{1}\{\text{sig}(\mathbf{w}^{T}[\mathbf{x}; \mathbf{1.0}]) > 1\}$ τ_0 } over all linear weights **w** and threshold τ_0 . Here, $sig(x) = 1/(1 + e^{-x})$ is the sigmoid function. Note that even though the set of all functions in our class can correspond to all weights s.t. $\|\mathbf{w}\| \leq R$ for some R, we still need to learn the correct weights that correspond to the partition functions, $\hat{I}_{\kappa}(.)$ in Assumption 5. As another example, if we select Gaussian Naive Bayes, the corresponding class of experts would consist of functions $\hat{\pi}(x) = \arg\max_{i \in \{0,1\}} \frac{1}{\sqrt{(2\pi)^{nd}|\Sigma|}} \exp^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu_i})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu_i})}$ where $\Sigma = \operatorname{Diag}(\boldsymbol{\sigma_i})$, over a range of values for $\{\boldsymbol{\mu_i}, \boldsymbol{\sigma_i}\}$'s. Again the correct parameters $\{\mu_i, \sigma_i\}$ that correspond to the partition functions in Assumption 5 need to be learned. Several reference implementations (e.g. sklearn in python) are available, which can be ported as a module into the basestation. Additional discussion (e.g. on complexity) is available in Appendix F.

H. Definition of Regret

The main objective is to minimize regret when compared to a genie strategy which knows the index function \mathcal{I} and the capacity regions \mathcal{C}_i 's. Let $\mathbf{r}(t)$ be the rate vector selected by a policy, at time t. Then the regret of the policy till time T is given by:

$$R(T) = \sum_{t=1}^{T} \left(f_{\mathbf{u}(t), \mathcal{I}(\mathbf{q}(t))} (\hat{c}_{\mathbf{u}(t), \mathcal{I}(\mathbf{q}(t))}) - \mathbb{E} \left[r_{\mathcal{I}(\mathbf{q}(t))} (\mathbf{r}(t), \mathbf{u}(t)) \right] \right)$$
(4)

where $\mathbf{q}(t), \mathbf{u}(t)$ are the channel-state vector and direction vector at time t, respectively. Note, that $f_{\mathbf{u}(t),\mathcal{I}(\mathbf{q}(t))}(\hat{c}_{\mathbf{u}(t),\mathcal{I}(\mathbf{q}(t))})$ is the maximum average rate that can be achieved in the direction $\mathbf{u}(t)$, by a *genie policy* that knows the capacity classes and the boundaries of the capacity regions. The regret measures the sub-optimality of the policy in question with respect to the *genie policy*, in an expected sense. The goal is to design a policy that yields R(T) that is sub-linear in T, for all times T large enough. This basically implies that the policy keeps learning the system as time progresses.

Remark 3: Although we have only used channel-state for clustering, other 5G parameters such as rank or angle of arrival of primary beam can also be used for clustering, by treating them as additional dimensions of the user channel-state. Moreover, we can explore parameters such as the power allocation among users by treating them as additional dimensions of the direction vectors.

IV. ALGORITHM

The algorithm is structured as an *epoch-greedy* strategy [13]. One key algorithmic idea is that if a rate vector **r** is scheduled for several different observed channel-state's **q**, then the

success notifications $Y(\mathbf{q}, \mathbf{r})$ provide useful information that can be leveraged using the class of binary experts Π to obtain a binary classifier that separates the channel-state space ${\cal P}$ into two regions \mathcal{P}_* and \mathcal{P}_*^c , where $\mathcal{P}_* = \{\mathbf{q} \in \mathcal{P} : \mathbf{r} \in \mathcal{C}_{\mathcal{I}(\mathbf{q})}\}.$ A carefully chosen set of rate points can then be used to form a group of binary classifying functions, which when composed together yields a mapping $\pi \colon \mathcal{P} \to [K]$, which is identical to $\mathcal{I}(\mathbf{q})$ with high probability.

The algorithm starts with an initialization phase and then proceeds in epochs. In initialization phase the algorithm constructs π by building a tree of binary classifiers which is then used to classify the channel-state points into K different classes. This stage is referred to as *initializing classifier* π . After building π , the algorithm runs in epochs similar to epoch-greedy policies for contextual bandits. At the beginning of each epoch, there is a class explore stage corresponding to improving the accuracy of classifier π . This is followed by a capacity explore stage aimed at learning the capacity regions of the K different channel-state partitions, in the given directions \mathcal{U} . The last stage in an epoch is the *exploitation* stage where we deduce the correct capacity class of the observed channel-state vector using π and then schedule the optimal rate vector according to the current belief about the boundary of the corresponding capacity region. An illustrative pseudo-code of our algorithm is shown in Algorithm 1, while a more detailed pseudo-code can be found as Algorithm 4. We will explain each of the stages/phases in more detail in subsequent sections.

Algorithm 1 Epoch-Greedy Algorithm for Online Capacity Class Learning and Rate Allocation

- 1: Initialize classifier π , by observing t_0 channel-state's, scheduling corresponding to carefully designed rate vectors and observing the notifications. (Initializing Classifier)
- 2: Epoch: l = 1. Time: $t = t_0$.
- 3: while $t \leq T$ do
- Update the classifier π by observing channel-state q, scheduling a carefully chosen rate point r, and using the notification $Y(\mathbf{q}, \mathbf{r})$. This is repeated K-1 times. (Class explore)
- Learn the boundaries of the K capacity regions in the directions \mathcal{U} , by scheduling carefully chosen rate points and using the current π . A total of $\alpha(l)$ rate points are scheduled in this part of the epoch. (Capacity explore)
- Schedule next s(l) rate points optimally using π and the 6: learned boundaries. (Exploit)
- Let $t = t + K 1 + \alpha(l) + s(l)$ and l = l + 1.
- 8: end while

A. Initializing Classifier π

The first stage of the algorithm is to initialize the mapping (multi-class classifier) π used to classify the different channelstate's into the K different classes. This mapping consists of K-1 binary experts from our class of experts, which are composed together in a tree-like structure, in order to yield the mapping π .

The detailed pseudo-code for this phase is provided as Algorithm 2. In the beginning of this phase, for several time-slots the channel-state's are observed and stored, while not making any scheduling decisions (for instance, the scheduler is allowed to proceed in its default behavior). This process is continued until we observe n_0 distinct channel-state vectors, which are essentially n_0 distinct i.i.d random variables sampled from $f_{\mathcal{O}}$.

Algorithm 2 Initializing the Classifier Tree

1: Schedule arbitrary rate vectors for the first n_0 channel-state vectors observed. Let i = 0 and form a tree \mathcal{T} where the root contains the n_0 initial channel-state points. There are no other nodes in the tree.

```
2: while i < K - 1 do
      Randomly select a rate point r.
      S_i = \{\}
4:
      for l = 1 : l_0 do
5:
6:
         Let \mathbf{q} be the observed channel-state at time-step t.
          Schedule rate r. (Class Explore)
         Let y \in \{0, 1\} be the notification received. Add (\mathbf{q}, y)
          Set t = t + 1.
9:
10:
      Construct a binary classifier \hat{\pi}_i by empirical risk mini-
11:
```

- mization (ERM) over S_i , over the expert set Π .
- 12: for all leaves j of T do
- Classify the channel-state at leaf j according to the 13: classifier $\hat{\pi}_i$. Let n_i be the number of channel-state points at leaf j.
- if $\frac{n_0\beta}{2}$ < number of leaf channel-state classified as $0 < n_j - \frac{n_0 \beta}{2}$ then
- Make leaf j into a parent of two new leaves where 15: the left leaf has all the channel-state's classified as 1 and the right has all the channel-state's classified as 0.

```
i = i + 1
16:
17:
             Break
          end if
18:
19:
      end for
20: end while
```

Then we begin initializing the tree-structure which is detailed in steps 2-20 of Algorithm 2. Note that in each iteration of the while loop in step 2 of Algorithm 2, a rate point is randomly selected and then for the following l_0 time-slots irrespective of the channel-state observed, this rate point is scheduled. The feedback observed helps us in building a binary classification data-set that can be used to train a classifier $\hat{\pi} \in \Pi$ which can differentiate all $\mathbf{q} \in \mathcal{P}$ such that $\mathbf{r} \in \mathcal{C}_{\mathcal{I}_{\mathbf{q}}}$ from the rest. We assume that the classifiers are trained in step 11 using empirical risk minimization (ERM) with the 0-1loss function. Therefore, we have that:

$$\hat{\pi}_i = \underset{\hat{\pi} \in \hat{\Pi}}{\operatorname{argmin}} \frac{1}{|S_i|} \sum_{(\mathbf{q}, y) \in S_i} \mathbb{1}\{\hat{\pi}(\mathbf{q}) \neq y\}.$$

At any point in time, an internal node \mathcal{N}_i in the tree stores the triplet $(\hat{\pi}_i, \mathbf{r_i}, S_i)$ where $\hat{\pi}_i$ is the expert obtained by ERM over the examples $\{(\mathbf{q}, y)\}$ stored in S_i which were in turn obtained by scheduling the rate \mathbf{r}_i for l_0 time-slots. A leaf of the tree \mathcal{L}_i stores a subset of the initial n_0 channel-state points. In each

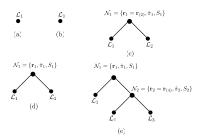


Fig. 2. Construction of a classification tree which represents the final initial classifier π that maps $\mathcal{P} \to [K]$ corresponding to channel-state class structure in Fig. 1.

iteration of the while loop, the classifier trained using the data collected by scheduling the current randomly chosen rate point, is only retained if it can split at least one of the current leaf nodes in the tree reliably into two distinct partitions. This is achieved by the check in step 14 of Algorithm 2. The while loop continues to iterate until the tree has K-1 internal nodes.

In order to illustrate this phase, let us consider a system as shown in Figure 1 with r=2 users, such that the channel-states can be partitioned into K=3 classes \mathcal{P}_1 , \mathcal{P}_2 and \mathcal{P}_3 with capacity regions \mathcal{C}_1 , \mathcal{C}_2 and \mathcal{C}_3 respectively.

For, simplicity let $\mathbf{r}_{(1)},...,\mathbf{r}_{(4)}$ be the first four rate points that are randomly chosen in step 3 of Algorithm 2, in that order (see Fig. 1). Since, $\mathbf{r}_{(1)}$ is a rate point that lies in all the capacity regions, the corresponding classifier $\hat{\pi}_1$ formed using that data collected in step 8, will classify most of the n_0 channel-state points as 1. Therefore, this will not split the current leaf node (the root node with n_0 initial channel-state vectors) into any partitions. Hence, the classifier and the rate point is discarded and the value of the iterator i remains unchanged. The tree remains the same with one leaf node as shown in Fig. 2(a)-(b).

In the next iteration of the while loop, the randomly chosen rate point is $\mathbf{r}_{(2)}$. The data collected using $\mathbf{r}_{(2)}$ is used to train a classifier $\hat{\pi}_1$, which classifies most points in class \mathcal{P}_2 as 1, while classifying most points outside of \mathcal{P}_2 as zero.⁶ This point splits the n_0 channel-state points in the current leaf node into two partitions. Therefore, the classifier is retained. An internal node $\mathcal{N}_1 = \{\mathbf{r}_1, \hat{\pi}_1, S_1\}$ is formed where $\mathbf{r}_1 = \mathbf{r}_{(2)}$. Moreover, two leaf nodes are formed where \mathcal{L}_1 is a leaf corresponding to all the n_0 channel-state vectors that are labeled as 1 by $\hat{\pi}_1$ and \mathcal{L}_2 contains the rest. This is illustrated in Fig. 2(c).

In the next iteration, the rate point $\mathbf{r}_{(3)}$ is chosen, which will effectively yield the same classifier as the one corresponding to $\mathbf{r}_{(2)}$. Therefore, this classifier will be insufficient to split any of the leaves in Fig. 2(c). Thus the value of i is unchanged and the tree remains the same as shown in Fig. 2(d).

Finally, the rate point $\mathbf{r}_{(4)}$ is chosen. The classifier $\hat{\pi}_2$ corresponding to this point ideally distinguishes between points lying in \mathcal{P}_1 from those outside of \mathcal{P}_1 . Thus, this new classifier can split the points in leaf \mathcal{L}_2 of the tree in Fig. 2(c), into two nodes, as shown in Fig. 2(d). This leads us to our final classifying tree π . Ideally (ignoring classification errors), a channel-state point belonging to $\mathcal{P}_1,\mathcal{P}_2$ and \mathcal{P}_3 will land in \mathcal{L}_3 , \mathcal{L}_1 and \mathcal{L}_2 respectively.

The parameters n_0 , l_0 have been chosen in order to ensure that w.h.p a correct classifying tree is obtained. The following lemma formalizes this claim.

 6 Note that this is just an initialization of the classifier and moreover the feedback received from scheduling is noisy. Therefore, the binary classifiers trained will not be fully accurate. However, n_0 and l_0 are designed to be large enough such that with high probability the tree structure is correct.

Lemma 1: Let
$$n_0 \geq \frac{24K}{\beta^2}\log\left(\frac{4\log(\frac{1}{\delta})+K}{\delta\lambda}\right)$$
 and l_0 is large enough such that $\frac{1}{1-2\rho}\sqrt{\frac{V}{l_0}}+\sqrt{\frac{2\log\left(\frac{l_0^2}{\delta}\right)}{l_0}}<\frac{\beta}{4K}$ and $l_0>\sqrt{\left(\frac{4\log(\frac{1}{\delta})+K-1}{K\lambda}\right)}$. Then with probability at least $1-3K\delta$, the loop in step 2 of Algorithm 2 is terminated after at most $\frac{4\log(\frac{1}{\delta})+K-1}{\lambda}$ iterations and further a correct classifying tree structure is obtained.

B. Class Explore

After the classification tree is initialized, the algorithm proceeds in epochs and the structure of the tree remains unchanged. The first few time-slots in each epoch are dedicated to improving the accuracy of the classifiers $\hat{\pi}_i$'s stored in the internal nodes of the tree \mathcal{N}_i 's. We name this part of an epoch class explore. The class explore phase in an epoch consists of K-1 time-steps t_1, \ldots, t_{K-1} . At time-step t_i , let the channel-state observed be \mathbf{q}_i . After the channel-state is observed, the rate vector \mathbf{r}_i stored in the internal node \mathcal{N}_i is scheduled and a notification y_i is received The data-sample (\mathbf{q}_i, y_i) is added to the set S_i and $\hat{\pi}_i$ is updated through ERM over the updated set S_i . This is performed for all $i = 1, 2, \dots, K - 1$. This phase is detailed in steps 7 -14 of Algorithm 4. The basic idea is to obtain one more training sample for each of the classifiers stored in the internal nodes, at the beginning of each epoch, thereby improving the classification accuracy of the global classier $\pi: \mathcal{P} \to [K]$. The following lemma provides an upper bound for the classification error of the global classifier $\hat{\pi}$ at the beginning of epoch l.

Lemma 2: At the end of the class explore phase in epoch l with probability at least $1-(K-1)\frac{\delta}{(l+l_0)^2}$ we have

$$\mathbb{P}(\pi(\mathbf{Q}) \neq \mathcal{I}(\mathbf{Q}))$$

$$\leq (K-1) \left(\left(\frac{1}{1-2\rho} \right) \sqrt{\frac{V}{l_0+l}} + \sqrt{\frac{2\log\left(\frac{(l_0+l)^2}{\delta}\right)}{l_0+l}} \right)$$

$$\triangleq (K-1)\epsilon(l_0+l,\delta),$$

where the probability is over the randomness in $\mathbf{Q} \sim f_{\mathcal{Q}}$ and the randomness in π due to the random samples in the training set

We have provided a proof of Lemma 2 in Appendix C.

Remark 4: Instance-dependent bounds of misclassification rates depend upon multiple variables including the precise channel-state distribution and noise model. The above bound holds for general system models and is useful to show a sublinear scaling of regret. To understand the impact of our algorithms in a practical setting, we have explored the performance benefits in Section V.

C. Capacity Explore

In each epoch, the class explore phase is followed by a few time-slots dedicated to *capacity explore*. This phase is described as steps 16-22 in Algorithm 4. It is aimed towards learning the boundaries of the K capacity classes in the directions \mathcal{U} . Note that there are K possible capacity classes and $D = |\mathcal{U}|$ direction vectors to explore. In the capacity explore phase of epoch l, for $\alpha(l, \delta)$ time-slots we observe the channel-state vectors, direction vectors and schedule carefully

Algorithm 3 Capacity Explore Update

```
1: for \forall k \in [K] and \mathbf{u} \in \mathcal{U} do
2: if m_{k,\mathbf{u}} > \frac{1}{2} then
3: C_{k,\mathbf{u}}[0] = \frac{C_{k,\mathbf{u}}[0] + C_{k,\mathbf{u}}[1]}{2}
4: else if m_{k,\mathbf{u}} < \frac{1}{2} then
5: C_{k,\mathbf{u}}[1] = \frac{C_{k,\mathbf{u}}[0] + C_{k,\mathbf{u}}[1]}{2}
6: end if
7: end for
```

designed rate vectors to learn the capacity region. We set $\alpha(l,\delta) = \frac{2D}{\beta} \left(\frac{16}{1-2\rho}\right)^2 \log\left(\frac{l^2}{\delta}\right).$ We initialize $C_{k,\mathbf{u}}[0] = 0$ and $C_{k,\mathbf{u}}[1] = C$ for all $k \in [K]$

We initialize $C_{k,\mathbf{u}}[0] = 0$ and $C_{k,\mathbf{u}}[1] = C$ for all $k \in [K]$ and $\mathbf{u} \in \mathcal{U}$ at the start of the algorithm. $C_{k,\mathbf{u}}[0]$ is a lower bound for $c_{\mathbf{u},i}^*$ and $C_{k,\mathbf{u}}[1]$ is an upper bound for $c_{\mathbf{u},i}^*$, and these values are updated after the capacity explore phase in every epoch.

Let $\tau_{l,k,\mathbf{u}}$ be the set of time-slots in which the channel-state \mathbf{q} observed is such that $\pi(\mathbf{q})=k$ and the direction vector specified is \mathbf{u} , in the capacity explore phase of epoch l. In all these time-slots, the rate $\frac{C_{k,\mathbf{u}}[0]+C_{k,\mathbf{u}}[1]}{2}\mathbf{u}$ is scheduled. $m_{k,\mathbf{u}}$ denotes the empirical mean of the success rates in scheduling the above rate vectors. The lower and upper bounds $C_{k,\mathbf{u}}[0]$ and $C_{k,\mathbf{u}}[1]$ are then updated depending on the value of $m_{k,\mathbf{u}}$ for all k,\mathbf{u} . The update procedure is detailed in Algorithm 3, which is similar to a traditional binary search procedure (or bisection method) for searching the boundary of the capacity regions in the given directions \mathcal{U} (see also [5]).

D. Exploitation

In every epoch, after the *exploration* phases, the overwhelming majority of time-slots are dedicated to *exploitation*. The *exploitation phase* in epoch l consists of $s(l) = O(\sqrt{l})$ time-slots. In each of these time-slots, a channel-state \mathbf{q} is observed and a direction vector \mathbf{u} is specified. The class $k = \pi(\mathbf{q})$ is identified according to our current global classifier and the rate vector $C_{k,\mathbf{u}}[0]\mathbf{u}$ is scheduled. This phase is detailed as steps 24 - 29 in Algorithm 4.

Remark 5: While Algorithm 4 satisfies the regret bound in Theorem 1, there are a few low-probability failure events for our algorithm. For instance, there is a small probability that the initial classifier tree-structure is incorrect (i.e., a failure event in Lemma 1). Similarly, at any epoch with a small probability, the capacity explore update (by Algorithm 3) can be incorrect (due to differences in empirical and true success probability). We have developed a robust version of our algorithm, which detects such low-probability failure events and corrects them online (please see Supplementary material). For the simulations in Section VI, we use the robust version of our algorithm.

V. REGRET BOUND

In this section, we provide our main theoretical result which provides a cumulative regret bound for Algorithm 4, when Assumptions 1-5 are satisfied.

Theorem 1: Under Assumptions 1-5, with probability at least $1 - \xi K D \delta$, Algorithm 4 achieves a regret bound of,

$$R(T) = \mathcal{O}\left(T^{2/3}\log\left(\frac{1}{\delta}\right)\left(D\log T + K + \sqrt{V}\right)\right),$$

at time T where $\xi < 13$.

```
Algorithm 4 Online Rate Allocation From Channel-State and Service Data
```

```
1: Initialize empty sets S_i = \{\} for i \in [K].
2: Initialize a single node tree \mathcal{T} where the node contains n_0
        different channel-state points.
3: Initialize capacity rate C_{k,\mathbf{u}}[0] = 0 and C_{k,\mathbf{u}}[1] = C for
        all k \in [K] and \mathbf{u} \in \mathcal{U}.
4: Initialize classifier \pi using Algorithm 2.
5: Set t = t_0 (time index) and l = 1 (epoch index).
6: while t \leq T do
        for i = 0 : K - 1 do
7:
8:
           \mathbf{r}_i is the rate vector stored in node \mathcal{N}_i.
           Let \mathbf{q} be the channel-state observed at time step t.
           Schedule rate \mathbf{r}_i. (Class Explore)
10:
11:
           Let y \in \{0,1\} be the notification received. Add
               (\mathbf{q},y) to S_i.
           Set t = t + 1.
12:
           Update the classifier \hat{\pi}_i in \mathcal{N}_i.
13:
        end for
14:
        Let the empirical means of success rate be m_{k,\mathbf{u}}=0
15:
           for all k \in [K] and \mathbf{u} \in \mathcal{U}.
        for s = 1 : \alpha(\delta, l) do
16:
           Observe (q, u).
17:
           Let k = \pi(\mathbf{q}).
18:
           Schedule rate vector \left(\frac{C_{k,\mathbf{u}}[0]+C_{k,\mathbf{u}}[1]}{2}\right) u. (Capacity
19:
           Update m_{k,\mathbf{u}} according to received notification y.
20:
           Set t = t + 1.
21:
        end for
22:
        Update C and \hat{S} according to Algorithm 3.
23:
24:
        for s = 1 : s(l) do
25:
           Observe (q, u).
           Let k = \pi(\mathbf{q}).
26:
           Schedule rate vector C_{k,\mathbf{u}}[0]\mathbf{u}. (Exploit)
27:
28:
           Let t = t + 1.
29:
        end for
        l = l + 1.
31: end while
```

Theorem 1 and its proof is available in greater detail in the Supplementary Material of this paper, where the explicit dependence on the various problem parameters has been specified.

Discussion: Theorem 1 states that the regret of Algorithm 4 scales as $\mathcal{O}(T^{2/3}\log T)$ as a function of time. The scaling is linear with respect to the number of classes K and the number of direction vectors D. It scales as \sqrt{V} in terms of the VC dimension of the class of experts. For a finite class of experts $\hat{\Pi}$, the VC dimensions is $\mathcal{O}(\log N)$, where $N=|\hat{\Pi}|$ is the number of experts.

It should be noted that *epoch-greedy* algorithms in bandit settings generally have a regret scaling of $O(T^{2/3})$ in the problem independent setting, because of explicit exploration. For instance, the epoch-greedy strategy in [13] has a similar regret scaling for the problem of stochastic contextual bandits with experts. However, we would like to highlight that our problem setting is significantly more complicated than the usual contextual bandits with experts problem, as in

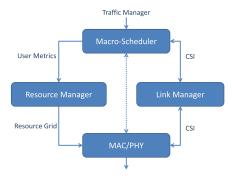


Fig. 3. Block diagram of scheduler module in WiNGS.

a contextual bandit setting when an arm is pulled under a context, we get a direct feedback about the reward of that arm under that context. However, in our problem setting when a channel-state is observed and a rate vector is scheduled, the received feedback only gives us a partial noisy feedback about the possible capacity class in which the channel-state belongs. The quality of the feedback also depends on the choice of the rate points. Further in our problem setting, even after the capacity classes are learned there is an additional task to recover the boundaries of the corresponding capacity regions. Therefore, the epoch-greedy algorithm proposed in this paper is a first step towards analyzing this setting, and we leave the study of algorithms with implicit exploration that can potentially yield $O(\sqrt{T})$ regret bound as future work.

VI. SIMULATION RESULTS

In this section we perform empirical simulation of our algorithm on the state-of-the-art Wireless Next-Generation Simulation (WiNGS) platform developed by AT&T Labs.

WiNGS includes a fully dynamic, event-driven system-level simulator which models large-scale cellular network deployments and the L3/L2/L1 protocols layers comprising the 5G New Radio (NR) air interface. Both planned and random deployments of base stations is supported, with users located indoors or outdoors generating traffic according to various finite and full-buffer traffic models. Packets generated by the traffic model pass through the Packet Data Convergence Protocol (PDCP), Radio Link Control (RLC), and Media Access Control (MAC) sublayers where functions including segmentation, re-transmissions, and HARQ (Hybrid Automatic RepeatreQuest) processes are implemented. The wireless channel is modeled with both long-term effects (e.g. log-normal shadowing, Line-of-Sight vs. Non Line-of-Sight pathloss) and short-term effects (e.g. fast fading due to the environmental scattering and user mobility). The physical layer functionality includes transport block formation based on link adaptation based on channel-state and ACK/NACK feedback. Codebook and channel reciprocity-based digital beamforming is used to generate linear precoders for both single user (SU) and multi-user MU-MIMO transmission on orthogonal or overlapping time/frequency resources. The probability of success for a transport block which is sent over the wireless channel is based on the post-MMSE receiver SINR to BLER (BLock Error Rate) mapping curves calculated from bit-precise link-level simulation.

Figure 3 provides a block diagram of the modules relevant to our simulations. The AT&T WiNGS scheduler runs in discrete time steps of size 1 ms. At the start of each time step, the traffic

manager sends the set of schedulable users and their user metrics to the macro-scheduler. The user metrics consist of information such as CQI (Channel Quality Index, generated every 10 ms), MIMO rank, NDI (New Data Indicator), set of pairable users (users that can be co-scheduled) and their corresponding MU-SINR, etc. A primary user is selected from the set of schedulable users by the macro-scheduler based on a round-robin policy. The macro-scheduler then picks a secondary user for MU-MIMO transmission from set of candidate pairable users. The macro-scheduler then passes the user-pair (primary user and secondary user) and their user-metrics to the resource manager. In the event the macro-scheduler is unable to find a secondary user to pair with primary user then only the primary user and its user metrics are sent to the resource manager. Furthermore it should be noted that for any failed MU transmission, the re-transmissions are sent in SU mode.⁸

For the selected user-pair (or user), the resource manager selects the MCS (Modulation and Coding Scheme) to be used according to the MU-SINR (or SU-SINR). As discussed in Remark 1, this map from channel quality (MU-SINR) to user rates (MCS) is adaptive, and is determined online by our algorithm. Specifically, for a given user the mapping of MU-SINR to MCS adapts according to the success/failure of the past transmissions. After selecting the MCS, resource manager fills the resource grid and sends it to the MAC/PHY layer to be scheduled.

To implement our algorithm, we have modified the default resource manager by selecting the MCS to be scheduled for the users according to our algorithm. We use a threshold value of $\frac{9}{10}$ instead of $\frac{1}{2}$ for $m_{k,\mathbf{u}}$ in Algorithm 3 since wireless carriers like AT&T requires probability of success to be greater than 90% for all users (desired service requirement). For all of our simulations, we use a robust version of our algorithm with few modifications. Specifically, the class of experts/classifying functions used in our simulations is the naive Bayes implemented in MATLAB. For building the classifier, we test 15 different MCS pairs (rate-vectors) for $l_0 = 50$. The value of n_0 is set to be 3000 (3 sec) and during this time we schedule according to AT&T resource manager. For sake of clarity in results we have excluded the first n_0 time steps from the results. Furthermore, the traffic model used in all simulations is full buffer i.e. all users have data to be transmitted at all time steps. We finally refer to Appendix F (in Supplementary Material) for more discussion on practical issues.

A. Clusters and Associated Capacity Regions

We first plot the result for a single run of our algorithm for a system with one base-station and n=4 users and d=5 direction vectors for each user-pair, i.e., 30 unique direction vectors. At any time-step the direction vector is selected uniformly at random from this set of 30 possible direction vectors. The probability of success and the throughput per resource block is plotted in Figure 4a and 4b respectively. We observe that our algorithm is able to converge to optimum within 10 sec and achieve probability of success higher then 90%. In Figure 4c and 4d we plot the probability of success and

⁷Users that have data to be transmitted.

⁸A failed MU transmission means either one or both users were unable to decode the packet. Every failed MU transmission can cause 1 or 2 SU re-transmissions based on if MU-transmission to one or both users was unsuccessful.

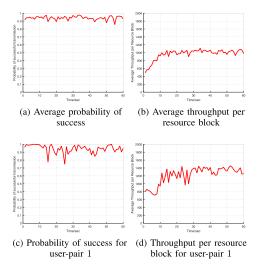


Fig. 4. Single run of our algorithm for a system with 4 users.

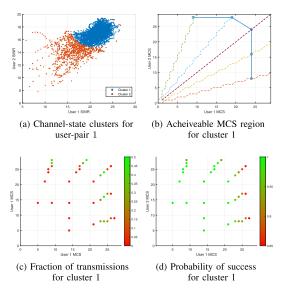


Fig. 5. Single cell results for our algorithm for user-pair 1.

throughput per resource block for user-pair 1 in the system. We observe that similar to before, our algorithm is able to learn the optimal throughput for user-pair 1 within 10 sec and achieve probability of success higher than 90%.

Figure 4a to 4d show that our algorithm operates in two distinct phases. The first phase, which we will refer to as the "explore" phase, consists of (*Initializing Classifier*) along with the first few epochs of our algorithm, where the algorithm explores for the majority of time steps. During the "explore" phase, the performance of the algorithm gradually improves, as it learns the system over time (first 10 sec for this scenario). The second phase, which we will refer to as "exploit" phase, consists of the epochs where our algorithm exploits for a majority of the time steps. The performance during the "exploit" phase remains relatively stable apart from small fluctuations due to channel variations (slow fading). It should be noted that due to averaging (over users), the variations observed in Figure 4a and 4b are smaller than the variations in Figure 4c and 4d during "exploit" phase.

In Figure 5a, we plot the channel-state clusters determined by our algorithm for user-pair 1. We observe that channel-state clusters determined by our algorithm are disjoint and cluster 1 contains better channel-states values than cluster 2. We plot

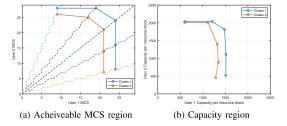


Fig. 6. Achievable MCS region and Capacity region for the different cluster of user-pair $1.\,$

the achievable MCS values of cluster 1 in Figure 5b where the dotted lines show the 5 different direction vectors. The fraction of time different MCS values are scheduled and their corresponding success rate are plotted in Figure 5c and 5d respectively. We observe that while our algorithm explores several different MCS values for any given direction vector, the MCS value that lies on the boundary of achievable MCS region⁹ is selected most often.

In Figure 6a, we plot the achievable MCS region for the two channel-state clusters. We observe that the achievable MCS region belonging to cluster 2 is subset of achievable MCS region of cluster 1 because cluster 1 contains better channel-state values as compared to cluster 2. In Figure 6b we plot the capacity regions of both clusters. It should be noted that the capacity region are slightly different from achievable MCS regions since the two users can have different probability of success for any given MCS value.

B. Single Cell

We consider a single base-station setting and analyze the performance of our algorithm against the state of the art AT&T scheduler. The parameter settings for the AT&T scheduler needs to be manually optimized for every scenario which is impractical in real deployments. Therefore, following practice in real deployments, for the simulations we consider 3 different parameter settings denoted as Policy 1, Policy 2 and Policy 3 respectively, where theses policies are hand-tuned to provide good results for a majority of scenarios. For our algorithm we have a set of d=11 direction vectors for each user-pair and at every time step, the direction vector selected by the scheduler is the direction vector which is closest to the MCS selected by the default AT&T scheduler. 10

For this setting we first consider 3 different scenarios with a single active base-station and n=4 users each. These scenarios differ in that the user locations are different (and thus, the channel and interference environments differ). For each scenario, we provide plot results for number of MU transmissions, probability of success, throughput per resource block and throughput per second. The "explore" phase typically lasts between 5 to 30 sec (depending on the number of user-pairs in the system); further, the performance of our algorithm during the "exploit" phase and that of the AT&T policies is relatively stable over time. Therefore we only display the first 1 minute of the scenario, even though a typical scenario in a practical setting can last for several minutes or longer.

The results for scenario 1 are shown in Figure 7. We observe that our algorithm is able to learn the channel-state clusters and suitable MCS values for different direction vectors within

⁹Those having a probability of success exceeding 90%.

¹⁰The MCS value selected by the AT&T scheduler is only used to select the direction vector. Our algorithm determines the MCS value to be scheduled along the selected direction vector.

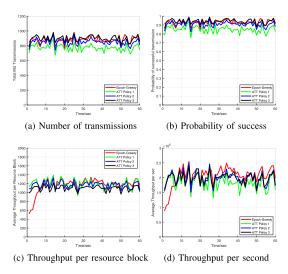


Fig. 7. Performance of different algorithms for Scenario 1 with 4 users. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more then 9.5%.

 ${\bf TABLE~I}$ Performance of Different Policies in Scenario 1 for Last 30 Sec

	Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3	
Average number of MU transmissions	887.7	767.8	845.9	892.6	
Average probability of success (%)	93.57	84.74	90.78	94.10	
Average throughput per resource block	1221	1184	1173	1100	
Average throughput per second (×10 ⁶)	2.170	1.824	1.988	1.976	
and the state of t					

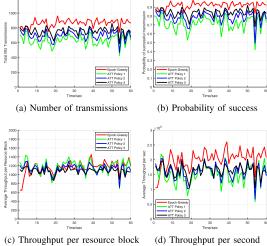


Fig. 8. Performance of different algorithms for Scenario 2 with 4 users. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more then 30%.

10 seconds. Furthermore when compared to the best performing AT&T policy (Policy 2), our algorithm is able to achieve 9.5% more throughput per second and improve the probability of success from 90.8% to 93.6%. Table I summarizes the performance from 30 to 60 sec for the different algorithms in scenario 1.

The results for the scenario 2 are shown in Figure 8. Similar to scenario 1 our algorithm is able to learn the channel-state

TABLE II
PERFORMANCE OF DIFFERENT POLICIES IN SCENARIO 2 FOR LAST 30 SEC

Average number of MU transmissions 876.0 634.0 698.5 756.8 Average probability of success (%) 92.82 71.68 78.37 83.85 Average throughput per resource block 1219 1172 1127 1083 Average throughput per resource block 2.130 1.500 1.500 1.646		Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3
of success (%) 92.82 /1.08 /8.3/ 83.85 Average throughput per resource block Average throughput		876.0	634.0	698.5	756.8
per resource block 1219 1172 1127 1083		92.82	71.68	78.37	83.85
Average throughput 2 120 1 500 1 500		1219	1172	1127	1083
per second ($\times 10^6$) 2.139 1.508 1.588 1.646	Average throughput per second $(\times 10^6)$	2.139	1.508	1.588	1.646

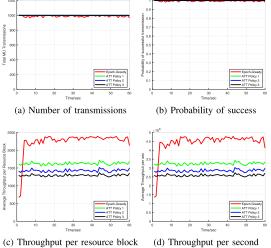


Fig. 9. Performance of different algorithms for Scenario 3 with 4 users. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more then 41%.

TABLE III
PERFORMANCE OF DIFFERENT POLICIES IN SCENARIO 3 FOR LAST 30 SEC

	Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3
Average number of MU transmissions	988.2	999.7	999.9	1000
Average probability of success (%)	99.4	100	100	100
Average throughput per resource block	2335	1633	1435	1299
Average throughput per second ($\times 10^6$)	4.615	3.265	2.869	2.597

clusters and optimum MCS values for different direction vectors within 10 sec. However unlike scenario 1 the AT&T policy which provides the best average throughput per second is Policy 3. Once again our algorithm outperforms the best AT&T Policy and achieves 30% more throughput per second and improve probability of success from 83.9% to 92.8%. Table II summarizes the performance from 30 to 60 sec for different algorithms in scenario 2.

Figure 9 provides the results for the scenario 3. Similar to previous two scenarios, our algorithm is quickly able to learn the channel-state clusters and optimum MCS values for different direction vectors within 10 sec. For this scenario the AT&T policy which provides the best average throughput per second is Policy 1. Similar to the previous two scenarios, our algorithm significantly outperforms Policy 1 and achieves 41% more throughput per second while still providing

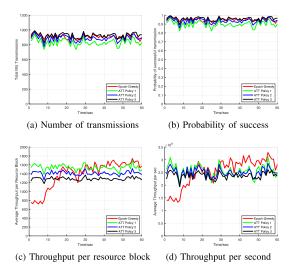


Fig. 10. Performance of different algorithms for Scenario 4 with 7 users. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more then 14.5%.

TABLE IV $\label{eq:performance} \textbf{PERFORMANCE OF DIFFERENT POLICIES IN SCENARIO 4 FOR LAST 30 SEC}$

	Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3
Average number of MU transmissions	908.4	829.8	886.4	926.9
Average probability of success (%)	94.91	89.62	93.52	96.01
Average throughput per resource block	1603	1530	1408	1287
Average throughput per second ($\times 10^6$)	2.913	2.543	2.497	2.386

probability of success of 99.4%. Table III summarizes the performance of different algorithms in scenario 3 for last 30 sec.

For scenario 4 we consider a single base-station with n=7 users. The results for the scenario 4 are shown in Figure 10. For this scenario, our algorithm takes longer time to learn the channel-state clusters and optimum MCS values for different direction vectors since there are 21 user-pairs in scenario 4 (as compared to 6 for scenarios 1-3). However our algorithm is able to learn optimum value of MCS for different direction vectors within 30 sec. Among the static policies, AT&T Policy 1 provides the best average throughput per second for scenario 4, however our algorithm is able to achieves 14.5% more throughput per second and improve the probability of success from 89.6% to 94.9%. Table IV summarizes the performance of different algorithms in scenario 4 during last 30 sec.

C. Multiple Cell

We now consider the multiple base-station setting and analyze the performance of our algorithm against the state of the art AT&T scheduler. For this setting we activate 4 base-stations in the neighborhood of primary base-station. For the AT&T scheduler, we have 3 new parameter setting 11 denoted as Policy 1, Policy 2 and Policy 3. For the simulations with AT&T scheduler we use the same policy on all 5 base-stations. For simulations with our algorithm we only run our

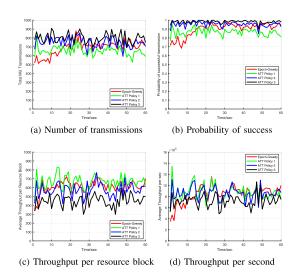


Fig. 11. Performance of different algorithms for Scenario 5 with 4 users connected to primary base-station. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more then 5%.

TABLE V
PERFORMANCE OF DIFFERENT POLICIES IN SCENARIO 5 FOR LAST 30 SEC

	Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3
Average number of MU transmissions	718.9	641.8	731.7	791.0
Average probability of success (%)	93.74	86.17	94.97	98.08
Average throughput per resource block	629.2	674.5	589.3	461.5
Average throughput per second $(\times 10^5)$	9.041	8.601	8.604	7.286

algorithm on primary base-station and run AT&T Policy 2 on the 4 neighboring base-stations. In the following, we compare the performance of primary base-station for the different algorithms. Furthermore for our algorithm we have set of d=11 direction vectors for each user-pair, where the direction vectors are selected using the same method as described for single cell simulations.

We present results for 3 different scenarios denoted by scenario 5, scenario 6 and scenario 7 with 4, 4 and 3 users connected to the primary base-station respectively. For all these scenarios we observe that the average throughput is significantly reduced due to inter-cell interference as compared to average throughput for single cell scenarios.

The results for scenario 5 are shown in Figure 11. We observe that our algorithm is able to learn the channel-state clusters and learn optimum MCS values for different direction vectors within 15 sec. Furthermore, the AT&T policy which provides the best average throughput per second is Policy 2, however our algorithm is able to achieve 5% more throughput per second and have high probability of success of 93.74%. Table V summarizes the performance from 30 to 60 sec for the different algorithms in scenario 5.

The results for scenario 6 are presented in Figure 12. We observe that similar to scenario 5, our algorithm is able to learn the channel-state clusters and optimum MCS values for different direction vectors within 15 sec. Furthermore our algorithm is able to achieve 90% more throughput per second than the best performing AT&T policy while having a high probability of success of 92.35%. Table VI summarizes the

¹¹These policies are different from the previous policies since there is interference from the neighboring base-stations.

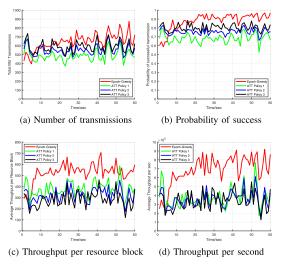


Fig. 12. Performance of different algorithms for Scenario 6 with 4 users connected to primary base-station. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more then 90%.

TABLE VI
PERFORMANCE OF DIFFERENT POLICIES IN SCENARIO 6 FOR LAST 30 SEC

	Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3
Average number of MU transmissions	694.0	501.6	558.9	598.5
Average probability of success (%)	92.35	69.19	75.92	80.87
Average throughput per resource block	548.0	385.6	331.1	290.3
Average throughput per second $(\times 10^6)$	7.598	3.925	3.733	3.502

TABLE VII
PERFORMANCE OF DIFFERENT POLICIES IN SCENARIO 7 FOR LAST 30 SEC

	Epoch- greedy	AT&T Policy 1	AT&T Policy 2	AT&T Policy 3
Average number of MU transmissions	525.5	503.6	563.5	580.6
Average probability of success (%)	93.35	88.97	97.31	99.27
Average throughput per resource block	212.2	199.3	146.3	97.5
Average throughput per second ($\times 10^6$)	2.243	2.025	1.657	1.139

performance from 30 to 60 sec for the different algorithms in scenario 6.

Finally we present the results for scenario 7 in Figure 13. Unlike scenario 5 and scenario 6 our algorithm is able to learn the channel-state clusters and optimum MCS values for different direction vectors within 10 sec because there are fewer users connected to primary base-station. For this scenario the AT&T policy which provides the best average throughput per second is Policy 2, however our algorithm achieves 10% more throughput per second and have a high probability of success of 93.35%. The performance from 30 to 60 sec for the different algorithms in scenario 7 is summarized in Table VII.

D. Summary of Results and Discussion

We observe that in all scenarios our algorithm is able to learn the channel-state clusters and their corresponding

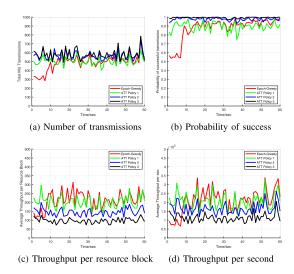


Fig. 13. Performance of different algorithms for Scenario 7 with 3 users connected to primary base-station. Our algorithm outperforms the current state-of-art and improves the average throughput per sec by more than 10%.

capacity regions in a short amount of time (10 to 15 sec for a 4 user system, 30 sec for a 7 user system). Furthermore, for both single cell and multi cell scenarios, our algorithm is able to match if not outperform the current state of art MU scheduling algorithm. It should be noted that the best performing static policy for the AT&T scheduler to maximize the throughput per sec differs across scenarios. Since we expect the scenario to dynamically change over time, AT&T policies would need to periodically be hand-tuned if one wishes to achieve scenario-specific optimal results. However, our algorithm is able to learn the correct channel-state to MCS mapping for different user-pairs, which allows it to match if not outperform the current state-of-art policy for all scenarios (even in cases like scenario 3 and scenario 6 shown in Figures 9 and 12, where all AT&T policies failed to achieve good throughput). We however emphasize that the real benefit of our approach is not in beating the best AT&T policy (after all, one could create a new hand-tuned policy), but to have an adaptive policy that does not require expensive hand-tuning for each scenario.

For a practical setting, an important metric used by wireless carriers like AT&T is the probability of success of a transmission, where an ideal policy should achieve more then 90% probability of success for all users. Our algorithm is able to ensure a high probability of success within 92% to 95% for all scenarios. Scenario 3 is an outlier with a high probability of success at 99.4% because the channel-state of users is very good and we are able to transmit the maximum value of MCS successfully with a high probability. Further, we can choose the desired probability of success of the system, by controlling the threshold value for $m_{k,\mathbf{u}}$ in Algorithm 3. For any fixed AT&T policy the probability of success not only varies significantly for different scenarios but is also below 90% for several scenarios. One could mitigate this by using a static policy (like Policy 3) which has better probability of success compared to other AT&T policies, however, the trade-off is that the throughput will be significantly lowered for some scenarios. Moreover, there are cases like scenario 2 or scenario 6 where even this is not possible, as the probability of success for all three AT&T policies is significantly below 90%. In conclusion, our algorithm not only achieves high throughput but also meets

the desired service requirement (probability of success) goal for all scenarios.

Finally, we comment on the applicability of our algorithm in real settings. Recall that our algorithm starts cold, and requires up to 30 seconds to ramp up, to match or outperform static policies. If users remain in the system for a very short amount of time (say, order of few seconds), then it would likely be more efficient to simply use a static policy. However, a typical scenario in a cellular network can potentially last from several minutes to much longer. For instance, if the same user has different sessions (e.g. multiple video streams viewed in sequence, or a user browsing multiple web-pages for an extended period of time), the learning "transfers" across these sessions, and the costs due to the early phase is amortized over time (see also Appendix F). Further, by also clustering users based upon their similar characteristics, we would be able to further reduce the complexity of online learning (and the "explore" phase duration). Studying such user-level clustering opens up interesting directions for future work.

ACKNOWLEDGMENT

The authors would like to sincerely thank Dr. Arunabha Ghosh for discussions on the model as well as providing extensive access to the AT&T WiNGS platform. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

REFERENCES

- I. Tariq, R. Sen, G. D. Veciana, and S. Shakkottai, "Online channel-state clustering and multiuser capacity learning for wireless scheduling," in *Proc. IEEE INFOCOM*, Paris, France, Apr. 2019, pp. 136–144.
- [2] Y. Du and G. de Veciana, "'Wireless networks without edges': Dynamic radio resource clustering and user scheduling," in *Proc. IEEE INFO-COM*, Apr. 2014, pp. 1–9.
- [3] R. Srikant and L. Ying, Communication Networks: An Optimization, Control, and Stochastic Networks Perspective. Cambridge, U.K.: Cambridge Univ. Press, 2014.
- [4] A. Agarwal, M. Dudík, S. Kale, J. Langford, and R. Schapire, "Contextual bandit learning with predictable rewards," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2012, pp. 19–26.
- [5] P. W. Goldberg and S. Kwek, "The precision of query points as a resource for learning convex polytopes with membership queries," in *Proc. Conf. Learn. Theory (COLT)*, 2000, pp. 225–235.
- [6] D. Avis and K. Fukuda, "A pivoting algorithm for convex hulls and vertex enumeration of arrangements and polyhedra," *Discrete Comput. Geometry*, vol. 8, no. 3, pp. 295–313, Sep. 1992.
- [7] G. Wunder, M. Kasparick, A. Stolyar, and H. Viswanathan, "Self-organizing distributed inter-cell beam coordination in cellular networks with best effort traffic," in *Proc. 8th Int. Symp. Modeling Optim. Mobile, Ad Hoc Wireless Netw. (WiOpt)*, May/Jun. 2010, pp. 295–302.
- [8] W. Yu, T. Kwon, and C. Shin, "Multicell coordination via joint scheduling, beamforming, and power spectrum adaptation," *IEEE Trans. Wireless Commun.*, vol. 12, no. 7, pp. 1–14, Jul. 2013.
- [9] T. Yoo and A. Goldsmith, "On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 3, pp. 528–541, Mar. 2006.
- [10] X. Xie and X. Zhang, "Scalable user selection for MU-MIMO networks," in *Proc. IEEE INFOCOM*, Apr. 2014, pp. 808–816.
- [11] D. Gesbert, M. Kountouris, R. W. Heath, C.-B. Chae, and T. Salzer, "From single user to multiuser communications: Shifting the MIMO paradigm," *IEEE Signal Process. Mag.*, vol. 24, no. 5, pp. 36–46, Sep. 2007.
- [12] A. Slivkins, "Contextual bandits with similarity information," in *Proc.* 24th Annu. Conf. Learn. Theory, 2011, pp. 679–702.

- [13] J. Langford and T. Zhang, "Epoch-Greedy algorithm for multi-armed bandits with side information," in *Proc. Adv. Neural Inf. Process. Syst.*, 2008, pp. 817–824.
- [14] S. Bubeck and N. Cesa-Bianchi, "Regret analysis of stochastic and non-stochastic multi-armed bandit problems," *Found. Trends Mach. Learn.*, vol. 5, no. 1, pp. 1–122, 2012.
- [15] A. Agarwal, D. Hsu, S. Kale, J. Langford, L. Li, and R. Schapire, "Taming the monster: A fast and simple algorithm for contextual bandits," in *Proc. Int. Conf. Mach. Learn.*, 2014, pp. 1638–1646.
- [16] M. Dudik et al., "Efficient optimal learning for contextual bandits," 2011, arXiv:1106.2369. [Online]. Available: http://arxiv.org/abs/ 1106.2369
- [17] S. Kwek and L. Pitt, "PAC learning intersections of halfspaces with membership queries," *Algorithmica*, vol. 22, nos. 1–2, pp. 53–75, Sep. 1998.
- [18] V. N. Vapnik and A. Y. Chervonenkis, "On the uniform convergence of relative frequencies of events to their probabilities," in *Measures of Complexity*. Cham, Switzerland: Springer, 2015, pp. 11–30.
- [19] A. K. Menon, B. van Rooyen, and N. Natarajan, "Learning from binary labels with instance-dependent corruption," 2016, arXiv:1605.00751. [Online]. Available: http://arxiv.org/abs/1605.00751



Isfar Tariq (Graduate Student Member, IEEE) received the B.S. degree in electrical engineering from the Lahore University of Management Sciences. He is currently pursuing the Ph.D. degree with the ECE Department, The University of Texas at Austin. His research interests include algorithm design for resource allocation, online learning, and wireless communication networks.



Rajat Sen (Member, IEEE) received the Ph.D. degree in electrical and computer engineering from The University of Texas at Austin in 2019. He is currently a Research Scientist at Google Research, Mountain View, CA, USA. His current research interests include online learning, bandit algorithms, time-series forecasting, and black-box optimization.



Thomas Novlan (Member, IEEE) received the B.S. degree (Hons.), M.S., and Ph.D. degrees in electrical engineering from The University of Texas at Austin in 2007, 2009, and 2012, respectively. In 2016, he joined AT&T, working on 5G technologies, such as dynamic spectrum sharing (DSS) and integrated access and backhaul (IAB). He is currently a Principal Member of Technical Staff with the Advanced Wireless Technology Group, AT&T Labs, Austin, TX, USA. He is working on 5G and beyond air-interface research and system design, including

applications of machine learning embedded in the radio access network.



Salam Akoum (Member, IEEE) received the Bachelor of Engineering degree in computer and communications engineering with a minor in mathematics (Hons.) from the American University of Beirut, Lebanon, in 2006, the Master of Science degree in electrical engineering from The University of Utah in 2008, and the Ph.D. degree in electrical engineering from The University of Texas at Austin in 2012. She is currently a Principal Member of Technical Staff with the Advanced Radio Technology Group, AT&T Labs. She is also an Active

Contributor to 3GPP 5G standardization, responsible for specifications of the physical-layer of the radio interface. She works on algorithm development and system architecture and optimization for current and next-generation wireless networks.



Milap Majmundar (Senior Member, IEEE) received the M.S.E.E. degree from Virginia Polytechnic Institute and State University and the M.B.A. degree from The University of Texas at Austin. He leads the Advanced Wireless Technology Group, AT&T Labs, responsible for research and development of techniques for next generation wireless networks. He holds about 70 granted patents in the area of wireless communications. He is currently interested in a range of research topics, including resource allocation, flexible network architectures, and spectrum usage.



Sanjay Shakkottai (Fellow, IEEE) received the Ph.D. degree from the ECE Department, University of Illinois at Urbana-Champaign in 2002. He is with The University of Texas at Austin, where he is currently the Temple Foundation Endowed Professor No. 4 and a Professor with the Department of Electrical and Computer Engineering. His research interests lie at the intersection of algorithms for resource allocation, statistical learning and networks, with applications to wireless communication networks, and online platforms. He received the NSF CAREER Award in 2004.



Gustavo de Veciana (Fellow, IEEE) received the B.S., M.S., and Ph.D. degrees in electrical engineering from the University of California at Berkeley in 1987, 1990, and 1993, respectively. He is currently a Professor and the Associate Chair of the Department of Electrical and Computer Engineering. From 2003 to 2007, he served as the Director and the Associate Director of the Wireless Networking and Communications Group (WNCG), The University of Texas at Austin. His research focuses on the design, analysis and control networks, information theory,

and applied probability. His current interests include measurement, modeling and performance evaluation, wireless and sensor networks, architectures and algorithms to design reliable computing, and network systems. In 2009, he was a Designated IEEE Fellow for his contributions to the analysis and design of communication networks. He was a recipient of the Cockrell Family Regents Chair in engineering and the National Science Foundation CAREER Award 1996. He was a co-recipient of the six best paper awards, including the IEEE William McCalla Best ICCAD Paper Award for 2000 and the Best Paper in ACM Transactions on Design Automation of Electronic Systems from January 2002 to 2004. He has been an Associate Editor and an Editor at large for the IEEE/ACM TRANSACTIONS ON NETWORKING. He currently serves on the board of trustees of IMDEA Networks, Madrid.