

An Efficient XOR-Scheduling Algorithm for Erasure Codes Encoding

Jianqiang Luo, Lihao Xu

Department of Computer Science, Wayne State University
5143 Cass Avenue, Detroit, MI, 48202, (313) 577-0730
jianqiang@wayne.edu, lihao@cs.wayne.edu

James S. Plank

Department of Electrical Engineering and Computer Science, University of Tennessee
203 Claxton Complex, Knoxville, TN, 37996-3450, (865) 974-4397
plank@cs.utk.edu

Abstract

In large storage systems, it is crucial to protect data from loss due to failures. Erasure codes lay the foundation of this protection, enabling systems to reconstruct lost data when components fail. Erasure codes can however impose significant performance overhead in two core operations: Encoding, where coding information is calculated from newly written data, and decoding, where data is reconstructed after failures.

This paper focuses on improving the performance of encoding, the more frequent operation. It does so by scheduling the operations of XOR-based erasure codes to optimize their use of cache memory. We call the technique XOR-scheduling and demonstrate how it applies to a wide variety of existing erasure codes. We conduct a performance evaluation of scheduling these codes on a variety of processors and show that XOR-scheduling significantly improves upon the traditional approach. Hence, we believe that XOR-scheduling has great potential to have wide impact in storage systems.

1 Introduction

As the amount of data increases exponentially in large data storage systems, it is crucial to protect data from loss when storage devices fail to work. Recently, both academic and industrial storage systems have addressed this issue by relying on erasure codes to tolerate component failures. Examples include projects such as OceanStore [20], GFS [12], RAIF [18], and RAIN [7], and companies like Network Appliance [22], HP [29], IBM [13], Cleversafe [9] and Allmydata [2], employing erasure codes such as RDP [10], B Code [31] and Reed-Solomon Codes [27, 6].

In an erasure coded system, a total of $n = k + m$ disks are employed, of which k hold data and m hold coding information. The act of *encoding* calculates the coding information from the data, and *decoding* reconstructs the data from surviving disks following one or more failures. Storage systems typically employ *Maximum Distance Separable (MDS)* codes [5], which ensure that the data can always be reconstructed as long as there are at least k disks that survive the failures.

Encoding is an operation performed constantly as new data is written to the system. Therefore, its performance overhead is crucial to overall system performance. There are two classes of MDS codes – *Reed-Solomon* codes [27] and *XOR* codes (e.g. RDP [10] and X Code [30]). Although there are storage systems based on both types of codes, the XOR codes outperform the others significantly [25] and form the basis of most recent storage systems [7, 9, 18].

This paper addresses the issue of optimizing the encoding performance of XOR codes. XOR codes are described by equations that specify how the coding information is calculated from the data, and most implementations are constructed directly and naively from this specification. However, the order of XOR operations is flexible and can impact the cache behavior of the codes significantly. Using this observation, we detail an algorithm for performing these operations with their cache behavior in mind, and give a comprehensive demonstration of how this improves performance on a variety of XOR codes and processing environments. Additionally, we assess the impact of encoding performance on the overall performance of a storage system, considering the entire memory hierarchy from cache to DRAM to disk.

The bottom line is that our algorithm improves the raw performance of encoding by 21 to 47 percent on various machine architectures. Moreover, the improvement applies to

all XOR codes and is therefore not code specific. Considering the system as a whole, the improvement is mitigated by the disparity in performance between the CPU and the disk and is less drastic — from 1.6 to 8.3 percent. However, considering the importance and constant use of storage systems, this improvement is significant as well. Hence, our algorithm will be very useful for practical storage systems.

2 Nomenclature and Erasure Codes

A storage system is composed of an **array** of n disks, each of which is the same size. Of these n disks, k of them hold **data** and the remaining m hold **coding** information, often termed **parity**, which is calculated from the data. We label the data disks D_0, \dots, D_{k-1} and the coding disks C_0, \dots, C_{m-1} . A typical system is pictured in Figure 1.

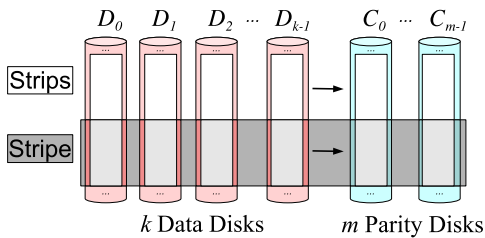


Figure 1. A typical storage system with erasure coding.

When encoding, one partitions each disk into **strips** of a fixed size. Each coding strip is encoded using one strip from each data disk, and the collection of $k + m$ strips that encode together is called a **stripe**. Thus, as in Figure 1, one may view each disk as a collection of strips, and one may view the entire system as a collection of stripes. Note that stripes are each encoded independently, and therefore if one desires to rotate the data and parity among the n disks for load balancing, one may do so by switching the disks' identities for each stripe.

Let us focus on a single stripe. Each XOR code has a parameter w , often termed the *word size* that further defines the code. This parameter is typically constrained by k and by the code. For example, for RDP [10], w must be less than or equal to k , and $w + 1$ must be a prime number. For the X Code and Liberation codes [30, 24], w must be a prime number less than or equal to k .

Each strip is partitioned into exactly w contiguous regions of bytes, called *packets*, labeled $D_{i,0}, \dots, D_{i,w-1}$ and $C_{j,0}, \dots, C_{j,w-1}$ for data and coding drives D_i and C_j respectively. Each packet is the same size, called the *packet size*. Therefore, strip sizes are defined by the product of w

and the packet size. An example of a stripe where $k = 4$, $m = 2$ and $w = 4$ is displayed in Figure 2.

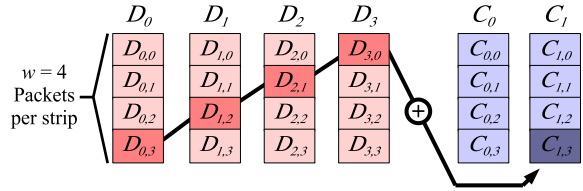


Figure 2. An example of one stripe where $k = 4$, $m = 2$ and $w = 4$.

For the purposes of defining a code, a packet size of one bit is convenient – coding bits are defined to be the XOR of collections of data bits. For example, in Figure 2, when the packet size is one, we can define the bit $C_{1,3}$ as being the XOR of bits $D_{3,0}, D_{2,1}, D_{1,2}$ and $D_{0,3}$, as it is in RDP coding [10]. However, for real implementations, packet sizes should be at least as big as a machine word, since CPU's can XOR two words in one machine instruction. Moreover, to improve cache behavior, larger packet sizes are often preferable [25].

Codes are defined by specifying how each coding packet is constructed from the data packets. This may be done by listing equations, as in RDP [10], EVENODD [3] and X Code [30], or it may be done by employing a Generator Matrix. To describe a code with a Generator Matrix, let us assume that the packet size is one bit. Therefore each data and coding strip is a w -bit word and that their concatenation, which we label $(D|C)$ is a wn -bit word called the *codeword*. The Generator Matrix G has wk rows and wn columns and a specific format: $G = (I|H)$, where I is a $wk \times wk$ identity matrix. Encoding adheres to the following equation:

$$(D|C) = D * G = D * (I|H)$$

Thus, specifying the matrix H is sufficient to specify the code. Codes such as Liberation codes [24], Blaum-Roth codes [4] and Cauchy Reed-Solomon codes [6] are specified in this manner.

Regardless of the specification technique, XOR codes boil down to lists of equations that construct coding packets as the XOR of collections of data packets. To be efficient, the total number of XOR operations should be minimized. Some codes, like RDP and the X Code achieve a lower bound of $(k - 1)$ XOR operations per coding word, while others like Liberation codes, EVENODD and the STAR code are just above this lower bound.

2.1 Erasure Codes

We do not attempt to summarize all research concerning XOR codes. However, we detail the codes that are relevant to this paper. We only consider MDS codes. Cauchy Reed-Solomon codes [6] are general-purpose XOR codes that may be defined for any values of k and m . The word size w is constrained such that $w \geq 2^n$, and the resulting Generator Matrix is typically dense, resulting in a large number of XOR operations for encoding. Plank and Xu describe how to produce sparser matrices [26], and these represent the best performing general-purpose erasure codes.

When m is constrained to equal two, the storage system is a RAID-6 system, and the two coding drives are typically labeled P and Q . There are many XOR codes designed for RAID-6, and these outperform Cauchy Reed-Solomon codes significantly. As stated above, RDP [10] achieves optimal encoding performance, and Liberation [24], Blaum-Roth [4] and EVENODD coding [3] are slightly worse. The STAR code extends EVENODD coding for $m = 3$ and like EVENODD greatly outperforms Cauchy Reed-Solomon coding. The X Code [30] and B Codes [31] are RAID-6 codes that also achieve optimal encoding performance, but require the coding information to be distributed evenly across all $(k + m)$ drives, and are therefore less flexible than the others.

Both Huang [16] and Hafner [14] provide techniques for grouping common XOR operations in certain codes to reduce their number. These techniques are especially effective for decoding Liberation and Blaum-Roth codes. In contrast to this work, we do not improve performance by reducing the number of XOR operations, but instead by improving how the order of XOR operations affects the cache.

Finally, in 2007, Plank released an open source erasure coding library called Jerasure [23] which implements both Reed-Solomon and XOR erasure codes. The XOR codes must use a Generator Matrix specification, with Cauchy Reed-Solomon, Liberation and Blaum-Roth codes included as basic codes. The XOR reduction technique of Hafner [14] is included to improve the performance of decoding. The library is implemented in C and has demonstrated excellent performance compared to other open-source erasure coding implementations [25].

3 CPU Cache

All modern processors have at least one CPU cache that lies between CPU and main memory. Its purpose is to bridge the performance gap between fast CPU's and relatively slow main memories [21]. Caches store recently referenced data, and when working effectively, reduce the number of accesses from CPU to main memory, each of which takes several instruction cycles and stalls the CPU.

When an access to a piece of data is satisfied by the cache, it is called a cache hit; otherwise, it is called a cache miss. Many algorithms that optimize performance do so by reducing the number of cache misses.

There are three types of cache miss: compulsory miss, capacity miss, and conflict miss [15]. Compulsory misses happen when a piece of data is accessed for the first time. Capacity misses occur because of LRU replacement policies of limited-size caches, and conflict misses occur in A-way associative caches when two pieces of data map to the same cache location.

To reduce cache misses, an algorithm needs to have enough locality in time, space, or both. Temporal locality is when the time period between two consecutive accesses to the same data is very short. Spatial locality is when the space difference between the data in a series of accesses is very small. Good temporal locality reduces capacity misses, and good spatial locality reduces both compulsory and conflict misses.

In section 2 above, we mention that large packet sizes are desirable. This is to reduce compulsory misses: when the first byte of a packet is read into the cache, its following bytes are also read into cache because of standard cache prefetch mechanisms. Then, when these bytes are used sequentially, their compulsory misses are avoided.

In this paper, we propose using a well-known optimization technology called *loop fusion* [1, 19, 21, 28] to improve temporal locality when performing XOR operations. In a piece of code, if there are two consecutive loops which access some shared data, it is better to combine these two loops into one loop to improve temporal locality. The reason is when the shared data is accessed in two consecutive loops, it might be moved out from CPU cache after the first loop is finished. Then in the second loop, the access to the shared data will incur cache misses. By combining the two loops, these cache misses are avoided. Employing loop fusion technology in our XOR-scheduling algorithm is the reason why our algorithm is faster than the traditional one.

4 XOR-Scheduling Algorithms

We motivate our scheduling algorithm with a toy example erasure code for $k = 2, m = 2$. In this code, $w = 1$, and our code is not an MDS code (since it is simply an example). We specify the code by assuming our packet size is one bit and listing the equations for the coding bits:

$$\begin{aligned} C_0 &= D_0 + D_1 \\ C_1 &= D_0 + D_1 \end{aligned}$$

Alternatively, we can specify our code with the following Generator Matrix:

$$G = (I|H) = \left(\begin{array}{cc|cc} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{array} \right)$$

Now, let us assume that we are implementing this code, and our packet size is two machine words. Let the two words on data disk D_i be $d_{i,0}$ and $d_{i,1}$, and let the coding words be labeled similarly. The encoding of the entire system is straightforward, and shown in Table 1.

$$\begin{array}{l} \overline{c_{0,0} = d_{0,0} + d_{1,0}} \\ \overline{c_{0,1} = d_{0,1} + d_{1,1}} \\ \overline{c_{1,0} = d_{0,0} + d_{1,0}} \\ \overline{c_{1,1} = d_{0,1} + d_{1,1}} \end{array}$$

Table 1. Encoding the Toy Example when the packet size equals two.

4.1 Traditional XOR-scheduling Algorithm

Traditional XOR-scheduling performs the encoding in the order specified by the encoding equations or Generator Matrix / Data Vector product. In our example, the coding words $c_{0,0}$, $c_{0,1}$, $c_{1,0}$ and $c_{1,1}$ are calculated in that order. To ease the explanation, we are going to assume that each coding word starts with a value of zero, and its calculation requires two XORs to update it. It will be clear how to remove this assumption later.

The schedule of XOR operations generated by the traditional algorithm in our toy example is listed in Table 2. In the table, each row contains two columns. The second column is the XOR operation, and the first column is the ID for that XOR operation. Throughout our examples, the same XOR operation may appear in different positions in different scheduling algorithms, but its ID will remain constant.

The first characteristic of the traditional algorithm is that it processes coding packets one by one, to completion. In our toy erasure code, there are two coding packets: C_0 and C_1 . In the traditional algorithm, the contents of C_1 will not be calculated until all of C_0 is finished. The second characteristic of the traditional algorithm is that when processing a coding packet, its associated data words are accessed in a strict order. For example, in Table 2, the execution sequence for coding packet C_0 is: $\{d_{0,0}, d_{0,1}, d_{1,0}, d_{1,1}\}$. In other words, when calculating C_0 , none of D_1 is accessed until all of D_0 has been accessed.

A pseudocode description of the traditional XOR-scheduling algorithm is shown in Algorithm 1.

ID	XOR
1	$c_{0,0} += d_{0,0}$
2	$c_{0,1} += d_{0,1}$
3	$c_{0,0} += d_{1,0}$
4	$c_{0,1} += d_{1,1}$
5	$c_{1,0} += d_{0,0}$
6	$c_{1,1} += d_{0,1}$
7	$c_{1,0} += d_{1,0}$
8	$c_{1,1} += d_{1,1}$

Table 2. The Result of Traditional Algorithm on the Toy Example

Algorithm 1 The Traditional Scheduling Algorithm

```

procedure XOR_scheduling()
1: for each coding packet  $C_j$  do
2:   for each data packet  $D_i$  in the calculation of  $C_j$  do
3:     for  $x=0;x<\text{packet size};x++$  do
4:       output  $c_{j,x} += d_{i,x}$ ;
5:     end for
6:   end for
7: end for

```

Algorithm 1 shows that the traditional scheduling algorithm has good spatial locality because it accesses the words in a packet sequentially.

4.2 A New XOR-scheduling Algorithm

The traditional XOR-scheduling algorithm follows the intuitive idea that coding words should be produced one by one. Instead, we can reorder the schedule so that it consumes *data* words one by one. Our new XOR-scheduling algorithm is based on this idea, and its characteristics are as follows:

1. The order of XOR operations is guided by the order of data words instead of coding words.
2. Each data word is used for all of its coding calculations before moving onto the next data word in the same packet.

The result of the new XOR-scheduling algorithm for the toy example is shown in Table 3.

The first characteristic requires that data packets are processed sequentially. As a result, in Table 3, all of the equations involving D_0 appear before those involving D_1 . The second characteristic requires that each data word is used for all coding calculations before moving onto the next data word in the same packet. This can be seen in Table 3

ID	XOR
1	$c_{0,0+} = d_{0,0}$
5	$c_{1,0+} = d_{0,0}$
2	$c_{0,1+} = d_{0,1}$
6	$c_{1,1+} = d_{0,1}$
3	$c_{0,0+} = d_{1,0}$
7	$c_{1,0+} = d_{1,0}$
4	$c_{0,1+} = d_{1,1}$
8	$c_{1,1+} = d_{1,1}$

Table 3. The Result of New Algorithm on the Toy Example

where $d_{0,0}$ is used to calculate both $c_{0,0}$ and $c_{1,0}$ before $d_{0,1}$ is touched. A pseudocode description of this algorithm is shown in Algorithm 2.

Algorithm 2 The New Scheduling Algorithm

procedure *XOR_scheduling*()

- 1: **for** each data packet D_i **do**
 - 2: **for** $x=0;x<\text{packet size};x++$ **do**
 - 3: **for** each coding packet C_j that uses D_i **do**
 - 4: output $c_{j,x} += d_{i,x}$;
 - 5: **end for**
 - 6: **end for**
 - 7: **end for**
-

Algorithm 2 differs with Algorithm 1 in the order of their loops. Below are the reasons for why Algorithm 2 has better data locality than Algorithm 1.

1. The innermost loop of Algorithm 2 in line 4 does not change the value of x , so it has good temporal locality of accessing the same $d_{i,x}$. This optimization is loop fusion.
2. The middle loop of Algorithm 2 in line 3 still preserves good spatial locality to $c_{j,x}$. This is because in most reliable storage systems, k will be bigger than m , and most coding words will remain in the cache across iterations of the innermost loop. Thus, since $d_{i,x}$ and $c_{j,x}$ should remain in the cache across an iteration of the inner loop, $d_{i,x+1}$ and $c_{j,x+1}$ should remain there too.

Loop fusion can not be as effectively employed in Algorithm 1 when k is much larger than m . For example, a typical RAID-6 system will have $k \geq 10$, while $m = 2$. Thus, were we to switch the order of the two inner loops in Algorithm 1, we would be more likely to lose spatial benefits, because it is less likely that the data words will remain in the cache across iterations of the inner loop. Although

we do not show the results, we did test loop fusion in Algorithm 1, and it did not improve performance.

We note that loop fusion technology reduces some spatial locality in Algorithm 2. In some environments with large packet sizes, this reduction can cause Algorithm 2 to perform worse than Algorithm 1. However, the peak performance achieved by Algorithm 2 is always better than that of Algorithm 1.

It is worth noting that the loop fusion technology used in our algorithm is a new view of encoding, and is not a simple code transformation that could be implemented by a compiler. In some respects, it is reminiscent of [21], which provides some examples of how manual optimization can greatly improve an algorithm’s performance.

4.3 XOR-Scheduling Implementation Details

Clearly, generating and using schedules that are as detailed as those in Figures 2 and 3 take too much space. Fortunately, it is not necessary, since a list of each data packet’s associated coding packets may be constructed simply from the packet’s row of the Generator Matrix. Once that list is generated, it may be used for every word in the packet.

Another important detail of our scheduling algorithm is how to initialize coding packets. We cannot simply initialize them to zero, since that increases the number of XOR operations. Instead, we copy the first data packets that are used for each coding packet, instead of XOR-ing them. Implementationally, this is simple in Algorithm 1, where we may identify the first data packets at the beginning of Line 2. In Algorithm 2, it is more difficult, since we need to know at Line 4 whether D_i is C_j ’s first data packet. Fortunately, all the codes that we address have regular structures which makes this determination straightforward.

A final detail involves the codes that have extra information in addition to their Generator Matrices – EVENODD, STAR and RDP. EVENODD employs a temporary packet, S that must be XOR’d with every packet in C_1 . To handle this, both algorithms perform two passes. In the first pass, the data packets are used to calculate C_0 , C_1 and S , and then a final pass XOR’s S into the packets of C_1 . The STAR code, which is an extension to EVENODD, is handled similarly. In RDP, all but one of the packets in C_1 are calculated using packets in C_0 in addition to data packets. For that reason, we also perform two passes in RDP: a first one that calculates C_0 and C_1 without the C_0 packets, and a second one that XOR’s the C_0 packets into C_1 .

5 Performance Evaluation

To test the two XOR-scheduling algorithms, we have conducted experiments that apply both algorithms to

Machine Name	CPU Speed	L1 Cache	L2 Cache	Cores	Bit	CPU Description	Memory
P4	3.0GHz	16KB	1MB	1	32-bit	Pentium 4 (520)	1 GB
Pd	2.8GHz	2*16KB	2*1MB	2	64-bit	Pentium Dual Core (D820)	1 GB
Pc2d	2.1GHz	2*32KB	3MB	2	64-bit	Pentium Core 2 Duo (T8100)	2 GB
Pc2q	2.4GHz	4*32KB	2*4MB	4	64-bit	Pentium Core 2 Quad (Q6600)	2 GB

Table 4. Details of the Test Platforms

many popular XOR-based erasure codes: the Liberation Codes [24], EVENODD [3], RDP [10], the X Code [30] and the STAR [17]. All but the STAR code are RAID-6 codes. Since CPU cache behavior is complicated, we have run our experiments on various machines to get a comprehensive view of our algorithm. There are four platforms used in our experiments, and they range from low-end to high-end. The detailed information on these platforms are shown in Table 4.

All these platforms run Linux, and the 64-bit machines run the 64-bit version of Linux. Since 64-bit machines perform XOR operations on 64-bit words, their encoding performance is as much as a factor of two faster than their 32-bit counterparts [25].

Our code is written in C and compiled using `gcc` with the `-O2` optimization flag set. `-O2` is the common optimization flag for installed packages and for the kernel of Linux [11]. While there may be more intricate compiler optimizations available on different machines, using `-O2` gives us the performance results that will match standard installations. Matching open source implementations, the code runs only on one thread, and thus does not take advantage of multiple cores.

For the Liberation codes, we use the Jerasure open source coding library [23] as the implementation for the traditional scheduling algorithm, as that is exactly what Jerasure implements. For all other codes and for implementing the new scheduling algorithm, we wrote our own implementations, starting with Jerasure as the base code.

5.1 Experiment Setup

We set up an experimental framework to evaluate XOR-scheduling algorithms’ performance. In our framework, all operations are performed in memory with no disk I/O involved. The reason is that employing the disk makes it very difficult to assess encoding performance accurately [25]. We evaluate the impact when disk operations are considered in Section 5.5 below.

The size of total input user data is 300 MB, which we assume will be shared among k data devices and encoded onto m coding devices. Because the input data is too large to be allocated in main memory sequentially, we only al-

locate 20 MB in main memory as a data buffer. We then allocate a coding buffer of the appropriate size, and proceed in phases until all 300 MB have been encoded. To simulate reading the data, we fill the buffer randomly.

We evaluate encoding performance using a metric of *Encoding Bandwidth*, calculated with the equation below:

$$\text{Encoding Bandwidth} = \frac{\text{Encoding Time}}{\text{Total Input Data Size}}$$

We use Encoding Bandwidth so that the performance evaluation is not dependent on the size of the data. It represents how quickly one can turn data into coding information with a given code, algorithm and machine. Since all but one of the codes evaluated are RAID-6 codes, their performance may be evaluated directly. The STAR code will encode with a lower bandwidth, since it creates three drives’ worth of coding information rather than two.

The encoding time is calculated using the `gettimeofday()` system call. Each data point is the average of three test runs.

5.2 Experiments on the Liberation Codes

We first focus on a performance comparison of the Liberation Codes. We do this because it is the one code for which we have an open source implementation. The parameters of our first experiment are $k = 11$, $m = 2$ and $w = 11$, as these represent a typical RAID-6 system. We modify the packet size, as [25] has demonstrated that varying packet sizes can impact performance. Figure 3 shows the performance comparison on machine P4.

Figure 3 allows us to make the following observations:

1. As packet sizes increase from a small value, the performance of both scheduling algorithms improves significantly.
2. For all packet sizes, the new scheduling algorithm achieves much better encoding performance than the traditional one.
3. The peak performance of the new scheduling algorithm is much higher than that of the traditional one.

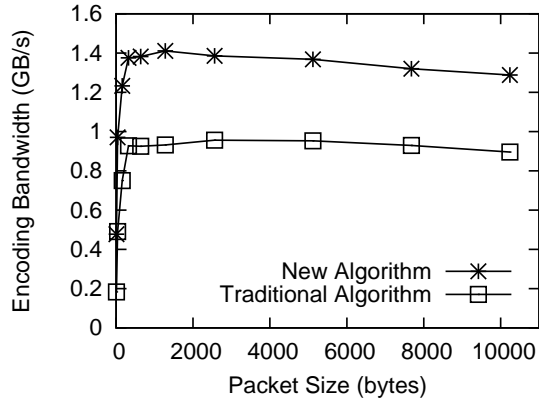


Figure 3. Encoding Performance On P4

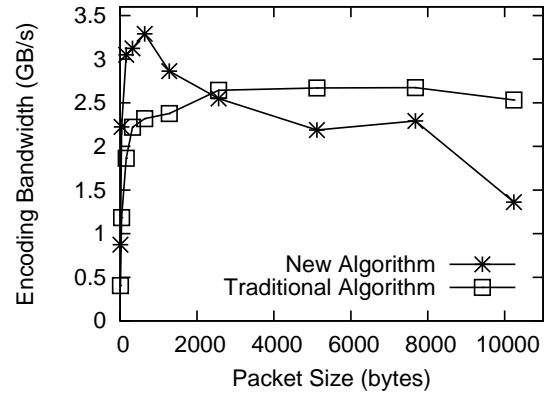


Figure 5. Encoding Performance on Pc2d

4. The new scheduling algorithm achieves its peak encoding performance with a relatively small packet size.

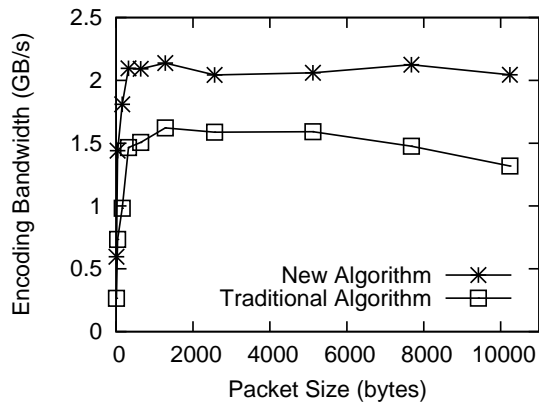


Figure 4. Encoding Performance on Pd

The first observation mirrors the observations in [25] which reflects greater spatial locality with larger packet sizes. The reason for the second and third observations is the new scheduling algorithm has better temporal locality than the traditional algorithm, matching our analysis in Section 4.2. The last observation is important and will recur in the data below. The traditional algorithm generates completed coding words quickly and sequentially. The new algorithm in contrast only completes generating coding data at the end of each stripe. Therefore, we may prefer smaller packets with this algorithm, and it is important that it performs well with relatively small packet sizes.

The results on the Pd machine are displayed in Figure 4. The four above observations on P4 apply also to Pd. The algorithms encode faster though, since Pd is a 64-bit machine.

Figure 5 shows the comparison results on machine Pc2d.

Although the machine has a slower processor than the others, it has larger caches and achieves better peak performance. We also see a new trend on this machine – after reaching peak performance with a rather small packet size, the performance of the new algorithm decreases to a point where it performs worse than the traditional algorithm. In Section 4.2, we discuss how the new scheduling algorithm can reduce spatial locality compared with traditional algorithm. This may be the reason for why performance decreases after obtaining peak performance.

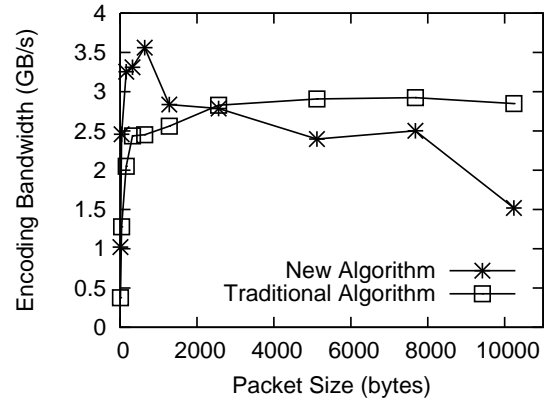


Figure 6. Encoding Performance on Pc2q

Nonetheless, we still see that the two most important observations hold: 1) the peak performance of new scheduling algorithm is much higher than that of traditional one; 2) the peak performance is obtained with small packet size.

Figure 6 shows the results on machine Pc2q. The performance of Pc2q is better than the others, but the trends are very similar to Pc2d, and the observations for Pc2d are also valid for Pc2q.

The peak performance of two scheduling algorithms on four machines are summarized in Table 5. It clearly displays

that the new scheduling algorithm achieves significant performance improvement, from 21% to 47%.

Machine	New	Traditional	Improvement
P4	1.411 GB/s	0.956 GB/s	47%
Pd	2.125	1.621	31%
Pc2d	3.290	2.673	23%
Pc2q	3.560	2.925	21%

Table 5. Comparison of Peak Encoding Performance of Liberation Codes for $k = w = 11$ and $m = 2$.

5.2.1 Modifying k and w

To observe potential sensitivity to the number of data devices and the number of packets per stripe, in Figure 7 we display the encoding performance of the Liberation codes for $k = w = 5$ and $k = w = 17$. In the top row of graphs, k and w equal five, and in the lower row they equal 17. The results mirror the results for $k = 11$, and the observations that held for $k = 11$ hold for the smaller and larger values of k . Thus, the new scheduling algorithm is stable for this code among this range of parameters.

5.3 Experiments on Other Erasure Codes

To test the algorithms on other codes, we tested EVENODD, RDP, X Code and STAR coding in our framework. The various codes impose different constraints on k and w , so we selected values that would match most closely with the Liberation code example. These values are summarized in Table 6.

Code	k	w	m
EVENODD	11	10	2
RDP	10	10	2
X Code	11	11	2
STAR	11	10	3

Table 6. Encoding parameters of various codes

The results are in Figure 8. In terms of peak performance, the codes match expectations. When $k = w$, the Liberation codes’ performance is theoretically identical to EVENODD [24], and this is reflected in the results. RDP and the X Code should encode the fastest, and they do.

STAR’s performance is worse than the others because it encodes an extra coding device.

In terms of the trends, the codes’ performance matches the Liberation codes, and all the observations for the Liberation codes hold for these codes as well. The results show that the algorithm is applicable to a wide variety of codes.

5.4 Encoding Performance of RDP

In [10], Corbett provides encoding performance of RDP Code, and it is the best reported results of RDP Code in the literature. The platform used in paper [10] contains two Pentium 4 CPUs of 2.8GHz of which one CPU is dedicated for RDP encoding. They do not report the cache sizes. The most similar platform in our tests is machine P4, but P4 only has one CPU. The performance comparison of our scheduling algorithm with their reported performance is given in Table 7.

Machine	Ours	Reported in [10]
P4	1.349 GB/s	1.55 GB/s

Table 7. Encoding Performance for RDP Code

Table 7 shows that our encoding performance is a little worse than their reported performance. Without knowing the exact details, the reasons can be manifold. However, a simple conclusion to draw is that with our new scheduling algorithm, our implementation can achieve encoding performance that is comparable to an implementation used in commercial products.

5.5 Data Write Performance

When erasure codes are used in a storage system, the overall performance of data write operations will be determined by a combination of encoding performance and disk write performance. A common view is that the encoding performance has no real impact on the overall performance of data write operations, because encoding performance is an order of magnitude faster than writing to disk. However, it still has a significant effect, shown in Table 8, mainly because of the existence of write buffer in main memory. The unit for all performance is GB/s.

In this table, for each machine we show the encoding performance of the new scheduling algorithm, the encoding performance of the traditional scheduling algorithm, and the disk write performance as measured by running `bf bonnie` [8] on machine P4, which has a 7200 RPM hard disk and runs the `ext2` file system. The data write perfor-

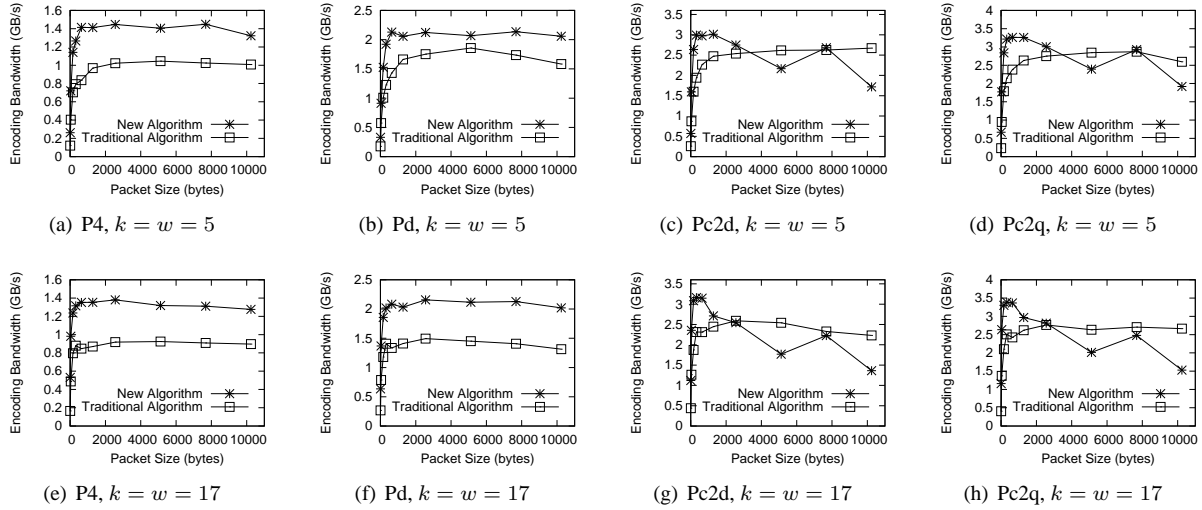


Figure 7. Encoding Performance for Liberation Code, varying k and w from 5 to 17

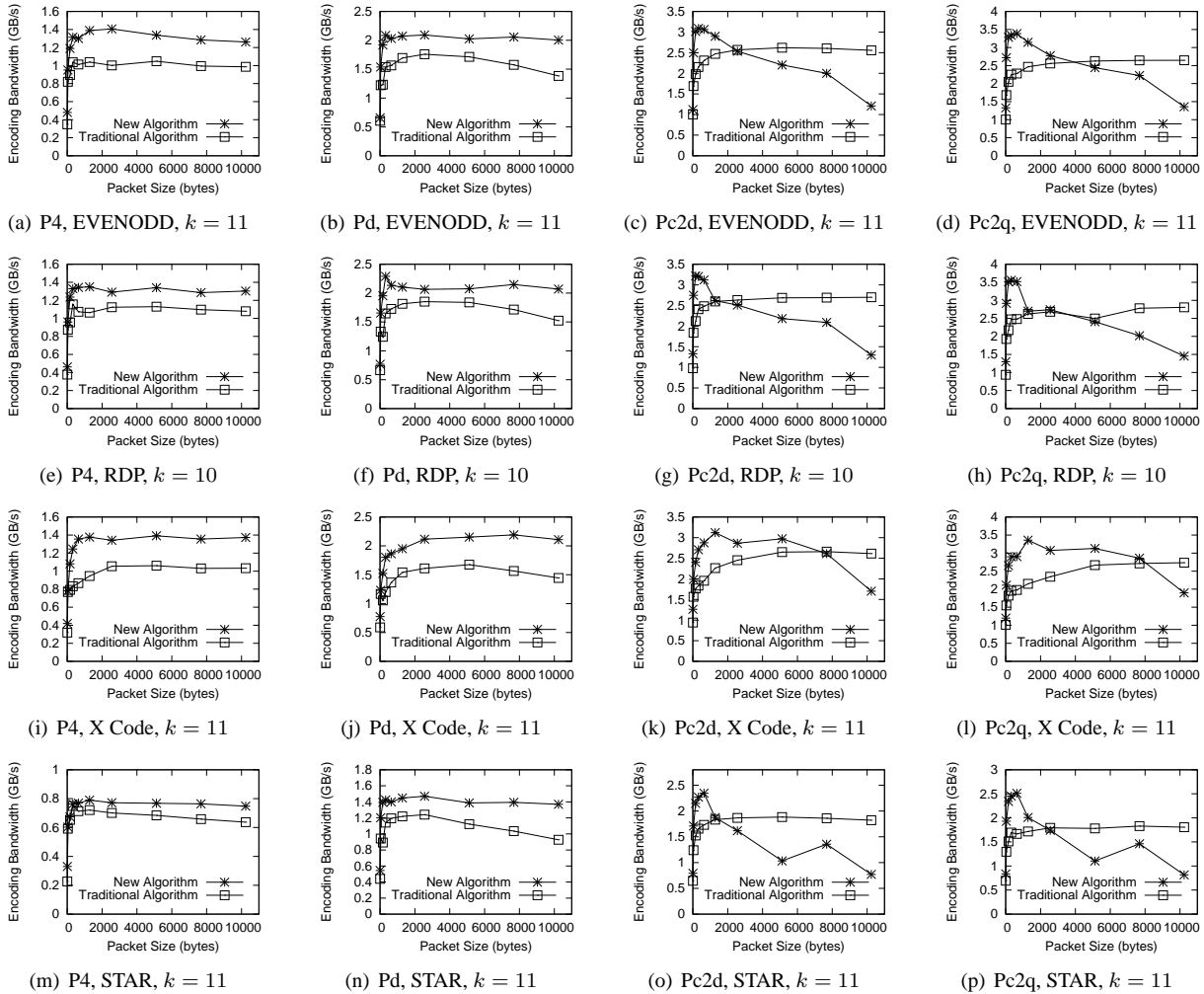


Figure 8. Encoding Performance for Various Erasure Codes

Machine	New	Traditional	Write	Improvement
P4	1.411	0.956	0.296	8.3%
Pd	2.125	1.621	0.296	3.8%
Pc2d	3.29	2.673	0.296	1.9%
Pc2q	3.56	2.925	0.296	1.6%

Table 8. Encoding Performance and Disk Write Performance

mance improvement in the last column is calculated using the straightforward equation below:

$$\text{Performance Improvement} = \frac{\frac{1}{\text{Traditional}} + \frac{1}{\text{Write}}}{\frac{1}{\text{New}} + \frac{1}{\text{Write}}} - 1$$

From the last column of Table 8, we can see that the performance improvement ranges from 1.6% to 8.3% on various platforms. In practical systems, this value can be even higher. The reason is that disk write performance can improve because of larger main memories and faster hard drives, but encoding performance may decrease because of limited cache sizes. If we assume that the encoding performance of two scheduling algorithms decreases by P percent, and disk write performance simultaneously increases by P percent, then we get the performance improvement trend presented in Figure 9. While this is merely speculative, it shows that encoding performance can have further impact on the overall performance of a storage system.

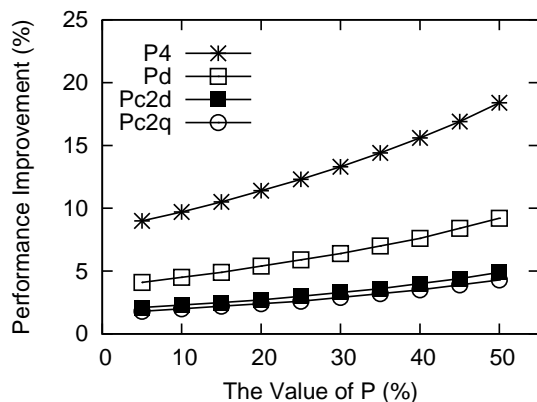


Figure 9. Potential Performance of Storage System

6 Conclusions

This paper proposes a new XOR-scheduling algorithm to improve the encoding performance of XOR-based erasure codes. For these erasure codes, the encoding performance is determined by two primary factors: the number of XOR operations and the cache behavior. This paper proposes a new scheduling algorithm which is able to efficiently utilize CPU cache and thus achieves much better encoding performance than the traditional algorithm. In a performance evaluation on some widely known erasure codes on a variety of platforms, we show that the encoding performance obtained by our scheduling algorithm significantly outperforms that of the traditional algorithm.

Two clear areas of future work are to incorporate our algorithm into **Jerasure** so that those who use this open source tool may take advantage of our work, and to assess the impact of this technique on decoding. Since decoding is a more complex activity than encoding and occurs more rarely, we have not yet performed this assessment. It will be a necessary step to providing a complete evaluation of the scheduling algorithm.

Additional future work is to put our new scheduling algorithm into a real storage system. Although we have analyzed how our scheduling algorithm improves the whole performance of storage systems, the data used in our analysis is synthetic and may not be representative. We plan to implement a reliable storage system [32] and use various scheduling algorithms in it to find how our scheduling algorithm can improve this system's performance. Another future work is to find similar optimization technologies for Reed-Solomon Codes. Although they do not perform as well as the XOR based codes in this paper, Reed-Solomon Codes are used widely not only in storage systems but also in network systems. Thus the performance improvements of our techniques will have wider impact.

References

- [1] R. Allen and K. Kennedy. Optimizing compilers for modern architectures. San Francisco, California, USA, 2001. Morgan Kaufmann.
- [2] Allmydata. Unlimited Online Backup, Storage, and Sharing. 2008.
- [3] M. Blaum, J. Brady, J. Bruck, and J. Menon. EVENODD: An Efficient Scheme for Tolerating Double Disk Failures in RAID Architectures. *IEEE Trans. on Computers*, 44(2):192–202, Feb. 1995.
- [4] M. Blaum and R. M. Roth. New Array Codes for Multiple Phased Burst Correction. *IEEE Trans. on Information Theory*, 39(1), Jan. 1993.
- [5] M. Blaum and R. M. Roth. On lowest density MDS codes. *IEEE Transactions on Information Theory*, 45(1):46–59, January 1999.

- [6] J. Blomer, M. Kalfane, M. Karpinski, R. Karp, M. Luby, and D. Zuckerman. An XOR-based Erasure-Resilient Coding Scheme. *Technical Report TR-95-048, International Computer Science Institute*, August 1995.
- [7] V. Bohossian, C. C. Fan, P. S. LeMahieu, M. D. Riedel, J. Bruck, and L. Xu. Computing in the RAIN: A Reliable Array of Independent Nodes. *IEEE Transaction on Parallel and Distributed Systems*, 12(2):99–114, 2001.
- [8] Tim Bray. <http://www.textuality.com/bonnie>.
- [9] Cleversafe, Inc. Cleversafe Dispersed Storage, Open source code distribution: <http://www.cleversafe.org/downloads>. 2008.
- [10] P. Corbett, B. English, A. Goel, T. Grcanac, S. Kleiman, J. Leong, and S. Sankar. Row-Diagonal Parity for Double Disk Failure Correction. *Proc. of 3rd USENIX Conference on File and Storage Technologies (FAST '04)*, March 2004.
- [11] gentoo wiki. <http://en.gentoo-wiki.com/wiki/cflags>.
- [12] S. Ghemawat, H. Gobioff, and S. T. Leung. The Google File System. *Proc. of the 19th ACM Symposium on Operating Systems Principles (SOSP '03)*, 2003.
- [13] J. L. Hafner. WEAVER Codes: Highly Fault Tolerant Erasure Codes for Storage Systems. *Proceedings of the Fourth USENIX Conference on File and Storage Technologies (FAST '05)*, 2005.
- [14] J. L. Hafner, V. Deenadhayalan, KK Rao, and J. A. Tomlin. Matrix Methods for Lost Data Reconstruction in Erasure Codes. *Proc. of 4th USENIX Conference on File and Storage Technologies (FAST '05)*, March 2005.
- [15] J. L. Hennessy, D. A. Patterson, and D. Goldberg. *Computer Architecture: A Quantitative Approach*, chapter Appendix C. Morgan Kaufmann, May 2002.
- [16] C. Huang, J. Li, and M. Chen. On Optimizing XOR-Based Codes for Fault-Tolerant Storage Applications. *Proc. IEEE Information Theory Workshop (ITW 2007)*, September 2007.
- [17] C. Huang and L. Xu. STAR: An Efficient Coding Scheme for Correcting Triple Storage Node Failures. *Proc. of FAST 2005*, Dec. 2005.
- [18] N. Joukov, A. M. Krishnakumar, C. Patti, A. Rai, S. Satnur, A. Traeger, and E. Zadok. RAIF: Redundant Array of Independent Filesystems. *Proceedings of 24th IEEE Conference on Mass Storage Systems and Technologies (MSST 2007)*, pages 199–212, September 2007.
- [19] M. Kowarschik. An overview of cache optimization techniques and cache-aware numerical algorithms. In *In Algorithms for Memory Hierarchies, volume 2625 of LNCS*, pages 213–232. Springer, 2003.
- [20] J. Kubiawicz, D. Bindel, Y. Chen, S. Czerwinski, P. Eaton, D. Geels, R. Gummadi, S. Rhea, H. Weatherspoon, C. Wells, and B. Zhao. OceanStore: An Architecture for Global-scale Persistent Storage. *Proceedings of the ninth international conference on Architectural support for programming languages and operating systems*, December 2000.
- [21] A. R. Lebeck and D. A. Wood. Cache Profiling and the SPEC Benchmarks: A Case Study. *IEEE Computer*, 27(10), September 1994.
- [22] B. Nisbet. FAS storage systems: Laying the foundation for application availability. *Network Appliance white paper: <http://www.netapp.com/us/library/analyst-reports/ar1056.html>*, February 2008.
- [23] J. S. Plank. Jerasure: A library in C/C++ facilitating erasure coding for storage applications. *Tech. Rep. CS-07-603, University of Tennessee*, September 2007.
- [24] J. S. Plank. The RAID-6 Liberation Codes. *FAST-2008: 6th Usenix Conference on File and Storage Technologies*, February 2008.
- [25] J. S. Plank, J. Luo, C. D. Schuman, L. Xu, and Z. Wilcox-O'Hearn. A Performance Evaluation and Examination of Open-Source Erasure Coding Libraries For Storage. to appear in *Proc. of FAST-2009: 7th Usenix Conference on File and Storage Technologies*, February 2009.
- [26] J. S. Plank and L. Xu. Optimizing Cauchy Reed Solomon Codes for Fault-Tolerant Network Storage Applications. *Proc. of The 5th IEEE International Symposium on Network Computing and Applications (IEEE NCA06)*, July 2006.
- [27] I. S. Reed and G. Solomon. Polynomial Codes over Certain Finite Fields. *Journal of the Society for Industrial and Applied Mathematics*, 8 (10):300–304, 1960.
- [28] J. Warren. A hierarchical basis for reordering transformations. In *POPL'84: Proceedings of the 11th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 272–282. ACM, 1984.
- [29] J. J. Wylie and R. Swaminathan. Determining fault tolerance of xor-based erasure codes efficiently. In *DSN '07: Proceedings of the 37th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, pages 206–215, Washington, DC, USA, 2007. IEEE Computer Society.
- [30] L. Xu. X-Code: MDS Array Codes with Optimal Encoding. *IEEE Trans. on Information Theory*, 45 (1):272–276, Jan. 1999.
- [31] L. Xu and J. Bruck. Low Density MDS Code and Factors of Complete Graphs. *IEEE Transactions on Information Theory*, Sep. 1999.
- [32] Lihao Xu. Hydra: A platform for survivable and secure data storage systems. *Proc. of StorageSS 2005, Fair Fax, Virginia*, Nov. 2005.