

# Adaptive Opportunistic Routing for Wireless Ad Hoc Networks

Abhijeet A. Bhorkar, Mohammad Naghshvar, *Student Member, IEEE*, Tara Javidi, *Member, IEEE*, and Bhaskar D. Rao, *Fellow, IEEE*

**Abstract**—A distributed adaptive opportunistic routing scheme for multihop wireless ad hoc networks is proposed. The proposed scheme utilizes a reinforcement learning framework to opportunistically route the packets even in the absence of reliable knowledge about channel statistics and network model. This scheme is shown to be optimal with respect to an expected average per-packet reward criterion. The proposed routing scheme jointly addresses the issues of learning and routing in an opportunistic context, where the network structure is characterized by the transmission success probabilities. In particular, this learning framework leads to a stochastic routing scheme that optimally “explores” and “exploits” the opportunities in the network.

**Index Terms**—Opportunistic routing, reward maximization, wireless ad hoc networks.

## I. INTRODUCTION

OPPORTUNISTIC routing for multihop wireless ad hoc networks has seen recent research interest to overcome deficiencies of conventional routing [1]–[6] as applied in wireless setting. Motivated by classical routing solutions in the Internet, conventional routing in ad hoc networks attempts to find a fixed path along which the packets are forwarded [7]. Such fixed-path schemes fail to take advantage of broadcast nature and opportunities provided by the wireless medium and result in unnecessary packet retransmissions. The opportunistic routing decisions, in contrast, are made in an online manner by choosing the next relay based on the actual transmission outcomes as well as a rank ordering of neighboring nodes. Opportunistic routing mitigates the impact of poor wireless links by exploiting the broadcast nature of wireless transmissions and the path diversity.

The authors in [1] and [6] provided a Markov decision theoretic formulation for opportunistic routing. In particular, it is shown that the optimal routing decision at any epoch is to select the next relay node based on a distance-vector summarizing the expected-cost-to-forward from the neighbors to the

destination. This “distance” is shown to be computable in a distributed manner and with low complexity using the probabilistic description of wireless links. The study in [1] and [6] provided a unifying framework for almost all versions of opportunistic routing such as SDF [2], Geographic Random Forwarding (GeRaF) [3], and ExOR [4], where the variations in [2]–[4] are due to the authors’ choices of cost measures to optimize. For instance, an optimal route in the context of ExOR [4] is computed so as to minimize the expected number of transmissions (ETX), while GeRaF [3] uses the smallest geographical distance from the destination as a criterion for selecting the next-hop.

The opportunistic algorithms proposed in [1]–[6] depend on a precise probabilistic model of wireless connections and local topology of the network. In a practical setting, however, these probabilistic models have to be “learned” and “maintained.” In other words, a comprehensive study and evaluation of any opportunistic routing scheme requires an integrated approach to the issue of probability estimation. Authors in [8] provide a sensitivity analysis for the opportunistic routing algorithm given in [6]. However, by and large, the question of learning/estimating channel statistics in conjunction with opportunistic routing remains unexplored.

In this paper, we first investigate the problem of opportunistically routing packets in a wireless multihop network when zero or erroneous knowledge of transmission success probabilities and network topology is available. Using a reinforcement learning framework, we propose a distributed adaptive opportunistic routing algorithm (d-AdaptOR) that minimizes the expected average per-packet cost for routing a packet from a source node to a destination. This is achieved by both sufficiently exploring the network using data packets and exploiting the best routing opportunities.

Our proposed reinforcement learning framework allows for a low-complexity, low-overhead, distributed asynchronous implementation. The significant characteristics of d-AdaptOR are that it is oblivious to the initial knowledge about the network, it is distributed, and it is asynchronous.

The main contribution of this paper is to provide an opportunistic routing algorithm that: 1) assumes no knowledge about the channel statistics and network, but 2) uses a reinforcement learning framework in order to enable the nodes to adapt their routing strategies, and 3) optimally exploits the statistical opportunities and receiver diversity. In doing so, we build on the Markov decision formulation in [6] and an important theorem in Q-learning proved in [9]. There are many learning-based routing solutions (both heuristic or analytically driven) for conventional routing in wireless or wired networks [10]–[15]. None of these

Manuscript received June 27, 2009; revised February 09, 2010; August 25, 2010; and February 28, 2011; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor C. Westphal; accepted May 17, 2011. Date of publication July 18, 2011; date of current version February 15, 2012. This work was supported in part by UC Discovery Grant #com07-10241, Ericsson, Intel Corporation, QUALCOMM, Inc., Texas Instruments, Inc., CWC at UCSD, and NSF CAREER Award CNS-0533035.

The authors are with the Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA 92093 USA (e-mail: abhorkar@ucsd.edu; naghshvar@ucsd.edu; tjavidi@ucsd.edu; brao@ucsd.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNET.2011.2159844

solutions exploits the receiver diversity gain in the context of opportunistic routing. However, for the sake of completeness, we provide a brief overview of the existing approaches. The authors in [10]–[14] focus on heuristic routing algorithms that adaptively identify the least congested path in a wired network. If the network congestion, hence delay, were to be replaced by time-invariant quantities,<sup>1</sup> the heuristics in [10]–[14] would become a special case of d-AdaptOR in a network with deterministic channels and with no receiver diversity. In this light, Theorem 1 in Section IV provides analytic guarantees for the heuristics obtained in [10]–[14]. In [15], analytic results for ant routing are obtained in wired networks without opportunism. Ant routing uses ant-like probes to find paths of optimal costs such as expected hop count, expected delay, and packet loss probability.<sup>2</sup> This dependence on ant-like probing represents a stark difference with our approach where d-AdaptOR relies solely on data packet for exploration.

The rest of the paper is organized as follows. In Section II, we discuss the system model and formulate the problem. Section III formally introduces our proposed adaptive routing algorithm, d-AdaptOR. We then state and prove the optimality theorem for d-AdaptOR in Section IV. In Section V, we present the implementation details and practical issues of d-AdaptOR. We perform simulation study of d-AdaptOR in Section VI. Finally, we conclude the paper and discuss future work in Section VII.

## II. SYSTEM MODEL

We consider the problem of routing packets from a source node 0 to a destination node  $d$  in a wireless ad hoc network of  $d+1$  nodes denoted by the set  $\Theta = \{0, 1, 2, \dots, d\}$ . The time is slotted and indexed by  $n \geq 0$  (this assumption is not technically critical and is only assumed for ease of exposition). A packet indexed by  $m \geq 1$  is generated at the source node 0 at time  $\tau_s^m$  according to an arbitrary distribution with rate  $\lambda > 0$ .

We assume a fixed transmission cost  $c_i > 0$  is incurred upon a transmission from node  $i$ . Transmission cost  $c_i$  can be considered to model the amount of energy used for transmission, the expected time to transmit a given packet, or the hop count when the cost is set to unity. We consider an opportunistic routing setting with no duplicate copies of the packets. In other words, at a given time only one node is responsible for routing any given packet. Given a successful packet transmission from node  $i$  to the set of neighbor nodes  $S$ , the next (possibly randomized) routing decision includes: 1) retransmission by node  $i$ ; 2) relaying the packet by a node  $j \in S$ ; or 3) dropping the packet altogether. If node  $j$  is selected as a relay, then it transmits the packet at the next slot, while other nodes  $k \neq j, k \in S$ , expunge that packet.

We define the termination event for packet  $m$  to be the event that packet  $m$  is either received at the destination or is dropped by a relay before reaching the destination. We denote this termination action by  $T$ . We define termination time  $\tau_T^m$  to be the stopping time when packet  $m$  is terminated. We discriminate

among the termination events as follows. We assume that upon the termination of a packet at the destination (successful delivery of a packet to the destination), a fixed and given positive delivery reward  $R$  is obtained, while no reward is obtained if the packet is terminated before it reaches the destination. Let  $r_m$  denote this random reward obtained at the termination time  $\tau_T^m$ , i.e., either zero if the packet is dropped prior to reaching the destination node or  $R$  if the packet is received at the destination.

Let  $i_{n,m}$  denote the index of the node which at time  $n$  transmits packet  $m$ , and accordingly let  $c_{i_{n,m}}$  denote the cost of transmission (equal to zero if at time  $n$  packet  $m$  is not transmitted). The routing scheme can be viewed as selecting a (random) sequence of nodes  $\{i_{n,m}\}$  for relaying packets  $m = 1, 2, \dots$ .<sup>3</sup> As such, the expected average per-packet reward associated with routing packets along a sequence of  $\{i_{n,m}\}$  up to time  $N$  is

$$J_N = \mathbf{E} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \quad (1)$$

where  $M_N$  denotes the number of packets terminated up to time  $N$  and the expectation is taken over the events of transmission decisions, successful packet receptions, and packet generation times.

**Problem (P):** Choose a sequence of relay nodes  $\{i_{n,m}\}$  in the absence of knowledge about the network topology such that  $J_N$  is maximized as  $N \rightarrow \infty$ .

In Section III, we propose the d-AdaptOR algorithm, which solves Problem (P). The nature of the algorithm allows nodes to make routing decisions in distributed, asynchronous, and adaptive manner.

**Remark 1:** The problem of opportunistic routing for multiple source–destination pairs, without loss of generality, can be decomposed to the single source–destination problem described above [Problem (P) is solved for each distinct flow].

## III. DISTRIBUTED ALGORITHM

Before we proceed with the description of d-AdaptOR, we provide the following notations. Let  $\mathcal{N}(i)$  denote the set of neighbors of node  $i$  including node  $i$  itself. Let  $\mathfrak{S}^i$  denote the set of potential reception outcomes due to a transmission from node  $i \in \Theta$ , i.e.,  $\mathfrak{S}^i = \{S : S \subseteq \mathcal{N}(i), i \in S\}$ . We refer to  $\mathfrak{S}^i$  as the state space for node  $i$ 's transmission. Furthermore, let  $\mathfrak{S} = \cup_{i \in \Theta} \mathfrak{S}^i$ . Let  $A(S) = S \cup \{T\}$  denote the space of all allowable actions available to node  $i$  upon successful reception at nodes in  $S$ . Finally, for each node  $i$ , we define a reward function on states  $S \in \mathfrak{S}^i$  and potential decisions  $a \in A(S)$  as

$$g(S, a) = \begin{cases} -c_a, & \text{if } a \in S \\ R, & \text{if } a = T \text{ and } d \in S \\ 0, & \text{if } a = T \text{ but } d \notin S. \end{cases}$$

### A. Overview of d-AdaptOR

As discussed before, the routing decision at any given time is made based on the reception outcome and involves retransmission, choosing the next relay, or termination. Our proposed

<sup>1</sup>The delay and congestion are highly time-varying quantities.

<sup>2</sup>Here, we note that unlike congestion or instantaneous delay, the expected delay under a stable and stationary routing algorithm is indeed time-invariant, and allow for similar mathematically sound treatment.

<sup>3</sup>Packets are indexed according to the termination order.

TABLE I  
NOTATIONS USED IN THE DESCRIPTION OF THE ALGORITHM

Symbol	Definition
$S_n^i$	Nodes receiving the transmission from node $i$ at time $n$
$a_n^i$	Decision taken by node $i$ at time $n$
$A(S)$	Set of available actions when nodes in $S$ receive a packet
$\mathcal{N}(i)$	Neighbors of node $i$ including node $i$
$g(S, a)$	Reward obtained by taking decision $a$ when set $S$ of nodes receive a packet
$\nu_n(i, S, a)$	Number of times up to time $n$ , nodes $S$ have received a packet from node $i$ and decision $a$ is taken
$N_n(i, S)$	Number of times up to time $n$ , nodes $S$ have received a packet from node $i$
$\Lambda_n(i, S, a)$	Score for node $i$ at time $n$ , when nodes $S$ have received the packet and decision $a$ is taken
$\Lambda_{max}^i$	Estimated best score for node $i$

scheme makes such decisions in a distributed manner via the following three-way handshake between node  $i$  and its neighbors  $\mathcal{N}(i)$ .

- 1) At time  $n$ , node  $i$  transmits a packet.
- 2) The set of nodes  $S_n^i$  who have successfully received the packet from node  $i$ , transmit acknowledgment (ACK) packets to node  $i$ . In addition to the node's identity, the acknowledgment packet of node  $k \in S_n^i$  includes a control message known as *estimated best score* (EBS) and denoted by  $\Lambda_{max}^k$ .
- 3) Node  $i$  announces node  $j \in S_n^i$  as the next transmitter or announces the termination decision  $T$  in a forwarding (FO) packet.

The routing decision of node  $i$  at time  $n$  is based on an adaptive (stored) score vector  $\Lambda_n(i, \cdot, \cdot)$ . The score vector  $\Lambda_n(i, \cdot, \cdot)$  lies in space  $\mathbb{R}^{v_i}$ , where  $v_i = \sum_{S \in \mathcal{S}^i} |A(S)|$ , and is updated by node  $i$  using the EBS messages  $\Lambda_{max}^k$  obtained from neighbors  $k \in S_n^i$ . Furthermore, node  $i$  uses a set of counting variables  $\nu_n(i, S, a)$  and  $N_n(i, S)$  and a sequence of positive scalars  $\{\alpha_n\}_{n=1}^\infty$  to update its score vector at time  $n$ . The counting variable  $\nu_n(i, S, a)$  is equal to the number of times neighbors  $S$  have received (and acknowledged) the packets transmitted from node  $i$  and routing decision  $a \in A(S)$  has been made up to time  $n$ . Similarly,  $N_n(i, S)$  is equal to the number of times the set of nodes  $S$  has received (and acknowledged) packets transmitted from node  $i$  up to time  $n$ . Lastly,  $\{\alpha_n\}_{n=1}^\infty$  is a fixed sequence of numbers available at all nodes.

Table I provides notations used in the description of the algorithm, while Fig. 1 gives an overview of the components of the algorithm. Next, we present further details.

### B. Detailed Description of d-AdaptOR

The operation of d-AdaptOR can be described in terms of initialization and four stages of transmission, reception and acknowledgment, relay, and adaptive computation as shown in Fig. 1. For simplicity of presentation, we assume a sequential timing for each of the stages. We use  $n^+$  to denote some

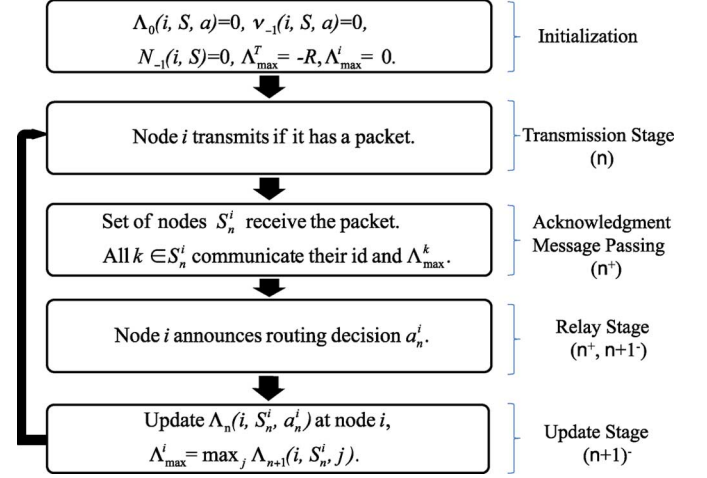


Fig. 1. Flow of the algorithm. The algorithm follows a four-stage procedure: transmission, acknowledgment, relay, and update.

(small) time after the start of  $n$ th slot and  $(n+1)^-$  to denote some (small) time before the end of  $n$ th slot such that  $n < n^+ < (n+1)^- < n+1$ .

#### 0) Initialization:

For all  $i \in \Theta, S \in \mathcal{S}^i, a \in A(S)$ , initialize  $\Lambda_0(i, S, a) = \nu_{-1}(i, S, a) = N_{-1}(i, S) = \Lambda_{max}^i = 0$ , while  $\Lambda_{max}^i = -R$ .

#### 1) Transmission Stage:

Transmission stage occurs at time  $n$  in which node  $i$  transmits if it has a packet.

#### 2) Reception and acknowledgment Stage:

Let  $S_n^i$  denote the (random) set of nodes that have received the packet transmitted by node  $i$ . In the reception and acknowledgment stage, successful reception of the packet transmitted by node  $i$  is acknowledged to it by all the nodes in  $S_n^i$ . We assume that the delay for the acknowledgment stage is small enough (not more than the duration of the time slot) such that node  $i$  infers  $S_n^i$  by time  $n^+$ .

For all nodes  $k \in S_n^i$ , the ACK packet of node  $k$  to node  $i$  includes the EBS message  $\Lambda_{max}^k$ .

Upon reception and acknowledgment, the counting random variable  $N_n$  is incremented as follows:

$$N_n(i, S) = \begin{cases} N_{n-1}(i, S) + 1, & \text{if } S = S_n^i \\ N_{n-1}(i, S), & \text{if } S \neq S_n^i. \end{cases}$$

#### 3) Relay Stage:

Node  $i$  selects a routing action  $a_n^i \in A(S_n^i)$  according to the following (randomized) rule parameterized by  $\epsilon_n(i, S) = \frac{1}{N_n(i, S) + 1}$ .

- With probability  $(1 - \epsilon_n(i, S_n^i))$

$$a_n^i \in \arg \max_{j \in A(S_n^i)} \Lambda_n(i, S_n^i, j)$$

is selected.<sup>4</sup>

<sup>4</sup>In case of ambiguity, node with the smallest index is chosen.

- With probability  $\epsilon_n(i, S_n^i)$

$$a_n^i \in A(S_n^i)$$

is selected uniformly with probability  $\frac{\epsilon_n(i, S_n^i)}{|A(S_n^i)|}$ .

Node  $i$  transmits FO, a control packet that contains information about routing decision  $a_n^i$  at some time strictly between  $n^+$  and  $(n+1)^-$ . If  $a_n^i \neq T$ , then node  $a_n^i$  prepares for forwarding in the next time slot, while nodes  $j \in S_n^i, j \neq a_n^i$  expunge the packet. If termination action is chosen, i.e.,  $a_n^i = T$ , all nodes in  $S_n^i$  expunge the packet.

Upon selection of routing action, the counting variable  $\nu_n$  is updated

$$\nu_n(i, S, a) = \begin{cases} \nu_{n-1}(i, S, a) + 1, & \text{if } (S, a) = (S_n^i, a_n^i) \\ \nu_{n-1}(i, S, a), & \text{if } (S, a) \neq (S_n^i, a_n^i). \end{cases}$$

#### 4) Adaptive Computation Stage:

At time  $(n+1)^-$ , after being done with transmission and relaying, node  $i$  updates score vector  $\Lambda_n(i, \cdot, \cdot)$  as follows.

- For  $S = S_n^i, a = a_n^i$

$$\Lambda_{n+1}(i, S, a) = \Lambda_n(i, S, a) + \alpha_{\nu_n(i, S, a)} \times (-\Lambda_n(i, S, a) + g(S, a) + \Lambda_{\max}^a). \quad (2)$$

- Otherwise

$$\Lambda_{n+1}(i, S, a) = \Lambda_n(i, S, a). \quad (3)$$

Furthermore, node  $i$  updates its EBS message  $\Lambda_{\max}^i$  for future acknowledgments as

$$\Lambda_{\max}^i = \max_{j \in A(S_n^i)} \Lambda_{n+1}(i, S_n^i, j).$$

### C. Computational Issues

The computational complexity and control overhead of d-AdaptOR is low.

1) *Complexity*: To execute stochastic recursion (2), the number of computations required per packet is order of  $O(\max_{i \in \Theta} |\mathcal{N}(i)|)$  at each time slot. The space complexity of d-AdaptOR is exponential in the number of neighbors, i.e.,  $O(\max_{i \in \Theta} 2^{|\mathcal{N}(i)|})$  for each node. The reduction in storage requirement using approximation techniques in [16] is left as future work.

2) *Control Overhead*: The number of acknowledgments per packet is order of  $O(\max_{i \in \Theta} |\mathcal{N}(i)|)$ , independent of network size.

3) *Exploration Overhead*: The adaptation to the optimal performance in the network is guaranteed via a controlled randomized routing strategy that can be viewed as cost of exploration. The cost of exploration is proportional to the total number of packets whose routes deviates from the optimal path. In proof of Theorem 1, we show that this cost increases sublinearly with the number of delivered packets, hence the per-packet exploration cost diminishes as the number of delivered packets grows. Additionally, communication of  $\Lambda_{\max}$  adds a very modest overhead

to the genie-aided or greedy-based schemes such as ExOR or SR.

### IV. ANALYTIC OPTIMALITY OF D-ADAPTOR

We will now state the main result establishing the optimality of the proposed d-AdaptOR algorithm under the assumptions of a time-invariant model of packet reception and reliable control packets. More precisely, we have the following assumptions.

*Assumption 1*: The probability of successful reception of a packet transmitted by node  $i$  at set  $S \subseteq \mathcal{N}(i)$  of nodes is  $P(S|i)$ , independent of time and all other routing decisions.

The probabilities  $P(\cdot)$  in Assumption 1 characterize a packet reception model that we refer to as *local broadcast model*. Note that for all  $S \neq S'$ , successful reception at  $S$  and  $S'$  are mutually exclusive and  $\sum_{S \subseteq \Theta} P(S|i) = 1$ . Furthermore, logically node  $i$  is always a recipient of its own transmission, i.e.,  $P(S|i) = 0$  iff  $i \notin S$ .

*Assumption 2*: The successful reception at set  $S$  due to transmission from node  $i$  is acknowledged perfectly to node  $i$ .

*Remark 2*: Assumption 1 is in line with the experimentally tested state of the art routing protocols MORE [17] and ExOR [4]. These studies seem to indicate that reasonably simple probabilistic models provide good abstractions of media access control (MAC) and physical (PHY) layers at the routing layer.

*Remark 3*: In practice, Assumption 2 is hard to satisfy. But as we will see in Section VI, when the rates and power of the control packets are set to maximize the reliability, the impact of violating this assumption can be kept extremely low.

*Remark 4*: In Section VI, we address the severity as well as the implications of Assumptions 1 and 2. In particular, via a set of QualNet simulations, we will show that d-AdaptOR exhibits many of its desirable properties in a realistic setup despite the relaxation of the analytical assumptions.

Given Assumptions 1 and 2, we are almost ready to present Theorem 1 regarding the optimality of d-AdaptOR among the class of policies that are oblivious to the network topology and/or channel statistics. More precisely, let a *distributed routing policy* be a collection  $\phi = \{\phi^i\}_{i \in \Theta}$  of routing decisions taken at nodes  $i \in \Theta$ , where  $\phi^i$  denotes a sequence of random actions  $\phi^i = \{a_0^i, a_1^i, \dots\}$  for node  $i$ . The policy  $\phi$  is said to be (P)-*admissible* if for all nodes  $i \in \Theta, S \in \mathfrak{S}^i, a \in A(S)$ , the event  $\{a_n^i = a\}$  belongs to the  $\sigma$ -field  $\mathcal{H}_n^i$  generated by the observations at node  $i$ , i.e.,  $\bigcup_{j \in \mathcal{N}(i)} \{S_0^j, a_0^j, \dots, S_{n-1}^j, a_{n-1}^j, S_n^j\}$ . Let  $\Phi$  denote the set of such (P)-admissible policies. Theorem 1 states that d-AdaptOR, denoted by  $\phi^* \in \Phi$ , is an optimal (P)-admissible policy.

*Theorem 1*: Suppose  $\sum_{n=0}^{\infty} \alpha_n = \infty, \sum_{n=0}^{\infty} \alpha_n^2 < \infty$ , and Assumptions 1 and 2 hold. Then, for all  $\phi \in \Phi$

$$\begin{aligned} \lim_{N \rightarrow \infty} \mathbf{E}^{\phi^*} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_t^m-1} c_{i_n, m} \right\} \right] \\ \geq \limsup_{N \rightarrow \infty} \mathbf{E}^{\phi} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_t^m-1} c_{i_n, m} \right\} \right] \end{aligned}$$

where  $\mathbf{E}^{\phi^*}$  and  $\mathbf{E}^{\phi}$  are the expectations taken with respect to policies  $\phi^*$  and  $\phi$ , respectively.<sup>5</sup>

Next, we prove the optimality of d-AdaptOR in two steps. In the first step, we show that  $\Lambda_n$  converges in an almost sure sense. In the second step, we use this convergence result to show that d-AdaptOR is optimal for Problem (P).

#### A. Convergence of $\Lambda_n$

Let  $U : \prod_i \mathbb{R}^{v_i} \rightarrow \prod_i \mathbb{R}^{v_i}$  be an operator on vector  $\Lambda$  such that

$$(U\Lambda)(i, S, a) = g(S, a) + \sum_{S' \in \mathcal{S}^a} P(S'|a) \max_{j \in A(S')} \Lambda(a, S', j). \quad (4)$$

Let  $\Lambda^* \in \prod_i \mathbb{R}^{v_i}$  denote the fixed point of operator  $U$ ,<sup>6</sup> i.e.,

$$\Lambda^*(i, S, a) = g(S, a) + \sum_{S' \in \mathcal{S}^a} P(S'|a) \max_{j \in A(S')} \Lambda^*(a, S', j). \quad (5)$$

The following lemma establishes the convergence of recursion (2) to the fixed point of  $U$ ,  $\Lambda^*$ .

*Lemma 1:* Let:

- J1)  $\Lambda_0(\cdot, \cdot, \cdot) = 0$ ,  $\Lambda_{\max}^T = -R$ ,  $\Lambda_{\max}^i = 0$  for all  $i \in \Theta$ ;
- J2)  $\sum_{n=0}^{\infty} \alpha_n = \infty$ ,  $\sum_{n=0}^{\infty} \alpha_n^2 < \infty$ .

Then, the sequence  $\Lambda_n$  obtained by the stochastic recursion (2) converges to  $\Lambda^*$  almost surely.

The proof uses known results on the convergence of a certain recursive stochastic process as presented by Fact 2 in Appendix-A.

#### B. Proof of Optimality

Using the convergence of  $\Lambda_n$ , we show that the expected average per-packet reward under d-AdaptOR is equal to the optimal expected average per-packet reward obtained for a genie-aided system where the local broadcast model is known perfectly. In other words, we take cue from known results associated with a closely related Auxiliary Problem (AP). In this Auxiliary Problem (AP), there exists a centralized controller with full knowledge of the local broadcast model  $P(\cdot|\cdot)$  as well as the transmission outcomes across the network [1], [6]. The objective in the Auxiliary Problem (AP) is a single-packet variation of that in Problem (P): the reward

$$\mathbf{E} \left[ r_m - \sum_{n=0}^{\tau_{T^m}^m - 1} c_{i_{n,m}} \right]$$

for routing a single packet  $m$  from the source to the destination is maximized over a set  $\Pi$  of (AP)-admissible policies, where this set  $\Pi$  of (AP)-admissible policies is a superset of (P)-admissible policies  $\Phi$  that also includes all topology-aware and

centralized policies. This Auxiliary Problem (AP) has been extensively studied in [1], [6], and [19], where a Markov decision formulation provides the following important result.

*Fact 1 [6, Theorem 2.1]:* Consider the unique solution  $V^* : \Theta \cup \{T\} \rightarrow \mathbb{R}^+$  to the following fixed-point equation:

$$\begin{aligned} V^*(d) &= R \\ V^*(i) &= \max \left( \left\{ -c_i + \sum_{S'} P(S'|i) \left( \max_{j \in S'} V^*(j) \right) \right\}, 0 \right) \end{aligned} \quad (6)$$

There exists an optimal topology-aware and centralized admissible policy  $\pi^* \in \Pi$  such that

$$\mathbf{E}^{\pi^*} \left[ r_m - \sum_{n=0}^{\tau_{T^m}^m - 1} c_{i_{n,m}} \right] = V^*(0). \quad (8)$$

Lemma 2 states the relationship between the solution of Problem (P) and that of the Auxiliary Problem (AP). More specifically, Lemma 2 shows that  $V^*(0)$  is an upper bound for the solution to Problem (P).

*Lemma 2:* For any (P)-admissible policy  $\phi \in \Phi$  for Problem (P) and for all  $N = 1, 2, \dots$

$$\mathbf{E}^{\phi} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_{T^m}^m - 1} c_{i_{n,m}} \right\} \right] \leq V^*(0).$$

The proof is given in Appendix-B. Intuitively, the result holds because the set of (P)-admissible policies is a subset of (AP)-admissible policies, i.e.,  $\Phi \subset \Pi$ .

Lemma 3 gives the achievability proof by showing that the expected average per-packet reward of d-AdaptOR is lower-bounded by  $V^*(0)$ .

*Lemma 3:* For any  $\delta > 0$

$$\liminf_{N \rightarrow \infty} \mathbf{E}^{\phi^*} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_{T^m}^m - 1} c_{i_{n,m}} \right\} \right] \geq V^*(0) - \delta.$$

The proof is given in Appendix-C. Lemmas 2 and 3 imply that  $\phi^*$  [which is (P)-admissible by construction] is an optimal policy under which

$$\lim_{N \rightarrow \infty} \mathbf{E}^{\phi^*} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_{T^m}^m - 1} c_{i_{n,m}} \right\} \right]$$

exists and is equal to  $V^*(0)$  establishing the proof of Theorem 1.

*Corollary 1:* When  $c_i = 1$ ,  $i \in \Theta$ , the network is connected, and  $R$  is greater than the worst-case routing cost,<sup>7</sup> d-AdaptOR minimizes

$$D_N = \mathbf{E}^{\pi} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \{ \tau_{T^m}^m - \tau_s^m \} \right] \quad (9)$$

the expected per-packet delivery time as  $N \rightarrow \infty$ .

<sup>5</sup>This is a strong notion of optimality and implies that the proposed algorithm's expected average reward is greater than the best-case performance ( $\limsup$ ) of all policies [18, p. 344].

<sup>6</sup>Existence and uniqueness of  $\Lambda^*$  is provided in Appendix-A.

<sup>7</sup>The worst-case routing cost can be determined by taking supremum over ETX metrics for all source-destination pairs.

This is because when  $c_i = 1$ ,  $R$  is sufficiently large, and the network is connected

$$V^*(0) = R - \inf_{\pi \in \Pi} \mathbf{E}^\pi \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_n, m} \right\} \right] \\ = R - \inf_{\pi \in \Pi} D_N, \text{ as } N \rightarrow \infty.$$

## V. PROTOCOL DESIGN AND IMPLEMENTATION ISSUES

In this section, we describe an 802.11 compatible implementation for d-AdaptOR.

### A. 802.11 Compatible Implementation

The implementation of d-AdaptOR, analogous to any opportunistic routing scheme, involves the selection of a relay node among the candidate set of nodes that have received and acknowledged a packet successfully. One of the major challenges in the implementation of an opportunistic routing algorithm in general, and the d-AdaptOR algorithm in particular, is the design of an 802.11 compatible acknowledgment mechanism at the MAC layer. We propose a practical and simple way to implement acknowledgment architecture.

The transmission at any node  $i$  is done according to an 802.11 CSMA/CA mechanism. Specially, before any transmission, transmitter  $i$  performs channel sensing and starts transmission after the backoff counter is decremented to zero. For each neighbor node  $j \in \mathcal{N}(i)$ , the transmitter node  $i$  then reserves a virtual time slot of duration  $T_{\text{ACK}} + T_{\text{SIFS}}$ , where  $T_{\text{ACK}}$  is the duration of the acknowledgment packet and  $T_{\text{SIFS}}$  is the duration of Short InterFrame Space (SIFS) [20]. Transmitter  $i$  then piggybacks a priority ordering of nodes  $\mathcal{N}(i)$  with each data packet transmitted. The priority ordering determines the virtual time slot in which the candidate nodes transmit their acknowledgment. Nodes in the set  $S^i$  that have successfully received the packet then transmit acknowledgment packets sequentially in the order determined by the transmitter node.

After a waiting time of  $T_{\text{wait}} = |\mathcal{N}(i)|(T_{\text{ACK}} + T_{\text{SIFS}})$  during which each node in the set  $S^i$  has had a chance to send an ACK, node  $i$  transmits a FOwarding control packet (FO). The FO packets contain the identity of the next forwarder, which may be node  $i$  again or any node  $j \in S^i$ . If  $T_{\text{wait}}$  expires and no FO packet is received (FO packet reception is unsuccessful), then the corresponding candidate nodes drop the received data packet. If the transmitter  $i$  does not receive any acknowledgment, node  $i$  retransmits the packet. The backoff window is doubled after every retransmission. Furthermore, the packet is dropped if the retry limit (set to 7) is reached.

In addition to the acknowledgment scheme, d-AdaptOR requires modifications to the 802.11 MAC frame format. Fig. 2 shows the modified MAC frame formats required by d-AdaptOR. The reserved bits in the type/subtype fields of the frame control field of the 802.11 MAC specification are used to indicate whether the rest of the frame is a d-AdaptOR data frame, a d-AdaptOR ACK, or a FO.<sup>8</sup> The data frame contains

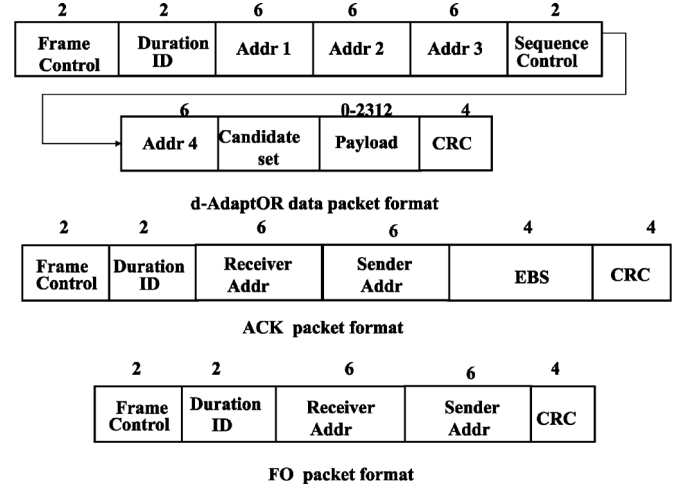


Fig. 2. Frame structure of the data packets, acknowledgment packets, and FO packets.

the candidate set in priority order, the payload, and the 802.11 Frame Check Sequence. The acknowledgment frame includes the data frame sender's address and the feedback EBS  $\Lambda_{\text{max}}$ . The FO packet is exactly the same as a standard 802.11 short control frame that uses different subtype value.

### B. d-AdaptOR in a Realistic Setting

1) *Loss of ACK and FO Packets:* Interference or low signal-to-noise ratio (SNR) can cause loss of ACK and FO packets. Loss of an ACK packet results in an incorrect estimation of nodes that have received the packet, and thus affects the performance of the algorithm. Loss of FO packet negatively impacts the throughput performance of the network. In particular, loss of an FO packet can result in the drop of data packets at all the potential relays, reducing the throughput performance. Hence, in our design, FO packets are transmitted at lower rates to ensure a reliable transmission.

2) *Increased Overhead:* As it is the case with any opportunistic scheme, d-AdaptOR adds a modest additional overhead to the standard 802.11 due to the added acknowledgment/handshake structure. This overhead increases linearly with the number of neighbors. Assuming a 802.11b physical layer operating at 11 Mb/s with an SIFS time of 10  $\mu\text{s}$ , preamble duration of 20  $\mu\text{s}$ , Physical Layer Convergence Protocol (PLCP) header duration of 4  $\mu\text{s}$ , and 512-B frame payloads, Table II compares the overhead in the data packet due to piggybacking and the control overhead due to ACK and FO packets for unicast 802.11, genie-aided opportunistic scheme, and d-AdaptOR. d-AdaptOR requires communication overhead of 4 extra bytes (for EBS) per ACK packet compared to the genie-aided opportunistic scheme, while unicast 802.11 does not require such overhead.

Note that the overhead cost can be reduced by restricting the number of nodes in the candidate list of MAC header to a given number, MAX-NEIGHBOUR. The unique ordering for the nodes in the candidate set is determined by prioritizing the nodes with respect to  $\Lambda_n(i, \{i, j\}, j), j \in \mathcal{N}(i)$  and then

<sup>8</sup>This enables the d-AdaptOR to communicate and be fully compatible with other 802.11 devices.

TABLE II  
OVERHEAD COMPARISONS

	Data Frame	Control packets	Total
802.11	397 $\mu$ s	40 $\mu$ s (ACK)	437 $\mu$ s
Genie-aided opportunistic scheme	400 $\mu$ s	115 $\mu$ s + 40 $\mu$ s (ACK+FO)	555 $\mu$ s
d-AdaptOR	400 $\mu$ s	124 $\mu$ s + 40 $\mu$ s (ACK+FO)	564 $\mu$ s

choosing the MAX-NEIGHBOUR highest priority nodes.<sup>9</sup> Such a limitation will sacrifice the diversity gain and, hence, the performance of any opportunistic routing algorithm for lower overhead. In practice, we have seen that limiting the neighbor set to 4 provides most of the diversity gain.

## VI. SIMULATIONS

In this section, we provide simulation studies in realistic wireless settings where the theoretical assumptions of our study do not hold. These simulations not only demonstrate a robust performance gain under d-AdaptOR in a realistic network, but also provide significant insight in the appropriate choice of the design parameters such as damping sequence  $\{\alpha_n\}$ , delivery reward  $R$ , etc. We first investigate the performance of d-AdaptOR with respect to the design parameters and network parameters in a grid topology of 16 nodes. We then use a realistic topology of 36 nodes with random placement to demonstrate robustness of d-Adaptor to the violation of the analytic Assumptions 1 and 2.

### A. Simulation Setup

In Sections VI-B and VI-C, using the appropriate choice of the design parameters, we compare the performance of d-AdaptOR against suitably chosen candidates. As a benchmark, when appropriate, we have compared the performance against a genie-aided policy that relies on full network topology information when selecting routes. This is nothing but  $\pi^*$  discussed in Section IV-B. We also compare against Stochastic Routing (SR) [1] (SR is the distributed implementation of policy  $\pi^*$ ) and ExOR [4] (an opportunistic routing policy with ETX metric) in which the empirical probabilistic structure of the network is used to implement opportunistic routing algorithms. As a result, their performance will be highly dependent on the precision of empirical probability associated with link  $p_{ij}$ . To provide a fair comparison, we have considered simple greedy versions of SR and ExOR. These algorithms adapt  $\{p_{ij}\}$  to the history of packet reception outcomes and rely on the updates to make routing decisions assuming error-free  $\{p_{ij}\}$ . We have also compared our performance against a conventional routing SRCR [21] with full knowledge of topology. In this setting, a conventional route is selected with perfect knowledge of link success probability at any given node. This comparison in effect provides a simple benchmark for all learning-based conventional routing policies in the literature such as Q-routing [10] and predictive Q-routing [12] when congestion is taken to be small enough (such that finding least congested paths coincides with finding the path with minimum expected number of transmissions).

Our simulations are performed in QualNet. We consider two sets of topologies in our experimental study.

- 1) Grid Topology: In Section VI-B, we study a grid topology consisting of 16 indoor nodes such that the nearest neighbors are separated by distance  $L$  meters. If unspecified,  $L$  is chosen to be 25 m. The source and the destination are chosen at the maximal distance (on diagonal) from each other.
- 2) Random Topology: In Section VI-C, we study a random topology consisting of 36 indoor nodes placed in an area of  $150 \times 150$  m<sup>2</sup>. Here, we investigate the performance under a multisource multidestination setting as the number of flows in the network is varied and each flow is specified via a randomly selected pair of source and destination.

The nodes are equipped with 802.11b radios placed in indoor environment transmitting at 11 Mb/s with transmission power 15 dBm. Note that the choice of indoor environment is motivated by the findings in [22], where opportunistic routing is found to provide significant diversity gains. The wireless medium model includes Rician fading with K-factor of 4 and log-normal shadowing with mean 4 dB. The path loss follows the two-ray model in [23] with path exponent of 3. The acknowledgment packets are short packets of length 24 B transmitted at 11 Mb/s, while FO packets are of length 20 B and transmitted at a lower rate of 1 Mb/s to ensure reliability. If unspecified, packets are generated according to a constant bit rate (CBR) source with rate 20 packets/s. The packets are assumed to be of length 512 B equipped with simple cyclic redundancy check (CRC) error detection. The cost of transmission is assumed to be one unit, and the reward  $R$  is set to 40. We have chosen  $\alpha_n^1 = \frac{1}{\sqrt{n} \log n}$  as the exploration parameter of choice.

### B. Effects of Design and Network Parameters

Here, we investigate the role and criticality of various design parameters of d-AdaptOR with respect to the expected number of transmission criterion. Let us start with design parameters  $\{\alpha_n\}$  and  $R$

1) *Exploration Parameter Sequence  $\{\alpha_n\}$* : The convergence rate of stochastic recursion (2) depends strongly on the choice of sequence  $\{\alpha_n\}$ . Convergence is slower with a faster decreasing sequence  $\{\alpha_n\}$  and results in less variance in the estimates of  $\Lambda_n$ , while with a slow decreasing sequence of  $\{\alpha_n\}$ , convergence is fast but results in large variance in the estimates of  $\Lambda_n$ . In Fig. 3, we have plotted the effect of the choice of  $\{\alpha_n\}$  sequence by comparing two sequences  $\{\alpha_n^1 = \frac{1}{\sqrt{n} \log n}\}$  and  $\{\alpha_n^2 = \frac{1}{n \log n}\}$ . Note that under sequence  $\{\alpha_n^2 = \frac{1}{n \log n}\}$ , d-AdaptOR is slower to adapt to the optimal performance while it shows a slightly smaller variance. This is because the choice of  $\{\alpha_n\}$  controls the rate with which greedy versus (randomly chosen) exploration actions are utilized. The optimization of the choice of  $\{\alpha_n\}$  is an interesting topic of study in stochastic approximation [24], [25], far beyond the scope of this work.

2) *Per-Packet Delivery Reward  $R$* : To ensure an acceptable performance of d-AdaptOR, the value of delivery reward,  $R$ , must be chosen sufficiently high. This would ensure the existence of routes under which the value of delivering a packet (as represented in  $R$ ) is worth (i.e., larger than) the cost of

<sup>9</sup>In case of ambiguity, the node with the smallest index is chosen.



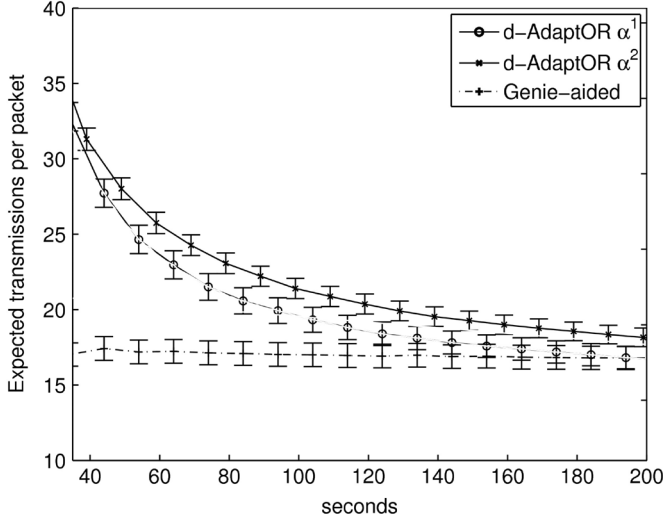


Fig. 3. Comparison for  $\alpha_n^1 = \frac{1}{\sqrt{n \log(n)}}$ ,  $\alpha_n^2 = \frac{1}{n \log(n)}$ .

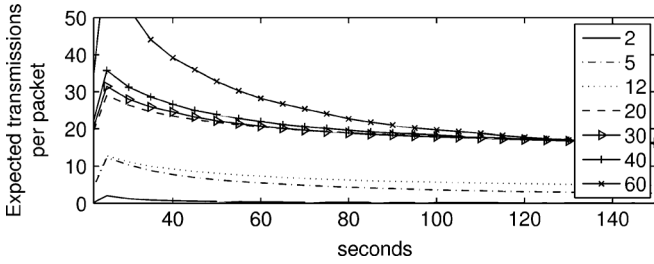


Fig. 4. Expected number of transmissions versus time as  $R$  is varied.

relaying and routing that packet. A reasonable choice of  $R$  is any value larger than the worst-case expected transmission cost. Increasing  $R$  beyond such a value does not affect the asymptotic optimality of the algorithm. Next, we study the performance of d-AdaptOR with respect to the convergence rate and delivery ratio.

Fig. 4 plots the expected number of transmissions rate as time progresses for various values of  $R$ . As seen in Fig. 4, if  $R$  increases beyond a threshold  $R_0$  (in the example provided here, this threshold is 18, but in general it depends on the network diameter), the expected number of transmissions per packet achieve the optimal value of  $R_0$ . In contrast, for  $R < R_0$ , the expected number of transmissions approaches zero as the packets not worth obtaining routing reward are dropped.<sup>10</sup> Fig. 4 also shows that the convergence rate of the expected number of transmissions for routing per packet under d-AdaptOR decreases as  $R$  increases. The slow convergence for  $R > R_0$  for large  $R$  is due to the flexibility of exploring longer paths. The slow convergence to zero for  $R < R_0$  near  $R_0$  is attributed to the fact that it takes a longer time for d-AdaptOR to realize that the packet is not worth relaying.

Fig. 5 plots the delivery ratio as  $R$  is varied. Fig. 5 shows that as  $R$  increases beyond a threshold  $R_0$ , the delivery ratio remains fixed. However, for sufficiently small  $R$ , nearly all the packets are dropped as the cost of transmission of the packet as well as relaying is not worth the obtained delivery reward. Due to very

<sup>10</sup>For  $R < R_0$ , we have plotted negative of the expected per-packet reward as the expected number of transmissions.

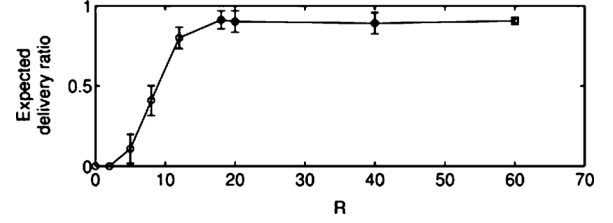


Fig. 5. Delivery ratio as  $R$  is varied.

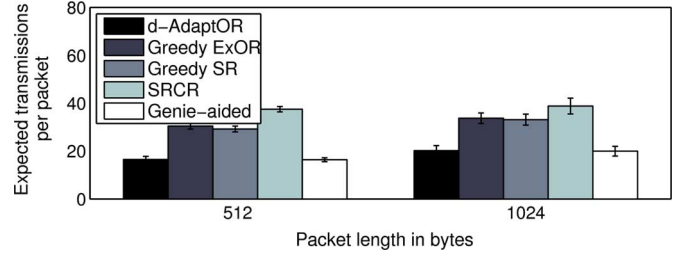


Fig. 6. d-AdaptOR performance as packet length is varied.

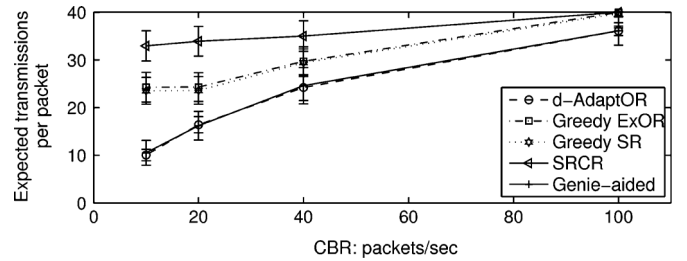


Fig. 7. Performance of d-AdaptOR as CBR traffic is varied.

slow convergence rate around  $R_0$  for  $R < R_0$ , we observe that a nonnegligible number of packets is delivered in the duration of experiment.

Next, we investigate the performance of d-AdaptOR with respect to other candidate protocols for the network parameters such as packet length, traffic rate, neighbor distance, and time-varying costs.

3) *Packet Length*: We have repeated our simulations for 1024-B packets. Fig. 6 plots the performance as the packet length is varied from 512 to 1024 B. Note that due to the decreasing packet transmission reliabilities, the expected routing cost per packet is increased with the packet size. However, the optimality of d-AdaptOR does not depend on the packet length.

4) *Traffic Rate*: Fig. 7 plots the mean number of transmissions versus CBR rate for candidate algorithms. Even though the performance gain for d-AdaptOR decreases somewhat with increase in the load, there is always a nonnegligible advantage over greedy solutions.

5) *Average Hop Length  $L$* : In an attempt to understand the performance gap between various opportunistic algorithms, specifically the gap between d-AdaptOR versus learning-based conventional routing algorithms [10]–[13] whose performance is bounded by SRCR, one needs to gain insight about the diversity gain achieved by opportunistic routing. Fig. 8 compares the expected transmission cost for the three opportunistic routing algorithms (d-AdaptOR, ExOR, and SR) and SRCR as the



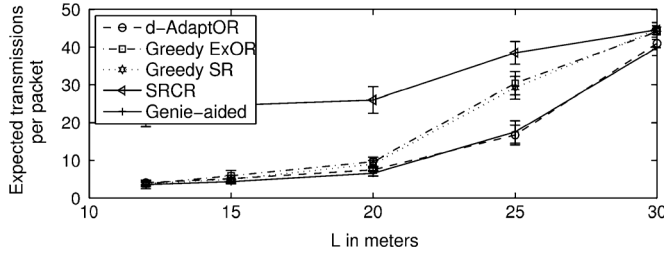


Fig. 8. Small hops provide significant receiver diversity gain.

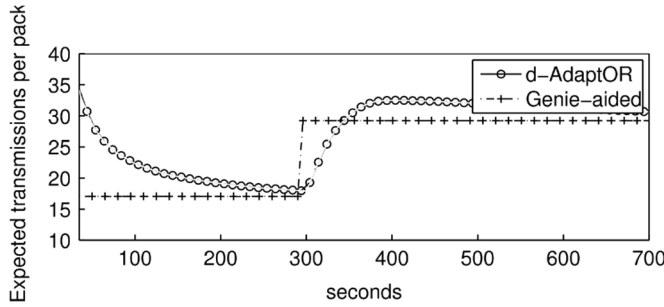


Fig. 9. Time-varying cost: Nodes go into sleep mode at time 300 s.

distance between the neighboring nodes in the grid topology, measured in  $L$  meters, is varied from 10 to 30 m. Note that for high values of  $L$ , the receiver diversity is low due to retransmission packet losses giving nearly similar performance for candidate protocols, while small  $L$  corresponds to a network with large receiver diversity gain. As expected, when  $L$  is small, all opportunistic routing schemes provide a significant improvement over conventional routing, but perhaps what is more interesting is the performance gain of learning-based d-AdaptOR over the greedy-based solutions in medium ranges.

6) *Time-Varying Cost*: In our analytical setup, we assume the transmission costs are fixed. Next, we discuss a simple scenario where the nodes have time-varying transmission costs. Consider a network in which nodes may go into an energy-saving mode when they do not participate in routing (e.g., to recharge their energy sources). Assume that upon entering the energy-saving mode, a node announces a high cost of transmission (100 instead of usual transmission cost of 1). Fig. 9 plots the expected average cost of d-AdaptOR when two nodes at the center of the grid move into an energy-saving mode. It shows that d-AdaptOR can track the genie-aided solution after the nodes move into the energy-saving mode.

### C. Case Study: Random Network

Here, we study a random network scenario consisting of 36 wireless nodes placed randomly, with the remaining parameters kept the same as the default parameters.

Fig. 10 plots the expected number of transmissions and the expected average per-packet reward for the candidate routing algorithms versus network operation time when a single flow is present in the random topology. We first note that, as expected, SRCR performs poorly compared to the opportunistic schemes as it fails to utilize the receiver diversity gain. This underlines our contribution over all existing learning-based solutions [10]–[13] that ignore receiver diversity. Furthermore,

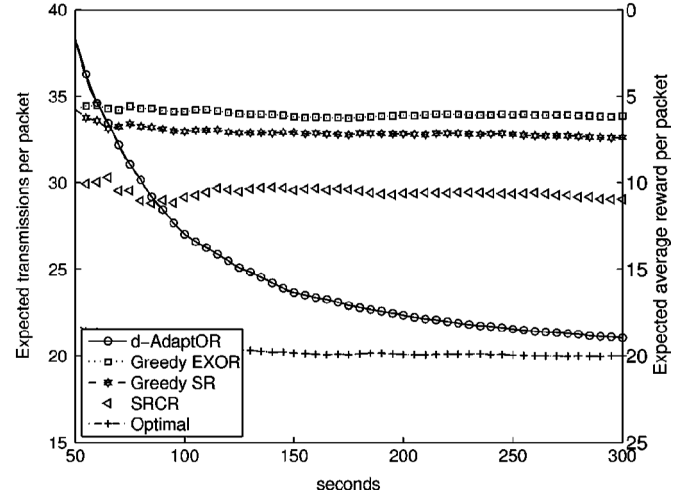


Fig. 10. Expected number of transmissions and average per-packet reward as function of operation time.

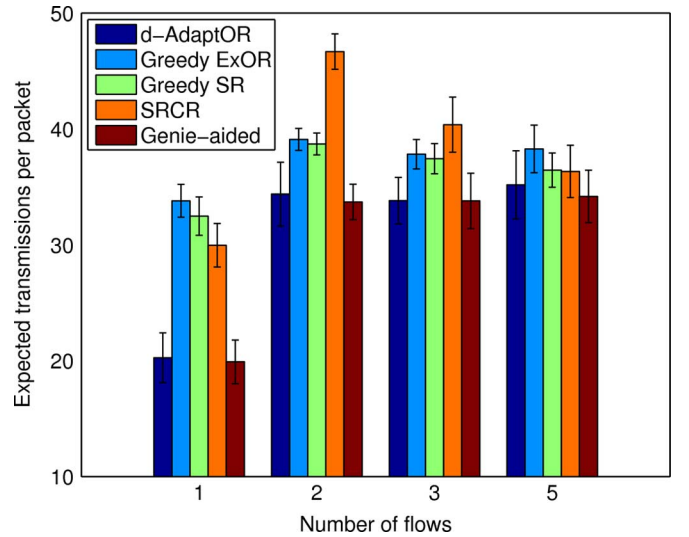


Fig. 11. d-AdaptOR versus distributed SR, ExOR, and SRCR performance for multiple flows.

Fig. 10 shows that the d-AdaptOR algorithm outperforms the greedy opportunistic schemes given sufficient number of packet deliveries. This is because the greedy versions of SR and ExOR fail to explore possible choices of routes and often result in strictly suboptimal routing policies. Fig. 10 also shows that the randomized routing decisions employed by d-AdaptOR work as a double-edged sword. On the one hand, they form a mechanism through which network opportunities are exhaustively explored until the globally optimal decisions are constructed, resulting in an improved long-term performance while these randomized decisions lead to a short-term performance loss. This, in fact, is reminiscent of the well-known exploration/exploitation tradeoff in stochastic control and learning literature.

Next, we study the performance of d-AdaptOR as the number of flows in the network is varied, where each flow is specified via a randomly selected pair of source and destination. Fig. 11 plots the expected number of transmissions and expected average reward for the candidate routing algorithms for the random topology. As seen in Fig. 11, d-AdaptOR maintains an

optimal performance. However, Fig. 11 also shows that the gap between d-AdaptOR and the greedy version of SR significantly decreases with an increase in number of flows where the natural pattern of traffic flow renders the (randomized) exploration phase less critical. In other words, while Fig. 11 is consistent with the Remark 1 in Section II regarding the decomposition of multiple-flow scenario to multiple single-flow scenarios, it also suggests that a joint design in which the multiplicity of flows provide a natural (and greedy) exploration of the network might be beneficial with regard to the transient/short-term performance measures of interest.

## VII. CONCLUSION AND FUTURE WORK

In this paper, we proposed d-AdaptOR, a distributed, adaptive, and opportunistic routing algorithm whose performance is shown to be optimal with zero knowledge regarding network topology and channel statistics. More precisely, under idealized assumptions, d-AdaptOR is shown to achieve the performance of an optimal routing with perfect and centralized knowledge about network topology, where the performance is measured in terms of the expected per-packet reward. Furthermore, we show that d-AdaptOR allows for a practical distributed and asynchronous 802.11 compatible implementation, whose performance was investigated via a detailed set of QualNet simulations under practical and realistic networks. Simulations show that d-AdaptOR consistently outperforms existing adaptive routing algorithms in practical settings.

The long-term average reward criterion investigated in this paper inherently ignores the short-term performance. To capture the performance of various adaptive schemes, however, it is desirable to study the performance of the algorithms over a finite horizon. One popular way to study this is via measuring the incurred “regret” over a finite horizon. Regret is a function of horizon  $N$  that quantifies the loss of the performance under a given adaptive algorithm relative to the performance of the topology-aware optimal one. More specifically, our results so far implies that the optimal rate of growth of regret is strictly sublinear in  $N$ , but fails to provide a conclusive understanding of the short-term behavior of d-AdaptOR. An important area of future work comprises developing adaptive algorithms that ensure optimal growth rate of regret.

The design of routing protocols requires a consideration of congestion control along with the throughput performance [26], [27]. Our work, however, does not consider this closely related issue. Incorporating congestion control in opportunistic routing algorithms to minimize expected delay without the topology and the channel statistics knowledge is an area of future research.

## APPENDIX

We start this section with a note on the notations used. On the probability space  $(\Omega, \mathcal{F}, P)$ , we use notation  $I : \Omega \rightarrow \{0, 1\}$  to denote the indicator random variable (with respect to  $\mathcal{F}$ ), such that for all  $\omega \in \Omega$ ,  $A \in \mathcal{F}$ ,  $I(A) = 1$  for all  $\omega \in A$ , and  $I(A) = 0$  for all  $\omega \notin A$ . For a vector  $x \in \mathbb{R}^D$ ,  $D \geq 1$ , we use  $x(l)$  to denote the  $l$ th element of the vector. Let  $\|x\|_v$  denote the weighted max-norm with positive weight vector  $v$ , i.e.,  $\|x\|_v = \max_l \frac{|x(l)|}{v(l)}$ . We denote the vector in  $\mathbb{R}^D$  with all

components equal to 1 by  $\mathbf{1}$ . We also use the notation  $X^n$  to represent the first  $n$  random elements of the random sequence  $\{X_k\}_{k=1}^\infty$ .

### A. Proof of Lemma 1

*Lemma 1:* Let:

J1)  $\Lambda_0(\cdot, \cdot, \cdot) = 0$ ,  $\Lambda_{\max}^T = -R$ ,  $\Lambda_{\max}^i = 0$  for all  $i \in \Theta$ ;

J2)  $\sum_{n=0}^\infty \alpha_n = \infty$ ,  $\sum_{n=0}^\infty \alpha_n^2 < \infty$ .

Then, the sequence  $\Lambda_n$  obtained by the stochastic recursion (2)

$$\begin{aligned} \Lambda_{n+1}(i, S, a) &= \Lambda_n(i, S, a) \\ &\quad + \alpha_{\nu_n(i, S, a)} (-\Lambda_n(i, S, a) + g(S, a) + \Lambda_{\max}^a) \end{aligned}$$

converges to  $\Lambda^*$  almost surely.

To prove Lemma 1, we note that the adaptive computation given by (2) utilizes a stochastic approximation algorithm to solve the MDP associated with Problem (AP). To study the convergence properties of this stochastic approximation, we appeal to known results in the intersection of learning and stochastic approximation given below.

In particular, consider a set of stochastic sequences on  $\mathbb{R}^D$ , denoted by  $\{x_n, \bar{\alpha}_n, \mathcal{M}_{n+1}\}$ , and the corresponding filtration  $\mathcal{G}_n$ , i.e., the increasing  $\sigma$ -field generated by  $\{x_n, \bar{\alpha}_n, \mathcal{M}_{n+1}\}$ , satisfying the following recursive equation:

$$x_{n+1} = x_n + \bar{\alpha}_n [U(x'_n) - x_n + \mathcal{M}_{n+1}]$$

where  $U$  is a mapping from  $\mathbb{R}^D$  into  $\mathbb{R}^D$  and  $x'_n = (x_{n_1}(1), x_{n_2}(2), \dots, x_{n_D}(D))$ ,  $0 \leq n_j \leq n$ ,  $j \in \{1, 2, \dots, D\}$ , is a vector of possibly delayed components of  $x_n$ . If no information is outdated, then  $n_j = n$  for all  $j$  and  $x'_n = x_n$ . The following important result on the convergence of  $x_n$  is provided in [9].

*Fact 2 [9, Theorem 2]:* Assume  $\{x_n, \bar{\alpha}_n, \mathcal{M}_{n+1}\}$  and  $U$  satisfy the following conditions.

- G1) For all  $n \geq 0$  and  $1 \leq l \leq D$ ,  $0 \leq \bar{\alpha}_n(l) \leq 1$  a.s.;  
for  $1 \leq l \leq D$ ,  $\sum_{n=0}^\infty \bar{\alpha}_n(l) = \infty$  a.s.;  
for  $1 \leq l \leq D$ ,  $\sum_{n=0}^\infty \bar{\alpha}_n^2(l) < \infty$  a.s.
- G2)  $\mathcal{M}_n$  is a martingale difference with finite second moment, i.e.,  $\mathbf{E}\{\mathcal{M}_{n+1}|\mathcal{G}_n\} = 0$ , and there exist constants  $A$  and  $B$  such that  $\mathbf{E}\{\mathcal{M}_{n+1}^2|\mathcal{G}_n\} \leq A + B(\max_{n' \leq n} \|x_{n'}\|)^2$ .
- G3) There exists a positive vector  $v$ , scalars  $\beta \in [0, 1)$  and  $C \in \mathbb{R}^+$ , such that

$$\|U(x)\|_v \leq \beta \|x\|_v + C.$$

- G4) Mapping  $U : \mathbb{R}^D \rightarrow \mathbb{R}^D$  satisfies the following properties.

- 1)  $U$  is componentwise monotonically increasing.
- 2)  $U$  is continuous.
- 3)  $U$  has a unique fixed point  $x^* \in \mathbb{R}^D$ .
- 4)  $U(x) - r\mathbf{1} \leq U(x - r\mathbf{1}) \leq U(x + r\mathbf{1}) \leq U(x) + r\mathbf{1}$ , for any  $r \in \mathbb{R}^+$ .

- G5) For any  $j$ ,  $n_j \rightarrow \infty$  as  $n \rightarrow \infty$ .

Then, the sequence of random vectors  $x_n$  converges to the fixed point  $x^*$  almost surely.

Let  $\mathcal{G}_n$  be the increasing  $\sigma$ -field generated by random vectors  $(\Lambda_n, S_n^i, a_n^i, \nu_n)$ . Let  $x_n = \Lambda_n$  be the random vector of

dimension  $D = \sum_{i \in \Theta} \sum_{S \in \mathfrak{S}^i} |A(S)|$ , generated via recursive equation (2). Furthermore

$$(U\Lambda_n)(i, S, a) = g(S, a) + \sum_{S' \in S^a} P(S'|a) \max_j \Lambda_n(a, S', j) \\ \bar{\alpha}_n(i, S, a) = \alpha_{\nu_n(i, S, a)} I(S_n^i = S, a_n^i = a).$$

Let  $\{\mathcal{M}_{n+1}\}$  be a random vector whose  $(i, S, a)$ th element is constructed as follows:

$$\mathcal{M}_{n+1}(i, S, a) = \max_j \Lambda_{n_a}(a, S_{n_a}^a, j) \\ - \sum_{S' \in \mathfrak{S}^a} P(S'|a) \max_j \Lambda_{n_a}(a, S', j)$$

where  $0 \leq n_a \leq n$ , and  $S_{n_a}^a$  is the most recent state visited by node  $a$ .

Now, we can rewrite (2) and (3) as in the form investigated in Fact 2, i.e.,

$$\Lambda_{n+1}(i, S, a) = \Lambda_n(i, S, a) + \bar{\alpha}_n(i, S, a) ((U\Lambda_{n_a})(i, S, a) \\ - \Lambda_n(i, S, a) + \mathcal{M}_{n+1}(i, S, a)).$$

The remaining steps of the proof reduce to verifying statements G1–G5. This is verified in Lemma 4.

*Lemma 4:*  $(\Lambda_n, \bar{\alpha}_n, \mathcal{M}_{n+1})$  satisfy conditions G1–G5.

*Proof:*

- (G1): It is shown in Lemma 6 that algorithm d-AdaptOR guarantees that every state-action is attempted infinitely often (i.o.). Hence

$$\sum_{n=0}^{\infty} \bar{\alpha}_n(i, S, a) = \sum_{n=0}^{\infty} \alpha_{\nu_n(i, S, a)} I(S_n^i = S, a_n^i = a) \\ \geq I((i, S, a) \text{ visited i.o.}) \left( \sum_{n=0}^{\infty} \alpha_n \right) = \infty.$$

However

$$\sum_{n=0}^{\infty} \bar{\alpha}_n^2(i, S, a) \leq \sum_{i, S, a} \sum_{n=0}^{\infty} \alpha_{\nu_n(i, S, a)}^2 I(S_n^i = S, a_n^i = a) \\ \leq \sum_{i \in \Theta} |\mathfrak{S}^i| d + 1 \sum_{n=0}^{\infty} \alpha_n^2 < \infty.$$

- (G2):

$$\mathbf{E}[\mathcal{M}_{n+1}|\mathcal{G}_n, n_a] = \mathbf{E}_{S^a} \left[ \max_j \Lambda_{n_a}(a, S^a, j) \right] \\ - \sum_{S'} P(S'|a) \max_j \Lambda_{n_a}(a, S', j) \\ = 0.$$

$$\mathbf{E}[\mathcal{M}_{n+1}|\mathcal{G}_n] = \mathbf{E}_{n_a}[\mathbf{E}[\mathcal{M}_{n+1}|\mathcal{G}_n, n_a]] = 0.$$

$$\mathbf{E}[\mathcal{M}_{n+1}^2|\mathcal{G}_n, n_a] \leq \mathbf{E}_{S^a} \left[ \left( \max_j \Lambda_{n_a}(a, S^a, j) \right)^2 \right] \\ \leq \max_{S^a} \max_j (\Lambda_{n_a}(a, S^a, j))^2 \\ \leq \|\Lambda_{n_a}\|^2.$$

$$\mathbf{E}[\mathcal{M}_{n+1}^2|\mathcal{G}_n] = \mathbf{E}_{n_a}[\mathbf{E}[\mathcal{M}_{n+1}^2|\mathcal{G}_n, n_a]] \\ \leq \mathbf{E}_{n_a}[\|\Lambda_{n_a}\|^2] \\ \leq \max_{n' \leq n} \|\Lambda_{n'}\|^2.$$

Thus Assumption (G2) of Fact 2 is satisfied.

- (G3): Let  $Z_d = \{S : d \in S, S \in \{\mathfrak{S}^i\}_{i \in \Theta}\}$  denote the set of states that contain the destination node  $d$ . Moreover, let  $Z_d^i = \{S : d \in S, i \in \Theta, S \in \mathfrak{S}^i\}$ . Let  $\tau_{Z_d}^\pi$  be the hitting time associated with set  $Z_d$  and policy  $\pi \in \Pi$ , i.e.,  $\tau_{Z_d}^\pi = \min\{n > 0 : \exists S \in Z_d, S \in \{S_n^i\}_{i \in \Theta}\}$ . Policy  $\pi$  is said to be proper if  $\text{Prob}(\tau_{Z_d}^\pi < \infty | \mathcal{F}_0) = 1$ . Let us now fix a proper deterministic stationary policy  $\pi \in \Pi$ . Existence of such a policy is guaranteed from the connectivity between 0 and  $d$ . Let  $F$  be the termination state that is reached after taking the termination action  $T$ . Let us define a policy dependent operator  $\mathcal{L}^\pi$

$$(\mathcal{L}^\pi \Lambda)(i, S, a) = g(S, a) + \sum_{S' \notin Z_d^a \cup F} P(S'|a) \Lambda(a, S', \pi(S')). \quad (10)$$

We then consider a Markov chain with states  $(i, S, a)$  and with the following dynamics: From any state  $(i, S, a)$ , we move to state  $(a, S', \pi(S'))$ , with probability  $P(S'|a)$ . Thus, subsequent to the first transition, we are always at a state of the form  $(i, S, \pi(S))$ , and the first two components of the state evolve according to policy  $\pi$ . As  $\pi$  is assumed proper, it follows that the system with states  $(i, S, a)$  also evolves according to a proper policy. We construct a matrix  $Q$  with each entry corresponding to the transition from state  $(i, S)$  to  $(\pi(S), S')$  with value equal to  $P(S'|\pi(S))$  for all  $S \notin Z_d^i \cup F, S' \notin Z_d^{\pi(S)} \cup F$  for all  $i$ . Since policy  $\pi$  is proper, the maximum eigenvalue of matrix  $Q$  is strictly less than 1. As  $Q$  is a nonnegative matrix, Perron Frobenius theorem guarantees the existence of a positive vector  $w$  with components  $w_{(i, S, a)}$  and some  $\beta \in [0, 1)$  such that

$$\sum_{S' \notin Z_d^a \cup F} P(S'|a) w_{(\pi(S), S', \pi(S'))} \leq \beta w_{(i, S, a)}. \quad (11)$$

From (11), we have a positive vector  $v$  such that  $\|(\mathcal{L}^\pi \Lambda) - \Lambda^\pi\|_v \leq \beta \|\Lambda - \Lambda^\pi\|_v$ , where  $\Lambda^\pi$  is the fixed point of equation  $\Lambda = \mathcal{L}^\pi \Lambda$ .

From the definition of  $U$  (4) and  $\mathcal{L}^\pi$  (10), we have  $|(U\Lambda)(\cdot, \cdot, \cdot)| \leq |(\mathcal{L}^\pi \Lambda)(\cdot, \cdot, \cdot)|$ . Using this and the triangle inequality, we obtain

$$\|U\Lambda\|_v \leq \|\mathcal{L}^\pi \Lambda\|_v \\ \leq \|\mathcal{L}^\pi \Lambda - \mathcal{L}^\pi \Lambda^\pi\|_v + \|\mathcal{L}^\pi \Lambda^\pi\|_v \\ \leq \beta \|\Lambda - \Lambda^\pi\|_v + \|\Lambda^\pi\|_v \\ \leq \beta \|\Lambda\|_v + 2\|\Lambda^\pi\|_v$$

establishing the validity of (G3).

- (G4): Assumption (G4) is satisfied by operator  $U$  using the following fact:

*Fact 3 [19, Proposition 4.3.1]:*  $U$  is monotonically increasing, continuous, and satisfies  $U(\Lambda) - \mathbf{r}1 \leq U(\Lambda -$

$r1) \leq U(\Lambda + r1) \leq U(\Lambda) + r1$ ,  $r > 0$ .  $\Lambda^*$  is a fixed point of  $U$ . From (5) and (6), we obtain

$$\max_{j \in A(S)} V^*(j) = \max_{j \in A(S)} \Lambda^*(i, S, j) + R. \quad (12)$$

Furthermore, using (5) and (12), for all  $i \in \Theta$

$$\Lambda^*(i, S, a) = g(S, a) + \sum_{S'} P(S'|a) \max_{j \in A(S')} V^*(j) - R. \quad (13)$$

The existence of fixed point  $\Lambda^*$  follows from (13), while uniqueness of  $\Lambda^*$  follows from uniqueness of  $V^*$  (Fact 1).

- (G5): Suppose  $n_a \rightarrow \infty$  as  $n \rightarrow \infty$ . Therefore, there exists  $N$  such that  $n_a < N$  for all  $n$ . This means that the number of times that node  $a$  has transmitted a packet is bounded by  $N$ . However, this contradicts Lemma 6, which says that each state-action pair  $(S, a)$  is visited i.o. Therefore,  $n_a \rightarrow \infty$  as  $n \rightarrow \infty$  for all  $a$ , and condition (G5) holds.

Thus, Assumptions (G1)–(G5) are satisfied. Hence, from Fact 2, our iterate (2) converges almost surely to  $\Lambda^*$ , the unique fixed point of  $U$ . ■

**Lemma 5:** If policy  $\phi^*$  is followed, then action  $a \in A(S)$  is selected i.o. if state  $S \in \mathfrak{S}$  is visited i.o.

*Proof:* Define the random variable  $K_n = I(S_n^i = S)$  for any  $i \in \Theta|\phi^*$ . Let  $\mathcal{K}_n$  be the  $\sigma$ -field generated by  $(K_1, K_2, \dots, K_n)$ . Let  $A_n = \{\omega : a_n^i = a, S_n^i = S \text{ for any } i \in \Theta|\phi^*\}$ .

From the construction of the algorithm, it is clear that  $A_n$  is  $\mathcal{K}_n$  measurable. Now, it is clear that under policy  $\phi^*$ ,  $A_{n+1}$  is independent of  $K^{n-1}$  given  $K_n$  and  $N_n(i, S)$ ,  $i \in \Theta$ . Define

$$P(A_{n+1}|K_n, N_n(i, S) \text{ for all } i \in \Theta) \geq \begin{cases} 0, & \text{if } K_n = 0 \\ \frac{\min_{i \in \Theta} \epsilon_n(i, S)}{|A(S)|}, & \text{if } K_n = 1 \end{cases} \quad (14)$$

$$\begin{aligned} & \sum_{n=0}^{\infty} \text{Prob}(A_{n+1}|\mathcal{K}_n) \\ & \geq \sum_{n=0}^{\infty} \text{Prob}(A_{n+1}|K_n, N_n(i, S) \text{ for all } i \in \Theta) \\ & \geq I(S \text{ is visited i.o.}) \sum_{n=0}^{\infty} \min_{i \in \Theta} \frac{\epsilon_n(i, S)}{|A(S)|} \\ & \geq \frac{I(S \text{ is visited i.o.})}{|A(S)|} \sum_{n=0}^{\infty} \frac{1}{\sum_{i \in \Theta} N_n(i, S) + 1} \\ & \geq \frac{I(S \text{ is visited i.o.})}{|A(S)|} \sum_{n=0}^{\infty} \frac{1}{n(d+1) + 1} = \infty. \end{aligned} \quad (15)$$

The next step of the proof is based on the following fact.

**Fact 4 [28, Corollary 5.29] (Extended Borel–Cantelli Lemma):** Let  $\mathcal{K}_k$  be an increasing sequence of  $\sigma$ -fields and let  $A_k$  be  $\mathcal{K}_k$ -measurable. If  $\sum_{k=1}^{\infty} \text{Prob}(A_k|\mathcal{K}_{k-1}) = \infty$ , then  $P(A_k \text{ i.o.}) = 1$ .

Thus, from Fact 4,  $a \in A(S)$  is visited i.o. if  $S$  is visited i.o. ■

**Lemma 6:** If policy  $\phi^*$  is followed, then each state-action  $(S, a)$  is visited infinitely often.

*Proof:* We say states  $S, S' \in \mathfrak{S}$  communicate if there exists a sequence of actions  $\{a_1, \dots, a_k, k < \infty\}$  such that probability of reaching state  $S'$  from state  $S$  following the sequence of actions  $\{a_1, \dots, a_k\}$  is greater than zero. Using Lemma 5, if state  $S \in \mathfrak{S}$  is visited i.o., then every action  $a \in A(S)$  is chosen i.o. as the set  $A(S)$  is finite. Hence, states  $S'$  such that  $P(S'|a) > 0$ ,  $S' \in \mathfrak{S}$ , are visited i.o. if  $S$  is visited i.o. By Lemma 5, every action  $a' \in A(S')$  is also visited i.o. Following similar argument and repeated application of Lemma 5, every state  $S'' \in \mathfrak{S}$  that communicates with state  $S$  and actions  $a \in A(S'')$  is visited i.o.

Under the assumption of the packet generation process in Section II, a packet is generated i.o. at the source node 0. Thus, state  $\{0\}$  is reached i.o. The construction of set  $\mathfrak{S}$  is such that every state  $S \in \mathfrak{S}$  communicates with state  $\{0\}$ . Thus, each  $(S, a)$  is visited i.o. since  $|\mathfrak{S}|$  is finite. ■

## B. Proof of Lemma 2

**Lemma 2:** For any (P)-admissible policy  $\phi \in \Phi$  for Problem (P) and for all  $N = 1, 2, \dots$

$$\mathbf{E}^\phi \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \leq V^*(0).$$

*Proof:* To prove the lemma, we refer to the Auxiliary Problem (AP). In this problem we have assumed the existence of a centralized controller with full knowledge of the local broadcast model. Mathematically speaking, let  $\mathbb{P}$  be the sample space of the random probability measures for the local broadcast model. Specifically,  $\mathbb{P} := \{p \in \mathbb{R}^{2^d} \times \mathbb{R}^d : p \text{ is a nonsquare left stochastic matrix}\}$ . Moreover, let  $\mathcal{P}_P$  be the trivial  $\sigma$ -field generated by the local broadcast model  $P \in \mathbb{P}$  (sample point in  $\mathbb{P}$ ), i.e.,  $\mathcal{P}_P = \{P, \mathbb{P} \setminus P, \emptyset, \mathbb{P}\}$ .<sup>11</sup> Recall that  $S_n^i$  denotes the set of nodes that have received the packet due to transmission from node  $i$  at time  $n$ , while  $a_n^i$  denotes the corresponding routing decision node  $i$  takes at time  $n$ .<sup>12</sup> For Auxiliary Problem (AP), a routing policy is a collection  $\pi = \{\pi^i\}_{i \in \Theta}$  of routing decisions taken for nodes  $i \in \Theta$  at the centralized controller, where  $\pi^i$  denotes a sequence of random actions  $\pi^i = \{a_0^i, a_1^i, \dots\}$  for node  $i$ . The routing policy  $\pi$  is said to be (AP)-admissible for Auxiliary Problem (AP) if the event  $\{a_n^i = a\}$  belongs to the product  $\sigma$ -field  $\mathcal{P}_P \times \prod_i \mathcal{H}_n^i$  [29].

From Fact 1, since  $\pi^*$  is the optimal policy for one packet, for each packet  $m$  and for any feasible policy  $\phi \in \Phi$

$$\begin{aligned} V^*(0) &= \mathbf{E}^{\pi^*} \left[ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} | \mathcal{F}_0 \right] \\ &\geq \mathbf{E}^\phi \left[ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right] \end{aligned}$$

<sup>11</sup> $\sigma$ -field captures the knowledge of the realization of local broadcast model and assumes a well-defined prior on these models.

<sup>12</sup> $S_n^i = \emptyset$ ,  $a_n^i = T$  if node  $i$  does not transmit at time  $n$ .

where the inequality follows from the fact that  $\Phi \subseteq \Phi$ . The remaining steps are straightforward

$$\begin{aligned} \mathbf{E}^\phi & \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \\ & \leq \mathbf{E}^\phi \left( \frac{1}{M_N} \sum_{m=1}^{M_N} V^*(0) \right) \\ & = V^*(0). \end{aligned}$$

### C. Proof of Lemma 3

*Lemma 3:* For any  $\delta > 0$

$$\liminf_{N \rightarrow \infty} \mathbf{E}^{\phi^*} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \geq V^*(0) - \delta.$$

*Proof:* From (5), (6), and (12), we obtain the following equality for all  $i \in \Theta$ ,  $S \in \mathfrak{S}^i$ :

$$\arg \max_{j \in A(S)} V^*(j) = \arg \max_{j \in A(S)} \Lambda^*(i, S, j). \quad (16)$$

Let

$$b = \min_{i \in \Theta} \min_{S \in \mathfrak{S}^i} \min_{\substack{j, k \in A(S) \\ \Lambda^*(i, S, j) \neq \Lambda^*(i, S, k)}} |\Lambda^*(i, S, j) - \Lambda^*(i, S, k)|.$$

Lemma 1 implies that, in an almost sure sense, there exists packet index  $m_1 < \infty$  such that for all  $n > \tau_s^{m_1}$ ,  $i \in \Theta$ ,  $S \in \mathfrak{S}^i$ ,  $a \in A(S)$

$$|\Lambda_n(i, S, a) - \Lambda^*(i, S, a)| \leq b/2.$$

In other words, from time  $\tau_s^{m_1}$  onwards, given any node  $i \in \Theta$  and set  $S \in \mathfrak{S}^i$ , the probability that d-AdaptOR chooses an action  $a \in A(S)$  such that  $\Lambda^*(i, S, a) \neq \max_{j \in A(S)} \Lambda^*(i, S, j)$  is upper-bounded by  $\epsilon_n(i, S)$ . Furthermore, since  $N_n(i, S) \rightarrow \infty$  (Lemma 6), for a given  $\gamma > 0$ , with probability 1, there exists a packet index  $m_2 < \infty$  such that for all  $n > \tau_s^{m_2}$ ,  $\max_{i, S} \epsilon_n(i, S) < \gamma$ .

Let  $m_0 = \max\{m_1, m_2\}$ . For all packets with index  $m \leq m_0$ , the overall expected reward is upper-bounded by  $m_0 R < \infty$  and lower-bounded by  $-\frac{m_0}{\lambda} d \max_i c_i > -\infty$ , hence their presence does not impact the expected average per-packet reward. Consequently, we only need to consider the routing decisions of policy  $\phi^*$  for packets  $m > m_0$ .

Consider the  $m$ th packet generated at the source. Let  $B_k^m$  be an event for which there exist  $k$  instances when d-AdaptOR routes packet  $m$  differently from the possible set of optimal actions. Mathematically speaking, event  $B_k^m$  occurs iff there exist instances  $\tau_s^m \leq n_1^m \leq n_2^m \dots n_k^m \leq \tau_T^m$  such that for all  $l = 1, 2, \dots, k$

$$\Lambda^*(i_{n_l^m, m}, S_{n_l^m}, a_{n_l^m}) \neq \max_{j \in A(S_{n_l^m})} \Lambda^*(i_{n_l^m, m}, S_{n_l^m}, j)$$

where  $S_{n_l^m}$  is the set of nodes that have successfully received packet  $m$  at time  $n_l^m$  due to transmission from node  $i_{n_l^m, m}$ . We call event  $B_k^m$  a misrouting of order  $k$ . For  $m > m_0$

$$\text{Prob}(B_k^m) \leq (\max_{i, S} \epsilon_n(i, S))^k \leq \gamma^k.$$

Now for packets  $m > m_0$ , let us consider the expected differential reward under policies  $\pi^*$  and  $\phi^*$

$$\begin{aligned} \mathbf{E}^{\pi^*} & \left[ \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \middle| \mathcal{F}_0 \right\} \right] - \mathbf{E}^{\phi^*} \left[ \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \\ & = V^*(0) - \mathbf{E}^{\phi^*} \left[ \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \\ & = \sum_{k=0}^{\infty} \mathbf{E}^{\phi^*} \left[ V^*(0) - \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \middle| B_k^m \right] \\ & \quad \times \text{Prob}(B_k^m) \\ & \leq \sum_{k=0}^{\infty} k R \text{Prob}(B_k^m) \\ & \leq R \sum_{k=1}^{\infty} k \gamma^k \\ & = \delta, \end{aligned} \quad (17)$$

where  $\delta = \frac{\gamma R}{(1-\gamma)^2}$ . Inequality (17) is obtained by noticing that maximum loss in the reward occurs if algorithm d-AdaptOR decides to drop packet  $m$  (no reward) while there exists a node  $j$  in the set of potential forwarders such that  $V^*(j) \approx R$ .

Thus, for all  $\delta > 0$ , the expected average per-packet reward under policy  $\phi^*$  is bounded as

$$\begin{aligned} \liminf_{N \rightarrow \infty} \mathbf{E}^{\phi^*} & \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} \left\{ r_m - \sum_{n=\tau_s^m}^{\tau_T^m-1} c_{i_{n,m}} \right\} \right] \\ & \geq \liminf_{N \rightarrow \infty} \mathbf{E}^{\phi^*} \left[ \frac{1}{M_N} \sum_{m=1}^{M_N} (V^*(0) - \delta) \right] \\ & = V^*(0) - \delta. \end{aligned}$$

### ACKNOWLEDGMENT

The authors would like to thank A. Plymoth and P. Johansson for the valuable discussions. They are grateful to the anonymous reviewers who provided thoughtful comments and constructive critique of the paper.

### REFERENCES

- [1] C. Lott and D. Teneketzis, "Stochastic routing in ad hoc wireless networks," in *Proc. 39th IEEE Conf. Decision Control*, 2000, vol. 3, pp. 2302–2307, vol. 3.
- [2] P. Larsson, "Selection diversity forwarding in a multihop packet radio network with fading channel and capture," *Mobile Comput. Commun. Rev.*, vol. 2, no. 4, pp. 47–54, Oct. 2001.
- [3] M. Zorzi and R. R. Rao, "Geographic random forwarding (GeRaF) for ad hoc and sensor networks: Multihop performance," *IEEE Trans. Mobile Comput.*, vol. 2, no. 4, pp. 337–348, Oct.–Dec. 2003.

- [4] S. Biswas and R. Morris, "ExOR: Opportunistic multi-hop routing for wireless networks," *Comput. Commun. Rev.*, vol. 35, pp. 33–44, Oct. 2005.
- [5] S. Jain and S. R. Das, "Exploiting path diversity in the link layer in wireless ad hoc networks," in *Proc. 6th IEEE WoWMoM*, Jun. 2005, pp. 22–30.
- [6] C. Lott and D. Teneketzis, "Stochastic routing in ad hoc networks," *IEEE Trans. Autom. Control*, vol. 51, no. 1, pp. 52–72, Jan. 2006.
- [7] E. M. Royer and C. K. Toh, "A review of current routing protocols for ad hoc mobile wireless networks," *IEEE Pers. Commun.*, vol. 6, no. 2, pp. 46–55, Apr. 1999.
- [8] T. Javidi and D. Teneketzis, "Sensitivity analysis for optimal routing in wireless ad hoc networks in presence of error in channel quality estimation," *IEEE Trans. Autom. Control*, vol. 49, no. 8, pp. 1303–1316, Aug. 2004.
- [9] J. N. Tsitsiklis, "Asynchronous stochastic approximation and Q-learning," in *Proc. 32nd IEEE Conf. Decision Control*, Dec. 1993, vol. 1, pp. 395–400.
- [10] J. Boyan and M. Littman, "Packet routing in dynamically changing networks: A reinforcement learning approach," in *Proc. NIPS*, 1994, pp. 671–678.
- [11] J. W. Bates, "Packet routing and reinforcement learning: Estimating shortest paths in dynamic graphs," 1995, unpublished.
- [12] S. Choi and D. Yeung, "Predictive Q-routing: A memory-based reinforcement learning approach to adaptive traffic control," in *Proc. NIPS*, 1996, pp. 945–951.
- [13] S. Kumar and R. Miikkulainen, "Dual reinforcement Q-routing: An on-line adaptive routing algorithm," in *Proc. Smart Eng. Syst., Neural Netw., Fuzzy Logic, Data Mining, Evol. Program.*, 2000, pp. 231–238.
- [14] S. S. Dhillon and P. Van Mieghem, "Performance analysis of the AntNet algorithm," *Comput. Netw.*, vol. 51, no. 8, pp. 2104–2125, 2007.
- [15] P. Purkayastha and J. S. Baras, "Convergence of Ant routing algorithm via stochastic approximation and optimization," in *Proc. IEEE Conf. Decision Control*, 2007, pp. 340–354.
- [16] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*. Belmont, MA: Athena Scientific, 1996.
- [17] S. Chachulski, M. Jennings, S. Katti, and D. Katabi, "Trading structure for randomness in wireless opportunistic routing," in *Proc. ACM SIGCOMM*, 2007, pp. 169–180.
- [18] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. New York: Wiley, 1994.
- [19] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Belmont, MA: Athena Scientific, 1997.
- [20] W. Stallings, *Wireless Communications and Networks*, 2nd ed. Upper Saddle River, NJ: Prentice-Hall, 2004.
- [21] J. Bicket, D. Aguayo, S. Biswas, and R. Morris, "Architecture and evaluation of an unplanned 802.11b mesh network," in *Proc. ACM MobiCom*, Cologne, Germany, 2005, pp. 31–42.
- [22] M. Kurth, A. Zubow, and J. P. Redlich, "Cooperative opportunistic routing using transmit diversity in wireless mesh networks," in *Proc. IEEE INFOCOM*, Apr. 2008, pp. 1310–1318.
- [23] J. Doble, *Introduction to Radio Propagation for Fixed and Mobile Communications*. Boston, MA: Artech House, 1996.
- [24] S. Russel and P. Norvig, *Artificial Intelligence: A Modern Approach*, 2nd ed. Upper Saddle River, NJ: Prentice-Hall, 2003.
- [25] R. Parr and S. Russell, "Reinforcement learning with hierarchies of machines," in *Proc. NIPS*, 1998, pp. 1043–1049.
- [26] P. Gupta and T. Javidi, "Towards throughput and delay optimal routing for wireless ad-hoc networks," in *Proc. Asilomar Conf.*, Nov. 2007, pp. 249–254.
- [27] M. J. Neely, "Optimal backpressure routing for wireless networks with multi-receiver diversity," in *Proc. CISS*, Mar. 2006, pp. 18–25.
- [28] L. Breiman, *Probability*. Philadelphia, PA: SIAM, 1992.
- [29] S. Resnick, *A Probability Path*. Boston, MA: Birkhuser, 1998.



**Abhijeet A. Bhorkar** received the B.Tech. and M.Tech. degrees in the electrical engineering from the Indian Institute of Technology, Bombay, India, both in 2006, and is currently pursuing the Ph.D. degree in electrical and computer engineering at the University of California, San Diego.

His research interests are primarily in the areas of stochastic control and estimation theory, information theory, and their applications in the optimization of wireless communication systems.



**Mohammad Naghshvar** (S'10) received the B.S. degree in electrical engineering from Sharif University of Technology, Tehran, Iran, in 2007, and is currently pursuing the M.S./Ph.D. degrees in electrical and computer engineering at the University of California, San Diego.

His research interests include stochastic control theory, network optimization, and wireless communication.



**Tara Javidi** (S'96–M'02) studied electrical engineering at the Sharif University of Technology, Tehran, Iran, from 1992 to 1996. She received the M.S. degrees in electrical engineering (systems) and applied mathematics (stochastics) and Ph.D. degree in electrical engineering and computer science from the University of Michigan, Ann Arbor, in 1998, 1999, and 2002, respectively.

From 2002 to 2004, she was an Assistant Professor with the Electrical Engineering Department, University of Washington, Seattle. She joined the University of California, San Diego, in 2005, where she is currently an Associate Professor of electrical and computer engineering. Her research interests are in communication networks, stochastic resource allocation, and wireless communications.

Dr. Javidi was a Barbour Scholar during the 1999–2000 academic year and received an NSF CAREER Award in 2004.



**Bhaskar D. Rao** (S'80–M'83–SM'91–F'00) received the B.Tech. degree in electronics and electrical communication engineering from the Indian Institute of Technology, Kharagpur, India, in 1979, and the M.S. and Ph.D. degrees from the University of Southern California, Los Angeles, in 1981 and 1983, respectively.

Since 1983, he has been with the University of California, San Diego, where he is currently a Professor with the Department of Electrical and Computer Engineering. His interests are in the areas of digital signal processing, estimation theory, and optimization theory, with applications to digital communications, speech signal processing, and human–computer interactions.

Dr. Rao has been a Member of the Statistical Signal and Array Processing Technical Committee of the IEEE Signal Processing Society. He is currently a Member of the Signal Processing Theory and Methods Technical Committee.