

AN EXPLORATION OF NON-ASYMPTOTIC LOW-DENSITY, PARITY CHECK ERASURE CODES FOR WIDE-AREA STORAGE APPLICATIONS

JAMES S. PLANK

*Department of Computer Science, University of Tennessee
Knoxville, Tennessee 37996, United States*

and

MICHAEL G. THOMASON

*Department of Computer Science, University of Tennessee
Knoxville, Tennessee 37996, United States*

Received November 2006

Revised January 2007

Communicated by Guest Editors

ABSTRACT

As peer-to-peer and widely distributed storage systems proliferate, the need to perform efficient erasure coding, instead of replication, is crucial to performance and efficiency. *Low-Density Parity-Check (LDPC)* codes have arisen as alternatives to standard erasure codes, such as Reed-Solomon codes, trading off vastly improved decoding performance for inefficiencies in the amount of data that must be acquired to perform decoding. The scores of papers written on LDPC codes typically analyze their *collective* and *asymptotic* behavior. Unfortunately, their practical application requires the generation and analysis of *individual* codes for *finite* systems.

This paper attempts to illuminate the practical considerations of LDPC codes for peer-to-peer and distributed storage systems. The three main types of LDPC codes are detailed, and a huge variety of codes are generated, then analyzed using simulation. This analysis focuses on the performance of *individual* codes for *finite* systems, and addresses several important heretofore unanswered questions about employing LDPC codes in real-world systems.

Keywords: Erasure codes, storage systems, fault-tolerance, peer-to-peer, low-density parity (LDPC) codes, Tornado codes

1. Introduction

Wide-area file systems typically employ replication to improve both the performance and fault-tolerance of file access. Specifically, consider a file system composed of storage nodes distributed across the wide area, and consider multiple clients, also distributed across the wide area, who desire to access a large file. The standard strategy that file systems employ is one where the file is partitioned into k blocks of a fixed size, and these blocks are replicated and distributed throughout the system. Such a scenario is depicted in Fig. 1, where a single file is partitioned into eight

blocks numbered one through eight, and each block is replicated four of eight storage servers. Three separate clients are shown accessing the file in its entirety by attempting to download each of the eight blocks from a nearby server.

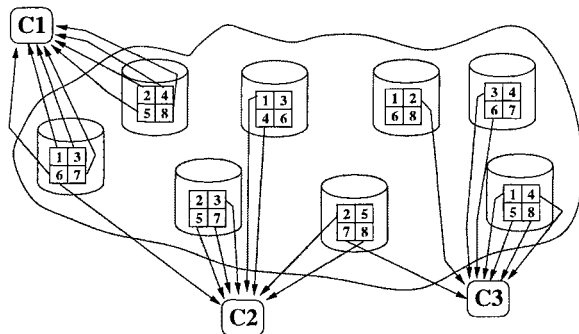


Fig. 1. A widely distributed file system hosting a file partitioned into eight blocks, each block replicated four times. Three clients are depicted accessing the file from different network locations.

Replicated systems such as these provide both fault-tolerance and improved performance over non-replicated storage systems. However, the costs are high. First, each block must be replicated m times to tolerate the failure of any $m - 1$ servers. Second, clients must find close copies of each of the file's blocks, which can be difficult, and the failure or slow access of any particular block can hold up the performance of the entire file's access [1].

Erasure encoding schemes improve both the fault-tolerance and downloading performance of replicated systems [31,33,7]. For example, with Reed-Solomon erasure encoding, instead of storing the blocks of the files themselves, $k + m$ encodings of the blocks are calculated, and these are stored instead. Now the clients need only download *any* k blocks, and from these, the k blocks of the file may be calculated. Such a scenario is depicted in Fig. 2, where 32 encoding blocks, labeled **A** through **Z** and **a** through **f** are stored, and the clients need only access the eight closest blocks to compute the file.

Reed-Solomon coding has been employed effectively in distributed storage systems [12,22], and in related functionalities such as fault-tolerant data structures [13], disk arrays [3] and checkpointing systems [18]. However, it is not without costs. Specifically, encoding involves breaking each block into words, and each word is calculated as the dot product of two length- k vectors under Galois Field arithmetic, which is more expensive than regular arithmetic. Decoding involves the inversion of an $k \times k$ matrix, and then each of the file's blocks is calculated with dot products as in encoding. Thus, as k grows, the costs of Reed-Solomon coding induce too much overhead [4].

In 1997, Luby *et al* published a landmark paper detailing a coding technique that thrives where Reed-Solomon coding fails [16]. Their codes, later termed "Tornado Codes," calculate m coding blocks from the k file blocks in linear time using only cheap exclusive-or (parity) operations. Decoding is also performed in linear time

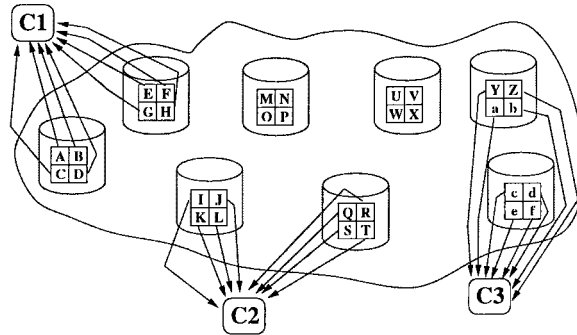


Fig. 2. The same system as Fig. 1, employing Reed-Solomon coding instead of replication. Again the file is partitioned into eight blocks, but now 32 encoding blocks are stored so that clients may employ *any* eight blocks to calculate the file.

using parity; however, rather than requiring any k blocks for decoding as in Reed-Solomon coding, they require fk blocks, where f is an *overhead factor* that is greater than one, but approaches one as k approaches infinity. A content-distribution system called “Digital Fountain” was built on Tornado Code technology, and in 1998 its authors formed a company of the same name [9].

Tornado Codes are instances of a class of codes called *Low-Density Parity-Check (LDPC)* codes, which have a long history dating back to the 60’s [10], but have received renewed attention since the 1997 paper. Since 1998, the research on LDPC codes has taken two paths – Academic research has resulted in many publications about LDPC codes [24,32,28,23], and Digital Fountain has both published papers [5,14,25] and received patents on various aspects of coding techniques.*

LDPC codes are based on bipartite graphs, which are employed to define codes based solely on parity operations. Nearly all published research on LDPC codes has had the same mission – to define codes that approach “channel capacity” asymptotically. In other words, they define codes where the overhead factor, f , approaches one as k approaches infinity. It has been shown [16] that codes based on regular graphs – those where each node has a constant incoming and outgoing cardinality – do not have this property. Instead, the “best” codes are based on randomly generated irregular graphs. A class of irregular graphs is defined, based on probability distributions of node cardinalities, and then properties are proven about the ensemble characteristics of this class. The challenge then becomes to design probability distributions that generate classes of graphs that approach channel capacity. Hundreds of such distributions have been published in the literature and on the web (see Table 1 for 80 examples).

Although the probabilistic method [2] with random graphs leads to powerful characterizations of LDPC ensembles, generating individual graphs from these probability distributions is a non-asymptotic, non-ensemble activity. In other words, while the properties of infinite collections of infinitely sized graphs is known, and

*U.S. Patents #6,073,250, #6,081,909, #6,163,870, #6,195,777, #6,320,520 and #6,373,406.

while there has been some work in finite-length analysis [8,19], the properties of individual, finite-sized graphs, especially for small values of k , have not been explored to date. Moreover, these properties have profound practical consequences.

Addressing aspects of these practical consequences is the goal of this paper. Specifically, we detail how three types of LDPC graphs are generated from given probability distributions and describe a method of simulation to analyze individual LDPC graphs. Then we generate a wide variety of LDPC graphs and analyze their performance in order to answer the following five practical questions:

- (i) What kind of overhead factors (f) can we expect for LDPC codes for small and large values of k ?
- (ii) Are the three types of codes equivalent, or do they perform differently?
- (iii) How do the published distributions fare in producing good codes for finite values of k ?
- (iv) Is there a great deal of random variation in code generation from a given probability distribution?
- (v) What effect does cascading have on the Simple codes?

In answering each question, we pose a challenge to the community to perform research that helps systems researchers make use of these codes. It is our hope that this paper will spur researchers on LDPC codes to include analyses of the non-asymptotic properties of individual graphs based on their research.

2. Three Types of LDPC Codes

Three distinct types of LDPC codes have been described in the academic literature. All are based on bipartite graphs that are randomly generated from probability distributions. We describe them briefly here. For detailed presentations on these codes, and standard encoding/decoding algorithms, please see other sources [16,11,27,23,32].

The graphs have $L+R$ nodes, partitioned into two sets – the *left* nodes, l_1, \dots, l_L , and the *right* nodes, r_1, \dots, r_R . Edges only connect left nodes to right nodes. A class of graphs G is defined by two probability distributions, Λ and P . These are vectors composed of elements $\Lambda_1, \Lambda_2, \dots$ and P_1, P_2, \dots such that $\sum_i \Lambda_i = 1$ and $\sum_i P_i = 1$. Let g be a graph in G . Λ_i is the probability that a left node in g has exactly i outgoing edges, and similarly, P_i is the probability that a right node in g has exactly i incoming edges.[†]

Given L , R , Λ and P , generating a graph g is in theory a straightforward task [16]. We describe our generation algorithm in section below. For this section, it suffices that given these four inputs, we can generate bipartite graphs from them.

To describe the codes below, we assume that we have k equal-sized blocks of data, which we wish to encode into $k+m$ equal-sized blocks, which we will distribute on the network. The nodes of LDPC graphs hold such blocks of data, and therefore we will use the term “node” and “block” interchangeably. Nodes can either initialize

[†]An alternate and more popular definition is to define probability distributions of the edges rather than the nodes using two vectors λ and ρ . The definitions are interchangeable since (Λ, P) may be converted easily to (λ, ρ) and vice versa.

their block's values from data, or they may calculate them from other blocks. The only operation used for these calculations is parity, as is common in RAID Level 5 disk arrays [6]. Each code generation method uses its graph to define an encoding of the k data blocks into $k + m$ blocks that are distributed on the network.

To decode, we assume that we download the blocks in order of proximity. Each time we download a block, we perform the standard decoding operation of removing the block's node and edges from the code's graph. If decoding yields the value of an undownloaded block, then that is equivalent to that node being downloaded, and we continue removing nodes and edges from the graph. If a block is downloaded and its node has already been removed from the graph, it is discarded (even though this useless download has affected performance). When all nodes of the graph are gone, the data has been recovered, and we do not need to download any more blocks.

2.1. Simple Codes

With Simple codes, $L = k$ and $R = m$. Each left node l_i holds the i -th data block, and each right node r_i is calculated to be the exclusive-or of all the left nodes that are connected to it. A very simple example is depicted in Fig. 3(a).

Simple codes can *cascade*, by employing $d > 1$ levels of bipartite graphs, g_1, \dots, g_d , where the right nodes of g_i are also the left nodes of g_{i+1} . The graph of level 1 has $L = k$, and those nodes contain the k data blocks. The remaining blocks of the encoding are right-hand nodes of the d graphs. Thus, $m = \sum_{i=1}^d R_i$. A simple three-level cascaded Simple code is depicted in Fig. 3(b).

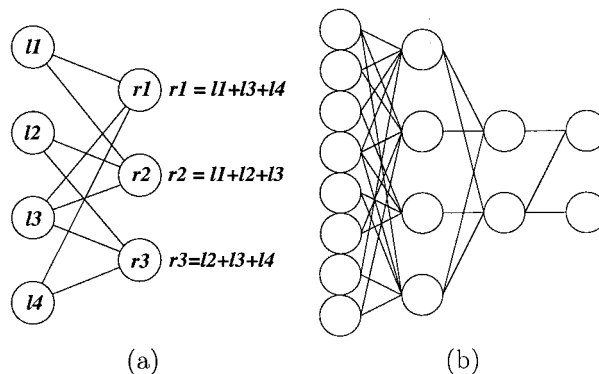


Fig. 3. (a) Example 1-level Simple code for $k = 4$, $m = 3$. (b) Example 3-level Simple code for $k = 8$, $m = 8$.

Encoding and decoding of both regular and cascading Simple codes are straightforward operations and are both linear time operations in the number of edges in the graph.

2.2. Gallager Codes

Gallager codes were introduced in the early 1960's [10] and are *unsystematic*,

as the $k + m$ blocks stored do not have to include the original k data blocks. The other two codes in this paper are systematic. Gallager codes employ a different kind of graph called a *Tanner* graph, where $L = k + m$, and $R = m$. The first step of creating a Gallager code is to use g to generate a $(k + m) \times n$ matrix M . This is employed to calculate the $k + m$ encoding blocks from the original k data blocks. These blocks are stored in the left nodes of g . The right nodes of g do not hold data, but instead are *constraint* nodes — each r_i has the property (guaranteed by the generation of M) that the exclusive-or of all nodes incident to it is zero. A simple Gallager code is depicted in Fig. 4(a).

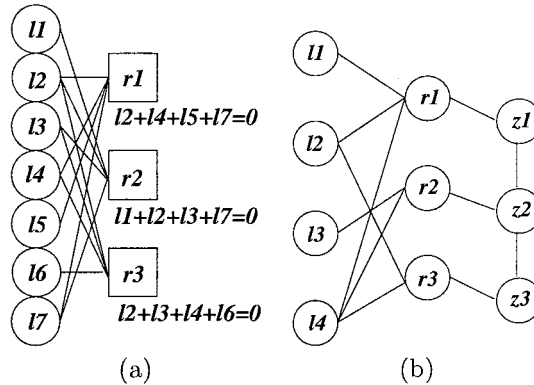


Fig. 4. (a) Example Gallager code for $k = 4$, $m = 3$. Note that the right nodes define constraints between the left nodes, and do not store encoding blocks. (b) Example IRA code for $k = 4$, $m = 3$. The left and accumulator nodes are stored as the encoding blocks. The right nodes are just used for calculations.

Encoding is an expensive operation, involving the generation of M , and calculation of the encoding blocks. Fortunately, if the graph is low density (i.e. the average cardinality of the nodes is small), M is a sparse matrix, and its generation and use for encoding and decoding is not as expensive as a dense matrix (as is the case with Reed-Solomon coding). Decoding is linear in the number of edges in the graph. Fortunately, M only needs to be generated once per graph, and then it may be used for all encoding/decoding operations.

2.3. IRA Codes

Irregular Repeat-Accumulate (IRA) Codes are systematic codes, as $L = k$ and $R = m$, and the information blocks are stored in the left nodes. However, an extra set of m nodes, z_1, \dots, z_m , are added to the graph in the following way. Each node r_i has an edge to z_i . Additionally, each node z_i has an edge to z_{i+1} , for $i < m$. These extra nodes are called *accumulator* nodes. For encoding, only blocks in the left and accumulator nodes are stored — the nodes in R are simply used to calculate the encodings and decodings, and these calculations proceed exactly as in the Simple codes. An example IRA graph is depicted in Fig. 4(b).

3. Asymptotic Properties of Codes

All three classes of LDPC codes have undergone asymptotic analyses that proceed as follows. A rate $\mathcal{R} = \frac{k}{k+m}$ is selected, and then Λ and P vectors are designed. From these, it may be proven that graphs generated from the distributions in Λ and P can asymptotically *achieve capacity*. In other words, they may be successfully decoded with fk downloaded blocks, where f approaches 1 from above as k approaches ∞ .

Unfortunately, in the real world, developers of wide-area storage systems cannot break up their data into infinitely many pieces. Limitations on the number of physical devices, plus the fact that small blocks of data do not transmit as efficiently as large blocks, dictate that k may range from single digits into the thousands. Therefore, a major question about LDPC codes (addressed by **Question 1** above) is how well they perform when k is in this range.

4. Assessing Performance

Our experimental methodology is as follows. For each of the three LDPC codes, we have written a program to randomly generate a bipartite graph g that defines an instance of the code, given k, m, Λ, P , and a seed for a random number generator. The generation follows the methodology sketched in [16]:

Table 1. The 80 published probability distributions (Λ and P) used to generate codes.

| Name | Source | # of Codes | Λ_{max} | P_{max} | Developed for | Rate: $[\frac{1}{3}, \frac{1}{2}, \frac{2}{3}]$ |
|------|--------|------------|-----------------|-----------|---------------|---|
| L97A | [16] | 2 | 1,048,577 | 30,050 | Simple | [0,1,1] |
| L97B | [16] | 8 | 8-47 | 16-28 | Simple | [0,4,4] |
| S99 | [26] | 19 | 2-3298 | 6-13 | Gallager | [4,7,8] |
| SS00 | [28] | 3 | 9-12 | 7-16 | Gallager | [0,3,0] |
| M00 | [17] | 14 | 2-20 | 3-8 | IRA | [0,6,8] |
| WK03 | [32] | 6 | 11-50 | 8-11 | Gallager* | [0,6,0] |
| RU03 | [23] | 2 | 8-13 | 6-7 | Gallager | [0,2,0] |
| R03 | [24] | 8 | 100 | 8 | IRA* | [0,8,0] |
| U03 | [30] | 22 | 6-100 | 6-19 | Gallager | [6,9,7] |

For each left node l_i , its number of outgoing edges ξ_i is chosen randomly from Λ , and for each right node r_i , its number of incoming edges ι_i is chosen randomly from P . This yields two total numbers of edges, $T_L = \sum_{i=1}^L \xi_i$ and $T_R = \sum_{i=1}^R \iota_i$ which may well differ by $D > 0$. Suppose $T_L > T_R$. To rectify this difference, we select a “shift” factor s such that $0 \leq s \leq 1$. Then we subtract sD edges randomly from the left nodes (modifying each ξ_i accordingly), and add $(1-s)D$ edges randomly to the right nodes (modifying each ι_i accordingly). This yields a total of T total edges coming from the left nodes and going to the right nodes.

Now, we define a new graph g' with T left nodes, T right nodes and a random matching of T edges between them. We use g' to define g , by having the first ξ_1 edges of g' define the edges in g coming from l_1 . The next ξ_2 edges in g' define the edges coming from l_2 , and so on. The right edges of g are defined similarly by the right edges of g' and ι_i .

At the end of this process, there is one potential problem with g — there may be duplicate edges between two nodes, which serve no useful purpose in coding or

decoding. We deal with this problem by deleting duplicate edges. An alternative method is to swap edges between nodes until no duplicate edges exist. We compared these two methods and found that neither outperformed the other, so we selected the edge deletion method since it is more efficient.

We evaluate each random graph by performing a Monte Carlo simulation of over 1000 random downloads, and calculating the average number of blocks required to successfully reconstruct the data. This is reported as the overhead factor f above. In other words, if $k = 100$, $m = 100$, and our simulation reports that $f = 1.10$, then on average, 110 random blocks of the 200 total blocks are required to reconstruct the 100 original blocks of data from the graph in question.

Theoretical work on LDPC codes typically calculates the *percentage of capacity* of the code, which is $\frac{1}{f}100\%$. We believe that for storage applications, the overhead factor is a better metric, since it quantifies how many block downloads are needed on average to acquire a file.

5. Experiments

5.1. Code Generation

The theoretical work on LDPC codes gives little insight into how the Λ and P vectors that they design will perform for smaller values of k . Therefore we have performed a rather wide exploration of LDPC code generation. First, we have employed 80 different sets of Λ and P from published papers on asymptotic codes. We call the codes so generated *published* codes. These are listed in Table 1, along with the codes and rates for which they were designed. The WK03 distributions are for Gallager codes on AWGN (Additive White Gaussian Noise) channels, and the R03 distributions are for IRA codes on AWGN and binary symmetric channels. In other words, neither is designed for the type of storage applications for which we employ them. We included the former as a curiosity and discovered that they performed very well. We included the latter because distributions for IRA codes are scarce.

Second, we have written a program that generates random Λ and P vectors, determines the ten best pairs that minimize f , and then goes through a process of picking random Λ 's for the ten best P 's and picking random P 's for the ten best Λ 's. This process is repeated, and the ten best Λ/P pairs are retained for subsequent iterations. Such a methodology is suggested by Luby *et al* [15]. We call the codes generated from this technique *Monte Carlo* codes.

Third, we observed that picking codes from some probability distributions resulted in codes with an extremely wide range of overhead factors (see section below). Thus, our third mode of attack was to take the best performing instances of the published and Monte Carlo codes, and use their left and right node cardinalities to define new Λ 's and P 's. For example, the Simple code in Fig. 3(a) can be generated from any number of probability distributions. However, it defines a probability distribution where $\Lambda = \langle 0, 0.75, 0.25 \rangle$ and $P = \langle 0, 0, 1 \rangle$. These new Λ 's and P 's may then be employed to generate new codes. We call the codes so generated *derived* codes.

5.2. Tests

The range of potential tests to conduct is colossal. As such, we limited it in the following way. We focus on three rates: $\mathcal{R} \in \{\frac{1}{3}, \frac{1}{2}, \frac{2}{3}\}$. In other words, $m = 2k$, $m = k$, and $m = \frac{k}{2}$. These are the rates most studied in the literature. For each of these rates, we generated the three types of codes from each of the 80 published distributions for all even k between 2 and 150, and for $k \in \{250, 500, 1250, 2500, 5000, 12500, 25000, 50000, 125000\}$ †. For Simple codes, we tested cascading levels from one to six.

For Monte Carlo codes, we tested all three codes with all three rates for even $k \leq 50$. As shown in section below, this code generation method is only useful for small k .

Finally, for each value of k , we used distributions derived the best current codes for all three coding methods (and all six cascading levels of Simple codes) to generate codes for the ten nearest values of k with the same rate. The hope is that good codes for one value of k can be employed to generate good codes for nearby values of k .

In sum, this makes for over 100,000 different data points, each of which was repeated with over 100 different random number seeds. The optimal code and overhead factor for each data point was recorded and the data is digested in the following section.

6. Results

Our computational engine is composed of 160 machines (Sun workstations running Solaris, Dell Pentium workstations running Linux, and a Macintosh PowerBook running OSX) which ran tests continuously for many months. We organize our results by answering each of the questions presented in Section above.

6.1. Question 1

What kind of overhead factors can we expect for LDPC codes for small and large values of k ?

All of our data is summarized in Fig. 5. For each value of k and m , the coding and generation method that produces the smallest overhead factor is plotted.

All three curves of Fig. 5 follow the same pattern. The overhead factor starts at 1 when $m = 1$ or $k = 1$, and the Simple codes become replication/parity codes with perfect performance. Then the factor increases with k until k reaches roughly twenty at which point it levels out until k increases to roughly 100. At that point, the factor starts to decrease as k increases, and it appears that it indeed goes to one as k gets infinitely large.

Although we only test three rates, it certainly appears that the overhead factor grows as the rate approaches zero. This is intuitive. At one end, any code with a rate of one will have an overhead factor of one. At the other, consider a one-level Simple code with $k = 3$ and $m = \infty$. There are only seven combinations of the left

†One exception is $k = 125000$ for $\mathcal{R} = \frac{1}{3}$, due to the fact that these graphs often exceeded the physical memory of our machines.

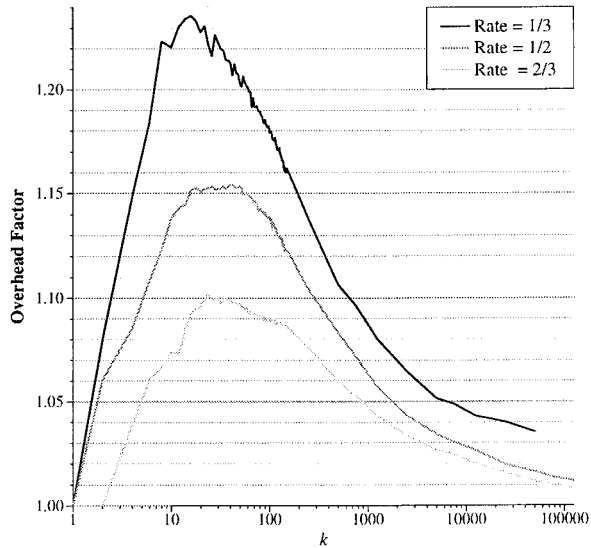


Fig. 5. The best codes for all generation methods for $1 \leq k \leq 125,000$, and $\mathcal{R} = \frac{1}{2}, \frac{2}{3}$.

nodes to which a right node may be connected. Therefore, the right nodes will be partitioned into at most $\frac{k}{7}$ groups, where each node in the group is equivalent. In other words, any download sequence that contains more than one block from a node group will result in overhead. Clearly, this argues for a higher overhead factor.

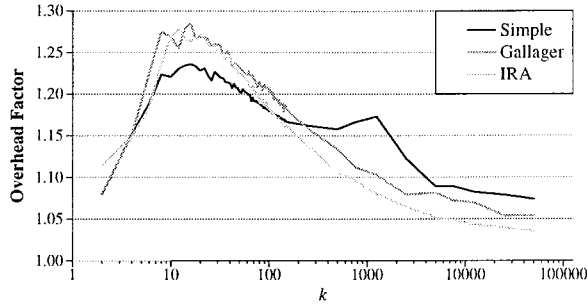
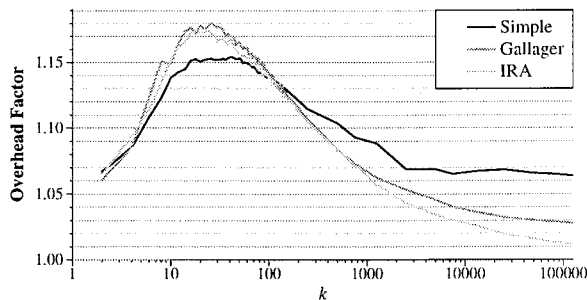
Challenge for the Community: The shape of the curves in Fig. 5 suggests that there is a lower bound for overhead factor as a function of k and m (or alternatively as a function of k and \mathcal{R}). It is a challenge to the theoretical community to quantify this lower bound for *finite* values of k and m , and then to specify exact methods for generating optimal or near optimal codes.

6.2. Question 2

Are the three types of codes equivalent, or do they perform differently?

They perform differently. Figs 6 - 8 show the best performing of the three different codes for the three rates. All three show a similar pattern – for small values of k , Simple codes perform the best. However, when k roughly equals 100, the IRA codes start to outperform the others, and the Gallager codes start to outperform the Simple codes. This trend continues to the maximum values of k .

Unfortunately, since the theoretical work on LDPC codes describes only asymptotic properties, little insight can be given as to why this pattern occurs. One curious point is the relationship between one-level Simple codes and Gallager codes. It is a trivial matter to convert a one-level Simple code into an equivalent Gallager code by adding m left nodes, l_{k+1}, \dots, l_{k+m} to the Simple graph, and m edges of the form (l_{k+i}, r_i) for $1 \leq i \leq m$. An example is the Simple code in Fig. 9(a),


 Fig. 6. Comparing methods, $\mathcal{R} = \frac{1}{3}$

 Fig. 7. Comparing methods, $\mathcal{R} = \frac{1}{2}$

which is equivalent to the Gallager code in Fig. 9(b). Both have overhead factors of 1.11. This fact would seem to imply that overhead factors for one-level Simple codes would be similar to, or worse than Gallager codes. However, when $k < 50$, the one-level Simple codes vastly outperform the others; the Gallager codes perform the worst. A clue to this behavior can be seen in Fig. 9(c). This is a Gallager code whose nodes have the same cardinalities as the code in Fig. 9(b), and thus would be generated by the same values of Λ and P . However, its overhead factor is 1.21!

To hammer this point home further, we performed the same conversion on a Simple graph where $k = m = 20$, and the overhead factor is 1.16. The node cardinalities of the equivalent Gallager graph were then used to generate values of Λ and P , which in turn were used to generate 500 new Gallager graphs with the exact same node cardinalities. The *minimum* overhead factor of these graphs was 1.31 (the average was 1.45, and the maximum was 1.58). What this suggests is that for smaller graphs, perhaps Λ and P need to be augmented with some other metric so that optimal codes can be generated easily.

Challenge to the community: A rigorous comparison of the practical utility of the three coding methods needs to be performed. In particular, a computation-

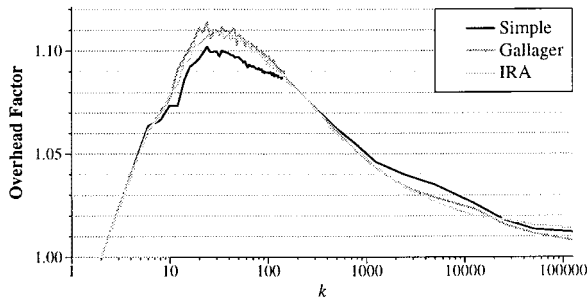


Fig. 8. Comparing methods, $\mathcal{R} = \frac{2}{3}$

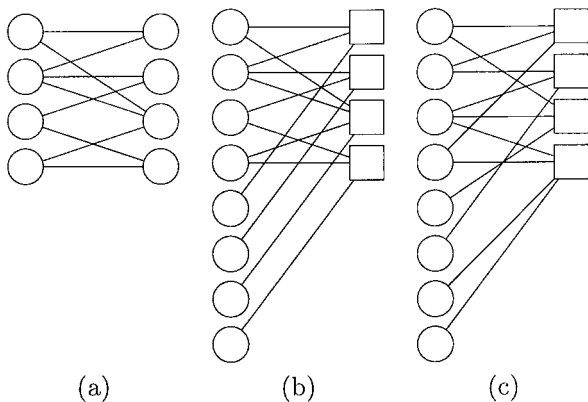


Fig. 9. (a) Example 1-level Simple code for $k = 4, m = 4$. (b) Equivalent Gallager code. (c) Gallager code generated from the same Λ and P as (b).

ally attractive method that yields (near) optimal codes for finite k would be exceptionally useful. This is highlighted by the fact that one-level Simple codes vastly outperform Gallager codes for small k , even though equivalent Gallager codes may be constructed from the Simple codes.

6.3. Question 3

How do the published distributions fare in producing good codes for finite values of k ?

In the next two graphs, we limit our scope to $R = \frac{1}{2}$, as the results for the other two rates are similar. First, we present the performance of the three code generation methods for the three coding methods for small k in Fig. 10. As in the other graphs, the best performing instance for each value of k is plotted.

In all coding methods, the Monte Carlo generation method produces better

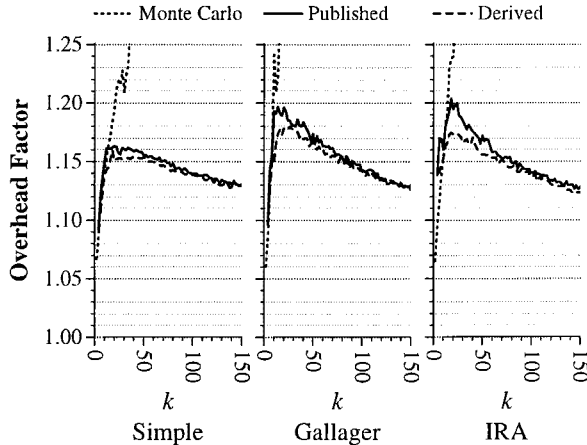


Fig. 10. Performance of various codes for $k \leq 150$ when $\mathcal{R} = \frac{1}{2}$.

codes than the published distributions when k is roughly less than 15. At that point, the exponential number of possible Λ/P combinations drastically reduces the effectiveness of Monte Carlo code generation. From that point until k is in the high double-digits, the performance of the published codes is worse than the derived codes. As k grows past 100, the derived and published codes perform roughly equally. Thus, for small $k (< 100)$, the answer to Question 3 is clearly *inadequately*.

Fig. 11 addresses which published distributions perform well in generating small codes. Each graph plots four curves – the best codes generated from distributions designed for the particular code and rate, the best codes generated from distributions designed for the particular code, but not for the rate, the best codes generated from distributions designed for other codes, and a reference curve showing the best codes from Fig. 10.

Table 2. Distributions that generated the best Simple codes for $k \leq 150$

| Source | Designed for | Rate | Percentage |
|--------|--------------|------|------------|
| [17] | IRA | 2/3 | 46.6% |
| [23] | Gallager | 1/2 | 24.0% |
| [30] | Gallager | 1/3 | 12.3% |
| [26] | Gallager | 2/3 | 11.0% |
| [26] | Gallager | 1/2 | 4.1% |
| [30] | Gallager | 1/2 | 2.0% |

In all three graphs, the worst codes were generated from distributions designed for the particular code, but for a different rate. In both the Gallager and IRA codes, the best codes were generated from distributions designed for the code and rate; and in the Simple codes, the best codes were clearly derived from distributions designed for *other* codes. Probing further, Table 2 shows the breakdown of which

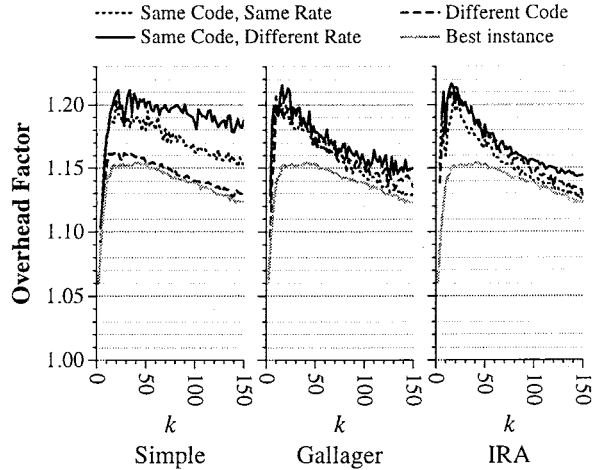


Fig. 11. Performance of published distributions for $k \leq 150$ when $\mathcal{R} = \frac{1}{2}$.

distributions produced the best Simple codes. The significance of this is none other than the fact that the derivation of good Simple codes for small k is clearly not well understood at this point.

For large k , we plot the best published and derived codes for all rates and coding methods in Fig. 12. Note that in each graph, the y-axis has a different scale. There are several interesting features of these graphs. In the middle graph, where $\mathcal{R} = \frac{1}{2}$, the published distributions perform best relative to the derived distributions. This is not surprising, since the bulk of the published distributions (46 of the 80) are for $\mathcal{R} = \frac{1}{2}$. For $\mathcal{R} = \frac{2}{3}$, all three coding methods perform similarly in their best instances. For $\mathcal{R} = \frac{1}{3}$, it is not surprising that the published distributions fare poorly in relation to the derived distributions, since only 10 of the 80 published distributions are for $\mathcal{R} = \frac{1}{3}$, and these are only for Gallager codes. It is interesting that given this fact, the derived IRA codes significantly outperform the others. It is also interesting that the published IRA codes for $\mathcal{R} = \frac{2}{3}$ perform so poorly in comparison to the derived codes.

As in the results on small k , (Fig. 2), in analyzing which distributions produce good graphs for large k , we find that for IRA and Gallager codes, the best codes are produced by distributions designed specifically for the code and rate. For Simple codes, the best codes are produced by distributions for *other* codes. We omit the data here for brevity. It may be obtained in [20].

In summary, for large k , our answer to Question 3 has to be that the published distributions perform poorly in relation to the derived codes. To explore this point further, we decided to experiment with deriving IRA codes without *any* basis in the published codes. Fig. 13 displays the results. In this graph, the dashed lines show our original best IRA codes, derived from the published distributions. The solid lines show the performance of codes derived solely from the Monte Carlo codes for

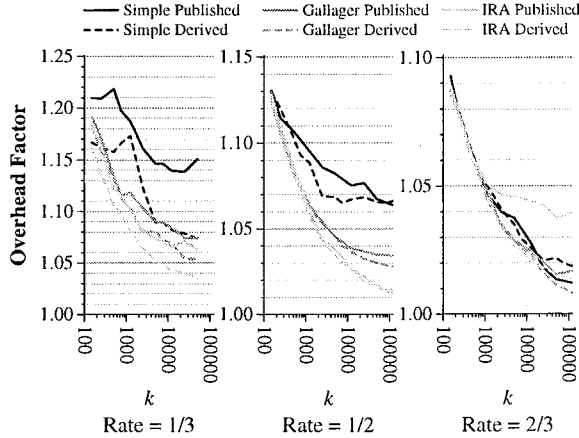
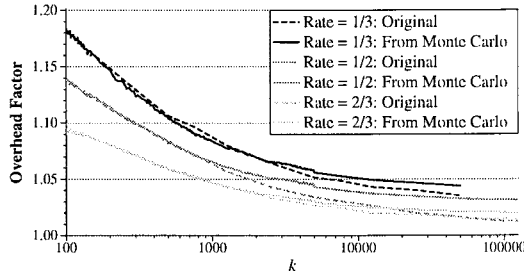

 Fig. 12. Performance of all codes and rates for large k .


Fig. 13. Derived IRA Codes.

even $k \leq 26$. This is just 13 starting points for each rate. To help in the derivation, we generate codes for the following values of k :

- Even values less than 150
- Multiples of 10 less than 1000
- Multiples of 50 less than 5000
- 7500, 12500, 25000, 50000, 75000 and 125000

The generations proceeded on 90 of our machines for 11 days. For all rates, when $k < 1000$, the codes derived from published distributions perform no better than the ones derived from the 13 starting points. As k progresses higher, the codes derived from published distributions outperform the others, although by less than 0.01 for $\mathcal{R} = \frac{1}{3}$ and $\frac{2}{3}$.

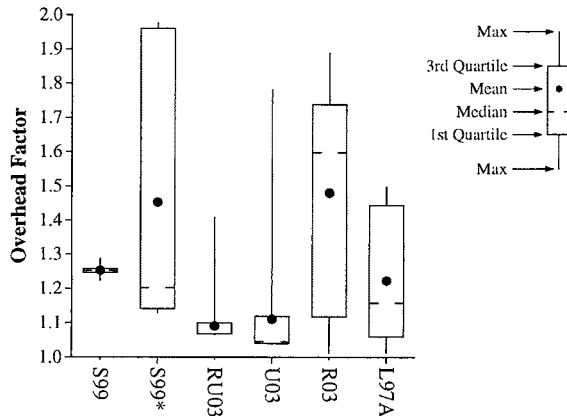
6.4. *Question 4*

Is there a great deal of random variation in code generation from a given probability distribution?

Obviously, this depends on the distribution, and how the distribution is utilized. In Table 3, we select six probability distributions in order to test their variation in code generation. For each of the distributions, we generated over 1000 random codes for $k = 125,000$, and present a digest of the results in Fig. 14. For each distribution we draw a Tukey plot [29], which shows the quartiles for the data and its mean.

Table 3. Range of code generation for given probability distributions.

| Source | Code | Rate Designed | Rate Used | Λ range | P range |
|--------|----------|---------------|-----------|-----------------|-----------|
| S99 | Gallager | 2/3 | 2/3 | 2 | 6 |
| S99* | Gallager | 2/3 | 1/2 | 2 | 6 |
| RU03 | Gallager | 1/2 | 1/2 | 2-13 | 7 |
| U03 | Gallager | 1/2 | 1/2 | 2-100 | 10-11 |
| R03 | IRA | 1/2 | 1/2 | 2-100 | 8 |
| L97A | Simple | 1/2 | 2/3 | 3-1M | 11-30K |

Fig. 14. The variation in code generation for six selected distributions, $k = 125,000$.

The first distribution, S99, from [26], is for a regular graph, where the left nodes each have two outgoing edges, and the right nodes have six incoming edges. As such, we expect little random deviation, which is borne out by the experiments. (We do expect some, because of the random nature of graph generation and of the downloading simulation).

S99* uses the same distribution, but for a different rate. As described in Section , when the total number of edges generated by the left and right nodes do not match, edges are added to or subtracted from random nodes until they do match. Thus, even a regular distribution such as this one, when employed for the wrong rate as in this instance, can generate a wide variety of graphs. It is interesting that this distribution produces better codes for the wrong rate, both in the best and

median case, than the rate for which it is developed. It is also interesting that this regular graph, which theoretically should achieve an asymptotic overhead factor of $\frac{1}{0.68731} = 1.45$ for $\mathcal{R} = \frac{2}{3}$ [26], in actuality achieves a far better one for both rates.

The next two distributions, RU03 and U03, are for Gallager graphs with rate $\frac{1}{2}$. RU03 is *right regular*, meaning all right-hand nodes have the same number of incoming edges, which is a desirable property, because it simplifies code analysis and distribution generation [26,23]. U03 is nearly right regular. Both distributions generate codes with a large spread in performance; however, both have the desirable quality that their medians are very close to their minimum values. In other words, one does not have to generate many codes to get one that performs optimally or near optimally.

The next distribution, for IRA graphs, is also right regular, but has far less desirable generation properties, as it has a very large range of overhead factors, and its median is extremely high. The last distribution, for two-level Simple codes, is one whose nodes have an exceptionally large range of cardinalities – over a million for left nodes (although with $k = 125,000$, the range is reduced to 32,769), and over 30,000 for right nodes. Interestingly, though, its range of overhead factors is less than R03, although it is still a large range.

While more distributions can be displayed, the bottom line remains the same — some distributions yield good codes with only a few iterations of code generation. Others require a longer time to generate good codes. Clearly, one must generate multiple instances of codes to find one that performs well for given values of k and m .

Challenge To The Community: Besides asymptotic performance, some measure of how quickly a distribution yields good codes in practice should be developed. While distributions such as R03 for IRA graphs and L97A for Simple graphs do produce excellent codes, they only do so in relatively rare cases, and thus are difficult to utilize.

6.5. Question 5

What effect does cascading have on Simple codes?

We tested up to six levels of cascading for each value of k and m . Fig. 15 plots the best level for each data point. Drawing general conclusions from Fig. 15 is difficult, but it appears that the optimal number of levels increases as $k \rightarrow \infty$. This increase is very gradual, and appears to be slower for $\mathcal{R} = \frac{1}{2}$ than for $\mathcal{R} = \frac{1}{3}$. The plot for $\mathcal{R} = \frac{1}{3}$ does not exhibit this trend — instead, for $14 \leq k \leq 80$, there are times when 2-Level codes outperform 1-Level codes. This is the point where the codes perform the worst (see Fig. 6), which may be significant. The failure of the codes for $\mathcal{R} = \frac{1}{3}$ to follow the trend of the others may also be due to the fact that we have no published probability distributions designed for this rate.

Finally, we suspect that to generate better multi-level Simple codes, each level may need its own distinct probability distribution. We did not test this hypothesis due to time constraints.

Challenge to the Community: Since the major researchers on Simple codes are focusing on corporate research, there is less theoretical research on them than the

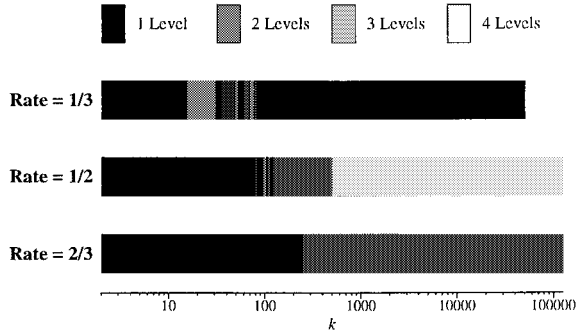


Fig. 15. Best performing Simple graphs, broken down by number of levels.

others. Regardless, understanding the properties of these codes and the implications of cascading behavior would be a welcome addition to the body of knowledge of these codes.

7. Conclusion

This paper has performed a practical exploration of the utility of LDPC codes for wide-area network storage applications. While the asymptotic properties of these codes have been well studied, we have attempted to illuminate their performance for finite systems by addressing five questions, whose answers we summarize below:

Question 1: The overhead factor of LDPC codes, while asymptotically approaching 1, reaches its maximum value when k is in the range of 10 to 100. This maximum value increases as the rate decreases, and may be roughly summarized as 1.20 for $\mathcal{R} = \frac{1}{3}$, 1.15 for $\mathcal{R} = \frac{1}{2}$ and 1.10 for $\mathcal{R} = \frac{2}{3}$.

Question 2: The three types of codes perform differently. Simple codes perform the best for $k < 100$. IRA codes perform the best for $k \geq 100$.

Question 3: Codes derived adaptively from other codes perform better than those derived from published distributions. Simple codes in particular do not perform well from the distributions designed for them.

Question 4: Some distributions produce codes that vary widely in performance. Others produce codes that are more consistent with one another. Concepts like right-regularity do not appear to make a difference.

Question 5: Cascading has an increasing effect on Simple codes as k grows. The different levels may require different probability distributions.

Bottom Line: From these questions, we can draw the following bottom line conclusions:

- While Λ and P suffice for deriving codes with asymptotically good performance, their use as generators of finite codes and indicators of finite code performance is lacking. In the course of this research, we have compiled a mass of codes which perform well, but these have come about by brute force Monte

Carlo techniques. The theoretical community is challenged to derive more effective techniques for generating finite codes, and the experimental community is challenged to explore other heuristics such as pseudo-codewords, stopping sets, trapping sets and girth as guideposts to generating good finite-length codes.

- Clearly, LDPC codes, even suboptimal ones, are very important alternatives to Reed-Solomon codes. A more thorough analysis comparing the performance of these two types of codes needs to be performed, with the goal of providing storage system users with recommendations for the optimal coding technique, value of k , and value of m , given their system's performance and failure parameters.

One limitation of the LDPC codes in this paper is that they have not been designed to adjust to different rates. It is easy to envision a situation where a file already broken into k blocks is spread among $k+m$ storage servers, and then m' new servers are added to the system. If the coding method that stores the original $k+m$ blocks can adapt efficiently to a rate of $\frac{k}{k+m+m'}$, then adding new coding blocks to the system is a straightforward and efficient operation. However, if the coding technique must be altered to accommodate the new rate, then old coding blocks must be discarded, and new ones calculated in their place, which will be inefficient. Reed-Solomon codes have the feature that they adapt to any rate, although they are very slow compared to LDPC, especially when k is large [21,7]. apply. New codes called LT codes and Raptor codes, that adapt to any rate with optimal asymptotic performance have been developed by Luby and Shokrollahi [14,25]. It is a subject of future work to perform a practical analysis of these codes.

Acknowledgements

This material is based upon work supported by the National Science Foundation under grants CNS-0615221, ACI-0204007 and ANI-0222945 and the Department of Energy under grant DE-FC02-01ER25465.

References

- [1] M. S. Allen and R. Wolski. The Livny and Plank-Beck Problems: Studies in data movement on the computational grid. In *SC2003*, Phoenix, November 2003.
- [2] N. Alon, J. W. Spencer, and P. Erdos. *The Probabilistic Method*. John Wiley & Sons, New York, 1992.
- [3] W. A. Burkhard and J. Menon. Disk array storage system reliability. In *23rd International Symposium on Fault-Tolerant Computing*, pages 432–441, Toulouse, France, June 1993.
- [4] J. Byers, M. Luby, M. Mitzenmacher, and A. Rege. A digital fountain approach to reliable distribution of bulk data. In *ACM SIGCOMM '98*, pages 56–67, Vancouver, August 1998.
- [5] J. W. Byers, M. Luby, and M. Mitzenmacher. Accessing multiple mirror sites in parallel: Using tornado codes to speed up downloads. In *IEEE INFOCOM*, pages 275–283, New York, NY, March 1999.
- [6] P. M. Chen, E. K. Lee, G. A. Gibson, R. H. Katz, and D. A. Patterson. RAID: High-performance, reliable secondary storage. *ACM Computing Surveys*, 26(2):145–185, June

- 1994.
- [7] R. L. Collins and J. S. Plank. Assessing the performance of erasure codes in the wide-area. In *DSN-05: International Conference on Dependable Systems and Networks*, Yokohama, Japan, 2005. IEEE.
 - [8] C. Di, D. Proietti, I. E. Telatar, T. J. Richardson, and R. L. Urbanke. Finite-length analysis of low-density parity-check codes on the binary erasure channel. *IEEE Transactions on Information Theory*, 48:1570–1579, June 2002.
 - [9] Digital Fountain, Inc. Next generation data transfer: the meta-content revolution. A Digital Fountain White Paper, www.digitalfountain.com, 2002.
 - [10] R. G. Gallager. *Low-Density Parity-Check Codes*. MIT Press, Cambridge, MA, 1963.
 - [11] H. Jin, A. Khandekar, and R. McEliece. Irregular repeat-accumulate codes. In *2nd International Symposium on Turbo codes and Related Topics*, Brest, France, September 2000.
 - [12] J. Kubiatiowicz, D. Bindel, Y. Chen, P. Eaton, D. Geels, R. Gummadi, S. Rhea, H. Weatherspoon, W. Weimer, C. Wells, and B. Zhao. Oceanstore: An architecture for global-scale persistent storage. In *Proceedings of ACM ASPLOS*, pages 190–201, Cambridge, MA, November 2000. ACM.
 - [13] W. Litwin and T. Schwarz. Lh*rs: a high-availability scalable distributed data structure using Reed Solomon codes. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, pages 237–248, 2000.
 - [14] M. Luby. LT codes. In *IEEE Symposium on Foundations of Computer Science*, 2002.
 - [15] M. Luby, M. Mitzenmacher, and A. Shokrollahi. Analysis of random processes via and-or tree evaluation. In *9th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 364–373, San Francisco, CA, January 1998. ACM.
 - [16] M. Luby, M. Mitzenmacher, A. Shokrollahi, D. Spielman, and V. Stemann. Practical loss-resilient codes. In *29th Annual ACM Symposium on Theory of Computing*, pages 150–159, El Paso, TX, 1997. ACM.
 - [17] R. J. McEliece. Achieving the Shannon Limit: A progress report. Plenary Talk, 38th Allerton Conference, October 2000.
 - [18] J. S. Plank. Improving the performance of coordinated checkpointers on networks of workstations using RAID techniques. In *15th Symposium on Reliable Distributed Systems*, pages 76–85, October 1996.
 - [19] J. S. Plank, A. L. Buchsbaum, R. L. Collins, and M. G. Thomason. Small parity-check erasure codes - exploration and observations. In *DSN-05: International Conference on Dependable Systems and Networks*, Yokohama, Japan, 2005. IEEE.
 - [20] J. S. Plank and M. G. Thomason. On the practical use of LDPC erasure codes for distributed storage applications. Technical Report CS-03-510, University of Tennessee, September 2003.
 - [21] J. S. Plank and M. G. Thomason. A practical analysis of low-density parity-check erasure codes for wide-area storage applications. In *DSN-2004: The International Conference on Dependable Systems and Networks*, pages 115–124, Florence, Italy, June 2004. IEEE.
 - [22] S. Rhea, C. Wells, P. Eaton, D. Geels, B. Zhao, H. Weatherspoon, and J. Kubiatiowicz. Maintenance-free global data storage. *IEEE Internet Computing*, 5(5):40–49, 2001.
 - [23] T. Richardson and R. Urbanke. Modern coding theory. Draft from lthcwww.epfl.ch/papers/ics.ps, August 2003.
 - [24] A. Roumy, S. Guemghar, G. Caire, and S. Verdu. Design methods for irregular repeat accumulate codes. In *IEEE International Symposium on Information Theory*, Yokohama, Japan, 2003.
 - [25] A. Shokrollahi. Raptor codes. Technical Report DR2003-06-001, Digital Fountain,

- 2003.
- [26] M. A. Shokrollahi. New sequences of linear time erasure codes approaching the channel capacity. In *Proceedings of AAEC-13, Lecture Notes in CS 1719*, pages 65–76, New York, 1999. Springer-Verlag.
 - [27] M. A. Shokrollahi. Codes and graphs. *Lecture Notes in Computer Science*, 1770, 2000.
 - [28] M. A. Shokrollahi and R. Storn. Design of efficient erasure codes with differential evolution. In *IEEE International Symposium on Information Theory*, Sorrento, Italy, 2000.
 - [29] E. R. Tufte. *The Visual Display of Quantitative Information*. Graphics Press, Cheshire, Connecticut, 1983.
 - [30] R. Urbanke *et al.* LdcpOpt - a fast and accurate degree distribution optimizer for LPDC ensembles. <http://lthcwww.epfl.ch/research/ldcpopt/index.php>, 2003.
 - [31] H. Weatherspoon and J. Kubiatowicz. Erasure coding vs. replication: A quantitative comparison. In *First International Workshop on Peer-to-Peer Systems (IPTPS)*, March 2002.
 - [32] S. B. Wicker and S. Kim. *Fundamentals of Codes, Graphs, and Iterative Decoding*. Kluwer Academic Publishers, Norwell, MA, 2003.
 - [33] Z. Zhang and Q. Lian. Reperasure: Replication protocol using erasure-code in peer-to-peer storage network. In *21st IEEE Symposium on Reliable Distributed Systems (SRDS'02)*, pages 330–339, October 2002.