

Probability of Partially Decoding Network-Coded Messages

Jessica Claridge and Ioannis Chatzigeorgiou

Abstract—In the literature there exists analytical expressions for the probability of a receiver decoding a transmitted source message that has been encoded using random linear network coding. In this work, we look into the probability that the receiver will decode at least a fraction of the source message, and present an exact solution to this problem for both non-systematic and systematic network coding. Based on the derived expressions, we investigate the potential of these two implementations of network coding for information-theoretic secure communication and progressive recovery of data.

Index Terms—Random linear network coding, rank-deficient decoding, probability analysis, information-theoretic security.

I. INTRODUCTION

Random linear network coding (RLNC) is the process of constructing coded packets, which are random linear combinations of source packets over a finite field [1]. If k source packets are considered, decoding at a receiving node starts after k linearly independent coded packets have been collected. The probability of recovering all of the k source packets when at least k coded packets have been received has been derived in [2]. However, the requirement for a large number of received coded packets before decoding can introduce undesirable delays at the receiving nodes. In an effort to alleviate this problem, *rank-deficient* decoding was proposed in [3] for the recovery of a subset of source packets when fewer than k coded packets have been obtained. Whereas the literature on network coding defines *decoding success* as the recovery of 100% of the source packets with a certain probability, the authors of [3] presented simulation results that measured the *fraction of decoding success*, that is, the recovery of a percentage of the source packets with a certain probability.

The fundamental problem that has motivated our work is the characterization of the probability of recovering some of the k source packets when n coded packets have been retrieved, where n can be smaller than, equal to or greater than k . This idea was considered in [4] for random network communications over a matroid framework. The authors show that partial decoding is highly unlikely. This problem has also been explored in the context of secure network coding, e.g., [5], [6]. Strict information-theoretic security can be achieved if and only if the mutual information between the packets available to an eavesdropper and the source packets is zero [7]. When network coding is used, *weak* security can be achieved if the eavesdropper cannot obtain k linearly independent

coded packets and, hence, cannot recover any meaningful information about the k source packets [5]. The authors of [5] obtained bounds on the probability of RLNC being weakly secure and showed that the adoption of large finite fields improves security. A different setting but a similar problem was investigated in [6]. Intermediate relay nodes between transmitting and receiving nodes were treated as potentially malicious, and criteria for characterizing the algebraic security of RLNC were defined. The authors demonstrated that the probability of an intermediate node recovering a strictly positive number of source packets tends to zero as the field size and the number of source packets go to infinity.

This paper revisits the aforementioned problem and obtains an exact expression for the probability that a receiving node will recover at least x of the k source packets if n coded packets are collected, for $x \leq n$. The derived expression can be seen as a generalization of [2, eq. (7)]. The paper also looks at the impact of transmitting source packets along with coded packets, known as *systematic* RLNC, as opposed to transmitting only coded packets, referred to as *non-systematic* RLNC.

In the remainder of the paper, Section II formulates the problem, Section III obtains the probability of recovering a fraction of a network-coded message, Section IV presents results and Section V summarizes the conclusions of this work.

II. SYSTEM MODEL AND PROBLEM FORMULATION

We consider a receiving network node, which collects n packets and attempts to reconstruct a message that consists of k source packets. The n packets could have been broadcast by a single transmitting node or could have been originated from multiple nodes that possess the same message.

In the case of *non-systematic* communication, transmitted packets are generated from the k source packets using RLNC over \mathbb{F}_q [1], where q is a prime power and \mathbb{F}_q denotes the finite field of q elements. In the case of *systematic* RLNC, a sequence of n_T transmitted packets consists of the k source packets and $n_T - k$ coded packets that have been generated as in the non-systematic case. In both cases, a coding vector of length k , which contains the weighting coefficients used in the generation of a packet, is transmitted along with each packet. At the receiving node, the coding vectors of the n successfully retrieved packets form the rows of a matrix $\mathbf{M} \in \mathbb{F}_q^{n \times k}$, where $\mathbb{F}_q^{n \times k}$ denotes the set of all $n \times k$ matrices over \mathbb{F}_q . The k source packets can be recovered from the n received packets if and only if k of the n coding vectors are linearly independent, implying that $\text{rank}(\mathbf{M}) = k$ for $n \geq k$. The probability that the $n \times k$ random matrix \mathbf{M} has rank k and, thus, the receiving node can reconstruct the entire message is given in [2] for non-systematic RLNC and [8] for systematic RLNC.

The objective of this paper is to derive the probability that a receiving node will reconstruct at least $x \leq k$ source packets

Data created during this work are openly available from the Lancaster University data archive at <http://dx.doi.org/10.17635/lancaster/researchdata/150>.

J. Claridge is with the Department of Mathematics, Royal Holloway, University of London, Egham, Surrey TW20 0EX, United Kingdom (e-mail: jessica.claridge.2013@live.rhul.ac.uk).

I. Chatzigeorgiou is with the School of Computing and Communications, Lancaster University, Lancaster LA1 4WA, United Kingdom (e-mail: i.chatzigeorgiou@lancaster.ac.uk).

upon reception of n network-coded packets. To formulate this problem, let \mathbf{e}_i denote the i -th unit vector of length k . A coding vector, or a row of \mathbf{M} , equal to \mathbf{e}_i represents the i -th source packet. Let X be the set of indices corresponding to the unit vectors contained in the rowspace of \mathbf{M} , denoted by $\text{Row}(\mathbf{M})$, so that $X = \{i : \mathbf{e}_i \in \text{Row}(\mathbf{M})\}$. We write $|X|$ to denote the cardinality of random variable X . Furthermore, we define random variables R and N to give the rank of \mathbf{M} and the number of rows in \mathbf{M} , respectively. The considered problem has been decomposed into the following two tasks:

- 1) Obtain the probability of recovering at least x source packets, provided that r out of the n received packets are linearly independent, for $x \leq r \leq k$. This is equivalent to finding the probability of $\text{Row}(\mathbf{M})$ containing at least x unit vectors, given \mathbf{M} has n rows and rank r . We denote this probability by $P(|X| \geq x | R = r, N = n)$.
- 2) Obtain the probability of recovering at least x source packets, provided that $n \geq x$ packets have been collected. We write $P(|X| \geq x | N = n)$ to refer to this probability.

Derivation of the probabilities $P(|X| \geq x | R = r, N = n)$ and $P(|X| \geq x | N = n)$ is the focus of the following section.

III. PROBABILITY ANALYSIS

The analysis presented in this section relies on the well-known Principle of Inclusion and Exclusion [9, Prop. 5.2.2], which is repeated below for clarity.

Lemma 1. Principle of inclusion and exclusion. *Given a set A , let f be a real valued function defined for all sets $S, J \subseteq A$. If $g(S) = \sum_{J: J \supseteq S} f(J)$ then $f(S) = \sum_{J: J \supseteq S} (-1)^{|J \setminus S|} g(J)$.*

For non-negative integers m and d , we denote by $\binom{m}{d}$ the binomial coefficient, which gives the number of d -element sets of an m -element set. The q -analog of the binomial coefficient, known as the *Gaussian binomial coefficient* and denoted by $\begin{bmatrix} m \\ d \end{bmatrix}_q$, enumerates all d -dimensional subspaces of an m -dimensional space over \mathbb{F}_q [9, p. 125].

Given \mathbf{M} has rank r , let $P(|X| = x | R = r, N = n)$ denote the probability of recovering exactly $x \leq r$ source packets or, equivalently, the probability of $\text{Row}(\mathbf{M})$ containing exactly $x \leq r$ unit vectors. The following theorem obtains an expression for $P(|X| = x | R = r, N = n)$, which is then used in the derivation of $P(|X| \geq x | R = r, N = n)$.

Theorem 1. *Given a random $n \times k$ matrix \mathbf{M} of rank r , the probability that the rowspace of \mathbf{M} contains exactly $x \leq r$ unit vectors is given by*

$$P(|X| = x | R = r, N = n) = \frac{\binom{k}{x}}{\begin{bmatrix} k \\ r \end{bmatrix}_q} \sum_{j=0}^{k-x} (-1)^j \binom{k-x}{j} \begin{bmatrix} k-x-j \\ r-x-j \end{bmatrix}_q. \quad (1)$$

Proof: For $S \subseteq J \subseteq \{1, \dots, k\}$, let $g(S)$ be the probability that $\{\mathbf{e}_i : i \in S\} \subseteq \text{Row}(\mathbf{M})$, that is, the probability that $S \subseteq X$. This is just the probability that $\text{Row}(\mathbf{M})$ contains a fixed $|S|$ -dimensional subspace, namely the space $V = \text{Span}\{\mathbf{e}_i : i \in S\}$. We see that, by considering the quotient space \mathbb{F}_q^k/V , there is a direct correspondence

between r -dimensional subspaces of \mathbb{F}_q^k containing V , and $(r - |S|)$ -dimensional subspaces of a $(k - |S|)$ -dimensional space. Hence, there are $\begin{bmatrix} k-|S| \\ r-|S| \end{bmatrix}_q$ r -dimensional subspaces of \mathbb{F}_q^k containing V . The probability that $\text{Row}(\mathbf{M})$ contains the space V is equal to

$$g(S) = \frac{\begin{bmatrix} k-|S| \\ r-|S| \end{bmatrix}_q}{\begin{bmatrix} k \\ r \end{bmatrix}_q} \quad (2)$$

where the denominator in (2) enumerates the r -dimensional subspaces of \mathbb{F}_q^k . Now, let $f(S)$ be the probability that $S = X$, that is, the probability that $\{\mathbf{e}_i : i \in S\} \subseteq \text{Row}(\mathbf{M})$ and $\mathbf{e}_i \notin \text{Row}(\mathbf{M})$ for $i \notin S$. It follows that $g(S) = \sum_{J \supseteq S} f(J)$. Invoking the Principle of Inclusion and Exclusion (Lemma 1) and using (2), we can write $f(S) = \sum_{J \supseteq S} (-1)^{|J \setminus S|} \cdot g(J)$ and expand it to

$$f(S) = \sum_{J \supseteq S} (-1)^{|J \setminus S|} \cdot \frac{\begin{bmatrix} k-|J| \\ r-|J| \end{bmatrix}_q}{\begin{bmatrix} k \\ r \end{bmatrix}_q} = \frac{1}{\begin{bmatrix} k \\ r \end{bmatrix}_q} \sum_{J' \subseteq \{1, \dots, k\} \setminus S} (-1)^{|J'|} \begin{bmatrix} k-|S|-|J'| \\ r-|S|-|J'| \end{bmatrix}_q \quad (3)$$

$$= \frac{1}{\begin{bmatrix} k \\ r \end{bmatrix}_q} \sum_{j=0}^{k-|S|} (-1)^j \binom{k-|S|}{j} \begin{bmatrix} k-|S|-j \\ r-|S|-j \end{bmatrix}_q \quad (4)$$

where (3) follows by setting $J' = J \setminus S$, and (4) follows since there are $\binom{k-|S|}{j}$ sets J' of size j . Considering that $f(S)$ is the probability that $X = S$, we can write

$$P(|X| = x | R = r, N = n) = \sum_{S: |S|=x} f(S) = \binom{k}{x} f(S') \quad (5)$$

where S' is any subset of $\{1, \dots, k\}$ of size x . The second equality in (5) holds since there are $\binom{k}{x}$ sets $S \subseteq \{1, \dots, k\}$ of size x . Substituting (4) in (5) gives the result. ■

Remark 1. Theorem 1 can be seen as a special case of [4, Proposition 6]. Whereas the proof in [4] uses elements of matroid theory, our paper proposes an alternative and more intuitive proof strategy.

Corollary 1. *Given a random $n \times k$ matrix \mathbf{M} of rank r , the probability that the rowspace of \mathbf{M} contains at least $x \leq r$ unit vectors is given by*

$$P(|X| \geq x | R = r, N = n) = \frac{1}{\begin{bmatrix} k \\ r \end{bmatrix}_q} \sum_{i=x}^r \binom{k}{i} \cdot \sum_{j=0}^{k-i} (-1)^j \binom{k-i}{j} \begin{bmatrix} k-i-j \\ r-i-j \end{bmatrix}_q. \quad (6)$$

Proof: By definition, $P(|X| \geq x | R = r, N = n)$ is equal to $\sum_{i=x}^r P(|X| = i | R = r, N = n)$. Substituting in (1) gives the result. ■

Note that, although \mathbf{M} is an $n \times k$ matrix, the probabilities in (1) and (6) hold for any value of $n \geq r$. Having obtained an expression for $P(|X| \geq x | R = r, N = n)$, we now proceed to the derivation of $P(|X| \geq x | N = n)$. This probability is denoted by $P_{\text{ns}}(|X| \geq x | N = n)$ and $P_s(|X| \geq x | N = n)$ for non-systematic and systematic RLNC, respectively. Expressions for each case are derived in the following two propositions.

Proposition 1. *If a receiving node collects n random linear combinations of k source packets, the probability that at least $x \leq k$ source packets will be recovered is*

$$P_{\text{ns}}(|X| \geq x | N = n) = \frac{1}{q^{nk}} \cdot \sum_{r=x}^{\min(n,k)} \left(\sum_{i=x}^r \binom{k}{i} \sum_{j=0}^{k-i} (-1)^j \binom{k-i}{j} \begin{bmatrix} k-i-j \\ r-i-j \end{bmatrix}_q \right) \prod_{\ell=0}^{r-1} (q^n - q^\ell). \quad (7)$$

Proof: Let $P(R = r | N = n)$ denote the probability that the $n \times k$ matrix \mathbf{M} has rank r . This is equivalent to the probability that r out of the n collected packets are linearly independent. The probability that at least x of the k source packets will be recovered can be obtained from

$$P_{\text{ns}}(|X| \geq x | N = n) = \sum_{r=x}^{\min(n,k)} P(R = r | N = n) P(|X| \geq x | R = r, N = n). \quad (8)$$

The probability $P(R = r | N = n)$ is equal to [10, Sec. II.A]

$$P(R = r | N = n) = \frac{1}{q^{nk}} \begin{bmatrix} n \\ r \end{bmatrix}_q \prod_{\ell=0}^{r-1} (q^k - q^\ell). \quad (9)$$

Substituting (6) and (9) into (8) and taking into account that

$$\frac{\begin{bmatrix} n \\ r \end{bmatrix}_q}{\begin{bmatrix} k \\ r \end{bmatrix}_q} \prod_{\ell=0}^{r-1} (q^k - q^\ell) = \prod_{\ell=0}^{r-1} (q^n - q^\ell) \quad (10)$$

leads to (7). \blacksquare

Proposition 2. *If k source packets and $n_{\text{T}} - k$ random linear combinations of those k source packets are transmitted over single-hop links, the probability that a receiving node will recover at least $x \leq k$ source packets from $n \leq n_{\text{T}}$ received packets is*

$$P_{\text{s}}(|X| \geq x | N = n) = \frac{1}{\binom{n_{\text{T}}}{n}} \cdot \sum_{r=x}^{\min(n,k)} \sum_{h=h_{\min}}^r \left(\binom{k}{h} \binom{n_{\text{T}}-k}{n-h} q^{-(n-h)(k-h)} \prod_{\ell=0}^{r-h-1} (q^{n-h} - q^\ell) \cdot \sum_{i=x_{\min}}^{r-h} \binom{k-h}{i} \sum_{j=0}^{k-h-i} (-1)^j \binom{k-h-i}{j} \begin{bmatrix} k-h-i-j \\ r-h-i-j \end{bmatrix}_q \right) \quad (11)$$

where $h_{\min} = \max(0, n - n_{\text{T}} + k)$ and $x_{\min} = \max(0, x - h)$.

Proof: Let us assume that some or none of the k transmitted source packets have been received and let $X' \subseteq X$ be the set of indices of the remaining source packets that can be recovered from the received coded packets. If n' of the $n_{\text{T}} - k$ coded packets have been received and k' source packets remain to be recovered, the respective coding vectors will form an $n' \times k'$ random matrix \mathbf{M}' . The probability that $r' \leq \min(k', n')$ coding vectors are linearly independent and at least $x' \leq r'$ source packets can be recovered is given by

$$P(|X'| \geq x', R' = r' | N' = n') = P(R' = r' | N' = n') P(|X'| \geq x' | R' = r', N' = n')$$

where the two terms of the product can be obtained from (9) and (6), respectively. The random variables N' and R'

denote the number of received coded packets and the rank of matrix \mathbf{M}' , respectively. If n of the n_{T} transmitted packets are received, the probability that h of them are source packets and the remaining $n - h$ are coded packets is

$$P(N' = n - h | N = n) = \binom{k}{h} \binom{n_{\text{T}}-k}{n-h} / \binom{n_{\text{T}}}{n}. \quad (12)$$

The coding vectors of the n received packets compose a matrix of rank r , based on which x or more source packets can be recovered when h of the n received packets are source packets. Parameters x' , r' , k' and n' , which are concerned with the received *coded* packets only, can be written as $x - h$, $r - h$, $k - h$ and $n - h$, respectively. The probability of recovering at least x source packets for all valid values of r and h is

$$P_{\text{s}}(|X| \geq x | N = n) = \sum_{r=x}^{\min(n,k)} \sum_{h=h_{\min}}^r P(N' = n - h | N = n) \cdot P(|X'| \geq \max(0, x - h), R' = r - h | N' = n - h) \quad (13)$$

which expands into (11). Note that $\max(0, x - h)$ ensures that the value of $|X'|$ is a non-negative integer when $h > x$. \blacksquare

Remark 2. In systematic RLNC, if the receiving node attempts to recover source packets as soon as the transmission is initiated, i.e., $n_{\text{T}} \leq k$, at least x source packets will certainly be recovered when $n \geq x$ source packets are received, that is,

$$P_{\text{s}}(|X| \geq x | N = n) = \begin{cases} 1, & \text{if } n_{\text{T}} \leq k \text{ and } x \leq n \\ 0, & \text{if } n_{\text{T}} \leq k \text{ and } x > n. \end{cases} \quad (14)$$

IV. RESULTS AND DISCUSSION

In order to demonstrate the exactness of the derived expressions, simulations that generated 60000 realisations of an $n \times k$ random matrix \mathbf{M} over \mathbb{F}_2 were carried out for $n = 1, \dots, 30$ and $k = 20$. In each case, matrix \mathbf{M} was converted into reduced row echelon form using Gaussian elimination. Then, the rows that correspond to unit vectors \mathbf{e}_i , which represent recoverable source packets, were counted and averaged over all realisations. Fig. 1(a) and Fig. 1(b) show that measurements obtained through simulations match the calculations obtained from (7) and (11) for non-systematic RLNC and systematic RLNC, respectively. In general, simulation results match analytical predictions for any finite field \mathbb{F}_q of order $q \geq 2$.

Fig. 2 considers the simple case of RLNC transmission over a broadcast erasure channel. If the transmission of n_{T} packets is modeled as a sequence of n_{T} Bernoulli trials whereby ε signifies the probability that a transmitted packet will be erased, the probability that a receiving node shall recover *at least* x of the k source packets can be expressed as

$$P(|X| \geq x) = \sum_{n=x}^{n_{\text{T}}} \binom{n_{\text{T}}}{n} (1 - \varepsilon)^n \varepsilon^{n_{\text{T}}-n} P(|X| \geq x | N = n). \quad (15)$$

The probability $P(|X| \geq x | N = n)$ is equal to (7) for non-systematic RLNC and (11) or (14), depending on the value of n_{T} , for systematic RLNC.

Fig. 2(a) focuses on non-systematic RLNC and depicts $P(|X| \geq x)$ in terms of n_{T} for $x \in \{2, 4, 10, 16, 20\}$ when $k = 20$, and for $x \in \{3, 6, 15, 24, 30\}$ when $k = 30$. Results

have been obtained for $q \in \{2, 8\}$ and $\varepsilon = 0.2$. For $q = 2$, the transmission of only a few additional coded packets can increase the fraction of the recovered message from at least $x/k = 0.1$ to $x/k = 1$. However, for q as low as 8, the range of n_T values for which a receiving node will proceed from recovering a small portion of the transmitted message to recovering the whole message gets very narrow. Furthermore, for $q = 2$, segmentation of the message into $k = 20$ source packets permits a receiving node to recover the same fraction (x/k) of the message with a higher probability than dividing the same message into $k = 30$ source packets.

Systematic RLNC is considered in Fig. 2(b). Besides the reduced decoding complexity [11], we observe that systematic RLNC enables a receiving node to gradually reveal an increasingly larger portion of the message as more packets are transmitted. However, a large number of source packets or a high order finite field impairs the progressive recovery of the message for $n_T > k$. This is because source packets are transmitted for $n_T \leq k$ but coded packets are sent for $n_T > k$; the decoding behaviour of a receiving node changes at $n_T = k$ and causes a change in the slope of $P(|X| \geq x)$ for $x/k = 0.8$.

The results show that, if information-theoretic security is required, non-systematic RLNC over finite fields of size 8 or larger can be used to segment each message into a large number of source packets. The number of transmitted packets can then be adjusted to the channel conditions to achieve a balance between the probability of legitimate nodes reconstructing the message and the probability of eavesdroppers being unable to decode even a portion of the message. If the objective of the system is to maximize the number of nodes that will recover at least a large part of a message, systematic RLNC over small finite fields can be used to divide data into source packets. If the receiving nodes do not suffer from limited computational capabilities, the size of the finite field can be increased to improve the probability of recovering the entire message.

V. CONCLUSIONS

This paper derived exact expressions for the probability of decoding a fraction of a source message upon reception of an arbitrary number of network-coded packets. Results unveiled the potential of non-systematic network coding in offering weak information-theoretic security, even when operations are over small finite fields. On the other hand, systematic network coding allows for the progressive recovery of the source message as the number of received packets increases, especially when the size of the finite field is small.

VI. ACKNOWLEDGMENTS

Jessica Claridge has been supported by an EPSRC PhD studentship. Both authors appreciate the support of the COST Action IC1104 and thank Simon R. Blackburn for his advice.

REFERENCES

[1] T. Ho, M. Médard, R. Koetter, D. R. Karger, M. Effros, J. Shi, and B. Leong, "A random linear network coding approach to multicast," *IEEE Trans. Inf. Theory*, vol. 52, no. 10, pp. 4413–4430, Oct. 2006.

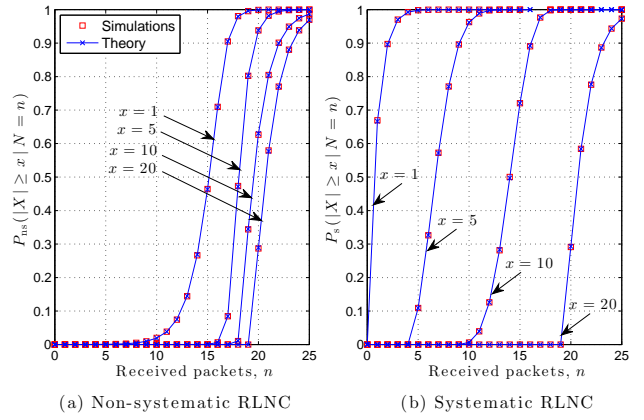


Fig. 1: Simulation results and theoretical values for (a) non-systematic RLNC and (b) systematic RLNC. The probability of recovering at least x source packets has been plotted for $q = 2$, $k = 20$, $x = 1, 5, 10, 20$ and $n_T = 30$.

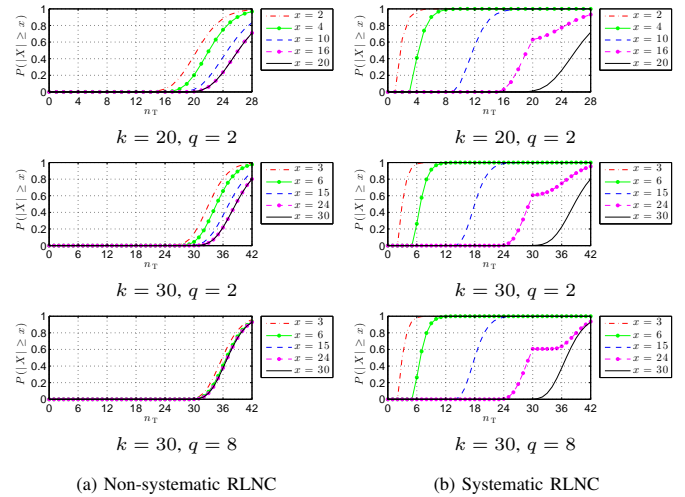


Fig. 2: Depiction of the probability of recovering at least x source packets when n_T packets have been transmitted over a packet erasure channel with $\varepsilon = 0.2$ using (a) non-systematic RLNC and (b) systematic RLNC.

[2] O. Trullols-Cruces, J. Barcelo-Ordinas, and M. Fiore, "Exact decoding probability under random linear network coding," *IEEE Commun. Lett.*, vol. 15, no. 1, pp. 67–69, Jan. 2011.

[3] Z. Yan, H. Xie, and B. W. Suter, "Rank deficient decoding of linear network coding," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing*, Vancouver, BC, May 2013, pp. 5080–5084.

[4] M. Gadouleau and A. Goupil, "A matroid framework for noncoherent random network communications," *IEEE Trans. Inf. Theory*, vol. 57, no. 2, pp. 1031–1045, Feb. 2011.

[5] K. Bhattad and K. R. Narayanan, "Weakly secure network coding," in *Proc. 1st Workshop on Network Coding, Theory and Applications*, Riva Del Garda, Italy, Apr. 2005.

[6] L. Lima, M. Médard, and J. Barros, "Random linear network coding: A free cipher?" in *Proc. IEEE Int. Symp. Inform. Theory*, Nice, France, Jun. 2007, pp. 546–550.

[7] N. Cai and R. W. Yeung, "Secure network coding," in *Proc. IEEE Int. Symp. on Inform. Theory*, Lausanne, Switzerland, Jun. 2002, p. 323.

[8] B. Shraider and N. M. Jones, "Systematic wireless network coding," in *Proc. IEEE Military Commun. Conf.*, Boston, MA, Oct. 2009.

[9] P. J. Cameron, *Combinatorics: Topics, techniques, algorithms*. Cambridge University Press, 1994.

[10] M. Gadouleau and Z. Yan, "Constant-rank codes and their connection to constant-dimension codes," *IEEE Trans. Inf. Theory*, vol. 56, no. 7, pp. 3207–3216, Jul. 2010.

[11] D. E. Lucani, M. Médard, and M. Stojanovic, "On coding for delay – Network coding for time-division duplexing," *IEEE Trans. Inf. Theory*, vol. 58, no. 4, pp. 2330–2348, Apr. 2012.