

A Practical Scheme for Wireless Network Operation

Radhika Gowaikar, *Student Member, IEEE*, Amir F. Dana, *Student Member, IEEE*, Babak Hassibi, and Michelle Effros, *Senior Member, IEEE*

Abstract—In many problems in wireline networks, it is known that achieving capacity on each link or subnetwork is optimal for the entire network operation. In this paper, we present examples of wireless networks in which decoding and achieving capacity on certain links or subnetworks gives us lower rates than other simple schemes, like forwarding. This implies that the separation of channel and network coding that holds for many classes of wireline networks does not, in general, hold for wireless networks. Next, we consider Gaussian and erasure wireless networks where nodes are permitted only two possible operations: nodes can either decode what they receive (and then re-encode and transmit the message) or simply forward it. We present a simple greedy algorithm that returns the optimal scheme from the exponential-sized set of possible schemes. This algorithm will go over each node at most once to determine its operation, and hence, is very efficient. We also present a decentralized algorithm whose performance can approach the optimum arbitrarily closely in an iterative fashion.

Index Terms—Forward/decode scheme, separation principle, wireless networks.

I. INTRODUCTION

IN A WIRELINE network having a single source and a single destination, we can think of information flow in the same terms as fluid flow, and obtain a max-flow min-cut result to get capacity. This treatment closely follows that of the Ford–Fulkerson [1] algorithm to give us a neat capacity result. This has been well understood for many years. However, until recently, similar min-cut capacity results were not known for any other class of network problems. Before we describe the recent results obtained in network problems, let us understand the general network problem. This can be stated in the context of a multiterminal network [2] as follows. We have a set of nodes, and the “channel” between these is specified by a probability transition function which governs the relationship between the signals transmitted by the nodes and how these are received by the other nodes. Every node can have messages that it wants to send to every other node. Because of the generality of this model, it

can be tailored to describe many practical systems easily. For instance, several wireless, as well as wireline, systems, (stationary) ad hoc and sensor networks, etc., can be modeled by choosing a suitable probability transition function.

In recent years, large ad hoc networks have received a lot of attention, starting with the work of Gupta and Kumar [3]. Most results involving these networks use relaying as a tool and consider issues like throughput, power efficiency, distortion. In addition, cooperation is a technique that has been shown to be very effective [4]. However, these methods study asymptotically large networks and give scaling laws, rather than exact results, for the performance measures that they study. In fact, finding the exact capacity region in this general setting is extremely challenging. In [2], outer bounds on the capacity region can be found. These have the form of “min-cut” upper bounds. Such an upper bound formalizes the intuitively satisfying notion that the rate from node a to node b cannot exceed the rate that any cutset of edges from a to b can support. However, determining whether schemes of network operation that reach this upper bound exist or not has proved to be very difficult. Even in simple relay networks, i.e., networks having one source node, one destination node, and a single other node (called the relay node), the answer to this question is not known, in general [2]. Only in special cases of the probability transition function (defined as “degraded” distributions) do we know schemes that can reach the upper bounds and thus attain capacity.

In this context, the results in [5] and [6] are remarkable. They say that in a wireline network setting, we can indeed achieve the min-cut upper bounds for a special case of a certain class of problems called multicast problems. In this problem, we have one source node and several sink nodes that want to receive the same message from the source. It turns out that using network coding techniques, we can achieve the min-cut capacity of the network. Further, [7] put this problem in an algebraic framework and presented *linear* schemes that also achieved this capacity. In addition, for some more general multicast problems, capacity has been shown to be achievable using linear network coding [7]. The work of [8]–[10] demonstrates the strengths of this algebraic approach.

We introduce the work presented in this paper by first examining a feature of the recent results in wireline networks and trying to determine if this feature is applicable in more general networks, viz., that in all the capacity-achieving schemes we have referred to above, the min-cut upper bounds are reached through separate channel and network coding. This means that there exists an optimal strategy in which each link in the wireline network can be made error-free by means of channel coding and network coding can be employed separately on top of this to determine which messages should be transmitted

Paper approved by Y. Fang, the Editor for Wireless Networks of the IEEE Communications Society. Manuscript received July 12, 2004; revised December 12, 2005 and July 10, 2006. This work was supported in part by the National Science Foundation under Grant CCR-0133818, in part by the Office of Naval Research under Grant N00014-02-1-0578, and in part by Caltech’s Lee Center for Advanced Networking. This paper was presented in part at the 41st Annual Allerton Conference on Communication, Control, and Computing, 2003, and in part at the Asilomar Conference, 2003.

R. Gowaikar and A. F. Dana were with the Department of Electrical Engineering, California Institute of Technology, Pasadena, CA 91125 USA. They are now with Qualcomm, Inc., San Diego, CA 92121 USA.

B. Hassibi and M. Effros are with the Department of Electrical Engineering, California Institute of Technology, Pasadena, CA 91125 USA (e-mail: hassibi@caltech.edu; effros@caltech.edu).

Digital Object Identifier 10.1109/TCOMM.2007.892448

on which link. This is quite unexpected, and leads us to wonder if such a separation can be optimal in more general network settings.

In the early sections of this paper, we will present simple wireless networks where this principle of separation fails. Thus we will show that operating wireless networks in a multihop manner, where each relay node decodes the message it receives, is not necessarily the right approach. This observation was first made in [11] and [12]. We will also suggest some schemes of operation that will outperform those that require the ability of relay nodes to decode.

We will focus attention on two specific wireless network models. The important features that characterize a wireless network are broadcast and interference. We will look at Gaussian wireless networks (GWNs) and erasure wireless networks (EWNs). The former has Gaussian channels as links and incorporates broadcast as well as interference. The second model has erasure channels as links and incorporates broadcast, but not interference. For these models, we will show that making links error-free can sometimes degrade the performance. In fact, asking nodes to simply forward their data rather than decoding it is sometimes more advantageous. This tells us that wireless networks need to be understood differently from wireline networks. We will see some explanations as to why this is the case later in the paper.

In our study of wireless networks, we propose a scheme of network operation that permits nodes only two operations. One is decoding to get the original data and then resending the same message as the source. The other is forwarding the data as received. Since each node has two options, we have an exponential-sized set of possible operations. We will present an algorithm that goes over each node at most once to find the optimal operation among this set of restricted operations. This will be a greedy algorithm that avoids searching over the exponential-sized set of possible operation allocations. We also present an algorithm that can approach the best rate arbitrarily closely in an iterative manner. This will be a “decentralized” algorithm, in the sense that each node needs only one bit of information from the destination in every iteration and no knowledge of the rest of the network in order to determine its own operation.

The organization of this paper is as follows. In Section II, we present two wireless network models. These will be the GWNs and EWNs. In Section III, we show that with these wireless models, making links or subnetworks error-free can be suboptimal. In Section IV, we will formally state the two operations that nodes will be permitted to perform. With this setup, we will state our problem of allocating appropriate operations in Section V. In Section VI, we will see how rates are calculated for all nodes in the network, and how asking certain nodes to decode and others to forward can affect the rate of the network. In Section VII, we will state our algorithm to find the optimal policy. In Section VIII, we will prove optimality of the algorithm. We will see some examples in Section IX that will show that the gap between the “all nodes decode” strategy and our method can be significant. In Section X, we will discuss the decentralized algorithm. We present upper bounds on the rate achievable by our scheme in Section XI. Conclusions and further questions are presented in Section XII.

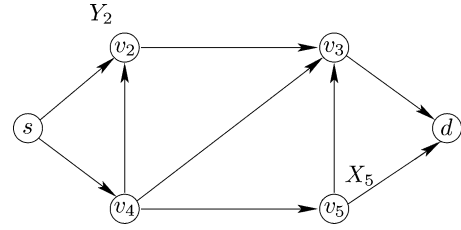


Fig. 1. Example of a network.

II. TWO WIRELESS NETWORK MODELS

In this section, we formalize two wireless network models. These are GWNs and EWNs. In both cases, the network consists of a directed, acyclic graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where \mathcal{V} is the set of vertices, and \mathcal{E} is the set of directed edges where each edge is a communication channel. We will denote $|\mathcal{V}| = V$ and $|\mathcal{E}| = E$. Also, we will have $\mathcal{V} = \{v_1, \dots, v_V\}$ and $\mathcal{E} = \{(v_i, v_j) : (v_i, v_j) \text{ is an edge}\}$. We will assume, without loss of generality, that $s = v_1$ is the source node and $d = v_V$ is the destination. The remaining nodes are the relay nodes which must aid communication between s and d . We will assume that every edge is on some directed path from s to d . If we have edges other than these, we remove them, and what remains is our graph \mathcal{G} . We will denote the message transmitted by vertex v_i by $X(v_i)$ and that received by node v_j by $Y(v_j)$. Fig. 1 represents a network with six vertices and nine edges, where v_1 is the source s and v_6 is the destination d . $X(v_5)$ is the message transmitted by v_5 and $Y(v_2)$ is that received by v_2 .

Gaussian Wireless Networks: In these networks, each edge (v_i, v_j) of the network is a Gaussian channel with some fixed attenuation factor $h_{i,j}$ associated with it. In a practical system, this may be some path loss that depends on the physical distances between the nodes. We will assume $h_{i,j}$ to be a nonnegative constant. We will assume that nodes broadcast messages, i.e., a node transmits the same message on all outgoing edges. Assuming that Fig. 1 represents a GWN, $X(v_5)$ is the message transmitted on edges (v_5, v_3) and (v_5, v_6) . We will also assume interference, i.e., the received signal at node v_i is the sum of all the signals transmitted on edges coming in to it and additive white Gaussian noise n_i of variance σ_i^2 . Therefore, in general, we have

$$Y(v_i) = n_i + \sum_{v_j: (v_j, v_i) \in \mathcal{E}} h_{j,i} X(v_j).$$

All n_i 's are assumed independent of each other as well as the messages. For Fig. 1, this implies that $Y(v_2) = h_{1,2}X(v_1) + h_{4,2}X(v_4) + n_2$. We will assume that all transmitting nodes have a power constraint of P .

Erasure Wireless Networks: In these networks, each edge (v_i, v_j) of the network is a binary erasure channel with erasure probability $\epsilon_{i,j}$. In addition, we assume that nodes (other than the source node) can transmit erasures, and they are received as erasures with probability 1. Denoting erasure by $*$, this assumption means that edges can also take $*$ as input, and this is always received as $*$. In short, the channel for edge (v_i, v_j) (for $v_i \neq s$) is modified as in Fig. 2. We incorporate broadcast in

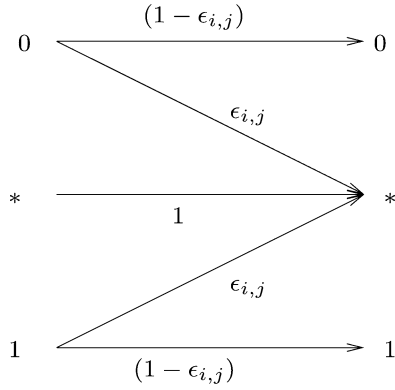


Fig. 2. Modified erasure channel.

the model, i.e., each transmitting node must send out the same signal on each outgoing edge. Now assuming that Fig. 1 represents a wireless erasure network, v_5 transmits $X(v_5)$ on edges (v_5, v_3) and (v_5, v_6) . However, we do not permit interference. This means that a node having several incoming edges sees messages from each edge without their interfering with each other. In general, if v_i has $\gamma_I(i)$ incoming edges, it will see $\gamma_I(i)$ messages that do not interfere with each other.¹ In Fig. 1, we see that $Y(v_2)$ consists of two received messages, the message coming in on edge (v_1, v_2) (which is $X(v_1)$ with some bits erased) and the message coming in on edge (v_4, v_2) (which is $X(v_4)$ with some bits erased). Finally, we mention that instead of the regular binary erasure channel, we can consider a channel with any finite alphabet \mathcal{A} as the input alphabet and get a more general EWN model. Our results go through for this also, but for simplicity, we restrict ourselves to binary inputs.

For both networks, we will assume instantaneous transmission on all links.

III. OPTIMIZING OVER SUBNETWORKS DOES NOT WORK

Theorem 1: For the wireless networks described in Section II, making subnetworks error-free can be suboptimal.

Proof: We give some examples to demonstrate this.

Gaussian Relay Networks: Consider a Gaussian parallel relay network consisting of two relay nodes and one source–destination pair. See Fig. 3(a). All four channel coefficients are assumed to be 1. The relay nodes v_2 and v_3 are solely to aid communication from source to destination. We assume that the noise power at each receiver is σ^2 and the transmit power at each node is P . Let $\rho \triangleq P/\sigma^2$ be the signal-to-noise ratio (SNR).

One way to view the network is as a cascade of a broadcast channel (from s to $\{v_2, v_3\}$) and a multiple-access channel (from $\{v_2, v_3\}$ to d). This is equivalent to assuming that the relays decode their messages correctly and code them again and transmit. If the relays are receiving independent information at rates R_1 and R_2 , we have $R_1 + R_2 \leq \log(1 + \rho)$ as the capacity

¹There exist network models in the physical layer that incorporate interference, which when abstracted to an erasure network model, act similarly to the interference-free model we have described here. For instance, simple division multiple-access schemes, such as TDMA, FDMA, or CDMA can be used to eliminate the interference.

region. These rate pairs (R_1, R_2) can be supported by the multiple-access channel, and hence, the maximum rate from s to d is no greater than $\log(1 + \rho)$. If the relays are receiving exactly the same information from the source, the maximum rate of this is $\log(1 + \rho)$. In this case, the multiple-access channel is used for correlated information, and can support rates up to $\log(1 + 4\rho)$. In either case, asking the relay nodes to decode limits the rate from s to d to $\log(1 + \rho)$. (We note also that the broadcast subnetwork is the bottleneck in both cases.)

Now consider another strategy in which the relay nodes do not decode, but only normalize their received signal to meet the power constraint and transmit it to the destination. In this case, the received signal at the destination is

$$Y(v_4) = \sqrt{\frac{P}{P + \sigma^2}} (2X(v_1) + n_2 + n_3) + n_4$$

where $X(v_1)$, $Y(v_4)$, n_2 , n_3 , n_4 are, respectively, the transmitted signal from the source, the received signal at the destination, and the noises introduced at v_2 , v_3 and d . Thus, the signal received by d is a scaled version of $X(v_1)$ with additive Gaussian noise. The maximum achievable rate, denoted by R_f , is

$$R_f = \log \left(1 + \frac{\frac{4P^2}{P + \sigma^2}}{\sigma^2 + \frac{2P\sigma^2}{P + \sigma^2}} \right) = \log \left(1 + \frac{4\rho^2}{3\rho + 1} \right)$$

where ρ is as before. Here, the subscript f stands for *forwarding*. We note that decoding at one of the relay nodes and forwarding at the other is always suboptimal.

In general, if we have $k (\geq 2)$ relay nodes in parallel rather than two, it can be easily checked that

$$R_d = \log(1 + \rho) \quad \text{and} \quad R_f = \log \left(1 + \frac{k^2 \rho^2}{(k + 1)\rho + 1} \right).$$

With this, we get a critical value of $\rho = 1/(k^2 - k - 1)$ below which decoding is better, and above which forwarding is better. Clearly, this goes to zero for large k . Therefore, in the limit of $k \rightarrow \infty$, it is always favorable to forward.

It turns out that this fact is also true for Gaussian relay networks in the presence of fading. The work of [11] shows that for fading Gaussian relay networks with n nodes, the asymptotic capacity achievable with the relay nodes decoding (and re-encoding) scales like $O(\log \log n)$, whereas with the forward scheme, it scales like $O(\log n)$.

Similar problems are considered in [13] and [14]. The former considers bounds and achievable rates for the Gaussian network with two parallel links, and the latter considers a network with a single source and destination and the other nodes acting as relays. The second result shows that the maximum rate achievable is $O(\log n)$. This is the same as that achieved by forwarding in our scheme.

Erasure Relay Network: Consider, once again, the network of Fig. 3(a), where now, each link represents an erasure channel with erasure probability $\epsilon_{i,j} = e$. Since we have broadcast, node s transmits the same messages to relay nodes v_2 and v_3 . If the relay nodes decode and re-encode, the rate is bounded by the

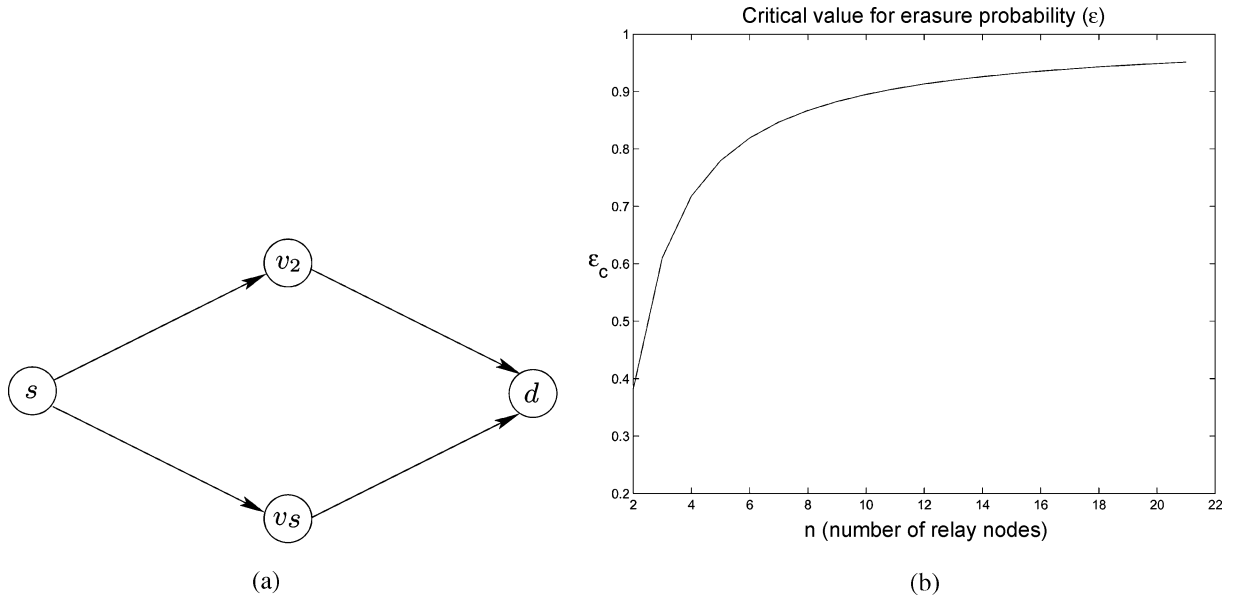


Fig. 3. Proof of *Theorem 1*. (a) Graph representation of a relay network with two relay nodes. (b) Critical value of erasure probability for k relay nodes.

sum-rate capacity of the broadcast system, which gives $R_d = 1 - e$.

If the relay nodes simply forward what they receive, it is easy to see that the destination sees an effective erasure probability of $(1 - (1 - e)^2)$. (We will spell out how to do this calculation for a general network in Section VI.) Forwarding erasures is possible since we are assuming the modified erasure channel of Fig. 2. With this, we have $R_f = 1 - (1 - (1 - e)^2)^2$. Comparing R_f and R_d , we can see that $e = (3 - \sqrt{5})/2$ is a critical value, above which decoding and re-encoding is better, and below which forwarding is better.

Thus we see that for this network also, making the broadcast subnetwork error-free is not always optimal.

In general, if we have k relay nodes in parallel rather than two, we have

$$R_d = 1 - e \quad \text{and} \quad R_f = 1 - (1 - (1 - e)^2)^k$$

and the critical value of e is as plotted in Fig. 3(b). Below this, forwarding is better, and above this, decoding is better. In the limit of large k , it is always better to forward.

From this, we see that making links or subnetworks error-free does not ensure optimal network operation. It can sometimes be provably suboptimal. \square

In this proof, a simple operation like forwarding the received data proved to be better than decoding it. We understand this as follows. Because of the broadcast present in wireless networks, the same data naturally gets passed on to the destination along many different paths. Therefore, some nodes receive better versions of the data on incoming links than other nodes, and are automatically in a better position to decode. Forcing all the nodes to decode and be error-free only imposes additional bottlenecks on the rate. Therefore, it is beneficial to carefully check the quality of the effective signal that various nodes get to see, and then decide whether to ask them to decode or not.

IV. A POSSIBLE SET OF NETWORK OPERATIONS

It follows from the previous discussions that to obtain the optimum rate over wireless networks, the nodes must perform operations other than just decoding. Determining what the optimum operation at each node should be, especially for a general wireless network, appears to be a daunting task. We shall therefore simplify the problem by allowing one of only two operations at every node. One will be the decode and re-encode operation as before. The other is the far simpler operation of forwarding the received data as is. The first operation, decode and re-encode, is typically the only operation used in multihop networks and many wireline networks. In effect, we are attempting to attain higher rates by introducing the additional operation of forwarding.

We will assume that the network operates in blocks of length n . We assume that the source s has a set of message indices

$$\Omega = \{1, 2, \dots, 2^{\lfloor nR \rfloor}\}$$

and an encoding function $f : \Omega \rightarrow \mathcal{X}^n$, where \mathcal{X} is \mathbb{R} for the GWN and $\{0, 1\}$ for the EWN. To transmit message $i \in \Omega$, the source transmits $f(i)$. With this, the source operates at rate R . $\{f(1), f(2), \dots, f(2^{\lfloor nR \rfloor})\}$ is the set of codewords or possible transmitted messages. This set is called the codebook, and is denoted by \mathcal{C} . We assume that all nodes have the codebook. For the Gaussian network, we will assume that the codebook meets the power constraint, i.e., $E\|f(i)\|^2 \leq P$.

In this paper, we restrict the relay nodes to two operations. These have been introduced in the examples of Section III, “forward” and “decode and re-encode.” We now state them formally.

Decode and Re-encode: This operation implies that when node v_i receives message $Y(v_i)$, it performs maximum-likelihood (ML) decoding of $Y(v_i)$ to determine which message

index was transmitted by s . Since it has the codebook, it re-encodes the message using the same codeword that the source s would have used, and transmits the same codeword. In short, it should act like a copy of the source.

However, for this to happen, we need that the decoding be error-free. This implies that the rate R at which the source operates should be no greater than the maximum rate at which node v_i can decode. We will see the relevance of this constraint in Section V.

Forward: We will describe this operation separately for the two network models. In the Gaussian network, node v_i receives message $Y(v_i)$ given by

$$Y(v_i) = n_i + \sum_{v_j: (v_j, v_i) \in \mathcal{E}} h_{j,i} X(v_j). \quad (1)$$

“Forwarding” implies that the node normalizes this signal to meet the power constraint and then transmits the message. Therefore, it transmits $X(v_i)$ given by

$$X(v_i) = \sqrt{\frac{P}{E \|Y(v_i)\|^2}} Y(v_i).$$

We will assume that $E \|Y(v_i)\|^2$ is known to v_i .

For the erasure network, nodes either decode without error and transmit the original codeword or “forward” the received data. Consider node v_i which sees data coming in on several edges, in the form of n -length blocks of bits and erasures. For the b th bit of such a block, it either sees erasures on every edge (and this sees an *effective* erasure), or gets to see the bit on at least one incoming edge. (It cannot happen that the node sees 1 on a particular edge and 0 on another edge for the b th position. This is because of our assumption that whenever an earlier node decodes, it does so without error.) Therefore, in our interference-free model, every relay node sees an effective erasure channel from the source, i.e., it sees the codeword transmitted by the source with some bits erased. “Forwarding” means broadcasting this sequence of bits and erasures.

Note that the effective erasure probability seen by node v_i is a function of the network topology and parameters, $\epsilon_{i,j}$. We will see in Section VI-C how this effective erasure probability can be calculated.

By restricting ourselves to only two operations, we have ensured that all nodes in the network see a Gaussian channel (with some effective SNR) or erasure channel (with some effective erasure probability) with respect to the transmitted codeword. Therefore, they can do ML decoding or typical set decoding if R is no greater than the rate that they can support. We will always ensure that R satisfies this constraint.

We can think of both operations as specific forms of network coding. In both networks and with both operations, all the information coming in at a node on different edges gets pooled together. This happens automatically in the Gaussian network and is done by the node itself in the erasure network. But the node has the choice of trying to decode, thus imposing a rate

constraint, or can simply forward the information, hoping that some other node would have a better chance of decoding.

Having described the two operations permitted to the relay nodes in the two networks, we are now ready to formally state the problem.

V. PROBLEM STATEMENT

Since we allow only two operations to nodes, decode and re-encode and forward, and every relay node must perform one of these, it is enough to specify the set of relay nodes that decode and re-encode in order to completely specify the working of the network. The source and destination will always be excluded from this set.

If a set $D \subseteq \mathcal{V} - \{s, d\}$ is the set of nodes that decode and re-encode, we will call D a **policy** for network operation.

Under policy D , each node of the network sees an effective (Gaussian or erasure) channel from the source. Let the effective SNR that node v_i sees under policy D be denoted by $\rho_D(v_i)$ for Gaussian networks. For erasure networks, we denote the effective erasure probability seen by node v_i under policy D by $e_D(v_i)$. Therefore, the rate that node v_i can support under policy D is $\log(1 + \rho_D(v_i))$ or $(1 - e_D(v_i))$ for Gaussian or erasure networks, respectively. In general, we will call this $R_D(v_i)$. Nodes in D as well as the destination must be able to perform error-free decoding. This means that the rate at which the source transmits must be no greater than the rates at which these nodes can decode. This tells us that under policy D , the rate R at which we can operate the network is constrained by

$$R \leq \min_{v_i \in D \cup \{d\}} R_D(v_i). \quad (2)$$

We denote this minimum by R_D

$$R_D = \min_{v_i \in D \cup \{d\}} R_D(v_i). \quad (3)$$

Intuitively, asking some nodes to decode means that there are more copies of the source in the network, and hence, the rate which the destination can support increases. On the other hand, asking a node to decode introduces a constraint on the rate R . This is the tradeoff for any policy D . For instance, in Fig. 1, consider nodes v_2 and v_4 . If v_4 forwards, node v_2 sees an effective erasure probability of $\epsilon_{4,2}\epsilon_{1,2} + \epsilon_{1,4}\epsilon_{1,2}(1 - \epsilon_{4,2})$. (We will see how this has been calculated in Section VI-C.) On the other hand, if v_4 decodes, node v_2 is at an advantage, since it sees a lower effective erasure probability, $\epsilon_{1,2}\epsilon_{4,2}$. However, asking v_4 to decode puts a constraint on the rate as seen by (2), since the rate that v_4 can support is only $(1 - \epsilon_{1,4})$. This constraint is $R_D \leq 1 - \epsilon_{1,4}$.

Our problem is to find the policy that gives the best rate, i.e., to find D such that R_D is maximized

$$\max_D \min_{v_i \in D \cup \{d\}} R_D(v_i).$$

First we need to address the question of finding $R_D(v_i)$, i.e., of finding the rate at node v_i under policy D . Recall that $X(v_i)$

and $Y(v_i)$ are the transmitted and received messages at node v_i . If we are using policy D , we will denote these by $X_D(v_i)$ and $Y_D(v_i)$. We may drop the subscript D if it is clear which policy we are referring to. Note that for the source, the transmitted message is $X(v_1)$, irrespective of the policy.

VI. DETERMINING THE RATE AT A NODE, $R_D(v_i)$

In this section, we describe a method to find the rate at an arbitrary node v_i when the set of decoding nodes is given by D . Therefore, we need to find the effective SNR or erasure probability of the received signal $Y_D(v_i)$. In order to do that, we need the concept of a partial ordering on the nodes.

A. Partial Ordering of Nodes

Consider two distinct nodes v_i and v_j of the network. Exactly one of the following will occur.

- 1) There is a directed path from v_i to v_j . In this case, we will say that $v_i < v_j$.
- 2) There is a directed path from v_j to v_i . In this case, we will say that $v_j < v_i$.
- 3) There is no directed path from v_i to v_j or from v_j to v_i . In this case, we will say that v_j and v_i are incomparable.

Note that since we assume acyclic networks, we cannot have directed paths both from v_i to v_j and from v_j to v_i . Thus, we have a partial ordering for nodes in the network. For example, in Fig. 1, we have $v_4 < v_3$, but v_2 and v_5 are incomparable. Note that the partial ordering gives us a (nonunique) sequence of nodes starting with s , such that for every v_i , all the nodes v_j that satisfy $v_j < v_i$ are before it in the sequence [15]. Call such a sequence \mathcal{S} . A possible sequence \mathcal{S} for Fig. 1 is $(s, v_4, v_2, v_5, v_3, d)$.

Next we address the issue of determining the rate under a particular policy. We discuss this separately for GWNs and EWNs.

B. Finding the Rate in GWNs

Recall that $Y_D(v_j)$ is the received signal at v_j under policy D . Once we know $Y_D(v_j)$, we can determine the signal power and the noise power in it. Denote these by $P_D(v_j)$ and $N_D(v_j)$, respectively. Consider node v_j . If it is decoding, $X_D(v_j) = X(v_1)$. If it is forwarding

$$\begin{aligned} X_D(v_j) &= \sqrt{\frac{P}{E \|Y_D(v_j)\|^2}} Y_D(v_j) \\ &= \sqrt{\frac{P}{P_D(v_j) + N_D(v_j)}} Y_D(v_j). \end{aligned}$$

We now outline a method for finding the rate for all the nodes by proceeding in the order given by \mathcal{S} . Without loss of generality, assume that the nodes are already numbered according to a partial ordering. Therefore, $\mathcal{S} = (v_1 = s, v_2, \dots, v_V = d)$. Then, for v_2 , we only have an edge coming in from s , and hence

$$Y_D(v_2) = h_{1,2}X(v_1) + n_2.$$

Let our induction hypothesis be that we know $Y_D(v_j)$ for $j = 1, \dots, i - 1$. For $Y_D(v_i)$, we now have

$$\begin{aligned} Y_D(v_i) &= n_i + \sum_{v_j:(v_j,v_i) \in \mathcal{E}} h_{j,i} X_D(v_j) \\ &= n_i + \sum_{v_j:(v_j,v_i) \in \mathcal{E}, v_j \in D \cup \{s\}} h_{j,i} X(v_1) \\ &\quad + \sum_{v_j:(v_j,v_i) \in \mathcal{E}, v_j \notin D \cup \{s\}} h_{j,i} X_D(v_j) \\ &= n_i + \sum_{v_j:(v_j,v_i) \in \mathcal{E}, v_j \in D \cup \{s\}} h_{j,i} X(v_1) \\ &\quad + \sum_{v_j:(v_j,v_i) \in \mathcal{E}, v_j \notin D \cup \{s\}} h_{j,i} \\ &\quad \times \sqrt{\frac{P}{P_D(v_j) + N_D(v_j)}} Y_D(v_j). \end{aligned} \quad (4)$$

By our hypothesis, we know all the $Y_D(v_j)$ that occur in the last summation. Substituting for these, we get $Y_D(v_i)$. Careful observation indicates that this will be a linear combination of $X(v_1)$ and the noise terms n_2, \dots, n_i .

In general, if this linear combination is given by

$$Y_D(v_i) = a_D X(v_1) + \sum_{j=2}^i a_{D,j}(v_i) n_j$$

we have $P_D(v_i) = a_D^2 P$ and $N_D(v_i) = \sum_{j=2}^i a_{D,j}^2(v_i) \sigma_j^2$. Once these are known, the SNR is simply $\rho_D(v_i) = P_D(v_i)/N_D(v_i)$, and the rate can be calculated as $R_D(v_i) = \log(1 + \rho_D(v_i))$. Clearly, the complexity of this procedure is $O(V)$.

C. Finding Rate in EWNs

We first put this problem in a graph-theoretic setting. We are given a directed, acyclic graph where certain nodes act as sources. For us, the set $D \cup \{s\}$ is the set of source nodes. All the edges of the graph have certain probabilities of failing, i.e., of being absent. For us, these are the erasure probabilities of the channel. With this setup, for every node v in the network (excluding s , but including those in D), we need to find the probability that there exists at least one directed path from some source node to this node. This is the network reliability problem in one of its most general formulations [16], [17]. This is a well-studied problem and is known to be $\#P$ -hard [17]. Although no polynomial-time algorithms to solve the problem are known, efficient algorithms for special graphs are known. An overview of the network reliability problem can be found in [18]. In the rest of this section, we propose two straightforward methods to compute the probabilities of connectivity that we are interested in. We will also mention some techniques that can reduce the computation involved in these methods.

Assume we have a policy D . Consider a node v_i of the network. To find $R_D(v_i)$, we need to find $e_D(v_i)$. A bit is erased at node v_i if it is erased on all incoming links. With each

edge (v_i, v_j) in the graph, associate a channel random variable $z(i, j)$. This takes the value 0 when a bit is erased, and the value 1 when a bit is not erased. Thus, it is a Bernoulli random variable with probability $(1 - \epsilon_{i,j})$.

Consider all the directed paths from s to v_i . Let there be k_i paths. Denote the paths by B_1, \dots, B_{k_i} . Let path B_j consist of l_j edges. We specify path B_j by writing in order the edges it traverses, i.e., with the sequence $((v_{j_1}, v_{j_2}), (v_{j_2}, v_{j_3}), \dots, (v_{j_{l_j}}, v_{j_{l_j+1}}))$. We know that $s = v_{j_1}$ and $v_i = v_{j_{l_j+1}}$. Consider the set of vertices excluding v_i that are on path v_j , i.e., $\{v_{j_i} : i = 1, \dots, l_j\}$. Some nodes in this set may belong to D , i.e., they are decoding nodes. In this case, we know that they transmit the original codeword exactly. Let t be the largest index in this set such that v_{j_t} decodes. Therefore, v_i will not receive bit b along path B_j only if an erasure occurs on an edge that comes after v_{j_t} in the path. We associate with path B_j the product of the random variables that affect this

$$Z_j = z(j_t, j_{t+1}) \cdot z(j_{t+1}, j_{t+2}) \cdots z(j_{l_j}, j_{l_j+1}).$$

This product is zero if one of the z random variables takes value zero, which, in turn, means that an erasure occurred on that edge.

Now, v_i sees an erasure only when none of the paths from s to itself manage to transmit the bit to it. Therefore, v_i sees an erasure when $Z_j = 0$ for *all* the paths B_j , $j = 1, \dots, k_i$. Therefore, we have

$$\begin{aligned} R_D(v_i) &= 1 - e_D(v_i) \\ &= 1 - P\left(\bigcap_{j=1}^{k_i} (Z_j = 0)\right) \\ &= P\left(\bigcup_{j=1}^{k_i} (Z_j \neq 0)\right). \end{aligned}$$

One way to evaluate this is by checking all possible combinations of values that the z variables can take and finding the total probability of those combinations that satisfy $\bigcup_{j=1}^{k_i} (Z_j \neq 0)$. This procedure has complexity $O(2^E)$. One observation that can make this procedure more efficient is that if we know that setting a certain subset of the z variables to 1 is enough to make the event $\bigcup_{j=1}^{k_i} (Z_j \neq 0)$ happen, then for every superset of this subset, setting all the z variables in that superset to 1 is also enough to make the event $\bigcup_{j=1}^{k_i} (Z_j \neq 0)$ happen. With this, we may have to check out fewer than the 2^E possible combinations of values for the z variables and reduce the complexity.

Another way to evaluate this is by using the inclusion-exclusion principle [15]. This gives us

$$P\left(\bigcup_{j=1}^{k_i} Z_j \neq 0\right) = \sum_{r=1}^{k_i} \sum_{1 \leq j_1 < \dots < j_r \leq k_i} (-1)^{r+1} P \times (Z_{j_1} \neq 0, \dots, Z_{j_r} \neq 0).$$

Since we have k_i paths, the above expression has $2^{k_i} - 1$ terms. A general term of the form $P(Z_{j_1} \neq 0, \dots, Z_{j_r} \neq 0)$ can be evaluated by first listing all the z variables that occur in at least one of the r terms. Say these are $z(i_1, j_1), \dots, z(i_q, j_q)$. Now $P(Z_{j_1} \neq 0, \dots, Z_{j_r} \neq 0)$ is given by the product $(1 - \epsilon_{i_1, j_1}) \times$

$\dots \times (1 - \epsilon_{i_q, j_q})$. This procedure has complexity $O(E2^k)$, where k is the $\max_i k_i$. In this procedure, the complexity of listing all the variables in a certain set of r terms can be reduced by storing the lists that one makes for sets of $(r - 1)$ terms and simply adding on the z terms from the r th term to the appropriate list.

VII. ALGORITHM TO FIND OPTIMUM POLICY

In general, since we have $V - 2$ relay nodes and each node has two options, forwarding and decoding and re-encoding, we have 2^{V-2} policies. To find the optimum policy, we can analyze the rate for each of these policies and determine the one that gives us the best rate. This strategy of exhaustive search requires us to analyze 2^{V-2} policies.

Here, we propose a greedy algorithm that finds the optimum policy D which maximizes the rate. This algorithm requires us to analyze at most $V - 2$ policies. In the next section, we will give a proof of correctness for this algorithm.

-
- 1) Set $D = \emptyset$.
 - 2) Compute $R_D(v_i)$ for all $v_i \in \mathcal{V}$. (Use techniques of Section VI.)
Find $R_D = \min_{v_i \in D \cup \{d\}} R_D(v_i)$.
 - 3) Find $M = \{v_i | v_i \notin \{s, d\} \cup D, R_D \leq R_D(v_i)\}$.
 - 4) If $M = \emptyset$, terminate. D is the optimal strategy.
 - 5) If $M \neq \emptyset$, find the largest $D' \subseteq M$ such that $\forall v \in D'$,
 $R_D(v) = \max_{v_i \in M} R_D(v_i)$.
Let $D = D \cup D'$.
-
- Return to 2.

At each stage of the algorithm, we look for nodes that are seeing a rate as good as or better than the current rate of network operation. If there are no such nodes, the algorithm terminates. If there are such nodes, we choose the best from among them. Thus, in every iteration, the nodes we add are such that they do not put additional constraints on the rate of the network. Therefore, the rate of the network can only increase in successive iterations.

Note that since we assume a finite network, this algorithm is certain to terminate. Also, since D cannot have more than $(V - 2)$ nodes, the algorithm cycles between steps 2–5 at most $(V - 2)$ times. This is significantly faster than the strategy of exhaustive search that requires us to analyze 2^{V-2} policies.

The complexity of the algorithm depends on how fast the computation of $R_D(v_i)$ can be done. We have seen techniques for this computation in Section VI.

VIII. ANALYSIS OF THE ALGORITHM

We first prove a lemma regarding the effect of decoding at a particular node on the rates supportable at other nodes.

Lemma 1: When node v is added to the decoding set D , the only nodes v_i that may see a change in rate are $v_i > v$. This change can only be an increase in rate, i.e., $\forall v_i$ such that $v_i > v$, we have $R_D(v_i) \leq R_{D \cup \{v\}}(v_i)$. Every other node v_j is unaffected, i.e., $R_D(v_j) = R_{D \cup \{v\}}(v_j)$.

Proof: We give a proof for the Gaussian network. We omit the proof for erasure networks, since it uses the same ideas.

Gaussian Network: Recall the computation of $\rho_D(v_i)$ described in Section VI-B. The computation for $Y_D(v_i)$ depends

only on (some of) the $Y_D(v_j)$ where (v_j, v_i) is an edge. Therefore, inductively, it is clear that $Y_D(v_i)$ (and hence, $\rho_D(v_i)$) depends only on the nodes v where $v < v_i$. Therefore, the only nodes that are affected when v changes its operation (from forwarding to decoding and re-encoding) are $v_i > v$. The rest are unaffected.

Consider one of the $X_D(v_j)$ terms in (4). Note that each of these are of power P of which some power is the signal power and the rest is the noise power. If v_j changes its operation from forwarding to decoding, $X_D(v_j) = X(v_1)$, i.e., the signal power increases to P and the noise power goes to 0. If v_j is forwarding, $X_D(v_j)$ is only a scaled version of $Y_D(v_j)$. Since it is always of power P , if the SNR at node v_j increases, the signal power in $X_D(v_j)$ increases while the noise power decreases. From (4), we see that in both these cases, there is an increase in the signal power of $Y_D(v_i)$ and a decrease in the noise power. This implies an increase in the SNR.

Therefore, when v is added to D , by induction, for all nodes $v_i > v$, the SNR, if affected, can only undergo an increase. Naturally, we have the same conclusion for the rate. \square

This lemma tells us that adding nodes to the set of decoding nodes can only increase the rate to other nodes. While this sounds like a good thing, it also puts a constraint on the rate, as indicated by (2). It is this tradeoff that our algorithm seeks to resolve by finding the optimal set of decoding nodes.

A. Proof of Optimality

Theorem 2: The algorithm of Section VII gives us an optimal set of decoding nodes.

Proof: Let S be an optimal set of decoding nodes. Let D be the set returned by the algorithm. We will prove that $R_D \geq R_S$. Then, since S is optimal, we will have $R_D = R_S$.

We prove $R_D \geq R_S$ in two steps. First we show that $R_{S \cup D} \geq R_S$. Then we show that $S \cup D - D = \emptyset$, i.e., $S \cup D = D$. This will complete the proof.

Step 1: In every iteration, the algorithm finds subsets D' and adds them to D . Denote by D_i the subset that is added to D in the i th iteration. Assuming the algorithm goes through m iterations, we have $D = D_1 \cup \dots \cup D_m$ where the union is over disjoint sets. In the algorithm, when D_i is added to D , all the nodes in it are decoding at the same rate, which is $R_{D_1 \cup \dots \cup D_{i-1}}(v)$ for $v \in D_i$. We will call this rate $R_{\text{algo},i}$. Consider the smallest i such that $D_i \not\subseteq S$, i.e., D_i is not already entirely in S .

Claim: Adding D_i to S does not decrease the rate, i.e., $R_{S \cup D_i} \geq R_S$.

Proof: Because of the acyclic assumption on the graph, we will have some nodes $v \in S$ such that $\forall u (\neq v) \in S$, we either have $v < u$, or v and u are incomparable. Let L be the set of all such nodes v . Note that by *Lemma 1*, node v supports a rate $R_S(v) = R_\emptyset(v)$. By (3), for every $v \in L$ we have the necessary condition

$$R_S \leq R_S(v) = R_\emptyset(v). \quad (5)$$

Also note that D_1, \dots, D_{i-1} are all in S , and by the definition of L and *Lemma 1*, we have

$$R_\emptyset(v) = R_{D_1 \cup \dots \cup D_{i-1}}(v). \quad (6)$$

We now consider two cases.

- If for some $w \in L$, we also have $w \in D_i$, then from (5) and (6), we have $R_S \leq R_S(w) = R_\emptyset(w) = R_{D_1 \cup \dots \cup D_{i-1}}(w) = R_{\text{algo},i}$.
- On the other hand, if none of the nodes in L are in D_i , pick any node $v \in L$. We have $v \notin D_i$. We now consider two subcases.
 - i) Let $v \notin D_1, \dots, D_{i-1}$. We note from Steps 3 and 5 of the algorithm that it picks out from the set of nodes not in D , all nodes with the best rate. Since v does not get picked, we have $R_{\text{algo},i} > R_{D_1 \cup \dots \cup D_{i-1}}(v)$. This, along with (5) and (6), gives us $R_S \leq R_{D_1 \cup \dots \cup D_{i-1}}(v) < R_{\text{algo},i}$.
 - ii) The other possibility is that $v \in D_1 \cup \dots \cup D_{i-1}$. Since the D_i 's are disjoint, there is a unique j such that $v \in D_j$. Since $v \in L$, by *Lemma 1*, $R_{\text{algo},j} = R_{D_1 \cup \dots \cup D_{j-1}}(v)$. With the same argument as that for (6), we have $R_\emptyset(v) = R_{D_1 \cup \dots \cup D_{j-1}}(v)$. But since the algorithm never decreases rate from one iteration to the next, we have $R_{\text{algo},i} \geq R_{\text{algo},j}$. Putting these together, we get $R_{\text{algo},i} \geq R_{\text{algo},j} = R_{D_1 \cup \dots \cup D_{j-1}}(v) = R_\emptyset(v)$. With (5), this gives us $R_S \leq R_S(v) = R_\emptyset(v) \leq R_{\text{algo},i}$.

Therefore, in every case, we have shown that $R_S \leq R_{\text{algo},i}$. This implies that adding the rest of the nodes from D_i to S will not put additional constraints on R_S , and hence, cannot decrease the rate. Therefore, we have $R_{S \cup D_i} \geq R_S$. \square

Since S is optimal, this proves that $S \cup D_i$ also achieves optimal rate. We can now call this set S , and for the next value of i such that $D_i \not\subseteq S$, we can prove that $S \cup D_i$ has optimal rate. Continuing like this, we have that $S \cup D$ is optimal, or, in other words, $R_{S \cup D} \geq R_S$.

Step 2: Next, we wish to show that $S \subseteq D$, i.e., $S \cup D - D = \emptyset$. Let us assume the contrary. Let $T = S \cup D - D$. Therefore, $T \cap D = \emptyset$ but $T \subseteq S$. Thus, $D \cup S = D \cup T$ where D and T are disjoint. Consider $v \in T$ such that $\forall u (\neq v) \in T$, we either have $v < u$ or v and u are incomparable. We have $R_{D \cup T}(v) = R_{D \cup S}(v)$. By *Lemma 1*, $R_{D \cup T}(v) = R_D(v)$. Also, the constraint of (2) tells us that $R_{D \cup S} \leq R_{D \cup S}(v)$. Finally, note that since the algorithm terminates without adding v to D , we have $R_D > R_D(v)$. Putting these inequalities together, we have $R_D > R_D(v) = R_{D \cup T}(v) = R_{D \cup S}(v) \geq R_{D \cup S}$. But this contradicts the fact that $S \cup D$ is optimal. Thus, we have $S \subseteq D$, i.e., $S \cup D = D$.

From Steps 1 and 2, we have $R_D \geq R_S$. But since S was an optimal policy, D is also an optimal policy. This proves that the algorithm does indeed return an optimal set of decoding nodes.

The only case in which this proof does not go through is when the algorithm returns $D = \emptyset$ and $S \neq \emptyset$. In this case, consider node $v \in L \subseteq S$, where L is as defined earlier. Since the algorithm does not pick up v , we have $R_\emptyset > R_\emptyset(v)$. But $R_S \leq R_S(v) = R_\emptyset(v)$ from (5). Thus, $R_S < R_\emptyset$. But this contradicts the optimality of S . Therefore, if there exists an optimal, nonempty S , the algorithm cannot return an empty D . \square

Corollary 1: The algorithm of Section VII returns the largest optimal policy D .

Proof: In the proof above, we have shown that for any optimal policy S , we have $S \subseteq D$. This implies that D is the largest optimal policy. \square

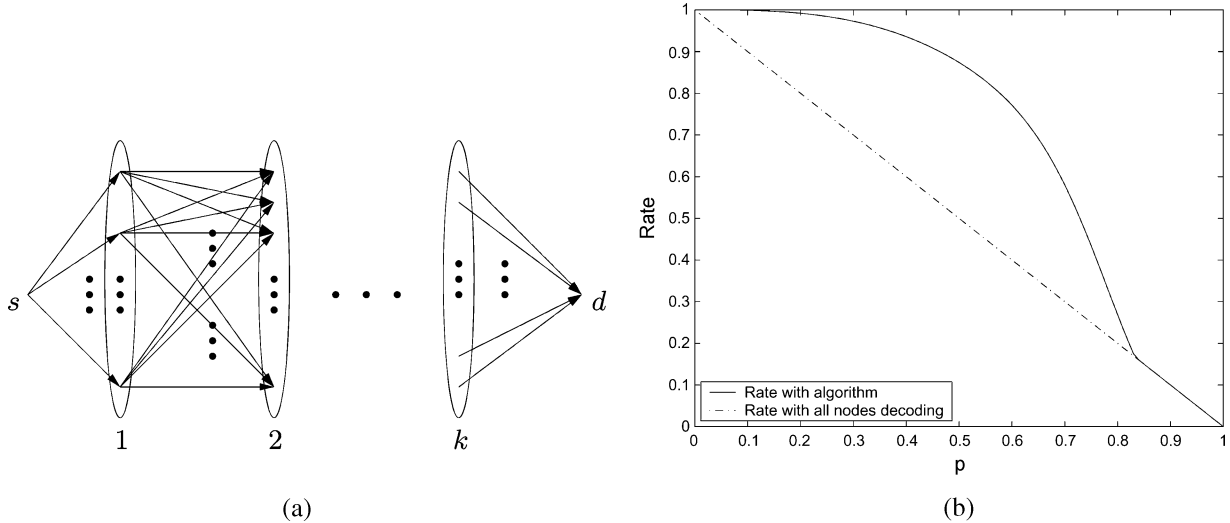


Fig. 4. Multistage relay network. (a) Model of a multistage relay network. (b) Rate for the multistage erasure relay network as given by the algorithm (solid curve) is much higher than that with all nodes decoding (dashed curve).

IX. EXAMPLES

In this section, we present some examples of networks and show how the algorithm runs on them.

A. Multistage Erasure Relay Networks

In Fig. 4(a), we have depicted a multistage relay network. In this, we have a single source and destination and k layers of relay nodes. The i th layer consists of l_i nodes. Between the i th and the $(i+1)$ th layer, we have a complete bipartite graph where all the edges are directed from the i th layer to the $(i+1)$ th. We assume that each of these edges has erasure probability ϵ_i . The source is connected to all the nodes in the first layer by erasure channels with erasure probability ϵ_0 , and all the nodes in the k th layer are connected to the destination by erasure channels with erasure probability ϵ_k . We will also call d the $(k+1)$ th layer and $l_{k+1} = 1$.

Because of the structure of this network, finding the rate under a particular policy is easier than indicated in Section VI-C. Denote by $Q_{i,j}$ the probability that in layer i there are j nodes that do not see an erasure. This defines $Q_{i,j}$ for $i = 1, 2, \dots, (k+1)$ and $j = 0, 1, \dots, l_i$. With this, for $i = 1$ we obtain

$$Q_{1,k} = \binom{l_1}{k} \epsilon_0^{l_1-k} (1 - \epsilon_0)^k. \quad (7)$$

For $i > 1$, we can show the recursion

$$Q_{i,k} = \binom{l_i}{k} \sum_{t=0}^{l_i-1} \epsilon_{i-1}^{t(l_i-k)} (1 - \epsilon_{i-1}^t)^k Q_{i-1,t}. \quad (8)$$

Denote by e_i the probability that the at least one node in the i th layer does not see an erasure. We can show that

$$e_i = \sum_{k=0}^{l_i} Q_{i,k} \left(1 - \frac{k}{l_i}\right).$$

Note that by symmetry, whenever a node decides to decode, all the nodes in that layer decode. When layer i decides to decode, we set $Q_{i,l_i} = 1$ and $Q_{i,j} = 0$ for $j \neq l_i$ and continue with the recursion of (8) for the other layers. This also extends to the case when more than one layer decodes.

Now, our algorithm proceeds as before, but operates on layers rather than nodes, and the effective erasure probability at layer i is e_i . As an explicit example, consider a multistage relay network with four layers between the source and destination. Let $l_1 = 3, l_2 = 6, l_3 = 4, l_4 = 5$, and $\epsilon_0 = p, \epsilon_1 = p^2, \epsilon_2 = p, \epsilon_3 = p^3, \epsilon_4 = p$ where p is any number in the interval $[0,1]$. For a fixed value of p , we can find the optimum policy for the network, and this will give us the optimal rate. Fig. 4(b) shows this optimal rate for the parameter p going from 0 to 1 (solid curve). This is not a smooth curve. The point where the right and left derivatives do not match is where either the optimum policy or the rate-determining layer changes. The rate with all nodes decoding has also been plotted (dashed curve). This rate is $1 - p$, and we see that the algorithm gives us dramatically higher rates.

B. Multistage Gaussian Relay Networks

We consider a multistage network similar to the one of the previous section, but in which the links represent Gaussian channels with fading coefficients h_i and with additive noise σ_i^2 at layer i . The indexing is identical to that in the erasure network (see Fig. 5).

Because of the structure of the network, it is easy to compute SNRs. Let $\rho(i)$ denote the SNR at layer i . Then, in the situation where all the nodes are forwarding, the following recursion gives us the SNR. We initialize the recursion as follows:

$$a(1) = h_0^2 P \quad b(1) = \sigma_1^2 \quad \rho(1) = \frac{a(1)}{b(1)}.$$

For the rest of the layers, i.e., $i \geq 2$, we have

$$\begin{aligned} a(i) &= a(i-1) \frac{h_i^2 l_i^2}{1 + \frac{1}{\rho(i-1)}} \\ b(i) &= b(i-1) \frac{h_i^2 l_i}{1 + \frac{1}{\rho(i-1)}} + \sigma_i^2 \\ \rho(i) &= \frac{a(i)}{b(i)}. \end{aligned}$$

As with the erasure relay network, whenever a node decides to decode, all the nodes in that layer decode. If some layers

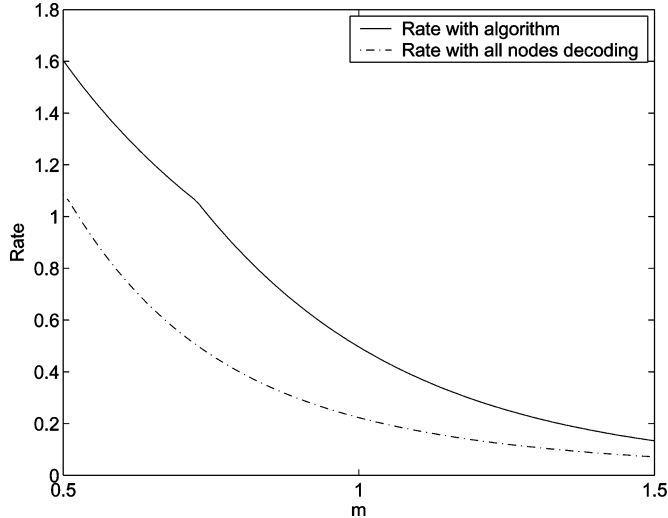


Fig. 5. Rate for the multistage Gaussian relay network as given by the algorithm (solid curve) is much higher than that with all nodes decoding (dashed curve).

decide to decode, a simple modification of the above recursion gives us the new rates. If i is the smallest number such that the i th layer decodes, then, clearly, the above recursion gives us rates for layers l_1 to l_i . For l_{i+1} , we set $a(i+1) = h_{i+1}^2 l_i^2 P$ and $b(i+1) = \sigma_{i+1}^2$. We have $\rho(i+1) = a(i+1)/b(i+1)$ as before, and we can continue with the recursion above for layers $(i+2)$, etc. We repeat this modification for each layer that decodes.

Once the SNR at a layer is known, the rate is given by $\log(1 + \rho)$ as usual. With this procedure for calculating rates, we use the algorithm of Section VII. It now operates on layers rather than nodes.

As an explicit example, consider a multistage relay network with three layers between the source and destination. Each node is restricted to using power $P = 1$. Let $l_1 = 2, l_2 = 5, l_3 = 3$, and $h_0 = 0.7, h_1 = 10, h_2 = 0.1, h_3 = 1$. We will have $\sigma_1^2 = m^2, \sigma_2^2 = m, \sigma_3^2 = m^3, \sigma_4^2 = m^2$ where m can be any positive real number. For a fixed value of m , we can find the optimum policy for the network and this will give us the optimal rate. Fig. 4(b) shows this optimal rate for the parameter m going from 0.5 to 1.5 (solid curve). As with the multistage erasure network, the curve is not smooth at points where the optimum policy or the rate-determining layer changes. We also see the advantage compared with the case when all nodes decode (dashed curve).

C. Erasure Network With Four Relay Nodes

Consider the relay network of Fig. 6(a). All the links have the same erasure probability p , where p is any number between 0 and 1. For this range of p , the algorithm has been used to find the optimum rates and policies. The rate is plotted in Fig. 6(b) (solid curve). Throughout, the optimum policy is $D = \{v_2, v_3, v_5\}$. The rate with all nodes decoding is $1 - p$ and is also plotted (dashed curve). As expected, the algorithm outperforms the all-decoding scheme.

D. Gaussian Network With Three Relay Nodes

In Fig. 7(a), we see a Gaussian network with three relay nodes. We assume that each node is restricted to use power $P = 1$. Let the additive noise variances be $\sigma_2^2 = m, \sigma_3^2 = m^3$,

$\sigma_4^2 = m^2, \sigma_5^2 = m^1$, where m can be an arbitrarily chosen real number. In Fig. 7(b), we see the rate returned by the algorithm for the optimal policy for $m \in [0.5, 1.5]$ (solid curve). The rate with all nodes decoding is also plotted (dashed curve). In the region $m \in [0.5, 0.58]$, we see that the optimal policy is, in fact, that of decoding at all nodes and the two curves match. After that, the optimal policy changes, and hence, we see that the optimal rate curve is not smooth.

E. Gaussian Network With Four Relay Nodes

In Fig. 8(a), we see a Gaussian network with four relay nodes. Each node, including the source, is restricted to using power $P = 1$. The attenuation factors associated with the edges are $h_{1,2} = 1, h_{1,4} = 2, h_{4,2} = 3, h_{2,3} = 4, h_{4,3} = 5, h_{4,5} = 1, h_{3,6} = 3, h_{5,3} = 2, h_{5,6} = 4$. The additive noise variances associated with the nodes are $\sigma_2^2 = m, \sigma_3^2 = m^3, \sigma_4^2 = m^2, \sigma_5^2 = m, \sigma_6^2 = m^3$, where m can be any positive real number. In Fig. 8, we see the rate returned by the algorithm for the optimal policy for $m \in [1, 3]$ (solid curve). The rate with all nodes decoding is also plotted (dashed curve). We see that the forward/decode scheme gives us significant improvements in the rate.

X. A DISTRIBUTED ALGORITHM FOR THE OPTIMAL POLICY

The algorithm as proposed in Section VII requires that the network parameters (noise variances or erasure probabilities) be known before the network operation begins, so that the optimum policy is known beforehand. With the algorithm in its current form, the nodes cannot determine for themselves if they should decode or forward. In this section, we propose a scheme that can permit nodes to determine their own operation.

The algorithm works iteratively to converge to a rate. In each iteration, the rate of operation of the network is incremented or decremented, depending on whether the previous transmission was successful or not. In every iteration, all the nodes get to decide their operation for themselves.

Let R^* be the maximum rate of the network. This is not known beforehand. We assume that parameters R, δ , and N are known to all the nodes beforehand. The block length n is also predetermined and known to all the nodes. In addition, we require that the nodes have a common source of randomness, so that they can generate the *same* random codebook individually. With this, consider the following algorithm.

- 1) All nodes generate the (same) codebook for rate R . They all set $k = 0$.
- 2) s transmits a randomly chosen codeword $X(v_1)$.
- 3) Every relay node v_i attempts to decode the received message $Y(v_i)$.
If it can decode without error, it transmits the decoded codeword.²
Else, it forwards the received message (with appropriate scaling, for the Gaussian network).
- 4) The destination attempts to decode the received message.
If it decodes without error, it sends back bit 1 to all the other nodes to indicate successful decoding.

²One method of error detection is for a node to perform typical set decoding, and assume an error if it finds more than one codeword that is jointly typical with the received message. Other methods of error detection are the introduction of cyclic redundancy checks (CRCs) or an automatic repeat request protocol, e.g. [19].

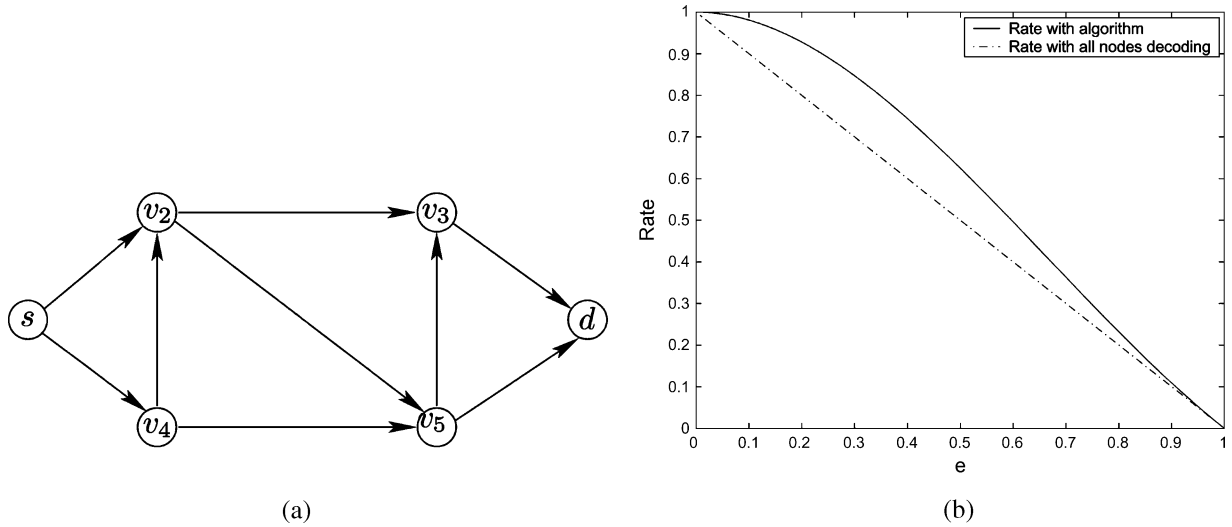


Fig. 6. Erasure network with four relay nodes. (a) Erasure network with four relay nodes. (b) Rate for the erasure network with four relay nodes as given by the algorithm (solid curve) is much higher than that with all nodes decoding (dashed curve).

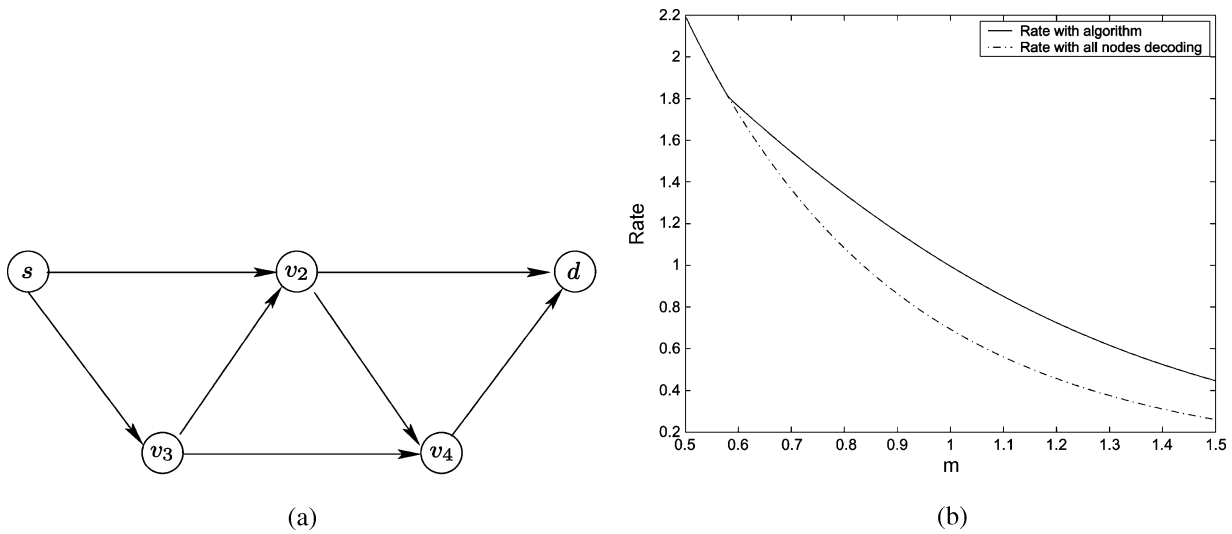


Fig. 7. Gaussian network with three relay nodes. (a) Gaussian network with three relay nodes. (b) Rate for the Gaussian relay network with three relay nodes as given by the algorithm (solid curve) is much higher than that with all nodes decoding (dashed curve).

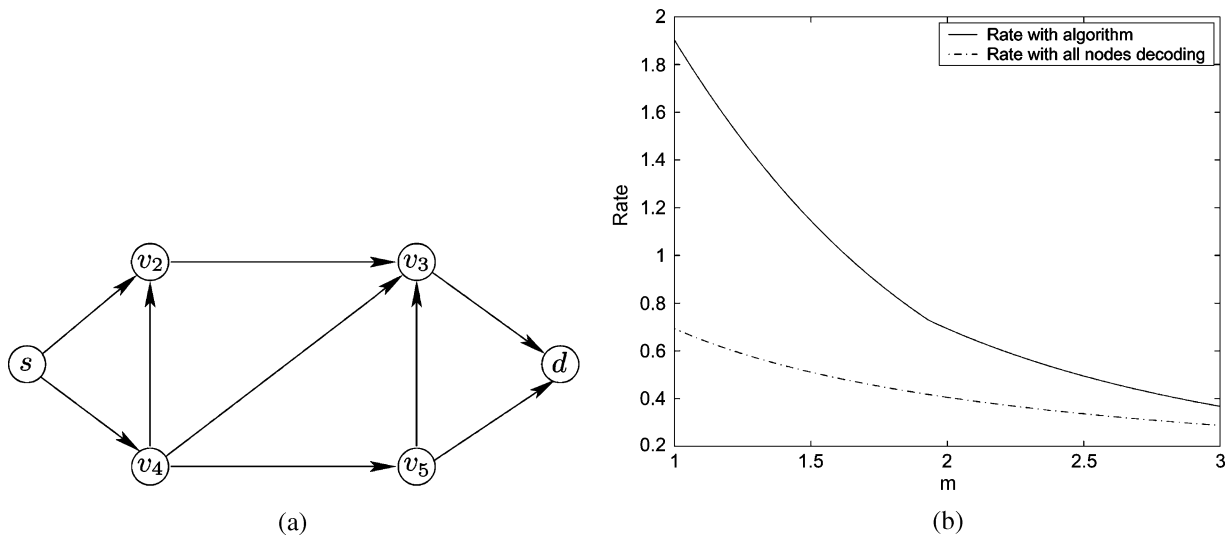


Fig. 8. Gaussian network with three relay nodes. (a) Gaussian network with four relay nodes. (b) Rate for the Gaussian relay network with four relay nodes as given by the algorithm (solid curve) is much higher than that with all nodes decoding (dashed curve).

Else, it sends back bit 0 to all other nodes.

5) All nodes increment k . $k = k + 1$.

If transmitted bit was zero, all nodes set $R = R - \delta/2^k$.

If transmitted bit was one, all nodes set $R = R + \delta/2^k$.

6) While $k \leq N$, go to Step 1.

Theorem 3: If the maximum rate of the network R^* is in the range $[R - \delta, R + \delta]$, the previous algorithm converges to it with an accuracy of $\delta/2^N$.

Proof: The source starts by transmitting at rate R . Each relay node receives messages on all incoming links and decodes the message if it can. If it cannot, it simply forwards what it has received. With this procedure, nodes decide their own operation. (The order in which they decide this is a partial order in the sense defined in Section VI-A.) After the destination receives all its incoming messages, it tries to decode. If $R > R^*$, the destination will definitely not be able to decode. If $R \leq R^*$, we claim that the destination will be able to decode. This is because when a node decodes, it only improves the rates for other nodes. Also, note that an arbitrary node v decides whether to decode or not only after all the nodes before it in the partial order have already determined if the rate they can support is greater or smaller than R . Since, by *Lemma 1*, these are the only nodes that affect the rate for v and they decode whenever they can, node v always gets to see the best situation it can, as far as rate R is concerned. This is true for the destination also.

Therefore, depending on whether the destination can decode or not, we can say if R^* is greater or smaller than R . If this bit of information is transmitted back to the source and other nodes, they can accordingly decide whether to increase or decrease the rate for the next transmission. Thus, we have a decision tree of rates such that the ability or inability of the decoder tells us which path to traverse in that tree we can finally converge on a rate sufficiently close to the actual rate R . \square

This algorithm provides a very natural mode of network operation that obviates the need for a central agent to know the entire network and decide the optimum policy. Although some communication from the destination to the source and other nodes is required, this is minimal and should be easily possible in a practical network setting.

We mention that the algorithm we present can be made more sophisticated, such that it works for all values of R^* , rather than just those in the interval $[R - \delta, R + \delta]$. We omit the details in the interests of brevity.

XI. UPPER BOUNDS ON THE MAXIMUM RATE

The algorithms of Section VII as well as Section X converge to the maximum rate possible with the decode/forward scheme, but we have no way of simply looking at the network and saying what this maximum rate will be. In this section, we present upper bounds on the rate achievable with the limited operations that we use in this paper.

A. Definitions

An $s - d$ cut is defined as a partition of the vertex set \mathcal{V} into two subsets \mathcal{V}_s and $\mathcal{V}_d = \mathcal{V} - \mathcal{V}_s$ such that $s \in \mathcal{V}_s$ and $d \in \mathcal{V}_d$. Clearly, an $s - d$ cut is determined simply by \mathcal{V}_s . For the $s - d$

cut given by \mathcal{V}_s , let the *cutset* $\mathcal{E}(\mathcal{V}_s)$ be the set of edges defined as

$$\mathcal{E}(\mathcal{V}_s) = \{(v_i, v_j) | (v_i, v_j) \in \mathcal{E}, v_i \in \mathcal{V}_s, v_j \in \mathcal{V}_d\}.$$

Finally, we define $X(\mathcal{V}_s)$ and $Y(\mathcal{V}_s)$ as

$$\begin{aligned} X(\mathcal{V}_s) &= \{v_i | (v_i, v_j) \in \mathcal{E}(\mathcal{V}_s)\} \\ Y(\mathcal{V}_s) &= \{v_j | (v_i, v_j) \in \mathcal{E}(\mathcal{V}_s)\}. \end{aligned}$$

Thus $X(\mathcal{V}_s)$ and $Y(\mathcal{V}_s)$ denote the nodes transmitting and receiving messages across the cut, respectively.

B. Upper Bound for Gaussian Networks

For Gaussian networks, it is evident that making the additive noise zero at certain nodes can only increase the maximum rate available at d . In particular, let us make the additive noise zero at all nodes except $Y(\mathcal{V}_s)$. Therefore, the received messages (and the transmitted messages) at all nodes in \mathcal{V}_s are exactly the same as that transmitted by the source. Now, if we permit the nodes in $Y(\mathcal{V}_s)$ to decode cooperatively, the rate at which they can decode will give us an upper bound on the rate that the destination can get.

Note that the SNR at node $v_j \in Y(\mathcal{V}_s)$ is

$$\frac{P}{\sigma_j^2} \left(\sum_{v_i: (v_i, v_j) \in \mathcal{E}(\mathcal{V}_s)} h_{i,j} \right)^2.$$

Since our codebook and noise are Gaussian distributed, the optimum scheme for decoding cooperatively is taking a suitable linear combination of received messages and then decoding that. For optimal decoding, we find the linear combination that gives us the best SNR. It is easy to show that the best SNR possible is the sum of the SNRs seen by each node in $Y(\mathcal{V}_s)$.

Therefore, an upper bound on the rate is

$$R \leq \log \left(1 + \sum_{v_j \in Y(\mathcal{V}_s)} \frac{P}{\sigma_j^2} \left(\sum_{v_i: (v_i, v_j) \in \mathcal{E}(\mathcal{V}_s)} h_{i,j} \right)^2 \right)$$

for every cut \mathcal{V}_s .

C. Upper Bound for Erasure Networks

As in the above section, we can obtain an upper bound on the rate for erasure networks by making certain links perfect, or free of erasures. Therefore, we can obtain an upper bound on the rate by making all edges other than those in $\mathcal{E}(\mathcal{V}_s)$ perfect. With this, all the received (and transmitted) messages in \mathcal{V}_s are exactly the same as the codeword transmitted by the source. Now, it is clear that the rate at which the nodes in $Y(\mathcal{V}_s)$ can decode cooperatively is an upper bound on the rate available at the destination.

Clearly, the effective erasure probability seen by the set of nodes $Y(\mathcal{V}_s)$ is $\prod_{(v_i, v_j) \in \mathcal{E}(\mathcal{V}_s)} \epsilon_{i,j}$. This gives us an upper bound on the rate. We have

$$R \leq 1 - \prod_{(v_i, v_j) \in \mathcal{E}(\mathcal{V}_s)} \epsilon_{i,j}$$

for every cut \mathcal{V}_s .

Note that in [20], a different min-cut upperbound is proposed and is shown to be achievable. This gives the capacity of the network under the assumption that the destination has perfect side-information regarding erasure locations from across the network. This is very different from the setup of this paper.

XII. CONCLUSIONS AND FURTHER QUESTIONS

To summarize, we have shown that making each link error-free in a wireless network is suboptimal. Thus, a multihop approach, in which every relay node decodes the received message, is not necessarily the correct approach for wireless networks. We have proposed a scheme for network operation that is of use in practical networks, and in which operations performed by a node are restricted to decoding and forwarding, both of which are common operations performed in a network setting. We have suggested an algorithm that finds the optimum policy without exhaustive search over an exponential number of policies, and also proposed a method to converge to the correct policy without having a central decision-making agent.

The algorithm of Section VII can find the maximum rate and optimum policy for any GWN or EWN. In addition, the bounds presented in Section XI give us some idea of what sort of optimal rates to expect. However, we still do not know what sort of policies are optimal in what ranges of erasure probabilities or SNR. The examples of Section III suggest that when the links are poor (high erasure probabilities or low SNR), it is better to decode. It would be interesting to know if this is true for general networks, and what thresholds exist below which a certain operation is always preferred.

Also, *Corollary 1* tells us that the algorithm returns the largest decoding set. Since decoding is the more costly of the two operations considered here, an algorithm that finds the smallest decoding set such that the maximum rate is obtained is of interest.

We note that both operations can be thought of as specific ways of doing network coding. We can also imagine a larger set of operations and optimal choice of operation from among these. The most general form of this would be when every node is free to use any function to encode the received data. This puts the problem in an information-theoretic setting, and a general solution for erasure networks is proposed in [20], where network coding techniques are used to obtain the precise capacity region for several multicast settings in erasure networks, assuming certain side-information. Naturally, this capacity region is an upper bound for the rates we have obtained in the absence of this side-information. Finding practical schemes that reach this capacity is an interesting avenue for future work.

REFERENCES

- [1] J. L. R. Ford and D. R. Fulkerson, *Flows in Networks*. Princeton, NJ: Princeton Univ. Press, 1962.
- [2] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [3] P. Gupta and P. R. Kumar, "The capacity of wireless networks," *IEEE Trans. Inf. Theory*, vol. 46, no. 3, pp. 388–404, Mar. 2000.
- [4] A. Sendonaris, E. Erkip, and B. Aazhang, "User cooperation diversity. Part I—System description. Part II—Implementation aspects and performance analysis," *IEEE Trans. Commun.*, vol. 51, no. 11, pp. 1927–1948, Nov. 2003.

- [5] R. Ahlswede, N. Cai, S.-Y. R. Li, and R. W. Yeung, "Network information flow," *IEEE Trans. Inf. Theory*, vol. 46, no. 4, pp. 1204–1216, Apr., 2000.
- [6] S.-Y. R. Li, R. W. Yeung, and N. Cai, "Linear network coding," *IEEE Trans. Inf. Theory*, vol. 49, no. 2, pp. 371–381, Feb., 2003.
- [7] R. Koetter and M. Médard, "An algebraic approach to network coding," *IEEE/ACM Trans. Netw.*, vol. 11, no. 5, pp. 782–795, Oct. 2003.
- [8] T. Ho, R. Koetter, M. Médard, D. Karger, and M. Effros, "The benefits of coding over routing in randomized setting," in *Proc. IEEE Int. Symp. Inf. Theory*, 2003, p. 442.
- [9] M. Effros, M. Médard, T. Ho, S. Ray, D. Karger, and R. Koetter, "Linear network codes: A unified framework for source channel, and network coding," in *Proc. DIMACS Workshop Netw. Inf. Theory*, 2003, invited paper, CD-ROM.
- [10] S. Jaggi, P. Sanders, P. A. Chou, M. Effros, S. Egner, K. Jain, and L. Tolhuizen, "Polynomial time algorithms for multicast network code construction," *IEEE Trans. Inf. Theory*, vol. 51, no. 6, pp. 1973–1982, Jun. 2005.
- [11] A. F. Dana, M. Sharif, B. Hassibi, and M. Effros, "Is broadcast plus multiaccess optimal for Gaussian wireless networks with fading?," in *Proc. 37th Asilomar Conf. Signals, Syst. Comput.*, Nov. 2003, vol. 2, pp. 1748–1752.
- [12] A. F. Dana, R. Gowaikar, B. Hassibi, M. Effros, and M. Mard, "Should we break a wireless network into sub-networks?," in *Proc. Annu. Allerton Conf. Commun., Control, Comput.*, 2003, CD-ROM.
- [13] B. Schein and R. Gallager, "The Gaussian parallel relay network," in *Proc. IEEE Int. Symp. Inf. Theory*, Jun. 2000, p. 22.
- [14] M. Gastpar and M. Vetterli, "On the capacity of wireless networks: The relay case," in *Proc. IEEE INFOCOM*, Jun. 2002, vol. 3, pp. 1577–1586.
- [15] J. H. van Lint and R. M. Wilson, *A Course in Combinatorics*. Cambridge, U.K.: Cambridge Univ. Press, 2001.
- [16] J. S. Provan and M. O. Ball, "The complexity of counting cuts and of computing the probability that a graph is connected," *SIAM J. Comput.*, vol. 12, no. 4, pp. 384–393, 1983.
- [17] L. G. Valiant, "The complexity of enumeration and reliability problems," *SIAM J. Comput.*, vol. 8, pp. 410–421, 1979.
- [18] H. L. Bodlaender and T. Wolle, "A note on the complexity of network reliability problems," 2003 [Online]. Available: <http://www.cs.uu.nl/research/techreps/aut/thomasw.html>
- [19] G. Caire and D. Tuninetti, "The throughput of hybrid-ARQ protocols for the Gaussian collision channel," *IEEE Trans. Inf. Theory*, vol. 47, no. 7, pp. 1971–1988, Jul. 2001.
- [20] A. F. Dana, R. Gowaikar, R. Palanki, B. Hassibi, and M. Effros, "Capacity of erasure wireless networks," *IEEE Trans. Inf. Theory*, vol. 52, no. 3, pp. 789–804, Mar. 2006.



Radhika Gowaikar (S'03) received the B.Tech. degree from the Indian Institute of Technology, Bombay, India, in 2001, and the M.S. and Ph.D. degrees from the California Institute of Technology, Pasadena, in 2002 and 2006, respectively, all in electrical engineering.

She is currently with Qualcomm, Inc., San Diego, CA. Her research interests include sensor and ad hoc networks, network coding for wireless networks, and decoding in multiple-antenna systems.



Amir F. Dana (S'98) was born in Tehran, Iran, in 1979. He received the B.S. degree in electrical engineering from the Sharif University of Technology, Tehran, Iran, in 2001, and the M.S. and Ph.D. degrees in electrical engineering from the California Institute of Technology, Pasadena, in 2002 and 2006, respectively.

He is currently with Qualcomm, Inc., San Diego, CA. His research interests are in the areas of information theory, and wireless communications. In particular, he has been working on power efficiency of sensor and ad hoc networks, network coding for wireless networks, and capacity of broadcast channels.



Babak Hassibi was born in Tehran, Iran, in 1967. He received the B.S. degree from the University of Tehran, Tehran, Iran, in 1989, and the M.S. and Ph.D. degrees from Stanford University, Stanford, CA, in 1993 and 1996, respectively, all in electrical engineering.

From October 1996 to October 1998, he was a Research Associate at the Information Systems Laboratory, Stanford University, and from November 1998 to December 2000, he was a Member of the Technical Staff, Mathematical Sciences Research Center, Bell

Laboratories, Murray Hill, NJ. Since January 2001, he has been with the Department of Electrical Engineering, California Institute of Technology, Pasadena, where he is currently an Associate Professor. He has also held short-term appointments at Ricoh California Research Center, the Indian Institute of Science, and Linköping University, Sweden. His research interests include wireless communications, robust estimation and control, adaptive signal processing and linear algebra. He is the coauthor of the books *Indefinite Quadratic Estimation and Control: A Unified Approach to H^2 and H^∞ Theories* (New York: SIAM, 1999) and *Linear Estimation* (Englewood Cliffs, NJ: Prentice-Hall, 2000).

Dr. Hassibi is a recipient of an Alborz Foundation Fellowship, the 1999 O. Hugo Schuck Best Paper Award of the American Automatic Control Council, the 2002 National Science Foundation Career Award, the 2002 Okawa Foundation Research Grant for Information and Telecommunications, the 2003 David and Lucille Packard Fellowship for Science and Engineering and the 2003 Presidential Early Career Award for Scientists and Engineers (PECASE). He has been a Guest Editor for the IEEE TRANSACTIONS ON INFORMATION THEORY special issue on "space-time transmission, reception, coding and signal processing" and an Associate Editor for Communications of the IEEE TRANSACTIONS ON INFORMATION THEORY during 2003–2006.



Michelle Effros (S'93–M'95–SM'03) received the B.S. degree with distinction in 1989, the M.S. degree in 1990, and the Ph.D. degree in 1994, all in electrical engineering, from Stanford University, Stanford, CA.

During the summers of 1988 and 1989, she was with Hughes Aircraft Company. She joined the faculty at the California Institute of Technology, Pasadena, in 1994, and is currently a Professor of Electrical Engineering. Her research interests include information theory, network coding, data compression, communications, pattern recognition,

and image processing.

Dr. Effros received Stanford's Frederick Emmons Terman Engineering Scholastic Award (for excellence in engineering) in 1989, the Hughes Masters Full-Study Fellowship in 1989, the National Science Foundation Graduate Fellowship in 1990, the AT&T Ph.D. Scholarship in 1993, the NSF CAREER Award in 1995, the Charles Lee Powell Foundation Award in 1997, the Richard Feynman-Hughes Fellowship in 1997, an Okawa Research Grant in 2000, and was cited by *Technology Review* as one of the world's top 100 young innovators in 2002. She is a member of Tau Beta Pi, Phi Beta Kappa, Sigma Xi, and the IEEE Information Theory, Signal Processing, and Communications societies. She served as the Editor of the *IEEE Information Theory Society Newsletter* from 1995 to 1998, and as a Member of the Board of Governors of the IEEE Information Theory Society from 1998 to 2003. She has served on the IEEE Signal Processing Society Image and Multi-Dimensional Signal Processing (IMDSP) Technical Committee since 2001. She was an Associate Editor for the joint special issue on Networking and Information Theory in the IEEE TRANSACTIONS ON INFORMATION THEORY and the IEEE/ACM TRANSACTIONS ON NETWORKING and is currently Associate Editor for Source Coding for the IEEE TRANSACTIONS ON INFORMATION THEORY.