# Joint Optimization of Relay Strategies and Resource Allocations in Cooperative Cellular Networks

Truman Chiu-Yam Ng and Wei Yu, *Member, IEEE*

*Abstract*— This paper considers a wireless cooperative cellular data network with a base station and many subscribers in which the subscribers have the ability to relay information for each other to improve the overall network performance. For a wireless network operating in a frequency-selective slow-fading environment, the choices of relay node, relay strategy, and the allocation of power and bandwidth for each user are important design parameters. The design challenge is compounded further by the need to take user traffic demands into consideration. This paper proposes a centralized utility maximization framework for such a network. We show that for a cellular system employing orthogonal frequency-division multiple-access (OFDMA), the optimization of physical-layer transmission strategies can be done efficiently by introducing a set of pricing variables as weighting factors. The proposed solution incorporates both user traffic demands and the physical channel realizations in a cross-layer design that not only allocates power and bandwidth optimally for each user, but also selects the best relay node and best relay strategy (i.e. decode-and-forward vs. amplify-and-forward) for each source-destination pair.

*Index Terms*— Cooperative communication, Lagrangian duality theory, network utility maximization, orthogonal frequency-division multiplex (OFDM), relay channel, spectrum optimization, wireless cellular networks.

## I. INTRODUCTION

**I**N a wireless network with many source-destination pairs, cooperative transmission by relay nodes has the potential to improve the overall network performance. When a relay or a group of relays are physically located between the source and the destination, the relays may facilitate transmission by first decoding the transmitted codeword, then forwarding the decoded codeword to the destination. This strategy is known as decode-and-forward (DF). Alternatively, a relay may simply amplify its received signal and employ a so-called amplify-and-forward (AF) strategy. In both cases, the use of relay has been shown to improve the overall transmission rate [1] and/or the diversity [2]–[7] of the wireless network.

This paper is motivated by the following questions: In a cooperative wireless network, which node should act as a relay? What relay strategy should be used? When, and in which frequency should relaying be employed? Clearly, the answers to these questions depend on the topology of the network. For example, consider a single-relay channel
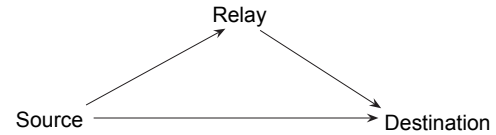
Fig. 1. A single-relay channel.

as shown in Fig. 1. When a relay is located closer to the source than to the destination, DF appears to be a natural choice. On the other hand, when the relay is located closer to the destination, its received signal-to-noise ratio (SNR) may not be high enough to allow decoding, in which case AF is more suited. However, the optimal operation of relays is more complicated than that provided by a rule-of-thumb. This is because the choices of relay node and relay strategy also depend on the amount of transmit power available at the source and at the relay, and further on the user traffic patterns. This is particularly true in a power- and bandwidth-limited network in which each node may act as a source/destination or relay simultaneously. In this case, the partitioning of the power and bandwidth between the transmission of one's own data vs. the relaying of other users' traffic becomes crucial. The optimal power and bandwidth allocation is further coupled with the choice of relay and the choice of relay strategies.

This paper takes a *system* view of the cooperative network, and aims to jointly optimize relay strategies and physical-layer resources in a network. The focus of this work is on a cellular data network with a single base station and many subscribers in each cell, where each subscriber has the ability to relay information for each other. This is called a multihop cellular network with peer-to-peer relaying [8]–[11]. However, the term multi-hop cellular usually refers to networks in which the relays simply route traffic from node to node. In this paper, we consider cooperative networks that use physical-layer relaying strategies, which take advantage of the broadcast nature of wireless channel, and allow the destination to cooperatively "combine" signals sent by both the source and the relay.

The target application of this paper is a fixed broadband access network in which the wireless channels are relatively stationary and channel estimations are feasible, and where centralized power control and bandwidth partitioning can be implemented. The use of cooperative strategies in a cellular network has many perceived benefits, e.g., throughput increase, coverage extension, power saving, and interference mitigation, etc. This paper presents an optimization framework to quantify some of these benefits.

In a cellular data network, bidirectional upstream and downstream transmissions are maintained for each subscriber. In a

traditional system with power control, subscribers closer to the base station transmit at a lower power level so they do not cause excess interference to other users. The users at the edge of the cell transmit at its maximum power limit. The main idea of a cooperative cellular network is to take advantage of the extra power available at users close to the base station by allowing them to act as a relay for users at the edge of the cell. The objective of this paper is to provide an efficient procedure for a joint optimization of the choice of relay and the choice of relay strategy for each transmission pair, and the allocation of power and bandwidth for each transmitting node in the network. The rest of this section contains a literature survey and highlights the main contributions of this paper.

### A. Related Work

Many well-known relay strategies can trace their roots to the information-theoretical study of the relay channel by Cover and El Gamal [12]. Cover and El Gamal showed that the capacity of a degraded relay channel is achieved by a block-Markov scheme with an infinite number of blocks and a decode-and-forward strategy at the relay. They further proposed a quantize-and-forward strategy for general relay channels. Amplify-and-forward can be interpreted as an analog counterpart to quantize-and-forward. Depending on the channel condition and performance criterion, DF may outperform AF, or vice versa [13]–[15].

In practical systems, the restriction that the relay node may not transmit and receive at the same time in the same frequency band is often imposed, giving rise to a so-called cheap relay channel [16], [17]. Moreover, the block-Markov scheme is often simplified into a two-time-slot strategy. These simplifications, although strictly suboptimal from an information-theoretical perspective, nevertheless make relay operations much easier to implement. The rest of the paper focuses on two-time-slot implementations of DF and AF strategies for cheap relay nodes [1], [18].

This paper envisions a type of relay-assisted network called cooperative network where user nodes act as relays for one another. We further assume that the cooperative network employs orthogonal frequency-division multiple-access (OFDMA), so that each user node may simultaneously act as a source, a destination, or as a relay, but at different frequency tones. For frequency-selective fading channels, power allocation among transmitting nodes and across frequency tones can greatly affect network performance. (The same channel model is also applicable to frequency-flat fading channels with adaptive modulations across fading states.) The issue of power control for relay-assisted networks has been dealt with in numerous studies in the literature. These studies all point to the importance of power allocation and they differ in performance criteria (rates [19] vs. probabilities of outage or error [20]), power constraints (separate- [21] vs. joint- [22] power constraints between the source and the relay), power control schemes (centralized [23] vs. distributed [24]), and system setups (single- [25] vs. multi-stage-relay [26], single- vs. multi-parallel-relay [27], three-terminal single-relay [28] vs. multi-nodes (ad-hoc, cellular) [29]). Along with power allocation, some studies listed above (e.g. [20], [27]) also jointly consider relay-node selection.

Besides power allocation and relay operations, bandwidth allocations for transmission nodes are equally important. Note that in some relay channel literature such as [19], [23], [30]–[33], bandwidth allocations refer to the assignments of orthogonal frequencies for transmissions of the source and the relay(s). However, this is *not* the type of bandwidth allocation that is investigated here. In this paper, each frequency tone is occupied by only one source-destination pair (with the possibility of relay participation), and both the source and the relay always transmit in the same frequency tone. Thus, bandwidth allocations in this paper refer to the assignments of source-destination pair to each frequency tone.

Despite the rich literature on relay-assisted networks, to the best of our knowledge, none of the previous work deals with the joint optimization of relay selection, relay-strategy selection, power and bandwidth allocation in a cooperative network employing physical-layer relaying strategies. This is in part due to the fact that such an optimization problem looks formidable at a first glance as so many different combinations of resource allocations and relay operations are possible. In addition, as the joint optimal scheme should also account for users' traffic demands, the optimization problem represents a *cross-layer* design. The main contribution of this paper is that the seemingly difficult joint optimization problem can be solved globally and efficiently within a network utility maximization framework. A key technique in this study is the use of pricing to determine the optimal relay strategies. Pricing has been used in earlier studies of both multi-hop hotspot network [34] and ad-hoc relay network [35], where prices are used to give user nodes an incentive to relay. However, in this paper, we introduce a centralized optimization framework, in which the network has control over the behavior of the user nodes. Therefore, the pricing variables in this paper only serve as weighting factors in the regulation of system resources.

### B. Summary of Contribution

This paper proposes a centralized optimization framework for the maximization of total system utility. The target system is a cooperative cellular network where subscribers act as relays for one another using physical-layer relaying strategies (e.g. DF, AF). The physical layer uses OFDM so that both power and bandwidth may be freely allocated among all nodes (but with the restriction that in each frequency tone, only one set of source and destination, plus possibly a relay, may be active). The main contributions of the paper are as follows:

- We show that the cross-layer problem may be decomposed into two subproblems: a utility maximization problem in the application layer and a joint relay-strategy and relay-node selection and power- and bandwidth-allocation problem in the physical layer. A set of dual variables coordinates the application-layer *demand* and physical-layer *supply* of rates.
- We further show that the physical-layer joint relay-strategy and relay-node selection and resource allocation problem may be solved globally and efficiently using another set of dual variables that accounts for the cost of power expended at each node.
- Together, these two dual decomposition steps solve the overall utility maximization problem globally with a

complexity that is linear in the number of frequency tones and the number of relay strategies, and quadratic in the number of relays.

### C. Outline of the Paper

The remainder of the paper is organized as follows. In Section II, we present a utility maximization framework for the joint resource allocation and relay-strategy and relay-node selection problem. We show that the overall cross-layer problem can be decomposed into two subproblems. Section III answers the question of which relay strategy is optimal for each source-destination, while taking both power constraints and users' traffic demands into consideration. Section IV contains simulation results that quantify the merit of relaying and illustrate the allocations of resource and operations of relays in an optimized cooperative network. Finally, conclusions are drawn in Section V.

### D. Notation

Boldface lower-case letters are used for column vectors. Vectors $\mathbf{0}$ and $\mathbf{1}$ denote the all-zero and all-one column vectors respectively. For two vectors of the same length, "$\succeq$" and "$\preceq$" are used to denote component-wise inequalities. Lower-case letter $x_i$ denotes the $i^{th}$ entry of vector $\boldsymbol{x}$. Boldface upper-case letters are used to denote matrices. For a matrix $\boldsymbol{X}$, $X(i,j)$ represents the entry on the $i^{th}$ row and the $j^{th}$ column. $X(i,:)$ denotes the $i^{th}$ row of $\boldsymbol{X}$, and $X(:,j)$ denotes the $j^{th}$ column of $\boldsymbol{X}$. The superscripts $^T$, $^H$, and $^*$ denote the transpose, the Hermitian, and the optimal value of the variable respectively.

## II. OPTIMIZATION USING PHYSICAL-LAYER RELAYING

This paper focuses on physical-layer relaying strategies where the source and the relay transmit in the same frequency, but in two different time-slots. Physical-layer relaying makes use of the broadcast nature of wireless channels, so that at each frequency tone, the destination receives and processes data sent by the source and the relay in both time-slots. Both DF and AF strategies are considered.

### A. Utility Maximization Framework

This paper adopts a utility maximization framework in which each data stream has an associated utility function, and the objective is to maximize the sum utility in the network. This network utility maximization framework is originated from the work of Kelly [36], and has been recently applied to many physical-layer design problems (e.g. [37], [38]). A utility function is a concave and increasing function of data rates that reflects user satisfaction. The choice of utility function depends on the underlying application (e.g. data, voice, video). By maximizing the sum utility, a network operator maintains a balance between competing resource demands by different users in a network. This paper assumes that the base station dictates when and how the user nodes may cooperate, which is realistic in a cellular environment with centralized control.

The cooperative cellular network of interest is shown in Fig. 2. The network consists of a base station and $K$ user nodes, all of which are equipped with a single antenna. Let $\mathcal{K} =$
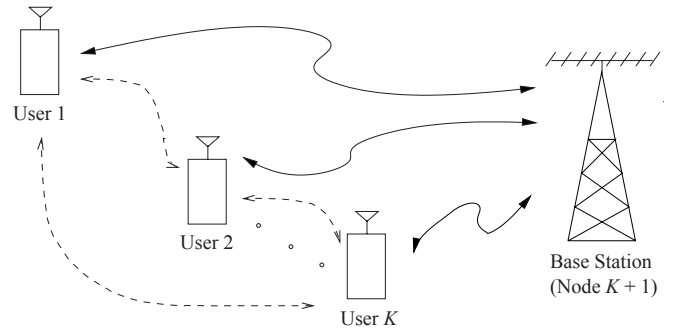


Fig. 2.   A cooperative cellular network with a base station and $K$ users.

$\{1, 2, ..., K\}$ be the set of user nodes. Denote the base station as node $K + 1$. Let $\mathcal{K}_+ = \{1, 2, ..., K + 1\}$ be the extended set of nodes. Each of the $K$ user nodes has both downstream and upstream communications with the base station. Let $(s, d)$ be a source-destination pair, or data stream. Then, the set of all data streams $\mathcal{M}$ is $\{(1, K + 1), (2, K + 1), ..., (K, K + 1), (K + 1, 1), (K + 1, 2), ..., (K + 1, K)\}$. The cardinality of $\mathcal{M}$ is $2K$.

We assume that the cooperative network employs an OFDMA physical-layer with $N$ tones. Let $\mathcal{N} = \{1, ..., N\}$ denote the set of tones. We further assume that the OFDM frames are synchronized throughout the network, so that cooperative transmission takes place in a tone-by-tone basis, and transmissions in different frequency tones do not interfere with each other. The wireless channel is modelled as a frequency-selective fading channel with coherence bandwidth in the order of the bandwidth of a few tones. This implies that fading between tones far away from each other is un-correlated. We assume that the network operates in a slow fading environment, so that channel estimation is possible and full channel side-information (CSI) is available at the base station. Although this paper focuses on a frequency-selective environment employing OFDMA, with suitable modifications, the optimization framework is also applicable to a frequency-flat fading environment if rate adaptation across fading states is possible.

To prevent inter-stream interference, the restriction that only one data stream can be active in each tone is imposed. However, each of the $2K$ data streams can use more than one frequency tones. We also impose the restriction that the active stream in each tone can use at most one relay (where the relay node can be any one of the other $K - 1$ user nodes). Moreover, the source-relay, source-destination, and relay-destination links all use the same frequency tone[1]. Note that upstream and downstream transmissions take place simultaneously in the network, so that the base station is both a source and a destination (but in different tones). However, the base station can never be a relay. On the other hand, all user nodes in the network can simultaneously be a source, a relay, and a destination (again in different tones).

The main challenge in designing such a network is to allocate power and bandwidth optimally across all frequency

---

[1]The generalization of this setup to the case where source-relay and relay-destination communications may occur in different tones is nontrivial. Computational complexity of the optimization increases considerably.

tones and across all nodes. Let $\boldsymbol{P}$ be a $(K + 1) \times N$ matrix such that $P(i, n)$ denotes the power expended by node $i$ in tone $n$ ($i \in \mathcal{K}_+, n \in \mathcal{N}$). Because only one data stream is active in each tone and only the source and the relay may transmit, the column vector $P(:, n)$ ($n \in \mathcal{N}$) has at most two non-zero entries. Similarly, let $\boldsymbol{R}$ be a $2K \times N$ matrix with $R(m, n)$ denoting the data rate of stream $m$ in tone $n$ ($m \in \mathcal{M}$, $n \in \mathcal{N}$). Since only one stream can be active in each tone, the column vector $R(:, n)$ ($n \in \mathcal{N}$) has at most one non-zero entry. Power and rate are related by the *achievable rate region*, denoted as $\boldsymbol{R} \in \mathcal{C}(\boldsymbol{P})$. The achievable rate region is the set of all possible rates achievable at a given power level. It implicitly accounts for the best possible use of relay strategies.

By definition, $(\boldsymbol{P1})_i$ ($i \in \mathcal{K}_+$), i.e. the row sum of $\boldsymbol{P}$, is the total power expended at node $i$, summing across all tones. For practical purposes, each node has a separate power constraint. Let $\boldsymbol{p^{max}} = [p_1^{max}, p_2^{max}, ..., p_{K+1}^{max}]^T$, where $p_i^{max}$ is the individual power constraint for node $i$. Similarly, $(\boldsymbol{R1})_m$ ($m \in \mathcal{M}$), i.e. the row sum of $\boldsymbol{R}$, gives the total data rate of stream $m$, summing across all tones.

Let $U_m$ be the utility function of data stream $m$, where $U_m$ is a function of achievable rate of stream $m$, i.e. $(\boldsymbol{R1})_m$. The objective of the optimization problem is to optimally choose the active data stream and allocate power in each tone, while at the same time selecting the best relay node and the best relay strategy, in order to maximize the network sum utility. Expressed succinctly, the optimization problem is

$$\max_{\boldsymbol{P}, \boldsymbol{R}} \quad \sum_{m \in \mathcal{M}} U_m((\boldsymbol{R1})_m) \qquad (1)$$
$$\text{s.t.} \quad \boldsymbol{P1} \preceq \boldsymbol{p^{max}}, \quad \boldsymbol{R} \in \mathcal{C}(\boldsymbol{P})$$

As mentioned earlier, this paper assumes a slow-fading environment with full CSI in which adaptive power- and bit-loading may be implemented in a centralized fashion. In a practical system, channel estimation needs to be done at the receivers and fedback to the base station, which then must solve the above optimization problem and inform all users the appropriate power levels and cooperative strategies. When the wireless fading channel has both fast-fading and slow-fading components, optimization may also be done in a statistical sense by adapting power allocation to the average SNR (or the slow-fading component of the channel). Online implementation may be further simplified by constructing a look-up table that maps the channel gains to the optimal cooperative strategies.

### B. Cross-layer Optimization via Dual Decomposition

Finding the optimal solution of (1) involves a search over all possible power and bandwidth allocations and over all possible relays and relay strategies. So, the optimization problem (1) is a mixed integer programming problem. However, in an OFDMA system with many narrow subcarriers, the optimal solution of (1) is always a convex function of $\boldsymbol{p^{max}}$, because the time-sharing of two transmission strategies can always be implemented via frequency-division multiplexing across frequency tones. The idea is that if two sets of rates using two different resource allocations and relay strategies are achievable individually, then their linear combination is also achievable by a frequency-division multiplex of the two sets of strategies. This is possible when the coherence bandwidth is larger than the width of a few tones, so that adjacent tones have similar channel conditions, and time-sharing can be achieved in the frequency domain. The idea of getting convexity through frequency-sharing is discussed earlier in [39] for a spectrum balancing problem. In particular, using the duality theory of [39], the following is true:

*Proposition 1:* The optimization problem (1) has zero duality gap in the limit as the number of OFDM tones goes to infinity. This is true even as discrete selections of data stream, relay, relaying strategy, and bit rate are involved.

The paragraph leading to the proposition contains the main idea of the proof. A detailed proof can be constructed along a line of argument as in [39]. The zero-duality-gap result opens the door for using *convex* optimization techniques for solving the utility maximization problem (1). The rest of this section develops a Lagrangian dual decomposition approach to solve (1).

First, introduce a new variable $\boldsymbol{t} = [t_{(1,K+1)}, t_{(2,K+1)}, ..., t_{(K,K+1)}, t_{(K+1,1)}, t_{(K+1,2)}, ..., t_{(K+1,K)}]^T$, and rewrite (1) as

$$\max_{\boldsymbol{P}, \boldsymbol{R}, \boldsymbol{t}} \quad \sum_{m \in \mathcal{M}} U_m(t_m) \qquad (2)$$
$$\text{s.t.} \quad \boldsymbol{P1} \preceq \boldsymbol{p^{max}}, \quad \boldsymbol{R1} \succeq \boldsymbol{t}, \quad \boldsymbol{R} \in \mathcal{C}(\boldsymbol{P}).$$

Because $U_m$ is an increasing function, when the objective of (2) is maximized, $\boldsymbol{t}$ must be equal to $\boldsymbol{R1}$. Therefore, (1) and (2) must have the same solution. The key step in dual decomposition is to relax the constraint $\boldsymbol{R1} \succeq \boldsymbol{t}$. The Lagrangian becomes

$$L(\boldsymbol{P}, \boldsymbol{R}, \boldsymbol{t}, \boldsymbol{\lambda}) = \sum_{m \in \mathcal{M}} U_m(t_m) + \boldsymbol{\lambda^T} \left( \boldsymbol{R1} - \boldsymbol{t} \right)$$
$$= \sum_{m \in \mathcal{M}} \left( U_m(t_m) + \lambda_m \left( \sum_{n \in \mathcal{N}} R(m, n) - t_m \right) \right), \quad (3)$$

where $\boldsymbol{\lambda} = [\lambda_{(1,K+1)}, \lambda_{(2,K+1)}, ..., \lambda_{(K,K+1)}, \lambda_{(K+1,1)}, \lambda_{(K+1,2)}, ..., \lambda_{(K+1,K)}]^T$, with each element $\lambda_m$ ($m \in \mathcal{M}$) being a dual variable corresponding to stream $m$. Observe that the dual function

$$g(\boldsymbol{\lambda}) = \begin{cases} \max_{\boldsymbol{P}, \boldsymbol{R}, \boldsymbol{t}} & L(\boldsymbol{P}, \boldsymbol{R}, \boldsymbol{t}, \boldsymbol{\lambda}) \\ \text{s.t.} & \boldsymbol{P1} \preceq \boldsymbol{p^{max}}, \quad \boldsymbol{R} \in \mathcal{C}(\boldsymbol{P}) \end{cases} \qquad (4)$$

consists of two sets of variables: application-layer variable $\boldsymbol{t}$, and physical-layer variables $\boldsymbol{P}$ and $\boldsymbol{R}$. Moreover, $g(\boldsymbol{\lambda})$ can be separated into two maximization subproblems, namely a utility maximization problem, corresponding to a rate adaptation problem in the application layer,

$$g_{appl}(\boldsymbol{\lambda}) = \max_{\boldsymbol{t}} \sum_{m \in \mathcal{M}} \left( U_m(t_m) - \lambda_m t_m \right), \qquad (5)$$

and a joint relay-strategy and relay-node selection and power and bandwidth allocation problem in the physical layer,

$$g_{phy}(\boldsymbol{\lambda}) = \begin{cases} \max_{\boldsymbol{P}, \boldsymbol{R}} & \sum_{m \in \mathcal{M}} \lambda_m \sum_{n \in \mathcal{N}} R(m, n) \\ \text{s.t.} & \boldsymbol{P1} \preceq \boldsymbol{p^{max}}, \quad \boldsymbol{R} \in \mathcal{C}(\boldsymbol{P}) \end{cases}. \qquad (6)$$

The optimization framework provides a layered approach to the sum utility maximization problem. The interaction between the layers is controlled through the use of the dual variable $\boldsymbol{\lambda}$ as a set of weighting factors, which centrally coordinates the application-layer *demand* and physical-layer *supply* of rates. The dual variable $\boldsymbol{\lambda}$ has a pricing interpretation: $\lambda_m$ ($m \in \mathcal{M}$) represents dollars per unit of bit rate. As a supplier of data rates, the physical layer attempts to maximize the total revenues by adaptively allocating power and bandwidth and selecting the best choice of relay and relaying scheme in each tone. A higher value of $\lambda_m$ induces the physical layer to allocate more resources to the corresponding data stream. As a consumer of data rates, the application layer aims to maximize the sum utility discounted by the total cost of data rates. A higher value of $\lambda_m$ indicates that the cost of rate for the corresponding data stream is high, thus inducing the application layer to lower traffic demand. Finally, because the sum utility maximization problem (2) has zero duality gap, it can be solved by minimizing the dual objective:

$$\text{minimize} \quad g(\boldsymbol{\lambda}) \tag{7}$$
$$\text{subject to} \quad \boldsymbol{\lambda} \succeq \mathbf{0}$$

One way to solve this dual problem is to update $\boldsymbol{\lambda}$ using a subgradient method [40] as follows:

*Subroutine 1:* Subgradient-based method for solving (7)

1) Initialize $\boldsymbol{\lambda}^{(0)}$.
2) Given $\boldsymbol{\lambda}^{(l)}$, solve the two subproblems (5) and (6) separately to obtain the optimal values $\boldsymbol{t}^*$, $\boldsymbol{P}^*$, and $\boldsymbol{R}^*$.
3) Perform a subgradient update for $\boldsymbol{\lambda}$:

$$\boldsymbol{\lambda}^{(l+1)} = \left[ \boldsymbol{\lambda}^{(l)} + (\boldsymbol{\nu}^{(l)})^T (\boldsymbol{t}^* - \boldsymbol{R}^* \mathbf{1}) \right]^+ \tag{8}$$

4) Return to step 2 until convergence. $\square$

The subgradient method above is guaranteed to converge to the optimal dual variable, if the step sizes $\boldsymbol{\nu}^{(l)}$ are chosen following a diminishing step size rule [41]. From the optimal dual variables, the optimal primal variables can then be found easily.

### C. Solutions of individual subproblems

This section describes efficient methods to solve each of the two individual subproblems corresponding to the application layer and physical layer respectively. In particular, the physical-layer joint relay-strategy and relay-node selection and resource allocation problem can be solved efficiently by introducing another set of pricing variables that accounts for the cost of power. This second decomposition at the physical layer, together with the first decomposition described in the previous section, solves the overall utility maximization problem globally and efficiently.

*1) Application Layer Subproblem:* Finding the optimal solution of $g_{appl}(\boldsymbol{\lambda})$ as described in (5) is relatively easy. Note that (5) can be solved by maximizing each of the summation terms separately. This means that the system searches for the optimal user traffic demand $t_m^*$ independently for each stream $m$, balancing the stream's utility with the cost of rate. Specifically, since $U_m$ is a concave function of $t_m$, $(U_m(t_m) - \lambda_m t_m)$ is also a concave function of $t_m$ ($m \in$
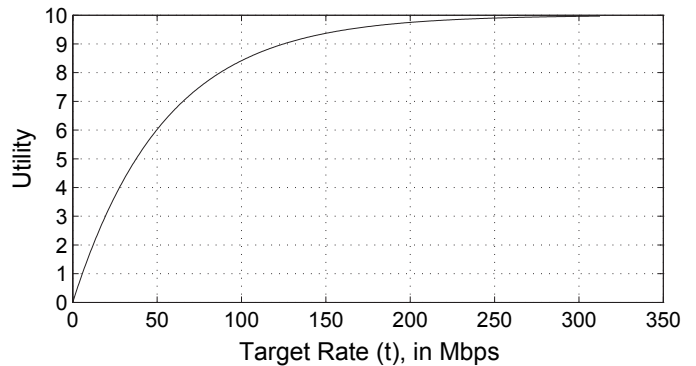


Fig. 3. An example of a utility function $U(t) = 10(1 - e^{-1.8421 \times 10^{-8} t})$.

$\mathcal{M}$). Therefore, $t_m^*$ can be found by taking the derivative of $(U_m(t_m) - \lambda_m t_m)$ with respect to $t_m$ and setting it to zero.

Many different choices of utility functions are possible depending on the application. Example 1 shows a class of utility function which will be used in simulation later in Section IV.

*Example 1:* Let $t$ be the data rate (in Mbps), and define the utility function

$$U(t) = \begin{cases} a\left(1 - e^{-bt}\right), & \text{if } t \geq 0 \\ -\infty, & \text{if } t < 0 \end{cases}, \tag{9}$$

where $a$ and $b$ are strictly positive real numbers. Fig. 3 shows an example where $a = 10$, and $b = 1.8421 \times 10^{-8}$. The variable $a$ represents the upper limit of the utility function. The value of $b$ is chosen so that at some target rate $c$, the utility obtained is equal to $0.9a$. For example, in Fig. 3, $c = 125$Mbps. Given $c$, $b = \frac{\ln(0.1)}{-c}$.

In the application layer subproblem, the per-stream maximization is of the form $(U(t) - \lambda t)$, where $\lambda$ is a constant. By calculus,

$$t^* = \max\left( 0, -\frac{1}{b} \ln \frac{\lambda}{ab} \right), \tag{10}$$

where $t^*$ is the value of $t$ which maximizes $(U(t) - \lambda t)$. $\square$

*2) Physical Layer Subproblem:* The physical-layer subproblem (6) is the more difficult of the two. Finding the optimal $\boldsymbol{R}$ involves selecting the best data stream, power allocation, relay node and relaying strategy in each tone. However, the per-node power constraint implies coupling across tones. This section introduces a second decomposition step that relaxes the power constraint. This removes the coupling across tones, resulting in a procedure that is *linear* in the number of tones. The main technique here is reminiscent of the weighted sum-rate maximization problem considered in [42].

The main step is to relax the power constraint $\boldsymbol{P1} \preceq \boldsymbol{p}^{max}$ by introducing prices into the objective function of (6):

$$\begin{aligned} Q &= \sum_{m \in \mathcal{M}} \lambda_m \sum_{n \in \mathcal{N}} R(m, n) + \boldsymbol{\mu}^T (\boldsymbol{p}^{max} - \boldsymbol{P1}) \\ &= \sum_{m \in \mathcal{M}} \lambda_m \sum_{n \in \mathcal{N}} R(m, n) + \\ &\quad \sum_{i \in \mathcal{K}_+} \mu_i \left( p_i^{max} - \sum_{n \in \mathcal{N}} P(i, n) \right), \end{aligned} \tag{11}$$

where $\boldsymbol{\mu} = [\mu_1, \mu_2, ..., \mu_{K+1}]^T$ is a vector dual variable. The key observation is that the dual function of the physical-layer subproblem

$$q(\boldsymbol{\mu}) = \begin{cases} \max\limits_{\boldsymbol{P}, \boldsymbol{R}} & Q(\boldsymbol{P}, \boldsymbol{R}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ \text{s.t.} & \boldsymbol{R} \in \mathcal{C}(\boldsymbol{P}) \end{cases} \quad (12)$$

can now be decoupled into $N$ per-tone maximization subproblems:

$$\max\limits_{P(:,n), R(:,n)} \sum_{m \in \mathcal{M}} \lambda_m R(m, n) - \sum_{i \in \mathcal{K}_+} \mu_i P(i, n) \quad (13)$$
$$\text{s.t.} \qquad R(:, n) \in \mathcal{C}(P(:, n))$$

In each tone, there is only one active stream. Further, only the source $\mathcal{S}$ and relay $\mathcal{R}$ expend power. Therefore, an alternative expression of (13) is

$$\max \quad \lambda_m R(m, n) - (\mu_{\mathcal{S}} P(\mathcal{S}, n) + \mu_{\mathcal{R}} P(\mathcal{R}, n)) \quad (14)$$
$$\text{s.t.} \quad R(:, n) \in \mathcal{C}(P(:, n)),$$

where the maximization is over $m$, choice of $\mathcal{R}$, $R(m, n)$, $P(\mathcal{S}, n)$, and $P(\mathcal{R}, n)$. The relaxation of $\boldsymbol{P1} \preceq \boldsymbol{p^{max}}$ removes the per-node power constraint. Hence, the per-tone problem (14) is easier to solve than the original physical-layer subproblem. Moreover, the physical-layer subproblem is decoupled across tones. So, the computational complexity becomes linear in the number of tones $N$. Again, using a pricing interpretation, $\mu_i$ ($i \in \mathcal{K}_+$) represents dollars per unit of power at node $i$. A higher value of $\mu_i$ discourages node $i$ from expending power, while a lower value of $\mu_i$ does the opposite. Thus, the physical-layer subproblem tries to maximize total revenue of data rates discounted by the cost of power. Together, $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ act as a set of weighting factors to centrally determine the optimal resource allocations and relay strategy at each tone.

A critical requirement for the decomposition of the physical-layer subproblems into $N$ per-tone subproblems is its convexity structure. Using the same argument as in Section II-B, the physical-layer subproblem (6) can always be made a convex function of $\boldsymbol{p^{max}}$ if time- or frequency-sharing can be implemented. Therefore, the physical-layer subproblem also has zero duality gap, and it can be solved optimally via the dual problem

$$\text{minimize} \quad q(\boldsymbol{\mu}) \quad (15)$$
$$\text{subject to} \quad \boldsymbol{\mu} \succeq \boldsymbol{0}$$

Again, a subgradient approach with appropriate step sizes $\epsilon^{(l)}$ may be used to solve the dual problem. From the optimal dual variables, the optimal primal variables can then be found relatively easily.

*Subroutine 2:* Subgradient-based method for solving (15)
1) Initialize $\boldsymbol{\mu^{(0)}}$.
2) Given $\boldsymbol{\mu^{(l)}}$, solve the $N$ per-tone maximization problems (14) separately, then combine the results to obtain $\boldsymbol{P^*}$ and $\boldsymbol{R^*}$.
3) Perform subgradient updates for $\boldsymbol{\mu}$:

$$\boldsymbol{\mu^{(l+1)}} = \left[ \boldsymbol{\mu^{(l)}} + (\boldsymbol{\epsilon^{(l)}})^T \left( \boldsymbol{P^*1} - \boldsymbol{p^{max}} \right) \right]^+ \quad (16)$$

4) Return to step 2 until convergence. □

### D. Overview of the Algorithm

The two subroutines presented in this section are interconnected hierarchically. The overall sum utility maximization problem (2) can be solved optimally in the dual domain using Subroutine 1. Step 2 of Subroutine 1 requires the solutions to both the application-layer subproblem (which is trivial) and the physical-layer subproblem (which requires Subroutine 2). Subroutines 1 and 2 use two nested subgradient update loops to search for the optimal dual variable. Because the overall optimization problem has zero duality gap, the following theorem holds:

*Theorem 1:* The algorithm summarized in the preceding paragraph always converges, and it converges to the global optimal value of the utility maximization problem (2). This is true whenever the time-sharing property holds, which is the case in the limit as the number of OFDM tones goes to infinity.

Computational experience suggests that the complexity of subgradient updates is polynomial in the dimension of the dual problem, (which is $K$ for $g(\boldsymbol{\lambda})$ and $2K$ for $q(\boldsymbol{\mu})$). Since the subroutines are connected hierarchically, the complexity of all subgradient updates of the proposed algorithm is polynomial in $K$.

## III. OPTIMAL RELAY STRATEGY SELECTION

The previous section shows that in a cooperative cellular network employing OFDMA with per-node power constraints, the physical-layer subproblem (6) can be solved globally and efficiently in the dual domain. However, this hinges upon efficient solutions to the per-tone problem (14), which is required in Step 2 of Subroutine 2. For each tone $n$ ($n \in \mathcal{N}$) the per-tone problem is a maximization of $\lambda_m R(m, n) - \mu_{\mathcal{S}} P(\mathcal{S}, n) - \mu_{\mathcal{R}} P(\mathcal{R}, n)$, where $\lambda_m$, $\mu_{\mathcal{S}}$ and $\mu_{\mathcal{R}}$ are fixed dual variables. The goal of this section is to show that this per-tone maximization problem can be solved efficiently via an exhaustive search.

The main idea is to consider $R(m, n)$ as the optimizing variable and to express $P(\mathcal{S}, n)$ and $P(\mathcal{R}, n)$ as functions of $R(m, n)$. As digital transmission is always implemented with a finite constellation, bit rate is always discrete. Consequently, the per-tone problem can be solved by an exhaustive search over a finite set defined by:
- Active data stream: $m$ ($m \in \mathcal{M}$) (which also implicitly determines $\mathcal{S}$ and $\mathcal{D}$.)
- Relay node: $\mathcal{R}$ ($\mathcal{R} \in \mathcal{K}, \mathcal{R} \neq \mathcal{S}$ or $\mathcal{D}$)
- Relaying strategy: {direct channel, decode-and-forward, amplify-and-forward}
- Bit rate: $R(m, n)$

The size of the search set is the product of the number of data streams, potential relays, relay strategies, and possible bit rates. An exhaustive search over such a discrete set is often feasible for a practical network.

However, expressing $P(\mathcal{S}, n)$ and $P(\mathcal{R}, n)$ as functions of $R(m, n)$ for each relay strategy is not entirely trivial. This is because with the participation of a relay, an entire range of power allocations at the source and at the relay is possible to achieve a fixed bit rate. The rest of this section shows that by using an extra optimization step that accounts for the pricing structure of the power availability, the optimal power

allocations at both the source and the relay can be readily found.

This paper focuses on two-time-slot implementation of DF and AF strategies. We also impose the restriction that $\mathcal{S}$ can only transmit in the first of two time-slots. This simplifies rate expressions significantly and results only in negligible performance loss as verified by simulation results.

For both relaying schemes, during the first time-slot, $\mathcal{S}$ transmits while both $\mathcal{R}$ and $\mathcal{D}$ receive. The difference between AF and DF is in the operation of $\mathcal{R}$. In AF, $\mathcal{R}$ amplifies the signal it receives in the first time-slot, and sends it out in the second time-slot. In DF, $\mathcal{R}$ attempts to decode its received signal in the first time-slot. If decoding is unsuccessful, $\mathcal{R}$ will remain silent in the second time-slot. Otherwise, $\mathcal{R}$ will re-encode the decoded data and then transmit it in the second time-slot. AF and DF relaying strategies can outperform each other depending on channel conditions and power allocations at $\mathcal{S}$ and $\mathcal{R}$. Although the achievable rate for the two-time-slot cheap relay channel as a function of transmit power have been previously derived in [1], [18], what is required here is the opposite, which is the optimal transmit power allocations as a function of data rate.

As mentioned earlier, perfect knowledge of channel gains and noise variances is assumed. Since transmission takes place in two time-slots, both actual power and actual data rate should be halved. Subscripts 1 and 2 are used to denote the first and the second time-slot respectively. Let $x_{\mathcal{S}1}$ be the symbol sent by $\mathcal{S}$, $y_{\mathcal{D}1}$ and $y_{\mathcal{D}2}$ be the received symbols at $\mathcal{D}$, and $y_{\mathcal{R}1}$ be the received symbol at $\mathcal{R}$. We assume that all nodes have one antenna each. The complex channel gains from $\mathcal{S}$ to $\mathcal{D}$, $\mathcal{S}$ to $\mathcal{R}$, and $\mathcal{R}$ to $\mathcal{D}$ are denoted by $h_{\mathcal{S}\mathcal{D}}$, $h_{\mathcal{S}\mathcal{R}}$, and $h_{\mathcal{R}\mathcal{D}}$ respectively. The channel gains are assumed to be identical in both time-slots. Moreover, $n_{\mathcal{D}1}$, $n_{\mathcal{D}2}$, and $n_{\mathcal{R}1}$ are circularly symmetric complex Gaussian noises $\mathcal{CN}(0, N_oW)^2$.

### A. Direct Channel (DC)

In DC, the source ($\mathcal{S}$) transmits directly to the destination ($\mathcal{D}$). The channel can be modelled as

$$y_{\mathcal{D}} = \sqrt{P(\mathcal{S},n)}h_{\mathcal{S}\mathcal{D}}x_{\mathcal{S}} + n_{\mathcal{D}} \qquad (17)$$

The achievable rate of DC is found using the well-known formula (in b/s/Hz):

$$R(m,n) \leq I(x_{\mathcal{S}}; y_{\mathcal{D}}) = \log_2\left(1 + \frac{P(\mathcal{S},n)|h_{\mathcal{S}\mathcal{D}}|^2}{\Gamma N_oW}\right), \quad (18)$$

where $\Gamma$ is the gap to capacity. For discrete bit-loading,

$$P^*(\mathcal{S},n) = (2^{R(m,n)} - 1)\frac{\Gamma N_oW}{|h_{\mathcal{S}\mathcal{D}}|^2}. \qquad (19)$$

### B. Decode-and-forward (DF)

The DF relay channel is shown in Fig. 4. In the first time-slot, $\mathcal{R}$ attempts to decode $x_{\mathcal{S}1}$. Assuming that decoding is
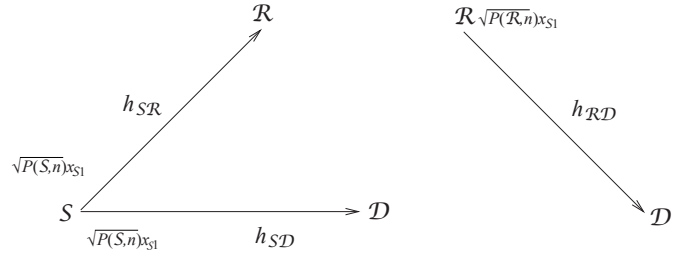
Fig. 4. The DF relay channel, where the first time-slot is shown on the left, and the second time-slot is shown on the right.

successful, $\mathcal{R}$ transmits $x_{\mathcal{S}1}$ in the second time-slot with power $P(\mathcal{R}, n)$. The channel equations are

$$y_{\mathcal{D}1} = \sqrt{P(\mathcal{S},n)}h_{\mathcal{S}\mathcal{D}}x_{\mathcal{S}1} + n_{\mathcal{D}1}, \qquad (20)$$
$$y_{\mathcal{R}1} = \sqrt{P(\mathcal{S},n)}h_{\mathcal{S}\mathcal{R}}x_{\mathcal{S}1} + n_{\mathcal{R}1}, \qquad (21)$$
$$y_{\mathcal{D}2} = \sqrt{P(\mathcal{R},n)}h_{\mathcal{R}\mathcal{D}}x_{\mathcal{S}1} + n_{\mathcal{D}2}. \qquad (22)$$

Successful decoding of $x_{\mathcal{S}1}$ at $\mathcal{R}$ requires

$$R(m,n) \leq I(x_{\mathcal{S}1}; y_{\mathcal{R}1}) = \log_2\left(1 + \frac{P(\mathcal{S},n)|h_{\mathcal{S}\mathcal{R}}|^2}{\Gamma N_oW}\right), \qquad (23)$$

or equivalently,

$$P(\mathcal{S},n) \geq (2^{R(m,n)} - 1)\frac{\Gamma N_oW}{|h_{\mathcal{S}\mathcal{R}}|^2}. \qquad (24)$$

Successful decoding at $\mathcal{D}$ requires

$$R(m,n) \leq I(x_{\mathcal{S}1}; y_{\mathcal{D}1}, y_{\mathcal{D}2})$$
$$= \log_2\left(1 + \frac{P(\mathcal{S},n)|h_{\mathcal{S}\mathcal{D}}|^2 + P(\mathcal{R},n)|h_{\mathcal{R}\mathcal{D}}|^2}{\Gamma N_oW}\right). \qquad (25)$$

Note that $\mathcal{D}$ receives two scaled and noisy versions of $x_{\mathcal{S}1}$ across two time-slots. A maximum-ratio-combining formula[3] is used to derive (25). In order for DF to work, the achievable rate of $x_{\mathcal{S}1}$ at $\mathcal{R}$ has to be higher than at $\mathcal{D}$. This means that the expression on the right hand side of (23) has to be greater than or equal to that of (25). Notice that this can happen only if $|h_{\mathcal{S}\mathcal{R}}| > |h_{\mathcal{S}\mathcal{D}}|$, which is intuitive.

Rearranging the terms of (25), it is clear that the minimum $P(\mathcal{R}, n)$ is

$$P^*(\mathcal{R},n) = \frac{\left(2^{R(m,n)} - 1\right)\Gamma N_oW - P(\mathcal{S},n)|h_{\mathcal{S}\mathcal{D}}|^2}{|h_{\mathcal{R}\mathcal{D}}|^2} \qquad (26)$$

If $P^*(\mathcal{R}, n) \leq 0$, then DF is not a suitable relaying scheme. Now, it remains to optimize $P(\mathcal{S}, n)$, which is not immediate since at a fixed rate $R(m, n)$, decreasing $P(\mathcal{S}, n)$ would increase $P^*(\mathcal{R}, n)$. Recall that the objective of the per-tone optimization problem for each tone $n$ as expressed in (14) is

$$\max_{P(\mathcal{S},n)} \lambda_m R(m,n) - \mu_{\mathcal{S}}P(\mathcal{S},n) - \mu_{\mathcal{R}}P^*(\mathcal{R},n) \qquad (27)$$
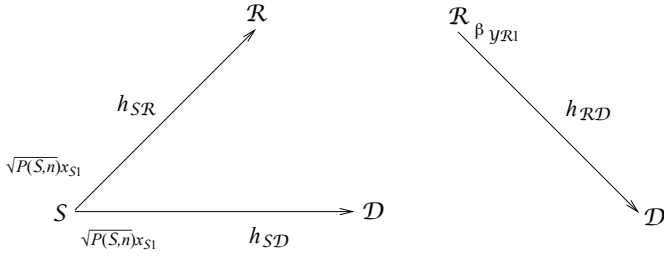
Fig. 5. The AF relay channel, where the first time-slot is shown on the left, and the second time-slot is shown on the right.

From (26), $P^*(\mathcal{R}, n)$ may be obtained as a function of $P(\mathcal{S}, n)$, so that the objective of the per-tone optimization problem (27) is a function of $P(\mathcal{S}, n)$ only. Since (27) is now unconstrained, it can be solved by taking the first-order derivative of its objective function (called $f$ below):

$$\frac{df}{dP(\mathcal{S}, n)} = -\mu_{\mathcal{S}} + \mu_{\mathcal{R}} \frac{|h_{\mathcal{SD}}|^2}{|h_{\mathcal{RD}}|^2} \tag{28}$$

Note that because $f$ is a linear function of $P(\mathcal{S}, n)$, the derivative is a constant. This implies that if $\frac{df}{dP(\mathcal{S}, n)} > 0$, then $P(\mathcal{S}, n)$ should be maximized. But, from (26), the maximizing $P(\mathcal{S}, n)$ makes $P^*(\mathcal{R}, n) = 0$. This means DF is unnecessary. On the other hand, if $\frac{df}{dP(\mathcal{S}, n)} \leq 0$, then $P^*(\mathcal{S}, n)$ should be minimized. This means that the expression for $P(\mathcal{S}, n)$, i.e. (24), should be satisfied with equality. Note that by (26), such a $P^*(\mathcal{S}, n)$ guarantees that $P^*(\mathcal{R}, n)$ is positive. The DF mode power allocation procedure is summarized below:

*Subroutine 3:* Optimal power allocation for a fixed $R(m, n)$ in the DF relay mode:

1) If $|h_{\mathcal{SR}}| <= |h_{\mathcal{SD}}|$, set $P^*(\mathcal{S}, n) = P^*(\mathcal{R}, n) = \infty$,
2) else if $-\mu_{\mathcal{S}} + \mu_{\mathcal{R}} \frac{|h_{\mathcal{SD}}|^2}{|h_{\mathcal{RD}}|^2} > 0$, then set $P^*(\mathcal{S}, n) = P^*(\mathcal{R}, n) = \infty$,
3) else set $P^*(\mathcal{S}, n)$ and $P^*(\mathcal{R}, n)$ according to (24) and (26), respectively, with equality.
4) Divide $R(m, n)$, $P^*(\mathcal{S}, n)$, and $P^*(\mathcal{R}, n)$ by 2. $\square$

As mentioned earlier, the last step is necessary since communication takes place in two time-slots.

### C. Amplify-and-forward (AF)

The AF relay channel is shown in Fig. 5. The relay $\mathcal{R}$ receives $x_{\mathcal{S}1}$ in the first time-slot and transmits an amplified version of $x_{\mathcal{S}1}$ in the second time-slot with power $P(\mathcal{R}, n)$. The channel equations are

$$y_{\mathcal{D}1} = \sqrt{P(\mathcal{S}, n)} h_{\mathcal{SD}} x_{\mathcal{S}1} + n_{\mathcal{D}1}, \tag{29}$$

$$y_{\mathcal{R}1} = \sqrt{P(\mathcal{S}, n)} h_{\mathcal{SR}} x_{\mathcal{S}1} + n_{\mathcal{R}1}, \tag{30}$$

$$y_{\mathcal{D}2} = \beta y_{\mathcal{R}1} h_{\mathcal{RD}} + n_{\mathcal{D}2}, \tag{31}$$

where

$$\beta = \sqrt{\frac{P(\mathcal{R}, n)}{P(\mathcal{S}, n)|h_{\mathcal{SR}}|^2 + N_oW}} \tag{32}$$

is the power amplification factor at $\mathcal{R}$. The AF scheme is a suitable choice when the relay does not have a sufficiently large SNR to decode the transmitted symbol. However, the AF scheme suffers from noise amplification.

To analyze the power requirement for AF, recognize that in order for the destination $\mathcal{D}$ to decode the signal $x_{\mathcal{S}1}$, which is sent across two time-slots, the following must hold:

$$R(m, n) \leq I(x_{\mathcal{S}1}; y_{\mathcal{D}1}, y_{\mathcal{D}2})$$
$$= \log_2 \left( 1 + \frac{1}{\Gamma} \left[ \frac{P(\mathcal{S}, n)|h_{\mathcal{SD}}|^2}{N_oW} + \frac{\frac{P(\mathcal{R}, n)P(\mathcal{S}, n)|h_{\mathcal{RD}}|^2|h_{\mathcal{SR}}|^2}{P(\mathcal{S}, n)|h_{\mathcal{SR}}|^2 + N_oW}}{N_oW \left( 1 + \frac{P(\mathcal{R}, n)|h_{\mathcal{RD}}|^2}{P(\mathcal{S}, n)|h_{\mathcal{SR}}|^2 + N_oW} \right)} \right] \right) \tag{33}$$

The variables of (33) are $P(\mathcal{S}, n)$ and $P(\mathcal{R}, n)$. Rearranging the terms gives

$$P^*(\mathcal{R}, n) = \frac{(c_1 P(\mathcal{S}, n) + c_2)(c_3 P(\mathcal{S}, n) + c_4)}{c_5 P(\mathcal{S}, n) + c_6}, \tag{34}$$

where

$$c_1 = |h_{\mathcal{SD}}|^2, \quad c_2 = -(2^{R(m,n)} - 1)\Gamma N_oW,$$
$$c_3 = |h_{\mathcal{SR}}|^2, \quad c_4 = N_oW,$$
$$c_5 = |h_{\mathcal{RD}}|^2(-|h_{\mathcal{SD}}|^2 - |h_{\mathcal{SR}}|^2),$$
$$c_6 = (2^{R(m,n)} - 1)\Gamma N_oW|h_{\mathcal{RD}}|^2.$$

Now, observe that $(c_3 P(\mathcal{S}, n) + c_4) > 0$. Thus, to ensure $P^*(\mathcal{R}, n) > 0$, the terms $(c_1 P(\mathcal{S}, n) + c_2)$ and $(c_5 P(\mathcal{S}, n) + c_6)$ must either be both greater than zero or both less than zero. It is not hard to see that a valid solution is obtained only when both terms are less than zero, leading to a feasible region for $P(\mathcal{S}, n)$ as

$$P_{\min}(\mathcal{S}, n) < P(\mathcal{S}, n) < P_{\max}(\mathcal{S}, n), \tag{35}$$

where

$$P_{\min}(\mathcal{S}, n) = \frac{(2^{R(m,n)} - 1)\Gamma N_oW}{|h_{\mathcal{SD}}|^2 + |h_{\mathcal{SR}}|^2},$$
$$P_{\max}(\mathcal{S}, n) = \frac{(2^{R(m,n)} - 1)\Gamma N_oW}{|h_{\mathcal{SD}}|^2}.$$

Notice that (35) contains strict inequalities. This is because if $P(\mathcal{S}, n)$ is equal to $P_{\min}(\mathcal{S}, n)$, then $P^*(\mathcal{R}, n)$ will be equal to infinity. On the other hand, $P(\mathcal{S}, n)$ cannot be equal to $P_{\max}(\mathcal{S}, n)$ either, because this implies $P^*(\mathcal{R}, n) = 0$, making AF relaying unnecessary.

Now, it remains to choose the optimal power allocations for a fixed $R(m, n)$. Similar to the analysis in the DF mode, the per-tone objective is as expressed in (14):

$$\max_{P(\mathcal{S}, n)} \lambda_m R(m, n) - \mu_{\mathcal{S}} P(\mathcal{S}, n) - \mu_{\mathcal{R}} P^*(\mathcal{R}, n). \tag{36}$$

Denote the objective function above $f$. The first step is to show that $f$ is a concave function of $P(\mathcal{S}, n)$. Compute

$$\frac{df}{dP(\mathcal{S}, n)} = -\mu_{\mathcal{S}} - \mu_{\mathcal{R}} \frac{dP^*(\mathcal{R}, n)}{dP(\mathcal{S}, n)}, \tag{37}$$

$$\frac{d^2 f}{dP(\mathcal{S}, n)^2} = -\mu_{\mathcal{R}} \frac{d^2 P^*(\mathcal{R}, n)}{dP(\mathcal{S}, n)^2} \tag{38}$$

It can be verified by algebra that for $P(\mathcal{S}, n)$ within the feasible region (35), $\frac{dP^*(\mathcal{R}, n)}{dP(\mathcal{S}, n)} < 0$ and $\frac{d^2 P^*(\mathcal{R}, n)}{dP(\mathcal{S}, n)^2} > 0$. Substituting the results into (37) and (38) shows that within the feasible region of $P(\mathcal{S}, n)$, $\frac{d^2 f}{dP(\mathcal{S}, n)^2} < 0$, and thereby proving the concavity of $f$.

The concavity of $f$ and the observation that $f$ is continuous ensure that there is a unique optimal value of $f$ within the
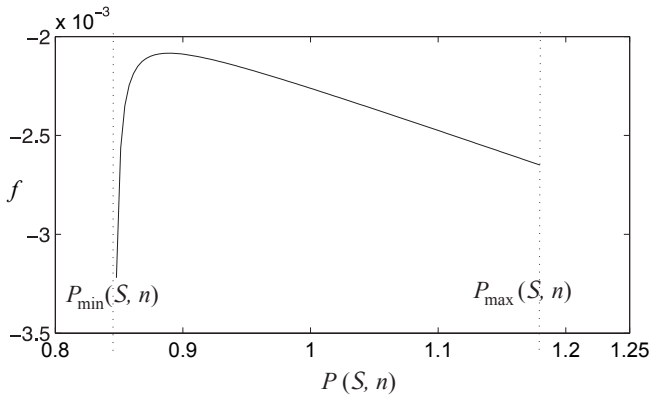
Fig. 6. An example of $f$ as a function of $P(\mathcal{S}, n)$ for the AF relay channel. Notice that there is a unique maximum value of $f$.

feasible region of $P(\mathcal{S}, n)$. The optimal value can be found by solving for the root of $\frac{df}{dP(\mathcal{S},n)}$. This can be done with a root-finding method such as the Newton's method. An example of $f$ as a function of $P(\mathcal{S}, n)$ is shown in Fig. 6. The AF mode power allocation procedure is summarized below:

*Subroutine 4:* Optimal power allocation for a fixed $R(m, n)$ in the AF relay mode:

1) Solve the equation $\frac{df}{dP(\mathcal{S},n)} = 0$ to obtain $P^*(\mathcal{S}, n)$, using either bisection or Newton's method, within the feasible region of $P(\mathcal{S}, n)$ (35).
2) Set $P^*(\mathcal{R}, n)$ according to (34).
3) Divide $R(m, n)$, $P^*(\mathcal{S}, n)$, and $P^*(\mathcal{R}, n)$ by 2. □

### D. Computational Complexity of Algorithm

All four subroutines presented in this paper are connected hierarchically to solve the overall utility maximization problem (2). The starting point of the proposed algorithm is Subroutine 1, which solves (2) optimally in the dual domain. Subroutine 1 requires the solutions to both the trivial application-layer subproblem and the more difficult physical-layer subproblem. Subroutine 2 solves the physical-layer subproblem in the dual domain, and requires solutions to the per-tone problem (14). The per-tone problem can be solved efficiently by searching over a discrete set. The discrete search requires the expressions of power at the source and the relay as a function of the bit rate for different relaying schemes. Subroutines 3 and 4 describe the associated procedures for DF and AF relaying respectively.

The subroutines together enable the cross-layer optimization problem to be solved globally and efficiently. The computational complexity of subgradient updates in Subroutines 1 and 2 together is polynomial in $K$, as discussed in Section II-D. The complexity of the per-tone problem is analyzed as follows. Let $A$ be the number of relaying schemes. In this paper, the relaying schemes considered are DC, DF, and AF, so $A = 3$. Let $B$ be the number of possible bit rate for each physical transmission, which is usually small in a practical system employing finite-sized constellations. Moreover, there are $2K$ data streams, $K - 1$ potential relay nodes for each frequency tone, and $N$ frequency tones. Thus, the complexity of solving the $N$ per-tone maximization problems is $O(ABK^2N)$. The quadratic dependency with $K$ is due to the fact that there are

$O(K)$ candidate relay nodes for each of the $2K$ data streams. In a practical system, if geographical information is available, the search over the suitable relay nodes can be easily reduced to $O(1)$ for each stream by restricting attentions to a fixed number of relays located roughly between the source and the destination. This reduces the overall complexity to be linear in $K$.

Note that the complexity of the proposed algorithm is linear in $N$, which is attractive in an OFDMA implementation. Although Theorem 1 requires the number of OFDM tones to go to infinity, this is only necessary to ensure the frequency-sharing property, which guarantees zero duality gap. In actual simulations with a finite number of tones $N$, we observe that the duality gap is very close to zero most of the time.

## IV. SIMULATIONS

This section presents simulation results for the proposed algorithm summarized in Section III-D, which solves the utility maximization problem that uses physical-layer relaying (2). We simulate two networks with 2 and 4 user nodes respectively. In both networks, the total system bandwidth is set to be 80MHz, and the number of OFDM tones $N = 256$. This corresponds to a tone width $W = 312.5$kHz. The channel gain between two nodes at each tone can be decomposed into a small-scale Rayleigh fading component and a large-scale path loss component with path loss exponent of 4. To simulate a frequency-selective environment, small-scale fading is i.i.d. across tones. Moreover, the gap to capacity $\Gamma$ is set to 1, which corresponds to perfect coding. The type of utility function chosen is as described in Example 1 with parameters $a = 10$ and $c = 125$Mbps for downstream communications, and $a = 1$ and $c = 12.5$Mbps for upstream communications. These parameters are chosen to reflect a preference for downstream communications in most applications.

### A. Network with 2 User Nodes

The first example is a wireless network with a base station and $K = 2$ user nodes. The base station (node $K + 1$) is fixed at $(0,0)$ and node 2 is fixed at $(10,0)$ in a two-dimensional plane. The location of node 1 changes as this simulation proceeds, but node 1 is always restricted to be closer to the base station. This is without loss of generality because for situations where node 1 is further away, the analysis would apply by simply switching the roles of node 1 and node 2. The power constraints of all nodes are the same such that $\frac{p_i^{max}}{N_oW} = 23$dB ($i \in \mathcal{K}_+$). This corresponds to a medium-SNR environment. There are a total of 4 data streams because there are 2 user nodes, and each of them can have downstream and upstream transmissions.

*1) Fix Node 1 at $(5, 0)$:* As a first example, node 1 is fixed at $(5, 0)$. Fig. 7 shows the location of nodes and describes whether a link can be used for direct transmissions, relaying transmissions, or both. The upper limit of system utility for this network is 22, which is calculated by summing the value of the parameter $a$ for all data streams. Using optimization methods presented in the previous section, it is found that by allowing relaying, the maximized sum utility increases from 17.11 to 18.79, quantifying the merit of relaying. Keep in
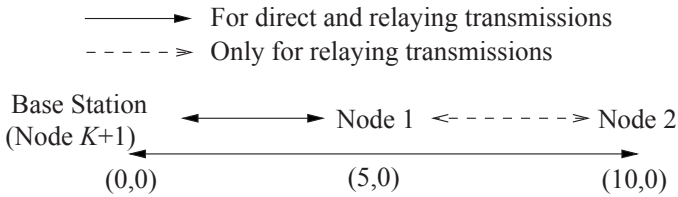
------→ For direct and relaying transmissions

- - - - -≻ Only for relaying transmissions

Base Station      ←-----→ Node 1 ←- - - - -≻ Node 2
(Node $K+1$)

(0,0)            (5,0)          (10,0)

Fig. 7. Topology of a cooperative network with 2 user nodes.

TABLE I

RATES OF VARIOUS DATA STREAMS IN THE 2-USER NETWORK

| Stream | No Relay | Allow Relay | Percentage Change |
|---|---|---|---|
| $(K+1,1)$ | 130.0Mbps | 115.9Mbps | $-10.8\%$ |
| $(K+1,2)$ | 50.8Mbps | 88.8Mbps | $74.8\%$ |
| $(1,K+1)$ | 27.0Mbps | 19.4Mbps | $-28.2\%$ |
| $(2,K+1)$ | 16.2Mbps | 15.8Mbps | $-2.5\%$ |

TABLE II

PERCENTAGE OF POWER EXPENDED AS A RELAY IN THE 2-USER NETWORK

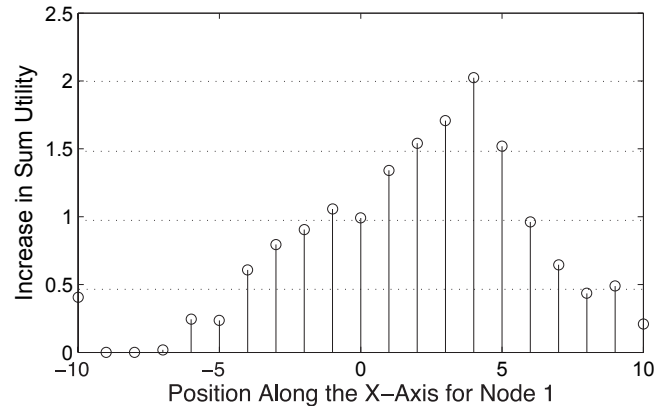| Node | Percentage of power expended as a relay |
|---|---|
| 1 | 47.6% |
| 2 | 0% |



Fig. 8. The benefit of relaying as a function of relay location. The increase in sum utility is plotted against various positions of node 1 in the 2-user network.

mind that the choice of utility function is highly dependent on the type of applications. Although the increase in utility in this example may not look impressive, it translates into significant increase in data rate of stream $(K+1,2)$, as shown in Table I. This result is logical because node 2 is farther away from the base station, and the choice of utility function favors downstream over upstream communications.

Each user node expends part of its total power to transmit its own upstream data, and the rest of its power to act as a relay for other nodes' transmissions. The optimization algorithm proposed in the previous two sections allows each user node to find the optimal power and bandwidth division between the two roles. However, each node has a finite power constraint, and the network has a finite bandwidth. Therefore, the increase in data rate of one stream has to come from decreases in data rates of other streams. As shown in Table I, the data rate of stream $(1,K+1)$ decreases by $28.2\%$ when relaying is allowed. This suggests that in order to maximize sum utility, node 1 has to sacrifice its own upstream data rate in return for the higher downstream data rate of node 2. In fact, node 1 expends $47.6\%$ of its power in the relay mode (see Table II).

*2) Effectiveness of Relaying vs User Location:* To illustrate how the location of user nodes affects the effectiveness of relaying, the position of node 1 is varied along the x-axis. The base station (node $K+1$) is kept at $(0,0)$, while node 2 is kept at $(10,0)$. Fig. 8 shows the increase in sum utility when relaying is allowed. Observe that the increase in performance is most significant when node 1 is at $(4,0)$, and gradually decreases as node 1 moves away along the x-axis in both directions.

*3) Relaying Scheme vs User Location:* To illustrate that the proposed optimization framework selects the optimal relaying scheme, Fig. 9 shows the dominant relaying scheme(s) for stream $(K+1,2)$ for different positions of node 1. The results of Fig. 9 nicely follow the rule-of-thumb that when $\mathcal{R}$ is close to $\mathcal{S}$ such that $\mathcal{R}$ can decode the received data, DF is the preferred relaying scheme. As $\mathcal{R}$ moves further away from $\mathcal{S}$, AF relaying scheme is preferred. However, the simple rule-of-thumb, without taking into consideration other factors, does not indicate where the transitions from DF to AF occur,

whereas Fig. 9 shows exactly where the transitions are. This is possible because the presented algorithm jointly optimizes relaying strategy, power allocation, and user traffic patterns.

*B. Network with 4 user nodes*

The second set of simulation results illustrates the optimized performance of a network with a base station and $K=4$ user nodes. Fig. 10 shows the location of different nodes and describes the allowable transmissions in each link. Note that in this example, both nodes 1 and 2 can potentially help the transmissions of nodes 3 and 4 by acting as relays. The proposed optimization algorithm selects the best relay in accordance with channel realizations and the availability of power and bandwidth.

The power constraints of all nodes are such that $\frac{p_i^{max}}{N_oW} = 20$dB ($i \in \mathcal{K}_+$). This corresponds to a medium-SNR environment. Using the proposed optimization algorithm, it is found that by allowing relaying, the maximized sum utility increases from $34.38$ to $37.20$, which again quantifies the merit of cooperative relaying. Note that each source-destination pair selects both the best relay and the best relay strategy in each frequency tone. This is done in a globally optimal way.

Table III shows the rates of various data streams with and without relaying. Again, as downstream is preferred by the virtue of the choice of the utility function, both streams $(K+1,3)$ and $(K+1,4)$ benefit tremendously from relaying. In fact, in the optimal solution, only streams $(K+1,3)$ and $(K+1,4)$ use relays. For both of these streams, node 1 or node 2 act as the relay, depending on the tone. This result illustrates that there is negligible performance loss (in this case, there is no loss at all) if we restrict attentions to relays located roughly between the source and the destination. In Table IV, we see that nodes 1 and 2 expend over $90\%$ of its power to act as
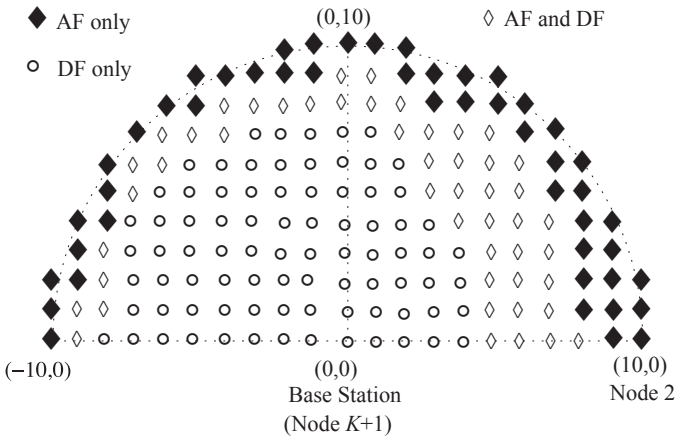
Fig. 9. The dominant relay mode (AF or DF) for stream $(K + 1, 2)$ for various positions of node 1 in the 2-user network. "AF and DF" indicates that AF and DF are both likely to occur, depending on the channel realizations in each tone. The DC mode is not shown here, but it occurs frequently for all positions of node 1, especially when node 1 is far away from the base station and/or node 2.
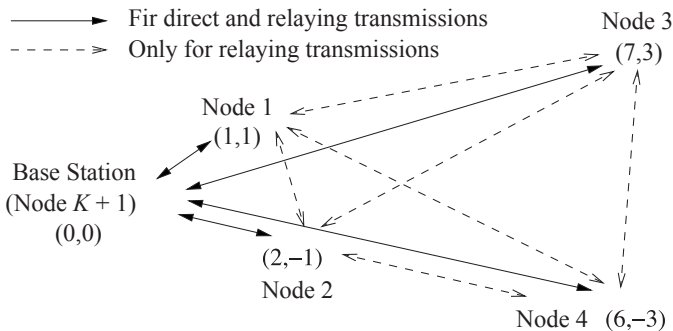


Fig. 10. Topology of a cooperative network with 4 user nodes.

a relay. These results illustrate that in a system with optimal allocations of bandwidth and power that truly maximizes the sum utility, nodes 1 and 2 would sacrifice their own data rates for the benefits of nodes 3 and 4.

## V. CONCLUSION

This paper proposes a utility maximization framework for the joint optimization of the best relay node, the best relay strategy, and the best power, bandwidth and rate allocations in a cellular network. Both amplify-and-forward and decode-and-forward relaying strategies are considered. By using a dual optimization technique for OFDMA systems, the seemingly difficult system optimization problem is solved efficiently and globally using a pricing structure. Specifically, the pricing structure decomposes the optimization problem into application-layer and physical-layer subproblems. The physical-layer subproblem can be further decomposed into per-tone optimization problems, which makes the overall computational complexity linear in the number of tones. This paper illustrates that by adopting a cross-layer approach that takes into account both the power and bandwidth availability and the traffic demand of each user, cooperative relaying in the physical layer has the potential to significantly improve the overall system performance.

TABLE III
RATES OF VARIOUS DATA STREAMS IN THE 4-USER NETWORK

| Stream | No Relay | Allow Relay | Percentage Change |
|---|---|---|---|
| $(K + 1, 1)$ | 152.8Mbps | 148.4Mbps | $-2.9\%$ |
| $(K + 1, 2)$ | 135Mbps | 129.4Mbps | $-4.1\%$ |
| $(K + 1, 3)$ | 46.6Mbps | 71.3Mbps | $53.0\%$ |
| $(K + 1, 4)$ | 54.1Mbps | 80.5Mbps | $48.8\%$ |
| $(1, K + 1)$ | 18.8Mbps | 16.6Mbps | $-11.7\%$ |
| $(2, K + 1)$ | 18.8Mbps | 16.3Mbps | $-13.3\%$ |
| $(3, K + 1)$ | 11.9Mbps | 13.8Mbps | $16.0\%$ |
| $(4, K + 1)$ | 14.1Mbps | 13.4Mbps | $-5.0\%$ |

TABLE IV
PERCENTAGE OF POWER EXPENDED AS A RELAY IN THE 4-USER
NETWORK

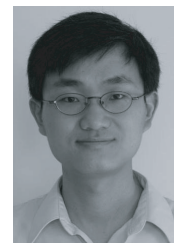| Node | Percentage of power expended as a relay |
|---|---|
| 1 | 94.9% |
| 2 | 92.2% |
| 3 | 0% |
| 4 | 0% |

## REFERENCES

[1] R. U. Nabar, H. Bölcskei, and F. W. Kneubühler, "Fading relay channels: performance limits and space-time signal design," *IEEE J. Sel. Areas Commun.*, vol. 22, no. 6, pp. 1099–1109, Aug. 2004.

[2] A. Nosratinia, T. E. Hunter, and A. Hedayat, "Cooperative communication in wireless networks," *IEEE Commun. Mag.*, vol. 42, no. 10, pp. 74–80, Oct. 2004.

[3] J. N. Laneman and G. W. Wornell, "Distributed space-time coded protocols for exploiting cooperative diversity in wireless networks," *IEEE Trans. Inf. Theory*, vol. 49, no. 10, pp. 2415–2425, Oct. 2003.

[4] J. N. Laneman, D. N. C. Tse, and G.W.Wornell, "Cooperative diversity in wireless networks: efficient protocols and outage behavior," *IEEE Trans. Inf. Theory*, vol. 50, no. 12, pp. 3062–3080, Dec. 2004.

[5] A. Sendonaris, E. Erkip, and B. Aazhang, "User cooperation diversity–part I: system description," *IEEE Trans. Commun.*, vol. 51, no. 11, pp. 1927–1938, Nov. 2003.

[6] ——, "User cooperation diversity–part II: implementation aspects and performance analysis," *IEEE Trans. Commun.*, vol. 51, no. 11, pp. 1939–1948, Nov. 2003.

[7] A. Stefanov and E. Erkip, "Cooperative coding for wireless networks," *IEEE Trans. Commun.*, vol. 52, no. 9, pp. 1470–1476, Sept. 2004.

[8] R. Pabst, B. Walke, D. Schultz, P. Herhold, H. Yanikomeroglu, S. Mukherjee, H. Viswanathan, M. Lott, W. Zirwas, M. Dohler, H. Aghvami, D. Falconer, and G. Fettweis, "Relay-based deployment concepts for wireless and mobile broadband radio," *IEEE Commun. Mag.*, vol. 42, no. 9, pp. 80–89, Sept. 2004.

[9] H. Viswanathan and S. Mukherjee, "Performance of cellular networks with relays and centralized scheduling," *IEEE Trans. Wireless Commun.*, vol. 4, no. 5, pp. 2318–2328, Sept. 2005.

[10] J. Cho and Z. J. Haas, "On the throughput enhancement of the downstream channel in cellular radio networks through multihop relaying," *IEEE J. Sel. Areas Commun.*, vol. 22, no. 7, pp. 1206–1219, Sept. 2004.

[11] H.-Y. Hsieh and R. Sivakumar, "On using peer-to-peer communication in cellular wireless data networks," *IEEE Trans. Mobile Comput.*, vol. 3, no. 1, pp. 57–72, Jan. 2004.

[12] T. M. Cover and A. A. E. Gamal, "Capacity theorems for the relay channel," *IEEE Trans. Inf. Theory*, vol. 25, no. 5, pp. 572–584, Sept. 1979.

[13] V. Sreng, H. Yanikomeroglu, and D. Falconer, "Relay selection strategies in cellular networks with peer-to-peer relaying," in *Proc. IEEE Veh. Technol. Conf. (VTC-Fall)*, pp. 1949–1953.

[14] M. Yu and J. Li, "Is amplify-and-forward practically better than decode-and-forward or vice versa?" in *Proc. IEEE Inter. Conf. Acoustics, Speech, and Signal Processing, (ICASSP)*, vol. 3, Mar. 2005, pp. 365–368.

[15] B. Can, H. Yomo, and E. D. Carvalho, "A novel hybrid forwarding scheme for OFDM based cooperative relay networks," in *Proc. IEEE Inter. Conf. Commun. (ICC) 2006.*

[16] M. A. Khojastepour, A. Sabharwal, and B. Aazhang, "On the capacity of 'cheap' relay networks," in *Proc. 37th Annu. Conf. Information Sciences and Systems (CISS) 2003*, pp. 12–14.

[17] ——, "On capacity of Gaussian 'cheap' relay channel," in *Proc. IEEE Global Telecommun. Conf. (GLOBECOM) 2003*, pp. 1776–1780.

[18] A. Wittneben and B. Rankov, "Impact of cooperative relays on the capacity of rank-deficient MIMO channels," in *Proc. 12th IST Summit on Mobile Wireless Communications 2003*, pp. 421–425.

[19] M. Dohler, A. Gkelias, and H. Aghvami, "A resource allocation strategy for distributed MIMO multi-hop communication systems," *IEEE Commun. Lett.*, vol. 8, no. 2, pp. 99–101, Feb. 2004.

[20] J. Luo, R. S. Blum, L. Cimini, L. Greenstein, and A. Haimovich, "Power allocation in a transmit diversity system with mean channel gain information," *IEEE Commun. Lett.*, vol. 9, no. 7, pp. 616–618, July 2005.

[21] A. Host-Madsen and J. Zhang, "Capacity bounds and power allocation for wireless relay channels," *IEEE Trans. Inf. Theory*, vol. 51, no. 6, pp. 2020–2040, June 2005.

[22] M. O. Hasna and M.-S. Alouini, "Optimal power allocation for relayed transmissions over Rayleigh-fading channels," *IEEE Trans. Wireless Commun.*, vol. 3, no. 6, pp. 1999–2004, Nov. 2004.

[23] I. Maric and R. D. Yates, "Bandwidth and power allocation for cooperative strategies in Gaussian relay networks," in *Proc. Thirty-Eighth Asilomar Conference on Signals, Systems and Computers*, vol. 2, pp. 1907–1911.

[24] J. Adeane, M. R. D. Rodrigues, and I. J. Wassell, "Centralised and distributed power allocation algorithms in cooperative networks," in *Proc. IEEE 6th Workshop on Signal Processing Advances in Wireless Communications (SPAWC) 2005*, pp. 333–337.

[25] X. Deng and A. M. Haimovich, "Power allocation for cooperative relaying in wireless networks," *IEEE Commun. Lett.*, vol. 9, no. 11, pp. 994–996, Nov. 2005.

[26] A. Reznik, S. R. Kulkarni, and S. Verdú, "Degraded Gaussian multirelay channel: capacity and optimal power allocation," *IEEE Trans. Inf. Theory*, vol. 50, no. 12, pp. 3037–3046, Dec. 2004.

[27] M. Chen, S. Serbetli, and A. Yener, "Distributed power allocation for parallel relay networks," in *Proc. IEEE Global Telecommun. Conf. (GLOBECOM) 2005*, vol. 3, Nov. 2005, pp. 1177–1181.

[28] E. G. Larsson and Y. Cao, "Collaborative transmit diversity with adaptive radio resource and power allocation," *IEEE Commun. Lett.*, vol. 9, no. 6, pp. 511–513, June 2005.

[29] S. Serbetli and A. Yener, "Optimal power allocation for relay assisted F/TDMA ad hoc networks," in *Proc. Inter. Conf. Wireless Networks, Commun. and Mobile Comp.*, vol. 2, pp. 1319–1324.

[30] Y. Liang and V. V. Veeravalli, "Resource allocation for wireless relay channels," in *Proc. 38th Asilomar Conference on Signals, Systems and Computers*, vol. 2, pp. 1902–1906.

[31] ——, "Gaussian orthogonal relay channels: optimal resource allocation and capacity," *IEEE Trans. Inf. Theory*, vol. 51, no. 9, pp. 3284–3289, Sept. 2005.

[32] G.-D. Yu, Z.-Y. Zhang, Y. Chen, S. Chen, and P.-L. Qiu, "Power allocation for non-regenerative OFDM relaying channels," in *Proc. Inter. Conf. Wireless Commun., Networking and Mobile Comp. (WCNC) 2005*, vol. 1, pp. 185–188.

[33] K. Lee and A. Yener, "On the achievable rate of three-node cognitive hybrid wireless networks," in *Proc. International Conference on Wireless*

[34] K. Chen, Z. Yang, C. Wagener, and K. Nahrstedt, "Market models and pricing mechanisms in a multihop wireless hotspot network," in *Proc. Second Annual Inter. Conf. on Mobile and Ubiquitous Systems: Networking and Services (MobiQuitous) 2005*, pp. 73–82.

[35] O. Ileri, S.-C. Mau, and N. Mandayam, "Pricing for enabling forwarding in self-configuring ad hoc networks," *IEEE J. Sel. Areas Commun.*, vol. 23, no. 1, pp. 151–162, Jan. 2005.

[36] F. P. Kelly, A. Maulloo, and D. Tan, "Rate control for communication networks: shadow prices, proportional fairness and stability," *J. Operations Research Society*, vol. 49, no. 3, pp. 237–252, 1998.

[37] G. Song and Y. G. Li, "Cross-layer optimization for OFDM wireless networks–part I: theoretical framework," *IEEE Trans. Wireless Commun.*, vol. 4, no. 2, pp. 614–624, Mar. 2005.

[38] ——, "Cross-layer optimization for OFDM wireless networks–part II: algorithm development," *IEEE Trans. Wireless Commun.*, vol. 4, no. 2, pp. 625–634, Mar. 2005.

[39] W. Yu and R. Lui, "Dual methods for nonconvex spectrum optimization of multicarrier systems," *IEEE Trans. Commun.*, vol. 54, no. 7, pp. 1310–1322, July 2006.

[40] N. Z. Shor, *Minimization Methods for Non-differentiable Functions.* Springer, 1985.

[41] S. Boyd, L. Xiao, and A. Mutapcic, "Subgradient methods," lecture notes of EE392o, Stanford University, Autumn Quarter 2003-2004.

[42] R. Cendrillon, W. Yu, M. Moonen, J. Verlinden, and T. Bostoen, "Optimal spectrum balancing for digital subscriber lines," *IEEE Trans. Commun.*, vol. 54, no. 5, pp. 922–933, May 2006.

**Truman Chiu-Yam Ng** received the B.A.Sc. degree in computer engineering from the University of Waterloo, Waterloo, ON, Canada, in 2004, and the M.A.Sc. degree in electrical engineering from the University of Toronto, Toronto, ON, Canada, in 2006. His areas of research are wireless relay channels, information theory, and convex optimization techniques.

**Wei Yu** (S'97-M'02) received the B.A.Sc. degree in Computer Engineering and Mathematics from the University of Waterloo, Waterloo, Ontario, Canada in 1997 and M.S. and Ph.D. degrees in Electrical Engineering from Stanford University, Stanford, CA, in 1998 and 2002, respectively. Since 2002, he has been an Assistant Professor with the Electrical and Computer Engineering Department at the University of Toronto, Toronto, Ontario, Canada, where he holds a Canada Research Chair. His main research interests are multiuser information theory, optimization, wireless communications and broadband access networks.

Prof. Yu is currently an Editor of the *IEEE Transactions on Wireless Communications*. He was a Guest Editor of the *IEEE Journal on Selected Areas in Communications* (Special Issue on Nonlinear Optimization on Communication Systems) in 2006, and a Guest Editor of *EURASIP Journal on Applied Signal Processing* (Special Issue on Advanced Signal Processing for Digital Subscriber Lines) in 2005.