

Integrated Radio Resource Allocation for Multihop Cellular Networks With Fixed Relay Stations

Y. Liu, *Student Member, IEEE*, R. Hoshyar, *Member, IEEE*, X. Yang, and R. Tafazolli, *Member, IEEE*

Abstract—Recently, the notion that a logical next step towards future mobile radio networks is to introduce multihop relaying into cellular networks, has gained wide acceptance. Nevertheless, due to the inherent drawbacks of multihop relaying, e.g., the requirement for extra radio resources for relaying hops, and the sensitivity to the quality of relaying routes, multihop cellular networks (MCNs) require a well-designed radio resource allocation strategy in order to secure performance gains. In this paper, the optimal radio resource allocation problem in MCNs, with the objective of throughput maximization, is formulated mathematically and proven to be NP-hard. Considering the prohibitive complexity of finding the optimal solution for such an NP-hard problem, we propose an efficient heuristic algorithm, named integrated radio resource allocation (IRRA), to find suboptimal solutions. The IRRA is featured as a low-complexity algorithm that involves not only base station (BS) resource scheduling, but also routing and relay station (RS) load balancing. Specifically, a load-based scheme is developed for routing. A mode-aware BS resource-scheduling scheme is proposed for handling links in different transmission modes, i.e., direct or multihop. Moreover, a priority-based RS load balancing approach is presented for the prevention of the overloading of RSs. Within the framework of the IRRA, the above three functions operate periodically with coordinated interactions. To prove the effectiveness of the proposed IRRA algorithm, a case study was carried out based on enhanced uplink UMTS terrestrial radio access/frequency-division duplex with fixed RSs. The IRRA is evaluated through system level simulations, and compared with two other cases: 1) nonrelaying and 2) relaying with a benchmark approach. The results show that the proposed algorithm can ensure significant gains in terms of cell throughput.

Index Terms—Cellular networks, fixed relay stations (RSs), multihop, radio resource allocation (RRA).

I. INTRODUCTION

MULTIHOP RELAYING has traditionally been studied in the context of ad hoc wireless networks mainly as a means of enabling the network operation without any infrastructure. In recent years, an upsurge of interest has been observed in the application of multihop relaying in cellular networks in order to create multihop cellular networks (MCNs) [1].

The propagation attenuation of a radio signal is proportional to the link distance raised to the power of a loss exponent. This

Manuscript received October 1, 2005; revised May 1, 2006. This work was supported in part by the European Union under IST Project WINNER (IST-2003-507581).

Y. Liu, R. Hoshyar, and R. Tafazolli are with the Center for Communication Systems Research (CCSR), University of Surrey, Guildford, Surrey GU2 7XH, U.K. (e-mail: y.liu@surrey.ac.uk; r.hoshyar@surrey.ac.uk; r.tafazolli@surrey.ac.uk).

X. Yang was with the Center for Communication Systems Research (CCSR), University of Surrey, Guildford, Surrey GU2 7XH, U.K. He is now with the Mobile Network Solutions Division, NEC Europe Ltd., London W3 6BL, U.K., (e-mail: xinjie.yang@uk.necur.com).

Digital Object Identifier 10.1109/JSAC.2006.881603

exponent can be up to 5.0 in shadowed urban cellular radio environments [13]. Therefore, by breaking a long-distance path into several segments, much lower path loss at each path segment can be achieved. Moreover, heavily shadowed users can employ multihop relaying to bypass obstacles, thereby gaining improved radio channel conditions. Due to the above, multihop relaying provides an opportunity for performance improvements in cellular systems. Nevertheless, multihop relaying has inherent drawbacks, e.g., the requirement for extra radio resources for relaying hops, and the sensitivity to the quality of relaying routes. Therefore, well-designed radio resource allocation (RRA) algorithms are crucial in MCNs, in order to effectively exploit the benefits of relaying, while minimizing its disadvantages.

For RRA in MCNs, aside from the scheduling of the conventional radio resources (this is termed as radio resource scheduling in this paper), transmission route selection should also be considered. In the literature, these two issues have generally been addressed separately.

In the mobile ad hoc network (MANET) research community, many routing algorithms have been proposed [17]. Some examples are optimized link state routing (OLSR), ad hoc on-demand distance vector (AODV), and dynamic source routing (DSR). In essence, these algorithms are designed with network infrastructureless in mind, and their main objective is to establish/maintain network connectivity, rather than to maximize system capacity. As a result, these algorithms are not suitable for MCNs. In recent years, routing in the context of MCNs has become a research issue, and a few algorithms have been proposed so far, e.g., location-based routing [4], path-loss-based routing [7], transmission-power-based routing [8], and congestion-based routing [3]. These algorithms can, to some extent, solve the route selection problem in MCNs. However, the selected routes are not necessarily optimal in terms of the system resource utilization efficiency.

In conventional cellular networks, the radio resource scheduling algorithms have been investigated thoroughly. Currently, the widely used algorithms in practical packet-based cellular networks are usually based on user prioritization combined with greedy resource loading [12]. However, these algorithms are all designed based on the assumption that every user in the system is directly connected to the BS, whereas this assumption is no longer valid in MCNs. Some work on resource scheduling in multihop systems has been published recently. In [2], a fractional bandwidth and power allocation algorithm is proposed for orthogonal regenerative frequency-division multihop communication systems, and in [9] a centralized downlink scheduling scheme is proposed for cellular networks utilizing small

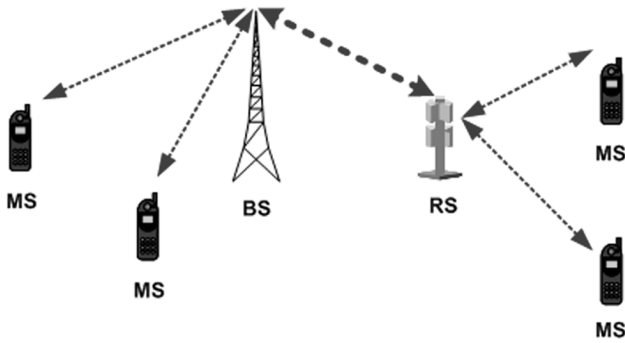


Fig. 1. An example of connectivity of network nodes in the system scenario.

number of relays. However, these algorithms generally do not take the capacity limits of RSs into account, and therefore during the resource scheduling, no consideration is given to the prevention of RS overloading.

The coordination between routing and resource scheduling in MCNs is also crucial and warrants careful investigation, especially since strong interdependency between the two functions is envisioned in MCNs. The strategy for effective coordination of routing and packet scheduling in packet-based MCNs, however, has not yet been fully investigated.

In this paper, transmission routes are treated as one extra dimension of the resource space of conventional cellular networks. The route selection and conventional radio resource scheduling are then integrated into one problem in the context of MCNs. The throughput maximization RRA problem is mathematically formulated and proven to be NP-hard. Considering the prohibitive complexity of searching for the optimal solution to such an NP-hard problem, we propose instead an efficient heuristic algorithm, named integrated radio resource allocation (IRRA) algorithm to find suboptimal solutions. To prove the effectiveness of the proposed algorithm, a case study was carried out based on enhanced uplink UMTS terrestrial radio access/frequency division duplex (UTRA/FDD) with fixed relay stations (RSs). As shown by the simulation results, the proposed algorithm can ensure significant throughput gains compared with the nonrelaying case and to a benchmark relaying approach.

The rest of this paper is organized as follows. The system scenario is described in Section II. In Section III, the throughput maximization RRA problem is formulated and proven to be NP-hard. In Section IV, a discussion on the design of an efficient algorithm for the NP-hard problem is presented. Next, based on the conclusions, the IRRA algorithm is proposed. A case study of the IRRA follows in Section V. Finally, Section VI concludes this paper.

II. SYSTEM SCENARIO

A multicell cellular radio scenario is considered. In each cell, a number of fixed RSs are deployed. A user terminal can either connect to the BS directly, or via an RS. Between RSs and the BS, good channel conditions (e.g., line-of-sight) are assumed, and hence direct transmissions are envisioned. An RS is assumed to be able to act as a relay for multiple users and has large enough buffer to hold the data being relayed. An example

TABLE I
POSSIBLE FREQUENCY-DIVISION RELAYING CASES

	BS-MS links	RS-MS links	BS-RS links
Case A	$fb A1$	$fb A1$	$fb A2$
Case B	$fb B1$	$fb B2$	$fb B1$
Case C	$fb C1$	$fb C2$	$fb C3$

of connectivity of network nodes in our system scenario is illustrated in Fig. 1.

In order to avoid the self-interference of RSs, the transmission and reception of an RS should not take place on the same frequency at the same time. Consequently, two basic relaying strategies are envisioned for multihop systems: frequency-division relaying and time-division relaying [2]. Under these schemes, the transmission and reception of an RS are either separated in the frequency-domain, e.g., by using different carrier or subcarrier frequencies, or in the time-domain, e.g., by employing different time slots or medium access control (MAC) frames. The study in this work is based on frequency-division relaying.

In order to reduce system complexity, in our work, it is assumed that each type of link has a fixed frequency band assigned to it. Therefore, three possible frequency-division relaying cases are envisioned, as shown in Table I.

In case A, BS-MS and RS-MS links share the same frequency band, and thus a user terminal does not need to switch frequency bands during a change of transmission/reception mode (from direct to multihop or *vice versa*). From the perspective of terminal complexity and mode switching latency, this is an advantage compared with cases B and C, especially when the respective frequency bands used in BS-MS and RS-MS links are spectrally distant.

In cases A and C, one frequency band is dedicated solely to BS-RS links. This configuration can effectively exploit the capacity of these links, which normally have high resource utilization efficiency due to very good radio channel conditions, as well as high traffic volume due to the fact that one RS can serve multiple users.

In this work, case A is assumed. Nevertheless, cases B and C will be studied and evaluated in our future work.

It is worth noting that the algorithm in this work is derived from frequency-division relaying case A. However, it can also be applied to the similar scenario of time-division relaying, in which BS-RS links have dedicated time resources but BS-MS and RS-MS links share the same time resources.

III. PROBLEM FORMULATION AND PROOF OF NP-HARDNESS

In this section, we formulate the throughput maximization radio resource allocation (TM-RRA) problem and prove it to be NP-hard.

We assume that the area of interest consists of M cells, each of which is overlaid by N_m RSs and serves K_m users, where $\sum_{m=1}^M N_m = N$ and $\sum_{m=1}^M K_m = K$.

We define a radio resource unit in the underlying cellular network (e.g., the combination of a data rate and a time slot) as a conventional resource unit (CRU), and the combination of a CRU and a transmission route as a general resource unit (GRU).

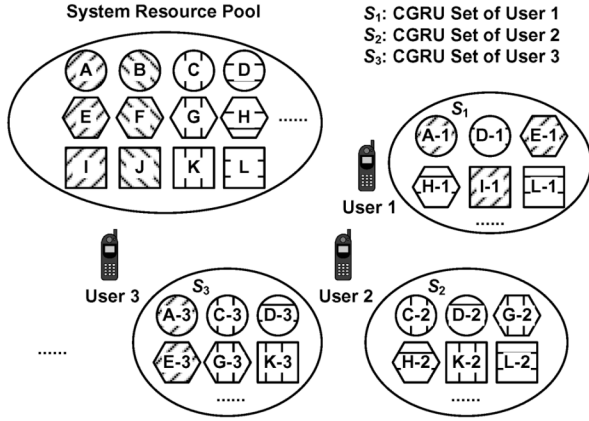


Fig. 2. System resource pool and user CGRU sets.

A GRU is essentially one allocation choice for a user, and it can be characterized by multiple attributes, which might vary depending on the nature of the cellular system considered. For example, in a TDMA-based MCN, one GRU may be characterized by three attributes: a time slot, a data rate, and a route. By contrast, in an FDMA-based MCN, a GRU might be characterized by a frequency, a data rate, and a route.

GRUs in the system can be easily determined based on a preexisting knowledge of the system’s conventional resource and the network topology. An example of the system resource pool is shown in Fig. 2, where the GRUs are indexed by A, B, C , and so on. Different shapes or fill patterns of GRUs represent different attributes. For example, we let the round, hexangular, and square shapes correspond to rates of 32, 64, and 128 kb/s, respectively, whereas “backward diagonal line,” “forward diagonal line,” “vertical line,” and “horizontal line” fill patterns represent routes 1–4, respectively. Then, GRUs $A = (32 \text{ kb/s, route 1})$, $B = (32 \text{ kb/s, route 2})$, and so on.

In our work, we assume that a user can only be allocated one GRU, i.e., one combination of a CRU and a route, at any allocation instant in time. In the case where multiple CRUs or multiple routes can be used for a user simultaneously, to keep the above assumption valid, we regard each combination of the multiple CRUs or routes as one extended CRU or route. GRUs containing normal or extended CRUs/routes will all be considered so that the RRA algorithm still only produces at most one GRU per user at any one allocation instant in time.

During resource allocation, not all the GRUs in the system resource pool would be valid candidates for each user. Some of them may not be required by a user, for instance if the data rate of the GRU is too high for a user given its queue size. We define these kinds of GRUs as unneeded GRUs (UNGRUs). In other cases, GRUs may not be usable by a user, for instance if the data rate is too high for the user to reach given its power limit. We define these kinds of GRUs as unusable GRUs (UUGRUs). We regard all other resource units, which are not UNGRUs or UUGRUs, as candidate GRUs (CGRUs), which are to be scheduled by the RRA algorithm.

Let S_k denote the CGRU set (CGRUS) of user k , which includes all the CGRUs of that user. It can be easily precalculated based on the knowledge of the system resource pool, as well as on the system status and the user’s profile (e.g., queue size and

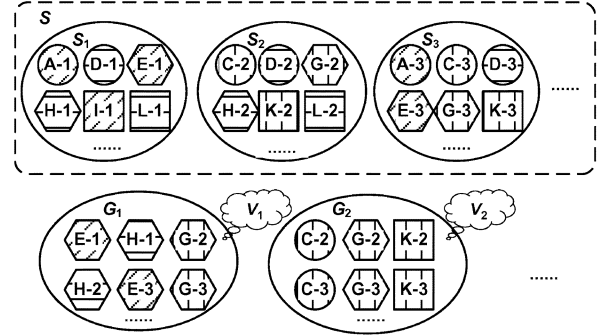


Fig. 3. User CGRU sets and relevant CGRU groups.

power headroom). As shown by the example in Fig. 2, GRUs A, D, E, H, I, L, \dots are valid candidates for user 1, and therefore $S_1 = \{A - 1, D - 1, E - 1, H - 1, I - 1, L - 1, \dots\}$. Each of the elements represents a resource allocation choice for user 1.

It is worth noting that one CGRU essentially contains two pieces of information: the GRU and the associated user. For example, $\text{CGRU } A - 1 = (\text{GRU } A, \text{user } 1)$. CGRUs in different CGRUSs have different associated users, hence should be regarded as different items. Consequently, $S_h \cap S_k = \Phi$ if $h \neq k$.

We let S denote the union of all the CGRUSs from S_1 to S_K , as shown in Fig. 3. We define the action selecting an element (CGRU) j from set S as having the meaning of allocating the corresponding GRU to the associated user. Each such action achieves a certain amount of throughput, while consuming some of the system’s capacity. We let x_j denote a Boolean variable indicating whether element j is selected (1 for “Yes” and 0 for “No”), and $d_j, q_{j,BSm}$, and $q_{j,RSn}$ denote the achieved throughput, the consumed capacity of the BSm and the consumed capacity of the RSn , respectively, if element j is selected.

We note that in S , some elements might share exclusive attributes. For example, in a time-division multiple-access (TDMA) system, a time slot is normally treated as an exclusive resource, i.e., if assigned to a user, it should not be used by others. In order to take this issue into account, we assume that set S includes totally I relevant CGRU groups (RCGRUGs) (G_1, \dots, G_I), each of which shares a certain attribute (e.g., a data rate or a route). As shown by the example in Fig. 3, all the items in G_1 share the attribute of hexangular shape (i.e., 64 kb/s), whereas all the items in G_2 share the attribute of having “vertical line” fill pattern (i.e., route 3). An element of S can be a member of multiple groups, such as $G - 2$, which is in both G_1 and G_2 . To introduce the aforementioned resource exclusivity, we assume that for $G_i, (i = 1, \dots, I)$, at most v_i elements can be selected for assigning to their associated users. The value of v_i is known depending on the attribute that the group is sharing. For instance, it is infinite for a group sharing a certain data rate, i.e., a data rate can be given to infinite number of users simultaneously. By contrast, in a TDMA system, it is 1 for the group sharing a particular time slot.

Assuming the above, the TM-RRA problem can be described as follows: given a particular system snapshot (system capacity, user profiles, etc.), select elements from the set S so as to achieve maximum system throughput, while adhering to the

following three constraints: 1) for each BS or RS, the aggregate capacity consumption by the selected elements of S is less than the corresponding capacity limit c ; 2) for each CGRUS (S_1, \dots, S_K), only one element is selected; and 3) for each RCGRUG (G_1, \dots, G_I), at most v_i elements are selected.

Mathematically, the TM-RRA problem can be formulated as

$$\text{Maximize } z = \sum_{j \in S} d_j x_j \quad (1)$$

Subject to

$$\sum_{j \in S} q_{j, \text{BS}_m} x_j \leq c_{\text{BS}_m} \quad \forall m = 1, \dots, M \quad (2)$$

$$\sum_{j \in S} q_{j, \text{RS}_n} x_j \leq c_{\text{RS}_n} \quad \forall n = 1, \dots, N \quad (3)$$

$$\sum_{j \in S_k} x_j = 1 \quad k = 1, \dots, K \quad (4)$$

$$\sum_{j \in G_i} x_j \leq v_i \quad i = 1, \dots, I \quad (5)$$

$$x_j = 0 \text{ or } 1 \quad j \in S. \quad (6)$$

Theorem 1: The TM-RRA problem is NP-hard.

Proof: We restrict the TM-RRA problem to the special case, where M is 1, N is 0, and v_i is infinite for all i from 1 to I . In this case, the constraints (3) and (5) are removed, and (2) is changed to the following form with the subscript m being dropped:

$$\sum_{j \in S} q_{j, \text{BS}} x_j \leq c_{\text{BS}}. \quad (7)$$

For CGRU j in set S , the achievable throughput d_j is known and the BS capacity consumption $q_{j, \text{BS}}$ can be easily precalculated in this case based on the interference status of the GRU, and on the channel gains and QoS requirements, etc. Moreover, cell capacity c_{BS} is normally a pre-known hardware limitation (e.g., maximum transmission power) or a value preset by the network operator (e.g., maximum load or rise over thermal (RoT) [11]). Furthermore, as mentioned earlier, all CGRUSs (S_1 to S_K) can be precalculated, and hence, it is clear that the only unknown in this special case is x_j .

Now, we recall a well-known NP-hard problem, the multiple choice knapsack problem (MCKP) [15], as follows: Given a knapsack, an item set B including b items, and a partition of the item set B_1, B_2, \dots, B_r , with $p_j =$ profit of item j , $w_j =$ weight of item j , and $c_{\text{kp}} =$ capacity of the knapsack, select items from set B so as to maximize the aggregate profit from all the selected items, while adhering to the following constraints: 1) the aggregate weights of the selected items should be less than the (weight bearing) capacity of the knapsack and 2) for each item subset (B_1 to B_r), only one item is selected.

Note that any arbitrary instance of the MCKP can be viewed as an instance of the aforementioned special case of the TM-RRA problem, by applying the following mappings: the item set B is mapped to the set S (including the mappings from p_j and w_j of the items in set B to d_j and q_j of the elements in set S , respectively), the item subsets B_1, \dots, B_r are mapped to the CGRUSs S_1, \dots, S_k , respectively, and finally, the knapsack

capacity c_{kp} is mapped to the BS capacity c_{BS} . The item selection for that MCKP instance is in one-to-one correspondence with the CGRU selection for the mapped instance of the special case of the TM-RRA problem. Thus, the MCKP can be regarded as a restricted version of the TM-RRA problem. Since the MCKP has already been proven to be NP-hard [15], we conclude that the TM-RRA problem is NP-hard.

IV. KEY ISSUES OF THE RRA DESIGN IN MCNS AND A HEURISTIC ALGORITHM

Due to the NP-hardness of the TM-RRA problem, the time required to find its optimal solution increases exponentially with the size of the problem. This prohibitive processing complexity is intolerable in practical systems.

In this section, to aid with the design of a practical RRA algorithm for MCNs, we discuss routing, radio resource scheduling, and their interactions. Based on the conclusions, an efficient heuristic algorithm, named IRRA, is derived.

A. Routing

From the architectural perspective, two basic routing strategies are envisioned in relay-based systems: centralized and distributed. Under the centralized strategy, routing is performed in a central controller, which normally possesses powerful processing capabilities so that sophisticated routing algorithms may be adopted to optimize the system performance. However, this strategy requires extensive information gathering from the distributed network nodes to the central controller, which inevitably results in signaling overheads and processing delay. With the distributed strategy, all network nodes from the source to the destination jointly perform route determination. This strategy can function when no central controller is reachable, but its performance is normally limited by network nodes' processing capabilities and knowledge of the network status.

In our system scenario, a BS could serve as a central controller in a cell, and, since maximally two hops are considered in our study, the signaling overhead as well as the signaling delay would be acceptable. Therefore the centralized strategy is chosen in our work in order to optimize the system performance.

As for the algorithm of routing, take a close look at the TM-RRA problem in (1)–(6), we can easily conclude that the optimal route of a user is the one that has highest achievable throughput with unit amount of induced system load. In our work, such a load-based routing algorithm is proposed, and its routing cost function is as follows:

$$\zeta_{\text{MS}_k - \text{CXX}} = \frac{C_{\text{MS}_k - \text{CXX}}}{R_{\text{MS}_k - \text{CXX}}} \quad (8)$$

where $\zeta_{\text{MS}_k - \text{CXX}}$ is the *load cost indicator (LCI)* of the route from user k to its connected receiver/transmitter CXX (a BS or an RS), and $C_{\text{MS}_k - \text{CXX}}$ and $R_{\text{MS}_k - \text{CXX}}$ represent the consumed system capacity of the route and the data rate on the route, respectively.

This cost function reflects the consumed system capacity when delivering unit amount of traffic on a particular route. Given a certain system capacity constraint, if every user employs the route with the least LCI from among the possible routes, the system throughput can be maximized.

It is worth noting that in the route-selection process for a user, in order to guarantee a fair comparison, it is important to calculate the LCIs of all of its routes based on the same data rate. However, the choice of this data rate normally does not affect the routing result as long as it is the same for all routes. This is shown in the case study in Section V, where after further derivation of (8), the data rate is removed completely from the cost function. However, if the result of routing does depend on the choice of the data rate in (8), the routing process should be performed based on the current data rate, or on the one that is assigned by the network and thus will be employed in the forthcoming transmission time intervals (TTIs). If these rates are zero, the minimum nonzero rate can be taken as an approximation.

B. Radio Resource Scheduling

As for the radio resource scheduling in MCNs, two basic strategies are envisioned: centralized and hierarchical. With the centralized strategy, the radio resource scheduling is mainly performed by the BS. By contrast, in the case of the hierarchical strategy, a layered algorithm architecture is foreseen: at the higher layer, the BS splits the radio resource space into fragments for individual RSs and itself, respectively. At the lower layer, the actual radio resource scheduling is performed by RSs (for multihop users) and the BS (for direct transmission users) according to their given resource spaces.

In general, the centralized strategy is simpler than the hierarchical one. The latter, on the other hand, is able to perform faster radio resource scheduling in order to quickly adapt to system dynamics. For the sake of simplicity, and easier coordination with the centralized routing, the centralized (BS-based) radio resource scheduling strategy is employed in our work. Nevertheless, the hierarchical strategy is an interesting topic for our future study.

In our system scenario, the resource scheduling for BS-RS links can still use the conventional algorithm due to the fact that these links are always in direct transmission mode and have dedicated frequency band. This issue will not be discussed further in this paper.

The scheduling for BS-MS and RS-MS links, on the other hand, is more complicated due to the sharing of the frequency spectrum. In this case, when allocating resources to users, not only the capacity constraint of the BS, but also those of RSs should be respected. Towards that end, we propose the following strategy: first, perform BS resource scheduling based on the capacity constraint of the BS, and then carry out RS load balancing to fine-tune the resource assignments in order to prevent the overloading of RSs. There are two issues requiring attention under this strategy.

First, BS resource scheduler should be aware of users' transmission/reception modes, and treat users in different modes differently. The main reason is that multihop users are not directly connected with BSs, hence the estimation of their consumed BS capacities is different from the conventional case for direct transmission users.

Second, the RS load balancing function should have a user prioritization mechanism in order to decide whose resource as-

signments should be tuned when an RS is overloaded. For instance, the following priority function can be used:

$$x_{\text{MS}_k} = \frac{\zeta_{\text{MS}_k\text{-BS}}}{\zeta_{\text{MS}_k\text{-RS}_n}} \quad (9)$$

where x_{MS_k} is the priority of user k , and RS n is the tentative RS (overloaded) of user k . Apparently, the higher the x_{MS_k} , the greater benefits this user can potentially bring to the system by using its tentative relaying route (instead of the direct transmission route). Therefore, if we use this priority function and start the overloading relief process from the user with the lowest priority, the benefit of multihop relaying can be preserved as much as possible.

C. Interactions Between Routing and Radio Resource Scheduling

In MCNs, multihop transmissions normally consume less system capacity if routing is appropriately performed, and the radio resource scheduler should promptly capture the saved resources and assign them to others who suffer a deficit. Consequently, radio resource scheduling should be based on the results of routing. On the other hand, radio resource scheduling affects the system interference/loading pattern, which in turn might affect the decisions of user route selection.

To perform RRA in the context of the above strongly interdependent scenario, the best mechanism is to perform global cross-optimization with lots of iterations. But as proven in Section III, the problem is NP-hard. Thus, the optimization algorithms are hardly to be feasible in practical systems.

A more practical solution would be to perform routing and resource scheduling within one algorithm framework, where resource scheduling always takes place immediately after routing. This allows the saved resources from multihop relaying to be captured promptly. Moreover, RS load balancing is considered within resource scheduling so that overloading of the RSs can be effectively avoided. This algorithm can be run periodically to quickly adapt resource assignments as required by the system dynamics.

D. Integrated Radio Resource Allocation (IRRA)

Based on the above discussions, we propose an IRRA algorithm, and it comprises three entities. The first entity, named, Load-Based Route Manager (LBRM), is responsible for the selection of user routes based on (8). The second, named, Base Station Resource Scheduler (BSRS) is for the scheduling of user resource assignments based on the available BS capacity. Users in different modes should be treated differently during this scheduling process. Finally, the third entity, named, Relay Station Load Balancer (RSLB), is responsible for the fine-tuning of user assigned resources/routes, based on a priority function, e.g., the one shown in (9), in order to avoid the overloading of RSs.

The IRRA is executed periodically as illustrated in Fig. 4, and the basic procedure in each execution is as follows.

- 1) The LBRM is run first, and outputs user routes.
- 2) Based on the user routes from the LBRM, the BSRS is executed and outputs user resource assignments (e.g., data rates) based on the available BS capacity.

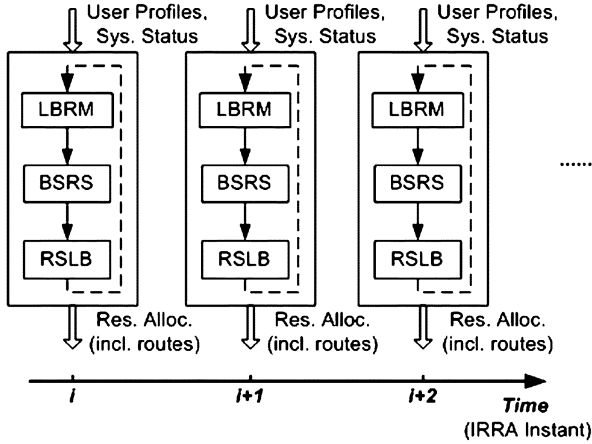


Fig. 4. The overall operating procedure of the IRRA.

- 3) Finally, the RS load balancing is carried out based on the user routes and resource assignments from the LBRM and the BSRS, respectively. The outputs of the RSLB are fine-tuned user resource allocations (including user routes).

The IRRA is generally a low complexity algorithm. Nevertheless, it is worth mentioning that at one execution instant, the above basic procedure can also be run iteratively, i.e., after step 3), return to step 1) to do the routing again based on the updated system status, then step 2), and so on, until the output converges or a predefined maximum number of iterations is reached. However, if the iterative procedure is used, the computational complexity of the algorithm increases.

V. A CASE STUDY IN ENHANCED UPLINK UTRA-FDD WITH FIXED RSS

To prove the effectiveness of the IRRA algorithm, in this section, we apply it to enhanced uplink UTRA-FDD [11] with fixed RSSs. The entities of the IRRA are described in greater detail, the algorithm complexity is estimated, and finally, the system performance is evaluated through system level simulations.

A. System Scenario

In this case study, to apply the frequency-division relaying explained in Section II, we introduce an extra carrier frequency for RS-BS links, whereas MS-RS links share the original carrier frequency with MS-BS links.

The resource scheduling in enhanced uplink UTRA-FDD is essentially transport format combination (TFC) selection [11]. Note that each TFC corresponds to a certain transmission data rate, and hence transmission rates are the major type of resource in enhanced uplink UTRA-FDD.

It is worth mentioning that the allocation of spreading codes is not considered in enhanced uplink UTRA-FDD, due to the fact that each user has a uniquely assigned scrambling sequence, thus the spreading code resources in use by a user do not affect those available for others.

B. Application of the IRRA

1) *Load-Based Route Manager (LBRM)*: In the LBRM, the route selection of each user is based on the LCIs of all of its possible routes. For MCNs with different underlying cellular sys-

tems, the LCI formula (8) can be derived into different concrete forms, depending on the interpretation of the system capacity. In the context of this case study, (8) can be derived further as follows.

Since the capacity of enhanced uplink UTRA-FDD is interference-limited, we define the consumed system capacity of a route as the total induced system interference by the transmitter

$$\begin{aligned}
 C_{MS_k-CRX} &= \sum_{m=1}^M P_{rx,MS_k-BS_m} + \sum_{n=1}^N P_{rx,MS_k-RS_n} \\
 &= P_{tx,MS_k} * \left(\sum_{m=1}^M \frac{1}{L_{MS_k-BS_m}} + \sum_{n=1}^N \frac{1}{L_{MS_k-RS_n}} \right) \\
 &= P_{tx,MS_k} * \rho_{MS_k} \quad (10)
 \end{aligned}$$

where P_{tx,MS_k} is the transmission power of user k , P_{rx,MS_k-BS_m} and P_{rx,MS_k-RS_n} represent the received power of user k at BS m and at RS n , respectively, $L_{MS_k-BS_m}$ and $L_{MS_k-RS_n}$, respectively, represent the end-to-end transmission losses between user k and BS m , and between user k and RS n , and finally, ρ_{MS_k} is a user-specific factor largely determined by the user's location.

Similar as the load factor in [14, Sec. 8.2.2.1], the load of user k at any receiver of interest, η_{MS_k-RX} , is defined as follows in this case study:

$$\eta_{MS_k-RX} = P_{rx,MS_k-RX} / I_{tot,RX} \quad (11)$$

where P_{rx,MS_k-RX} is the received power of user k at the receiver of interest, and $I_{tot,RX}$ represents the total interference at the receiver. In order to guarantee the quality of signal reception, the signal-to-interference ratio (SIR) of the user at its connected/targeted receiver should be maintained around a certain target value, which normally is a pre-known value depending on the user data rate and the quality-of-service (QoS) (e.g., block error rate) requirement. Therefore, the load of user k at its connected receiver, η_{MS_k-CRX} , can be calculated as follows [12], [14]:

$$\begin{aligned}
 \eta_{MS_k-CRX} &= \frac{P_{rx,MS_k-CRX}}{I_{tot,CRX}} = \frac{SIR_{MS_k-CRX}}{1 + SIR_{MS_k-CRX}} \\
 &\approx \frac{SIR_{t,MS_k-CRX}}{1 + SIR_{t,MS_k-CRX}} \quad (12)
 \end{aligned}$$

where SIR_{MS_k-CRX} and SIR_{t,MS_k-CRX} represent the actual received SIR and the SIR target of user k at its connected receiver, respectively.

Then, P_{tx,MS_k} can be obtained from the following equation:

$$\begin{aligned}
 P_{tx,MS_k} &= P_{rx,MS_k-CRX} * L_{MS_k-CRX} \\
 &= \eta_{MS_k-CRX} * I_{tot,CRX} * L_{MS_k-CRX} \quad (13)
 \end{aligned}$$

where η_{MS_k-CRX} can be calculated with (12), and $I_{tot,CRX}$ and L_{MS_k-CRX} can normally be estimated based on system measurements.

Considering (8), (10), (12), and (13) all together, we obtain

$$\begin{aligned}
 \zeta_{MS_k-CRX} &\approx \frac{SIR_{t,MS_k-CRX} * I_{tot,CRX} * L_{MS_k-CRX} * \rho_{MS_k}}{(1 + SIR_{t,MS_k-CRX}) * R_{MS_k-CRX}} \quad (14)
 \end{aligned}$$

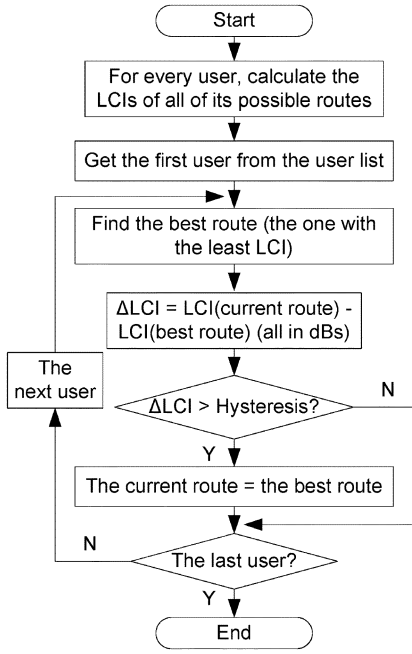


Fig. 5. The operating procedure of the LBRM.

Since during the route determination, the comparison of LCIs is always performed between all the possible routes of a particular user, ρ_{MS_k} , R_{MS_k-CRX} and SIR_{t,MS_k-CRX} can be considered to be fixed values. Therefore, (14) can be redefined, without affecting the routing results, as follows:

$$\zeta_{MS_k-CRX} \approx I_{tot,CRX} * L_{MS_k-CRX}. \quad (15)$$

Based on (15), the LBRM can easily determine and select the route with the highest load efficiency for each user. However, in order to avoid “ping-pong” effects causing unnecessary signaling overheads and oscillations, an LCI hysteresis (e.g., 2 dB) should be applied when updating user routes. The operating procedure of the LBRM is illustrated in Fig. 5.

It is worth noting that in the LBRM, the route selections for individual users are performed independently. Consequently, a number of users may happen to choose the same RS, which then becomes overloaded. Therefore, the routes selected in the LBRM might require fine-tuning by the subsequent RSLB. The reason for such a design, instead of performing the fine-tuning in the LBRM itself, is that the RSLB has full knowledge of all users’ routes and resource assignments as calculated by the LBRM and the BSRS, respectively, and therefore, it is able to make better decisions.

2) *Base Station Resource Scheduler (BSRS)*: The operating procedure of the BSRS is shown in Fig. 6. Note that the rate allocated to a user could be zero, which would indicate that the user is not allowed to transmit in the next IRRA period. Load estimation needs to be performed for both the current and the new rates, in order to see how much extra load will be incurred at the BS. Unlike conventional scheduling algorithms in enhanced uplink UTRA-FDD [12], the BSRS employs different load estimation approaches for users in different transmission modes.

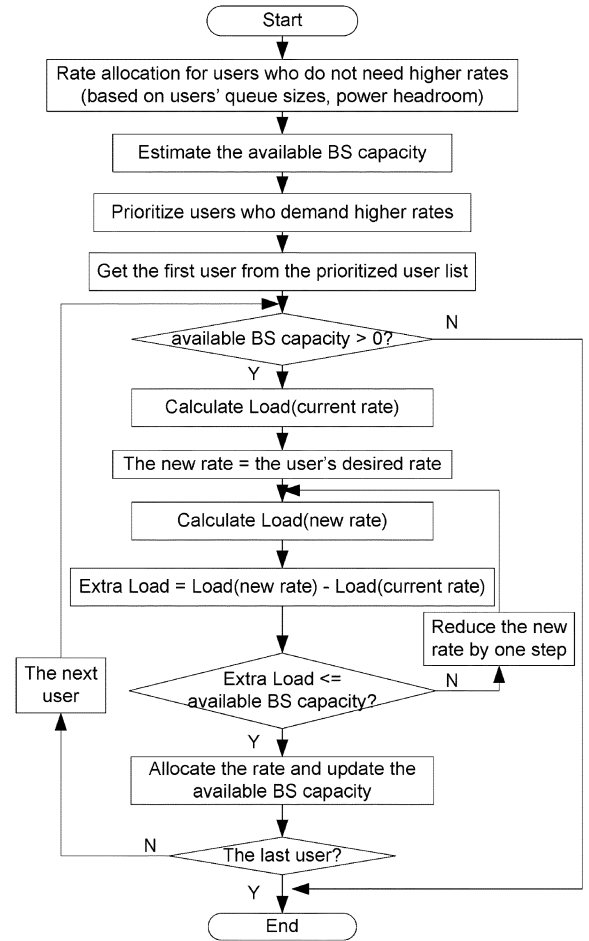


Fig. 6. The operating procedure of the BSRS.

For a user in direct transmission mode, the load of the user at the BS can be easily calculated with (12). On the other hand, for a user in multihop mode, since its connected receiver is a RS, not the BS, (12) can only be employed to calculate its load at the connected RS. Nevertheless, based on (11) and (13), it is not hard to see that the load of user k at its interfered receiver, η_{MS_k-IRX} , can be derived from the load at its connected receiver η_{MS_k-CRX} with the following equation:

$$\eta_{MS_k-IRX} = \frac{\eta_{MS_k-CRX} * I_{tot,CRX} * L_{MS_k-CRX}}{I_{tot,IRX} * L_{MS_k-IRX}}. \quad (16)$$

Therefore, the load at the BS of a user in multihop mode should be calculated based on (12) and (16) together.

3) *Relay Station Load Balancer (RSLB)*: As illustrated in Fig. 7, the operating procedure of the RSLB is: first, estimate the loads of RSs with (12) and (16) based on the tentative user transmission routes and rates from the LBRM and the BSRS, respectively, then check whether any RSs would be overloaded, and if so, commence the overloading-relief process for those overloaded RSs in order of descending load.

During the overloading-relief process of a particular RS, as mentioned in Section IV-B, a user prioritization mechanism is required in order to decide from which user the process starts, and with what sequence it proceeds. The priority function employed in this case study is provided in (9).

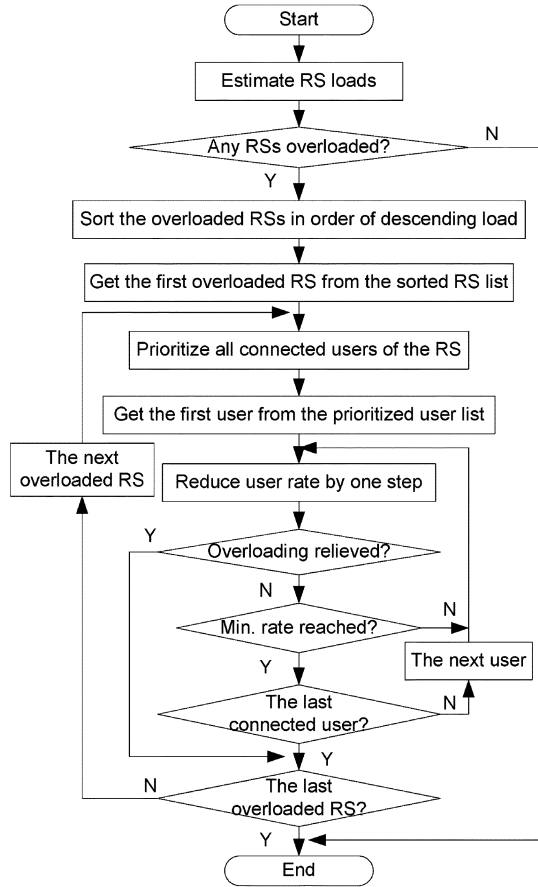


Fig. 7. The operating procedure of the RSLB.

As for the overloading-relief approaches, two basic options are envisioned: switching the chosen user from the overloaded RS to an underloaded one, or reducing its data rate. The first scheme is generally applicable in scenarios where the RS capacity is limited by transmission power or by time slots, etc., instead of by interference. In these cases, if we switch a user from an overloaded RS to an underloaded one, the load caused by that user is then completely shifted from the old RS to the new one. However, if the capacity of the RS is limited by interference, the first scheme is not effective, since the load (i.e., interference in this case) cannot be completely shifted between RSs, and if the underloaded RS is too far away, the user might generate even more load on the old RS after being switched to the new one. In this case study, the second approach is employed.

C. Algorithm Complexity

We consider the cell of interest to be cell m , which includes one BS, N_m RSs, and K_m users. We assume that each user has D candidate rates $R_1 > R_2 > \dots > R_D$, and that a user can only be relayed by RSs in the same cell, and thus the total number of possible routes for a user is $N_m + 1$.

As for the RBLM, the LCI calculation for a particular user requires $N_m + 1$ multiplications based on (15), the searching for its best route requires N_m comparisons, and the subsequent user route updating with respect to the given hysteresis at most requires two operations. Hence, at most, $K_m(2N_m + 3)$ operations are required in one execution of the LBRM.

It is not hard to see that the worst case of the BSRS is when all users are in multihop mode and are demanding higher rates. In this case, the BSRS involves two steps only. The first one is to prioritize all the users, and this requires at most $(K_m - 1)K_m/2$ operations [18]. The second step is to allocate rates to individual users. We assume that user k 's current rate is rate number f_k ($1 \leq f_k \leq D$). Then, at most, f_k rates (R_1, \dots, R_{f_k}) need to be tested for this user. In each of these tests, six operations are required to calculate the new load based on (12) and (16), and two operations are required to calculate the extra load and compare it with the available BS capacity. Hence, together with the six operations required by the estimation of a user's current load, at most, $\sum_{k=1}^{K_m} (8f_k + 6)$ operations are required in the second step. Therefore, the BSRS involves at most $K_m^2/2 + 11K_m/2 + \sum_{k=1}^{K_m} 8f_k$ operations in each execution.

The worst case of the RSLB is when all users are connected to a single (overloaded) RS so that the complexity for prioritizing the connected users is maximized. In this case, the RSLB involves three steps. The first step is to estimate the load of the RS, and this requires load estimations of K_m users based on (12), and then a summation of K_m load values. These require totally $3K_m - 1$ operations. The second step is to prioritize the connected users of the RS, and this requires at most $(K_m - 1)K_m/2$ operations [18]. The final step is the overloading-relief for the RS. We assume that user k is assigned rate number f'_k ($1 \leq f'_k \leq D$) by the BSRS. Then, at most, $D - f'_k$ rates ($R_{f'_k+1}, \dots, R_D$) need to be tested for this user in the relief process. Each test involves two operations for calculating the load of the new rate based on (12), and two operations for checking whether or not the resulting load reduction is sufficient for relieving the overloading. Hence, at most, $\sum_{k=1}^{K_m} 4(D - f'_k)$ operations are required in the final step. Therefore, the RSLB involves at most $K_m^2/2 + K_m(5/2 + 4D) - 1 - \sum_{k=1}^{K_m} 4f'_k$ operations.

Note that in the worst case of the BSRS, rate number f_k is finally assigned to user k ($k = 1$ to K_m). Therefore, the f'_k used in the RSLB should be the same as the f_k employed in the BSRS. Consequently, the total number of operations in the three entities is at most $K_m^2 + K_m(2N_m + 4D + 11) - 1 + \sum_{k=1}^{K_m} 4f_k$. It is not hard to see that the global worst case occurs when $f_k = D$ ($k = 1$ to K_m). Therefore, the worst-case algorithm complexity can be approximated as $O(K_m^2 + 2N_m K_m + 8DK_m + 11K_m)$.

It is clear that this quadratic complexity is far lower than the complexities of optimal algorithms which grow exponentially with the size of the problem.

D. Simulation Results

By system level simulation, the IRRA is evaluated and compared with the nonrelaying case, as well as with a benchmark relaying approach that adopts pathloss-based routing [7] for user route determination and a conventional scheduling algorithm [12] for rate scheduling. The simulation parameters and settings are listed in Table II.

Fig. 8 demonstrates the obtained throughput with different session arrival rates. We observe that when the offered traffic load (reflected by the session arrival rate) is light, without the help of relaying, traffic can still be delivered satisfactorily. However, when the offered traffic load becomes heavy, systems with

TABLE II
 SIMULATION PARAMETERS

Parameters	Settings/Explanations
Cellular layout	Hexagonal grid, omni-directional sites, 3 tiers with wrap around
Cell radius R (km)	1.8
Propagation model (dB)	$128.1 + 37.6 \log_{10}(R)$
Channel type	3GPP Pedestrian A
Std. deviation of slow fading (dB)	8.0
Correlation distance of slow fading (m)	50
BS antenna gain plus cable loss (dBi)	14
User antenna gain (dBi)	0
Maximum user power (dBm)	21
Maximum RS power (dBm)	24
Links with closed-loop power control	MS-BS, MS-RS, RS-BS
Closed-loop power control step size (dB)	1
User TFCS (Transport Format Combination Set) (kbps)	8,16,32,64,128,256,384
RS TFCS (kbps)	8,16,32,64,128,256,384,768,1000
TTI (Transmission Time Interval) (ms)	10
Scheduling period / IRRA period (ms)	100
Priority function for BS scheduling	Proportional fairness[12]
Traffic model	Near real time video[10]
Session arrival distribution model	Poisson
Session arrival rate (session/cell/s)	0.25; 0.5; 0.75; 1.0; 1.25
Session duration distribution model	Shifted exponential [10]
Minimum session duration (s)	20
Mean session duration (s)	40
Number of fixed RSs per cell	6, symmetrically located on the perimeter of a circle
RS-to-BS distance r	0.65 $\ast R$; 0.75 $\ast R$; 0.85 $\ast R$
BS/RS load threshold	70%

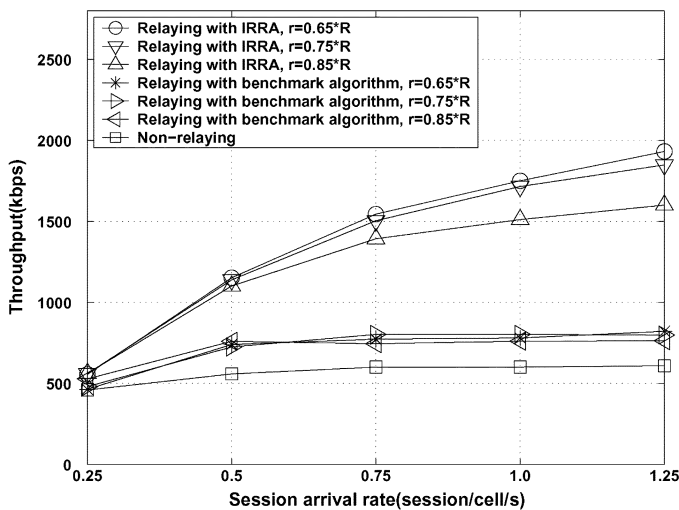


Fig. 8. Average cell throughput versus session arrival rate.

relaying perform significantly better than that without. Moreover, in the case of relaying with the IRRA, the ratio of relayed cell throughput to total cell throughput can be up to 75%, as indicated in Fig. 9.

The IRRA can achieve much higher cell throughput than the benchmark relaying approach. This is mainly due to the following reasons: first, the routing criterion used in the IRRA, i.e., LCI, can well reflect the most beneficial route for a certain user in terms of uplink resource/load efficiency. Second, the benchmark relaying approach has no mechanism to effectively avoid

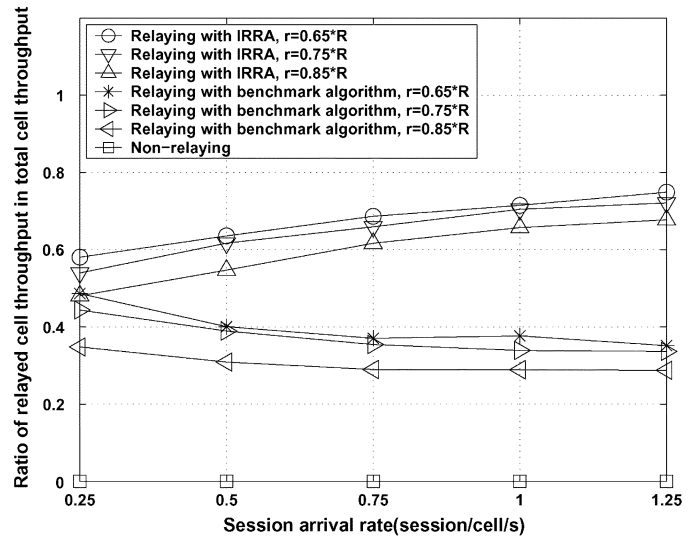


Fig. 9. Ratio of relayed cell throughput to total cell throughput versus session arrival rate.

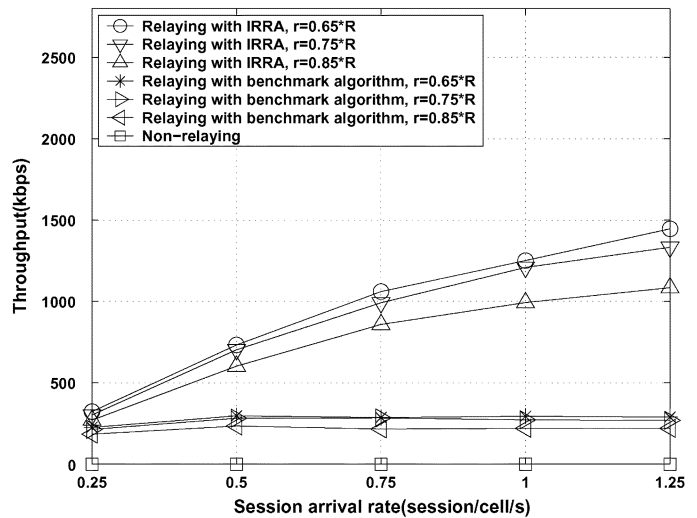


Fig. 10. Average relayed cell throughput versus session arrival rate.

the overloading of RSs. As a result, RSs are very likely to be overloaded when offered traffic load becomes heavy. For example, as indicated by Fig. 9, in the case of relaying with the benchmark algorithm, ratios of relayed cell throughput to total cell throughput decrease dramatically when the session arrival rate increases. Third, in the IRRA, the radio resource scheduling is well coordinated with routing, and therefore the benefits of relaying can be promptly captured and translated into throughput improvements. Finally, the IRRA is run periodically, hence the user transmission rates and routes are able to effectively adapt to system dynamics, e.g., traffic bursts.

It can also be observed that the gain of the proposed IRRA algorithm varies with the RS-to-BS distance r , and the maximum cell throughput gain is about 215% in comparison to the nonrelaying case, obtained when r equals 65% of the cell radius. This is because currently, we assume users can only be relayed by the RSs in the same cell. Consequently, when r becomes bigger, the number of eligible users for relaying will decrease. This is indicated in Fig. 10, where the average relayed cell throughput decreases as r increases. This reveals that the actual performance

of a MCN depends not only on the RRA algorithm characteristics, but also on how many users employ relaying. Therefore, the deployment of RSs should try to match user distributions in order to increase the number of users around RSs, thereby potentially increasing the number of users using relaying.

VI. CONCLUSION

We studied the RRA problem in MCNs. The optimization problem with the objective of system throughput maximization was mathematically formulated and proven to be NP-hard. Considering the prohibitive processing complexity of finding the optimal solution for such an NP-hard problem, we proposed an efficient heuristic algorithm, named IRRA, to provide suboptimal RRA in practical systems. To prove the effectiveness of the IRRA algorithm, a case study was carried out based on enhanced uplink UTRA-FDD with fixed RSs. As shown by the simulation results of the case study, the IRRA can ensure significant gains in terms of cell throughput compared with the nonrelaying case and to a benchmark relaying approach.

ACKNOWLEDGMENT

The authors would like to thank L. Hanzo II for his valuable comments. The authors would also like to express their gratitude to the anonymous reviewers for the insightful comments.

REFERENCES

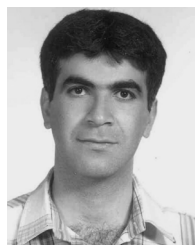
- [1] R. Pabst, B. H. Walke, D. C. Schultz, P. Herhold, H. Yanikomeroglu, S. Mukherjee, H. Viswanathan, M. Lott, W. Zirwas, M. Dohler, H. Aghvami, D. D. Falconer, and G. P. Fettweis, "Relay-based deployment concepts for wireless and mobile broadband radio," *IEEE Commun. Mag.*, vol. 42, no. 9, pp. 80–89, Sep. 2004.
- [2] M. Dohler, A. Gkelias, and H. Aghvami, "Resource allocation for FDMA-based regenerative multihop links," *IEEE Trans. Wireless Commun.*, vol. 3, no. 6, pp. 1989–1993, Nov. 2004.
- [3] T. Rouse, S. McLaughlin, and I. Band, "Congestion-based routing strategies in multihop TDD-CDMA networks," *IEEE J. Sel. Areas Commun.*, vol. 23, no. 3, pp. 668–681, Mar. 2005.
- [4] Z. Dawy, S. Davidovic, and I. Oikonomidis, "Coverage and capacity enhancement of CDMA cellular systems via multihop transmission," in *Proc. IEEE GLOBECOM*, Dec. 2003, vol. 2, pp. 1147–1151.
- [5] A. Fujiwara, S. Takeda, H. Yoshino, and T. Otsu, "Area coverage and capacity enhancement by multihop connection of CDMA cellular network," in *Proc. IEEE Veh. Technol. Conf.-Fall*, Sep. 2002, vol. 4, pp. 2371–2374.
- [6] Y. Yamao, T. Otsu, A. Fujiwara, H. Murata, and S. Yoshida, "Multihop radio access cellular concept for fourth-generation mobile communications system," in *Proc. IEEE PIMRC*, Sep. 2002, vol. 1, pp. 59–63.
- [7] V. Sreng, H. Yanikomeroglu, and D. D. Falconer, "Relayer selection strategies in cellular networks with peer-to-peer relaying," in *Proc. IEEE Veh. Technol. Conf.-Fall*, Oct. 2003, vol. 3, pp. 1949–1953.
- [8] A. A. N. A. Kusuma and L. L. H. Andrew, "Minimum power routing for multihop cellular networks," in *Proc. IEEE GLOBECOM*, Nov. 2002, vol. 1, pp. 37–41.
- [9] H. Viswanathan and S. Mukherjee, "Performance of cellular networks with relays and centralized scheduling," in *Proc. IEEE Veh. Technol. Conf.-Fall*, Oct. 6–9, 2003, vol. 3, pp. 1923–1928.
- [10] H. Nyberg, C. Johansson, and B. Olin, "A streaming video traffic model for the mobile access network," in *Proc. IEEE Veh. Technol. Conf.-Fall*, Oct. 2001, vol. 1, pp. 423–427.
- [11] 3GPP, "Feasibility study of enhanced uplink for UTRA FDD," (TR 25.896) V6.0.0, Mar. 2004. [Online]. Available: www.3gpp.org
- [12] Qualcomm Europe, "Reference node-B scheduler for EUL (3GPP R1-031246)," Nov. 2003. [Online]. Available: www.3gpp.org
- [13] T. S. Rappaport and T. Rappaport, *Wireless Communications Principles and Practice*. Englewood Cliffs, NJ: Prentice-Hall, 2001.
- [14] H. Holma and A. Toskala, *WCDMA for UMTS: Radio Access for Third Generation Mobile Communications*. New York: Wiley, 2000.
- [15] S. Martello and P. Toth, *Knapsack Problems, Algorithm and Computer Implementation*. New York: Wiley, 1990.

- [16] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. New York: Freeman, 1979.
- [17] IETF, "Mobile ad-hoc networks (MANET) charter," [Online]. Available: <http://www.ietf.org/html.charters/manet-charter.html>
- [18] H. S. Wilf, *Algorithms and Complexity*. Englewood Cliffs, NJ: Prentice-Hall, 1986.



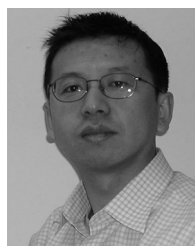
Yajian Liu (S'05) received the B.Eng. and M. Eng. degrees from the School of Electrical and Information Engineering, Hunan University, Hunan, China, in 1998 and 2001, respectively. He is currently working towards the Ph.D. degree at the Center for Communication Systems Research (CCSR), University of Surrey, Surrey, U.K.

From 2001 to 2003, he was a full-time Researcher on mobile communication systems in Shanghai R&D Center, Huawei Technologies Company, Ltd., China. Since 2003, he has been actively involved in IST European projects MUMOR, WINNER, and FIREWORKS. His research interests include radio resource management for multihop cellular networks, cooperative relaying for OFDM/OFDMA systems, and packet scheduling for 3G and beyond systems.



Reza Hoshyar (M'03) received the B.S. degree in communications engineering, and the M.S. and Ph.D. degrees both in mobile communications from Tehran University, Tehran, Iran, in 1991, 1996, and 2001, respectively.

From 2002 to June 2006, he has been a Research Fellow in the Mobile Communications Research Group, Center for Communication Systems Research (CCSR), University of Surrey, Surrey, U.K. During this period, he has been actively involved in IST European projects MUMOR, STRIKE, and WINNER. Currently, he is a Senior Research Fellow at CCSR and participating as the technical work package leader in IST European project FIREWORKS, focusing on advanced digital signal processing, coding and modulation, and scheduling techniques for relaying, and cooperative communications in OFDM/OFDMA systems. He received the top and second rank awards for his B.S. and M.S. degrees from Tehran University.



Xinjie Yang received the B.S. and M.S. degrees both in telecommunications engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 1993 and 1998, respectively, and the Ph.D. degree in mobile communications from the University of Surrey, Surrey, U.K., in 2003.

He is currently a 3G Product Manager at NEC Europe, Ltd., London, U.K. He worked as a Research Fellow at the Center for Communication Systems Research (CCSR), University of Surrey, between 2001 and 2005. He holds several patents and has published more than 20 technical papers. His main research interests are RRM and mobility management for 3G and beyond systems.



Rahim Tafazolli (M'89) is a Professor of Mobile/Personal Communications and Head of Mobile Communications Research at the Center for Communication Systems Research (CCSR), University of Surrey, Surrey, U.K. He is the editor of *Technologies for the Wireless Future* (Vol. 1, 2004 and Vol. 2, 2006). He is nationally and internationally known in the field of mobile communications and acts as External Examiner for British Telecom M.Sc. Course. He has been active in research for over 20 years and has authored and coauthored more than 300 papers in refereed international journals and conferences.

Prof. Tafazolli is a consultant to many mobile companies, has Lectured, Chaired and has been invited as Keynote Speaker to a number of IEE and IEEE workshops and conferences. He has been Technical Advisor to many mobile companies, and the European Union all in the field of mobile/wireless communications. He is the Founder and past Chairman of IEE International Conference on Third-Generation Mobile Communications. He is Chairman of EU Expert Group on Mobile Technology Platform, E-Mobility as well as Chairman of Working Group on Post-IP.