# Joint Chance-Constrained Predictive Resource Allocation for Energy-Efficient Video Streaming

Ramy Atawia, *Student Member, IEEE*, Hatem Abou-zeid, *Member, IEEE*,
Hossam S. Hassanein, *Senior Member, IEEE*, and Aboelmagd Noureldin, *Senior Member, IEEE*

*Abstract*—Predictive resource allocation (PRA) techniques that exploit knowledge of the future signal strength along roads have recently been recognized as promising approaches to save base station (BS) energy and improve user quality of service (QoS). Recent studies on human mobility patterns and wireless signal strength measurements along buses and trains have indeed supported the practical potential of PRA. An unresolved challenge, however, is modeling the *uncertainty* in the predictions, and developing real-time robust solutions that incorporate *probabilistic* QoS guarantees. This is of paramount importance in PRA due to the prediction *time horizon* that adds considerable complexity and increases the rate uncertainty in the problem. With these developments in mind, this paper addresses energy-efficient PRA applied to stored video streaming using *chance constrained* programming. The proposed solution incorporates: 1) uncertainty in predicted user rates; 2) a joint level of probabilistic constraint satisfaction over a time horizon; and 3) both optimal gradient-based and real-time guided heuristic solutions. Our framework fundamentally differs from previous PRA work in the literature where nonstochastic approaches with assumptions of perfect prediction were primarily used to demonstrate the potential energy savings and QoS gains. Numerical simulations based on a standard compliant long term evolution (LTE) system are provided to examine and compare the developed solution. Unlike existing energy-efficient PRA, the proposed framework achieves the desired QoS level under imperfect channel predictions. This robustness is attained without compromising the energy-efficiency compared to opportunistic schedulers, and thus supports PRA implementation in practice.

*Index Terms*—Channel state prediction, energy efficiency, kalman filter, radio access networks, resource allocation, robustness, video streaming.

## I. Introduction

THROUGHOUT the past decade, global environmental changes have been driving policy makers to enforce stringent regulations on the wireless industry. To this end, research in green wireless communications is gaining momentum to reduce the electrical power consumption of wireless networks

[1]. Among the network elements, the Radio Access Network (RAN) accounts for more than 50% of the network energy consumption [2]. As such, designing novel energy-efficient RAN frameworks is paramount to reducing the network carbon footprint while satisfying the increasing Telecom market demands. This includes techniques such as efficient Power Amplifier (PA) design [1], cell switch off [3], [4], and traffic-aware scheduling, among others. A more efficient RAN is also beneficial for operators as it can postpone investment in equipment installations and new spectrum. Thus, in addition to minimizing the energy-related operational expenditures (OpEx), the capital expenditures (CapEx) can also be reduced since radio equipment installations can make up to 70% of CapEx [5]. Concurrently, mobile video traffic is experiencing tremendous growth and forecasted to account for about two thirds of the consumer traffic in 2017 [6]. To address these recent developments, this paper presents a novel approach toward *energy-efficient* wireless video streaming.

Radio signal measurement studies indicate that cellular network users moving along the same path will typically experience similar signal strength variations as reported in [7], [8]. Advanced navigation and channel prediction techniques [9]–[11] will also enable accurate calculation of user traces and future channel rates. *Predictive* resource allocation that exploits these patterns of signal strength over a time horizon has recently been recognized as a promising approach to improve video streaming QoS [12], [13], and transmission energy [14], [15]. This is accomplished by leveraging the knowledge of the future link capacity users are expected to experience, and then performing *long-term* predictive Resource Allocation (RA) plans over several seconds. By doing so, BSs can prioritize users headed to poor channel conditions, or delay transmission until a user reaches better channel conditions. Stored video content such as YouTube and Netflix is well suited for such approaches as it can be strategically prebuffered and stored on the local cache of the User Equipment (UE).

The potential energy saving gains of PRA reported in recent literature [12]–[16] are very encouraging, and demonstrate the need for further investigation. The initial works on PRA primarily used ideal predicted data rates, and demonstrated the potential gains of such proactive mechanisms. However, such formulations depend on the *average* value of future data rates, and thus, are not robust to channel variations. Therefore, QoS satisfaction is not guaranteed under uncertain channel predictions. In this paper, we address this problem and developed a *robust* framework that provides *probabilistic* QoS guarantees. This enables network operators to prioritize users and

applications by offering a mechanism to control the probability of constraint satisfaction. Furthermore, by modeling uncertainty, the framework can strike a balance between providing high energy efficiency gains when predictions are accurate, and minimizing the risks associated with erroneous predictions during periods of uncertainty. The main contributions of the paper can be summarized in the following:

- We develop a *robust* PRA framework that accounts for rate uncertainties and provide QoS guarantees over a *time horizon*, with the objective of minimizing energy consumption. As recent practical and theoretical findings indicate that the variations in predicted rates can be modeled as multivariate normal random numbers [17], we employ probabilistic Joint Chance Constrained Programming (JCCP) to formulate the problem mathematically and then obtain its deterministic closed form.

- We then show that the resultant JCCP formulation is non-convex and apply proportional risk allocation for joint chance constraints. The problem is decomposed into two convex sub-problems, where the first stage optimizes the *individual* risk levels at each time slot, which are subsequently used by the second stage to solve the robust RA problem. By applying such a *non-uniform* risk allocation, we generalize the solution to achieve less conservative (i.e, energy-efficient) and more practical QoS aware RA decisions.

- Although the resultant two stage formulation is proven to be convex and double differentiable, it's solution is accomplished at a high complexity that may not meet real-time requirements. Thus, we develop an efficient low complexity guided search heuristic that guarantees the satisfaction of joint QoS levels.

- Extensive measurements in [8] indicate that the *variance* in the predicted channel measurements changes with the time of the day and geographical area. Hence, modeling the random rates with constant variances would result in sub-optimal results. Due to the inconsistency in the rate variance over time and location, we adopt Kalman Filter (KF) to accurately track such variations, providing an additional degree of robustness to the statistical parameters. With such a framework, QoS guarantees can be ensured during high variance while energy minimization is achieved during low varying cases.

In the following section, we provide a background to JCCP and review the related literature. Section III presents the problem statement and notations, while Section IV discusses the JCCP formulation for the robust PRA, as applied to energy-efficient video streaming. In Section V, we develop the solution methods using both a gradient-based approach and a guided heuristic for real-time allocations. The design of the Kalman filer is thereafter presented in Section VI, and numerical results are discussed in Section VII. Finally, we conclude the paper in Section VIII.

## II. BACKGROUND AND RELATED WORK

In wireless channel prediction, the future rates can not be perfectly predicted and thus typically modeled as random variables. The traditional approach in wireless PRA strategies

[12]–[15] is to replace each of these random variables by the expected (average) value and solve the resulting deterministic optimization problem. However, this approach results in *non-robust* and suboptimal allocations as the probability of experiencing a lower or higher data rate than the expected values are totally ignored. In particular, experiencing a lower rate than the average value will make the allocated resources insufficient to deliver the future demand, causing QoS dissatisfaction. On the other hand, if users experience higher values than the average, excessive resources may have been allocated to satisfy the demand, resulting in suboptimal resource utilization and energy savings. With this trade-off between QoS satisfaction and resource utilization, handling the errors in predicted rate and the channel variations during resource allocation is very challenging. To this end, robust stochastic optimization techniques have been introduced in which the predicted rate is modelled as a random variable rather than its mean value [18]. The variance and the probability density function (PDF) account for the cases in which the actual rate fluctuates above or below the mean value. The formulation in this case incorporates Chance Constrained Programming (CCP) [19] that can guarantee the satisfaction of user QoS at a certain level $\beta \in [0, 1]$. In essence, a chance constraint can be formulated as

$$Pr\{F(x_t, \eta_t) \geq D_t\} \geq \beta, \qquad \forall t \in \mathcal{T}, \qquad (1)$$

where $x_t$ is the resource allocation variable at time slot $t$, and $\eta_t$ denotes the random data rate. The function $F(x_t, \eta_t)$ models the relation between $x_t$, $\eta_t$ and the demand $D_t$ for each time slot $t$ in the time horizon $\mathcal{T}$. The above formulation guarantees that the allocation at each time slot satisfies the corresponding demand with at least probability $\beta$. This represents the QoS level, where a higher value results in allocating more resources (i.e., more energy consumption) to ensure demand satisfaction.

However, such form of chance constraint can only guarantee the QoS satisfaction level during each time slot, and does not model the satisfaction over the *time horizon*. In particular, allocating less resources in one time slot will result in the demand dissatisfaction in both the current and the future instances. Thus, satisfying $\beta\%$ of the demand of one time slot will not guarantee the same satisfaction degree in the coming time slot, since the latter does not account for the partial satisfaction in the prior slots. This is because the demand across the time slots is cumulative and allocation should be able to recover from outages in the previous slots. To avoid the propagation of such outages, allocation of all the time slots in the horizon should be jointly considered. This is typically done using Joint Chance Constrained Programming (JCCP) [20] and expressed mathematically as follows

$$Pr\{F(x_t, \eta_t) \geq D_t, \quad \forall t \in \mathcal{T}\} \geq \beta. \qquad (2)$$

JCCP has been successfully adopted in the literature to solve numerous networking problems where the decision made on one constraint affects the satisfaction of the others. Among these, application to routing and bandwidth assignment, and uplink resource allocation in OFDM networks [21] where the QoS satisfaction of one user might affect the others. In such models, the chance constraints are found to be independent and

their joint probability is simply the product of their individual probabilities. However, such an independence is not applicable in PRA since the constraints are no longer independent due to the cumulative demand at each time slot. Due to the difficulty of obtaining the pairs of joint probabilities, Boole's inequality [22] can be used to bound this joint probability. However, applying such a bound is very conservative and can result in suboptimal allocations that deteriorate the network optimization objective. Therefore, the individual probabilities of each constraint should be optimized to result in less conservative solutions. Example of applications that apply time dependent JCCP are model predictive control [23], [24] and the unit commitment in power generation systems [25] in which the demand is cumulative among the time slots and therefore joint satisfaction is needed. Individual probabilities of chance constraints can be determined optimally if the RA problem with unknown individual probabilities remains affine or convex, as in [26]. Otherwise, both individual probabilities and RA decisions are jointly determined using simulation based or iterative search techniques as in [25]. In summary, the joint chance constraint solves for two decision vectors: 1) the individual probabilities of each time slot QoS constraint, and 2) the resource allocation among the users. The former is subjected to Boole's inequality while the latter is subjected to user QoS satisfaction at each time slot in order to satisfy the overall QoS level over the time horizon.

The common challenge in both types of CCP is that the problem does not have a closed form solution when expressed in the form in (1) or (2). As such, the problem is either solved using simulation based approaches or analytical methods. In the former type, realizations of the random component are generated [18] and allocation is decided to satisfy $\beta^{th}$ percentile of the scenarios. On the other hand, analytical methods replace the chance constraints with the cumulative distribution function of the random variable. These methods are found to provide better accuracy when the inverse cumulative density function is invertible, unimodal and results in affine or convex optimization. Nevertheless, the simulation based methods remain as an alternative to provide an acceptable solution when the analytical approximation fails.

The previous *non-robust* works in PRA expressed the rate by its average value and did not include any form of uncertainty modeling or probabilistic QoS constraints [12]–[16]. Such a limitation is highlighted by recent practical studies on channel predictions [8], [17], [27] {which indicate that the errors associated with future rates should not be ignored and require exclusive handling by PRA. In particular, the authors in [28] consider maximizing the spectral efficiency while satisfying minimum user demands. However, the objective of minimizing energy consumption and ensuring joint QoS satisfaction over the time horizon was not addressed. A constant PDF was also assumed, and real-time mechanisms to track the channel variance were not developed. An initial study on modeling uncertainty for *energy-efficient* PRA was introduced in our prior work in [29]. Therein, fuzzy-based optimization was adopted which has the advantage of implementation simplicity, but resulted in conservative solutions that over satisfy the constraints at the expense of higher energy [30].

In this paper, we present a joint chance constraint based approach for *energy-efficient* PRA applied to stored video streaming where the satisfaction of the cumulative demand at all the time slots are jointly considered. As discussed in the introduction, the proposed solution incorporates: 1) uncertainty in predicted user rates, 2) a joint level of constraint satisfaction over a time horizon, 3) both gradient-based optimal and real-time guided heuristic solutions, and 4) adaptive tracking of variations in modeled random rates. It should be noted that the mobility-based rate prediction itself and statistical error modeling as reported in [8], [17], [27] are not the main objective of this paper.

## III. SYSTEM OVERVIEW

### A. Preliminaries

We use the following notational conventions throughout the paper: $\mathcal{X}$ denotes a set and its cardinality is denoted by $X$. Matrices are denoted with subscripts, e.g. $\mathbf{x} = (x_{a,b} : a \in \mathbb{Z}_+, b \in \mathbb{Z}_+)$, and matrix transpose and inverse are denoted as $\mathbf{x}'$ and $\mathbf{x}^{-1}$ respectively. $Pr\left(\bigcap_{\forall \mathcal{S}} s_i\right)$ and $Pr\left(\bigcup_{\forall \mathcal{S}} s_i\right)$ denote the joint and disjoint probabilities of all events in set $\mathcal{S}$. The gradient and Hessian of function $\mathbf{f}(\cdot)$ are denoted by $\nabla \mathbf{f}(\cdot)$ and $\nabla^2 \mathbf{f}(\cdot)$ in order. $\tilde{r}$ represents a random variable, whose probability density function follows normal distribution, while its cumulative density function is the Q function denoted as $Q$. The $n^{th}$ percentile of a zero mean and unit variance normally distributed random variable is denoted by $Q_{1-n}^{-1}$. $E[\cdot]$ denotes the expectation of a random variable.

### B. Problem Definition

Consider a BS with an active user set $\mathcal{M}$, where an arbitrary user is denoted by $i \in \mathcal{M}$. Users request stored video content from video servers. We assume that the wireless link is the bottleneck and the video content is always available at the BS. We assume that user's mobility trace is known for the next $T$ seconds, called the prediction window, and at a per second granularity. This results in a total of $T$ time slots within the prediction window, which we denote by the set $\mathcal{T} = \{1, 2, \ldots, T\}$. At each time slot $t \in \mathcal{T}$, the BS resources (airtime fractions) are shared among the active users. We define the resource allocation matrix $\mathbf{x} = (x_{i,t} \in [0, 1] : i \in \mathcal{M}, t \in \mathcal{T})$ which gives the fraction of time slot $t$ that the BS's bandwidth is assigned to user $i$. The average available rate for user $i$ at time slot $t$ is denoted as $\bar{r}_{i,t}$, which is calculated by mapping the predicted user traces to the Radio Environment Map (REM) at the service provider. The main objective of the proposed predictive resource allocation scheme is to minimize the energy consumed by the BS in transmitting the video content to the users while satisfying their QoS level. To achieve this, we incorporate the following energy and QoS models.

*1) Energy Minimization:* Studies on BS energy consumption and sleeping strategies [4], [31], reveal that the energy consumption $E$ is approximately linearly proportional to the airtime fraction of the BS [15], [32]. This is commonly referred

to as time duty-cycling. In essence, $E = P \times \Delta T$ where $P$ is the total transmitted power by the BS and $\Delta T$ is the time during which the BS was switched ON. The dominant part of the power is that transmitted over the wireless channel, which is largely constant as downlink power control is not employed in the current LTE 3GPP standards [31], [33]. Accordingly, the energy consumption can be expressed in terms of the airtime $\Delta T$ to avoid dependencies on the constant power fraction that varies with BS type [32]. Therefore, as in [13], [15], we minimize the energy consumption by minimizing the total time air fractions $x_{i,t}$ allocated to all the users.

*2) QoS Satisfaction:* To achieve energy savings under QoS satisfaction, the BS should use the minimum resources needed to guarantee the video delivery at the target user quality over a time horizon. Existing energy-efficient RA approaches reveal that playback interruptions, due to buffer underrun, are among the primary sources of user dissatisfaction with video delivery services [14], [34]. In essence, video freezing occurs when the allocated airtime up to time slot $t$ results in delivering a total amount of video less than the corresponding cumulative streaming demand. This demand can be denoted as $D_{i,t} = V_i \times t$, where $V_i$ is the fixed streaming rate of user $i$ corresponding to the requested video quality. The number of video stops can therefore provide a sound QoS metric when modeling RA to optimize the trade-off between energy-minimization and QoS satisfaction.

A promising predictive green video delivery strategy was introduced in [15] where the airtime is minimized to decrease energy consumption, while ensuring that the total amount of data is greater than the minimum cumulative streaming rate to avoid video freezing. This is formulated as follows

$$\underset{\mathbf{x}}{\text{minimize}} \quad \sum_{t=1}^{T} \sum_{i=1}^{M} x_{i,t} \qquad (3)$$

subject to:

$$\text{C1:} \quad \sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'} \geq D_{i,t}, \quad \forall i \in \mathcal{M}, t \in \mathcal{T},$$

$$\text{C2:} \quad \sum_{i=1}^{M} x_{i,t} \leq 1, \qquad \forall t \in \mathcal{T},$$

$$\text{C3:} \quad x_{i,t} \geq 0 \qquad \forall i \in \mathcal{M}, t \in \mathcal{T}.$$

The QoS constraint C1 ensures that the cumulative video content requirement is not violated at each time slot, and C2 expresses the resource limitation at each base station. It ensures that the sum of the airtime of all users is less than 1 s, which is the time slot duration. Finally, C3 is the non-negativity of airtime fraction.

The objective function in (3) refers to the minimization of the total allocated airtime which allows the base station to use time duty cycling and go into sleeping mode [33] and thus saves energy. As such, the base station transmission is set to only consume power either to satisfy the minimum user's demand or to strategically prebuffer the video during the user's peak radio conditions; to avoid future allocations at the cell edge. This allocation is unlike the Maximum Throughput (MT) which aims to

exploit the total resources and thus does not allow the BS to switch to sleeping mode.

In summary, the allocation in (3) achieves both energy minimization and QoS satisfaction under perfect future channel knowledge. More details of such perfect knowledge based PRA can be found in [14], [15]. However as discussed previously, the above formulation depends on the *average* value of future data rates and thus it is not robust to any channel variations. Consequently, QoS satisfaction cannot be guaranteed under practical considerations.

In this paper we formulate a *robust* predictive allocation strategy that calculates the airtime fraction $x_{i,t}$ for every user, at each time slot, using JCCP in which future rates are modeled as random numbers.

### C. Framework Overview

Our framework for robust PRA is based on Joint Chance Constrained Programming (JCCP) to provide long term QoS satisfaction at desired level $\beta$. By offering a mechanism to control the value of $\beta$, operators may achieve a balance between prediction gains and the risks associated with erroneous predictions. The main components of the framework are summarized below.

1) The JCCP model presented in detail in Section IV performs the robust predictive airtime allocation in two consecutive stages:
   - *Risk Allocation.* This stage determines the probability of constraint satisfaction at each time slot such that the total QoS level over the time horizon is achieved. This step is performed once at the beginning of the prediction window. The main challenge is to distribute such probabilities in a way that optimizes the allocation of the next stage in terms of energy consumption.
   - *Robust PRA.* Here the actual airtime fraction for each user is allocated such that the total energy consumption is minimized while satisfying the QoS levels. The allocation is determined using the calculated values from the previous stage, the user demand in addition to the average and variance of the future random rates.

2) In order to provide an additional element of robustness, the time varying variance of the predicted user rates is also estimated. The initial values of future variances are either increased or decreased based on the previous channel measurements by the user and their correlation coefficient with the current measurements. The updated variance is then provided to the robust PRA where QoS guarantees can be achieved during high variance, and energy minimization during low variance.

The formulation and implementation details of framework components are presented in following three sections.

## IV. ROBUST PROBLEM FORMULATION USING CHANCE CONSTRAINTS

In this section, we first model the robust PRA framework for video streaming using traditional *individual* chance constraints

which is found to be a convex optimization problem. Thereafter, the problem is extended to the non-convex *joint* chance constraint model to enable QoS satisfaction of the cumulative demand over the time horizon. To provide a tractable solution, the problem is then decomposed into two convex stages that can be optimally solved individually.

In what follows, we adopt the Gaussian distribution error model for the predicted rate introduced in [27], as in recent robust RA works [28], [29]. In particular, predicting the future rates using autoregressive filters, resulted in a Gaussian distributed error model compared to the actual set of collected data [27]. This is supported by the same distribution attained while applying the 3GPP correlated shadowing on the average value of predicted rates [35]. In our model the rate is predicted at a 1 s granularity, which is generally deduced from a large number of samples due to the small feedback interval (1 ms) of the users participating in channel prediction [31]. Such a scenario supports the Central limit theorem (CLT) which approximates the PDF of users' predicted rate as a Gaussian distribution [28]. Nevertheless, all the introduced formulations are applicable for other error models with closed form and invertible CDF.[1]

### A. Individual Chance Constraint Programming

The robust equivalent of the PRA in (3) is attained by replacing the QoS constraint C1 with the individual chance constraint, where predicted rates are replaced by random variables, and a probabilistic constraint is developed as follows

$$\underset{\mathbf{x}}{\text{minimize}} \quad \sum_{t=1}^{T} \sum_{i=1}^{M} x_{i,t} \tag{4}$$

subject to:

C1: $\quad Pr\left\{\sum_{t'=0}^{t} \tilde{r}_{i,t'} x_{i,t'} \geq D_{i,t}\right\} \geq \beta, \quad \forall i \in \mathcal{M}, t \in \mathcal{T},$

C2: $\quad \sum_{i=1}^{M} x_{i,t} \leq 1, \qquad\qquad \forall t \in \mathcal{T},$

C3: $\quad x_{i,t} \geq 0 \qquad\qquad\qquad \forall i \in \mathcal{M}, t \in \mathcal{T}.$

Herein, the predicted data rate $\tilde{r}_{i,t'}$ is modeled as a random variable following a normal distribution: $\tilde{r}_{i,t'} \sim N(\bar{r}_{i,t}, \sigma_{i,t}^2)$, and $\beta \in [0, 1]$ is the QoS satisfaction level.

Accordingly, the summation of the normally distributed random data rates in C1 of (4) is a multivariate normal distribution whose mean is the summation of means of all single random variables, which we denote as $\mu$. The corresponding variance is the covariance matrix denoted by $\Sigma$, and can be

evaluated as follows

$$\mu = \sum_{t'=0}^{t} \bar{r}_{i,t}, \qquad \Sigma = \begin{bmatrix} \sigma_{i,0}^2 & \cdots & \sigma_{i,0,t} \\ \cdots & \sigma_{i,1}^2 & \cdots \\ \sigma_{i,t,0} & \cdots & \sigma_{i,t}^2 \end{bmatrix}, \quad (5)$$

where $\sigma_{i,t,h} = E[(\tilde{r}_{i,t} - \bar{r}_{i,t})(\tilde{r}_{i,h} - \bar{r}_{i,h})]$ and $\sigma_{i,t}^2 = \sigma_{i,t,h}, \forall t = h$.

The deterministic closed form of (4) can be expressed using the multivariate random variables and normal cumulative distribution function as shown below.

$$Q\left(\frac{D_{i,t} - \sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'}}{\sqrt{\sum_{t'=0}^{t} \sum_{h=0}^{t} x_{i,t'}^2 \sigma_{i,t',h}}}\right) \geq \beta, \forall i \in \mathcal{M}, t \in \mathcal{T},$$

$$\sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'} + Q_\beta^{-1} \sqrt{\sum_{t'=0}^{t} \sum_{h=0}^{t} x_{i,t'}^2 \sigma_{i,t',h}} \geq D_{i,t}. \tag{6}$$

The independence between the realizations of random predicted channel rate at each time slot implies that $\sigma_{i,t',h} = 0, \forall t' \neq h$. Accordingly, the chance constraint is represented as follows

$$\sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'} + Q_\beta^{-1} \sqrt{\sum_{t'=0}^{t} x_{i,t'}^2 \sigma_{i,t'}^2} \geq D_{i,t}, \tag{7}$$

$$\forall i \in \mathcal{M}, t \in \mathcal{T}.$$

The above constraint representation is a second order cone programming (SOCP) model which is convex for $\beta > 0.5$ [37] and results in a negative value for the inverse of the Q-function. Finally, the deterministic closed form of (4) using individual chance constraint with the preceding assumptions can be summarized below

$$\underset{\mathbf{x}}{\text{minimize}} \quad \sum_{t=1}^{T} \sum_{i=1}^{M} x_{i,t} \tag{8}$$

subject to:

C1: $\quad \sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'} + Q_\beta^{-1} \sqrt{\sum_{t'=0}^{t} x_{i,t'}^2 \sigma_{i,t'}^2} \geq D_{i,t},$

$$\forall i \in \mathcal{M}, t \in \mathcal{T},$$

C2: $\quad \sum_{i=1}^{M} x_{i,t} \leq 1, \qquad\qquad \forall t \in \mathcal{T},$

C3: $\quad x_{i,t} \geq 0 \qquad\qquad\qquad \forall i \in \mathcal{M}, t \in \mathcal{T}.$

As mentioned in Section II, this type of chance constraint formulation ensures that the QoS is satisfied each time slot at a certain level $\beta$. However, it does not model the cumulative satisfaction for each user over the time horizon in which the per slot demand satisfaction is dependent on the total data delivered in the preceding time slots. In order to avoid future buffer starvation, the allocation in each time slot should compensate

---

[1]It has to be noted that the total probability of negative realizations for the normally distributed random rate has a non-significant value ($\approx 0$). This is attributed to the high average rate values that maintain a positive distribution under typical variances in the 3GPP models and standards [29], [35], [31], [36].

the unsatisfied previous demands. This is why the joint chance constraint model is needed.

### B. Joint Chance Constraint Programming

The joint chance constraint form for the C1 constraint in (3) can be expressed as follows

$$Pr\left\{\bigcap_{\forall t \in \mathcal{T}} \sum_{t'=0}^{t} \tilde{r}_{i,t'} x_{i,t'} \geq D_{i,t}\right\} \geq \beta, \forall i \in \mathcal{M}. \quad (9)$$

We denote the event of individual QoS satisfaction by $S_{i,t} \triangleq \left\{\sum_{t'=0}^{t} \tilde{r}_{i,t'} x_{i,t'} \geq D_{i,t}\right\}$. Similarly, the event of individual QoS dissatisfaction is denoted by $S_{i,t}^c$. The probability of joint satisfaction of event $S_{i,t}$ is represented as the complement of disjoint probability of the dissatisfaction event as in (10)

$$Pr\left\{\bigcap_{\forall t \in \mathcal{T}} S_{i,t}\right\} = 1 - Pr\left\{\bigcup_{\forall t \in \mathcal{T}} S_{i,t}^c\right\}, \forall i \in \mathcal{M}. \quad (10)$$

According to Boole's inequality, the disjoint probability is tightly bounded from above by the total probability of all individual events [22] as follows

$$Pr\left\{\bigcup_{\forall t \in \mathcal{T}} S_{i,t}^c\right\} \leq \sum_{\forall t \in \mathcal{T}} Pr\left\{S_{i,t}^c\right\}, \forall i \in \mathcal{M}. \quad (11)$$

The joint probability of the QoS satisfaction event is therefore bounded as below

$$Pr\left\{\bigcap_{\forall t \in \mathcal{T}} S_{i,t}\right\} \geq 1 - \sum_{\forall t \in \mathcal{T}} Pr\left\{S_{i,t}^c\right\}, \forall i \in \mathcal{M},$$

$$Pr\left\{\bigcap_{\forall t \in \mathcal{T}} S_{i,t}\right\} \geq \beta, \forall i \in \mathcal{M}, \quad (12)$$

$$\sum_{\forall t \in \mathcal{T}} Pr\left\{S_{i,t}^c\right\} \leq 1 - \beta, \forall i \in \mathcal{M}.$$

The above equation implies that the joint probability is satisfied if the summation of individual probabilities of the compliment event is kept below the probability of QoS dissatisfaction (i.e., $1 - \beta$). Accordingly, the joint chance constraint in (9) can be replaced by the two constraints in (13) and (14)

$$Pr\left\{\sum_{t'=0}^{t} \tilde{r}_{i,t'} x_{i,t'} < D_{i,t}\right\} < \zeta_{i,t}, \forall i \in \mathcal{M}, t \in \mathcal{T}. \quad (13)$$

$$\sum_{\forall t \in \mathcal{T}} \zeta_{i,t} \leq 1 - \beta, \forall i \in \mathcal{M}. \quad (14)$$

where $\zeta_{i,t}$ is denoted as the probability for not satisfying the individual QoS constraint (i.e., $Pr\left\{S_{i,t}^c\right\}$) and is called the probability of *risk* [23].

Each probabilistic constraint in (13) will have the same deterministic equivalent form as the individual chance constraint but with $\beta$ replaced by $\zeta_{i,t}$. After incorporating (13) and (14), this JCCP formulation becomes a function of both variables: $\zeta_{i,t}$ and $x_{i,t}$ as summarized below

$$\underset{\mathbf{x}, \zeta}{\text{minimize}} \quad \sum_{t=1}^{T} \sum_{i=1}^{M} x_{i,t} \quad (15)$$

subject to:

$$C1: \quad \sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'} + Q_{1-\zeta_{i,t}}^{-1} \sqrt{\sum_{t'=0}^{t} x_{i,t'}^2 \sigma_{i,t'}^2} \geq D_{i,t},$$

$$\forall i \in \mathcal{M}, t \in \mathcal{T},$$

$$C2: \quad \sum_{i=1}^{M} x_{i,t} \leq 1, \qquad\qquad \forall t \in \mathcal{T},$$

$$C3: \quad x_{i,t} \geq 0 \qquad\qquad \forall i \in \mathcal{M}, t \in \mathcal{T},$$

$$C4: \quad \sum_{\forall t \in \mathcal{T}} \zeta_{i,t} \leq 1 - \beta, \qquad \forall i \in \mathcal{M}.$$

Indeed the above formulation is no longer convex and thus the optimal solution can not be guaranteed by traditional optimization techniques. A proof of its non-convexity is provided in Appendix A. Therefore, to provide a tractable solution, the above formulation is split into two stages: *Risk Allocation* and *Robust PRA*. The first stage determines the optimal values for each risk level (i.e., solves for $\zeta_{i,t}$), while the second stage solves the PRA problem given the calculated QoS satisfaction levels in the prior stage (i.e., solves for $x_{i,t}$).

**Stage A: Risk Allocation** In this stage, the value of risk probabilities for each constraint is determined such that Boole's inequality (14) is satisfied to guarantee the joint probability satisfaction of (9). An initial feasible solution is to uniformly distribute the probability of risk $(1 - \beta)$ over all the time horizon. In other words, assign an equal risk probability $\zeta_{i,t}$ among all the time slots of each user as below

$$\zeta_{i,t} = \frac{1 - \beta}{T}, \quad \forall i \in \mathcal{M}. \quad (16)$$

However, such equal risk allocation was proven to be very conservative [23] and results in suboptimal resource allocation that compromises the energy savings of the PRA obtained in the second stage. Hence, optimal risk allocation is applied to consider the optimality of the second stage in addition to the Boole's inequality constraint C4 in (15).

Note that lower risk probability $\zeta_{i,t}$ results in higher airtime $x_{i,t}$, and that $x_{i,t}$ is inversely proportional to its corresponding average rate $\bar{r}_{i,t}$ as depicted in (3). Therefore, the risk of each time slot is allocated proportionally to the corresponding average rate $\bar{r}_{i,t}$ in order to minimize the energy consumption during the resource allocation stage. In other words, time slots with low average data rate will suffer from high airtime for QoS satisfaction. Thus, assigning low risk probability to these slots will result in additional airtime. To that end, the following risk allocation optimization is introduced in (17) to achieve the

optimality of the second stage as well

$$\text{minimize} \quad \sum_{t=1}^{T} \left(\frac{\hat{r}_i}{\bar{r}_{i,t}}\right)^n y_{i,t} \quad \forall i \in \mathcal{M}, \qquad (17)$$

$$\text{subject to:} \quad \sum_{\forall t \in \mathcal{T}} Q(y_{i,t}) \le 1 - \beta, \forall i \in \mathcal{M}.$$

where: $y_{i,t} = Q_{\zeta_{i,t}}^{-1}$ to represent the constraint in a differentiable form, $\hat{r}_i = \max_t \bar{r}_{i,t}$ and $n$ is the risk proportionality parameter whose value is positive. The value of $n$ captures the trade-off between the risk of not satisfying the QoS at a certain time slot and the energy savings. For very small values of $n$, the risk is fairly distributed among the time slots and the user will not suffer from successive video degradations. On the other hand, more energy savings are obtained when the value of $n$ increases since high risk is allowed at low data rate values. The mobile operator then may tune $n$ based on the maximum allowable consecutive degradation, or the desired energy savings. The above problem is convex given that $\beta \ge 0.5$, which is valid for practical considerations. A proof of this convexity is provided in Appendix B.

**Stage B: Robust PRA**

After solving the first stage in (17), and determining the risk probabilities $\zeta_{i,t}$ for each constraint, the problem in (15) can be solved without constraint C4. The resulting formulation preserves the form of SOCP, which is still convex due to the positiveness of the calculated risk probabilities.

## V. GRADIENT BASED AND GUIDED HEURISTIC SOLUTION METHODS

After decomposing the joint chance constraint programming into two convex optimization stages, the solution methods for each stage are introduced in this section.

### A. Risk Allocation Solution

The constrained proportional risk allocation in (17) is solved by calculating the Lagrange formulation and then using Newton's method to search for the saddle points that satisfy the Karush–Kuhn–Tucker (KKT) optimality conditions as follows

$$\mathcal{L}(y, \lambda) = \sum_{t=1}^{T} \left(\frac{\hat{r}_i}{\bar{r}_{i,t}}\right)^n y_{i,t} - \lambda \left(\sum_{\forall t \in \mathcal{T}} Q(y_{i,t}) - (1 - \beta)\right),$$

$$(18)$$

$$\forall i \in \mathcal{M}.$$

where $\lambda \ge 0$ is the Lagrange multiplier associated with the constraint in (17).

Since the above problem is optimized for each user separately and performed only once at the beginning of the time horizon, optimal path searching methods provide acceptable performance. We therefore apply Newton's method as summarized in Algorithm 1. The algorithm starts with the uniform risk allocation and then iteratively searches for the saddle

---

**Algorithm 1.** Newton's Method for Proportional Risk Allocation

---

**Input:** Time Horizon: $\mathcal{T}_i$, Average Predicted Rates: $\bar{r}_i$, QoS Level: $\beta$ and Risk Proportionality Factor: $n$

**Output:** $y_i$

Initialization : $\zeta_{i,t} = \frac{1-\beta}{\mathcal{T}_i}$, $y_{i,t} = Q_{\zeta_{i,t}}^{-1}$, $\forall t \in \mathcal{T}$, $\lambda = \lambda_0$, $\epsilon = 0.001$, $\Delta y_i = \Delta y_0$ and $\mathcal{L} = [y_i \ \lambda]^T$

1: **while** $\Delta y_i \ge \epsilon$ **do**

2: $\quad \frac{\partial \mathcal{L}(y_{i,t}, \lambda)}{\partial y_{i,t}} = \left(\frac{\hat{r}_i}{\bar{r}_{i,t}}\right)^n + \lambda \frac{1}{\sqrt{2\pi}} e^{\frac{-y_{i,t}^2}{2}}$

$\quad \frac{\partial \mathcal{L}(y_{i,t}, \lambda)}{\partial \lambda} = -\left(\sum_{\forall} t' \in \mathcal{T} Q(y_{i,t}) - (1 - \beta)\right)$

$\quad \frac{\partial^2 \mathcal{L}(y_{i,t}, \lambda)}{\partial y_{i,t}^2} = -\lambda \frac{1}{\sqrt{2\pi}} y_{i,t} e^{\frac{-y_{i,t}^2}{2}}$

$\quad \frac{\partial^2 \mathcal{L}(y_{i,t}, \lambda)}{\partial y_{i,t} \partial \lambda} = \frac{1}{\sqrt{2\pi}} e^{\frac{-y_{i,t}^2}{2}}$

3: $\quad$ *Construct:* $\nabla \mathcal{L}(y_i, \lambda)$ and $\nabla^2 \mathcal{L}(y_i, \lambda)$

4: $\quad$ *Calculate* $(\nabla^2 \mathcal{L}(y_i, \lambda))^{-1}$

5: $\quad \Delta \mathcal{L} = -(\nabla^2 \mathcal{L}(y_i, \lambda))^{-1} \nabla \mathcal{L}(y_i, \lambda)$

6: $\quad \mathcal{L} = \mathcal{L} + \Delta \mathcal{L}$

7: $\quad \Delta y_i = \Delta \mathcal{L}(1 : T)$

8: $\quad \Delta \lambda = \Delta \mathcal{L}(T + 1)$

9: $\quad y_i = y_i + \Delta y_i$

10: $\quad \lambda = \lambda + \Delta \lambda$

11: **end while**

12: **return** $y_i$

---

points along the gradient while the step size is guided by the Hessian matrix. The calculated step value $\Delta \mathcal{L}$ contains the change in both the decision vector $y_i$ and the Lagrange multiplier $\lambda$ which are denoted as $\Delta y_i$ and $\Delta \lambda$, respectively. In each iteration, both decision vectors are updated using the calculated step, and the algorithm stops when the iterations no longer result in a significant enhancement, denoted by $\epsilon$.

### B. Robust Predictive Resource Allocation

The calculated risk probabilities for each user at every time slot are now readily available to the robust PRA stage from the risk allocation solution. The objective of this stage is to solve for the airtime allocation formulated in (15). The solution of this stage is much more complex compared to the risk allocation since here the airtime is determined jointly for all the users over the total time horizon. Based on the users' feedback, this stage is recomputed every $\tau$ seconds according to the received amount of data. To address the resulting impractical complexity, a guided heuristic is also introduced to provide a real-time resource allocation solution, while the derivative based and line search methods are used to provide benchmark solutions.

*1) Interior Point With Barrier Based Solution:* The interior point method (IPM) with barrier function was proven to satisfy the KKT optimality conditions, and thus achieves the optimal solution for SOCP formulation [38]. The unconstrained log-barrier formulation for the optimization problem in (15) is

expressed below

$$\underset{\mathbf{x}}{\text{minimize}} \quad B_\gamma(x) = F(x) + \gamma \Phi(x) \qquad (19)$$

where:

$$F(x) = \sum_{t=1}^{T} \sum_{i=1}^{M} x_{i,t},$$

$$\Phi(x) = -\sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{M}} log \left( \sum_{t'=0}^{t} \bar{r}_{i,t'} x_{i,t'} \right.$$

$$\left. + y_{i,t} \sqrt{\sum_{t'=0}^{t} x_{i,t'}^2 \sigma_{i,t'}^2} - D_{i,t} \right)$$

$$- \sum_{t \in \mathcal{T}} log \left( 1 - \sum_{i=1}^{M} x_{i,t} \right),$$

where $\gamma$ is the barrier parameter that controls the convergence of the solution method.

Path following algorithms can be applied to obtain the optimal solution such as Newton's method discussed previously in Algorithm 1. The dimension of the Hessian matrix depends on both the number of users and the allocation time slots of each user. Its inverse matrix will be of high computational complexity, thereby limiting its use for real-time resource allocation.

*2) Guided Heuristic Solution:* To provide a low complex alternative solution, a guided search algorithm is introduced that exploits the problem's features rather than the direct gradient based iterative search. The algorithm first calculates the minimum allocation for the users to ensure constraint satisfaction (i.e., satisfy C1 in (15)) given the calculated risk probabilities and the requested demands. In case of radio resource limit violations (i.e., C2 in (15)), airtime reallocation of users is done by granting the excess user requirement in other time slots. In order to achieve energy minimization, users are allocated the residual airtime when they reach the peak average rate location. Residual airtime is the remaining airtime after satisfying the QoS constraints (first step) for all users. The heuristic is summarized in Algorithm 2.

**Minimal airtime allocation:** To ensure the satisfaction of QoS constraint, C1 in (15) is turned to equality in the quadratic form $ax^2 + bx + c = 0$ and solved using (20) (Lines 4-11) of Algorithm 2. This is achieved as follows

$$x_{i,t'} = \frac{-b_{i,t'} + \sqrt{b_{i,t'}^2 - 4a_{i,t'} c_{i,t'}}}{2a_{i,t'}}, \qquad (20)$$

$$\text{Where:} \quad a_{i,t'} = \bar{r}_{i,t'}^2 - (y_{i,t'} \sigma_{i,t'})^2,$$

$$b_{i,t'} = -2K_{i,t'} \bar{r}_{i,t'}^2,$$

$$c_{i,t'} = K_{i,t'}^2 - (y_{i,t'} L_{i,t'})^2,$$

$$K_{i,t'} = D_{i,t'} - \sum_{h=0}^{t'-1} x_{i,t'} \bar{r}_{i,t'},$$

$$L_{i,t'} = \sum_{h=0}^{t'-1} x_{i,t'}^2 \bar{r}_{i,t'}^2.$$

---

**Algorithm 2.** Guided Heuristic Robust Green Allocation

**Input:** Users: $\mathcal{M}$, Time Horizon: $\mathcal{T}_i$, Mean of Predicted Rates: $\bar{R}$, Rate Variances: $\Sigma$, Risk Levels: $Y$ and Demand: $D$

**Output:** $X$

1: *Initialization* : $x = \emptyset$,
    $t_i^{(p)} = \underset{t \in \mathcal{T}}{\text{argmax}} \{\bar{R}_i\}, \forall i \in \mathcal{M}$   /* time slot with maximum average rate (cell center) */
2: **for all** $t \in \mathcal{T}$
3:    $\tau_t = 0$  /* total airtime fraction allocated in time slot $t$ */
4:    **for all** $i \in \mathcal{M}$ **do**
5:       **if** $t < t_i^{(p)}$ **then**
6:          Calculate $x_{i,t}$ using (20)   /* minimal airtime allocation*/
7:          $\tau_t = \tau_t + x_{i,t}$
8:       **else**
9:          $M := M \setminus i$   /* remove user from minimal allocation after reaching cell center*/
10:      **end if**
11:   **end for**
12:   **if** $\tau_t > 1$ **then**
13:      $i^{(*)} := \underset{i \in \mathcal{M}}{\text{argmax}} \{x_{i,t}\}$,   /*choose the user with maximum airtime violating the constraint*/
14:      $\delta x_{i^*,t} = \tau_t + x_{i^*,t} - 1$   /*violating airtime excess fraction*/
15:      **for** $n := t - 1$ **to** $0$ **do**
16:         **if** $\tau_n + \delta x_{i^*,t} < 1$ **then**
17:            $x_{i^*,n} := x_{i^*,n} + \delta x_{i^*,t}$   /*Repair the solution*/
18:            $\tau_n := \tau_n + \delta x_{i^*,t}$
19:            $\tau_t = 1$
20:         **end if**
21:      **end for**
22:   **end if**
23: **end for**
24: **for all** $i \in \mathcal{M}$ **do**
25:   *AllocatePeaks* $(\tau_t, t_i^p)$
26: **end for**
27: **return** $X$

---

**Allocation Repair:** The total allocated airtime to all users in each time slot is calculated and the radio resource limitation constraint, C2 in (15), is checked. In case of any violations, the excess airtime is allocated in other time slots with unused resources. Particularly, the heuristic compensates (recovers) any time slot $t \in T$ with a total allocated airtime fractions (i.e., $\tau_t = \sum_{\forall i \in I} x_{i,t} \forall t$) more than the slot duration (1 sec.) which occurs due to 1) an increased number of users, 2) high traffic per user or, 3) high QoS level ($\beta$). The heuristic solves this case by iteratively picking the user with the maximum airtime fraction in this time slot and prebuffering his video content in advance to ensure airtime minimization under demand satisfaction (Lines 12-21) in Algorithm 2.

**Peak Average Rate Allocation:** The above allocation strategy guarantees the satisfaction of both QoS and resource constraints. Thus, it continues until the peak data rate time slot is reached. The allocation strategy is then changed (Line 24)

to allocate the demand of the future time slots in advance, to minimize the airtime. This follows the following steps for each user $i$

- Calculate the residual demand for user $i$: $\delta D_{i,t'} = D_{i,T} - \sum_{t=0}^{t=t'} D_{i,t}$
- Repeat the allocation strategy in step 1 until either the total residual demand is allocated or the peak rate time slot is full.
- In case of remaining demand while the peak rate time slot is fully loaded, the second peak average rate with remaining airtime is selected and the above procedure continues.
- In each iteration, the residual demand is decremented by $x_{i,t'} \times (\bar{r}_{i,t'} - y_{i,t'}\sigma_{i,t'})$, which is a conservative method since it assumes the worst case channel capacity of the current rate.
- The algorithm terminates when all users received their total demand denoted as $D_{i,T}$.

Both the feasibility and optimality of the obtained resource allocation solution are highly sensitive to the variance $\sigma^2$. Applying the second stage with low variance does not guarantee the constraint satisfaction since less airtime will be allocated to the user according to (20), especially during low data rates when high risk probability is allowed.

On the other hand, using a large variance $\sigma^2$ results in a conservative solution that allocates too much airtime especially in relatively high data rate time slots when low risk is applied. Due to the fluctuation of $\sigma^2$ with the user location and time of the day as previously mentioned in Section I, a fixed value of $\sigma^2$ becomes suboptimal. We therefore propose to adaptively track the variance $\sigma^2$ based on the user's previous measurements. The tracking procedure is implemented using Kalman Filter (KF) described in detail in the following section.

## VI. KALMAN FILTER BASED VARIANCE ESTIMATION

The variance of the random predicted rates are updated using the channel measurements by the user in the previous time slot. The measured rate variance by user $i$ during the previous time slot $t-1$ is denoted as $\bar{\sigma}_{i,t-1}^2$ and calculated as follows

$$\bar{\sigma}_{i,t-1}^2 = (\bar{r}_{i,t-1} - \mathbf{\bar{r}}_{i,t-1})^2, \qquad (21)$$

where $\mathbf{\bar{r}}_{i,t-1}$ is the average measured data rate by user $i$ during the previous time slot $t-1$. We then denote $\bar{\delta}\sigma_{i,t}^2$ as the ratio between the measured and the initial theoretical variances denoted as $\bar{\sigma}_{i,t-1}^2$ and $\sigma_{i,t-1}^2$, respectively. Although the variance ratio represents the actual deviations from the initial variance, the former still varies from one time slot to another. Accordingly, the change in the variance over time is modeled as a Gaussian process and thus can be accurately estimated using Kalman Filter, which is known to be the optimal linear estimator in the mean square error sense. It is composed of two stages as summarized below [39]:

**Prediction Phase:**

$$\mathcal{X}_t^- = \Phi_t \mathcal{X}_{t-1}^+. \qquad (22)$$
$$\mathcal{P}_t^- = \Phi_t \mathcal{P}_{t-1}^+ \Phi_t' + \mathcal{Q}. \qquad (23)$$

**Measurement Phase:**

$$\mathcal{K}_t = \mathcal{P}_t^- H_t'(H_t \mathcal{P}_t^- H_t' + \mathcal{R})^{-1}. \qquad (24)$$
$$\mathcal{X}_t^+ = \mathcal{X}_t^- + \mathcal{K}_t(z_t - H_t \mathcal{X}_t^-). \qquad (25)$$
$$\mathcal{P}_t^+ = \mathcal{P}_t^- - K_t H_t \mathcal{P}_t^-. \qquad (26)$$

where $\mathcal{X}_t^-$ and $\mathcal{X}_t^+$ are the priori and posterior state vectors respectively. $\mathcal{P}_t^-$ and $\mathcal{P}_t^+$ are the priori and posterior state estimation covariance matrices respectively. $H$ and $\Phi$ are the observation (design) and state transition matrices respectively, while $\mathcal{K}$ is the KF gain. $\mathcal{Q}$ and $\mathcal{R}$ are the process and the measurement noise covariance matrices respectively.

The Kalman filter performs state vector estimation using two phases: Prediction and Measurement. In the first phase, the predicted state value $\mathcal{X}_t^-$ is calculated using the previously estimated value $\mathcal{X}_{t-1}^+$ in time slot $t-1$ as indicated in (22). In the measurement phase, the new state is calculated using a weighted difference between the observed measurements $z_t$ and the predicted state (25). This weighting is done using Kalman gain $\mathcal{K}_t$ calculated in (24), that is dependent on both the measurement noise covariance $\mathcal{R}$ and the predicted state estimation covariance $\mathcal{P}_t^-$ in (24).

In our problem, the priori state $\mathcal{X}_t^-$ represents the variance ratio $\delta\sigma_{i,t}^2$ and equals the corrected state of the previous time epoch $\mathcal{X}_{t-1}^+$ by setting the state transition to unity. The observation $z_t$ represents the measured variance ratio $\bar{\delta}\sigma_{i,t}^2$ shown in (21). The observed measurements $z_t$ and the predicted state $\mathcal{X}_t^-$ represent different values for the same quantity (i.e., variance ratio), and therefore the state observation matrix $H$ is set to unity. The updated ratio $\delta\sigma_{i,t}^{2+}$ will be then used to update the predicted variances in the remaining time slots, denoted as $\sigma_{i,t+\delta t}^{2+}$, while simultaneously considering their correlation with the current measurement as follows

$$\sigma_{i,t+\delta t}^{2+} = \left(1 + \rho_{t,t+\delta t}(\delta\sigma_{i,t+\delta t}^{2+} - 1)\right)\sigma_{i,t+\delta t}^2, \forall\, \delta t \in [1, T-t], \qquad (27)$$

where $\rho_{t,t+\delta t}$ is the channel correlation coefficient between the channel fading at time $t$ and $t + \delta t$.

According to (27), in case of high correlation (i.e., $\rho_{t,t+\delta t} \approx 1$), the future variance will be multiplied by the value of current updated ratio and the term in the brackets becomes 1. On the other hand, very low correlation results in no updates of the future variance. In our model, we calculate the correlation coefficient using an exponentially decaying function with the correlation distance $d_{cor}$ according to the 3GPP slow fading model [35].

## VII. PERFORMANCE EVALUATION

### A. Simulation Set-up

The presented robust PRA techniques are simulated for an LTE network using the Network Simulator (ns-3) and Gurobi optimizer based environment in [40], with model parameters and KF initial values (i.e., $P_0$, $Q$, $R$ and $\delta\sigma_0$) as indicated in Table I. The 3GPP correlated slow fading model and its parameters [35] are incorporated in the received UE power and

TABLE I
SUMMARY OF MODEL PARAMETERS

| Parameter | Value |
|---|---|
| BS transmit power | 43 dBm |
| Bandwidth | 5 MHz |
| Time Horizon $T$ | 60 s |
| Streaming rate V | 0.5, 1, 1.5 [Mbps] |
| Bit Error Rate | $5 \times 10^{-5}$ |
| Shadow correlation distance($d_{cor}$) [35] | $50m$ |
| Shadow standard deviation($\sigma$) [35] | $6 dB$ |
| Velocity | From 25 km/h to 60 km/h |
| $P_0$ | 1 |
| $Q$ | 0.1 |
| $R$ | 1 |
| $\delta\sigma_0$ | 1 |
| Risk Proportionality Factor $n$ | 4 |
| Feedback interval $\tau$ | $5s.$ |
| Packet size | $10^3$ [bytes] |
| Packet rate (from core network to BS) | $10^3 s^{-1}$ |
| Total number of packets | $7.5 \times 10^3$ |
| Buffer size | $10^9$ [bits] |

thus provide predicted rate variations. Simulation results are averaged over 50 runs for statistical validation. Users follow different predefined paths within the cell at varying velocities from 25 to 60 Km/h and request a video stream at a fixed quality. Although the allocation is done at each base station separately, neighbouring BSs are placed at an inter-cell distance of 600 m for practical calculation of SINR and channel rates. The actual transmission rate by the network (i.e., transport block size) during the allocated airtime varies according to the selected Modulation and Coding Scheme (MCS) which is based on the reported CQI, after SINR mapping, by the users [31], [41].

## B. Evaluation Metrics and Scheme Notations

In order to assess the introduced Robust Predictive Resource Allocation (R-PRA) framework, we use the two metrics previously discussed in Section III. The first is the percentage of videos stops which reflects the user QoS level. Mathematically, it is calculated as the percentage of time slots in which constraint C1 in (3) is violated. A maximum allowable degradation level is defined as the boundary for the metric, and is equal to $(1 - \beta) \times 100\%$. The second metric is the average BS airtime which is used to measure the energy consumption in the network. During resource allocation, both the BS and UE consume energy in transmission and reception of data. Therefore, minimizing airtime reduces the energy consumption proportionally [32]. The objective function in (3) is used to quantitatively measure this metric.

In this evaluation study, we denote the proposed optimal ICCP and JCCP, and their corresponding heuristics with the following abbreviations:

- **Optimal-ICCP:** refers to formulation in (8) whose solution is obtained using the IPM of Section V implemented in Gurobi.
- **Heuristic-ICCP:** refers to formulation in (8) whose solution is obtained using the guided heuristic in Algorithm 2.
- **Optimal-JCCP:** based on the original non-convex JCCP formulation in (15) and solved using the sequential

quadratic programming in MATLAB for a global optimal risk and airtime allocations.

- **Optimal-ERA-JCCP:** uses the two stage JCCP in which the first stage solution is obtained with equal risk values (16) and the second stage (8) is solved using the IPM implemented in Gurobi.
- **Heuristic-ERA-JCCP:** similar to the Optimal-ERA-JCCP but the second stage is solved using the guided heuristic in Algorithm 2.
- **Optimal-PRA-JCCP:** similar to the Optimal-ERA-JCCP with first stage formulated as in (17) and solved with Lagrangian Newton in Algorithm 1.
- **Heuristic-PRA-JCCP:** similar to the Heuristic-ERA-JCCP with the first stage formulated as in (17) and solved with Lagrangian Newton in Algorithm 1.

The optimal techniques are used to 1) evaluate the robustness of the introduced framework, and 2) assess the developed real-time guided heuristic in Algorithm 2. The non-convex Optimal-JCCP is used to evaluate the feasibility of the decomposed two-stage JCCP.

## C. Simulation Results

*1) Comparison With Existing Non-Predictive and Non-Robust RA:* The first simulated scenario is for one user moving across the cell from one edge to the other. Both the predicted average and the actual experienced rates are shown in Fig. 1(a). We consider three typical classes of RA:

- **NP-RA:** refers to opportunistic Non-predictive Resource Allocation and the widely used Proportional Fairness will be adopted as a type of this class.
- **NR-PRA:** refers to the existing energy-efficient Non-Robust Predictive Resource Allocation in [15], which assumed perfect prediction and represented the future rate by its average value.
- **R-PRA:** refers to the energy-efficient Robust Predictive Resource Allocation introduced in this paper in its two main forms (ICCP and JCCP).

The NR-PRA assumes perfect prediction of the future channel rates and results in the minimum energy consumption compared to both the NP-RA and the R-PRA as illustrated in Fig. 2(a). This is because, NR-PRA strategically allocates the minimal airtime that satisfies the demand based on the average predicted rate until the user reaches the cell center. On the other hand, the introduced R-PRA conservatively allocates more airtime than the NR-PRA to guarantee QoS satisfaction under rate variations. The NP-RA, however, greedily assigns all the available resources and thus delivers the video to the user during the initial low rates regardless of the future high rates as shown in Fig. 1(b). On the other hand, Fig. 2(b) shows that the low-energy NR-PRA failed to satisfy the QoS demand as we can see that the user suffered from a large percentage of video stops. On the contrary, the proposed R-PRA (ICCP and JCCP) was able to compensate the variations by strategically allocating more airtime and the result is much fewer video stops. The traditional non-energy aware NP-RA filled the buffer of the user in the first few seconds, resulting in the highest QoS satisfaction
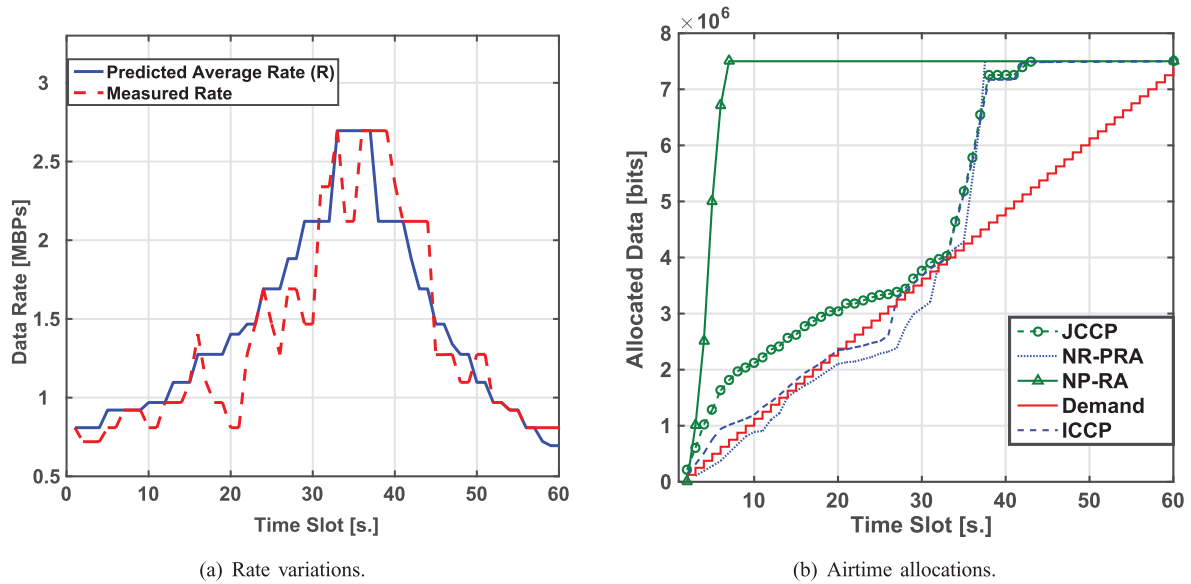
(a) Rate variations.



(b) Airtime allocations.

Fig. 1. Illustrative allocation and rate variations examples for the considered techniques.



(a) BS airtime.



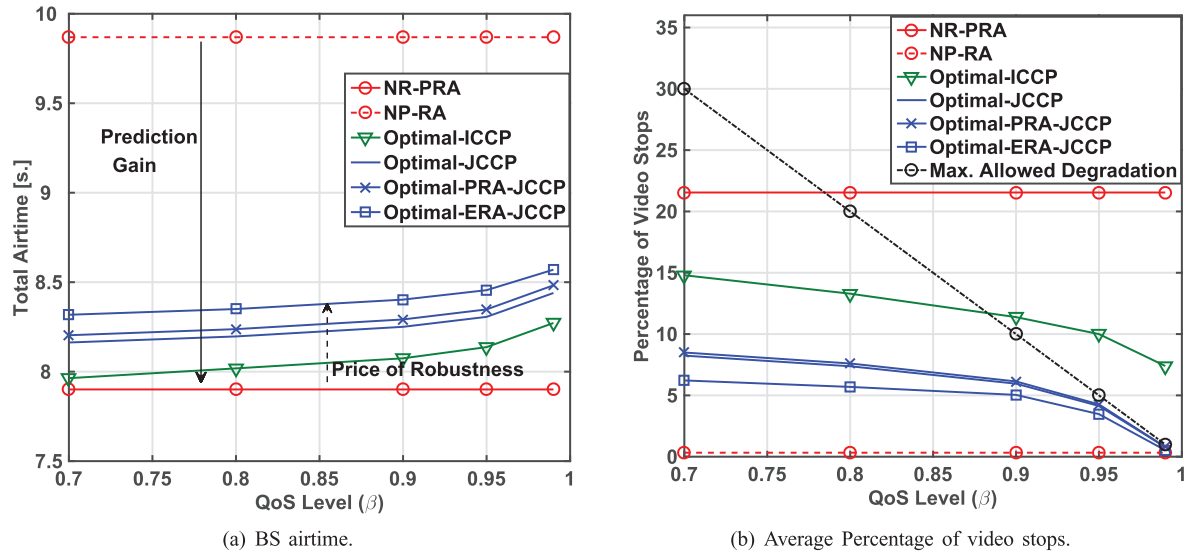(b) Average Percentage of video stops.

Fig. 2. Percentage of video stops and average BS airtime for varying QoS degrees $\beta$ for 1 user experiencing rate variations.

with a negligible number of stops, but at the cost of high energy consumption.

To summarize, the NR-PRA previously introduced in [15] provides large energy savings, denoted as the *Prediction Gain*, compared to the NP-RA. However, this gain was achieved with unacceptable QoS violations under imperfect predictions. To overcome this limitation, the introduced R-PRA is designed to simultaneously satisfy the QoS requirements and energy minimization. This comes at the cost of slightly decreasing the prediction gain by an amount referred to as the *Price of Robustness* that accounts for rate variations. The above conclusions can also be drawn from the higher load scenario in Fig. 4, and indicate that robust PRA can provide significant gains under practical considerations of imperfect predictions. These results are obtained for the optimal forms of the introduced R-PRA (i.e., Optimal-ICCP and Optimal-JCCP) to assess

their performance bounds, and the developed real-time heuristic which will be assessed separately. We first compare the performance of the optimal ICCP and JCCP.

*2) Performance of R-PRA: ICCP and JCCP:* Under the aforementioned low load scenario, the Optimal-ICCP violates the maximum allowable video degradation in case of large QoS levels (i.e., $\beta \geq 0.9$) as shown in Fig. 2(b). This is attributed to the ignored dependency between the allocations in the time slots. More specifically, the demand violation occurred at $t = 20\ s$ in Fig. 1(b) due to the low rate (shown in Fig. 1(a)), resulting in cumulative degradations in the following time slots. This is because the potential outage was not accounted for before hand. We can see that the buffer occupancy remained below the demand from $t = 20\ s$ to $t = 25\ s$ in Fig. 1(b) until the reallocation is done and the unmet demand is compensated. This violation was avoided for lower values of $\beta$ due to the

(a) Average BS airtime.
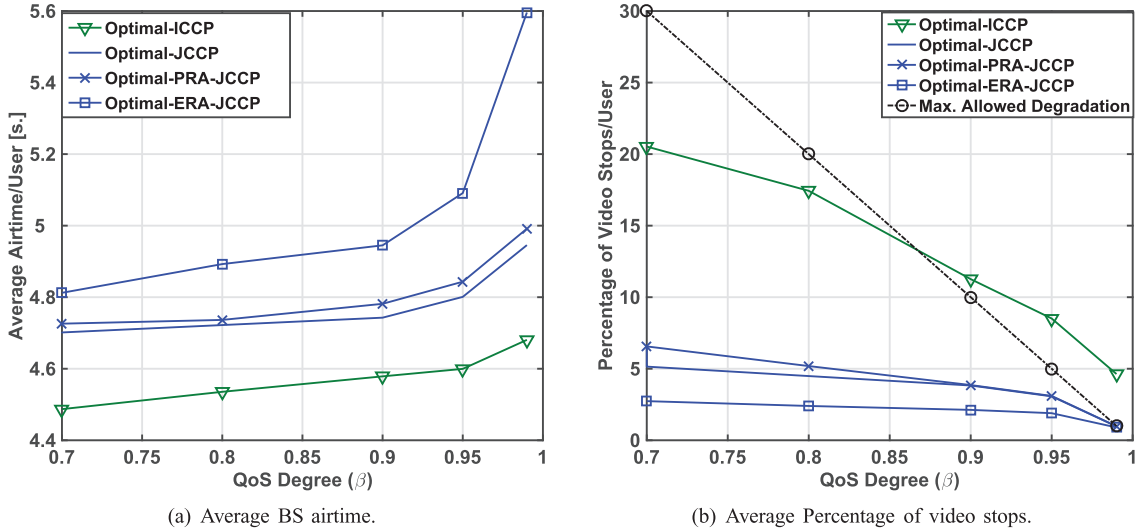


(b) Average Percentage of video stops.

Fig. 3.   Percentage of video stops and average BS airtime for varying QoS degrees $\beta$ for 4 Users experiencing slow fading with imperfect predictions.



(a) Users experiencing imperfect predictions, $\beta = 0.95$ and $V = 0.5 Mbps$



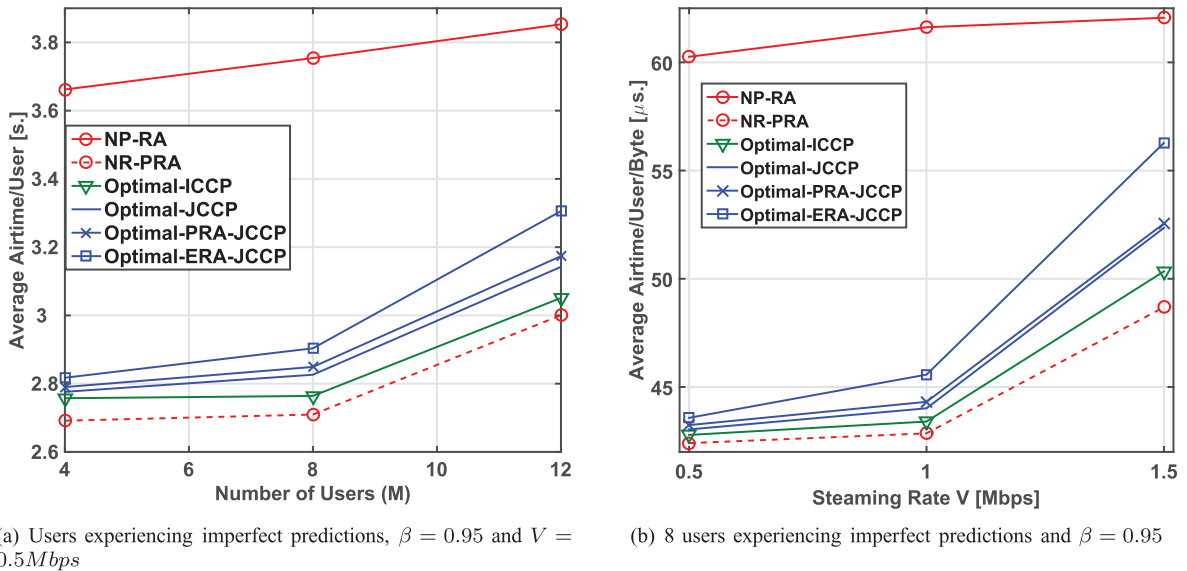(b) 8 users experiencing imperfect predictions and $\beta = 0.95$

Fig. 4.   Performance of Robust PRA for different simulation scenarios.

continuous feedback from the user every $\tau$ seconds that enabled the network to recover video outages.

On the other hand, all the JCCP forms: Optimal-JCCP, Optimal-ERA-JCCP and Optimal-PRA-JCCP were able to avoid the above propagation of video stops and thus did not violate the maximum allowed degradation at all QoS levels as shown in Fig. 2(b). This was done at the expense of energy savings (i.e., a higher price of robustness) compared to ICCP as depicted in Fig. 2(a). The results also demonstrate the ability of the decomposed convex forms of JCCP (Optimal-ERA-JCCP and Optimal-PRA-JCCP) to obtain a solution that satisfies the QoS level. However, compared to the global optimal solution, the Optimal-PRA-JCCP was able to satisfy the QoS level with less energy compared to the Optimal-ERA-JCCP. This result emphasizes the importance of optimizing the risk values over

TABLE II
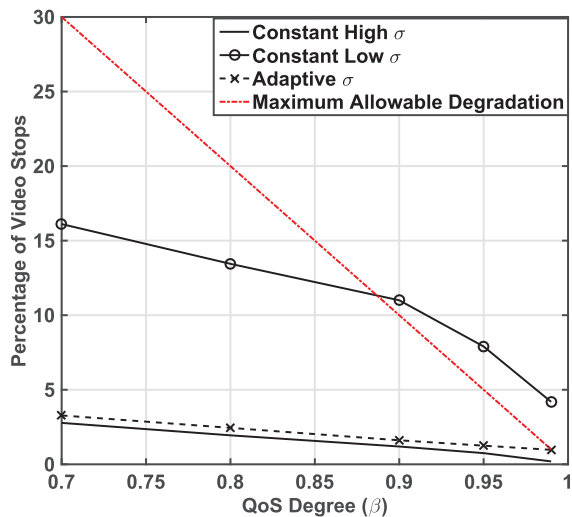OPTIMALITY GAP OF HEURISTIC ALGORITHMS

| Technique | Optimality Gap | | | |
|---|---|---|---|---|
| | 1 User | 4 Users | 8 Users | 12 Users |
| Heuristic-ICCP | 0.1 % | 0.15 % | 0.25 % | 0.3 % |
| Heuristic-ERA-JCCP | 0.1 % | 0.2 % | 0.5 % | 1.2 % |
| Heuristic-PRA-JCCP | 0.1 % | 0.15 % | 0.32 % | 0.45 % |

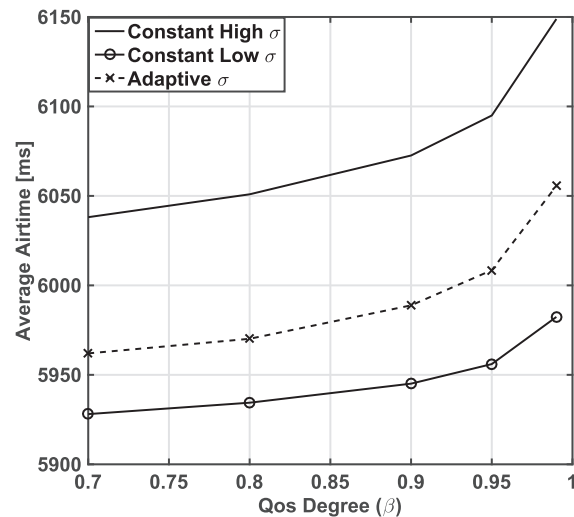the time horizon to control the conservatism of JCCP, especially when the user is located near the cell edge.

The performance results also indicates that the energy saving gap between the Optimal-PRA-JCCP and Optimal-ERA-JCCP increases with higher QoS levels ($\beta$), number of users and higher streaming rates as shown in Fig. 3(a), Fig. 4(a) and Fig. 4(b), respectively. In particular, as $\beta$ increases, lower risk

TABLE III
COMPLEXITY MEASURES FOR INTRODUCED ROBUST TECHNIQUES

| Technique | Order of Magnitude | Execution Time 1 User | Execution Time 12 Users |
|---|---|---|---|
| Optimal-ICCP | $O(\sqrt{2TM + T}(MT)^2(MT(T+1)/2 + M + 1)) \approx O(\sqrt{MT}(M^3T^4))$ | 90 *s.* | 980 *s.* |
| Heuristic-ICCP | $O(MT + T^2)$ | $< 1ms.$ | $< 1ms.$ |
| Optimal JCCP | $O(\sqrt{2TM + 2T}(2MT)^2(MT(T+1)/2 + M + 1 + T)) \approx O(\sqrt{MT}(M^3T^4))$ | 140*s.* | 1560*s.* |
| Optimal-ERA-JCCP | $O(\sqrt{2TM + T}(MT)^2(MT(T+1)/2 + M + 1)) \approx O(\sqrt{MT}(M^3T^4))$ | 90*s.* | 980*s.* |
| Heuristic-ERA-JCCP | $O(MT + T^2)$ | $< 1ms.$ | $< 1ms.$ |
| Optimal-PRA-JCCP | $O(\sqrt{2TM + T}(MT)^2(MT(T+1)/2 + M + 1) + M(T+1)^3) \approx O(\sqrt{MT}(M^3T^4))$ | 90*s.* | 980*s.* |
| Heuristic-PRA-JCCP | $O(MT + T^2 + M(T+1)^3) \approx O(M(T)^3)$ | $< 1ms.$ | $< 1ms.$ |



(a) Average Percentage of video stops.



(b) Average BS airtime.

Fig. 5. Percentage of video stops and average BS airtime for varying QoS degrees $\beta$ for 4 Users rate variations. Allocation is done using Heuristic-PRA-JCCP.

values are attained and the value of the inverse Q-function decreases exponentially which results in more airtime to satisfy C1 in (15). Similarly, increasing the number of users or streaming rate will result in more conservative RA for the cell edge users which decreases the BS airtime available for the cell center users to pre-buffer the video. It should be noted that the range of airtime varies across the scenarios since users follow different paths and velocities in each case.

*3) Optimality and Complexity Analysis:* In order to evaluate the introduced guided heuristic, the optimality gap $Z$ is measured between the heuristic based solutions and the optimal results as $Z = \frac{M(\mathbf{x}) - M(\mathbf{x}^*)}{M(\mathbf{x}^*)} \times 100$, where $M(\mathbf{x})$ and $M(\mathbf{x}^*)$ are the values of objective functions corresponding to the heuristic and optimal solutions, respectively. A small optimality gap indicates that the heuristic solution is very close to the optimal one.

From Table II we observe that the heuristic solutions can provide the energy savings with small optimality gaps. This performance degrades with an increased competition at the cell center due to either a large number of users located in the cell peak during the same slot or few residual airtime due to conservative allocation of cell edge users (the case of ERA-JCCP). In particular, increasing the number of users at the cell peak will increase the optimality gap since the residual resources

(after allocating the cell edge users) need to be proportionally allocated while considering the future rates. This was not handled by the heuristic algorithm to maintain its low complexity. Instead, the heuristic performs a greedy allocation to the users with the maximum rates. As for QoS satisfaction, the guided heuristic solutions follow the same performance trends as their corresponding optimal counterparts, i.e., the ICCP forms fail to satisfy the maximum degradation at high QoS levels while the JCCP succeed for all values.

We next analyze the computational complexity of the different allocation strategies. For SOCP formulations, the optimal solution, using interior point method, require a maximum of $O(\sqrt{K})$ iterations [37] where $K$ is the number of constraints. Each iteration has a complexity of $O(m^2 \sum_{i=1}^{K} n_i)$ [42], where $m$ denotes the total number of decision variables and $n_i$ is the dimension of the $i^{th}$ constraint. For the Newton's method, the main complexity lies in the calculation of the Hessian matrix inverse with a dimension $m \times m$. This gives a complexity of $O(m^3)$ for each step in Newton's method. Table III summarizes the two complexity measures for all the considered techniques as a function of the problem dimensions, i.e., number users $M$ and time slots $T$. For the heuristic in Algorithm 2, the QoS satisfaction has a complexity of $O(MT)$. The peak allocations and solution repairing have complexities of $O(M(T - t_p))$

and $O(MT)$, respectively. We also report the execution time measured within the simulation environment on a Quad Core i7-Processor, 3.2 GHz machine. These results highlight the incapability of the optimal solution methods to facilitate real-time implementation. It should be noted that increasing the number of users does not result in a proportional increase in execution time since the algorithms can be executed on multiple threads when there are multiple users. Moreover, the complexity of Newton's method which was executed for each user individually completes in less than 1 ms.

*4) Adaptive Variance Estimation:* The simulations were extended to test the robustness of the PRA framework to the *variations* in the channel variance. Such variations in the rate variance are typically observed in practical measurements due to the different landscapes and degrees of urbanization [8]. A conservative approach to tackle such variabilities is to optimize with a constantly large value (highest value revealed in simulations) for the rate variance. This will ensure meeting the QoS satisfaction level using JCCP as in Fig. 5(a). However, it compromises the energy efficiency as shown in Fig. 5(b). On the other hand, starting with a fixed lower value (smallest value revealed in simulations) of variance will result in less energy consumption but at the expense of QoS degradation even when JCCP is applied. The KF based tracking algorithm starts with an arbitrary value of variance, and then continuously adapts its value based on the error between the channel measurements and initial values. It is therefore able to satisfy the QoS for all values of $\beta$, and with a lower airtime compared to the high variance case. In this scenario, the evaluation is based on the Heuristic-PRA-JCCP since it has a practical complexity and results in more energy savings compared to the Heuristic-ERA-JCCP as highlighted previously.

## VIII. CONCLUSION

In this paper, we addressed the problem of *predictive* resource allocation for energy efficient video streaming. In contrast to previous efforts [12]–[16], we developed a *robust*-PRA framework with uncertainty in mind that provides *joint probabilistic* QoS guarantees. By offering a mechanism to control the probability constraint satisfaction, operators may control the trade-off between energy savings and the risks associated with erroneous predictions. Furthermore, in order to facilitate practical deployment, near-optimal real-time solutions coupled with a channel variation tracking technique were developed.

The performance evaluation results demonstrated the strong ability of the R-PRA to avoid the QoS violations exhibited by existing non-robust PRA approaches, which do not consider the impact of imperfect predictions. This QoS satisfaction was achieved while still providing significant energy savings compared to non-predictive RA. With regards to user satisfaction at low QoS levels, both the ICCP and JCCP were able to maintain the percentage of video stops below the maximum allowed value. However, the ICCP dominates the energy saving of JCCP as the latter appeared to be more conservative than needed. On the other hand, at high QoS levels, the demand accumulation over time slots hinders the user satisfaction by the low-energy ICCP, and the JCCP is the only feasible strategy.

The increased conservatism of JCCP necessitated optimal allocation of the QoS-risk over the time horizon to limit the excess energy consumptions. This robustness was at the expense of a slight increase in the complexity needed for calculating the risk of each time slot. The computational complexity was handled through guided heuristic which resulted in the same QoS performance without significant increase in complexity especially in case of conservative JCCP and high load scenarios. As PRA is an emerging resource allocation paradigm there are several directions for future work. This includes 1) extending the framework to solve non-invertible cumulative density functions that model rate uncertainty with irregular probability density functions, 2) cooperative channel variance estimation between the users to achieve a faster convergence to the optimal value, and 3) joint optimization of video quality and energy-efficiency.

## APPENDIX A

The objective function and all constraints in (15) are linear except the QoS one. The convexity of this first constraint will be checked using the Hessian matrix, which should be positive semidefinite [37]. Let the QoS constraint for the first user ($i = 0$) at time $t = 1$ be denoted as $f(x_{0,0}, x_{0,1}, \zeta_{0,1})$. In the standard form, the constraint is represented as follows

$$f(x_{0,0}, x_{0,1}, \zeta_{0,1}) = -\sum_{t'=0}^{1} \bar{r}_{0,t'} x_{0,t'} - Q^{-1}_{1-\zeta_{0,1}} \sqrt{\sum_{t'=0}^{1} x_{0,t'}^2 \sigma_{0,t'}^2} \tag{28}$$

For the ease of representation, let $f(x_{0,0}, x_{0,1}, \zeta_{0,1})$, $x_{0,0}$, $x_{0,1}$ and $\zeta_{0,1}$ be denoted as $\mathcal{F}$, $x_0$, $x_1$ and $\zeta$ respectively. The Hessian matrix $H$ can then be defined as follows

$$H = \nabla^2 \mathcal{F} = \begin{bmatrix} \frac{\partial^2 \mathcal{F}}{\partial x_0^2} & \frac{\partial^2 \mathcal{F}}{\partial x_0 \partial x_1} & \frac{\partial^2 \mathcal{F}}{\partial x_0 \partial \zeta} \\ \frac{\partial^2 \mathcal{F}}{\partial x_1 \partial x_0} & \frac{\partial^2 \mathcal{F}}{\partial x_1^2} & \frac{\partial^2 \mathcal{F}}{\partial x_1 \partial \zeta} \\ \frac{\partial^2 \mathcal{F}}{\partial \zeta \partial x_0} & \frac{\partial^2 \mathcal{F}}{\partial \zeta \partial x_1} & \frac{\partial^2 \mathcal{F}}{\partial \zeta^2} \end{bmatrix} \tag{29}$$

$$\frac{\partial^2 \mathcal{F}}{\partial x_0^2} = -Q^{-1}_{1-\zeta} \sigma_{0,0}^2 \frac{x_1^2 \sigma_{0,1}^2}{\left(\sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2}\right)^3} \tag{30}$$

$$\frac{\partial^2 \mathcal{F}}{\partial x_1^2} = -Q^{-1}_{1-\zeta} \sigma_{0,1}^2 \frac{x_0^2 \sigma_{0,0}^2}{\left(\sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2}\right)^3} \tag{31}$$

$$\frac{\partial^2 \mathcal{F}}{\partial \zeta^2} = -\frac{\partial^2 Q^{-1}_{1-\zeta}}{\partial \zeta^2} \sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2} \tag{32}$$

$$\frac{\partial^2 \mathcal{F}}{\partial x_0 \partial x_1} = \frac{\partial^2 \mathcal{F}}{\partial x_1 \partial x_0} = Q^{-1}_{1-\zeta} \sigma_{0,0}^2 \sigma_{0,1}^2 \frac{x_0 x_0, \sigma_{0,0}^2 \sigma_{0,1}^2}{\left(\sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2}\right)^3} \tag{33}$$

$$\frac{\partial^2 \mathcal{F}}{\partial x_0 \partial \zeta} = \frac{\partial^2 \mathcal{F}}{\partial \zeta \partial x_0} = -\frac{\partial Q^{-1}_{1-\zeta}}{\partial \zeta} \frac{x_0 \sigma_{0,0}^2}{\sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2}} \tag{34}$$

$$\frac{\partial^2 \mathcal{F}}{\partial x_1 \partial \zeta} = \frac{\partial^2 \mathcal{F}}{\partial \zeta \partial x_1} = -\frac{\partial Q^{-1}_{1-\zeta}}{\partial \zeta} \frac{x_1 \sigma_{0,1}^2}{\sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2}} \tag{35}$$

The function ($\mathcal{F}$) is convex if the Hessian matrix is positive semidefinite. In particular, all the principle minors should be positive or zero. The value of satisfaction degree of individual chance constraint (i.e., $\zeta$) should be less than 0.5 to satisfy the constraint (summation of $\zeta$) for $\beta > 0.5$. Accordingly, the inverse of Q function $Q_{1-\zeta}^{-1}$ is less than 0. Thus, all the first order principle minors are positive. The first second-order principle minor (constructed by deleting the third row and column) is always positive for all the values of $x_0$ and $x_1$, $\sigma_{0,0}$ and $\sigma_{0,0}$. However, this is not the case for the other second order principle minors whose positiveness depend on the actual values of $x_0$ and $x_1$, $\sigma_{0,0}$ and $\sigma_{0,0}$. For illustration, the value of a second order principle (constructed by deleting the second row and column) is calculated as follows

$$
\Delta_5 = \left( -Q_{1-\zeta}^{-1} \sigma_{0,0}^2 \frac{x_1^2 \sigma_{0,1}^2}{\left( \sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2} \right)^3} \right)
$$
$$
\left( -\frac{\partial^2 Q_{1-\zeta}^{-1}}{\partial \zeta^2} \sqrt{\sum_{t'=0}^{1} x_{t'}^2 \sigma_{0,t'}^2} \right) -
$$
$$
\left( -\frac{\partial Q_{1-\zeta}^{-1}}{\partial \zeta} \frac{x_0 \sigma_{0,0}^2}{\sqrt{\sum_{t'=0}^{1} x_{0,t'}^2 \sigma_{0,t'}^2}} \right)^2 \tag{36}
$$

It can be observed that $\Delta_5$ is only positive for specific values of allocation decisions and the variance. For instance, by assuming the variance $\sigma_0$ is greater than the variance $\sigma_1$, the second term will be greater than the first term, and thus $\Delta_5 < 0$. Accordingly, the Hessian matrix is neither positive nor negative semidefinite and hence the problem is non-convex.

## APPENDIX B

All the equations in (17) are linear and thus convex except the second constraint whose convexity is checked as follows

$$
F(y, \beta) = \sum_{\forall t \in \mathcal{T}} Q(y_{i,t}) - 1 + \beta
$$
$$
\nabla^2 F(y, \beta) = Q''(y_{i,t}) = \frac{1}{\sqrt{2\pi}} y_{i,t} e^{\frac{-y_{i,t}^2}{2}}.
$$

Since we assume $\beta \geq 0.5$ for practical QoS levels, the constraint holds iff $\sum_{\forall t \in \mathcal{T}} Q(y_{i,t}) \leq 0.5$. This implies that $Q(y_{i,t}) \leq 0.5$ which occurs when $y_{i,t} \geq 0$. The Hessian matrix is a diagonal matrix of positive entries that represents its eigenvalues. Accordingly, the Hessian matrix is positive semidefinite and this proves the convexity of function.

## REFERENCES

[1] K. Davaslioglu and E. Ayanoglu, "Quantifying potential energy efficiency gain in green cellular wireless networks," *IEEE Commun. Surveys Tuts.*, vol. 16, no. 4, pp. 2065–2091, May 2014.

[2] L. Correia *et al.*, "Challenges and enabling technologies for energy aware mobile radio networks," *IEEE Commun. Mag.*, vol. 48, no. 11, pp. 66–72, Nov. 2010.

[3] M. A. Marsan, L. Chiaraviglio, D. Ciullo, and M. Meo, "Optimal energy savings in cellular access networks," in *Proc. IEEE Int. Conf. Commun. (ICC) Workshops*, 2009, pp. 1–5.

[4] E. Oh and B. Krishnamachari, "Energy savings through dynamic base station switching in cellular wireless access networks," in *Proc. IEEE GLOBECOM*, 2010, pp. 1–5.

[5] iGR. (2013). *U.S. Regional and Small Operator Network Infrastructure Capex and Opex Forecast, 2012–2017* [Online]. Available: https://igr-inc.com/, accessed on Mar. 29, 2015.

[6] CISCO. (2014). *Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2013–2018* [Online]. Available: http://www.cisco.com/c/en/us/solutions/service-provider/visual-networking-index-vni/index.html, accessed on Mar. 29, 2015.

[7] A. Schulman *et al.*, "Bartendr: A practical approach to energy-aware cellular data scheduling," in *Proc. ACM Mobicom*, 2010, pp. 85–96.

[8] H. Abou-Zeid, H. S. Hassanein, Z. Tanveer, and N. AbuAli, "Evaluating mobile signal and location predictability along public transportation routes," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, 2015, pp. 1195–1200.

[9] M. Elazab, A. Noureldin, and H. Hassanein, "Integrated cooperative localization for connected vehicles in urban canyons," in *Proc. IEEE Globecom*, 2015, pp. 1–6.

[10] A. Mahmoud, A. Noureldin, and H. Hassanein, "VANETs positioning in urban environments: A novel cooperative approach," in *Proc. IEEE Veh. Technol. Conf. (VTC-Fall)*, 2015, pp. 1–7.

[11] C. Brunner and D. Flore, "Generation of pathloss and interference maps as SON enabler in deployed UMTS networks," in *Proc. IEEE Veh. Technol. Conf. (VTC-Spring)*, 2009, pp. 1–5.

[12] H. Abou-Zeid, H. S. Hassanein, and S. Valentin, "Energy-efficient adaptive video transmission: Exploiting rate predictions in wireless networks," *IEEE Trans. Veh. Technol.*, vol. 63, no. 5, pp. 2013–2026, Jun. 2014.

[13] H. Abou-Zeid and H. S. Hassanein, "Toward green media delivery: Location-aware opportunities and approaches," *IEEE Wireless Commun.*, vol. 21, no. 4, pp. 38–46, Aug. 2014.

[14] Z. Lu and G. de Veciana, "Optimizing stored video delivery for mobile networks: The value of knowing the future," in *Proc. IEEE INFOCOM*, 2013, pp. 2806–2814.

[15] H. Abou-Zeid and H. S. Hassanein, "Predictive green wireless access: Exploiting mobility and application information," *IEEE Wireless Commun.*, vol. 20, no. 5, pp. 92–99, Oct. 2013.

[16] R. Margolies *et al.*, "Exploiting mobility in proportional fair cellular scheduling: Measurements and algorithms," in *Proc. IEEE INFOCOM*, 2014, pp. 1339–1347.

[17] N. Bui and J. Widmer, "Modelling throughput prediction errors as Gaussian random walks," in *Proc. KuVS*, 2014, pp. 1–3.

[18] P. Kali and S. W. Wallace, *Stochastic Programming*. New York, NY, USA: Springer, 1994.

[19] A. Charnes and W. W. Cooper, "Chance-constrained programming," *Manage. Sci.*, vol. 6, no. 1, pp. 73–79, 1959.

[20] B. L. Miller and H. M. Wagner, "Chance constrained programming with joint constraints," *Oper. Res.*, vol. 13, no. 6, pp. 930–945, 1965.

[21] B. Nunez, P. Adasme, I. Soto, J. Cheng, M. Letournel, and A. Lisser, "A chance constrained approach for uplink wireless OFDMA networks," in *Proc. 9th Int. Symp. Commun. Syst. Netw. Digit. Signal Process. (CSNDSP)*, 2014, pp. 754–757.

[22] S. S. Venkatesh, *The Theory of Probability: Explorations and Applications*. Cambridge, U.K.: Cambridge Univ. Press, 2012.

[23] M. Ono and B. C. Williams, "Iterative risk allocation: A new approach to robust model predictive control with a joint chance constraint," in *Proc. 47th IEEE Conf. Decision Control (CDC)*, 2008, pp. 3427–3432.

[24] F. Oldewurtel, C. N. Jones, and M. Morari, "A tractable approximation of chance constrained stochastic MPC based on affine disturbance feedback," in *Proc. IEEE Conf. Decision Control (CDC)*, 2008, pp. 4731–4736.

[25] U. A. Ozturk, M. Mazumdar, and B. A. Norman, "A solution to the stochastic unit commitment problem using chance constrained programming," *IEEE Trans. Power Syst.*, vol. 19, no. 3, pp. 1589–1598, Aug. 2004.

[26] M. Ono and B. C. Williams, "Decentralized chance-constrained finite-horizon optimal control for multi-agent systems," in *Proc. IEEE Conf. Decision Control (CDC)*, 2010, pp. 138–145.

[27] N. Bui, F. Michelinakis, and J. Widmer, "A model for throughput prediction for mobile users," in *Proc. Eur. Wireless (EW)*, 2014, pp. 1–6.

[28] N. Y. Soltani, S.-J. Kim, and G. B. Giannakis, "Chance-constrained optimization of OFDMA cognitive radio uplinks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 3, pp. 1098–1107, Mar. 2013.

[29] R. Atawia, H. Abou-Zeid, H. Hassanein, and A. Noureldin, "Robust resource allocation for predictive video streaming under channel uncertainty," in *Proc. IEEE GLOBECOM*, Dec. 2014, pp. 4683–4688.

[30] S. Chen, "Comparing probabilistic and fuzzy set approaches for designing in the presence of uncertainty," Ph.D. thesis, Virginia Polytechnic Institute and State University, Blacksburg, VA, USA, 2000.

[31] 3GPP, "LTE; Evolved universal terrestrial radio access (E-UTRA); Physical layer procedures," Technical Specification TS 36.213 v12.5.0, 2015.

[32] C. Desset, B. Debaillie, V. Giannini, A. Fehske, G. Auer, and H. Holtkamp, "Flexible power modeling of LTE base stations," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, Apr. 2012, pp. 2858–2862.

[33] H. Holtkamp, G. Auer, S. Bazzi, and H. Haas, "Minimizing base station power consumption," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 2, pp. 297–306, Feb. 2014.

[34] A. ParandehGheibi, M. Médard, A. Ozdaglar, and S. Shakkottai, "Avoiding interruptions—A QoE reliability function for streaming media applications," *IEEE J. Sel. Areas Commun.*, vol. 29, no. 5, pp. 1064–1074, May 2011.

[35] 3GPP, "LTE; Evolved universal terrestrial radio access (E-UTRA); Further advancements for E-UTRA physical layer aspects," Technical Report TR 36.814 V9.0.0, 2010.

[36] 3GPP, "E-UTRA; Base station (BS) radio transmission and reception (release 10)," Technical Specification TS 36.104 V10.2.0, Dec. 2011.

[37] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.

[38] Z.-Q. Luo and W. Yu, "An introduction to convex optimization for communications and signal processing," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 8, pp. 1426–1438, Aug. 2006.

[39] A. Noureldin, T. B. Karamat, and J. Georgy, *Fundamentals of Inertial Navigation, Satellite-Based Positioning and Their Integration*. New York, NY, USA: Springer, 2013.

[40] H. Abou-Zeid, H. S. Hassanein, and R. Atawia, "Towards mobility-aware predictive radio access: Modeling; Simulation; and evaluation in LTE networks," in *Proc. ACM Int. Conf. Model. Anal. Simul. Wireless Mobile Syst. (MSWiM)*, 2014, pp. 109–116.

[41] R. Kwan, C. Leung, and J. Zhang, "Multiuser scheduling on the downlink of an LTE cellular system," *Res. Lett. Commun.*, vol. 2008, no. 3, pp. 1–4, 2008.

[42] M. S. Lobo, L. Vandenberghe, S. Boyd, and H. Lebret, "Applications of second-order cone programming," *Linear Algebra Appl.*, vol. 284, no. 1, pp. 193–228, 1998.

**Ramy Atawia** (S'12) received the B.Sc. and M.Sc. degrees (both with Hons.) in communication engineering from the German University in Cairo (GUC), Cairo, Egypt, in 2012 and 2013, respectively. He is currently pursuing the Ph.D. degree at Queen's University, Kingston, ON, Canada. He is also a Research and Teaching Assistant with Queen's University. His research interests include context-aware radio resource management, robust optimization, and network planning. He is serving as a TPC member and Reviewer in the IEEE flagship conferences and journals. He is currently an R&D Intern at Alcatel-Lucent (Nokia), Belgium.

**Hatem Abou-zeid** (S'04–M'14) received the Ph.D. degree in electrical and computer engineering from Queen's University, Kingston, ON, Canada, in 2014 (where his research was recognized with a medal nomination for innovation). He has been a Software Engineer with the CISCO R&D Center, Ottawa, Canada, since 2015. He has contributed to the development of scalable traffic engineering and IP routing protocols for service provide and data center networks. He also held the position of Postdoctoral Fellow prior to joining CISCO. His research interests include SDNs, context-aware radio access networks, and adaptive video delivery for vehicular communications. He is also an experienced Lecturer and has been granted several teaching fellowships with Queen's University to instruct freshman and senior-level engineering courses. He is also a Technical Reviewer for several prestigious conferences and journals.

**Hossam S. Hassanein** (S'86–M'90–SM'05) is a leading authority in the areas of broadband, wireless and mobile networks architecture, protocols, control and performance evaluation. His record spans more than 500 publications in journals, conferences, and book chapters, in addition to numerous keynotes and plenary talks in flagship venues. He is the Founder and the Director of the Telecommunications Research Laboratory (TRL), Queen's University School of Computing, with extensive international academic and industrial collaborations. He is a Former Chair of the IEEE Communication Society Technical Committee on Ad hoc and Sensor Networks (TC AHSN). He is the Distinguished Speaker of the IEEE Communications Society (Distinguished Lecturer 2008–2010). He has received several recognitions and best paper awards at top international conferences.

**Aboelmagd Noureldin** (S'98–M'02–SM'08) received the B.Sc. degree in electrical engineering and the M.Sc. degree in engineering physics from Cairo University, Giza, Egypt, in 1993 and 1997, respectively, and the Ph.D. degree in electrical and computer engineering from the University of Calgary, Calgary, AB, Canada, in 2002. He is a Professor with the Department of Electrical and Computer Engineering, Royal Military College of Canada (RMCC), Kingston, Ontario, Canada. He is also the Founder and the Director of the Navigation and Instrumentation Research Group, RMCC. His research interests include GPS, wireless location and navigation, indoor positioning and multisensor fusion. He has authored more than 200 papers published in journals and conference proceedings. His research work led to several patents in the area of position, location and navigation.